



CS-570

Statistical Signal Processing

Lecture 7: Dictionary Learning

Spring Semester 2019

Grigorios Tsagkatakis

Today's Objectives

NO CLASS ON WEDNESDAY

Topics:

- Dictionary Learning

Disclaimer: Material used:

Zhang, Zheng, et al. "A survey of sparse representation: algorithms and applications." *IEEE access* 3 (2015): 490-530.

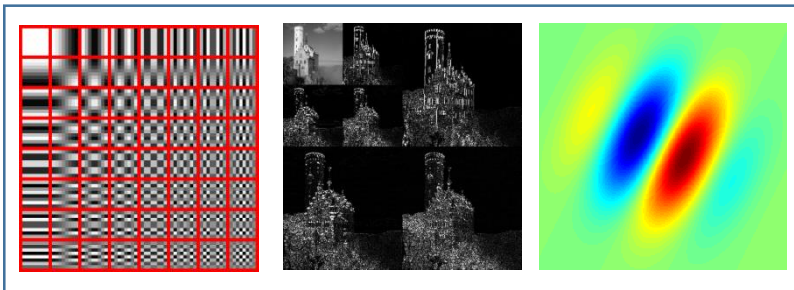


Sparse Signal Modeling

Key idea $\min \|\mathbf{y} - \mathbf{D}\mathbf{s}\|_2 \quad \text{s.t.} \quad \|\mathbf{s}\|_0 \leq K$

Greedy $\|\mathbf{s}\|_1$

Dictionary learning



$$\min \|\mathbf{Y} - \mathbf{D}\mathbf{S}\|_F$$
$$\text{s.t.} \quad \|\mathbf{S}_i\|_1 \leq K, \|\mathbf{D}_i\|_2 \leq 1$$

Transform

- DFT
- DCT
- Wavelets

Learned from examples

- KSVD

A special type of dictionary: two-ortho case

- Motivation for over-complete dictionary: many signals are mixtures of diverse phenomena; no single basis can describe them well

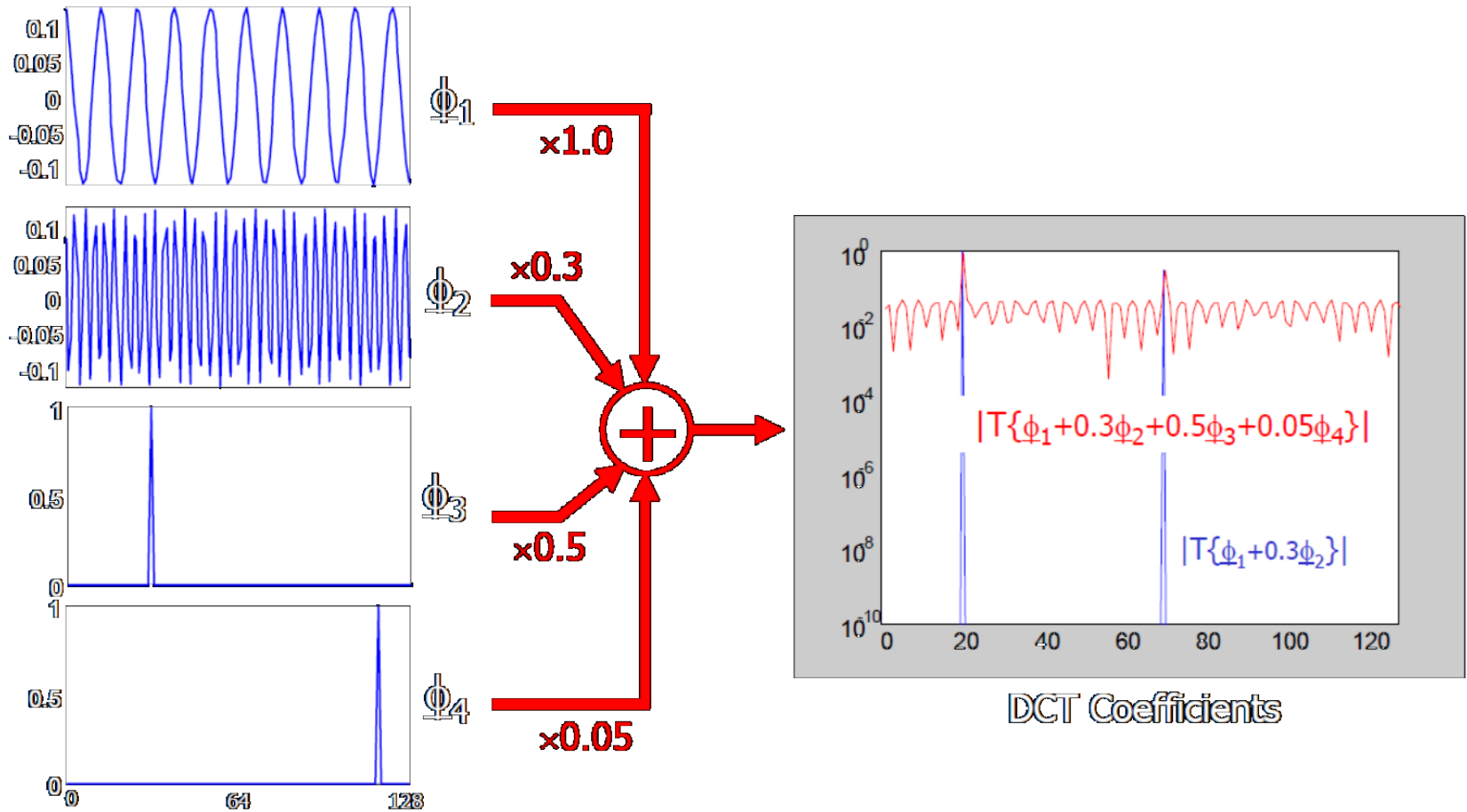
Two-ortho case: A is a concatenation of 2 orthonormal matrices

$$A = [\Psi, \Phi] \quad \text{where } \Psi\Psi^* = \Psi^*\Psi = \Phi\Phi^* = \Phi^*\Phi = I$$

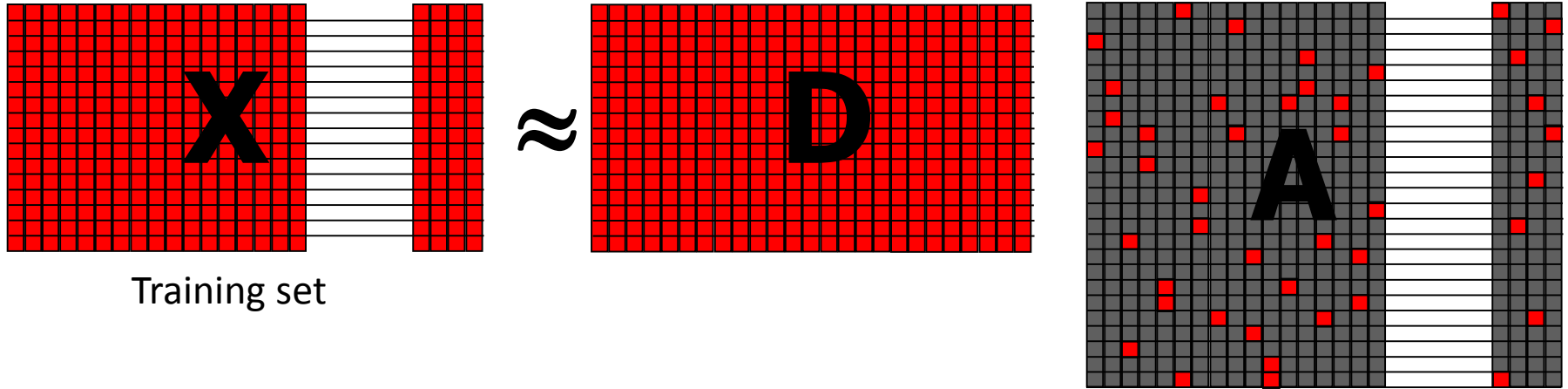
- A classical example: $A = [I, F]$ (F : Fourier matrix)
representing a signal y as a superposition of spikes and sinusoids



