

# Dimensionality Reduction by PCA

# Pattern Matrix

- ▶ Statistics and machine learning typically starts from data given in the form of **observations**, **feature vectors** or **patterns**
- ▶ Feature vectors (in some  $m$ -dimensional Euclidean space)

$$\mathbf{x}_i \in \mathcal{X} \subseteq \mathbb{R}^m, \quad i = 1, \dots, n$$

- ▶ Patterns can be summarized into the **pattern matrix**

$$\mathbf{X} \in \mathbb{R}^{n \times m}, \quad \mathbf{X} = \begin{pmatrix} \mathbf{x}'_1 \\ \dots \\ \mathbf{x}'_i \\ \dots \\ \mathbf{x}'_n \end{pmatrix} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1m} \\ \dots & \dots & \dots & \dots \\ x_{i1} & x_{i2} & \dots & x_{im} \\ \dots & \dots & \dots & \dots \\ x_{n1} & x_{n2} & \dots & x_{nm} \end{pmatrix}$$

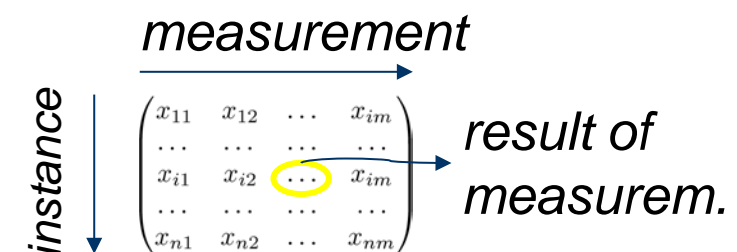
*transposed i-th pattern*

# Examples: Pattern Matrices

$$\mathbf{X} \in \mathbb{R}^{n \times m}$$

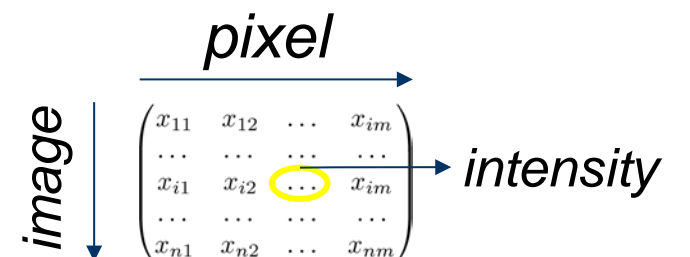
## ► Measurement vectors

- $i$ : instance number, e.g. a house
- $j$ : measurement, e.g. the area of a house



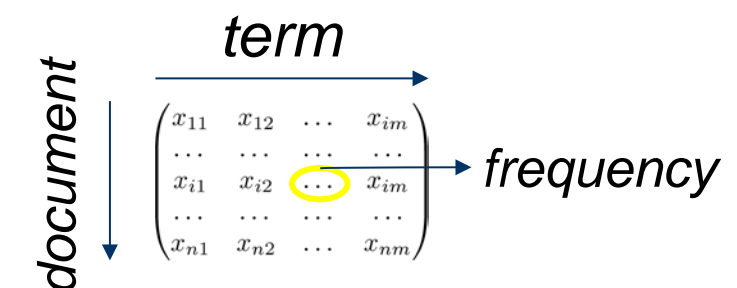
## ► Digital images as gray-scale vectors

- $i$ : image number
- $j$ : pixel value at location  $j=(k,l)$



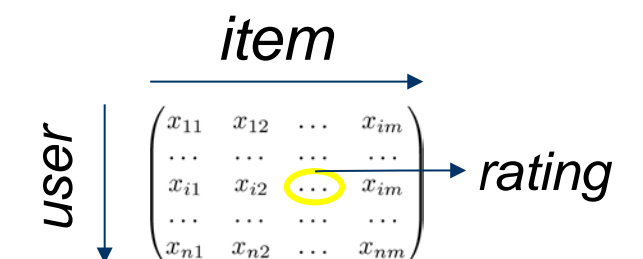
## ► Text documents in bag-of-words representation

- $i$ : document number
- $j$ : term (word or phrase) in a vocabulary

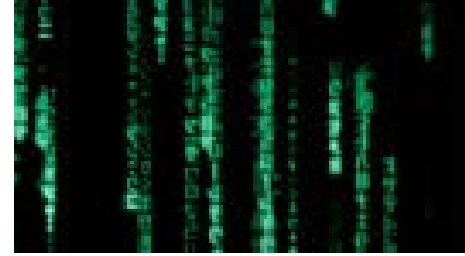


## ► User rating data

- $i$ : user number
- $j$ : item (book, movie)



