

4

Fast rank and subspace tracking

In this chapter, we introduce FRAHST (*Fast Rank-Adaptive recursive row-Householder Subspace Tracking*), our extension that adapts both projection matrix and subspace rank on the fly. We provide details on how the original complexity of $\mathcal{O}(Nr + r^3)$ can be reduced to $\mathcal{O}(Nr + r^2)$. Subsequently, we show how the method is used for anomaly detection and how it can be integrated in a real-time monitoring solution.

4.1

Tracking the principal subspace basis

Strobach's algorithm presented in Section 3.3.4 is the state-of-the-art in the low complexity class subspace trackers but currently no adaptation has been put forward to cope with the rank tracking problem. Throughout our research, we experimented with several alternatives (such as YAST from Badeau et al. [2008] and FDP from Doukopoulos and Moustakides [2008]) and found this particular method to have the best properties: provides excellent subspace estimates very efficiently in terms of computational and memory requirements, and has no special parameters to tune. The algorithm is an implementation of the orthogonal principle where the Q and S matrices (3-4) are directly updated through a recursive Householder reduction. During execution, the algorithm reads in new data snapshots $\mathbf{z}(t) \in \mathbb{R}^N$ from N streams at time t and tracks the principal subspace projection matrix \mathbf{Q} with dominant complexity of $3Nr$ flops, which is the lower bound for this class of algorithms where r is the *known and constant* subspace dimension.

We adopt the algorithm as listed in Figure 3.3 using the approximation of $\psi = 0$, consequently (3-23c) can be omitted and (3-23d) simplified. We can save these few computation steps because we are only interested in subspace tracking and these simplifications are barely noticeable in practice (as seen in Figure 5.3). This corresponds to the *a posteriori* projection approximation criterion described in Section 3.3.3. More concretely, we use the Householder reduction to recursively update the tracked subspace basis $\mathbf{Q}(t)$ of dimension $N \times r$ considering the new data point:

$$\begin{bmatrix} \mathbf{S}(t) \\ 0 \dots 0 \end{bmatrix} = \mathbf{H}(t) \begin{bmatrix} \alpha \mathbf{S}(t-1) + \mathbf{h}(t)\mathbf{h}^\top(t) \\ Z^{1/2}(t)\mathbf{h}^\top(t) \end{bmatrix} \quad (4-1)$$

$$\begin{bmatrix} \mathbf{Q}(t) & \mathbf{z}_q(t) \end{bmatrix} = \begin{bmatrix} \mathbf{Q}(t-1) & \bar{\mathbf{z}}_\perp(t) \end{bmatrix} \mathbf{H}(t) \quad (4-2)$$

where $\bar{\mathbf{z}}_\perp(t)$ is the unit-norm version of the complement of the orthogonal projection of $\mathbf{z}(t)$ onto the *principal* subspace spanned by the column vectors of $\mathbf{Q}(t-1)$ and $Z^{1/2}$ is the norm of this complement. The S -matrix is strongly dominated by its diagonal elements, whose values approximate the principal eigenvalues of the covariance, so we can interpret the transformation as a compressor that keeps most of the energy in the principal eigenvalues given the new input information. This is not an arbitrary heuristic, but follows from the orthogonal iteration principal described in Section 3.3.2.

4.2

Tracking the principal subspace rank

Most subspace tracking algorithms have the dimensionality r of the principal subspace given as a parameter. However, an algorithm that can automatically infer the latent dimensions is widely applicable, especially for monitoring changes in the underlying data stream.

We propose to adapt r , so that we maintain a high percentage f_E of the energy $E(t)$. Energy thresholding is a common method to estimate the number of principal components [Jolliffe, 2002] and is also adopted in SPIRIT [Papadimitriou et al., 2005, Papadimitriou, 2005, Sun, 2007], which has already shown to be useful for detecting unusual patterns in data center measurements [Hoke et al., 2006]. The term energy is common in the signal processing literature and $E(t)$ is defined as

$$E(t) = \sum_{t=1}^T \|\mathbf{z}(t)\|^2 = \sum_{t=1}^T \sum_{i=1}^N z_i^2(t). \quad (4-3)$$

We define T as the total number of intervals read so far. We let $\tilde{\mathbf{z}}(t)$ be the reconstruction of $\mathbf{z}(t)$ based on the previously learnt projection matrix, defined as follows:

$$\tilde{\mathbf{z}}(t) = \mathbf{Q}(t-1)\mathbf{Q}^T(t-1)\mathbf{z}(t) = \mathbf{Q}(t-1)\mathbf{h}(t). \quad (4-4)$$

Similarly, we have the energy $\tilde{E}(t)$ of the reconstruction $\tilde{\mathbf{z}}$ defined as

$$\tilde{E}(t) = \sum_{t=1}^T \|\tilde{\mathbf{z}}(t)\|^2. \quad (4-5)$$

Lemma 3 *Assuming the basis vectors \mathbf{q}_i , $1 \leq i \leq r$ are orthonormal, we have*

$$\tilde{E}(t) = \frac{1}{T} \sum_{t=1}^T \|\mathbf{h}(t)\|^2 = \frac{T-1}{T} \tilde{E}(t-1) + \frac{1}{T} \|\mathbf{h}(t)\|^2.$$