Example



Dictionary learning



Training set

$$\text{Min}_{\mathbf{D}, \mathbf{A}} \sum_{j=1}^P \|\mathbf{D}\underline{\alpha}_j - \underline{x}_j\|_2^2 \quad \text{s.t.} \quad \forall j, \|\underline{\alpha}_j\|_0 \leq L$$

Each example is a linear combination of atoms from \mathbf{D}

Each example has a sparse representation with no more than L atoms

Divergence – Matrix Rank

The **rank** of a matrix M is the size of the largest collection of linearly independent columns of M (the **column rank**) or the size of the largest collection of linearly independent *rows of M* (the **row rank**)

- Row Echelon Form

$$\begin{bmatrix} 1 & 2 & 1 \\ -2 & -3 & 1 \\ 3 & 5 & 0 \end{bmatrix} \xrightarrow{R_2 \rightarrow 2r_1+r_2} \begin{bmatrix} 1 & 2 & 1 \\ 0 & 1 & 3 \\ 3 & 5 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 1 \\ 0 & 1 & 3 \\ 3 & 5 & 0 \end{bmatrix} \xrightarrow{R_3 \rightarrow -3r_1+r_3} \begin{bmatrix} 1 & 2 & 1 \\ 0 & 1 & 3 \\ 0 & -1 & -3 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 1 \\ 0 & 1 & 3 \\ 0 & -1 & -3 \end{bmatrix} \xrightarrow{R_3 \rightarrow r_2+r_3} \begin{bmatrix} 1 & 2 & 1 \\ 0 & 1 & 3 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 1 \\ 0 & 1 & 3 \\ 0 & 0 & 0 \end{bmatrix} \xrightarrow{R_1 \rightarrow -2r_2+r_1} \begin{bmatrix} 1 & 0 & -5 \\ 0 & 1 & 3 \\ 0 & 0 & 0 \end{bmatrix} \quad \left. \vphantom{\begin{bmatrix} 1 & 0 & -5 \\ 0 & 1 & 3 \\ 0 & 0 & 0 \end{bmatrix}} \right\} \text{Rank}=2$$

A matrix is in **row echelon form** if

- (i) all nonzero rows are above any rows of all zeroes
- (ii) The leading coefficient of a nonzero row is always strictly to the right of the leading coefficient of the row above it



Matrix Rank

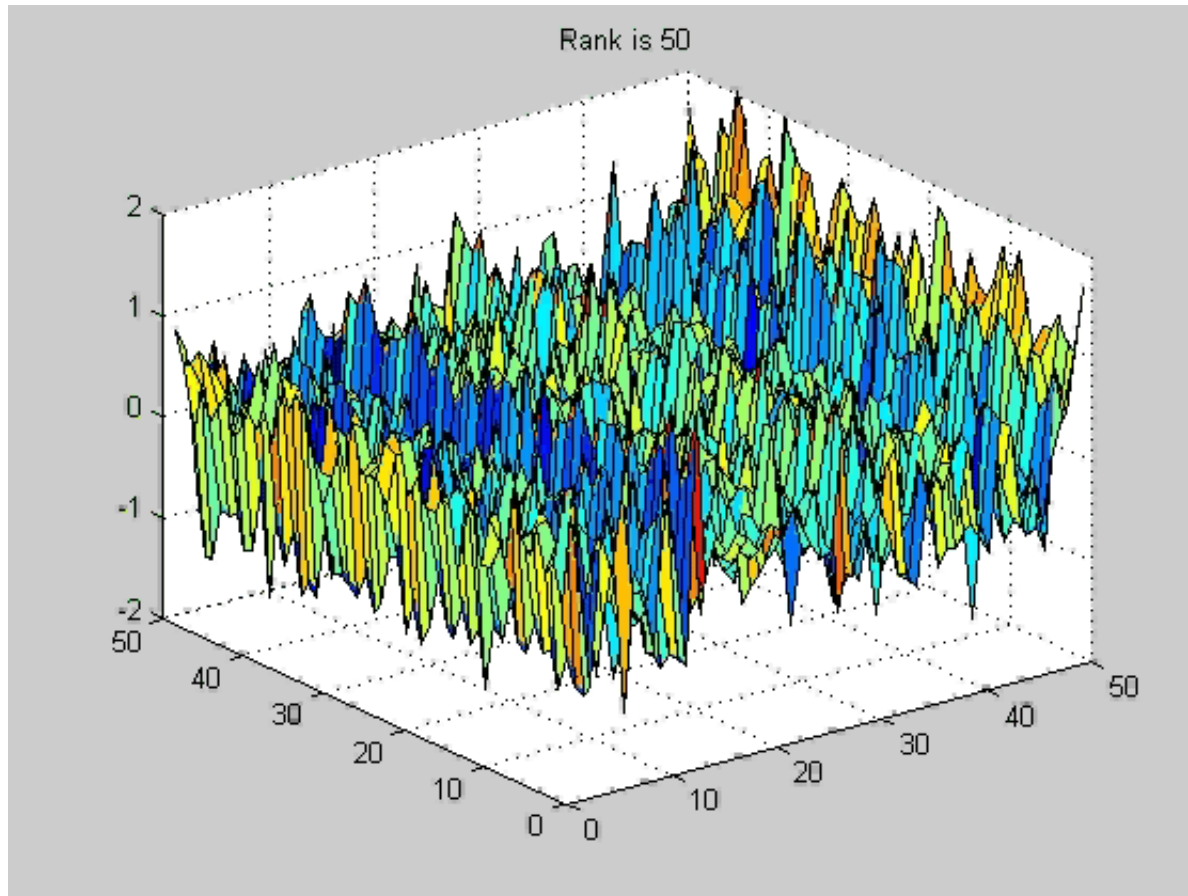
- The rank of an $m \times n$ matrix is a nonnegative integer and cannot be greater than either m or n . That is, $\text{rank}(M) \leq \min(m, n)$.
- A matrix that has a rank as large as possible is said to have **full rank**; otherwise, the matrix is **rank deficient**.

$$\text{rank}(AB) \leq \min(\text{rank } A, \text{rank } B).$$

$$\text{rank}(A^T A) = \text{rank}(A A^T) = \text{rank}(A) = \text{rank}(A^T)$$



Matrix Rank



Singular Value Decomposition (SVD)

