**NC STATE**
UNIVERSITY

## CSC-591: Foundations of Data Science
## T/Th. 12:50-2:05pm. EBI-1005.

Ranga Raju Vatsavai

Chancellors Faculty Excellence Associate Professor in Geospatial Analytics
Department of Computer Science, North Carolina State University (NCSU)
Associate Director, Center for Geospatial Analytics, NCSU
&
Joint Faculty, Oak Ridge National Laboratory (ORNL)

W11: 10/27/15-10/29/15

---

**NC STATE**
UNIVERSITY

# Admin: Changes in grading

- Midterm-1: 15%
- Midterm-2: 20%
- Final: 35%

- Midterm-2 Topics
  - Regression (all topics covered in the class)
  - Information theory, attribute selection
  - Dimensionality reduction
  - Nonparametric hypothesis testing

# Dimensionality Reduction

- Previously, feature (or attribute) selection
  - Preserves original (reduced) attribute set
    - $X_1$, $X_2$, $X_3$, …$X_{d-1}$, $X_d$
- Dimensionality reduction
  - Preserve as much structure as possible
  - Structure: relationships that affects class separability
  - New feature (transformed) space

10/26/15     © Raju Vatsavai     CSC-591. 3

# Principal Component Analysis

- Principal component analysis (PCA) is an orthogonal transformation original (correlated data) variables into a set of values of linearly uncorrelated variables called principal components.
  - $1^{st}$ component: direction of greatest variability in the data
  - $2^{nd}$ component: orthogonal to $1^{st}$, greatest variability of what's left
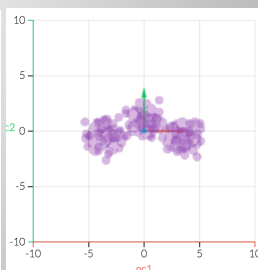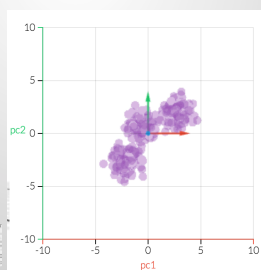  - … and so on until d (original dimensionality)
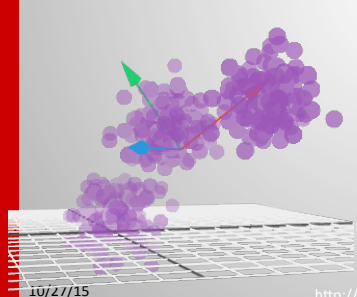
10/26/15     © Raju Vatsavai     CSC-591. 4

# PCA

- Choose first "m" components which become m new dimensions (or new coordinate system)
  - Project data onto new dimensions (i.e., change coordinate of every data point to these new dimensions)
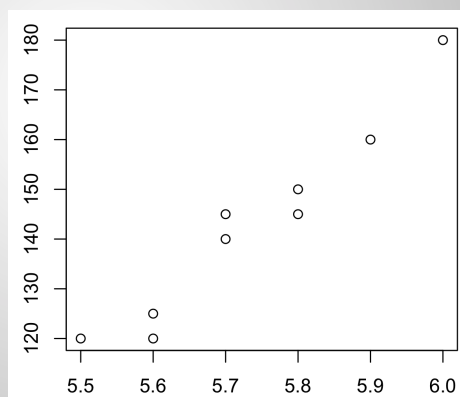


10/27/15     http://setosa.io/ev/principal-component-analysis/     CSC-591. 5

# Greatest Variability

- 2-d example
- Find z that maximize variability, why?
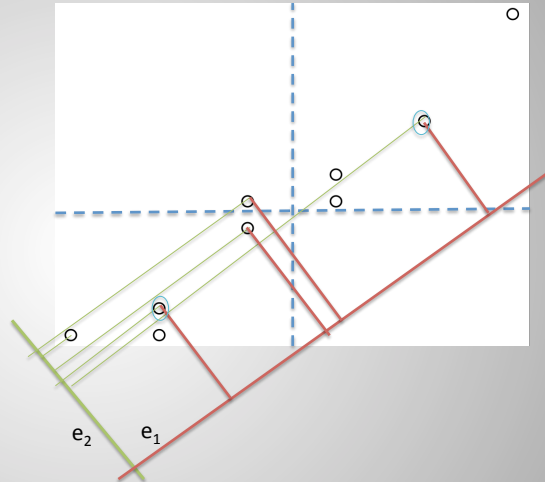