# Latent Structure

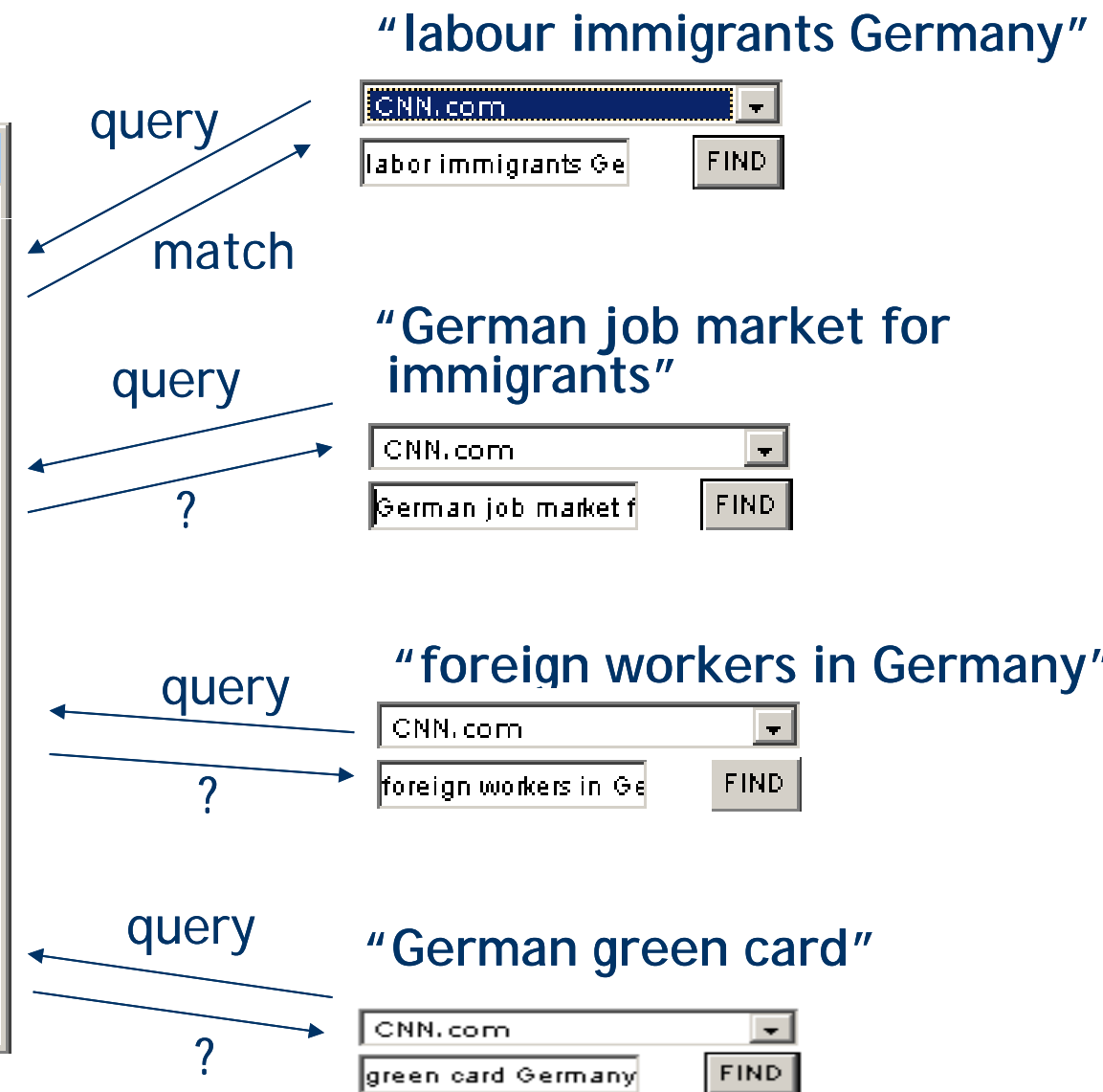


- ▶ Given a matrix that “encodes” data ...
- ▶ Potential problems

- too large
- too complicated
- missing entries
- noisy entries
- lack of structure
- ...

$$\mathbf{A} = \begin{pmatrix} a_{11} & \dots & a_{1j} & \dots & a_{1m} \\ \dots & \dots & \dots & \dots & \dots \\ a_{i1} & \dots & a_{ij} & \dots & a_{im} \\ \dots & \dots & \dots & \dots & \dots \\ a_{n1} & \dots & a_{nj} & \dots & a_{nm} \end{pmatrix}$$

# Vocabulary Mismatch & ~~Robustness~~





# Document-Term Matrix

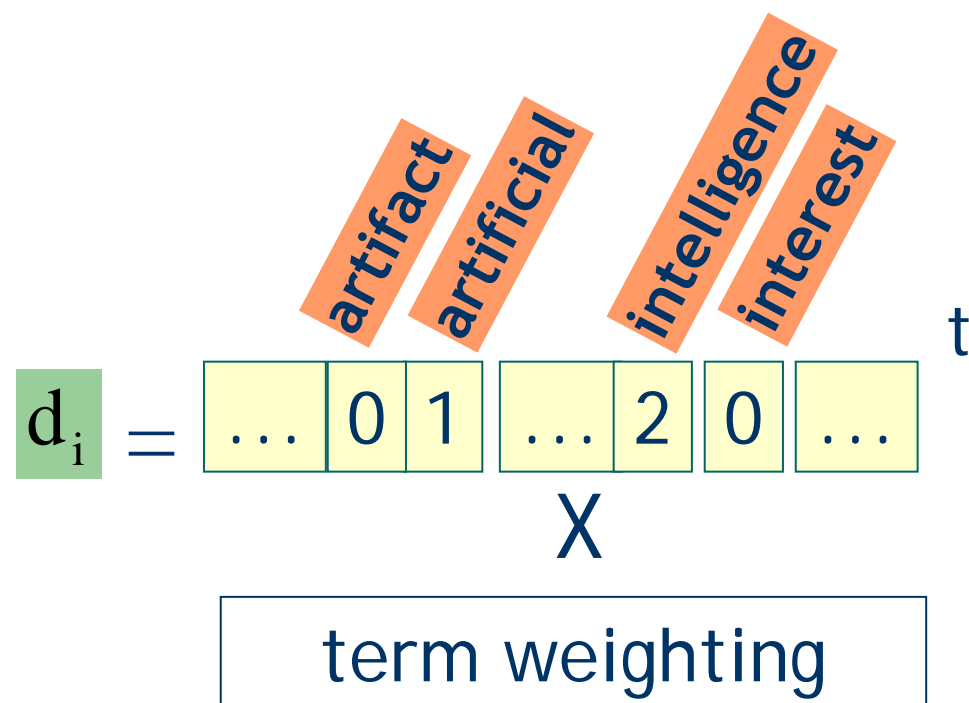
$D$  = Document collection

$W$  = Lexicon/Vocabulary

intelligence

$w_j$

*Texas Instruments said it has developed the first 32-bit computer chip designed specifically for artificial intelligence applications [...]*