Given any $m \times n$ matrix \mathbf{M} , algorithm to find matrices \mathbf{U} , $\mathbf{\Sigma}$, and \mathbf{V} such that $\mathbf{M} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$

- \mathbf{U} : left singular vectors (orthonormal)
- $\mathbf{\Sigma}$: diagonal containing singular values
- \mathbf{V} : right singular vectors (orthonormal)

$$\mathbf{M} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$$

$m \times m$ $m \times n$ V is $n \times n$

$$\begin{pmatrix} M \end{pmatrix} = \begin{pmatrix} U \end{pmatrix} \begin{pmatrix} s_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & s_n \end{pmatrix} \begin{pmatrix} V \end{pmatrix}^T$$



Singular Value Decomposition (SVD)

Properties

- The s_i are called the singular values of \mathbf{M}
- If \mathbf{M} is singular, some of the s_i will be 0
- In general $\text{rank}(\mathbf{M}) = \text{number of nonzero } s_i$
- SVD is mostly unique (up to permutation of SV)



M-term approximation

- SVD can be used to compute optimal **low-rank approximations**.
- Approximation problem: Find \mathbf{A}_k of rank k such that

$$\mathbf{A}_k = \min_{X: \text{rank}(X)=k} \|\mathbf{A} - X\|_F \longleftarrow \text{Frobenius norm}$$

$$\|\mathbf{A}\|_F \equiv \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2}.$$

\mathbf{A}_k and X are both $m \times n$ matrices.

Typically, want $k \ll r$.

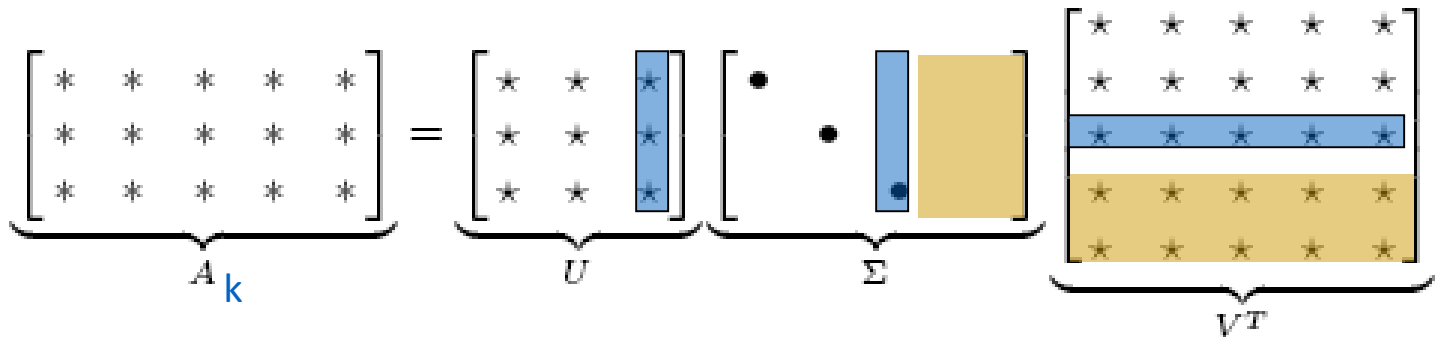


Low-rank Approximation

- Solution via SVD

$$A_k = U \operatorname{diag}(\sigma_1, \dots, \sigma_k, \underbrace{0, \dots, 0}) V^T$$

*set smallest $r-k$
singular values to zero*



$$A_k = \sum_{i=1}^k \sigma_i u_i v_i^T$$

*column notation: **sum**
of rank 1 matrices*

Method of optimal directions (K. Engan and S. Husoy 1999)

$$\min_{D \in \mathcal{D}, X \in \mathcal{X}_\Omega} \|Y - DX\|_F^2.$$

MOD: least squares

1 Fix D , solve for X :

$$\min_{X \in \mathcal{X}_\Omega} \|Y - DX\|_F^2.$$

2 Fix X , solve for D :

$$\min_D \|Y - DX\|_F^2.$$

3 (Optional) Normalization:

$$D_{:,i} = D_{:,i} / \|D_{:,i}\|_2.$$

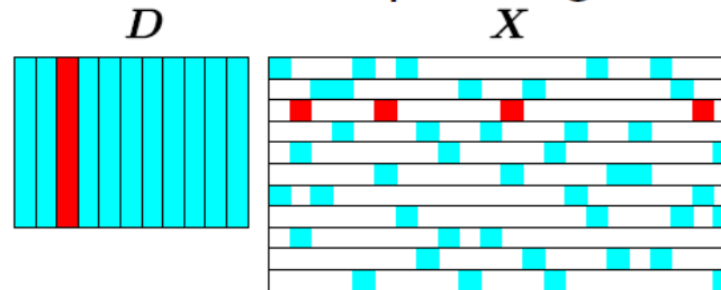


The K-SVD (M. Aharon, et al. 2006)

$$\min_{D \in \mathcal{D}, X \in \mathcal{X}_\Omega} \|Y - DX\|_F^2.$$

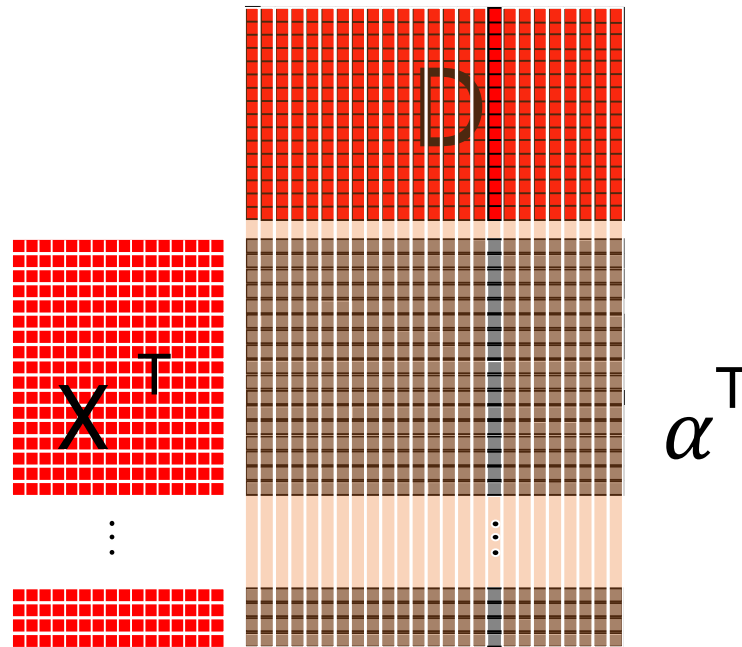
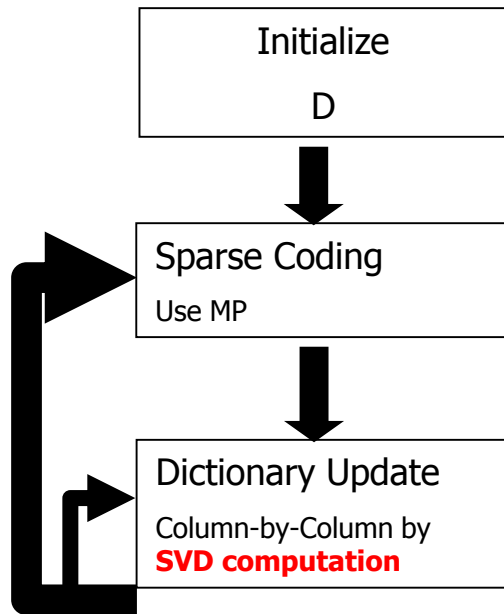
For each column:

Update: this column in D & the corresponding row in X .



