

Dictionary Learning

CS 754

Problem statement

- Images or small patches from images can often be expressed as a linear combination of columns (or “atoms”) of a matrix.

$$\mathbf{y}_i = \mathbf{D}\mathbf{s}_i = \sum_{k=1}^K \mathbf{d}_k s_{ik}$$

Index for image/image-patch = i

- This matrix is called a **dictionary** and the coefficients of the linear combination are called **dictionary coefficients**.

Problem statement

- We have seen previously that these columns could belong to the DCT or wavelet basis matrix – these are called universal bases.
- But we can go one step beyond – we can infer (or learn) the dictionary from the data. This is called **dictionary learning**.

Problem statement

- The dictionary can be forced to obey a variety of constraints:
 1. it may be orthonormal
 2. it may be overcomplete (i.e. the number of columns exceeds the number of rows)
 3. it may be non-negative.
- The learned dictionary can be useful in a variety of applications – compression, denoising, deblurring, inpainting as well as compressive recovery.

Dictionary learning algorithms

- We will study the following dictionary learning algorithms:
 1. Principal Components Analysis (PCA)
 2. Non-negative sparse coding
 3. Method of optimal directions (MOD)
 4. Union of orthonormal bases
 5. K-SVD

PCA (Principal Components Analysis)

Principal Components Analysis (PCA)

- Consider N vectors (or points) \mathbf{x}_i , each containing d elements, each represented as a **column vector**.
- We say: $\forall i, \mathbf{x}_i \in R^d$. d could be very large – like 250,000 or more.
- Our aim is to extract some k features from each \mathbf{x}_i , $k \ll d$.
- Effectively we are projecting the original vectors from a d -dimensional space to a k -dimensional space. This is called **dimensionality reduction**. PCA is **one method** of dimensionality reduction.

PCA

- How do we pick the ‘ k ’ appropriate features?
- We look into a notion of “compressibility” – how much can the data be compressed, allowing for some small errors.

PCA: Algorithm

1. Compute the mean of the given points:

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i, \mathbf{x}_i \in R^d, \bar{\mathbf{x}} \in R^d$$

2. Deduct the mean from each point:

$$\bar{\mathbf{x}}_i = \mathbf{x}_i - \bar{\mathbf{x}}$$

3. Compute the covariance matrix of these mean-deducted points:

$$\mathbf{C} = \frac{1}{N-1} \sum_{i=1}^N \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T = \frac{1}{N-1} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T, \text{Note : } \mathbf{C} \in R^{d \times d}$$

Note: \mathbf{C} is a symmetric matrix, and it is positive - semidefinite



PCA: algorithm

4. Find the eigenvectors of \mathbf{C} :

$$\mathbf{C}\mathbf{V} = \mathbf{V}\Lambda, \mathbf{V} \in R^{d \times d}, \Lambda \in R^{d \times d}$$

\mathbf{V} – matrix of eigenvectors (each column is an eigenvector),

Λ – diagonal matrix of eigenvalues

Note: \mathbf{V} is an orthonormal matrix (i.e. $\mathbf{V}\mathbf{V}^T = \mathbf{V}^T\mathbf{V} = \mathbf{I}$), as \mathbf{C} is a covariance matrix and hence it is symmetric.

Note: Λ contains non - negative values on the diagonal (eigen - values)

5. Extract the k eigenvectors corresponding to the k largest eigenvalues. This is called the extracted eigenspace:

$$\hat{\mathbf{V}}_k = \mathbf{V}(:,1:k)$$

There is an implicit assumption here that the first k indices indeed correspond to the k largest eigenvalues. If that is not true, you would need to pick the appropriate indices.

PCA: algorithm

6. Project each point onto the eigenspace, giving a vector of k **eigen-coefficients** for that point.

$$\alpha_{ik} = \hat{V}_k^T \bar{x}_i, \alpha_{ik} \in R^k; \alpha_i = V^T \bar{x}_i, \alpha_i \in R^d$$

As V is orthonormal, we have

$$\bar{x}_i = V\alpha_i = V(:,1)\alpha_i(1) + V(:,2)\alpha_i(2) + \dots + V(:,d)\alpha_i(d)$$

$$\approx \hat{V}_k \alpha_{ik} = \hat{V}_k(:,1)\alpha_{ik}(1) + \hat{V}_k(:,2)\alpha_{ik}(2) + \dots + \hat{V}_k(:,k)\alpha_{ik}(k)$$

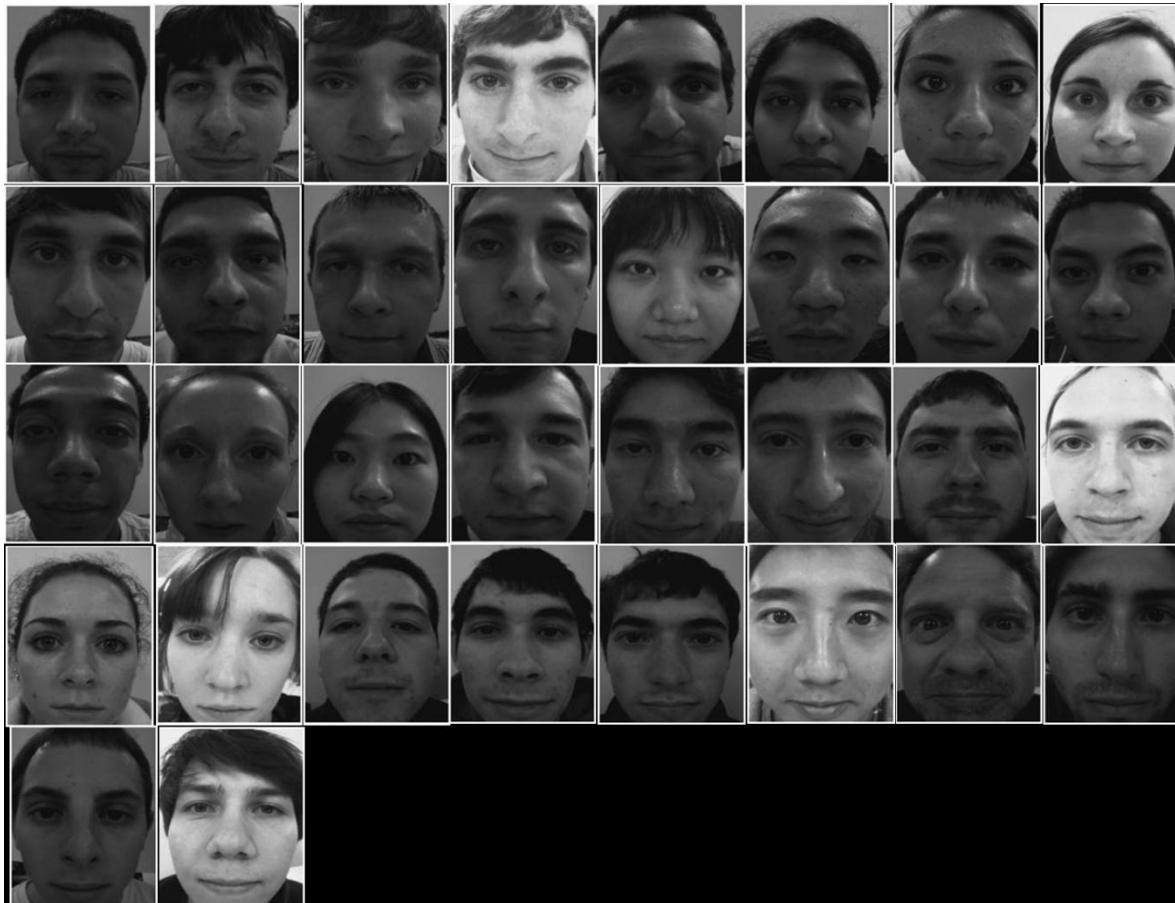
We are representing each face as a linear combination of the k eigenvectors corresponding to the k largest eigenvalues. The coefficients of the linear combination are the eigen-coefficients.

Note that α_{ik} is a vector of the eigencoefficients of the i -th sample point, and it has k elements. The j -th element of this vector is denoted as $\alpha_{ik}(j)$.

PCA and Face Recognition: Eigen-faces

- Consider a database of cropped, frontal face images (which we will assume are aligned and under the same illumination). These are the **gallery images**.
- We will reshape each such image (a 2D array of size $H \times W$ after cropping) to form a column vector of $d = HW$ elements. Each image will be a vector \mathbf{x}_i , as per the notation on the previous two slides.
- And then carry out the six steps mentioned before.
- The eigenvectors that we get in this case are called **eigenfaces**. Each eigenvector has d elements. If you reshape those eigenvectors to form images of size $H \times W$, those images **look like (filtered!) faces**.

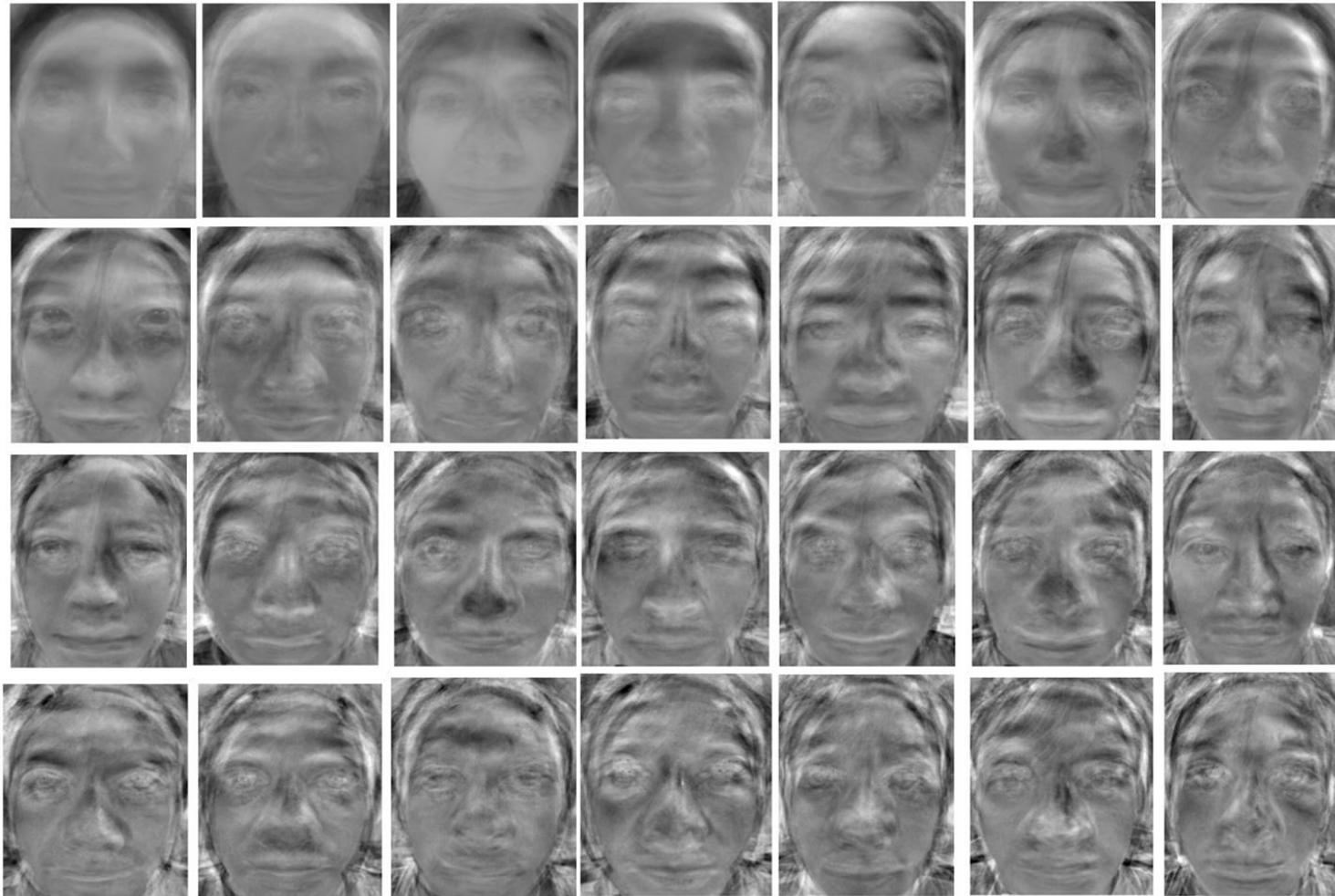
Example 1



A face database

http://people.ece.cornell.edu/land/courses/ece4760/FinalProjects/s2011/bjh78_caj65/bjh78_caj65/

Top 25 Eigen-faces for this database!



http://people.ece.cornell.edu/land/courses/ece4760/FinalProjects/s2011/bjh78_caj65/bjh78_caj65/

PCA and Face recognition: Eigenfaces

- For each gallery image, you compute the eigen-coefficients. You then store the eigen-coefficients and the identity of the person in a database.
- You also store \hat{V}_k, \bar{x} in the database.
- During the testing phase, you are given a probe image (say) z_p in the form of a column vector of HW elements.
- You deduct the mean image from z_p :

$$\bar{z}_p = z_p - \bar{x}$$

PCA and Face recognition: Eigenfaces

- You then project the mean-deducted face image onto the eigen-space:

$$\alpha_p = \hat{V}_k^T \bar{z}_p \quad \xrightarrow{\text{blue arrow}} \text{Eigen-coefficients of the probe image } z_p.$$

- Now, compare α_p with all the α_{ik} (eigen-coefficients of the gallery images) in the database.
- Find the closest match in terms of the squared distance between the eigen-coefficients. That gives you the identity (see next slide).

PCA and Face recognition: Eigenfaces

$$j_p = \arg \min_l \|\alpha_p - \alpha_l\|_2^2$$

↓ ↓
Eigen-coefficients Eigen-coefficients
of the probe of the l -th gallery
image \mathbf{z}_p . image \mathbf{x}_l .

Note: other distance measures (different from sum of squared differences) may also be employed. One example is sum of absolute differences, given as follows: $\|\alpha_p - \alpha_l\|_1$

Another could be normalized dot product (and this distance measure should be maximized!):
$$\frac{\alpha_p \cdot \alpha_l}{\|\alpha_p\|_2 \|\alpha_l\|_2}$$

PCA and Face recognition: eigenfaces

- The eigen-face images contain more and more high frequency information as the corresponding eigen-values decrease.
- Although PCA is a technique known for a long time, its application in face recognition was pioneered by Turk and Pentland in a classic paper in 1991.

M. Turk and A. Pentland (1991). Eigenfaces for recognition, *Journal of Cognitive Neuroscience*, 3(1): 71–86.

PCA and Face recognition: eigenfaces

- We can regard the k eigenfaces as **key signatures**.
- We express each face image as a **linear combination** of these eigenfaces, i.e. the average face + (say) 3 times eigenface 1 + (say) 5 times eigenface 2 + (say) -1 times eigenface 3 and so on. (note: 3,5,-1 here are the eigen-coefficients, and some of them can be negative).

$$\bar{\mathbf{x}}_i = \mathbf{V}\mathbf{a}_i = \mathbf{V}(:,1)\mathbf{a}_i(1) + \mathbf{V}(:,2)\mathbf{a}_i(2) + \dots + \mathbf{V}(:,d)\mathbf{a}_i(d)$$

$$\approx \hat{\mathbf{V}}_k \mathbf{a}_{ik} = \hat{\mathbf{V}}_k(:,1)\mathbf{a}_{ik}(1) + \hat{\mathbf{V}}_k(:,2)\mathbf{a}_{ik}(2) + \dots + \hat{\mathbf{V}}_k(:,d)\mathbf{a}_{ik}(k)$$

One word of caution: Eigen-faces

- The algorithm described earlier is computationally infeasible for face images (or eigenfaces), as it requires storage of a $d \times d$ Covariance matrix (d – the number of image pixels - could be more than 10,000). And the computation of the eigen-vectors of such a matrix is a $O(d^3)$ operation!
- We will study a modification to this that will bring down the computational cost drastically.

Eigen-faces: reducing computational complexity.

- Consider the covariance matrix:

$$\mathbf{C} = \frac{1}{N-1} \sum_{i=1}^N \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T = \frac{1}{N-1} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T, \text{ Note: } \mathbf{C} \in R^{d \times d}$$

Note: \mathbf{C} is a symmetric matrix, and it is positive - semidefinite

- It will require too much memory if d is large, and computing its eigenvectors will be a horrendous task!
- Consider the case when N is much less than d . This is very common in face recognition applications. The number of training images is usually much smaller than the size of the image.

Eigen-faces: reducing computational complexity.

- In such a case, the rank of \mathbf{C} is at the most $N-1$. So \mathbf{C} will have at the most $N-1$ non-zero eigenvalues.



- We can write \mathbf{C} in the following way:

$$\mathbf{C} = \frac{1}{N-1} \sum_{i=1}^N \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^T \propto \mathbf{X} \mathbf{X}^T, \text{ where}$$

$$\mathbf{X} = [\bar{\mathbf{x}}_1 \mid \bar{\mathbf{x}}_2 \mid \dots \mid \bar{\mathbf{x}}_N] \in R^{d \times N}$$

Back to Eigen-faces: reducing computational complexity.

- Consider the matrix $\mathbf{X}^T \mathbf{X}$ (size $N \times N$) instead of $\mathbf{X} \mathbf{X}^T$ (size $d \times d$). Its eigenvectors are of the form:

$$\mathbf{X}^T \mathbf{X} \mathbf{w} = \lambda \mathbf{w}, \mathbf{w} \in R^N$$

$$\rightarrow \mathbf{X} \mathbf{X}^T (\boxed{\mathbf{X} \mathbf{w}}) = \lambda (\boxed{\mathbf{X} \mathbf{w}}) \text{ [pre - multiplying by } \mathbf{X}]$$



$\mathbf{X} \mathbf{w}$ is an eigenvector of $\mathbf{C} = \mathbf{X} \mathbf{X}^T$! Computing all eigenvectors of \mathbf{C} will now have a complexity of only $O(N^3)$ for computation of the eigenvectors of $\mathbf{X}^T \mathbf{X} + O(N \times dN)$ for computation of $\mathbf{X} \mathbf{w}$ from each \mathbf{w} = total of $O(N^3 + dN^2)$ which is much less than $O(d^3)$. Note that \mathbf{C} has at most only $\min(N-1, d)$ eigenvectors corresponding to non-zero eigen-values (why?).

Eigenfaces: Algorithm ($N \ll d$ case)

1. Compute the mean of the given points:

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i, \mathbf{x}_i \in R^d, \bar{\mathbf{x}} \in R^d$$

2. Deduct the mean from each point:

$$\bar{\mathbf{x}}_i = \mathbf{x}_i - \bar{\mathbf{x}}$$

3. Compute the following matrix:

$$\mathbf{L} = \mathbf{X}^T \mathbf{X}, \mathbf{L} \in R^{N \times N}, \mathbf{X} = [\bar{\mathbf{x}}_1 | \bar{\mathbf{x}}_2 | \dots | \bar{\mathbf{x}}_N] \in R^{d \times N}$$

Note : \mathbf{L} is a symmetric matrix, and it is positive - semidefinite

Eigen-faces: Algorithm ($N \ll d$ case)

4. Find the eigenvectors of \mathbf{L} :

$$\mathbf{LW} = \mathbf{W}\Gamma, \mathbf{W} - \text{eigenvectors}, \Gamma - \text{eigenvalues},$$

$$\mathbf{WW}^T = \mathbf{I}$$

5. Obtain the eigenvectors of \mathbf{C} from those of \mathbf{L} :

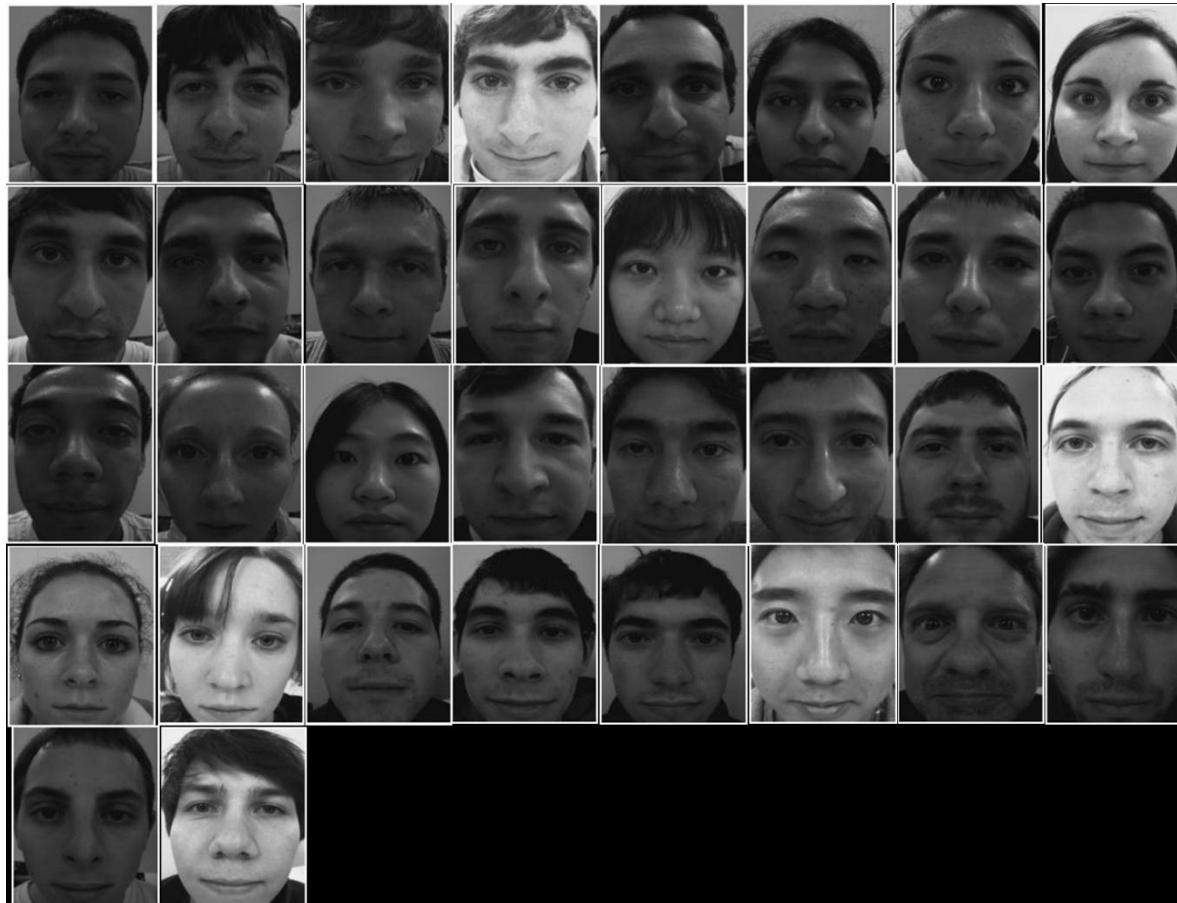
$$\mathbf{V} = \mathbf{XW}, \mathbf{X} \in R^{d \times N}, \mathbf{W} \in R^{N \times N}, \mathbf{V} \in R^{N \times N}$$

6. Unit-normalize the columns of \mathbf{V} .

7. \mathbf{C} will have at most only N eigenvectors corresponding to non-zero eigen-values*. Out of these you pick the top k ($k < N$) corresponding to the largest eigen-values.

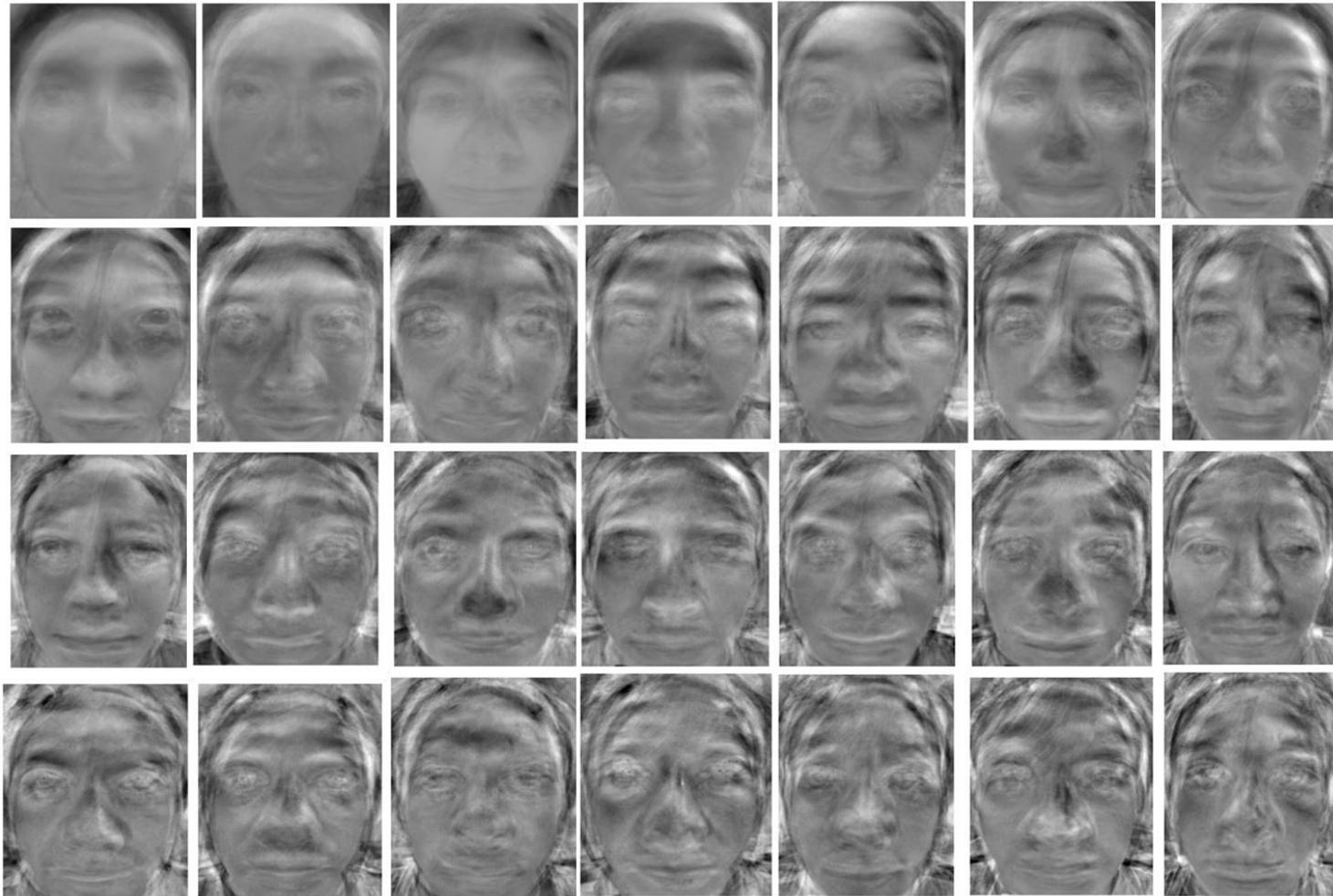
* Actually this number is at most $N-1$ – this is due to the mean subtraction, else it would have been at most N .

