# Principal Component Analysis

STAT5241 Section 2

Statistical Machine Learning

Xiaofei Shi

# Data Visualization

- How can we visualize data?

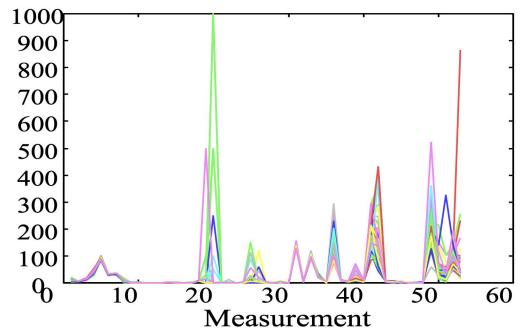- Example: 53 features for 65 people?

- 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…
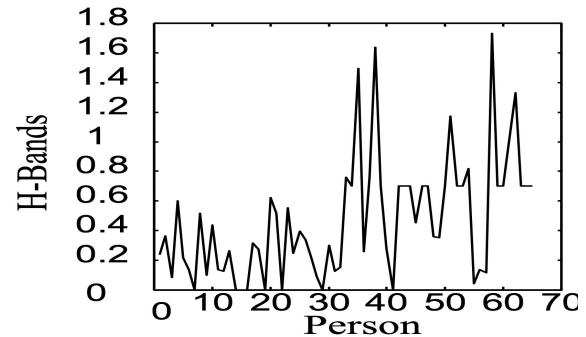
COLUMBIA UNIVERSITY

- • Curves (65 curves, one for each person)



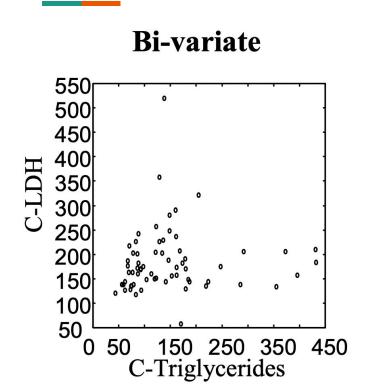Difficult to compare the different patients…
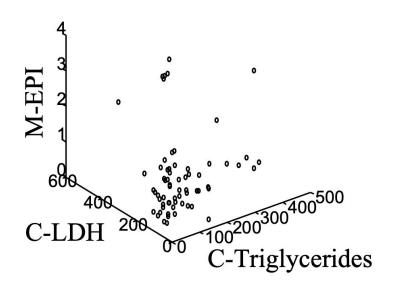
- Curves (53 pictures, one for each feature)



Difficult to see the correlations between the features...
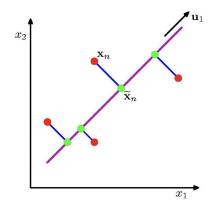
# Bi-variate



# Tri-variate

# Moreover...

- 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



**PCA:**

- Orthogonal projection of the data onto a lower-dimension linear space that <u>equivalently</u>...
    1. *minimizes* the mean squared distance between
        - data points (red points) and projections (green points)
        - i.e. sum of squares of blue line lengths
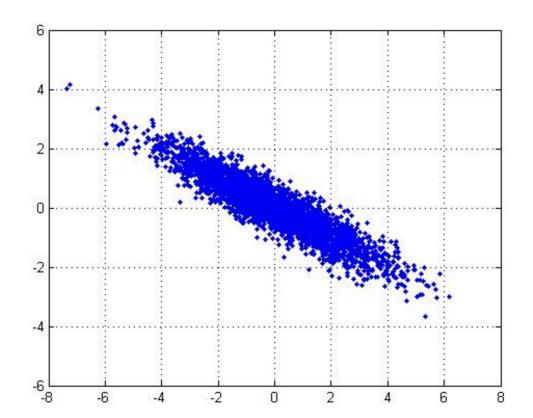    2. *maximizes* variance of projected data (green points)

# PCA

**Idea:**

❑ Given data points in a N-dimensional space, project them into a lower dimensional space while preserving as much information as possible.

