

Chapter 4

Column and Column- Row Decompositions



Column Decompositions

Motivation

- SVD is often hard to interpret and yields dense factorizations
 - NMF tries to address these problems with varying success
- But if original data is sparse & easy to interpret, why not use it in the decompositions?

The CX decomposition

- In the **CX decomposition** we are given a matrix **A** and a rank k , and we need to select k **columns** of **A** into matrix **C** and build matrix **X** s.t. we minimize $\|\mathbf{A} - \mathbf{CX}\|_{\xi}$
- ξ is either F or 2
- A.k.a. **column subset selection problem** (CSSP)

Why CX?

- The columns of **C** preserve the original interpretation of columns of **A**
 - Even complex constraints are satisfied if the original data satisfied them
- Feature selection
 - Selects the columns that can be used to explain the rest
 - Compare to the dimensionality reduction

Alternative target function

- Building \mathbf{C} is the hard part of CX decompositions
- Given \mathbf{A} and \mathbf{C} , \mathbf{X} can be computed with the pseudo-inverse
 - $\mathbf{X} = \mathbf{C}^+ \mathbf{A}$
 - Alternative target function for CX:
minimize $\|\mathbf{A} - \mathbf{C}\mathbf{C}^+ \mathbf{A}\|_{\xi} = \|\mathbf{A} - \mathbf{P}_{\mathbf{C}} \mathbf{A}\|_{\xi}$

How to select C?

- Exhaustive: try all $\binom{m}{k}$ subsets of columns
 - Not very scalable
- Try to select the columns in a clever way
 - But how?
- Sample columns w.r.t. carefully selected probabilities
 - Avoids deterministic worst-case scenarios

Related idea: RRQR

- The **rank-revealing QR** (RRQR) factorization of matrix \mathbf{A} is

$$\mathbf{A}\mathbf{\Pi} = \mathbf{Q}\mathbf{R} = \mathbf{Q} \begin{pmatrix} \mathbf{R}_{11} & \mathbf{R}_{12} \\ 0 & \mathbf{R}_{22} \end{pmatrix}$$

n -by- n orthogonal
 Permutation matrix that satisfies
 n -by- m
 k -by- k upper-triangular w/ positive diagonal
 k -by- $(m-k)$
 $(n-k)$ -by- $(m-k)$
 k th singular value of \mathbf{A}

$$\frac{\sigma_k(\mathbf{A})}{p_1(k, m)} \leq \sigma_{\min}(\mathbf{R}_{11}) \leq \sigma_k(\mathbf{A})$$

$$\sigma_{k+1}(\mathbf{A}) \leq \sigma_{\max}(\mathbf{R}_{22}) \leq p_2(k, m)\sigma_{k+1}(\mathbf{A})$$

Some polynomial on k and m

CX and RRQR

- Let $\mathbf{A}\mathbf{\Pi} = \mathbf{Q}\mathbf{R}$ and let $\mathbf{\Pi}_k$ be the first k columns of $\mathbf{\Pi}$ and $\mathbf{C} = \mathbf{A}\mathbf{\Pi}_k$ some k columns of \mathbf{A}
- Now $\|\mathbf{A} - \mathbf{P}_c\mathbf{A}\|_{\xi} = \|\mathbf{R}_{22}\|_{\xi}$, $\xi = F$ or 2
- In particular $\|\mathbf{A} - \mathbf{P}_c\mathbf{A}\|_2 \leq p_2(k, m)\|\mathbf{A} - \mathbf{A}_k\|_2$
 - $\mathbf{A}_k = \mathbf{U}_k\mathbf{\Sigma}_k\mathbf{V}_k^T$ (truncated SVD)
 - $\mathbf{C}\mathbf{X}$ is $p_2(k, m)$ -approximation to SVD