Example 1



A face database

Top 25 Eigen-faces for this database!



http://people.ece.cornell.edu/land/courses/ece4760/FinalProjects/s2011/bjh78_caj65/bjh78_caj65/

Example 2



The Yale Face database



Top 25 eigenfaces from
the previous database



Reconstruction of a face image
using the top $1, 8, 16, 32, \dots, 104$
eigenfaces (i.e. k varied from 1 to
104 in steps of 8)

$$\mathbf{x}_i = \bar{\mathbf{x}} + \hat{\mathbf{V}}\mathbf{a}_i \approx \bar{\mathbf{x}} + \sum_{l=1}^k \hat{\mathbf{V}}(:, l)\mathbf{a}_{ik}(l)$$

What if both N and d are large?

- This can happen, for example, if you wanted to build an eigenspace for face images of all people in Mumbai.
- Divide people into coherent groups based on some visual attributes (eg: gender, age group etc) and build separate eigenspaces for each group.

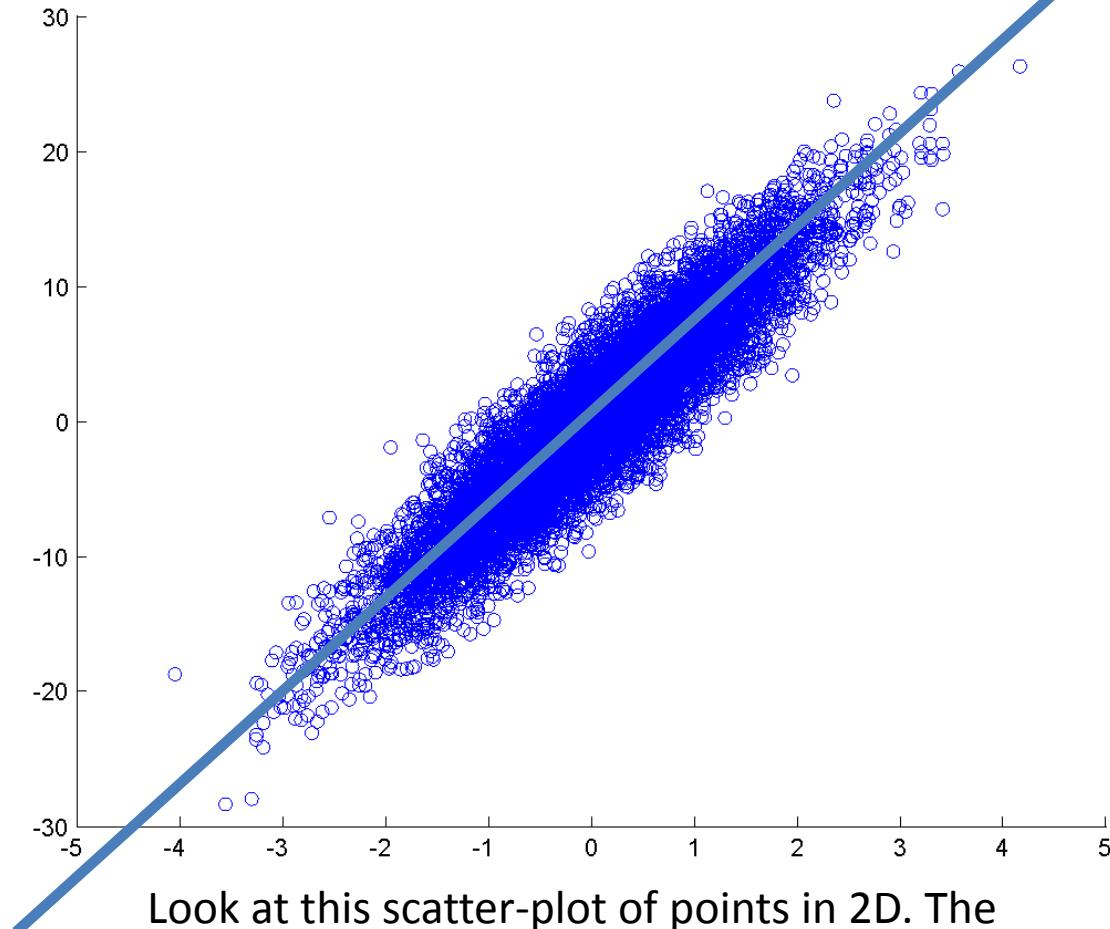
PCA: A closer look

- PCA has many applications – apart from face/object recognition – in image processing/computer vision, statistics, econometrics, finance, agriculture, and you name it!
- Why PCA? What's special about PCA? See the next slides!

PCA: what does it do?

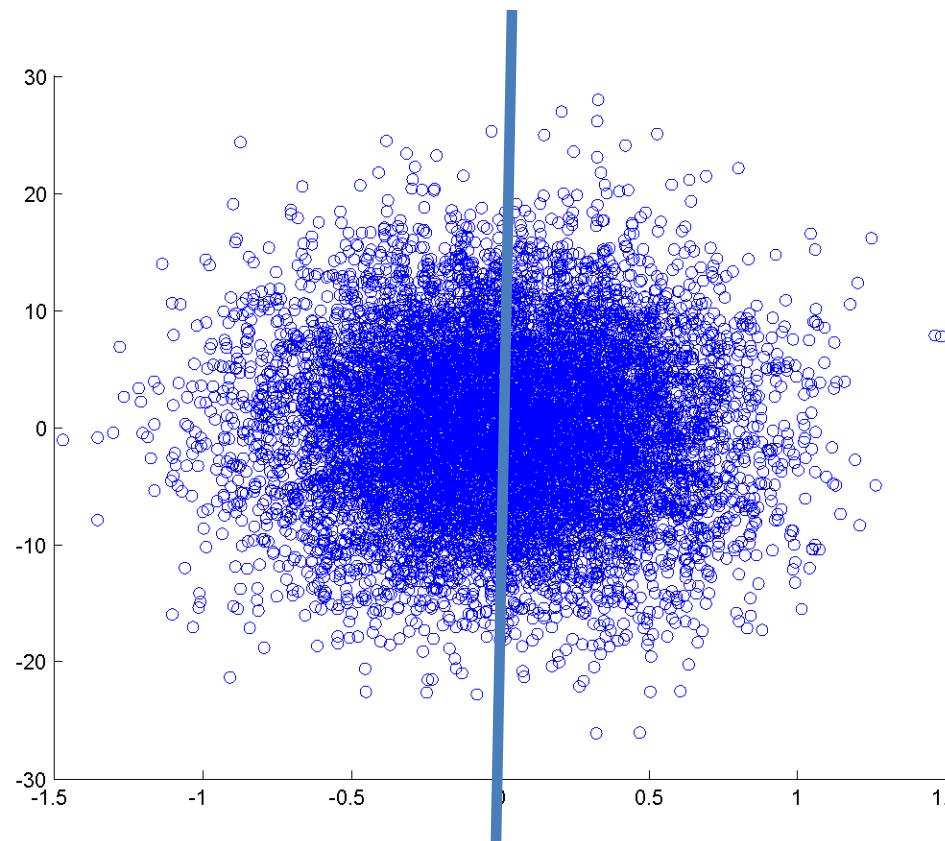
- It finds ' k ' perpendicular directions (all passing through the mean vector) such that the original data are approximated as **accurately** as possible when projected onto these ' k ' directions.
- We will see soon why these ' k ' directions are eigenvectors of the covariance matrix of the data!

PCA



Look at this scatter-plot of points in 2D. The points are highly spread out in the direction of the light blue line.

PCA



This is how the data would look if they were rotated in such a way that the major axis of the ellipse (the light blue line) now coincided with the Y axis. As the spread of the X coordinates is now relatively insignificant (**observe the axes!**), we can approximate the rotated data points by their projections onto the Y-axis (i.e. their Y coordinates alone!). This was not possible prior to rotation!

PCA

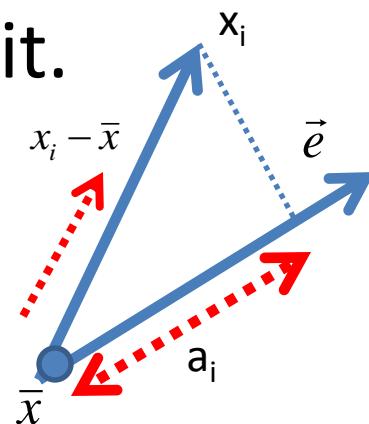
- As we could ignore the X-coordinates of the points post rotation and represent them just by the Y-coordinates, we have performed some sort of lossy data compression – or dimensionality reduction.
- The job of PCA is to perform such a rotation as shown on the previous two slides!

PCA

- **Aim of PCA:** Find the line \vec{e} passing through the sample mean (i.e. \bar{x}), such that the projection of any mean-deducted point $x_i - \bar{x}$ onto \vec{e} , most accurately approximates it.

Projection of $x_i - \bar{x}$ onto \vec{e} is $= a_i \vec{e}, a_i \in R,$

$$a_i = \vec{e}^T (x_i - \bar{x})$$



$$\text{Error of approximation} = \| a_i \vec{e} - (x_i - \bar{x}) \|^2$$

Note: Here \vec{e} is a **unit vector**.

PCA

- Summing up over all points, we get:

$$\begin{aligned} \text{Sum total error of approximation} &= J(\vec{\mathbf{e}}) = \sum_{i=1}^N \| (a_i \vec{\mathbf{e}}) - (\mathbf{x}_i - \bar{\mathbf{x}}) \|^2 \\ &= \sum_{i=1}^N \| (a_i \vec{\mathbf{e}}) \|^2 + \sum_{i=1}^N \| \mathbf{x}_i - \bar{\mathbf{x}} \|^2 - 2 \sum_{i=1}^N a_i \vec{\mathbf{e}}^T (\mathbf{x}_i - \bar{\mathbf{x}}) \\ &= \sum_{i=1}^N \| (a_i \vec{\mathbf{e}}) \|^2 + \sum_{i=1}^N \| \mathbf{x}_i - \bar{\mathbf{x}} \|^2 - 2 \sum_{i=1}^N a_i \vec{\mathbf{e}}^T (\mathbf{x}_i - \bar{\mathbf{x}}) \\ &= \sum_{i=1}^N a_i^2 + \sum_{i=1}^N \| \mathbf{x}_i - \bar{\mathbf{x}} \|^2 - 2 \sum_{i=1}^N a_i \vec{\mathbf{e}}^T (\mathbf{x}_i - \bar{\mathbf{x}}) \\ &= \sum_{i=1}^N a_i^2 + \sum_{i=1}^N \| \mathbf{x}_i - \bar{\mathbf{x}} \|^2 - 2 \sum_{i=1}^N a_i^2, (\because a_i = \vec{\mathbf{e}}^T (\mathbf{x}_i - \bar{\mathbf{x}})) \\ &= -\sum_{i=1}^N a_i^2 + \sum_{i=1}^N \| \mathbf{x}_i - \bar{\mathbf{x}} \|^2 \end{aligned}$$

PCA

Sum total error of approximation = $J(\vec{e}) = -\sum_{i=1}^N a_i^2 + \sum_{i=1}^N \|\mathbf{x}_i - \bar{\mathbf{x}}\|^2$

$$= -\sum_{i=1}^N (\vec{e}^t (\mathbf{x}_i - \bar{\mathbf{x}}))^2 + \sum_{i=1}^N \|\mathbf{x}_i - \bar{\mathbf{x}}\|^2$$

This term is proportional to the variance of the data points when projected onto the direction \mathbf{e} .

$$= -\sum_{i=1}^N \vec{e}^t (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^t \vec{e} + \sum_{i=1}^N \|\mathbf{x}_i - \bar{\mathbf{x}}\|^2$$

$$= -\vec{e}^t \mathbf{S} \vec{e} + \sum_{i=1}^N \|\mathbf{x}_i - \bar{\mathbf{x}}\|^2 \quad (\text{where } \mathbf{S} = (N-1)\mathbf{C})$$

PCA

$$J(\vec{e}) = -\vec{e}^t \mathbf{S} \vec{e} + \sum_{i=1}^N \|\mathbf{x}_i - \bar{\mathbf{x}}\|^2$$

Independent of the direction e

Minimizing $J(\vec{e})$ w.r.t. \vec{e} is equivalent to maximizing $\vec{e}^t \mathbf{S} \vec{e}$ w.r.t. \vec{e} .

We use the method of Lagrange multipliers to do so, while simultaneously
imposing the constraint that $\vec{e}^t \vec{e} = 1$. See appendix for details

So we have to take the derivative of the following modified function
w.r.t. \vec{e} (and set it to 0)

$$\tilde{J}(\vec{e}) = \vec{e}^t \mathbf{S} \vec{e} - \lambda (\vec{e}^t \vec{e} - 1)$$

Taking derivative of $\tilde{J}(\vec{e})$ w.r.t. \vec{e} and setting it to 0, we get

$$\mathbf{S} \vec{e} = \lambda \vec{e},$$

so \vec{e} is an eigen - vector of \mathbf{S} .

As $\vec{e}^t \mathbf{S} \vec{e} = \lambda$ and we wish to maximize $\vec{e}^t \mathbf{S} \vec{e}$, we choose \vec{e} to be the
eigen - vector corresponding to the maximum eigenvalue of \mathbf{S} .

PCA

- PCA thus projects the data onto that direction that minimizes the total squared difference between the data-points and their respective projections along that direction.
- This **equivalently** yields the direction along which the spread (or variance) will be maximum.
- Why? Note that the eigenvalue of a covariance matrix tells you the variance of the data when projected along that particular eigenvector:

$$\vec{S}\vec{e} = \lambda\vec{e} \rightarrow \vec{e}^t \vec{S}\vec{e} = \lambda$$

$$\vec{e}^t \vec{S}\vec{e} = \sum_{i=1}^N (\vec{e}^t (\mathbf{x}_i - \bar{\mathbf{x}}))^2$$

This term is proportional to the variance of the data when projected along \mathbf{e} .

PCA

- But for most applications (including face recognition), just a single direction is absolutely insufficient!
- We will need to project the data (from the high-dimensional, i.e. d -dimensional space) onto k ($k \ll d$) different mutually perpendicular directions.
- What is the criterion for deriving these directions?
- We seek those k directions for which the total reconstruction error of all the N images when projected on those directions is minimized.

PCA

- We seek those k directions for which the total reconstruction error of all the N images when projected on those directions is minimized.

$$J(\{\mathbf{e}_j\}_{j=1}^k) = \sum_{i=1}^N \left\| (\mathbf{x}_i - \bar{\mathbf{x}}) - \sum_{j=1}^k (\mathbf{e}_j^t (\mathbf{x}_i - \bar{\mathbf{x}})) \mathbf{e}_j \right\|_2^2$$

- One can prove that these k directions will be the eigenvectors of the \mathbf{S} matrix (equivalently covariance matrix of the data) corresponding to the k -largest eigenvalues. These k directions form the eigen-space.
- If the eigenvalues of \mathbf{S} are distinct, these k directions are defined uniquely (up to a sign factor)

PCA

- One can prove that these k directions will be the eigenvectors of the \mathbf{S} matrix (equivalently covariance matrix of the data) corresponding to the k -largest eigenvalues. These k directions form the eigen-space.



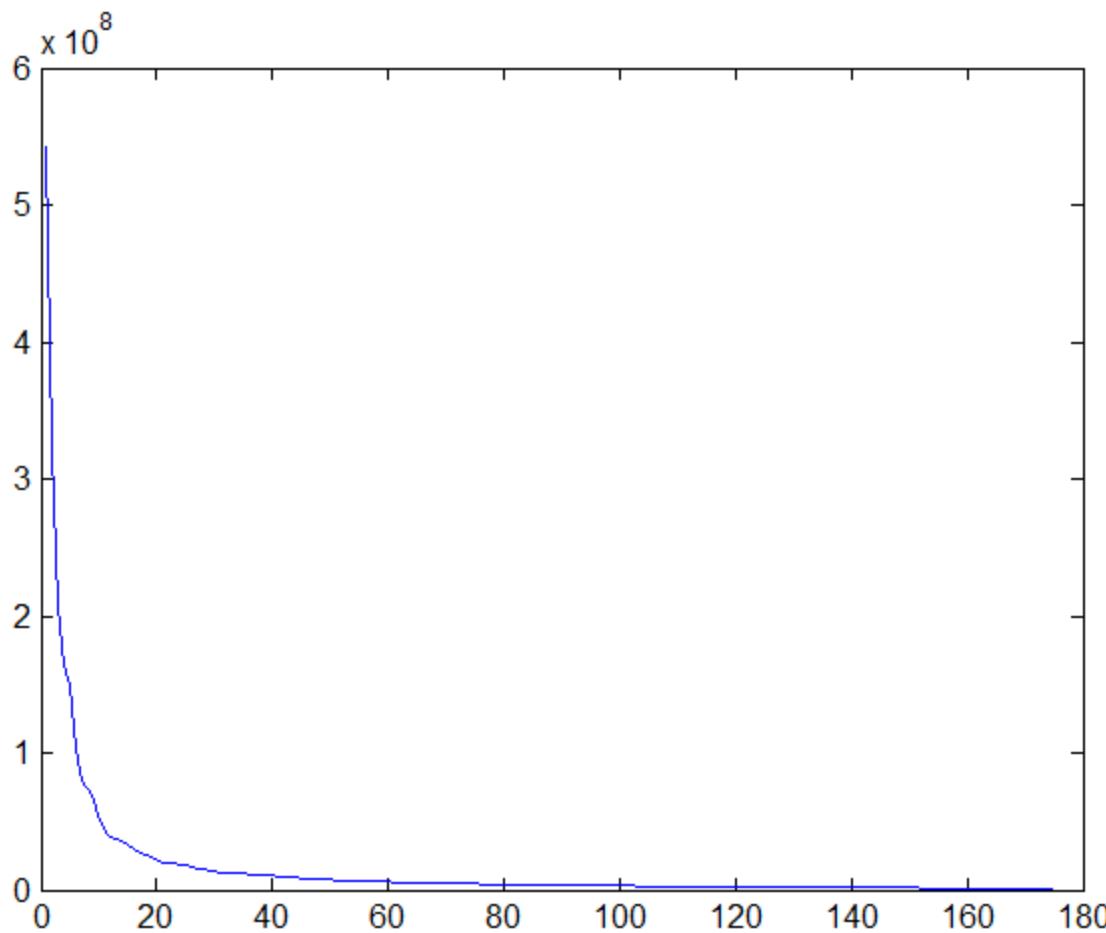
- **Sketch of the proof:**
 - ✓ Assume we have found \mathbf{e}_1 and are looking for \mathbf{e}_2 (where \mathbf{e}_2 is perpendicular to \mathbf{e}_1 and \mathbf{e}_2 has unit magnitude).
 - ✓ Write out the objective function with the two constraints.
 - ✓ Take the derivative of the objective function, set it to 0, and do some algebra to see that \mathbf{e}_2 is the eigenvector of \mathbf{S} with the second largest eigenvalue.
 - ✓ Proceed similarly for other directions.

Some observations about PCA for face images

- The matrix \mathbf{V} (with all columns) is orthonormal. Hence the squared error between any image and its approximation using just top k eigenvectors is given by:

$$\begin{aligned}\|\mathbf{x}_i - \tilde{\mathbf{x}}_i\|^2 &= \|V(\mathbf{a}_i - \tilde{\mathbf{a}}_i)\|^2 (\text{why?}) \\ &= \|\mathbf{a}_i - \tilde{\mathbf{a}}_i\|^2 (\text{why?}) \\ &= \sum_{j=k+1}^d \mathbf{a}_i^2 (\text{why?})\end{aligned}$$

This error is small ***on an average for a well-aligned group of face images*** – we will see why on the next slide.



The eigenvalues of the covariance matrix typically decay fast in value (if the faces were properly normalized). Note that the j -th eigenvalue is proportional to the variance of the j -th eigencoefficient, i.e.

$$\lambda_j = \mathbf{e}_j^t \mathbf{S} \mathbf{e}_j = \sum_{i=1}^N \mathbf{e}_j^t (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^t \mathbf{e}_j = (N-1)E(\alpha_{ij}^2)$$

What this means is that the data have low variance when projected along most of the eigenvectors, i.e. effectively the data are concentrated in a lower-dimensional subspace of the d -dimensional space.

PCA: Compression of a set of images

- Consider a database of N images that are “similar” (eg: all are face images, all are car images, etc.)
- Build an eigen-space from some subset of these images (could be all images, as well)
- We know that these images can often be reconstructed very well (i.e. with low error) using just a few eigenvectors.

PCA: Compression of a set of images

- Use this fact for image compression.
- Original data storage = d pixels $\times N$ images = Nd bytes (assume one byte per pixel intensity) = $8Nd$ bits.
- After PCA: Nk \times number of bits to store each eigen-coefficient = $32Nk$ bits (remember $k \ll d$, example: $d \sim 250,000$ and $k \sim 100$).
- Plus storage of eigenvectors = $32dk$ bits (remember $k \ll M$ as well).
- Plus mean image = $8d$ bits.
- Total: $32(N+d)k + 8d$ bits

PCA: Compression of a set of images

- Example: $N= 5000, d = 250000, k = 100$
- Original size/(size after PCA compression) $\sim 12.2.$
- Note: we allocated 32 bits for every element of the eigen-vector. This is actually very conservative and you can have further savings using several tricks.

PCA: Compression of a set of images

- This differs a lot from JPEG compression.
- JPEG uses discrete cosine transform (DCT) as the basis. PCA tunes the basis to the underlying training data! If you change the training data, the basis changes!
- The performance of the PCA compression algorithm will depend on how compactly the eigen-space can represent the data.
- We will study image compression using JPEG and other standards later on in the course.

Which is the best orthonormal basis?

Relationship between PCA and DCT

- Consider a set of M data-points (e.g. image patches in a vectorized form) represented as a linear combination of column vectors of an ortho-normal basis matrix:

$$\mathbf{q}_i = \mathbf{U}\boldsymbol{\theta}_i, \mathbf{q}_i \in R^{N \times 1}, \boldsymbol{\theta}_i \in R^{N \times 1}, \mathbf{U}\mathbf{U}^T = \mathbf{U}^T\mathbf{U} = \mathbf{I}$$

- Suppose we reconstruct each patch using only a subset of some k coefficients as follows:

$\tilde{\mathbf{q}}_i^{(k)} = \mathbf{U}\tilde{\boldsymbol{\theta}}_i$, where $\tilde{\boldsymbol{\theta}}_i$ is obtained by setting to 0 all except k coefficients (the same k coefficients are retained for all patches)

Which is the best orthonormal basis?

Relationship between PCA and DCT

- For which orthonormal basis \mathbf{U} is the following error the lowest:

$$E(\mathbf{U}) = \sum_{i=1}^M \left\| \tilde{\mathbf{q}}_i^{(k)} - \mathbf{q}_i \right\|^2$$

Which is the best orthonormal basis?

Relationship between PCA and DCT

- The answer is the PCA basis, i.e. the set of k eigenvectors of the correlation matrix \mathbf{C} , corresponding to the k largest eigen-values.
Here is \mathbf{C} is defined as:

$$\mathbf{C} = \frac{1}{M-1} \sum_{i=1}^M \mathbf{q}_i \mathbf{q}_i^T,$$

$$C_{kl} = \frac{1}{M-1} \sum_{i=1}^M q_{ikl} q_{ilk}$$

PCA: separable 2D version

- Find the correlation matrix \mathbf{C}_R of row vectors from the patches.
- Find the correlation matrix \mathbf{C}_C of column vectors from the patches.
- The final PCA basis is the Kronecker product of the individual bases:

$$\mathbf{C}_R = \frac{1}{M-1} \sum_{i=1}^M \sum_{j=1}^n q_i(j,:)^T q_i(j,:);$$

$[\mathbf{V}_R, \mathbf{D}_R] = eig(\mathbf{C}_R); q_i(j,:) \in R^{1 \times n} - j^{th}$ row vector of \mathbf{q}_i

$$\mathbf{C}_C = \frac{1}{M-1} \sum_{i=1}^M \sum_{j=1}^n q_i(:,j) q_i(:,j)^T;$$

$[\mathbf{V}_C, \mathbf{D}_C] = eig(\mathbf{C}_C); q_i(:,j) \in R^{n \times 1} - j^{th}$ column vector of \mathbf{q}_i

$$\mathbf{V} = \mathbf{V}_R \otimes \mathbf{V}_C; \mathbf{V} \mathbf{V}^T = I; \mathbf{V} \in R^{n^2 \times n^2}, \mathbf{V}_R \in R^{n \times n}, \mathbf{V}_C \in R^{n \times n}, \mathbf{q}_i \in R^{n \times n}$$

But PCA is not used in JPEG, because...

- It is image-dependent, and the basis matrix would need to be computed afresh for each image.
- The basis matrix would need to be **stored** for each image.
- It is **expensive** to compute – $O(n^3)$ for a vector with n elements.

The DCT is used instead!

DCT and PCA

- DCT can be computed very fast using fft.
- It is universal – no need to store the DCT bases explicitly.
- DCT has very good energy compaction properties, only slightly worse than PCA.

Experiment

- Suppose you extract $M \sim 100,000$ small-sized (8×8) patches from a set of images.
- Compute the column-column and row-row correlation matrices.

$$\mathbf{C}_C = \frac{1}{M-1} \sum_{i=1}^M \mathbf{P}_i \mathbf{P}_i^T = \frac{1}{M-1} \sum_{i=1}^M \sum_{j=1}^8 P_i(:, j) P_i(:, j)';$$

See code [here](#).

$$\mathbf{C}_R = \frac{1}{M-1} \sum_{i=1}^M \mathbf{P}_i^T \mathbf{P}_i = \frac{1}{M-1} \sum_{i=1}^M \sum_{j=1}^8 P_i(j, :)' P_i(j, :);$$

- Compute their eigenvectors \mathbf{V}_R and \mathbf{V}_C .
- The eigenvectors will be very similar to the columns of the 1D-DCT matrix! (as evidenced by dot product values).
- Now compute the Kronecker product of \mathbf{V}_R and \mathbf{V}_C and call it \mathbf{V} . Reshape each column of \mathbf{V} to form an image. These images will appear very similar to the DCT bases.