Proof: It is straightforward from applying Pythagorean theorem and using the orthonormality of the vectors \mathbf{q}_i , $1 \leq i \leq r$

$$\begin{aligned} \|\tilde{\mathbf{z}}(t)\|^2 &= \|h_1(t)\mathbf{q}_1(t-1) + \dots + h_r(t)\mathbf{q}_r(t-1)\|^2 \\ &= h_1^2(t)\|\mathbf{q}_1(t-1)\|^2 + \dots + h_r^2(t)\|\mathbf{q}_r(t-1)\|^2 \\ &= h_1^2(t) + \dots + h_r^2(t) = \|\mathbf{h}(t)\|^2 \end{aligned}$$

The result follows from summing over T . \square

This underlies the importance of a subspace tracker that guarantees orthonormal estimates for the projection matrix. For instance, SPIRIT requires a re-orthonormalization step for Lemma 3 to hold, and this amounts to an extra of $\mathcal{O}(Nr^2)$ flops¹. Our algorithm does not need this extra step, and updates the Q -matrix maintaining its orthonormality by construction.

In terms of parameters, we have a low-energy f_E and a high-energy F_E thresholds, and we keep enough number of latent variables r so that the retained energy is within the range $[f_E E(t), F_E E(t)]$. Whenever it gets outside these bounds, we increase or decrease by one unit accordingly. The main steps for the online rank estimation are:

- Estimate the high-dimensional data energy $E(t)$ incrementally from the sum of squares of \mathbf{z}_i , for all $1 \leq i \leq N$.
- Estimate the energy $\tilde{E}(t)$ of the current $r(t)$ latent variables.
- Adapt $r(t+1)$ if necessary. We add a new basis if the current latent variables capture too little energy, i.e., $\tilde{E}(t) < f_E E(t)$, or drop one dimension if the maintained energy is too high, i.e., $\tilde{E}(t) > F_E E(t)$.

The following lemma proves that the above steps guarantee the relative reconstruction error is within the specified interval $[f_E E(t), F_E E(t)]$.

Lemma 4 *The relative squared error of the reconstruction satisfies*

$$1 - F_E \leq \frac{\sum_{t=1}^T \|\mathbf{z}(t) - \tilde{\mathbf{z}}(t)\|^2}{\sum_{t=1}^T \|\mathbf{z}(t)\|^2} \leq 1 - f_E$$

¹This computational complexity is necessary for all orthonormalization procedures such as the classical and stabilized Gram-Schmidt as well as QR decomposition algorithms.

Proof: The measure $\mathbf{z}_\perp(t) = \mathbf{z}(t) - \tilde{\mathbf{z}}(t)$ has been mentioned before in Section 3.3.4 where $\mathbf{z}(t) \cdot \mathbf{z}_\perp(t) = 0$. Therefore, we have

$$\begin{aligned} \|\mathbf{z}(t) - \tilde{\mathbf{z}}(t)\|^2 &= \tilde{\mathbf{z}}^\top(t)[\mathbf{z}(t) - \tilde{\mathbf{z}}(t)] - \underbrace{\mathbf{z}^\top(t)[\mathbf{z}(t) - \tilde{\mathbf{z}}(t)]}_{=0 \text{ (by orthogonality)}} \\ &= \|\mathbf{z}(t)\|^2 - \|\tilde{\mathbf{z}}(t)\|^2 \\ &= \|\mathbf{z}(t)\|^2 - \|\mathbf{h}(t)\|^2 \quad (\text{by Lemma 3}). \end{aligned}$$

The results follows from the summing over T and from the definitions of E (4-3) and \tilde{E} (4-5). \square

In Section 5.3.4, we demonstrate that FRAHST's reconstruction error lies in the given range as expected: particularly when setting a range $[f_E, F_E] = [0.96, 0.98]$, our algorithm correctly maintains a relative reconstruction error between 2% and 4% across various datasets.

4.2.1

Exponential forgetting

We can better exploit the temporal nature of the data by using an exponential forgetting factor, which allows us to follow trend drifts over time. We use the same α from (3-1) to properly keep track of the energy, discounting it with the same rate. The update at each step is then:

$$E(t+1) = \frac{t-1}{t}\alpha E(t) + \frac{1}{t}\|\mathbf{z}(t)\|^2 \quad (4-6)$$

$$\tilde{E}(t+1) = \frac{t-1}{t}\alpha \tilde{E}(t) + \frac{1}{t}\|\tilde{\mathbf{z}}(t)\|^2. \quad (4-7)$$

The value of α does not affect the computational cost of our method and is therefore more appealing than a sliding window to compute sample statistics, as the latter has buffering requirements.

4.2.2

Adapting data structures

Since r may change over time, our algorithm needs to account for resizing of the two internal matrices Q and S matrices with dimensions $N \times r$ and $r \times r$ respectively. In order to best preserve the properties of each matrix and subsequently minimize losses in the quality of the tracked subspace, we devise two simple heuristics to restructure both matrices. When r is incremented, we append the normalized orthogonal error given by the current updated

projection matrix:

$$\mathbf{z}_{\perp}^{\dagger}(t) = \mathbf{z}(t) - \mathbf{Q}(t)\mathbf{Q}^{\top}(t)\mathbf{z}(t) \quad (4-8)$$

$$\mathbf{Q}(t+1) = \begin{bmatrix} \mathbf{Q}(t) & \mathbf{z}^{\dagger}(t)/\|\mathbf{z}^{\dagger}(t)\| \end{bmatrix} \quad (4-9)$$

where \mathbf{q}_{r+1} now serves as an instantaneous estimate for the new basis able to capture interesting information in the new direction. Even though we relaxed the upper-right triangular shape constraint of the S -matrix, in the case of an increase in r we restructure the matrix by adding the following estimate to its new diagonal element:

$$\mathbf{S}(t+1) = \begin{bmatrix} \mathbf{S}(t) & 0 \\ 0 & \|\mathbf{z}_{\perp}^{\dagger}(t)\|^2 \end{bmatrix}. \quad (4-10)$$