- Find best planar approximation of 3D data
- Find best 12-D approximation of 10,000-D data

❑ In particular, choose **<u>linear</u>** projection that minimizes *squared error* in reconstructing the original data.
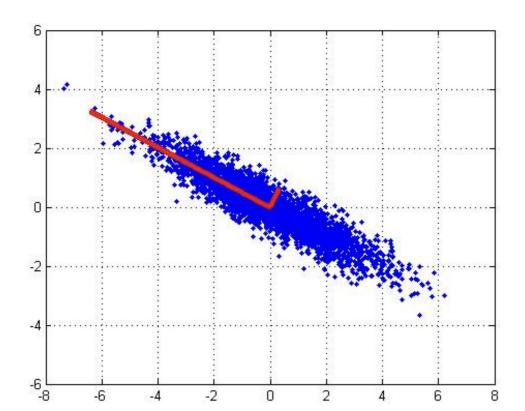
# Example:

# Example:

# Algorithm: sequential

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

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

To find $\mathbf{w_1}$, maximize the variance of projection of $\mathbf{x}$
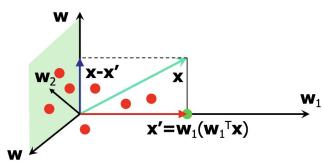
# Algorithm: sequential

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

1st PCA vector

To find $\mathbf{w}_1$, maximize the variance of projection of $\mathbf{x}$

$$\mathbf{w}_2 = \arg\max_{\|\mathbf{w}\|=1} \frac{1}{m} \sum_{i=1}^{m} \{[\mathbf{w}^T (\mathbf{x}_i - \underbrace{\mathbf{w}_1 \mathbf{w}_1^T \mathbf{x}_i})]^2\}$$

2nd PCA vector

$\mathbf{x}'$ projection onto w_1

To find $\mathbf{w}_2$, we maximize the **variance** of the projection in the **residual** subspace



w

$\mathbf{w}_2$    x-x'    x

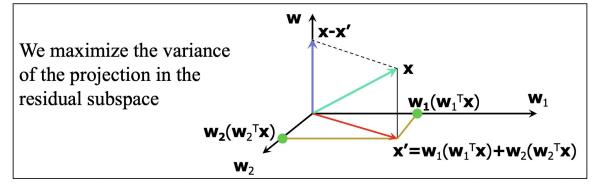$\mathbf{w}_1$

x'=$\mathbf{w}_1(\mathbf{w}_1^T\mathbf{x})$

w

# Algorithm: sequential

Given $\mathbf{w_1}, \ldots, \mathbf{w_{k-1}}$, we calculate $\mathbf{w_k}$ principal vector as before:

Maximize the variance of projection of $\mathbf{x}$

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

$k^{\text{th}}$ PCA vector

$\mathbf{x'}$ projection onto previous directions

We maximize the variance of the projection in the residual subspace



$\mathbf{w}$
$\mathbf{x\text{-}x'}$
$\mathbf{x}$
$\mathbf{w_1(w_1{}^Tx)}$
$\mathbf{w_1}$
$\mathbf{w_2(w_2{}^Tx)}$
$\mathbf{x'=w_1(w_1{}^Tx)+w_2(w_2{}^Tx)}$
$\mathbf{w_2}$

# Algorithm: sample covariance matrix

- Given data $\{\mathbf{x}_1, \ldots, \mathbf{x}_m\}$, compute covariance matrix $\Sigma$

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

- **PCA** basis vectors = the eigenvectors of $\Sigma$

- Larger eigenvalue $\Rightarrow$ more important eigenvectors

# Algorithm: sample covariance matrix

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

    % $\underline{\mathbf{X}}$ = N × m data matrix, <u>N is number of features</u>
    % … each data point $\mathbf{x}_i$ = column vector, i=1..m

- $\underline{\mathbf{x}} = \dfrac{1}{m} \displaystyle\sum_{i=1}^{m} \mathbf{x}_i$

- $\mathbf{X} \leftarrow$ subtract mean $\underline{\mathbf{x}}$ from each column vector $\mathbf{x}_i$ in $\underline{\mathbf{X}}$

- $\Sigma \leftarrow \mathbf{X}\mathbf{X}^T$ … covariance matrix of $\mathbf{X}$

- { $\lambda_i$, $\mathbf{u}_i$ }$_{i=1..N}$ = eigenvectors/eigenvalues of $\Sigma$
  where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_N$
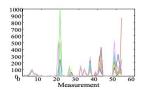
- Return { $\lambda_i$, $\mathbf{u}_i$ }$_{i=1..k}$
  % top $k$ PCA components

