



Jonatas dos Santos Grosman

LER: Anotação e classificação automática de entidades e relações

Dissertação de Mestrado

Dissertação apresentada como requisito parcial para obtenção do grau de Mestre pelo Programa de Pós-graduação em Informática da PUC-Rio.

Orientador: Prof. Hélio Côrtes Vieira Lopes

Rio de Janeiro
Abril de 2017



Jonatas dos Santos Grosman

LER: Anotação e classificação automática de entidades e relações

Dissertação apresentada como requisito parcial para obtenção do grau de Mestre pelo Programa de Pós-graduação em Informática da PUC-Rio. Aprovada pela Comissão Examinadora abaixo assinada.

Prof. Hélio Côrtes Vieira Lopes

Orientador

Departamento de Informática – PUC-Rio

Prof^a. Simone Diniz Junqueira Barbosa

Departamento de Informática – PUC-Rio

Prof^a. Maria Cláudia de Freitas

Departamento de Letras – PUC-Rio

Prof. Marcus Vinícius Soledade Poggi de Aragão

Departamento de Informática – PUC-Rio

Prof. Márcio da Silveira Carvalho

Coordenador Setorial do Centro Técnico Científico – PUC-Rio

Rio de Janeiro, 20 de Abril de 2017

Todos os direitos reservados. É proibida a reprodução total ou parcial do trabalho sem autorização da universidade, do autor e do orientador.

Jonatas dos Santos Grosman

Graduou-se em Sistemas de Informação pela Faculdade de Educação Tecnológica do Estado Rio de Janeiro.

Ficha Catalográfica

Grosman, Jonatas dos Santos

LER: Anotação e classificação automática de entidades e relações / Jonatas dos Santos Grosman; orientador: Hélio Côrtes Vieira Lopes. – 2017.

v., 196 f: il. color. ; 30 cm

Dissertação (mestrado) - Pontifícia Universidade Católica do Rio de Janeiro, Departamento de Informática.

Inclui bibliografia

1. Informatica – Teses. 2. Processamento de linguagem natural. 3. Aprendizado automático. 4. Extração de informação. 5. Ontologias. 6. Curadoria de dados. I. Lopes, Hélio Côrtes Vieira. II. Pontifícia Universidade Católica do Rio de Janeiro. Departamento de Informática. III. Título.

Agradecimentos

Começarei pelo clichê.

Palavras não bastam para descrever tamanha gratidão que tenho pelas pessoas que fizeram parte do processo que culminou neste trabalho.

Agora ao menos clichê.

Eu poderia separar em três grupos as pessoas à agradecer, as que me fizeram entrar no mestrado, as que me fizeram continuar nele e, por fim, as que me fizeram concluí-lo. Inicialmente tentei transcrever os agradecimentos segundo essa divisão, mas me peguei em tantas intersecções que decidi escrever na ordem que encontrei algum sentido.

Pois bem, eu poderia começar o agradecimento regredindo muito no tempo, lembrando dos incentivos dos meus professores e colegas de faculdade, já que foi nesse tempo que minha paixão pela computação nasceu, mas só não o farei por intuito de brevidade, preferindo começar pelos meus primeiros passos fora dela. Foi nesse tempo que conheci as pessoas que mais contribuíram para eu entender onde estava me metendo e que rumo tomar. Começo pelo estágio, onde conheci um grande amigo, Lucas Bastos, que a seu jeito me fazia querer saber mais sobre computação. Indo para a minha primeira experiência científica fora da faculdade, no tempo que fiquei no Laboratório Nacional de Computação Científica, conheci pessoas brilhantes como Anderson Menezes, Antônio Tadeu, Artur Ziviani, Bruno Bastos, Bruno Correa, Iuri Malinoski, Thiago Cardozo e Vinicius Moreira. Eles podem nem fazer ideia do quanto tenho a agradecer pelas lições aprendidas, então escrevo aqui para que tenham ciência. Na minha passagem pelo Observatório Nacional, mesmo que breve, já que tive de abandoná-los para poder fazer o mestrado, conheci pessoas como Pedro Rocha e Selma Junqueira, que me ajudaram com o empurrão que faltava para querer enveredar pelo caminho que tomei. Tenho muito a agradecer (muito mesmo) aos meus grandes amigos Fábio Albuquerque, Felipe Gomes, Wesley Hinsch e Wendson Chaper, pessoas com quem trabalhei e me ajudaram a manter a sanidade em alguns momentos difíceis.

Entrando agora nos domínios da PUC-Rio, gostaria de falar sobre os amigos que sofreram comigo os pesares da pós-graduação, Bruno Pontes, Djalma Lúcio, Grazi Kapps, João Magela, Luiz Felipe Netto, Renato Moraes e Victor Thomaz, pois com eles o percurso foi, de qualquer modo, duro, porém alegre. Tenho de falar também dos amigos de laboratório, Cássio Almeida, Jefry

Sastre, Sonia Fiol e William Fernandes, sempre dispostos a ajudar e rir quando eu estava em apuros. Agradeço também a todos os professores que influenciaram este trabalho, em especial Simone Diniz, Maria Cláudia e Marcus Poggi (não é por acaso a participação deles na banca examinadora). Reservo aqui um espaço para falar sobre duas pessoas que foram fundamentais para a conclusão deste trabalho, Pedro Furtado e Hélio Lopes, o primeiro um grande amigo, diria irmão, que fiz no departamento e o segundo outro grande amigo e por acaso também meu orientador. Se este trabalho merece aplausos, eu os divido com os ambos.

Por fim, porém não menos importante, já que se a ordem fosse determinante esta parte deveria vir na capa, gostaria de agradecer a minha família por todo apoio que me deram, dão e sei que darão. Cabe aqui fazer uma transcrição mais detalhada dos culpados. Aos meus pais Maria e Paulo, detentores da minha admiração e respeito por todo o sacrifício que fizeram para que eu pudesse realizar meus sonhos, ao meu irmão Felipe por seu apoio e admiração (gostaria de dizer "pelos conselhos", mas estes é melhor não seguir a risca). Agradeço também a minha esposa Júlia pelo amor, paciência e apoio dado durante o processo de construção deste trabalho (ainda espero conseguir que me libere para o doutorado), e minha cunhada Juliana e meus sogros Aníbal e Kátia pelo apoio e sincera admiração que me deram.

E depois de tantos nomes, era de se esperar que ao cabo pedisse desculpas por ter esquecido algum (o que não seria nenhum espanto, já que tenho lá minhas dúvidas se tão poucas linhas seriam suficientes para cobrir minha gratidão), então peço aqui minha absolvição, e digo que se isso de fato aconteceu, não me julgue mal. Se o seu nome fosse para estar aqui, é sinal de que me conhece bem e logo sabe o quão desmemoriado sou.

Resumo

Grosman, Jonatas dos Santos; Lopes, Hélio Côrtes Vieira. **LER: Anotação e classificação automática de entidades e relações**. Rio de Janeiro, 2017. 196p. Dissertação de Mestrado – Departamento de Informática, Pontifícia Universidade Católica do Rio de Janeiro.

Diversas técnicas para extração de informações estruturadas de dados em linguagem natural foram desenvolvidas e demonstraram resultados muito satisfatórios. Entretanto, para obterem tais resultados, requerem uma série de atividades que geralmente são feitas de modo isolado, como a anotação de textos para geração de corpora, etiquetamento morfossintático, engenharia e extração de atributos, treinamento de modelos de aprendizado de máquina etc., o que torna onerosa a extração dessas informações, dado o esforço e tempo a serem investidos. O presente trabalho propõe e desenvolve uma plataforma em ambiente web, chamada LER (Learning Entities and Relations) que integra o fluxo necessário para essas atividades, com uma interface que visa a facilidade de uso. Outrossim, o trabalho mostra os resultados da implementação e uso da plataforma proposta.

Palavras-chave

Processamento de linguagem natural; Aprendizado automático; Extração de informação; Ontologias; Curadoria de dados.

Abstract

Grosman, Jonatas dos Santos; Lopes, Hélio Côrtes Vieira (Advisor). **LER: Annotation and automatic classification of entities and relations**. Rio de Janeiro, 2017. 196p. Dissertação de Mestrado – Departamento de Informática, Pontifícia Universidade Católica do Rio de Janeiro.

Many techniques for the structured information extraction from natural language data have been developed and have demonstrated their potentials yielding satisfactory results. Nevertheless, to obtain such results, they require some activities that are usually done separately, such as text annotation to generate corpora, Part-Of-Speech tagging, features engineering and extraction, machine learning models' training etc., making the information extraction task a costly activity due to the effort and time spent on this. The present work proposes and develops a web based platform called LER (Learning Entities and Relations), that integrates the needed workflow for these activities, with an interface that aims the ease of use. The work also shows the platform implementation and its use.

Keywords

Natural language processing; Automatic learning; Information extraction; Ontologies; Data curation.

Sumário

1	Introdução	20
2	Revisão bibliográfica	24
2.1	Mineração de textos e processamento de linguagem natural	24
2.2	Anotação de dados textuais	28
3	A plataforma	37
3.1	Arquitetura	40
3.2	Controle de usuários	43
3.3	Gerência de projetos	44
3.4	O ERAS	46
3.4.1	Gerência dos dados	48
3.4.2	Anotação	55
3.4.3	Estatísticas	59
3.5	Aprendizado automático	68
3.6	Publicação de serviços	77
4	Experimentos	81
4.1	Experimento de anotação	81
4.1.1	Metodologia	81
4.1.2	Resultados	82
4.2	Experimento de aprendizado automático	85
4.2.1	Dados	86
4.2.2	Metodologia	95
4.2.3	Resultados	97
4.2.3.1	NER	97
4.2.3.2	RE	102
4.2.3.3	Uso dos modelos finais	105
5	Conclusão	111
	Referências bibliográficas	114
A	Experimento de anotação: Guia de anotação	118
A.1	FERRAMENTA DE ANOTAÇÃO	118
A.1.1	PROCESSO DE ANOTAÇÃO	120
A.1.2	ATALHOS	127
A.2	TAREFA DE ANOTAÇÃO	127
A.2.1	RÓTULOS	127
A.2.1.1	Actor	128
A.2.1.2	Event	129
A.2.1.3	Location	131
A.2.1.4	Time	132
A.2.1.5	Tipos primitivos	132
A.2.2	RELAÇÕES	134

A.2.3	CONECTORES	138
A.2.4	EXEMPLOS	139
A.3	CONSIDERAÇÕES FINAIS	141
B	Experimento de anotação: Comentários	143
C	Experimento de anotação: Dados dos participantes	146
C.1	User 01(A,0,5)	146
C.2	User 02(A,0,5)	146
C.3	User 03(B,0,5)	147
C.4	User 04(B,0,5)	148
C.5	User 05(A,0,10)	149
C.6	User 06(A,0,10)	149
C.7	User 07(B,0,10)	150
C.8	User 08(B,0,10)	151
C.9	User 09(A,2,5)	151
C.10	User 10(A,2,5)	152
C.11	User 11(B,2,5)	153
C.12	User 12(A,2,10)	154
C.13	User 13(A,2,10)	154
C.14	User 14(B,2,10)	155
C.15	User 15(A,4,5)	156
C.16	User 16(A,4,5)	156
C.17	User 17(B,4,5)	157
C.18	User 18(A,4,10)	158
C.19	User 19(A,4,10)	159
C.20	User 20(B,4,10)	160
D	Experimento de anotação: Tabelas e gráficos	162
E	Experimento de aprendizado automático: Tabelas e gráficos	179

Lista de figuras

1.1	Exemplos de dificuldades para ferramentas computacionais convencionais para tratamento de textos, ao lidarem com <i>tweets</i> : (a) abreviações de palavras e referência geográfica (Linha Amarela) usando um identificador de uma conta do <i>Twitter</i> (@LinhaAmarelaRJ); (b) horários e <i>hyperlinks</i> ; (c) abreviações especiais do domínio do <i>Twitter</i> , como “RT” para representar “ <i>retweet</i> ”.	21
2.1	Classes da ontologia TEDO e <i>object properties</i> (<i>datatype properties</i> omitidas para legibilidade) (6)	24
2.2	Resultado, em grafo, do fluxo proposto em (6) [Conforme figura do mesmo artigo, com adição das <i>tags</i> das relações, uma vez que no trabalho haviam números que faziam referência a uma tabela não apresentada aqui.]	27
2.3	Resultado, em RDF, do fluxo proposto em (6) [Figura obtida no referido trabalho.]	28
2.4	Processo de extração de informações do sistema TwitIE [Figura retirada de (22)]	31
2.5	Tela do ambiente de anotação do BRAT para uma aplicação específica de extração de eventos biomédicos.	32
2.6	Interface do WebAnno para configuração de projetos e definição do conjunto de <i>tags</i> . [Figura retirada de (23)]	33
2.7	Interface do WebAnno para curadoria de dados. [Figura retirada de (23)]	33
2.8	Interface do WebAnno para monitoração dos projetos. [Figura retirada de (23)]	34
3.1	Fluxo geral de estruturação de dados	38
3.2	Fluxo de curadoria dos dados	38
3.3	Fluxo de aprendizado automático	39
3.4	Fluxo de implantação dos modelos	39
3.5	LER	40
3.6	Arquitetura do LER	41
3.7	Exemplo de escalonamento vertical do LER	41
3.8	Exemplo de escalonamento horizontal do LER	42
3.9	Exemplo de escalonamento horizontal do LER utilizando cluster	42
3.10	Tela inicial da plataforma	43
3.11	Resumo do modelo de dados do LER	43
3.12	Menu da plataforma por níveis de acesso: (a) Administrador; (b) Usuário; (c) Colaborador	44
3.13	Criação de projetos	45
3.14	Descrição de verbos em português no Freeling 4.0	45
3.15	Área de projetos	46
3.16	Anotações usuário A: (a) anotação do documento 01.txt seguido de sua re-anotação; (b) anotação do documento 02.txt seguido de sua re-anotação	47

3.17 Anotações usuário B: (a) anotação do documento 01.txt seguido de sua re-anotação; (b) anotação do documento 02.txt seguido de sua re-anotação	47
3.18 Ontologia: Entidades	48
3.19 Ontologia: Relações	48
3.20 Área de gerência de dados	49
3.21 Criação de pacotes	49
3.22 Gerenciamento de colaboradores	50
3.23 Comentários	50
3.24 Perspectivas do documento	51
3.25 Perspectiva de validação: (a) estado inicial do documento; (b) estado após expansão da sentença	52
3.26 Perspectiva de validação com concordância total: (a) estado inicial do documento; (b) estado após expansão da sentença	53
3.27 Perspectiva de ações do colaborador: (a) estado inicial do documento; (b) estado após desfazer algumas ações do colaborador	54
3.28 Formato exportado	55
3.29 Tela de anotação	55
3.30 Guia de anotação	56
3.31 Processo de anotação de entidade: (a) área desejada selecionada; (b) acionando menu de contexto para seleção de entidade; (c) estado final da anotação	56
3.32 Processo de anotação de entidade com múltiplas palavras: (a) área desejada selecionada; (b) acionando menu de contexto para seleção de entidade; (c) estado final da anotação	57
3.33 Processo de anotação de relação: (a) entidade origem selecionada; (b) entidade destino alcançada; (c) estado final da anotação	57
3.34 Exemplo de tentativa de anotação de relação inválida	57
3.35 Exemplo de anotação de relação quando há múltiplas relações possíveis entre as entidades	58
3.36 Processo de anotação de conector: (a) área desejada selecionada; (b) acionando menu de contexto para seleção de ação Connect; (c) relação desejada alcançada; (d) estado final da anotação	58
3.37 Processo de remoção de anotação: (a) menu de contexto acionado no elemento desejado para executar ação Remove; (b) estado final da remoção	58
3.38 Processo de remoção de anotação em uma área: (a) menu de contexto acionado após selecionada a área desejada para executar ação Remove; (b) estado final da remoção	59
3.39 Adicionando comentários no documento	59
3.40 Estatísticas de status	59
3.41 Estatísticas de cobertura dos tokens	60
3.42 Estatísticas de distribuição dos tokens nos documentos	60
3.43 Estatísticas de entidades por colaborador, também disponíveis para relações e conectores.	60
3.44 Estatísticas de distribuição das entidades no GSA, também disponíveis para relações e conectores.	61
3.45 Nuvem de palavras do GSA	61

3.46	Tabela de ocorrências das palavras no GSA	62
3.47	Curva de frequência acumulada das palavras no GSA	62
3.48	Definição de entidades para cálculo de concordância: (a) anotação do usuário; (b) perspectiva criada para cálculo de concordância	63
3.49	Definição de relações para cálculo de concordância: (a) anotação do usuário; (b) perspectiva criada para cálculo de concordância	64
3.50	Definição de conectores para cálculo de concordância: (a) anotação do usuário; (b) perspectiva criada para cálculo de concordância	64
3.51	Tabela de classificação dos usuários A e B para os exemplos dados no capítulo	65
3.52	Mapa de concordância para o exemplo dado no capítulo: (a) entidades, relações e conectores; (b) apenas entidades; (c) apenas relações; (d) apenas conectores	66
3.53	Curva de concordância com o GSA ao longo do tempo	67
3.54	Curva de auto-concordância ao longo do tempo	67
3.55	Estatísticas básicas sobre colaboradores e documentos	68
3.56	Configurações gerais das tarefas	68
3.57	Configurações dos dados das tarefas	69
3.58	Engenharia de atributos da tarefas	71
3.59	Atributos disponíveis: (a) NER; (b) RE	71
3.60	Exemplo de atributos para a tarefa de NER: (a) documento de entrada; (b) exemplo de uma matriz criada utilizando o documento de entrada	72
3.61	Exemplo de atributos para a tarefa de RE: (a) documento de entrada; (b) perspectiva do documento de entrada criada para geração da matriz de atributos; (c) exemplo simplificado de uma matriz criada utilizando o documento de entrada; (d) exemplo de uma matriz criada utilizando o documento de entrada	73
3.62	Configurações dos classificadores	74
3.63	Execução de tarefas: (a) estado inicial da tarefa; (b) tarefa em execução; (c) tarefas executando em paralelo	75
3.64	Log de execução da tarefa	75
3.65	Resultados da tarefa	76
3.66	Scores do modelo treinado	76
3.67	Criação do serviço	78
3.68	Serviços disponíveis	78
3.69	Teste do serviço	79
3.70	Teste externo do serviço	79
4.1	Visões da ontologia usada na anotação dos dados: (a) entidades; (b) relações	82
4.2	Visões da ontologia usada na anotação dos dados para o experimento de aprendizado automático: (a) entidades; (b) relações	85
4.3	Estatísticas de tokens nos dados utilizados no experimento: (a) Cobertura dos tokens; (b) Histograma de distribuição de documentos por quantidade de tokens	87
4.4	Estatísticas de entidades nos dados utilizados no experimento: (a) Distribuição de entidades nos dados; (b) Histograma de distribuição de documentos por quantidade de entidades	88

4.5	Estatísticas de relações nos dados utilizados no experimento: (a) Distribuição de relações nos dados; (b) Histograma de distribuição de documentos por quantidade de relações	89
4.6	Estatísticas de conectores nos dados utilizados no experimento: (a) Distribuição de conectores nos dados; (b) Histograma de distribuição de documentos por quantidade de conectores	90
4.7	Estatísticas de auto-concordância ano longo do tempo nos dados utilizados no experimento: (a) auto-concordância com entidades, relações e conectores; (b) auto-concordância com entidades; (c) auto-concordância com relações; (d) auto-concordância com conectores	91
4.8	Nuvem de palavras associadas às entidades nos dados utilizados no experimento	91
4.9	Curva de frequência acumulada de palavras associadas às entidades: (a) todas entidades; (b) excluindo a entidade com mais ocorrências, neste caso <i>Location</i>	93
4.10	Nuvem de palavras associadas aos conectores nos dados utilizados no experimento	94
4.11	Curva de frequência acumulada de palavras associadas aos conectores: (a) todos conectores; (b) excluindo conector com mais ocorrências, neste caso <i>hasEvent</i>	95
4.12	Generalização e remoções usadas para redução das classes	101
4.13	Teste com resposta razoável: (a) <i>tweet</i> original; (b) anotação predita; (c) anotação esperada	106
4.14	Teste com resposta razoável: triplas RDF retornadas	106
4.15	Teste com boa resposta: (a) <i>tweet</i> original; (b) anotação predita; (c) anotação esperada	107
4.16	Teste com boa resposta: triplas RDF retornadas	108
4.17	Teste com resposta perfeita: (a) <i>tweet</i> original; (b) anotação predita; (c) anotação esperada	109
4.18	Teste com resposta perfeita: triplas RDF retornadas	109
A.1	Ferramenta de anotação	118
A.2	Resumo da ontologia de anotação	119
A.3	Adicionando comentários	119
A.4	Rotulação, passo 1	120
A.5	Rotulação, passo 2	121
A.6	Rotulação, resultado final	121
A.7	Relacionamento entre os rótulos (relação inválida)	122
A.8	Relação inválida de <i>#Car</i> para <i>#Accident</i> , porém válida de <i>#Accident</i> para <i>#Car</i>	122
A.9	Relacionamento entre os rótulos (relação válida)	122
A.10	Relacionamento entre os rótulos, resultado final	122
A.11	Relações entre sentenças	123
A.12	Criação de conectores, passo 1	123
A.13	Criação de conectores, passo 2	124
A.14	Criação de conectores, passo 3	124
A.15	Criação de conectores, resultado final	124
A.16	Remoção de rótulos, relações e conectores	125

A.17 Remoção por área de de rótulos, relações e conectores, passo 1	125
A.18 Remoção por área de de rótulos, relações e conectores, passo 2	126
A.19 Remoção por área de de rótulos, relações e conectores, resultado final	126
A.20 Hierarquia de classes da ontologia de eventos de trânsito	128
A.21 Exemplo de rotulação, Actor	128
A.22 Exemplo de rotulação, Car e Motorcycle	129
A.23 Exemplo de rotulação, Interdiction e Event	129
A.24 Exemplo de rotulação, Protest e Interdiction	129
A.25 Exemplo de rotulação, RoadWork	130
A.26 Exemplo de rotulação, WeatherEvent	130
A.27 Exemplo de rotulação, Solution e Accident	130
A.28 Exemplo de rotulação, Breakdown	130
A.29 GoodTrafficSituation	130
A.30 HeavyTrafficSituation	131
A.31 SlowTrafficSituation	131
A.32 Exemplo 1 de rotulação, Location	131
A.33 Exemplo 2 de rotulação, Location	131
A.34 Exemplo 1 de rotulação, Time	132
A.35 Exemplo 2 de rotulação, Time	132
A.36 Tipos primitivos	132
A.37 wayEffect:BothDirections	133
A.38 wayEffect:OneDirection	133
A.39 wayEffect:Partially	133
A.40 xsd:unsignedInt	133
A.41 xsd:string	134
A.42 causes	135
A.43 flowsTo	135
A.44 hasActor	135
A.45 hasEvent	136
A.46 hasSupporter	136
A.47 hasTime	136
A.48 isAlternativeFor	137
A.49 isEdgeFor	137
A.50 isReferenceFor	137
A.51 isRestrictedTo	137
A.52 hasNumericQuantity	138
A.53 hasStringQuantity	138
A.54 hasWayEffect	138
A.55 conectores	139
A.56 exemplo 1	139
A.57 exemplo 2	139
A.58 exemplo 3	140
A.59 exemplo 4	140
A.60 exemplo 5	140
A.61 exemplo 6	140
A.62 exemplo 7	141
A.63 exemplo 8	141

A.64	Tweet de evento ferroviário	142
A.65	Tweet de evento climático	142
A.66	Tweet de informação geral	142
A.67	Tweet fora de contexto	142
D.1	Experimento de anotação, tokens	163
D.2	Experimento de anotação, entidades	163
D.3	Experimento de anotação: relações	164
D.4	Experimento de anotação, conectores	164
D.5	Experimento de anotação, concordância com entidades, relações e conectores	165
D.6	Experimento de anotação, concordância com entidades	165
D.7	Experimento de anotação, concordância com relações	166
D.8	Experimento de anotação, concordância com conectores	166
D.9	Experimento de anotação, concordância com GSA ao longo do tempo do grupo (0,5): (a) entidades, relações e conectores; (b) entidades; (c) relações; (d) conectores	167
D.10	Experimento de anotação, auto-concordância ao longo do tempo do grupo (0,5): (a) entidades, relações e conectores; (b) entidades; (c) relações; (d) conectores	168
D.11	Experimento de anotação, concordância com GSA ao longo do tempo do grupo (0,10): (a) entidades, relações e conectores; (b) entidades; (c) relações; (d) conectores	169
D.12	Experimento de anotação, auto-concordância ao longo do tempo do grupo (0,10): (a) entidades, relações e conectores; (b) entidades; (c) relações; (d) conectores	170
D.13	Experimento de anotação, concordância com GSA ao longo do tempo do grupo (2,5): (a) entidades, relações e conectores; (b) entidades; (c) relações; (d) conectores	171
D.14	Experimento de anotação, auto-concordância ao longo do tempo do grupo (2,5): (a) entidades, relações e conectores; (b) entidades; (c) relações; (d) conectores	172
D.15	Experimento de anotação, concordância com GSA ao longo do tempo do grupo (2,10): (a) entidades, relações e conectores; (b) entidades; (c) relações; (d) conectores	173
D.16	Experimento de anotação, auto-concordância ao longo do tempo do grupo (2,10): (a) entidades, relações e conectores; (b) entidades; (c) relações; (d) conectores	174
D.17	Experimento de anotação, concordância com GSA ao longo do tempo do grupo (4,5): (a) entidades, relações e conectores; (b) entidades; (c) relações; (d) conectores	175
D.18	Experimento de anotação, auto-concordância ao longo do tempo do grupo (4,5): (a) entidades, relações e conectores; (b) entidades; (c) relações; (d) conectores	176
D.19	Experimento de anotação, concordância com GSA ao longo do tempo do grupo (4,10): (a) entidades, relações e conectores; (b) entidades; (c) relações; (d) conectores	177

D.20 Experimento de anotação, auto-concordância ao longo do tempo do grupo (4,10): (a) entidades, relações e conectores; (b) entidades; (c) relações; (d) conectores

178

Lista de tabelas

3.1	Concordâncias observadas e esperadas ao acaso no exemplo dado	65
4.1	Concordância com entidades, relações e conectores nos grupos	83
4.2	Concordância com entidades nos grupos	83
4.3	Concordância com relações nos grupos	83
4.4	Concordância com conectores nos grupos	83
4.5	Concordância com entidades, relações e conectores nos tipos de anotador	84
4.6	Concordância com entidades nos tipos de anotador	84
4.7	Concordância com relações nos tipos de anotador	84
4.8	Concordância com conectores nos tipos de anotador	84
4.9	Estatísticas de tempo do conjunto de dados usado no experimento de aprendizado automático	86
4.10	Resultado NER passo 5	97
4.11	Gridsearch no NER para o SVC no passo 5	97
4.12	Resultado NER passo 6	98
4.13	Gridsearch no NER para o Random Forest no passo 8	98
4.14	Gridsearch no NER para o Stochastic Gradient Descent no passo 8	98
4.15	Gridsearch no NER para o SVC no passo 8	99
4.16	Resultado NER passo 8	99
4.17	Resultado final detalhado: Random Forest (validação)	99
4.18	Resultado final detalhado: Random Forest (teste)	100
4.19	Resultado NER final	101
4.20	Resultado final detalhado: Random Forest (validação com redução de classes)	102
4.21	Resultado final detalhado: Random Forest (teste com redução de classes)	102
4.22	Resultado RE passo 5	103
4.23	Gridsearch no RE para o Frank Wolfe SSVM no passo 8	103
4.24	Gridsearch no RE para o Structured Perceptron no passo 8	103
4.25	Resultado RE passo 8	103
4.26	Resultado final detalhado: Frank Wolfe SSVM (validação)	103
4.27	Resultado final detalhado: Frank Wolfe SSVM (teste)	104
4.28	Resultado RE final	104
4.29	Resultado final detalhado: Frank Wolfe SSVM (validação com redução de classes)	104
4.30	Resultado final detalhado: Frank Wolfe SSVM (teste com redução de classes)	105
C.1	Resumo das respostas dos participantes ao questionário do experimento de anotação	161
D.1	Resumo colaborações	162
E.1	Training scores: model-fscore-TEDO-NER-RF-0	179

E.2	Revalidation scores: model-fscore-TEDO-NER-RF-0	180
E.3	Training scores: model-fscore-TEDO-NER-SGD-0	180
E.4	Revalidation scores: model-fscore-TEDO-NER-SGD-0	181
E.5	Training scores: model-fscore-TEDO-NER-SVC-0	181
E.6	Revalidation scores: model-fscore-TEDO-NER-SVC-0	182
E.7	Training scores: model-fscore-TEDO-NER-RF-1	182
E.8	Revalidation scores: model-fscore-TEDO-NER-RF-1	183
E.9	Training scores: model-fscore-TEDO-NER-SGD-1	183
E.10	Revalidation scores: model-fscore-TEDO-NER-SGD-1	184
E.11	Best parameters: model-fscore-TEDO-NER-SVC-1	184
E.12	Training scores: model-fscore-TEDO-NER-SVC-1	185
E.13	Revalidation scores: model-fscore-TEDO-NER-SVC-1	185
E.14	Best parameters: model-fscore-TEDO-NER-RF-2	186
E.15	Training scores: model-fscore-TEDO-NER-RF-2	186
E.16	Revalidation scores: model-fscore-TEDO-NER-RF-2	187
E.17	Best parameters: model-fscore-TEDO-NER-SGD-2	187
E.18	Training scores: model-fscore-TEDO-NER-SGD-2	188
E.19	Revalidation scores: model-fscore-TEDO-NER-SGD-2	188
E.20	Best parameters: model-fscore-TEDO-NER-SVC-2	189
E.21	Training scores: model-fscore-TEDO-NER-SVC-2	189
E.22	Revalidation scores: model-fscore-TEDO-NER-SVC-2	190
E.23	Training scores: model-fscore-TEDO-NER-RF-FINAL	190
E.24	Test scores: model-fscore-TEDO-NER-RF-FINAL	191
E.25	Training scores: model-fscore-TEDO-NER-RF-FINAL-AGRUPADO	191
E.26	Test scores: model-fscore-TEDO-NER-RF-FINAL-AGRUPADO	192
E.27	Training scores: model-fscore-TEDO-RE-SP-0	192
E.28	Revalidation scores: model-fscore-TEDO-RE-SP-0	192
E.29	Training scores: model-fscore-TEDO-RE-FWSSVM-0	193
E.30	Revalidation scores: model-fscore-TEDO-RE-FWSSVM-0	193
E.31	Best parameters: model-fscore-TEDO-RE-SP-1	193
E.32	Training scores: model-fscore-TEDO-RE-SP-1	193
E.33	Revalidation scores: model-fscore-TEDO-RE-SP-1	194
E.34	Best parameters: model-fscore-TEDO-RE-FWSSVM-1	194
E.35	Training scores: model-fscore-TEDO-RE-FWSSVM-1	194
E.36	Revalidation scores: model-fscore-TEDO-RE-FWSSVM-1	194
E.37	Training scores: model-fscore-TEDO-RE-FWSSVM-FINAL	195
E.38	Test scores: model-fscore-TEDO-RE-FWSSVM-FINAL	195
E.39	Training scores: model-fscore-TEDO-RE-FWSSVM-FINAL- AGRUPADO	195
E.40	Test scores: model-fscore-TEDO-RE-FWSSVM-FINAL-AGRUPADO	196

Call me Ishmael.

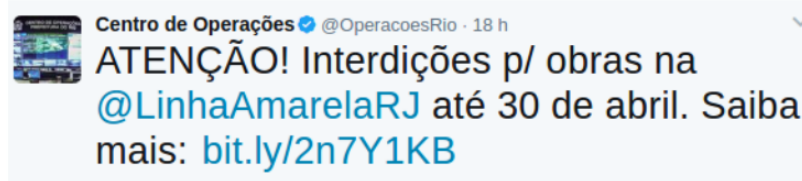
Herman Melville, *Moby-Dick; or, The Whale.*

1 Introdução

Um trabalho de 2014 (1) indicou o potencial de retorno econômico de técnicas de mineração de dados em geral. Estimou-se, por exemplo, que um melhor aproveitamento de dados com as técnicas apropriadas de mineração textual e não-textual permitiria ao sistema de saúde americano a criação de mais de US\$ 300 bilhões em valor anualmente. Igualmente, o uso eficiente de informações dos dados para melhorar operações e detectar fraudes poderia gerar para a administração do setor público europeu até US\$ 250 bilhões de valor potencial anual.

O volume de dados não-estruturados disponíveis em formato digital tem crescido de modo intenso nas últimas décadas. Empresas de diversas áreas percebem, cada vez mais claramente, o potencial econômico do uso inteligente dos dados de grandes repositórios, sobretudo de redes sociais (2). Algumas companhias têm usado ferramentas de mídia social como *Facebook* e *Twitter* para prover serviços e interagir com clientes, tendo como resultado uma enorme quantidade de conteúdo gerado por usuários disponível sem custos (3).

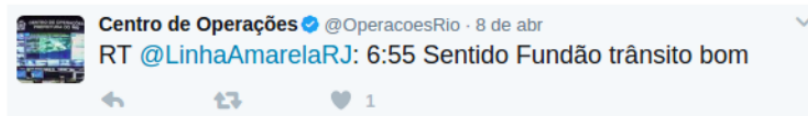
Há, portanto, um grande universo de dados não-estruturados, com elevado valor, passíveis de aquisição com custos relativamente baixos. O problema dos dados textuais, típicos de redes sociais e microblogs em geral, reside em sua intratabilidade computacional em seus formatos originais. Há diferenças de dificuldades mesmo entre as fontes. O *Twitter*, por exemplo, limitado a 140 caracteres, representa um grande desafio para as tecnologias convencionais de mineração de dados textuais em razão de conter linguagem informal, *emojis*, *hashtags* e outras formas diferenciadas de escrita e expressão (4). A figura 1.1 ilustra essa realidade em *tweets* em português relacionados a questões de trânsito.



(a)



(b)



(c)

Figura 1.1: Exemplos de dificuldades para ferramentas computacionais convencionais para tratamento de textos, ao lidarem com *tweets*: (a) abreviações de palavras e referência geográfica (Linha Amarela) usando um identificador de uma conta do *Twitter* (@LinhaAmarelaRJ); (b) horários e *hyperlinks*; (c) abreviações especiais do domínio do *Twitter*, como “RT” para representar “*retweet*”.

Diversas técnicas de *Text Analytics* (TA) foram desenvolvidas e demonstraram capacidade de extrair informações (2). Uma delas faz uso de ontologias de domínio como estrutura para a interpretação computacional de dados textuais, o que despertou um interesse diferenciado em razão da união entre dois grandes campos de pesquisa: ontologias e processamento de linguagem natural (NLP). Essa união tende a gerar bons resultados, como demonstrado por resultados de algumas pesquisas (5), (6).

Entretanto, os bons resultados gerados nestes trabalhos requerem uma série de atividades que geralmente são feitas de modo isolado, o que se torna uma dificuldade, tendo em vista o esforço e tempo a serem investidos. É um forte exemplo a etapa de anotação dos textos para geração dos conjuntos de dados de treinamento e teste: existem ferramentas de anotação disponíveis, grátis e até mesmo em código aberto, apresentando vantagens e desvantagens. Uma desvantagem comum é a falta de suporte ao uso de ontologias como geradores de *tags*, muito menos que levem em conta as propriedades (*object properties* e *datatype properties*) descritas nos arquivos de ontologias. Algumas ferramentas se mostram pouco atrativas no que tange à usabilidade, requerendo muitos e complicados passos para anotação, o que induz a erros. Além disso, a maioria das ferramentas de anotação não está integrada, em um mesmo ambiente, às ferramentas de modelagem e aprendizado de máquina. Não se tem um sistema geral que integre as etapas de coleta e armazenamento de dados; caracterização do texto; anotação colaborativa com base em ontologias; treinamento, teste e persistência de modelos e disponibilização final destas

etapas como serviços automáticos de estruturação de dados textuais em ambiente de nuvem.

Assim, a contribuição deste trabalho é a criação deste sistema integrado no formato *web*, acessível via *browser*, com uma interface de anotação de uso mais simples, mantendo as principais vantagens dos mais conhecidos sistemas de anotação, integração com um conjunto de algoritmos de aprendizado de máquina e, como produto final, a hospedagem das etapas desenvolvidas como serviços de geração dados no formato RDF para os mais diversos domínios. Este sistema tem uma forma totalmente modular, no formato de serviços, o que permitirá a conexão, no futuro, de outras tarefas além de NER e RE.

Além do desenvolvimento do sistema, são implementadas alternativas para melhoria dos resultados já alcançados em trabalhos semelhantes (5), (7) pelo uso de outras estratégias de criação de *features* relacionadas às funções morfossintáticas das palavras e pela avaliação do sistema como um todo em um problema já tratado anteriormente (5), (6) com alguma melhoria na ontologia de eventos usada.

Os próximos capítulos deste trabalho tratarão da revisão bibliográfica (capítulo 2), da descrição da plataforma construída (capítulo 3), dos experimentos conduzidos para avaliação do processo de anotação e aprendizado automático (capítulo 4) e por fim das conclusões e sugestões para trabalhos futuros (capítulo 5).

Ignorance is the parent of fear.

Herman Melville, *Moby-Dick; or, The Whale*.

2

Revisão bibliográfica

Este capítulo apresenta a revisão bibliográfica dos trabalhos mais relevantes sobre mineração de textos (seção 2.1) e anotação de dados textuais (seção 2.2)

2.1

Mineração de textos e processamento de linguagem natural

Ittoo e colegas (2) discutem acerca dos principais desafios e tendências no campo de TA, termo usado para a união dos conceitos de *Text Mining* e NLP. São apresentados os resultados recentes mais importantes nas pesquisas acadêmicas em TA sob diversos aspectos. Uma das aplicações que merecem destaque, segundo os autores, é a descrita em (6), onde é proposta uma metodologia para a interpretação de *tweets* relacionados a eventos de trânsito na cidade do Rio de Janeiro. O artigo propõe e usa uma ontologia de eventos de trânsito denominada TEDO (*Traffic Event Domain Ontology*) para modelar as situações de trânsito como eventos, compostos por atores, locais e horários. A TEDO se baseia nas noções de evento (8) e relações entre eventos (9), (10). A figura 2.1 mostra o esquema resumido desta ontologia.

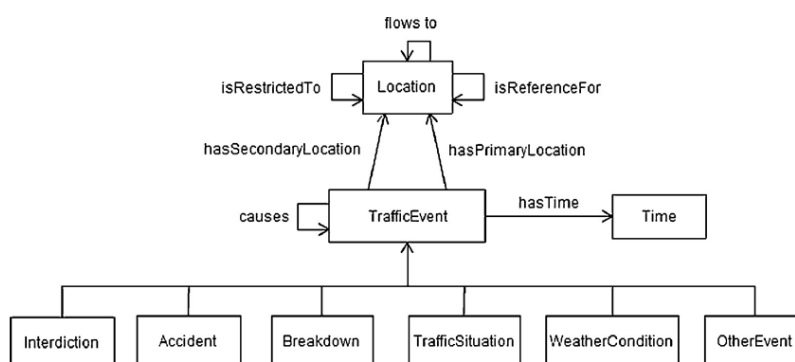


Figura 2.1: Classes da ontologia TEDO e *object properties* (*datatype properties* omitidas para legibilidade) (6)

A interpretação dos *tweets* envolve um conjunto de operações em série: (1) Tokenização dos *tweets* e etiquetagem morfosintática (*Part-of-speech tagging*); (2) Reconhecimento de entidades (*Named Entity Recognition* NER); (3) Geolocalização; (4) Extração de relações entre entidades (*Relation Extraction* RE); (5) Geração de triplas RDF (*Resource Description Framework*).

A etapa de tokenização e POS *tagging* gera atributos para o aprendizado de máquina que expressam as funções sintáticas dos *tokens* nas sentenças

analisadas, usando para isso um algoritmo proprietário (formato *Web service*) chamado F-EXT (11). As informações dos *tokens* geradas por este algoritmo e usadas na construção de atributos no referido artigo são:

- Token (W_i): o conteúdo do *token* X_i ;
- Simple Token (SMW_i): o conteúdo simplificado do *token* X_i (em caixa baixa, sem quaisquer caracteres especiais ou pontuações);
- Simplified Token (SW_i): o conteúdo simplificado do *token* X_i (em caixa baixa, sem quaisquer caracteres especiais, também removendo letras, pontuação e números de tamanho 1);
- Part-of-speech (POS_i): a etiqueta morfossintática do *token* X_i ;
- Stemmed Word (STW_i): a raiz da palavra presente no *token* X_i . Exemplo: se X_i é *blocked*, STW_i é *block*.

Outros atributos foram construídos (*feature engineering*) e, unidos ao POS *tagging*, fornecem informações úteis à aplicação de algoritmos de aprendizado de máquina. A primeira etapa de aprendizado é a de NER, onde as entidades descritas na TEDO são usadas como classes em uma tarefa de classificação. Os atributos construídos para esta tarefa são:

- CurrT (X): o *token* da posição atual, X;
- PrevT (X, N): o *token* N posições antes do *token* X;
- NextT (X, N): o *token* N posições depois do *token* X;
- CurrWSC (X): indica se X começa em letra maiúscula e se há apenas uma letra maiúscula em todo o *token*;
- LocType (X): indica se a representação de X em letras minúsculas pertence a um determinado conjunto de palavras que designam localidades. Exemplos: “avenida”, “av.”, “rua”, “estrada”.

O classificador usado nesta tarefa é uma implementação SMO (*Sequential Minimal Optimization*) (12) da família de métodos SVM (*Support Vector Machine*) do pacote Weka 3.6.5 (13). Uma das conclusões do trabalho é a de que o método SVM apresenta os melhores resultados para NER, frente a outros algoritmos testados.

Cita-se também uma etapa de geolocalização determina as coordenadas das entidades nomeadas identificadas, usando para isso o algoritmo *SmartGeocode* (7). Em seguida, as relações entre entidades são extraídas pela etapa de RE, que cria uma árvore de dependências G_T a partir de um *tweet* T com entidades nomeadas. Para a extração de G_T , são computados atributos dos pares de elementos textuais K e L, que são usados por um algoritmo para o aprendizado da natureza das (possíveis) relações entre K e L. A implementação de RE do referido trabalho usa um perceptron estruturado de margem larga (14) que estabelece um peso para cada aresta de um grafo dirigido completo (todas as relações entre todos os nós e em todas as direções). Em seguida, um algoritmo de *Maximum Spanning Tree* define, no grafo completo, a árvore final com as relações mais importantes.

Os atributos de nó usados no trabalho são:

- Word (W): indica as palavras do nó K;
- Simplified Word(SW): as palavras simplificadas do nó K;
- Ruler Entity (RE): a entidade nomeada em K, em se tratando de entidade relevante. Se não relevante, usa-se a função morfosintática, POS;
- Named Entity (NE): a entidade nomeada em K, ignorada se não for uma entidade relevante;
- Punctuation (PUNCT): indica se há uma pontuação no texto em K.

Com estes atributos de nó, são construídos os atributos das relações, que são:

- PossRel: indica se os nós K e L podem ter uma relação;
- ConcT (Y): as palavras concatenadas dos nós K e L;
- BetT (Y): a concatenação de todos os *tokens* entre K e L;
- Near (Y, N): a concatenação de todos os N *tokens* antes de K e N *tokens* depois de L;
- AbsLocPair: indica se K e L têm uma relação, onde a entidade em K é <restriction> e em L é <location-name>. Se não há relação, nenhum valor é usado;
- MetaWithLoc: indica se K e L têm uma relação, onde a entidade em K é <reference> ou <direction> e em L é <location-name>. Se não há relação, nenhum valor é usado.

Ao fim desse processo, a informação de entidades e relações é traduzida como triplas RDF e, dessa forma, o *tweet* em texto livre ganha uma estrutura semântica baseada na ontologia.

Como exemplo, apresenta-se um *tweet* em sua versão original e os resultados, respectivamente nas figuras 2.2 e 2.3, do fluxo discutido acima em formato de grafo e em RDF:

Tweet: Acidente entre 2 carros na Av das Américas na pista sentido Grota Funda próximo ao número 19880

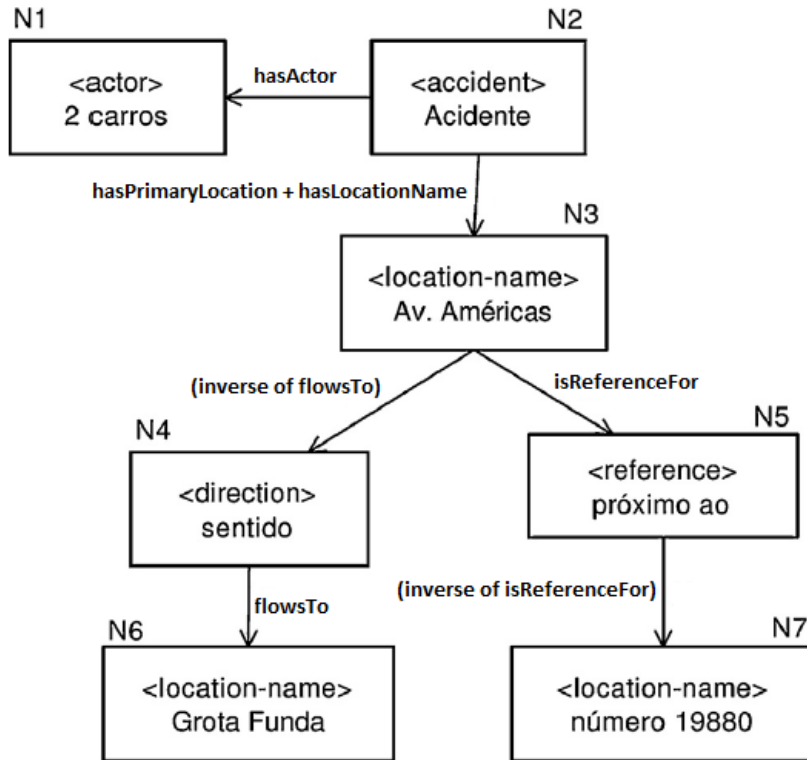


Figura 2.2: Resultado, em grafo, do fluxo proposto em (6) [Conforme figura do mesmo artigo, com adição das *tags* das relações, uma vez que no trabalho haviam números que faziam referência a uma tabela não apresentada aqui.]

```

1. <tedo:Accident
   rdf:about='http://www.ld.inf.puc-rio.br/tedo/trafficevent/100''>
2.   <tedo:hasActor>2 carros</tedo:hasActor >
3.   <tedo:hasPrimaryLocation
   id="http://www.ld.inf.puc-rio.br/tedo/location/101">
4.   <tedo:hasSecondaryLocation
   id="http://www.ld.inf.puc-rio.br/tedo/location/102">
5.   <tedo:hasSecondaryLocation
   id="http://www.ld.inf.puc-rio.br/tedo/location/103">
6.   <tedo:hasTime id="http://www.ld.inf.puc-rio.br/time/104">
7. </tedo:Accident>
8. <!-- Primary Location -->
9. <tedo:Location
   rdf:about="http://www.ld.inf.puc-rio.br/location/101">
10.  <tedo:hasLocationName>Av. Das Amricas</tedo:hasLocationName>
11.  <tedo:hasCoordinates>-22.998864, -43.365984</tedo:hasCoordinates>
12.  <tedo:flowsTo
   rdf:resource="http://www.ld.inf.puc-rio.br/location/102">
13.  <tedo:isReferenceFor
   rdf:resource="http://www.ld.inf.puc-rio.br/location/103">
14. </tedo:Location>
15. <!-- Secondary Location -->
16. <tedo:Location
   rdf:about="http://www.ld.inf.puc-rio.br/location/102">
17.  <tedo:hasLocationName>Grotta Funda</tedo:hasLocationName>
18.  <tedo:hasCoordinates>-23.015379, -43.521634</tedo:hasCoordinates>
19. </tedo:Location> \\
20. <!-- Secondary Location -->\\
21. <tedo:Location
   rdf:about="http://www.ld.inf.puc-rio.br/location/103">
22.  <tedo:hasLocationName>Nmero 19880</tedo:hasLocationName>
23.  <tedo:hasCoordinates> -23.016279, -43.514426</tedo:hasCoordinates>
24. </tedo:Location>
25. <!-- Time -->
26. <tedo:Time
   rdf:about="http://www.ld.inf.puc-rio.br/time/104">
27.  <tedo:hasPublicationTime>05/03/2012 07:07:01</tedo:hasPublicationTime>
28. </tedo:Time>

```

Figura 2.3: Resultado, em RDF, do fluxo proposto em (6) [Figura obtida no referido trabalho.]

2.2

Anotação de dados textuais

Sistemas de NLP se baseiam em técnicas de aprendizado supervisionado, que para bons resultados dependem de uma grande quantidade de dados anotados manualmente por especialistas do domínio a ser modelado. Em alguns casos, a falta de dados anotados para treino é um obstáculo ao desenvolvimento de ferramentas de NLP baseadas em aprendizado de máquina. O processo de anotação é altamente custoso em termos financeiros e de tempo. Contudo, a falta de dados anotados manualmente submetem os sistemas de NLP ao que se denomina “gargalo de aquisição de conhecimento” (15). Em anos recentes, um dos principais caminhos adotados para facilitar a aquisição destes dados é o uso de técnicas de *crowdsourcing* via *web*, isto é, pelo uso de anotadores não especializados recrutados na internet. A título de exemplo, em trabalho recente na Indonésia (16), propõe-se uma inovadora ferramenta móvel para a anotação colaborativa, avaliando-a em um experimento envolvendo 15 estudantes indonésios, que anotaram 1500 dados textuais através de seus *smartphones*.

Hovy e Lavid (17) discutem sobre o fato da qualidade dos modelos treinados estar diretamente ligada à qualidade dos dados utilizados. Descrevem ainda a necessidade de haver no mínimo dois anotadores atuando sobre um mesmo documento, para que ao cabo da anotação, a concordância entre eles seja calculada, e sendo esta insatisfatória, há indícios de problemas na definição da tarefa, ou até mesmo, indique o fato de que a tarefa em questão seja muito difícil. Segundo os autores, as etapas que devem ser seguidas em um processo de anotação são:

1. Identificar e preparar uma seleção dos textos representativos como material de partida para o “corpus de treinamento”;
2. Instanciar uma determinada teoria linguística ou conceito linguístico, especificar o conjunto de tags a usar, suas condições de aplicabilidade, etc. Esta etapa inclui o início da criação do guia de anotação;
3. Anotar parte do corpus de treinamento, a fim de determinar a viabilidade tanto da instanciação quanto do guia de anotação;
4. Medir os resultados (comparar as decisões dos anotadores) e decidir quais as medidas apropriadas, e como devem ser aplicadas;
5. Determinar que nível de concordância deve ser considerado satisfatório (baixa concordância significa pouca consistência na anotação para permitir que os algoritmos de aprendizado de máquina sejam treinados com sucesso). Se a concordância não é (ainda) satisfatória, o processo se repete a partir do passo 2, com as mudanças apropriadas à teoria, sua instanciação, o guia e as instruções do anotador. Caso contrário, o processo continua para o passo 6.
6. Anotar uma grande parte do corpus, possivelmente ao longo de vários meses ou anos, com muitas verificações intermediárias, melhorias, etc.
7. Quando um material suficiente for anotado, treinar um modelo de aprendizado automático em parte do corpus, e posteriormente medir o desempenho do modelo na parte restante dos dados;
8. Se o desempenho do modelo for satisfatório, este pode ser utilizado para anotação automática de documentos de mesma natureza dos utilizados no treino. Se o desempenho não for satisfatório, o processo se repete, possivelmente a partir do passo 2, ou a partir do passo 6, se forem necessários mais dados de treinamento.

Para a execução do processo a pouco descrito, se faz necessário o uso de uma ferramenta de anotação. Há diversas ferramentas disponíveis, algumas de uso grátis. Uma das mais conhecidas e usadas, a GATE Teamware, faz parte de um programa de pesquisa de mais de 20 anos e é um dos muitos componentes da GATE (18), uma arquitetura geral para engenharia de textos (19), (20). Após anos de desenvolvimento e contribuições de pesquisadores e usuários de diversas áreas, a GATE se tornou um grande ecossistema, cobrindo de forma diferenciada todo o ciclo de vida de desenvolvimento de sistemas

de análises de textos. Essa família de ferramentas tem crescido ao longo dos anos e hoje compreende uma aplicação *desktop* para desenvolvedores, uma aplicação *web* baseada em trabalho colaborativo, bibliotecas em Java, uma arquitetura e um processo bem definidos. A GATE Teamware pode ser baixada e instalada como um serviço para gestão de todo o projeto de anotação, desde o cadastro dos anotadores, passando por todo o processo de monitoração dos comportamentos dos anotadores e gerando um conjunto final de dados anotados, corrigidos e selecionados por um curador. Essa ferramenta parece ter um foco maior em integrar componentes para anotação automática e menor em prover interfaces para anotação manual. Também se trata de uma ferramenta complexa para usuários não especializados, requerendo um considerável tempo de treinamento, o que de certa forma limita a ideia de *crowdsourcing* para a anotação. Outro ponto importante a se destacar é que esse sistema permite o uso de ontologias para a anotação de entidades, mas não das relações entre elas, o que é fundamental para o caso tratado em (6). Mesmo com essas limitações da aplicação de anotação, a família GATE como um todo fornece muitas ferramentas e bibliotecas úteis para trabalhos futuros neste tema. Um exemplo interessante e no contexto de (6) é o do sistema TwitIE (22), dedicado ao clássico reconhecimento de entidades nomeadas (NER) e à extração de informações (IE) justamente em *tweets*, que pelo uso de módulos da GATE implementa um processo semelhante ao dos *tweets* de trânsito (ver figura 2.4). Cabe ressaltar que essa aplicação não se baseia em ontologias de domínio, como no caso do TEDO, mas busca por entidades nomeadas em geral em um ambiente reconhecidamente difícil e ruidoso, como são os microblogs em geral.

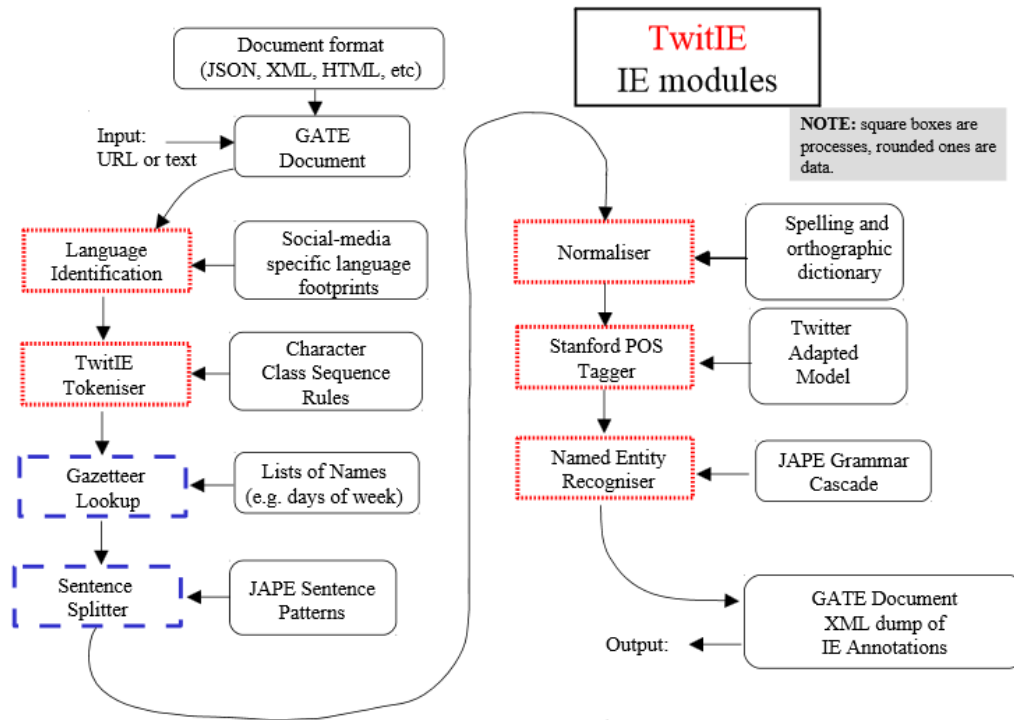


Figura 2.4: Processo de extração de informações do sistema TwitIE [Figura retirada de (22)]

Outra ferramenta de anotação é a BRAT (*Brat Rapid Annotation Tool*) (21), que foi a primeira ferramenta em código aberto para anotação totalmente baseada em *web* e suportando anotações colaborativas de múltiplas camadas simultaneamente em uma única cópia de documento. Ela apresenta algumas limitações, como a lentidão para processamento de documentos com mais de 100 sentenças, poucos formatos de arquivo suportados e, principalmente, impossibilidade de carregamento direto de um arquivo de ontologia para uso na identificação de entidades e relações. A figura 2.5 mostra um exemplo da tela do ambiente de anotação do BRAT para uma aplicação específica de extração de eventos biomédicos.

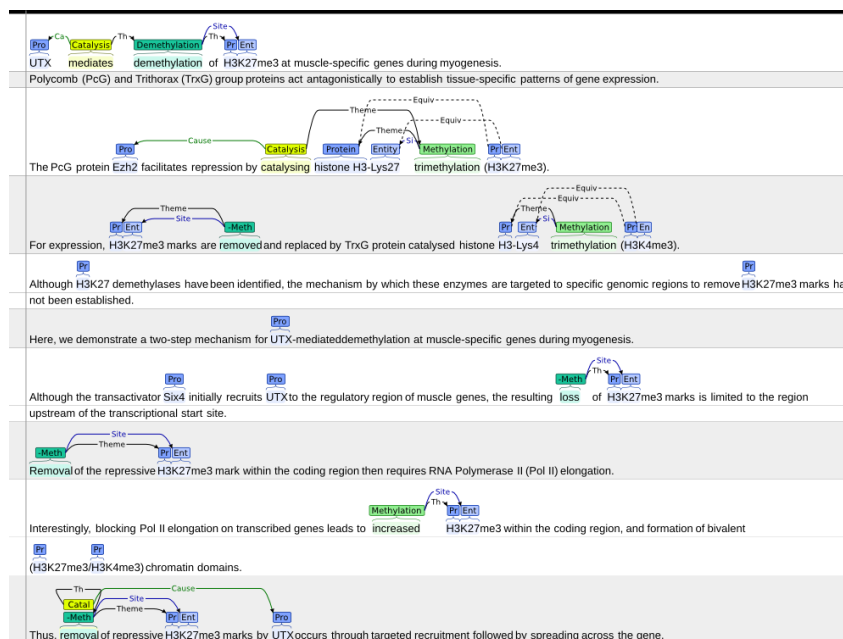


Figura 2.5: Tela do ambiente de anotação do BRAT para uma aplicação específica de extração de eventos biomédicos.

Há também a ferramenta WebAnno (23), que utiliza a BRAT como interface de *frontend*, mas substitui sua camada *server* para adicionar o suporte ao gerenciamento de usuários e da qualidade da anotação com os procedimentos de curadoria. A implementação da WebAnno também disponibiliza uma interface modificada da BRAT para facilitação da anotação em *crowdsourcing*. Semelhante às outras ferramentas, ela permite a anotação de relações de dependência, onde duas anotações de POS tag são conectadas por uma relação direta dentro de opções já definidas no sistema, anotações de co-referências, mas não permite o carregamento de uma ontologia para a anotação das entidades e suas relações. O mais próximo que esta chega de uma flexibilidade para os tipos de entidades é a possibilidade de se fornecer um conjunto de *tags* (*tagset*) a serem usadas na anotação. As figuras 2.6, 2.7 e 2.8 mostram telas do WebAnno para configuração de projeto, curadoria e monitoração de projeto.

