# Dimensionality Reduction by PCA

# Pattern Matrix

$$d$$

$$\mathbf{X} = \{ \mathbf{x}^k_i \}$$

$$\mathbf{x}^k$$

$$N$$

$$\mathbf{y} = \{ \mathbf{y}^k \}$$

$$\alpha$$

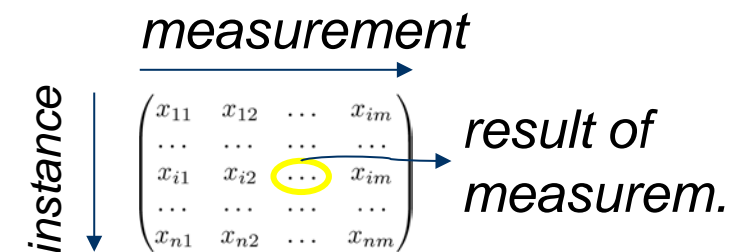$$\mathbf{w}$$

# Examples: Pattern Matrices

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

- ## Measurement vectors

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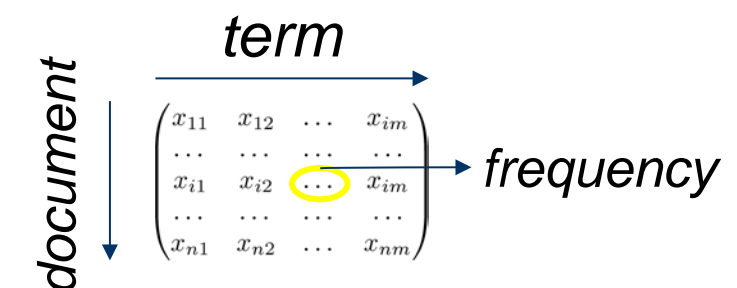
    - $j$: measurement, e.g. the area of a house

$$\begin{array}{c}\text{measurement} \\ \text{instance} \downarrow \begin{pmatrix} x_{11} & x_{12} & \dots & x_{im} \\ \dots & \dots & \dots & \dots \\ x_{i1} & x_{i2} & \dots & x_{im} \\ \dots & \dots & \dots & \dots \\ x_{n1} & x_{n2} & \dots & x_{nm} \end{pmatrix} \rightarrow \begin{array}{l}\text{result of} \\ \text{measurem.}\end{array}\end{array}$$

- ## Digital images as gray-scale vectors

    - $i$: image number

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

$$\begin{array}{c}\text{pixel} \\ \text{image} \downarrow \begin{pmatrix} x_{11} & x_{12} & \dots & x_{im} \\ \dots & \dots & \dots & \dots \\ x_{i1} & x_{i2} & \dots & x_{im} \\ \dots & \dots & \dots & \dots \\ x_{n1} & x_{n2} & \dots & x_{nm} \end{pmatrix} \rightarrow \text{intensity}\end{array}$$
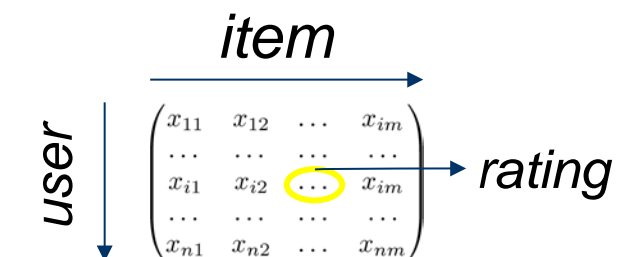
- ## Text documents in bag-of-words representation

    - $i$: document number

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

$$\begin{array}{c}\text{term} \\ \text{document} \downarrow \begin{pmatrix} x_{11} & x_{12} & \dots & x_{im} \\ \dots & \dots & \dots & \dots \\ x_{i1} & x_{i2} & \dots & x_{im} \\ \dots & \dots & \dots & \dots \\ x_{n1} & x_{n2} & \dots & x_{nm} \end{pmatrix} \rightarrow \text{frequency}\end{array}$$

