

Machine Learning

Prof. Barbara Caputo

Dip. Ingegneria Informatica, Automatica e Gestionale, Roma



SAPIENZA
UNIVERSITÀ DI ROMA

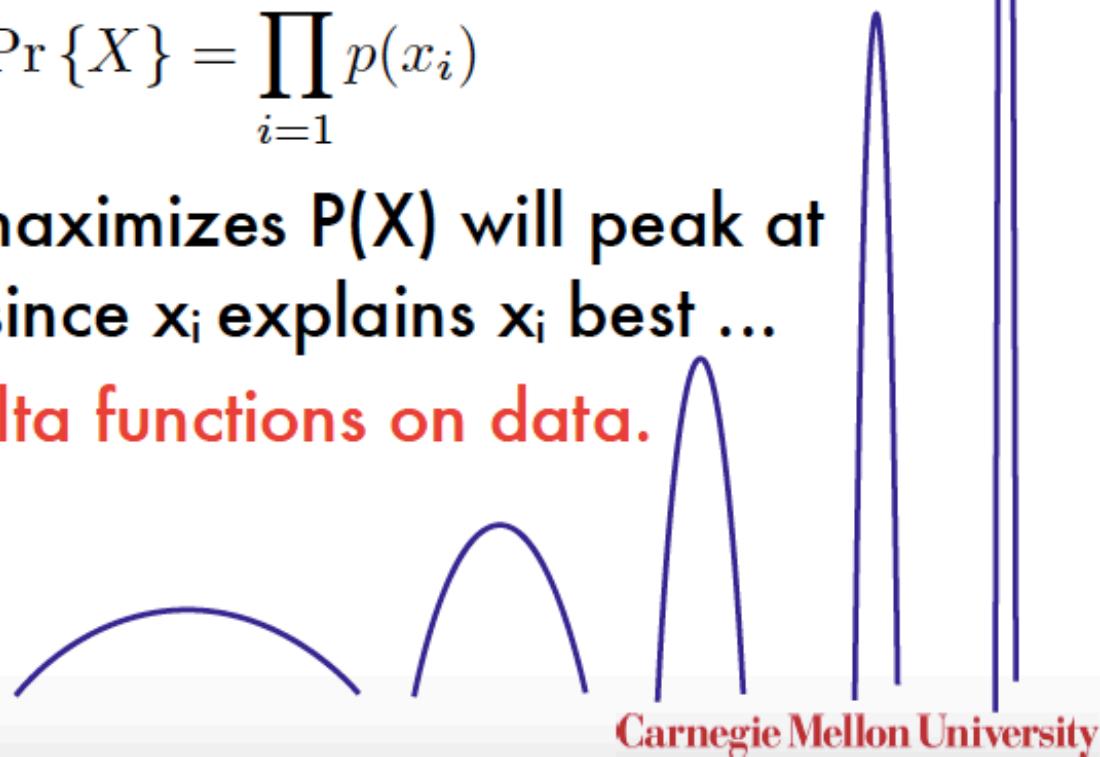
Model Selection

Maximum Likelihood

- Need to measure how well we do
- For density estimation we care about

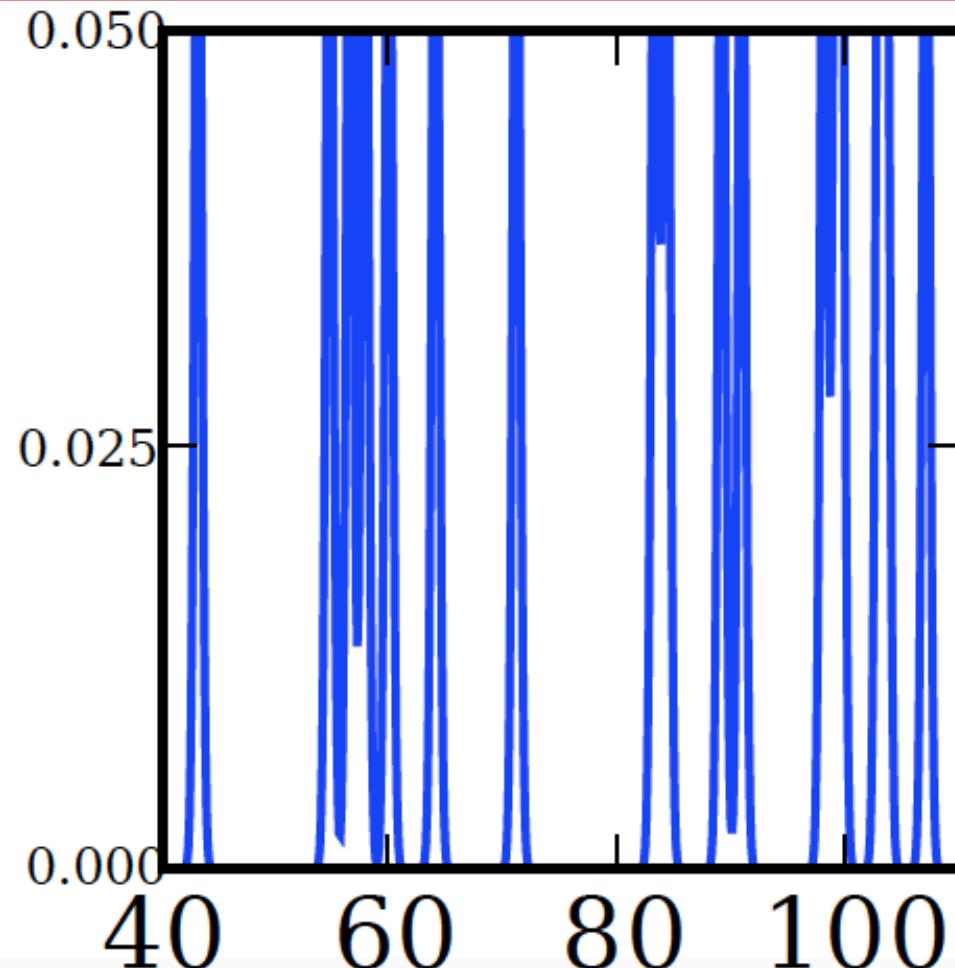
$$\Pr \{X\} = \prod_{i=1}^m p(x_i)$$

- Finding α that maximizes $P(X)$ will peak at all data points since x_i explains x_i best ...
- Maxima are delta functions on data.
- Overfitting!



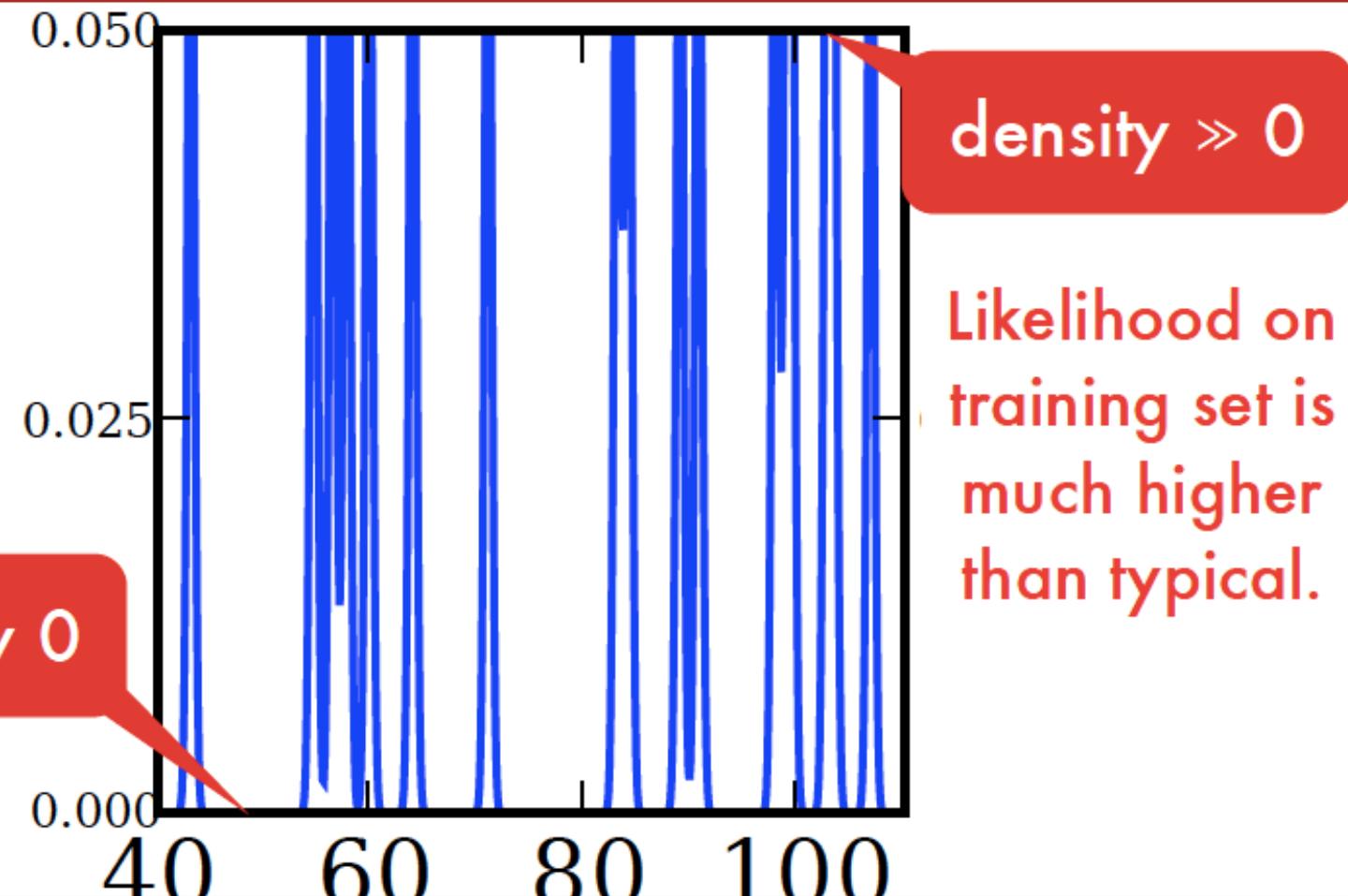
Carnegie Mellon University

Overfitting



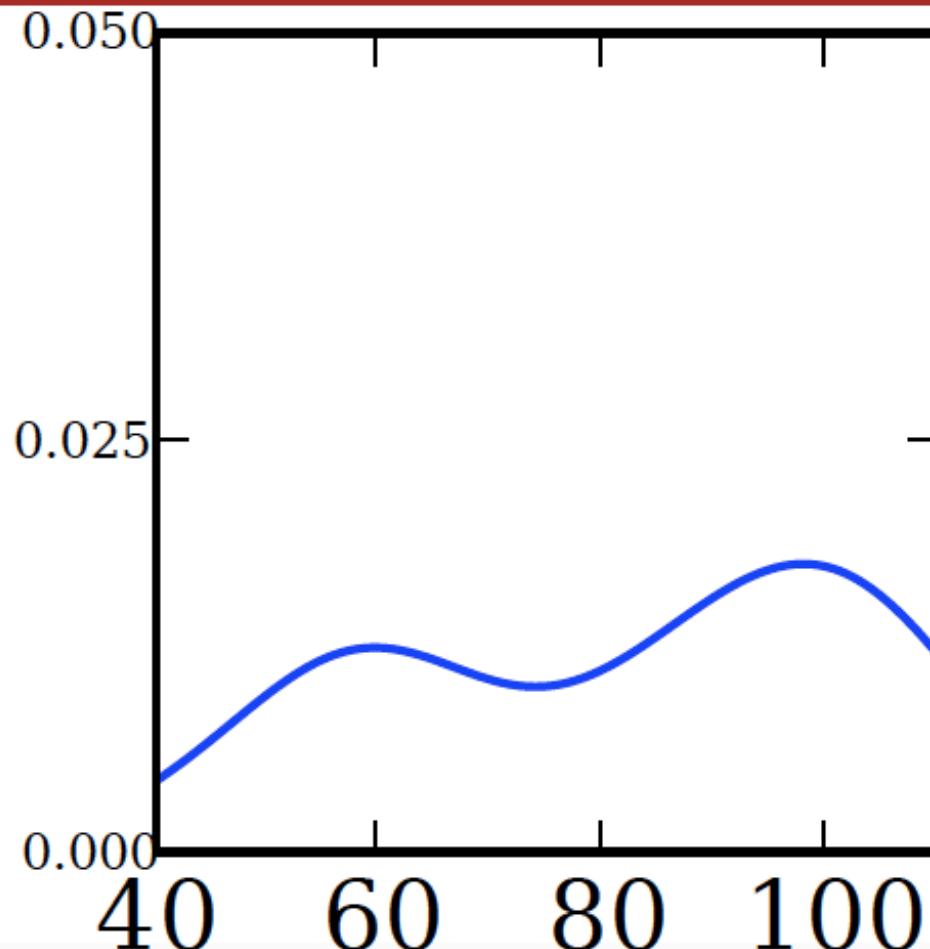
Likelihood on
training set is
much higher
than typical.

Overfitting



Carnegie Mellon University

Underfitting



Likelihood on
training set is
very similar to
typical one.

Too simple.

Model Selection

- Validation

- Use some of the data to estimate density.
- Use other part to evaluate how well it works
- Pick the parameter that works best

$$\mathcal{L}(X'|X) := \frac{1}{n'} \sum_{i=1}^{n'} \log \hat{p}(x'_i)$$

• • •

Model Selection

- Validation

- Use some of the data to estimate density.
- Use other part to evaluate how well it works
- Pick the parameter that works best

$$\mathcal{L}(X'|X) := \frac{1}{n'} \sum_{i=1}^{n'} \log \hat{p}(x'_i)$$

- Learning Theory

- Use data to build model
- Measure complexity and use this to bound

$$\frac{1}{n} \sum_{i=1}^n \log \hat{p}(x_i) - \mathbf{E}_x [\log \hat{p}(x)]$$

Carnegie Mellon University

Model Selection

- Validation

- Use some of the data to estimate density.
- Use other part to evaluate how well it works
- Pick the parameter that works best

easy

$$\mathcal{L}(X'|X) := \frac{1}{n'} \sum_{i=1}^{n'} \log \hat{p}(x'_i)$$

- Learning Theory

- Use data to build model
- Measure complexity and use this to bound

$$\frac{1}{n} \sum_{i=1}^n \log \hat{p}(x_i) - \mathbf{E}_x [\log \hat{p}(x)]$$

Model Selection

- Validation

- Use some of the data to estimate density.
- Use other part to evaluate how well it works
- Pick the parameter that works best

easy

$$\mathcal{L}(X'|X) := \frac{1}{n'} \sum_{i=1}^{n'} \log \hat{p}(x'_i)$$

wasteful

- Learning Theory

- Use data to build model
- Measure complexity and use this to bound

$$\frac{1}{n} \sum_{i=1}^n \log \hat{p}(x_i) - \mathbf{E}_x [\log \hat{p}(x)]$$

Carnegie Mellon University

Model Selection

- Validation

- Use some of the data to estimate density.
- Use other part to evaluate how well it works
- Pick the parameter that works best

easy

$$\mathcal{L}(X'|X) := \frac{1}{n'} \sum_{i=1}^{n'} \log \hat{p}(x'_i)$$

wasteful

- Learning Theory

- Use data to build model
- Measure complexity and use this to bound

difficult

$$\frac{1}{n} \sum_{i=1}^n \log \hat{p}(x_i) - \mathbf{E}_x [\log \hat{p}(x)]$$

Carnegie Mellon University

Model Selection

- Leave-one-out Crossvalidation

Model Selection

- Leave-one-out Crossvalidation
 - Use **almost all** data to estimate density.