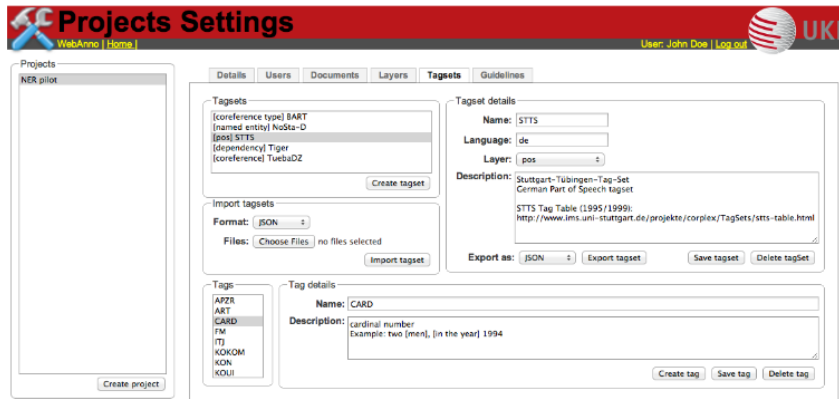


Figura 2.6: Interface do WebAnno para configuração de projetos e definição do conjunto de tags. [Figura retirada de (23)]

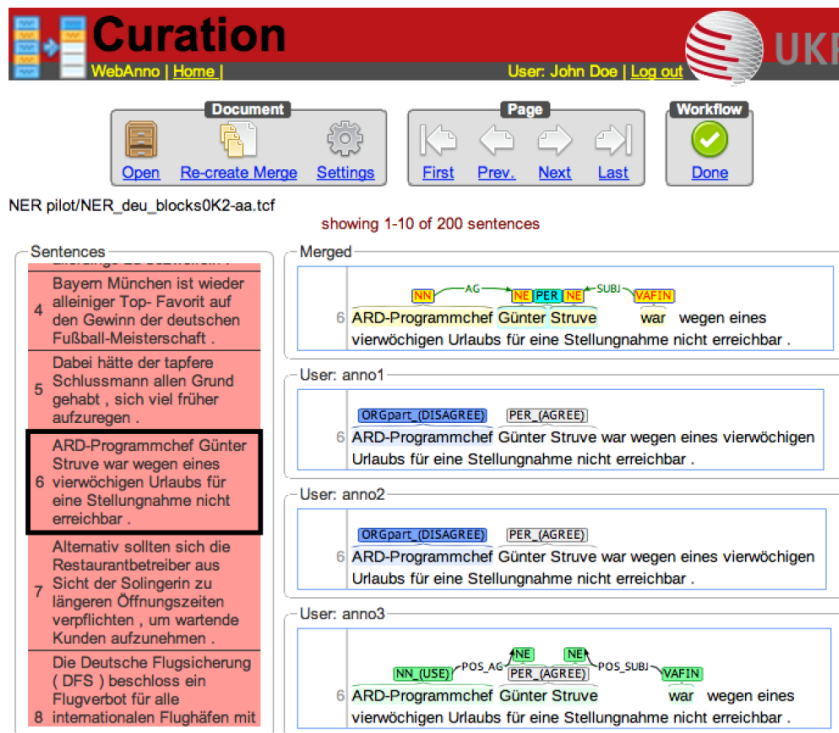


Figura 2.7: Interface do WebAnno para curadoria de dados. [Figura retirada de (23)]

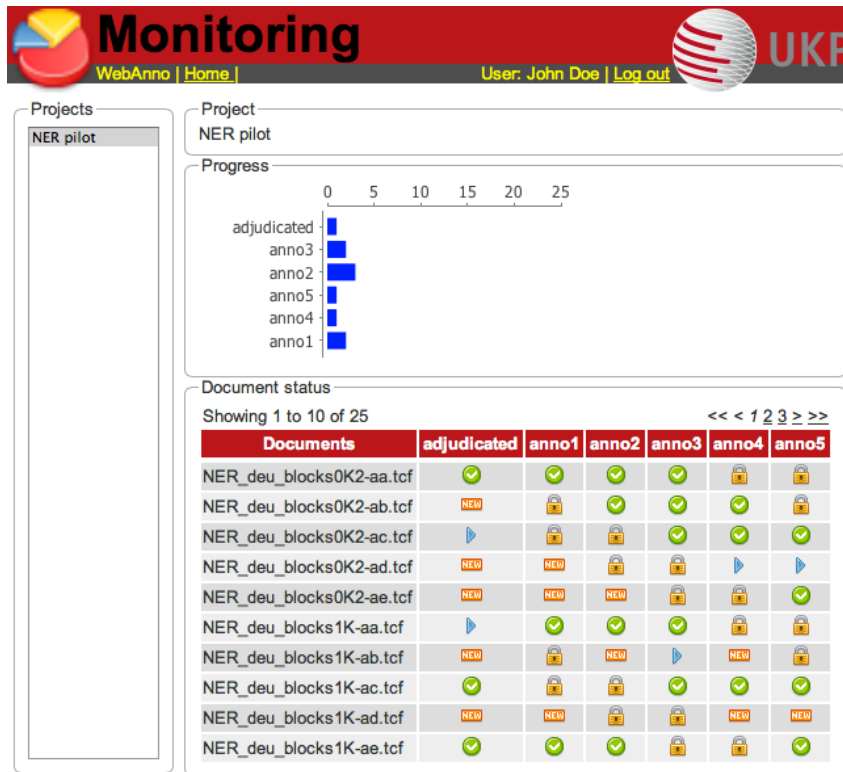


Figura 2.8: Interface do WebAnno para monitoração dos projetos. [Figura retirada de (23)]

Pode-se citar também ferramentas para anotações de propósito genérico, como a MMAX2 (24) ou a WordFreak (25), que como muitas outras, não são baseadas em *web* e não fornecem um ambiente para o gerenciamento do projeto de anotação, curadoria etc.

A análise das características destas ferramentas, seus pontos fortes e fracos, unida às demandas que a metodologia descrita em (6) levanta em termos de sistema, leva às seguintes conclusões sobre as características mais importantes em um sistema de preparação de dados e geração de modelos de NLP para reconhecimento de entidades e extração de informações com base em ontologias de domínio:

1. Deve ser uma ferramenta em ambiente *web*, para facilitação de acesso e centralização de dados;
2. Deve fornecer uma interface simples e direta para a anotação, onde pessoas não especializadas não tenham maiores dificuldades;
3. Deve fornecer um sistema robusto de gerenciamento de projetos de anotação;
4. Deve monitorar e registrar todas as ações de anotação, para uso destes dados em estudos posteriores e cálculos de métricas de anotação;
5. Deve permitir o carregamento direto de arquivos de ontologia (OWL, por exemplo) para seu uso como base de anotação tanto de entidades quanto de relações;

6. Deve permitir o carregamento de arquivos e documentos que guiem as anotações (*annotation guidelines*);
7. Deve permitir o uso direto, também em ambiente *web*, dos dados anotados e selecionados como entradas para um sistema de construção de modelos de aprendizado de máquina.

Com vistas a executar todos os passos da metodologia de (6) em um só ambiente, usando as melhores e mais pertinentes filosofias dos sistemas de anotação estudados, contando também com um sistema que use os dados anotados e forneça uma interface simples de construção dos modelos de aprendizado de máquina, os próximos capítulos mostram a proposta, implementação, uso e avaliação de um novo sistema, totalmente baseado em nuvem.

*It is not down on any map; true places never
are.*

Herman Melville, *Moby-Dick; or, The Whale*.

3

A plataforma

O objetivo do presente trabalho é, a partir da ideia descrita e testada em (6), desenvolver um sistema em nuvem, baseado em módulos de serviço, acessível via *browser*, que comporte em um único ambiente (do ponto de vista do usuário) todas as ferramentas necessárias para:

1. a manipulação e persistência de conjuntos de documentos textuais (fontes de dados);
2. o uso de bibliotecas de NLP com tokenizadores (*tokenizers*) e etiquetadores morfológicos (POS *taggers*) para a caracterização dos textos;
3. a anotação dos textos por diversos usuários cadastrados, com base em ontologias de domínio em formato OWL, carregáveis como *uploads* convencionais;
4. a curadoria dos documentos, por parte do “dono” do projeto, pela comparação das anotações e monitoração de parâmetros que indiquem os comportamentos dos anotadores;
5. o treinamento, a avaliação e a persistência de modelos de aprendizado de máquina para a tarefa de NER, que visam ao reconhecimento automático das entidades descritas nas ontologias;
6. o treinamento, a avaliação e a persistência de modelos de aprendizado de máquina para a tarefa de RE, que visam à construção automática de grafos de relações entre entidades, também baseada nas relações entre entidades descritas nas ontologias;
7. a construção e a disponibilização de serviços de leitura e estruturação, em formato RDF, de dados textuais com base em modelos de NER e RE gerados previamente.

As figuras 3.1, 3.2, 3.3 e 3.4 apresentam diagramas que descrevem alguns fluxos da plataforma, onde: as setas, ilustram a direção e fluxo de execução; os retângulos, os processos; as elipses, entradas e saídas; as nuvens, representam saídas instanciáveis.

Em um nível mais geral, como ilustrado na figura 3.1, o sistema deve ser composto por 3 grandes etapas: curadoria (gestão e preparação de dados), aprendizado de máquina (construção, treinamento e avaliação de modelos, tendo como base os dados da etapa anterior) e distribuição de serviços (compostos pelos modelos gerados na etapa anterior).

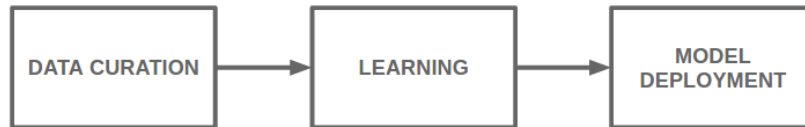


Figura 3.1: Fluxo geral de estruturação de dados

A figura 3.2 ilustra em maiores detalhes o fluxo da etapa de curadoria. Ela recebe dois *inputs* fundamentais: os dados textuais e a ontologia do domínio a ser modelado. O conjunto de textos alimentados passa por uma etapa de tokenização e POS *Tagging*, quando cada *token* separado é caracterizado quanto à sua função morfosintática. Os dados tokenizados são então disponibilizados em um ambiente de anotação, onde a ontologia alimentada serve de base para a disponibilização das *tags* (a partir de suas classes) e das relações (a partir de suas *object properties* e *datatype properties*). As anotações de todos os participantes são então validadas por um curador, cuja responsabilidade é comparará-las, corrigir os erros e definir o conjunto de anotações corretas, isto é, o *Gold Standard Annotation* (GSA). A GSA é, portanto, o conjunto de dados a ser usado na etapa de aprendizado de máquina.

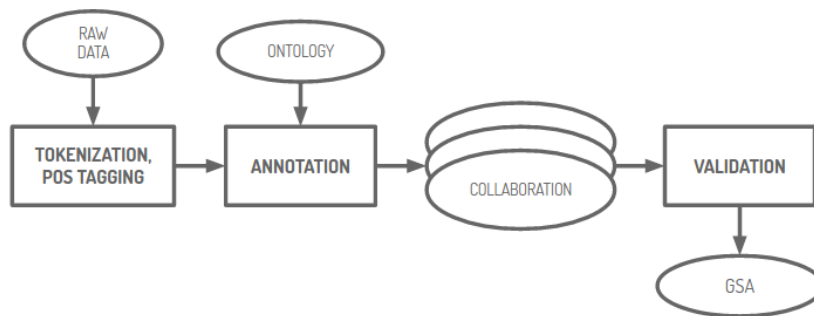


Figura 3.2: Fluxo de curadoria dos dados

A figura 3.3 ilustra o fluxo da etapa de aprendizado de máquina. Os dados GSA entram e alimentam uma etapa de engenharia de atributos, onde servem como base para a criação, automática e manual, de atributos para os classificadores. Os atributos são então aplicados sobre os dados GSA e os transformam em dados numéricos, isto é, o formato esperado pelos algoritmos de classificação. O aprendizado acontece, então, na etapa seguinte, onde os modelos treinados são treinados e avaliados, fornecendo, como saída, os modelos definitivos capazes de processar automaticamente novos dados textuais puros.

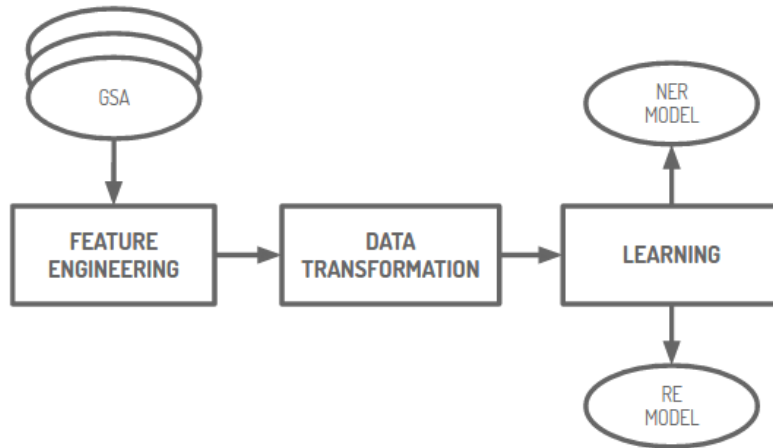


Figura 3.3: Fluxo de aprendizado automático

A terceira e última grande etapa, ilustrada em detalhes pela figura 3.4, recebe os modelos de NER e RE definidos na etapa anterior e são combinados em uma cadeia NER → RE. Essa cadeia de modelos é então disponibilizada como um serviço *online* que, recebendo dados textuais puros, identifica entidades e relações (pertencentes à ontologia usada) e exporta os resultados na forma de triplas RDF.

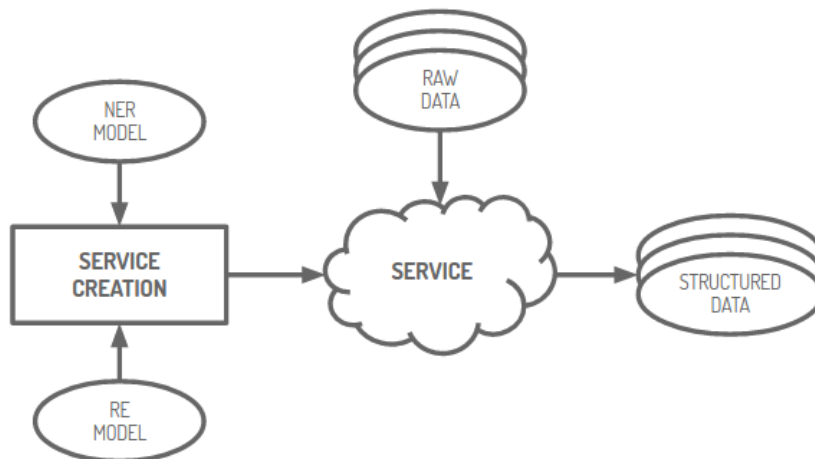


Figura 3.4: Fluxo de implantação dos modelos

Definidos os fluxos gerais, desenvolveu-se uma plataforma *web* denominada LER (*Learning Entities and Relations*) com módulos que reproduzem os fluxos descritos acima. Uma vez que a etapa de curadoria é uma das mais onerosas e importantes do fluxo (26), ela foi implementada como um subsistema do LER, sendo denominada ERAS (*Entities and Relations Annotation System*). Desse modo, a estrutura do sistema implementado por ser entendida pelo diagrama da figura 3.5.

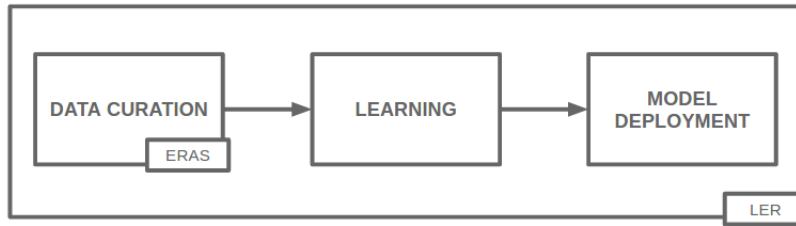


Figura 3.5: LER

3.1 Arquitetura

Para o desenvolvimento da plataforma, optou-se por um formato modularizado seguindo o modelo arquitetural REST, onde a comunicação entre os módulos se dá por meio do protocolo HTTP. O controle de acesso aos módulos da plataforma é feito por meio de tokens seguindo o padrão JWT (27), sendo estes transmitidos no cabeçalho das requisições entre os módulos. As principais tecnologias utilizadas para o desenvolvimento do LER foram:

- *Python*: Linguagem de programação utilizada para a construção do *backend*;
- *MongoDB*: SGBD (Sistema de Gerenciamento de Banco de Dados) utilizado para armazenamento dos dados gerados pela plataforma;
- *Scikit-learn*: Biblioteca utilizada para execução dos algoritmos de aprendizado automático;
- *PyStruct*: Biblioteca utilizada para execução dos algoritmos de aprendizado automático estruturado;
- *Javascript*: Linguagem de programação utilizada para a construção do *frontend*;
- *AngularJS*: Biblioteca responsável pelo controle dos componentes do *frontend* bem como as chamadas aos serviços do *backend*

Como pode ser visto na figura 3.6, a plataforma foi dividida em 5 módulos distintos:

- *Authentication*: Módulo responsável pelo controle de acesso dos usuários na plataforma, sendo o único módulo com acesso às credenciais destes e o responsável pela geração e verificação da sua validade;
- *Data*: Módulo responsável pelo controle de todos os dados da plataforma (com exceção dos dados de usuários, que estão no escopo do módulo *Authentication*);
- *Learning*: Módulo responsável pela execução de todos os métodos de aprendizado de automático da plataforma (*Tokenization*, *POS Tagging*, *NER* e *RE*);

- *Services*: Módulo responsável pelo gerenciamento dos serviços que recebem textos puros e geram dados estruturados em triplas RDF através do uso dos modelos de NER e RE criados no módulo *Learning*;
- *Client*: Módulo que provê uma interface web da plataforma para o usuário.

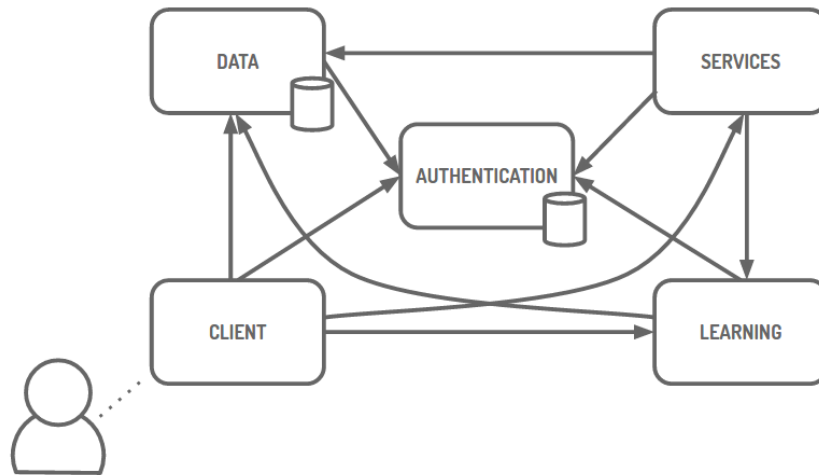


Figura 3.6: Arquitetura do LER

A vantagem desta modularização é a facilitação do escalonamento vertical e horizontal da plataforma em ambientes com alta demanda. Um escalonamento vertical, como exemplificado na figura 3.7, é apenas um *upgrade* nas configurações da máquina hospedeira da aplicação. Um aumento horizontal de escala pode ser obtido tanto separando os módulos em máquinas diferentes, podendo, por exemplo, deixar um módulo com maior demanda de processamento em uma máquina independente (ver figura 3.8), ou até mesmo criando um *cluster* de máquinas para aumentar o poder de processamento de quaisquer módulos (ver figura 3.9).

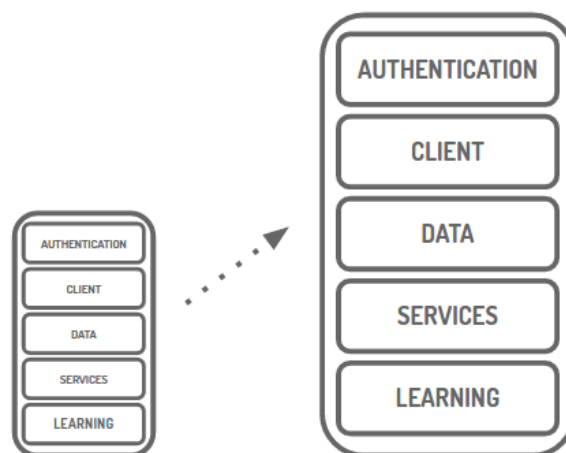


Figura 3.7: Exemplo de escalonamento vertical do LER

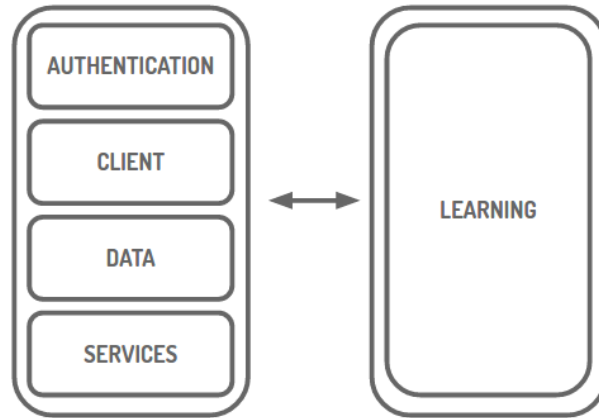


Figura 3.8: Exemplo de escalonamento horizontal do LER

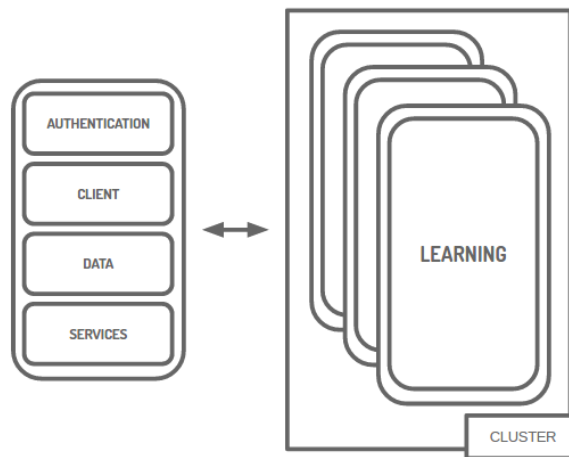


Figura 3.9: Exemplo de escalonamento horizontal do LER utilizando cluster

As principais funções da plataforma serão descritas nas próximas seções deste trabalho. Um exemplo do resultado final, no que diz respeito à interface de usuário, pode ser visto na figura 3.10. Um resumo do modelo de dados da plataforma na figura 3.11.

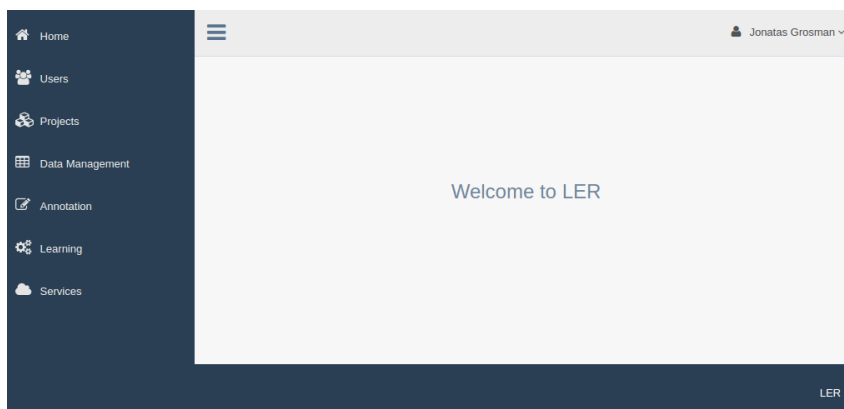


Figura 3.10: Tela inicial da plataforma

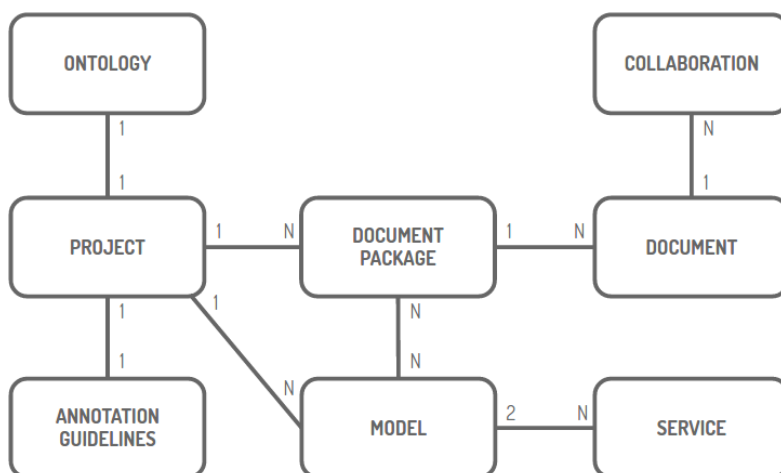


Figura 3.11: Resumo do modelo de dados do LER

PUC-Rio - Certificação Digital Nº 1421597/CA

3.2 Controle de usuários

O controle de usuários na plataforma é feito por meio de login e senha. Há 3 tipos distintos de usuários, com diferentes níveis de acesso à plataforma, como pode ser visto na figura 3.12. Os níveis de acesso disponíveis são:

- Administrador: Tem acesso a todas as áreas da plataforma;
- Usuário comum: Tem acesso a todas as áreas da plataforma com exceção da área de controle de usuários;
- Colaborador: Tem acesso apenas à tela inicial e à área de anotação.

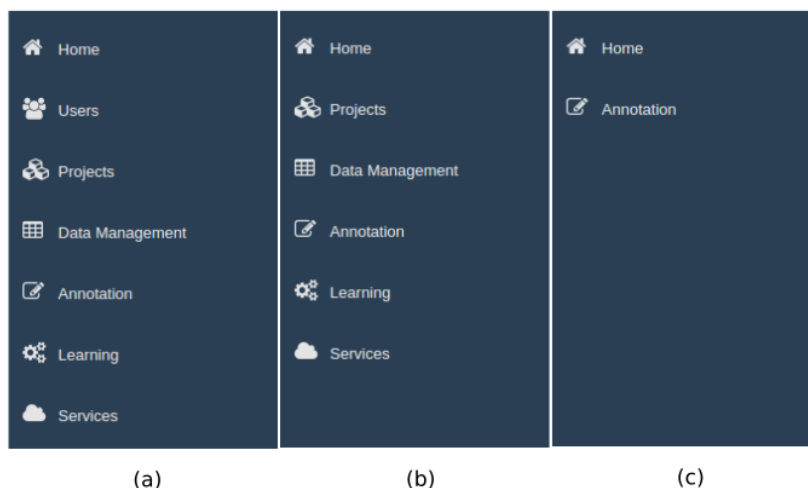


Figura 3.12: Menu da plataforma por níveis de acesso: (a) Administrador; (b) Usuário; (c) Colaborador

3.3

Gerência de projetos

A primeira etapa do fluxo de uso da plataforma trata da criação do projeto (ver figura 3.13), que congregará em si as tarefas de anotação e construção de modelos NER e RE. Ainda nesta etapa se definem, por *upload*, a ontologia no formato OWL (obrigatória) e o guia de anotação no formato PDF (opcional). É também obrigatória a escolha de um idioma para todo o projeto. Um elemento básico de diferenciação entre a presente implementação e a que foi adotada em (6) é justamente o POS *Tagger* utilizado. O LER adota o sistema Freeling 4.0 (28), explicado em resumo para uma versão anterior em (29). Esse POS *Tagger* apresenta algumas vantagens em relação à ferramenta F-EXT (11), que não está mais disponível para uso público e comercial, mas apenas para fins de pesquisa. Em primeiro lugar, o Freeling tem suporte para 12 idiomas além de português e inglês, únicos suportados pelo F-EXT. Além disso, o Freeling codifica a informação morfológica das POS *tags* tendo como base uma estrutura de etiquetas de tamanhos variáveis onde cada caractere corresponde a uma característica morfológica específica, sendo o primeiro sempre relativo à categoria (verbo, adjetivo etc.), de modo que um *token* que no F-EXT seria descrito simplesmente como um verbo, no Freeling em português, por exemplo, seria caracterizado em 6 níveis adicionais (ver figura 3.14). Essa abordagem multidimensional da função morfosintática do *token* possibilita a criação de mais dimensões nos atributos de POS, o que pode melhorar a capacidade de separação dos classificadores na etapa de aprendizado de máquina.

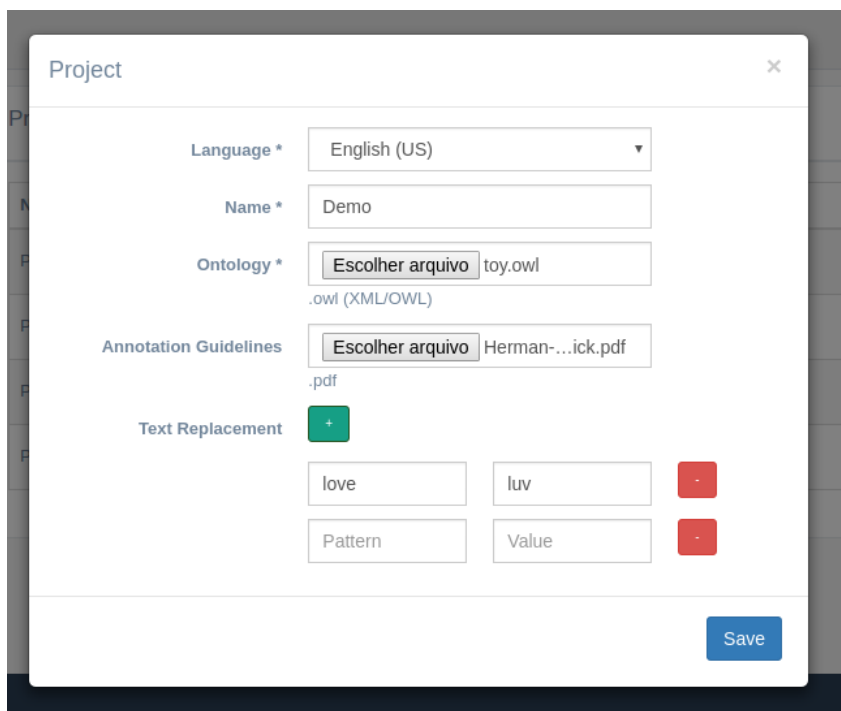


Figura 3.13: Criação de projetos

PUC-Rio - Certificação Digital Nº 1421597/CA

Position	Attribute	Values
0	category	<i>V:verb</i>
1	type	M:main; A:auxiliary; S:semiauxiliary
2	mood	<i>I:indicative; S:subjunctive; M:imperative; P:pastparticiple; G:gerund; N:infinitive</i>
3	tense	<i>P:present; I:imperfect; F:future; S:past; C:conditional; M:plusquamperfect</i>
4	person	1:1; 2:2; 3:3
5	num	S:singular; P:plural
6	gen	F:feminine; M:masculine; C:common; N:neuter

Figura 3.14: Descrição de verbos em português no Freeling 4.0

A figura 3.15 mostra um exemplo de tela com vários projetos criados, onde o administrador ou usuário comum pode editá-los ou excluí-los.

Name	Language	Ontology	Annotation Guidelines	
Project 1	pt-br	[Add] [Delete] [Edit]	[Add] [Delete] [Edit]	[Edit] [Delete]
Project 2	pt-br	[Add] [Delete] [Edit]	[Add]	[Edit] [Delete]
Project 3	pt-br	[Add] [Delete] [Edit]	[Add]	[Edit] [Delete]
Project 4	pt-br	[Add] [Delete] [Edit]	[Add]	[Edit] [Delete]
Demo	en-us	[Add] [Delete] [Edit]	[Add] [Delete] [Edit]	[Edit] [Delete]

Figura 3.15: Área de projetos

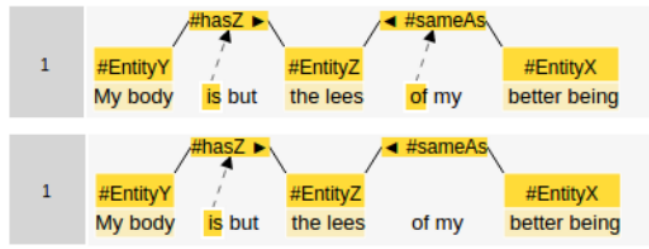
3.4 O ERAS

Como comentado anteriormente, a tarefa de anotação e curadoria de dados é uma das mais custosas em tempo no fluxo de estruturação de dados, merecendo um sistema dedicado especialmente a ela. As seções a seguir descreverão como o sistema ERAS lida este processo, nas perspectivas do curador e dos colaboradores.

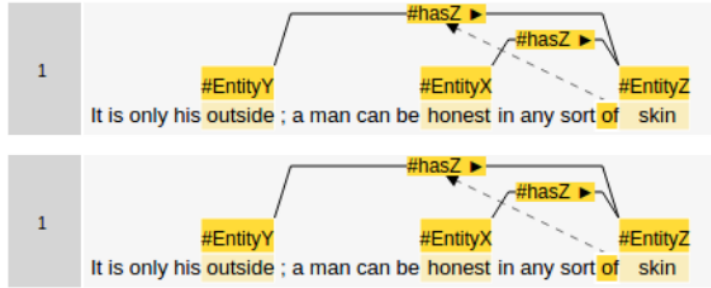
Com o intuito de tornar a descrição mais clara nos exemplos apresentados a seguir, será definido um cenário hipotético onde 2 usuários, A e B, realizaram a anotação e re-anotação das seguintes frases extraídas de (30):

- 01.txt: *My body is but the lees of my better being;*
- 02.txt: *It is only his outside; a man can be honest in any sort of skin.*

As figuras 3.16 e 3.17 mostram os resultados do cenário descrito acima, e as figuras 3.18 e 3.19 mostram as visões da ontologia que fôra utilizada para anotar os textos.

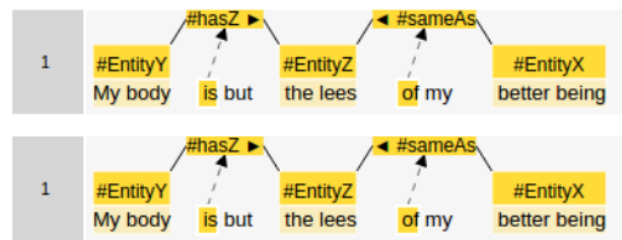


(a)

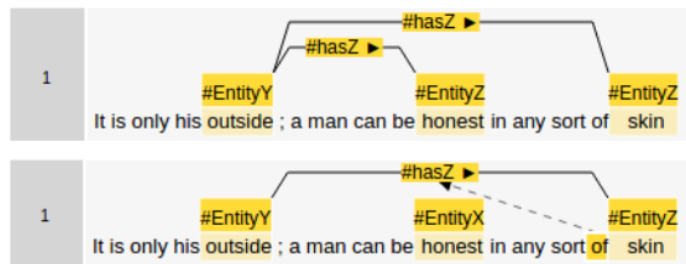


(b)

Figura 3.16: Anotações usuário A: (a) anotação do documento 01.txt seguido de sua re-anotação; (b) anotação do documento 02.txt seguido de sua re-anotação



(a)



(b)

Figura 3.17: Anotações usuário B: (a) anotação do documento 01.txt seguido de sua re-anotação; (b) anotação do documento 02.txt seguido de sua re-anotação

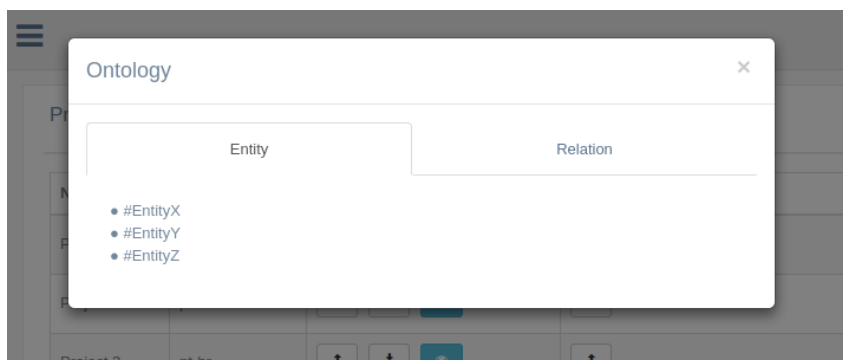


Figura 3.18: Ontologia: Entidades

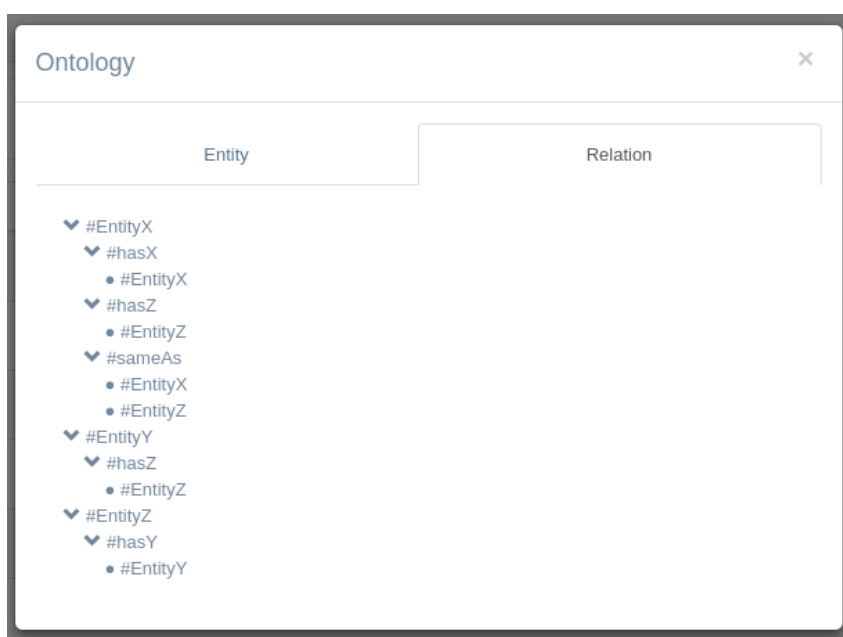


Figura 3.19: Ontologia: Relações

3.4.1

Gerência dos dados

Cada documento tem um estado dentre 3 possíveis: *UNCHECKED*, *PRECHECKED* e *CHECKED*. O estado *UNCHECKED* é o estado original de qualquer documento que é carregado no sistema e significa que ele está disponível para anotação. O estado *PRECHECKED* é um estado intermediário em que o documento não está disponível para anotação ou uso no aprendizado automático. O estado *CHECKED* significa que o documento faz parte do GSA a ser disponibilizado para uso na etapa de aprendizado de máquinas. A figura 3.20 mostra os documentos 01.txt e 02.txt com estados *PRECHECKED* e *CHECKED*, respectivamente.

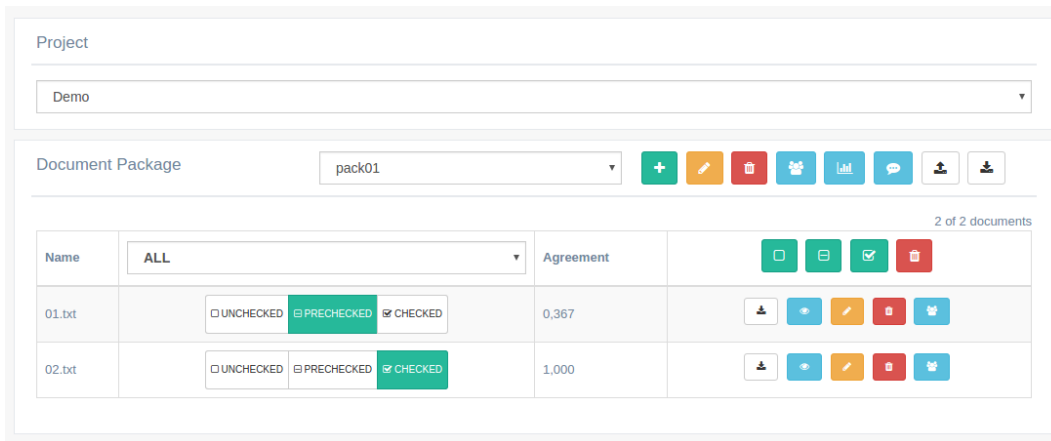


Figura 3.20: Área de gerência de dados

Cada documento precisa fazer parte de um determinado pacote. A criação dos pacotes também acontece na interface representada pela figura 3.20. A figura 3.21 apresenta os parâmetros da criação de um pacote. O *Precheck agreement threshold* é um parâmetro que indica o mínimo valor de concordância entre 2 anotadores diferentes sobre um mesmo documento para que este seja automaticamente colocado no estado *PRECHECKED*, isto é, sem a necessidade de que todos os anotadores daquele pacote precisem anotar o tal documento. Esse *threshold* visa agilizar o processo de anotação pela eliminação da necessidade de que todos os anotadores anotem todos os documentos do pacote. Os botões seguintes indicam quais itens entre *tag*, *relation* e *connector* são considerados para o cálculo da concordância para o *threshold*.

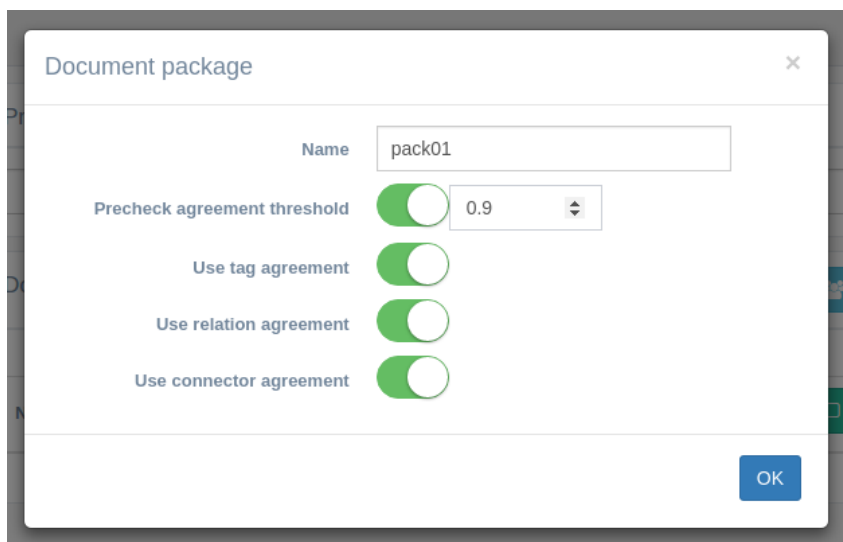


Figura 3.21: Criação de pacotes

A figura 3.22 apresenta a tela de gerenciamento dos colaboradores, onde o gestor do projeto constrói grupos aos quais cada anotador será alocado.

Um anotador só pode fazer parte de um único grupo. Nesta tela também são definidos os parâmetros de re-anotação *Warm up size*, que corresponde ao número de documentos a serem anotados para mero “aquecimento” sem registro, e *Reannotation step*, que indica o número de anotações finalizadas após o qual o documento mais antigo depois da última re-anotação será repetido para o cálculo da auto-concordância, isto é, concordância do anotador consigo mesmo.

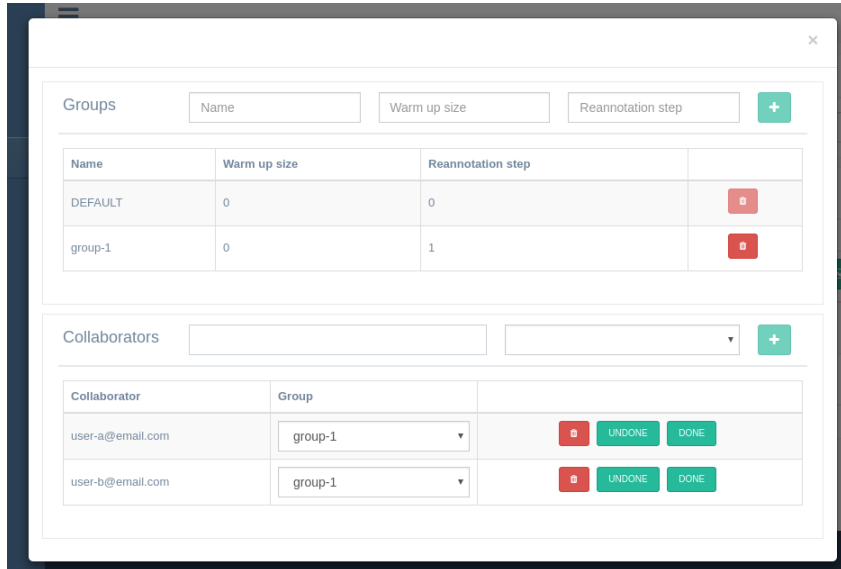


Figura 3.22: Gerenciamento de colaboradores

A figura 3.23 mostra uma interface onde o curador pode consultar todos os comentários de cada colaborador em cada documento. Esses comentários são colocados pelos colaboradores na interface de anotação.

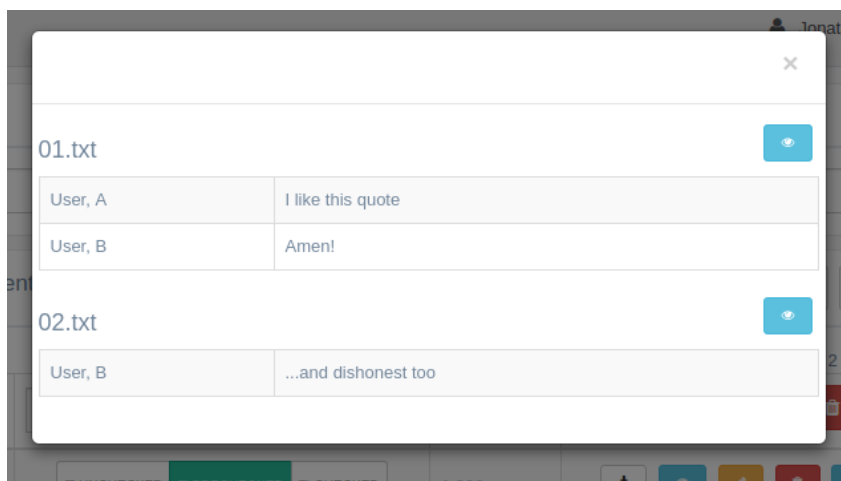


Figura 3.23: Comentários

A figura 3.24 apresenta um *menu* onde o curador pode selecionar qual perspectiva enxergar em determinado documento, podendo escolher o resul-

tado de uma anotador específico ou uma visão conjunta de todas as anotações realizadas, pela seleção de *Validation*.

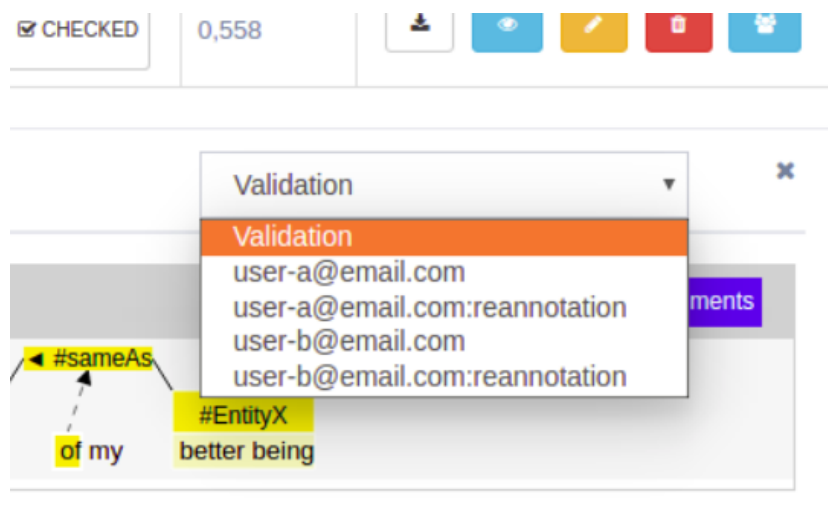
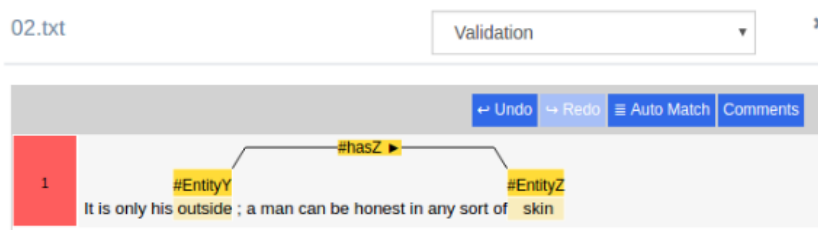
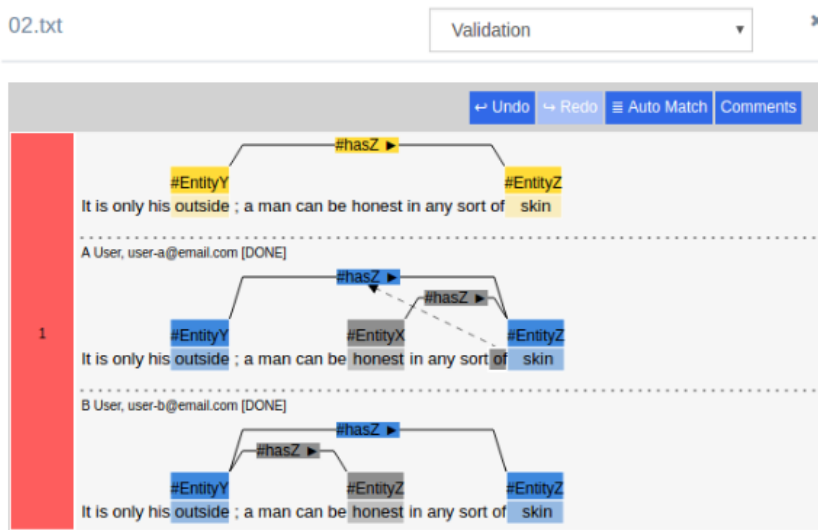


Figura 3.24: Perspectivas do documento

A figura 3.25 mostra o caso da seleção de *Validation* na figura 3.24. A cor avermelhada indica a existência de discordâncias entre os anotadores do documento em questão. O clique sobre esta região abre o quadro do documento colocando todas as anotações juntas, para que o curador possa, também via cliques e deleções, determinar a configuração final da anotação do tal documento. Há também um botão *Auto Match* que, se pressionado pelo curador, anotará o documento definitivo automaticamente com todas as *tags* e relações em que todos os anotadores concordaram.



(a)



(b)

Figura 3.25: Perspectiva de validação: (a) estado inicial do documento; (b) estado após expansão da sentença

A figura 3.26 mostra um exemplo de documento sem a região avermelhada, indicando a concordância total entre os anotadores.



Figura 3.26: Perspectiva de validação com concordância total: (a) estado inicial do documento; (b) estado após expansão da sentença

A figura 3.27 mostra o log de ações de cada anotador sobre determinado documento.

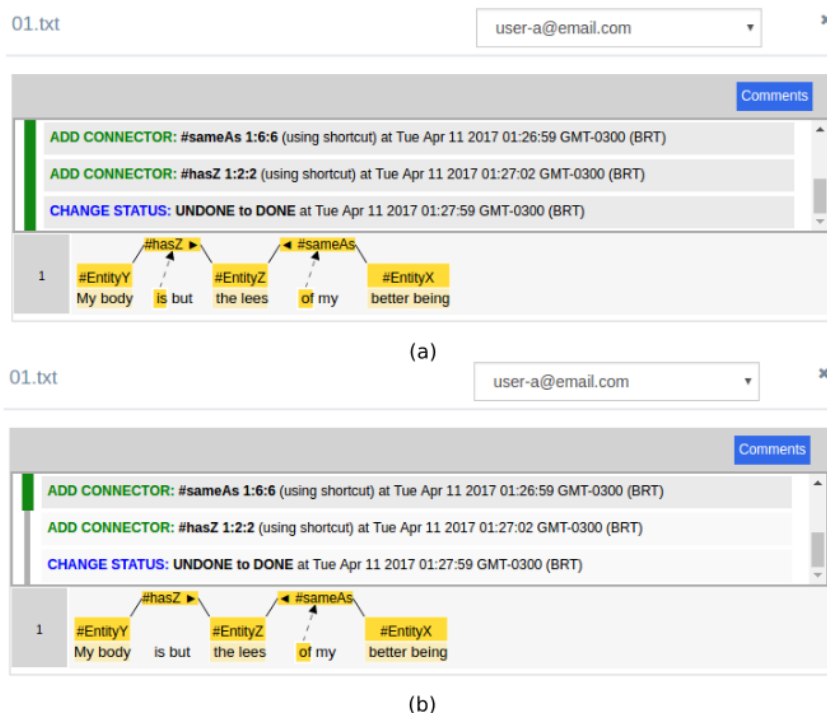


Figura 3.27: Perspectiva de ações do colaborador: (a) estado inicial do documento; (b) estado após desfazer algumas ações do colaborador

Por último, na perspectiva do curador, a figura 3.28 apresenta o formato do documento anotado, quando exportado por *download*. Este formato segue o padrão descrito em (31), com algumas adaptações.

```
#METADATA
author Herman Melville

#TEXT
My body is but the lees of my better being

#COLLABORATORS
user-a@email.com
user-b@email.com

#STATUS
PRECHECKED
DONE,R-DONE
DONE,R-DONE

#DESCRIPTION
ID FORM PROB POS LEMMA TAG RELATION CONNECTOR
0 My 1.0 PRP$ my B-#EntityY;B-#EntityY,R-B-#EntityY;B-#EntityY
1 body 1.0 NN body I-#EntityY;I-#EntityY,R-I-#EntityY;I-#EntityY
2 is 1.0 VBZ be O;O;O O;O;O B-0-0-0-4-#hasZ;B-0-0-0-4-#ha
3 but 0.992177 CC but O;O;O O;O;O O;O;O
4 the 1.0 DT the B-#EntityY;B-#EntityY,R-B-#EntityY;B-#EntityY

1491887273 DEFAULT,OPEN

#LOG-0
TIME ACTION
1491883829 DEFAULT,OPEN
1491883927 DEFAULT,ADD,TAG,0-8-8-#EntityX
1491883954 SHORTCUT,UNDO,ADD,TAG,0-8-8-#EntityX
REMOVE,TAG,0-8-8-#EntityX
1491883987 DEFAULT,ADD,TAG,0-8-9-#EntityX
1491884048 DEFAULT,OPEN
1491884083 DEFAULT,ADD,TAG,0-8-1-#EntityY
1491884173 DEFAULT,CHANGE,STATUS,UNDO-R-DONE

#COMMENTS-0
I like this quote

#COMMENTS-1
Amen!
```

Figura 3.28: Formato exportado

3.4.2 Anotação

Os colaboradores enxergam apenas a interface de anotação, ilustrada pela figura 3.29.

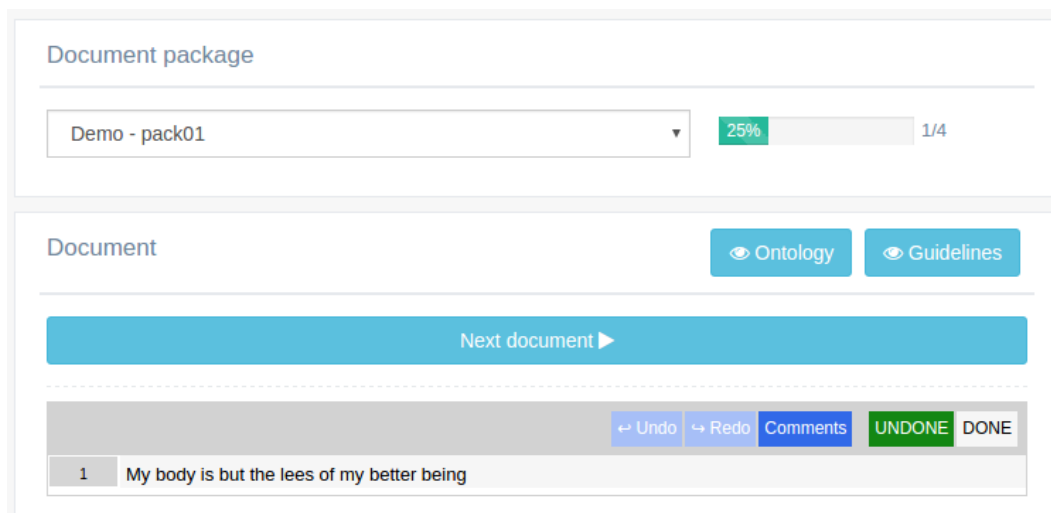


Figura 3.29: Tela de anotação

O colaborador pode consultar a estrutura da ontologia e o guia durante a anotação, como ilustrado pelas figuras 3.18, 3.19 e 3.30.

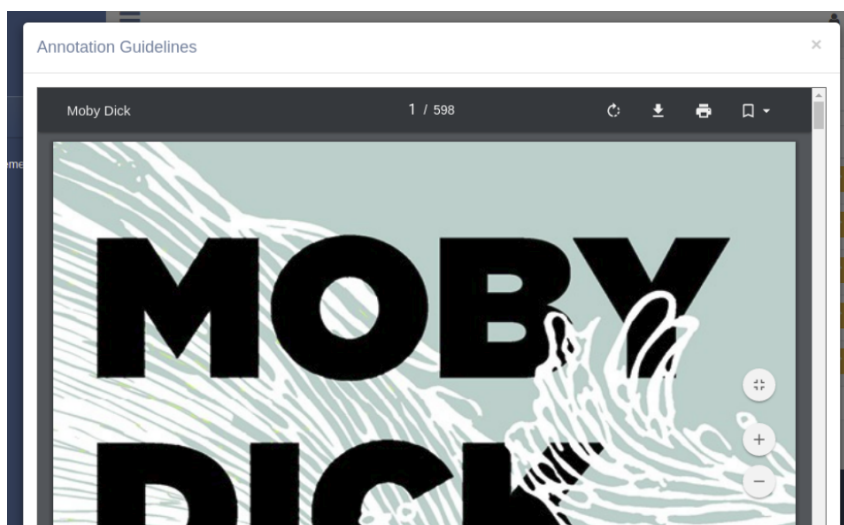


Figura 3.30: Guia de anotação

As figuras 3.31, 3.32, 3.33, 3.34, 3.35, 3.36, 3.37 e 3.38 apresentam as ações de anotação disponíveis para o colaborador.

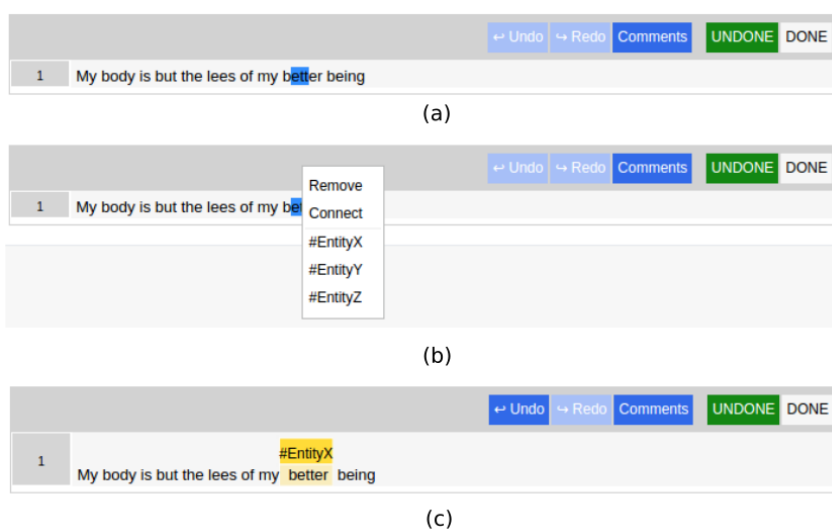


Figura 3.31: Processo de anotação de entidade: (a) área desejada selecionada; (b) acionando menu de contexto para seleção de entidade; (c) estado final da anotação



Figura 3.32: Processo de anotação de entidade com múltiplas palavras: (a) área desejada selecionada; (b) acionando menu de contexto para seleção de entidade; (c) estado final da anotação

PUC-Rio - Certificação Digital Nº 1421597/CA

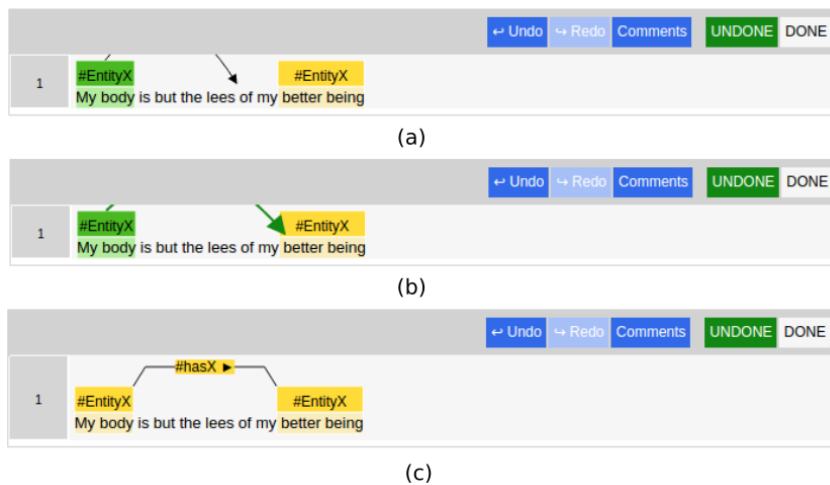


Figura 3.33: Processo de anotação de relação: (a) entidade origem selecionada; (b) entidade destino alcançada; (c) estado final da anotação



Figura 3.34: Exemplo de tentativa de anotação de relação inválida

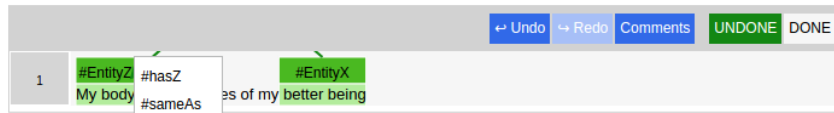
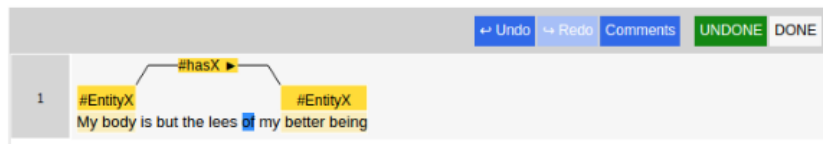
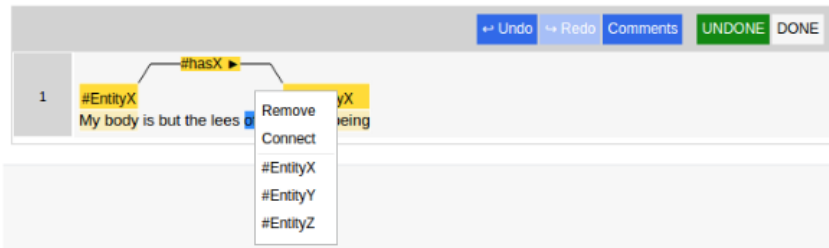


Figura 3.35: Exemplo de anotação de relação quando há múltiplas relações possíveis entre as entidades



(a)



(b)

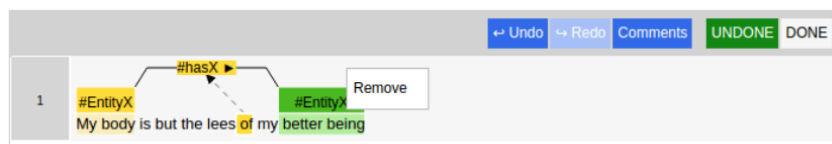


(c)



(d)

Figura 3.36: Processo de anotação de conector: (a) área desejada selecionada; (b) acionando menu de contexto para seleção de ação Connect; (c) relação desejada alcançada; (d) estado final da anotação



(a)



(b)

Figura 3.37: Processo de remoção de anotação: (a) menu de contexto acionado no elemento desejado para executar ação Remove; (b) estado final da remoção

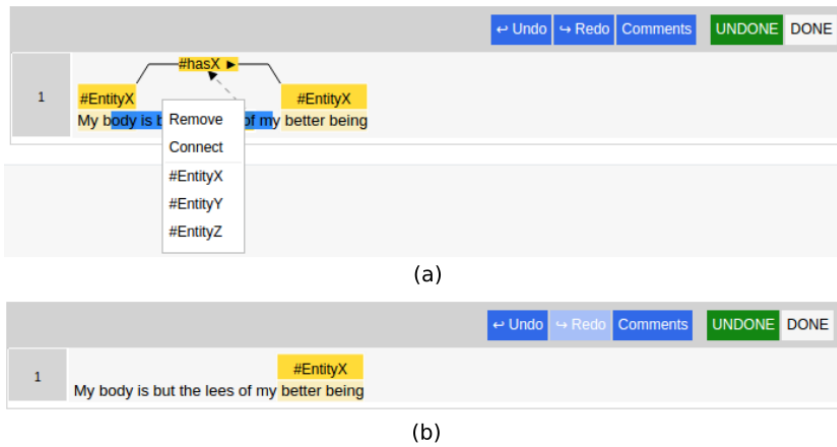


Figura 3.38: Processo de remoção de anotação em uma área: (a) menu de contexto acionado após selecionada a área desejada para executar ação Remove; (b) estado final da remoção

O colaborador pode adicionar comentários (ver figura 3.39) em cada documento, conforme já comentado na seção anterior.