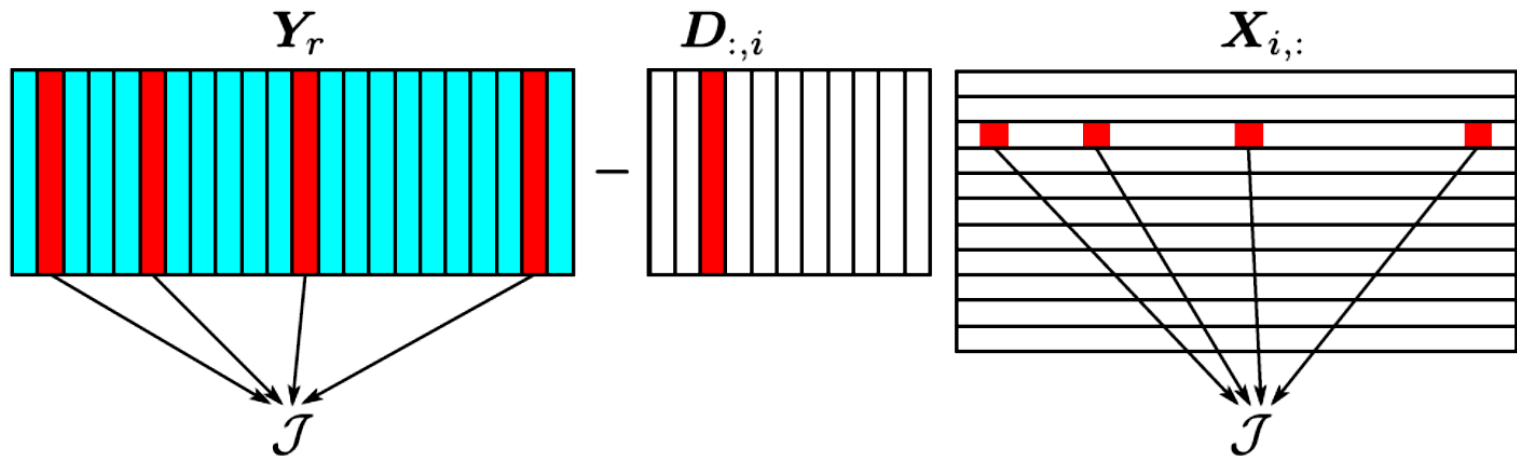
Fix: other columns in D & the corresponding rows in X .

K-SVD algorithm



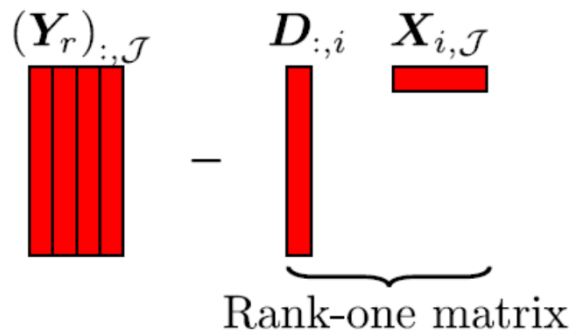
K-SVD details

$$\begin{aligned}
 & \|Y - DX\|^2 \\
 &= \|Y - D_{:,j \neq i} X_{j \neq i,:} - D_{:,i} X_{i,:}\|^2 \\
 &= \|Y_r - D_{:,i} X_{i,:}\|^2 \\
 &= \left\| (Y_r)_{:,J} - D_{:,i} X_{i,J} \right\|^2 + c
 \end{aligned}$$



K-SVD details

$$\begin{aligned}
 & \|Y - DX\|^2 \\
 &= \|Y - D_{:,j \neq i} X_{j \neq i,:} - D_{:,i} X_{i,:}\|^2 \\
 &= \|Y_r - D_{:,i} X_{i,:}\|^2 \\
 &= \left\| (Y_r)_{:, \mathcal{J}} - D_{:,i} X_{i, \mathcal{J}} \right\|^2 + c
 \end{aligned}$$



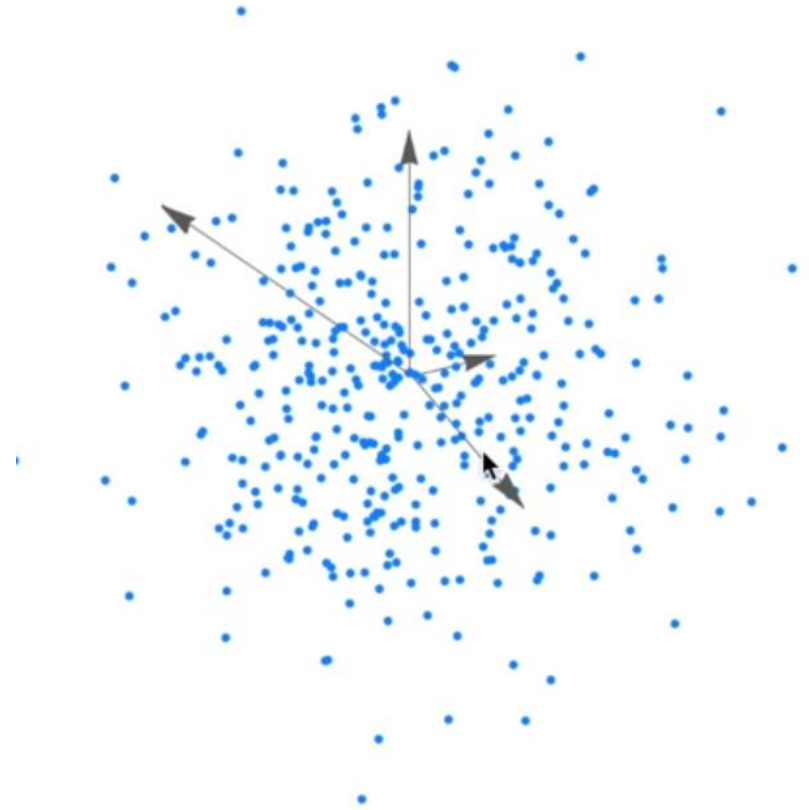
SVD: optimal rank-one matrix approximation.

$$\begin{aligned}
 A &= \sum \lambda_i \mathbf{u}_i \mathbf{v}_i^T & \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \\
 &\approx \lambda_1 \mathbf{u}_1 \mathbf{v}_1^T
 \end{aligned}$$

K-SVD algorithm

Here is three-dimensional data set, spanned by over-complete dictionary of four vectors.

What we want is to update each of these vector to better represent the data.