Model Selection

- Leave-one-out Crossvalidation
 - Use **almost all** data to estimate density.
 - Use single instance to estimate how well it works

Model Selection

- Leave-one-out Crossvalidation
 - Use **almost all** data to estimate density.
 - Use **single instance** to estimate how well it works

$$\log p(x_i|X \setminus x_i) = \log \frac{1}{n-1} \sum_{j \neq i} k(x_i, x_j)$$

—

Model Selection

- Leave-one-out Crossvalidation
 - Use **almost all** data to estimate density.
 - Use single instance to estimate how well it works

$$\log p(x_i | X \setminus x_i) = \log \frac{1}{n-1} \sum_{j \neq i} k(x_i, x_j)$$

- This has huge variance

Model Selection

- Leave-one-out Crossvalidation
 - Use **almost all** data to estimate density.
 - Use single instance to estimate how well it works

$$\log p(x_i|X \setminus x_i) = \log \frac{1}{n-1} \sum_{j \neq i} k(x_i, x_j)$$

- This has huge variance
- Average over estimates for all training data

Model Selection

- Leave-one-out Crossvalidation
 - Use **almost all** data to estimate density.
 - Use single instance to estimate how well it works

$$\log p(x_i|X \setminus x_i) = \log \frac{1}{n-1} \sum_{j \neq i} k(x_i, x_j)$$

- This has huge variance
- Average over estimates for all training data
- Pick the parameter that works best

Model Selection

- Leave-one-out Crossvalidation
 - Use almost all data to estimate density.
 - Use single instance to estimate how well it works

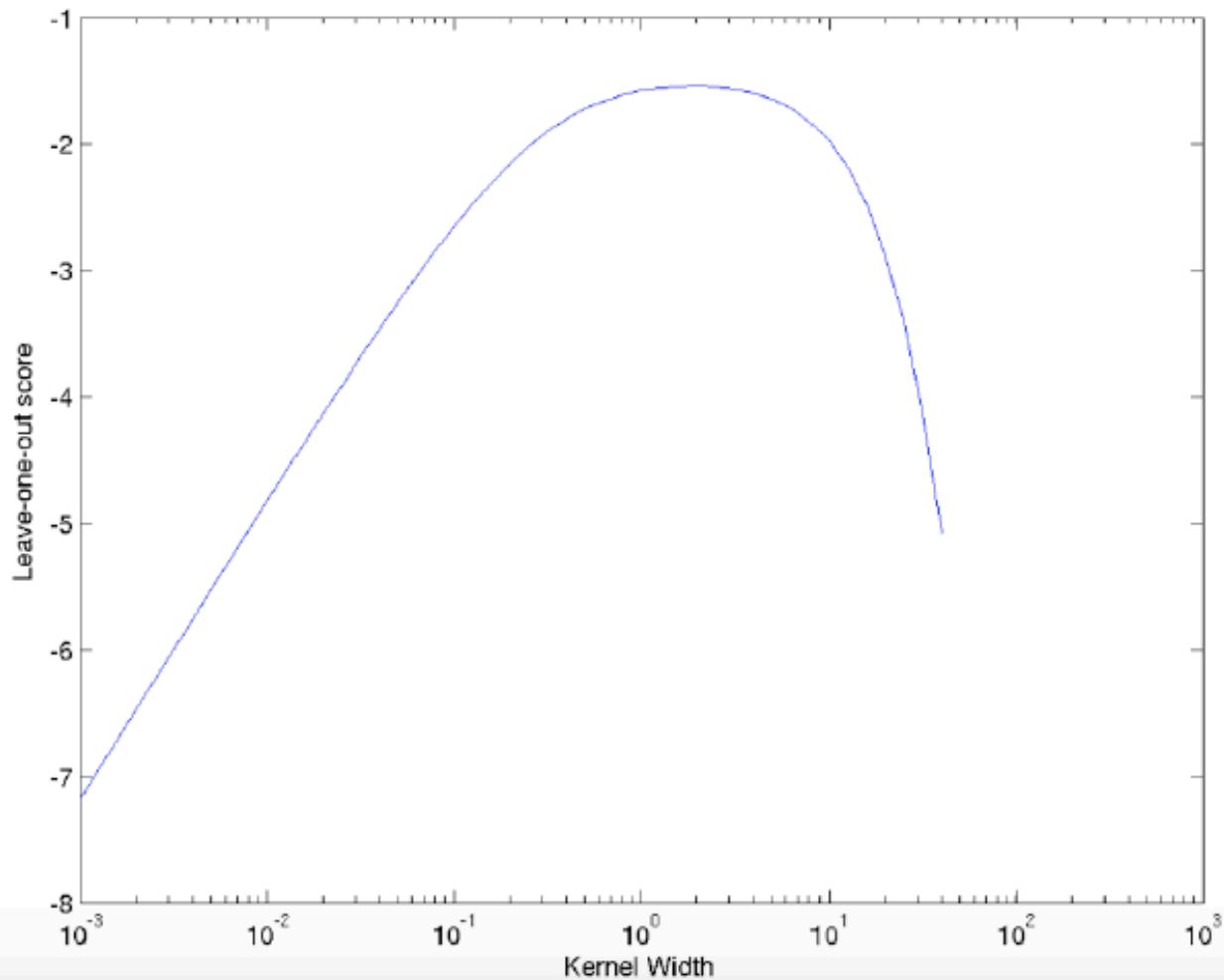
$$\log p(x_i|X \setminus x_i) = \log \frac{1}{n-1} \sum_{j \neq i} k(x_i, x_j)$$

- This has huge variance
- Average over estimates for all training data
- Pick the parameter that works best
- Simple implementation

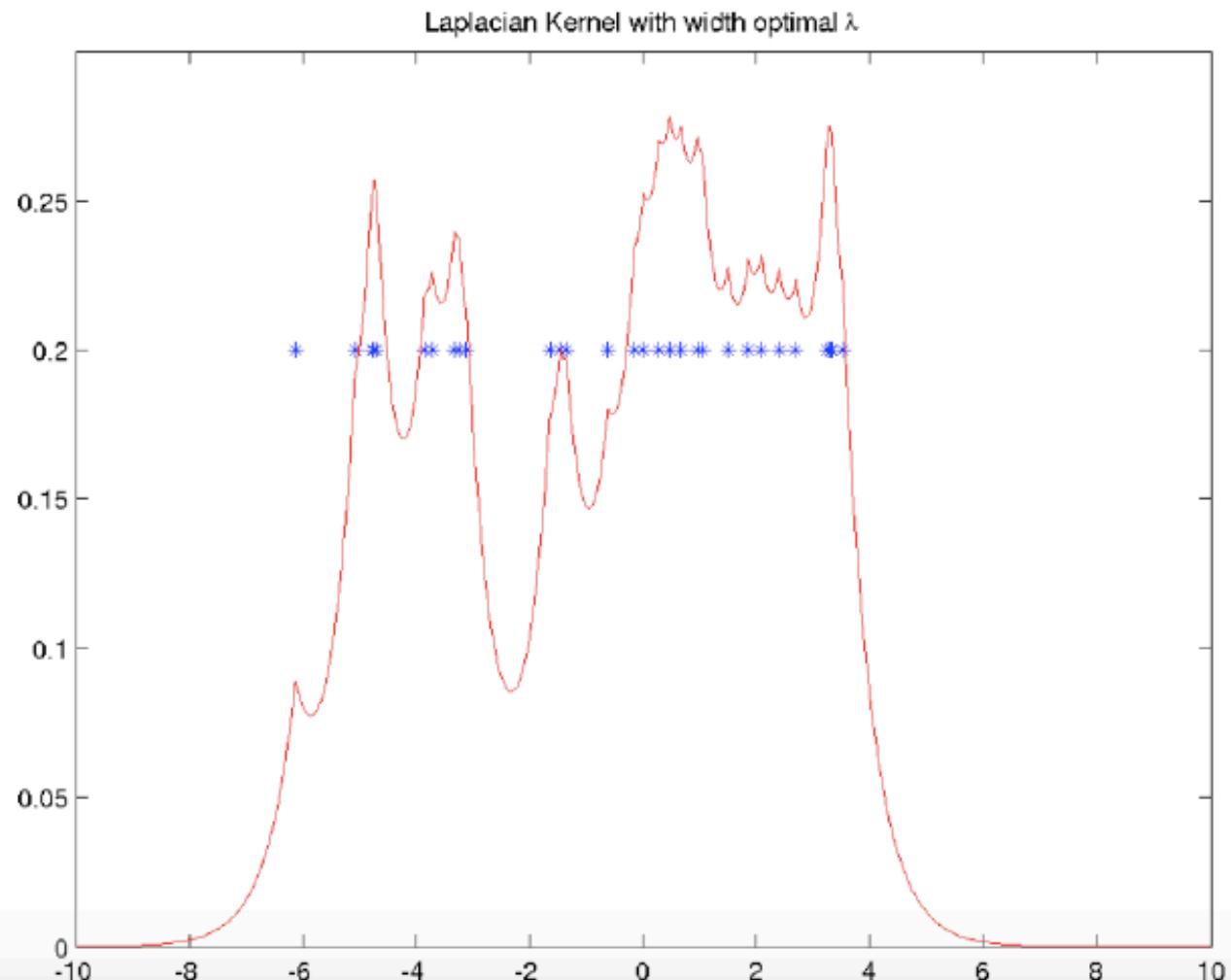
$$\frac{1}{n} \sum_{i=1}^n \log \left[\frac{n}{n-1} p(x_i) - \frac{1}{n-1} k(x_i, x_i) \right] \text{ where } p(x) = \frac{1}{n} \sum_{i=1}^n k(x_i, x)$$

Carnegie Mellon University

Leave-one out estimate



Optimal estimate



Model Selection

- k-fold Crossvalidation

Model Selection

- k-fold Crossvalidation
 - Partition data into k blocks (typically 10)

Model Selection

- k-fold Crossvalidation
 - Partition data into k blocks (typically 10)
 - Use all but one block to compute estimate

Model Selection

- k-fold Crossvalidation
 - Partition data into k blocks (typically 10)
 - Use all but one block to compute estimate
 - Use remaining block as validation set

Model Selection

- k-fold Crossvalidation
 - Partition data into k blocks (typically 10)
 - Use all but one block to compute estimate
 - Use remaining block as validation set
 - Average over all validation estimates

Model Selection

- k-fold Crossvalidation
 - Partition data into k blocks (typically 10)
 - Use all but one block to compute estimate
 - Use remaining block as validation set
 - Average over all validation estimates

$$\frac{1}{k} \sum_{i=1}^k l(p(X_i | X \setminus X_i))$$

Model Selection

- k-fold Crossvalidation
 - Partition data into k blocks (typically 10)
 - Use all but one block to compute estimate
 - Use remaining block as validation set
 - Average over all validation estimates
$$\frac{1}{k} \sum_{i=1}^k l(p(X_i | X \setminus X_i))$$
- Almost unbiased (e.g. via Luntz and Brailovski, 1969)
(error is for $(k-1)/k$ sized set)

Model Selection

- k-fold Crossvalidation
 - Partition data into k blocks (typically 10)
 - Use all but one block to compute estimate
 - Use remaining block as validation set
 - Average over all validation estimates
$$\frac{1}{k} \sum_{i=1}^k l(p(X_i | X \setminus X_i))$$
- Almost unbiased (e.g. via Luntz and Brailovski, 1969)
(error is for $(k-1)/k$ sized set)
- Pick best parameter

Contents

- ❑ Motivation
- ❑ PCA algorithms
- ❑ Applications

Some of these slides are taken from

- Karl Booksh Research group
- Tom Mitchell
- Ron Parr

2

Motivation

PCA Applications

- Data Visualization
- Data Compression
- Noise Reduction

Data Visualization

Example:

- Given 53 blood and urine samples (features) from 65 people.
- How can we visualize the measurements?