0.3536	0.4904	0.4619	0.4157	0.3536	0.2778	0.1913	0.0975
0.3536	0.4157	0.1913	-0.0975	-0.3536	-0.4904	-0.4619	-0.2778
0.3536	0.2778	-0.1913	-0.4904	-0.3536	0.0975	0.4619	0.4157
0.3536	0.0975	-0.4619	-0.2778	0.3536	0.4157	-0.1913	-0.4904
0.3536	-0.0975	-0.4619	0.2778	0.3536	-0.4157	-0.1913	0.4904
0.3536	-0.2778	-0.1913	0.4904	-0.3536	-0.0975	0.4619	-0.4157
0.3536	-0.4157	0.1913	0.0975	-0.3536	0.4904	-0.4619	0.2778
0.3536	-0.4904	0.4619	-0.4157	0.3536	-0.2778	0.1913	-0.0975
0.3517	-0.4493	-0.4278	0.4230	0.3754	0.3247	-0.2250	-0.1245
0.3534	-0.4366	-0.2276	-0.0110	-0.3078	-0.4746	0.4732	0.2975
0.3543	-0.3101	0.1728	-0.4830	-0.3989	0.0498	-0.4299	-0.4109
0.3546	-0.1115	0.4799	-0.3005	0.3342	0.4102	0.1856	0.4761
0.3547	0.1141	0.4823	0.2944	0.3301	-0.4182	0.1745	-0.4771
0.3543	0.3104	0.1771	0.4834	-0.3977	-0.0322	-0.4308	0.4103
0.3535	0.4357	-0.2319	0.0143	-0.3009	0.4656	0.4851	-0.2975
0.3520	0.4468	-0.4328	-0.4204	0.3686	-0.3253	-0.2342	0.1261
0.3520	-0.4461	-0.4305	0.4224	0.3696	0.3247	0.2342	0.1283
0.3537	-0.4338	-0.2345	-0.0114	-0.3000	-0.4671	-0.4814	-0.3028
0.3545	-0.3086	0.1662	-0.4896	-0.4007	0.0359	0.4261	0.4102
0.3548	-0.1145	0.4763	-0.3031	0.3339	0.4198	-0.1800	-0.4713
0.3548	0.1056	0.4839	0.2926	0.3349	-0.4194	-0.1766	0.4733
0.3543	0.3043	0.1863	0.4833	-0.4028	-0.0354	0.4269	-0.4097
0.3532	0.4389	-0.2269	0.0180	-0.3008	0.4654	-0.4811	0.3037
0.3512	0.4562	-0.4300	-0.4126	0.3694	-0.3242	0.2335	-0.1319

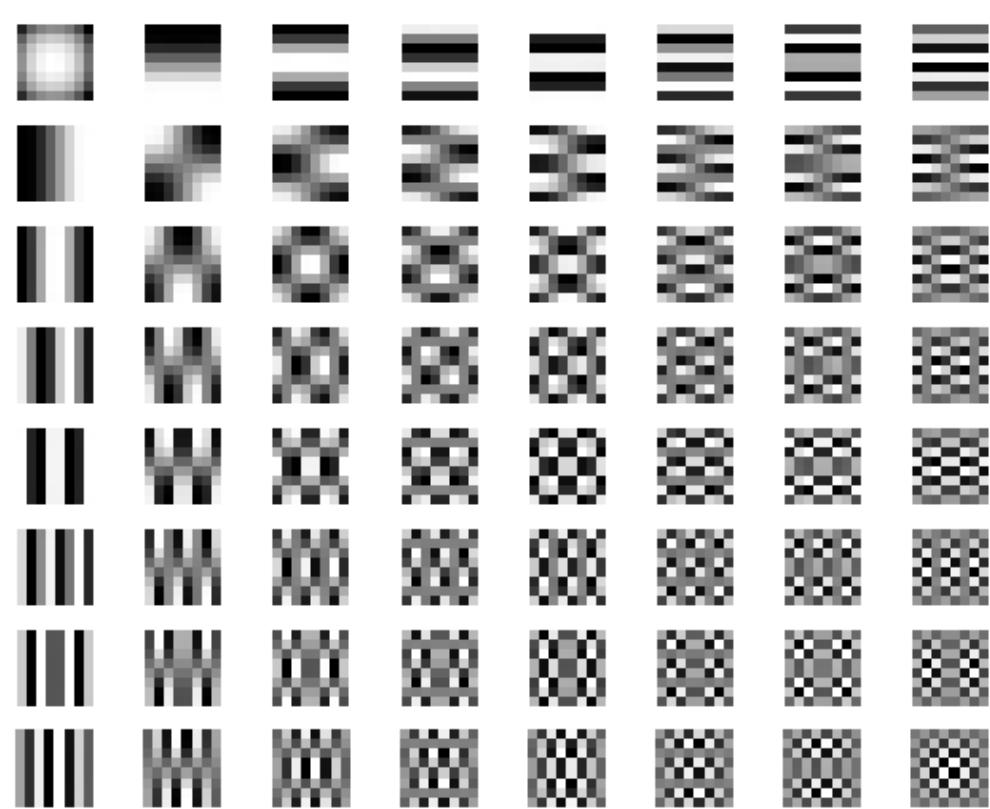
DCT matrix: `dctmtx`
command from MATLAB (see
code on website)

V_C: Eigenvectors of column-
column correlation matrix

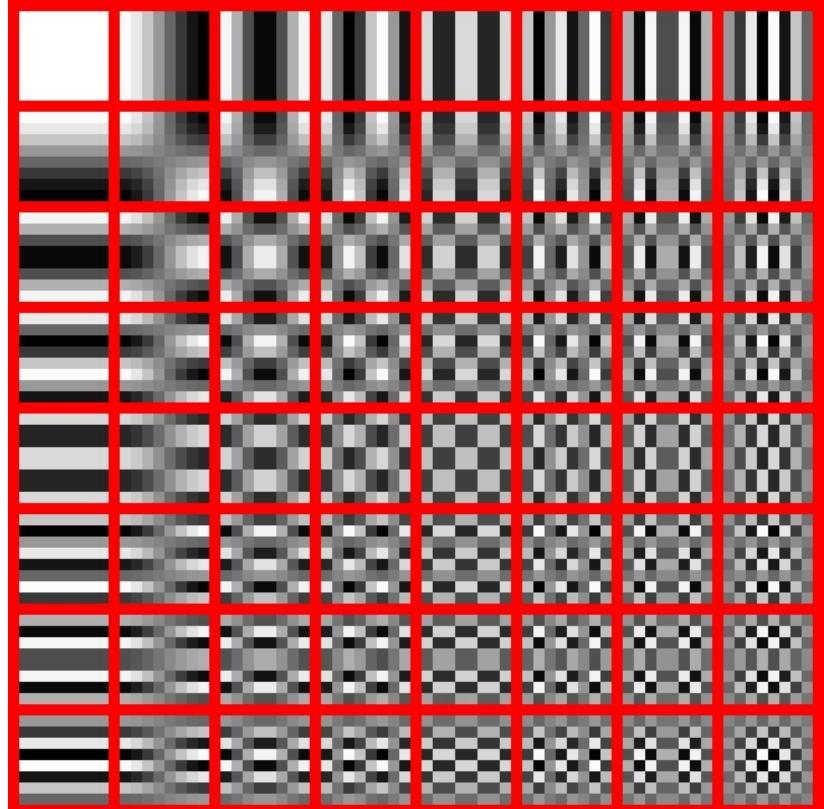
V_R: Eigenvectors of row-row
correlation matrix

Absolute value of dot products between the columns of DCT matrix and columns of V_R (left) and V_C (right)

1.0000	0.0007	0.0032	0.0002	0.0013	0.0001	0.0005	0.0000	1.0000	0.0002	0.0029	0.0001	0.0010	0.0000	0.0004	0.0000
0.0007	0.9970	0.0097	0.0689	0.0009	0.0322	0.0003	0.0110	0.0002	0.9965	0.0028	0.0766	0.0005	0.0314	0.0009	0.0107
0.0033	0.0106	0.9968	0.0118	0.0713	0.0004	0.0334	0.0025	0.0029	0.0025	0.9969	0.0046	0.0728	0.0017	0.0304	0.0013
0.0002	0.0718	0.0124	0.9926	0.0007	0.0927	0.0017	0.0276	0.0001	0.0795	0.0044	0.9923	0.0029	0.0916	0.0015	0.0243
0.0010	0.0001	0.0737	0.0004	0.9942	0.0008	0.0780	0.0010	0.0008	0.0003	0.0747	0.0026	0.9948	0.0061	0.0696	0.0004
0.0000	0.0261	0.0015	0.0962	0.0005	0.9934	0.0011	0.0569	0.0000	0.0246	0.0021	0.0949	0.0069	0.9940	0.0131	0.0452
0.0003	0.0007	0.0276	0.0021	0.0802	0.0010	0.9964	0.0013	0.0003	0.0004	0.0252	0.0003	0.0715	0.0137	0.9970	0.0002
0.0000	0.0076	0.0026	0.0227	0.0012	0.0596	0.0015	0.9979	0.0000	0.0076	0.0013	0.0207	0.0001	0.0476	0.0009	0.9986



64 columns of \mathbf{V} – each reshaped to form an 8 x 8 image, and rescaled to fit in the 0-1 range.
Notice the similarity between the DCT bases and the columns of \mathbf{V} . Again, \mathbf{V} is the Kronecker product of \mathbf{V}_R and \mathbf{V}_C .



DCT bases

DCT and PCA

- DCT is very close to PCA when the patches come from what is called as a **stationary first order Markov process**, i.e.

$$q_i = \rho q_{i-1} + \eta_i, \eta_i \sim N(0, \sigma_\eta^2), \rho < 1$$

$$E(q_i q_{i-1}) = \rho \sigma_q^2, E(q_i q_{i-2}) = \rho^2 \sigma_q^2, \dots, E(q_i q_{i-n+1}) = \rho^{n-1} \sigma_q^2,$$

$$\mathbf{C} = \sigma_q^2 \begin{pmatrix} 1 & \rho & \rho^2 & \cdot & \rho^{n-1} \\ \rho & 1 & \cdot & \cdot & \cdot \\ \rho^2 & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \rho^{n-1} & \cdot & \cdot & \cdot & 1 \end{pmatrix}, C_{ij} = E(q_i q_j)$$

DCT and PCA

- One can show that the eigenvectors of the covariance matrix of the form seen on the previous slide are the DCT basis vectors!
- Natural images approximate this first order Markov model, and hence DCT is almost as good as PCA for compression of a large ensemble of image patches.
- DCT has the advantage of being a universal basis and also the DCT coefficients are more efficiently computable than PCA coefficients (because DCT computation uses FFT).

More results from the previous experiment. See code [here](#).

1.0000	0.9902	0.9795	0.9733	0.9682	0.9639	0.9604	0.9570
0.9902	1.0005	0.9908	0.9795	0.9734	0.9684	0.9643	0.9605
0.9795	0.9908	1.0010	0.9908	0.9796	0.9735	0.9689	0.9646
0.9733	0.9795	0.9908	1.0005	0.9904	0.9793	0.9735	0.9686
0.9682	0.9734	0.9796	0.9904	1.0004	0.9903	0.9794	0.9734
0.9639	0.9684	0.9735	0.9793	0.9903	1.0001	0.9903	0.9793
0.9604	0.9643	0.9689	0.9735	0.9794	0.9903	1.0004	0.9904
0.9570	0.9605	0.9646	0.9686	0.9734	0.9793	0.9904	1.0002

CR/CR(1,1) -
 Notice it can be
 approximated by the
 form shown two slides
 before, with $\rho \sim 0.99$

1.0000	0.9888	0.9770	0.9704	0.9648	0.9599	0.9554	0.9510
0.9888	1.0004	0.9891	0.9768	0.9703	0.9646	0.9596	0.9548
0.9770	0.9891	1.0004	0.9886	0.9764	0.9698	0.9640	0.9587
0.9704	0.9768	0.9886	0.9994	0.9878	0.9755	0.9687	0.9627
0.9648	0.9703	0.9764	0.9878	0.9986	0.9870	0.9746	0.9676
0.9599	0.9646	0.9698	0.9755	0.9870	0.9978	0.9861	0.9734
0.9554	0.9596	0.9640	0.9687	0.9746	0.9861	0.9967	0.9847
0.9510	0.9548	0.9587	0.9627	0.9676	0.9734	0.9847	0.9951

CC/CC(1,1) -
 Notice it can be
 approximated by the
 form shown two slides
 before, with $\rho \sim 0.9888$

Non-negative matrix factorization (NMF) and non-negative sparse coding (NNSC)

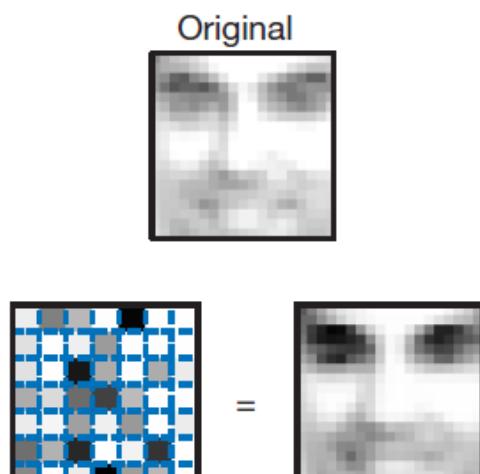
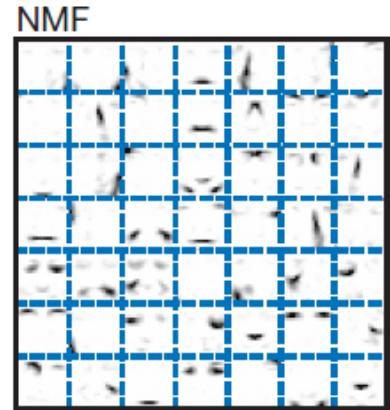
NMF

- Consider a matrix $\mathbf{Y} \in \mathbb{R}^{n \times N}$, each column of which is a vectorised patch of n elements, i.e. $\mathbf{Y} = [\mathbf{y}_1 | \mathbf{y}_2 | \dots | \mathbf{y}_N]$.
- Then NMF seeks to factorize \mathbf{Y} into the product of matrices \mathbf{W} (size $n \times r$) and \mathbf{H} (size $r \times N$) which are both non-negative, such that $\mathbf{Y} \approx \mathbf{WH}$.
- Note: non-negativity means that every element of \mathbf{W} , \mathbf{H} must be non-negative.

NMF

- The columns of \mathbf{W} are basis vectors and the values in the i -th column of \mathbf{H} are coefficients.
- We seek to approximate each vector \mathbf{y}_i with a linear combination of the columns of \mathbf{W} .
- Usually, r is chosen to be smaller than n or N , which gives good model compression.
- The speciality of NMF is that it allows only additive combinations of the columns of \mathbf{W} to approximate \mathbf{Y} .
- This is in tune with the intuition of combining different parts to form a whole.

PCA versus NMF bases



Source:

Lee and Seung, "Learning the parts of objects by non-negative matrix factorization",

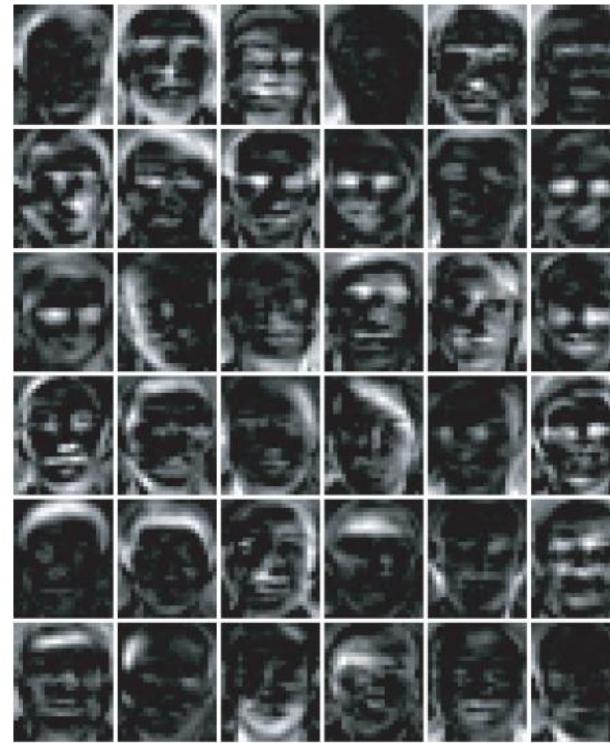
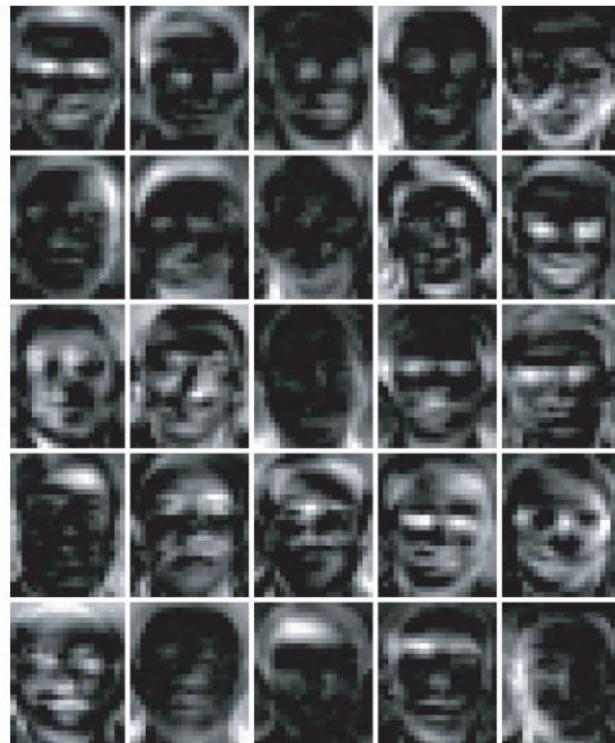
<http://lsa.colorado.edu/LexicalSemantics/seung-nonneg-matrix.pdf>

NMF bases are non-negative unlike eigenfaces. In that sense, each basis image from NMF is more like a model face than from the eigenfaces.

In practice they do not appear as sparse as seen in this figure, but more like what you will see on the next slide.



Source:
Wang et al, NMF framework
for face recognition,
<http://www.cs.ucsb.edu/~mturk/pubs/IJPRAI05.pdf>



NMF: Objective function

- The objective function to be minimized is:

$$E(\mathbf{W}, \mathbf{H}) = \|\mathbf{Y} - \mathbf{WH}\|_F^2 \text{ such that } \mathbf{W} \geq \mathbf{0}, \mathbf{H} \geq \mathbf{0}$$

- Method of alternating optimization using projected gradient descent with adaptive step size:

Fix $\mathbf{H}, \mathbf{W} \leftarrow \mathbf{W} + \delta 2(\mathbf{Y} - \mathbf{WH})(\mathbf{H}^T)$ such that $\mathbf{W} \geq \mathbf{0}$

Fix $\mathbf{W}, \mathbf{H} \leftarrow \mathbf{H} + \delta 2\mathbf{W}^T(\mathbf{Y} - \mathbf{WH})$ such that $\mathbf{H} \geq \mathbf{0}$

NMF: Multiplicative updates

- The objective function to be minimized is:

$$E(\mathbf{W}, \mathbf{H}) = \|\mathbf{Y} - \mathbf{WH}\|_F^2 \text{ such that } \mathbf{W} \geq \mathbf{0}, \mathbf{H} \geq \mathbf{0}$$

- However there are also multiplicative updates of the following form:

$$H_{ab} = H_{ab} \frac{(W^T Y)_{ab}}{(W^T W H)_{ab}}, 1 \leq a \leq r, 1 \leq b \leq N$$

$$W_{ia} = W_{ia} \frac{(Y H^T)_{ia}}{(W H H^T)_{ia}}, 1 \leq a \leq r, 1 \leq b \leq N$$

NMF: Multiplicative updates

- The multiplicative updates maintain the non-negativity of \mathbf{W} and \mathbf{H} and converge to a stationary point of $E(\mathbf{W}, \mathbf{H})$.
- The multiplicative updates also guarantee decrease of the cost function. They are based on majorization-minimization.

Non-negative sparse coding

- An extension to NMF with sparseness penalty imposed on the coefficients.
- The objective function to be minimized is:

$$E(\mathbf{W}, \mathbf{H}) = \|\mathbf{Y} - \mathbf{WH}\|_F^2 + \lambda \|\mathbf{H}\|_1$$

such that $\mathbf{W} \geq 0, \mathbf{H} \geq 0, \mathbf{w}_i^t \mathbf{w}_i = 1$ for $1 \leq i \leq r$

- Need to be careful to avoid multiple solutions – also impose unit norm constraint on all columns of \mathbf{W} .

Non-negative sparse coding

- Algorithm:

Fix H ,

$$W \leftarrow W + \delta 2(Y - WH)(H^T)$$

Set negative values of W to 0 and

unit normalize its columns so that $w_i^t w_i = 1$ for $1 \leq i \leq r$

Fix W ,

$$H \leftarrow H + \delta 2W^T(Y - WH) + \lambda$$

Set negative values of H to 0

Results

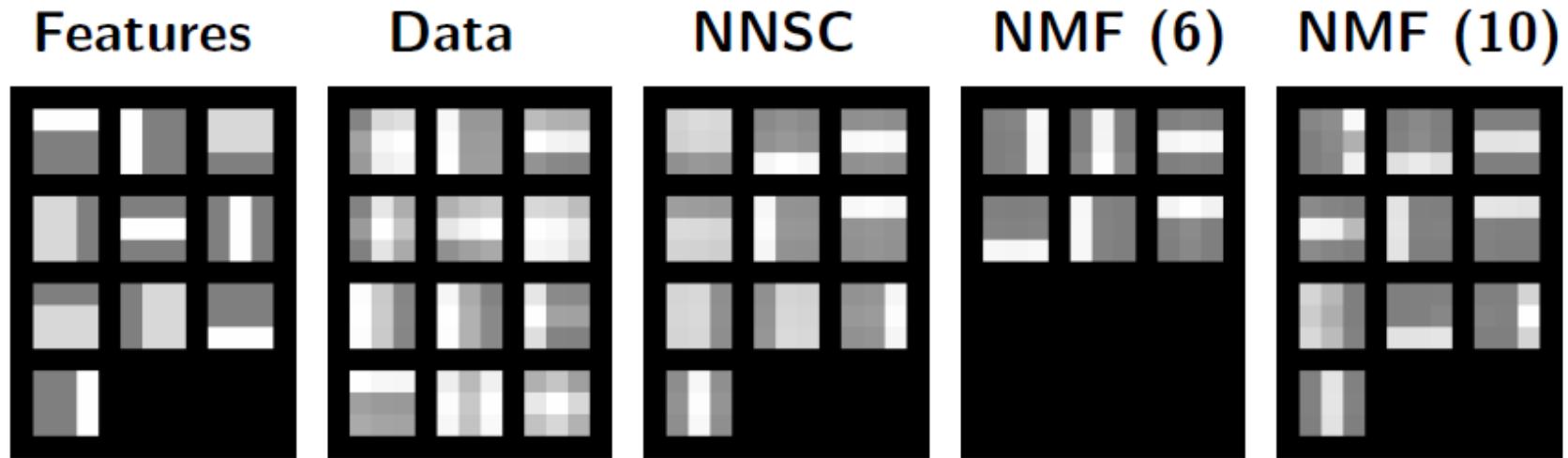


Figure 1: Experiments on bars data. Features: The 10 original features that were used to construct the dataset. Data: A random sample of 12 data vectors. These constitute superpositions of the original features. NNSC: Features learned by NNSC, with dimensionality of the hidden representation equal to 10, starting from random initial values. NMF (6): Features learned by NMF, with dimensionality 6. NMF (10): Features learned by NMF, with dimensionality 10. See main text for discussion.

NNSC and Poisson Data

- The noise affecting real-world images is known to follow the Poisson distribution.
- The Poisson distribution is given as follows:

$$P(X = i) = \frac{\lambda^i}{i!} e^{-\lambda}$$

X is a Poisson random variable, i is a non-negative integer, λ is the mean of the Poisson distribution

- For such a distribution, the mean is equal to the variance.

NNSC and Poisson Data

- Consider an image I , a Poisson-corrupted version of this image is obtained as $J_i \sim \text{Poisson}(I_i)$ where I_i is the average intensity (photon flux) at pixel i .
- Note that the noisy intensity at any pixel location i is always a non-negative integer, drawn from a Poisson distribution having mean (hence variance) I_i .