- ## User rating data

    - $i$: user number

    - $j$: item (book, movie)

$$\begin{array}{c}\text{item} \\ \text{user} \downarrow \begin{pmatrix} x_{11} & x_{12} & \dots & x_{im} \\ \dots & \dots & \dots & \dots \\ x_{i1} & x_{i2} & \dots & x_{im} \\ \dots & \dots & \dots & \dots \\ x_{n1} & x_{n2} & \dots & x_{nm} \end{pmatrix} \rightarrow \text{rating}\end{array}$$

# 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 [...]*

$d_i$

**Document-Term Matrix**

$W$

| | $w_1$ | ... | $w_j$ | ... | $w_J$ |
|---|---|---|---|---|---|
| $d_1$ | | | | | |
| ... | | | ... | | |
| $d_i$ | | ... | $c(d_i, w_j)$ | ... | |
| ... | | | ... | | |
| $d_I$ | | | | | |

$D$

artifact   artificial   intelligence   interest

$t$

$d_i = $ | ... | 0 | 1 | ... | 2 | 0 | ... |

$X$

term weighting

5

# 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

0

1

0

2

# Vocabulary Mismatch & ~~Robustness~~

7

# 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 [...]*

$d_i$

**Document-Term Matrix**

$W$

artifact  artificial  intelligence  interest

$d_i = $ | ... | 0 | 1 | ... | 2 | 0 | ... |  t

X

term weighting

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

Clustering rows?            Clustering columns?

# Latent Structure

- Given a matrix that "encodes" data …

- Potential problems
  - too large
  - too complicated
  - missing entries
  - noisy entries
  - lack of structure
  - …

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

- 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?

9

# 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**"



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 $f(z)$: $\quad [x_1, x_2] = z * \underline{v} = z * [v_1, v_2]$
- $\underline{v}$ 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: $a^{(l)}$ is the projection of each point $x^{(l)}$ onto v
  - v chosen to minimize the residual variance
  - Equivalently, v is the direction of maximum variance
  - Extensions: best two dimensions: xi= ai*v + bi*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 $C = X^T X$
4. Compute k largest eigenvectors of C

$$C = VDV^T$$

# Singular Value Decomposition

- Alternative method to calculate (still subtract mean 1$^{st}$)
- 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*S matrix provides coefficients
  - Example $x_I = U_{I,1} S_{11} v_1 + U_{I,2} S_{22} v_2 + \ldots$

- Gives the least-squares approximation to X of this form

$$\underset{N \, x \, D}{X} \approx \underset{N \, x \, K}{U} \quad \underset{K \, x \, K}{S} \quad \underset{K \, x \, D}{V^T}$$

# Glorious SVD

$$X = USV^T$$

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

# Collaborative Filtering (Netflix)

From Y. Koren
of BellKor team

users

|  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 |  | 3 |  | ? | 5 |  |  | 5 |  | 4 |  |
| 2 |  |  | 5 | 4 |  |  | 4 |  |  | 2 | 1 | 3 |
| 3 | 2 | 4 |  | 1 | 2 |  | 3 |  | 4 | 3 | 5 |  |
| 4 |  | 2 | 4 |  | 5 |  |  | 4 |  |  | 2 |  |
| 5 |  |  | 4 | 3 | 4 | 2 |  |  |  |  | 2 | 5 |
| 6 | 1 |  | 3 |  | 3 |  |  | 2 |  |  | 4 |  |

movies

$$X_{N \times D} \approx U_{N \times K} \quad S_{K \times K} \quad V^T_{K \times D}$$

# Latent Space Models

Model ratings matrix as "user" and "movie" positions

Infer values from known ratings

Extrapolate to unranked

# Latent Space Models

**serious**

Braveheart

The Color
Purple

Amadeus

Lethal Weapon

Sense and
Sensibility

**"Chick
flicks"?**

Ocean's 11

The Lion King

Dumb and
Dumber

The Princess
Diaries

Independence
Day

**escapist**

# "Eigen-faces"

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



$$X$$
$$N \times D$$

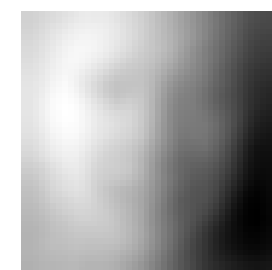# "Eigen-faces"