# Singular value decomposition (SVD)

Singular Value Decomposition of the **centered** data matrix **X.**

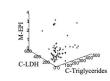$$\mathbf{X}_{\text{features} \times \text{samples}} = \mathbf{USV}^{\mathrm{T}}$$
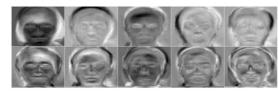
1. Data visualization (blood example)

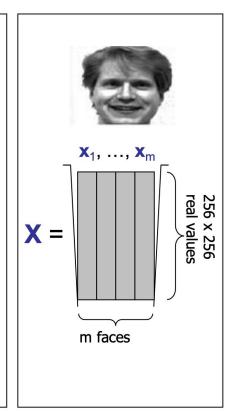2. Noise reduction (eigenfaces)

3. Data compression (image example)

# Application: eigenface

☐ Example data set: Images of faces
  - Eigenface approach
    [Turk & Pentland], [Sirovich & Kirby]
☐ Each face **x** is …

  - 256 × 256 values (luminance at location)
  - **x** in $\Re^{256 \times 256}$ (view as 64K dim vector)

☐ Form **X** = [ $\mathbf{x}_1$ , …, $\mathbf{x}_m$ ] **centered** data matrix

☐ Compute $\Sigma = \mathbf{XX}^\top$

☐ Problem: $\Sigma$ is 64K × 64K … HUGE!!!
   (34 GB in memory)



$\mathbf{x}_1$ , …, $\mathbf{x}_m$

**X** =

256 x 256 real values

m faces

# Computational complexity

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

❑ Given $N{\times}N$ covariance matrix $\Sigma$, can compute
- all $N$ eigenvectors/eigenvalues in $O(N^3)$
- first $k$ eigenvectors/eigenvalues in $O(k N^2)$

❑ But if $N=64K$, EXPENSIVE!

# Computational complexity: how about...

- Note that  $m << 64K$
- Use $L = X^T X$ instead of $\Sigma = XX^T$
- If $v$ is eigenvector of $L$
  then $Xv$ is eigenvector of $\Sigma$
- $O(Nm^2) + O(km^2)$
- 64M vs 42,000M operations

COLUMBIA
UNIVERSITY

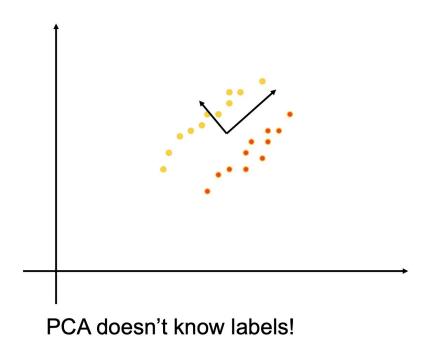# (Dis)advantages

- Required carefully handled data. Sensitive to data preprocessing quality.

- Completely knowledge-free!

# (Dis)advantages



PCA doesn't know labels!
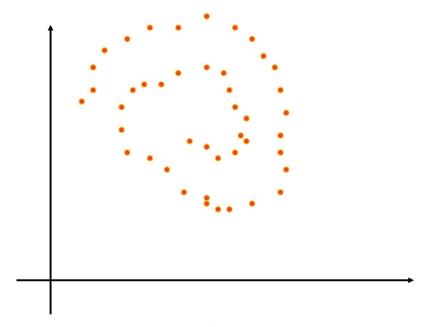
# (Dis)advantages

PCA cannot capture NON-LINEAR structure!

# Takeaways

- PCA:
  - Finds orthonormal basis for data
  - Sorts dimensions in order of "importance"
  - Usually discard unimportant dimensions
- Applications:
  - Visualization
  - Data compression / compact representation
  - Remove noise to improve classification (hopefully)
- Disadvantages:
  - Doesn't know class labels
  - Can only capture linear variations
- One of many ways to reduce dimensionality!

COLUMBIA
UNIVERSITY

# References

- Trevor Hastie, Robert Tibshirani, Jerome Friedman: The Elements of Statistical Learning: Data Mining, Inference and Prediction, Chapter 14.5

- Ziv Bar-Joseph, Tom Mitchell, Pradeep Ravikumar and Aarti Singh: CMU 10-701

COLUMBIA
UNIVERSITY