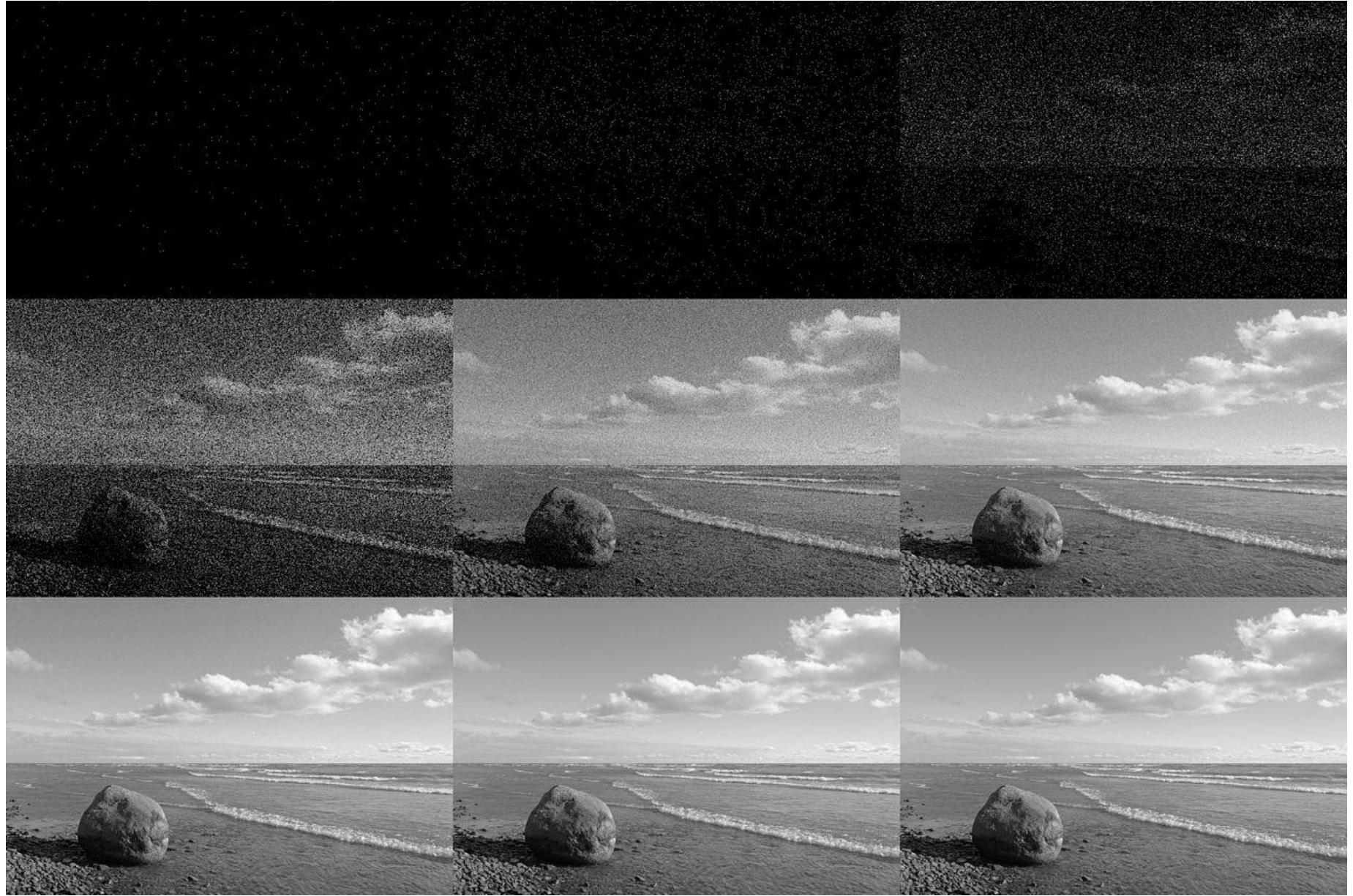


Image quality at 9 different average rates of photon emission.

https://en.wikipedia.org/wiki/Shot_noise

NNSC and Poisson denoising

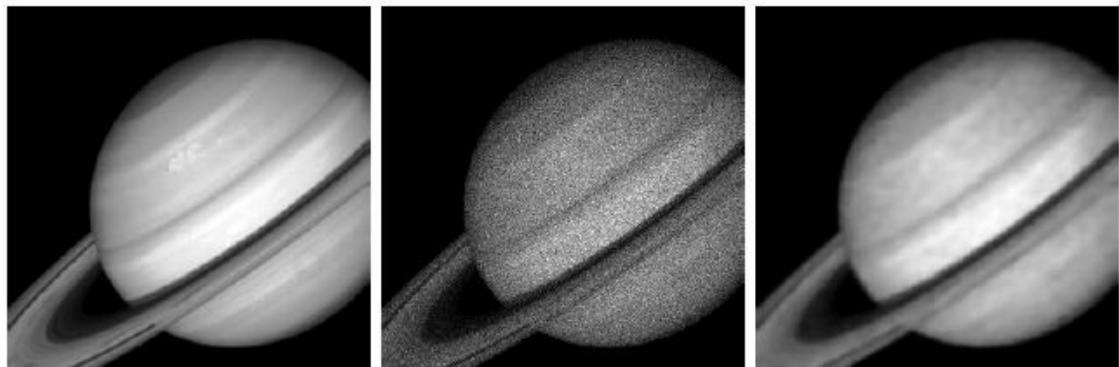
- NNSC lends itself naturally to Poisson denoising due to the non-negativity constraints – on measurements as well as data (mean of Poisson = photon rate).
- Consider $\mathbf{Y} \sim \text{Poisson}(\mathbf{WH})$, i.e. $y_{kl} \sim \text{Poisson}([\mathbf{WH}]_{kl})$, $1 \leq k \leq n$, $1 \leq l \leq N$. We have $[\mathbf{WH}]_{kl} = [\mathbf{W}\mathbf{h}_l]_k$.
- The objective function to be minimized is as follows:

$$E(W, H) = \sum_{k=1}^n \sum_{l=1}^N -y_{kl} \log(W\mathbf{h}_l)_k + (W\mathbf{h}_l)_k + \rho \sum_{l=1}^N \sum_{c=1}^r h_{cl}$$

such that $W \geq 0, H \geq 0, w_i^t w_i = 1$ for $1 \leq i \leq r$

Non-negative
Poisson likelihood
where y_{kl} is a noisy
pixel

Sparsity on
coefficients



(a) Original image

(b) Noisy image, PSNR
 $= 21.39$

(c) Our algorithm,
PSNR = 32.88

Figure 3.1: Saturn image, Peak = 50



(a) Original image

(b) Noisy image, PSNR
 $= 19.71$

(c) Our algorithm,
PSNR = 28.63

Figure 3.2: House image, Peak = 50



(a) Original image

(b) Noisy image, $\text{PSNR} = 10.37$ (c) Our algorithm, $\text{PSNR} = 21.40$

Figure 3.3: Pepper image, Peak = 5



(a) Original image

(b) Noisy image, $\text{PSNR} = 11.68$ (c) Our algorithm, $\text{PSNR} = 23.15$

A comment on image processing under Poisson noise

- The relative error due to Poisson noise is equal to standard deviation divided by mean = $(\lambda)^{-0.5}$.
- At low photon counts , i.e. low λ , the noise variance is less, but its relative effect is more, and hence the problem is more challenging!
- This is unlike the signal independent additive iid Gaussian noise!

Learning Overcomplete Dictionaries

Introduction: Complete and over-complete dictionaries

- Signals are often represented as a linear combination of basis functions (e.g. Fourier or wavelet representation).
- The basis functions always have the same dimensionality as the (discrete) signals they represent.
- The number of basis vectors is traditionally the same as the dimensionality of the signals they represent.
- These bases may be orthonormal (Fourier, wavelet, PCA) or may not be orthonormal (in fact any full rank matrix can be used to represent signals – albeit not compactly).

Introduction: Complete and over-complete bases

- A more general representation for signals uses so called “over-complete dictionaries”, where the number of basis functions is MORE than the dimensionality of the signals.
- Complete and over-complete bases:

$$\mathbf{x} = \mathbf{As}$$

$$\mathbf{x} \in R^n, \mathbf{A} \in R^{n \times n}$$

$$\mathbf{s} \in R^m$$

$$\mathbf{x} = \mathbf{As}$$

$$\mathbf{x} \in R^n, \mathbf{A} \in R^{n \times m}$$

$$\mathbf{s} \in R^m (m > n)$$



Introduction: Construction of over-complete bases

- Over-complete bases can be created by union of multiple sets of complete bases.
- Example 1: A signal with n values can be represented using a union of $n \times n$ Fourier and $n \times n$ Haar wavelet bases, yielding a $n \times 2n$ basis matrix.
- Example 2: A signal with n values can be represented by adding sinusoids of more frequencies to an existing Fourier basis matrix with n vectors.

Introduction: uniqueness? ☹

- With complete bases, the representation of the signal is always unique.
- Example: Signals are uniquely defined by their wavelet or Fourier transforms, the eigen-coefficients of any signal (given PCA basis) are uniquely defined.
- This uniqueness is LOST with over-complete basis.

Introduction: compactness! ☺

- Advantage: over-complete bases afford much greater compactness in signal representation.
- Example: Consider two types of audio-signals – whistles and claps. Signals of either type can be represented in a complete Fourier or wavelet basis (power-law of compressibility will apply).
- BUT: imagine two complete bases respectively learned for whistles and claps – B_1 and B_2 .

Introduction: compactness! ☺

- Suppose \mathcal{B}_1 and \mathcal{B}_2 are such that a whistle (resp. clap) signal will likely have a compact representation in the whistle (resp. clap) basis and not in the other one.
- A whistle+clap signal will NOT have a compact representation in either basis - \mathcal{B}_1 or \mathcal{B}_2 !
- But the whistle+clap signal WILL have a compact representation in the overcomplete basis $\mathcal{B} = [\mathcal{B}_1 \ \mathcal{B}_2]$.
- Another (and simpler) example: cosines and spikes!

More problems 😞

- Since a signal can have many representations in an over-complete basis, which one do we pick?
- **Pick the sparsest one**, i.e. the one with least number of non-zero elements which either perfectly reconstructs the signal, or reconstructs the signal up to some error.
- Finding the sparsest representation for a signal in an over-complete basis is an NP-hard problem (i.e. the best known algorithm for this task has exponential time complexity 😞)

More problems 😞

- In other words, the following problem is NP-hard:

$$x = As$$

$$x \in R^n, A \in R^{n \times m}$$

$$s \in R^m (m > n)$$

$$s^* = \min \|s\|_0 \text{ subject to } x = As$$

Solving (?) those problems

- The NP-hard problem has several methods for approximation – basis pursuit (BP), matching pursuit (MP), orthogonal matching pursuit (OMP) and many others.
- None of them will give you the sparsest solution always – but they will yield the sparsest solutions under certain conditions on the overcomplete dictionary \mathbf{A} .

Learning the dictionaries

- So far we assumed that the basis (i.e. \mathbf{A}) was fixed, and optimized for the sparse representation.
- Now, we need to learn \mathbf{A} as well – along with the sparse codes!
- We've learned about PCA – but it gives an orthonormal basis.

Learning the dictionary: analogy with K-means

- In K-means, we start with a bunch of K cluster centers and assign each point in the dataset to the nearest cluster center.
- The cluster centers are re-computed by taking the mean of all points assigned to a cluster.
- The assignment and cluster-center computation problems are iterated until a convergence criterion is met.

Learning the bases: analogy with K-means

- K-means is a special sparse coding problem where each point is represented by strictly one of K dictionary elements.
- Our dictionary (or bases) learning problem is more complex: we are trying to express each point as a linear combination of a **subset** of dictionary elements (or a **sparse** linear combination of dictionary elements).

Learning the bases: Method 1

(Olshausen and Field, 1996)

$x_k = k^{th}$ signal; $x_k = As_k + \text{noise}$

$$\begin{aligned}\{A, s_k\} &= \arg \max_{A, s_k} \prod_k p(A, s_k | x_k) = \arg \max_A \prod_k p(x_k | A, s_k) p(A) p(s_k | A) \\ &= \arg \min_{A, s_k} \sum_k \log p(s_k | A) + \log p(x_k | A, s_k)\end{aligned}$$

Keeping s_k fixed, optimize on A using the foll. estimator

$$\arg \min_A \sum_{k=1}^N (\lambda \|s_k\|_1 + \|x_k - As_k\|^2)$$

Gradient descent

$$A^{(t+1)} = A^{(t)} - \eta \sum_{k=1}^N (A^{(t)} s_k - x_k) s_k^T$$

$$p(s_k) \propto e^{-\lambda \|s_k\|_1}; p(x_k | A, s_k) \propto e^{-\|x_k - As_k\|^2 / (2\sigma^2)}$$

Uniform prior on A

Gaussian likelihood, i.e.
assuming Gaussian noise

Laplacian prior on s_k
(given any A)

Two-step iterative procedure

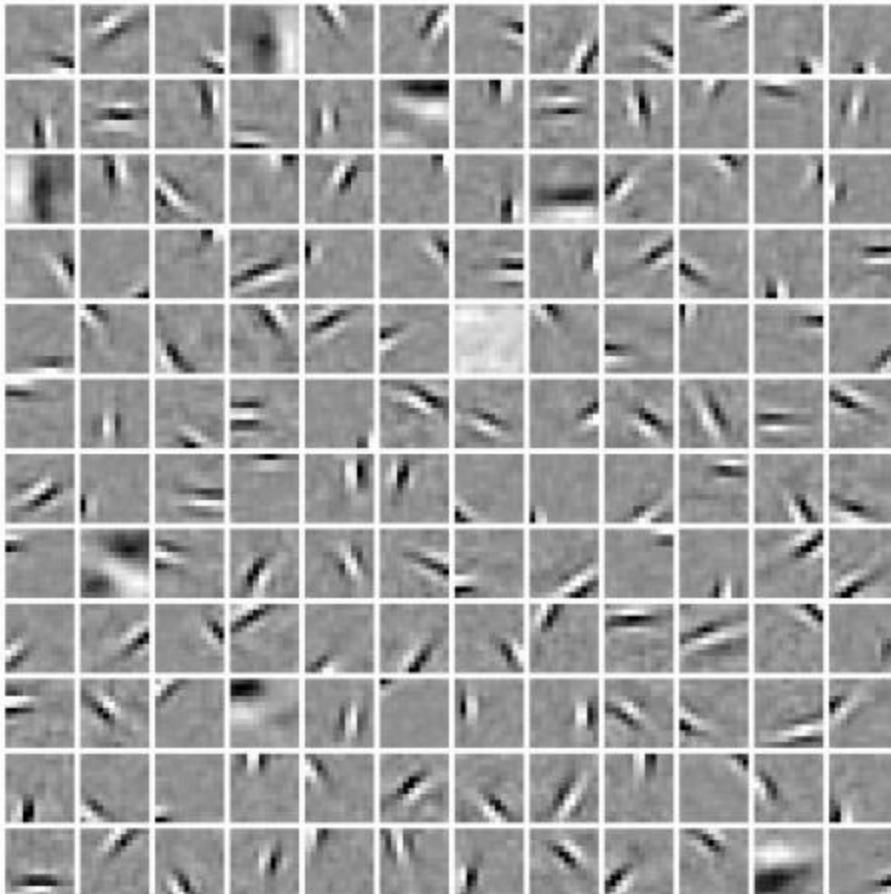
$$E(\mathbf{A}, \{\mathbf{s}_k\}) = \sum_{k=1}^N (\lambda \|\mathbf{s}_k\|_1 + \|\mathbf{x}_k - \mathbf{A}\mathbf{s}_k\|^2)$$

Gradient descent

$$\mathbf{A}^{(t+1)} = \mathbf{A}^{(t)} - \eta \sum_{k=1}^N (\mathbf{A}^{(t)} \mathbf{s}_k - \mathbf{x}_k) \mathbf{s}_k^t$$

- Fix the basis \mathbf{A} and obtain the sparse coefficients for each signal using MP, OMP or BP (some papers – like the one by Olshausen and Field - use gradient descent for this step!).
- Now fix the coefficients, and update the basis vectors (using various techniques, one of which was described on the previous slide).
- Normalize each basis vector to unit norm (so we may use *projected* gradient descent).
- Repeat the previous two steps until some error criterion is met.

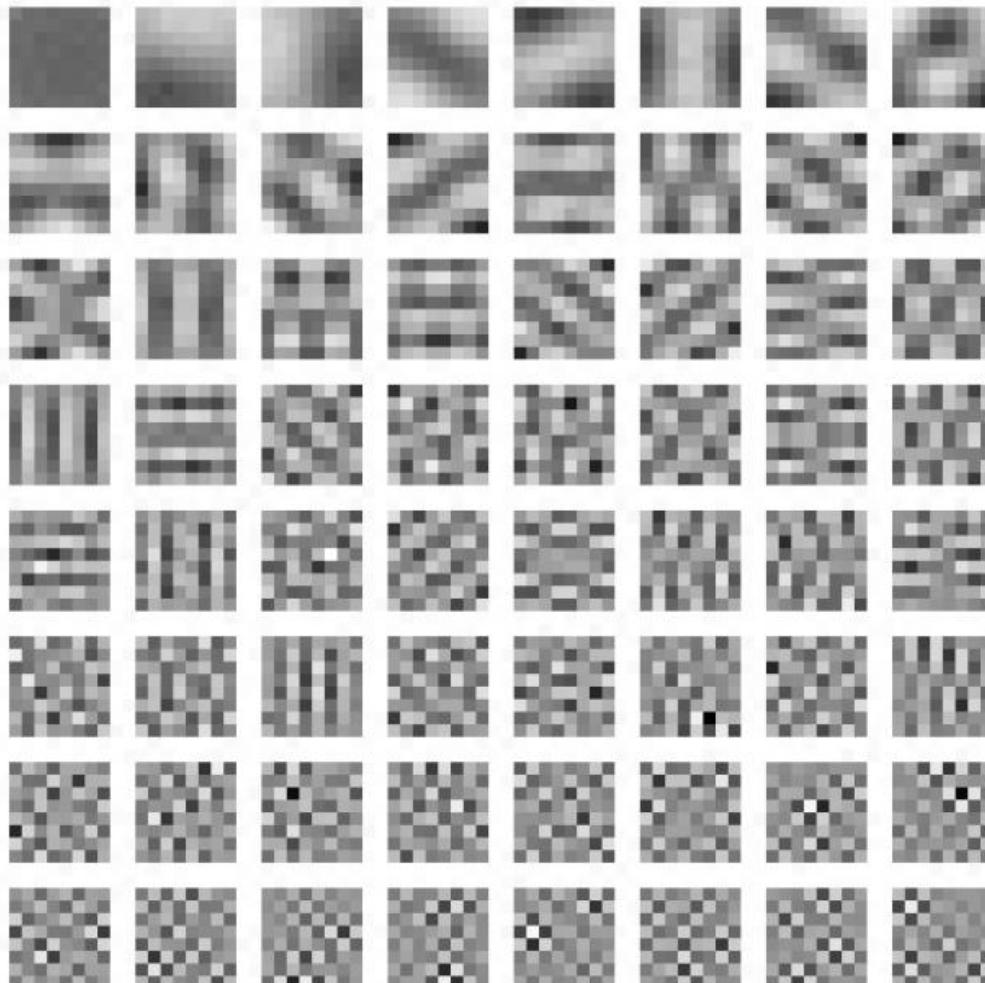
Toy Experiment 1



Result of basis learning (dictionary with 144 elements) with sparsity constraints on the codes. Training performed on 12 x 12 patches extracted from natural images.

It is conjectured that the simple cells of the V1 area of the visual cortex in the brain have responses similar in shape to these learned dictionary elements.

Ref: Olshausen and Field, “Natural image statistics and efficient coding”
<https://pdfs.semanticscholar.org/e309/e441a38ccee6456bd02e0f1e894e44180d53.pdf>



These are the results one obtains upon performing PCA on the 12 x 12 patches.

Ref: Olshausen and Field, “Natural image statistics and efficient coding”

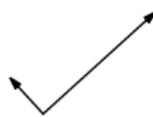
<https://pdfs.semanticscholar.org/e309/e441a38ccee6456bd02e0f1e894e44180d53.pdf>

Toy Experiment 2

- Data-points generated as a (Laplacian/super-Laplacian) random linear combination of some arbitrarily chosen basis vectors. Note: the coefficients are chosen from a (super)-Laplacian distribution.
- In the no-noise situation, the aim is to extract the basis vectors and the coefficients of the linear combination.
- The fitting results are shown on the next slide. The true and estimated directions agree quite well.

Ref: Lewicki and Sejnowski, “Learning overcomplete representations”

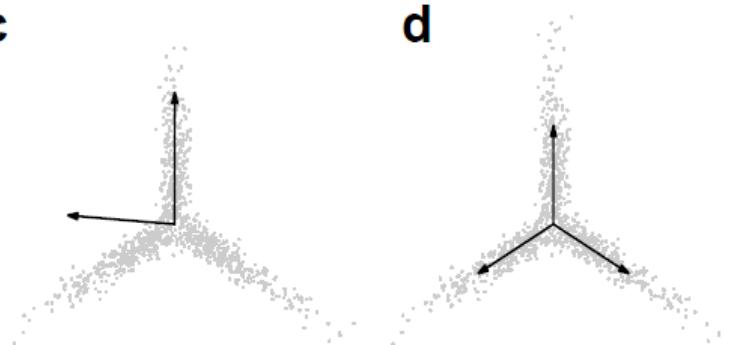
a



b



c



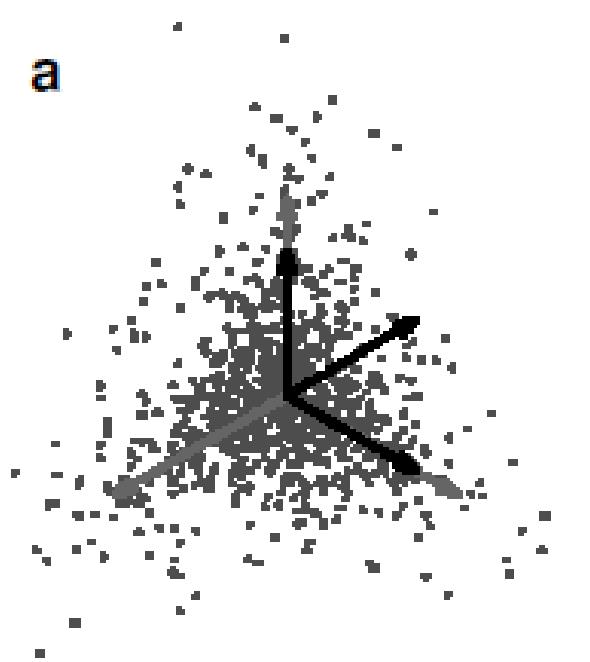
d



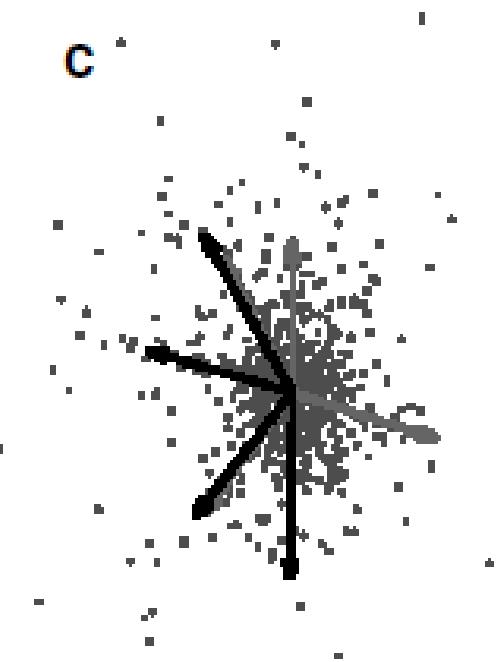
Figure 2: Fitting two-dimensional data with different bases. (a) PCA makes the assumption that the data have a Gaussian distribution. The optimal basis vectors are orthogonal and are not efficient at representing non-orthogonal density distributions. (b) ICA does not require that the vectors be orthogonal and can fit more general types of densities. (c) Some data distributions, however, cannot be modeled adequately by either PCA or ICA. (d) Allowing the basis to be overcomplete allows the three-armed distribution to be fit with simple and independent distributions for the basis vector coefficients.

Ref: Lewicki and Sejnowski, “[Learning overcomplete representations](#)”

a



c



Ref: Lewicki and Sejnowski, “[Learning overcomplete representations](#)”

Learning the Bases: Method of Optimal Directions (MOD) - Method 2

- Given a fixed dictionary \mathbf{A} , assume sparse codes for every signal are computed using OMP, MP etc.
- The overall error for dictionary update is now given as

$$E(\mathbf{A}) = \sum_{i=1}^n \|\mathbf{y}_i - \mathbf{As}_i\|^2 = \|\mathbf{Y} - \mathbf{AS}\|_F^2$$

- We want to find dictionary \mathbf{A} that minimizes this error (assuming fixed sparse codes).

Ref: Engan et al, "Method of optimal directions for frame design"

Learning the Bases: Method of Optimal Directions (MOD) - Method 2

- Take the derivative of $E(\mathbf{A})$ w.r.t. \mathbf{A} and set it to 0 . This gives us the following update:

$$(\mathbf{Y} - \mathbf{AS})\mathbf{S}^T = \mathbf{0} \Rightarrow \mathbf{A}^{(t+1)} = \mathbf{YS}^{(t)^T} (\mathbf{S}^{(t)}\mathbf{S}^{(t)^T})^{-1}$$

- Following the update of \mathbf{A} , each column in \mathbf{A} is independently rescaled to unit norm.
- The updates of \mathbf{A} and \mathbf{S} alternate with each other till some convergence criterion is reached.
- This method is more efficient than the one by Olshausen and Field.

Learning the Bases: Method 3- Union of Orthonormal Bases

- Like before, we represent a signal in the following way:

$$\mathbf{X} = \mathbf{AS} + \varepsilon$$

$$(A, S) = \min_{A, S} \|X - AS\|^2 + \lambda \|S\|_1$$

- A is an over-complete dictionary, but let us assume that it is a union of ortho-normal bases, in the form

$$A = [A_1 | A_2 | \dots | A_M]$$

$$\forall i, 1 \leq i \leq M, A_i A_i^T = I$$

Learning the Bases: Method 3- Union of Ortho-normal Bases

- The coefficient matrix \mathbf{S} can now be written as follows (M subsets, each corresponding to a single orthonormal basis):

$$\mathbf{S} = \begin{pmatrix} \mathbf{S}_1 \\ \mathbf{S}_2 \\ \vdots \\ \vdots \\ \mathbf{S}_M \end{pmatrix}$$

- Assuming we have fixed bases stored in \mathbf{A} , the coefficients in \mathbf{S} can be estimated using block coordinate descent, described on the following slide.

Learning the Bases: Method 3- Union of Ortho-normal Bases

```
for  $t = 1:T\{$ 
```

$$\lambda_t = \lambda_0(1 - t/T);$$

```
Update  $\mathbf{S}$  using BCR with parameter  $\lambda_t$ 
```

```
\}
```

```
BCR( $\lambda_t, \mathbf{S}\{$ 
```

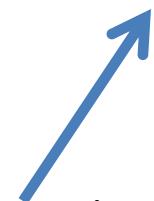
```
for  $m = 1:M\{$ 
```

$$\mathbf{X}_m = \mathbf{X} - \sum_{j \neq m} \mathbf{A}_j \mathbf{S}_j$$

There is a quick way of performing this optimization given an ortho-normal basis – SOFT THRESHOLDING (could be replaced by hard thresholding if you had a stronger sparseness prior than an L1 norm)

$$\text{Update } \mathbf{S}_m \text{ as follows : } \mathbf{S}_m = \arg \min_{\mathbf{S}^*} \|\mathbf{X}_m - \mathbf{A}_m \mathbf{S}^*\|^2 + \lambda_t \|\mathbf{S}^*\|_1$$

```
\}
```



Learning the Bases: Method 3- Union of Ortho-normal Bases

- Given the coefficients, we now want to update the dictionary which is done as follows:

for $m = 1 : M$

{

$$\mathbf{X}_m = \mathbf{X} - \sum_{j \neq m} \mathbf{A}_j \mathbf{S}_j;$$

$$\mathbf{S}_m \mathbf{X}_m^T = \mathbf{U} \Lambda \mathbf{V}^T;$$

$$\mathbf{A}_m = \mathbf{V} \mathbf{U}^T;$$

}

Why are we doing this? It is related to the so-called orthogonal Procrustes problem - a well-known application of SVD. We will see this on the next slide.