In the case when r is decremented, the two matrices are simply truncated from $r \times r$ to $(r-1) \times (r-1)$ and we discard the values from the previous entries. The motivation for these updates comes mainly from similar ideas employed in [Yang, 1995b] and while other heuristics were experimented, including very simple ones such as initializing with random orthonormal vectors, they all severely damage the subspace estimates. This completes the necessary steps to our extension which we summarise below in Figure 4.1.

4.3

Exception handling

The zero-input case requires the operation of the algorithm in ‘idle mode’, which enables the algorithm to handle cases of vanishing inputs. It can be seen in (4-11a) and (4-11b) that a null input will zero both $Z(t)$ and $\mathbf{h}(t)$. Hence no updating of the basis estimate in $\mathbf{Q}(t)$ is necessary because there is no innovation in a zero input. We add a check to test the condition $Z(t) < \sigma$ before (4-11c) to bypass all remaining computations for that step. The singularity threshold σ is machine and platform dependent but is a positive constant very close to zero.

Fortunately, our rank estimation routine will already guarantee that the S -matrix is not rank deficient, hence we do not need to handle special conditions elsewhere.

4.4

Achieving lower asymptotic computational complexity

Our rank estimation is dominated by $2Nr$ flops from computing the reconstruction error based on the updated projection matrix (4-11n), hence

parameters: $0 < f_E < F_E < 1$ and $0 < \alpha < 1$
 $r(0) \leftarrow 1$
 $\mathbf{Q}(0) \leftarrow$ random orthonormal
 $\mathbf{S}(0) \leftarrow \sigma \mathbf{I}$, where σ is a small positive constant
 for $t = 1, 2, \dots$

read $\mathbf{z}(t) \in \mathbb{R}^N$

$\mathbf{h}(t) = \mathbf{Q}^\top(t-1)\mathbf{z}(t)$ (4-11a)

$Z(t) = \mathbf{z}^\top(t)\mathbf{z}(t) - \mathbf{h}^\top(t)\mathbf{h}(t)$ (4-11b)

$\mathbf{X}(t) = \alpha \mathbf{S}(t-1) + \mathbf{h}(t)\mathbf{h}^\top(t)$ (4-11c)

$\mathbf{X}^\top(t)\mathbf{b}(t) = Z^{1/2}\mathbf{h}(t) \xrightarrow{\text{solve}} \mathbf{b}(t)$ (4-11d)

$\varphi^2(t) = \frac{1}{2} + \frac{1}{\sqrt{4(\mathbf{b}^\top(t)\mathbf{b}(t) + 1)}}$ (4-11e)

$\delta(t) = \frac{\varphi(t)}{Z^{1/2}(t)}$ (4-11f)

$\mathbf{v}(t) = \left[\frac{1 - 2\varphi^2(t)}{2\varphi(t)} \right] \mathbf{b}(t)$ (4-11g)

$\mathbf{S}(t) = \mathbf{X}(t) - \frac{1}{\delta(t)}\mathbf{v}(t)\mathbf{h}^\top(t)$ (4-11h)

$\mathbf{e}(t) = \delta(t)\mathbf{z}(t) - \mathbf{Q}(t-1)[\delta(t)\mathbf{h}(t) - \mathbf{v}(t)]$ (4-11i)

$\mathbf{Q}(t) = \mathbf{Q}(t-1) - 2\mathbf{e}(t)\mathbf{v}^\top(t)$ (4-11j)

$E(t) = \alpha E(t-1) + \|\mathbf{z}\|^2$ (4-11k)

$\tilde{E}(t) = \alpha \tilde{E}(t-1) + \|\mathbf{h}\|^2$ (4-11l)

if $\tilde{E}(t) < f_E E(t)$ and $r(t) < N$ (4-11m)

$\mathbf{z}_\perp^\dagger(t) = \mathbf{z}(t) - \mathbf{Q}(t)\mathbf{Q}^\top(t)\mathbf{z}(t)$ (4-11n)

$\mathbf{Q}(t+1) = [\mathbf{Q}(t) \quad \mathbf{z}_\perp^\dagger(t)/\|\mathbf{z}_\perp^\dagger(t)\|]$ (4-11o)

$\mathbf{S}(t+1) = \begin{bmatrix} \mathbf{S}(t) & \mathbf{0} \\ \mathbf{0} & \|\mathbf{z}_\perp^\dagger(t)\|^2 \end{bmatrix}$ (4-11p)

$r(t+1) = r(t) + 1$ (4-11q)

else if $\tilde{E}(t) > F_E E(t)$ and $r(t) > 1$ (4-11r)

$\mathbf{Q}(t+1) =$ delete the rightmost column of $\mathbf{Q}(t)$ (4-11s)

$\mathbf{S}(t+1) =$ delete bottom row and rightmost column of $\mathbf{S}(t)$ (4-11t)

$r(t+1) = r(t) - 1$ (4-11u)

Figure 4.1: FRAHST algorithm

the entire algorithm has time complexity of $5Nr + \mathcal{O}(N + r^3)$ per update. The space complexity $\mathcal{O}(Nr)$ is similarly very low.