Data Visualization

- Matrix format (65x53)

	H-WBC	H-RBC	H-Hgb	H-Hct	H-MCV	H-MCH	H-MCHC
A1	8.0000	4.8200	14.1000	41.0000	85.0000	29.0000	34.0000
A2	7.3000	5.0200	14.7000	43.0000	86.0000	29.0000	34.0000
A3	4.3000	4.4800	14.1000	41.0000	91.0000	32.0000	35.0000
A4	7.5000	4.4700	14.9000	45.0000	101.0000	33.0000	33.0000
A5	7.3000	5.5200	15.4000	46.0000	84.0000	28.0000	33.0000
A6	6.9000	4.8600	16.0000	47.0000	97.0000	33.0000	34.0000
A7	7.8000	4.6800	14.7000	43.0000	92.0000	31.0000	34.0000
A8	8.6000	4.8200	15.8000	42.0000	88.0000	33.0000	37.0000
A9	5.1000	4.7100	14.0000	43.0000	92.0000	30.0000	32.0000

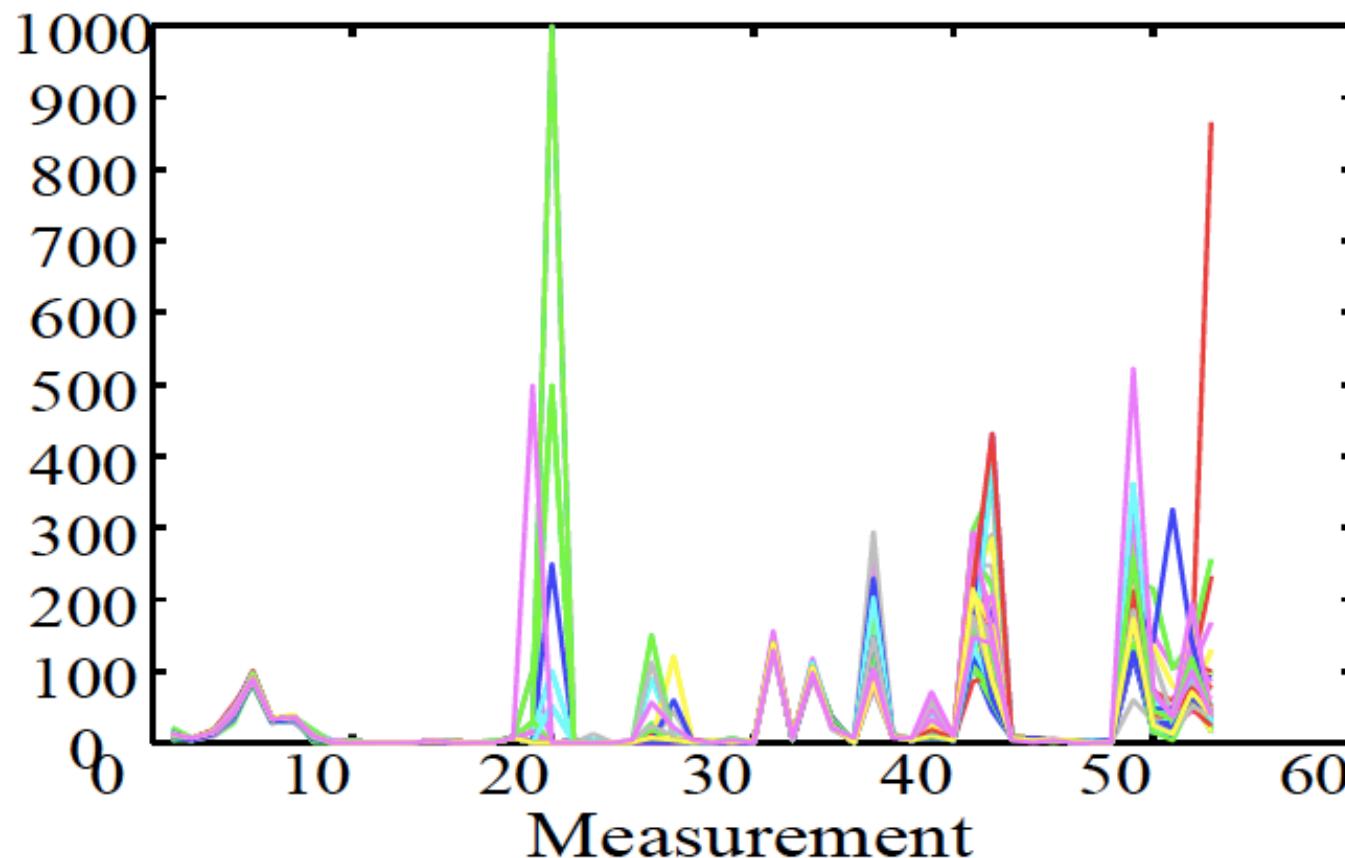
Instances

Features

Difficult to see the correlations between the features...

Data Visualization

- Spectral format (65 curves, one for each person)

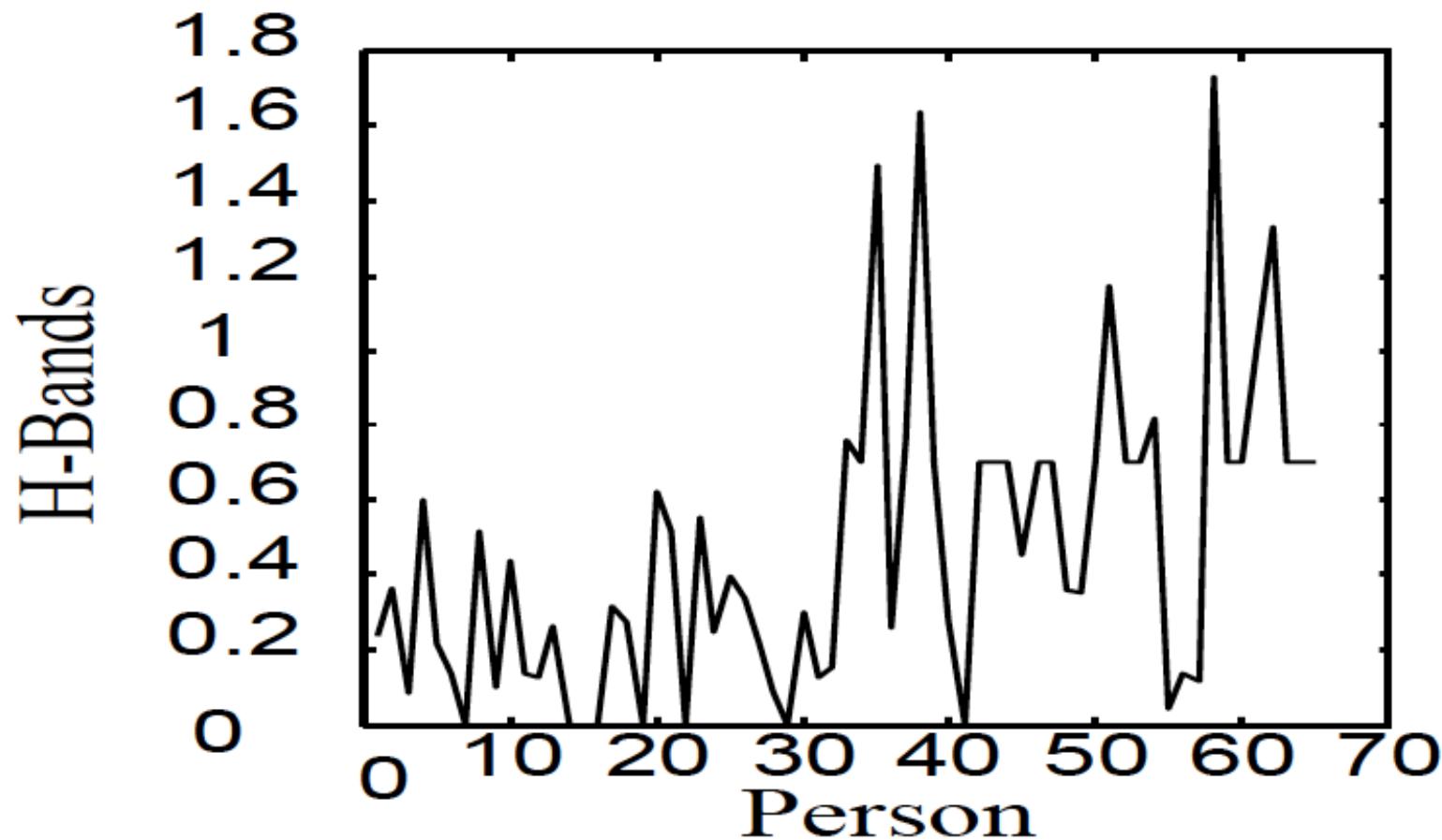


Difficult to compare the different patients...

7

Data Visualization

- Spectral format (53 pictures, one for each feature)

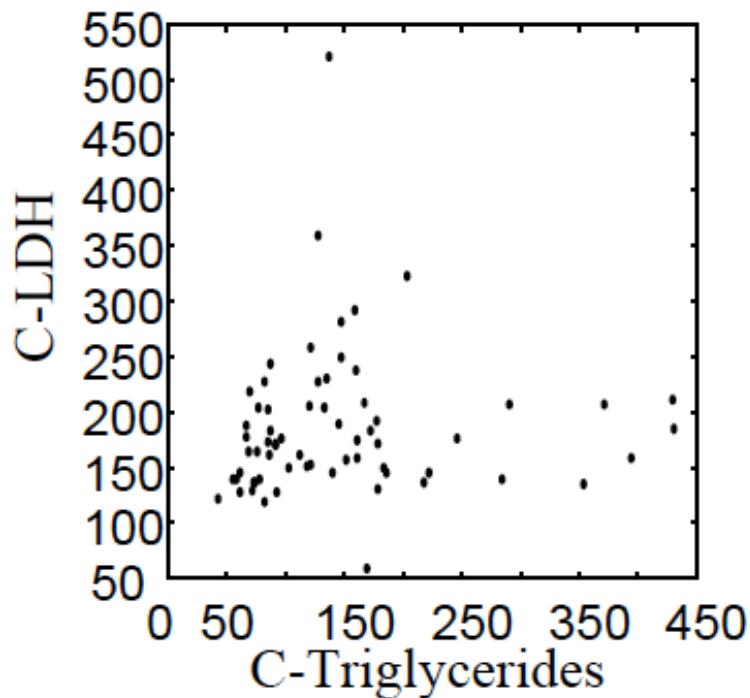


Difficult to see the correlations between the features...

8

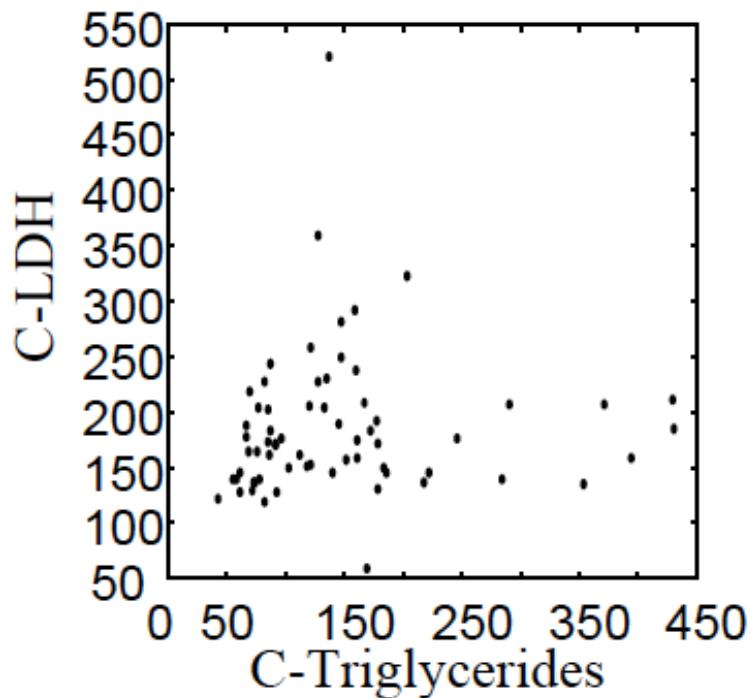
Data Visualization

Bi-variate

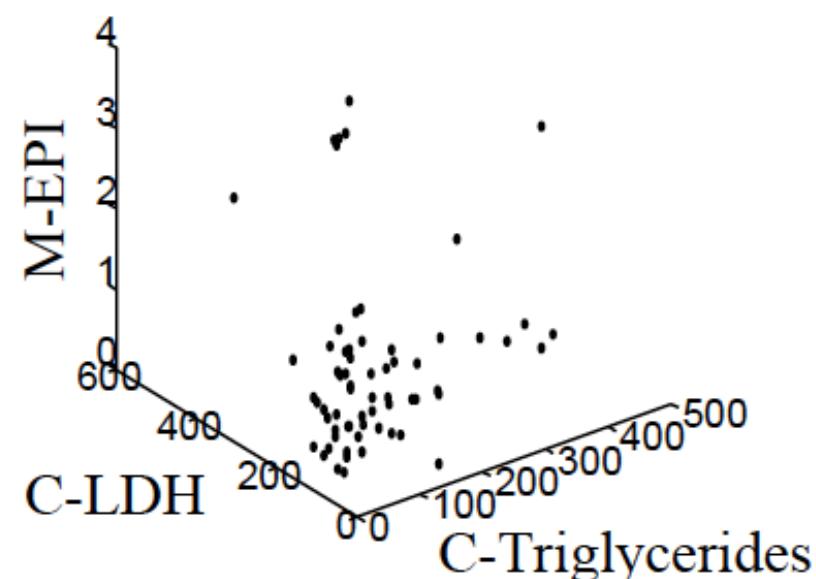


Data Visualization

Bi-variate

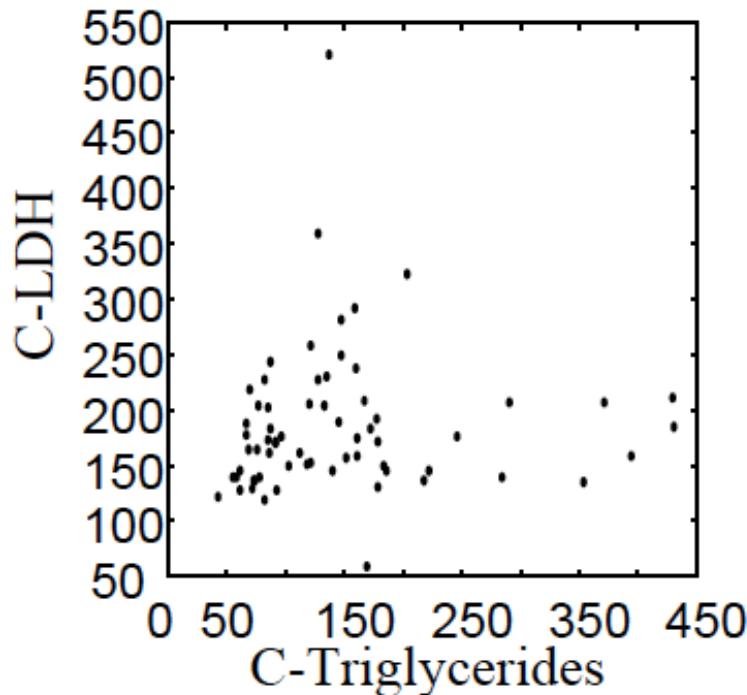


Tri-variate

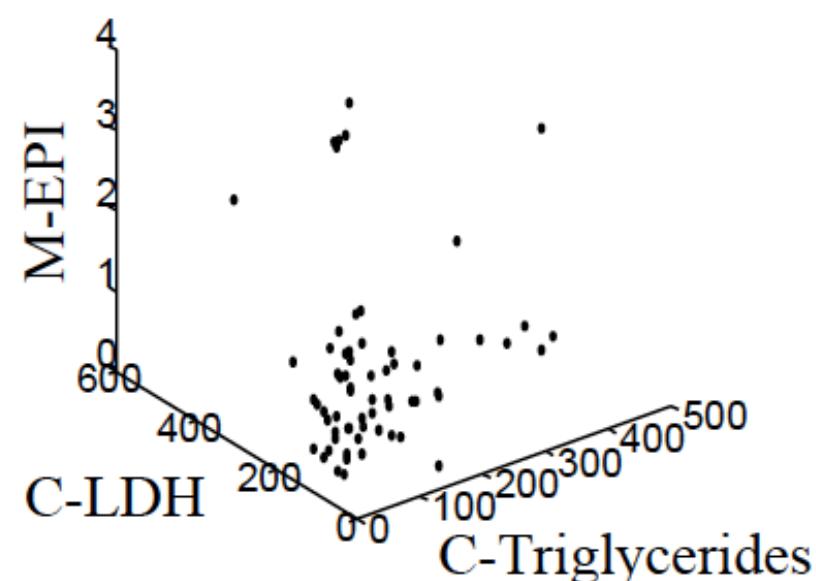


Data Visualization

Bi-variate



Tri-variate



How can we visualize the other variables???