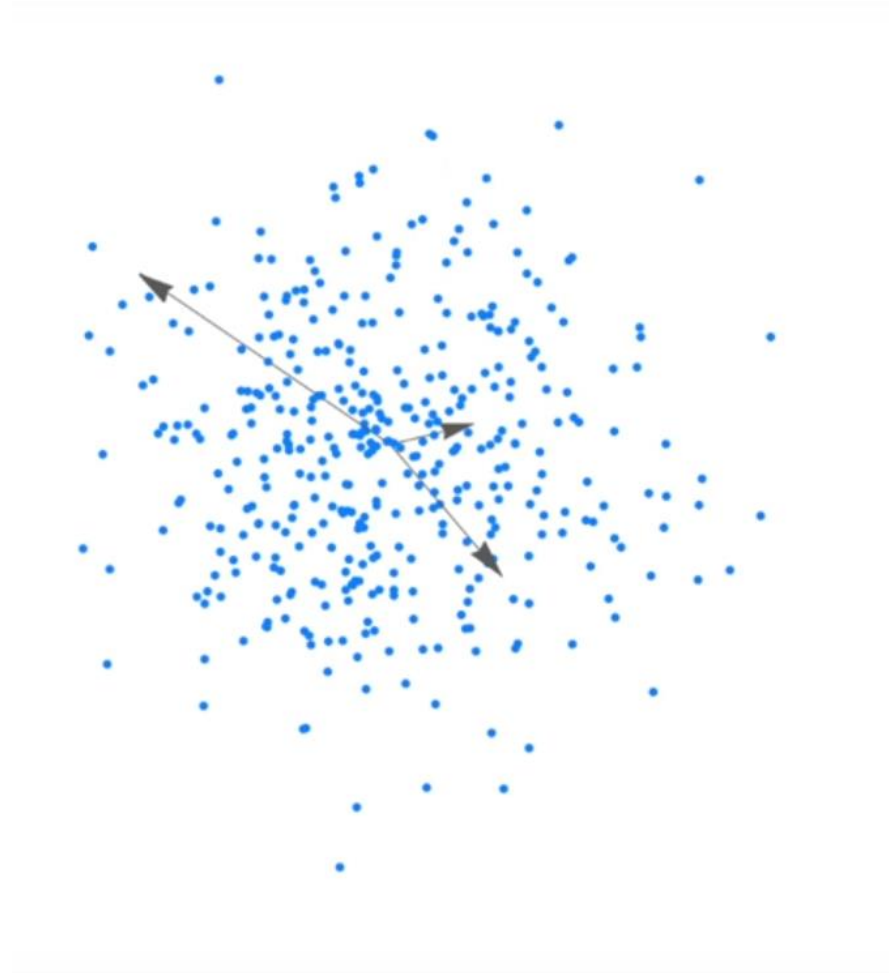
Spring Semester 2019



K-SVD algorithm

1. Remove one of these vector

If we do sparse coding using only three vectors, from the dictionary, we cannot perfectly represent the data.

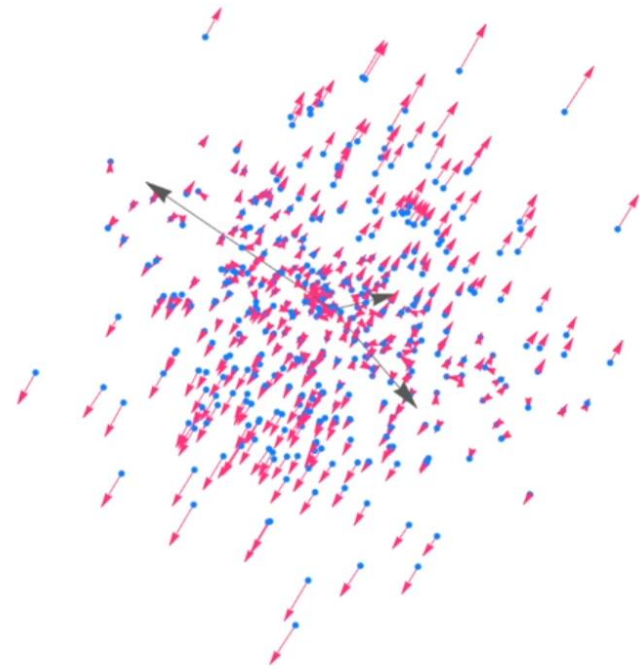


Spring Semester 2019



K-SVD algorithm

2. Find approximation error on each data point

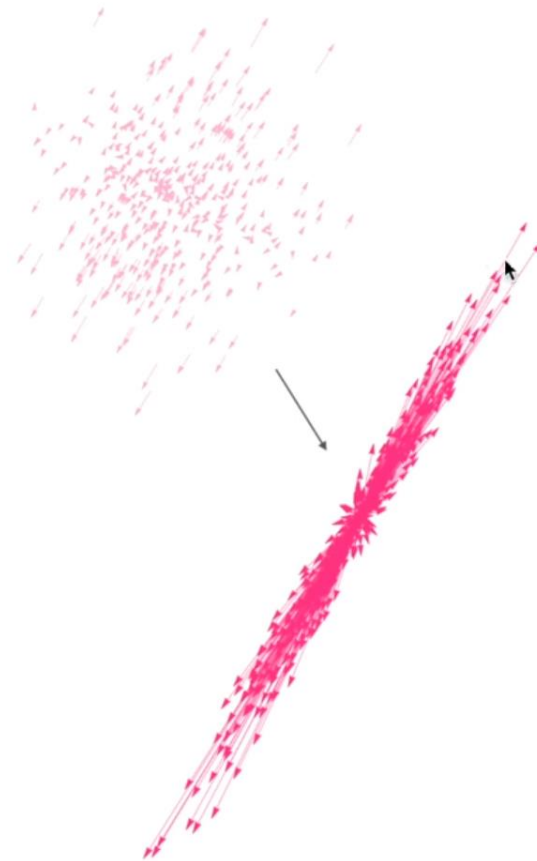


Spring Semester 2019



K-SVD algorithm

2. Find approximation error on each data point



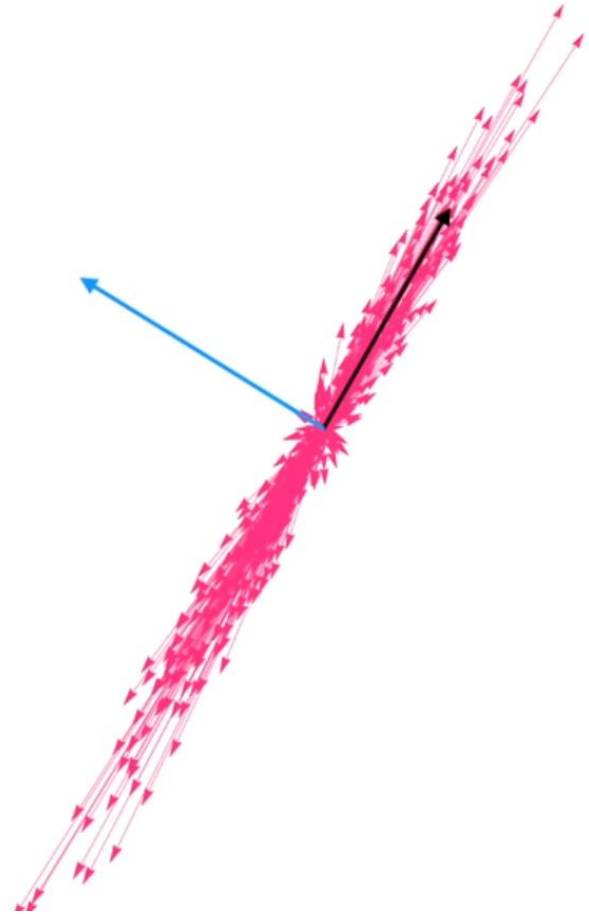
Spring Semester 2019



K-SVD algorithm

3. Apply SVD on error matrix

The SVD provides us a set of orthogonal basis vector sorted in order of decreasing ability to represent the variance error matrix.