Computing CX by sampling

- Let $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ be the input and its SVD and \mathbf{V}_k the truncated \mathbf{V}
- Sample columns of \mathbf{A} with replacement
 - Probability p_j for selecting column j is
$$p_j = \|(\mathbf{V}_k^T)_j\|_2^2/k$$
 - Sample $O(k^2 \log(1/\delta)/\epsilon^2)$ columns and repeat $\log(1/\delta)$ times returning the least-error sample

Notes on sampling

- We can prove that

$$\|\mathbf{A} - \mathbf{P}_C \mathbf{A}\|_F \leq (1 + \epsilon) \|\mathbf{A} - \mathbf{A}_k\|_F$$

with probability at least $1 - \delta$

- Notice that \mathbf{C} has much more than k columns
- $O(k^2 \log(1/\delta)/\epsilon^2)$ with large hidden constants

Why does sampling work?

- Intuitively, if \mathbf{A} is of low rank ($k \ll n$), \mathbf{A} should have many almost-similar columns
- If we sample many columns enough, we should get a representative for each set of similar columns
 - ⇒ We need to sample more columns than the rank
 - Or our error depends on the rank...

CX with exact k

- Construct larger-than- k CX decomposition as above for $c = O(k \log k)$ columns (and using rank- k truncated SVD)
 - Let $\mathbf{\Pi}_1$ be the m -by- c matrix that selects c columns s.t.
 $\mathbf{C} = \mathbf{A}\mathbf{\Pi}_1$
 - Let \mathbf{D}_1 be c -by- c diagonal s.t. if j th column is selected on round i , $(\mathbf{D}_1)_{jj} = (cp_j)^{-1/2}$
- Run RRQR algorithm for $\mathbf{V}_k^T \mathbf{\Pi}_1 \mathbf{D}_1$ to select exactly k columns of $\mathbf{V}_k^T \mathbf{\Pi}_1 \mathbf{D}_1$ with matrix $\mathbf{\Pi}_2$ (c -by- k)
 - return $\mathbf{C} = \mathbf{A}\mathbf{\Pi}_1\mathbf{\Pi}_2$

From $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$



Notes on the exact- k CX

- $\Pr[\|\mathbf{A} - \mathbf{P}_c \mathbf{A}\|_F \leq \Theta(k \log^{1/2} k) \|\mathbf{A} - \mathbf{A}_k\|_F] \geq 0.8$
- The sampling phase still requires really many columns (high hidden constants)
 - But in practice something like $c = 5k$ works
- Any RRQR algorithm can be used for the second step
 - But the analysis depends on the chosen algorithm

Non-Negative CX

Motivation

- If data is non-negative, so is \mathbf{C}
 - But \mathbf{X} can contain negative values in standard \mathbf{CX}
 - Non-negative \mathbf{X} yields “parts-of-whole” interpretation similar to NMF
 - Selected columns are “pure” while others are mixtures of the pure columns
 - Non-negativity also improves sparsity

The non-negative CX decomposition

- In the **non-negative CX decomposition** (NNCX) we are given a non-negative matrix \mathbf{A} and a rank k , and we need to select k **columns** of \mathbf{A} into matrix \mathbf{C} and build a non-negative matrix \mathbf{X} s.t. we minimize $\|\mathbf{A} - \mathbf{CX}\|_F$