... difficult to see in 4 or higher dimensional spaces...

Data Visualization

- Is there a representation better than the coordinate axes?

Data Visualization

- Is there a representation better than the coordinate axes?
- Is it really necessary to show all the 53 dimensions?
 - ... what if there are strong correlations between the features?

Data Visualization

- Is there a representation better than the coordinate axes?
- Is it really necessary to show all the 53 dimensions?
 - ... what if there are strong correlations between the features?
- How could we find the *smallest* subspace of the 53-D space that keeps the *most information* about the original data?

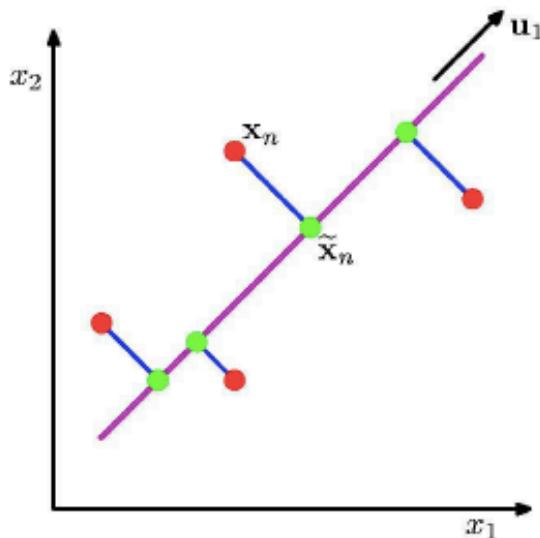
Data Visualization

- Is there a representation better than the coordinate axes?
- Is it really necessary to show all the 53 dimensions?
 - ... what if there are strong correlations between the features?
- How could we find the *smallest* subspace of the 53-D space that keeps the *most information* about the original data?
- A solution: **Principal Component Analysis**

PCA Algorithms

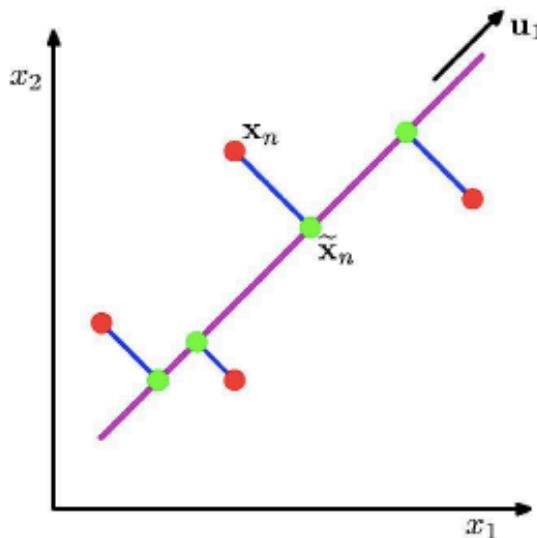
Principal Component Analysis

PCA:



Principal Component Analysis

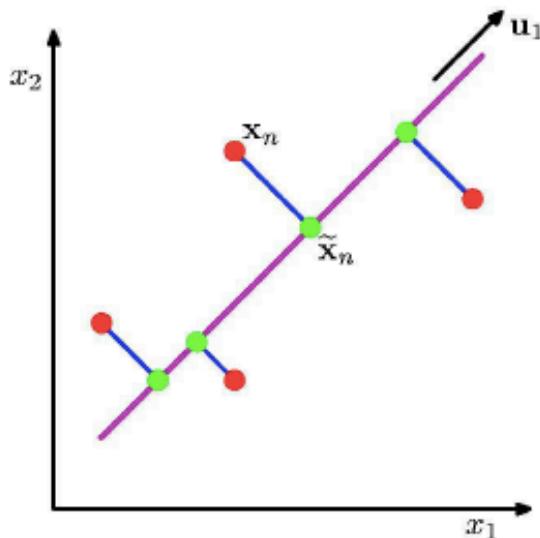
PCA:



Orthogonal projection of the data onto a lower-dimension linear space that...

- maximizes variance of projected data (purple line)

Principal Component Analysis



PCA:

Orthogonal projection of the data onto a lower-dimension linear space that...

- maximizes variance of projected data (purple line)
- minimizes mean squared distance between
 - data point and
 - projections (sum of blue lines)

12

Principal Component Analysis

Idea:

- Given data points in a d -dimensional space, project them into a **lower dimensional** space while **preserving as much information** as possible.
 - Find best planar approximation to 3D data
 - Find best 12-D approximation to 10^4 -D data

Principal Component Analysis

Idea:

- Given data points in a d -dimensional space, project them into a **lower dimensional** space while **preserving as much information** as possible.
 - Find best planar approximation to 3D data
 - Find best 12-D approximation to 10^4 -D data
- In particular, choose projection that **minimizes *squared error*** in reconstructing the original data.

Principal Component Analysis

□ **PCA Vectors** originate from the center of mass.

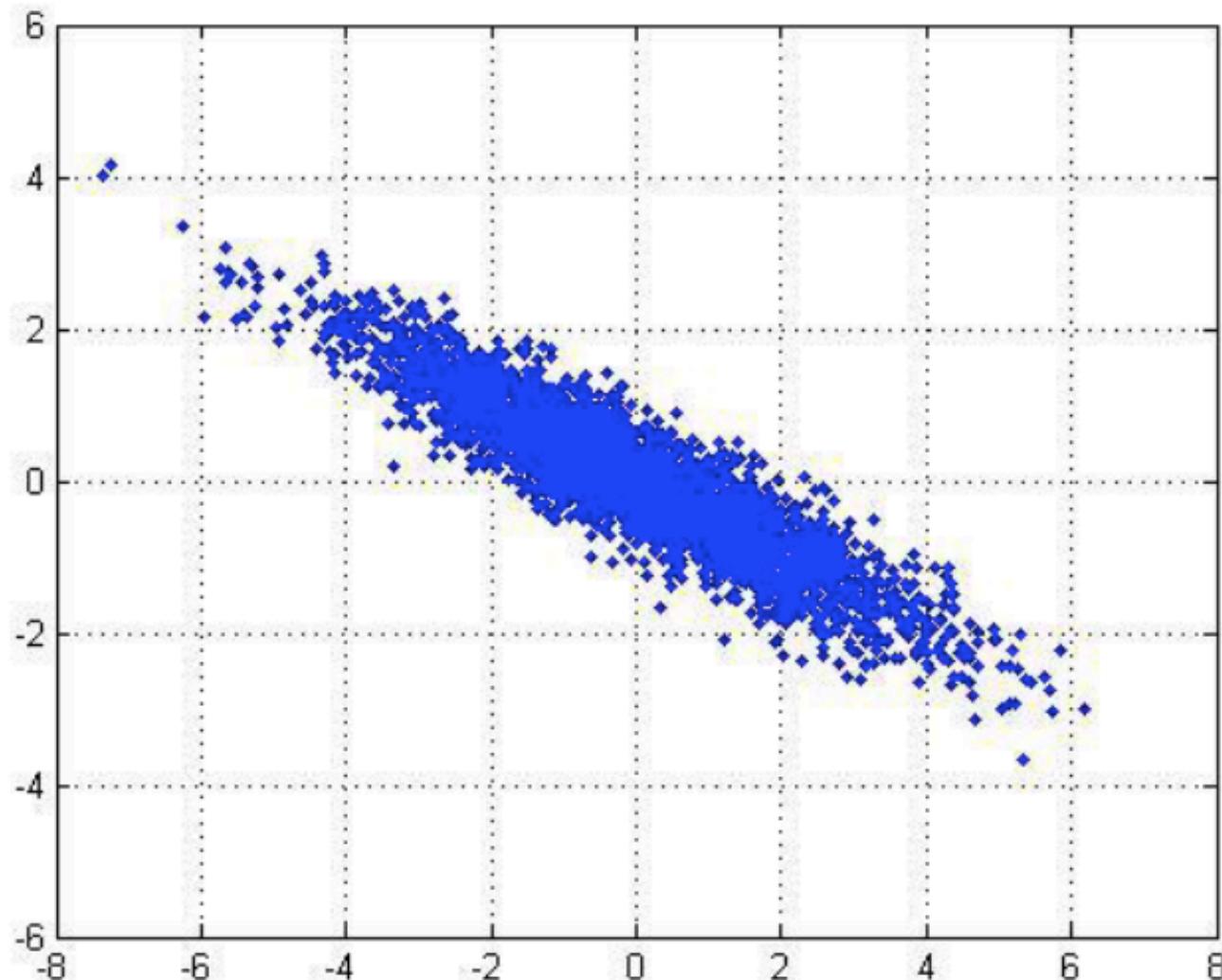
Principal Component Analysis

- **PCA Vectors** originate from the center of mass.
- Principal component #1: points in the direction of the **largest variance**.

Principal Component Analysis

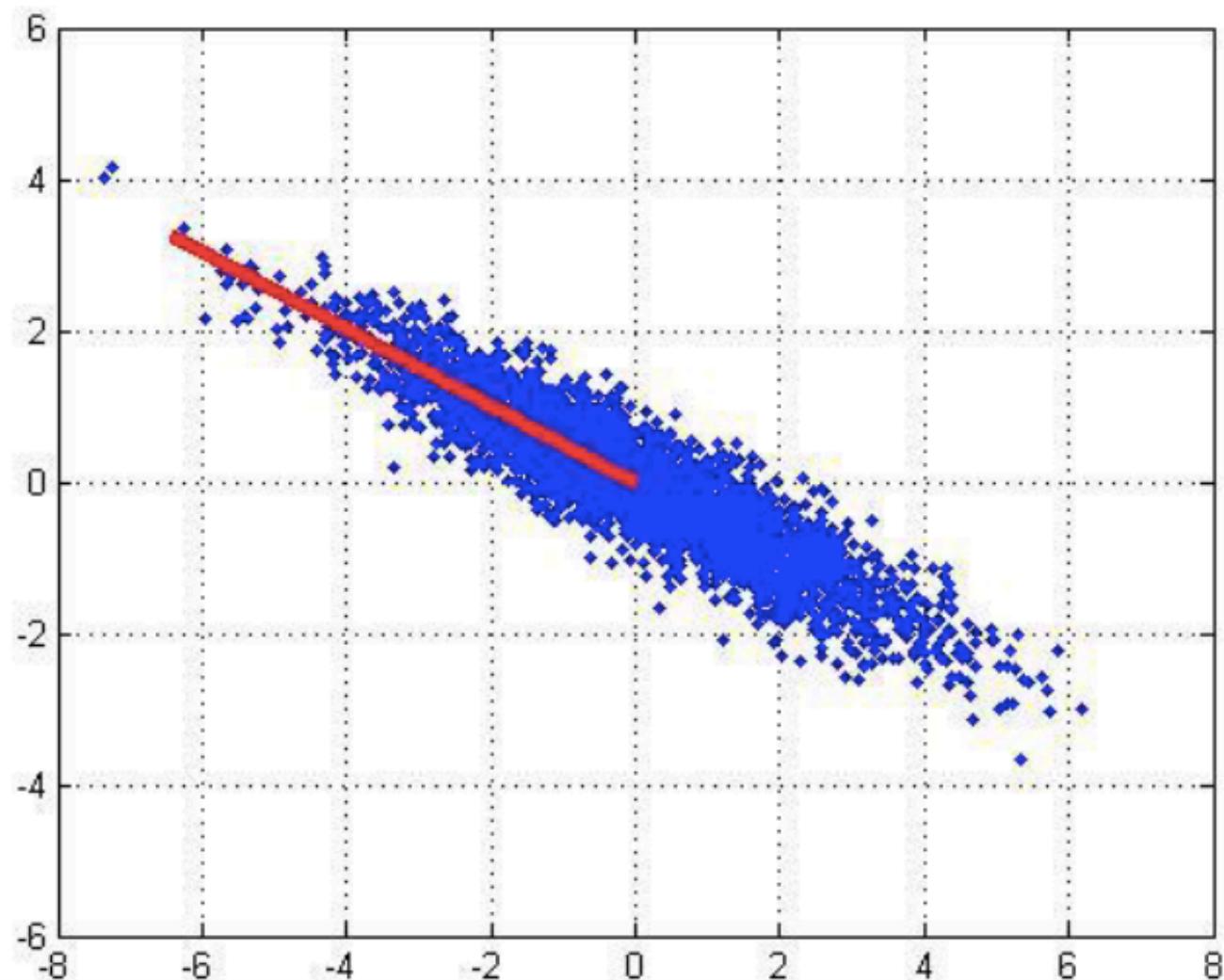
- **PCA Vectors** originate from the center of mass.
- Principal component #1: points in the direction of the **largest variance**.
- Each subsequent principal component
 - is **orthogonal** to the previous ones, and
 - points in the directions of the **largest variance of the residual subspace**