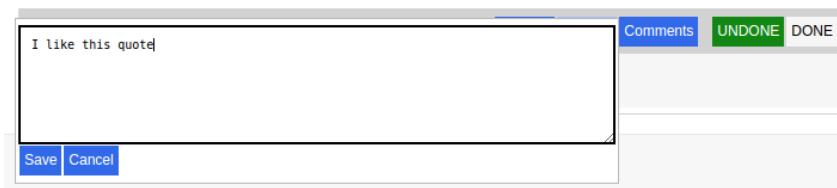


Figura 3.39: Adicionando comentários no documento

3.4.3 Estatísticas

O curador também tem acesso a estatísticas para cada pacotes, ilustradas nas figuras 3.40, 3.41, 3.42 e 3.43.

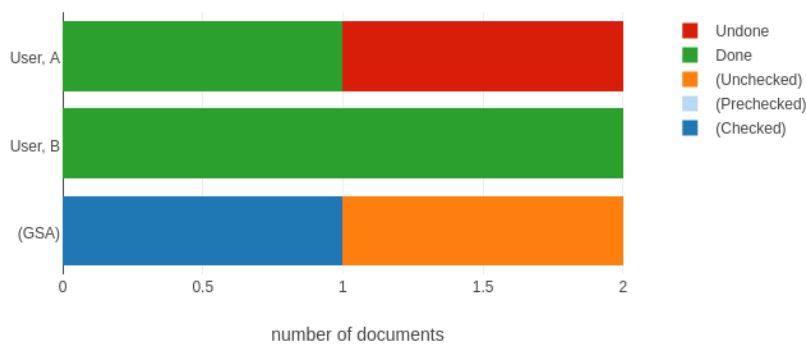


Figura 3.40: Estatísticas de status

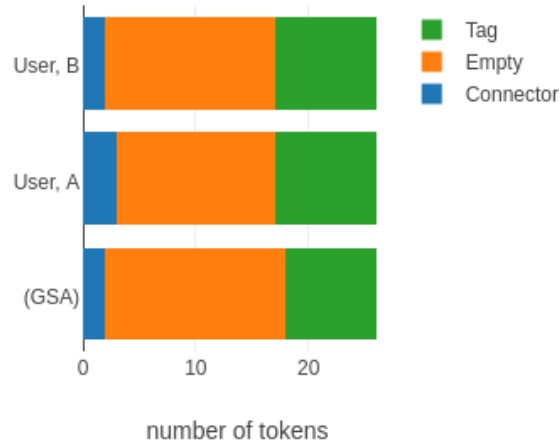


Figura 3.41: Estatísticas de cobertura dos tokens

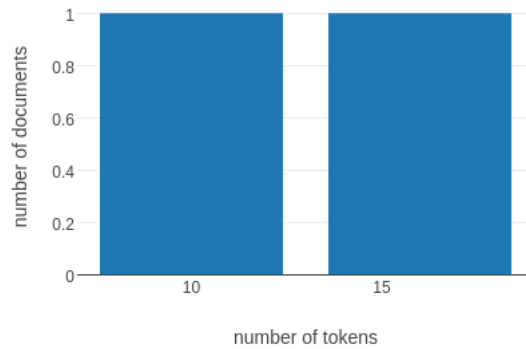


Figura 3.42: Estatísticas de distribuição dos tokens nos documentos

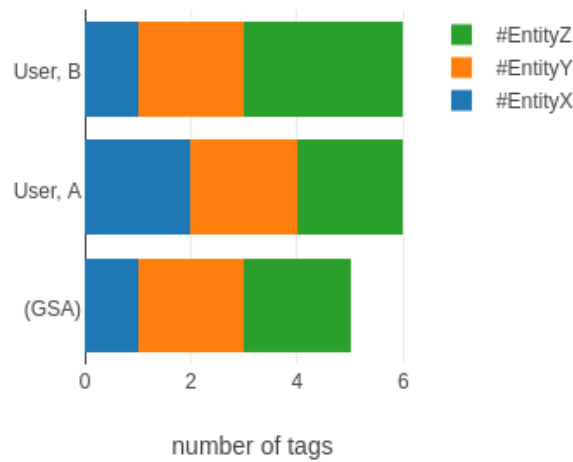


Figura 3.43: Estatísticas de entidades por colaborador, também disponíveis para relações e conectores.

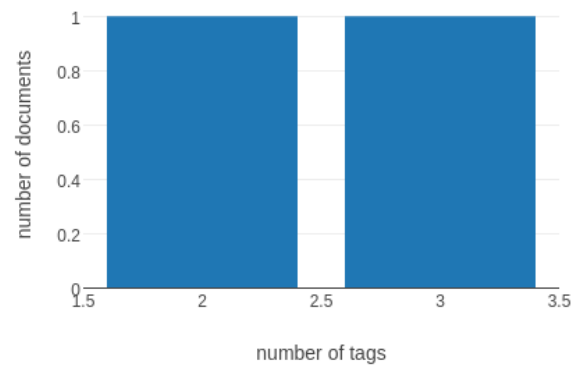


Figura 3.44: Estatísticas de distribuição das entidades no GSA, também disponíveis para relações e conectores.

Algumas funcionalidades importantes para a criação de atributos na etapa de aprendizado de máquina são ilustradas pelas figuras 3.45, 3.46 e 3.47. Essas visualizações estão disponíveis para as entidades, úteis para a tarefa NER, bem como para as relações, úteis para a tarefa RE.



Figura 3.45: Nuvem de palavras do GSA

Words	Occurrences
of	2
is	1

Figura 3.46: Tabela de ocorrências das palavras no GSA

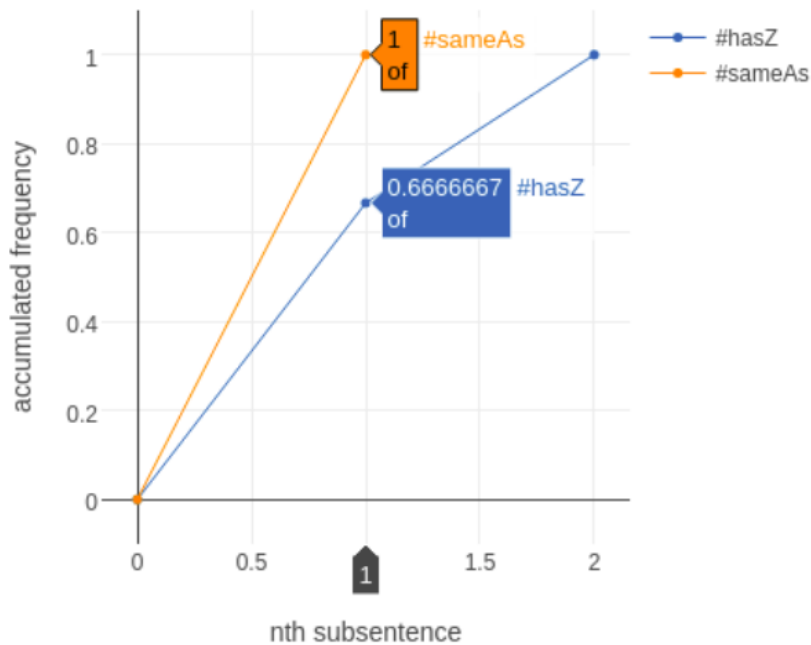


Figura 3.47: Curva de frequência acumulada das palavras no GSA

A concordância entre as anotações é calculada pelo índice kappa (32):

$$K = \frac{P_a + P_e}{1 - P_e} \tag{3-1}$$

Onde:

P_a — “Concordância observada”

P_e — “Probabilidade de concordância ao acaso”

O cálculo da probabilidade de concordância ao acaso é dado pela seguinte equação:

$$P_e = \frac{1}{|N|^2} \sum_{|K|} x_{k1} x_{k2} \tag{3-2}$$

Onde:

N — “Instâncias categorizadas”

K — “Categorias disponíveis”

x_{ki} — “Número de vezes que o observador i categorizou instâncias como sendo k ”

As instâncias que serão avaliadas para o cálculo da concordância, são criadas a partir de chaves únicas dadas a fragmentos das anotações, definidas segundos regras particulares para cada tipo de anotação:

Entidades — T:<identificador do documento>:<sentença do fragmento>:<índice do fragmento na sentença>, e.g., T:01.txt:0:0

Relações — R:<identificador do documento>:<sentença do fragmento origem>:<índice do fragmento origem na sentença>:<classe do fragmento origem>:<sentença do fragmento destino>:<índice do fragmento destino na sentença>:<classe do fragmento destino>, e.g., R:01.txt:0:0:#EntityY:0:4:#EntityZ

Conectores — C:<identificador do documento>:<sentença do fragmento origem>:<índice do fragmento origem na sentença>:<classe do fragmento origem>:<sentença do fragmento destino>:<índice do fragmento destino na sentença>:<classe do fragmento destino>:<sentença do fragmento>:<índice do fragmento na sentença>, e.g., C:01.txt:0:0:#EntityY:0:4:#EntityZ:0:2

A figura 3.48 mostra como são definidas as instâncias no caso das entidades. A figura 3.48(a) mostra o texto anotado pelo usuário, e a figura 3.48(b) a perspectiva que será levada em consideração para o cálculo da concordância. Os fragmentos são definidos *token a token* com o intuito de minimizar a penalização de pequenos erros na anotação.

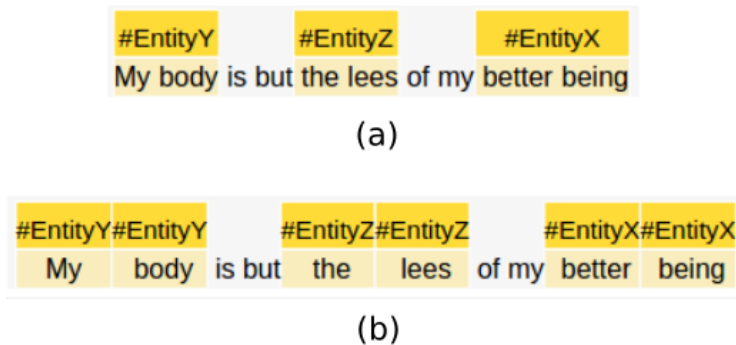
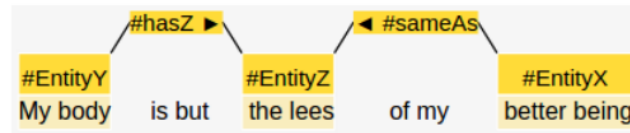
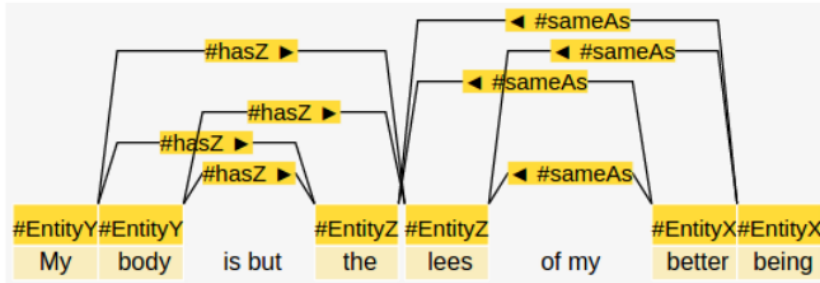


Figura 3.48: Definição de entidades para cálculo de concordância: (a) anotação do usuário; (b) perspectiva criada para cálculo de concordância

A figura 3.49 mostra como são definidas as instâncias no caso das relações. A figura 3.49(a) mostra o texto anotado pelo usuário, e a figura 3.49(b) a perspectiva que será levada em consideração para o cálculo da concordância. É possível notar nessa imagem a criação de novas relações. Isso se deve ao fato de ser levado em conta a mesma lógica de penalização de erros adotada para o cálculo de concordância nas entidades.



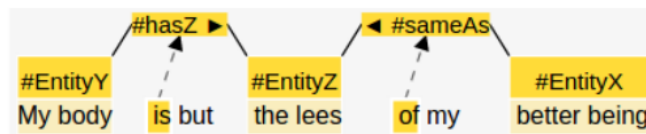
(a)



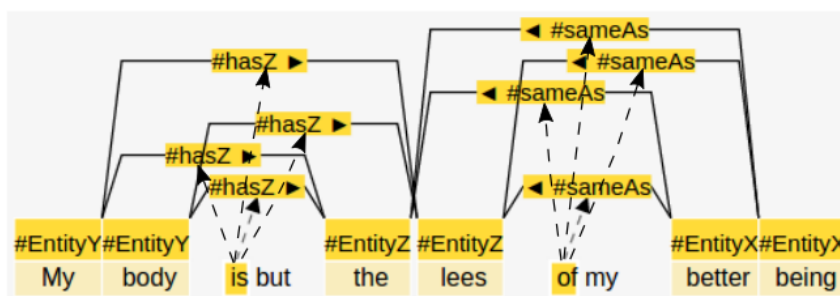
(b)

Figura 3.49: Definição de relações para cálculo de concordância: (a) anotação do usuário; (b) perspectiva criada para cálculo de concordância

A figura 3.50 mostra como são definidas as instâncias no caso dos conectores. A figura 3.50(a) mostra o texto anotado pelo usuário, e a figura 3.50(b) a perspectiva que será levada em consideração para o cálculo da concordância.



(a)



(b)

Figura 3.50: Definição de conectores para cálculo de concordância: (a) anotação do usuário; (b) perspectiva criada para cálculo de concordância

A classificação das instâncias definidas para o cálculo de concordância utiliza o tipo (T para entidades, R para relações e C para conectores) e rótulo dado. Na ausência de classificação por parte de um anotador a uma instância, um rótulo especial (EMPTY) é atribuído a esta. A figura 3.51 mostra como

ficaria uma tabela de classificação levando em consideração os exemplos de anotação mostrados nas figuras 3.16 e 3.17.

		USER-A										TOTAL
		T:#EntityX	T:#EntityY	T:#EntityZ	T:EMPTY	R:#hasZ	R:#sameAs	R:EMPTY	C:#hasZ	C:#sameAs	C:EMPTY	
USER-B	T:#EntityX	2	0	0	0	0	0	0	0	0	0	2
	T:#EntityY	0	3	0	0	0	0	0	0	0	0	3
	T:#EntityZ	1	0	3	0	0	0	0	0	0	0	4
	T:EMPTY	0	0	0	17	0	0	0	0	0	0	17
	R:#hasZ	0	0	0	0	5	0	1	0	0	0	6
	R:#sameAs	0	0	0	0	0	4	0	0	0	0	4
	R:EMPTY	0	0	0	0	1	0	0	0	0	0	1
	C:#hasZ	0	0	0	0	0	0	0	4	0	0	4
	C:#sameAs	0	0	0	0	0	0	0	0	4	0	4
	C:EMPTY	0	0	0	0	0	0	0	1	0	0	1
TOTAL	3	3	3	17	6	4	1	5	4	0	46	

Figura 3.51: Tabela de classificação dos usuários A e B para os exemplos dados no capítulo

Utilizando os dados mostrados na figura 3.51 é possível efetuar o cálculo dos valores de P_a e P_e para cada categoria, como pode ser visto na tabela 3.1.

Tabela 3.1: Concordâncias observadas e esperadas ao acaso no exemplo dado

	P_{ak}	P_{ek}
T:EntityX	0.0435	0.0028
T:EntityY	0.0652	0.0043
T:EntityZ	0.0652	0.0057
T:EMPTY	0.3696	0.1366
R:hasZ	0.1087	0.0170
R:sameAs	0.0870	0.0076
R:EMPTY	0.0000	0.0005
C:hasZ	0.0870	0.0095
C:sameAs	0.0870	0.0076
C:EMPTY	0.0000	0.0000
	0.9132	0.1916

Os valores finais para P_a e P_e no exemplo descrito são 0.9132 e 0.1916 respectivamente, o que leva a um valor de 0,8926 para o índice kappa.

Na figura 3.52 é possível ver os valores das concordâncias do modo que são visualizados na plataforma, no formato de um mapa de calor. Nota-se que além da concordância entre os anotadores A e B, há o cálculo da auto-concordância, essa comparação é feita utilizando os documentos re-anotados pelo próprio usuário caso o mesmo tenha um passo de re-anotação definido (ver seção 3.4.1), outro cálculo também feito é a concordância entre o anotador e a anotação ideal (que no exemplo aqui mostrado será o documento resultante das anotações que coincidiram entre ambos anotadores). No ERAS esses valores são calculados apenas com base nos documentos em estado *CHECKED*.

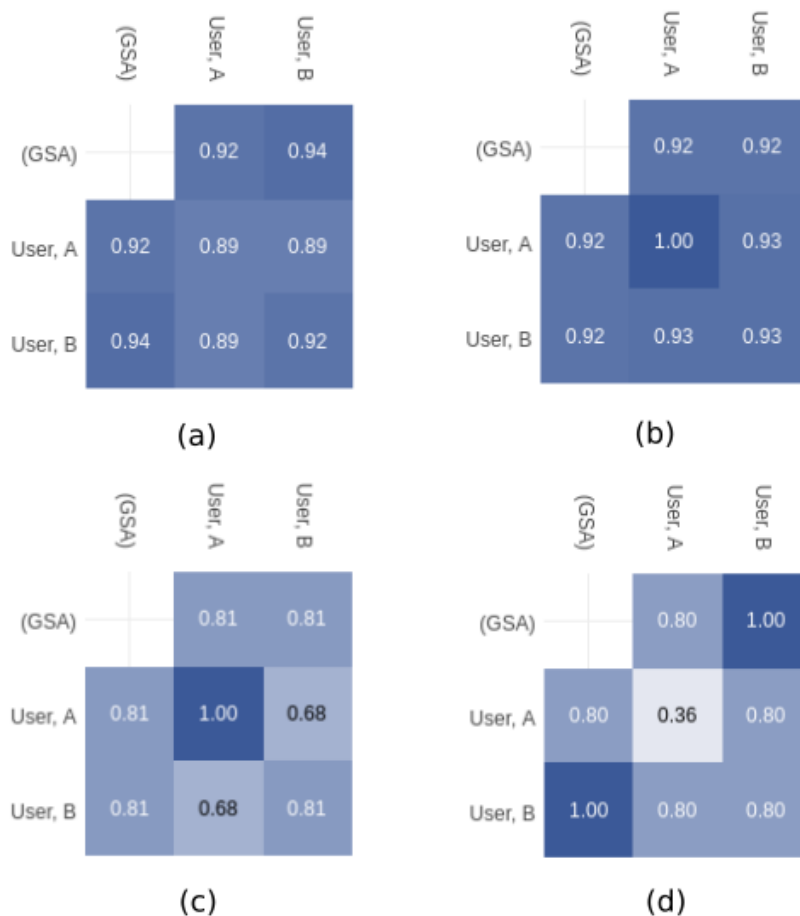


Figura 3.52: Mapa de concordância para o exemplo dado no capítulo: (a) entidades, relações e conectores; (b) apenas entidades; (c) apenas relações; (d) apenas conectores

Uma característica interessante presente no ERAS e, até onde pesquisado, não percebida em outros sistemas de anotação, é ilustrada pelas figuras 3.53 e 3.54. Trata-se da visualização da evolução temporal das concordâncias dos anotadores em relação ao GSA e em relação a eles mesmos.

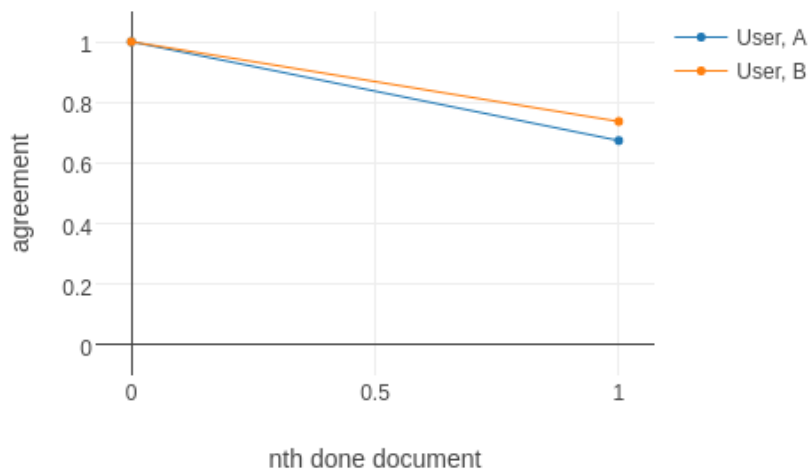


Figura 3.53: Curva de concordância com o GSA ao longo do tempo

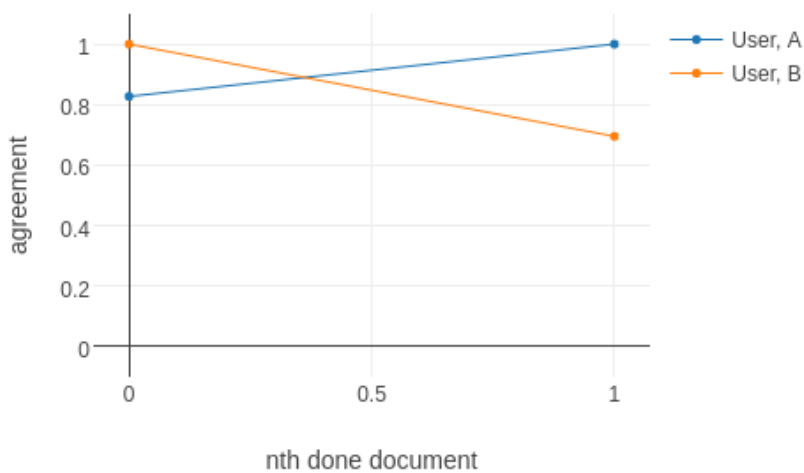


Figura 3.54: Curva de auto-concordância ao longo do tempo

Por fim, o curador tem uma visualização conjunta das estatísticas gerais de cada anotador, como ilustrado pela figura 3.55.

Collaborations ↓

Collaborator	Time spent	Time spent AVG	Time spent SD	# Views	# Views AVG	# Views SD
User, B	1m 37s	48s	1s	2,00	1,00	0,00
User, A	19m 32s	9m 46s	7m 50s	16,00	8,00	0,00

Documents ↓

Name	# Tokens	Time spent	Time spent AVG	Time spent SD	# Views	# Views AVG	# Views SD	
02.txt	13	2m 46s	1m 23s	33s	9	4,50	3,50	
01.txt	11	18m 23s	9m 11s	8m 24s	9	4,50	3,50	

Figura 3.55: Estatísticas básicas sobre colaboradores e documentos

3.5

Aprendizado automático

O módulo de aprendizado automático possibilita a criação de tarefas para o treinamento e o teste de classificadores para NER e RE. A figura 3.56 mostra o formulário principal de configuração da tarefa *General*, onde se define o seu nome, tipo e o projeto de onde consumir os dados GSA.

The screenshot shows a web interface with a 'General' tab selected. The form contains the following fields:

- Name:** demo-NER
- Type:** Named Entity Recognition
- Project:** Demo

A blue 'Save' button is positioned at the bottom right of the form area.

Figura 3.56: Configurações gerais das tarefas

A figura 3.57 mostra o formulário de configuração dos dados a serem ou não usados no treinamento e teste do classificador da tarefa. Uma barra superior indica a quantidade e o percentual de documentos para treino e teste. A parte inferior mostra uma árvore representativa da ontologia usada, onde se pode generalizar algumas classes pelas classes que as contêm ou mesmo excluir algumas entidades.

Package Name	# checked documents	DO NOT USE	TRAINING	TEST
pack01	2	DO NOT USE	TRAINING	TEST
pack02	2	DO NOT USE	TRAINING	TEST
pack03	0	DO NOT USE	TRAINING	TEST

Class	KEEP	GENERALIZE	REMOVE
#EntityX	KEEP	GENERALIZE	REMOVE
#EntityY	KEEP	GENERALIZE	REMOVE
#EntityZ	KEEP	GENERALIZE	REMOVE

Save

Figura 3.57: Configurações dos dados das tarefas

O formulário de atributos para o classificador, ilustrado pelas figuras 3.58 e 3.59, apresenta diversas opções para NER e RE. Estes atributos são baseados nos descritos por (6). Para NER, há 6 tipos disponíveis:

- FORM($X, step, regex$): aplica uma expressão regular ($regex$) sobre o conteúdo completo do parâmetro FORM (conforme resultado do POS *Tagger*) do *token* localizado a um número de passos ($step$, com valor nulo, positivo ou negativo) do *token* X . (Note-se que, se $step$ é nulo, o atributo representará a aplicação da $regex$ sobre a representação FORM do próprio *token* X .);
- LEMMA($X, step, regex$): o mesmo conceito de FORM($X, step, regex$), porém aplicado sobre a representação LEMMA;
- POS($X, step$): o POS *tag* do *token* localizado a $step$ passos do *token* X ;
- PROB($X, step$): o valor da probabilidade, segundo o POS *Tagger*, do conteúdo de POS do *token* localizado a $step$ passos do *token* X ;
- RANGE-FORM($X, step, regex$): aplica uma expressão regular sobre o conteúdo da *string* formada pela concatenação – com espaços – dos parâmetros FORM dos *tokens* em um *range* de $step$ a partir do *token* X ;
- RANGE-LEMMA($X, step, regex$): o mesmo conceito de RANGE-FORM($X, step, regex$), porém aplicado sobre a representação LEMMA.

Há uma diferença entre NER e RE no tocante aos tipos de instância a serem avaliados. Em NER, avalia-se *token a token*. Em RE, avalia-se as relações entre “nós”, isto é, entre entidades nomeadas segundo a ontologia, presentes no texto. São 12 os atributos disponíveis:

- RANGE-FORM($X, step, regex$): o mesmo conceito usado no NER, considerando, porém, que o parâmetro *step* se baseia em um dos nós da relação X , dependendo do sinal: se *step* é negativo, aplica-se à esquerda da relação; se positivo, à direita desta. Valores nulos de *step* não implementam esta *feature*;
- RANGE-LEMMA($X, step, regex$): o mesmo conceito de RANGE-FORM, descrito no item anterior, aplicado sobre LEMMA;
- POSITIONAL-FORM($X, step, regex$): o mesmo conceito usado no NER, considerando, porém, que o parâmetro *step* se baseia em um dos nós da relação X , dependendo do sinal: se *step* é negativo, aplica-se à esquerda da relação; se positivo, à direita desta. Valores nulos de *step* não implementam esta *feature*;
- POSITIONAL-LEMMA($X, step, regex$): o mesmo conceito de POSITIONAL-FORM, descrito no item anterior, aplicado sobre LEMMA;
- POS($X, step$): o mesmo conceito usado no NER, considerando, porém, que o parâmetro *step* se baseia em um dos nós da relação X , dependendo do sinal: se *step* é negativo, aplica-se à esquerda da relação; se positivo, à direita desta. Valores nulos de *step* não implementam esta *feature*;
- INTERIOR-RANGE-FORM($X, regex$): aplica uma expressão regular sobre o conteúdo da *string* formada pela concatenação, com espaços, dos parâmetros FORM dos *tokens* entre os nós da relação X ;
- INTERIOR-RANGE-LEMMA($X, regex$): o mesmo conceito de INTERIOR-RANGE-FORM, descrito no item anterior, aplicado sobre LEMMA;
- NODE-TO-NODE-DISTANCE(X): calcula a distância, em *tokens*, entre os nós da relação X . Para fins de redução da magnitude desse atributo a uma ordem de grandeza próxima dos outros, que são binários, o valor é dividido pela máxima quantidade de *tokens* encontrada nos documentos de treino do modelo;
- NODE-TO-NODE-DISTANCE-WITH-SIGNAL(X): mesmo conceito do atributo descrito acima, porém possibilitando valores negativos. A ideia é que este atributo, além de representar a distância entre nós, carregue também a informação do sentido da relação: números negativos traduzem relações da direita para a esquerda e, positivos, da esquerda para a direita;
- CLASS-NODE-FROM(X): classe de NER do nó de saída da relação X ;
- CLASS-NODE-TO(X): classe de NER do nó de chegada da relação X ;
- POSSIBLE-RELATION(X): indica se a relação é ou não possível, com base nos dados do conjunto de treinamento.

Há também a opção de configurar a leitura das POS *Tags* em dimensões ou níveis específicos, como já comentado anteriormente. Uma das características mais importantes está no botão *Auto generate*, que, com base nos parâmetros da parte superior do formulário, gera atributos automaticamente. Essa geração automática cria expressões regulares pelo uso das tabelas de frequências acumuladas de palavras ligadas às entidades e às relações (conectores).

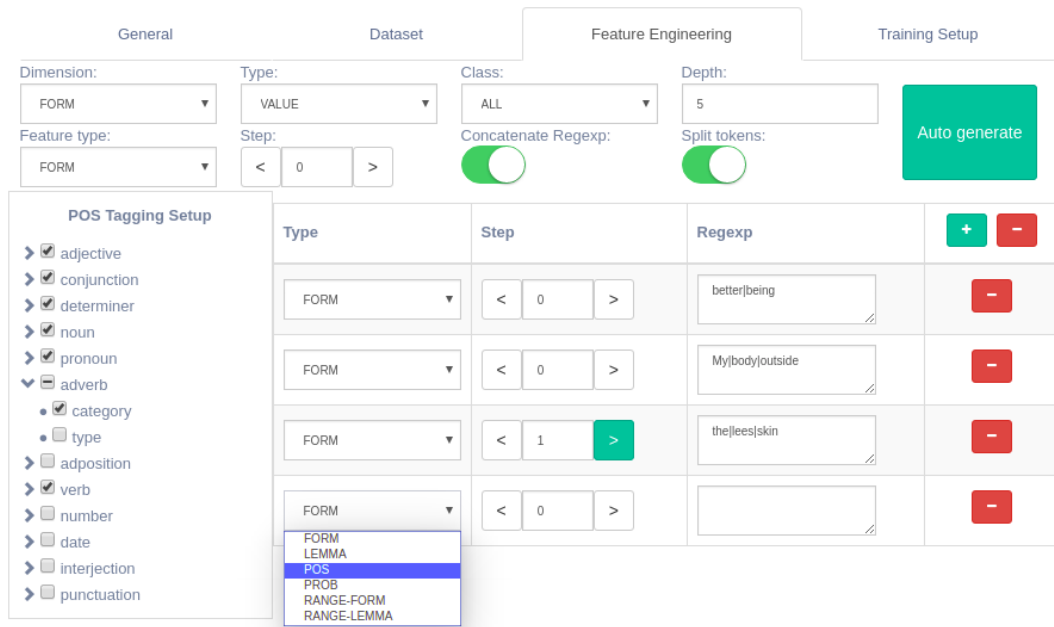


Figura 3.58: Engenharia de atributos da tarefas

PUC-Rio - Certificação Digital Nº 1421597/CA

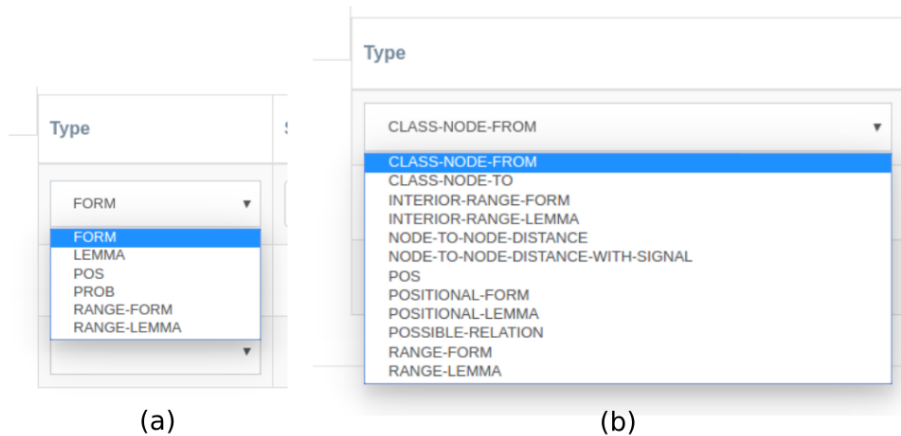
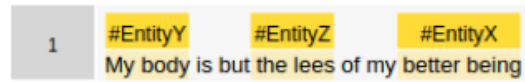


Figura 3.59: Atributos disponíveis: (a) NER; (b) RE

A partir dos atributos definidos para a tarefa, uma matriz numérica é criada com base nos dados de treino, um exemplo do formato desta matriz para a tarefa de NER pode ser vista na figura 3.60. Nota-se que cada amostra

de aprendizado equivale a um *token* do documento. Outro ponto a se observar na figura é que a classificação do *token* (coluna Y) não equivale apenas ao que este representa no domínio, mas também se o mesmo participa da descrição no início do bloco anotado ou não. Deve-se observar também que a coluna Y da referida imagem não está no formato numérico com intuito de facilitar a explicação desta, já que as classes devem ser codificadas para valores numéricos antes de submetida aos algoritmos de aprendizado automático.



(a)

	FORM(X, 0, "t")	FORM(X, -1, "of")	RANGE-FORM(X, -3, "is")	POS(X, 0)[verb]	Y
My	0	0	0	0	B-#EntityY
body	0	0	0	0	I-#EntityY
is	0	0	0	1	O
but	1	0	1	0	O
the	1	0	1	0	B-#EntityZ
lees	0	0	1	0	I-#EntityZ
of	0	0	0	0	O
my	0	1	0	0	O
better	1	0	0	0	B-#EntityX
being	0	0	0	1	I-#EntityX

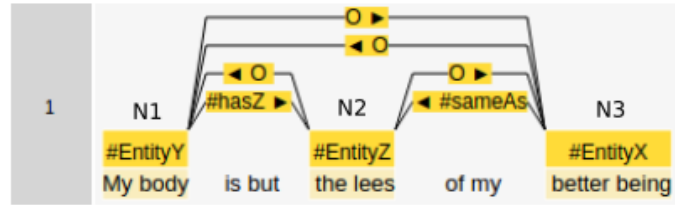
(b)

Figura 3.60: Exemplo de atributos para a tarefa de NER: (a) documento de entrada; (b) exemplo de uma matriz criada utilizando o documento de entrada

Como na figura acima, a figura 3.61 apresenta um exemplo do formato da matriz de treino, só que desta vez para a tarefa de RE. É possível notar nesta figura uma diferença crucial entre as instâncias das tarefas, que é o fato de cada amostra não ser mais apenas um *token*, porém um documento inteiro, sendo este descrito na forma de um grafo completo, onde os nós são as entidades anotadas, e o que será classificado na tarefa é a existência de uma aresta e seu rótulo. Esta definição do rótulo de uma aresta pelo modelo é uma melhoria ao que foi feito por (6), onde apenas era feita a classificação da existência das arestas, ou seja, o problema era tratado como uma classificação binária, sendo a rotulação feita por meio de regras pré-definidas para o domínio. Outra melhoria feita foi o fato de a abordagem utilizada no LER poder retornar grafos desconexos e até mesmo aceitar ciclos nestes. Na proposta anterior apenas uma árvore era dada como saída do classificador.



(a)



(b)

	NODE-TO-NODE-DISTANCE(X)	INTERIOR-RANGE-FORM(X, "is")	Y
N1 → N2	4	1	#hasZ
N1 → N3	8	1	O
N2 → N1	4	1	O
N2 → N3	4	0	O
N3 → N1	8	1	O
N3 → N2	4	0	#sameAs

(c)

	FEATURES	Y
01.txt	[[4, 1], [8, 1], [4, 1], [4, 0], [8, 1], [4, 0]]	[#hasZ, O, O, O, O, #sameAs]
02.txt	[...]	[...]

(d)

Figura 3.61: Exemplo de atributos para a tarefa de RE: (a) documento de entrada; (b) perspectiva do documento de entrada criada para geração da matriz de atributos; (c) exemplo simplificado de uma matriz criada utilizando o documento de entrada; (d) exemplo de uma matriz criada utilizando o documento de entrada

O último menu da criação da tarefa, *Training Setup*, trata das configurações de *gridsearch* (33) a serem executadas na tarefa, para a busca dos melhores parâmetros em relação a 4 diferentes índices: Accuracy; Precision; Recall; F1. Os algoritmos de aprendizado automático disponíveis na plataforma são os mesmos utilizados por (6) com acréscimo de alguns outros. Para a tarefa de NER: Support Vector Classification (34); Random Forest (35); Linear Stochastic Gradient Descent (36). Para a tarefa de RE: Structured Perceptron (37); Frank-Wolfe SSVM (38). Tais algoritmos foram adicionados apenas com o intuito de demonstrar que dada a forma como o LER fora implementando, é possível o uso de vários algoritmos para as tarefas de NER e RE. Espera-se que com a evolução da ferramenta e o uso desta em diversos domínios diferentes, mais algoritmos sejam incluídos.

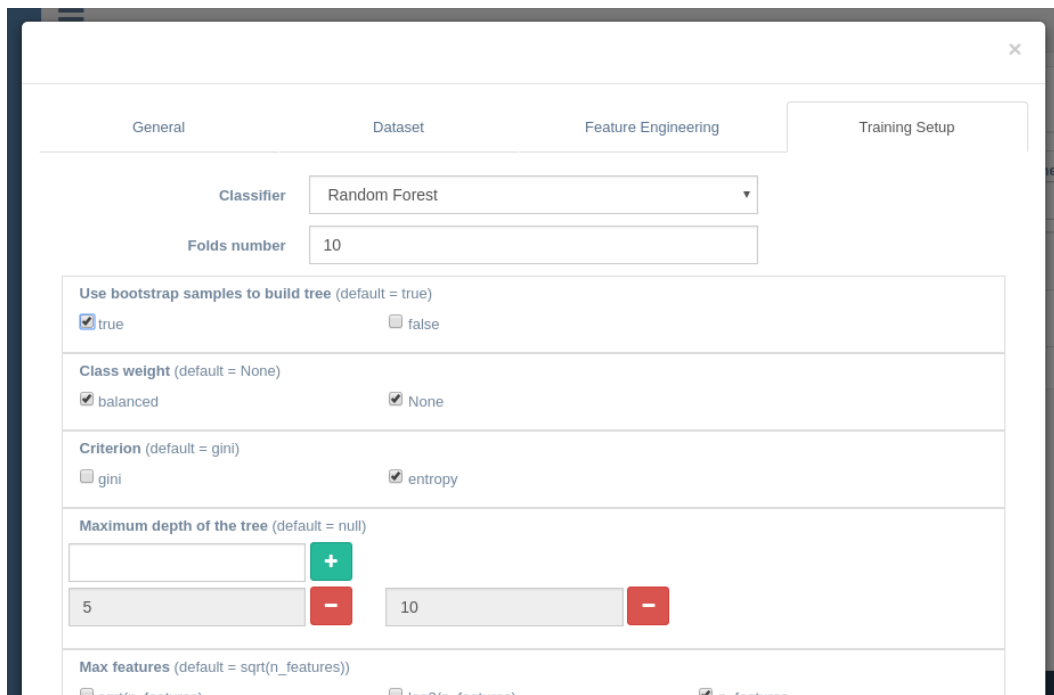








Figura 3.62: Configurações dos classificadores





Criadas as tarefas, elas podem ser iniciadas ou encerradas em modo paralelo e independente, como ilustrado pela figura 3.63.

Name	Type	Project	Status	Task to c
demo-NER	NER	Demo	off	   

(a)

Name	Type	Project	Status	Task to c
demo-NER	NER	Demo	running	 9s 

(b)

Name	Type	Project	Status	Task to c
demo-NER	NER	Demo	running	 10s 
demo-NER-clone	NER	Demo	running	 7s 

(c)

Figura 3.63: Execução de tarefas: (a) estado inicial da tarefa; (b) tarefa em execução; (c) tarefas executando em paralelo

Durante a execução, ou mesmo em caso de encerramento com sucesso ou erro, é possível visualizar o log de cada tarefa (ver figura 3.64), que apresenta em formato de texto as principais informações sobre seu processo de treinamento e teste.

```

INFO [2017-04-11 02:48:12,756]: Starting logger.
INFO [2017-04-11 02:48:12,879]: Master Thread configuration DONE.
INFO [2017-04-11 02:48:12,968]: Starting Master Thread.
INFO [2017-04-11 02:48:13,097]: Slave Process started.
INFO [2017-04-11 02:48:13,290]: Master sent data through PIPE conn to Slave.
INFO [2017-04-11 02:48:13,388]: Master sent msg to Slave to receive data.
INFO [2017-04-11 02:48:13,388]: Starting NER learning function.
INFO [2017-04-11 02:48:13,507]: Slave received data.
INFO [2017-04-11 02:48:13,507]: Collecting training documents
INFO [2017-04-11 02:48:14,117]: Collecting test documents
INFO [2017-04-11 02:48:14,532]: Building features
INFO [2017-04-11 02:48:14,639]: Building tag setup
    
```


Figura 3.64: Log de execução da tarefa

Por fim, em caso de rodada completa da tarefa e sem erros, é possível ler os resultados numéricos dos melhores parâmetros aplicados, em termos 4 índices alvo, como ilustrado pela figura 3.65 e 3.66.

Task elapsed time: 43s

- model-recall-demo-NER >
- model-precision-demo-NER >
- model-fscore-demo-NER >
- model-accuracy-demo-NER v

Classifier: Random Forest
Grid search metric: accuracy
Training elapsed CPU time: 0.0889750000000047 seconds

Best Parameters 

Parameter	Value
Max features	n_features
Class weight	None

- Training scores >
- Test scores >

Figura 3.65: Resultados da tarefa

Test scores v



Class	Support	Accuracy	Recall	Precision	F-score
O	18	0.884	1	0.857	0.923
B-#EntityX	1	0.961	1	0.5	0.666
B-#EntityY	2	0.961	1	0.666	0.8
B-#EntityZ	2	0.923	0	0	0
I-#EntityX	1	0.961	0	0	0
I-#EntityY	1	0.961	0	0	0
I-#EntityZ	1	0.961	0	0	0
Average, SD		0.945, 0.028	0.428, 0.494	0.289, 0.347	0.341, 0.400

Figura 3.66: Scores do modelo treinado

O sistema LER avalia os resultados de todos os modelos – de NER e RE – e escolhe os melhores parâmetros para cada um dos clássicos índices *Accuracy*, *Precision*, *Recall* e *F1* (39):

$$Accuracy = \frac{tp + tn}{tp + tn + fp + fn} \quad (3-3)$$

$$Precision = \frac{tp}{tp + fp} \quad (3-4)$$

$$Recall = \frac{tp}{tp + fn} \quad (3-5)$$

$$F1 = 2 \cdot \frac{Precision \cdot Recall}{Precision + Recall} \quad (3-6)$$

Sendo:

tp — “true positive”

fp — “false positive”

tn — “true negative”

fn — “false negative”

3.6

Publicação de serviços

Após todo o processo de curadoria e treinamento dos modelos, estes podem ser finalmente implantados para o uso na tarefa de estruturação dos dados textuais. No LER isso é feito por meio de serviços que disponibilizarão uma interface utilizando o protocolo HTTP. A figura 3.67 mostra a tela de criação de tais serviços, onde o usuário escolhe dentre todos os modelos treinados para determinado projeto, qual modelo ficará encarregado para o reconhecimento de entidades e qual será o modelo responsável pela extração de relações, bem como a definição de mapeamentos do texto classificado para propriedades das entidades.

Figura 3.67: Criação do serviço

Assim que o serviço é criado, ficará ativo e esperando requisições. A figura 3.68 mostra a tela que lista os serviços ativos na plataforma, nesta mesma imagem é possível notar a presença de um *token* de autenticação associado ao serviço, esta é a forma como as requisições a este serão controladas, apenas quem estiver de posse desse *token* poderá ter suas requisições aceitas.

Name	Project	NER model	RE model	Token	URL
demo-dervice	Demo	model-precision-demo-NER	model-fscore-demo-RE	013b8764-1e9c-4f9e-a31b-ca59a3b287d5	http://127.0.0.1:50003/services/58e572cf3249b0480fef1301/predict?token=013b8764-1e9c-4f9e-a31b-ca59a3b287d5

Figura 3.68: Serviços disponíveis

A resposta dada pelo serviço é estruturada em dois formatos, JSON (40) e triplas RDF (41), as figuras 3.69 e 3.70 mostram exemplos de retornos dados pelos serviços, onde a primeira mostra o resultado do teste que o próprio LER disponibiliza em sua interface, e a segunda o resultado do acesso aos serviços utilizando uma aplicação externa.

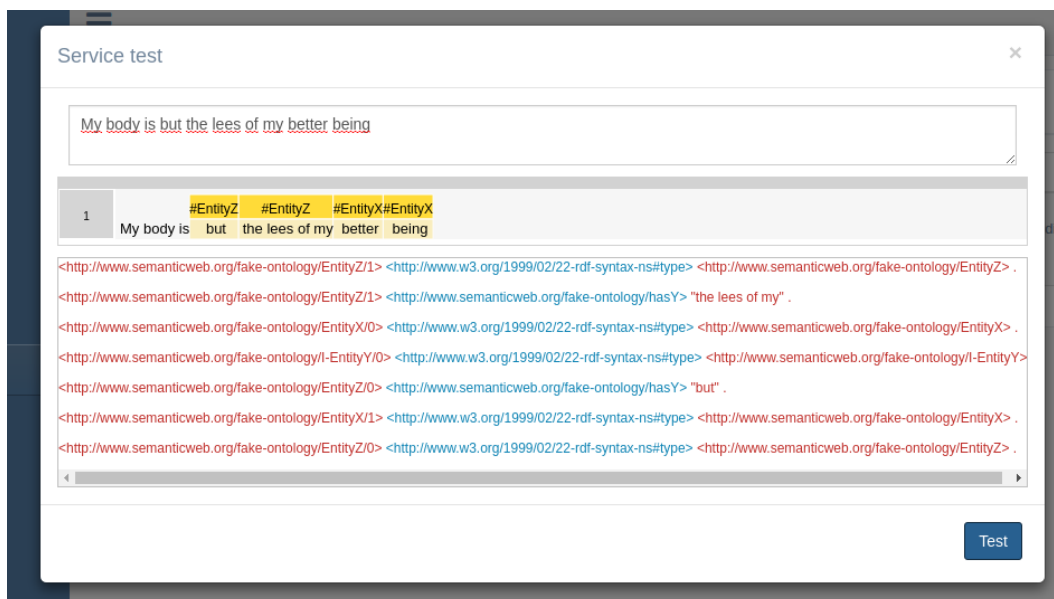


Figura 3.69: Teste do serviço

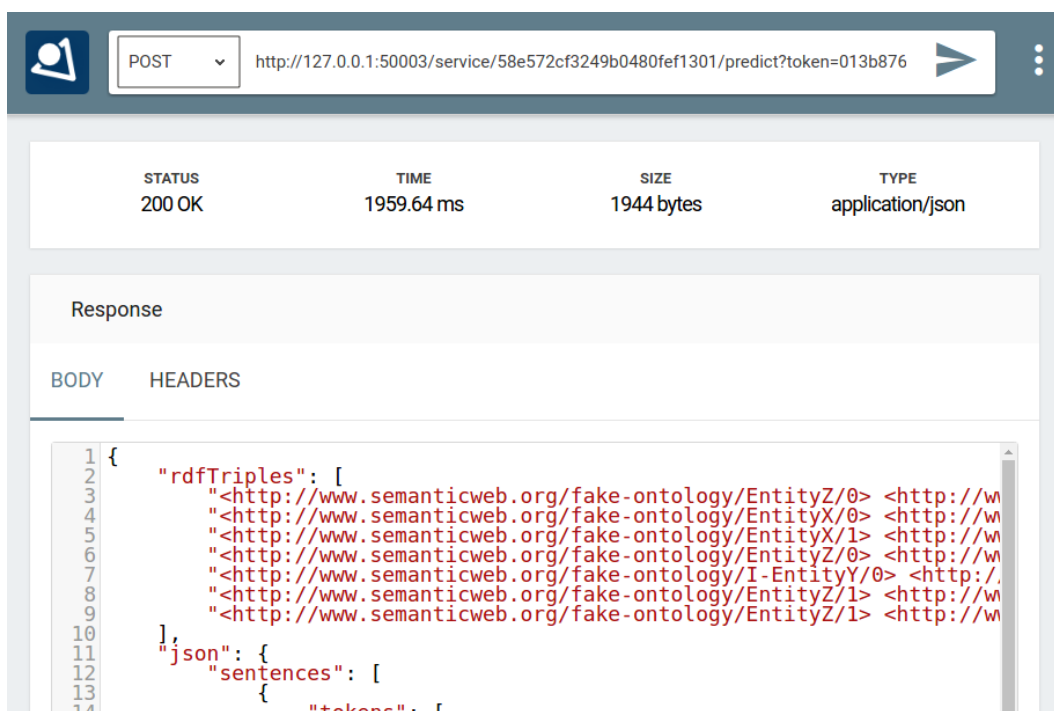


Figura 3.70: Teste externo do serviço

I try all things, I achieve what I can.

Herman Melville, *Moby-Dick; or, The Whale*.

4 Experimentos

Este capítulo apresenta as descrições e os resultados de dois tipos de experimentos, o primeiro relativo a uma avaliação do ERAS com a participação de um grupo de anotadores recrutados e o segundo relativo ao funcionamento do LER sobre dados de *tweets* de trânsito anotados de forma especializada, tendo uma versão modificada da TEDO como base.

4.1 Experimento de anotação

Um experimento de anotação foi conduzido com o intuito de avaliar o comportamento da plataforma para a tarefa de curadoria de dados, tarefa esta de domínio do subsistema ERAS. Optou-se por executar a tarefa de curadoria de anotações sobre o mesmo domínio de dados avaliados por (6), que é o de *tweets* de eventos de trânsito.

As próximas seções descrevem a metodologia empregada no experimento, seguida pelos resultados acompanhados de uma análise.

4.1.1 Metodologia

Foram coletados 100 *tweets* de eventos de trânsito, 50 da conta @operacoesRio e 50 da conta @odia24horas. Em seguida, foram separados 25 *tweets* de cada conta com o objetivo de criar o guia de anotação (o resultado pode ser visto no Apêndice A) e, também, realizar adaptações na ontologia TEDO com base nas dificuldades encontradas ao utilizar a mesma para anotar os textos. A Figura 4.1 mostra a ontologia resultante.



Figura 4.1: Visões da ontologia usada na anotação dos dados: (a) entidades; (b) relações

O restante dos *tweets* foi anotado por 20 participantes de forma a validar a efetividade do guia de anotação e a consistência da ontologia, com base nas anotações realizadas, concordâncias obtidas, comentários feitos nos documentos (ver Apêndice B) e respostas dadas ao questionário distribuído no final do experimento (ver Apêndice C). Os participantes foram separados em seis grupos com configurações distintas de aquecimento e passo de re-anotação. Além disso, os participantes foram separados em dois tipos. Uma parte deles anotou baseando-se apenas no guia de anotação – tipo A. A outra parte – tipo B – teve acesso ao guia de anotação e, também, receberam uma explicação presencial sobre a tarefa e tiraram dúvidas ao longo da anotação.

4.1.2 Resultados

A análise sobre a concordância dos participantes foi separada em duas etapas. Primeiramente avaliou-se a diferença de concordância entre os grupos. Num segundo momento, avaliou-se a diferença de concordância entre os diferentes tipos de participantes. A concordância foi medida de duas formas, primeiramente tomando como referência o GSA (que foi criado pelo autor deste trabalho) e depois comparando as anotações dos participantes de cada grupo entre eles mesmos que é a auto-concordância. Para cada agrupamento, foi calculado o valor médio de concordância acompanhado do seu desvio-padrão.

Os resultados da concordância entre grupos são apresentados nas Tabelas 4.1, 4.2, 4.3 e 4.4. Notou-se que a tarefa de identificação de relações parece mais

complexa, pois a média de concordância de todos os grupos é consideravelmente inferior quando comparada a tarefa de identificação de entidades. A tarefa de identificação de conectores também teve média de concordância baixa, o que já era esperado, pois os conectores têm dependência direta com as relações. Com os resultados obtidos, não é possível afirmar que a configuração de aquecimento e passo de re-anotação tenham alguma influência na melhoria da concordância, seja com o GSA, como também a própria auto-concordância.

Tabela 4.1: Concordância com entidades, relações e conectores nos grupos

Grupo (aquecimento, passo de re-anotação)	Concordância com o GSA	Auto-concordância
(0,5)	0.572 (σ 0.135)	0.767 (σ 0.116)
(2,5)	0.530 (σ 0.036)	0.800 (σ 0.125)
(4,5)	0.497 (σ 0.166)	0.777 (σ 0.133)
(0,10)	0.417 (σ 0.189)	0.480 (σ 0.228)
(2,10)	0.543 (σ 0.107)	0.773 (σ 0.023)
(4,10)	0.520 (σ 0.087)	0.720 (σ 0.066)

Tabela 4.2: Concordância com entidades nos grupos

Grupo (aquecimento, passo de re-anotação)	Concordância com o GSA	Auto-concordância
(0,5)	0.835 (σ 0.077)	0.918 (σ 0.066)
(2,5)	0.820 (σ 0.070)	0.887 (σ 0.055)
(4,5)	0.713 (σ 0.204)	0.847 (σ 0.110)
(0,10)	0.628 (σ 0.247)	0.643 (σ 0.321)
(2,10)	0.817 (σ 0.058)	0.900 (σ 0.020)
(4,10)	0.760 (σ 0.092)	0.790 (σ 0.085)

Tabela 4.3: Concordância com relações nos grupos

Grupo (aquecimento, passo de re-anotação)	Concordância com o GSA	Auto-concordância
(0,5)	0.273 (σ 0.278)	0.438 (σ 0.293)
(2,5)	0.243 (σ 0.137)	0.623 (σ 0.200)
(4,5)	0.143 (σ 0.315)	0.413 (σ 0.412)
(0,10)	0.040 (σ 0.268)	0.007 (σ 0.439)
(2,10)	0.170 (σ 0.243)	0.473 (σ 0.367)
(4,10)	0.150 (σ 0.185)	0.493 (σ 0.190)

Tabela 4.4: Concordância com conectores nos grupos

Grupo (aquecimento, passo de re-anotação)	Concordância com o GSA	Auto-concordância
(0,5)	0.088 (σ 0.260)	0.525 (σ 0.205)
(2,5)	0.020 (σ 0.026)	0.673 (σ 0.320)
(4,5)	0.097 (σ 0.167)	0.920 (σ 0.139)
(0,10)	0.060 (σ 0.081)	0.415 (σ 0.582)
(2,10)	0.123 (σ 0.150)	0.590 (σ 0.373)
(4,10)	0.093 (σ 0.075)	0.363 (σ 0.335)

Os resultados da concordância entre tipos de participantes são apresentados nas Tabelas 4.5, 4.6, 4.7 e 4.8. Como na concordância entre grupos, notou-se que as tarefas de identificação de relações e de conectores parecem mais complexas que a de identificação de entidades. Observou-se, porém, que há uma diferença significativa na concordância dos grupos. O tipo B, que recebeu uma

explicação presencial sobre a tarefa e teve dúvidas sanadas ao longo da anotação, teve uma média de concordância bem superior quando comparado com o tipo A, que apenas teve acesso ao guia de anotação. Assim pode-se supor que o guia de anotação pode ser importante, mas a presença de um especialista do domínio da tarefa de anotação foi determinante para um melhor desempenho dos participantes do experimento.

Tabela 4.5: Concordância com entidades, relações e conectores nos tipos de anotador

Tipo anotador	Concordância com o GSA	Auto-concordância
A	0.452 (σ 0.132)	0.681 (σ 0.190)
B	0.600 (σ 0.049)	0.752 (σ 0.134)

Tabela 4.6: Concordância com entidades nos tipos de anotador

Tipo anotador	Concordância com o GSA	Auto-concordância
A	0.700 (σ 0.172)	0.765 (σ 0.201)
B	0.848 (σ 0.033)	0.916 (σ 0.045)

Tabela 4.7: Concordância com relações nos tipos de anotador

Tipo anotador	Concordância com o GSA	Auto-concordância
A	0.069 (σ 0.239)	0.315 (σ 0.362)
B	0.318 (σ 0.102)	0.501 (σ 0.338)

Tabela 4.8: Concordância com conectores nos tipos de anotador

Tipo anotador	Concordância com o GSA	Auto-concordância
A	0.022 (σ 0.134)	0.681 (σ 0.345)
B	0.166 (σ 0.090)	0.404 (σ 0.347)

Os questionários respondidos ao final do experimento, que podem ser encontrados no Apêndice C, mostram que os participantes tiveram uma impressão positiva da ferramenta, destacando a sua facilidade de uso. Nos questionários, também foi possível ver que grande parte dos participantes do tipo B, que receberam treinamento presencial, não sentiram necessidade de ler o guia de anotação de forma detalhada. Notou-se também que um guia muito extenso não agradou aos participantes.

Com base nos resultados obtidos no experimento, que podem ser vistos de forma mais detalhada no Apêndice D, percebe-se que a concordância média obtida na tarefa de identificação de relações e de conectores foi insatisfatória. Para sanar este problema, um caminho possível seria o de modificar a ontologia de forma a facilitar o entendimento das relações.

4.2

Experimento de aprendizado automático

Os resultados obtidos no experimento de anotação indicaram a necessidade de algumas pequenas mudanças para a melhoria da ontologia. Em primeiro lugar, as classes *HeavyTrafficSituation* e *SlowTrafficSituation* causaram dúvidas em alguns anotadores, incluindo o autor desta dissertação. Qual seria a diferença entre um tráfego pesado e um tráfego lento? Destarte, as duas classes se transformaram em *BadTrafficSituation*, o exato oposto de *GoodTrafficSituation*. Foram também retiradas as relações *hasNumericQuantity* e *hasStringQuantity*, pois entendeu-se que apenas a informação sobre o tipo de *Actor* era realmente importante e não as quantidades, pelo que, doravante, “2 Carros” seria classificado simplesmente como *Car*. Outrossim, a classe *WeatherEvent* foi retirada, uma vez que em nenhum momento apareceu no conjunto de dados usado. A mudança mais importante diz respeito às relações com radical “is”, como *isReferenceFor*, ou *isEdgeFor*, posto que geraram muitas dúvidas quanto à direção de relação: uma *Location A* era algo em relação a outra *Location B* ou esta *Location B* tinha a *Location A* como alvo de uma relação? Para simplificar, todas as relações “is...” se tornaram “has...”, de modo que a *Location* principal do evento sempre tem (“has”) alguma relação que sai dela para *Locations* secundários, o que também torna os diagramas finais da anotação mais claros. A única exceção foi a relação *isRestrictedTo*, que foi retirada, pois diversas vezes confundia-se com *isReferenceTo*. Todas as mudanças descritas podem ser vistas quando as figuras 4.1 e 4.2 são comparadas, onde a primeira mostra a versão da ontologia utilizada no experimento de anotação e a segunda a versão utilizada para anotar os documentos para o experimento de aprendizado automático.