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


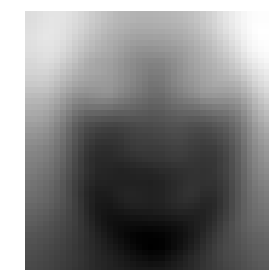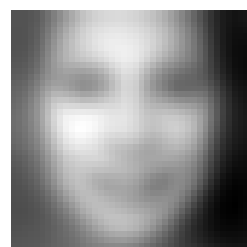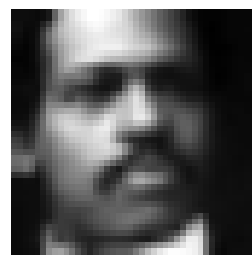
$$X_{N \times D} \approx U_{N \times K} \quad S_{K \times K} \quad V^T_{K \times D}$$

| 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          ...
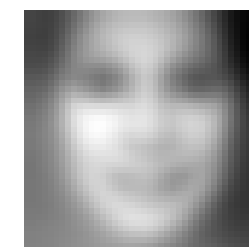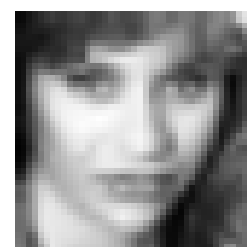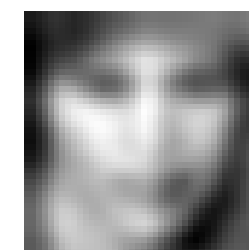