Document-Term Matrix

	$W$					
	$w_1$	...	$w_j$	...	$w_J$	
$D$	$d_1$					
	...					
	$d_i$			$c(d_i, w_j)$		
	...					
	$d_I$					

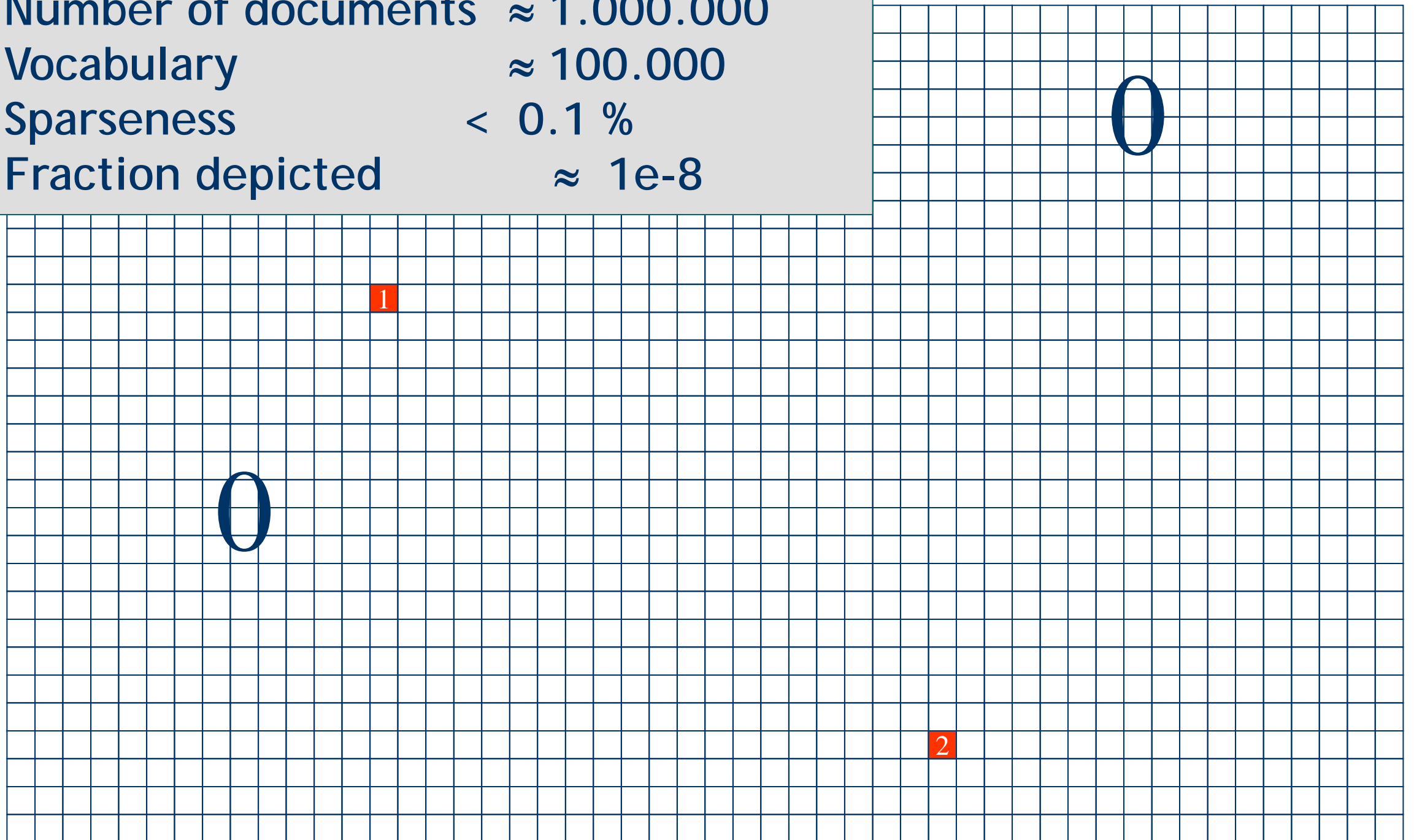
Clustering rows?

Clustering columns?

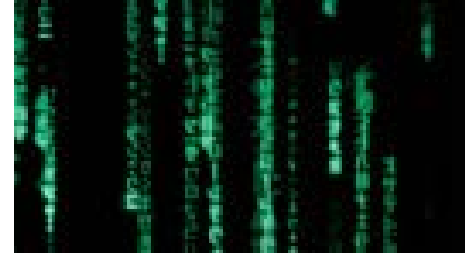
# A 100 Million<sup>ths</sup> of a Typical Document-term Matrix

Typical:

- Number of documents  $\approx 1.000.000$
- Vocabulary  $\approx 100.000$
- Sparseness  $< 0.1 \%$
- Fraction depicted  $\approx 1e-8$



# Latent Structure



- ▶ Given a matrix that “encodes” data ...

- ▶ Potential problems

- too large
- too complicated
- missing entries
- noisy entries
- lack of structure
- ...

$$\mathbf{A} = \begin{pmatrix} a_{11} & \dots & a_{1j} & \dots & a_{1m} \\ \dots & \dots & \dots & \dots & \dots \\ a_{i1} & \dots & a_{ij} & \dots & a_{im} \\ \dots & \dots & \dots & \dots & \dots \\ a_{n1} & \dots & a_{nj} & \dots & a_{nm} \end{pmatrix}$$

- ▶ Is there a **simpler** way to **explain** entries?
- ▶ There might be a **latent structure** underlying the data.
- ▶ How can we “find” or “reveal” this structure?