PUC-Rio - Certificação Digital Nº 1421597/CA

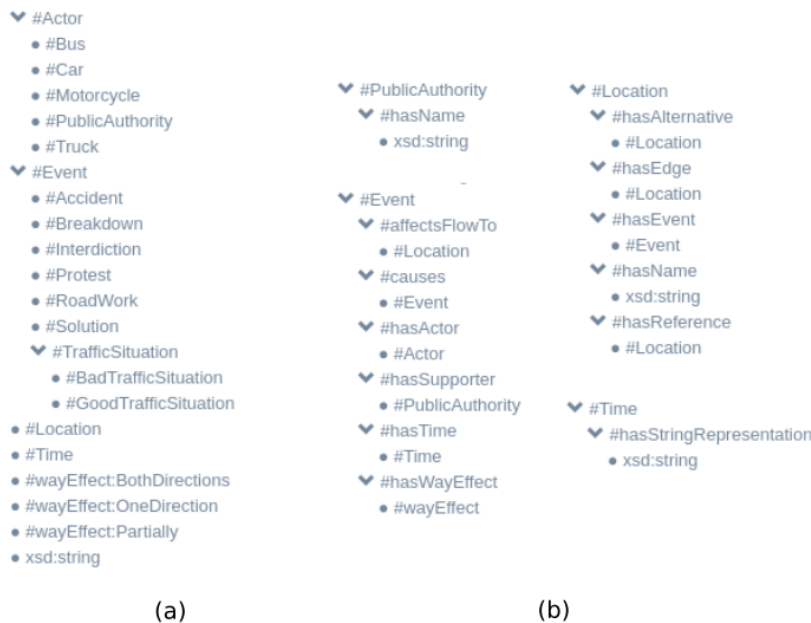


Figura 4.2: Visões da ontologia usada na anotação dos dados para o experimento de aprendizado automático: (a) entidades; (b) relações

4.2.1

Dados

Para os experimentos de aprendizado de máquina, 600 novos *tweets* foram coletados (metade pertence à conta @operacoesRio e a outra metade à conta @odia24horas) e anotados exclusivamente pelo autor deste trabalho. Usou-se uma re-anotação ativada por passo 10, ou seja, as anotações de 60 *tweets* foram repetidas para medições de auto-concordância. A tabela 4.9 apresenta os resultados de tempos de anotação.

Tabela 4.9: Estatísticas de tempo do conjunto de dados usado no experimento de aprendizado automático

Tempo gasto anotando	Tempo médio por documento
15h 52m	1m 35s (σ 1m 44s)

A figura 4.3(a) mostra, para o total de tokens em todo o conjunto de dados, a divisão entre não marcados, marcados como *Tag* ou marcados como *Connector*. Percebe-se a proximidade entre as quantidades de *tokens* participando como *Tag* e como *Connector*, e menos *tokens* não anotados, o que indica que boa parte do conteúdo dos documentos tinha informação relevante. Percebe-se também que um elevado percentual dos *tokens* oferece alguma informação importante para o aprendizado de máquinas. A parte (b) da figura mostra a distribuição do número de *tokens* entre os documentos, que segue um padrão normal com média próxima a 20. Em se tratando de *tweets*, o número máximo de *tokens* estará sempre limitado.

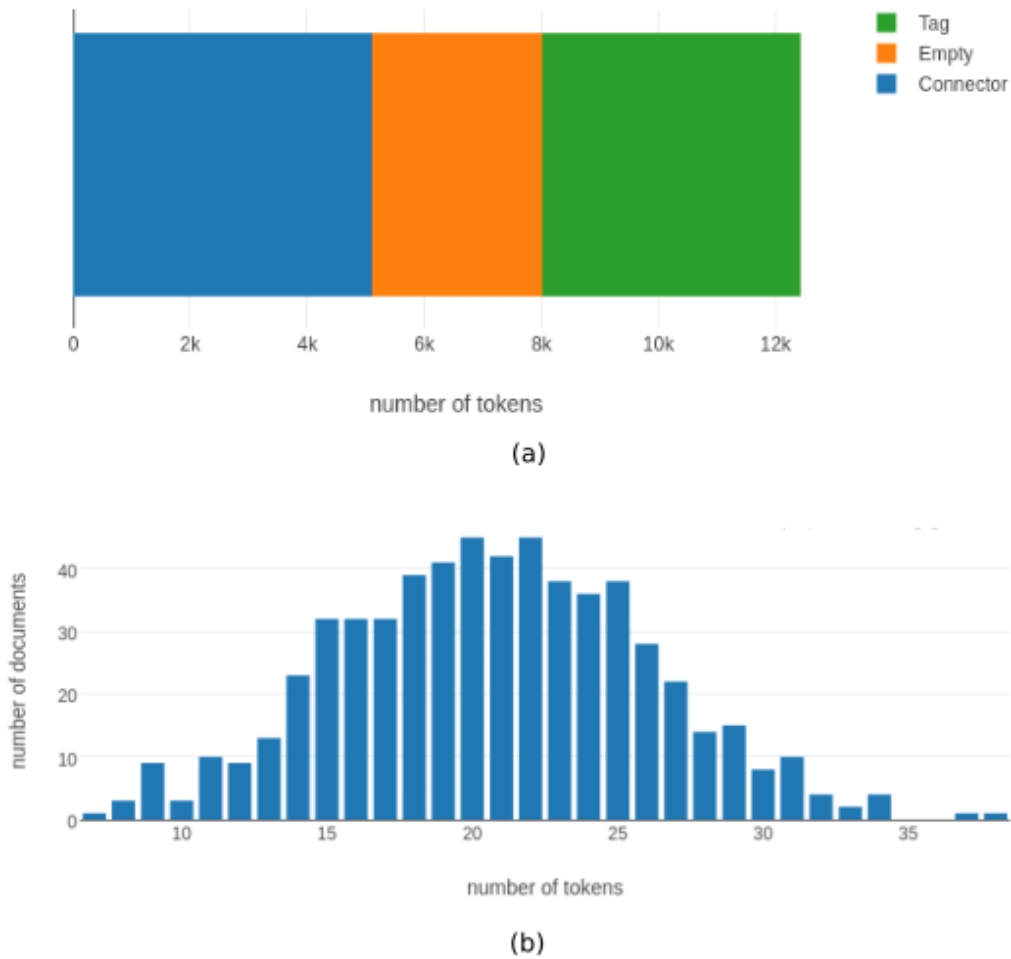
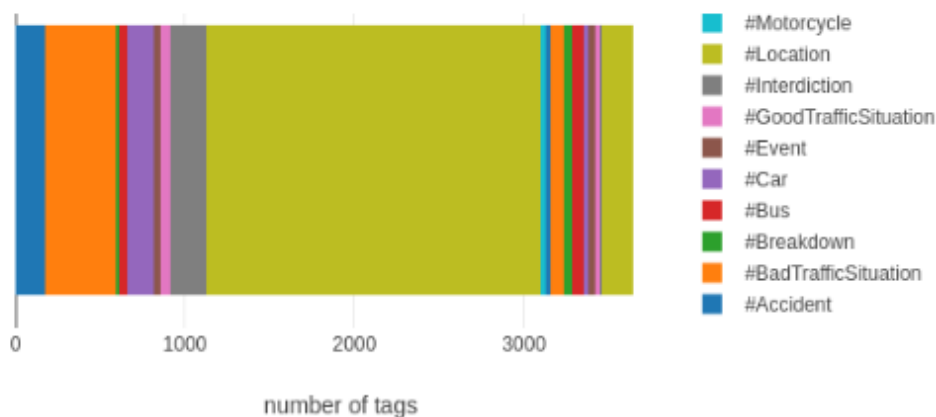
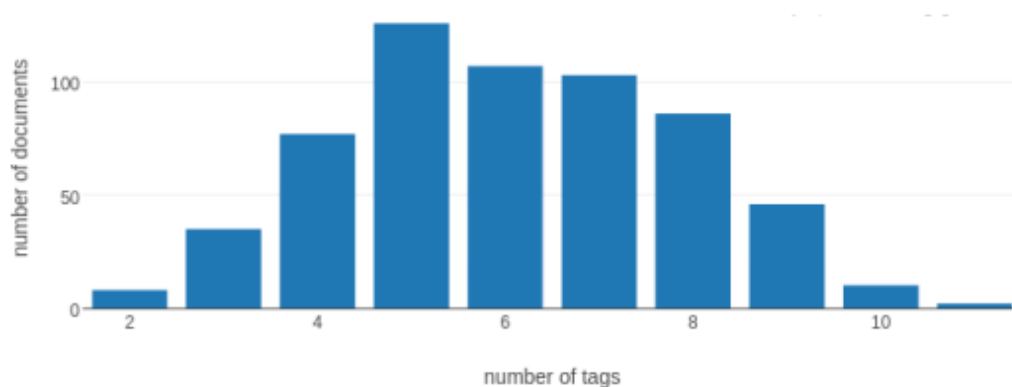


Figura 4.3: Estatísticas de tokens nos dados utilizados no experimento: (a) Cobertura dos tokens; (b) Histograma de distribuição de documentos por quantidade de tokens

A figura 4.4 apresenta os dados para *Tags*. A parte (a) da figura deixa evidente o desbalanceio entre as classes, com uma maioria de *Location* e algumas classes aparecendo muito raramente. A parte (b) apresenta a distribuição do número de *Tags* ou entidades identificadas nos documentos.



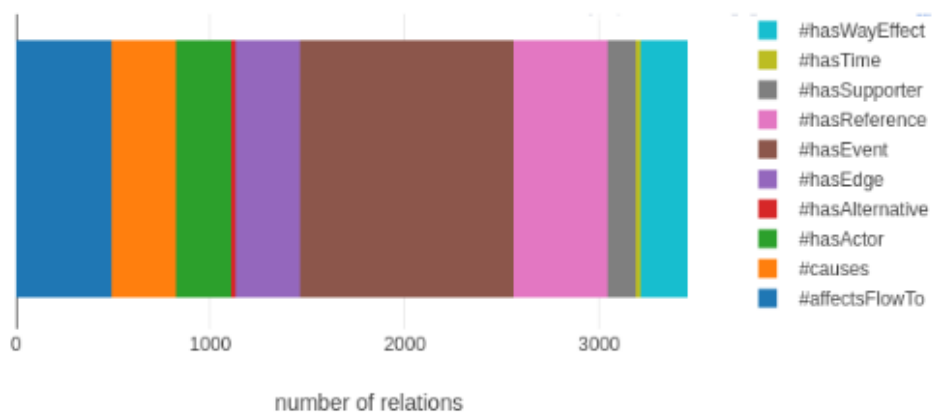
(a)



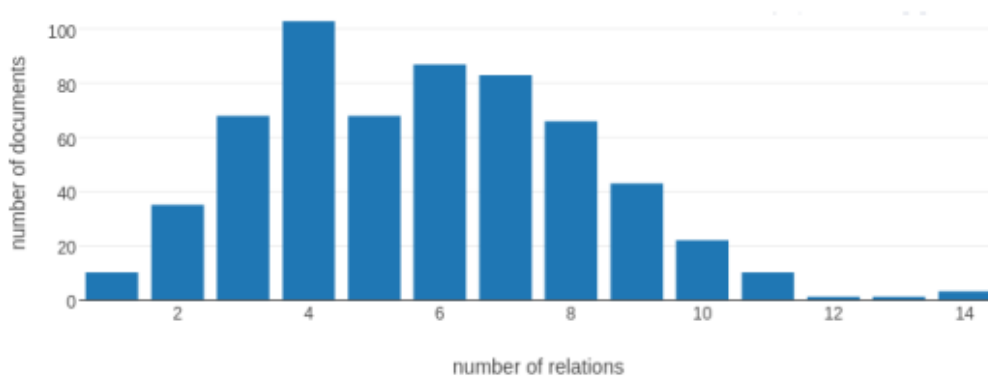
(b)

Figura 4.4: Estatísticas de entidades nos dados utilizados no experimento: (a) Distribuição de entidades nos dados; (b) Histograma de distribuição de documentos por quantidade de entidades

A figura 4.5 mostra a distribuição entre *relations*, onde também ocorre um desbalanço, com as relações envolvendo classes mais frequentes aparecendo, logicamente, mais vezes. O desbalanço, neste caso, é menor do que o observado para as entidades.



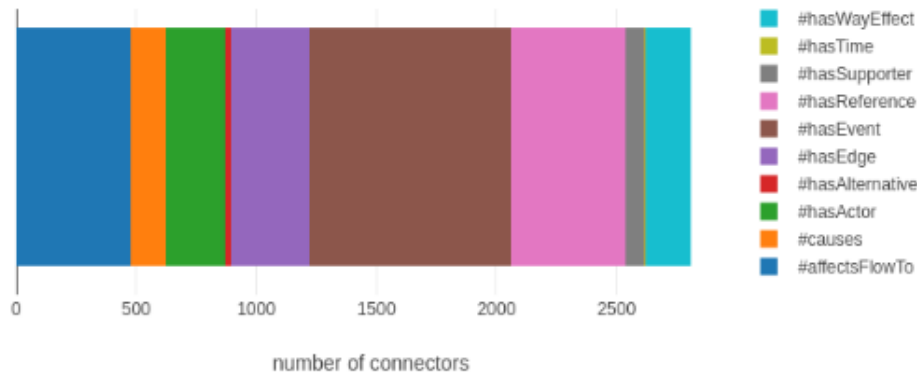
(a)



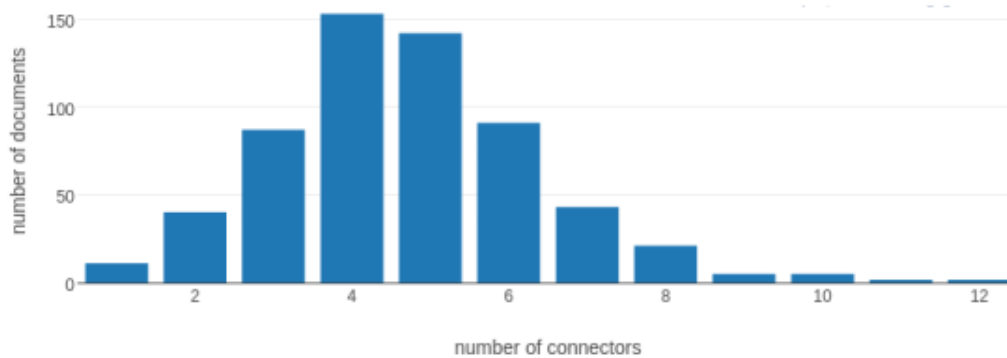
(b)

Figura 4.5: Estatísticas de relações nos dados utilizados no experimento: (a) Distribuição de relações nos dados; (b) Histograma de distribuição de documentos por quantidade de relações

A figura 4.6 apresenta a estatística quanto às participações dos *connectors* em cada relação, que tende a seguir a mesma distribuição verificada para as relações.



(a)



(b)

Figura 4.6: Estatísticas de conectores nos dados utilizados no experimento: (a) Distribuição de conectores nos dados; (b) Histograma de distribuição de documentos por quantidade de conectores

A figura 4.7, uma das mais importantes para o entendimento dos resultados posteriores, apresenta as estatísticas de auto-concordância, uma em relação a cada item da anotação. A parte (a), referente à concordância de todos os itens, mostra que em alguns documentos houve forte variação na re-anotação. A parte (b), referente à concordância apenas entre entidades (*Tags*), mostra que pouco da discordância geral se deve a este item. A parte (c), referente exclusivamente às relações, explica a maior parte da discordância geral. A parte (d), referente exclusivamente aos *connectors*, mostra as maiores discordâncias, mas é apenas uma consequência direta da discordância nas relações. Assim, a anotação das relações foi menos precisa que a das entidades, o que influenciará, de certa forma, na diferença entre os resultados de NER e RE posteriores.

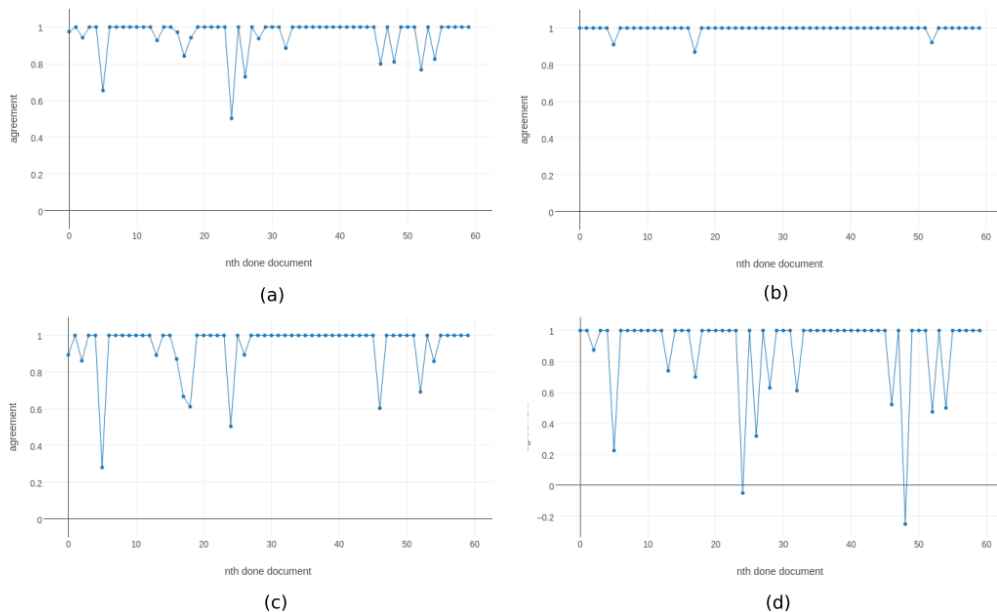


Figura 4.7: Estatísticas de auto-concordância ano longo do tempo nos dados utilizados no experimento: (a) auto-concordância com entidades, relações e conectores; (b) auto-concordância com entidades; (c) auto-concordância com relações; (d) auto-concordância com conectores

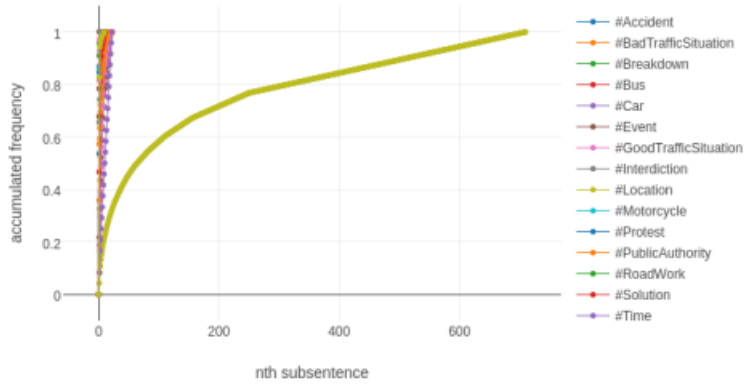
A figura 4.8 mostra a nuvem de palavras gerada pelo sistema relativa às entidades, caracterizando os principais termos ligados à identificação das classes.



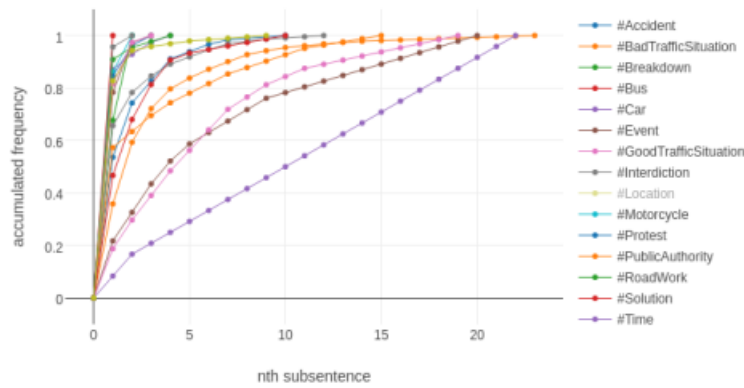
Figura 4.8: Nuvem de palavras associadas às entidades nos dados utilizados no experimento

A figura 4.9 apresenta as curvas de frequências acumuladas das palavras

associadas a cada entidade. A parte (b) se diferencia da (a) por eliminar o caso *Location*, que muda a escala e dificulta a visualização das outras classes. A parte linear da curva, quando acontece, representa uma taxa de aumento constante na frequência acumulada quando da adição de uma nova palavra, mostrando que a generalização automática das expressões para a identificação de determinada classe se faz mais difícil. A parte (a) tem um exemplo claro dessa situação nos caso de *Location*, onde a parte não linear diz respeito a palavras que indicam locais que aparecem mais frequentemente, como “Ponte Rio-Niterói” e “Centro”, enquanto a parte linear se refere a palavras menos frequentes. Nesses casos, apenas a parte não linear com a maior taxa de crescimento deve ser usada para uma geração automática de atributos que não implique em um aumento exagerado da quantidade total de atributos. A parte (b) tem outro exemplo na classe *Time*, onde a taxa de aumento constante acontece em quase toda a curva, o que é de se esperar, uma vez que cada evento terá um tempo diferente dos outros, com raros casos de repetição. Nesses casos, com tempos e outros tipos semelhantes, recomenda-se o uso de expressões regulares que generalizem a referida informação. A geração automática dos atributos usa como parâmetro principal um valor de profundidade (*depth*), que representa o valor da abscissa destes gráficos. O valor exato a ser considerado é algo subjetivo, variando de acordo com a quantidade máxima de atributos que cada um tolere. Para o presente trabalho, no caso da tarefa NER e tendo como base a parte (b) da figura, considerou-se 20 um valor razoável.



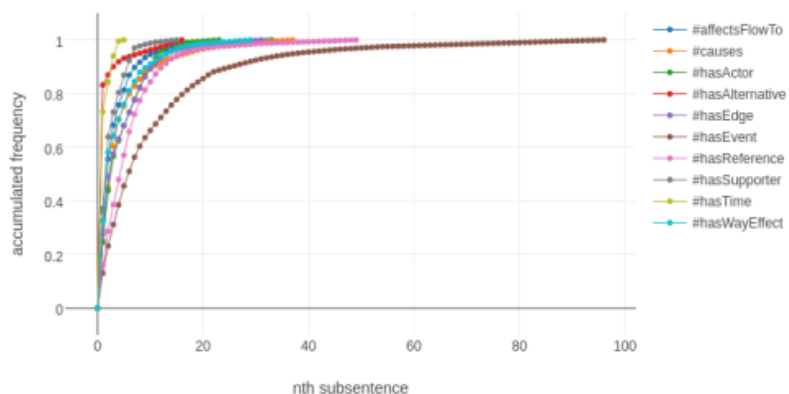
(a)



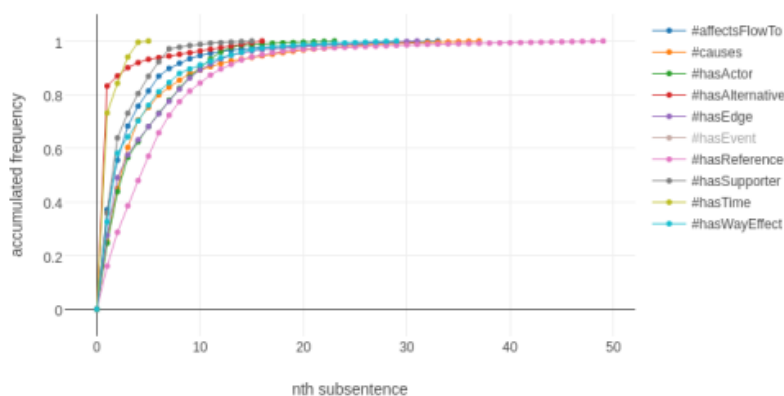
(b)

Figura 4.9: Curva de frequência acumulada de palavras associadas às entidades: (a) todas entidades; (b) excluindo a entidade com mais ocorrências, neste caso *Location*

A figura 4.10 apresenta a nuvem de palavras referente aos conectores, elucidando alguns casos típicos de indicação de sentido de fluxo, de locais, entre outros.



(a)



(b)

Figura 4.11: Curva de frequência acumulada de palavras associadas aos conectores: (a) todos conectores; (b) excluindo conector com mais ocorrências, neste caso *hasEvent*

4.2.2

Metodologia

Os experimentos de NER e RE, ambos adotando uma estratégia de *cross-validation* de 10-*fold*, foram conduzidos em alguns passos semelhantes:

1. Separação randômica de 20% dos dados de cada conta de *tweets* de trânsito para serem utilizados como conjunto de teste, isto é, destacados do fluxo do experimento para serem utilizados apenas na etapa final, sem iterações sobre eles para melhoria dos modelos;
2. Execução do mesmo procedimento do item anterior nos 80% restantes, para criação dos conjuntos de treino e revalidação;
3. Uso do gerador automático de atributos com base nos dados;
4. Adição de alguns atributos que não foram contemplados no item anterior;
5. Treinamento dos modelos utilizando todos os classificadores disponíveis com suas configurações padrão;

6. Apenas nas tarefas de NER, execução da mesma ação do item anterior com acréscimo de atributos;
7. Avaliação e seleção, com base no F1 médio, dos classificadores com resultados minimamente razoáveis e, apenas no caso do NER, verificação da melhoria ou não do resultado pelo acréscimo de atributos do item anterior;
8. Execução do *gridsearch* em todos os classificadores selecionados no item anterior;
9. Treino do modelo final com os melhores parâmetros encontrados no item anterior, usando para treino o conjunto de treino completo (treino e validação) e para o teste único e definitivo o conjunto de teste separado no primeiro item;
10. Execução do item anterior utilizando agrupamentos e remoções de entidades, para mero teste desta estratégia, sem objetivar qualquer melhoria de modelo.

Os atributos para NER foram os seguintes:

- LEMMA{0} automático do tipo VALUE com separação das palavras e sem concatenação das expressões regulares;
- RANGE-LEMMA{-4, 4} automático do tipo FROM-CONNECTOR, com separação das palavras e sem concatenação das expressões regulares;
- RANGE-LEMMA{-4, 4} automático do tipo TO-CONNECTOR, com separação das palavras e sem concatenação das expressões regulares;
- POS{-3, -2, -1, 0, 1, 2, 3} apenas com as categorias da árvore de *POS Tagging*

Os atributos para RE foram os seguintes:

- NODE-TO-NODE-DISTANCE;
- NODE-TO-NODE-DISTANCE-WITH-SIGNAL;
- CLASS-NODE-FROM;
- CLASS-NODE-TO;
- POSSIBLE-RELATION;
- INTERIOR-RANGE-LEMMA{-4, 4} automático do tipo CONNECTOR, com separação das palavras e sem concatenação das expressões regulares;

4.2.3

Resultados

Esta seção apresenta os resultados dos experimentos descritos, executados no ambiente LER, que funciona em um servidor Intel® Core™i5, com CPU 650 @ 3.20GHz e 12GB de RAM. Esses resultados estão descritos de forma mais detalhada no Apêndice E. As próximas seções estão divididas entre os experimentos executados para a tarefa de NER e RE respectivamente.

4.2.3.1

NER

O passo 5 da metodologia, cujos resultados podem ser vistos na tabela 4.10, indica que, com parâmetros *default*, o algoritmo Random Forest se mostra o mais rápido. Entretanto, o SGD apresenta o melhor resultado de F1. Percebe-se claramente o estranho comportamento do classificador SVC, que resultou em um F1 extremamente baixo. Isso provavelmente se deve ao fato do *kernel default* para o SVC ser do tipo RBF, que tende a piores resultados quando os valores dos atributos não são normalmente distribuídos (42), como é o caso aqui tratado, cuja distribuição é binomial.

Tabela 4.10: Resultado NER passo 5

Classificador	Tempo tarefa	Tempo CPU treino	F1 validação	F1 revalidação
Random Forest	4m 31s	0.514s	0.554 (σ 0.281)	0.682 (σ 0.353)
Stochastic Gradient Descent	4m 52s	1.672s	0.574 (σ 0.293)	0.752 (σ 0.310)
SVC	11m 47s	50.304s	0.055 (σ 0.205)	0.067 (σ 0.224)

Uma vez que os parâmetros *default* do SVC apresentaram tal problema, um *gridsearch* especial foi feito para este classificador ainda no passo 5, para uma avaliação sobre a possibilidade ou não de deixá-lo entre os avaliados nos passos seguintes. Os resultados podem ser vistos na tabela 4.11, onde o *kernel* linear se mostrou o mais adequado.

Tabela 4.11: Gridsearch no NER para o SVC no passo 5

Parâmetro	Valores testados	Melhor valor
Kernel	Radial Basis Function, Linear, Polynomial, Sigmoidal	Linear

Passando ao passo 6, onde a mudança de atributos ocorre apenas pelo aumento da quantidade de níveis do *POS Tagging* consideradas, agora com todos, nova rodada é realizada mantendo-se os parâmetros *default*, exceto pela mudança pontual do *kernel* do SVC de RBF para Linear. Os resultados podem ser vistos na tabela 4.12. Percebe-se que, como era esperado, os tempos de treinamento aumentaram, mas mantendo a mesma proporção relativa entre os classificadores, onde se vê que o SVC é muitíssimo mais lento que os demais. Não obstante, pelo uso de um *kernel* linear, também apresenta o melhor resultado de todos. Em termos de atributos, percebe-se que os resultados melhoraram em todos os classificadores, o que leva à conclusão de que o uso de mais níveis de POS tende a melhorar a qualidade da classificação no caso aqui estudado.

Tabela 4.12: Resultado NER passo 6

Classificador	Tempo tarefa	Tempo CPU treino	F1 validação	F1 revalidação
Random Forest	6m 19s	0.712s	0.567 (σ 0.243)	0.777 (σ 0.296)
Stochastic Gradient Descent	6m 54s	2.959s	0.609 (σ 0.284)	0.776 (σ 0.260)
SVC	1h 11m 43s	64.763s	0.731 (σ 0.226)	0.879 (σ 0.221)

Os passos 7 e 8, isto é, a seleção dos classificadores minimamente viáveis (neste caso, todos) e execução de *gridsearch* em todos estes para a escolha da melhor combinação classificador-parâmetros, têm seus resultados apresentados nas tabelas 4.13, 4.14, 4.15 e 4.16. O SVC apresenta o melhor desempenho no conjunto de treino, com um maior F1 médio e menor desvio padrão, sendo, de outra sorte, vencido pelo Random Forest no conjunto de revalidação, que por sua vez apresentou um valor de F1 médio próximo ao do SVC no conjunto de treino, mas com um desvio padrão um pouco maior. Por fim, escolheu-se o Random Forest e seus melhores parâmetros como a estratégia final para NER em razão do maior valor de F1 para a revalidação. É de se notar que o seu tempo de treino foi o dobro do observado no SVC, o que se deve ao elevado número de árvores adotado. Todavia, sabe-se que a complexidade do SVC na presente implementação cresce mais que quadraticamente com o número de instâncias (34), enquanto a do Random Forest depende mais do número de árvores - já fixado - e suas profundidades, do que da quantidade de dados de treino (35, 43), de onde se conclui que essa diferença de tempos deverá diminuir, ou mesmo se inverter, com o aumento da quantidade de dados.

Tabela 4.13: Gridsearch no NER para o Random Forest no passo 8

Parâmetro	Valores testados	Melhor valor
Class weight	balanced, None	None
Criterion	gini, entropy	gini
Max features	$\sqrt{n_features}$, $\log_2(n_features)$, $n_features$	$n_features$
Number of trees in the forest	10, 50, 100	50
Use out-of-bag samples to estimate accuracy	true, false	true
One vs Rest	true, false	true

Tabela 4.14: Gridsearch no NER para o Stochastic Gradient Descent no passo 8

Parâmetro	Valores testados	Melhor valor
Class weight	balanced, None	None
Learning rate	constant, optimal, invscaling	
Loss function	Hinge, Log, Modified huber, Epsilon insensitive, Squared epsilon insensitive, Squared hinge, Perceptron, Huber, Epsilon insensitive, Squared epsilon insensitive	Hinge
One vs Rest	true, false	true
Penalty	None, 11, 12, elasticnet	elasticnet

Tabela 4.15: Gridsearch no NER para o SVC no passo 8

Parâmetro	Valores testados	Melhor valor
Penalty parameter C	0.1, 1, 10, 100	1
Class weight	balanced, None	None
One vs Rest	true, false	true
Use Shrinking heuristic	true, false	true

Tabela 4.16: Resultado NER passo 8

Classificador	Tempo tarefa	Tempo CPU treino	F1 validação	F1 revalidação
Random Forest	1d 14h 11m 55s	218.456s	0.745 (σ 0.253)	0.906 (σ 0.213)
Stochastic Gradient Descent	11h 30m 37s	12.315s	0.676 (σ 0.233)	0.833 (σ 0.229)
SVC	13h 29m 28s	108.015s	0.751 (σ 0.192)	0.881 (σ 0.216)

As tabelas 4.17 e 4.18 apresentam os valores de F1 por classe para NER referentes à validação e teste final, respectivamente. Obteve-se, portanto, um classificador com F1-score médio de 0.821, com desvio padrão entre classes de 0.287, com valores de Precision e Recall razoavelmente próximos ao valor de F1. Nota-se que a maioria dos valores individuais de F1 foi 1 ou muito próxima de 1, restando algumas classes com valores de F1 nulos, todas ligadas ao Enumerator *wayEffect*.

Tabela 4.17: Resultado final detalhado: Random Forest (validação)

Class	Support	Accuracy, σ CV	Recall, σ CV	Precision, σ CV	F1, σ CV
O	6457	0.964, 0.016	0.980, 0.010	0.964, 0.019	0.972, 0.012
B-#Accident	141	0.999, 0.000	0.962, 0.053	0.991, 0.025	0.975, 0.028
B-#BadTrafficSituation	339	0.996, 0.002	0.948, 0.046	0.944, 0.044	0.945, 0.032
B-#Breakdown	17	0.999, 0.000	0.962, 0.104	0.944, 0.157	0.940, 0.115
B-#Bus	39	1.000, 0.000	1.000, 0.000	1.000, 0.000	1.000, 0.000
B-#Car	126	0.999, 0.000	0.983, 0.049	1.000, 0.000	0.990, 0.027
B-#Event	37	0.998, 0.000	0.667, 0.204	0.925, 0.114	0.749, 0.131
B-#GoodTrafficSituation	45	0.997, 0.002	0.546, 0.307	0.625, 0.370	0.560, 0.315
B-#Interdiction	165	0.999, 0.001	0.974, 0.043	0.995, 0.014	0.984, 0.026
B-#Location	1586	0.980, 0.007	0.960, 0.017	0.920, 0.036	0.939, 0.024
B-#Motorcycle	22	0.999, 0.000	0.944, 0.157	1.000, 0.000	0.962, 0.104
B-#Protest	17	0.999, 0.000	0.802, 0.275	1.000, 0.000	0.858, 0.208
B-#PublicAuthority	63	0.998, 0.001	0.854, 0.212	0.870, 0.204	0.839, 0.183
B-#RoadWork	32	0.999, 0.000	0.925, 0.138	1.000, 0.000	0.955, 0.083
B-#Solution	55	0.999, 0.000	0.916, 0.220	1.000, 0.000	0.937, 0.165
B-#Time	20	0.998, 0.001	0.60, 0.447	0.604, 0.447	0.586, 0.426
B-#Truck	30	1.000, 0.000	1.000, 0.000	1.000, 0.000	1.000, 0.000
B-#wayEffect:BothDirections	30	0.999, 0.000	0.980, 0.059	0.901, 0.153	0.929, 0.092
B-#wayEffect:OneDirection	4	0.998, 0.000	0.000, 0.000	0.000, 0.000	0.000, 0.000
B-#wayEffect:Partially	146	0.998, 0.001	0.840, 0.285	0.846, 0.299	0.838, 0.284
I-#BadTrafficSituation	105	0.997, 0.002	0.872, 0.167	0.904, 0.109	0.878, 0.124
I-#Event	8	0.998, 0.000	0.200, 0.400	0.200, 0.400	0.200, 0.400
I-#GoodTrafficSituation	43	0.997, 0.002	0.732, 0.329	0.825, 0.290	0.750, 0.288
I-#Location	396	0.976, 0.013	0.617, 0.206	0.754, 0.161	0.657, 0.141
I-#Protest	2	0.997, 0.000	0.500, 0.000	0.500, 0.000	0.500, 0.000
I-#PublicAuthority	9	0.998, 0.001	0.500, 0.408	0.666, 0.471	0.555, 0.415
I-#Time	8	0.995, 0.001	0.270, 0.270	0.450, 0.450	0.338, 0.338
I-#wayEffect:BothDirections	34	0.999, 0.000	0.800, 0.400	0.713, 0.379	0.748, 0.382
I-#wayEffect:Partially	2	0.998, 0.000	0.500, 0.000	1.000, 0.000	0.666, 0.000
Average, σ		0.996, 0.007	0.753, 0.263	0.812, 0.250	0.767, 0.252

Tabela 4.18: Resultado final detalhado: Random Forest (teste)

Class	Support	Accuracy	Recall	Precision	F1
O	1553	0.962	0.981	0.960	0.971
B-#Accident	36	0.999	0.972	1.000	0.985
B-#BadTrafficSituation	72	0.997	0.986	0.946	0.965
B-#Breakdown	5	1.000	1.000	1.000	1.000
B-#Bus	11	1.000	1.000	1.000	1.000
B-#Car	28	0.999	0.964	1.000	0.981
B-#Event	5	0.999	0.800	0.800	0.800
B-#GoodTrafficSituation	15	0.998	0.800	0.923	0.857
B-#Interdiction	43	0.999	0.976	1.000	0.988
B-#Location	390	0.976	0.938	0.917	0.927
B-#Motorcycle	7	1.000	1.000	1.000	1.000
B-#Protest	9	1.000	1.000	1.000	1.000
B-#PublicAuthority	18	0.997	0.777	0.933	0.848
B-#RoadWork	12	0.999	0.916	1.000	0.956
B-#Solution	20	0.997	0.800	0.941	0.864
B-#Time	4	0.999	0.750	1.000	0.857
B-#Truck	7	0.999	0.857	1.000	0.923
B-#wayEffect:BothDirections	8	0.999	1.000	0.888	0.941
B-#wayEffect:Partially	42	0.996	0.000	0.000	0.000
I-#BadTrafficSituation	24	0.999	0.857	0.947	0.900
I-#GoodTrafficSituation	13	0.998	0.958	0.958	0.958
I-#Location	115	0.976	0.846	0.916	0.879
I-#Time	5	0.997	0.626	0.827	0.712
I-#wayEffect:BothDirections	10	1.000	0.000	0.000	0.000
I-#wayEffect:Partially	2	0.998	0.200	0.250	0.222
Average, σ		0.995, 0.009	0.800, 0.289	0.848, 0.290	0.821, 0.287

Após a escolha do melhor classificador e teste final, estudou-se também o mesmo conjunto de parâmetros sobre um conjunto de dados com generalização das classes pertencentes à *Actor* e exclusão da classe *wayEffect:Partially*, que foi a mais problemática no teste já discutido. O estado final da árvore pode ser visto na figura 4.12.

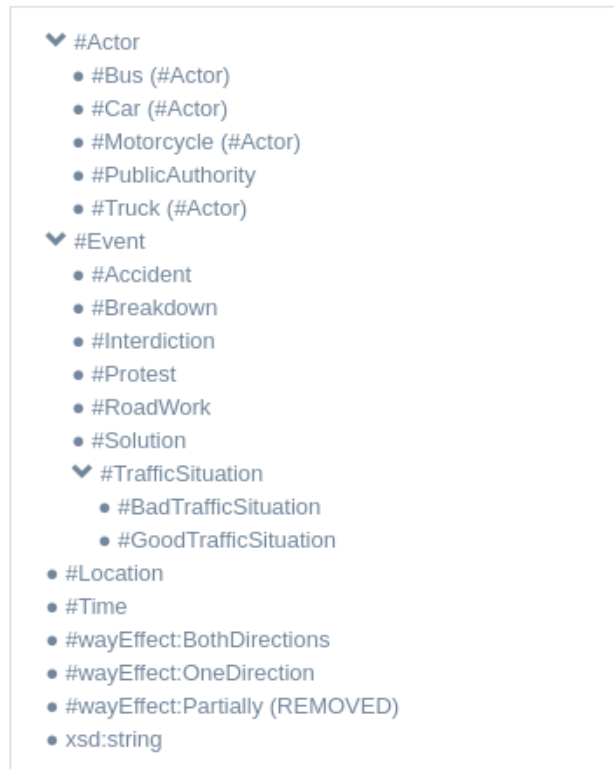


Figura 4.12: Generalização e remoções usadas para redução das classes

A tabela 4.19 mostra o resultado global com redução de classes comparado ao resultado original sem redução. Conclui-se que, neste caso específico, a redução do número de classes piorou o valor de F1, diminuindo, porém, o tempo de treinamento. Isso não significa que a estratégia de generalização não possa melhorar o desempenho de NER em outros casos. Detalhes dos resultados deste teste de generalizações são apresentados nas tabelas 4.20 e 4.21.

Tabela 4.19: Resultado NER final

Classificador	Tempo tarefa	Tempo CPU treino	F1 validação	F1 teste
Random Forest	3h 20m 7s	387.326s	0.767 (σ 0.252)	0.821 (σ 0.287)
Random Forest (com redução de classes)	1h 44m 22s	252.572s	0.734 (σ 0.273)	0.818 (σ 0.287)

Tabela 4.20: Resultado final detalhado: Random Forest (validação com redução de classes)

Class	Support	Accuracy, σ CV	Recall, σ CV	Precision, σ CV	F1, σ CV
O	6605	0.965, 0.016	0.981, 0.009	0.966, 0.019	0.973, 0.012
B-#Accident	141	0.999, 0.000	0.962, 0.053	0.984, 0.046	0.972, 0.033
B-#Actor	217	0.999, 0.000	0.983, 0.035	1.000, 0.000	0.991, 0.018
B-#BadTrafficSituation	339	0.995, 0.002	0.942, 0.053	0.937, 0.039	0.938, 0.031
B-#Breakdown	17	0.999, 0.000	0.962, 0.104	1.000, 0.000	0.977, 0.062
B-#Event	37	0.998, 0.001	0.636, 0.161	0.925, 0.114	0.733, 0.086
B-#GoodTrafficSituation	45	0.996, 0.002	0.531, 0.315	0.556, 0.351	0.517, 0.314
B-#Interdiction	165	0.999, 0.001	0.974, 0.043	0.995, 0.014	0.984, 0.026
B-#Location	1586	0.980, 0.007	0.957, 0.019	0.922, 0.033	0.939, 0.024
B-#Protest	17	0.999, 0.000	0.833, 0.288	1.000, 0.000	0.875, 0.216
B-#PublicAuthority	63	0.998, 0.001	0.854, 0.212	0.925, 0.159	0.857, 0.154
B-#RoadWork	32	0.999, 0.000	0.925, 0.138	1.000, 0.000	0.955, 0.083
B-#Solution	55	0.999, 0.000	0.916, 0.220	0.987, 0.033	0.930, 0.163
B-#Time	20	0.998, 0.001	0.433, 0.453	0.500, 0.500	0.458, 0.465
B-#wayEffect:BothDirections	30	0.999, 0.000	0.980, 0.059	0.889, 0.169	0.920, 0.102
B-#wayEffect:OneDirection	4	0.998, 0.000	0.000, 0.000	0.000, 0.000	0.000, 0.000
I-#BadTrafficSituation	105	0.997, 0.003	0.785, 0.303	0.822, 0.290	0.798, 0.290
I-#Event	8	0.998, 0.000	0.200, 0.400	0.133, 0.266	0.160, 0.320
I-#GoodTrafficSituation	43	0.997, 0.002	0.757, 0.339	0.851, 0.300	0.774, 0.299
I-#Location	396	0.976, 0.013	0.618, 0.202	0.777, 0.168	0.665, 0.138
I-#Protest	2	0.997, 0.000	0.500, 0.000	0.500, 0.000	0.500, 0.000
I-#PublicAuthority	9	0.998, 0.001	0.575, 0.422	0.538, 0.411	0.555, 0.415
I-#Time	8	0.995, 0.001	0.270, 0.270	0.460, 0.460	0.340, 0.340
I-#wayEffect:BothDirections	34	0.999, 0.000	0.900, 0.300	0.747, 0.295	0.808, 0.286
Average, σ		0.995, 0.008	0.728, 0.276	0.767, 0.278	0.734, 0.273

Tabela 4.21: Resultado final detalhado: Random Forest (teste com redução de classes)

Class	Support	Accuracy	Recall	Precision	F1
O	1597	0.962	0.979	0.964	0.971
B-#Accident	36	1.000	1.000	1.000	1.000
B-#Actor	53	0.999	0.981	1.000	0.990
B-#BadTrafficSituation	72	0.997	0.986	0.946	0.965
B-#Breakdown	5	1.000	1.000	1.000	1.000
B-#Event	5	0.999	0.800	1.000	0.888
B-#GoodTrafficSituation	15	0.997	0.800	0.857	0.827
B-#Interdiction	43	0.999	0.976	1.000	0.988
B-#Location	390	0.978	0.946	0.922	0.934
B-#Protest	9	1.000	1.000	1.000	1.000
B-#PublicAuthority	18	0.997	0.777	0.875	0.823
B-#RoadWork	12	0.999	0.916	1.000	0.956
B-#Solution	20	0.997	0.800	0.941	0.864
B-#Time	4	0.998	0.750	0.600	0.666
B-#wayEffect:BothDirections	8	0.999	1.000	0.888	0.941
I-#BadTrafficSituation	24	0.999	0.000	0.000	0.000
I-#GoodTrafficSituation	13	0.998	0.958	0.958	0.958
I-#Location	115	0.975	0.846	0.916	0.879
I-#Time	5	0.996	0.608	0.833	0.703
I-#wayEffect:BothDirections	10	1.000	0.000	0.000	0.000
Average		0.995 (σ 0.009)	0.806 (σ 0.289)	0.835 (σ 0.292)	0.818 (σ 0.287)

4.2.3.2

RE

Seguindo em RE o que foi feito para NER, os resultados do passo 5 podem ser vistos na tabela 4.22. O algoritmo Frank Wolfe SSVM apresentou, com seus parâmetros *default*, o melhor resultado de F1 para validação e revalidação, porém com um tempo de treinamento mais que 20 vezes maior que o Structured Perceptron.

Tabela 4.22: Resultado RE passo 5

Classificador	Tempo tarefa	Tempo CPU treino	F1 validação	F1 revalidação
Frank Wolfe SSVM	3h 41m 34s	335.189s	0.617 (σ 0.340)	0.674 (σ 0.356)
Structured Perceptron	2h 55m 23s	14.983s	0.552 (σ 0.294)	0.635 (σ 0.293)

Passando direto ao passo 8, execução de *gridsearch*, definiu-se os melhores parâmetros para cada algoritmo dentro dos critérios de busca de cada caso, conforme pode ser visto nas tabelas 4.23 e 4.24.

Tabela 4.23: Gridsearch no RE para o Frank Wolfe SSVM no passo 8

Parâmetro	Valores testados	Melhor valor
Penalty parameter C	0.1, 1, 10, 100	100
Use weight averaging	true, false	false

Tabela 4.24: Gridsearch no RE para o Structured Perceptron no passo 8

Parâmetro	Valores testados	Melhor valor
Average	true, false	true
Batch learning	true, false	false

Os algoritmos foram reavaliados, ainda sobre o conjunto de revalidação. Os resultados apresentados na tabela 4.25 indicam que, mais uma vez, o algoritmo Frank Wolfe SSVM alcança o melhor desempenho em termos de F1, novamente apresentando o maior tempo de treinamento. Seguindo o critério de maior F1 na revalidação, escolheu-se o Frank Wolfe SSVM, com seus melhores parâmetro, como estratégia final para teste.

Tabela 4.25: Resultado RE passo 8

Classificador	Tempo tarefa	Tempo CPU treino	F1 validação	F1 revalidação
Frank Wolfe SSVM	7h 42m 12s	212.895s	0.658 (σ 0.269)	0.725 (σ 0.223)
Structured Perceptron	3h 5m 47s	16.456s	0.634 (σ 0.291)	0.665 (σ 0.342)

Os resultados finais, apresentados nas tabelas 4.26 e 4.27, mostram que foi obtido um classificador com F1 médio de 0.701 sobre o conjunto de teste, com valores de Precision e Recall próximos a F1.

Tabela 4.26: Resultado final detalhado: Frank Wolfe SSVM (validação)

Class	Support	Accuracy, σ CV	Recall, σ CV	Precision, σ CV	F1, σ CV
O	13797	0.885, 0.019	0.917, 0.033	0.944, 0.017	0.930, 0.013
#affectsFlowTo	390	0.987, 0.003	0.826, 0.068	0.714, 0.069	0.763, 0.051
#causes	261	0.992, 0.003	0.804, 0.144	0.778, 0.072	0.782, 0.086
#hasActor	230	0.998, 0.001	0.965, 0.043	0.956, 0.078	0.958, 0.044
#hasAlternative	18	0.998, 0.001	0.166, 0.267	0.214, 0.364	0.185, 0.304
#hasEdge	279	0.979, 0.012	0.522, 0.325	0.438, 0.290	0.420, 0.240
#hasEvent	864	0.957, 0.008	0.707, 0.186	0.598, 0.104	0.626, 0.080
#hasReference	390	0.966, 0.012	0.314, 0.175	0.289, 0.198	0.292, 0.174
#hasSupporter	112	0.998, 0.001	0.923, 0.136	0.934, 0.074	0.919, 0.079
#hasTime	21	0.999, 0.000	0.900, 0.223	0.958, 0.093	0.904, 0.157
Average, σ		0.976, 0.033	0.704, 0.264	0.682, 0.271	0.678, 0.270

Tabela 4.27: Resultado final detalhado: Frank Wolfe SSVM (teste)

Class	Support	Accuracy	Recall	Precision	F1
O	3442	0.913	0.966	0.932	0.949
#affectsFlowTo	99	0.990	0.747	0.860	0.799
#causes	71	0.992	0.690	0.830	0.753
#hasActor	57	0.997	0.859	0.942	0.899
#hasAlternative	8	0.997	0.125	0.200	0.153
#hasEdge	52	0.990	0.615	0.640	0.627
#hasEvent	226	0.966	0.623	0.726	0.671
#hasReference	100	0.976	0.240	0.558	0.335
#hasSupporter	33	0.999	0.939	1.000	0.968
#hasTime	4	0.999	0.750	1.000	0.857
Average, σ		0.982, 0.025	0.655, 0.263	0.769, 0.237	0.701, 0.254

O teste da árvore da figura 4.12 foi realizado também para este classificador e os resultados comparativos para RE são apresentados na tabela 4.28, onde mais uma vez se vê que o resultado piorou com a redução de classes, com correspondente redução do tempo de treinamento, exatamente como o comportamento para NER. As tabelas 4.29 e 4.30 detalham os resultados da redução para cada classe de relação.

Tabela 4.28: Resultado RE final

Classificador	Tempo tarefa	Tempo CPU treino	F1 validação	F1 teste
Frank Wolfe SSVM	4h 24m 58s	322.663s	0.678 (σ 0.270)	0.701 (σ 0.254)
Frank Wolfe SSVM (com redução de classes)	2h 30m 42s	205.059s	0.647 (σ 0.304)	0.662 (σ 0.282)

Tabela 4.29: Resultado final detalhado: Frank Wolfe SSVM (validação com redução de classes)

Class	Support	Accuracy, σ CV	Recall, σ CV	Precision, σ CV	F1, σ CV
O	11817	0.863, 0.048	0.906, 0.081	0.927, 0.022	0.913, 0.037
#affectsFlowTo	390	0.986, 0.005	0.745, 0.107	0.790, 0.113	0.755, 0.062
#causes	261	0.990, 0.004	0.784, 0.146	0.771, 0.083	0.766, 0.084
#hasActor	230	0.998, 0.001	0.960, 0.051	0.965, 0.064	0.960, 0.043
#hasAlternative	18	0.997, 0.001	0.095, 0.233	0.071, 0.174	0.081, 0.199
#hasEdge	279	0.975, 0.011	0.446, 0.300	0.407, 0.269	0.358, 0.196
#hasEvent	864	0.946, 0.019	0.691, 0.180	0.614, 0.146	0.615, 0.078
#hasReference	390	0.961, 0.026	0.226, 0.198	0.305, 0.238	0.210, 0.170
#hasSupporter	112	0.999, 0.000	0.971, 0.057	0.908, 0.101	0.933, 0.059
#hasTime	21	0.999, 0.000	0.858, 0.224	0.958, 0.093	0.880, 0.152
Average, σ		0.971, 0.039	0.668, 0.293	0.672, 0.296	0.647, 0.304

Tabela 4.30: Resultado final detalhado: Frank Wolfe SSVM (teste com redução de classes)

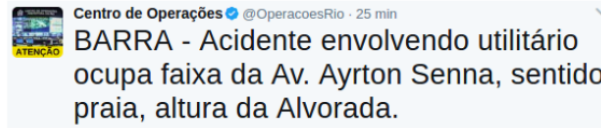
Class	Support	Accuracy	Recall	Precision	F1
O	2884	0.889	0.963	0.906	0.934
#affectsFlowTo	99	0.988	0.666	0.891	0.763
#causes	71	0.991	0.661	0.870	0.752
#hasActor	57	0.997	0.912	0.945	0.928
#hasAlternative	8	0.996	0.125	0.125	0.125
#hasEdge	52	0.983	0.673	0.454	0.542
#hasEvent	226	0.957	0.473	0.781	0.589
#hasReference	100	0.970	0.130	0.448	0.201
#hasSupporter	33	0.998	0.878	1.000	0.935
#hasTime	4	0.999	0.750	1.000	0.857
Average, σ		0.977, 0.032	0.623, 0.282	0.742, 0.281	0.662, 0.282

4.2.3.3

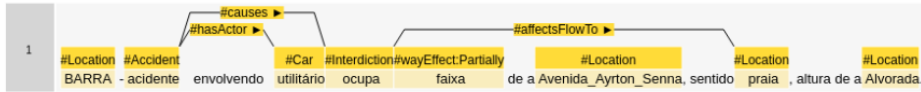
Uso dos modelos finais

Uma vez escolhidos os classificadores de NER e RE, estes foram estabelecidos para um serviço de exemplo, onde foram executados testes com dados totalmente novos e mais recentes retirados diretamente dos canais de *twitter*, isto é, dados que não participaram do conjunto de treinamento e tampouco do conjunto de teste.

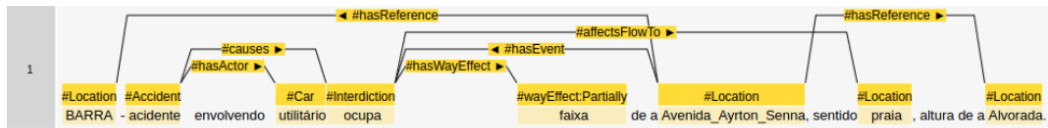
A figura 4.13 mostra um exemplo de resposta razoável, com triplas RDF correspondentes apresentadas na figura 4.14. Percebe-se claramente que as falhas se devem, neste caso, somente ao classificador de RE, pois as entidades de NER foram correta e totalmente identificadas. Isso provavelmente se deve ao que foi diagnosticado na figura 4.7, onde as quedas na auto-concordância nas anotações das relações implicaram quedas ainda maiores nas auto-concordâncias de conectores, o que provavelmente atrapalhou o processo posterior de geração automática de atributos, gerando expressões regulares ambíguas e reduzindo o desempenho de RE como um todo. Esse resultado mostra a importância da busca por um processo cada vez melhor e mais aperfeiçoado de anotação.



(a)



(b)



(c)

Figura 4.13: Teste com resposta razoável: (a) *tweet* original; (b) anotação predita; (c) anotação esperada

PUC-Rio - Certificação Digital Nº 1421597/CA

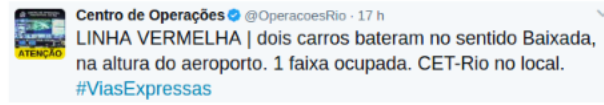
```

1 <http://inf.puc-rio.br/ontologies/TEDO/Car/0> <http://www.w3.org/1999/02/22-rdf-
  syntax-ns#type> <http://inf.puc-rio.br/ontologies/TEDO/Car> .
2 <http://inf.puc-rio.br/ontologies/TEDO/Accident/0> <http://www.w3.org/1999/02/22-rdf-
  syntax-ns#type> <http://inf.puc-rio.br/ontologies/TEDO/Accident> .
3 <http://inf.puc-rio.br/ontologies/TEDO/Accident/0> <http://inf.puc-rio.br/ontologies
  /TEDO/hasActor> <http://inf.puc-rio.br/ontologies/TEDO/Car/0> .
4 <http://inf.puc-rio.br/ontologies/TEDO/Location/2> <http://inf.puc-rio.br/ontologies
  /TEDO/hasName> "Avenida_Ayrton_Senna" .
5 <http://inf.puc-rio.br/ontologies/TEDO/Accident/0> <http://inf.puc-rio.br/ontologies
  /TEDO/causes> <http://inf.puc-rio.br/ontologies/TEDO/Interdiction/0> .
6 <http://inf.puc-rio.br/ontologies/TEDO/Interdiction/0> <http://inf.puc-rio.br/
  ontologies/TEDO/affectsFlowTo> <http://inf.puc-rio.br/ontologies/TEDO/Location
  /0> .
7 <http://inf.puc-rio.br/ontologies/TEDO/Location/1> <http://inf.puc-rio.br/ontologies
  /TEDO/hasName> "BARRA" .
8 <http://inf.puc-rio.br/ontologies/TEDO/Location/3> <http://inf.puc-rio.br/ontologies
  /TEDO/hasName> "Alvorada" .
9 <http://inf.puc-rio.br/ontologies/TEDO/Location/0> <http://www.w3.org/1999/02/22-rdf-
  syntax-ns#type> <http://inf.puc-rio.br/ontologies/TEDO/Location> .
10 <http://inf.puc-rio.br/ontologies/TEDO/Location/0> <http://inf.puc-rio.br/ontologies
  /TEDO/hasName> "praia" .
11 <http://inf.puc-rio.br/ontologies/TEDO/Interdiction/0> <http://www.w3.org
  /1999/02/22-rdf-syntax-ns#type> <http://inf.puc-rio.br/ontologies/TEDO/
  Interdiction> .
12 <http://inf.puc-rio.br/ontologies/TEDO/Location/3> <http://www.w3.org/1999/02/22-rdf-
  syntax-ns#type> <http://inf.puc-rio.br/ontologies/TEDO/Location> .
13 <http://inf.puc-rio.br/ontologies/TEDO/Location/1> <http://www.w3.org/1999/02/22-rdf-
  syntax-ns#type> <http://inf.puc-rio.br/ontologies/TEDO/Location> .
14 <http://inf.puc-rio.br/ontologies/TEDO/Location/2> <http://www.w3.org/1999/02/22-rdf-
  syntax-ns#type> <http://inf.puc-rio.br/ontologies/TEDO/Location> .
  
```

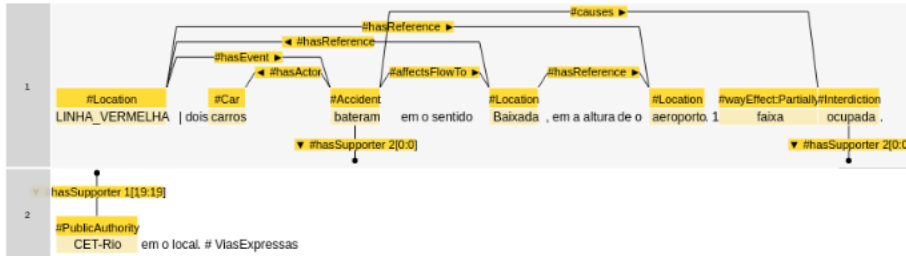
Figura 4.14: Teste com resposta razoável: triplas RDF retornadas

Já a figura 4.15, com triplas correspondentes na figura 4.16, mostra um exemplo de resposta melhor, os únicos equívocos da anotação foram a ausência

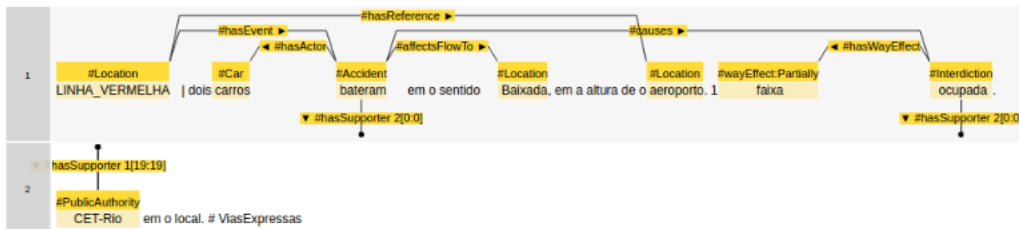
da relação *hasWayEffect* entre a entidade *Interdiction* e *wayEffect:Partially*, e a presença de algumas relações do tipo *hasReference* que não deveriam existir.