The specific problem we are solving is given below. Note that it cannot be solved using a pseudo-inverse as that will not impose orthonormality constraint, if there is noise in the data or if the coefficients are perturbed or thresholded.

$$\min_{\mathbf{A}} \|\mathbf{X}_m - \mathbf{A} \mathbf{S}_m\|^2 \text{ s.t. } \mathbf{A} \mathbf{A}^T = \mathbf{I}$$

$$A^* = \min_A \|X - AS\|_F^2 \text{ s.t. } A^T A = I$$

$$\min_A \|X - AS\|_F^2 = \min_A \text{trace}((X - AS)^T (X - AS))$$

$$= \min_A \text{trace}(X^T X - 2X^T AS + S^T S)$$

$$= \max_A \text{trace}(X^T AS)$$

$$= \max_A \text{trace}(ASX^T) \quad \because \text{trace}(FG) = \text{trace}(GF)$$

-- Let $SX^T = Q$, SVD of Q gives $Q = UDV^T$ --

$$= \max_A \text{trace}(AUDV^T)$$

$$= \max_A \text{trace}(V^T AUD)$$

$$= \max_A \text{trace}(Z(A)D) \text{ where } Z(A) = V^T AU$$

$$= \max_A \sum_i z_{ii} d_{ii} \leq \sum_i d_{ii} (\because Z(A)^T Z(A) = I)$$

The maximum is achieved for $Z(A) = I$, i.e.

$$V^T AU = I \Rightarrow A = VU^T$$

Learning the Bases: Method 3- Union of Ortho-normal Bases

- Keeping all bases in \mathbf{A} fixed, update the coefficients in \mathbf{S} using a known sparse coding technique.
- Keeping the coefficients in \mathbf{S} fixed, update the bases in \mathbf{A} using the aforementioned SVD-based method.
- Repeat the above two steps until a convergence criterion is reached.

Learning the Bases: Method 4 – K-SVD

- Recall: we want to learn a dictionary and *sparse* codes on that dictionary given some data-points:

$$\min_{A,S} \|Y - AS\|^2 \text{ subject to } \forall i, \|s_i\|_0 \leq T_0$$

- Starting with a fixed dictionary, sparse coding follows as usual – OMP, BP, MP etc. The criterion could be based on reconstruction error or L0-norm of the sparse codes.
- The dictionary is updated one column at a time.

$$\begin{aligned}
 \|Y - AS\|_F^2 &= \left\| Y - \sum_{j=1}^K a_j s^j \right\|_F^2 = \left\| Y - \sum_{j \neq k}^K a_j s^j - a_k s^k \right\|_F^2 \\
 &= \|E_k - a_k s^k\|_F^2
 \end{aligned}$$

Does NOT depend
on the k-th
dictionary column

Row 'k' (NOT
column) of
matrix S

How to find a_k , given the above expression? We have decomposed the original error matrix, i.e. $Y-AS$, into a sum of rank-1 matrices, out of which only the last term depends on a_k . So we are trying to find a rank-1 approximation for E_k , and this can be done by computing the SVD of E_k , and using the singular vectors corresponding to the largest singular value.

$$E_k = U \Lambda V^T, a_k = u_1, s^k = \Lambda(1,1)v_1^T$$

Problem! The dictionary codes may no more be sparse! SVD does not have any in-built sparsity constraint in it! So, we proceed as follows:

$$\omega_k = \{i \mid s^k(i) \neq 0\}$$

Consider the error matrix defined as follows:

$$\|Y_{\omega_k} - AS_{\omega_k}\|_F^2 = \left\| \left(Y - \sum_{j \neq k} a_j s^j \right)_{\omega_k} - a_k s_{\omega_k}^k \right\|_F^2 = \|E_{k, \omega_k} - (a_k s_{\omega_k}^k)\|_F^2$$

$$E_{k, \omega_k} = U \Lambda V^T$$

$$a_k \leftarrow u_1, s_{\omega_k}^k \leftarrow \Lambda(1,1)v_1^T$$

Considers only those columns of Y (i.e. only those data-points) that actually **USE** the k -th dictionary atom, effectively yielding a smaller matrix, of size p by $|\omega_k|$

This update affects the sparse codes of only those data-points that used the k -th dictionary element, i.e. those that are part of ω_k .

Task: Find the best dictionary to represent the data samples $\{\mathbf{y}_i\}_{i=1}^N$ as sparse compositions, by solving

$$\min_{\mathbf{D}, \mathbf{X}} \{ \|\mathbf{Y} - \mathbf{DX}\|_F^2 \} \quad \text{subject to } \forall i, \|\mathbf{x}_i\|_0 \leq T_0.$$

Initialization : Set the dictionary matrix $\mathbf{D}^{(0)} \in \mathbb{R}^{n \times K}$ with ℓ^2 normalized columns. Set $J = 1$.

Repeat until convergence (stopping rule):

- *Sparse Coding Stage:* Use any pursuit algorithm to compute the representation vectors \mathbf{x}_i for each example \mathbf{y}_i , by approximating the solution of

$$i = 1, 2, \dots, N, \quad \min_{\mathbf{x}_i} \{ \|\mathbf{y}_i - \mathbf{Dx}_i\|_2^2 \} \quad \text{subject to } \|\mathbf{x}_i\|_0 \leq T_0.$$

- *Codebook Update Stage:* For each column $k = 1, 2, \dots, K$ in $\mathbf{D}^{(J-1)}$, update it by
 - Define the group of examples that use this atom, $\omega_k = \{i \mid 1 \leq i \leq N, \mathbf{x}_T^k(i) \neq 0\}$.
 - Compute the overall representation error matrix, \mathbf{E}_k , by

$$\mathbf{E}_k = \mathbf{Y} - \sum_{j \neq k} \mathbf{d}_j \mathbf{x}_T^j.$$

- Restrict \mathbf{E}_k by choosing only the columns corresponding to ω_k , and obtain \mathbf{E}_k^R .
- Apply SVD decomposition $\mathbf{E}_k^R = \mathbf{U} \Delta \mathbf{V}^T$. Choose the updated dictionary column $\tilde{\mathbf{d}}_k$ to be the first column of \mathbf{U} . Update the coefficient vector \mathbf{x}_R^k to be the first column of \mathbf{V} multiplied by $\Delta(1, 1)$.
- Set $J = J + 1$.

http://www.cs.technion.ac.il/~elad/publications/journals/2004/32_KSV_D_IEEE_TSP.pdf

KSVD and K-means

- Limit the sparsity factor $T_0 = 1$.
- Enforce all the sparse codes to be either 0 or 1.
- Then you get the K-means algorithm! ☺

Implementation issues

- KSVD is a popular and effective method. But some implementation issues haunt K-SVD (life is never easy 😊).
- KSVD is susceptible to **local minima and over-fitting** if K is too large (just like you can get meaningless clusters if your number of cluster is too small, or get meaningless densities if the number of histogram bins is too many).
- KSVD convergence is not fully guaranteed. The dictionary updates given fixed sparse codes ensure the error decreases. However the sparse codes given fixed dictionary may not decrease the error – it is affected by the behaviour of the sparse coding approximation algorithms.
- You can speed up the algorithm by removing extremely “unpopular” dictionary elements, or removing duplicate (or near-duplicate) columns of the dictionary.

(Many) Applications of KSVD

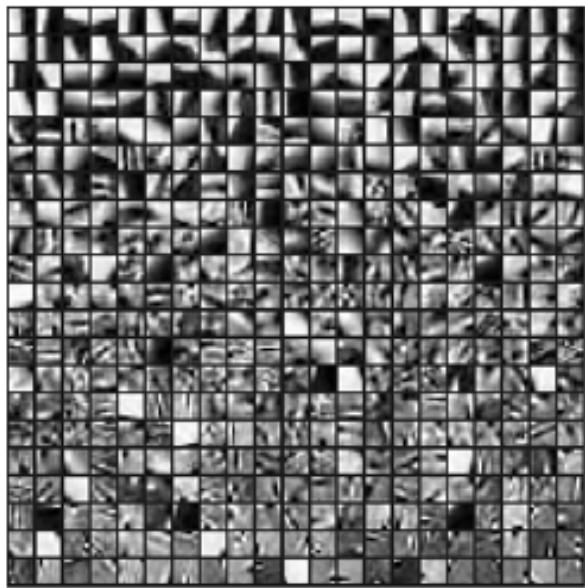
- Image denoising
- Image inpainting
- Image deblurring
- Blind compressive sensing
- Classification
- Compression
- And you can work on many more ☺

Application: Image Compression

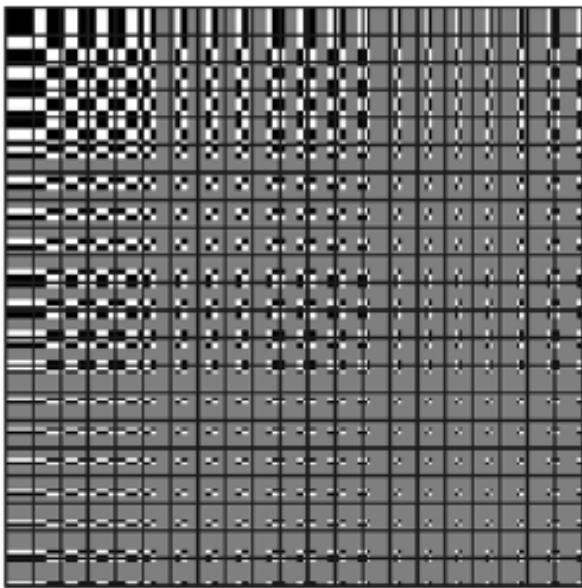
- Recall JPEG algorithm
- It divides the image to be compressed into non-overlapping patches of size 8×8 .
- It computes the 2D-DCT coefficients of each patch, quantizes them and stores only the non-zero coefficients – using Huffman encoding.
- KSVD can be used in image compression by replacing the orthonormal 2D-DCT dictionary by an overcomplete dictionary trained on a database of similar images.

Application: Image Compression (Training Phase)

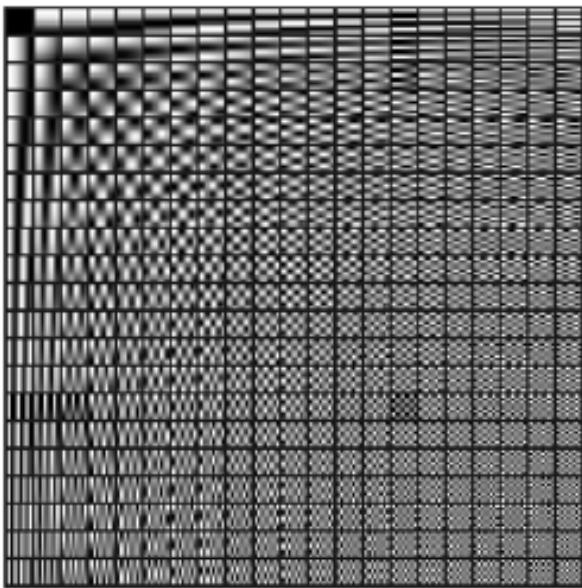
- Training set for dictionary learning: a set of 11000 patches of size 8×8 – taken from a face image database. Dictionary size $K = 441$ atoms (elements).
- OMP used in the sparse coding step during training – stopping criterion is a fixed number of coefficients $T_0 = 10$.
- Over-complete Haar and DCT dictionaries – of size 64×441 – and ortho-normal DCT basis of size 64×64 (JPEG), also used for comparison.



(a)



(b)



(c)

Fig. 5. (a) The learned dictionary. Its elements are sorted in an ascending order of their variance and stretched to maximal range for display purposes. (b) The overcomplete separable Haar dictionary and (c) the overcomplete DCT dictionary are used for comparison.

Application: Image Compression (Testing Phase)

- A lossy image compression algorithm is evaluated using an ROC curve – the X axis contains the average number of bits (BPP) to store the signal.
- The Y-axis is the compression error or PSNR. Normally, the acceptable error e is fixed and the number of bits is calculated.
- The test image is divided into non-overlapping patches of size 8×8 .
- Each patch is projected onto the trained dictionary and its sparse code is obtained using OMP given a fixed error e .

Application: Image Compression (Testing Phase)

- The encoded image contains the following:
 1. Sparse-code coefficients for each patch and the indices of each non-zero coefficient (in the dictionary).
 2. The number of coefficients used to represent each patch (different patches will need different number of coefficients).

Application: Image Compression (Testing Phase)

- The average number of bits per pixel (BPP) is calculated as:

$$BPP = \frac{a \bullet \# \text{ patches} + \# \text{ coeffs} \bullet (b + Q)}{\# \text{ pixels}}$$

Sum total of the number of coefficients representing each patch

Number of bits required to store the number of coefficients for each patch

Number of bits required to store the dictionary index for each coefficient

Number of bits required to code each coefficient (quantization level).
The quantization makes this a lossy compression algorithm.

Huffman encoding

The diagram illustrates the calculation of BPP. It shows the formula $BPP = \frac{a \bullet \# \text{ patches} + \# \text{ coeffs} \bullet (b + Q)}{\# \text{ pixels}}$. Blue arrows point from the terms in the formula to three boxes: the first term points to 'Number of bits required to store the number of coefficients for each patch'; the second term points to 'Number of bits required to code each coefficient (quantization level)'. A blue arrow also points from the term $\# \text{ coeffs} \bullet (b + Q)$ to the third box, which contains the text 'The quantization makes this a lossy compression algorithm.' Red arrows point from the bottom two boxes up to the word 'Huffman' in the 'Huffman encoding' box.

http://www.cs.technion.ac.il/~elad/publications/journals/2004/32_KSVD_IIEEE_TS_P.pdf

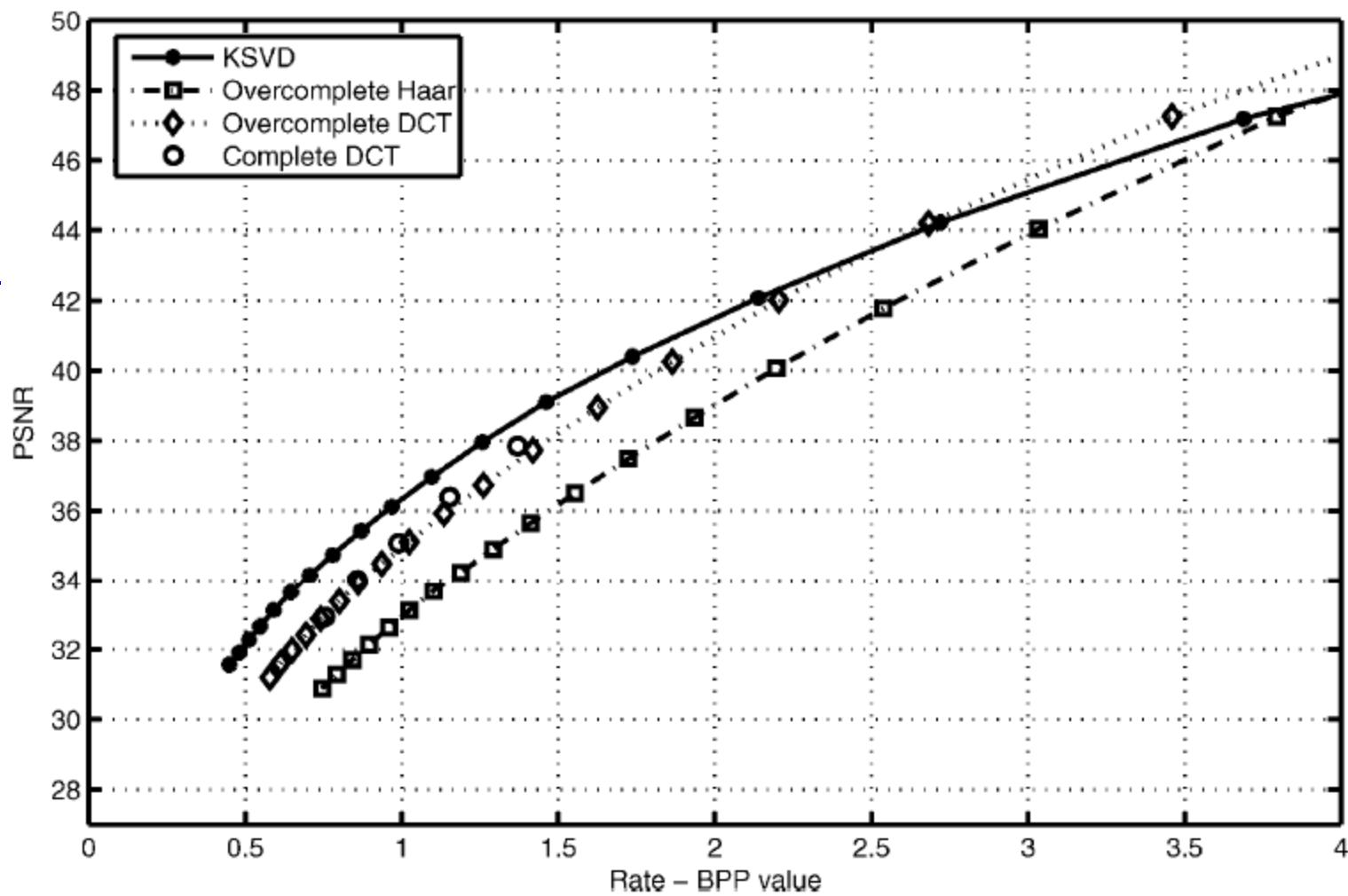


Fig. 8. Compression results: rate-distortion graphs.

Even if the error 'e' for the OMP was fixed, we need to compute the total MSE between the true and the compressed image. This is due to effects of quantization while storing the sparse coefficient values for each patch.

http://www.cs.technion.ac.il/~elad/publications/journals/2004/32_KSVD_IEEE_TSP.pdf

K-SVD dictionary
8 bits per coefficients
PSNR =34.1564
Rate = 0.70651 BPP

Overcomplete DCT dictionary
8 bits per coefficients
PSNR =32.4021
Rate = 0.69419 BPP

Complete DCT dictionary
8 bits per coefficients
PSNR =32.3917
Rate = 0.70302 BPP



Fig. 9. Sample compression results.

Application: Image Denoising

- KSVD for denoising seeks to minimize the following objective function:

$$(X, D, \{\alpha_{ij}\}) = \arg \min_{X, D, \{\alpha_{ij}\}} \lambda \|X - Y\|^2 + \sum_{i,j} \mu_{ij} \|\alpha_{ij}\|_0 + \sum_{i,j} \|D\alpha_{ij} - R_{ij} X\|^2$$

Y = noisy image

X = underlying clean image (to be estimated)

α_{ij} = sparse dictionary coefficients for patch at location (i,j)

D = dictionary

R_{ij} = matrix that extracts the patch x_{ij} from image X, i.e.

$x_{ij} = R_{ij} X$

Application: Image Denoising

- Note: The dictionary may be learned a priori from a corpus of image patches. The patches from the noisy image can then be denoised by mere sparse coding.
- The more preferable method is to train the dictionary directly on the noisy image in tandem with the sparse coding step (as in the previous slide).
- This avoids having to depend on the training set and allows for tuning of the dictionary to the underlying image structure (as opposed to the structure of some other images).

KSVD Algorithm for Denoising (Dictionary learned on the noisy image)

- Set $X = Y$, $D = \text{overcomplete DCT}$
- Until some “convergence” criterion is satisfied, repeat the following:
 1. Obtain the sparse codes for every patch (typically using OMP) as follows:

$$\forall i, j, \min \left\| \alpha_{ij} \right\|_0 \text{ s.t. } \left\| R_{ij} X - D \alpha_{ij} \right\|^2 \leq C \sigma^2$$

2. Perform the dictionary learning update typical for KSVD.
 - Estimate the final image X by averaging the reconstructed overlapping patches, OR estimate X given D and α_{ij} :

$$X = \arg \min_X \lambda \|X - Y\|^2 + \sum_{ij} \|D \alpha_{ij} - R_{ij} X\|^2$$

$$\Rightarrow X = (\lambda I + \sum_{ij} R_{ij}^T R_{ij})^{-1} (\lambda Y + \sum_{ij} R_{ij}^T D \alpha_{ij})$$

KSVD Algorithm for Denoising (Dictionary learned on the noisy image)

$$X = \arg \min_X \lambda \|X - Y\|^2 + \sum_{ij} \|D\alpha_{ij} - R_{ij}X\|^2$$

$$\Rightarrow X = (\lambda I + \sum_{ij} R_{ij}^T R_{ij})^{-1} (\lambda Y + \sum_{ij} R_{ij}^T D\alpha_{ij})$$

$$\lambda \text{ set to } \frac{30}{\sigma}$$

This equation is a mathematically rigorous way to show how X was reconstructed by averaging overlapping denoised patches together with the noisy image as well.

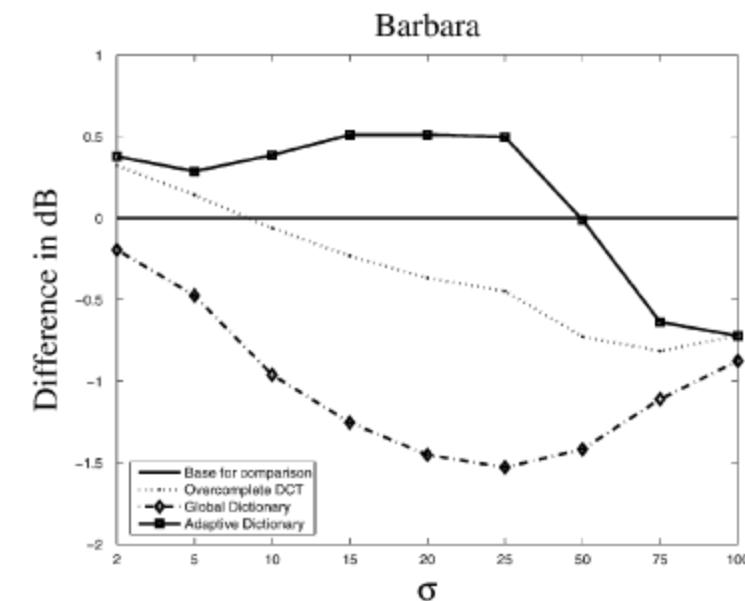
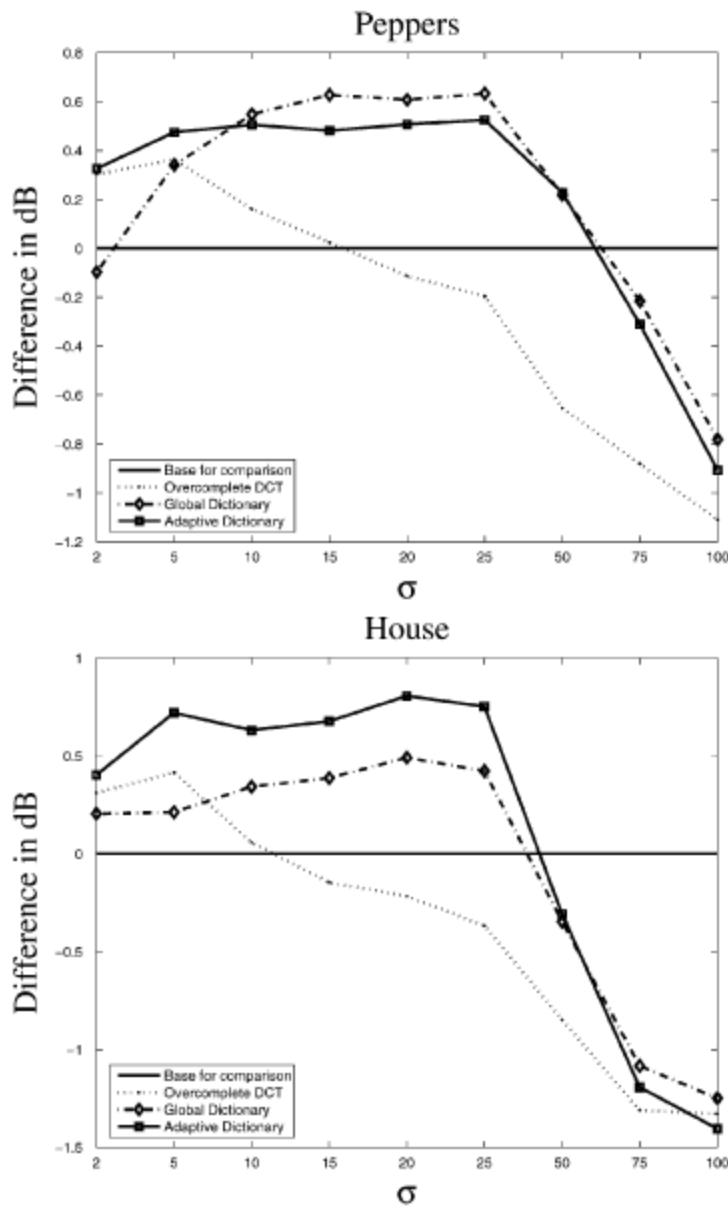


Fig. 4. Comparison between the three presented methods (overcomplete DCT, global trained dictionary, and adaptive dictionary trained on patches from the noisy image) and the results achieved recently in [23] for three test images.

Baseline for comparison: method by Portilla et al, "Image denoising using scale mixture of Gaussians in the wavelet domain", IEEE TIP 2003.

http://www.cs.technion.ac.il/~elad/publications/journals/2005/KSVD_Denoising_IEEE_TIP.pdf

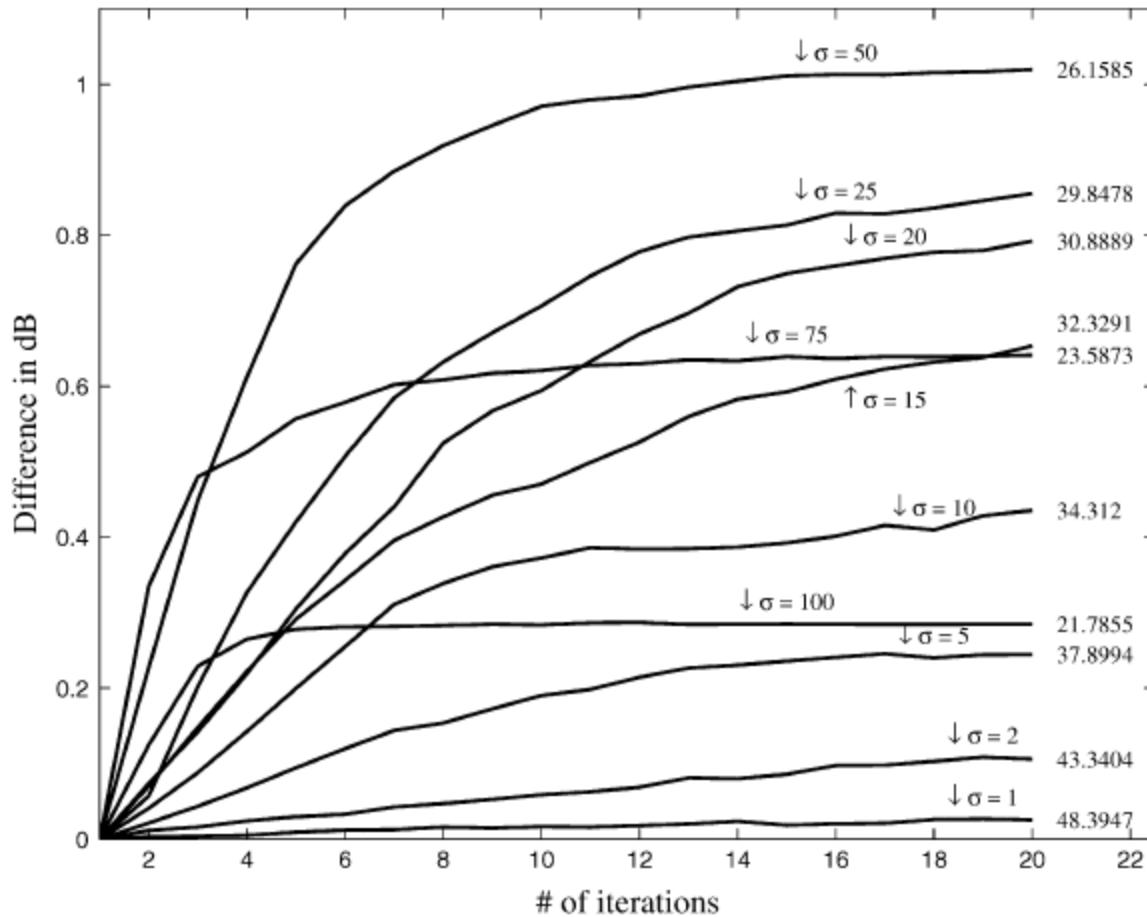


Fig. 5. Improvement in the denoising results after each iteration of the K-SVD algorithm, executed on noisy patches of the image ‘Peppers.’