# Matrix Decomposition

- Common approach: approximately **factorize** matrix

$$\mathbf{A} \approx \hat{\mathbf{A}} = \mathbf{L} \cdot \mathbf{R}$$

approximation      left factor      right factor

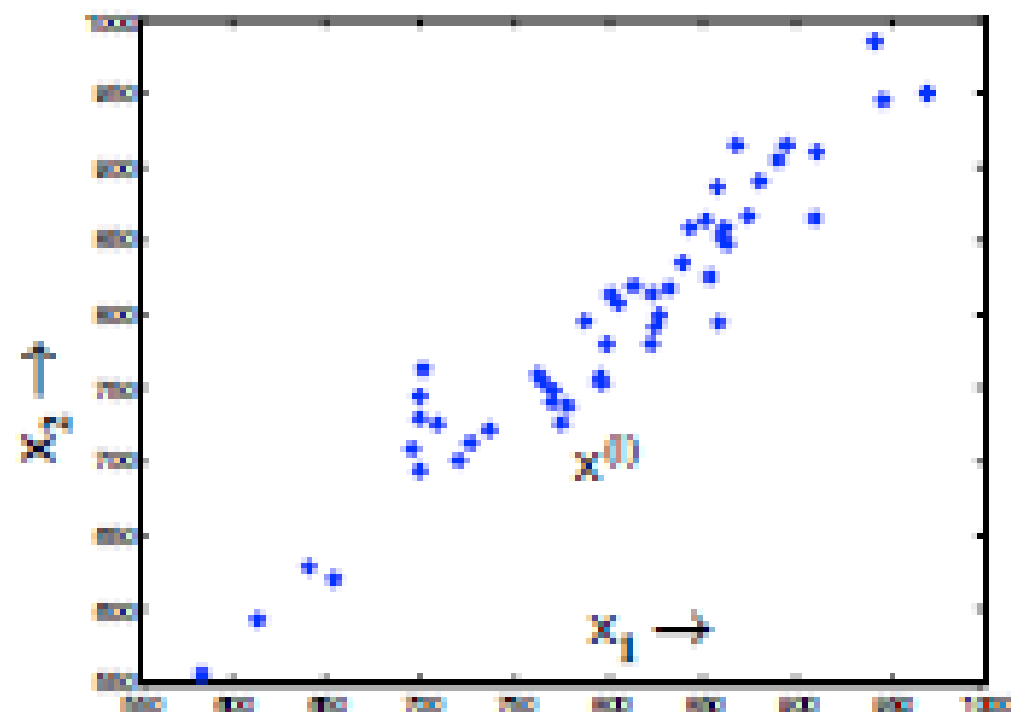
- Factors are typically constrained to be “thin”

$$\begin{array}{c} \overbrace{\hspace{2cm}}^m \\ \begin{array}{|c|} \hline \mathbf{A} \\ \hline \end{array} \\ \underbrace{\hspace{2cm}}_n \end{array} \approx \begin{array}{c} \underbrace{\hspace{2cm}}_n \\ \begin{array}{|c|} \hline \mathbf{L} \\ \hline \end{array} \end{array} \cdot \begin{array}{c} \overbrace{\hspace{2cm}}^m \\ \begin{array}{|c|} \hline \mathbf{R} \\ \hline \end{array} \\ \underbrace{\hspace{2cm}}_q \end{array}$$

reduction  
 $n \cdot m \gg n \cdot q + m \cdot q$   
 factors = latent structure (?)

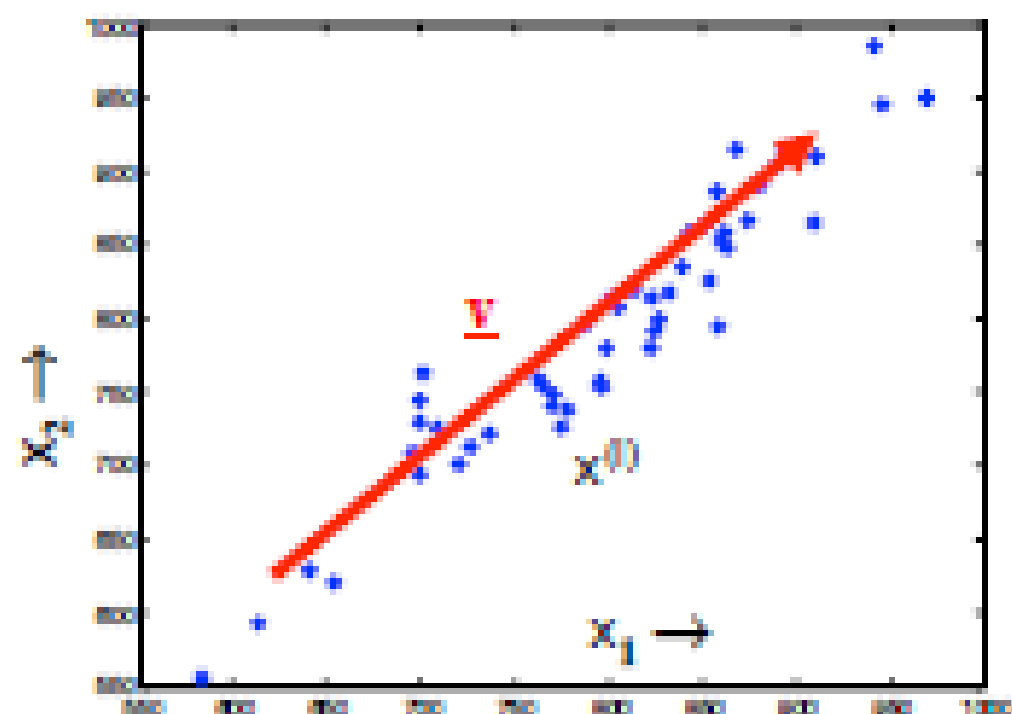
# Dimensionality reduction

- Ex: data with two real values  $[x_1, x_2]$
- We'd like to describe each point using only one value  $[z_1]$
- We'll communicate a "model" to convert:  $[x_1, x_2] \sim f(z_1)$