10/26/15     © Raju Vatsavai     CSC-591. 6

3

# Greatest Variability

- 2-d example
- Find z that maximize variability, why?



$e_2$    $e_1$

Reduce cases where two points are close in e-space but very far in (x,y)-space
Minimize distances between original points and their projections

10/27/15                    © Raju Vatsavai                    CSC-591. 7

# How to Get Principal Components

- Center the data at 0. That is, $x_{1,j} = x_{i,j} - \mu_1$
- Normalize attributes if needed



Reduce cases where two points are close in e-space but very far in (x,y)-space
Minimize distances between original points and their projections

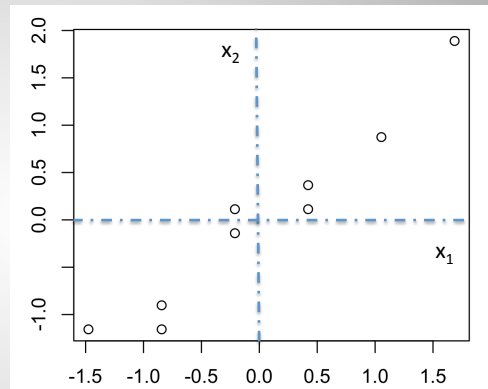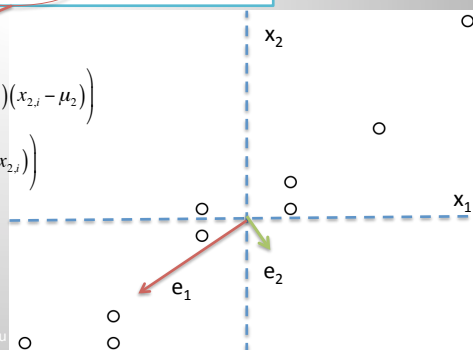10/27/15                    © Raju Vatsavai                    CSC-591. 8

# Principal Components

- Compute covariance matrix $\Sigma$
  - Do x1 and x2 tend to increase (or decrease) together?

|       | $x_1$     | $x_2$     |
|-------|-----------|-----------|
| $x_1$ | 1.0000000 | 0.9696882 |
| $x_2$ | 0.9696882 | 1.0000000 |

  - $Cov(x_1,x_2) = \frac{1}{n}\left(\sum_{i=1}^{n}(x_{1,i}-\mu_1)(x_{2,i}-\mu_2)\right)$

    $= \frac{1}{n}\left(\sum_{i=1}^{n}(x_{1,i})(x_{2,i})\right)$

10/27/15                                    © Raju

---

# Principal Components

- Multiply covariance matrix $\Sigma$ by a vector

$$\begin{pmatrix} 2.0 & 0.8 \\ 0.8 & 0.6 \end{pmatrix}\begin{pmatrix} -1 \\ +1 \end{pmatrix} = \begin{pmatrix} -1.2 \\ -0.2 \end{pmatrix}$$

$$\begin{pmatrix} 2.0 & 0.8 \\ 0.8 & 0.6 \end{pmatrix}\begin{pmatrix} -1.2 \\ -0.2 \end{pmatrix} = \begin{pmatrix} -2.5 \\ -1.08 \end{pmatrix}$$

...

