Unsupervised Learning with Mixtures of Gaussians

# EM Algorithm - continued

Machine Learning – CSE446

Carlos Guestrin

University of Washington

May 20, 2013

1

---

# Supervised Learning of Mixtures of Gaussians

*K Components*

*total:*
*K-1 +*
$K\left(\frac{m^2+m}{2}\right)$
*+km*

- Mixtures of Gaussians:
  - Prior class probabilities: P(y) ← *K-1 params (multinomial)*
  - Likelihood function per class: P(**x**|y=i) ← $N(\mu_i, \Sigma_i)$

  $\uparrow$ *symmetric m×m matrix*
  $\frac{m^2}{2} + \frac{m}{2}$ *params*

  *x → m dims*  *m*

- Suppose, for each data point, we know location **x** and class y
  - Learning is easy… ☺

  - For prior P(y) ↩  $P(y=i) = \dfrac{count(y=i) \text{ in data}}{N}$

  - For likelihood function:

    $P(X|y=i)$ ⌐  $\mu_i$ *is average of* $x^j$ *for points in class i*

    $\Sigma_i$ ←  $\sigma_{uv} = \dfrac{\sum_{j \text{ in cluster } i} (x_u^j - \mu_{iu})(x_v^j - \mu_{iv})}{num \text{ points in Cluster } i}$
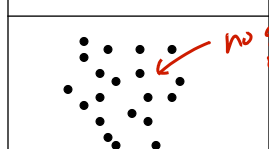
2

---

1

# Unsupervised Learning: not as hard as it looks

*we don't have $y^j$*

Sometimes easy

*← well separated*

*no good mixture of Gaussians*

Sometimes impossible

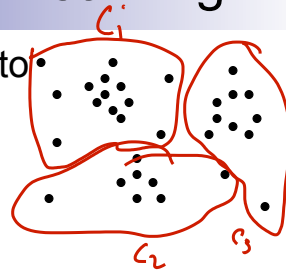and sometimes in between

*← overlapping but meaningful clusters*

*IN CASE YOU'RE WONDERING WHAT THESE DIAGRAMS ARE, THEY SHOW 2-d UNLABELED DATA (X VECTORS) DISTRIBUTED IN 2-d SPACE. THE TOP ONE HAS THREE VERY CLEAR GAUSSIAN CENTERS*

©Carlos Guestrin 2005-2013   3

---

# EM: "Reducing" Unsupervised Learning to Supervised Learning

- If we knew assignment of points to classes ➜ Supervised Learning!

  $C_1$

  $C_2$   $C_3$

- Expectation-Maximization (EM)
  - ☐ Guess assignment of points to classes   *or clusters*
  - ☐ Recompute model parameters
  - ☐ Iterate

©Carlos Guestrin 2005-2013   4

# Back to Unsupervised Learning of Mixtures of Gaussians – a simple version

A simple case:

We have unlabeled data $x_1$ $x_2$ … $x_m$

We know there are k classes

We know $P(y_1)$ $P(y_2)$ $P(y_3)$ … $P(y_k)$

We <u>don't</u> know $\mu_1$ $\mu_2$ .. $\mu_k$

We can write P( data | $\mu_1$…. $\mu_k$)

$$= p\left(x_1...x_m \mid \mu_1...\mu_k\right)$$

$$= \prod_{j=1}^{m} p\left(x_j \mid \mu_1...\mu_k\right)$$

$$= \prod_{j=1}^{m} \sum_{i=1}^{k} p\left(x^j \mid \mu_i\right) P\left(y=i\right)$$

$$\propto \prod_{j=1}^{m} \sum_{i=1}^{k} \exp\left(-\frac{1}{2\sigma^2}\left\|x^j - \mu_i\right\|^2\right) P\left(y=i\right)$$

---

# EM for simple version of Mixtures of Gaussians: The E-step

- If we know $\mu_1,\ldots,\mu_k$ → easily compute prob.

  point $x^j$ belongs to class y=i

$$p\left(y=i \mid x^j, \mu_1...\mu_k\right) \propto \exp\left(-\frac{1}{2\sigma^2}\left\|x^j - \mu_i\right\|^2\right) P\left(y=i\right)$$