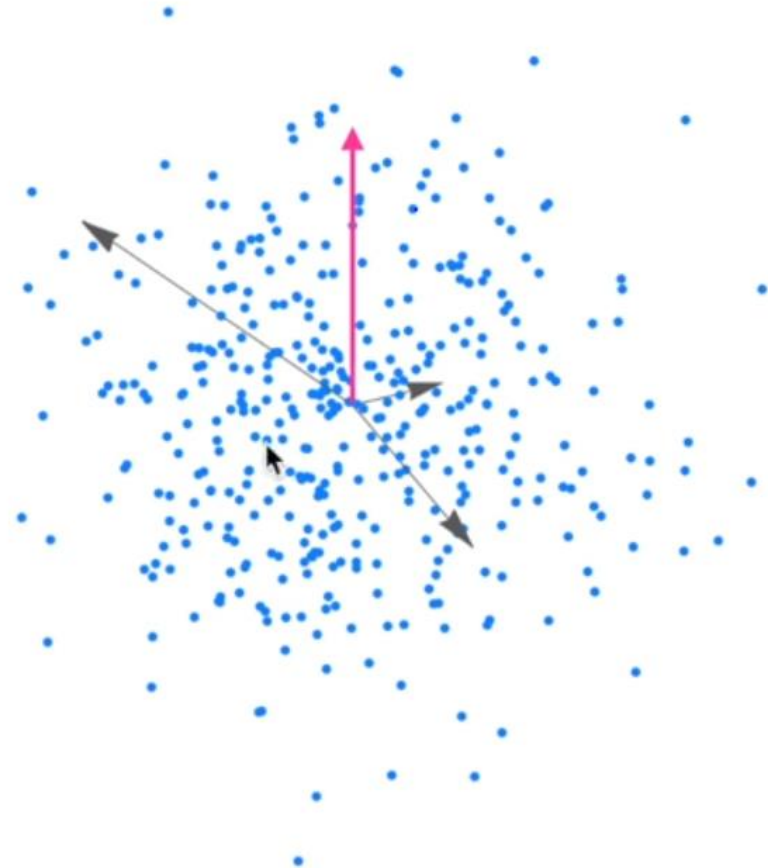
Spring Semester 2019



K-SVD algorithm

3. Replace the chosen vector with the first eigenvector of error matrix.

4. Do the same for other vectors.



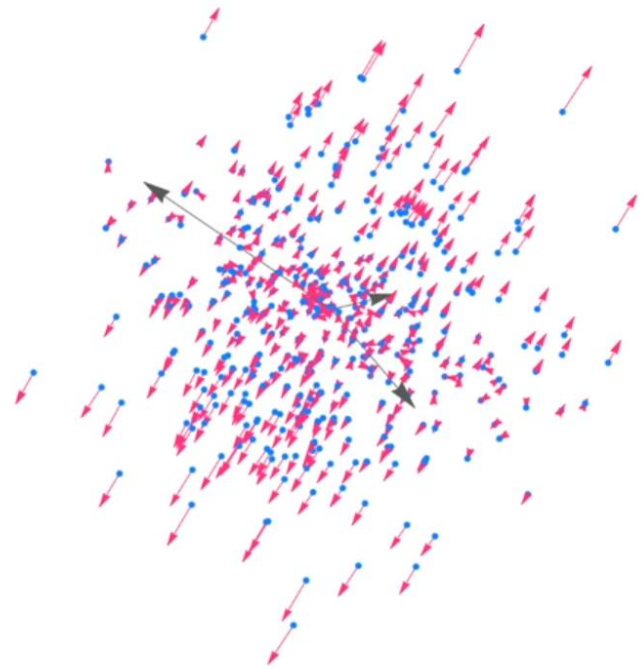
Spring Semester 2019



K-SVD algorithm

But, there is **not all**, but a **few** data points using the chosen vector.

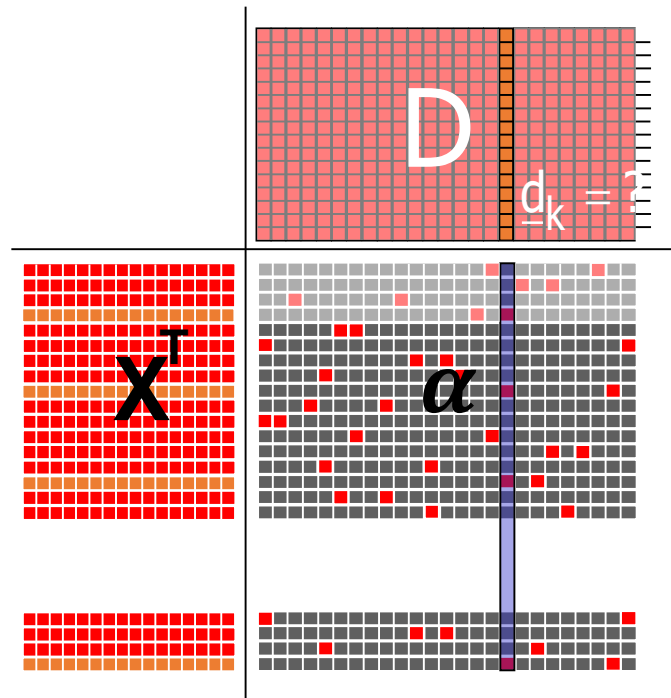
Then, it is not necessary to calculate error for all data points, but instead a few of them that are using the chosen vector.



Spring Semester 2019



K-SVD algorithm



Spring Semester 2019



K-SVD algorithm

1. Initialize the dictionary randomly
2. Using any pursuit algorithm to find a sparse coding α , for the input data X using dictionary D .
3. Update D :
 - a. Remove a basis vector d_k
 - b. Compute the approximation error E_k on data points that were actually using d_k
 - c. Take SVD of E_k
 - d. Update d_k .
4. Repeat to step 2 until convergence.