2D Gaussian dataset



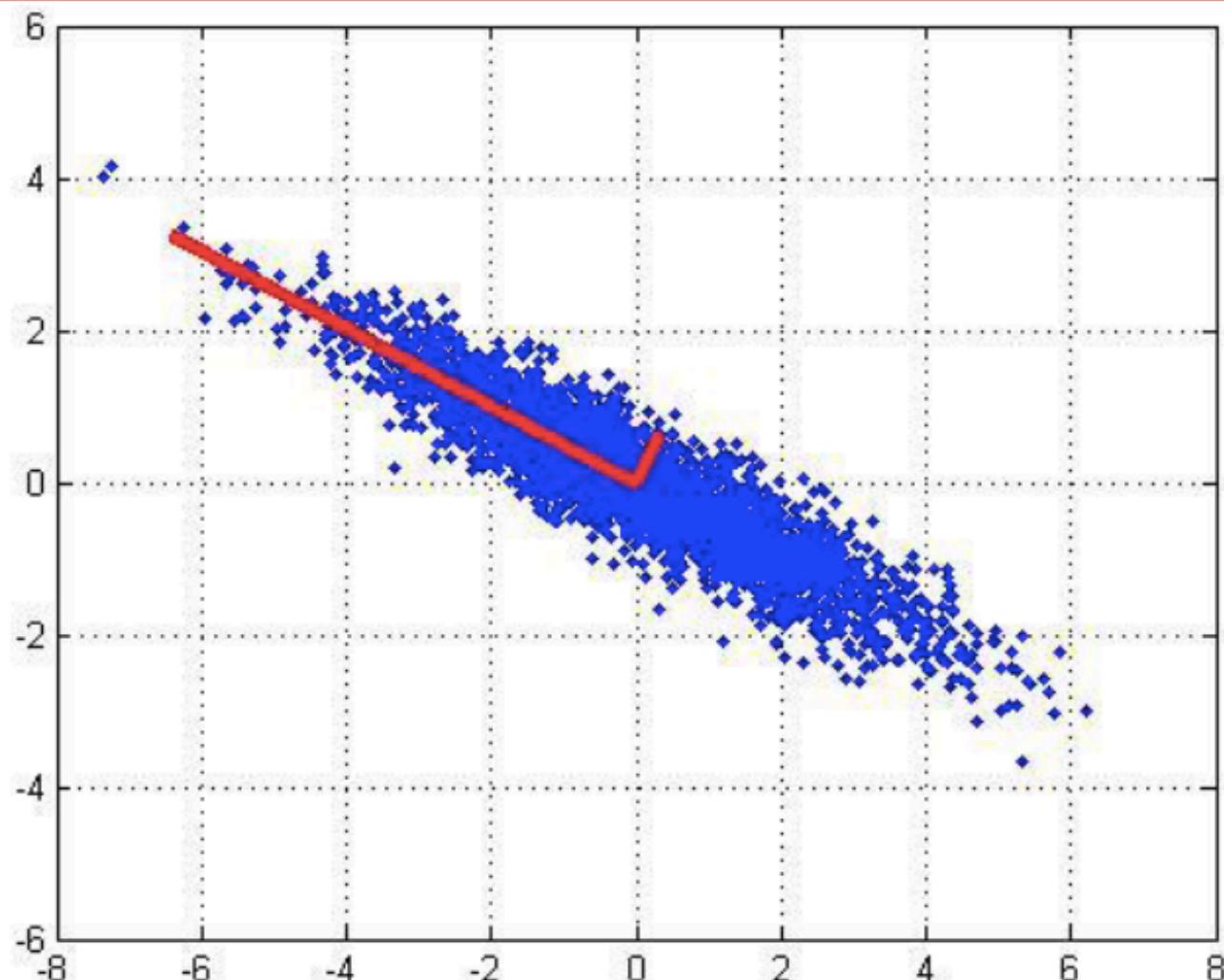
15

1st PCA axis



16

2nd PCA axis



17

PCA algorithm I (sequential)

Given the **centered** data $\{\mathbf{x}_1, \dots, \mathbf{x}_m\}$, compute the principal vectors:

PCA algorithm I (sequential)

Given the **centered** data $\{\mathbf{x}_1, \dots, \mathbf{x}_m\}$, compute the principal vectors:

$$\mathbf{w}_1 = \arg \max_{\|\mathbf{w}\|=1} \frac{1}{m} \sum_{i=1}^m \{(\mathbf{w}^T \mathbf{x}_i)^2\} \quad \text{1st PCA vector}$$

PCA algorithm I (sequential)

Given the **centered** data $\{\mathbf{x}_1, \dots, \mathbf{x}_m\}$, compute the principal vectors:

$$\mathbf{w}_1 = \arg \max_{\|\mathbf{w}\|=1} \frac{1}{m} \sum_{i=1}^m \{(\mathbf{w}^T \mathbf{x}_i)^2\} \quad \text{1st PCA vector}$$

We maximize the variance of projection of \mathbf{x}

$$\mathbf{w}_2 = \arg \max_{\|\mathbf{w}\|=1} \frac{1}{m} \sum_{i=1}^m \{[\mathbf{w}^T (\mathbf{x}_i - \underbrace{\mathbf{w}_1 \mathbf{w}_1^T \mathbf{x}_i}_{\mathbf{x}' \text{ PCA reconstruction}})]^2\} \quad k^{\text{th}} \text{ PCA vector}$$

PCA algorithm I (sequential)

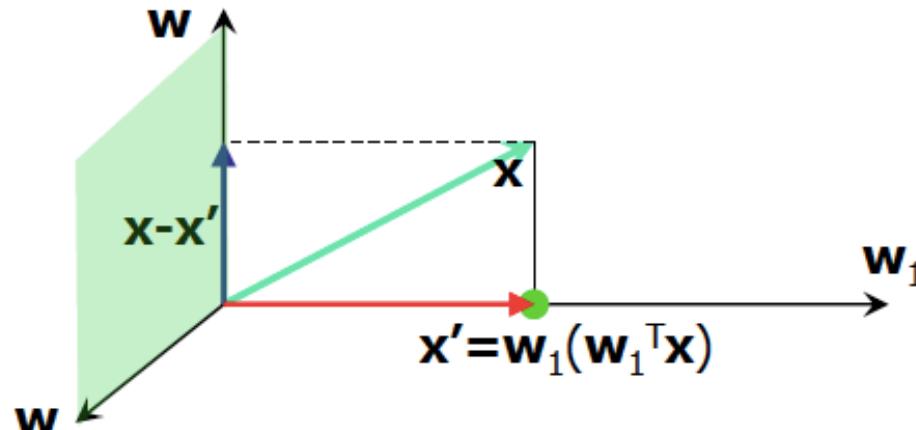
Given the **centered** data $\{\mathbf{x}_1, \dots, \mathbf{x}_m\}$, compute the principal vectors:

$$\mathbf{w}_1 = \arg \max_{\|\mathbf{w}\|=1} \frac{1}{m} \sum_{i=1}^m \{(\mathbf{w}^T \mathbf{x}_i)^2\} \quad \text{1st PCA vector}$$

We maximize the variance of projection of \mathbf{x}

$$\mathbf{w}_2 = \arg \max_{\|\mathbf{w}\|=1} \frac{1}{m} \sum_{i=1}^m \{[\mathbf{w}^T (\mathbf{x}_i - \underbrace{\mathbf{w}_1 \mathbf{w}_1^T \mathbf{x}_i}_{\mathbf{x}' \text{ PCA reconstruction}})]^2\} \quad k^{\text{th}} \text{ PCA vector}$$

We maximize the variance
of the projection in the
residual subspace



PCA algorithm I (sequential)

Given $\mathbf{w}_1, \dots, \mathbf{w}_{k-1}$, we calculate \mathbf{w}_k principal vector as before:

Maximize the variance of projection of \mathbf{x}

$$\mathbf{w}_k = \arg \max_{\|\mathbf{w}\|=1} \frac{1}{m} \sum_{i=1}^m \left\{ \underbrace{\left[\mathbf{w}^T (\mathbf{x}_i - \sum_{j=1}^{k-1} \mathbf{w}_j \mathbf{w}_j^T \mathbf{x}_i) \right]^2} \right\}$$

k^{th} PCA vector
 \mathbf{x}' PCA reconstruction

PCA algorithm I (sequential)

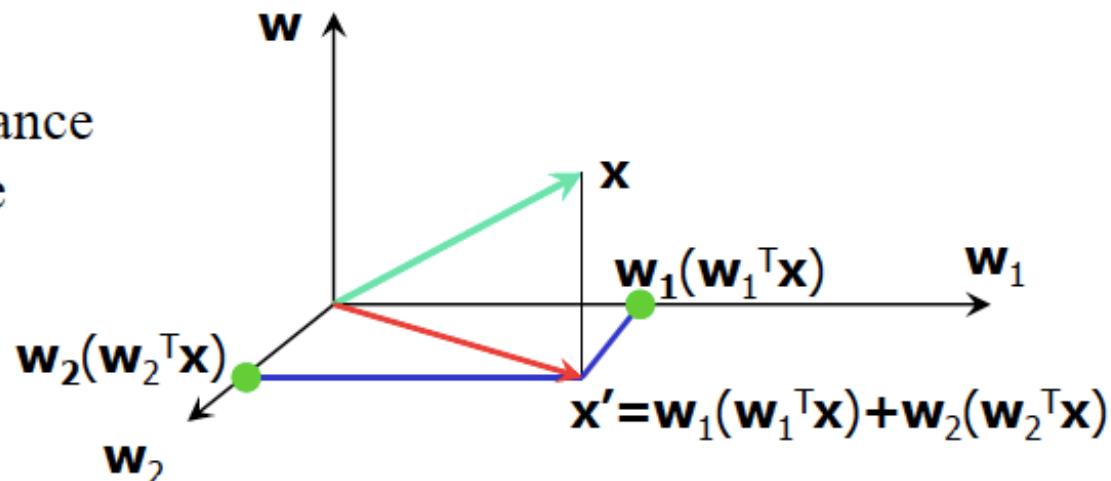
Given $\mathbf{w}_1, \dots, \mathbf{w}_{k-1}$, we calculate \mathbf{w}_k principal vector as before:

Maximize the variance of projection of \mathbf{x}

$$\mathbf{w}_k = \arg \max_{\|\mathbf{w}\|=1} \frac{1}{m} \sum_{i=1}^m \left\{ \underbrace{\left[\mathbf{w}^T (\mathbf{x}_i - \sum_{j=1}^{k-1} \mathbf{w}_j \mathbf{w}_j^T \mathbf{x}_i) \right]^2} \right\}$$

k^{th} PCA vector
 \mathbf{x}' PCA reconstruction

We maximize the variance of the projection in the residual subspace



PCA algorithm II (sample covariance matrix)

- Given data $\{\mathbf{x}_1, \dots, \mathbf{x}_m\}$, compute covariance matrix Σ

PCA algorithm II (sample covariance matrix)

- Given data $\{\mathbf{x}_1, \dots, \mathbf{x}_m\}$, compute covariance matrix Σ

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$

where

$$\bar{\mathbf{x}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$$

PCA algorithm II (sample covariance matrix)

- Given data $\{\mathbf{x}_1, \dots, \mathbf{x}_m\}$, compute covariance matrix Σ

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$

where

$$\bar{\mathbf{x}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$$

- PCA** basis vectors = the eigenvectors of Σ
- Larger eigenvalue \Rightarrow more important eigenvectors

PCA algorithm II (sample covariance matrix)

PCA algorithm(\mathbf{X} , k): top k eigenvalues/eigenvectors

```
%  $\mathbf{X}$  = N × m data matrix,  
% ... each data point  $\mathbf{x}_i$  = column vector,  $i=1..m$ 
```

PCA algorithm II (sample covariance matrix)

PCA algorithm(\mathbf{X} , k): top k eigenvalues/eigenvectors

- % \mathbf{X} = $N \times m$ data matrix,
- % ... each data point \mathbf{x}_i = column vector, $i=1..m$
- $\mathbf{\bar{x}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$

PCA algorithm II (sample covariance matrix)

PCA algorithm(\mathbf{X} , k): top k eigenvalues/eigenvectors

- % $\mathbf{X} = N \times m$ data matrix,
- % ... each data point $\mathbf{x}_i =$ column vector, $i=1..m$
- $\underline{\mathbf{x}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$
- $\mathbf{X} \leftarrow$ subtract mean $\underline{\mathbf{x}}$ from each column vector \mathbf{x}_i in \mathbf{X}

PCA algorithm II (sample covariance matrix)

PCA algorithm(\mathbf{X} , k): top k eigenvalues/eigenvectors

- % $\mathbf{X} = N \times m$ data matrix,
- % ... each data point $\mathbf{x}_i =$ column vector, $i=1..m$
- $\underline{\mathbf{x}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$
- $\mathbf{X} \leftarrow$ subtract mean $\underline{\mathbf{x}}$ from each column vector \mathbf{x}_i in \mathbf{X}
- $\Sigma \leftarrow \mathbf{X}\mathbf{X}^T$... covariance matrix of \mathbf{X}
- $\{ \lambda_i, \mathbf{u}_i \}_{i=1..N} =$ eigenvectors/eigenvalues of Σ
... $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_N$

PCA algorithm II (sample covariance matrix)

PCA algorithm(\mathbf{X} , k): top k eigenvalues/eigenvectors

- % $\mathbf{X} = N \times m$ data matrix,
- % ... each data point \mathbf{x}_i = column vector, $i=1..m$
- $\underline{\mathbf{x}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$
- $\mathbf{X} \leftarrow$ subtract mean $\underline{\mathbf{x}}$ from each column vector \mathbf{x}_i in \mathbf{X}
- $\Sigma \leftarrow \mathbf{X}\mathbf{X}^T$... covariance matrix of \mathbf{X}
- $\{ \lambda_i, \mathbf{u}_i \}_{i=1..N} =$ eigenvectors/eigenvalues of Σ
... $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_N$
- Return $\{ \lambda_i, \mathbf{u}_i \}_{i=1..k}$
% top k PCA components

22

PCA algorithm III (SVD of the data matrix)

Singular Value Decomposition of the **centered** data matrix \mathbf{X} .

PCA algorithm III (SVD of the data matrix)

Singular Value Decomposition of the **centered** data matrix \mathbf{X} .

$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_m] \in \mathbb{R}^{N \times m}, \quad \begin{array}{l} m: \text{number of instances,} \\ N: \text{dimension} \end{array}$

PCA algorithm III (SVD of the data matrix)

Singular Value Decomposition of the **centered** data matrix \mathbf{X} .