K=4

K=50
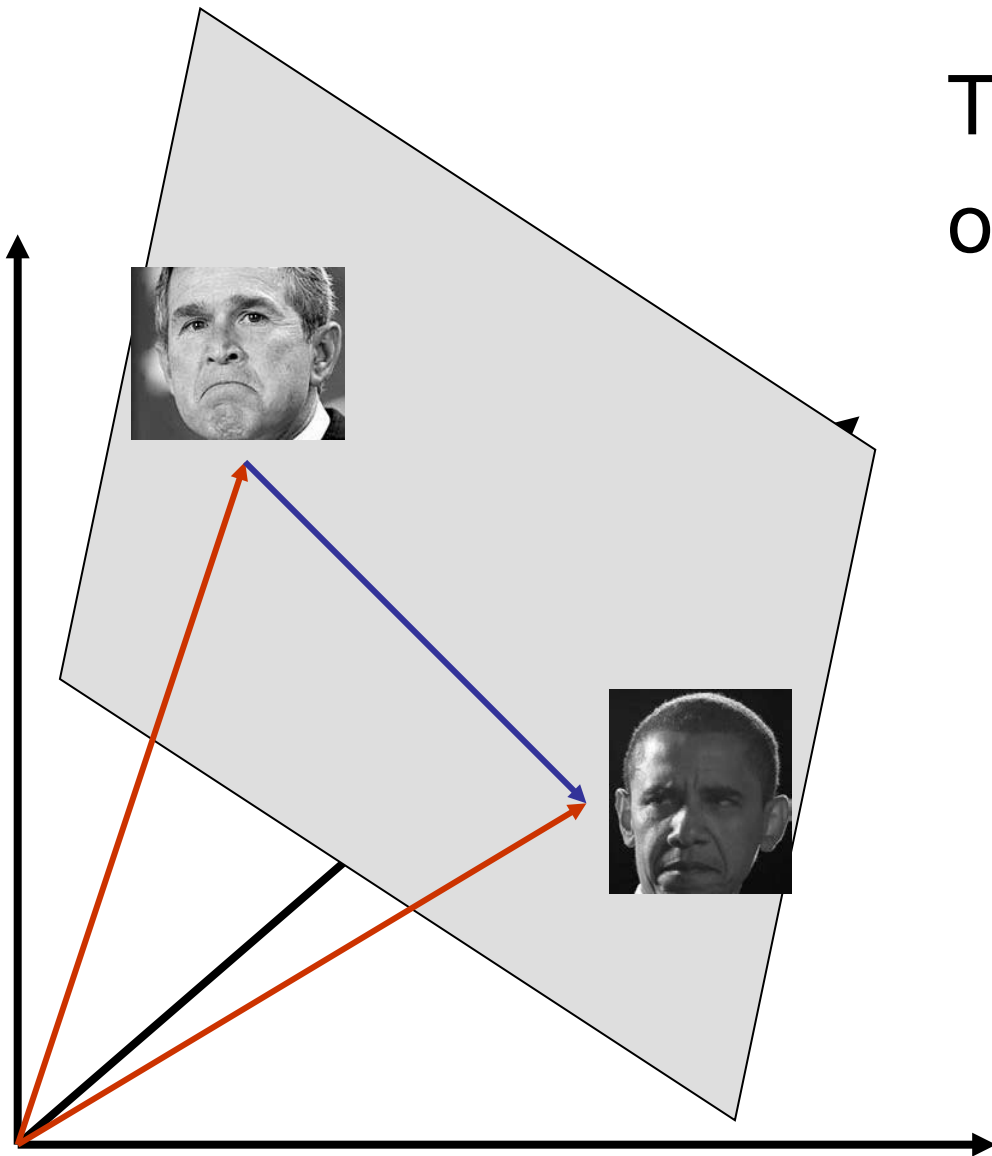
K=4

K=50

# The Face Subspace

The set of faces is a "subspace" of the set of images

- Suppose it is K dimensional

- We can find the best subspace using PCA

- This is like fitting a "hyper-plane" to the set of faces

  • spanned by vectors $\mathbf{v}_1, \mathbf{v}_2, ..., \mathbf{v}_K$

Any face: $\mathbf{x} \approx \overline{\mathbf{x}} + a_1\mathbf{v}_1 + a_2\mathbf{v}_2 + \ldots + a_k\mathbf{v}_k$
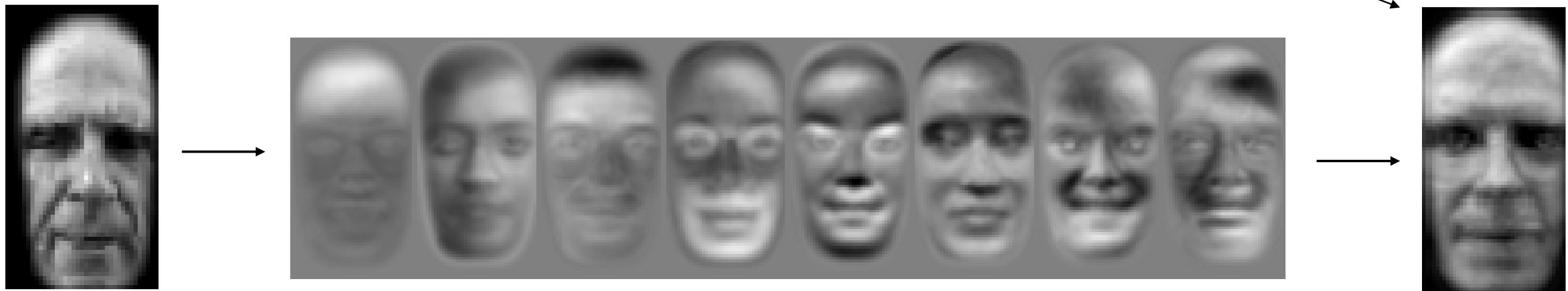
# Projecting onto the Eigenface Subspace