Example

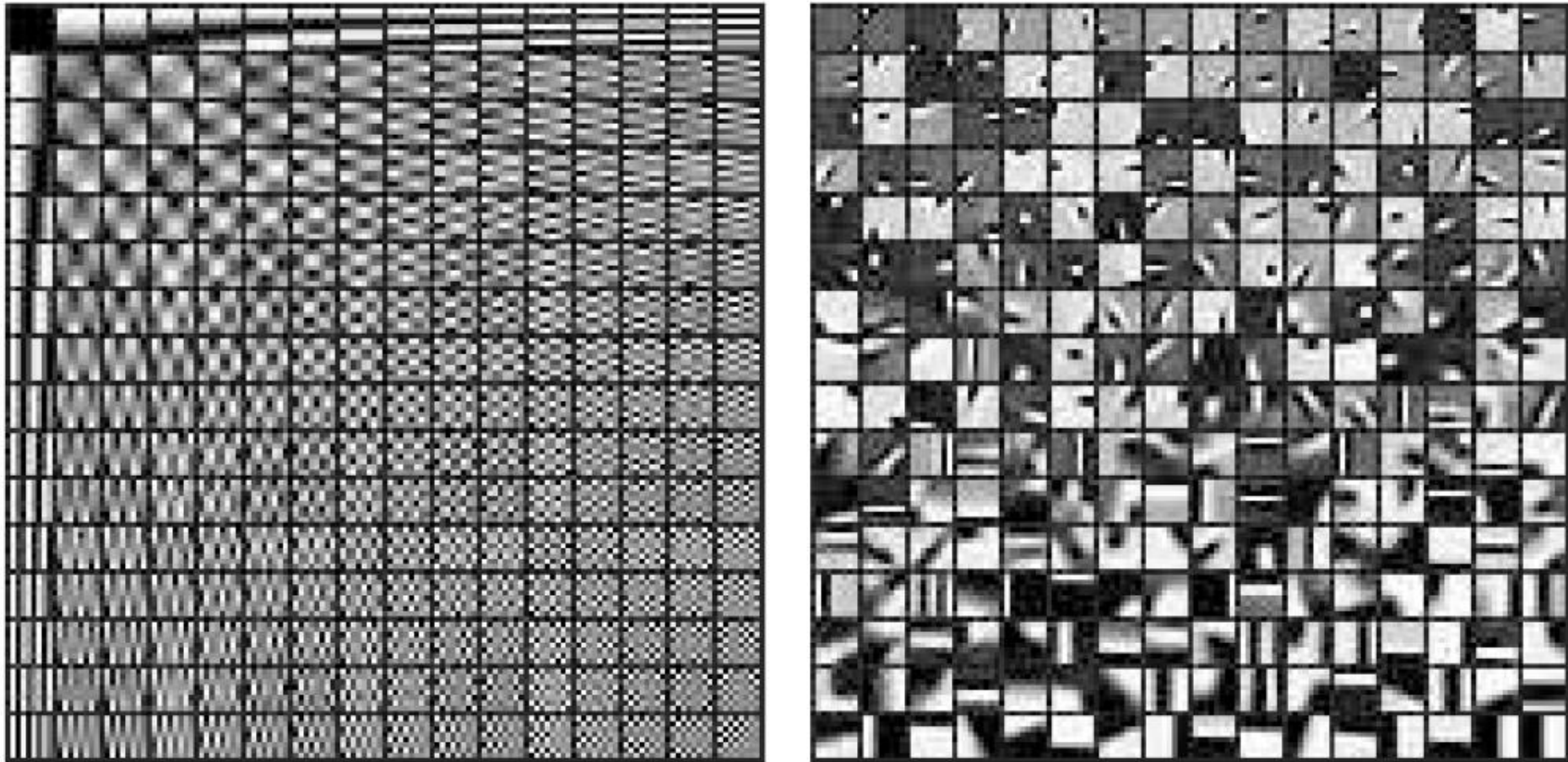
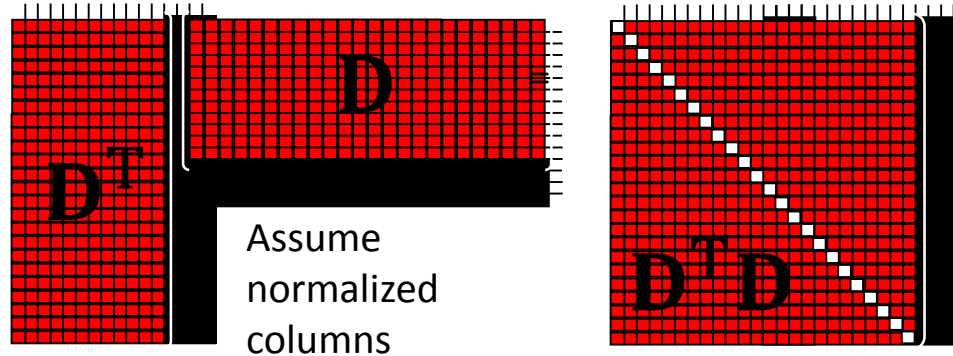


Fig. 2. Left: Overcomplete DCT dictionary. Right: Globally trained dictionary.

The Mutual Coherence

Compute



- The Mutual Coherence $\mu(\mathbf{D})$ is the largest off-diagonal entry in absolute value
- Other ways to characterize the dictionary
 - Restricted Isometry Property - RIP,
 - Exact Recovery Condition - ERC,
 - Spark

Basis pursuit success

Theorem: Given a noisy signal $y = \mathbf{D}\alpha + v$ where $\|v\|_2 \leq \varepsilon$ and α is sufficiently sparse,

$$\|\alpha\|_0 < \frac{1}{4} \left(1 + \frac{1}{\mu} \right)$$

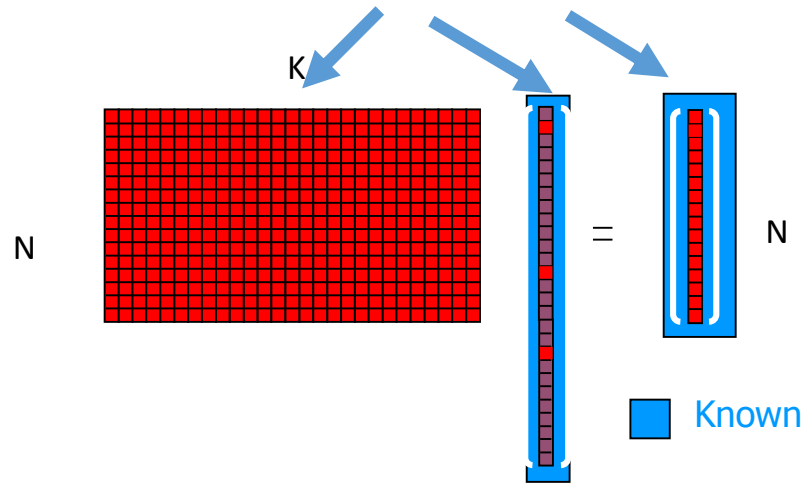
then Basis-Pursuit: $\min_{\alpha} \|\alpha\|_1$ s. t. $\|\mathbf{D}\alpha - y\|_2 \leq \varepsilon$
leads to a stable result: $\|\hat{\alpha} - \alpha\|_2^2 \leq \frac{4\varepsilon^2}{1 - \mu(4\|\alpha\|_0 - 1)}$