# EM for simple version of Mixtures of Gaussians: The M-step

- If we know prob. point $x^j$ belongs to class $y=i$

  $\rightarrow$ MLE for $\mu_i$ is weighted average

  - imagine k copies of each $x^j$, each with weight $P(y=i|x^j)$:

$$\mu_i = \frac{\sum_{j=1}^{m} P\left(y=i\middle|x^j\right) x^j}{\sum_{j=1}^{m} P\left(y=i\middle|x^j\right)}$$

7

# E.M. for Simple version of Mixtures of Gaussians

**E-step**

Compute "expected" classes of all datapoints for each class

$$p\left(y=i\middle|x^j,\mu_1...\mu_k\right) \propto \exp\left(-\frac{1}{2\sigma^2}\left\|x^j-\mu_i\right\|^2\right) P\left(y=i\right)$$

*Just evaluate a Gaussian at $x^j$*

**M-step**

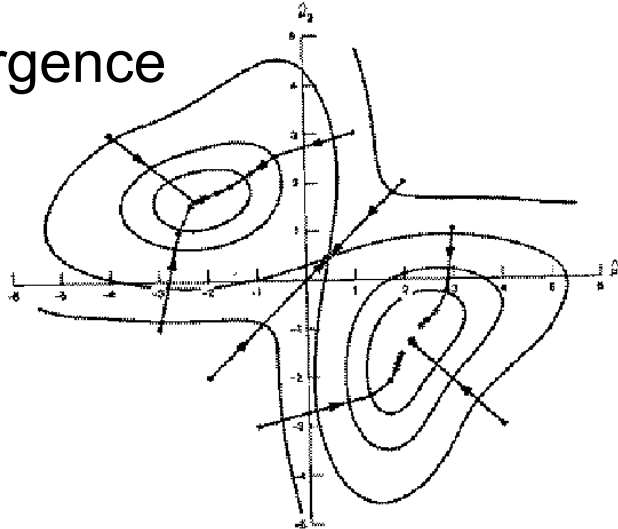Compute Max. like **μ** given our data's class membership distributions

$$\mu_i = \frac{\sum_{j=1}^{m} P\left(y=i\middle|x^j\right) x^j}{\sum_{j=1}^{m} P\left(y=i\middle|x^j\right)}$$

8

# E.M. Convergence

- EM is coordinate ascent on an interesting potential function
- Coord. ascent for bounded pot. func. ! convergence to a local optimum guaranteed



- This algorithm is REALLY USED. And in high dimensional state spaces, too. E.G. Vector Quantization for Speech Data

9

---

# E.M. for axis-aligned GMM

$$\Sigma = \begin{pmatrix} \sigma_1^2 & 0 & 0 & \cdots & 0 & 0 \\ 0 & \sigma_2^2 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \sigma_3^2 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & \sigma_{m-1}^2 & 0 \\ 0 & 0 & 0 & \cdots & 0 & \sigma_m^2 \end{pmatrix}$$

Iterate. On the $t$'th iteration let our estimates be

$\lambda_t = \{ \mu_1^{(t)}, \mu_2^{(t)} \ldots \mu_k^{(t)}, \Sigma_1^{(t)}, \Sigma_2^{(t)} \ldots \Sigma_k^{(t)}, p_1^{(t)}, p_2^{(t)} \ldots p_k^{(t)} \}$

$p_i^{(t)}$ is shorthand for estimate of prior $P(y=i)$ on t'th iteration

**E-step**

Compute "expected" classes of all datapoints for each class

$$P\left( y = i \big| x^j, \lambda_t \right) \propto p_i^{(t)} p\left( x^j \big| \mu_i^{(t)}, \Sigma_i^{(t)} \right)$$