- Repeat
  - Turns towards the direction of variance

10/27/15                                    © Raju

5

# Principal Components

- We saw that the slope of vector converging towards maximum variance, but length is growing faster
  - So we want a vector $e_i$'s that aren't turned by $\Sigma$
- Finding the basis of $\Sigma$
  - $\Sigma\, e = \lambda\, e$
  - e's are called eigenvectors
  - $\lambda$'s are corresponding eigenvalues
- Principal components = eigenvectors with largest eigenvalues

10/27/15       © Raju Vatsavai       CSC-591. 11

# Finding Principal Components

- Find eigenvalues
  - Solve: det($\Sigma$-$\lambda$I) = 0     $\det\begin{pmatrix} 2.0-\lambda & 0.8-0 \\ 0.8-0 & 0.6-\lambda \end{pmatrix} = 0$

$$\det\begin{pmatrix} 2.0-\lambda & 0.8-0 \\ 0.8-0 & 0.6-\lambda \end{pmatrix} = \left(2.0-\lambda\right)\left(0.6-\lambda\right) - \left(0.8\right)\left(0.8\right) = 0$$

$$\lambda^2 - 2.6\lambda + 0.56 = 0$$

$$\{\lambda_1, \lambda_2\} = \frac{1}{2}\left(2.6 \pm \sqrt{2.6^2 - 4*0.56}\right) = \{2.36, 0.23\}$$

10/27/15       © Raju Vatsavai       CSC-591. 12

## Finding Principal Components

- Find i$^{th}$ eigenvector by solving $\Sigma$ e$_i$ = $\lambda_i$ e$_i$

$$\begin{pmatrix} 2.0 & 0.8 \\ 0.8 & 0.6 \end{pmatrix} \begin{pmatrix} e_{1,1} \\ e_{1,2} \end{pmatrix} = 2.36 \begin{pmatrix} e_{1,1} \\ e_{1,2} \end{pmatrix}$$

$$2.0e_{1,1} + 0.8e_{1,2} = 2.36e_{1,1}$$

$$0.8e_{1,1} + 0.6e_{1,2} = 2.36e_{1,2} \qquad 0.8e_{1,2} = (2.36 - 2.0)e_{1,1}$$

$$e_{1,1} = \frac{0.8}{0.36}e_{1,2} = 2.2e_{1,2}$$

- Lots of vectors that satisfy this condition
- The simplest is $e_1 = \begin{bmatrix} 2.2 \\ 1.0 \end{bmatrix}$

10/27/15 © Raju Vatsavai CSC-591. 13

## Finding Principal Components

- To avoid multiple solutions, we want e$_i$'s to be unit vectors, i.e. $\|e_1\| = 1$

$$e_1 = \begin{bmatrix} 0.91 \\ 0.41 \end{bmatrix}$$

- Now solve for 2$^{nd}$ eigenvector

$$\begin{pmatrix} 2.0 & 0.8 \\ 0.8 & 0.6 \end{pmatrix} \begin{pmatrix} e_{1,1} \\ e_{1,2} \end{pmatrix} = 0.23 \begin{pmatrix} e_{1,1} \\ e_{1,2} \end{pmatrix}$$

$$e_2 = \begin{bmatrix} -0.41 \\ 0.91 \end{bmatrix}$$

10/27/15 © Raju Vatsavai CSC-591. 14

## Projection

- After finding principal components, we need to project our data onto new dimensions
- $e_1, \dots e_m$ are new (first m) dimensions
- We have instance $x = \{x_1, \dots x_d\}$ (original coordinates)
- We want new coordinates $x' = \{x'_1, \dots x'_m\}$
  - Center each instance: $(x' - \mu)$
  - Project to each dimension: $(x' - \mu)^T e_j$ for j=1,…,m

10/27/15        © Raju Vatsavai        CSC-591. 15

## Projection

- After finding principal components, we need to project our data onto new dimensions
- $e_1, \dots e_m$ are new (first m) dimensions
- We have instance $x = \{x_1, \dots x_d\}$ (original coordinates)
- We want new coordinates $x' = \{x'_1, \dots x'_m\}$
  - Center each instance: $(x' - \mu)$
  - Project to each dimension: $(x' - \mu)^T e_j$ for j=1,…,m