The value of r depends on the threshold range and the actual data read at each interval t . For example, in the specific datasets we experimented, r

rarely exceeds the value of 6 even for N larger than 200 when capturing 96% of the energy. Despite the fact that $r \ll N$ in most cases, in the worst case analysis we have $r = N$ and one can perceive the typical $\mathcal{O}(r^3)$ complexity for solving the systems of linear equations in step (4-11d) as a stumbling block. Gaussian elimination or LU decomposition are the typical² algorithms of choice and require approximately r^3 flops. It is theoretically possible to achieve $\mathcal{O}(r^{2.376})$ by the asymptotically fastest known matrix multiplication (but impractical) Coppersmith–Winograd algorithm [Coppersmith and Winograd, 1990, Robinson, 2005] or with the more realistic Strassen method that requires about $4.7r^{2.807}$ flops [Bailey et al., 1990]. But neither alternatives are optimal. In [Strobach, 2009a], the author suggests that the recursivity that is inherent in the subspace tracking algorithm can be exploited by recursive updating factorizations of inner matrices, but not much detail is provided. We follow this thread and show how to apply the same suggested ideas for our rank-adaptive extension and reduce the complexity of step (4-11d) to $\mathcal{O}(r^2)$.

4.4.1

Efficiently updating QR-factorizations

Instead of refactoring the X matrix on every update from scratch, it is much more efficient to update each factor as necessary and a procedure for such recurrent updates is widely known and require $\mathcal{O}(r^2)$ flops (Chapter 12.5 in Golub and Van Loan [1996] and Gill et al. [1972]). The idea is to work with the QR decompositions of the X - and S - $r \times r$ square matrices and recognize that (4-11c) and (4-11h) are ‘QR = QR + rank one’ updates that can be expressed as

$$\mathbf{X}(t) = \mathbf{X}^Q(t)\mathbf{X}^R(t) = \mathbf{S}^Q(t)\mathbf{S}^R(t) - \mathbf{h}^\top(t)\mathbf{h}(t) \quad (4-12)$$

$$\mathbf{S}(t) = \mathbf{S}^Q(t)\mathbf{S}^R(t) = \mathbf{X}^Q(t)\mathbf{X}^R(t) - 1/\delta(t)\mathbf{v}(t)\mathbf{h}^\top(t) \quad (4-13)$$

We illustrate the procedure considering the case $\mathbf{A}^Q\mathbf{A}^R = \mathbf{A} \in R^{r \times r}$, where we need to compute the QR factorization of $\mathbf{A} + \mathbf{bc}^\top = \mathbf{A}_1^Q\mathbf{A}_1^R$ and $\mathbf{b}, \mathbf{c} \in R^r$. The expression is written in the form $\mathbf{A} + \mathbf{bc}^\top = \mathbf{A}^Q(\mathbf{A}^R + \mathbf{wc}^\top)$ where $\mathbf{w} = \mathbf{A}^Q\mathbf{b}$. We use the procedure detailed in [Golub and Van Loan, 1996] to restore the QR factorization of \mathbf{A} which relies on computing two sequences

²Most software packages seem to call the DGESV from LAPACK to compute the solution of linear equations.

J and G of Givens rotations as shown below:

$$\text{QR-UPDATE}(\mathbf{A}^Q, \mathbf{A}^R, \mathbf{b}, \mathbf{c}) \quad (4-14)$$

$$\mathbf{J}_1^\top \dots \mathbf{J}_{m-1} \mathbf{w} = \pm \|\mathbf{w}\| \mathbf{e}_1 \quad (4-15)$$

$$\mathbf{J}_1^\top \dots \mathbf{J}_{m-1} (\mathbf{A}^R + \mathbf{w} \mathbf{c}^\top) = \mathbf{H} \pm \|\mathbf{w}\| \mathbf{e}_1 \mathbf{b}^\top = \mathbf{H}_1 \quad (4-16)$$

$$\mathbf{G}_{m-1}^\top \dots \mathbf{G}_1^\top \mathbf{H}_1 = \mathbf{A}_1^R \quad (4-17)$$

$$\mathbf{A}_1^Q = \mathbf{A}^Q \mathbf{J}_{m-1} \dots \mathbf{J}_1 \mathbf{G}_1 \dots \mathbf{G}_{m-1} \quad (4-18)$$

$$\text{output } \mathbf{A}_1^Q, \mathbf{A}_1^R \quad (4-19)$$

where $\mathbf{e}_1 \equiv [1 \ 0 \ \dots \ 0]^\top \in \mathbb{R}^r$. Each \mathbf{J}_k is a Given rotation in planes k and $k + 1$ [Algorithm 5.1.3 in Golub and Van Loan, 1996]. \mathbf{H}_1 can be shown to be an ‘almost upper triangular’ Hessenberg matrix and its QR factorization (4-17) can be computed by applying the G rotations in $\mathcal{O}(r^2)$ [Algorithm 5.2.3 in Golub and Van Loan, 1996]. The entire update requires about $26r^2$ flops. Once we have the QR factors, it is straight-forward to solve the linear system from (4-11d) with complexity $\mathcal{O}(r^2)$ by using back-substitution, which follows from pre-multiplying the transpose of \mathbf{X}^Q and from the fact that \mathbf{X}^R is an upper triangular matrix (step (4-20b)).