*Just evaluate a Gaussian at $x^j$*

M-step

Compute Max. like **μ** given our data's class membership distributions

$$\mu_i^{(t+1)} = \frac{\sum_j P\left( y = i \big| x^j, \lambda_t \right) x^j}{\sum_j P\left( y = i \big| x^j, \lambda_t \right)}$$

$$p_i^{(t+1)} = \frac{\sum_j P\left( y = i \big| x^j, \lambda_t \right)}{m}$$

$m$ = #records

10

5

# E.M. for General GMMs

Iterate. On the *t*'th iteration let our estimates be

$$\lambda_t = \{ \mu_1^{(t)}, \mu_2^{(t)} \dots \mu_k^{(t)}, \Sigma_1^{(t)}, \Sigma_2^{(t)} \dots \Sigma_k^{(t)}, p_1^{(t)}, p_2^{(t)} \dots p_k^{(t)} \}$$

**E-step**

Compute "expected" classes of all datapoints for each class

$$P\left( y = i \middle| x^j, \lambda_t \right) \propto p_i^{(t)} p\left( x^j \middle| \mu_i^{(t)}, \Sigma_i^{(t)} \right)$$

*Just evaluate a Gaussian at $x^j$*

M-step

Compute Max. like **μ** given our data's class membership distributions

$$\mu_i^{(t+1)} = \frac{\sum_j P\left( y = i \middle| x^j, \lambda_t \right) x^j}{\sum_j P\left( y = i \middle| x^j, \lambda_t \right)} \qquad \Sigma_i^{(t+1)} = \frac{\sum_j P\left( y = i \middle| x^j, \lambda_t \right) \left[ x^j - \mu_i^{(t+1)} \right]\left[ x^j - \mu_i^{(t+1)} \right]^T}{\sum_j P\left( y = i \middle| x^j, \lambda_t \right)}$$

$$p_i^{(t+1)} = \frac{\sum_j P\left( y = i \middle| x^j, \lambda_t \right)}{m}$$

*m* = #records

# Gaussian Mixture Example: Start



p=0.333

p=0.333

0.333

# After first iteration

# After 2nd iteration

7

# After 3rd iteration



p=0.34

p=0.307

# After 4th iteration



p=0.331

p=0.288

8

# After 5th iteration



ρ=0.322

ρ=0.285

17

# After 6th iteration



ρ=0.315

ρ=0.287

18

9

# After 20th iteration



p=0.234   p=0.334

**19**

# Some Bio Assay data



**20**

10

# GMM clustering of the assay data

21

# Resulting Density Estimator

22

11

## Three classes of assay
(each learned with it's own mixture model)

Compound =
IL-1
TNF
none

n u c l e u s

23

## Resulting Bayes Classifier

Compound =
IL-1
TNF
none

n u c l e u s

24

Resulting Bayes Classifier, using posterior probabilities to alert about ambiguity and anomalousness