- ## The eigenfaces $\mathbf{v_1}$, ..., $\mathbf{v_K}$ span the space of faces

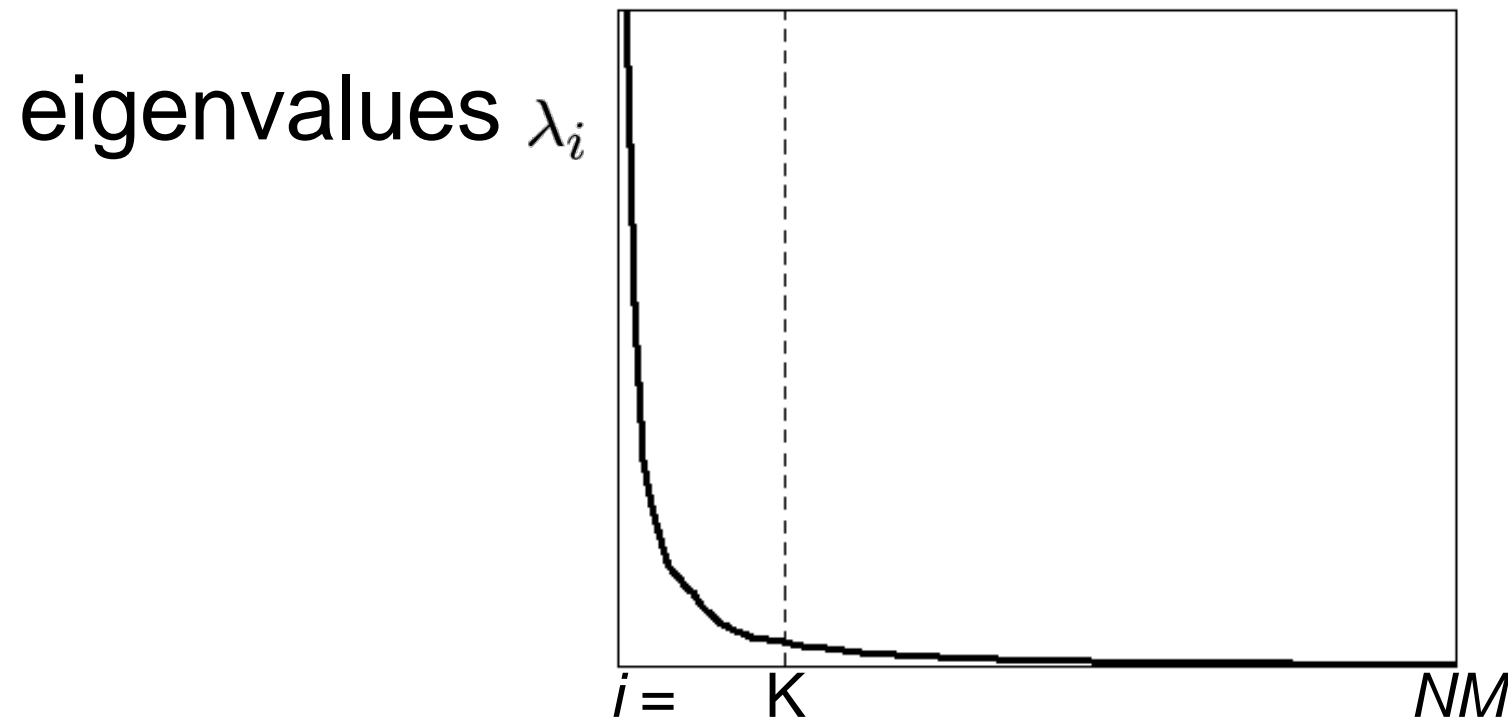  - A face is converted to eigenface coordinates by

$$\mathbf{x} \rightarrow (\underbrace{(\mathbf{x} - \bar{\mathbf{x}}) \cdot \mathbf{v_1}}_{a_1}, \; \underbrace{(\mathbf{x} - \bar{\mathbf{x}}) \cdot \mathbf{v_2}}_{a_2}, \ldots, \; \underbrace{(\mathbf{x} - \bar{\mathbf{x}}) \cdot \mathbf{v_K}}_{a_K})$$

$$\mathbf{x} \approx \bar{\mathbf{x}} + a_1 \mathbf{v_1} + a_2 \mathbf{v_2} + \ldots + a_K \mathbf{v_K}$$