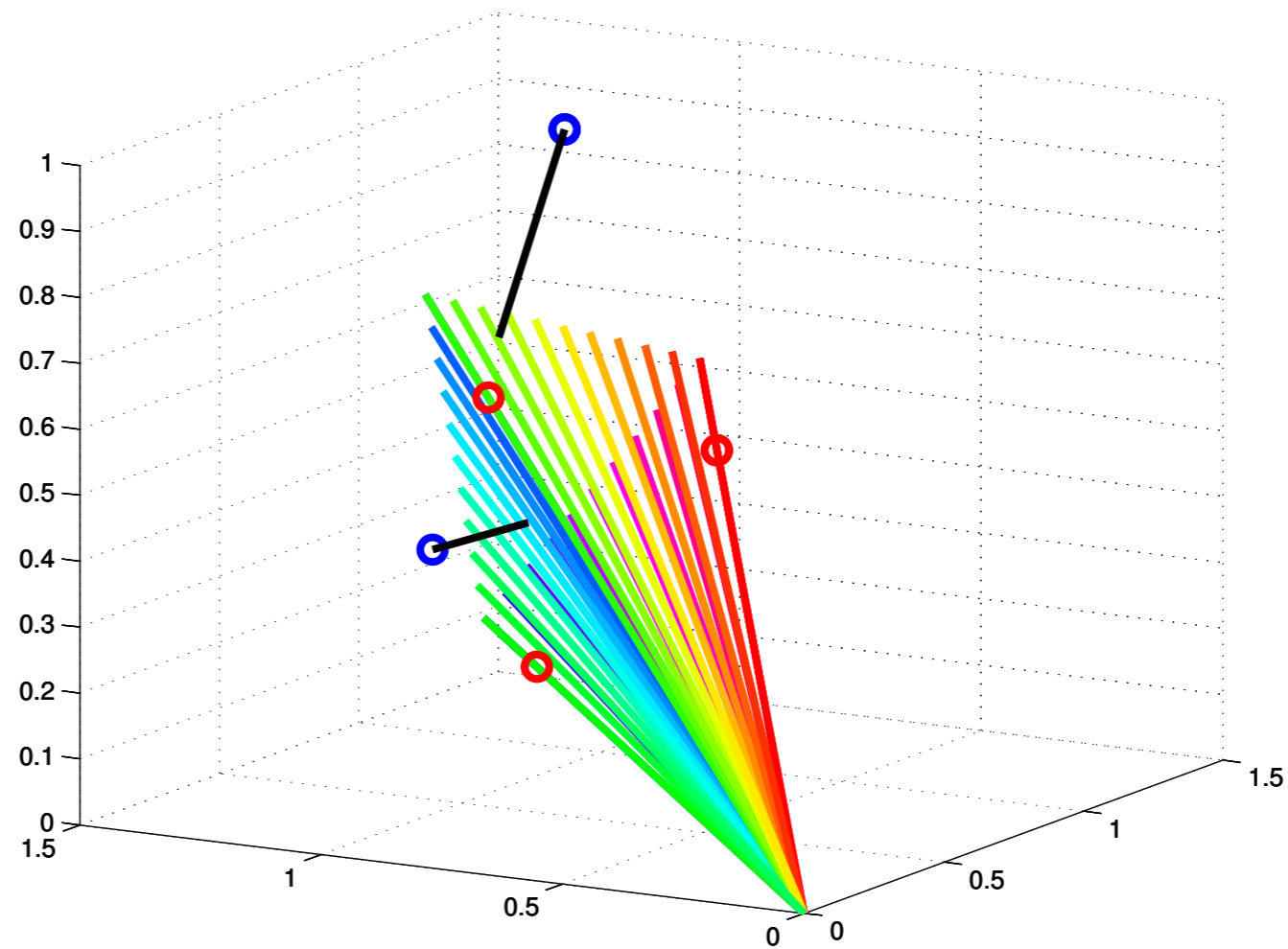
Geometry of NNCX

Columns in \mathbf{C}

Columns not in \mathbf{C}

Convex cone

Projections



Cones and columns

- Consider the cone spanned by columns of \mathbf{A} , $\text{cone}(\mathbf{A})$
 - If removing column j of \mathbf{A} changes the cone, that column is **extremal**
 - Otherwise it is **internal**
 - Selecting all extremal columns to \mathbf{C} gives us $\mathbf{A} = \mathbf{C}\mathbf{X}$ with nonnegative \mathbf{X}

Algorithm for NNCX

- When we cannot select all extremal columns, we must choose which of them to select
 - Our goal is to maximize the volume of the convex cone
 - Finding the extremal columns is not easy
 - Given the columns, we must compute the non-negative projection