Compound =
IL-1
TNF
none

Yellow means ANOMALOUS
Cyan means AMBIGUOUS

nucleus

**Yellow means anomalous**

**Cyan means ambiguous**
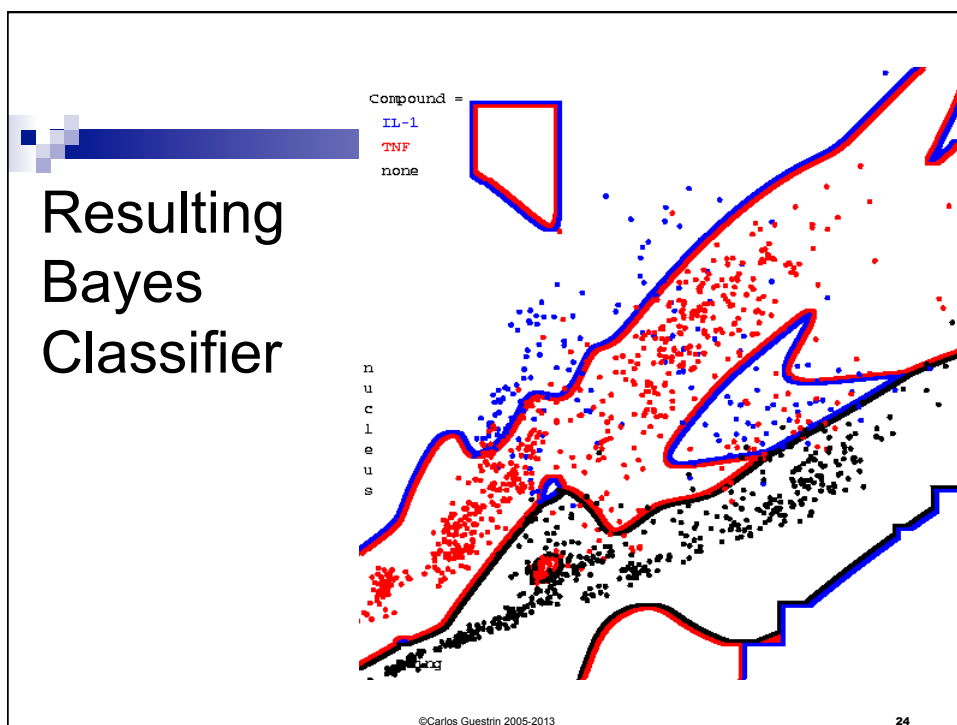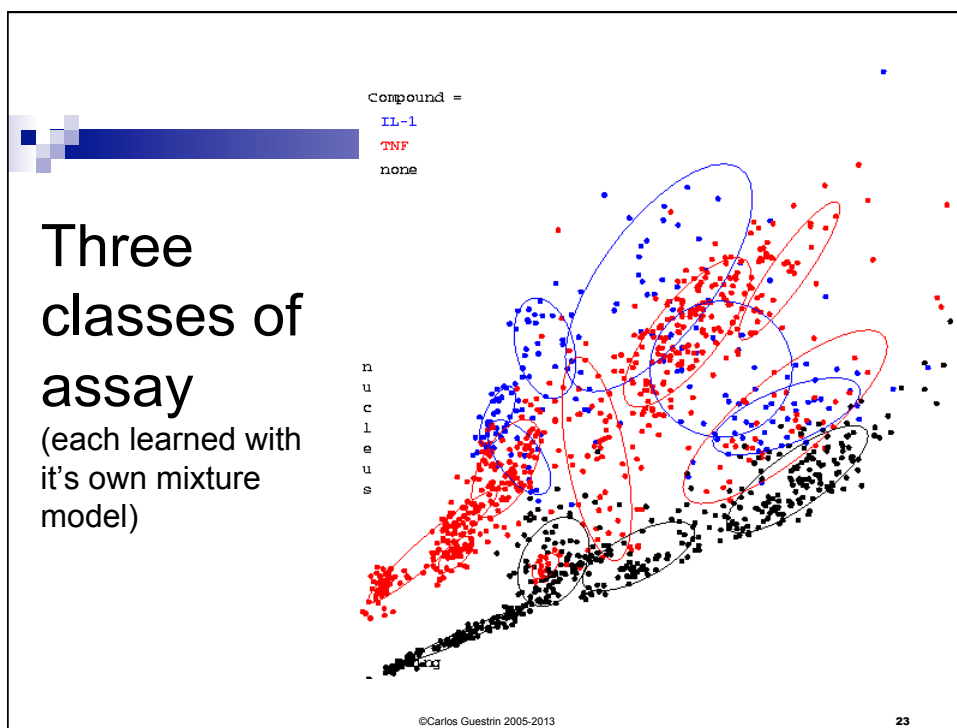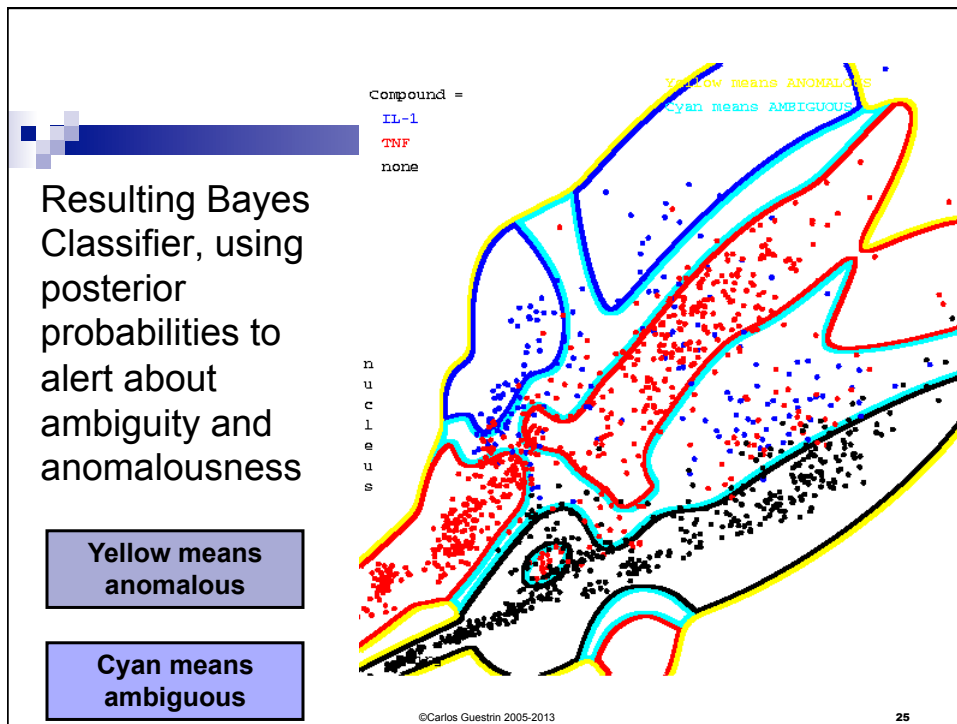
©Carlos Guestrin 2005-2013                                          25