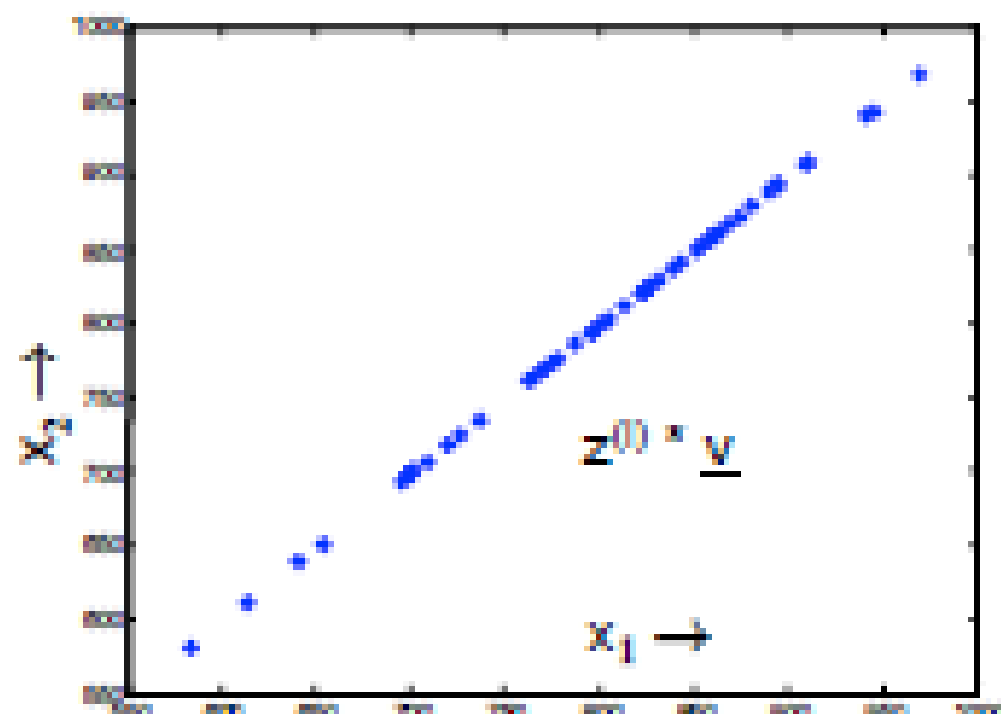
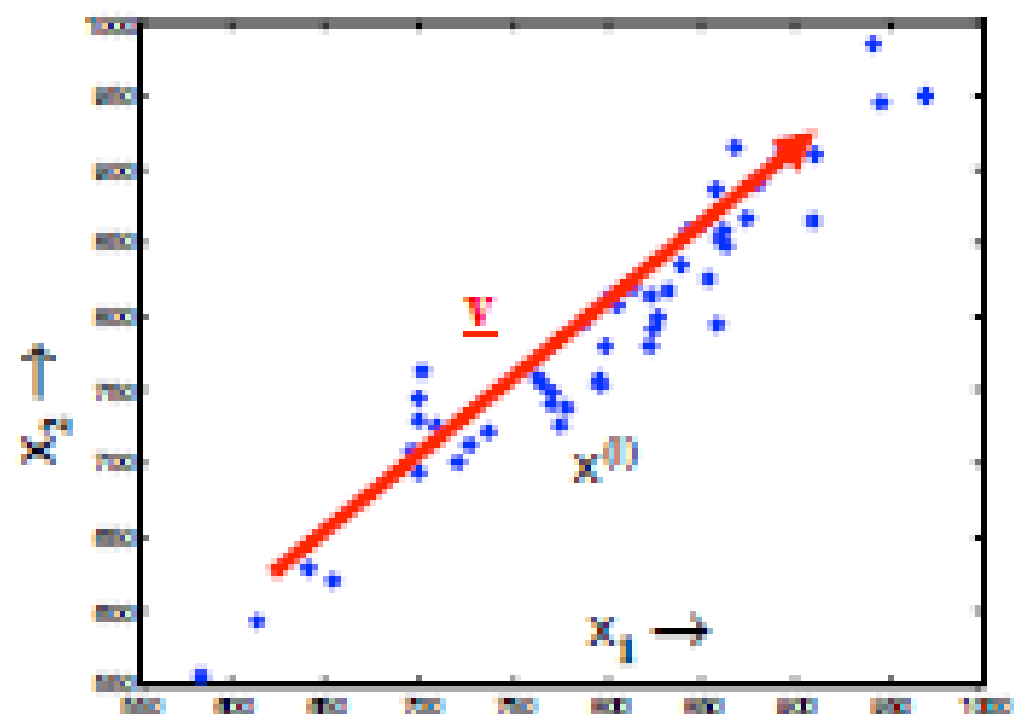
# Dimensionality reduction