We are now left with the task to account for the resizing of the QR factors of the S -matrix. Similarly, we use four procedures based on Givens rotations from [Golub and Van Loan, 1996, Hajek, 2009] for both appending and deleting column and row vector of the original matrix according to steps (4-11p) and (4-11t) whilst maintaining the QR structure with quadratic complexity. We do not give all details here, but all the routines QR-INSERT-ROW, QR-INSERT-COL, QR-DELETE-COL and QR-DELETE-ROW are all very much alike to the above QR-UPDATE and surgically manipulate the values of the upper triangular A^R -matrix via Givens rotations which is then used to produce a Hessenberg matrix to yield the final two rotations that are used to update the A^Q -matrix. We give the pseudo-code for the asymptotically faster modification of FRAHST in Figure 4.2.

Considering that each routine costs roughly a total of $26r^2$ flops, and adding r^2 operations from the back-substitution step totals about $5Nr + 157r^2$ which suggests that the asymptotically faster modified version only pays off when $r > 157$. In practice, any implementation can switch to this version when the dimension of the principal subspace is that high while using few operations for the most common lower rank case. Regarding the implementation, we note that MATLAB, for instance, only offers the QR-UPDATE function³ but all

³See the documentation for the *qrupdate* at <http://www.mathworks.com/>.

Initialize $\mathbf{S}^Q(0) = \mathbf{S}^R(0) = \sigma \mathbf{I}$

Same steps (4-11a)-(4-11u) except for the following.

Replace (4-11c) and (4-11d) with:

$$\mathbf{X}^Q(t)\mathbf{X}^R(t) \leftarrow \text{QR-UPDATE}(\alpha\mathbf{S}^Q(t), \alpha\mathbf{S}^R(t), \mathbf{h}, \mathbf{h}) \quad (4-20a)$$

$$\mathbf{X}^R\mathbf{b} = (\mathbf{X}^Q)^\top \mathbf{Z}(t)^{1/2} \mathbf{h}(t) \xrightarrow{\text{back-substitution}} \mathbf{b}(t). \quad (4-20b)$$

Replace (4-11h) with:

$$\mathbf{S}^Q(t)\mathbf{S}^R(t) \leftarrow \text{QR-UPDATE}(\mathbf{X}^Q(t), \mathbf{X}^R(t), 1/\delta(t)\mathbf{v}, \mathbf{h}). \quad (4-20c)$$

Replace (4-11p) with:

$$\mathbf{S}^Q(t+1)\mathbf{S}^R(t+1) \leftarrow \text{QR-INSERT-ROW}(\mathbf{S}^Q(t), \mathbf{S}^R(t), r(t+1), \mathbf{0}) \quad (4-20d)$$

$$\mathbf{S}^Q(t+1)\mathbf{S}^R(t+1) \leftarrow \text{QR-INSERT-COL}(\mathbf{S}^Q(t+1), \mathbf{S}^R(t+1), r(t+1), [0 \ 0 \ \dots \ \|\mathbf{z}_\perp^\dagger(t)\|^2]). \quad (4-20e)$$

Replace (4-11t) with:

$$\mathbf{S}^Q(t+1)\mathbf{S}^R(t+1) \leftarrow \text{QR-DELETE-COL}(\mathbf{S}^Q(t), \mathbf{S}^R(t), r(t)) \quad (4-20f)$$

$$\mathbf{S}^Q(t+1)\mathbf{S}^R(t+1) \leftarrow \text{QR-DELETE-ROW}(\mathbf{S}^Q(t+1), \mathbf{S}^R(t+1), r(t)). \quad (4-20g)$$

Figure 4.2: Asymptotically faster FRAHST algorithm with recurrent QR updates.

routines are found in the Fortran package from Hajek [2009] which is greatly optimized.

4.4.2 Efficiently updating LU-factorizations

We explore a second alternative to reduce the $\mathcal{O}(r^3)$ term from the original complexity in terms of fast Bennett LU-factor updating [Bennett, 1965], cited in [Strobach, 2009a] and the algorithm is detailed in [Stange et al., 2007] and it only requires $4r^2$ flops for updating. Similarly to the previous extension, this allows the linear system of equations to be solved in $\mathcal{O}(r^2)$ by exploiting the ‘LU = LU + rank one’ updates in each step. However, now we are dealing with the lower and upper triangular factors of the S - and X -matrices. Analogous to the previous section, the main goal is to update the LU decomposition after a rank-one update:

$$\underbrace{\mathbf{A}^L \mathbf{A}^U}_{\mathbf{A}} + \mathbf{bc}^\top = \underbrace{\mathbf{A}_1^Q \mathbf{A}_1^R}_{\text{New LU factors}} \quad (4-21)$$

which can be accomplished with an algorithm such as Bennett’s algo-

rithm, which is listed below [from Stange et al., 2007]:

LU-UPDATE($\mathbf{A}^L, \mathbf{A}^U, \mathbf{b}, \mathbf{c}$)
 for $1 \leq i \leq r$ (4-22)

$$\mathbf{A}_{ii}^U = \mathbf{A}_{ii}^U + b_i c_i \quad (4-23)$$

$$c_i = c_i / \mathbf{A}_{ii}^U \quad (4-24)$$

for $i + 1 \leq j \leq r$ (4-25)

$$b_j = b_j - b_i \mathbf{L}_{ji} \quad (4-26)$$

$$\mathbf{L}_{ji} = \mathbf{L}_{ji} + c_i b_j \quad (4-27)$$

$$\mathbf{U}_{ij} = \mathbf{U}_{ij} + b_i c_j \quad (4-28)$$

$$c_j = c_j - c_i \mathbf{U}_{ij} \quad (4-29)$$

Once we have the LU factors for the X - matrix, solving the linear system of equations from (4-11d) is a direct application of both forward- and back-substitutions. The pseudo-code for the asymptotically faster modification of FRAHST is given in Figure 4.3.