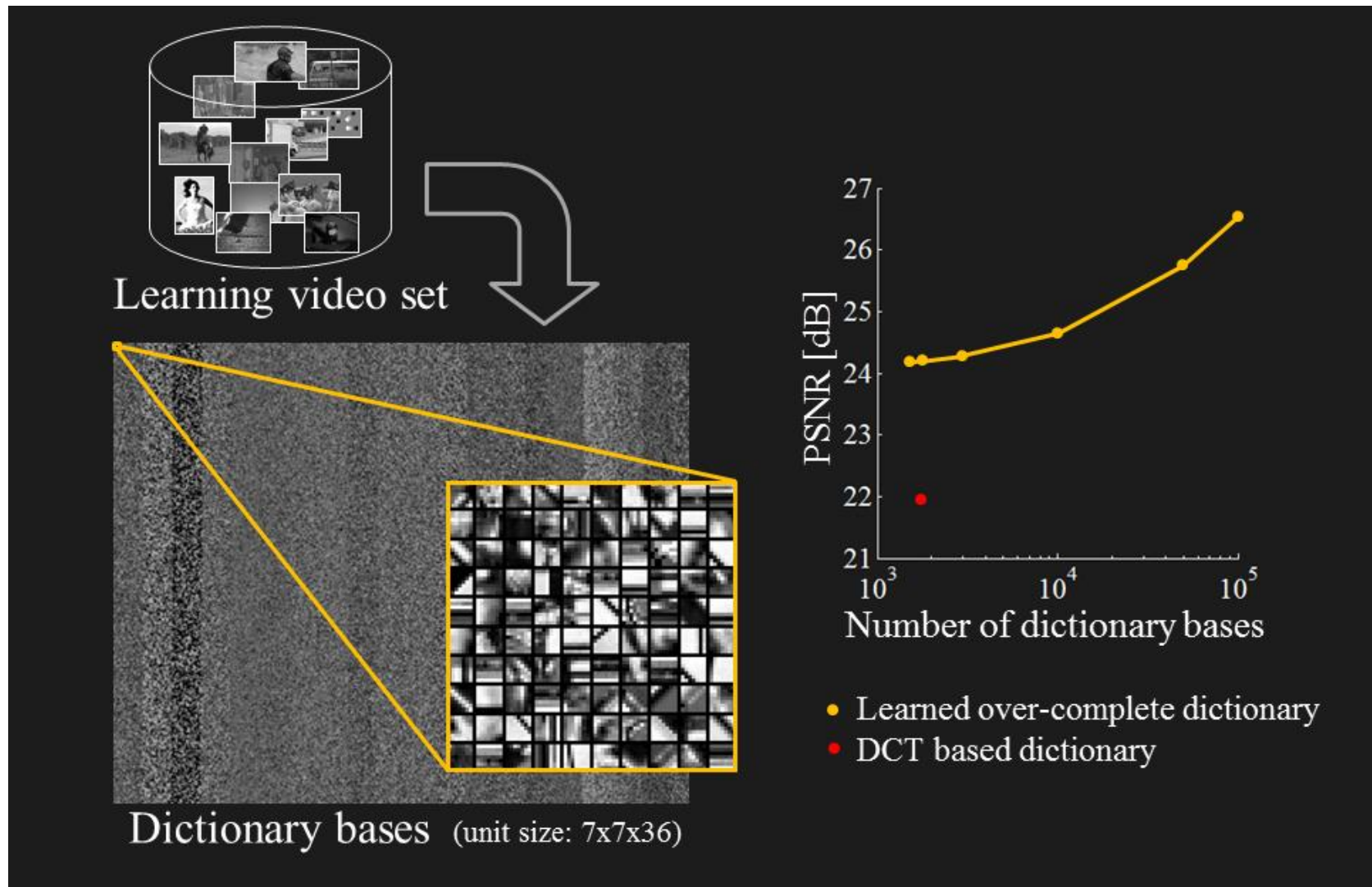
Dictionary Learning

- How to correctly choose the basis for representing the data ?

$$D\alpha = x$$



Accuracy increases with dictionary size



Y. Hitomi, J. Gu, M. Gupta, T. Mitsunaga, and S. K. Nayar. Video from a single coded exposure photograph using a learned over-complete dictionary. In IEEE Intl. Conf. Computer Vision, 2011

Denoising

Image Denoising [E. & Aharon ('06)]



The MAP estimator for denoising this image patch is built by solving

$$\hat{\alpha} = \arg \min_{\alpha} \|\alpha\|_0 \text{ subject to } \|\mathbf{D}\alpha - \mathbf{y}\|_2^2 \leq T \quad \longrightarrow \quad \hat{\mathbf{x}} = \mathbf{D}\hat{\alpha}$$

De-noising

- Learn a **patch dictionary**.
- For each patch, compute the **sparse representation** then use it to **reconstruct** the patch.

$$\mathbf{x}^* = \arg \min_x \|\mathbf{x}\|_1 + \lambda \|\mathbf{Ax} - \mathbf{b}\|_1$$

$$\mathbf{b} = \mathbf{Ax}^*$$



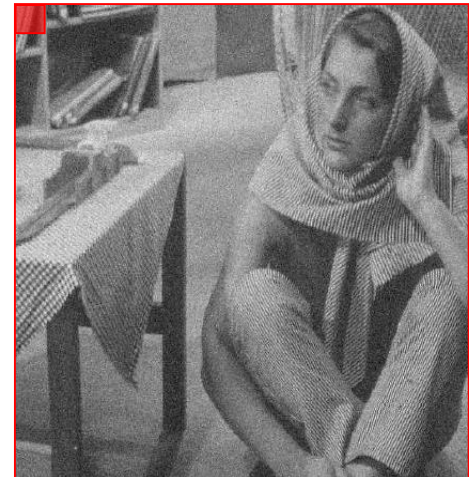
What data to train on?

Option 1:

- ❑ Use a database of images,
- ❑ Works fine (~ 0.5 -1dB below the SotA).

Option 2:

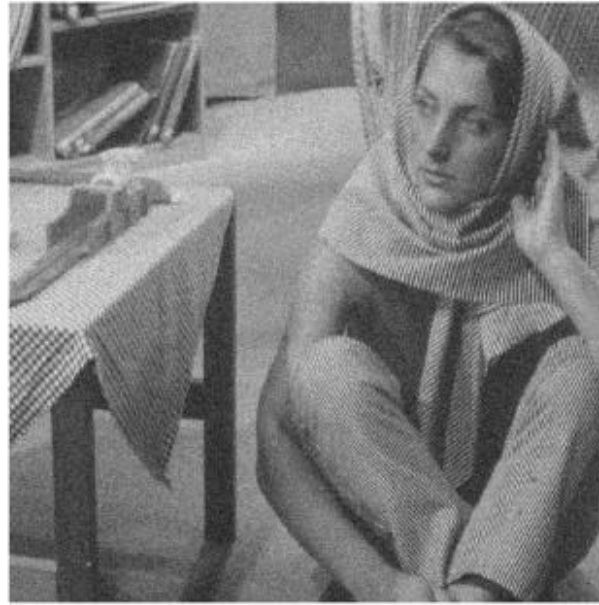
- ❑ Use the corrupted image itself !!
- ❑ Simply sweep through all patches of size N -by- N (overlapping blocks),
- ❑ Image of size 1000^2 pixels $\sim 10^6$ examples to use – more than enough.
- ❑ This works much better!



Original Image



Noisy Image (22.1307 dB, $\sigma=20$)



Denoised Image Using
Global Trained Dictionary (28.8528 dB)



Denoised Image Using
Adaptive Dictionary (30.8295 dB)



S Fig. 6. Example of the denoising results for the image “Barbara” with $\sigma = 20$ —the original, the noisy, and two restoration results.