# Choosing the Dimension K



eigenvalues $\lambda_i$

$i =$    K                    NM
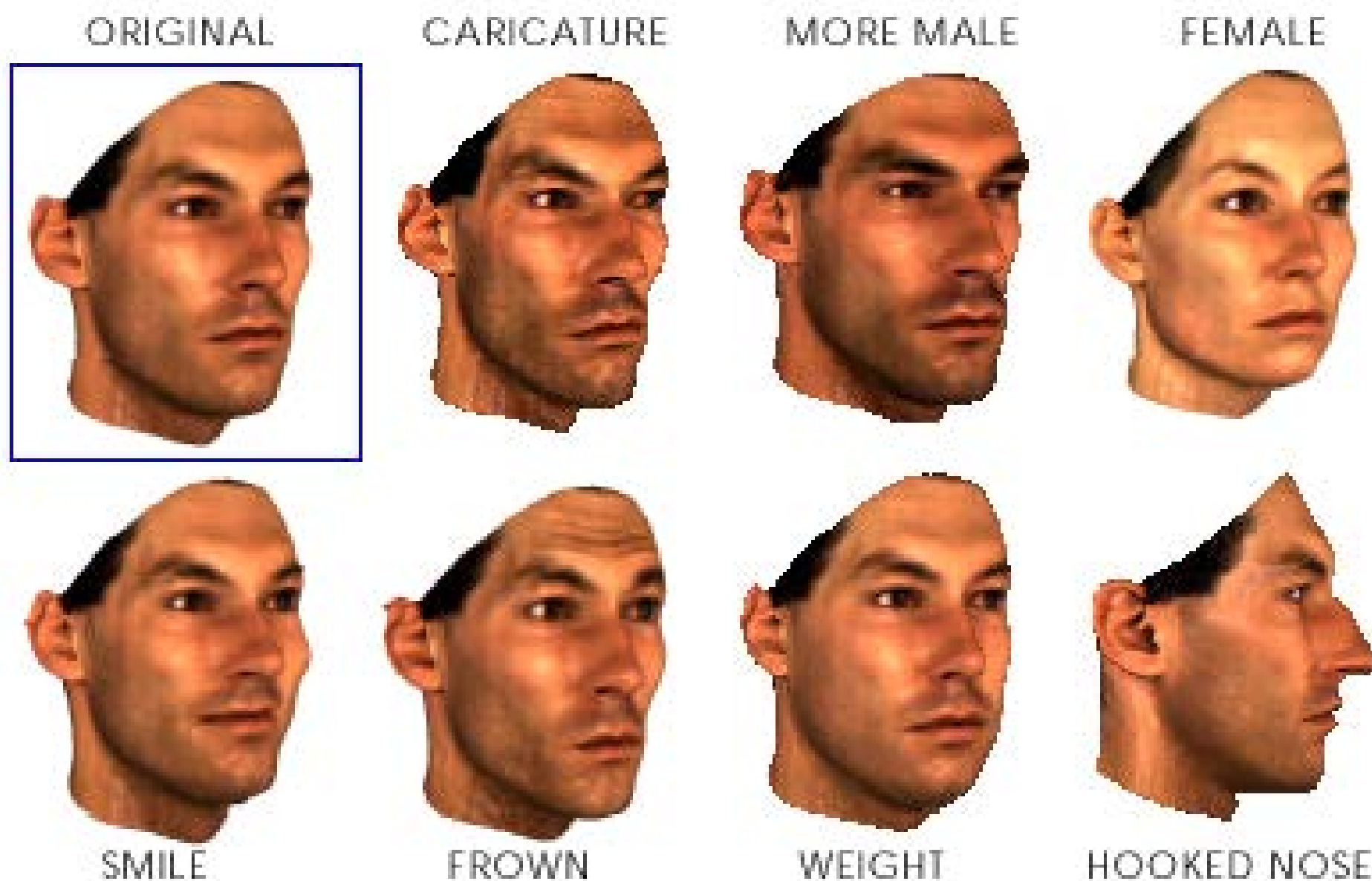
- How many eigenfaces to use?

- Look at the decay of the eigenvalues

  - the eigenvalue tells you the amount of variance "in the direction" of that eigenface
  - ignore eigenfaces with low variance

# PCA with depth data: Blinz & Vetter, 1999

# Non-linear Dimensionality Reduction

# Learn an embedding ("self-supervision")

[Collobert and Weston 2008]



house, where the professor lived without his wife and child; or so he said jokingly sometimes: "Here's where I live. My house." His daughter often added, without resentment, for the visitor's information, "It started out to be for me, but it's really his." And she might reach in to bring forth an inch-high table lamp with fluted shade, or a blue dish the size of her little fingernail, marked "Kitty" and half full of eternal milk; but she was sure to replace these, after they had been admired, pretty near exactly where they had been. The little house was very orderly, and just big enough for all it contained, though to some tastes the bric-à-brac in the parlor might seem excessive. The daughter's preference was for the store-bought gimmicks and ap little es, the toasters and carpet sweepers of Lilliput, but she knew that most adult visitors would