Initialize $\mathbf{S}^L(0) = \mathbf{S}^U(0) = \sigma \mathbf{I}$

Same steps (4-11a)-(4-11u) except for the following.

Replace (4-11c) and (4-11d) with:

$$\mathbf{X}^L(t) \mathbf{X}^U(t) \leftarrow \text{LU-UPDATE}(\alpha \mathbf{S}^L(t), \alpha \mathbf{S}^U(t), \mathbf{h}, \mathbf{h}) \quad (4-30a)$$

$$\mathbf{X}^L \mathbf{y}(t) = \mathbf{Z}(t)^{1/2} \mathbf{h}(t) \xrightarrow{\text{forward-substitution}} \mathbf{y}(t). \quad (4-30b)$$

$$\mathbf{X}^U \mathbf{b} = \mathbf{y}(t) \xrightarrow{\text{back-substitution}} \mathbf{b}(t). \quad (4-30c)$$

Replace (4-11h) with:

$$\mathbf{S}^L(t) \mathbf{S}^U(t) \leftarrow \text{LU-UPDATE}(\mathbf{X}^L(t), \mathbf{X}^U(t), 1/\delta(t) \mathbf{v}, \mathbf{h}). \quad (4-30d)$$

Replace (4-11p) with:

$$\mathbf{S}^L(t+1) = \begin{bmatrix} \mathbf{S}^L(t) & \mathbf{0} \\ \mathbf{0} & \|\mathbf{z}_{\perp}^{\dagger}(t)\| \end{bmatrix} \quad \mathbf{S}^U(t+1) = \begin{bmatrix} \mathbf{S}^U(t) & \mathbf{0} \\ \mathbf{0} & \|\mathbf{z}_{\perp}^{\dagger}(t)\| \end{bmatrix} \quad (4-30e)$$

Replace (4-11t) with:

$$\mathbf{S}^L(t+1) = \text{delete bottom row and right most column of } \mathbf{S}^L(t) \quad (4-30f)$$

$$\mathbf{S}^U(t+1) = \text{delete bottom row and right most column of } \mathbf{S}^U(t) \quad (4-30g)$$

Figure 4.3: Asymptotically faster FRAHST algorithm with recurrent LU updates.

The total complexity is about $5Nr + 10r^2$ due to two $4r^2$ LU updates and two r^2 substitutions for solving the linear system. This is the lowest complexity

we achieve for the FRAHST algorithm. However, Bennett’s algorithm is known to cause numerical stability problems during the updating procedure [Stange et al., 2007], and we indeed experience a loss of quality in the subspace estimates for this faster version in practice. Therefore, in general we are recommending the previous QR version for a more stable approach that performs equally well as the original. In Chapter 5, we compare the performance of all the proposed versions (see Figure 5.4(c)).

Summary of the algorithm

The latent variables $\mathbf{h}(t)$ give us a compact representation of the raw data $\mathbf{z}(t)$ and FRAHST guarantees high reconstruction accuracy (in terms of the relative squared error, which is less than $1 - f_E$). When our streams are highly correlated, as we often expect to be the case, the latent dimension r is much smaller than the number N of streams and the Q -matrix captures the linear combinations between them. Traditional batch methods for estimating the principal components require time that depends on the duration T , and most subspace trackers do not adapt the rank of the principal subspace in an online manner. Our algorithm does not need to store any past values and has *very* low computational requirements and can be applied to several important tasks, such as anomaly detection by monitoring the current tracked rank and forecasting by fitting auto-regressive models on the latent variables for the next interval, especially to handle missing data – the latter techniques were advocated in [Papadimitriou et al., 2005], but FRAHST has a lower dominant complexity $\mathcal{O}(Nr)$ instead of $\mathcal{O}(Nr^2)$ needed for SPIRIT to guarantee the correct expected reconstruction error. This makes our algorithm more attractive for the streaming scenario.

4.5

Real-time Anomaly Detection System

Now, we formalize the anomaly detection procedure based on the FRAHST algorithm, and elaborate on the architecture for implementing a real-time monitoring system.

4.5.1

Raising alarms

Our anomaly detection procedure follows directly from the rank-adaptive nature of FRAHST. We propose to raise alarms whenever there is an increase in rank, as laid out below.

```

last ← 0
for t = 1, 2, ...
    FRAHST(t) (4-31)
    if r(t) > r(t - 1) (4-32)
        if t > last + 1: raise an alarm at t (4-33)
        last ← t (4-34)

```

Figure 4.4: Anomaly detection routine

The control variable `last` is used to suppress alarms when there are rank increments in consecutive interval, which are likely to be false alarms. The intuition is that a change in the underlying streams might be so great, that more than one additional latent variables are necessary to achieve the same reconstruction error.