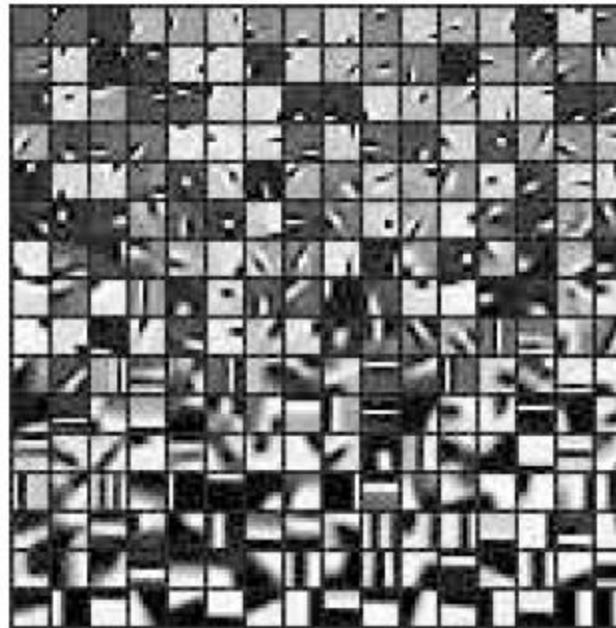
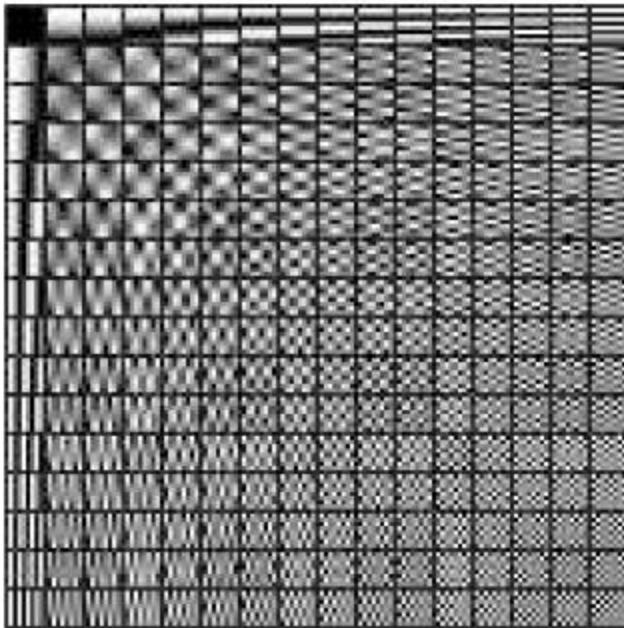
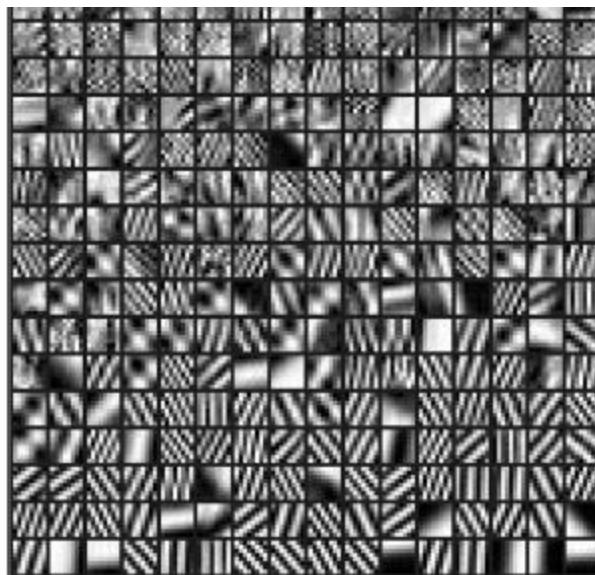


Fig. 2. Left: Overcomplete DCT dictionary. Right: Globally trained dictionary.



http://www.cs.technion.ac.il/~elad/Variou s/KSVD_Matlab_ToolBox.zip

Fig. 7. Example of the denoising results for the image “Barbara” with $\sigma = 20$ —the adaptively trained dictionary.

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)



<http://www.cs.technion.ac.il/~elad/publications/journals/2005/KSVD Denoising IEEE TIP.pdf>

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

Application of KSVD: Filling in Missing Pixels

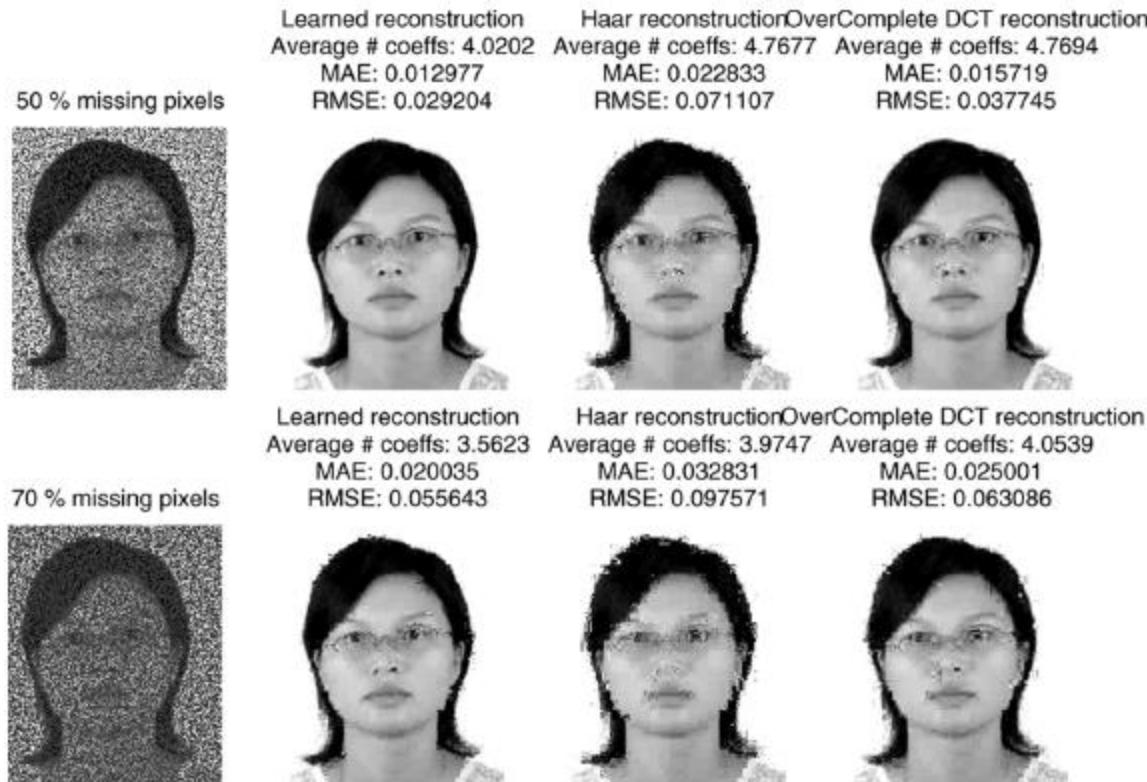


Fig. 7. The corrupted image (left) with the missing pixels marked as points and the reconstructed results by the learned dictionary, the overcomplete Haar dictionary, and the overcomplete DCT dictionary, respectively. The different rows are for 50% and 70% of missing pixels.

http://www.cs.technion.ac.il/~elad/publications/journals/2004/32_KSVD_IIEEE_TSP.pdf

Application of KSVD: Filling in Missing Pixels

- An over-complete dictionary is trained a priori on a set of face images.
- A test face image (not part of the training set) is synthetically degraded by masking out 50 to 70 percent of the pixels.
- Patches from the degraded image are sparse coded by projection onto the trained dictionary using OMP.
- OMP is modified so that only the non-degraded pixels are considered during any error computation (the dictionary elements are therefore appropriately re-scaled to unit norm).

Application of KSVD: Filling in Missing Pixels

$y_{ij} = n^2 \times 1$ patch in corrupted image at location (i, j)

$\tilde{y}_{ij} = /P \times 1$ vector containing the P non-corrupted intensities from y_{ij}

$\therefore \tilde{y}_{ij} = \Phi_{ij} D s_{ij}$ where $\Phi_{ij} \in R^{P \times n^2}, D \in R^{n^2 \times K}, s_{ij} \in R^{K \times 1}$

$$E(\{s_{ij}\}) = \sum_{i,j} \left\| \tilde{y}_{ij} - \Phi_{ij} D s_{ij} \right\|^2 \text{ s.t. } \|s_{ij}\|_0 \leq T_0$$

D learned a-priori

Sparse codes determined using OMP or ISTA

Application of KSVD: Filling in Missing Pixels

- The patches are reconstructed in the following manner:

$$y_{ij} = Ds_{ij}$$

NOTE: Dictionary elements without masking!!

KSVD Application: Coded Exposure Image

It is a coded superposition (i.e. summation) of T sub-frames within a unit integration time of the video camera.

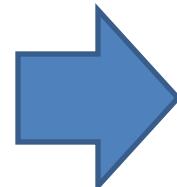
Coded exposure image
(captured in one unit
integration time of the
camera)

$$I(x, y) = \sum_{t=1}^T S(x, y, t) \bullet E(x, y, t)$$

Binary code at
time instant t

Sub-frame at
time instant t

Conventional capture
(simple integration across
time, without modulation
by binary codes)



$$\forall x, y, t S(x, y, t) = 1$$

KSVD Application: Sparse Coding in Video CS

$$I(x, y) = \sum_{t=1}^T S(x, y, t) \bullet E(x, y, t) \rightarrow \mathbf{y} = \Phi \mathbf{f}$$

$$\mathbf{f} = \Psi \boldsymbol{\theta}$$

$$\therefore \mathbf{y} = \Phi \Psi \boldsymbol{\theta}$$

\mathbf{y} is a vectorized form of the snapshot image \mathbf{I} of n pixels.

\mathbf{f} is a vectorized form of the underlying video with T frames and each frame having n pixels, hence \mathbf{f} has nT pixels in total.

Ψ is a 3D DCT basis – the basis can also be “learned offline on a representative set of videos”

Φ is a matrix of size $n \times nT$ containing values from the binary code \mathbf{S} . See below:

$$\Phi = (\text{diag}(S_1) | \text{diag}(S_2) | \dots | \text{diag}(S_T))$$

$$S_t(x, y) = S(x, y, t) \text{ for } t = 1 \text{ to } T$$

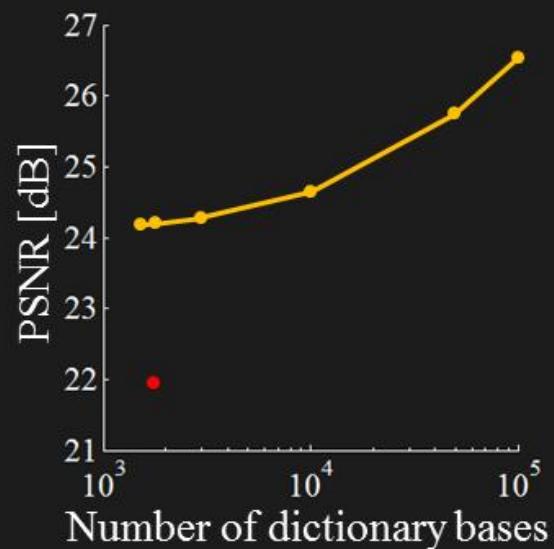
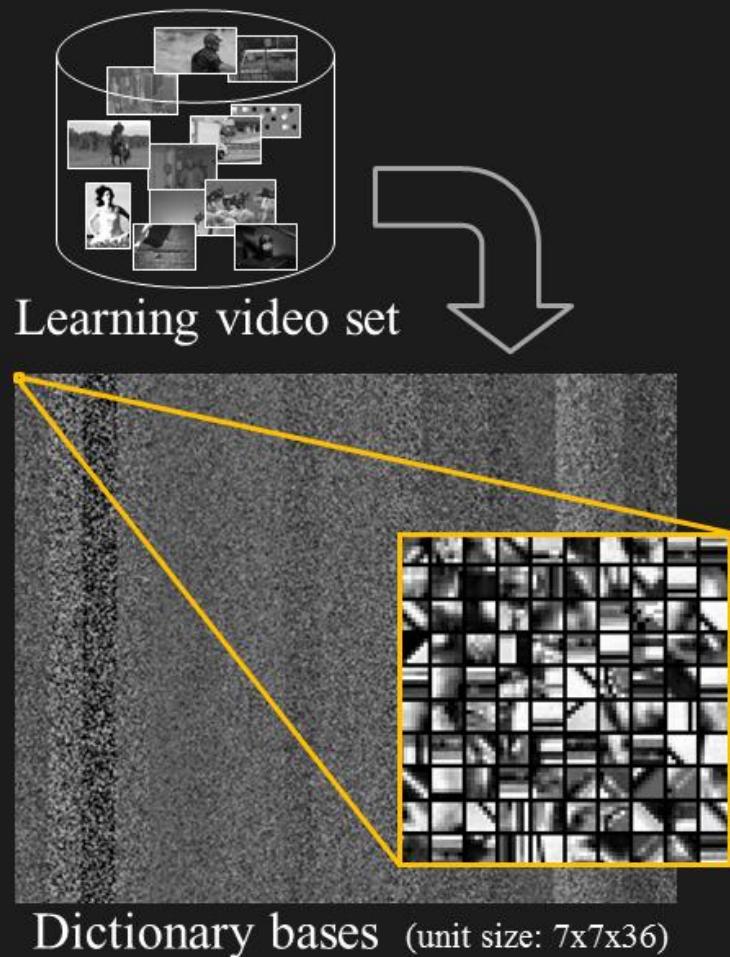
$\text{Diag}(S_t)$ represents a $n \times n$ diagonal matrix. The diagonal contains elements from the mask S_t .

Dictionary Learning

- Done offline – training set was 20 video sequences, each video rotated in 8 directions and played forward + backward = 320 videos.
- All videos had **target frame rate** (2000 fps, as we work with a 60 fps camera and want 36 fold gain).
- Video-patch size was $7 \times 7 \times 36 = 1764 \times 1$
- Offline learning: KSVD (*), $K = 100,000$ atoms
- Sparse coding done online (using OMP)

KSVD is a dictionary learning technique. Given a set of input patches in $n \times N$, it learns a set of vectors ($n \times K$), such that a sparse linear combination of these vectors approximates each patch as closely as possible. There are K coefficients per patch, out of which most are encouraged to be 0 (sparse).

Learned Over-complete Dictionary



- Learned over-complete dictionary
- DCT based dictionary

http://www.cs.columbia.edu/CAVE/projects/single_shot_video/

Video reconstruction results on real data available below:

<https://www.youtube.com/watch?v=JAYC0C3NIdY>

- Ideally, the videos whose patches are used for training must have the TARGET frame rate.
- Suppose we have a 60 fps video camera and want to increase its temporal resolution to say $60T$ fps. This would require coded versions of T consecutive sub-frames in the Hitomi camera (cf: slides on CS systems) to be superposed to produce a snapshot.
- Thus from a single snapshot with $n \times n$ pixels, we would like to reconstruct a video of size $n \times n \times T$. The dictionary needs to be trained on patches of size $m \times m \times T$ (m much less than n).
- Thus ideally, we would require the availability of videos acquired at a frame rate of $60T$ (or close to it).
- If you wish to design a system with a different value of T (say T'), you may ideally need to re-train the dictionary on patches taken from videos with appropriate frame rate.

Blind compressed sensing

- In standard compressed sensing, we assume knowledge of the basis (may be orthonormal or even overcomplete) in which the signal is sparse.
- However this knowledge may not be available to us always.
- For example, consider a remote sensing application with compressive acquisitions of water bodies followed by forests over mountain ranges.
- Images of the former are likely to be highly sparse in DCT whereas the latter may not be – and another basis may serve a better purpose.
- Blind compressed sensing is the task of inferring the signal basis as well as the signal coefficients directly from the compressed data.

Blind compressed sensing

- Consider N signals and their compressive measurements in the following form:

$$y_i = \Phi_i x_i + \eta_i, \Phi_i \in R^{m \times n}, x_i \in R^n, y_i \in R^m, m \ll n$$

- Our sparsity model for each signal is as follows:

$$x_i = D\theta_i, D \in R^{n \times K}, \theta_i \in R^K$$

- Combined model is:

$$y_i = \Phi_i D\theta_i + \eta_i$$

Blind compressed sensing

- The objective function to be minimized is:

$$E(D, \{\theta_i\}) = \sum_{i=1}^N \|y_i - \Phi_i D \theta_i\|^2 \text{ s.t. } \|\theta_i\|_0 \leq T, d_k^t d_k = 1$$

- The algorithm for this is alternating in nature: sparse coding assuming a fixed dictionary, and dictionary updates assuming sparse codes.
- Sparse coding: OMP

$$E_1(\theta_i) = \|y_i - \Phi_i D \theta_i\|^2 \text{ s.t. } \|\theta_i\|_0 \leq T$$

Blind compressed sensing

- The dictionary update is as follows:

$$\begin{aligned} E(d_k) &= \sum_{i=1}^N \left\| y_i - \Phi_i \sum_{l=1}^K d_l \theta_{il} \right\|^2 \\ &= \sum_{i=1}^N \left\| y_i - \Phi_i \sum_{l=1, l \neq k}^K d_l \theta_{il} - \Phi_i d_k \theta_{ik} \right\|^2 \\ &= \sum_{i=1}^N \left\| y_i^{(-k)} - \Phi_i d_k \theta_{ik} \right\|^2 \end{aligned}$$

Blind compressed sensing

- Taking derivative w.r.t. d_k and setting to 0 gives:

$$-2 \sum_{i=1}^N \Phi_i^t (y_i^{(-k)} - \Phi_i d_k \theta_{ik}) \theta_{ik} = 0$$

$$\rightarrow d_k = \left(\sum_{i=1}^N \Phi_i^t \Phi_i \theta_{ik}^2 \right)^{-1} \sum_{i=1}^N \Phi_i^t y_i^{(-k)} \theta_{ik}$$

$$d_k \leftarrow d_k / \|d_k\|_2$$

$$\theta_{ik} = \frac{d_k^t \Phi_i^t y_i^{(-k)}}{d_k^t \Phi_i^t \Phi_i d_k}$$

This matrix will be invertible with probability 1 only if N is such that $Nm \geq n$. Notice that it is important for every point to have a different sensing matrix, otherwise the dictionary cannot be inferred.

- Repeat for other columns iteratively.

Blind CS: Pseudocode

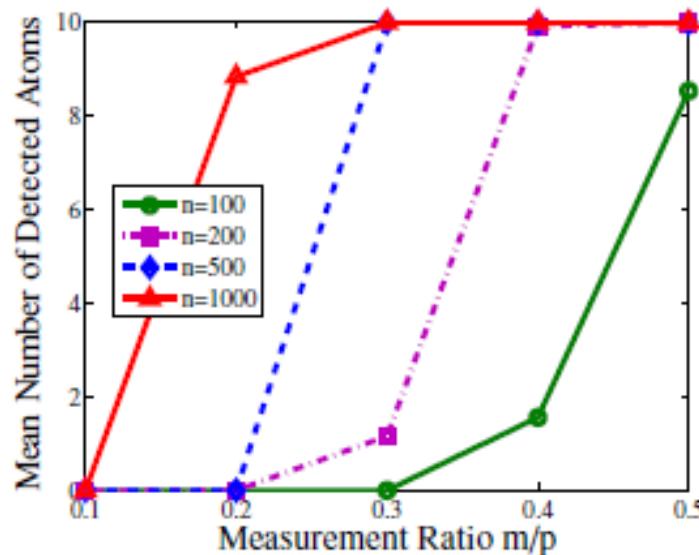
- Start with a random dictionary, infer the sparse codes for each data point.
- Given the sparse codes, infer each column of the dictionary keeping the others fixed. Infer the corresponding dictionary coefficients as well.
- Repeat the previous two steps till a convergence criterion is reached.

Blind compressed sensing: experiments

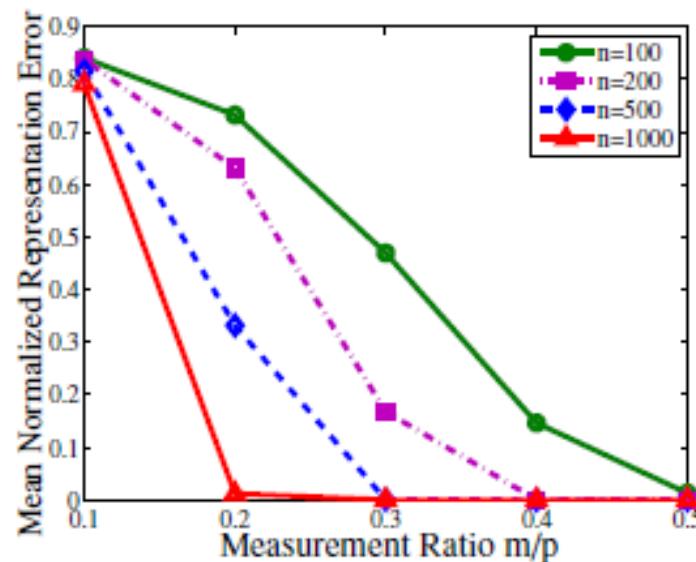
- Synthetic dictionary of size 50 by 10 ($K = 10$) generated iid from $N(0,1)$.
- Training signals generated using a linear combination of these dictionary columns, with coefficients generated from $N(0,80)$.
- Corresponding compressive measurements are generated using a different random projection for each training signal.
- The dictionary and coefficients are inferred from the compressed data, starting with random initial conditions.

Blind compressed sensing: experiments

- The inferred dictionary columns are compared to the true dictionary columns using dot products.
- The number of correctly identified (i.e. dot product of more than 0.95) dictionary columns is computed.
- Plots seen on next slide.



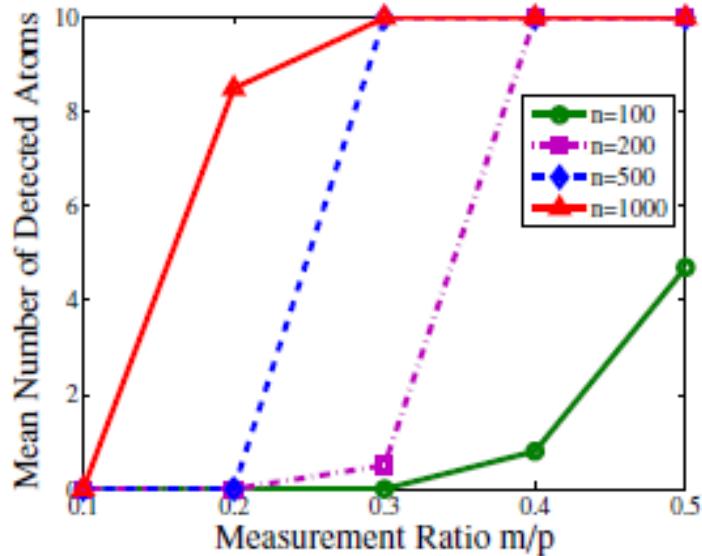
(a)



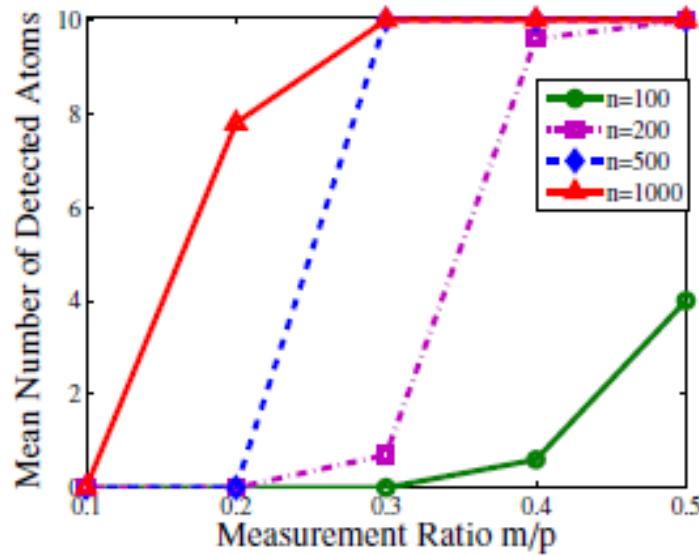
(b)

Fig. 2: Results for synthetic data. (a) Plot of the mean number of detected atoms in 50 trials for varying n and m/p . (b) Plot of mean normalized representation error of the training signals in 50 trials for varying n and m/p .

<https://arxiv.org/pdf/1504.01169.pdf>



(a)



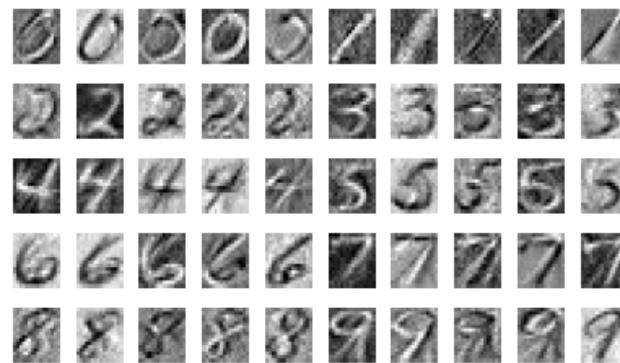
(b)

Fig. 3: Results for noisy CS measurements. Plot of mean number of detected atoms in 10 trials for the case that (a) $\text{SNR} = 50 \text{ dB}$, and (b) $\text{SNR} = 30 \text{ dB}$.