The convex_cone algorithm

- Set $\mathbf{R} \leftarrow \mathbf{A}$
- **repeat**
 - Select column \mathbf{c} with highest norm in the residual \mathbf{R}
 - Normalize \mathbf{c} to unit norm
 - Solve nonnegative \mathbf{x} that minimizes $\|\mathbf{R} - \mathbf{c}\mathbf{x}^T\|$
 - Set $\mathbf{R} \leftarrow \mathbf{R} - \mathbf{c}\mathbf{x}^T$
- **until** k columns are selected
- Set \mathbf{C} to the columns of \mathbf{A} corresponding to the selected \mathbf{c} and solve nonnegative \mathbf{X} minimizing $\|\mathbf{A} - \mathbf{C}\mathbf{X}\|_F$

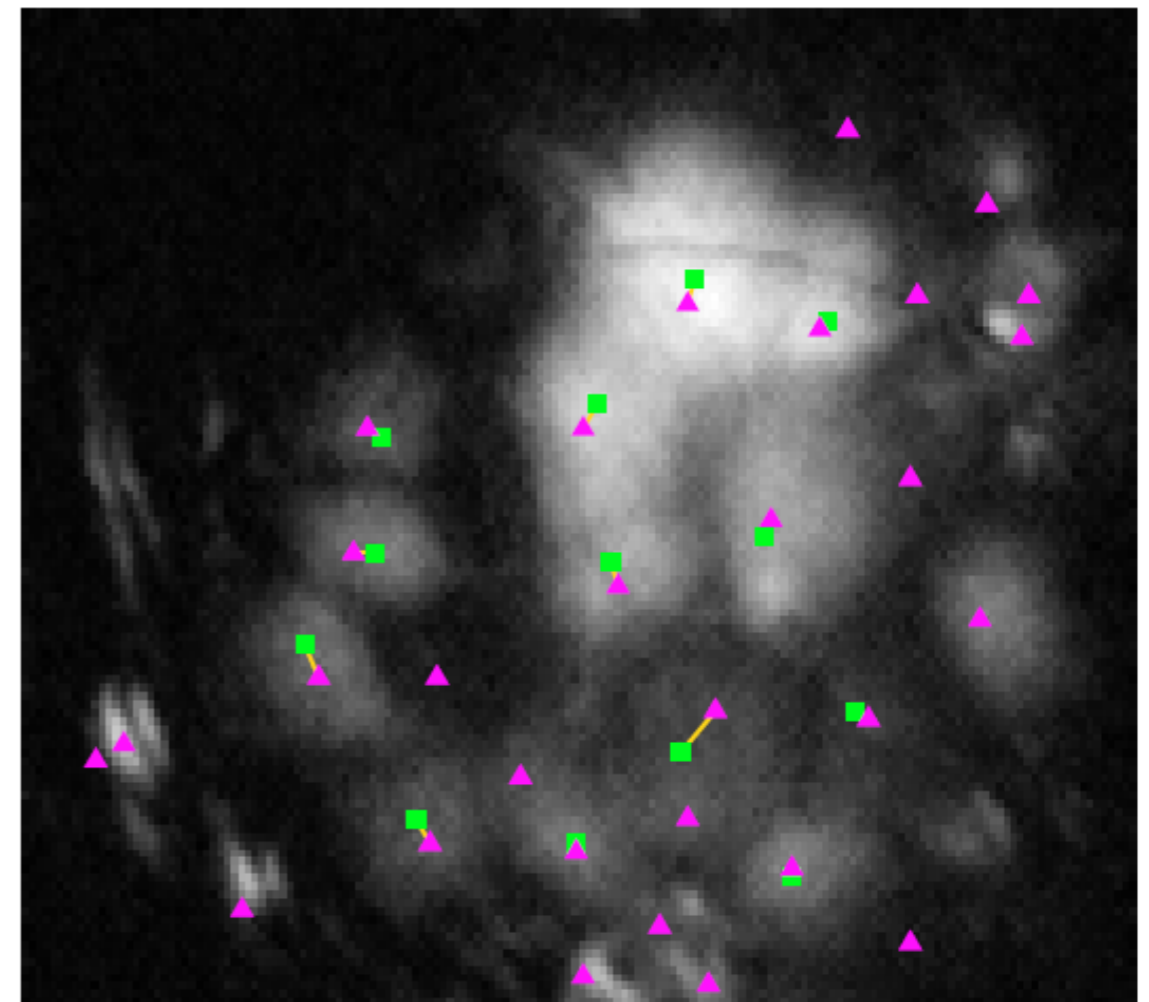
Solving for non-negative X

- Given C , finding non-negative X is the same as with NMF
- Convex optimization with linear constraints
- Or truncated-to-zero pseudo-inverse

Application: Neuroimaging

- Record brain cell activity over time
 - Every row is one frame
- Assume some columns contain the pure glomerulus signal
 - **C** identifies these signals
 - **X** explains how the signals are mixed in the brain images

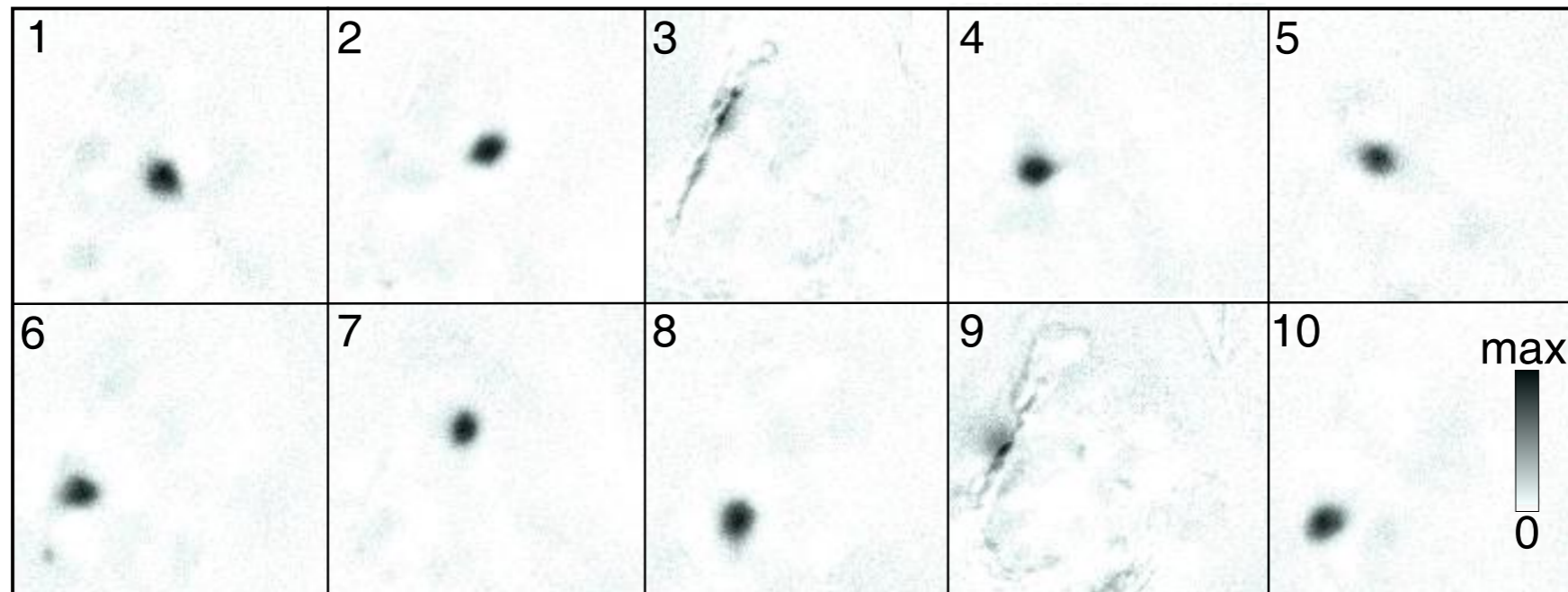
Movie frame with real and found locations marked