- Ex: data with two real values  $[x_1, x_2]$
- We'd like to describe each point using only one value  $[z_1]$
- We'll communicate a "model" to convert:  $[x_1, x_2] \sim f(z_1)$
- Ex: linear function  $f(z)$ :  $[x_1, x_2] = z * \underline{v} = z * [v_1, v_2]$
- $\underline{v}$  is the same for all data points (communicate once)
- $z$  tells us the closest point on  $v$  to the original point  $[x_1, x_2]$



# Dimensionality reduction

- Ex: data with two real values  $[x_1, x_2]$
- We'd like to describe each point using only one value  $[z_1]$
- We'll communicate a "model" to convert:  $[x_1, x_2] \sim f(z_1)$
- Ex: linear function  $f(z)$ :  $[x_1, x_2] = z * \underline{v} = z * [v_1, v_2]$
- $\underline{v}$  is the same for all data points (communicate once)
- $z$  tells us the closest point on  $v$  to the original point  $[x_1, x_2]$

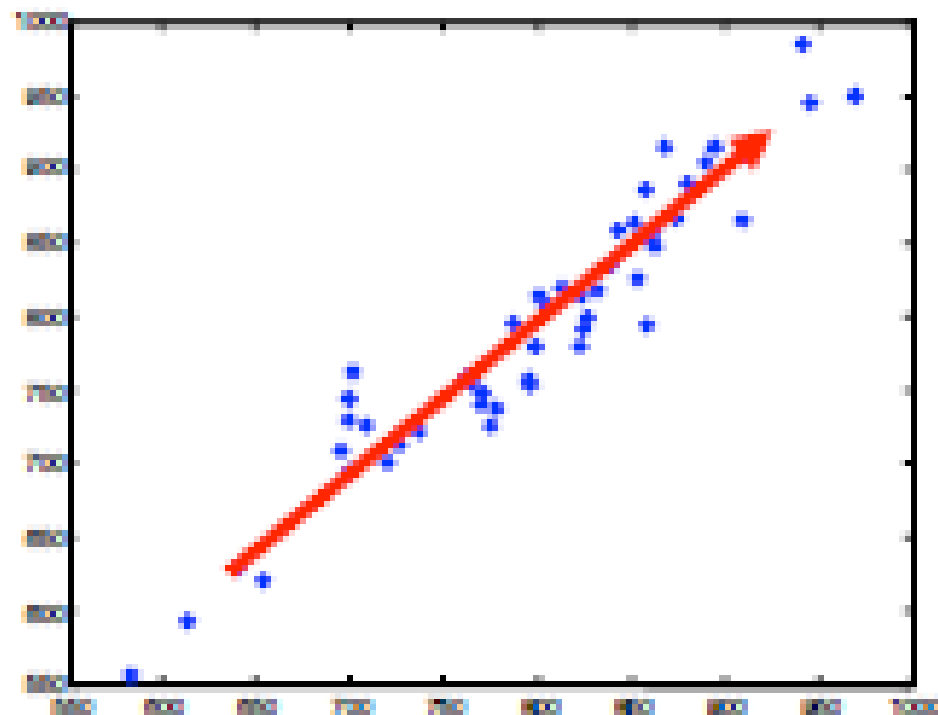


# Principal Components Analysis

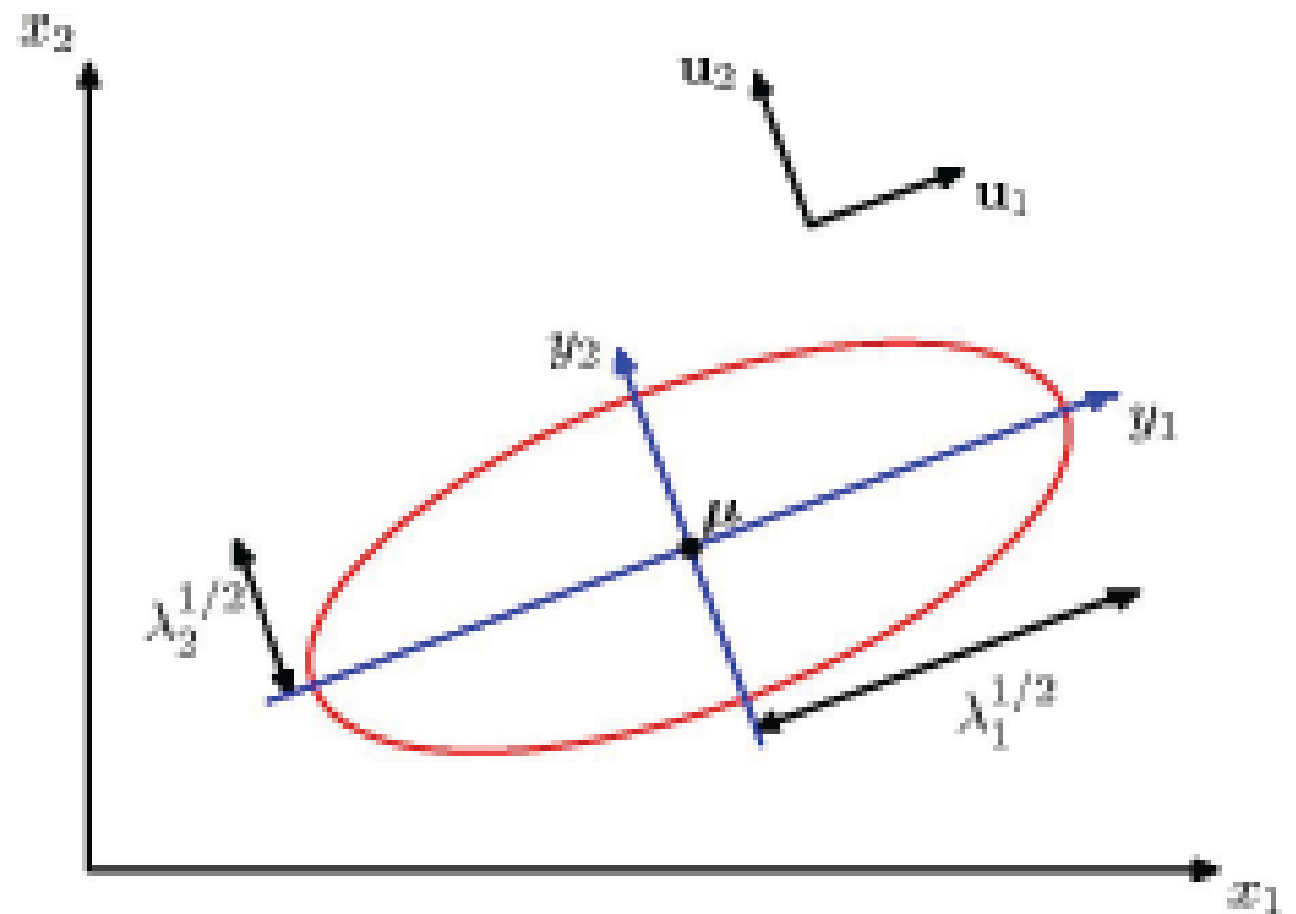
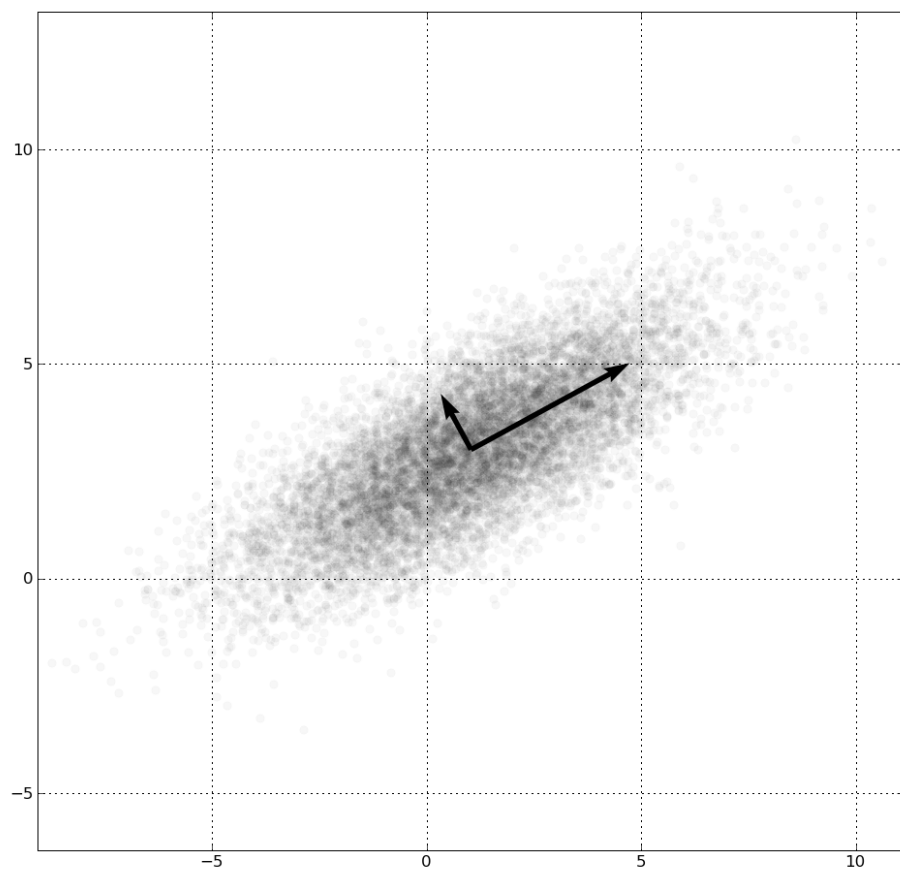
- What is the vector that would most closely reconstruct  $X$ ?

$$\min_{a,v} \sum_i (x^{(i)} - a^{(i)}v)^2$$