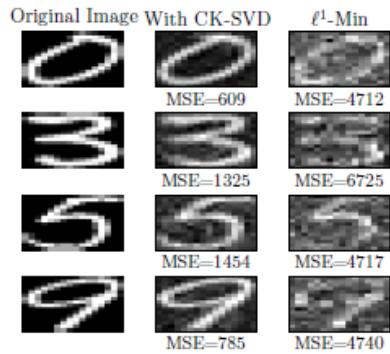
<https://arxiv.org/pdf/1504.01169.pdf>



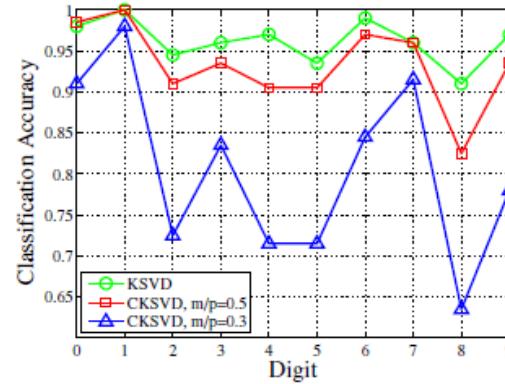
(a) Samples



(b) The Learned Dictionary



(c) CS Recovery



(d) Classification

Fig. 4: Results for the USPS dataset. (a) Samples from the USPS dataset. (b) Samples from the learned dictionary atoms using CK-SVD on ten classes of the USPS dataset for the measurement ratio $m/p = 0.5$. (c) Reconstruction accuracy comparison of OMP using the CK-SVD learned dictionary and ℓ^1 -minimization using Daubechies 8 wavelets. (d) Plot of classification accuracy for K-SVD on the data versus CK-SVD on the CS measurements.

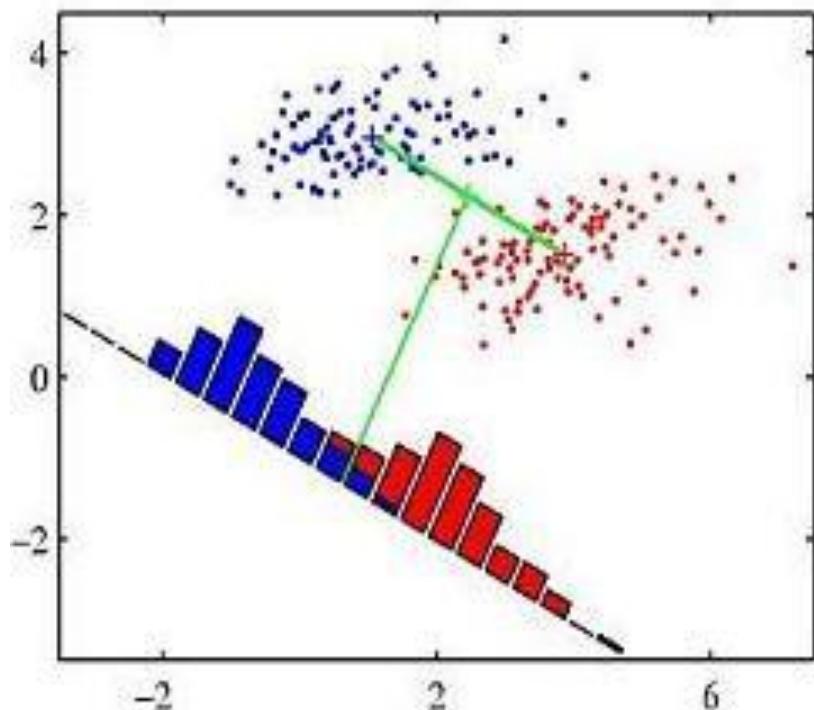
Transform Learning for Classification

Transform learning

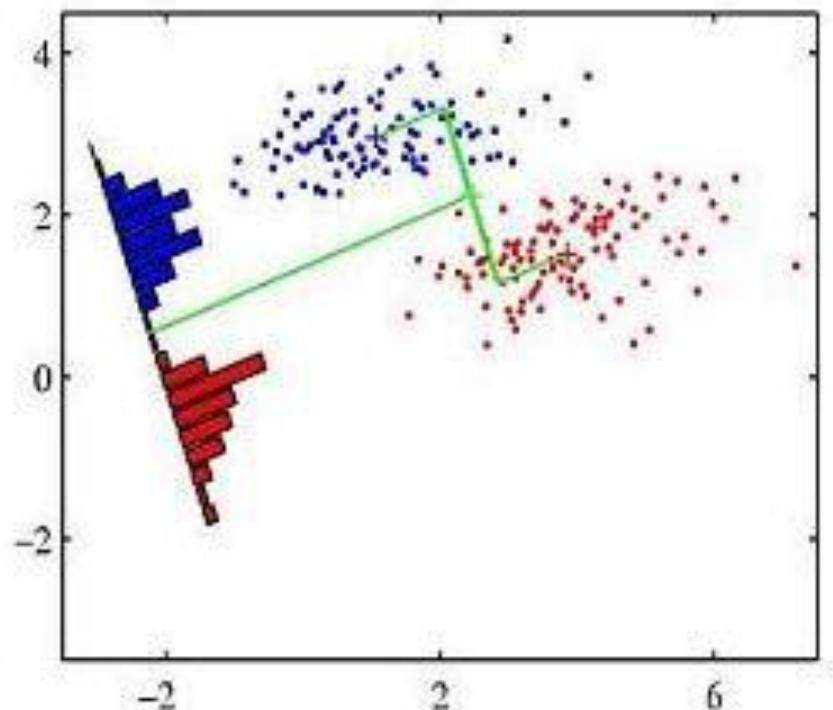
- Techniques like PCA, NMF, KSVD, MOD, etc. are design for (compact) image representation.
- They sometimes yield great results on image classification.
- But they are not designed for optimal image classification.

Transform learning

- Given input points $\{\mathbf{x}_i\}$, $i = 1$ to N , we will learn a transformation vector \mathbf{w} with the aim of optimal classification of the transformed points $y_i = \mathbf{w}^t \mathbf{x}_i$.
- Two techniques for this:
 - Fisher's linear discriminant
 - Mutual information based discriminant



Not so good \mathbf{w}



Good \mathbf{w}

<https://stackoverflow.com/questions/33844198/reproduce-fisher-linear-discriminant-figure>

Fisher's linear discriminant (FLD)

- Also called linear discriminant analysis (LDA).
- We start off considering $c = 2$ classes C1 and C2 and generalize later.
- We have:

$$\mathbf{m}_1 = \frac{\sum_{\mathbf{x}_i \in C_1}^N \mathbf{x}_i}{n_1}, \mathbf{m}_2 = \frac{\sum_{\mathbf{x}_i \in C_2}^N \mathbf{x}_i}{n_2}$$

Fisher's linear discriminant (FLD)

- Consider:

$$y_i = \mathbf{w}^t \mathbf{x}_i;$$

$$\tilde{m}_1 = \frac{\sum_{x_i \in C_1}^N \mathbf{w}^t \mathbf{x}_i}{n_1} = \mathbf{w}^t \mathbf{m}_1,$$

$$\tilde{m}_2 = \frac{\sum_{x_i \in C_2}^N \mathbf{w}^t \mathbf{x}_i}{n_2} = \mathbf{w}^t \mathbf{m}_2$$

Criterion for a good \mathbf{w} :

The transformed mean vectors should be as far apart as possible.

But that is not enough. We want the difference between the transformed mean vectors to be large RELATIVE to the variances of the transformed points of each class. **Otherwise, there may be conflation between the points of the two classes due to a large spread even if their means are far apart.**

Fisher's linear discriminant (FLD)

- Consider:

$$\tilde{s}_1^2 = \sum_{y_i \in C_1} (y_i - \tilde{m}_1)^2, \tilde{s}_2^2 = \sum_{y_i \in C_2} (y_i - \tilde{m}_2)^2$$

- The vector \mathbf{w} should be designed based on maximizing the following criterion:

$$J(\mathbf{w}) = \frac{\|\tilde{\mathbf{m}}_1 - \tilde{\mathbf{m}}_2\|^2}{\tilde{s}_1^2 + \tilde{s}_2^2} \rightarrow = \|\mathbf{w}^t (\mathbf{m}_1 - \mathbf{m}_2)\|^2 = \mathbf{w}^t (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^t \mathbf{w}$$

$= \mathbf{w}^t \mathbf{S}_B \mathbf{w}$

Between-class scatter matrix

$$J(\mathbf{w}) = \frac{\|\tilde{\mathbf{m}}_1 - \tilde{\mathbf{m}}_2\|^2}{\tilde{s}_1^2 + \tilde{s}_2^2} \rightarrow \begin{aligned} &= \|\mathbf{w}^t(\mathbf{m}_1 - \mathbf{m}_2)\|^2 = \mathbf{w}^t(\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^t \mathbf{w} \\ &= \mathbf{w}^t \mathbf{S}_B \mathbf{w} \end{aligned}$$

Between-class scatter matrix

$$\begin{aligned}\tilde{s}_1^2 &= \sum_{y_i \in C_1} \|y_i - \tilde{\mathbf{m}}_1\|^2 = \sum_{y_i \in C_1} \mathbf{w}^t(\mathbf{x}_i - \mathbf{m}_1)(\mathbf{x}_i - \mathbf{m}_1)^t \mathbf{w} = \mathbf{w}^t \mathbf{S}_I \mathbf{w} \\ \tilde{s}_2^2 &= \mathbf{w}^t \mathbf{S}_W \mathbf{w}\end{aligned}$$

$$\tilde{s}_1^2 + \tilde{s}_2^2 = \mathbf{w}^t (\mathbf{S}_I + \mathbf{S}_W) \mathbf{w} = \mathbf{w}^t \mathbf{S}_W \mathbf{w}$$

Within-class scatter matrix

$$\therefore J(\mathbf{w}) = \frac{\mathbf{w}^t \mathbf{S}_B \mathbf{w}}{\mathbf{w}^t \mathbf{S}_W \mathbf{w}}$$

Find \mathbf{w} which maximizes this criterion – often called the Fisher criterion

Fisher's linear discriminant

- Taking derivatives and setting to 0, we have:

$$J(\mathbf{w}) = \frac{\mathbf{w}^t \mathbf{S}_B \mathbf{w}}{\mathbf{w}^t \mathbf{S}_W \mathbf{w}}, \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = 0 \rightarrow$$

$$\mathbf{w}^t \mathbf{S}_W \mathbf{w} (2 \mathbf{S}_B \mathbf{w}) = \mathbf{w}^t \mathbf{S}_B \mathbf{w} (2 \mathbf{S}_W \mathbf{w}) \rightarrow$$

$$\mathbf{S}_B \mathbf{w} = \lambda \mathbf{S}_W \mathbf{w} \rightarrow \mathbf{S}_W^{-1} \mathbf{S}_B \mathbf{w} = \lambda \mathbf{w}$$

$$\lambda = \mathbf{w}^t \mathbf{S}_B \mathbf{w} / \mathbf{w}^t \mathbf{S}_W \mathbf{w}$$

Called a Generalized eigenvalue problem.

Assumes that \mathbf{S}_W is not singular.

Fisher's linear discriminant

- In the $c = 2$ case:

$$S_B = (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^t$$

$$\rightarrow S_B \mathbf{w} = \beta(\mathbf{m}_1 - \mathbf{m}_2) \text{ for some scalar } \beta = (\mathbf{m}_1 - \mathbf{m}_2)^t \mathbf{w}$$

$$\text{Recall: } S_W^{-1} S_B \mathbf{w} = \lambda \mathbf{w} \rightarrow S_W^{-1} \beta(\mathbf{m}_1 - \mathbf{m}_2) = \lambda \mathbf{w}$$

$$\rightarrow \mathbf{w} \text{ points in the same direction as } S_W^{-1}(\mathbf{m}_1 - \mathbf{m}_2)$$

- Thus FLD determines a direction that optimizes the criterion mentioned earlier.
- Since S_B is rank 1, only a single such direction is possible.

Fisher's linear discriminant: classification

- Thus FLD determines a direction that optimizes the criterion mentioned earlier.
- It reduces dimensionality from d to 1.
- The actual classification however proceeds by nearest neighbor method in the transformed space.

Fisher's linear discriminant: case of c classes

- In this case, we project the d -dimensional points \mathbf{x}_i onto a $c-1$ dimensional space (we assume $d \geq c$).
- Generalizing the within-class scatter matrix:

$$S_W = \sum_{i=1}^c S_i, S_i = \sum_{x_j \in C_i} (\mathbf{x}_j - \mathbf{m}_i)(\mathbf{x}_j - \mathbf{m}_i)^t,$$
$$\mathbf{m}_i = \sum_{x_j \in C_i} \mathbf{x}_j / n_i$$

Fisher's linear discriminant: case of c classes

- Define the total scatter matrix:

$$\mathbf{m} = \sum_{i=1}^c \sum_{x_j \in C_i} \mathbf{x}_j / N = \sum_{i=1}^c n_i \mathbf{m}_i / N,$$

$$S_T = \sum_{j=1}^N (\mathbf{x}_j - \mathbf{m})(\mathbf{x}_j - \mathbf{m})^t$$

Fisher's linear discriminant: case of c classes

$$\mathbf{m} = \sum_{i=1}^c \sum_{x_j \in C_i} \mathbf{x}_j / N = \sum_{i=1}^c n_i \mathbf{m}_i / N,$$

$$S_T = \sum_{j=1}^N (\mathbf{x}_j - \mathbf{m})(\mathbf{x}_j - \mathbf{m})^t = \sum_{i=1}^c \sum_{x_j \in C_i} (\mathbf{x}_j - \mathbf{m})(\mathbf{x}_j - \mathbf{m})^t$$

$$= \sum_{i=1}^c \sum_{x_j \in C_i} (\mathbf{x}_j - \mathbf{m}_i + \mathbf{m}_i - \mathbf{m})(\mathbf{x}_j - \mathbf{m}_i + \mathbf{m}_i - \mathbf{m})^t$$

$$= \sum_{i=1}^c \sum_{x_j \in C_i} (\mathbf{x}_j - \mathbf{m}_i)(\mathbf{x}_j - \mathbf{m}_i)^t + \sum_{i=1}^c \sum_{x_j \in C_i} (\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})^t$$

$$= S_W + \sum_{i=1}^c n_i (\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})^t = S_W + S_B$$

Fisher's linear discriminant: case of c classes

- We perform a dimensionality reduction from d dimensions to $c-1$ dimensions using a transform matrix \mathbf{W} of size $d \times (c-1)$ where:

$$\forall i, \mathbf{y}_i = \mathbf{W}^T \mathbf{x}_i$$

- We define scatter matrices in the transformed space:

$$\tilde{\mathbf{S}}_W = \mathbf{W}^T \mathbf{S}_W \mathbf{W}, \tilde{\mathbf{S}}_B = \mathbf{W}^T \mathbf{S}_B \mathbf{W}$$

Fisher's criterion: case of c classes

- The Fisher criterion which we maximize w.r.t. \mathbf{W} is:

$$J(\mathbf{W}) = \frac{|\mathbf{W}^T \mathbf{S}_B \mathbf{W}|}{|\mathbf{W}^T \mathbf{S}_W \mathbf{W}|}$$
$$\tilde{\mathbf{S}}_W = \mathbf{W}^T \mathbf{S}_W \mathbf{W}, \tilde{\mathbf{S}}_B = \mathbf{W}^T \mathbf{S}_B \mathbf{W}$$

The determinant of the scatter matrix is a scalar measure of the scatter. As the determinant is the product of eigenvalues, it is a product of variances in the principal directions and thus the square of the scattering volume of the hyperellipsoid.

- This is a generalized eigenvalue problem again.

Fisher's criterion: case of c classes

- The solution is a set of $c-1$ generalized eigenvectors of the following form:

$$\mathbf{S}_B \mathbf{w}_i = \lambda_i \mathbf{S}_W \mathbf{w}_i$$

- There are no more than $c-1$ such eigenvectors with non-zero eigenvalues as the rank of \mathbf{S}_B is at the most $c-1$. (As \mathbf{S}_B is the sum of c rank-one matrices.)

Fisher's linear discriminant: uniqueness

- Is the solution \mathbf{W} unique?
- Answer is no – it is unique only up to an arbitrary non-singular linear transformation because

$$J(\mathbf{W}) = \frac{|\mathbf{W}^T \mathbf{S}_B \mathbf{W}|}{|\mathbf{W}^T \mathbf{S}_W \mathbf{W}|} = \frac{|(\mathbf{WA})^T \mathbf{S}_B \mathbf{WA}|}{|(\mathbf{WA})^T \mathbf{S}_W (\mathbf{WA})|} = \frac{|A| |\mathbf{W}^T \mathbf{S}_B \mathbf{W}| |A|}{|A| |\mathbf{W}^T \mathbf{S}_W \mathbf{W}| |A|}$$

Fisherfaces: applications in face recognition

- Problem: the matrix S_w becomes singular.
- Why? Because the number of images is usually much smaller than the image dimensions!
- PCA to the rescue!

Fisherfaces: applications in face recognition

- Reduce the dimensionality of the d -dimensional points – of which there are N – using PCA.
- Project onto the largest $N-c$ eigenvectors – the former causes no loss of information. (This technique is called “Fisherfaces”- see equation 5 of [this](#).)
- Then perform FLD to reduce dimensionality to $c-1$. See equation on next slide.

Fisherfaces: applications in face recognition

- The criterion for this modified FLD is as follows:

$$W_{pca} = \arg \max_W |W^T S_T W|$$
$$W_{fld} = \arg \max_W \frac{|W^T W_{pca}^T S_B W_{pca} W|}{|W^T W_{pca}^T S_W W_{pca} W|}.$$

 d x (N-c)
matrix

Results with Fisherfaces

- Training data: Database of face images under different lighting conditions.
- Subsets 1, 2, 3, 4, 5 where the latitudinal and longitudinal angles of the light source direction are respectively within 15, 30, 45, 60, 75 degrees of the camera axis.
- Interpolation: Training on sets 1 and 5 – testing on the others.
- Extrapolation: training on set 1, testing on the others.

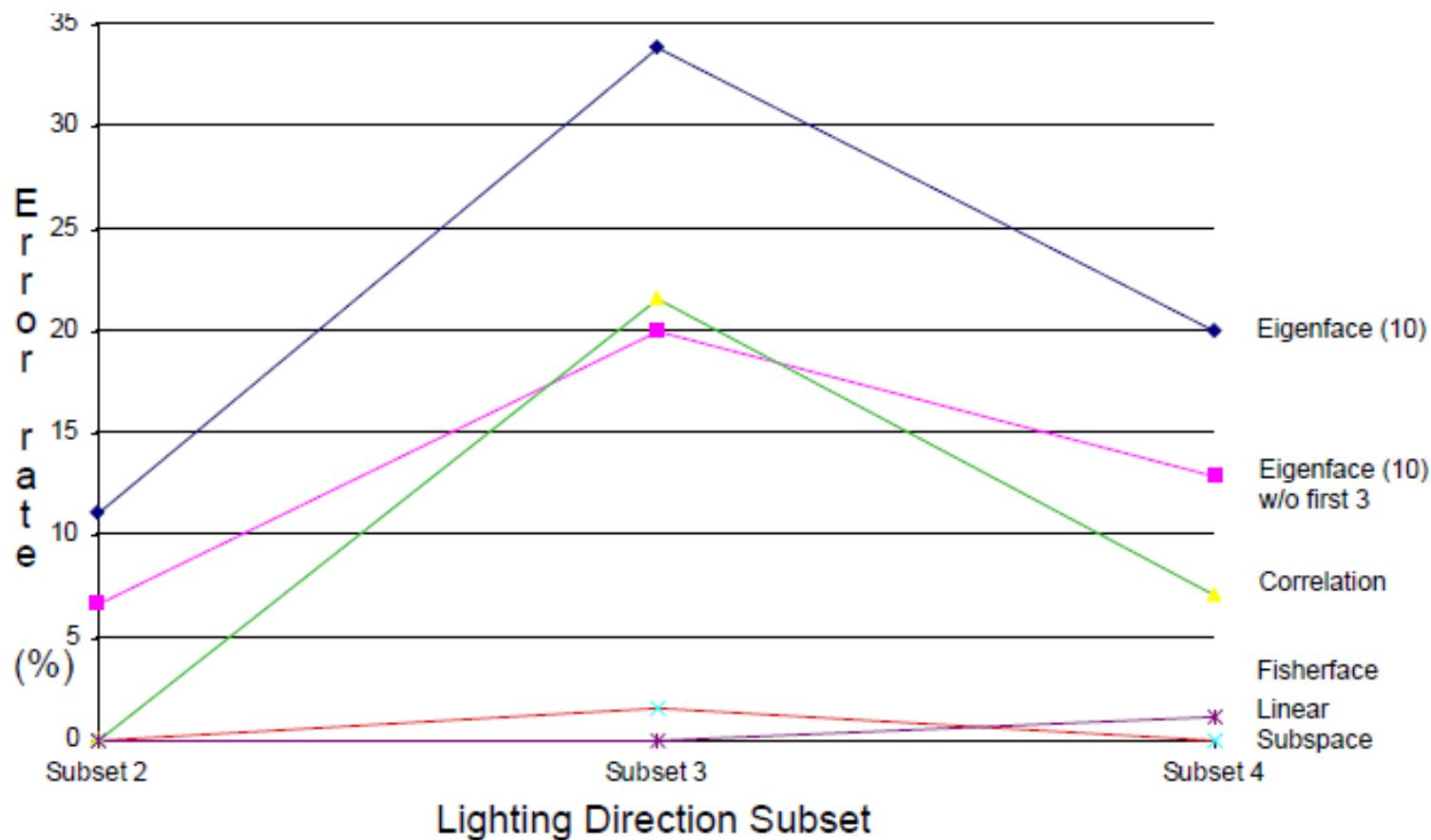


Fig. 4. Interpolation: When each of the methods is trained on images from both near frontal and extreme lighting (Subsets 1 and 5), the graph and corresponding table show the methods' relative performance under intermediate lighting conditions.

Fisherfaces produces better results!

<http://www.face-rec.org/algorithms/LDA/belhumeur96eigenfaces.pdf>

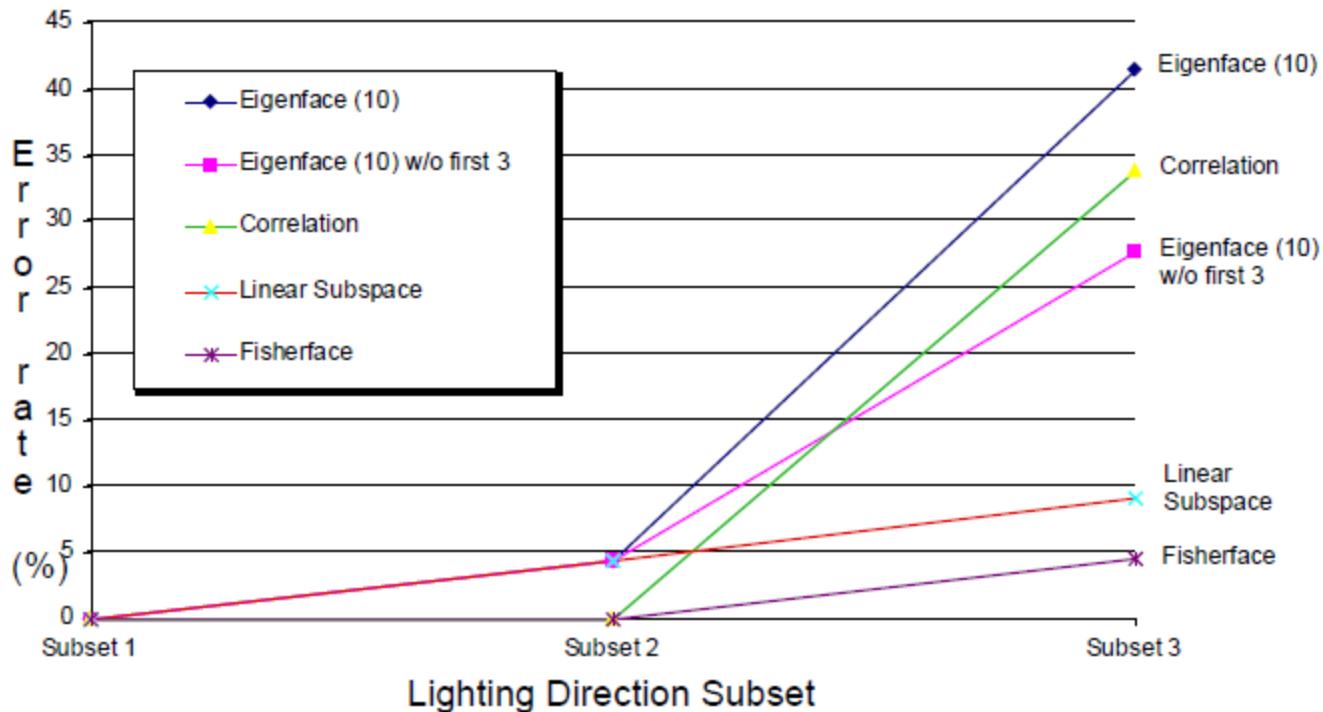


Fig. 3. Extrapolation: When each of the methods is trained on images with near frontal illumination (Subset 1), the graph and corresponding table show their relative performance under extreme light source conditions.

Fisherfaces produces better results!

<http://www.face-rec.org/algorithms/LDA/belhumeur96eigenfaces.pdf>

Classification using Mutual Information

- The FLD approach is constrained to reduce data dimensions to at most $c-1$ when there are c classes.
- It is desirable to remove this limitation and have more flexibility.
- Mutual information = quantity that characterizes the dependence between two random variables.

Mutual information: definition

- Given random variables X and Y , the mutual information between them $I(X,Y)$ is defined as:

$$I(X,Y) = \sum_{X=x, Y=y} p(x,y) \log\left(\frac{p(x,y)}{p(x)p(y)}\right)$$

- For continuous random variables, it is:

$$I(X,Y) = \iint p(x,y) \log\left(\frac{p(x,y)}{p(x)p(y)}\right) dx dy$$

Mutual information: definition

- It thus tells you how much information X and Y contain about each other – or how close they come to being independent of each other.
- MI has the property of being non-negative and symmetric.
- Zero-valued for independent random variables.