$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_m] \in \mathbb{R}^{N \times m}, \quad m: \text{number of instances},$
 $N: \text{dimension}$

$$\mathbf{X}_{\text{features} \times \text{samples}} = \mathbf{U} \mathbf{S} \mathbf{V}^T$$

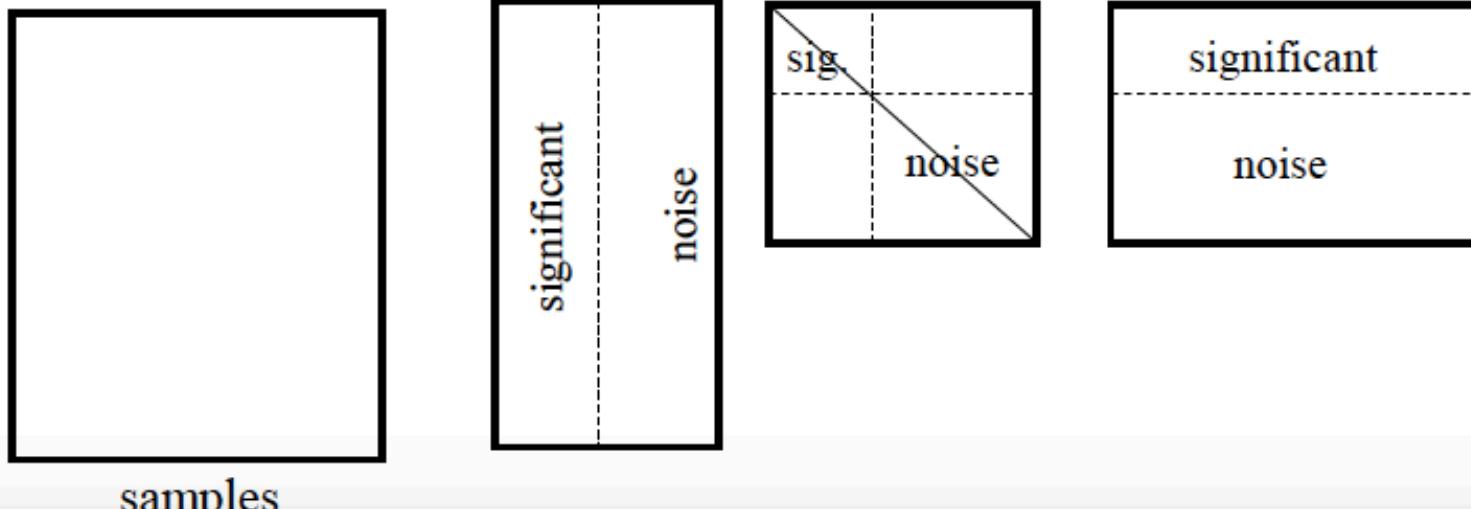
PCA algorithm III (SVD of the data matrix)

Singular Value Decomposition of the **centered** data matrix \mathbf{X} .

$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_m] \in \mathbb{R}^{N \times m}, \quad m: \text{number of instances},$
 $N: \text{dimension}$

$$\mathbf{X}_{\text{features} \times \text{samples}} = \mathbf{U} \mathbf{S} \mathbf{V}^T$$

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



PCA algorithm III

- **Columns of \mathbf{U}**

- the principal vectors, $\{ \mathbf{u}^{(1)}, \dots, \mathbf{u}^{(k)} \}$
- orthogonal and has unit norm – so $\mathbf{U}^T \mathbf{U} = \mathbf{I}$
- Can reconstruct the data using linear combinations of $\{ \mathbf{u}^{(1)}, \dots, \mathbf{u}^{(k)} \}$

PCA algorithm III

- **Columns of \mathbf{U}**

- the principal vectors, $\{ \mathbf{u}^{(1)}, \dots, \mathbf{u}^{(k)} \}$
- orthogonal and has unit norm – so $\mathbf{U}^T \mathbf{U} = \mathbf{I}$
- Can reconstruct the data using linear combinations of $\{ \mathbf{u}^{(1)}, \dots, \mathbf{u}^{(k)} \}$

- **Matrix \mathbf{S}**

- Diagonal
- Shows importance of each eigenvector

- **Columns of \mathbf{V}^T**

- The coefficients for reconstructing the samples

5 minutes break

Applications

Face Recognition

- Want to identify specific person, based on facial image
 - Robust to glasses, lighting,...
- ⇒ Can't just use the given 256 x 256 pixels



26

Applying PCA: Eigenfaces

Method A: Build a PCA subspace for each person and check which subspace can reconstruct the test image the best

Applying PCA: Eigenfaces

Method A: Build a PCA subspace for each person and check which subspace can reconstruct the test image the best

Method B: Build one PCA database for the whole dataset and then classify based on the weights.

Applying PCA: Eigenfaces

Method A: Build a PCA subspace for each person and check which subspace can reconstruct the test image the best

Method B: Build one PCA database for the whole dataset and then classify based on the weights.

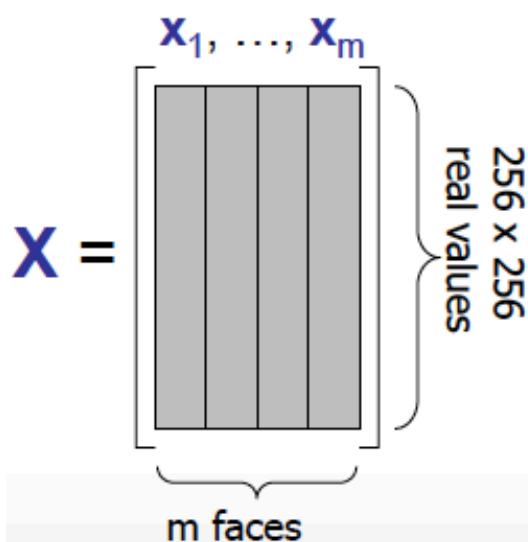


- Example data set: Images of faces
 - Famous Eigenface approach [Turk & Pentland], [Sirovich & Kirby]

Applying PCA: Eigenfaces

Method A: Build a PCA subspace for each person and check which subspace can reconstruct the test image the best

Method B: Build one PCA database for the whole dataset and then classify based on the weights.



- Example data set: Images of faces
 - Famous Eigenface approach [Turk & Pentland], [Sirovich & Kirby]
- Each face \mathbf{x} is ...
 - 256×256 values (luminance at location)
 - \mathbf{x} in $\mathbb{R}^{256 \times 256}$ (view as 64K dim vector)
- Form $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_m]$ **centered** data mtx
- Compute $\Sigma = \mathbf{X}\mathbf{X}^T$
- Problem: Σ is $64K \times 64K \dots$ **HUGE!!!**

Computational Complexity

- Suppose m instances, each of size N
 - Eigenfaces: $m=500$ faces, each of size $N=64K$

Computational Complexity

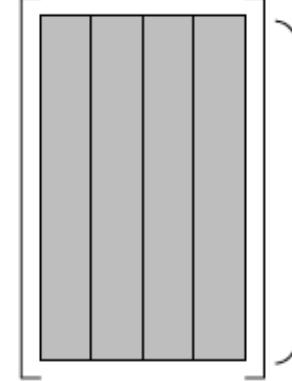
- Suppose m instances, each of size N
 - Eigenfaces: $m=500$ faces, each of size $N=64K$
- Given $N \times N$ covariance matrix Σ , can compute
 - all N eigenvectors/eigenvalues in $O(N^3)$
 - first k eigenvectors/eigenvalues in $O(k N^2)$
- But if $N=64K$, EXPENSIVE!

A Clever Workaround

- Note that $m << 64K$
- Use $\mathbf{L} = \mathbf{X}^T \mathbf{X}$ instead of $\Sigma = \mathbf{X} \mathbf{X}^T$
- If \mathbf{v} is eigenvector of \mathbf{L}
then \mathbf{Xv} is eigenvector of Σ

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1, \dots, \mathbf{x}_m \end{bmatrix}$$

256×256
 $\{\}$ real values
 $\{\}$ m faces



A Clever Workaround

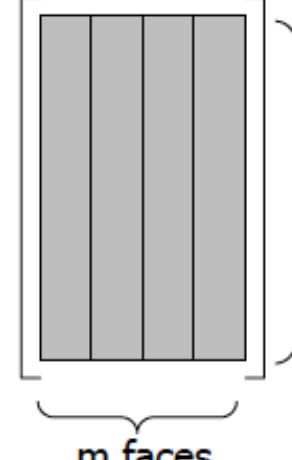
- Note that $m \ll 64K$
- Use $\mathbf{L} = \mathbf{X}^T \mathbf{X}$ instead of $\Sigma = \mathbf{X} \mathbf{X}^T$
- If \mathbf{v} is eigenvector of \mathbf{L}
then \mathbf{Xv} is eigenvector of Σ

Proof: $\mathbf{L} \mathbf{v} = \gamma \mathbf{v}$

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1, \dots, \mathbf{x}_m \end{bmatrix}$$

256 x 256
real values

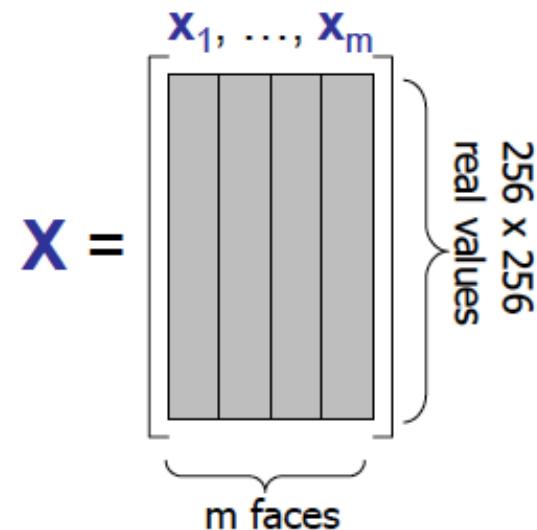
$\brace{m \text{ faces}}$



A Clever Workaround

- Note that $m \ll 64K$
- Use $\mathbf{L} = \mathbf{X}^T \mathbf{X}$ instead of $\Sigma = \mathbf{X} \mathbf{X}^T$
- If \mathbf{v} is eigenvector of \mathbf{L}
then \mathbf{Xv} is eigenvector of Σ

Proof: $\mathbf{L} \mathbf{v} = \gamma \mathbf{v}$
 $\mathbf{X}^T \mathbf{X} \mathbf{v} = \gamma \mathbf{v}$



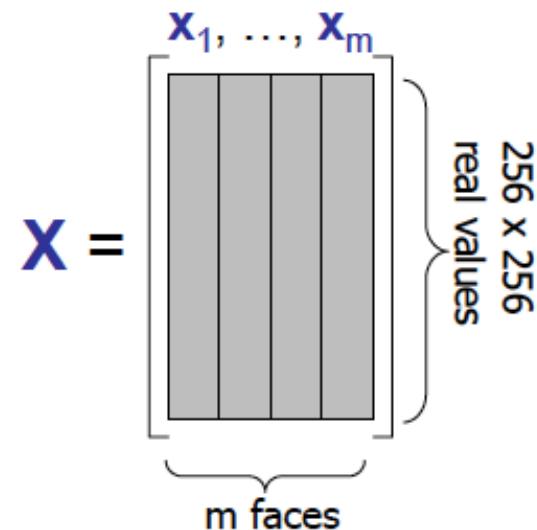
A Clever Workaround

- Note that $m << 64K$
- Use $\mathbf{L} = \mathbf{X}^T \mathbf{X}$ instead of $\Sigma = \mathbf{X} \mathbf{X}^T$
- If \mathbf{v} is eigenvector of \mathbf{L}
then \mathbf{Xv} is eigenvector of Σ

Proof: $\mathbf{L} \mathbf{v} = \gamma \mathbf{v}$

$$\mathbf{X}^T \mathbf{X} \mathbf{v} = \gamma \mathbf{v}$$

$$\mathbf{X} (\mathbf{X}^T \mathbf{X} \mathbf{v}) = \mathbf{X}(\gamma \mathbf{v}) = \gamma \mathbf{Xv}$$



A Clever Workaround

- Note that $m << 64K$
- Use $\mathbf{L} = \mathbf{X}^T \mathbf{X}$ instead of $\Sigma = \mathbf{X} \mathbf{X}^T$
- If \mathbf{v} is eigenvector of \mathbf{L}
then \mathbf{Xv} is eigenvector of Σ

Proof: $\mathbf{L} \mathbf{v} = \gamma \mathbf{v}$

$$\mathbf{X}^T \mathbf{X} \mathbf{v} = \gamma \mathbf{v}$$

$$\mathbf{X} (\mathbf{X}^T \mathbf{X} \mathbf{v}) = \mathbf{X}(\gamma \mathbf{v}) = \gamma \mathbf{Xv}$$

$$(\mathbf{X} \mathbf{X}^T) \mathbf{X} \mathbf{v} = \gamma (\mathbf{Xv})$$

$$\Sigma (\mathbf{Xv}) = \gamma (\mathbf{Xv})$$

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1, \dots, \mathbf{x}_m \end{bmatrix}$$

256×256
 m faces
real values

Principle Components (Method B)



30

Reconstructing... (Method B)



- ... faster if train with...
 - only people w/out glasses
 - same lighting conditions

Shortcomings

- Requires carefully controlled data:
 - All faces centered in frame
 - Same size
 - Some sensitivity to angle

Shortcomings

- Requires carefully controlled data:
 - All faces centered in frame
 - Same size
 - Some sensitivity to angle
- Method is completely knowledge free
 - (sometimes this is good!)
 - Doesn't know that faces are wrapped around 3D objects (heads)
 - Makes no effort to preserve class distinctions

Happiness subspace (method A)



33

Disgust subspace (method A)



34

Image Compression

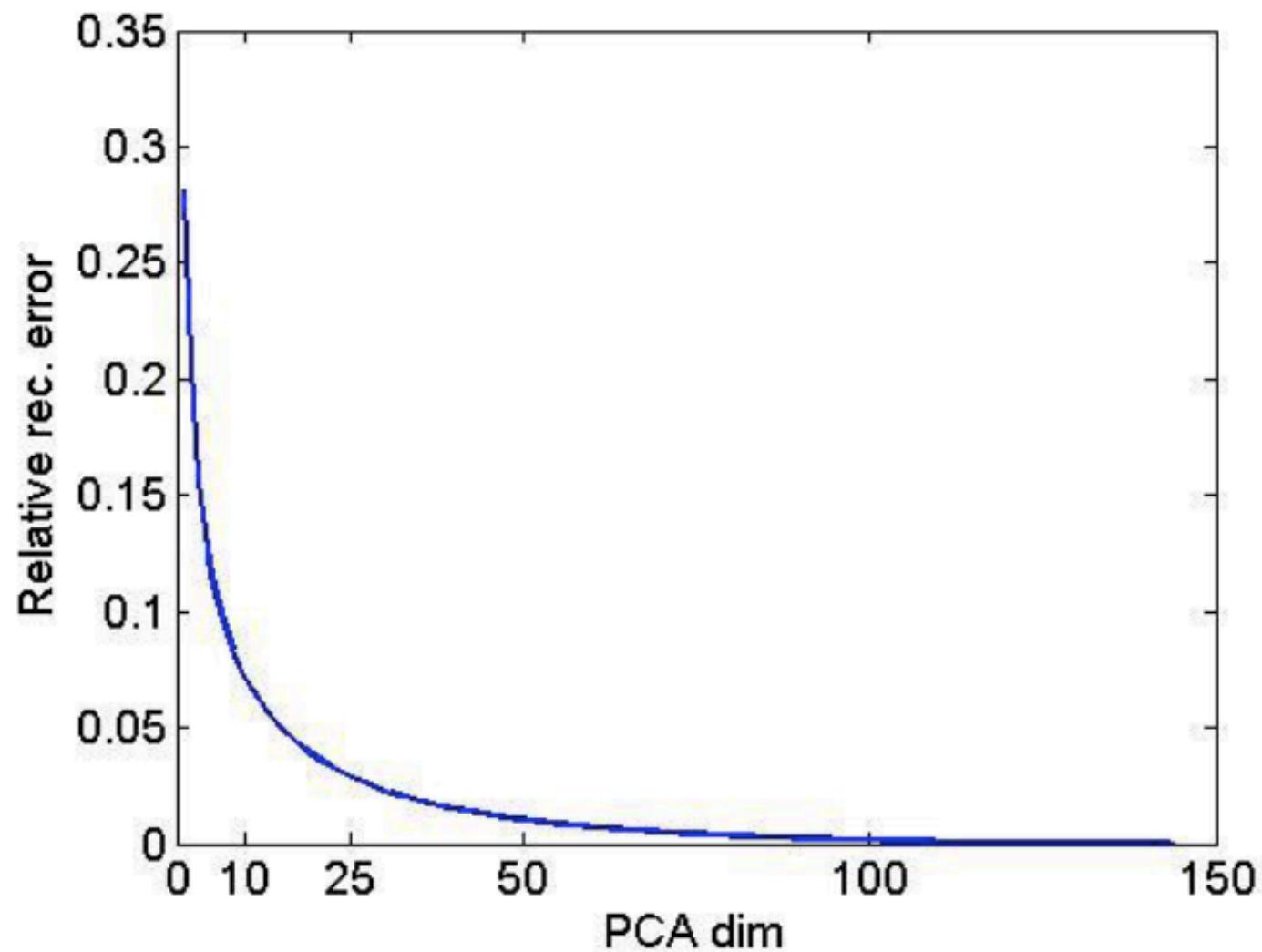
Original Image



- Divide the original 372x492 image into patches:
 - Each patch is an instance that contains 12x12 pixels on a grid
 - View each as a 144-D vector