4.5.2 Event-driven Architecture

We devise a loose-coupled event-driven architecture where all modules are *publishers* or *subscribers* to an event broker, also known as event/message bus [and widely popularized by Hohpe and Woolf, 2003]. The most important components in a data center can be monitored using the Simple Network Management Protocol (SNMP), which include routers, access servers, switches, bridges, hubs, temperature sensors and computer hosts. The variables accessible via SNMP are organized in hierarchical Management Information Bases (MIBs) and a manager can get information from agents using the UDP protocol, typically via port 161. Not all data is available via SNMP, such as statistics from a running MySQL database or an Apache web server. As illustrated in Figure 4.5, specific plugin modules are implemented to obtain all desired information which is then normalized into a JSON⁴ format which is sent to the bus to be consumed by different processes via remote calls. We also integrate with Nagios⁵ in order to reuse the existing plugins to enable a wide variety of data to be available into our system. This normalized message contains: timestamp, event type identifier and the collected variables along with their names. The user configures, in runtime via a web interface, a model that represents the entities in the monitored infrastructure and are mapped to this raw data accordingly. There are many challenges to implement efficient communication

⁴Javascript Object Notation

⁵A popular opensource monitoring system.

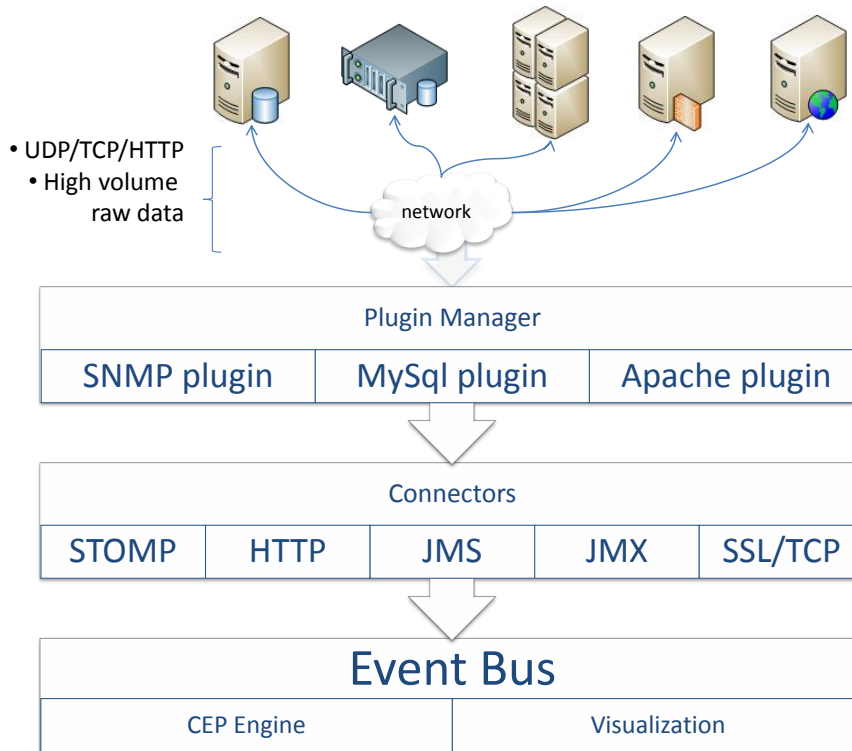


Figure 4.5: Raw asynchronous events are collected and normalized messages are sent to the event bus.

strategies to ensure minimum latency. For this reason most communication is done over a persistent TCP socket using the STOMP⁶ protocol.

Different implementations were considered, but we chose Esper⁷ as the event stream processing engine⁸ to power our solution. It is an excellent open source project, easy to embed and with the probably most expressive event processing query language available. A CEP engine is analogous to a ‘inverted’ traditional database, where instead of having queries being made over persistent tuples of data we now have persistent queries that run continuously over streaming data.

The idea is to monitor streams as defined by continuous queries over the incoming raw data. More specifically, the data center operator chooses a set of streams to be monitored together and a query can be placed to join the corresponding underlying streams in order to produce the input vector $\mathbf{z}(t)$ for our algorithm at periodic equidistant intervals – the join operator and the output rate of the query are features of the core stream processing engine. This is necessary due to the asynchronous nature of the original events. Figure

⁶Streaming Text Orientated Message Protocol

⁷More information at <http://www.espertech.com/>.

⁸The terms event stream processing (ESP), complex event processing (CEP) and data stream management system (DSMS) may be used interchangeably.