---

# E.M.: The General Case

- E.M. widely used beyond mixtures of Gaussians
  - The recipe is the same…

- Expectation Step:  Fill in missing data, given current values of parameters, $\theta^{(t)}$
  - If variable $y$ is missing (could be many variables)
  - Compute, for each data point $\mathbf{x}^j$, for each value $i$ of $y$:
    - $P(y=i|\mathbf{x}^j,\theta^{(t)})$

- Maximization step:  Find maximum likelihood parameters for (weighted) "completed data":
  - For each data point $\mathbf{x}^j$, create $k$ weighted data points
    - 
  - Set $\theta^{(t+1)}$ as the maximum likelihood parameter estimate for this weighted data

- Repeat

©Carlos Guestrin 2005-2013                                          26

13

# What you should know

- K-means for clustering:
  - □ algorithm
  - □ converges because it's coordinate ascent
- EM for mixture of Gaussians:
  - □ How to "learn" maximum likelihood parameters (locally max. like.) in the case of unlabeled data
- Be happy with this kind of probabilistic analysis
- Remember, E.M. can get stuck in local minima, and empirically it <u>DOES</u>
- EM is coordinate ascent

27

# Acknowledgements

- K-means & Gaussian mixture models presentation contains material from excellent tutorial by Andrew Moore:
  - □ http://www.autonlab.org/tutorials/
- K-means Applet:
  - □ http://www.elet.polimi.it/upload/matteucc/Clustering/tutorial_html/AppletKM.html
- Gaussian mixture models Applet:
  - □ http://www.neurosci.aist.go.jp/%7Eakaho/MixtureEM.html

28

# Dimensionality Reduction PCA

Machine Learning – CSE446

Carlos Guestrin

University of Washington

May 20, 2013

29

---

# Dimensionality reduction

- Input data may have thousands or millions of dimensions!
    - □ e.g., text data has
- **Dimensionality reduction**: represent data with fewer dimensions
    - □ easier learning – fewer parameters
    - □ visualization – hard to visualize more than 3D or 4D
    - □ discover "intrinsic dimensionality" of data
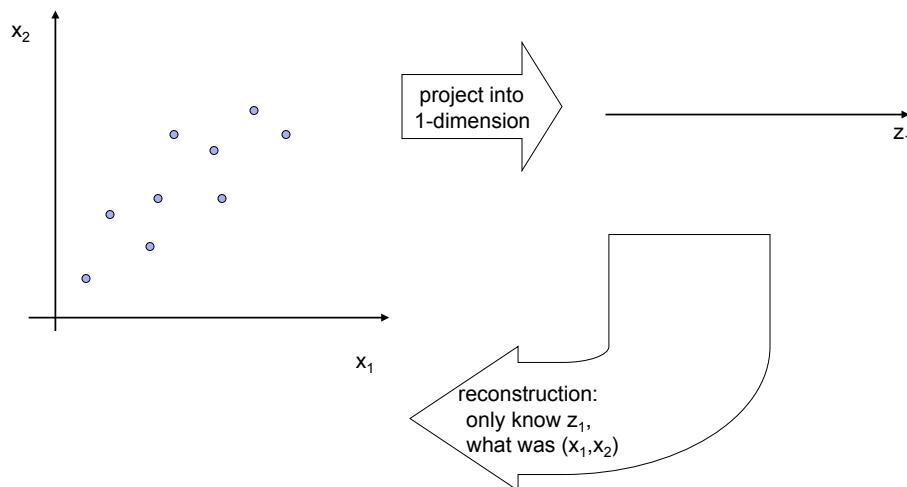        - ■ high dimensional data that is truly lower dimensional