39

L_2 error and PCA dim



40

PCA compression: 144D \Rightarrow 60D



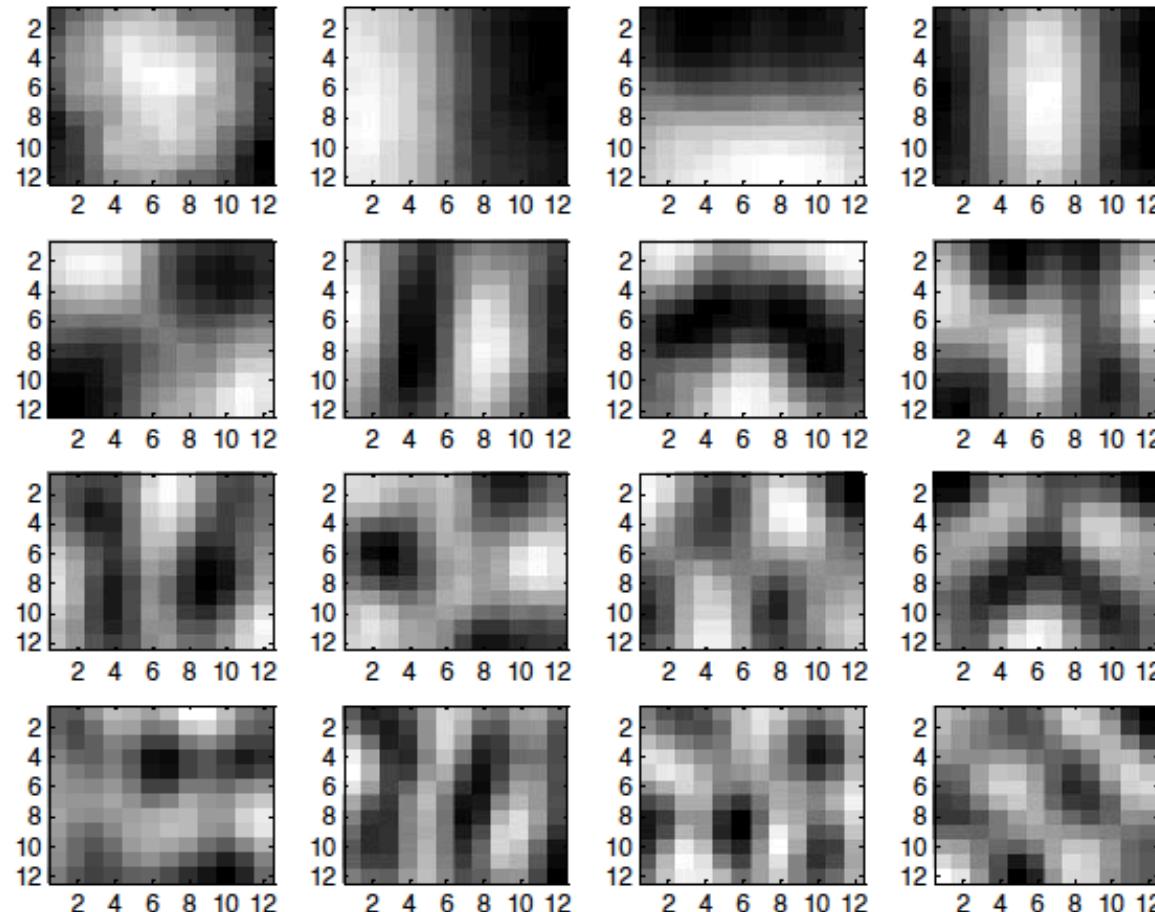
41

PCA compression: 144D \Rightarrow 16D



42

16 most important eigenvectors

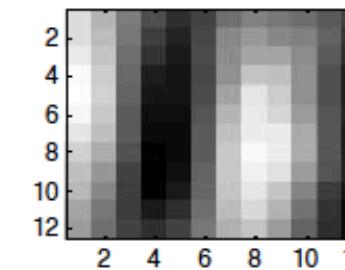
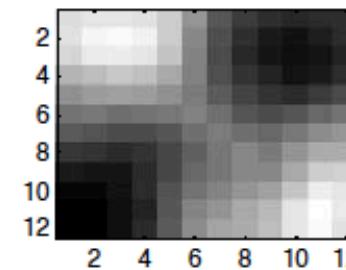
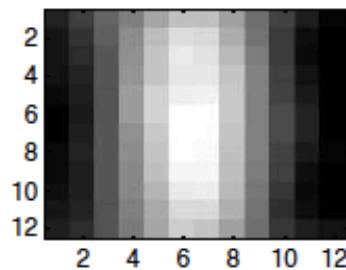
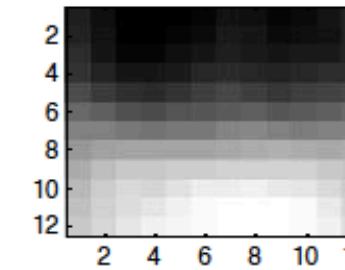
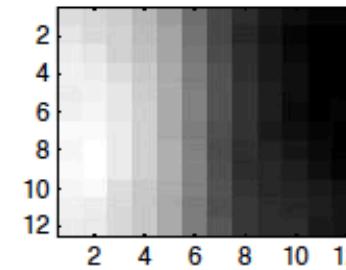
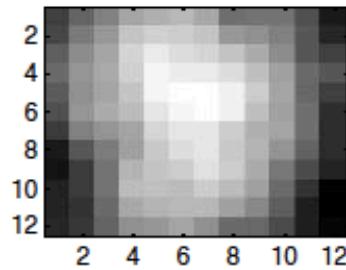


PCA compression: 144D \Rightarrow 6D



44

6 most important eigenvectors

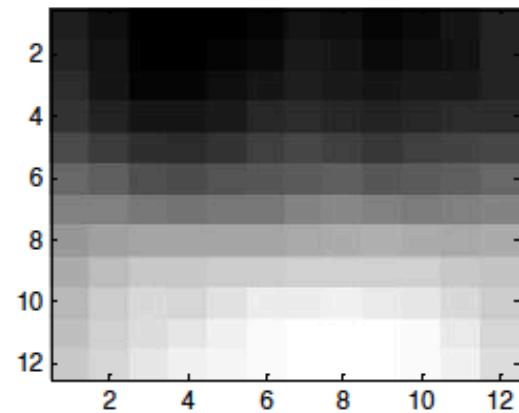
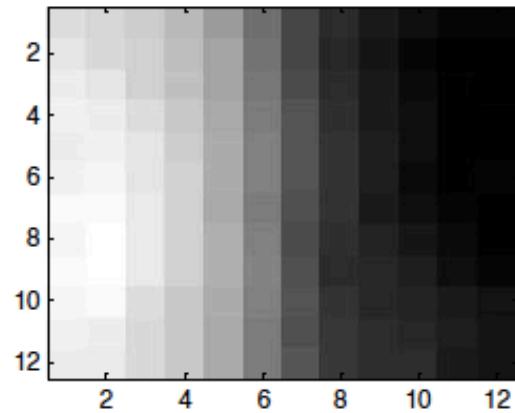
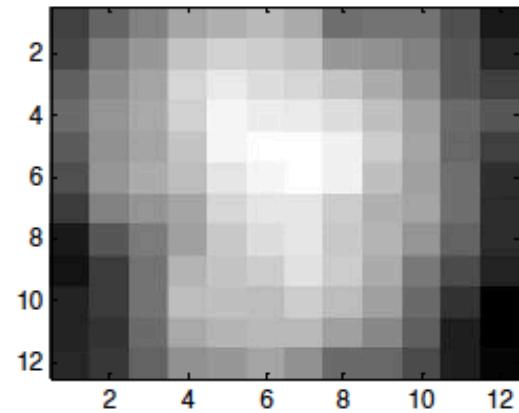