(a)



(b)



(c)

Figura 4.15: Teste com boa resposta: (a) *tweet* original; (b) anotação predita; (c) anotação esperada

```

1 <http://inf.puc-rio.br/ontologies/TEDO/Car/0> <http://www.w3.org/1999/02/22-rdf-
  syntax-ns#type> <http://inf.puc-rio.br/ontologies/TEDO/Car> .
2 <http://inf.puc-rio.br/ontologies/TEDO/Location/2> <http://inf.puc-rio.br/ontologies
  /TEDO/hasName> "Baixada" .
3 <http://inf.puc-rio.br/ontologies/TEDO/Interdiction/0> <http://inf.puc-rio.br/
  ontologies/TEDO/hasSupporter> <http://inf.puc-rio.br/ontologies/TEDO/
  PublicAuthority/0> .
4 <http://inf.puc-rio.br/ontologies/TEDO/PublicAuthority/0> <http://www.w3.org
  /1999/02/22-rdf-syntax-ns#type> <http://inf.puc-rio.br/ontologies/TEDO/
  PublicAuthority> .
5 <http://inf.puc-rio.br/ontologies/TEDO/Location/0> <http://inf.puc-rio.br/ontologies
  /TEDO/hasEvent> <http://inf.puc-rio.br/ontologies/TEDO/Accident/0> .
6 <http://inf.puc-rio.br/ontologies/TEDO/Interdiction/0> <http://www.w3.org
  /1999/02/22-rdf-syntax-ns#type> <http://inf.puc-rio.br/ontologies/TEDO/
  Interdiction> .
7 <http://inf.puc-rio.br/ontologies/TEDO/Location/0> <http://www.w3.org/1999/02/22-rdf
  -syntax-ns#type> <http://inf.puc-rio.br/ontologies/TEDO/Location> .
8 <http://inf.puc-rio.br/ontologies/TEDO/Location/2> <http://www.w3.org/1999/02/22-rdf
  -syntax-ns#type> <http://inf.puc-rio.br/ontologies/TEDO/Location> .
9 <http://inf.puc-rio.br/ontologies/TEDO/Accident/0> <http://inf.puc-rio.br/ontologies
  /TEDO/hasSupporter> <http://inf.puc-rio.br/ontologies/TEDO/PublicAuthority/0> .
10 <http://inf.puc-rio.br/ontologies/TEDO/Location/1> <http://www.w3.org/1999/02/22-rdf
  -syntax-ns#type> <http://inf.puc-rio.br/ontologies/TEDO/Location> .
11 <http://inf.puc-rio.br/ontologies/TEDO/Location/2> <http://inf.puc-rio.br/ontologies
  /TEDO/hasReference> <http://inf.puc-rio.br/ontologies/TEDO/Location/0> .
12 <http://inf.puc-rio.br/ontologies/TEDO/Location/0> <http://inf.puc-rio.br/ontologies
  /TEDO/hasName> "LINHA_VERMELHA" .
13 <http://inf.puc-rio.br/ontologies/TEDO/Location/0> <http://inf.puc-rio.br/ontologies
  /TEDO/hasReference> <http://inf.puc-rio.br/ontologies/TEDO/Location/1> .
14 <http://inf.puc-rio.br/ontologies/TEDO/Location/2> <http://inf.puc-rio.br/ontologies
  /TEDO/hasReference> <http://inf.puc-rio.br/ontologies/TEDO/Location/1> .
15 <http://inf.puc-rio.br/ontologies/TEDO/Accident/0> <http://inf.puc-rio.br/ontologies
  /TEDO/causes> <http://inf.puc-rio.br/ontologies/TEDO/Interdiction/0> .
16 <http://inf.puc-rio.br/ontologies/TEDO/Accident/0> <http://www.w3.org/1999/02/22-rdf
  -syntax-ns#type> <http://inf.puc-rio.br/ontologies/TEDO/Accident> .
17 <http://inf.puc-rio.br/ontologies/TEDO/Accident/0> <http://inf.puc-rio.br/ontologies
  /TEDO/hasActor> <http://inf.puc-rio.br/ontologies/TEDO/Car/0> .
18 <http://inf.puc-rio.br/ontologies/TEDO/Accident/0> <http://inf.puc-rio.br/ontologies
  /TEDO/affectsFlowTo> <http://inf.puc-rio.br/ontologies/TEDO/Location/2> .
19 <http://inf.puc-rio.br/ontologies/TEDO/Location/1> <http://inf.puc-rio.br/ontologies
  /TEDO/hasName> "aeroporto" .

```

Figura 4.16: Teste com boa resposta: triplas RDF retornadas

O terceiro e último exemplo, ilustrado na figura 4.17 e com triplas na figura 4.18, mostra um caso de extração perfeita das informações, combinando exatamente com o resultado de uma anotação manual em termos de NER e RE.

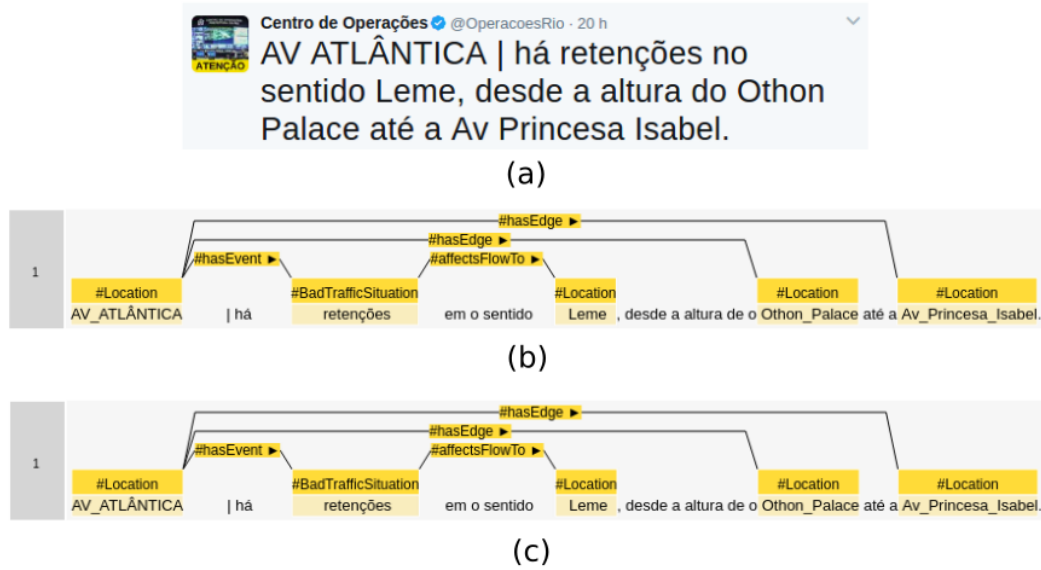


Figura 4.17: Teste com resposta perfeita: (a) *tweet* original; (b) anotação predita; (c) anotação esperada

PUC-Rio - Certificação Digital Nº 1421597/CA

```

1 <http://inf.puc-rio.br/ontologies/TEDO/Location/1> <http://www.w3.org/1999/02/22-rdf
  -syntax-ns#type> <http://inf.puc-rio.br/ontologies/TEDO/Location> .
2 <http://inf.puc-rio.br/ontologies/TEDO/Location/2> <http://www.w3.org/1999/02/22-rdf
  -syntax-ns#type> <http://inf.puc-rio.br/ontologies/TEDO/Location> .
3 <http://inf.puc-rio.br/ontologies/TEDO/BadTrafficSituation/0> <http://www.w3.org
  /1999/02/22-rdf-syntax-ns#type> <http://inf.puc-rio.br/ontologies/TEDO/
  BadTrafficSituation> .
4 <http://inf.puc-rio.br/ontologies/TEDO/Location/0> <http://inf.puc-rio.br/ontologies
  /TEDO/hasEdge> <http://inf.puc-rio.br/ontologies/TEDO/Location/1> .
5 <http://inf.puc-rio.br/ontologies/TEDO/Location/2> <http://inf.puc-rio.br/ontologies
  /TEDO/hasName> "Av_Princesa_Isabel" .
6 <http://inf.puc-rio.br/ontologies/TEDO/Location/0> <http://inf.puc-rio.br/ontologies
  /TEDO/hasName> "AV_ATLÂNTICA" .
7 <http://inf.puc-rio.br/ontologies/TEDO/Location/1> <http://inf.puc-rio.br/ontologies
  /TEDO/hasName> "Othon_Palace" .
8 <http://inf.puc-rio.br/ontologies/TEDO/BadTrafficSituation/0> <http://inf.puc-rio.br
  /ontologies/TEDO/affectsFlowTo> <http://inf.puc-rio.br/ontologies/TEDO/Location
  /3> .
9 <http://inf.puc-rio.br/ontologies/TEDO/Location/3> <http://inf.puc-rio.br/ontologies
  /TEDO/hasName> "Leme" .
10 <http://inf.puc-rio.br/ontologies/TEDO/Location/0> <http://www.w3.org/1999/02/22-rdf
  -syntax-ns#type> <http://inf.puc-rio.br/ontologies/TEDO/Location> .
11 <http://inf.puc-rio.br/ontologies/TEDO/Location/0> <http://inf.puc-rio.br/ontologies
  /TEDO/hasEdge> <http://inf.puc-rio.br/ontologies/TEDO/Location/2> .
12 <http://inf.puc-rio.br/ontologies/TEDO/Location/3> <http://www.w3.org/1999/02/22-rdf
  -syntax-ns#type> <http://inf.puc-rio.br/ontologies/TEDO/Location> .
13 <http://inf.puc-rio.br/ontologies/TEDO/Location/0> <http://inf.puc-rio.br/ontologies
  /TEDO/hasEvent> <http://inf.puc-rio.br/ontologies/TEDO/BadTrafficSituation/0> .
  
```

Figura 4.18: Teste com resposta perfeita: triplas RDF retornadas

*I know not all that may be coming, but be it
what it will, I'll go to it laughing.*

Herman Melville, *Moby-Dick; or, The Whale*.

5 Conclusão

Este trabalho propôs a implementação de um sistema *web*, o LER, para a execução das atividades de curadoria de dados, aprendizado de máquina e construção de serviços. A etapa de curadoria é de tal modo importante e complexa que inspirou a implementação de um subsistema especial, o ERAS, como subsistema do LER. A etapa de aprendizado de máquina realiza o treinamento de modelos para as tarefas de identificação de entidades (NER) e de extração de relações entre as entidades identificadas (RE). A última etapa, a de construção de serviços, combina classificadores de NER e RE previamente treinados com dados gerados na etapa de curadoria, disponibilizando serviços preparados para receberem dados textuais puros e exportarem a informação estruturada extraída do texto.

O ERAS foi implementado com o objetivo de ser uma ferramenta robusta e de simples uso por pessoal não especializado na tarefa de anotação. Para fins de avaliação do estado atual deste sistema, um grupo de anotadores de diferentes origens e perfis foi recrutado para uma tarefa de anotação de *tweets* de trânsito, tendo como referência uma versão levemente modificada da ontologia TEDO, proposta em (6). Diversas questões relacionadas ao processo de anotação foram analisadas e questionários foram preenchidos. Os resultados do experimento mostraram que a ferramenta atendeu satisfatoriamente o que fora proposto neste trabalho, pois ela se mostrou muito adequada à tarefa de anotação e recebeu muitos elogios por parte dos participantes do experimento. Durante a condução dos experimento tornou-se claro o fato de que a construção de uma ontologia não é um processo estático, e que a definição de uma que se adéque bem a dados textuais, é uma tarefa muito difícil e trabalhosa, porém essa etapa do processo é de suma importância e deve ser executada com atenção, já que afeta diretamente a qualidade dos dados anotados. Outro ponto que pôde-se notar nestes resultados foi que a tarefa de identificação de relações entre entidades no domínio de *tweets* de trânsito é complexa e por vezes ambígua, o que acaba gerando anotações finais com qualidade abaixo do esperado, porém notou-se também que ter um especialista da tarefa dando suporte aos anotadores é de suma importância.

O restante do sistema LER, isto é, as etapas de aprendizado automático e criação de serviços, foi testado pelo uso de um conjunto de 600 *tweets* de trânsito anotados exclusivamente pelo autor deste trabalho, tendo como base ontológica uma versão da TEDO otimizada a partir das experiências de anotação citadas acima. Para as tarefas de NER e RE, diferentes conjuntos de atributos, classificadores, algoritmos e parâmetros foram avaliados, selecionando os mais adequados usando os valores de F1 e tempos de treinamento, nesta ordem, como principais critérios. Para a tarefa NER, obteve-se um melhor resultado no conjunto de teste (F1 = 0.821, Precision = 0.848 e Recall =

0.800) pelo uso do classificador Random Forest com 50 árvores (entre outros parâmetro ótimos), atributos de contexto gerados automaticamente pela ferramenta e informações completas de POS *Tagging*. Para a tarefa RE, obteve-se um melhor resultado no conjunto de teste ($F1 = 0.701$, $Precision = 0.769$ e $Recall = 0.655$) utilizando o algoritmo Frank Wolfe SSVM, atributos relativos às entidades identificadas e atributos de contexto gerados automaticamente pela ferramenta (fazendo intenso uso dos conectores registrados na etapa de anotação). Alguns testes acerca da redução de classes por generalização baseada na ontologia foram feitos e não melhoraram os resultados de NER e RE, exceto em relação ao tempo de treinamento.

Os melhores modelos foram selecionados e alocados em série como um serviço, para a demonstração da automação complexa do fluxo de extração das informações em novos *tweets* de trânsito nunca “vistos” pelos classificadores. Os resultados são interessantes, mas melhorias podem ser feitas. Neste caso específico, a perda da eficiência do serviço reside no modelo para RE, que teve um pior F1 comparado ao obtido para NER. Contudo, novos trabalhos deverão visar a melhoria de ambos os classificadores.

Recomenda-se que trabalhos futuros atuem em diversas linhas:

- Realização de novos experimentos de anotação com um número maior de anotadores, por mais tempo e com mais documentos. A ferramenta ERAS está preparada para registrar diversos dados relativos ao comportamento dos anotadores, mesmo relativos à dinâmica das concordâncias. Os resultados do aprendizado de máquina deixam claro, como infere-se por pura lógica, que o processo de anotação, se otimizado, produz melhores classificadores.
- Uso de técnicas de *process mining* (44) nos *logs* de anotação, isto é, uma atividade derivada do que se diz no item acima. Este tipo de análise seria interessante, por exemplo, para descobrir quais fluxos de uso levam a maiores concordâncias, o que seria uma informação muito útil já que afeta diretamente a qualidade dos modelos treinados;
- Melhoria no *feedback* dado nas anotações. Atualmente o único tipo de retorno que o anotador pode dar via sistema se resume a uma caixa de texto livre em cada documento anotado. Abordagens onde o usuário também desse o *feedback* no nível das entidades, relações e conectores poderiam trazer mais informações úteis ao processo de curadoria. Outra função interessante seria a possibilidade do anotador atribuir um índice de incerteza após efetuar determinada anotação, o que facilitaria a separação dos documentos entre níveis de dificuldade e identificação de melhorias na ontologia definida para o domínio;
- Criar mecanismos que auxiliem o curador na construção da ontologia. Um curador poderia definir apenas algumas entidades e relações genéricas e o LER o ajudaria, por exemplo, a identificar quais classes deveriam ser especializadas ou generalizadas. Isso tornaria o sistema mais proativo e amigável para pessoas que não estão familiarizadas com a construção de ontologias;
- Teste do sistema LER em outros domínios, diferentes de *tweets* de trânsito. Um trabalho relativo a relatórios operacionais da área de

petróleo e gás (45) já demonstrou, com bons resultados, a viabilidade do uso do LER em um caso demasiadamente complexo e desafiador.

- Aplicação de técnicas de seleção automática ou indução de atributos (46), para melhoria da capacidade de classificação dos algoritmos atualmente usados.
- Aplicação de outros métodos para NER e RE, inclusive considerando-as como uma única tarefa (metodologia “end-to-end”) (47). Recomenda-se também testes do uso de *Deep Learning* para NER e RE, conforme avaliado por outros autores (48).

Referências bibliográficas

- [1] FILIPPOV, S.. **Mapping text and data mining in academic and research communities in Europe**. Lisbon Council, 2014.
- [2] ITTOO, A.; NGUYEN, L. M. ; VAN DEN BOSCH, A.. **Text analytics in industry: Challenges, desiderata and trends**. *Computers in Industry*, 78:96–107, 2016.
- [3] HE, W.; ZHA, S. ; LI, L.. **Social media competitive analysis and text mining: A case study in the pizza industry**. *International Journal of Information Management*, 33(3):464–472, 2013.
- [4] ATSERIAS, J.; CODINA, J.. **What is the text of a tweet?** In: *WORKSHOP PROGRAMME*, p. 29.
- [5] REDLICH, L. R.. **Modelagem de eventos de trânsito com base em clipping de grandes massas de dados da web**. Master's thesis, PUC-Rio, 2013.
- [6] ALBUQUERQUE, F. C.; CASANOVA, M. A.; LOPES, H.; REDLICH, L. R.; DE MACEDO, J. A. F.; LEMOS, M.; DE CARVALHO, M. T. M. ; RENSO, C.. **A methodology for traffic-related twitter messages interpretation**. *Computers in Industry*, 78:57–69, 2016.
- [7] DA COSTA ALBUQUERQUE, F.; CASANOVA, M. A.; DE MACEDO, J. A. F.; DE CARVALHO, M. T. M. ; RENSO, C.. **A proactive application to monitor truck fleets**. In: *MOBILE DATA MANAGEMENT (MDM), 2013 IEEE 14TH INTERNATIONAL CONFERENCE ON*, volumen 1, p. 301–304. IEEE, 2013.
- [8] SOWA, J. F.; OTHERS. **Knowledge representation: logical, philosophical, and computational foundations**, volumen 13. MIT Press, 2000.
- [9] KANEIWA, K.; IWAZUME, M. ; FUKUDA, K.. **An upper ontology for event classifications and relations**. In: *AUSTRALASIAN JOINT CONFERENCE ON ARTIFICIAL INTELLIGENCE*, p. 394–403. Springer, 2007.
- [10] WORBOYS, M.; HORNSBY, K.. **From objects to events: Gem, the geospatial event model**. In: *INTERNATIONAL CONFERENCE ON GEOGRAPHIC INFORMATION SCIENCE*, p. 327–343. Springer, 2004.
- [11] MOTTA, E. N.; FERNANDES, E. R. ; MILIDIÚ, R. L.. **F-ext-ws-2.0: A web service for natural language processing**. 2010.

- [12] PLATT, J. C.. **12 fast training of support vector machines using sequential minimal optimization**. Advances in kernel methods, p. 185–208, 1999.
- [13] HALL, M.; FRANK, E.; HOLMES, G.; PFAHRINGER, B.; REUTEMANN, P.; WITTEN, I. H.. **The weka data mining software: an update**. ACM SIGKDD explorations newsletter, 11(1):10–18, 2009.
- [14] FERNANDES, E. R.; DOS SANTOS, C. N. ; MILIDIÚ, R. L.. **Latent structure perceptron with feature induction for unrestricted coreference resolution**. In: JOINT CONFERENCE ON EMNLP AND CONLL-SHARED TASK, p. 41–48. Association for Computational Linguistics, 2012.
- [15] CHU, T. C.; HUANG, R.. **Sourcing the crowd for a few good ones: Event type detection**. In: 24TH INTERNATIONAL CONFERENCE ON COMPUTATIONAL LINGUISTICS, p. 1239, 2012.
- [16] MADLBERGER, L.; ROMADHONY, A. ; PURWARIANTI, A.. **Gotong royong in nlp research a mobile tool for collaborative text annotation in indonesia**. In: ASIAN LANGUAGE PROCESSING (IALP), 2016 INTERNATIONAL CONFERENCE ON, p. 99–102. IEEE, 2016.
- [17] HOVY, E.; LAVID, J.. **Towards a ‘science’of corpus annotation: a new methodological challenge for corpus linguistics**. International journal of translation, 22(1):13–36, 2010.
- [18] OF SHEFFIELD, T. U.. **Gate home**. <https://gate.ac.uk/>, 2017. Acesso em 07 de abril de 2017.
- [19] CUNNINGHAM, H.; MAYNARD, D.; BONTCHEVA, K. ; TABLAN, V.. **Gate: an architecture for development of robust hlt applications**. In: PROCEEDINGS OF THE 40TH ANNUAL MEETING ON ASSOCIATION FOR COMPUTATIONAL LINGUISTICS, p. 168–175. Association for Computational Linguistics, 2002.
- [20] CUNNINGHAM, H.; MAYNARD, D. ; BONTCHEVA, K.. **Text processing with gate**. Gateway Press CA, 2011.
- [21] STENETORP, P.; PYYSALO, S.; TOPIĆ, G.; OHTA, T.; ANANIADOU, S. ; TSUJII, J.. **Brat: a web-based tool for nlp-assisted text annotation**. In: PROCEEDINGS OF THE DEMONSTRATIONS AT THE 13TH CONFERENCE OF THE EUROPEAN CHAPTER OF THE ASSOCIATION FOR COMPUTATIONAL LINGUISTICS, p. 102–107. Association for Computational Linguistics, 2012.
- [22] BONTCHEVA, K.; DERCZYNSKI, L.; FUNK, A.; GREENWOOD, M. A.; MAYNARD, D. ; ASWANI, N.. **Twitie: An open-source information extraction pipeline for microblog text**. In: RANLP, p. 83–90, 2013.
- [23] YIMAM, S. M.; GUREVYCH, I.; DE CASTILHO, R. E. ; BIEMANN, C.. **We-banno: A flexible, web-based and visually supported system for distributed annotations**. In: ACL (CONFERENCE SYSTEM DEMONSTRATIONS), p. 1–6, 2013.

- [24] MÜLLER, C.; STRUBE, M.. **Multi-level annotation of linguistic data with mmax2**. *Corpus technology and language pedagogy: New resources, new tools, new methods*, 3:197–214, 2006.
- [25] MORTON, T.; LACIVITA, J.. **Wordfreak: an open tool for linguistic annotation**. In: PROCEEDINGS OF THE 2003 CONFERENCE OF THE NORTH AMERICAN CHAPTER OF THE ASSOCIATION FOR COMPUTATIONAL LINGUISTICS ON HUMAN LANGUAGE TECHNOLOGY: DEMONSTRATIONS-VOLUME 4, p. 17–18. Association for Computational Linguistics, 2003.
- [26] WISSLER, L.; ALMASHRAEE, M.; DÍAZ, D. M. ; PASCHKE, A.. **The gold standard in corpus annotation**. In: IEEE GSC, 2014.
- [27] CAMPBELL, B.; JONES, M. ; MORTIMORE, C.. **Json web token (jwt) profile for oauth 2.0 client authentication and authorization grants**. 2015.
- [28] CENTER, T. R.. **Welcome | Freeling Homepage**. <http://nlp.cs.upc.edu/freeling/node/1>, 2017. [Acesso em 20 de Março de 2017].
- [29] PADRÓ, L.; STANILOVSKY, E.. **Freeling 3.0: Towards wider multilinguality**. In: PROCEEDINGS OF THE LANGUAGE RESOURCES AND EVALUATION CONFERENCE (LREC 2012), Istanbul, Turkey, May 2012. ELRA.
- [30] MELVILLE, H.. **Moby-dick; or, the whale. 1851**. Ed. Harrison Hayford et al. Evanston: Northwestern UP and the Newberry Library, 1988.
- [31] UNIVERSALDEPENDENCIES. **Conll-u format**. <http://universaldependencies.org/format.html>, 2017. [Online; accessed 9-april-2017].
- [32] CARLETTA, J.. **Assessing agreement on classification tasks: the kappa statistic**. *Computational linguistics*, 22(2):249–254, 1996.
- [33] WIKIPEDIA. **Hyperparameter optimization**. https://en.wikipedia.org/w/index.php?title=Hyperparameter_optimization, 2017. [Acesso em 26 de Março de 2017].
- [34] CHANG, C.-C.; LIN, C.-J.. **Libsvm: a library for support vector machines**. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 2(3):27, 2011.
- [35] BREIMAN, L.. **Random forests**. *Machine learning*, 45(1):5–32, 2001.
- [36] BOTTOU, L.. **Large-scale machine learning with stochastic gradient descent**. In: PROCEEDINGS OF COMPSTAT'2010, p. 177–186. Springer, 2010.

- [37] HUANG, L.; FAYONG, S. ; GUO, Y.. **Structured perceptron with inexact search**. In: PROCEEDINGS OF THE 2012 CONFERENCE OF THE NORTH AMERICAN CHAPTER OF THE ASSOCIATION FOR COMPUTATIONAL LINGUISTICS: HUMAN LANGUAGE TECHNOLOGIES, p. 142–151. Association for Computational Linguistics, 2012.
- [38] LACOSTE-JULIEN, S.; JAGGI, M.; SCHMIDT, M. ; PLETSCHER, P.. **Block-coordinate frank-wolfe optimization for structural svms**. arXiv preprint arXiv:1207.4747, 2012.
- [39] POWERS, D. M.. **Evaluation: from precision, recall and f-measure to roc, informedness, markedness and correlation**. 2011.
- [40] CROCKFORD, D.. **The application/json media type for javascript object notation (json)**. 2006.
- [41] BECKETT, D.. **Rdf 1.1 n-triples**. W3C Recommendation, 2014.
- [42] SCIKIT LEARN. **Preprocessing data**. <http://scikit-learn.org/stable/modules/preprocessing.html>, 2017. [Online; accessed 9-april-2017].
- [43] SOLÉ, X.; RAMISA, A. ; TORRAS, C.. **Evaluation of random forests on large-scale classification problems using a bag-of-visual-words representation**. In: CCIA, p. 273–276, 2014.
- [44] AALST, W.; ADRIANSYAH, A.; MEDEIROS, A. K. A.; ARCIERI, F.; BAIER, T.; BLICKLE, T.; BOSE, J. C.; BRAND, P.; BRANDTJEN, R.; BUIJS, J. ; OTHERS. **Process mining manifesto**. In: BUSINESS PROCESS MANAGEMENT WORKSHOPS, p. 169–194. Springer, 2012.
- [45] FURTADO, T.. **Interpretação automática de relatórios de operação de equipamentos**. Master's thesis, PUC-Rio, 2017.
- [46] MOTTA, E. N.. **Indução e Seleção Incrementais de Atributos no Aprendizado Supervisionado**. PhD thesis, PUC-Rio, 2014.
- [47] MIWA, M.; BANSAL, M.. **End-to-end relation extraction using lstms on sequences and tree structures**. arXiv preprint arXiv:1601.00770, 2016.
- [48] CHIU, J. P.; NICHOLS, E.. **Named entity recognition with bidirectional lstm-cnns**. arXiv preprint arXiv:1511.08308, 2015.

A

Experimento de anotação: Guia de anotação

O documento contido nas seções seguintes é o guia de anotação que foi dado aos participantes do experimento de anotação descrito no capítulo 4 desta dissertação.

A.1

FERRAMENTA DE ANOTAÇÃO

Para a tarefa de anotação que será explicada na seção A.2 deste documento, o anotador utilizará uma ferramenta online implementada para tal fim. Esta ferramenta fará uso de uma ontologia para a tarefa de anotação, esta ontologia é basicamente uma série de regras que definem a estrutura de um domínio, ou seja, quais entidades fazem parte deste domínio, quais são as relações possíveis entre estas entidades e quais as restrições associadas à estas relações. A ferramenta que será utilizada no processo de anotação (Figura A.1) tem basicamente as seguintes características:

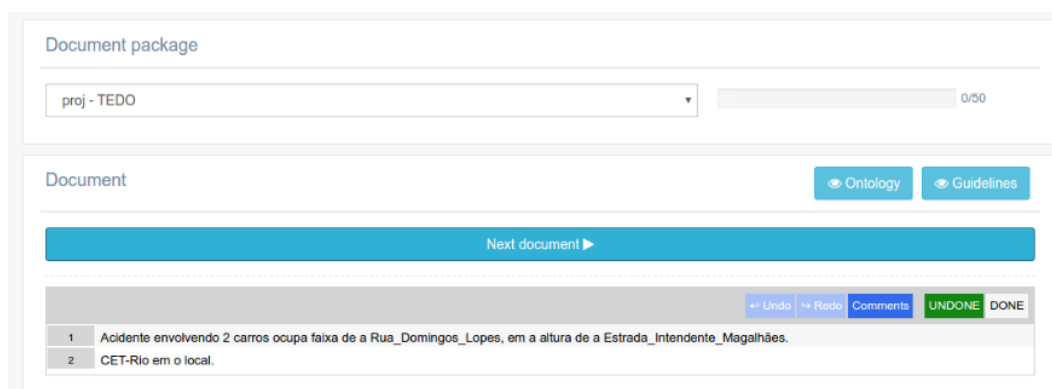


Figura A.1: Ferramenta de anotação

- Seletor de pacotes de documentos disponíveis para anotação;
- Indicador de progresso, mostrando quantos documentos estão disponíveis para anotação e quanto já foram finalizados;
- Apresentação de um resumo da ontologia que descreve o domínio da anotação (é mostrado quando se clica no botão “Ontology”, como pode ser visto na Figura A.2, este resumo está descrito no formato de árvore: no 1º nível estão todas as entidades disponíveis para rotulação, no 2º as relações possíveis que podem partir da entidade em questão e no 3º as entidades que podem ser atingidas pela relação);

- Guia de anotação (este documento), que é apresentado quando o botão “Guidelines” é acionado;
- Função de ir para o próximo documento (os documentos são apresentados aleatoriamente para anotação, porém não há obrigação de finalizar os documentos na ordem apresentada, é possível anotar parte de um documento e ir para outro caso queira, que em algum momento no futuro o documento não finalizado será reapresentado);
- Capacidade de adicionar comentários ao documento (acionando o botão “Comments”), como pode ser visto na Figura A.3;
- Desfazer (acionando o botão undo) ou refazer (acionando o botão redo) ações;
- Finalizar a anotação de um documento (acionando o botão “DONE”, uma vez finalizado o documento não será mais apresentado para o anotador);
- Possibilidade de rotulação das palavras, criação de relações entre os rótulos e associação de palavras às relações;
- Todas as ações que afetam o documento anotado são salvas automaticamente.

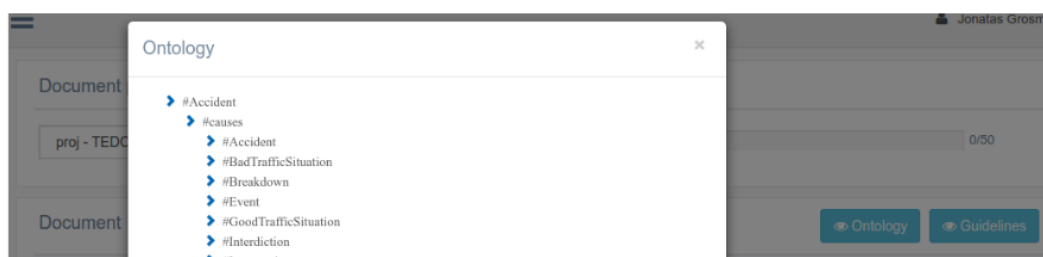


Figura A.2: Resumo da ontologia de anotação

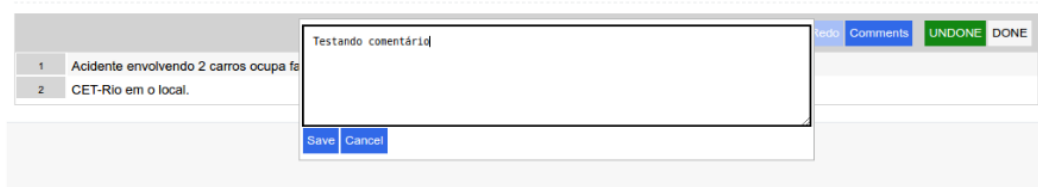


Figura A.3: Adicionando comentários

A.1.1

PROCESSO DE ANOTAÇÃO

O processo de anotação pode ser dividido basicamente em 3 partes, rotulação de palavras, criação de relacionamento entre rótulos e criação de conectores. O processo de rotulação é a associação de um bloco de texto a uma determinada entidade descrita na ontologia que modela o domínio do documento. Este processo segue os seguintes passos:

- Selecionar a palavra ou conjunto de palavras que serão rotuladas (Figura A.4). Não é necessário selecionar toda a palavra que será rotulada apenas parte dela já é o suficiente;
- Abrir o menu de contexto (clitando com o botão direito do mouse) e selecionar o rótulo desejado (Figura A.5), logo após isso o rótulo estará associado à(s) palavra(s) selecionada(s) (Figura A.6).

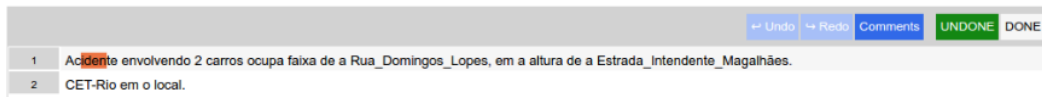


Figura A.4: Rotulação, passo 1

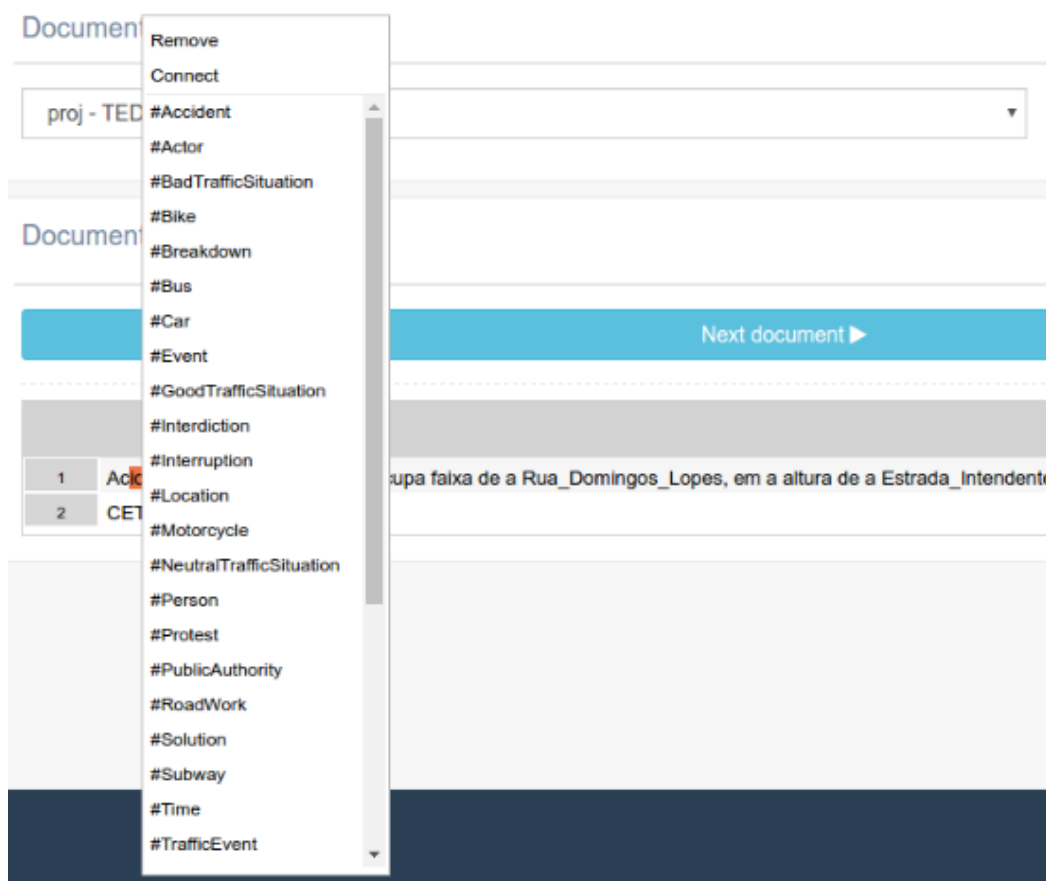


Figura A.5: Rotulação, passo 2

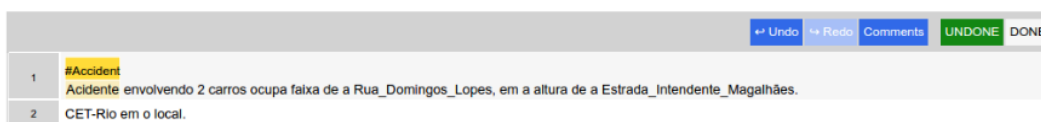


Figura A.6: Rotulação, resultado final

O processo de criação de relacionamento entre rótulos é a conexão de um bloco de texto rotulado à outro, de tal forma que descreva uma relação prevista na ontologia, esse processo segue os seguintes passos:

- Clicar com o botão esquerdo do mouse sobre o rótulo que será a origem da relação, após isso um arco será criado, ir em direção ao rótulo que será o destino da relação, caso a relação seja inválida o arco ficará vermelho (Figura A.7) e a relação não poderá ser criada (Isso ocorre por restrições definidas na ontologia, como é possível ver na Figura A.8, não existe relação possível partindo de #Car para a #Accident, porém o inverso existe), caso contrário ficará verde (Figura A.9) e para finalizar o processo

basta clicar sobre o rótulo destino que a relação estará feita (Figura A.10).

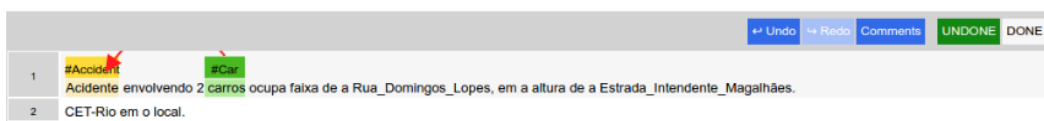


Figura A.7: Relacionamento entre os rótulos (relação inválida)



Figura A.8: Relação inválida de #Car para #Accident, porém válida de #Accident para #Car

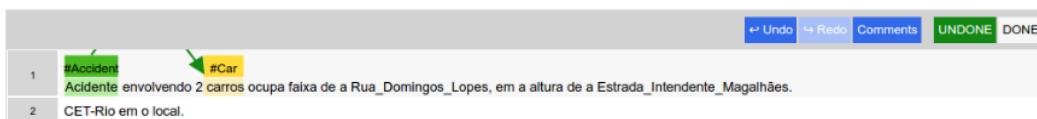


Figura A.9: Relacionamento entre os rótulos (relação válida)

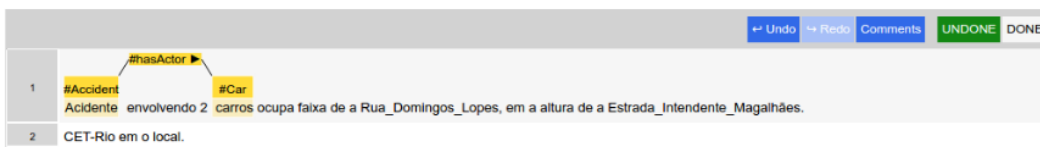


Figura A.10: Relacionamento entre os rótulos, resultado final

No processo de relacionamento não é necessário que os rótulos estejam na mesma sentença, como pode ser visto na Figura A.11, a descrição desse tipo

de relação contém a informação da sentença e palavras que estão na ponta oposta da relação. Tal informação está no formato SENTENÇA[PALAVRA INICIAL : PALAVRA FINAL], ou seja, a descrição 1[2:3] pode ser lida como “a ponta oposta da relação está ligada ao conjunto de palavras na sentença 1 do intervalo de 2 a 3”.

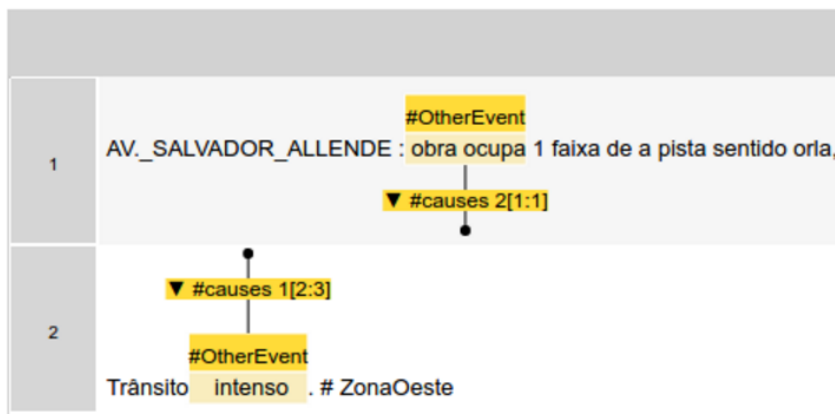


Figura A.11: Relações entre sentenças

O processo de criação de conectores é a conexão de um bloco de texto a uma relação, de tal forma que as palavras que compõem tal bloco descrevam a relação em questão. Esse processo segue os seguintes passos:

- Selecionar a palavra ou conjunto de palavras que serão associadas à relação (Figura A.12), não é necessário selecionar toda a palavra apenas parte dela já é o suficiente;
- Abrir o menu de contexto (clitando com o botão direito do mouse) e selecionar o item “Connect” (Figura A.13);
- Após aparecer o arco de conexão (Figura A.14), mova o mouse até a relação desejada e clique com o botão esquerdo, feito isso o bloco estará associado à relação (Figura A.15).

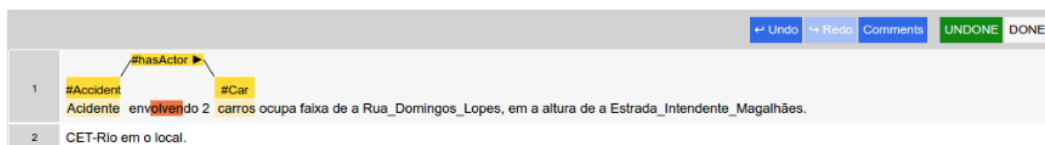


Figura A.12: Criação de conectores, passo 1

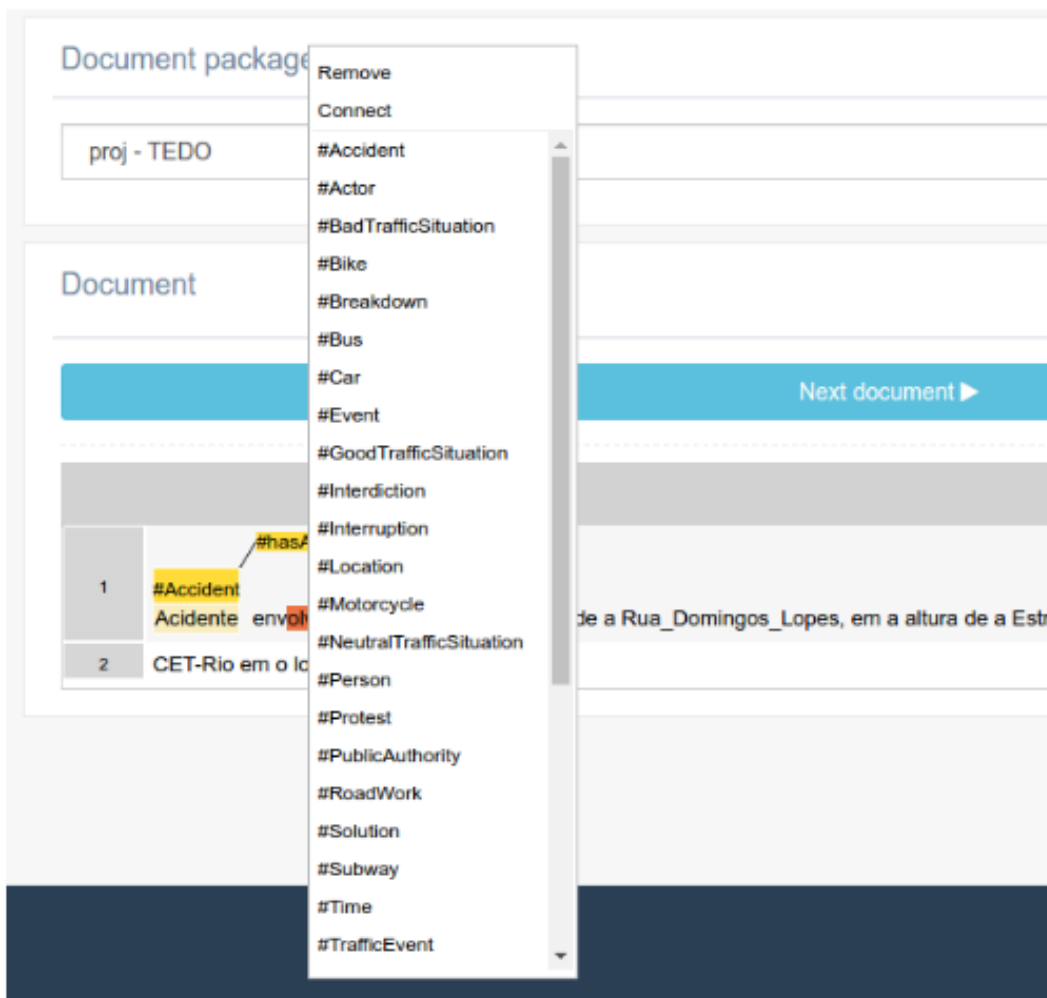


Figura A.13: Criação de conectores, passo 2

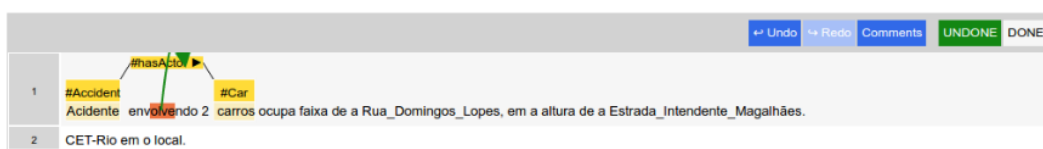


Figura A.14: Criação de conectores, passo 3

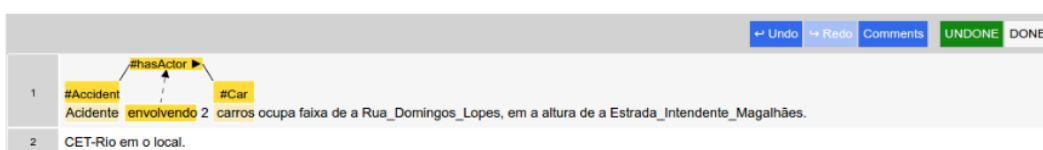


Figura A.15: Criação de conectores, resultado final

A remoção de rótulos, relações e conectores, pode ser feita clicando com o botão direito do mouse sobre o item em questão (Figura A.16) e clicando em “Remove” ou pode ser feita removendo os itens presentes em uma área do texto selecionada, basta apenas selecionar a área desejada (Figura A.17), abrir o menu de contexto com o botão direito do mouse (Figura A.18), e clicar na opção “Remove” que então toda a área selecionada será limpa (Figura A.19).

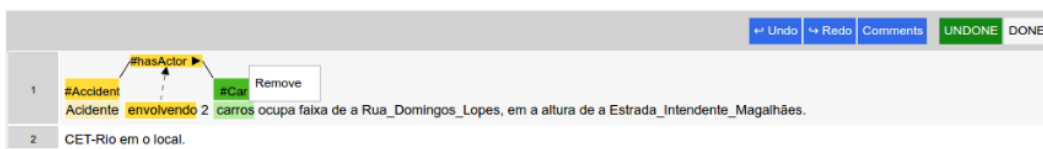


Figura A.16: Remoção de rótulos, relações e conectores

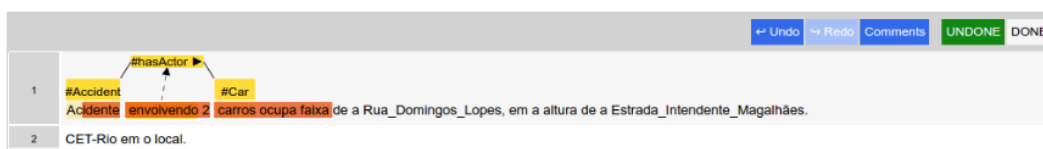


Figura A.17: Remoção por área de de rótulos, relações e conectores, passo 1

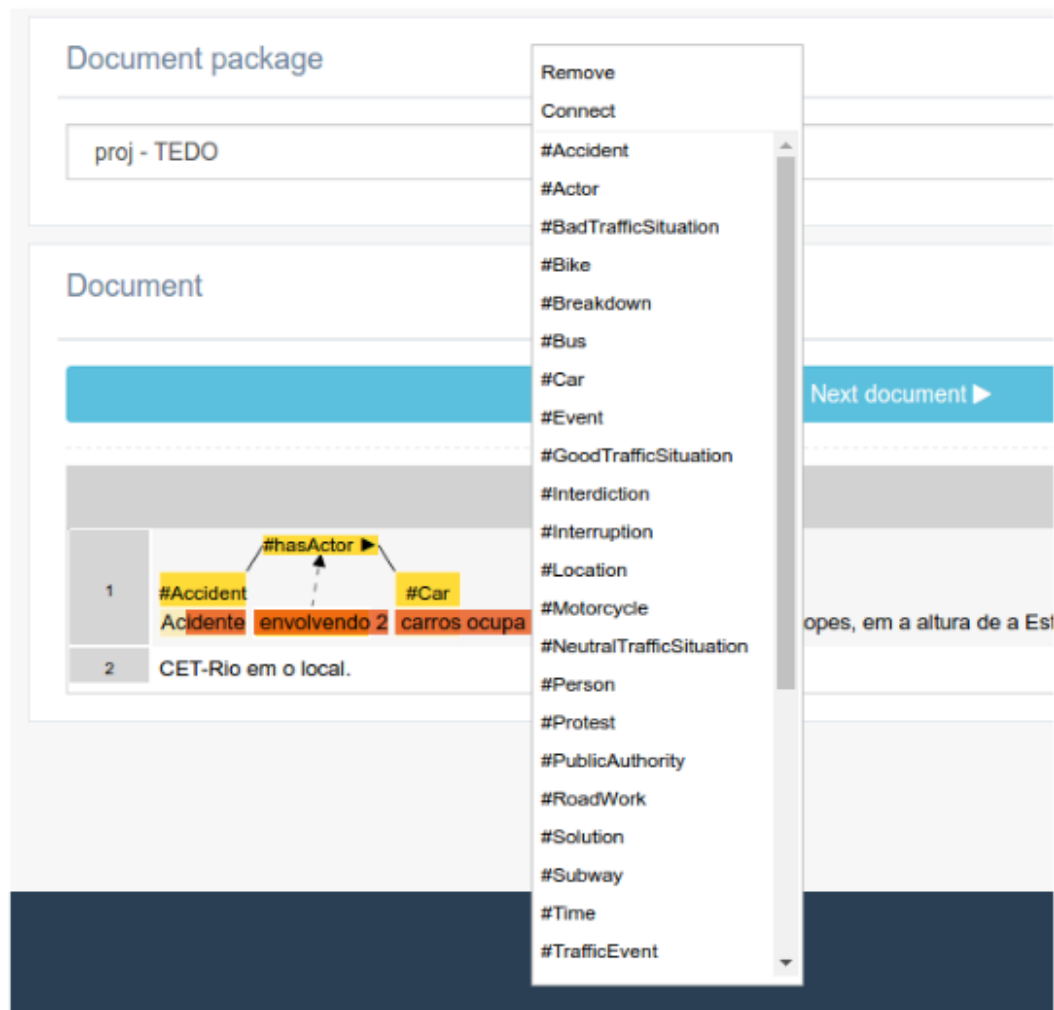


Figura A.18: Remoção por área de de rótulos, relações e conectores, passo 2

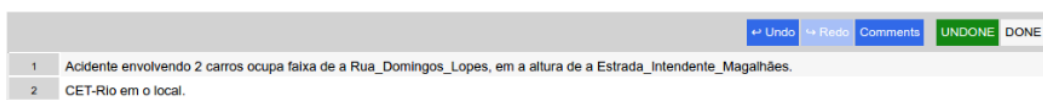


Figura A.19: Remoção por área de de rótulos, relações e conectores, resultado final

Assim que o processo de anotação estiver finalizado basta clicar no botão “DONE” que o documento será fechado para edição e não aparecerá mais para ser anotado, caso existam mais documentos para serem anotados, após finalizado o documento corrente, um novo documento será apresentado aleatoriamente (o mesmo que aconteceria se clicar no botão “Next document”), caso contrário a barra de status ficará em 100% e não aparecerão mais documentos para serem anotados.

A.1.2 ATALHOS

Com intuito de tornar o processo de anotação mais prático e ágil alguns atalhos estão disponíveis na ferramenta:

- Tecla T (Rotulação): Após selecionar um bloco de texto e acionar a tecla, o último rótulo utilizado (em qualquer parte do documento) será associado ao bloco selecionado. Este atalho é bastante útil quando há N ocorrências de uma mesma entidade da ontologia em um mesmo documento, neste cenário, basta apenas que se atribua um rótulo seguindo o fluxo padrão descrito na seção A.1.1 e então os demais podem ser atribuídos utilizando a tecla T;
- Tecla C (criação de conectores): Basta apenas selecionar a área desejada e acionar o atalho que o arco de conexão será criado e poderá associar o bloco de texto à relação desejada.
- Tecla R (remoção de uma área do texto): Basta apenas selecionar a área desejada e acionar o atalho que todos os itens associados a área selecionada serão removidos.
- Tecla Z (desfazer): Basta apenas acionar o atalho que a última ação executada será desfeita.
- Tecla Y (refazer): Basta apenas acionar o atalho que a última ação desfeita será refeita.

A.2 TAREFA DE ANOTAÇÃO

Como já foi dito na seção anterior, a tarefa de anotação está associada a uma ontologia que descreve o domínio dos dados a serem anotados, na tarefa em particular descrita por este documento o domínio são tweets que descrevem eventos de trânsito na cidade do Rio de Janeiro, esta seção tem como objetivo descrever especificamente como tais dados devem ser anotados, mostrando sempre que possível exemplos destas anotações.

A.2.1 RÓTULOS

Na Figura A.20 é possível ver todas as classes da ontologia que será utilizada na tarefa de anotação. Todas estas classes estão disponíveis para serem rótulos de partes dos tweets de trânsito que deverão ser anotados. Nas próximas seções serão apresentadas descrições de cada uma das classes que compõem a ontologia, bem como exemplos de uso destas nos tweets.

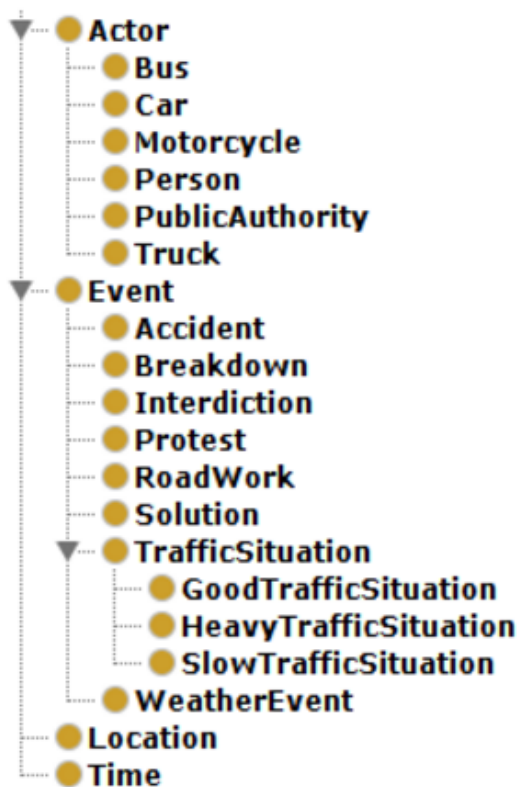


Figura A.20: Hierarquia de classes da ontologia de eventos de trânsito

A.2.1.1

Actor

Esta classe descreve os atores envolvidos nos eventos de trânsito, tal classe deve ser utilizada apenas quando não há especialização desta (Bus, Car, Motorcycle, Person, PublicAuthority e Truck) que consiga descrever o ator em questão, um exemplo (hipotético) dessa situação pode ser visto na Figura A.21, outro exemplo do uso da classe (no caso uma especialização desta) pode ser visto na Figura A.22.

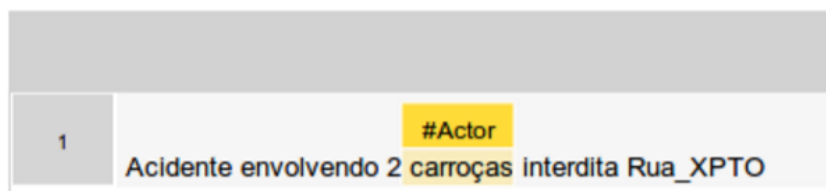


Figura A.21: Exemplo de rotulação, Actor

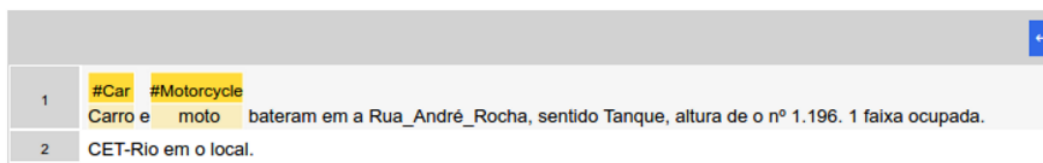


Figura A.22: Exemplo de rotulação, Car e Motorcycle

A.2.1.2 Event

Esta classe descreve os eventos descritos no tweet, e deve ser apenas utilizada para rotulação quando não há especialização desta que descreve o evento em questão, como pode ser visto na Figura A.23. As especializações da classe Event são Accident (a ocorrência de um acidente qualquer, no nosso caso, de um acidente de trânsito), Breakdown (uma falha qualquer, no nosso caso pode ser a falha de um carro por exemplo), Interdiction (a interdição de algo, no nosso caso de uma via), Protest (eventos como manifestações públicas, protestos e etc), RoadWork (uma obra acontecendo em uma via), Solution (Soluções para problemas causados por outros eventos, como remoção de carros acidentados, conserto de semáforos e etc), TrafficSituation (condições de trânsito, que podem ser especializadas em GoodTrafficSituation, HeavyTrafficSituation ou SlowTrafficSituation) e WeatherEvent (eventos climáticos, como chuva, ventanias e etc). As Figuras A.24, A.25, A.26, A.27, A.28, A.29, A.30 e A.31 mostram exemplos de uso das especializações da classe Event.

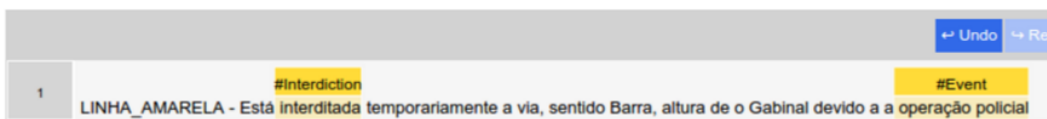


Figura A.23: Exemplo de rotulação, Interdiction e Event



Figura A.24: Exemplo de rotulação, Protest e Interdiction

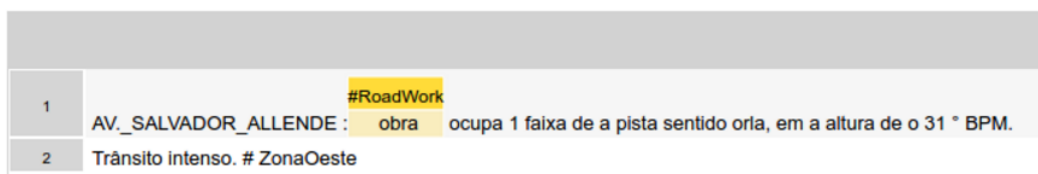


Figura A.25: Exemplo de rotulação, RoadWork



Figura A.26: Exemplo de rotulação, WeatherEvent

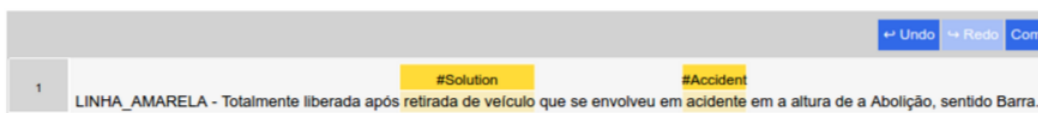


Figura A.27: Exemplo de rotulação, Solution e Accident

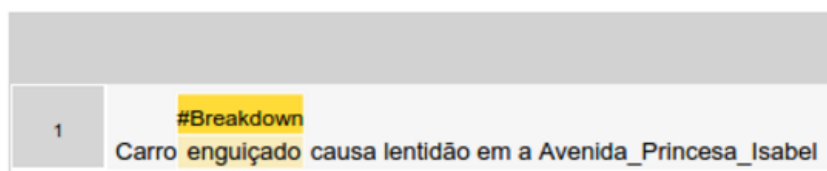


Figura A.28: Exemplo de rotulação, Breakdown

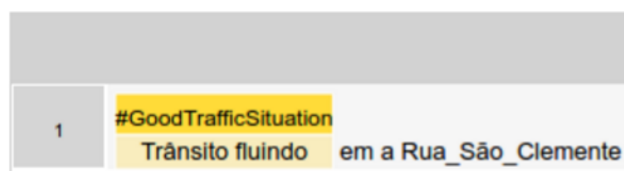


Figura A.29: GoodTrafficSituation

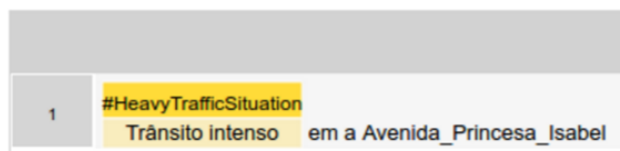


Figura A.30: HeavyTrafficSituation

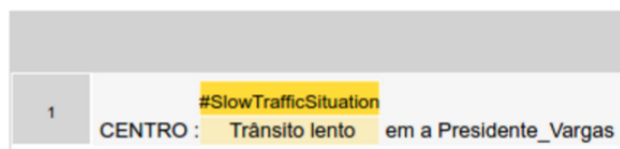


Figura A.31: SlowTrafficSituation

A.2.1.3 Location

Esta classe descreve a localização dos eventos descritos no tweet, essa classe não está necessariamente ligada apenas a endereços, como ruas, avenidas e etc, ela pode ser utilizada também em pontos de referência como estações de metrô, praças, pontos turísticos e etc. As Figuras A.32 e A.33 mostram exemplos de uso dessa classe.

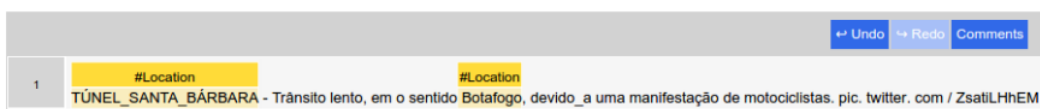


Figura A.32: Exemplo 1 de rotulação, Location



Figura A.33: Exemplo 2 de rotulação, Location

A.2.1.4

Time

Esta classe descreve o tempo em que os eventos descritos ocorreram ou ocorrerão. As Figuras A.34 e A.35 mostram exemplos de uso dessa classe.