■ Human expert

▲ Algorithm

Application cont'd



Top-10 rows of X from NNCX decomposition shows the shape and location of glomeruli

Column-Row Decompositions

The CUR decomposition

- In the **CUR decomposition** we are given matrix **A** and integers c and r , and our task is to select c columns of **A** to matrix **C** and r rows to matrix **R**, and build c -by- r matrix **U** minimizing $\|\mathbf{A} - \mathbf{CUR}\|_F$
- Often $c = r = k$

Why CUR?

- If selecting the actual columns in CX is good, selecting the actual columns and rows must be even better
- We find prototypical columns *and rows*
- \mathbf{U} is usually small, so if \mathbf{C} and \mathbf{R} are sparse, storing CUR takes little space

Solving CUR: general idea

- CUR is two-sided CX
- Simple algorithm idea:
 - Solve CX for \mathbf{A} and \mathbf{A}^T and solve for \mathbf{U} given \mathbf{C} and \mathbf{R}
 - $\mathbf{U} = \mathbf{C}^+ \mathbf{A} \mathbf{R}^+$
- Better algorithms take into account the columns selected to \mathbf{C} when computing \mathbf{R}

Simple CUR algorithm

- Sample columns proportional to their L_2 -norm
- Sample rows proportional to their L_2 -norm
- Build $\mathbf{W} = \mathbf{A}[R,C]$ (the sub-matrix of columns in \mathbf{R} and rows in \mathbf{C})
- Let $\mathbf{W} = \mathbf{X}\mathbf{\Sigma}\mathbf{Y}^T$ be an SVD of \mathbf{W} , and set
$$\mathbf{U} \leftarrow \mathbf{Y}(\mathbf{\Sigma}^+)^2\mathbf{X}^T$$

Fancier CUR algorithm

- Find \mathbf{C} similar to exact- k CX earlier
 - Sample $O(k/\epsilon)$ additional columns
- Find $\mathbf{Z} \in \text{span}(\mathbf{C})$, $\mathbf{Z}^T \mathbf{Z} = \mathbf{I}$, such that $\|\mathbf{A}^T - \mathbf{A}^T \mathbf{Z} \mathbf{Z}^T\|_F \leq (1 + O(\epsilon)) \|\mathbf{A} - \mathbf{C} \mathbf{X}^*\|_F$
 - Use \mathbf{Z} to get the probabilities for sampling $O(k \log k)$ rows of \mathbf{A} and reduce that to $O(k)$ rows
 - Sample $O(k/\epsilon)$ additional rows
- Set $\mathbf{U} = \mathbf{X}^* \mathbf{Z}^T \mathbf{A} \mathbf{R}^+$

Optimal X

Comments on the Boutsidis & Woodruff algorithm

- Slight variations of the above algorithm achieve:
 - selects the smallest number of rows and columns for $(1+\epsilon)$ approximation
 - matrix **U** has the smallest possible rank

CX and CUR summary

- Rows and columns of the original data should be interpretable
 - Also admit local constraints in the data
- CX and CUR decompositions are forms of feature selection
 - Applications when we need “prototypical” rows and columns

Literature

- Drineas, Mahoney & Muthukrishnan (2006): *Subspace sampling and relative-error matrix approximation: Column-based methods*. In APPROX/RANDOM '06
- Boutsidis, Mahoney & Drineas (2010): *An improved approximation algorithm for the column subset selection problem*. arXiv:0812.4293v2
- Boutsidis & Woodruff (2014): *Optimal CUR matrix decompositions*. arXiv:1405.7910v2
- Strauch (2014): *Column subset selection with applications to neuroimaging data*. PhD thesis, U. Konstanz