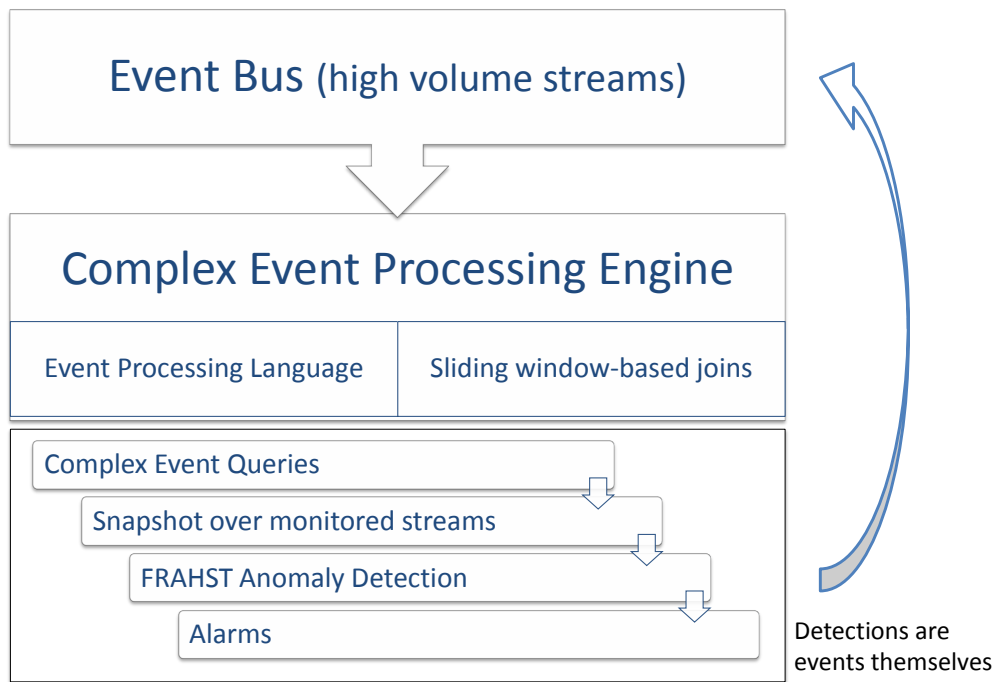


Figure 4.6: We illustrate the processing flow in the system: a query is used to join raw streams into a derived complex event that feeds our algorithm, which is then apt to detect anomalies in variables under surveillance.

4.6 overviews the flow of the anomaly detection system and it is important to realize that we have multiple instances of FRAHST running attached to different queries with different set of streams. For example, one instance will be detecting anomalies from a query that represents data regarding cpu and memory usage from a cluster A, while another will monitor cluster B. This solution is flexible and allows an expert user to leverage his previous knowledge while maintaining an interpretation over the projected latent variables.

We use the adapter pattern for the actual integration with FRAHST where a processor is attached to a continuous query and updates the algorithm upon new events. The step (4-33) in Figure 4.4 notifies a listener which consequently publishes an alarm event to the message broker so further actions are taken; such as update charts, send SMS or XMPP messages. It is interesting to observe that generated alarms can also be considered event streams, and hence can be monitored. The sequence diagram in Figure 4.7 better illustrates the interaction between the different subcomponents in the monitoring system. During the each step interval, data is collected from different sources and is aggregated by user-defined queries. Each query has an associated query listener, which is responsible for forming the input data vector and delegating

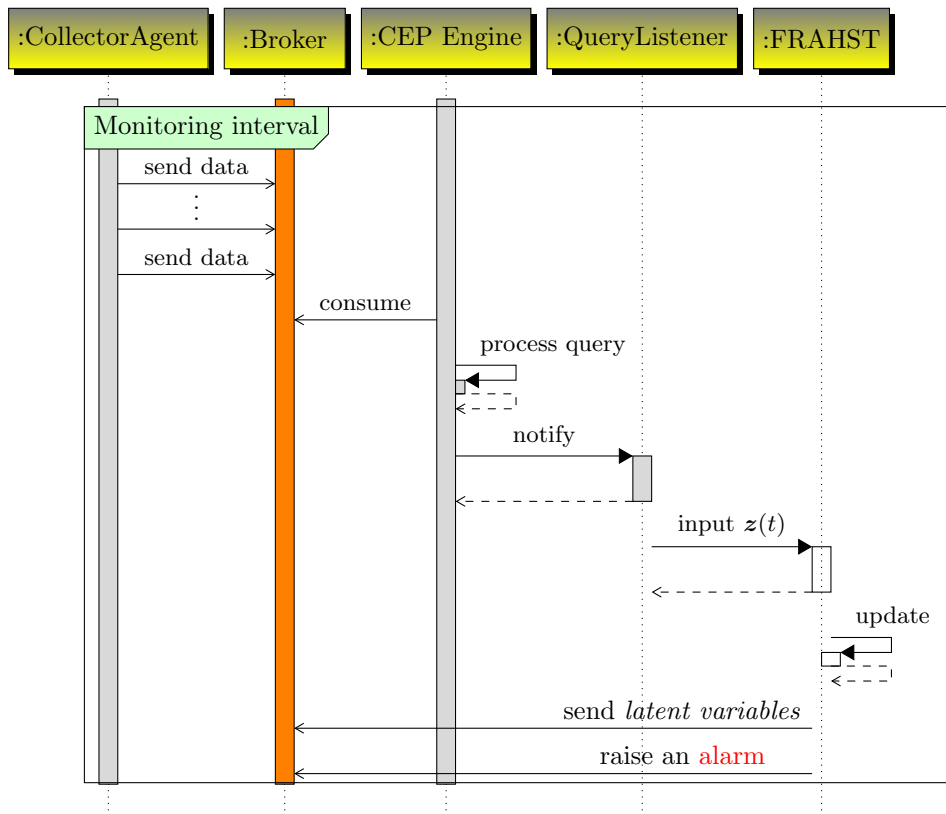


Figure 4.7: Simplified sequence diagram of the update step.

over an adapter to the FRAHST implementation. Although we do not depict it here, a user interface can consume both latent variables and alarm information from the broker for visualization purposes. Most of the system is implemented in the Java language, and ActiveMQ⁹ was used for the event broker.

Our solution has been adopted for a real data center, where it has been shown to achieve high throughput and low latency. This use case has been presented at an international forum recently [Clemente and Vieira, 2009].

⁹See <http://activemq.apache.org/>.