Figura A.34: Exemplo 1 de rotulação, Time

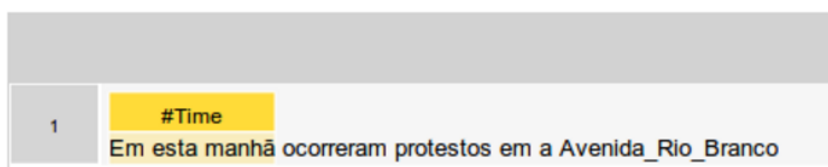


Figura A.35: Exemplo 2 de rotulação, Time

A.2.1.5

Tipos primitivos

Além das classes já citadas nas seções anteriores, como pode ser visto na Figura A.36, há tipos primitivos na ontologia que também estão disponíveis para serem utilizados como rótulos.

- #wayEffect:BothDirections
- #wayEffect:OneDirection
- #wayEffect:Partially
- xsd:string
- xsd:unsignedInt

Figura A.36: Tipos primitivos

O tipo `wayEffect:BothDirections` (Figura A.37), descreve o estado da via em que ambos os sentidos desta estão afetados por determinado evento, a mesma lógica se aplica aos tipos `wayEffect:OneDirection` (Figura A.38) e `wayEffect:Partially` (Figura A.39), onde o primeiro relata que um dos sentidos da via está comprometido e o último passa a ideia de que apenas parte da via foi afetada, porém mesmo assim o fluxo para ambos os sentidos continua acessível.

1	AV_NIEMEYER : via interditada em #wayEffect:BothDirections os dois sentidos devido_a um acidente entre ônibus e moto, alt
2	Opção : Lagoa-Barra. # alert

Figura A.37: wayEffect:BothDirections

1	Avenida_das_Américas está interditada, #wayEffect:OneDirection em o sentido São_Conrado, altura de o Pedra_de_Itaú
---	---

Figura A.38: wayEffect:OneDirection

1	AV_SALVADOR_ALLENDE : obra ocupa #wayEffect:Partially 1 faixa de a pista sentido orla, em a altura de o 31 ° BPM.
2	Trânsito intenso. # ZonaOeste

Figura A.39: wayEffect:Partially

Os tipos `xsd:unsignedInt` (Figura A.40) e `xsd:string` (Figura A.41) são utilizados apenas em uma situação específica, que é quando a quantidade de atores envolvidos no evento está presente, onde o primeiro tipo será utilizado quando esta quantidade é dada em formato numérico, e o segundo quando esta é dada em formato não numérico.

1	Acidente envolvendo xsd:unsignedInt 2 carros ocupa faixa de a Rua_Domingos_Lopes, e
2	CET-Rio em o local.

Figura A.40: xsd:unsignedInt

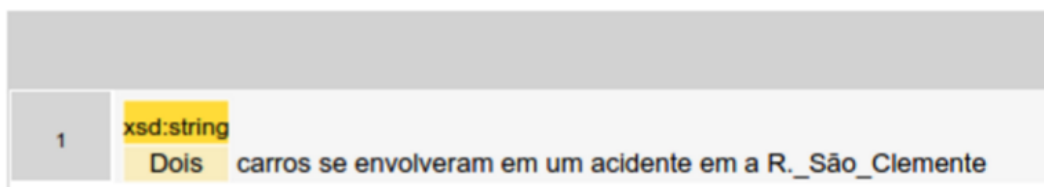


Figura A.41: xsd:string

A.2.2 RELAÇÕES

Além de classes e tipos primitivos, a ontologia de trânsito utilizada na anotação dos tweets fornece também a descrição das seguintes relações:

- causes: Indica que um determinado evento causa outro (Figura A.42);
- flowsTo: Relaciona duas localizações onde a fonte da relação descreve uma via e o alvo a direção do fluxo desta (Figura A.43). Os termos “sentido”, “direção” e “que vai para” são indicadores da presença desse tipo de relação;
- hasActor: Relaciona um evento qualquer aos atores envolvidos no mesmo (Figura A.44);
- hasEvent: Relaciona uma localização a eventos ocorridos nesta (Figura A.45);
- hasSupporter: Relaciona um evento qualquer a uma autoridade pública que está dando suporte ao mesmo (Figura A.46);
- hasTime: Relaciona um evento qualquer ao período em que este ocorreu (Figura A.47);
- isAlternativeFor: Relaciona 2 locais, onde um é uma via alternativa para o outro, geralmente esta relação ocorre quando há um evento de interdição (Figura A.48). Os termos “prefira a”, “siga pela” e “vá pela” são indicadores da presença desse tipo de relação;
- isEdgeFor: Trata a relação entre dois locais, onde a fonte de relação descreve um limite para o alvo, geralmente esta relação está presente quando se reporta um evento de tráfego em uma via e os limites desta afetados (Figura A.49). Os termos “a partir”, “até” e “entre” são indicadores da presença desse tipo de relação;
- isReferenceFor: Descreve a relação entre dois locais, onde a fonte da relação é uma referência que ajuda a definir com mais precisão a localização do alvo. Os termos “na altura”, “próximo” e “perto” são indicadores da presença desse tipo de relação (Figura A.50);
- isRestrictedTo: Descreve a relação entre dois locais, onde a fonte da relação é uma referência que define com exatidão a localização do alvo (Figura A.51);

- hasNumericQuantity: Relação que liga o tipo de ator envolvido em determinado evento e sua quantidade, quando descrita de forma numérica (Figura A.52);
- hasStringQuantity: Relação que liga o tipo de ator envolvido em determinado evento e sua quantidade, quando descrita de forma não numérica (Figura A.53);
- hasWayEffect: Esta relação liga um evento qualquer ao modo como a via foi afetada por ele (Figura A.54);

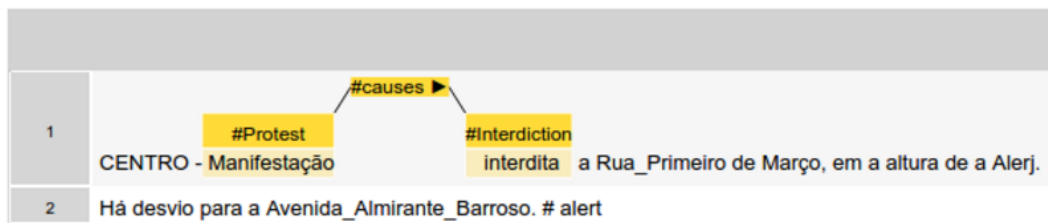


Figura A.42: causes

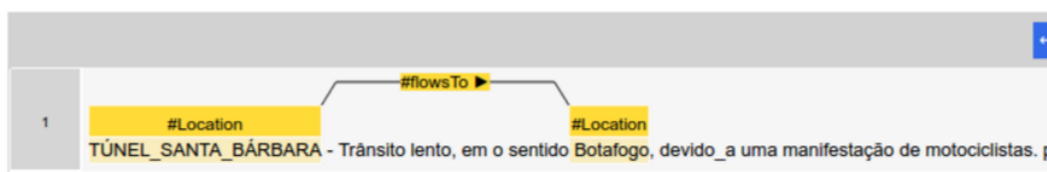


Figura A.43: flowsTo

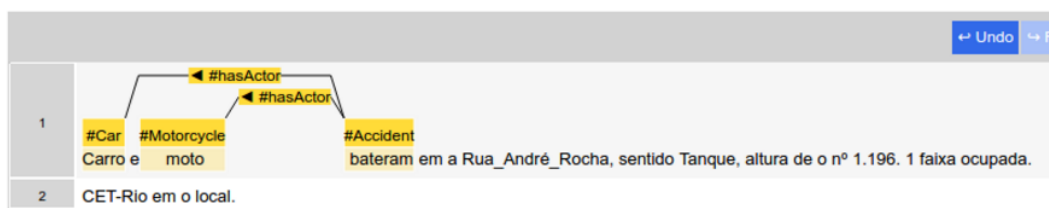


Figura A.44: hasActor

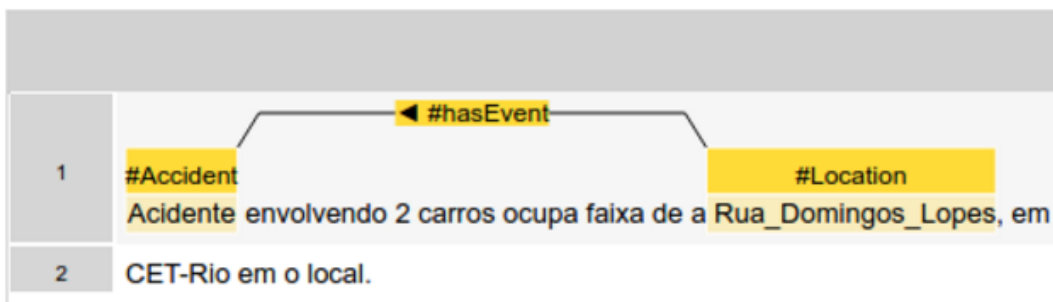


Figura A.45: hasEvent

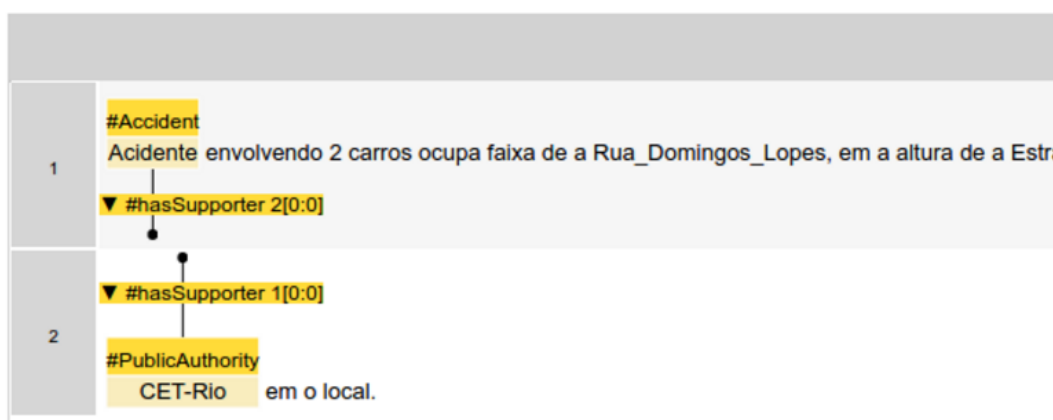


Figura A.46: hasSupporter



Figura A.47: hasTime

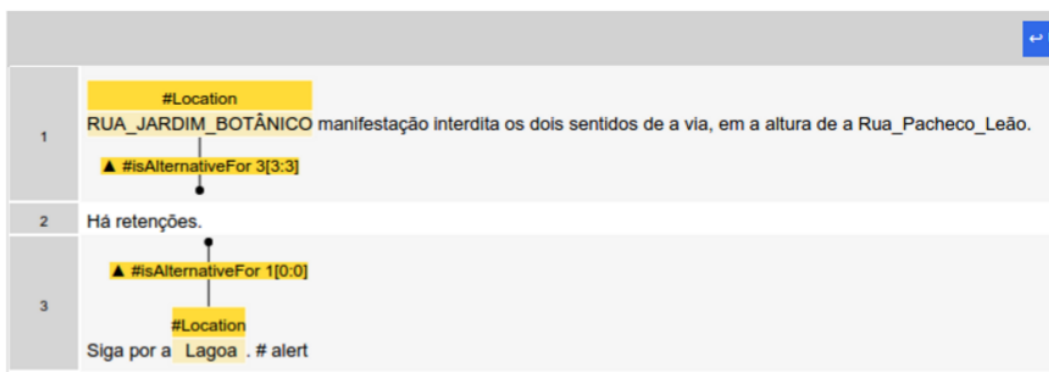


Figura A.48: isAlternativeFor

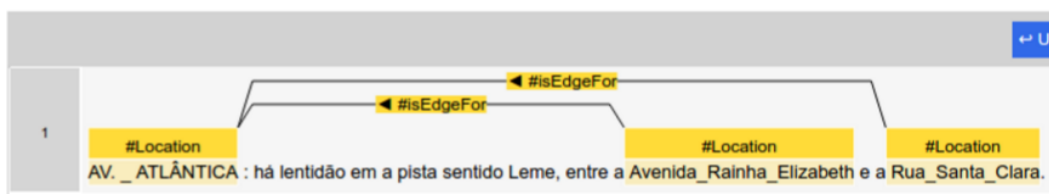


Figura A.49: isEdgeFor

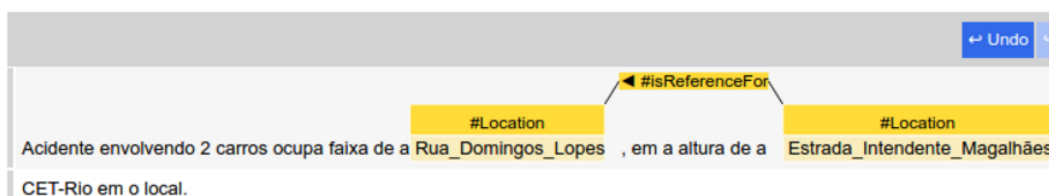


Figura A.50: isReferenceFor

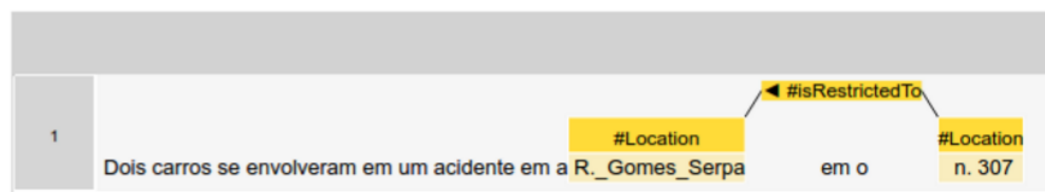


Figura A.51: isRestrictedTo

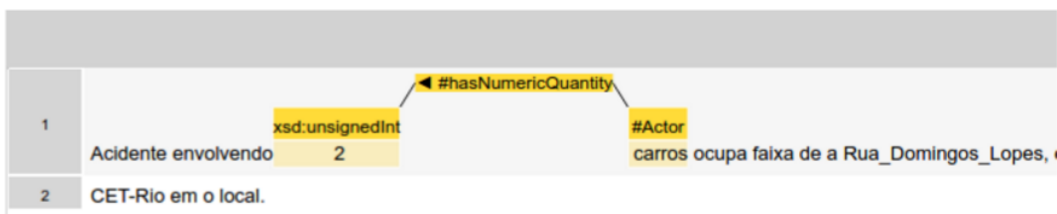


Figura A.52: hasNumericQuantity

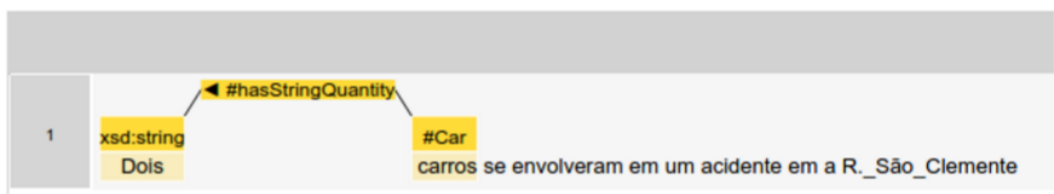


Figura A.53: hasStringQuantity

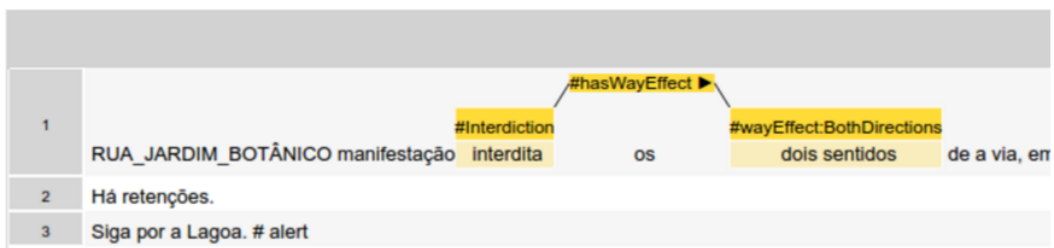


Figura A.54: hasWayEffect

PUC-Rio - Certificação Digital Nº 1421597/CA

A.2.3 CONECTORES

Como já foi dito na seção A.1.1, na ferramenta de anotação há a possibilidade de associação de blocos de texto a uma relação, neste caso, o anotador deve buscar associar partes do texto (que não estão rotulados) a relações que foram deduzidas justamente pela presença desta parte do texto, como pode ser visto na Figura A.55.

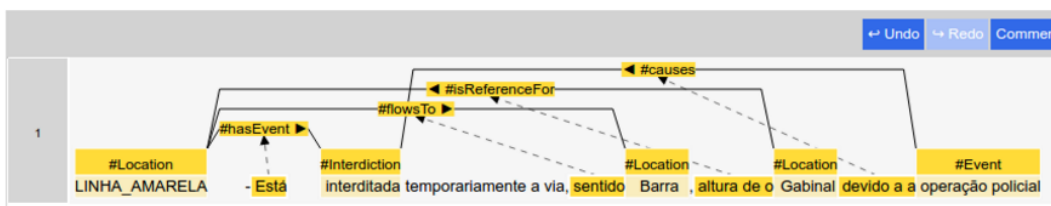


Figura A.55: conectores

A.2.4 EXEMPLOS

As Figuras A.56, A.57, A.58, A.59, A.60, A.61, A.62 e A.63 são exemplos de tweets completamente anotados.

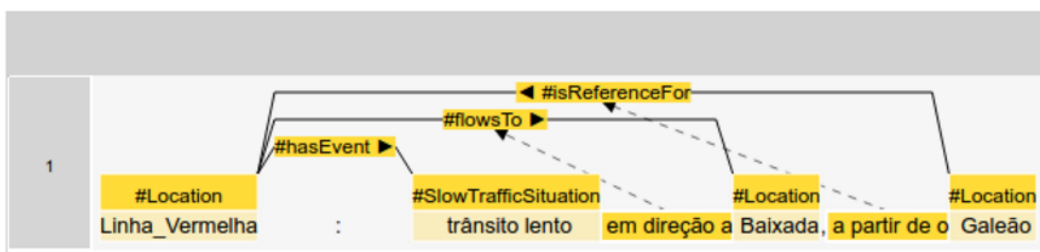


Figura A.56: exemplo 1

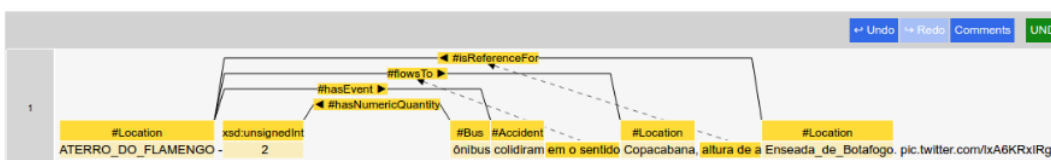


Figura A.57: exemplo 2

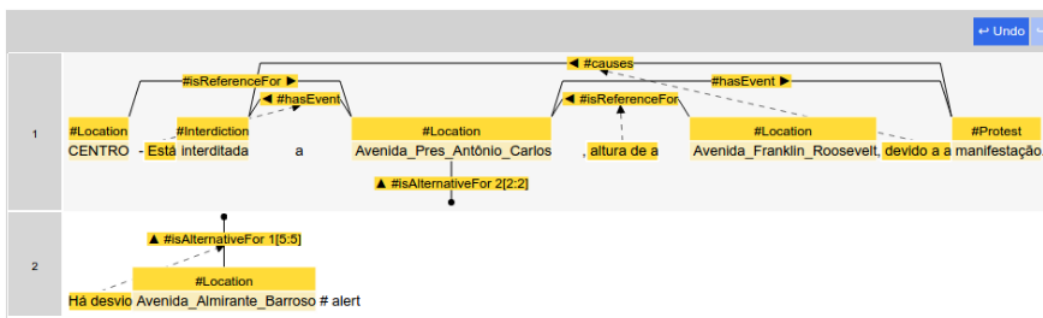


Figura A.58: exemplo 3

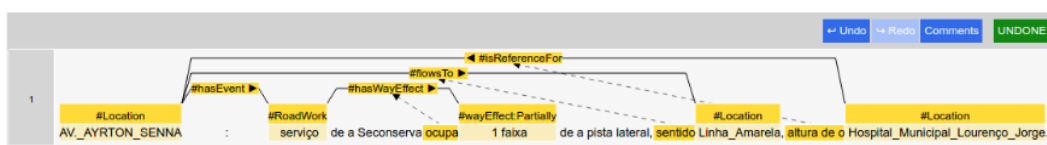


Figura A.59: exemplo 4

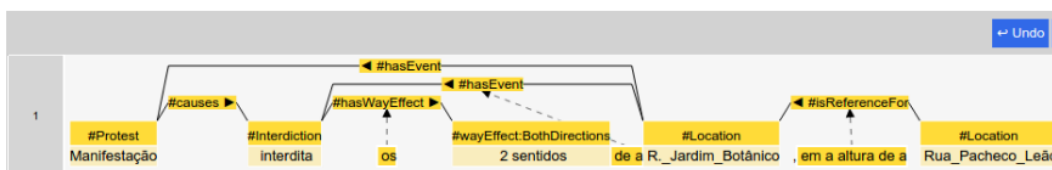


Figura A.60: exemplo 5

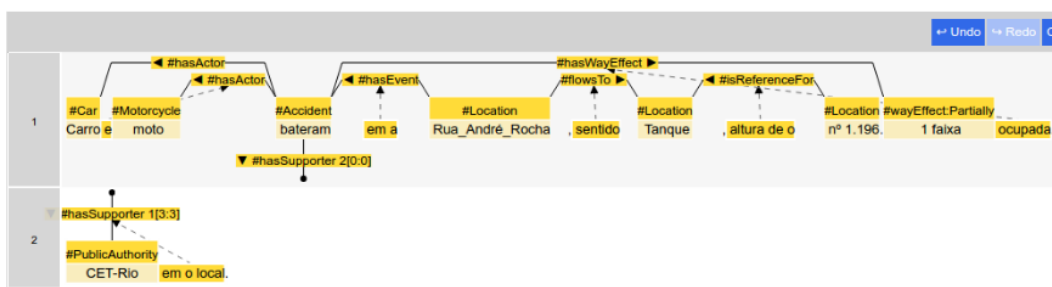


Figura A.61: exemplo 6

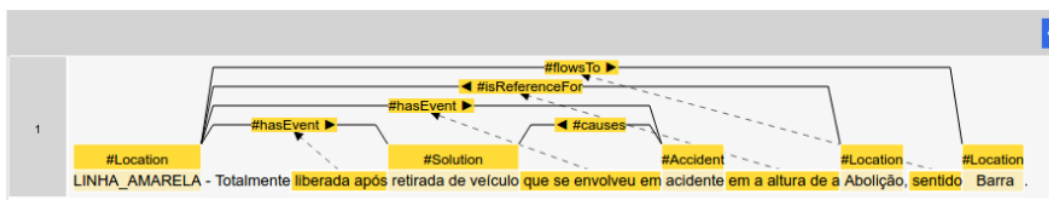


Figura A.62: exemplo 7

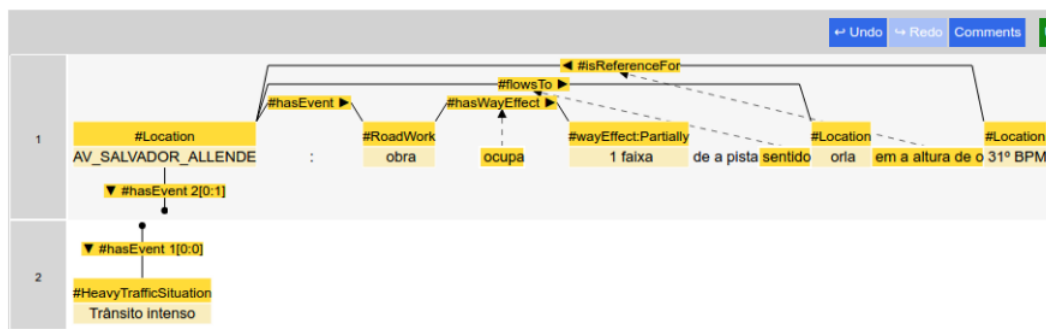


Figura A.63: exemplo 8

A.3 CONSIDERAÇÕES FINAIS

Durante o processo de anotação não hesite em consultar este guia sempre que necessário para sanar dúvidas, e caso este guia não seja o suficiente para saná-las, anote o tweet da forma que achar mais coerente e deixe um comentário (como foi explicado na seção A.1) descrevendo suas dificuldades, estes comentários são muito importantes para o enriquecimento deste guia e amadurecimento da ontologia utilizada na anotação. Um ponto importante, que deve ser levado em consideração no processo de anotação, é que por enquanto apenas tweets que dizem respeito diretamente a eventos rodoviários devem ser anotados, os demais devem ser finalizados sem realizar anotações, como por exemplo, documentos que descrevam eventos ferroviários (Figura A.64), eventos climáticos (sem relação direta com um evento rodoviário, Figura A.65), informações gerais (Figura A.66) e documentos fora de contexto (como respostas de outros tweets, Figura A.67).

1	Linha 2 de o metrô tem circulação interrompida entre Colégio e Pavuna.
2	Leia mais em @jornalodia http : // bit.ly/2c7WrXI

Figura A.64: Tweet de evento ferroviário

1	Aviso_de_Ressaca : de 10h, de o próximo_domingo, até a as 10h_de_segunda-feira (14/11).
2	Ondas de até 2,5 metros.
3	Fonte : Marinha_do_Brasil

Figura A.65: Tweet de evento climático

1	DOMINGO (06/11) TEM_IRONMAN (@ Brasillronman) NO_RIO.
2	Veja o esquema de trânsito para o evento em o link http : // bit.ly/ironman2016tra nsito ... pic.twitter.com/XPJiz3Z23V

Figura A.66: Tweet de informação geral

1	@ sramichelle Isso foi dentro de o metrô
---	--

Figura A.67: Tweet fora de contexto

B

Experimento de anotação: Comentários

Este apêndice contém todos os comentários feitos pelos participantes do experimento de anotação descrito no capítulo 4 desta dissertação. Os comentários estão separados por documentos, onde cada bloco de comentário contém o texto do *tweet* que representa o documento e os comentários feitos pelos usuários. Há de se notar que apenas alguns dos 50 documentos anotados estão descritos aqui, já que não haviam comentários em todos eles, e mesmo havendo, nem todos usuários realizaram esta ação.

Tweet: Dificuldades na Av. Presidente Antônio Carlos e também na R. Primeiro de Março.

- User, 10(A,2,5): *Fora de contexto.*

Tweet: Av. Borges de Medeiros com retenções em direção ao Rebouças, a partir da Rua Mário Ribeiro.

- User, 10(A,2,5): *Coloquei “retenções” como #Event pois não encontrei uma rotulação de descrevesse este evento. Pensei em colocar #Interdiction, porém interdição e retenção não são a mesma coisa.*

Tweet: Lentidão na Av. Epitácio Pessoa, sentido Rebouças, a partir do Corte do Cantagalo.

- User, 09(A,2,5): *Seguindo o exemplo da figura 56 do Guia, em conflito com a informação do uso do “isEdgeFor” nos termos “a partir”*

Tweet: Fluxo livre nos dois sentidos da Ponte, 13 minutos.

- User, 02(A,0,5): *Para esse caso, primeiro que apareceu para mim até o momento, seria interessante permitir um relação de duração de tempo. Semelhante seria o tempo para percorrer o Rebouças, ou ir do Leblon até o Centro (tem uma placa próximo a gávea indicando isso).*
- User, 13(A,2,10): *Em “13 minutos” fiquei na dúvida em usar #Time (relacionado a evento) ou xsd:unsignedInt (Associado a quantidade de atores) ou nenhum dos dois.*

Tweet: Via liberada!

- User, 13(A,2,10): *“Via liberada” seria #hasEvent (pelo que entendi), mas não tem #Location neste Tweet.*

- User, 18(A,4,10): *Neste não anotei nada e dei DONE. Apesar de ter evento rodoviário, não há contexto algum.*

Tweet: GASÔMETRO - Lentidão na via, sentido Centro, a partir do Into.

- User, 16(A,4,5): *comentario centro da cidade*
- User, 08(B,0,10): *SlowTrafficSituation e HeavyTrafficSituation parecem sinônimos para mim.*

Tweet: Obra interdita a Estrada da Gávea, sentido Barra, altura do nº 696. Desvio: Autoestrada Lagoa-Barra.

- User, 02(A,0,5): *Observação de que os atalhos deveriam ser desabilitados enquanto o usuário digita o comentário.*

Tweet: RT @_ecoponte : 06:30 - 15 min. Trânsito para o Rio apresenta retenção dos acessos à Mocanguê. [http:// ecoponte.com.br](http://ecoponte.com.br)

- User, 08(B,0,10): *Fiquei na dúvida na “retenção dos acessos a”. Marquei IsReferenceFor como um chute.*

Tweet: RT @_ecoponte : 05:23 - 13 min. Fluxo sentido Niterói segue normal. [http:// ecoponte.com.br](http://ecoponte.com.br)

- User, 16(A,4,5): *TESTANDO A BAGAÇA.....*

Tweet: RECREIO - Acidente envolvendo ônibus e utilitário ocupa faixa da Av. das Américas, altura do BRT Salvador Allende, sentido Barra.

- User, 08(B,0,10): *Acidente e interdição deveriam ser associados ao mesmo local?*

Tweet: Viaduto 31 de Março apresenta lentidão no sentido Centro. Mesmo panorama no Túnel Santa Bárbara.

- User, 02(A,0,5): *#flowsTo e #hasEvent estão sobrepostos no tweet inferior*

Tweet: GÁVEA - Atropelamento envolvendo moto ocupa faixa da Av. Rodrigo Otávio, altura da R. Mário Ribeiro.

- User, 08(B,0,10): *Acidente e interdição deveriam ser conectados ao mesmo local? Ou apenas interdição?*

Tweet: VILA DA PENHA - Trânsito lento na Av. Meriti, sentido Largo do Bicão, entre a Av. Vicente de Carvalho e a Rua Volta.

- User, 13(A,2,10): *Não sei se neste caso em que é descrita uma restrição de local usando ruas, Exemplo: Na Rua XPTO entre Rua A e Rua B. Entre Ruas A e B dá ideia de trecho específico, então parece que é correto usar #isRestrictedTo ao invés de #isReferenceFor, mas fiquei na dúvida. Então usei a seguinte regra, se aparece a palavra “Entre” em trechos “Rua tal e Rua Tal” estou usando #isRestrictedTo.*

Tweet: ANIL - Acidente envolvendo 2 carros ocupa faixa da Est. de Jacarepaguá, sentido Grajaú, altura da R. Sd. Genaro Pedro Lima.

- User, 08(B,0,10): *Quando uma via tem dois eventos, acidente e interdição, devemos ligá-la a ambos os eventos?*

Tweet: PARADA DE LUCAS - Acidente entre dois carros na pista central da Av. Brasil, sentido Centro. Lentidão no local. Uma faixa ocupada.

- User, 02(A,0,5): *Assumindo que 'em o local.' = no local*

Tweet: CIDADE NOVA - Acidente envolvendo moto e carro ocupa faixa da central da Av. Pres. Vargas, sentido Candelária, altura da Cidade Nova.

- User, 06(A,0,10): *tive que marca 2 palavras na mesma tag nao esta marcando so uma quando estao muito proximas.*

Tweet: RT @LinhaAmarelaRJ : 6h55 No sentido Barra trânsito lento do acesso 9, #AvBrasil , ao acesso 8, Bonsucesso, #FluxodeVeículos

- User, 17(B,4,5): *Esta sentença está bem estranha, não entendi muito bem.*

Tweet: AV BRÁS DE PINA Ônibus enguiçado ocupa faixa no sentido Olaria, altura do nº 585. CET-Rio na via. Lentidão.

- User, 08(B,0,10): *IsRestrictedTo e IsReferenceFor são sinônimos para mim.*

Tweet: AV BRASIL #IRAJÁ Caminhão enguiçado ocupa faixa da pista central sentido Zona Oeste. CET-Rio na via. Sem retenção.

- User, 08(B,0,10): *Breakdown e Interdiction deveriam ser associados a Location? Ou apenas Interdiction?*

C

Experimento de anotação: Dados dos participantes

Nas próximas seções, estão as informações gerais de cada participante do experimento, seguidas das respostas que deram sobre o questionário dado ao final do experimento de anotação descrito no capítulo 4. Após os questionários há uma tabela com o resumo das opiniões extraídas das respostas dos participantes.

C.1

User 01(A,0,5)

Recebeu treinamento presencial: Não

Aquecimento: 0

Passo de re-anotação: 5

Idade: 24

Escolaridade: Superior completo

Profissão: Analista de sistemas

Já participou de algum outro experimento de anotação de texto? Caso tenha participado, qual o grau de dificuldade da tarefa que acabou de executar em comparação às anteriores?

Não

Quais foram as principais dificuldades que encontrou durante a tarefa?

Nenhuma

Você tentou executar a tarefa lendo apenas parte do guia de anotação? Caso sim, acha que terminou de ler todo ele a partir de que documento finalizado?

Não

O que achou do guia de anotação? O que mudaria nele?

Nada

O que achou da ferramenta de anotação? O que mudaria nela?

Achei uma ótima ideia, diferente mas boa.

C.2

User 02(A,0,5)

Recebeu treinamento presencial: Não

Aquecimento: 0

Passo de re-anotação: 5

Idade: 29

Escolaridade: Pós-graduação (mestrado)

Profissão: Analista de sistemas

Já participou de algum outro experimento de anotação de texto? Caso tenha participado, qual o grau de dificuldade da tarefa que acabou de executar em comparação às anteriores?

Não, primeira vez que participo desse tipo de experimento.

Quais foram as principais dificuldades que encontrou durante a tarefa?

Quando tinha dúvidas procurava um exemplo da relação do que eu estava vendo, foram um ou dois casos que não encontrei e não consigo recordá-los agora. A maior dificuldade foi utilizar o guia no formato PDF através da página. Outra dificuldade foi adicionar comentários, que não desativava as teclas de atalho enquanto o texto era digitado.

Você tentou executar a tarefa lendo apenas parte do guia de anotação? Caso sim, acha que terminou de ler todo ele a partir de que documento finalizado?

Primeiro li todo documento, quando acabei que comecei a realizar as primeiras anotações. Para saber se estava realizando uma anotação corretamente, inicialmente procurei um exemplo similar e fui comparando até terminar.

O que achou do guia de anotação? O que mudaria nele?

O guia foi fundamental para a realização da tarefa, sendo claro e sucinto. Porém, dado o contexto da ferramenta, o guia podia ser uma página web interativa. Uma mudança imediata seria reabrir o guia na página em que eu estava lendo e não na primeira.

O que achou da ferramenta de anotação? O que mudaria nela?

De forma geral o uso da ferramenta foi muito fácil, o que acabou motivando a realização da tarefa de anotação. Sugiro as seguintes melhorias:

1 - Permitir cancelar a ação clicando em qualquer lugar da página e não somente na área de anotação. Por exemplo, quando escolho uma palavra e digito C para conecta-lá, para cancelar era necessário clicar na área dos tweets.

2 - Corrigir a sobreposição de rótulos entre relações de diferentes tweets.

3 - Melhorar o resize da tela. No meu notebook que tem uma pequena tela e costume dar zoom no Windows 10 (aumento do tamanho texto de aplicativos) não ficou legal.

C.3

User 03(B,0,5)

Recebeu treinamento presencial: Sim

Aquecimento: 0

Passo de re-anotação: 5

Idade: 29

Escolaridade: Doutorado em andamento

Profissão: Estudante

Já participou de algum outro experimento de anotação de texto? Caso tenha participado, qual o grau de dificuldade da tarefa que acabou de executar em comparação às anteriores?

Não, nunca havia participado.

Quais foram as principais dificuldades que encontrou durante a tarefa?

Talvez, de início, fazer a associação correta do termo ao subtipo específico. Depois de algumas anotações, isso passou a ser mais fácil.

Você tentou executar a tarefa lendo apenas parte do guia de anotação? Caso sim, acha que terminou de ler todo ele a partir de que documento finalizado?

Eu li apenas uma parte do guia e assiste à explicação realizada pelo autor.

O que achou do guia de anotação? O que mudaria nele?

Pelo que li, o guia estava bem feito.

O que achou da ferramenta de anotação? O que mudaria nela?

A ferramenta é excelente e, por isso, gostaria de dar os parabéns ao autor. Talvez tentar diminuir o tempo de espera entre a anotação de cada termo, mas não sei se é possível.

C.4

User 04(B,0,5)

Recebeu treinamento presencial: Sim

Aquecimento: 0

Passo de re-anotação: 5

Idade: 30

Escolaridade: Mestrado em Informática

Profissão: Estudante

Já participou de algum outro experimento de anotação de texto? Caso tenha participado, qual o grau de dificuldade da tarefa que acabou de executar em comparação às anteriores?

Não, foi a primeira vez.

Quais foram as principais dificuldades que encontrou durante a tarefa?

Os tweets não sempre tem uma unica interpretação, o q dificulta ser coerente na anotação.

Você tentou executar a tarefa lendo apenas parte do guia de anotação? Caso sim, acha que terminou de ler todo ele a partir de que documento finalizado?

Eu não li o guia de anotação. Só segui as suas instruções com a aula practica. A verdade, acho o guia muito comprido. Pode ser q a quantidade de páginas desanime ao pessoal. Acho melhor o q fizemos, aprender na hora com exemplos.

O que achou do guia de anotação? O que mudaria nele?

Depois de anotar voltei em ele e está muito explicativo, so q comprido.

O que achou da ferramenta de anotação? O que mudaria nela?

Gostei da ferramenta. O fato de ser simple e facil a anotação de atores, eventos, localizações, etc, ajuda muito. Gostei das setinhas em vermelho, pois mais de uma vez tentei fazer uma coisa que não estava certa e a ferramenta ajudou a não errar. Acho que o fato de ter shortcut para o ultimo label colocado, subliminarmente ajuda a ser organizado na anotação. Gostei do desenho da ferramenta pois é simple, limpo (não está sobrecarregado de opções) e amigavel. Agora, quando voce finaliza a anotação poderia sair algum resultado associado a seu desempenho (Ex: percentual de anotados corretamente, etc) Pois como está deixa uma sensação de inserteção.

C.5

User 05(A,0,10)

Recebeu treinamento presencial: Não

Aquecimento: 0

Passo de re-anotação: 10

Idade: 39

Escolaridade: Doutor em informática

Profissão: Analista

Já participou de algum outro experimento de anotação de texto? Caso tenha participado, qual o grau de dificuldade da tarefa que acabou de executar em comparação às anteriores?

Não

Quais foram as principais dificuldades que encontrou durante a tarefa?

Somente o fato de ser algo totalmente novo.

Você tentou executar a tarefa lendo apenas parte do guia de anotação? Caso sim, acha que terminou de ler todo ele a partir de que documento finalizado?

Sim. Não li todo.

O que achou do guia de anotação? O que mudaria nele?

O guia tem bastante informação. Usuários de sistema tem preguiça de ler manual. Se fosse eu colocaria os exemplos logo no início e as informações adicionais depois.

O que achou da ferramenta de anotação? O que mudaria nela?

A ferramenta é muito boa. Não mudaria nada.

C.6

User 06(A,0,10)

Recebeu treinamento presencial: Não

Aquecimento: 0

Passo de re-anotação: 10

Idade: 34

Escolaridade: Superior Completo

Profissão: Designer Gráfico

Já participou de algum outro experimento de anotação de texto? Caso tenha participado, qual o grau de dificuldade da tarefa que acabou de executar em comparação às anteriores?

Não

Quais foram as principais dificuldades que encontrou durante a tarefa?

Demorei um pouco para entender no começo, tive que voltar nos exemplos varias vezes ate pegar agilidade

Você tentou executar a tarefa lendo apenas parte do guia de anotação? Caso sim, acha que terminou de ler todo ele a partir de que documento finalizado?

Não

O que achou do guia de anotação? O que mudaria nele?

Poderia ser mais dinâmico, talvez um infográfico...

O que achou da ferramenta de anotação? O que mudaria nela?

Interessante, tentaria deixar ela mais didática, intuitiva...

C.7

User 07(B,0,10)

Recebeu treinamento presencial: Sim

Aquecimento: 0

Passo de re-anotação: 10

Idade: 49

Escolaridade: Doutorado em andamento

Profissão: Estatístico

Já participou de algum outro experimento de anotação de texto? Caso tenha participado, qual o grau de dificuldade da tarefa que acabou de executar em comparação às anteriores?

Nunca participei.

Quais foram as principais dificuldades que encontrou durante a tarefa?

Saber as rotulações mais adequadas e em alguns casos as relações.

Você tentou executar a tarefa lendo apenas parte do guia de anotação? Caso sim, acha que terminou de ler todo ele a partir de que documento finalizado?

Li quase todo o documento, mas durante as anotações consultei varias vezes.

O que achou do guia de anotação? O que mudaria nele?

O guia esta bom e foi fundamental. Poderia ter mais exemplos completos.

Um cheat sheet seria útil

O que achou da ferramenta de anotação? O que mudaria nela?

Ferramenta prática e fácil de operar. Não mudaria nada.

C.8**User 08(B,0,10)**

Recebeu treinamento presencial: Sim

Aquecimento: 0

Passo de re-anotação: 10

Idade: 30

Escolaridade: Doutorado em andamento

Profissão: Analista de sistemas

Já participou de algum outro experimento de anotação de texto? Caso tenha participado, qual o grau de dificuldade da tarefa que acabou de executar em comparação às anteriores?

Sim, participei de alguns experimentos de anotação no passado. Achei este experimento mais fácil de executar por conta da ferramenta e dos exemplos práticos dados pela pessoa que realizou o experimento antes de eu iniciar a anotação.

Quais foram as principais dificuldades que encontrou durante a tarefa?

Achei algumas tags ambíguas, tais como `IsReferenceFor/IsRestrictedTo` e `HeavyTraffic/SlowTraffic`.

Você tentou executar a tarefa lendo apenas parte do guia de anotação? Caso sim, acha que terminou de ler todo ele a partir de que documento finalizado?

Praticamente não li o guia de anotação. Ajudou muito a pessoa que realizou o experimento fazer alguns exemplos práticos e depois responder algumas dúvidas ao longo do processo de anotação. A anotação foi muito fluida.

O que achou do guia de anotação? O que mudaria nele?

Pelo fato de não ter quase olhado o guia, não tenho propriedade para avaliá-lo. Mas exemplos práticos ajudariam muito.

O que achou da ferramenta de anotação? O que mudaria nela?

A ferramenta de anotação está ótima. Não mudaria nada nela.

C.9**User 09(A,2,5)**

Recebeu treinamento presencial: Não

Aquecimento: 2

Passo de re-anotação: 5

Idade: 31

Escolaridade: Ensino Superior em andamento

Profissão: Militar

Já participou de algum outro experimento de anotação de texto? Caso tenha participado, qual o grau de dificuldade da tarefa que acabou de executar em comparação às anteriores?

Não.

Quais foram as principais dificuldades que encontrou durante a tarefa?

No inicio, associar todas as sentenças, mas logo a logica ajudou no raciocínio

Você tentou executar a tarefa lendo apenas parte do guia de anotação? Caso sim, acha que terminou de ler todo ele a partir de que documento finalizado?

Não, li todo Guia antes.

O que achou do guia de anotação? O que mudaria nele?

Bem detalhado, apenas colocaria mais de um exemplo por função para comparar diversas situações de emprego das sentenças

O que achou da ferramenta de anotação? O que mudaria nela?

Muito útil, caso possível um auto vinculo primário das sentenças possíveis para que ficasse apenas os ajustes finos por parte do usuário.

C.10

User 10(A,2,5)

Recebeu treinamento presencial: Não

Aquecimento: 2

Passo de re-anotação: 5

Idade: 28

Escolaridade: Mestre em Sistemas e Computação

Profissão: Programadora

Já participou de algum outro experimento de anotação de texto? Caso tenha participado, qual o grau de dificuldade da tarefa que acabou de executar em comparação às anteriores?

Nunca participei.

Quais foram as principais dificuldades que encontrou durante a tarefa?

Saber quando uma palavra era o nome de alguma rua, bairro ou outra coisa. Como não sou do Rio e não o frequento muito, eu não reconhecia alguns nomes de ruas e pontos de referência. Algumas vezes tive que pesquisar no Google para saber se era nome de Rua, ou apenas uma referência ou outra coisa. Como por exemplo a palavra Gasômetro.

Outra dificuldade foi saber o que anotar para determinada palavra. Por exemplo: “Trânsito sentido Gávea”, Eu não sabia se anotava como “one direction” ou como “flow to”. Acabei anotando como flow to em todos os sentidos, pois me baseei nos exemplos. Outro exemplo: A palavra “retenção”, para mim retenção não é a mesma coisa que interdição, então todas as vezes coloquei como “event”

Você tentou executar a tarefa lendo apenas parte do guia de anotação? Caso sim, acha que terminou de ler todo ele a partir de que documento finalizado?

Não. Eu li ele completo antes de começar a fazer. Porém, consultei ele diversas vezes durante todo o processo de anotação.

O que achou do guia de anotação? O que mudaria nele?

Achei o guia muito bom e muito bem explicado. A única coisa que mudaria nele é que colocaria uma sessão explicando alguns termos utilizados, como por exemplo: domínio, entidades, relações. No início estes termos já são apresentados para o leitor sem uma explicação prévia, e me senti meio perdida algumas vezes. Fui captar o que cada termo significava apenas quando analisei os exemplos.

O que achou da ferramenta de anotação? O que mudaria nela?

Muito boa, intuitiva, e de fácil utilização. Eu colocaria uma opção de revisão e edição das anotações finalizadas. Algumas vezes tive a impressão de que havia feito alguma anotação errada, porém não tinha como voltar para conferir. E no início não tenho certeza se fiz correto, pois fui pegar o jeito das anotações apenas depois de algumas já finalizadas. Então seria legal poder conferir e se houver alguma errada, poder editar.

C.11

User 11(B,2,5)

Recebeu treinamento presencial: Sim

Aquecimento: 2

Passo de re-anotação: 5

Idade: 26

Escolaridade: Superior completo

Profissão: Estudante

Já participou de algum outro experimento de anotação de texto? Caso tenha participado, qual o grau de dificuldade da tarefa que acabou de executar em comparação às anteriores?

Não, nunca tinha feito.

Quais foram as principais dificuldades que encontrou durante a tarefa?

A ferramenta mostra um processo de carga por cada passo feito, o por cada anotação. Poderia guardar so nu final de cada oração.

Você tentou executar a tarefa lendo apenas parte do guia de anotação? Caso sim, acha que terminou de ler todo ele a partir de que documento finalizado?

Não, a gente teve um tutorial personalizado, so usei para consultar os relacionamentos.

O que achou do guia de anotação? O que mudaria nele?

Ta bom. mais para melhorar poderia gravar suas proprias explicações, assim como as perguntas mais frequentes do pessoal que so tem o manual.

O que achou da ferramenta de anotação? O que mudaria nela?

A ferramenta e muito simples e cobre as funcionalidades esperadas dela, so foi muito chato que a cada anotação saia o simbolo de carregando e frizaba ate acabar a carga.

C.12**User 12(A,2,10)**

Recebeu treinamento presencial: Não

Aquecimento: 2

Passo de re-anotação: 10

Idade: 30

Escolaridade: Ensino superior incompleto

Profissão: Técnico em informática

Já participou de algum outro experimento de anotação de texto? Caso tenha participado, qual o grau de dificuldade da tarefa que acabou de executar em comparação às anteriores?

Não, foi a primeira vez.

Quais foram as principais dificuldades que encontrou durante a tarefa?

Foi fazer algumas conexões e classificar algumas palavras.

Você tentou executar a tarefa lendo apenas parte do guia de anotação? Caso sim, acha que terminou de ler todo ele a partir de que documento finalizado?

Não.

O que achou do guia de anotação? O que mudaria nele?

Bom, só achei que tava faltando algumas coisas do guia na ferramenta.

O que achou da ferramenta de anotação? O que mudaria nela?

Achei a ferramenta muito boa, ela é bem simples e fácil de usar.

C.13**User 13(A,2,10)**

Recebeu treinamento presencial: Não

Aquecimento: 2

Passo de re-anotação: 10

Idade: 32

Escolaridade: Doutorado em andamento

Profissão: Professor

Já participou de algum outro experimento de anotação de texto? Caso tenha participado, qual o grau de dificuldade da tarefa que acabou de executar em comparação às anteriores?

Nunca participei.

Quais foram as principais dificuldades que encontrou durante a tarefa?

Nas primeiras anotações que fiz foi complicado relacionar as anotações além do básico, que é veículo, evento manutenção da via, etc. Conforme fui consultando os exemplos novamente do Guia melhorou o meu entendimento e pude fazer anotações mais completas.

Você tentou executar a tarefa lendo apenas parte do guia de anotação? Caso sim, acha que terminou de ler todo ele a partir de que documento finalizado?

Eu li o guia todo uma vez. Depois quando tinha dúvida olhava os exemplos do guia. Foi muito útil os exemplos com todos os tipos de anotações. Ficou mais fácil de entender olhando estes exemplos completos.

O que achou do guia de anotação? O que mudaria nele?

Colocaria mais exemplos, principalmente aqueles que apresentam os vários tipos de anotações. Fica mais fácil e rápido de entender olhando os exemplos.

O que achou da ferramenta de anotação? O que mudaria nela?

Mudaria o “Menu suspenso” de modo que houvesse vários “sub menus” de acordo com os tipos de anotações mais gerais e mais específicas. Exemplo, “Carro”, “Moto”, “Ônibus”, seria um sub menu dentro de “Veículo”.

C.14

User 14(B,2,10)

Recebeu treinamento presencial: Sim

Aquecimento: 2

Passo de re-anotação: 10

Idade: 29

Escolaridade: Doutorado em andamento

Profissão: Estudante

Já participou de algum outro experimento de anotação de texto? Caso tenha participado, qual o grau de dificuldade da tarefa que acabou de executar em comparação às anteriores?

Não, esse foi meu primeiro experimento de anotação de texto.

Quais foram as principais dificuldades que encontrou durante a tarefa?

O conhecimento de alguns nomes de bairros (eu sou estrangeiro, ainda estou aprendendo Português)

Você tentou executar a tarefa lendo apenas parte do guia de anotação? Caso sim, acha que terminou de ler todo ele a partir de que documento finalizado?

Eu recebi uma aula de como usar a ferramenta.

O que achou do guia de anotação? O que mudaria nele?

Muito bem estruturado, bem feito. Não mudaria nada.

O que achou da ferramenta de anotação? O que mudaria nela?

É uma ferramenta bem feita. Gostei do fluxo para fazer as anotações, eu incluiria uma legenda com os atalhos (T, R e C), e tentaria fazê-la mais independente do servidor quando uma ação é feita... tipo enviar só quando a anotação estiver concluída, e ir guardando as mudanças num arquivo temporário se for preciso.

C.15**User 15(A,4,5)**

Recebeu treinamento presencial: Não

Aquecimento: 4

Passo de re-anotação: 5

Idade: 45

Escolaridade: Mestrado em Informática

Profissão: SysAdmin e Desenvolvedor

Já participou de algum outro experimento de anotação de texto? Caso tenha participado, qual o grau de dificuldade da tarefa que acabou de executar em comparação às anteriores?

Não.

Quais foram as principais dificuldades que encontrou durante a tarefa?

A principal dificuldade foi ter que ler um documento de 22 páginas. Onde, só no final do documento é que realmente fica claro como fazer as anotações.

Você tentou executar a tarefa lendo apenas parte do guia de anotação? Caso sim, acha que terminou de ler todo ele a partir de que documento finalizado?

Não. Consultei o documento a todo momento!

O que achou do guia de anotação? O que mudaria nele?

O guia de anotação é bom e bem informativo. Contudo, fica cansativo ter ler um monte de definições, para só depois, começar a parte de anotações. Na minha opinião, ficaria um pouco menos cansativo, se no início do guia as informações fossem menos detalhadas e em seguida, ter exemplos mais simples. E todos os detalhes ficariam para o final. Assim, o usuário, poderia já entenderia o que fazer e como fazer, já no início do documento.

O que achou da ferramenta de anotação? O que mudaria nela?

A ferramenta é muito boa. A única coisa que eu mudaria seria tirar a informação do número de documentos que faltam pra anotar. Só a porcentagem estaria boa. Esse número me deixou apreensivo, parecia que não ia terminar nunca :-)

C.16**User 16(A,4,5)**

Recebeu treinamento presencial: Não

Aquecimento: 4

Passo de re-anotação: 5

Idade: 31

Escolaridade: Ensino superior incompleto

Profissão: Analista de Sistemas

Já participou de algum outro experimento de anotação de texto? Caso tenha participado, qual o grau de dificuldade da tarefa que acabou de executar em comparação às anteriores?

Não.

Quais foram as principais dificuldades que encontrou durante a tarefa?

Algumas tags do experimento não tinham uma relação prevista na ontologia com outras, então “adaptei” com uma relação entre outras tags. Senti falta de mais informações na tela principal. Me enrolei um pouco com os atalhos.

Você tentou executar a tarefa lendo apenas parte do guia de anotação? Caso sim, acha que terminou de ler todo ele a partir de que documento finalizado?

Como a aparência da aplicação me pareceu bastante simples, tentei fazer de maneira intuitiva. Tive algumas dificuldades e, então, recorri a documentação. Depois de ler apenas uma vez, consegui finalizar tudo normalmente.

O que achou do guia de anotação? O que mudaria nele?

Nada. Estava ótimo! Apenas retiraria algumas informações do guia e colocaria na tela principal para deixar a aplicação ainda mais gerenciável pelo usuário.

O que achou da ferramenta de anotação? O que mudaria nela?

Apenas retiraria algumas informações do guia e colocaria na tela principal para deixar a aplicação ainda mais gerenciável pelo usuário.

C.17

User 17(B,4,5)

Recebeu treinamento presencial: Sim

Aquecimento: 4

Passo de re-anotação: 5

Idade: 20

Escolaridade: Superior em andamento

Profissão: Estudante

Já participou de algum outro experimento de anotação de texto? Caso tenha participado, qual o grau de dificuldade da tarefa que acabou de executar em comparação às anteriores?

Nunca participei.

Quais foram as principais dificuldades que encontrou durante a tarefa?

A subjetividade da tarefa foi algo que me trouxe bastante dificuldade, em alguns momentos ficava em dúvida de qual seria a melhor forma de realizar a anotação. Além disso, como eram muitas anotações, tive dificuldade em manter a atenção para não deixar passar nada.

Você tentou executar a tarefa lendo apenas parte do guia de anotação? Caso sim, acha que terminou de ler todo ele a partir de que documento finalizado?

Executei a tarefa lendo pouco do guia de anotação. Baseei-me, majoritariamente, na sua explicação e nos exemplos que fez junto com o pessoal.

O que achou do guia de anotação? O que mudaria nele?

Não li o suficiente do guia de anotação para opinar.

O que achou da ferramenta de anotação? O que mudaria nela?

Achei a ferramenta muito boa. A forma de classificar as palavras e de construir relações entre elas é bastante intuitiva e após realizar isso o desenho formado é uma ótima forma de verificar se o que fez está certo. Realmente não sei o que mudaria, para mim está ótimo do jeito que está.

C.18

User 18(A,4,10)

Recebeu treinamento presencial: Não

Aquecimento: 4

Passo de re-anotação: 10

Idade: 30

Escolaridade: Superior completo

Profissão: Estudante

Já participou de algum outro experimento de anotação de texto? Caso tenha participado, qual o grau de dificuldade da tarefa que acabou de executar em comparação às anteriores?

Não

Quais foram as principais dificuldades que encontrou durante a tarefa?

Na hora de finalizar a anotação de um tweet. Eu ficava em dúvida se eu realmente tinha anotado tudo correto e, por existir diversos tipos de rotulações/relacionamentos, achava que sempre podia estar faltando algo a ser rotulado/relacionado/conectado.

Você tentou executar a tarefa lendo apenas parte do guia de anotação? Caso sim, acha que terminou de ler todo ele a partir de que documento finalizado?

Tentei, na verdade fiz uma leitura rápida em algumas partes do guia de anotação e com certeza não li o documento todo no primeiro estudo do guia. Acho que sim.

O que achou do guia de anotação? O que mudaria nele?

Achei completo e bom para referência em todo momento que se tinha dúvida com relação à alguma anotação. Eu colocaria a parte dos exemplos mais pro início, pois o documento é extenso e no meu caso nem cheguei a usar todas os tipos de relação. Sempre que eu tinha alguma dúvida na anotação, procurava primeiro algum exemplo parecido completamente anotado. Acho que todas as figuras apresentadas podiam estar completamente anotadas, pois pode gerar confusão quando o anotador olha para uma figura que tem um tweet parcialmente anotado e acha que o mesmo é um exemplo completamente anotado. Para começar o guia e explicar relações/conexões/rotulações, eu colocaria tweets simples nas figuras, mas sempre completamente anotados.

O que achou da ferramenta de anotação? O que mudaria nela?

Gostei da ferramenta, achei a interface boa e separa bem o texto a medida que novas informações são inseridas, deixando claro o que está sendo feito e

visualmente legível, mesmo quando se tem muitas informações. Eu colocaria todas as opções de conexões/relações visíveis, sem a necessidade de clicar no botão direito do mouse para escolher, desta forma acredito que iria ajudar o anotador, pois com tudo visível ia ser mais fácil de decorar todas opções fazendo com que seja mais difícil o anotador esquecer de algo na anotação.

C.19

User 19(A,4,10)

Recebeu treinamento presencial: Não

Aquecimento: 4

Passo de re-anotação: 10

Idade: 28

Escolaridade: Superior Completo

Profissão: Analista de Sistemas

Já participou de algum outro experimento de anotação de texto? Caso tenha participado, qual o grau de dificuldade da tarefa que acabou de executar em comparação às anteriores?

Não, esta foi a primeira vez na qual participei de um experimento deste tipo.

Quais foram as principais dificuldades que encontrou durante a tarefa?

Encontrei algumas dificuldades nos casos que imagino serem os difíceis, sem atores bem claros, ou que continham muitas informações sobre sentidos e ruas. Não soube fazer alguns casos, onde não havia nada para ser anotado, até onde pude perceber.

Você tentou executar a tarefa lendo apenas parte do guia de anotação? Caso sim, acha que terminou de ler todo ele a partir de que documento finalizado?

Sim. Tentei fazer ao ler apenas as primeiras seções do documento e alguns exemplos básicos de anotações, mas não consegui ir muito longe com alguns casos difíceis. Por isso, reli a documentação, desta vez por completo e assim consegui fazer os demais.

O que achou do guia de anotação? O que mudaria nele?

O guia lhe dá um bom direcionamento, mas eu tentaria colocar talvez mais exemplos práticos. Talvez isto não seja possível, pois poderia influenciar no resultado do experimento, mas é algo a se pensar.

O que achou da ferramenta de anotação? O que mudaria nela?

Achei a ferramenta bem desenvolvida e bem dinâmica, mas talvez peque um pouco em usabilidade para aumentar a velocidade das anotações. Creio que deveria ser estudada alguma maneira além dos atalhos para se anotar mais rapidamente. Talvez poder, por exemplo, marcar mais de uma palavra e anotá-las todas juntas, nos casos em que todas forem de tipos iguais.

C.20**User 20(B,4,10)**

Recebeu treinamento presencial: Sim

Aquecimento: 4

Passo de re-anotação: 10

Idade: 19

Escolaridade: Ensino Superior Incompleto

Profissão: Estudante

Já participou de algum outro experimento de anotação de texto? Caso tenha participado, qual o grau de dificuldade da tarefa que acabou de executar em comparação às anteriores?

Nunca tinha participado. Achei razoavelmente tranquilo, a partir das primeiras 5/10 anotações fica bem mais fácil.

Quais foram as principais dificuldades que encontrou durante a tarefa?

Tive um problema com comentário, não podia escrever z, pois ele fazia um undo. Também me confundi inicialmente com as setas de relação entre locais. O sentido delas (ex: flowsTo, referenceFor) acabou me confundindo.

Você tentou executar a tarefa lendo apenas parte do guia de anotação? Caso sim, acha que terminou de ler todo ele a partir de que documento finalizado?

Sim, lendo apenas a parte que você mostrou ao explicar a tarefa. Nos geral as dúvidas que eu tinha podiam ser tiradas no guia, mas você acabou respondendo pessoalmente.

O que achou do guia de anotação? O que mudaria nele?

Não li inteiro para falar a verdade.

O que achou da ferramenta de anotação? O que mudaria nela?

Achei muito boa. Acho bem importante ela ser super rápida, pois qualquer tempo de espera ao montar as relações pode ser tornar um pouco tedioso. Principalmente se forem feitas um grande número de anotações. O visual está bem legal e o portal está bem intuitivo. Parabéns!