- Given  $v$ :  $a^{(i)}$  is the projection of each point  $x^{(i)}$  onto  $v$
- $v$  chosen to minimize the residual variance
- Equivalently,  $v$  is the direction of maximum variance
- Extensions: best two dimensions:  $x_i = a_i v + b_i w + m$



# Eigenvectors



# PCA

Given pattern matrix  $X$ ,

1. Subtract mean from each point
2. (sometimes) scale each dimension by its variance
3. Compute covariance matrix  $S = X^T X$
4. Compute  $k$  largest eigenvectors of  $S$

$$S = V D V^T$$

# Singular Value Decomposition

- Alternative method to calculate (still subtract mean 1<sup>st</sup>)
- Decompose  $X = U S V^T$ 
  - Orthogonal:  $X^T X = V S S V^T = V D V^T$
  - $X X^T = U S S U^T = U D U^T$
- $U \cdot S$  matrix provides coefficients
  - Example  $x_i = U_{i,1} S_{11} V_1 + U_{i,2} S_{22} V_2 + \dots$
- Gives the least-squares approximation to  $X$  of this form

$$\boxed{\begin{matrix} X \\ N \times D \end{matrix}} \approx \boxed{\begin{matrix} U \\ N \times K \end{matrix}} \boxed{\begin{matrix} S \\ K \times K \end{matrix}} \boxed{\begin{matrix} V^T \\ K \times D \end{matrix}}$$