**True Snippet**

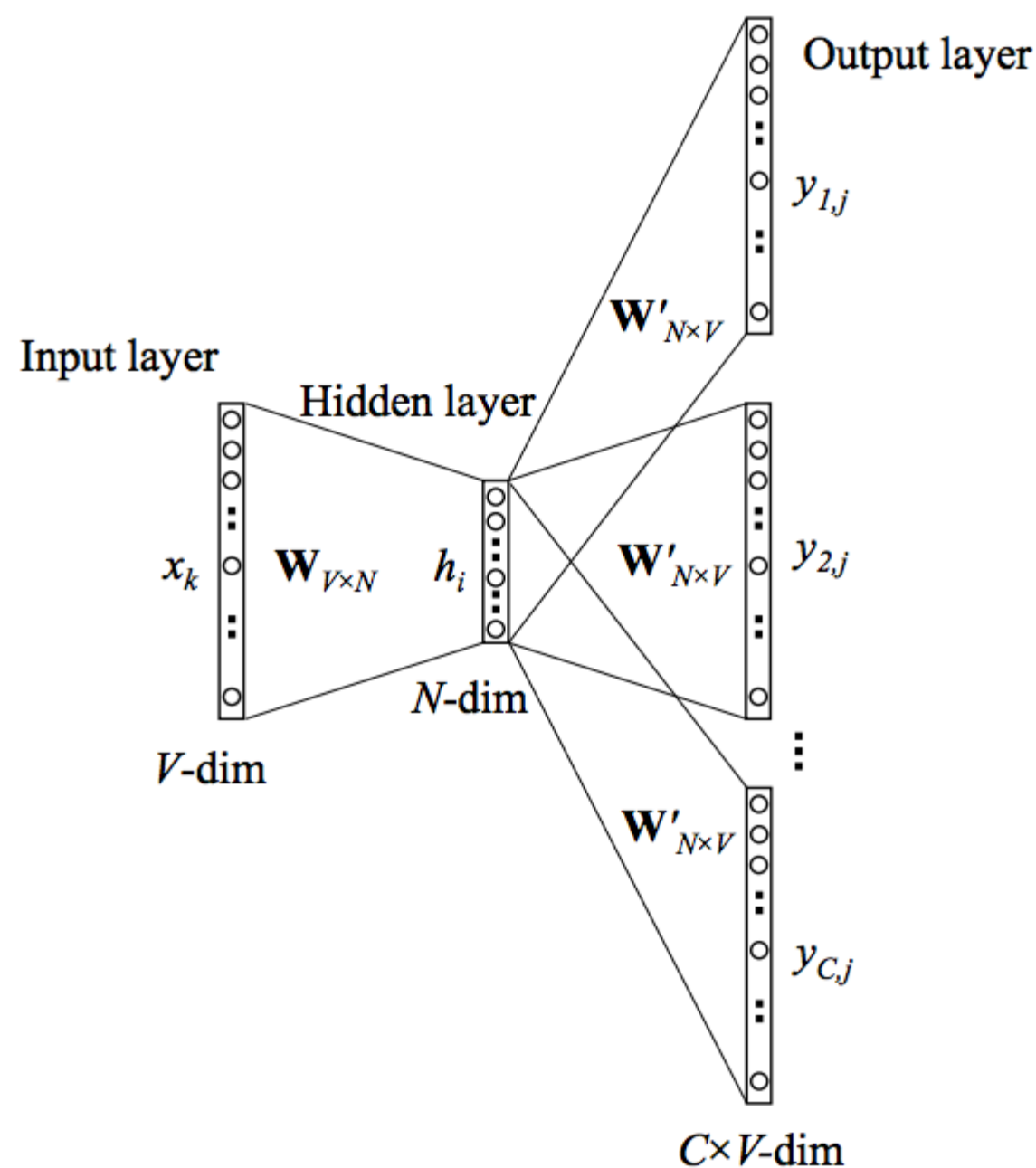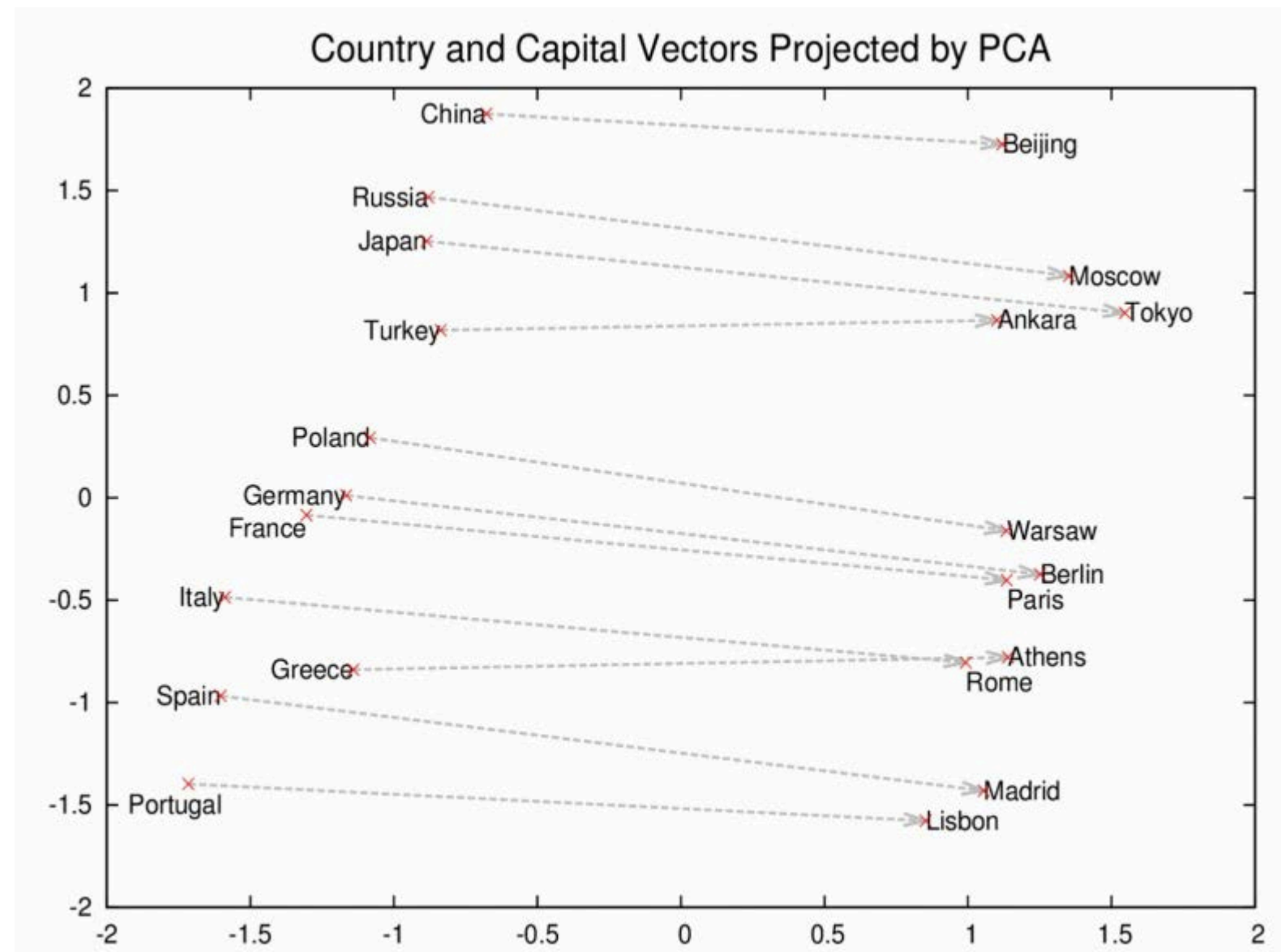**NN**

# Skip Gram (word2vec)

# Learning word2vec embeddings

| Word | Cosine distance |
| --- | --- |
| norway | 0.760124 |
| denmark | 0.715460 |
| finland | 0.620022 |
| switzerland | 0.588132 |
| belgium | 0.585835 |
| netherlands | 0.574631 |
| iceland | 0.562368 |
| estonia | 0.547621 |
| slovenia | 0.531408 |

## Country and Capital Vectors Projected by PCA

# Example

$$vec(\text{"man"}) - vec(\text{"king"}) + vec(\text{"woman"}) = vec(\text{"queen"})$$



http://deeplearning4j.org/word2vec.html

# Visual Context Task



Carl Doersch, Abhinav Gupta, and Alexei A. Efros. **Unsupervised Visual Representation Learning by Context Prediction.** In *ICCV 2015*