Participantes	Questões					
	Já participou de outro experimento de anotação	Expressou alguma dificuldade ao utilizar o LER	Expressou alguma dificuldade na tarefa de identificar entidades, relações ou conectores	Tentou executar a tarefa sem ler todo o guia de anotação	Propôs melhorias no guia	Propôs melhorias na ferramenta
User 01(A,0,5)	-	-	-	-	-	-
User 02(A,0,5)	-	✓	-	-	✓	✓
User 03(B,0,5)	-	-	✓	✓	-	✓
User 04(B,0,5)	-	-	✓	✓	✓	✓
User 05(A,0,10)	-	-	-	✓	✓	-
User 06(A,0,10)	-	-	✓	-	✓	✓
User 07(B,0,10)	-	-	✓	✓	✓	-
User 08(B,0,10)	✓	-	✓	✓	-	-
User 09(A,2,5)	-	-	✓	-	✓	✓
User 10(A,2,5)	-	-	✓	-	✓	✓
User 11(B,2,5)	-	✓	-	✓	✓	✓
User 12(A,2,10)	-	-	✓	-	-	-
User 13(A,2,10)	-	-	✓	-	✓	✓
User 14(B,2,10)	-	-	✓	✓	-	✓
User 15(A,4,5)	-	-	-	-	✓	✓
User 16(A,4,5)	-	✓	✓	✓	✓	✓
User 17(B,4,5)	-	-	✓	✓	-	-
User 18(A,4,10)	-	-	✓	✓	✓	✓
User 19(A,4,10)	-	-	✓	✓	✓	✓
User 20(B,4,10)	-	✓	✓	✓	-	✓

Tabela C.1: Resumo das respostas dos participantes ao questionário do experimento de anotação

D

Experimento de anotação: Tabelas e gráficos

Este apêndice contém gráficos e tabelas com informações adicionais sobre o experimento de anotação descrito no capítulo 4 desta dissertação.

Tabela D.1: Resumo colaborações

Colaborador	Tempo gasto	Tempo médio	# visualizações	# visualizações média
User, 01(A,0,5)	2h 18m 4s	2m 45s (σ 4m 3s)	63	1,26 (σ 0,52)
User, 02(A,0,5)	2h 51m 45s	3m 26s (σ 5m 40s)	126	2,52 (σ 1,86)
User, 03(B,0,5)	3h 12m 57s	3m 51s (σ 6m 30s)	55	1,10 (σ 0,30)
User, 04(B,0,5)	2h 42m 31s	3m 15s (σ 4m 16s)	51	1,02 (σ 0,14)
User, 05(A,0,10)	2h 9m 30s	2m 35s (σ 5m 20s)	61	1,22 (σ 0,58)
User, 06(A,0,10)	9h 55m 37s	11m 54s (σ 48m 22s)	57	1,14 (σ 0,40)
User, 07(B,0,10)	8h 4m 40s	9m 41s (σ 31m 29s)	53	1,06 (σ 0,24)
User, 08(B,0,10)	1h 36m 58s	1m 56s (σ 1m 50s)	50	1,00 (σ 0,00)
User, 09(A,2,5)	1d 12m 43s	29m 3s (σ 2h 55m 13s)	58	1,16 (σ 1,12)
User, 10(A,2,5)	2h 41m 31s	3m 13s (σ 4m 1s)	52	1,04 (σ 0,20)
User, 11(B,2,5)	2h 13m 33s	2m 40s (σ 3m 49s)	55	1,10 (σ 0,41)
User, 12(A,2,10)	5h 20m 20s	6m 24s (σ 34m 5s)	52	1,04 (σ 0,20)
User, 13(A,2,10)	2h 45m 46s	3m 18s (σ 3m 11s)	50	1,00 (σ 0,00)
User, 14(B,2,10)	2h 8m 25s	2m 34s (σ 5m 37s)	63	1,26 (σ 1,13)
User, 15(A,4,5)	2h 12m 7s	2m 38s (σ 6m 28s)	300	6,00 (σ 7,91)
User, 16(A,4,5)	3h 8m 45s	3m 46s (σ 8m 14s)	60	1,20 (σ 0,40)
User, 17(B,4,5)	2h 1s	2m 24s (σ 3m 24s)	52	1,04 (σ 0,20)
User, 18(A,4,10)	5h 3m 8s	6m 3s (σ 8m 33s)	134	2,68 (σ 2,60)
User, 19(A,4,10)	3h 26m 43s	4m 8s (σ 13m 13s)	61	1,22 (σ 0,41)
User, 20(B,4,10)	1h 50m 24s	2m 12s (σ 2m 31s)	51	1,02 (σ 0,14)

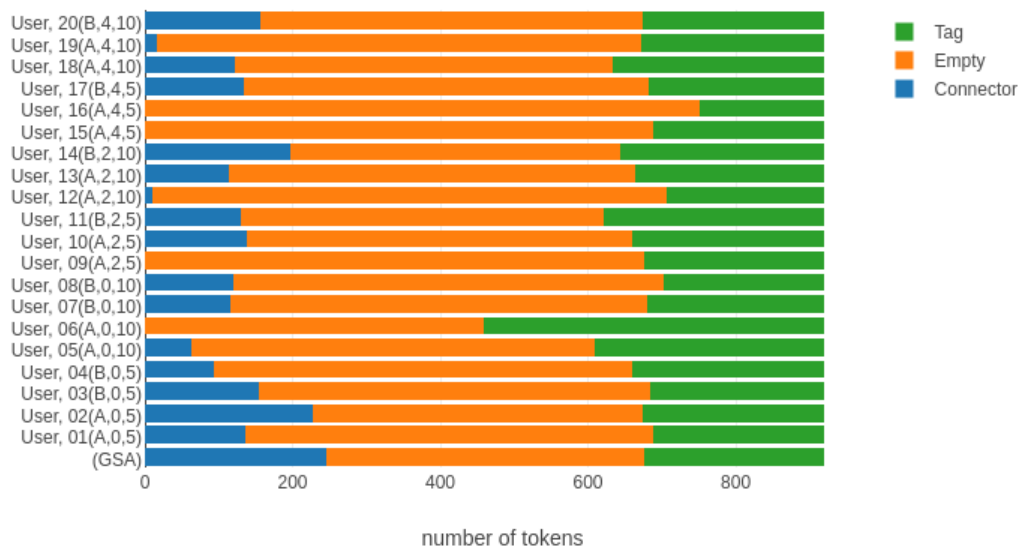


Figura D.1: Experimento de anotação, tokens

PUC-Rio - Certificação Digital Nº 1421597/CA

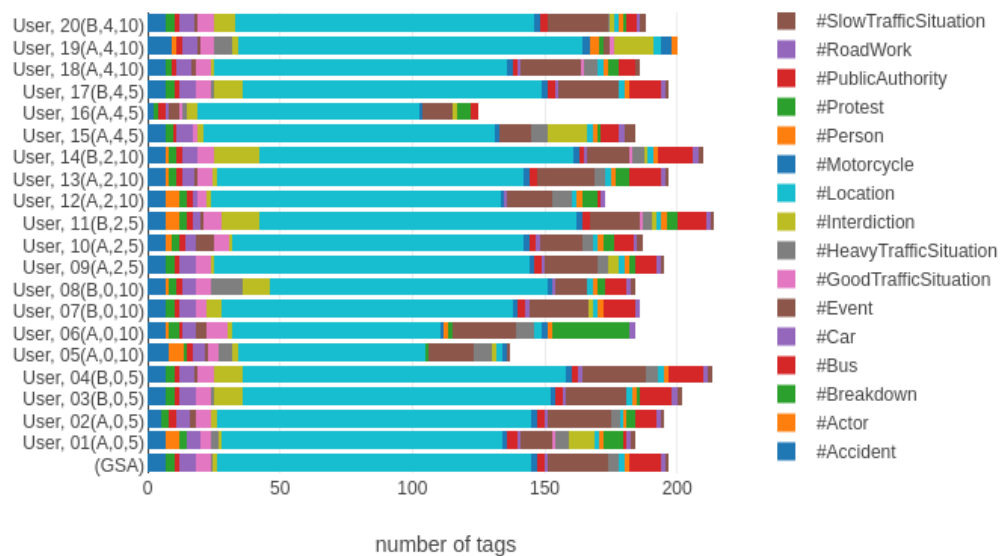


Figura D.2: Experimento de anotação, entidades

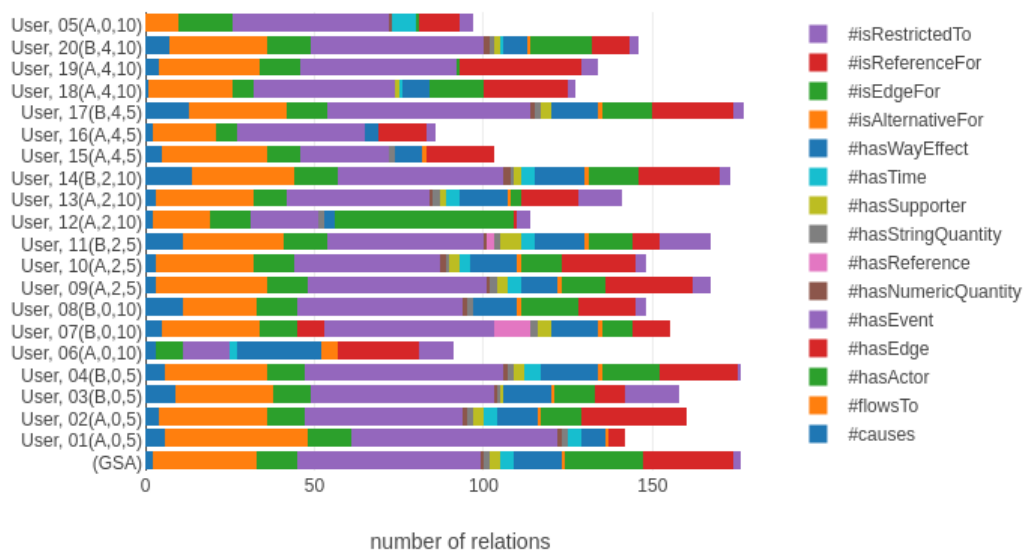


Figura D.3: Experimento de anotação: relações

PUC-Rio - Certificação Digital N° 1421597/CA

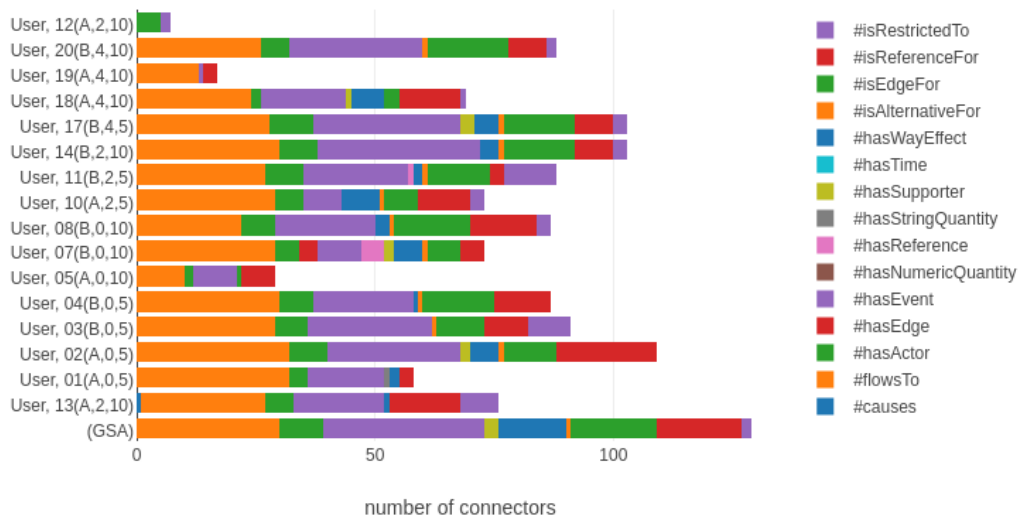
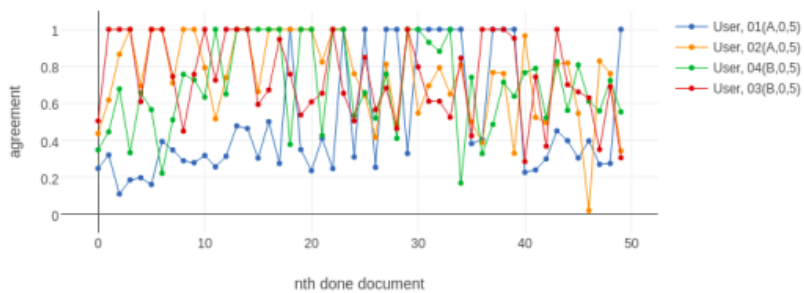
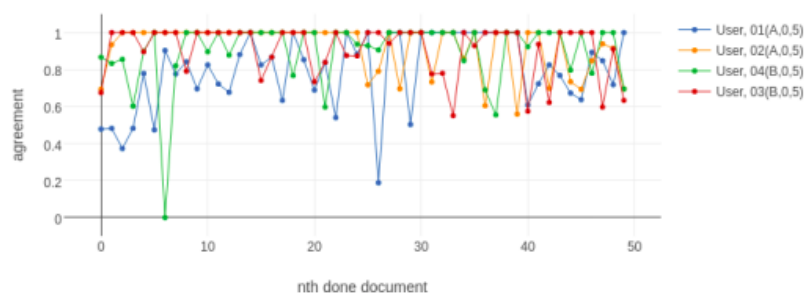


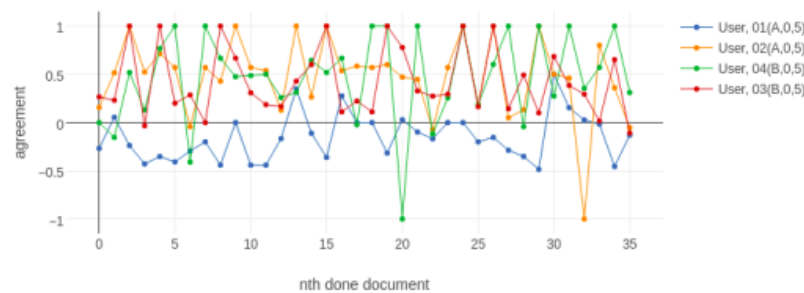
Figura D.4: Experimento de anotação, conectores



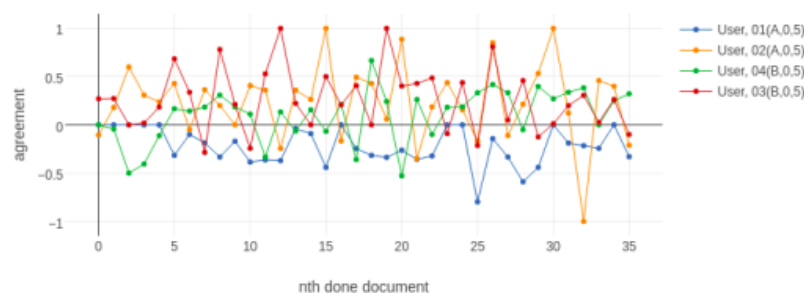
(a)



(b)

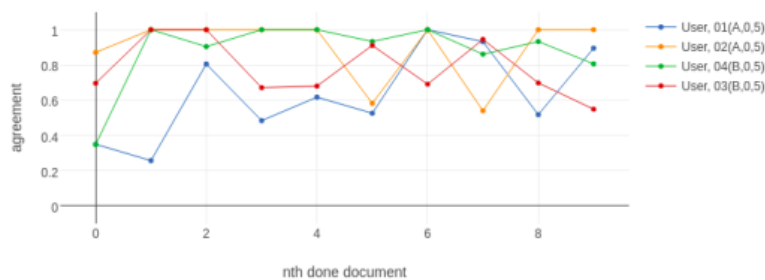


(c)

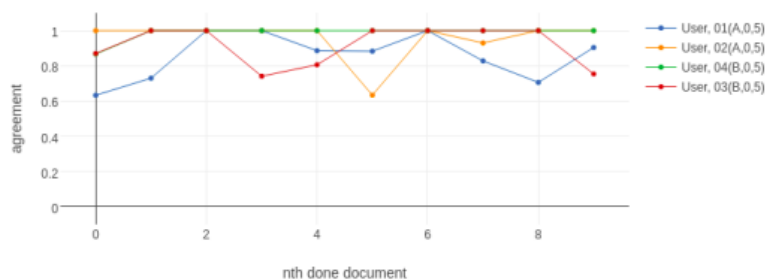


(d)

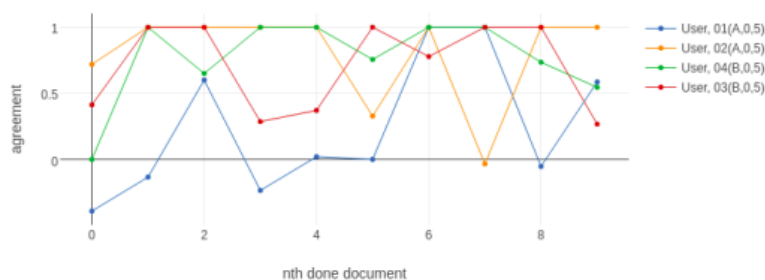
Figura D.9: Experimento de anotação, concordância com GSA ao longo do tempo do grupo (0,5): (a) entidades, relações e conectores; (b) entidades; (c) relações; (d) conectores



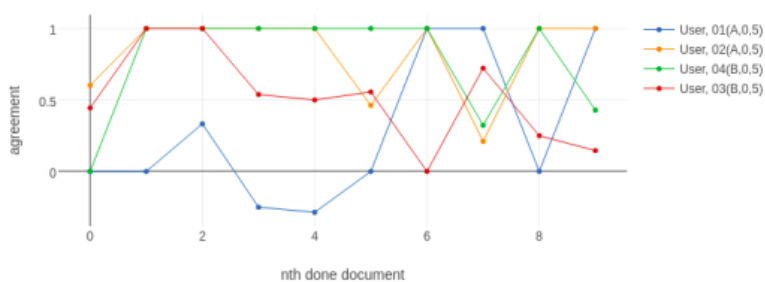
(a)



(b)

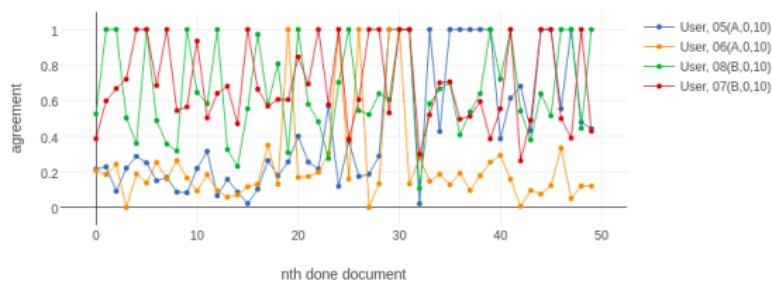


(c)

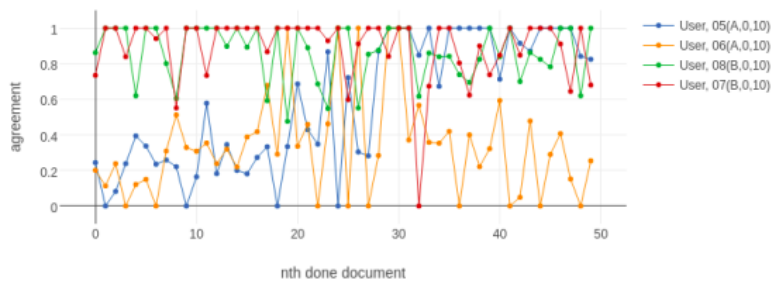


(d)

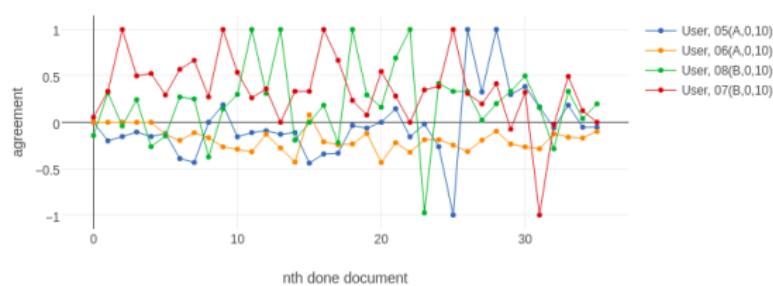
Figura D.10: Experimento de anotação, auto-concordância ao longo do tempo do grupo (0,5): (a) entidades, relações e conectores; (b) entidades; (c) relações; (d) conectores



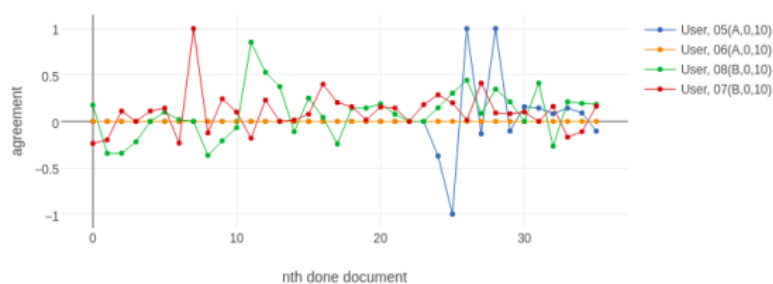
(a)



(b)

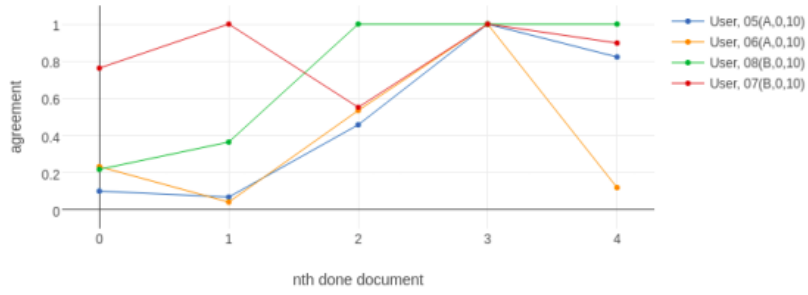


(c)

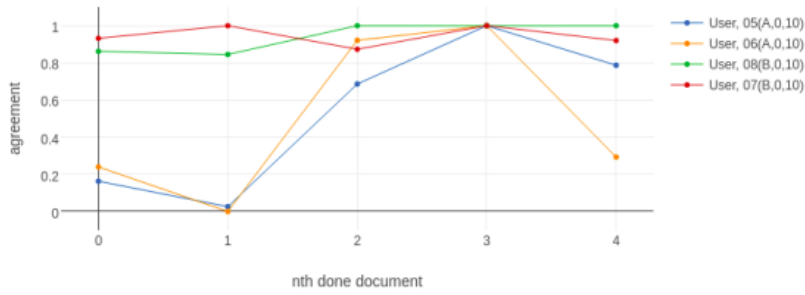


(d)

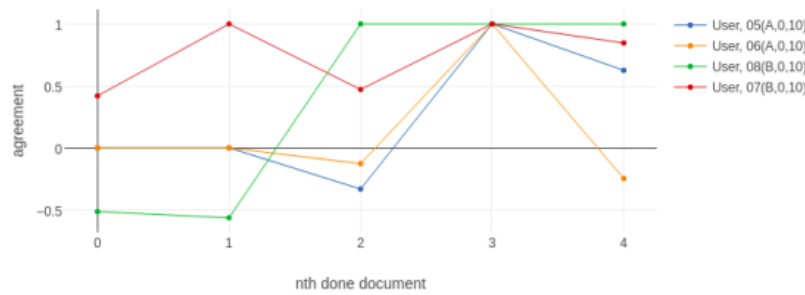
Figura D.11: Experimento de anotação, concordância com GSA ao longo do tempo do grupo (0,10): (a) entidades, relações e conectores; (b) entidades; (c) relações; (d) conectores



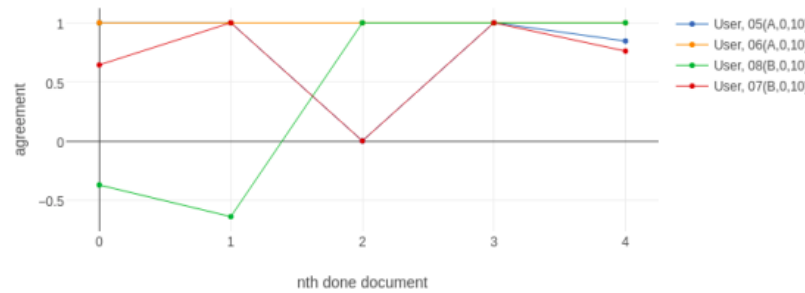
(a)



(b)

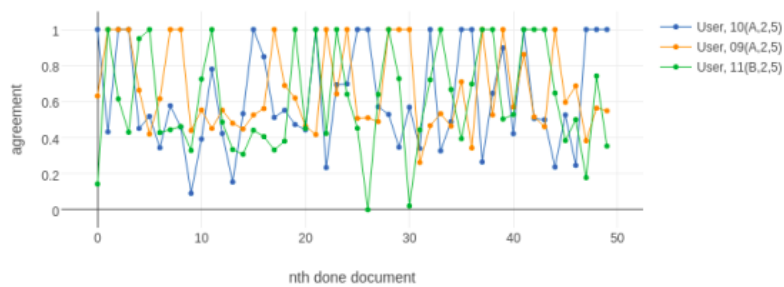


(c)

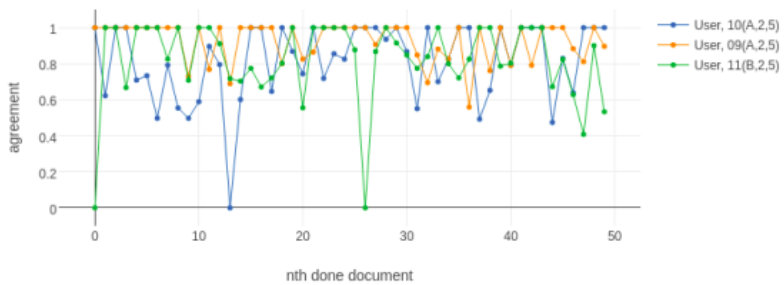


(d)

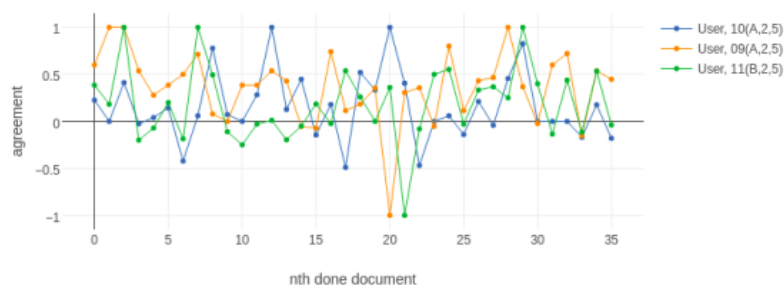
Figura D.12: Experimento de anotação, auto-concordância ao longo do tempo do grupo (0,10): (a) entidades, relações e conectores; (b) entidades; (c) relações; (d) conectores



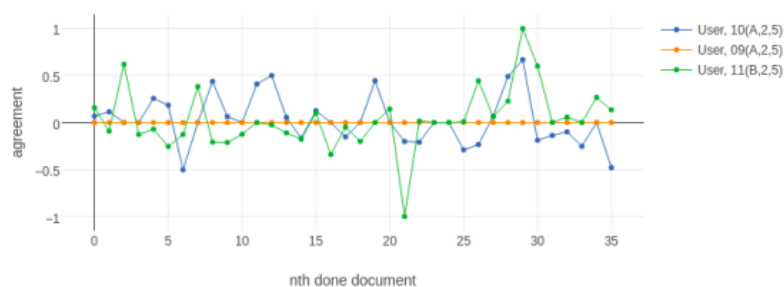
(a)



(b)

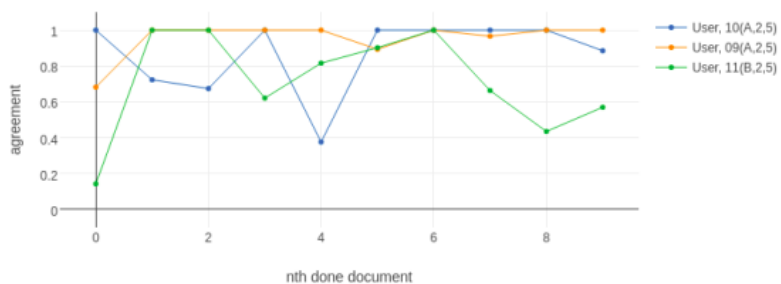


(c)

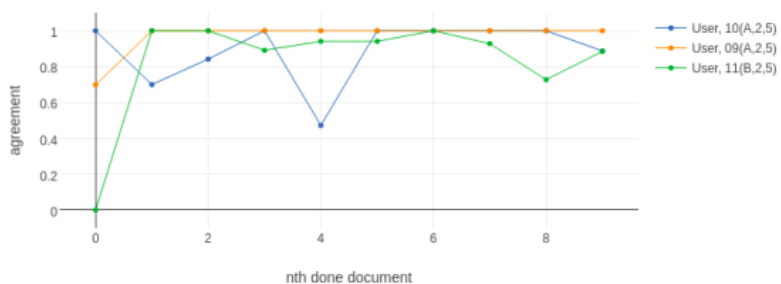


(d)

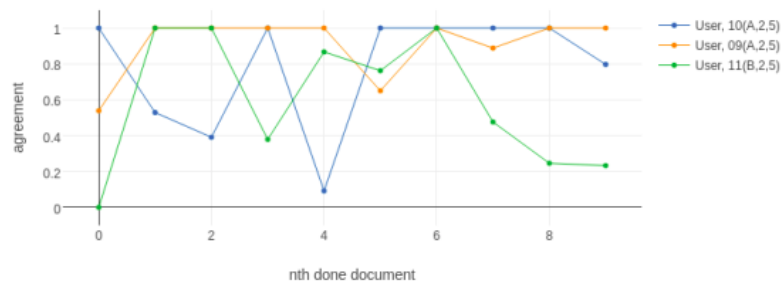
Figura D.13: Experimento de anotação, concordância com GSA ao longo do tempo do grupo (2,5): (a) entidades, relações e conectores; (b) entidades; (c) relações; (d) conectores



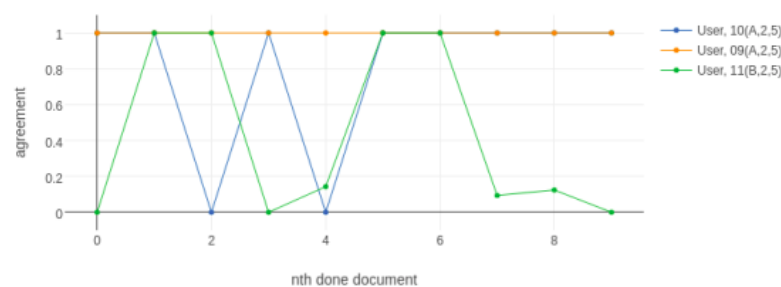
(a)



(b)

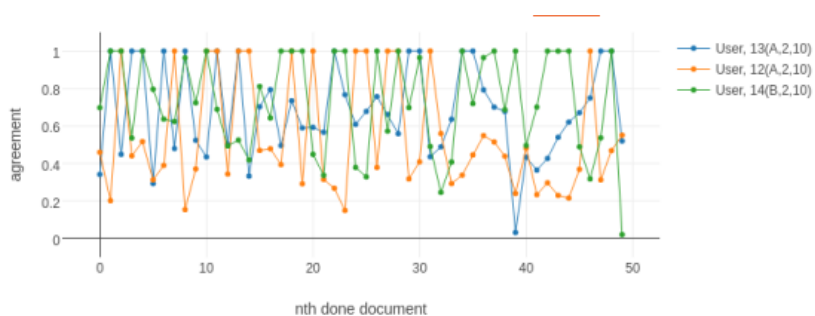


(c)

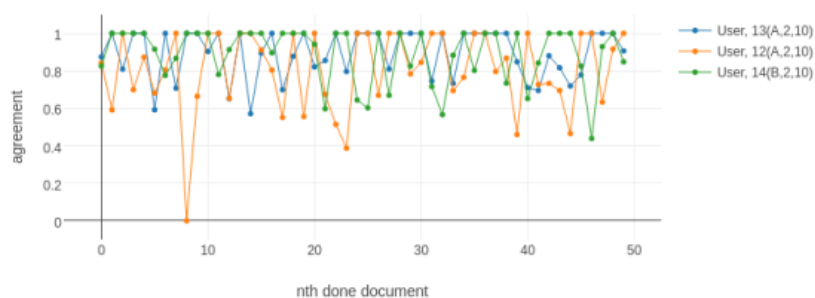


(d)

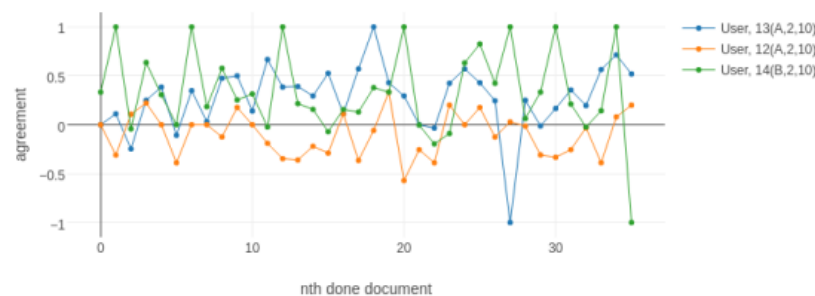
Figura D.14: Experimento de anotação, auto-concordância ao longo do tempo do grupo (2,5): (a) entidades, relações e conectores; (b) entidades; (c) relações; (d) conectores



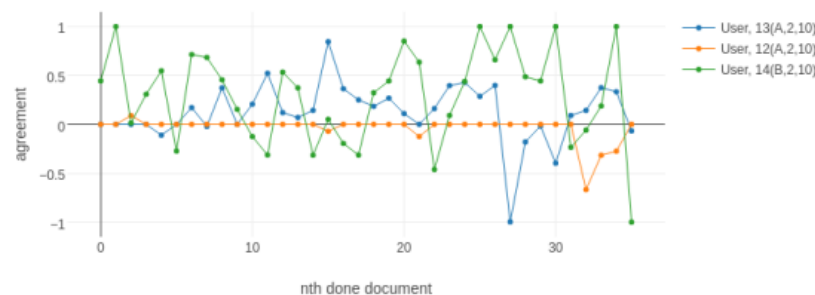
(a)



(b)

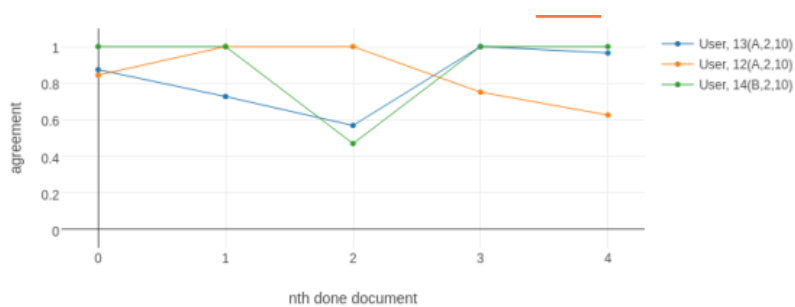


(c)

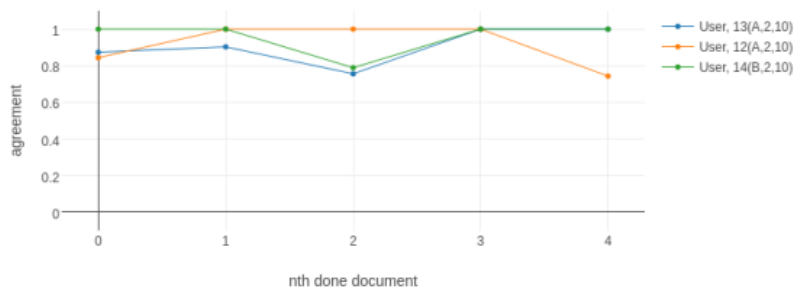


(d)

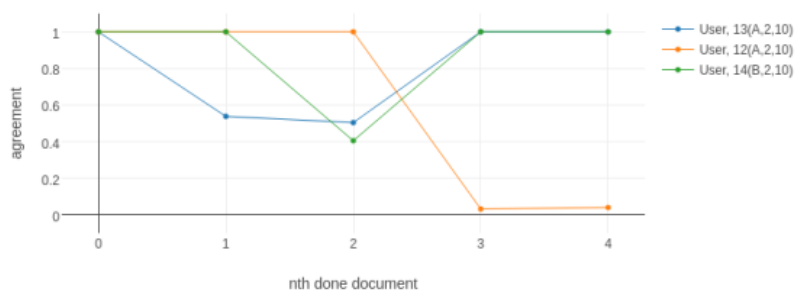
Figura D.15: Experimento de anotação, concordância com GSA ao longo do tempo do grupo (2,10): (a) entidades, relações e conectores; (b) entidades; (c) relações; (d) conectores



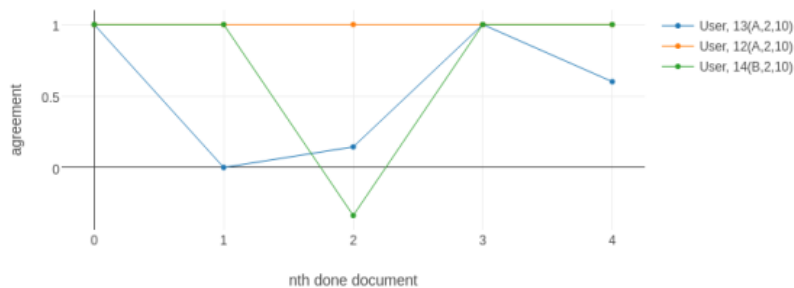
(a)



(b)

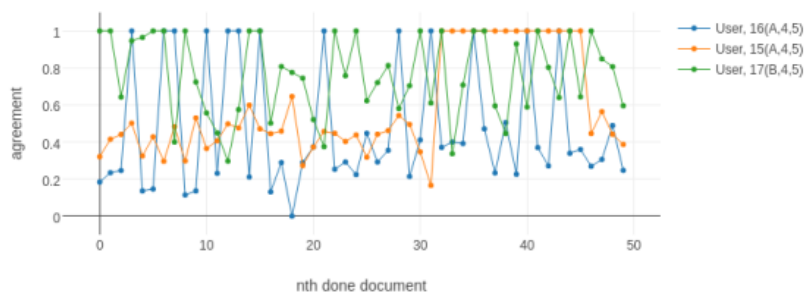


(c)

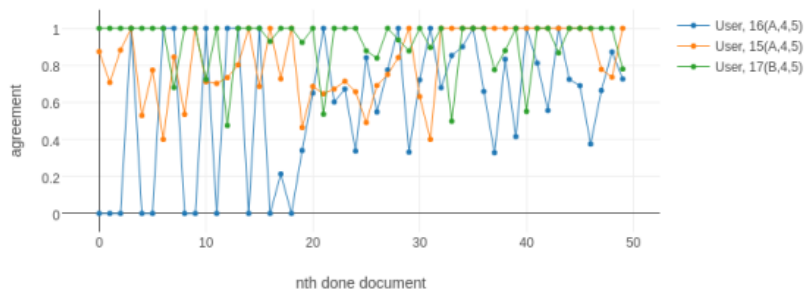


(d)

Figura D.16: Experimento de anotação, auto-concordância ao longo do tempo do grupo (2,10): (a) entidades, relações e conectores; (b) entidades; (c) relações; (d) conectores



(a)



(b)

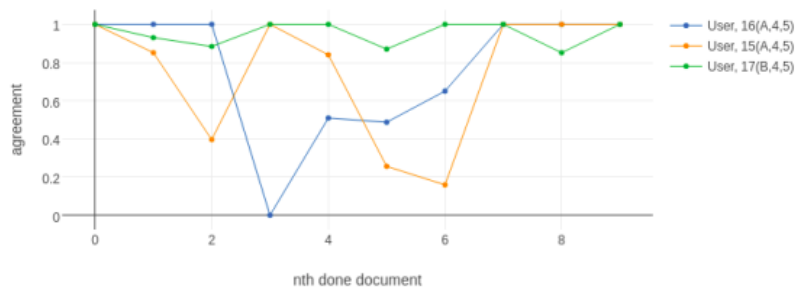


(c)

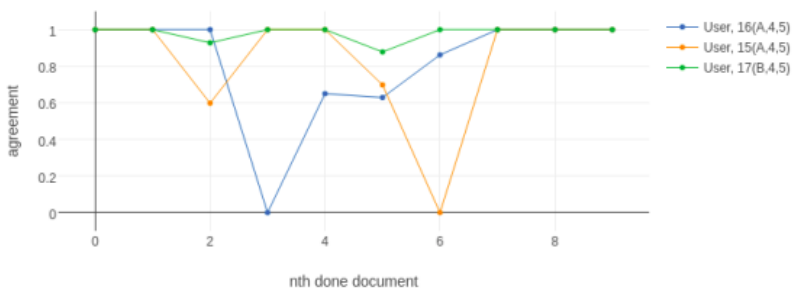


(d)

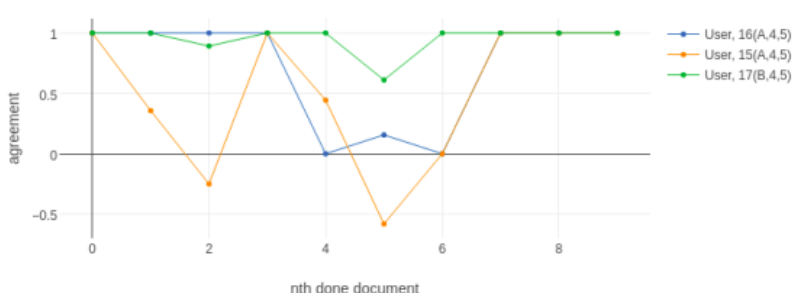
Figura D.17: Experimento de anotação, concordância com GSA ao longo do tempo do grupo (4,5): (a) entidades, relações e conectores; (b) entidades; (c) relações; (d) conectores



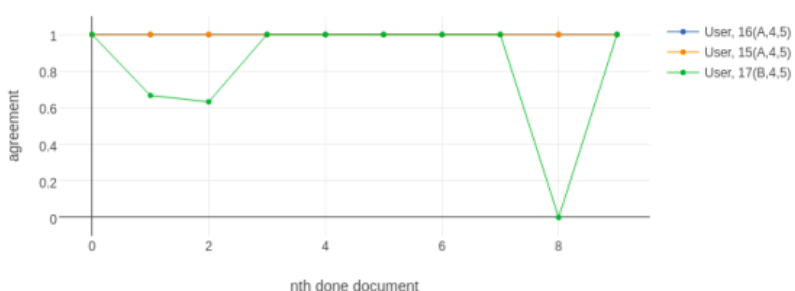
(a)



(b)

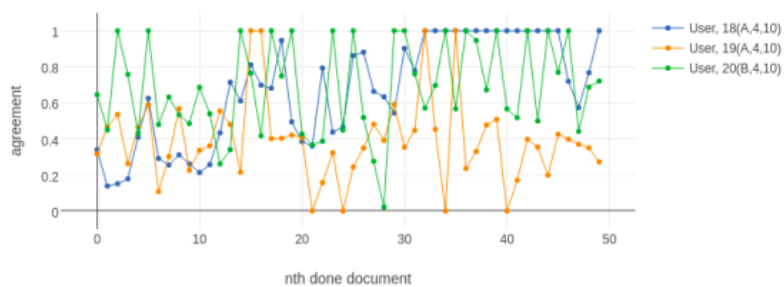


(c)

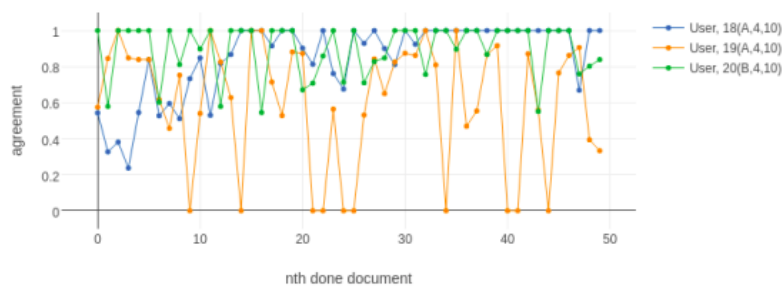


(d)

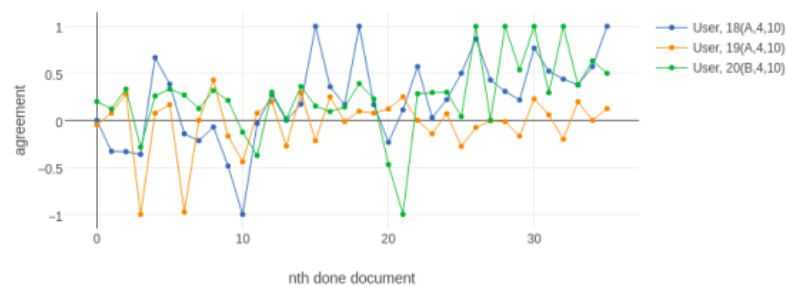
Figura D.18: Experimento de anotação, auto-concordância ao longo do tempo do grupo (4,5): (a) entidades, relações e conectores; (b) entidades; (c) relações; (d) conectores



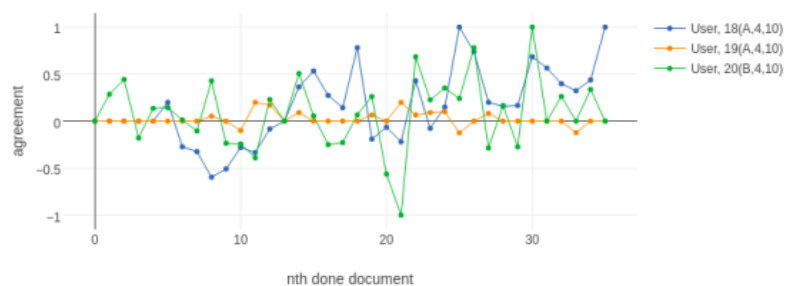
(a)



(b)

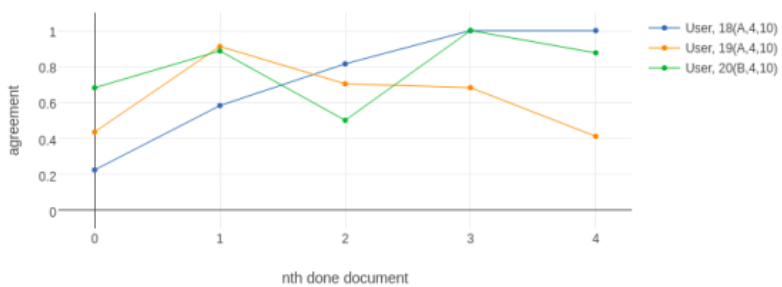


(c)

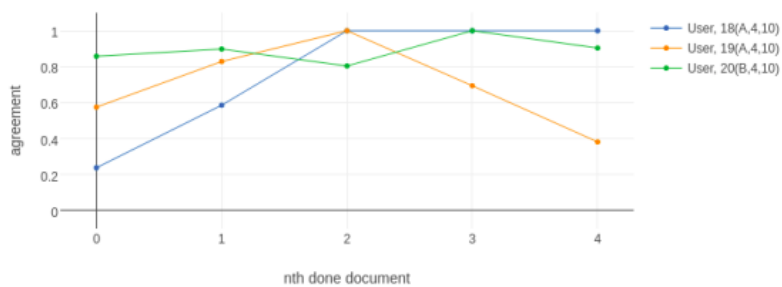


(d)

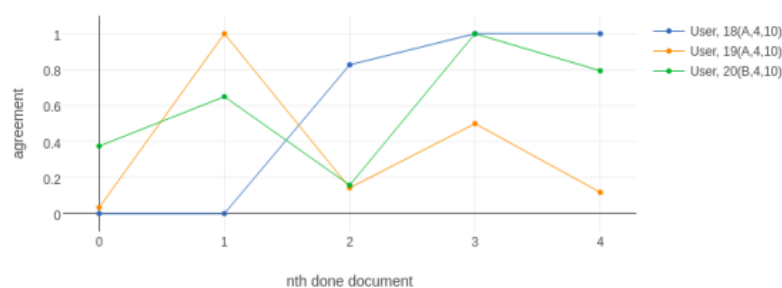
Figura D.19: Experimento de anotação, concordância com GSA ao longo do tempo do grupo (4,10): (a) entidades, relações e conectores; (b) entidades; (c) relações; (d) conectores



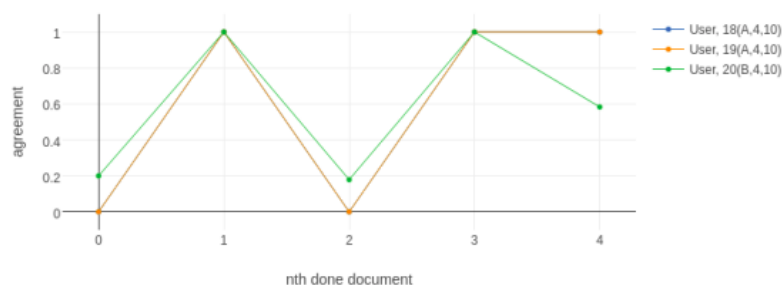
(a)



(b)



(c)



(d)

Figura D.20: Experimento de anotação, auto-concordância ao longo do tempo do grupo (4,10): (a) entidades, relações e conectores; (b) entidades; (c) relações; (d) conectores

E

Experimento de aprendizado automático: Tabelas e gráficos

Este apêndice contém gráficos e tabelas com informações adicionais sobre o experimento de aprendizado automático descrito no capítulo 4 desta dissertação.

Tabela E.1: Training scores: model-fscore-TEDO-NER-RF-0

Class	Support	Accuracy, σ CV	Recall, σ CV	Precision, σ CV	F1, σ CV
O	5262	0.925, 0.018	0.978, 0.011	0.912, 0.019	0.944, 0.014
B-#Accident	114	0.996, 0.002	0.795, 0.162	0.921, 0.096	0.843, 0.119
B-#BadTrafficSituation	273	0.991, 0.005	0.863, 0.093	0.868, 0.097	0.863, 0.084
B-#Breakdown	14	0.999, 0.000	0.531, 0.440	0.625, 0.484	0.565, 0.450
B-#Bus	31	0.997, 0.001	0.646, 0.313	0.777, 0.323	0.693, 0.306
B-#Car	103	0.997, 0.002	0.798, 0.284	0.826, 0.288	0.810, 0.284
B-#Event	32	0.995, 0.001	0.138, 0.314	0.125, 0.259	0.130, 0.282
B-#GoodTrafficSituation	29	0.996, 0.002	0.178, 0.299	0.350, 0.450	0.213, 0.312
B-#Interdiction	141	0.994, 0.003	0.631, 0.238	0.853, 0.290	0.720, 0.255
B-#Location	1270	0.956, 0.011	0.859, 0.083	0.856, 0.054	0.854, 0.054
B-#Motorcycle	18	0.998, 0.001	0.799, 0.258	0.831, 0.207	0.800, 0.218
B-#Protest	16	0.993, 0.002	0.000, 0.000	0.000, 0.000	0.000, 0.000
B-#PublicAuthority	55	0.996, 0.002	0.311, 0.327	0.500, 0.471	0.366, 0.355
B-#RoadWork	31	0.998, 0.001	0.568, 0.268	0.873, 0.311	0.672, 0.282
B-#Solution	45	0.996, 0.003	0.580, 0.270	0.875, 0.330	0.687, 0.282
B-#Time	17	0.998, 0.001	0.428, 0.494	0.428, 0.494	0.428, 0.494
B-#Truck	24	0.997, 0.002	0.320, 0.448	0.380, 0.468	0.322, 0.431
B-#wayEffect:BothDirections	19	0.999, 0.000	0.583, 0.478	0.600, 0.489	0.590, 0.483
B-#wayEffect:OneDirection	3	0.998, 0.000	0.000, 0.000	0.000, 0.000	0.000, 0.000
B-#wayEffect:Partially	123	0.997, 0.001	0.741, 0.263	0.887, 0.296	0.805, 0.273
I-#BadTrafficSituation	85	0.997, 0.002	0.867, 0.132	0.925, 0.065	0.891, 0.090
I-#Event	8	0.998, 0.000	0.300, 0.400	0.333, 0.421	0.314, 0.408
I-#GoodTrafficSituation	28	0.997, 0.002	0.300, 0.378	0.500, 0.500	0.350, 0.390
I-#Location	306	0.974, 0.014	0.399, 0.239	0.741, 0.277	0.500, 0.247
I-#Protest	2	0.998, 0.000	0.500, 0.000	1.000, 0.000	0.666, 0.000
I-#PublicAuthority	9	0.997, 0.001	0.576, 0.318	0.809, 0.269	0.655, 0.285
I-#Time	8	0.995, 0.001	0.000, 0.000	0.000, 0.000	0.000, 0.000
I-#wayEffect:BothDirections	22	0.999, 0.001	0.700, 0.400	0.766, 0.395	0.723, 0.391
I-#wayEffect:Partially	2	0.998, 0.000	0.500, 0.000	1.000, 0.000	0.666, 0.000
Average, σ		0.992, 0.015	0.513, 0.275	0.640, 0.309	0.554, 0.281

Tabela E.2: Revalidation scores: model-fscore-TEDO-NER-RF-0

Class	Support	Accuracy	Recall	Precision	F1
O	1195	0.940	0.986	0.923	0.954
B-#Accident	27	0.996	0.888	0.857	0.872
B-#BadTrafficSituation	66	0.990	0.909	0.845	0.875
B-#Breakdown	3	0.999	0.666	1.000	0.800
B-#Bus	8	0.997	0.625	0.833	0.714
B-#Car	23	0.998	0.869	1.000	0.930
B-#Event	5	0.996	0.000	0.000	0.000
B-#GoodTrafficSituation	16	0.994	0.312	1.000	0.476
B-#Interdiction	24	0.997	0.916	0.916	0.916
B-#Location	316	0.964	0.892	0.895	0.893
B-#Motorcycle	4	0.998	0.750	0.750	0.750
B-#Protest	1	0.999	0.000	0.000	0.000
B-#PublicAuthority	8	0.999	0.875	1.000	0.933
B-#RoadWork	1	1.000	1.000	1.000	1.000
B-#Solution	10	0.998	0.700	1.000	0.823
B-#Time	3	1.000	1.000	1.000	1.000
B-#Truck	6	0.996	0.166	0.333	0.222
B-#wayEffect:BothDirections	11	0.999	0.909	1.000	0.952
B-#wayEffect:OneDirection	1	0.999	0.000	0.000	0.000
B-#wayEffect:Partially	23	0.999	0.956	1.000	0.977
I-#BadTrafficSituation	20	0.999	0.950	1.000	0.974
I-#GoodTrafficSituation	15	0.996	0.000	0.000	0.000
I-#Location	90	0.969	0.600	1.000	0.749
I-#wayEffect:BothDirections	12	0.998	0.422	0.883	0.571
Average, σ		0.993, 0.014	0.641, 0.357	0.759, 0.367	0.682, 0.353

Tabela E.3: Training scores: model-fscore-TEDO-NER-SGD-0

Class	Support	Accuracy, σ CV	Recall, σ CV	Precision, σ CV	F1, σ CV
O	5262	0.923, 0.022	0.957, 0.013	0.928, 0.033	0.942, 0.015
B-#Accident	114	0.998, 0.001	0.888, 0.070	1.000, 0.000	0.939, 0.039
B-#BadTrafficSituation	273	0.991, 0.003	0.918, 0.075	0.858, 0.091	0.881, 0.046
B-#Breakdown	14	0.999, 0.000	0.687, 0.428	0.750, 0.433	0.708, 0.422
B-#Bus	31	0.998, 0.001	0.918, 0.153	0.902, 0.184	0.887, 0.130
B-#Car	103	0.998, 0.001	0.938, 0.065	0.983, 0.032	0.959, 0.035
B-#Event	32	0.996, 0.001	0.259, 0.312	0.412, 0.405	0.281, 0.281
B-#GoodTrafficSituation	29	0.996, 0.002	0.202, 0.318	0.310, 0.406	0.204, 0.277
B-#Interdiction	141	0.994, 0.004	0.837, 0.160	0.775, 0.235	0.788, 0.174
B-#Location	1270	0.950, 0.013	0.864, 0.081	0.840, 0.073	0.846, 0.038
B-#Motorcycle	18	0.999, 0.000	0.854, 0.204	0.934, 0.124	0.873, 0.145
B-#Protest	16	0.994, 0.001	0.327, 0.061	0.750, 0.353	0.419, 0.117
B-#PublicAuthority	55	0.996, 0.002	0.467, 0.394	0.592, 0.465	0.511, 0.411
B-#RoadWork	31	0.999, 0.001	0.698, 0.361	0.833, 0.314	0.728, 0.317
B-#Solution	45	0.997, 0.001	0.841, 0.184	0.925, 0.139	0.856, 0.112
B-#Time	17	0.997, 0.001	0.380, 0.451	0.309, 0.382	0.333, 0.398
B-#Truck	24	0.999, 0.000	0.760, 0.388	0.765, 0.388	0.760, 0.384
B-#wayEffect:BothDirections	19	0.999, 0.001	0.600, 0.489	0.525, 0.453	0.552, 0.461
B-#wayEffect:OneDirection	3	0.998, 0.000	0.000, 0.000	0.000, 0.000	0.000, 0.000
B-#wayEffect:Partially	123	0.996, 0.001	0.757, 0.293	0.851, 0.290	0.791, 0.278
I-#BadTrafficSituation	85	0.998, 0.001	0.779, 0.284	0.835, 0.290	0.800, 0.279
I-#Event	8	0.998, 0.000	0.459, 0.407	0.533, 0.452	0.492, 0.425
I-#GoodTrafficSituation	28	0.997, 0.002	0.393, 0.482	0.335, 0.431	0.355, 0.443
I-#Location	306	0.964, 0.014	0.359, 0.342	0.609, 0.356	0.382, 0.255
I-#Protest	2	0.998, 0.000	0.260, 0.000	0.272, 0.000	0.266, 0.000
I-#PublicAuthority	9	0.996, 0.001	0.089, 0.126	0.259, 0.366	0.133, 0.188
I-#Time	8	0.995, 0.001	0.000, 0.000	0.000, 0.000	0.000, 0.000
I-#wayEffect:BothDirections	22	0.998, 0.000	0.318, 0.322	0.386, 0.415	0.309, 0.322
I-#wayEffect:Partially	2	0.998, 0.000	0.500, 0.000	1.000, 0.000	0.666, 0.000
Average, σ		0.992, 0.016	0.562, 0.296	0.637, 0.294	0.574, 0.293

Tabela E.4: Revalidation scores: model-fscore-TEDO-NER-SGD-0

Class	Support	Accuracy	Recall	Precision	F1
O	1195	0.927	0.960	0.927	0.943
B-#Accident	27	0.998	0.925	0.961	0.943
B-#BadTrafficSituation	66	0.986	0.984	0.730	0.838
B-#Breakdown	3	0.999	0.666	1.000	0.800
B-#Bus	8	1.000	1.000	1.000	1.000
B-#Car	23	0.998	0.956	0.956	0.956
B-#Event	5	0.997	0.400	0.666	0.500
B-#GoodTrafficSituation	16	0.993	0.437	0.636	0.518
B-#Interdiction	24	0.999	0.958	1.000	0.978
B-#Location	316	0.961	0.936	0.850	0.891
B-#Motorcycle	4	0.999	1.000	0.800	0.888
B-#Protest	1	1.000	1.000	1.000	1.000
B-#PublicAuthority	8	1.000	1.000	1.000	1.000
B-#RoadWork	1	0.999	1.000	0.500	0.666
B-#Solution	10	0.998	0.900	0.818	0.857
B-#Time	3	1.000	1.000	1.000	1.000
B-#Truck	6	0.999	0.833	1.000	0.909
B-#wayEffect:BothDirections	11	0.995	0.181	1.000	0.307
B-#wayEffect:OneDirection	1	0.998	0.000	0.000	0.000
B-#wayEffect:Partially	23	0.999	1.000	0.958	0.978
I-#BadTrafficSituation	20	0.999	1.000	0.952	0.975
I-#GoodTrafficSituation	15	0.997	0.000	0.000	0.000
I-#Location	90	0.958	0.733	1.000	0.846
I-#wayEffect:BothDirections	12	0.995	0.144	0.866	0.247
Average, σ		0.991, 0.017	0.750, 0.342	0.817, 0.280	0.752, 0.310