10/27/15        © Raju Vatsavai        CSC-591. 16

## Key Properties

- Eigenvectors (e) maximizes the variance
- What is the variance along eigenvector
- Variance of projected points $(x^T e) = \lambda$

## How many dimensions

- Of all eigenvectors $e_1, \ldots, e_d$, we want $e_m$, m<<d
- We know, eigenvalue $\lambda_i$ = variance along ei
  - Sort eigenvector s.t. $\lambda_1 \leq \lambda_2 \ldots \leq \lambda_d$
  - Pick first m eigenvectors which explain 90% or (95%) of total variance
- Or, plot eigenvalues as function of dimensions
  - (like K-means)

# Example: Eigen Faces
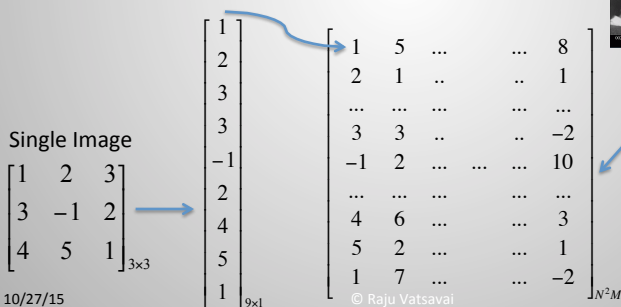
- Face recognition
  using PCA
  (Eigenfaces)

# Representation of Images

- Data: Set of N image of size KxK
- Convert the images into $K^2 \times N$ matrix
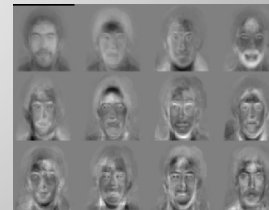- Each NxN image becomes a $N^2$ vector in the matrix

Single Image

$$\begin{bmatrix} 1 & 2 & 3 \\ 3 & -1 & 2 \\ 4 & 5 & 1 \end{bmatrix}_{3\times3} \rightarrow \begin{bmatrix} 1 \\ 2 \\ 3 \\ 3 \\ -1 \\ 2 \\ 4 \\ 5 \\ 1 \end{bmatrix}_{9\times1}$$

$$\begin{bmatrix} 1 & 5 & ... & ... & 8 \\ 2 & 1 & .. & .. & 1 \\ ... & ... & ... & ... & ... \\ 3 & 3 & .. & .. & -2 \\ -1 & 2 & ... & ... & ... & 10 \\ ... & ... & ... & ... & ... \\ 4 & 6 & ... & ... & 3 \\ 5 & 2 & ... & ... & 1 \\ 1 & 7 & ... & ... & -2 \end{bmatrix}_{N^2M}$$

# Representation of Images

1. Average dataset
2. Center the data
   1. Each face differs from the average by vector
3. Compute covariance matrix
4. Find eigenvectors and eigenvalues (set of M eigenvectors each $K^2$ dim)
5. Convert into images (take each column (Eigenvector) and convert it KxK image)

CSC-591. 21

# Nonlinear Dimensionality Reduction

- Its difficult to represent more than 3-d data
- Simplify by assuming that the data of interest lie on an embedded non-linear manifold within the higher-dimensional space. If the manifold is of low enough dimension, the data can be visualized in the low-dimensional space.
- Several approaches
  - ISOMAP
  - LLE

10/27/15       © Raju Vatsavai       CSC-591. 22

# Additional Resources and Acknowledgements

- PCA Tutorials/References
  - A Tutorial on Principal Component Analysis, by Jonathon Shlens, arXiv
  - Principal Component Analysis, by H. Abdi, L. Williams, Weily, 2010
  - (search on web for these articles)
- D. Mladenic