MI for classification

- Consider data $\{\mathbf{x}_i\}$ for $i = 1$ to N where each \mathbf{x}_i in \mathbb{R}^D .
- Each \mathbf{x}_i has label c_i which belongs to exactly one of N_c classes.
- Let X = R.V. for data, C = R.V. for class label.
- Infer transform vector \mathbf{w} such that the mutual information between $\{c_i\}$ and $\{\mathbf{y}_i = \mathbf{w}^t \mathbf{x}_i\}$ is maximized.

MI for classification

- Infer transform vector \mathbf{w} such that the mutual information between $\{c_i\}$ and $\{\mathbf{y}_i = \mathbf{w}^t \mathbf{x}_i\}$ is maximized.
- For our problem, we seek to maximize:

$$I(C, Y) = \sum_c \int p(c, y) \log \frac{p(c, y)}{P(c)p(y)} dy$$

- To compute MI, we need expressions for all the probabilities.

Probability expressions

- The expression for $P(c)$ is

$$P(c_p) = J_p / N, J_p = (\# \text{ points in class } p) / N$$

- Expression for $p(\mathbf{y})$:

$$p(\mathbf{y}) = \sum_{c=1}^{N_c} p(c, \mathbf{y}) = \sum_{c=1}^{N_c} p(\mathbf{y}/c) P(c)$$

$$p(\mathbf{y}|c_p) = \frac{1}{J_p} \sum_{j=1}^{J_p} G(\mathbf{y} - \mathbf{y}_{pj}, \sigma^2 I)$$

$$p(c_p, \mathbf{y}) = \frac{1}{N} \sum_{j=1}^{J_p} G(\mathbf{y} - \mathbf{y}_{pj}, \sigma^2 I), \quad p = 1, \dots, N_c$$

Segway: Kernel density estimation

- What motivates the specific expression for $p(\mathbf{y} | c_p)$?

$$p(\mathbf{y}|c_p) = \frac{1}{J_p} \sum_{j=1}^{J_p} G(\mathbf{y} - \mathbf{y}_{pj}, \sigma^2 I)$$

- This is a common way of **non-parametric** estimation of the probability density.
- Non-parametric: because it does not assume any families of densities (eg: Gaussian, Laplacian, Cauchy, etc.)

Segway: Kernel density estimation

- G in the expression below is called the kernel, which is often chosen to be Gaussian.

$$p(\mathbf{y}|c_p) = \frac{1}{J_p} \sum_{j=1}^{J_p} G(\mathbf{y} - \mathbf{y}_{pj}, \sigma^2 I)$$

- This estimator is often called the Parzen-Rosenblatt density estimator.
- Histogramming is another form of non-parametric density estimation – but produces discontinuous estimates.
- Note: Histogramming and kernel method are both **estimates** of a true underlying (unknown) probability density.

Segway: Kernel density estimation

- The choice of σ in this estimator is critical:

$$p(\mathbf{y}|c_p) = \frac{1}{J_p} \sum_{j=1}^{J_p} G(\mathbf{y} - \mathbf{y}_{pj}, \sigma^2 I)$$

- The optimal σ depends on the second derivatives of the underlying true density – which is unknown.
- A rule of thumb is to choose $\sigma = O(n^{-0.2})$ where n is the number of samples.

Probability expressions

- Expression for $p(\mathbf{y})$:

$$p(\mathbf{y}) = \sum_{c=1}^{N_c} p(c, \mathbf{y}) = \sum_{c=1}^{N_c} p(\mathbf{y}/c) P(c)$$

$$p(\mathbf{y}|c_p) = \frac{1}{J_p} \sum_{j=1}^{J_p} G(\mathbf{y} - \mathbf{y}_{pj}, \sigma^2 I)$$

$$p(c_p, \mathbf{y}) = \frac{1}{N} \sum_{j=1}^{J_p} G(\mathbf{y} - \mathbf{y}_{pj}, \sigma^2 I), \quad p = 1, \dots, N_c$$

$$p(\mathbf{y}) = \frac{1}{N} \sum_{p=1}^{N_c} \sum_{j=1}^{J_p} G(\mathbf{y} - \mathbf{y}_{pj}, \sigma^2 I) = \frac{1}{N} \sum_{i=1}^N G(\mathbf{y} - \mathbf{y}_i, \sigma^2 I)$$

Quadratic mutual information

- The mutual information expression involves log of summations – difficult to efficiently solve.
- Instead, we use the quadratic mutual information (QMI):

$$QMI(c, y) = \sum_c \int (p(c, y) - P(c)p(y))^2 dy$$

QMI expressions

$$I_T(C, Y) = \sum_c \int_{\mathbf{y}} p(c, \mathbf{y})^2 d\mathbf{y} + \sum_c \int_{\mathbf{y}} P(c)^2 p(\mathbf{y})^2 d\mathbf{y} - 2 \sum_c \int_{\mathbf{y}} p(c, \mathbf{y}) P(c) p(\mathbf{y}) d\mathbf{y}$$

$$V_{IN} \equiv \sum_c \int_{\mathbf{y}} p(c, \mathbf{y})^2 d\mathbf{y}$$

$$V_{ALL} \equiv \sum_c \int_{\mathbf{y}} P(c)^2 p(\mathbf{y})^2 d\mathbf{y}$$

$$V_{BTW} \equiv \sum_c \int_{\mathbf{y}} p(c, \mathbf{y}) P(c) p(\mathbf{y}) d\mathbf{y}$$

$$I_T(C, Y) = V_{IN} + V_{ALL} - 2V_{BTW} \quad \frac{\partial I_T}{\partial \mathbf{y}_i} = \frac{\partial V_{IN}}{\partial \mathbf{y}_i} + \frac{\partial V_{ALL}}{\partial \mathbf{y}_i} - 2 \frac{\partial V_{BTW}}{\partial \mathbf{y}_i}$$

QMI expressions

$$V_{IN} \equiv \sum_c \int_{\mathbf{y}} p(c, \mathbf{y})^2 d\mathbf{y} \quad I_T(C, Y) = V_{IN} + V_{ALL} - 2V_{BTW}$$

$$V_{ALL} \equiv \sum_c \int_{\mathbf{y}} P(c)^2 p(\mathbf{y})^2 d\mathbf{y}$$

$$V_{BTW} \equiv \sum_c \int_{\mathbf{y}} p(c, \mathbf{y}) P(c) p(\mathbf{y}) d\mathbf{y}$$

Within-class differences

$$V_{IN}(\{c_i, \mathbf{y}_i\}) = \sum_c \int_{\mathbf{y}} p(c, \mathbf{y})^2 d\mathbf{y} = \frac{1}{N^2} \sum_{p=1}^{N_c} \sum_{k=1}^{J_p} \sum_{l=1}^{J_p} G(\mathbf{y}_{pk} - \mathbf{y}_{pl}, 2\sigma^2 I)$$

$$V_{ALL}(\{c_i, \mathbf{y}_i\}) = \sum_c \int_{\mathbf{y}} P(c)^2 p(\mathbf{y})^2 d\mathbf{y} = \frac{1}{N^2} \left(\sum_{p=1}^{N_c} \left(\frac{J_p}{N} \right)^2 \right) \sum_{k=1}^N \sum_{l=1}^N G(\mathbf{y}_k - \mathbf{y}_l, 2\sigma^2 I)$$

$$V_{BTW}(\{c_i, \mathbf{y}_i\}) = \sum_c \int_{\mathbf{y}} p(c, \mathbf{y}) P(c) p(\mathbf{y}) d\mathbf{y} = \frac{1}{N^2} \sum_{p=1}^{N_c} \frac{J_p}{N} \sum_{j=1}^{J_p} \sum_{k=1}^N G(\mathbf{y}_{pj} - \mathbf{y}_k, 2\sigma^2 I) .$$

Differences of all pairs of samples
(regardless of class)

QMI expressions for derivatives

$$\frac{\partial I_T}{\partial \mathbf{y}_i} = \frac{\partial V_{IN}}{\partial \mathbf{y}_i} + \frac{\partial V_{ALL}}{\partial \mathbf{y}_i} - 2 \frac{\partial V_{BTW}}{\partial \mathbf{y}_i}$$

$$\frac{\partial}{\partial \mathbf{y}_i} G(\mathbf{y}_i - \mathbf{y}_j, 2\sigma^2 I) = G(\mathbf{y}_i - \mathbf{y}_j, 2\sigma^2 I) \frac{(\mathbf{y}_j - \mathbf{y}_i)}{2\sigma^2}$$

$$\frac{\partial}{\partial \mathbf{y}_{ci}} V_{IN} = \frac{1}{N^2 \sigma^2} \sum_{k=1}^{J_c} G(\mathbf{y}_{ck} - \mathbf{y}_{ci}, 2\sigma^2 I) (\mathbf{y}_{ck} - \mathbf{y}_{ci})$$

$$\frac{\partial}{\partial \mathbf{y}_{ci}} V_{ALL} = \frac{1}{N^2 \sigma^2} \left(\sum_{p=1}^{N_c} \left(\frac{J_p}{N} \right)^2 \right) \sum_{k=1}^N G(\mathbf{y}_k - \mathbf{y}_i, 2\sigma^2 I) (\mathbf{y}_k - \mathbf{y}_i)$$

$$\frac{\partial}{\partial \mathbf{y}_{ci}} V_{BTW} = \frac{1}{N^2 \sigma^2} \sum_{p=1}^{N_c} \frac{J_p + J_c}{2N} \sum_{j=1}^{J_p} G(\mathbf{y}_{pj} - \mathbf{y}_{ci}, 2\sigma^2 I) (\mathbf{y}_{pj} - \mathbf{y}_{ci})$$

QMI optimization

- Proceeds by projected gradient ascent with adaptive step-size and can reach a local maximum of the QMI:

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \eta \frac{\partial I}{\partial \mathbf{w}} = \mathbf{w}_t + \eta \sum_{i=1}^N \frac{\partial I}{\partial \mathbf{y}_i} \frac{\partial \mathbf{y}_i}{\partial \mathbf{w}}$$

- This was for optimizing over a single direction.
- For multiple directions, we need to optimize over a matrix \mathbf{W} (size $D \times d$, $d < D$) such that $\mathbf{y}_i = \mathbf{W}^T \mathbf{x}_i$.
- We need to impose orthonormality of \mathbf{W} – hence the need for a projection step.

QMI optimization: orthogonalization

- Let $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_d$ be the row vectors of \mathbf{W} – which is not at this point orthogonal.
- Define the following operator:

$$proj_u(v) = \frac{v \bullet u}{u \bullet u} u$$

- The next slide gives the orthogonalization procedure – called Gram-Schmidt O.P.

Gram Schmidt Orthogonalization

$\mathbf{u}_1 = \mathbf{w}_1, \mathbf{e}_1 = \mathbf{u}_1 / \|\mathbf{u}_1\|_2$ Procedure

This forces \mathbf{u}_2 to be orthogonal to \mathbf{u}_1 .

$\mathbf{u}_2 = \mathbf{w}_2 - proj_{\mathbf{u}_1}(\mathbf{w}_2), \mathbf{e}_2 = \mathbf{u}_2 / \|\mathbf{u}_2\|_2,$

...

$\mathbf{u}_d = \mathbf{w}_d - \sum_{k=1}^{d-1} proj_{\mathbf{u}_k}(\mathbf{w}_d), \mathbf{e}_d = \mathbf{u}_d / \|\mathbf{u}_d\|_2$

NOTE : $\mathbf{u}_2 \bullet \mathbf{u}_1 = (\mathbf{w}_2 - proj_{\mathbf{u}_1}(\mathbf{w}_2)) \bullet \mathbf{u}_1$

$$= (\mathbf{w}_2 - \frac{\mathbf{w}_2 \bullet \mathbf{u}_1}{\|\mathbf{u}_1\|^2} \mathbf{u}_1) \bullet \mathbf{u}_1 = \mathbf{w}_2 \bullet \mathbf{u}_1 - \mathbf{w}_2 \bullet \mathbf{u}_1 = 0$$

https://en.wikipedia.org/wiki/Gram%E2%80%93Schmidt_process

Results

- Three classes:
 - C1: Bimodal Gaussians with centers at $(1,0,0)$ and $(-1,0,0)$
 - C2: Bimodal Gaussians with centers at $(0,1,0)$ and $(0,-1,0)$
 - C3: Single Gaussian at $(0,0.7,1)$

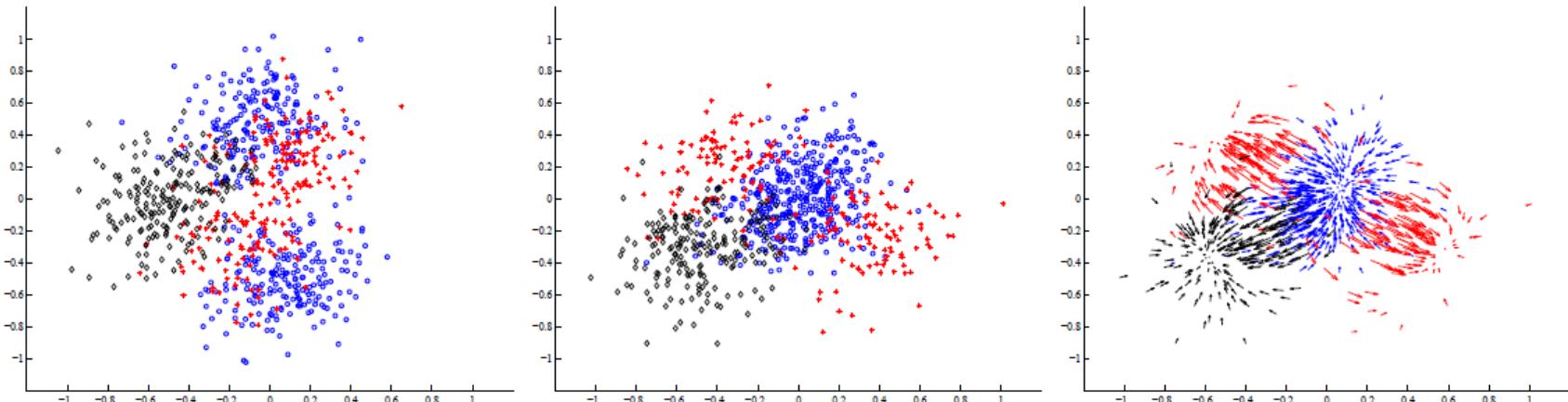


Figure 2: Three classes in three dimensions projected onto a two-dimensional subspace using both LDA (left) and MMI-based projection (center). Final information forces, with a narrow kernel, are shown on the right.

<http://www.jmlr.org/papers/volume3/torkkola03a/torkkola03a.pdf>

Compressed Sensing with Overcomplete Dictionaries

Towards CS with overcomplete dictionaries

- A **basis** in a vector space contains vectors that are all linearly independent.
- A **frame** generalizes this concept to vectors that may be linearly dependent.
- A set of vectors \mathbf{e}_k , $1 \leq k \leq n$, form a **frame** if there exist positive values A and B , $A \leq B$, such that for every vector \mathbf{v} , we have:

$$A\|\mathbf{v}\|^2 \leq \sum_k |\mathbf{v}^t \mathbf{e}_k|^2 \leq B\|\mathbf{v}\|^2$$

Towards CS with overcomplete dictionaries

- A **frame** is overcomplete or redundant if it not a basis (of the vector space).
- A frame is called a **tight frame** if $A = B$. For a tight frame, we have for any \mathbf{v} :

$$\mathbf{v} = \frac{1}{A} \sum_k (\mathbf{v}^t \mathbf{e}_k) \mathbf{e}_k$$

Towards CS with overcomplete dictionaries

- We mention a generalization of the restricted isometry property called the D-RIP.

Definition 1.3 (D-RIP) Let Σ_s be the union of all subspaces spanned by all subsets of s columns of D . We say that the measurement matrix A obeys the restricted isometry property adapted to D (abbreviated D-RIP) with constant δ_s if

$$(1 - \delta_s)\|v\|_2^2 \leq \|Av\|_2^2 \leq (1 + \delta_s)\|v\|_2^2$$

holds for all $v \in \Sigma_s$.

We point out that Σ_s is just the image under D of all s -sparse vectors.

Note that Σ_s is the image under D of all s -sparse vectors. The only difference is that D here has more columns than rows unlike the case of orthonormal matrices before. Also random Gaussian and Bernoulli matrices obey D-RIP with high probability given $O(s \log n)$ rows, just like they obey RIP.

CS with tight frames

- Let Ψ be a tight frame of size $n \times d$, and Φ be a $m \times n$ measurement matrix satisfying D-RIP with $\delta_{2s} < 0.08$ or $\delta_{7s} < 0.6$. Then the solution \mathbf{x}^* to the optimization problem P1* below satisfies:

$$\|\mathbf{x}^* - \mathbf{x}\|_2 \leq C_0 \varepsilon + C_1 \frac{\|\Psi^* \mathbf{x} - (\Psi^* \mathbf{x})_s\|_1}{\sqrt{s}}$$

P1*: $\min \|\Psi^* \mathbf{x}\|_1$ such that $\|\mathbf{y} - \mathbf{Ax}\|_2 \leq \varepsilon$;

$\mathbf{y} = \mathbf{Ax} + \boldsymbol{\eta}$; \mathbf{x} is the signal, \mathbf{y} is the measurement vector

$\Psi^* \mathbf{x}$ is the signal - coefficient vector

References

- Section 3.8 of “Pattern Classification” by Duda and Hart
- <http://en.wikipedia.org/wiki/Eigenface>
- M. Turk and A. Pentland (1991). ["Eigenfaces for recognition"](#). *Journal of Cognitive Neuroscience*

Appendix: Covariance matrix

- In probability and statistics, the expected value of a scalar random variable x is called its mean:

$$\mu = E(x) = \int xp(x)dx \approx \frac{1}{N} \sum_{i=1}^N x_i$$

Probability density function of x Sample values of the random variable

- The variance of a scalar random variable is the expected value of its squared deviation around the mean:

$$\sigma^2 = E((x-\mu)^2) = \int (x-\mu)^2 p(x)dx \approx \frac{1}{N-1} \sum_{i=1}^N (x_i - \mu)^2$$

Appendix: Covariance matrix

- The expected value of a vector random variable is called its mean vector:

$$\vec{\mu} = E(\vec{x}), \vec{x} = (x(1), x(2), \dots, x(d)), \vec{\mu} = (\mu(1), \mu(2), \dots, \mu(d))$$

$$\mu_k \approx \frac{1}{N} \sum_{i=1}^N x_{ki}, 1 \leq k \leq d$$

→ k-th sample value of \mathbf{x}_i

Appendix: Covariance matrix

- The variance is now replaced by a covariance matrix. Each entry contains the covariance between two elements of the vector:

$$\mathbf{C} = \begin{pmatrix} E((x(1) - \mu(1))^2) & E((x(1) - \mu(1))(x(2) - \mu(2))) & \dots & E((x(d) - \mu(d))(x(1) - \mu(1))) \\ E((x(2) - \mu(2))(x(1) - \mu(1))) & E((x(2) - \mu(2))^2) & \dots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ E((x(1) - \mu(1))(x(d) - \mu(d))) & \vdots & \ddots & E((x(d) - \mu(d))^2) \end{pmatrix} \in \mathbb{R}^{d \times d}$$

$$C_{kl} = E((x(k) - \mu(k))(x(l) - \mu(l))) \approx \frac{1}{N-1} \sum_{i=1}^N (x_i(k) - \mu(k))(x_i(l) - \mu(l)), 1 \leq k \leq d, 1 \leq l \leq d$$

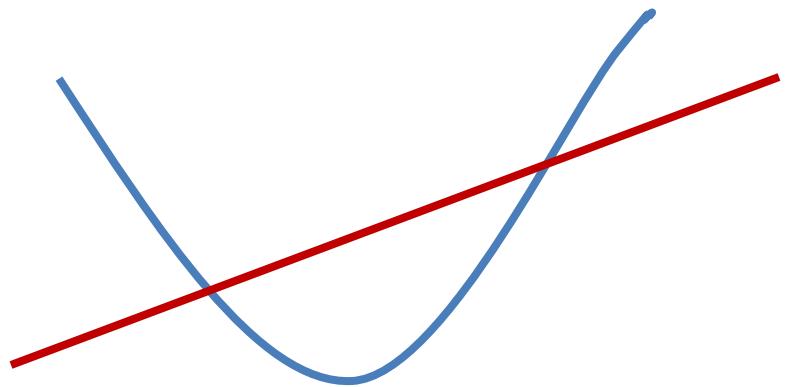
Note: $\mu(k)$ = k -th element of vector μ

Appendix: A word about Lagrange Multipliers

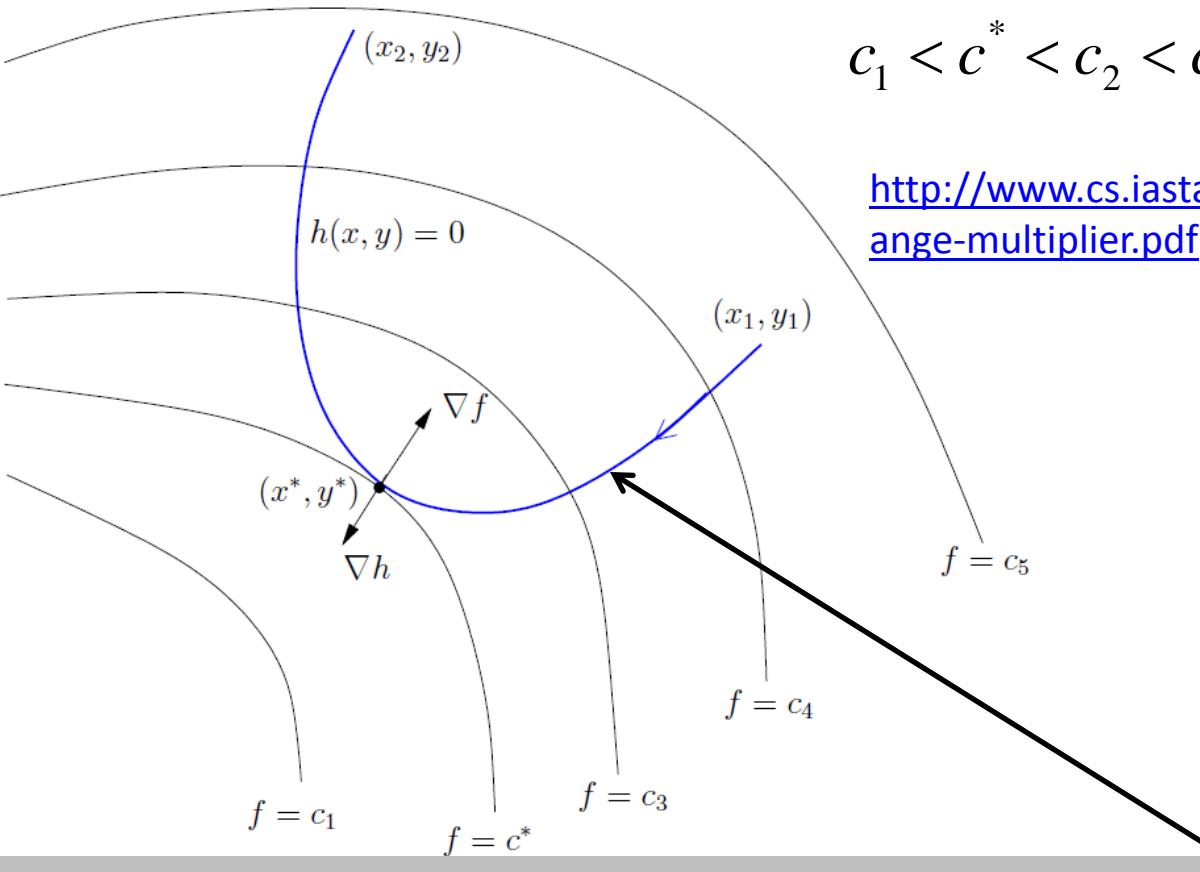
- Consider you want to find the minimum or maximum of $f(x,y)$.
- Normally, you take the derivative and set it to 0 and check the sign of the second derivative.
- Now you want to find the optimum subject to a constraint that $h(x,y) = 0$.

Appendix: A word about Lagrange Multipliers

- Physical analogy: if you drop a pebble into a parabola-shaped bowl, the pebble will fall to the bottom-most point.
- But if you placed a wooden plank into the bowl, the pebble will settle somewhere on the plank!



$$c_1 < c^* < c_2 < c_3 < c_4 < c_5$$



<http://www.cs.iastate.edu/~cs577/handouts/lagrange-multiplier.pdf>

Suppose we are walking from (x_1, y_1) to (x_2, y_2) along the constraint curve $h(x, y) = 0$. Initially, the tangent vector along $h(x, y) = 0$ has a component along $-\text{grad}(f)$, and hence there is a decrease in the value of $f(x, y)$ as we move along $h(x, y) = 0$. If the point (x^*, y^*) is a local minimum of $f(x, y)$, a small motion along $h(x, y) = 0$ from (x^*, y^*) will have no component along $\text{grad}(f)$, else it would cause an increase in the value of f . Hence the tangent to $h(x, y) = 0$ at (x^*, y^*) will be perpendicular to $\text{grad}(f)$. As we move further, motion along $h(x, y) = 0$ will have a component along $+\text{grad}(f)$ leading to an increase in the value of $f(x, y)$. At the minimum point, we have $\text{grad}(f)$ perpendicular to the tangent, i.e. $\text{grad}(f)$ and $\text{grad}(h)$ are collinear. Hence, there exists some value λ (the **Lagrange multiplier**) such that we have:

$$\nabla f(x^*, y^*) = \lambda \nabla h(x^*, y^*) \rightarrow \nabla(f(x^*, y^*) - \lambda h(x^*, y^*)) = 0$$

If (x', y') is the extremum of function $f(x, y)$, then a small perturbation to (x', y') will lead to no change in the value of f - in the limit when the magnitude of the perturbation is 0.

Hence we have $|\nabla f(x', y') = 0|$.

But here we are constrained to move only along the curve $h(x, y) = 0$.

Let \mathbf{s} be the tangent to this curve at a point (x', y') .

The magnitude of the change in f due to infinitesimal movement along the curve is given by $\nabla f(x', y') \bullet \mathbf{s}$.

If (x', y') is an extremum of f along this curve, then we must have

$$\nabla f(x', y') \bullet \mathbf{s} = 0.$$

But the normal to the curve is $\nabla h(x', y')$ and it is perpendicular to \mathbf{s} .

Hence we have $\nabla f(x', y')$ is collinear with $\nabla h(x', y')$.

$$\therefore \nabla f(x', y') = \lambda \nabla h(x', y') \text{ for some value } \lambda.$$

In our derivation for PCA, instead of $f(x, y)$ we have $\mathbf{e}^t \mathbf{S} \mathbf{e}$, and instead of $h(x, y) = 0$, we have $\mathbf{e}^t \mathbf{e} - 1 = 0$. Note that \mathbf{e} is a vector in d dimensional space.