Tabela E.5: Training scores: model-fscore-TEDO-NER-SVC-0

Class	Support	Accuracy, σ CV	Recall, σ CV	Precision, σ CV	F1, σ CV
O	5262	0.805, 0.017	0.959, 0.009	0.787, 0.020	0.865, 0.013
B-#Accident	114	0.985, 0.004	0.000, 0.000	0.000, 0.000	0.000, 0.000
B-#BadTrafficSituation	273	0.966, 0.008	0.000, 0.000	0.000, 0.000	0.000, 0.000
B-#Breakdown	14	0.997, 0.001	0.000, 0.000	0.000, 0.000	0.000, 0.000
B-#Bus	31	0.995, 0.001	0.000, 0.000	0.000, 0.000	0.000, 0.000
B-#Car	103	0.987, 0.004	0.000, 0.000	0.000, 0.000	0.000, 0.000
B-#Event	32	0.995, 0.001	0.000, 0.000	0.000, 0.000	0.000, 0.000
B-#GoodTrafficSituation	29	0.996, 0.002	0.000, 0.000	0.000, 0.000	0.000, 0.000
B-#Interdiction	141	0.982, 0.007	0.000, 0.000	0.000, 0.000	0.000, 0.000
B-#Location	1270	0.910, 0.012	0.876, 0.038	0.662, 0.048	0.753, 0.042
B-#Motorcycle	18	0.997, 0.001	0.000, 0.000	0.000, 0.000	0.000, 0.000
B-#Protest	16	0.993, 0.002	0.000, 0.000	0.000, 0.000	0.000, 0.000
B-#PublicAuthority	55	0.992, 0.005	0.000, 0.000	0.000, 0.000	0.000, 0.000
B-#RoadWork	31	0.995, 0.002	0.000, 0.000	0.000, 0.000	0.000, 0.000
B-#Solution	45	0.993, 0.005	0.000, 0.000	0.000, 0.000	0.000, 0.000
B-#Time	17	0.996, 0.001	0.000, 0.000	0.000, 0.000	0.000, 0.000
B-#Truck	24	0.997, 0.002	0.000, 0.000	0.000, 0.000	0.000, 0.000
B-#wayEffect:BothDirections	19	0.997, 0.001	0.000, 0.000	0.000, 0.000	0.000, 0.000
B-#wayEffect:OneDirection	3	0.998, 0.000	0.000, 0.000	0.000, 0.000	0.000, 0.000
B-#wayEffect:Partially	123	0.984, 0.005	0.000, 0.000	0.000, 0.000	0.000, 0.000
I-#BadTrafficSituation	85	0.989, 0.004	0.000, 0.000	0.000, 0.000	0.000, 0.000
I-#Event	8	0.998, 0.000	0.000, 0.000	0.000, 0.000	0.000, 0.000
I-#GoodTrafficSituation	28	0.996, 0.002	0.000, 0.000	0.000, 0.000	0.000, 0.000
I-#Location	306	0.962, 0.022	0.000, 0.000	0.000, 0.000	0.000, 0.000
I-#Protest	2	0.997, 0.000	0.000, 0.000	0.000, 0.000	0.000, 0.000
I-#PublicAuthority	9	0.996, 0.001	0.000, 0.000	0.000, 0.000	0.000, 0.000
I-#Time	8	0.995, 0.001	0.000, 0.000	0.000, 0.000	0.000, 0.000
I-#wayEffect:BothDirections	22	0.997, 0.001	0.000, 0.000	0.000, 0.000	0.000, 0.000
I-#wayEffect:Partially	2	0.997, 0.000	0.000, 0.000	0.000, 0.000	0.000, 0.000
Average, σ		0.982, 0.037	0.063, 0.232	0.050, 0.184	0.055, 0.205

Tabela E.6: Revalidation scores: model-fscore-TEDO-NER-SVC-0

Class	Support	Accuracy	Recall	Precision	F1
O	1195	0.798	0.963	0.773	0.857
B-#Accident	27	0.985	0.000	0.000	0.000
B-#BadTrafficSituation	66	0.965	0.000	0.000	0.000
B-#Breakdown	3	0.998	0.000	0.000	0.000
B-#Bus	8	0.995	0.000	0.000	0.000
B-#Car	23	0.987	0.000	0.000	0.000
B-#Event	5	0.997	0.000	0.000	0.000
B-#GoodTrafficSituation	16	0.991	0.000	0.000	0.000
B-#Interdiction	24	0.987	0.000	0.000	0.000
B-#Location	316	0.909	0.863	0.682	0.762
B-#Motorcycle	4	0.997	0.000	0.000	0.000
B-#Protest	1	0.999	0.000	0.000	0.000
B-#PublicAuthority	8	0.995	0.000	0.000	0.000
B-#RoadWork	1	0.999	0.000	0.000	0.000
B-#Solution	10	0.994	0.000	0.000	0.000
B-#Time	3	0.998	0.000	0.000	0.000
B-#Truck	6	0.996	0.000	0.000	0.000
B-#wayEffect:BothDirections	11	0.994	0.000	0.000	0.000
B-#wayEffect:OneDirection	1	0.999	0.000	0.000	0.000
B-#wayEffect:Partially	23	0.987	0.000	0.000	0.000
I-#BadTrafficSituation	20	0.989	0.000	0.000	0.000
I-#GoodTrafficSituation	15	0.992	0.000	0.000	0.000
I-#Location	90	0.952	0.000	0.000	0.000
I-#wayEffect:BothDirections	12	0.993	0.000	0.000	0.000
Average, σ		0.979, 0.042	0.076, 0.252	0.060, 0.201	0.067, 0.224

Tabela E.7: Training scores: model-fscore-TEDO-NER-RF-1

Class	Support	Accuracy, σ CV	Recall, σ CV	Precision, σ CV	F1, σ CV
O	5262	0.937, 0.017	0.982, 0.007	0.926, 0.020	0.953, 0.012
B-#Accident	114	0.997, 0.001	0.863, 0.068	0.916, 0.098	0.886, 0.069
B-#BadTrafficSituation	273	0.991, 0.002	0.867, 0.076	0.865, 0.057	0.865, 0.060
B-#Breakdown	14	0.999, 0.000	0.531, 0.440	0.625, 0.484	0.565, 0.450
B-#Bus	31	0.998, 0.000	0.785, 0.177	0.944, 0.157	0.831, 0.119
B-#Car	103	0.996, 0.002	0.746, 0.302	0.781, 0.288	0.759, 0.293
B-#Event	32	0.995, 0.000	0.253, 0.313	0.440, 0.419	0.288, 0.298
B-#GoodTrafficSituation	29	0.997, 0.002	0.274, 0.325	0.600, 0.489	0.342, 0.351
B-#Interdiction	141	0.993, 0.002	0.778, 0.129	0.872, 0.148	0.812, 0.107
B-#Location	1270	0.973, 0.009	0.910, 0.094	0.875, 0.092	0.892, 0.091
B-#Motorcycle	18	0.998, 0.001	0.703, 0.356	0.633, 0.342	0.656, 0.340
B-#Protest	16	0.994, 0.002	0.133, 0.188	0.333, 0.471	0.190, 0.269
B-#PublicAuthority	55	0.996, 0.001	0.398, 0.382	0.555, 0.496	0.456, 0.419
B-#RoadWork	31	0.998, 0.000	0.444, 0.316	0.750, 0.408	0.531, 0.326
B-#Solution	45	0.997, 0.001	0.727, 0.302	0.850, 0.327	0.775, 0.300
B-#Time	17	0.998, 0.001	0.380, 0.451	0.428, 0.494	0.399, 0.465
B-#Truck	24	0.997, 0.001	0.370, 0.388	0.530, 0.443	0.390, 0.353
B-#wayEffect:BothDirections	19	0.999, 0.000	0.500, 0.500	0.500, 0.500	0.500, 0.500
B-#wayEffect:OneDirection	3	0.998, 0.000	0.000, 0.000	0.000, 0.000	0.000, 0.000
B-#wayEffect:Partially	123	0.996, 0.002	0.773, 0.282	0.869, 0.296	0.815, 0.284
I-#BadTrafficSituation	85	0.996, 0.002	0.766, 0.275	0.861, 0.293	0.805, 0.274
I-#Event	8	0.998, 0.000	0.521, 0.434	0.420, 0.355	0.455, 0.372
I-#GoodTrafficSituation	28	0.997, 0.002	0.291, 0.375	0.350, 0.450	0.315, 0.403
I-#Location	306	0.976, 0.014	0.380, 0.336	0.681, 0.449	0.461, 0.353
I-#Protest	2	0.998, 0.000	0.500, 0.000	1.000, 0.000	0.666, 0.000
I-#PublicAuthority	9	0.997, 0.001	0.178, 0.127	0.583, 0.424	0.271, 0.191
I-#Time	8	0.995, 0.001	0.192, 0.038	0.672, 0.227	0.297, 0.069
I-#wayEffect:BothDirections	22	0.999, 0.000	0.592, 0.389	0.684, 0.389	0.608, 0.362
I-#wayEffect:Partially	2	0.998, 0.000	0.500, 0.000	1.000, 0.000	0.666, 0.000
Average, σ		0.993, 0.012	0.529, 0.256	0.674, 0.233	0.567, 0.243

Tabela E.8: Revalidation scores: model-fscore-TEDO-NER-RF-1

Class	Support	Accuracy	Recall	Precision	F1
O	1195	0.949	0.989	0.935	0.961
B-#Accident	27	0.997	0.925	0.925	0.925
B-#BadTrafficSituation	66	0.989	0.893	0.819	0.855
B-#Breakdown	3	0.999	0.666	1.000	0.800
B-#Bus	8	0.999	0.875	1.000	0.933
B-#Car	23	0.997	0.956	0.846	0.897
B-#Event	5	0.996	0.000	0.000	0.000
B-#GoodTrafficSituation	16	0.995	0.437	1.000	0.608
B-#Interdiction	24	0.997	0.916	0.880	0.897
B-#Location	316	0.970	0.911	0.911	0.911
B-#Motorcycle	4	0.998	0.500	1.000	0.666
B-#Protest	1	1.000	1.000	1.000	1.000
B-#PublicAuthority	8	1.000	1.000	1.000	1.000
B-#RoadWork	1	1.000	1.000	1.000	1.000
B-#Solution	10	0.998	0.700	1.000	0.823
B-#Time	3	1.000	1.000	1.000	1.000
B-#Truck	6	0.996	0.166	0.333	0.222
B-#wayEffect:BothDirections	11	1.000	1.000	1.000	1.000
B-#wayEffect:OneDirection	1	0.999	0.000	0.000	0.000
B-#wayEffect:Partially	23	0.999	0.956	1.000	0.977
I-#BadTrafficSituation	20	0.998	0.950	0.904	0.926
I-#GoodTrafficSituation	15	0.996	0.533	1.000	0.695
I-#Location	90	0.969	0.411	0.902	0.564
I-#wayEffect:BothDirections	12	1.000	1.000	1.000	1.000
Average, σ		0.993, 0.012	0.741, 0.321	0.852, 0.290	0.777, 0.296

Tabela E.9: Training scores: model-fscore-TEDO-NER-SGD-1

Class	Support	Accuracy, σ CV	Recall, σ CV	Precision, σ CV	F1, σ CV
O	5262	0.939, 0.014	0.974, 0.011	0.935, 0.019	0.954, 0.010
B-#Accident	114	0.997, 0.002	0.881, 0.127	0.938, 0.138	0.897, 0.110
B-#BadTrafficSituation	273	0.993, 0.002	0.896, 0.073	0.899, 0.066	0.894, 0.037
B-#Breakdown	14	0.999, 0.000	0.593, 0.466	0.600, 0.469	0.593, 0.462
B-#Bus	31	0.999, 0.000	0.955, 0.125	0.962, 0.104	0.950, 0.094
B-#Car	103	0.998, 0.001	0.829, 0.284	0.850, 0.291	0.836, 0.281
B-#Event	32	0.995, 0.001	0.244, 0.312	0.425, 0.438	0.274, 0.310
B-#GoodTrafficSituation	29	0.997, 0.001	0.329, 0.342	0.454, 0.433	0.343, 0.337
B-#Interdiction	141	0.995, 0.002	0.837, 0.146	0.903, 0.194	0.840, 0.144
B-#Location	1270	0.966, 0.006	0.882, 0.053	0.868, 0.102	0.871, 0.062
B-#Motorcycle	18	0.999, 0.000	0.896, 0.174	0.896, 0.174	0.870, 0.133
B-#Protest	16	0.993, 0.001	0.219, 0.086	0.777, 0.314	0.296, 0.052
B-#PublicAuthority	55	0.996, 0.002	0.229, 0.340	0.333, 0.471	0.266, 0.385
B-#RoadWork	31	0.999, 0.000	0.704, 0.319	0.841, 0.309	0.750, 0.295
B-#Solution	45	0.998, 0.000	0.825, 0.330	0.833, 0.333	0.820, 0.321
B-#Time	17	0.997, 0.001	0.523, 0.466	0.488, 0.437	0.493, 0.432
B-#Truck	24	0.999, 0.001	0.481, 0.424	0.570, 0.469	0.511, 0.430
B-#wayEffect:BothDirections	19	0.998, 0.001	0.433, 0.472	0.500, 0.500	0.450, 0.471
B-#wayEffect:OneDirection	3	0.998, 0.000	0.000, 0.000	0.000, 0.000	0.000, 0.000
B-#wayEffect:Partially	123	0.996, 0.002	0.846, 0.288	0.840, 0.293	0.841, 0.287
I-#BadTrafficSituation	85	0.997, 0.001	0.731, 0.342	0.809, 0.310	0.726, 0.309
I-#Event	8	0.998, 0.000	0.860, 0.195	0.860, 0.195	0.860, 0.195
I-#GoodTrafficSituation	28	0.996, 0.002	0.366, 0.411	0.382, 0.412	0.368, 0.407
I-#Location	306	0.973, 0.014	0.347, 0.372	0.622, 0.417	0.392, 0.352
I-#Protest	2	0.998, 0.000	0.304, 0.000	0.583, 0.000	0.400, 0.000
I-#PublicAuthority	9	0.996, 0.001	0.430, 0.156	0.809, 0.269	0.503, 0.128
I-#Time	8	0.995, 0.001	0.000, 0.000	0.000, 0.000	0.000, 0.000
I-#wayEffect:BothDirections	22	0.999, 0.000	0.666, 0.387	0.742, 0.408	0.673, 0.368
I-#wayEffect:Partially	2	0.998, 0.000	1.000, 0.000	1.000, 0.000	1.000, 0.000
Average, σ		0.993, 0.012	0.596, 0.297	0.680, 0.264	0.609, 0.284

Tabela E.10: Revalidation scores: model-fscore-TEDO-NER-SGD-1

Class	Support	Accuracy	Recall	Precision	F1
O	1195	0.945	0.975	0.940	0.957
B-#Accident	27	0.998	0.925	1.000	0.961
B-#BadTrafficSituation	66	0.991	1.000	0.804	0.891
B-#Breakdown	3	0.999	0.666	1.000	0.800
B-#Bus	8	1.000	1.000	1.000	1.000
B-#Car	23	0.998	0.956	0.956	0.956
B-#Event	5	0.994	0.400	0.200	0.266
B-#GoodTrafficSituation	16	0.995	0.625	0.833	0.714
B-#Interdiction	24	0.998	0.958	0.958	0.958
B-#Location	316	0.967	0.924	0.884	0.904
B-#Motorcycle	4	0.998	0.500	1.000	0.666
B-#Protest	1	0.998	1.000	0.333	0.500
B-#PublicAuthority	8	0.999	0.875	1.000	0.933
B-#RoadWork	1	1.000	1.000	1.000	1.000
B-#Solution	10	1.000	1.000	1.000	1.000
B-#Time	3	1.000	1.000	1.000	1.000
B-#Truck	6	0.999	1.000	0.857	0.923
B-#wayEffect:BothDirections	11	0.995	0.272	1.000	0.428
B-#wayEffect:OneDirection	1	0.999	0.000	0.000	0.000
B-#wayEffect:Partially	23	0.998	1.000	0.884	0.938
I-#BadTrafficSituation	20	0.998	0.900	0.947	0.923
I-#GoodTrafficSituation	15	0.995	0.400	1.000	0.571
I-#Location	90	0.968	0.377	0.918	0.535
I-#wayEffect:BothDirections	12	0.997	0.666	1.000	0.800
Average, σ		0.993, 0.013	0.767, 0.288	0.854, 0.266	0.776, 0.260

Tabela E.11: Best parameters: model-fscore-TEDO-NER-SVC-1

Parameter	Values	Best value
Kernel	Radial Basis Function, Linear, Polynomial, Sigmoidal	Linear

Tabela E.12: Training scores: model-fscore-TEDO-NER-SVC-1

Class	Support	Accuracy, σ CV	Recall, σ CV	Precision, σ CV	F1, σ CV
O	5262	0.948, 0.014	0.976, 0.008	0.946, 0.017	0.960, 0.010
B-#Accident	114	0.997, 0.002	0.862, 0.155	0.948, 0.081	0.898, 0.123
B-#BadTrafficSituation	273	0.995, 0.002	0.927, 0.054	0.930, 0.028	0.928, 0.034
B-#Breakdown	14	0.999, 0.000	0.718, 0.422	0.750, 0.433	0.732, 0.425
B-#Bus	31	1.000, 0.000	1.000, 0.000	1.000, 0.000	1.000, 0.000
B-#Car	103	0.999, 0.000	0.975, 0.039	0.941, 0.066	0.956, 0.042
B-#Event	32	0.996, 0.001	0.411, 0.324	0.603, 0.400	0.459, 0.327
B-#GoodTrafficSituation	29	0.997, 0.002	0.481, 0.331	0.800, 0.400	0.577, 0.335
B-#Interdiction	141	0.997, 0.002	0.890, 0.098	0.958, 0.054	0.920, 0.066
B-#Location	1270	0.976, 0.008	0.929, 0.029	0.919, 0.032	0.924, 0.025
B-#Motorcycle	18	0.999, 0.000	0.916, 0.166	1.000, 0.000	0.947, 0.108
B-#Protest	16	0.995, 0.003	0.419, 0.282	1.000, 0.000	0.537, 0.273
B-#PublicAuthority	55	0.998, 0.001	0.653, 0.374	0.757, 0.408	0.694, 0.381
B-#RoadWork	31	1.000, 0.000	0.888, 0.314	0.888, 0.314	0.888, 0.314
B-#Solution	45	0.998, 0.001	0.898, 0.165	0.975, 0.066	0.923, 0.106
B-#Time	17	0.998, 0.002	0.571, 0.494	0.542, 0.474	0.555, 0.482
B-#Truck	24	0.999, 0.000	0.650, 0.450	0.687, 0.451	0.659, 0.442
B-#wayEffect:BothDirections	19	0.999, 0.000	0.969, 0.062	1.000, 0.000	0.983, 0.033
B-#wayEffect:OneDirection	3	0.998, 0.000	0.500, 0.500	0.500, 0.500	0.500, 0.500
B-#wayEffect:Partially	123	0.997, 0.001	0.822, 0.281	0.802, 0.304	0.802, 0.282
I-#BadTrafficSituation	85	0.997, 0.001	0.869, 0.151	0.901, 0.143	0.876, 0.127
I-#Event	8	0.998, 0.000	0.386, 0.474	0.386, 0.474	0.386, 0.474
I-#GoodTrafficSituation	28	0.997, 0.001	0.539, 0.443	0.660, 0.439	0.555, 0.416
I-#Location	306	0.977, 0.010	0.542, 0.124	0.754, 0.209	0.611, 0.128
I-#Protest	2	0.998, 0.000	0.500, 0.000	1.000, 0.000	0.666, 0.000
I-#PublicAuthority	9	0.998, 0.001	0.666, 0.471	0.666, 0.471	0.666, 0.471
I-#Time	8	0.995, 0.001	0.000, 0.000	0.000, 0.000	0.000, 0.000
I-#wayEffect:BothDirections	22	0.999, 0.000	0.929, 0.112	0.951, 0.096	0.939, 0.101
I-#wayEffect:Partially	2	0.998, 0.000	0.500, 0.000	1.000, 0.000	0.666, 0.000
Average, σ		0.995, 0.010	0.703, 0.241	0.802, 0.226	0.731, 0.226

Tabela E.13: Revalidation scores: model-fscore-TEDO-NER-SVC-1

Class	Support	Accuracy	Recall	Precision	F1
O	1195	0.958	0.982	0.953	0.967
B-#Accident	27	0.997	0.925	0.925	0.925
B-#BadTrafficSituation	66	0.996	1.000	0.916	0.956
B-#Breakdown	3	0.999	0.666	1.000	0.800
B-#Bus	8	1.000	1.000	1.000	1.000
B-#Car	23	0.998	0.956	0.916	0.936
B-#Event	5	0.997	0.400	0.666	0.500
B-#GoodTrafficSituation	16	0.998	0.875	1.000	0.933
B-#Interdiction	24	1.000	1.000	1.000	1.000
B-#Location	316	0.976	0.911	0.944	0.927
B-#Motorcycle	4	0.998	0.750	0.750	0.750
B-#Protest	1	1.000	1.000	1.000	1.000
B-#PublicAuthority	8	0.999	1.000	0.888	0.941
B-#RoadWork	1	1.000	1.000	1.000	1.000
B-#Solution	10	0.999	0.900	1.000	0.947
B-#Time	3	1.000	1.000	1.000	1.000
B-#Truck	6	1.000	1.000	1.000	1.000
B-#wayEffect:BothDirections	11	1.000	1.000	1.000	1.000
B-#wayEffect:OneDirection	1	0.999	0.000	0.000	0.000
B-#wayEffect:Partially	23	0.998	1.000	0.920	0.958
I-#BadTrafficSituation	20	0.999	0.950	1.000	0.974
I-#GoodTrafficSituation	15	0.999	0.933	1.000	0.965
I-#Location	90	0.969	0.533	0.750	0.623
I-#wayEffect:BothDirections	12	1.000	1.000	1.000	1.000
Average, σ		0.995, 0.010	0.866, 0.238	0.901, 0.208	0.879, 0.221

Tabela E.14: Best parameters: model-fscore-TEDO-NER-RF-2

Parameter	Values	Best value
Class weight	balanced, None	None
Criterion	gini, entropy	gini
Max features	sqrt(n_features), log2(n_features), n_features	n_features
Number of trees in the forest	10, 50, 100	50
Use out-of-bag samples to estimate accuracy	true, false	true
One vs Rest	true, false	true

Tabela E.15: Training scores: model-fscore-TEDO-NER-RF-2

Class	Support	Accuracy, σ CV	Recall, σ CV	Precision, σ CV	F1, σ CV
O	5262	0.959, 0.017	0.976, 0.012	0.961, 0.017	0.968, 0.013
B-#Accident	114	0.998, 0.001	0.954, 0.062	0.945, 0.088	0.946, 0.058
B-#BadTrafficSituation	273	0.996, 0.001	0.955, 0.038	0.942, 0.051	0.947, 0.025
B-#Breakdown	14	1.000, 0.000	1.000, 0.000	1.000, 0.000	1.000, 0.000
B-#Bus	31	1.000, 0.000	1.000, 0.000	1.000, 0.000	1.000, 0.000
B-#Car	103	0.999, 0.000	0.990, 0.029	1.000, 0.000	0.994, 0.015
B-#Event	32	0.998, 0.001	0.670, 0.334	0.814, 0.318	0.701, 0.311
B-#GoodTrafficSituation	29	0.997, 0.002	0.336, 0.350	0.525, 0.453	0.378, 0.359
B-#Interdiction	141	0.999, 0.001	0.964, 0.060	0.994, 0.017	0.978, 0.036
B-#Location	1270	0.980, 0.007	0.965, 0.018	0.916, 0.030	0.940, 0.023
B-#Motorcycle	18	0.999, 0.000	0.944, 0.157	1.000, 0.000	0.962, 0.104
B-#Protest	16	0.998, 0.001	0.736, 0.289	1.000, 0.000	0.811, 0.221
B-#PublicAuthority	55	0.998, 0.001	0.795, 0.329	0.861, 0.314	0.812, 0.311
B-#RoadWork	31	0.999, 0.000	0.935, 0.122	1.000, 0.000	0.961, 0.072
B-#Solution	45	0.999, 0.000	0.875, 0.330	0.850, 0.327	0.861, 0.327
B-#Time	17	0.998, 0.001	0.571, 0.494	0.447, 0.438	0.484, 0.447
B-#Truck	24	1.000, 0.000	0.900, 0.300	0.900, 0.300	0.900, 0.300
B-#wayEffect:BothDirections	19	0.999, 0.000	0.983, 0.049	0.925, 0.160	0.943, 0.103
B-#wayEffect:OneDirection	3	0.998, 0.000	0.500, 0.500	0.250, 0.250	0.333, 0.333
B-#wayEffect:Partially	123	0.997, 0.001	0.791, 0.273	0.858, 0.292	0.820, 0.277
I-#BadTrafficSituation	85	0.997, 0.002	0.878, 0.162	0.928, 0.135	0.886, 0.121
I-#Event	8	0.998, 0.000	0.300, 0.400	0.400, 0.489	0.333, 0.421
I-#GoodTrafficSituation	28	0.997, 0.003	0.675, 0.419	0.590, 0.381	0.615, 0.380
I-#Location	306	0.975, 0.015	0.472, 0.149	0.768, 0.212	0.572, 0.155
I-#Protest	2	0.998, 0.000	0.500, 0.000	1.000, 0.000	0.666, 0.000
I-#PublicAuthority	9	0.997, 0.002	0.500, 0.408	0.666, 0.471	0.555, 0.415
I-#Time	8	0.995, 0.001	0.000, 0.000	0.000, 0.000	0.000, 0.000
I-#wayEffect:BothDirections	22	0.999, 0.000	0.769, 0.395	0.736, 0.383	0.747, 0.382
I-#wayEffect:Partially	2	0.997, 0.000	0.500, 0.000	0.500, 0.000	0.500, 0.000
Average, σ		0.995, 0.008	0.739, 0.255	0.785, 0.256	0.745, 0.253

Tabela E.16: Revalidation scores: model-fscore-TEDO-NER-RF-2

Class	Support	Accuracy	Recall	Precision	F1
O	1195	0.970	0.991	0.962	0.976
B-#Accident	27	1.000	1.000	1.000	1.000
B-#BadTrafficSituation	66	0.997	0.954	0.984	0.969
B-#Breakdown	3	0.999	0.666	1.000	0.800
B-#Bus	8	1.000	1.000	1.000	1.000
B-#Car	23	1.000	1.000	1.000	1.000
B-#Event	5	0.998	0.600	0.750	0.666
B-#GoodTrafficSituation	16	0.997	0.937	0.833	0.882
B-#Interdiction	24	1.000	1.000	1.000	1.000
B-#Location	316	0.979	0.939	0.936	0.938
B-#Motorcycle	4	1.000	1.000	1.000	1.000
B-#Protest	1	1.000	1.000	1.000	1.000
B-#PublicAuthority	8	1.000	1.000	1.000	1.000
B-#RoadWork	1	1.000	1.000	1.000	1.000
B-#Solution	10	1.000	1.000	1.000	1.000
B-#Time	3	1.000	1.000	1.000	1.000
B-#Truck	6	1.000	1.000	1.000	1.000
B-#wayEffect:BothDirections	11	1.000	1.000	1.000	1.000
B-#wayEffect:OneDirection	1	0.999	0.000	0.000	0.000
B-#wayEffect:Partially	23	0.998	0.956	0.916	0.936
I-#BadTrafficSituation	20	1.000	1.000	1.000	1.000
I-#GoodTrafficSituation	15	0.998	1.000	0.882	0.937
I-#Location	90	0.974	0.511	0.901	0.652
I-#wayEffect:BothDirections	12	1.000	1.000	1.000	1.000
Average, σ		0.996, 0.008	0.898, 0.230	0.923, 0.202	0.906, 0.213

Tabela E.17: Best parameters: model-fscore-TEDO-NER-SGD-2

Parameter	Values	Best value
Class weight	balanced, None	None
Learning rate	constant, optimal, invscaling	
Loss function	Hinge, Log, Modified huber, Epsilon insensitive, Squared epsilon insensitive, Squared hinge, Perceptron, Huber, Epsilon insensitive, Squared epsilon insensitive	Hinge
One vs Rest	true, false	true
Penalty	None, l1, l2, elasticnet	elasticnet

Tabela E.18: Training scores: model-fscore-TEDO-NER-SGD-2

Class	Support	Accuracy, σ CV	Recall, σ CV	Precision, σ CV	F1, σ CV
O	5262	0.941, 0.012	0.964, 0.021	0.946, 0.018	0.955, 0.009
B-#Accident	114	0.998, 0.001	0.902, 0.099	0.967, 0.060	0.929, 0.063
B-#BadTrafficSituation	273	0.994, 0.001	0.929, 0.070	0.897, 0.063	0.909, 0.035
B-#Breakdown	14	0.999, 0.000	0.843, 0.329	0.875, 0.330	0.857, 0.327
B-#Bus	31	1.000, 0.000	1.000, 0.000	1.000, 0.000	1.000, 0.000
B-#Car	103	0.998, 0.001	0.844, 0.287	0.871, 0.293	0.857, 0.289
B-#Event	32	0.995, 0.002	0.366, 0.343	0.486, 0.417	0.374, 0.309
B-#GoodTrafficSituation	29	0.997, 0.002	0.418, 0.358	0.633, 0.433	0.485, 0.369
B-#Interdiction	141	0.996, 0.002	0.879, 0.095	0.941, 0.082	0.905, 0.073
B-#Location	1270	0.969, 0.008	0.909, 0.054	0.871, 0.068	0.888, 0.053
B-#Motorcycle	18	0.999, 0.000	0.870, 0.176	0.901, 0.164	0.863, 0.137
B-#Protest	16	0.995, 0.002	0.344, 0.122	0.888, 0.157	0.468, 0.099
B-#PublicAuthority	55	0.996, 0.002	0.617, 0.407	0.650, 0.397	0.608, 0.375
B-#RoadWork	31	0.999, 0.001	0.807, 0.328	0.888, 0.314	0.837, 0.314
B-#Solution	45	0.997, 0.002	0.877, 0.130	0.850, 0.162	0.847, 0.111
B-#Time	17	0.997, 0.002	0.285, 0.451	0.250, 0.400	0.265, 0.421
B-#Truck	24	0.999, 0.001	0.673, 0.384	0.766, 0.395	0.709, 0.381
B-#wayEffect:BothDirections	19	0.999, 0.001	0.633, 0.458	0.700, 0.458	0.650, 0.450
B-#wayEffect:OneDirection	3	0.998, 0.000	0.500, 0.500	0.500, 0.500	0.500, 0.500
B-#wayEffect:Partially	123	0.996, 0.001	0.836, 0.289	0.860, 0.297	0.846, 0.290
I-#BadTrafficSituation	85	0.997, 0.001	0.807, 0.277	0.831, 0.282	0.818, 0.277
I-#Event	8	0.998, 0.000	0.575, 0.471	0.480, 0.410	0.518, 0.430
I-#GoodTrafficSituation	28	0.997, 0.002	0.570, 0.405	0.611, 0.419	0.586, 0.406
I-#Location	306	0.968, 0.014	0.444, 0.326	0.653, 0.376	0.484, 0.304
I-#Protest	2	0.997, 0.000	0.347, 0.000	0.533, 0.000	0.421, 0.000
I-#PublicAuthority	9	0.996, 0.001	0.444, 0.350	0.411, 0.340	0.383, 0.272
I-#Time	8	0.995, 0.001	0.269, 0.269	0.175, 0.175	0.212, 0.212
I-#wayEffect:BothDirections	22	0.999, 0.001	0.450, 0.471	0.444, 0.471	0.447, 0.471
I-#wayEffect:Partially	2	0.998, 0.000	1.000, 0.000	1.000, 0.000	1.000, 0.000
Average, σ		0.993, 0.012	0.669, 0.235	0.720, 0.225	0.676, 0.233

Tabela E.19: Revalidation scores: model-fscore-TEDO-NER-SGD-2

Class	Support	Accuracy	Recall	Precision	F1
O	1195	0.942	0.994	0.920	0.956
B-#Accident	27	0.998	0.925	1.000	0.961
B-#BadTrafficSituation	66	0.994	0.954	0.900	0.926
B-#Breakdown	3	0.999	0.666	1.000	0.800
B-#Bus	8	0.998	0.750	1.000	0.857
B-#Car	23	0.998	0.913	0.954	0.933
B-#Event	5	0.997	0.400	0.666	0.500
B-#GoodTrafficSituation	16	0.995	0.500	1.000	0.666
B-#Interdiction	24	0.999	0.958	1.000	0.978
B-#Location	316	0.970	0.889	0.933	0.910
B-#Motorcycle	4	0.998	0.750	0.750	0.750
B-#Protest	1	1.000	1.000	1.000	1.000
B-#PublicAuthority	8	0.998	0.625	1.000	0.769
B-#RoadWork	1	1.000	1.000	1.000	1.000
B-#Solution	10	1.000	1.000	1.000	1.000
B-#Time	3	1.000	1.000	1.000	1.000
B-#Truck	6	0.999	0.833	1.000	0.909
B-#wayEffect:BothDirections	11	1.000	1.000	1.000	1.000
B-#wayEffect:OneDirection	1	0.999	0.000	0.000	0.000
B-#wayEffect:Partially	23	0.999	1.000	0.958	0.978
I-#BadTrafficSituation	20	0.998	0.950	0.950	0.950
I-#GoodTrafficSituation	15	0.997	0.800	0.923	0.857
I-#Location	90	0.963	0.288	0.866	0.433
I-#wayEffect:BothDirections	12	0.998	0.750	1.000	0.857
Average, σ		0.993, 0.013	0.789, 0.257	0.909, 0.206	0.833, 0.229

Tabela E.20: Best parameters: model-fscore-TEDO-NER-SVC-2

Parameter	Values	Best value
Penalty parameter C	0.1, 1, 10, 100	1
Class weight	balanced, None	None
One vs Rest	true, false	true
Use Shrinking heuristic	true, false	true

Tabela E.21: Training scores: model-fscore-TEDO-NER-SVC-2

Class	Support	Accuracy, σ CV	Recall, σ CV	Precision, σ CV	F1, σ CV
O	5262	0.956, 0.013	0.975, 0.009	0.958, 0.016	0.966, 0.010
B-#Accident	114	0.999, 0.001	0.961, 0.049	0.975, 0.039	0.967, 0.038
B-#BadTrafficSituation	273	0.996, 0.001	0.961, 0.045	0.933, 0.035	0.946, 0.022
B-#Breakdown	14	0.999, 0.000	1.000, 0.000	0.975, 0.066	0.986, 0.036
B-#Bus	31	0.999, 0.000	1.000, 0.000	0.907, 0.107	0.948, 0.060
B-#Car	103	0.999, 0.000	0.884, 0.296	0.891, 0.298	0.887, 0.296
B-#Event	32	0.996, 0.001	0.490, 0.266	0.731, 0.313	0.570, 0.263
B-#GoodTrafficSituation	29	0.997, 0.002	0.459, 0.370	0.625, 0.436	0.509, 0.379
B-#Interdiction	141	0.998, 0.001	0.970, 0.047	0.968, 0.036	0.968, 0.028
B-#Location	1270	0.976, 0.009	0.943, 0.030	0.900, 0.038	0.921, 0.033
B-#Motorcycle	18	0.999, 0.000	0.938, 0.155	0.987, 0.035	0.953, 0.104
B-#Protest	16	0.998, 0.000	0.669, 0.239	1.000, 0.000	0.774, 0.194
B-#PublicAuthority	55	0.997, 0.002	0.758, 0.312	0.796, 0.331	0.758, 0.299
B-#RoadWork	31	1.000, 0.000	0.814, 0.355	0.888, 0.314	0.833, 0.333
B-#Solution	45	0.999, 0.000	0.950, 0.093	0.909, 0.118	0.925, 0.091
B-#Time	17	0.998, 0.001	0.523, 0.466	0.542, 0.474	0.526, 0.460
B-#Truck	24	0.999, 0.000	0.900, 0.300	0.783, 0.325	0.824, 0.303
B-#wayEffect:BothDirections	19	0.999, 0.001	0.900, 0.300	0.800, 0.331	0.833, 0.307
B-#wayEffect:OneDirection	3	0.998, 0.001	0.500, 0.500	0.500, 0.500	0.500, 0.500
B-#wayEffect:Partially	123	0.997, 0.001	0.847, 0.184	0.970, 0.047	0.889, 0.136
I-#BadTrafficSituation	85	0.998, 0.001	0.941, 0.082	0.926, 0.092	0.928, 0.063
I-#Event	8	0.998, 0.000	0.375, 0.460	0.400, 0.489	0.386, 0.474
I-#GoodTrafficSituation	28	0.997, 0.002	0.508, 0.402	0.685, 0.450	0.558, 0.401
I-#Location	306	0.977, 0.011	0.546, 0.260	0.674, 0.272	0.596, 0.262
I-#Protest	2	0.997, 0.000	0.347, 0.000	0.800, 0.000	0.484, 0.000
I-#PublicAuthority	9	0.997, 0.001	0.583, 0.311	0.777, 0.314	0.522, 0.109
I-#Time	8	0.995, 0.001	0.416, 0.006	0.540, 0.100	0.465, 0.034
I-#wayEffect:BothDirections	22	0.999, 0.000	0.706, 0.397	0.694, 0.396	0.687, 0.383
I-#wayEffect:Partially	2	0.998, 0.000	0.500, 0.000	1.000, 0.000	0.666, 0.000
Average, σ		0.995, 0.009	0.737, 0.220	0.811, 0.165	0.751, 0.192

Tabela E.22: Revalidation scores: model-fscore-TEDO-NER-SVC-2

Class	Support	Accuracy	Recall	Precision	F1
O	1195	0.958	0.984	0.952	0.967
B-#Accident	27	0.998	0.925	1.000	0.961
B-#BadTrafficSituation	66	0.996	0.969	0.941	0.955
B-#Breakdown	3	0.999	0.666	1.000	0.800
B-#Bus	8	1.000	1.000	1.000	1.000
B-#Car	23	1.000	1.000	1.000	1.000
B-#Event	5	0.997	0.600	0.600	0.600
B-#GoodTrafficSituation	16	0.996	0.687	0.916	0.785
B-#Interdiction	24	1.000	1.000	1.000	1.000
B-#Location	316	0.976	0.927	0.933	0.930
B-#Motorcycle	4	0.999	1.000	0.800	0.888
B-#Protest	1	1.000	1.000	1.000	1.000
B-#PublicAuthority	8	1.000	1.000	1.000	1.000
B-#RoadWork	1	1.000	1.000	1.000	1.000
B-#Solution	10	1.000	1.000	1.000	1.000
B-#Time	3	1.000	1.000	1.000	1.000
B-#Truck	6	1.000	1.000	1.000	1.000
B-#wayEffect:BothDirections	11	0.998	0.818	1.000	0.900
B-#wayEffect:OneDirection	1	0.999	0.000	0.000	0.000
B-#wayEffect:Partially	23	0.998	1.000	0.884	0.938
I-#BadTrafficSituation	20	1.000	1.000	1.000	1.000
I-#GoodTrafficSituation	15	0.998	1.000	0.882	0.937
I-#Location	90	0.968	0.477	0.767	0.589
I-#wayEffect:BothDirections	12	0.998	0.833	1.000	0.909
Average, σ		0.995, 0.010	0.870, 0.233	0.903, 0.211	0.881, 0.216

Tabela E.23: Training scores: model-fscore-TEDO-NER-RF-FINAL

Class	Support	Accuracy, σ CV	Recall, σ CV	Precision, σ CV	F1, σ CV
O	6457	0.964, 0.016	0.980, 0.010	0.964, 0.019	0.972, 0.012
B-#Accident	141	0.999, 0.000	0.962, 0.053	0.991, 0.025	0.975, 0.028
B-#BadTrafficSituation	339	0.996, 0.002	0.948, 0.046	0.944, 0.044	0.945, 0.032
B-#Breakdown	17	0.999, 0.000	0.962, 0.104	0.944, 0.157	0.940, 0.115
B-#Bus	39	1.000, 0.000	1.000, 0.000	1.000, 0.000	1.000, 0.000
B-#Car	126	0.999, 0.000	0.983, 0.049	1.000, 0.000	0.990, 0.027
B-#Event	37	0.998, 0.000	0.667, 0.204	0.925, 0.114	0.749, 0.131
B-#GoodTrafficSituation	45	0.997, 0.002	0.546, 0.307	0.625, 0.370	0.560, 0.315
B-#Interdiction	165	0.999, 0.001	0.974, 0.043	0.995, 0.014	0.984, 0.026
B-#Location	1586	0.980, 0.007	0.960, 0.017	0.920, 0.036	0.939, 0.024
B-#Motorcycle	22	0.999, 0.000	0.944, 0.157	1.000, 0.000	0.962, 0.104
B-#Protest	17	0.999, 0.000	0.802, 0.275	1.000, 0.000	0.858, 0.208
B-#PublicAuthority	63	0.998, 0.001	0.854, 0.212	0.870, 0.204	0.839, 0.183
B-#RoadWork	32	0.999, 0.000	0.925, 0.138	1.000, 0.000	0.955, 0.083
B-#Solution	55	0.999, 0.000	0.916, 0.220	1.000, 0.000	0.937, 0.165
B-#Time	20	0.998, 0.001	0.600, 0.447	0.604, 0.447	0.586, 0.426
B-#Truck	30	1.000, 0.000	1.000, 0.000	1.000, 0.000	1.000, 0.000
B-#wayEffect:BothDirections	30	0.999, 0.000	0.980, 0.059	0.901, 0.153	0.929, 0.092
B-#wayEffect:OneDirection	4	0.998, 0.000	0.000, 0.000	0.000, 0.000	0.000, 0.000
B-#wayEffect:Partially	146	0.998, 0.001	0.840, 0.285	0.846, 0.299	0.838, 0.284
I-#BadTrafficSituation	105	0.997, 0.002	0.872, 0.167	0.904, 0.109	0.878, 0.124
I-#Event	8	0.998, 0.000	0.200, 0.400	0.200, 0.400	0.200, 0.400
I-#GoodTrafficSituation	43	0.997, 0.002	0.732, 0.329	0.825, 0.290	0.750, 0.288
I-#Location	396	0.976, 0.013	0.617, 0.206	0.754, 0.161	0.657, 0.141
I-#Protest	2	0.997, 0.000	0.500, 0.000	0.500, 0.000	0.500, 0.000
I-#PublicAuthority	9	0.998, 0.001	0.500, 0.408	0.666, 0.471	0.555, 0.415
I-#Time	8	0.995, 0.001	0.270, 0.270	0.450, 0.450	0.338, 0.338
I-#wayEffect:BothDirections	34	0.999, 0.000	0.800, 0.400	0.713, 0.379	0.748, 0.382
I-#wayEffect:Partially	2	0.998, 0.000	0.500, 0.000	1.000, 0.000	0.666, 0.000
Average, σ		0.996, 0.007	0.753, 0.263	0.812, 0.250	0.767, 0.252

Tabela E.24: Test scores: model-fscore-TEDO-NER-RF-FINAL

Class	Support	Accuracy	Recall	Precision	F1
O	1553	0.962	0.981	0.960	0.971
B-#Accident	36	0.999	0.972	1.000	0.985
B-#BadTrafficSituation	72	0.997	0.986	0.946	0.965
B-#Breakdown	5	1.000	1.000	1.000	1.000
B-#Bus	11	1.000	1.000	1.000	1.000
B-#Car	28	0.999	0.964	1.000	0.981
B-#Event	5	0.999	0.800	0.800	0.800
B-#GoodTrafficSituation	15	0.998	0.800	0.923	0.857
B-#Interdiction	43	0.999	0.976	1.000	0.988
B-#Location	390	0.976	0.938	0.917	0.927
B-#Motorcycle	7	1.000	1.000	1.000	1.000
B-#Protest	9	1.000	1.000	1.000	1.000
B-#PublicAuthority	18	0.997	0.777	0.933	0.848
B-#RoadWork	12	0.999	0.916	1.000	0.956
B-#Solution	20	0.997	0.800	0.941	0.864
B-#Time	4	0.999	0.750	1.000	0.857
B-#Truck	7	0.999	0.857	1.000	0.923
B-#wayEffect:BothDirections	8	0.999	1.000	0.888	0.941
B-#wayEffect:Partially	42	0.996	0.000	0.000	0.000
I-#BadTrafficSituation	24	0.999	0.857	0.947	0.900
I-#GoodTrafficSituation	13	0.998	0.958	0.958	0.958
I-#Location	115	0.976	0.846	0.916	0.879
I-#Time	5	0.997	0.626	0.827	0.712
I-#wayEffect:BothDirections	10	1.000	0.000	0.000	0.000
I-#wayEffect:Partially	2	0.998	0.200	0.250	0.222
Average, σ		0.995, 0.009	0.800, 0.289	0.848, 0.290	0.821, 0.287

Tabela E.25: Training scores: model-fscore-TEDO-NER-RF-FINAL-AGRUPADO

Class	Support	Accuracy, σ CV	Recall, σ CV	Precision, σ CV	F1, σ CV
O	6605	0.965, 0.016	0.981, 0.009	0.966, 0.019	0.973, 0.012
B-#Accident	141	0.999, 0.000	0.962, 0.053	0.984, 0.046	0.972, 0.033
B-#Actor	217	0.999, 0.000	0.983, 0.035	1.000, 0.000	0.991, 0.018
B-#BadTrafficSituation	339	0.995, 0.002	0.942, 0.053	0.937, 0.039	0.938, 0.031
B-#Breakdown	17	0.999, 0.000	0.962, 0.104	1.000, 0.000	0.977, 0.062
B-#Event	37	0.998, 0.001	0.636, 0.161	0.925, 0.114	0.733, 0.086
B-#GoodTrafficSituation	45	0.996, 0.002	0.531, 0.315	0.556, 0.351	0.517, 0.314
B-#Interdiction	165	0.999, 0.001	0.974, 0.043	0.995, 0.014	0.984, 0.026
B-#Location	1586	0.980, 0.007	0.957, 0.019	0.922, 0.033	0.939, 0.024
B-#Protest	17	0.999, 0.000	0.833, 0.288	1.000, 0.000	0.875, 0.216
B-#PublicAuthority	63	0.998, 0.001	0.854, 0.212	0.925, 0.159	0.857, 0.154
B-#RoadWork	32	0.999, 0.000	0.925, 0.138	1.000, 0.000	0.955, 0.083
B-#Solution	55	0.999, 0.000	0.916, 0.220	0.987, 0.033	0.930, 0.163
B-#Time	20	0.998, 0.001	0.433, 0.453	0.500, 0.500	0.458, 0.465
B-#wayEffect:BothDirections	30	0.999, 0.000	0.980, 0.059	0.889, 0.169	0.920, 0.102
B-#wayEffect:OneDirection	4	0.998, 0.000	0.000, 0.000	0.000, 0.000	0.000, 0.000
I-#BadTrafficSituation	105	0.997, 0.003	0.785, 0.303	0.822, 0.290	0.798, 0.290
I-#Event	8	0.998, 0.000	0.200, 0.400	0.133, 0.266	0.160, 0.320
I-#GoodTrafficSituation	43	0.997, 0.002	0.757, 0.339	0.851, 0.300	0.774, 0.299
I-#Location	396	0.976, 0.013	0.618, 0.202	0.777, 0.168	0.665, 0.138
I-#Protest	2	0.997, 0.000	0.500, 0.000	0.500, 0.000	0.500, 0.000
I-#PublicAuthority	9	0.998, 0.001	0.575, 0.422	0.538, 0.411	0.555, 0.415
I-#Time	8	0.995, 0.001	0.270, 0.270	0.460, 0.460	0.340, 0.340
I-#wayEffect:BothDirections	34	0.999, 0.000	0.900, 0.300	0.747, 0.295	0.808, 0.286
Average, σ		0.995, 0.008	0.728, 0.276	0.767, 0.278	0.734, 0.273

Tabela E.26: Test scores: model-fscore-TEDO-NER-RF-FINAL-AGRUPADO

Class	Support	Accuracy	Recall	Precision	F1
O	1597	0.962	0.979	0.964	0.971
B-#Accident	36	1.000	1.000	1.000	1.000
B-#Actor	53	0.999	0.981	1.000	0.990
B-#BadTrafficSituation	72	0.997	0.986	0.946	0.965
B-#Breakdown	5	1.000	1.000	1.000	1.000
B-#Event	5	0.999	0.800	1.000	0.888
B-#GoodTrafficSituation	15	0.997	0.800	0.857	0.827
B-#Interdiction	43	0.999	0.976	1.000	0.988
B-#Location	390	0.978	0.946	0.922	0.934
B-#Protest	9	1.000	1.000	1.000	1.000
B-#PublicAuthority	18	0.997	0.777	0.875	0.823
B-#RoadWork	12	0.999	0.916	1.000	0.956
B-#Solution	20	0.997	0.800	0.941	0.864
B-#Time	4	0.998	0.750	0.600	0.666
B-#wayEffect:BothDirections	8	0.999	1.000	0.888	0.941
I-#BadTrafficSituation	24	0.999	0.000	0.000	0.000
I-#GoodTrafficSituation	13	0.998	0.958	0.958	0.958
I-#Location	115	0.975	0.846	0.916	0.879
I-#Time	5	0.996	0.608	0.833	0.703
I-#wayEffect:BothDirections	10	1.000	0.000	0.000	0.000
Average, σ		0.995, 0.009	0.806, 0.289	0.835, 0.292	0.818, 0.287

Tabela E.27: Training scores: model-fscore-TEDO-RE-SP-0

Class	Support	Accuracy, σ CV	Recall, σ CV	Precision, σ CV	F1, σ CV
O	11274	0.885, 0.013	0.957, 0.019	0.911, 0.017	0.933, 0.008
#affectsFlowTo	313	0.984, 0.005	0.603, 0.171	0.741, 0.063	0.647, 0.110
#causes	218	0.991, 0.002	0.692, 0.144	0.801, 0.149	0.717, 0.084
#hasActor	186	0.997, 0.002	0.975, 0.049	0.908, 0.094	0.936, 0.055
#hasAlternative	16	0.998, 0.000	0.000, 0.000	0.000, 0.000	0.000, 0.000
#hasEdge	226	0.984, 0.006	0.234, 0.202	0.441, 0.324	0.295, 0.240
#hasEvent	710	0.953, 0.008	0.347, 0.147	0.609, 0.168	0.428, 0.133
#hasReference	307	0.971, 0.009	0.234, 0.198	0.366, 0.254	0.230, 0.151
#hasSupporter	96	0.998, 0.001	0.809, 0.245	0.807, 0.280	0.773, 0.251
#hasTime	18	0.997, 0.002	0.545, 0.379	0.678, 0.398	0.565, 0.366
Average, σ		0.976, 0.033	0.539, 0.312	0.626, 0.270	0.552, 0.294

Tabela E.28: Revalidation scores: model-fscore-TEDO-RE-SP-0

Class	Support	Accuracy	Recall	Precision	F1
O	2523	0.888	0.982	0.894	0.936
#affectsFlowTo	77	0.987	0.558	0.934	0.699
#causes	43	0.991	0.441	0.904	0.593
#hasActor	44	0.996	0.886	0.906	0.896
#hasAlternative	2	0.999	0.500	0.500	0.500
#hasEdge	53	0.983	0.150	0.666	0.246
#hasEvent	154	0.957	0.292	0.703	0.412
#hasReference	83	0.969	0.084	0.304	0.132
#hasSupporter	16	0.999	0.875	1.000	0.933
#hasTime	3	1.000	1.000	1.000	1.000
Average, σ		0.977, 0.032	0.577, 0.325	0.781, 0.221	0.635, 0.293

Tabela E.29: Training scores: model-fscore-TEDO-RE-FWSSVM-0

Class	Support	Accuracy, σ CV	Recall, σ CV	Precision, σ CV	F1, σ CV
O	11274	0.905, 0.010	0.968, 0.007	0.922, 0.009	0.944, 0.006
#affectsFlowTo	313	0.987, 0.004	0.679, 0.108	0.805, 0.089	0.729, 0.076
#causes	218	0.992, 0.004	0.762, 0.126	0.809, 0.119	0.774, 0.093
#hasActor	186	0.998, 0.001	0.970, 0.059	0.928, 0.079	0.945, 0.041
#hasAlternative	16	0.998, 0.000	0.000, 0.000	0.000, 0.000	0.000, 0.000
#hasEdge	226	0.983, 0.008	0.423, 0.227	0.548, 0.201	0.451, 0.200
#hasEvent	710	0.962, 0.006	0.570, 0.088	0.710, 0.148	0.620, 0.077
#hasReference	307	0.977, 0.008	0.000, 0.000	0.000, 0.000	0.000, 0.000
#hasSupporter	96	0.998, 0.000	0.950, 0.124	0.872, 0.203	0.885, 0.154
#hasTime	18	0.999, 0.001	0.764, 0.306	1.000, 0.000	0.825, 0.237
Average, σ		0.980, 0.027	0.609, 0.347	0.659, 0.350	0.617, 0.340

Tabela E.30: Revalidation scores: model-fscore-TEDO-RE-FWSSVM-0

Class	Support	Accuracy	Recall	Precision	F1
O	2523	0.912	0.969	0.929	0.949
#affectsFlowTo	77	0.988	0.610	0.921	0.734
#causes	43	0.995	0.790	0.894	0.839
#hasActor	44	0.997	0.931	0.911	0.921
#hasAlternative	2	0.999	0.000	0.000	0.000
#hasEdge	53	0.986	0.660	0.614	0.636
#hasEvent	154	0.968	0.694	0.690	0.692
#hasReference	83	0.972	0.000	0.000	0.000
#hasSupporter	16	0.999	0.937	1.000	0.967
#hasTime	3	1.000	1.000	1.000	1.000
Average, σ		0.982, 0.025	0.659, 0.354	0.696, 0.367	0.674, 0.356

Tabela E.31: Best parameters: model-fscore-TEDO-RE-SP-1

Parâmetro	Valores testados	Melhor valor
Average	true, false	true
Batch learning	true, false	false

Tabela E.32: Training scores: model-fscore-TEDO-RE-SP-1

Class	Support	Accuracy, σ CV	Recall, σ CV	Precision, σ CV	F1, σ CV
O	11274	0.899, 0.009	0.961, 0.008	0.922, 0.013	0.941, 0.005
#affectsFlowTo	313	0.987, 0.003	0.700, 0.105	0.762, 0.070	0.726, 0.074
#causes	218	0.991, 0.003	0.762, 0.104	0.778, 0.122	0.763, 0.096
#hasActor	186	0.998, 0.001	0.979, 0.042	0.922, 0.086	0.947, 0.049
#hasAlternative	16	0.998, 0.000	0.047, 0.116	0.142, 0.349	0.071, 0.174
#hasEdge	226	0.984, 0.006	0.373, 0.185	0.600, 0.252	0.438, 0.183
#hasEvent	710	0.959, 0.006	0.549, 0.098	0.658, 0.089	0.589, 0.069
#hasReference	307	0.976, 0.009	0.138, 0.070	0.406, 0.252	0.198, 0.105
#hasSupporter	96	0.998, 0.000	0.912, 0.108	0.902, 0.151	0.892, 0.098
#hasTime	18	0.998, 0.002	0.735, 0.328	0.919, 0.160	0.780, 0.271
Average, σ		0.979, 0.029	0.616, 0.315	0.701, 0.246	0.634, 0.291

Tabela E.33: Revalidation scores: model-fscore-TEDO-RE-SP-1

Class	Support	Accuracy	Recall	Precision	F1
O	2523	0.903	0.964	0.924	0.944
#affectsFlowTo	77	0.987	0.701	0.794	0.744
#causes	43	0.995	0.790	0.894	0.839
#hasActor	44	0.997	0.954	0.893	0.923
#hasAlternative	2	0.998	0.000	0.000	0.000
#hasEdge	53	0.983	0.396	0.552	0.461
#hasEvent	154	0.966	0.610	0.706	0.655
#hasReference	83	0.968	0.072	0.272	0.114
#hasSupporter	16	0.999	0.937	1.000	0.967
#hasTime	3	1.000	1.000	1.000	1.000
Average, σ		0.980, 0.028	0.642, 0.352	0.703, 0.317	0.665, 0.342

Tabela E.34: Best parameters: model-fscore-TEDO-RE-FWSSVM-1

Parâmetro	Valores testados	Melhor valor
Penalty parameter C	0.1, 1, 10, 100	100
Use weight averaging	true, false	false

Tabela E.35: Training scores: model-fscore-TEDO-RE-FWSSVM-1

Class	Support	Accuracy, σ CV	Recall, σ CV	Precision, σ CV	F1, σ CV
O	11274	0.885, 0.015	0.929, 0.036	0.935, 0.023	0.931, 0.010
#affectsFlowTo	313	0.986, 0.007	0.769, 0.108	0.743, 0.121	0.746, 0.086
#causes	218	0.991, 0.004	0.762, 0.136	0.768, 0.111	0.751, 0.084
#hasActor	186	0.998, 0.001	0.983, 0.038	0.952, 0.070	0.965, 0.040
#hasAlternative	16	0.997, 0.001	0.226, 0.341	0.309, 0.440	0.242, 0.353
#hasEdge	226	0.980, 0.009	0.478, 0.273	0.411, 0.202	0.411, 0.222
#hasEvent	710	0.957, 0.006	0.649, 0.206	0.647, 0.162	0.598, 0.092
#hasReference	307	0.968, 0.016	0.218, 0.213	0.303, 0.316	0.199, 0.196
#hasSupporter	96	0.998, 0.001	0.962, 0.089	0.874, 0.145	0.905, 0.100
#hasTime	18	0.999, 0.001	0.792, 0.293	0.966, 0.066	0.833, 0.234
Average, σ		0.976, 0.033	0.677, 0.268	0.691, 0.249	0.658, 0.269

Tabela E.36: Revalidation scores: model-fscore-TEDO-RE-FWSSVM-1

Class	Support	Accuracy	Recall	Precision	F1
O	2523	0.815	0.799	0.976	0.879
#affectsFlowTo	77	0.982	0.870	0.614	0.720
#causes	43	0.994	0.813	0.833	0.823
#hasActor	44	0.997	0.954	0.893	0.923
#hasAlternative	2	0.998	1.000	0.400	0.571
#hasEdge	53	0.962	0.886	0.305	0.454
#hasEvent	154	0.938	0.974	0.455	0.621
#hasReference	83	0.927	0.554	0.202	0.296
#hasSupporter	16	0.999	0.937	1.000	0.967
#hasTime	3	1.000	1.000	1.000	1.000
Average, σ		0.961, 0.055	0.879, 0.127	0.668, 0.293	0.725, 0.223

Tabela E.37: Training scores: model-fscore-TEDO-RE-FWSSVM-FINAL

Class	Support	Accuracy, σ CV	Recall, σ CV	Precision, σ CV	F1, σ CV
O	13797	0.885, 0.019	0.917, 0.033	0.944, 0.017	0.930, 0.013
#affectsFlowTo	390	0.987, 0.003	0.826, 0.068	0.714, 0.069	0.763, 0.051
#causes	261	0.992, 0.003	0.804, 0.144	0.778, 0.072	0.782, 0.086
#hasActor	230	0.998, 0.001	0.965, 0.043	0.956, 0.078	0.958, 0.044
#hasAlternative	18	0.998, 0.001	0.166, 0.267	0.214, 0.364	0.185, 0.304
#hasEdge	279	0.979, 0.012	0.522, 0.325	0.438, 0.290	0.420, 0.240
#hasEvent	864	0.957, 0.008	0.707, 0.186	0.598, 0.104	0.626, 0.080
#hasReference	390	0.966, 0.012	0.314, 0.175	0.289, 0.198	0.292, 0.174
#hasSupporter	112	0.998, 0.001	0.923, 0.136	0.934, 0.074	0.919, 0.079
#hasTime	21	0.999, 0.000	0.900, 0.223	0.958, 0.093	0.904, 0.157
Average, σ		0.976, 0.033	0.704, 0.264	0.682, 0.271	0.678, 0.270

Tabela E.38: Test scores: model-fscore-TEDO-RE-FWSSVM-FINAL

Class	Support	Accuracy	Recall	Precision	F1
O	3442	0.913	0.966	0.932	0.949
#affectsFlowTo	99	0.990	0.747	0.860	0.799
#causes	71	0.992	0.690	0.830	0.753
#hasActor	57	0.997	0.859	0.942	0.899
#hasAlternative	8	0.997	0.125	0.200	0.153
#hasEdge	52	0.990	0.615	0.640	0.627
#hasEvent	226	0.966	0.623	0.726	0.671
#hasReference	100	0.976	0.240	0.558	0.335
#hasSupporter	33	0.999	0.939	1.000	0.968
#hasTime	4	0.999	0.750	1.000	0.857
Average, σ		0.982, 0.025	0.655, 0.263	0.769, 0.237	0.701, 0.254

Tabela E.39: Training scores: model-fscore-TEDO-RE-FWSSVM-FINAL-AGRUPADO

Class	Support	Accuracy, σ CV	Recall, σ CV	Precision, σ CV	F1, σ CV
O	11817	0.863, 0.048	0.906, 0.081	0.927, 0.022	0.913, 0.037
#affectsFlowTo	390	0.986, 0.005	0.745, 0.107	0.790, 0.113	0.755, 0.062
#causes	261	0.990, 0.004	0.784, 0.146	0.771, 0.083	0.766, 0.084
#hasActor	230	0.998, 0.001	0.960, 0.051	0.965, 0.064	0.960, 0.043
#hasAlternative	18	0.997, 0.001	0.095, 0.233	0.071, 0.174	0.081, 0.199
#hasEdge	279	0.975, 0.011	0.446, 0.300	0.407, 0.269	0.358, 0.196
#hasEvent	864	0.946, 0.019	0.691, 0.180	0.614, 0.146	0.615, 0.078
#hasReference	390	0.961, 0.026	0.226, 0.198	0.305, 0.238	0.210, 0.170
#hasSupporter	112	0.999, 0.000	0.971, 0.057	0.908, 0.101	0.933, 0.059
#hasTime	21	0.999, 0.000	0.858, 0.224	0.958, 0.093	0.880, 0.152
Average, σ		0.971, 0.039	0.668, 0.293	0.672, 0.296	0.647, 0.304

Tabela E.40: Test scores: model-fscore-TEDO-RE-FWSSVM-FINAL-AGRUPADO

Class	Support	Accuracy	Recall	Precision	F1
O	2884	0.889	0.963	0.906	0.934
#affectsFlowTo	99	0.988	0.666	0.891	0.763
#causes	71	0.991	0.661	0.870	0.752
#hasActor	57	0.997	0.912	0.945	0.928
#hasAlternative	8	0.996	0.125	0.125	0.125
#hasEdge	52	0.983	0.673	0.454	0.542
#hasEvent	226	0.957	0.473	0.781	0.589
#hasReference	100	0.970	0.130	0.448	0.201
#hasSupporter	33	0.998	0.878	1.000	0.935
#hasTime	4	0.999	0.750	1.000	0.857
Average, σ		0.977, 0.032	0.623, 0.282	0.742, 0.281	0.662, 0.282