PCA compression: 144D \Rightarrow 3D

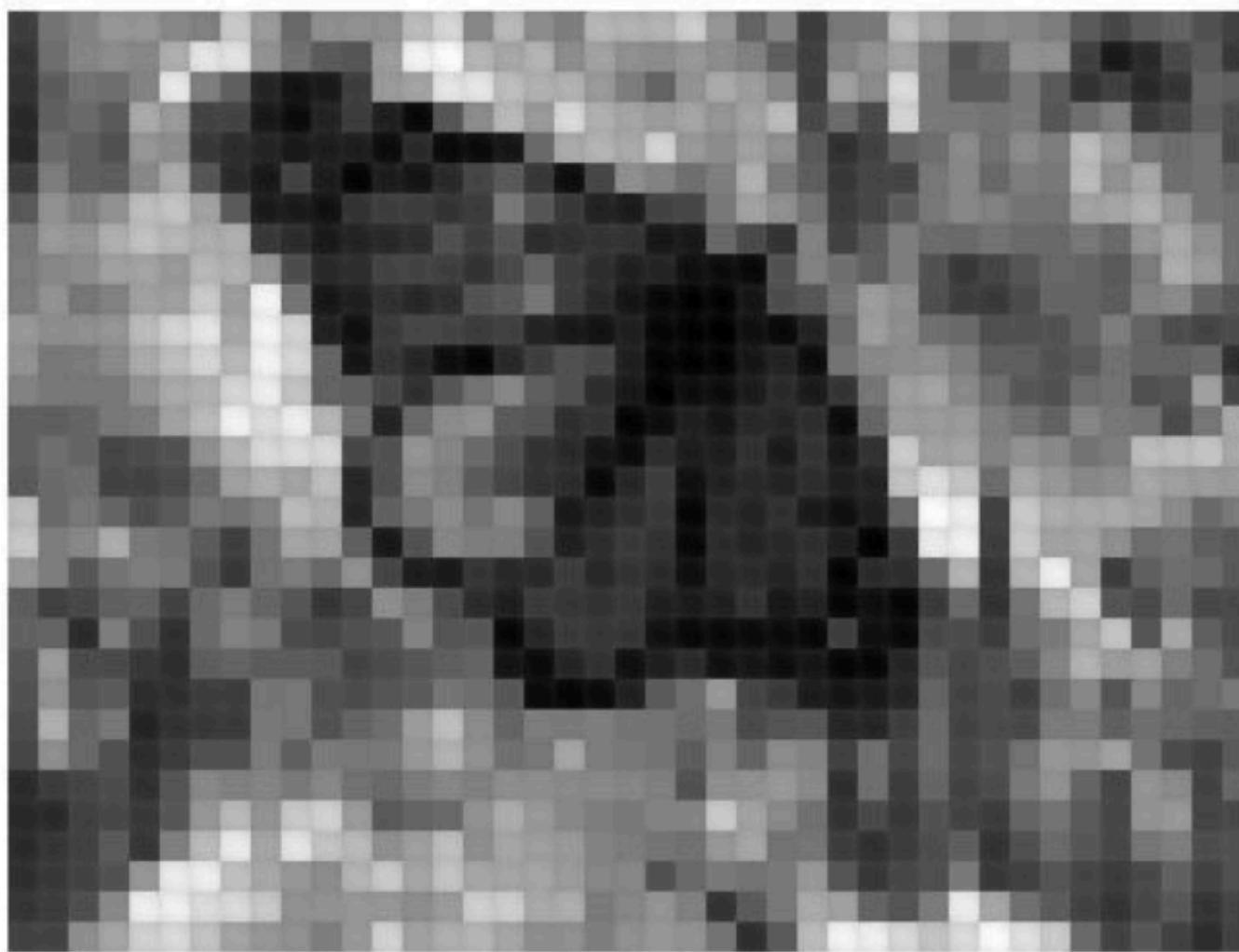


46

3 most important eigenvectors

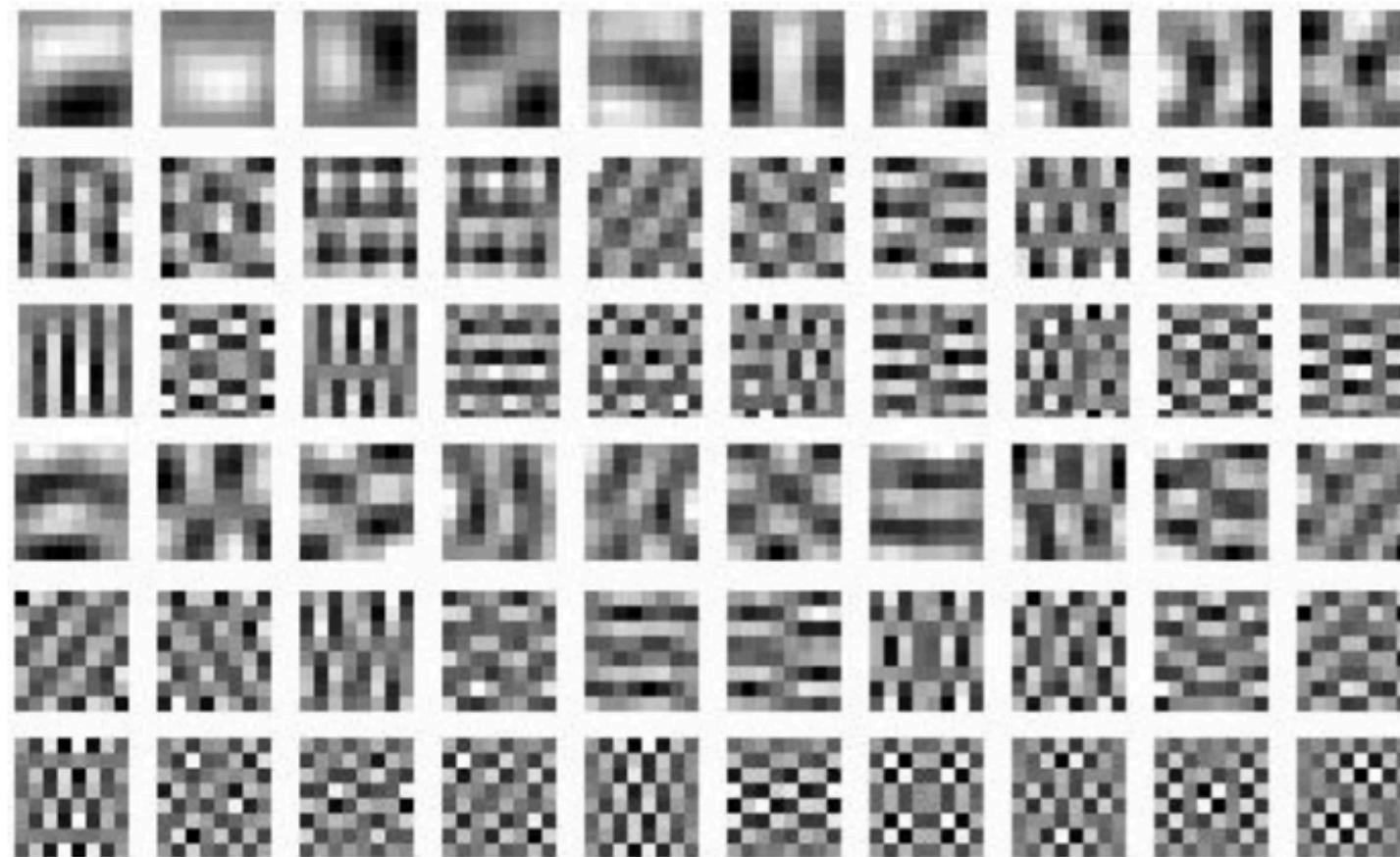


PCA compression: 144D \Rightarrow 1D



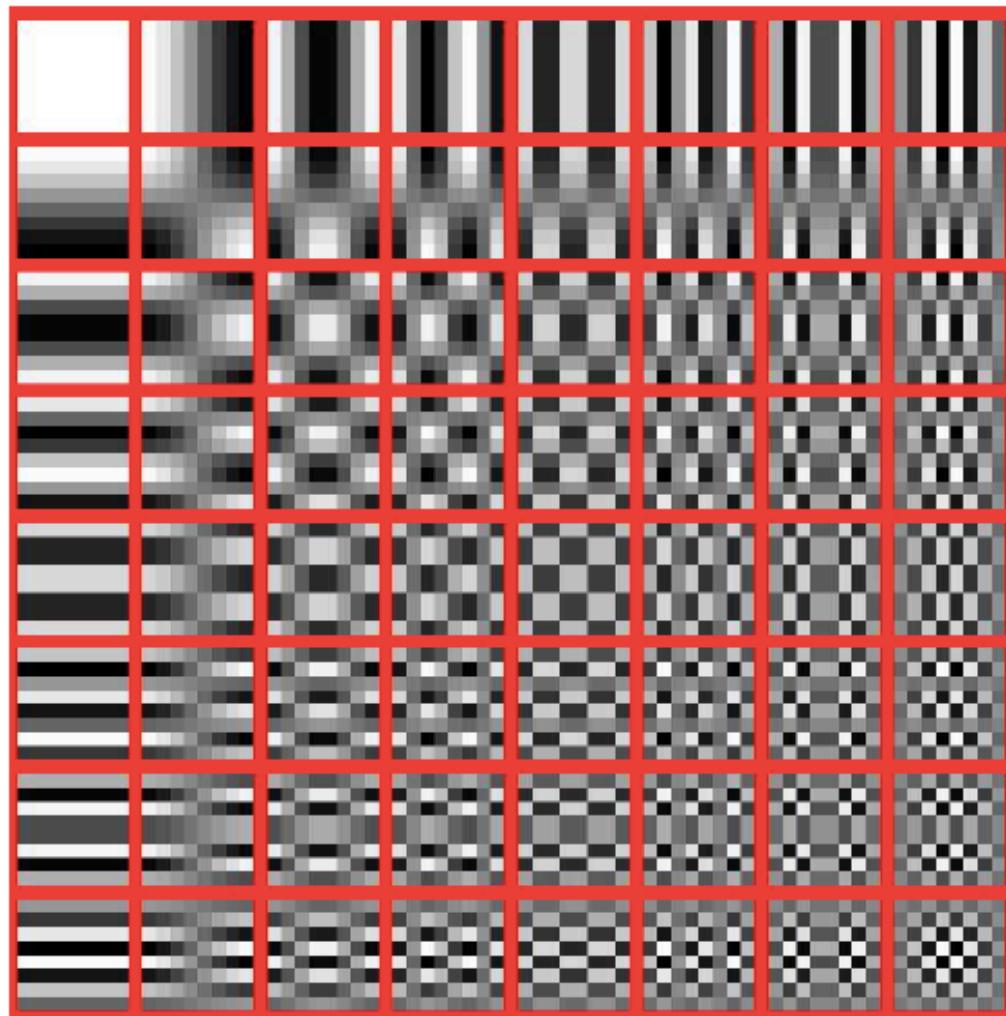
48

60 most important eigenvectors



Looks like the discrete cosine bases of JPG!...

2D Discrete Cosine Basis

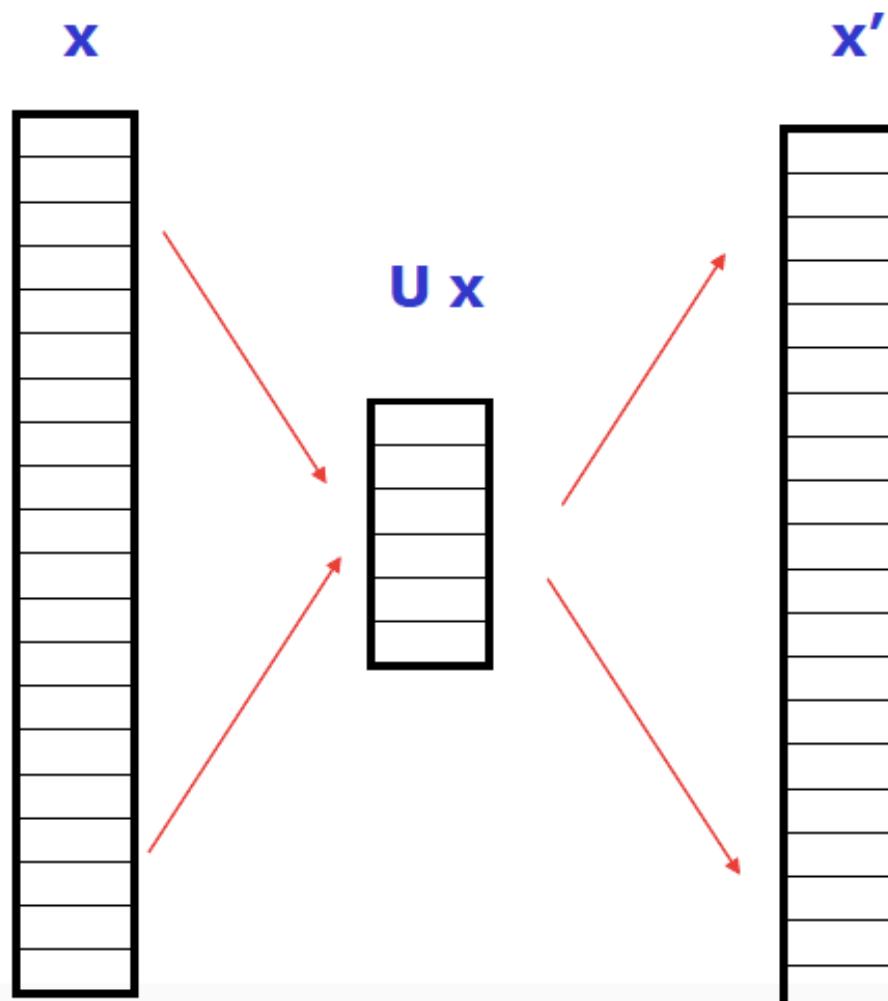


http://en.wikipedia.org/wiki/Discrete_cosine_transform

50

Noise Filtering

Noise Filtering



52

Noisy image

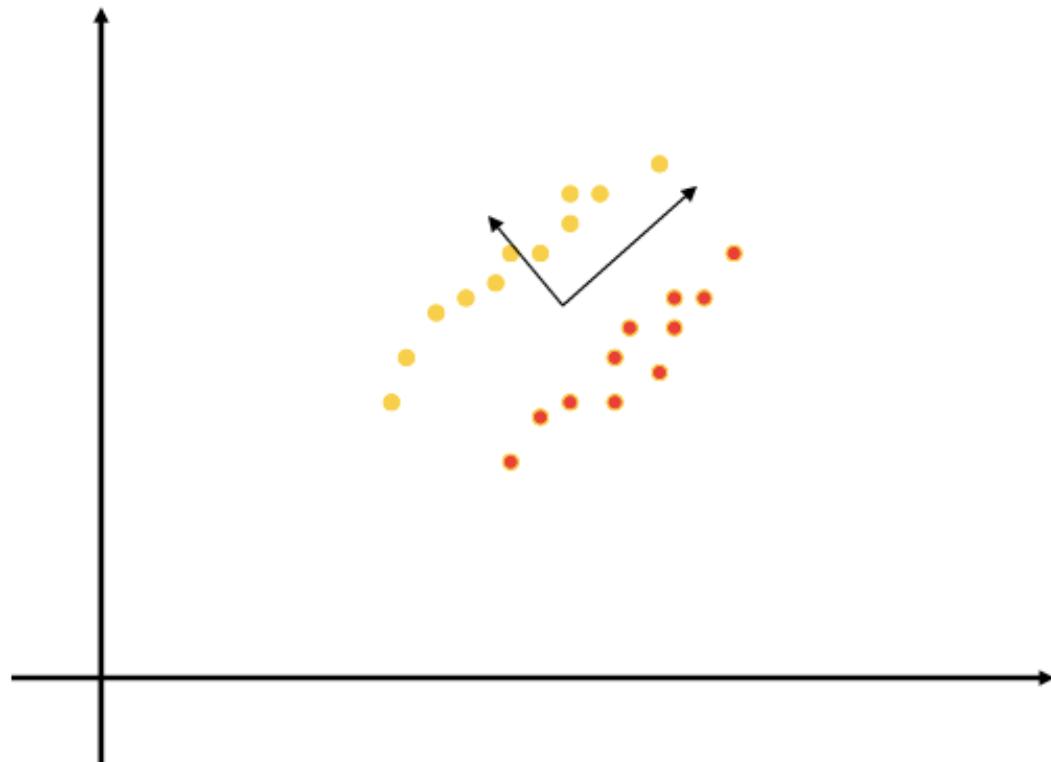


Denoised image using 15 PCA components



PCA Shortcomings

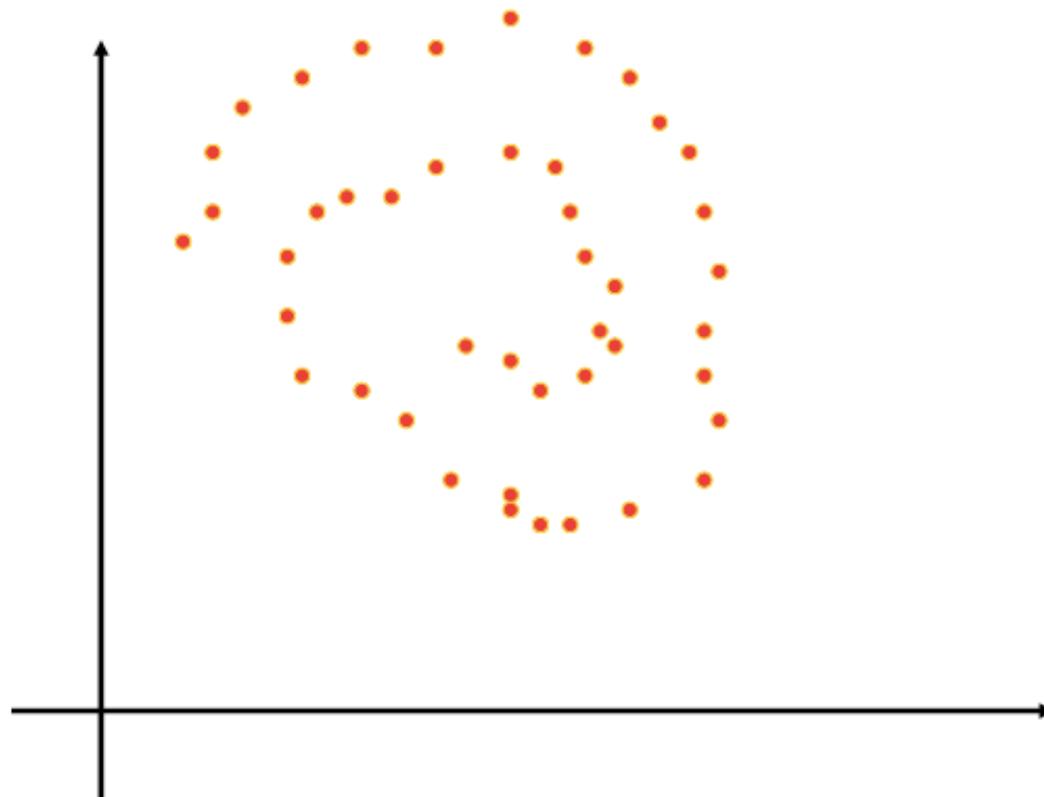
Problematic Data Set for PCA



PCA doesn't know labels!

56

Problematic Data Set for PCA



PCA cannot capture NON-LINEAR structure!

58

PCA Conclusions

❑ PCA

- finds orthonormal basis for data
- Sorts dimensions in order of “importance”
- Discard low significance dimensions

❑ Uses:

- Get compact description
- Ignore noise
- Improve classification (hopefully)

❑ Not magic:

- Doesn't know class labels
- Can only capture linear variations

❑ One of many tricks to reduce dimensionality!

That's all!