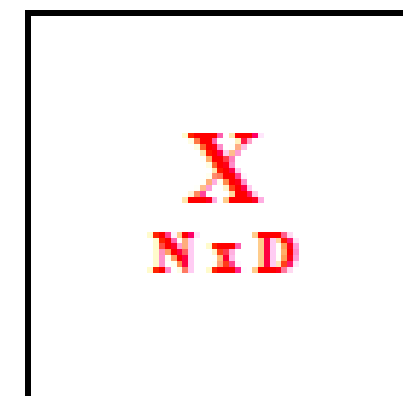
# Glorious SVD

$$X = USV^T$$

- $XX^T$  and  $X^T X$  share the same eigenvalues
- Even better: their eigenvectors are related
  - $Xv_i$  is an eigenvector of  $XX^T$

# “Eigen-faces”

- “Eigen-X” = represent X using PCA
- Ex: Viola Jones data set
  - 24x24 images of faces = 576 dimensional measurements



# “Eigen-faces”

- “Eigen-X” = represent X using PCA
- Ex: Viola Jones data set
  - 24x24 images of faces = 576 dimensional measurements
  - Take first K PCA components

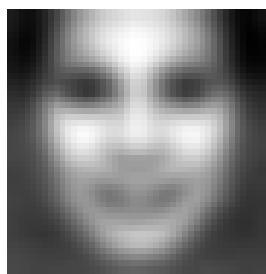


+

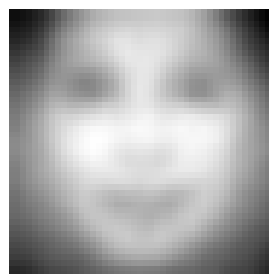


⋮

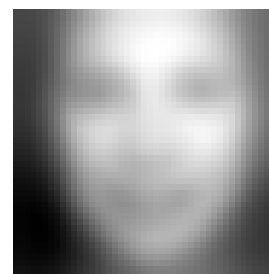
$$\begin{matrix} X \\ N \times D \end{matrix} \approx \begin{matrix} U \\ N \times K \end{matrix} \begin{matrix} S \\ K \times K \end{matrix} \begin{matrix} V^T \\ K \times D \end{matrix}$$



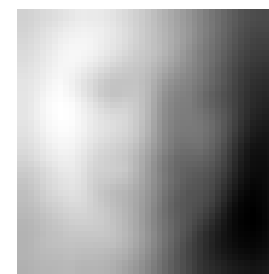
Mean



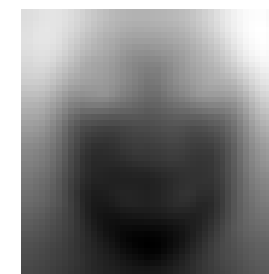
V(1,:)



V(2,:)



V(3,:)

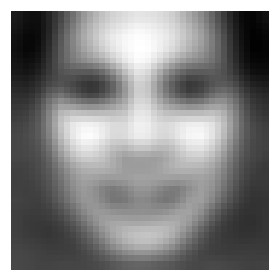


V(4,:)

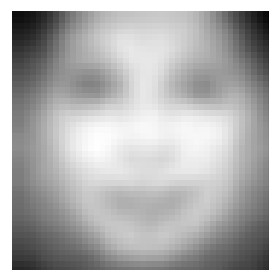
...

# “Eigen-faces”

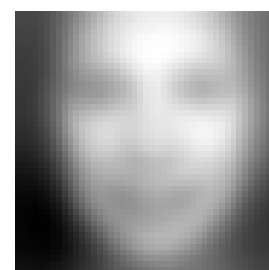
- “Eigen-X” = represent X using PCA
- Ex: Viola Jones data set
  - 24x24 images of faces = 576 dimensional measurements
  - Take first K PCA components



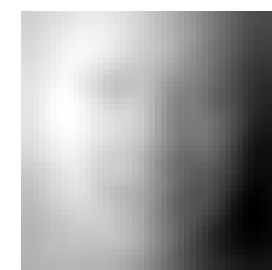
Mean



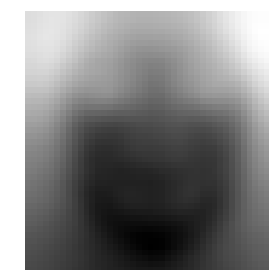
Dir 1



Dir 2

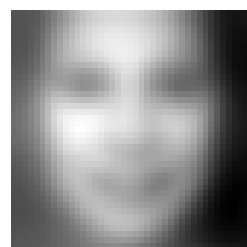
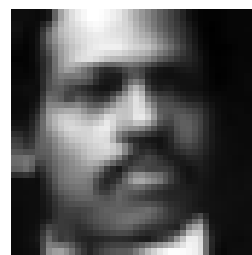


Dir 3

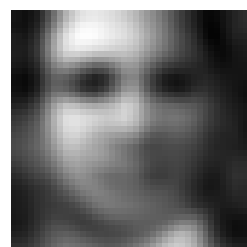


Dir 4

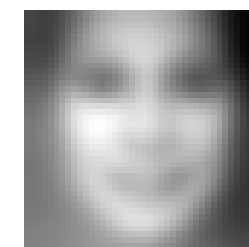
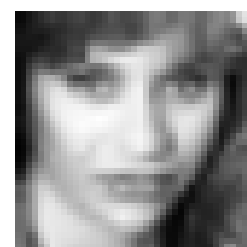
...



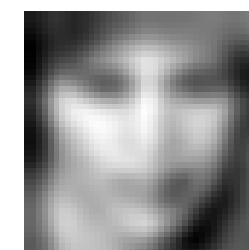
K=4



K=50



K=4



K=50