# Lower dimensional projections

- Rather than picking a subset of the features, we can new features that are combinations of existing features

- Let's see this in the unsupervised setting
  - just **X**, but no Y

# Linear projection and reconstruction



$x_2$

project into 1-dimension

$z_1$

reconstruction: only know $z_1$, what was $(x_1, x_2)$

$x_1$

# Principal component analysis – basic idea

- Project n-dimensional data into k-dimensional space while preserving information:
  - e.g., project space of 10000 words into 3-dimensions
  - e.g., project 3-d into 2-d

- Choose projection with minimum reconstruction error

# Linear projections, a review

- Project a point into a (lower dimensional) space:
  - **point**: $\mathbf{x} = (x_1,\ldots,x_d)$
  - **select a basis** – set of basis vectors – $(\mathbf{u}_1,\ldots,\mathbf{u}_k)$
    - we consider orthonormal basis:
      - $\mathbf{u}_i \bullet \mathbf{u}_i = 1$, and $\mathbf{u}_i \bullet \mathbf{u}_j = 0$ for $i \neq j$
  - **select a center** – $\bar{\mathbf{x}}$, defines offset of space
  - **best coordinates** in lower dimensional space defined by dot-products: $(z_1,\ldots,z_k)$, $z_i = (\mathbf{x}-\bar{\mathbf{x}})\bullet\mathbf{u}_i$
    - minimum squared error

# PCA finds projection that minimizes reconstruction error
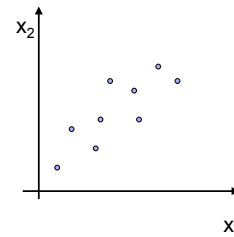
- Given m data points: $\mathbf{x}^i = (x_1^i,\ldots,x_d^i)$, i=1…N
- Will represent each point as a projection:

  □ $\hat{\mathbf{x}}^i = \bar{\mathbf{x}} + \sum_{j=1}^{k} z_j^i \mathbf{u}_j$  where:  $\bar{\mathbf{x}} = \dfrac{1}{N}\sum_{i=1}^{N}\mathbf{x}^i$  and  $z_j^i = (\mathbf{x}^i - \bar{\mathbf{x}})\cdot\mathbf{u}_j$

- PCA:
  - □ Given k<<d, find $(\mathbf{u}_1,\ldots,\mathbf{u}_k)$ minimizing reconstruction error:

    $error_k = \sum_{i=1}^{N}(\mathbf{x}^i - \hat{\mathbf{x}}^i)^2$



$x_2$

$x_1$

---

# Understanding the reconstruction error

$$\hat{\mathbf{x}}^i = \bar{\mathbf{x}} + \sum_{j=1}^{k} z_j^i \mathbf{u}_j$$

$$z_j^i = (\mathbf{x}^i - \bar{\mathbf{x}})\cdot\mathbf{u}_j$$

□ Given k<<d, find $(\mathbf{u}_1,\ldots,\mathbf{u}_k)$ minimizing reconstruction error:

$$error_k = \sum_{i=1}^{N}(\mathbf{x}^i - \hat{\mathbf{x}}^i)^2$$

- Note that $\mathbf{x}^i$ can be represented exactly by d-dimensional projection:

  $\mathbf{x}^i = \bar{\mathbf{x}} + \sum_{j=1}^{d} z_j^i \mathbf{u}_j$

- Rewriting error:

18

# Reconstruction error and covariance matrix

$$error_k = \sum_{i=1}^{N} \sum_{j=k+1}^{d} [\mathbf{u}_j \cdot (\mathbf{x}^i - \bar{\mathbf{x}})]^2 \qquad \Sigma = \frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}^i - \bar{\mathbf{x}})(\mathbf{x}^i - \bar{\mathbf{x}})^T$$

# Minimizing reconstruction error and eigen vectors

- Minimizing reconstruction error equivalent to picking orthonormal basis ($\mathbf{u}_1,\ldots,\mathbf{u}_d$) minimizing:

$$error_k = N \sum_{j=k+1}^{d} \mathbf{u}_j^T \Sigma \mathbf{u}_j$$

- Eigen vector:


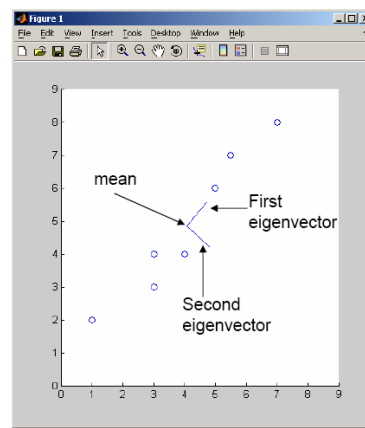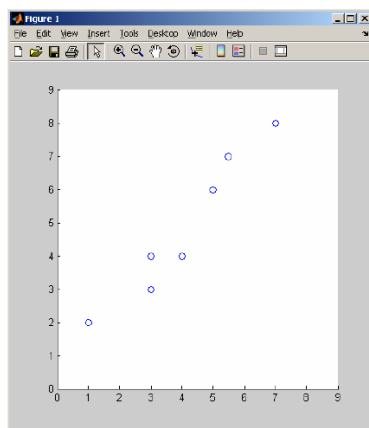- Minimizing reconstruction error equivalent to picking ($\mathbf{u}_{k+1}$, …,$\mathbf{u}_d$) to be eigen vectors with smallest eigen values

# Basic PCA algoritm

- Start from m by n data matrix **X**
- **Recenter**: subtract mean from each row of **X**
  - $X_c \leftarrow X - \overline{X}$
- **Compute covariance matrix**:
  - $\Sigma \leftarrow 1/N\ X_c^T\ X_c$
- Find **eigen vectors and values** of $\Sigma$
- **Principal components:** k eigen vectors with highest eigen values

---

# PCA example

$$\hat{x}^i = \bar{x} + \sum_{j=1}^{k} z_j^i u_j$$

# PCA example – reconstruction

$$\hat{\mathbf{x}}^i = \bar{\mathbf{x}} + \sum_{j=1}^{k} z_j^i \mathbf{u}_j$$

only used first principal component

# Eigenfaces [Turk, Pentland '91]

- Input images:

- Principal components:

# Eigenfaces reconstruction

- Each image corresponds to adding 8 principal components:

---

# Scaling up

- Covariance matrix can be really big!
  - $\Sigma$ is d by d
  - Say, only 10000 features
  - finding eigenvectors is very slow…

- Use singular value decomposition (SVD)
  - finds to k eigenvectors
  - great implementations available, e.g., R or Matlab svd

# SVD

- Write $X = W S V^T$
  - $X \leftarrow$ data matrix, one row per datapoint
  - $W \leftarrow$ weight matrix, one row per datapoint – coordinate of $x^i$ in eigenspace
  - $S \leftarrow$ singular value matrix, diagonal matrix
    - in our setting each entry is eigenvalue $\lambda_j$
  - $V^T \leftarrow$ singular vector matrix
    - in our setting each row is eigenvector $v_j$

# PCA using SVD algoritm

- Start from m by n data matrix $X$
- **Recenter**: subtract mean from each row of $X$
  - $X_c \leftarrow X - \overline{X}$
- Call SVD algorithm on $X_c$ – ask for k singular vectors
- **Principal components:** k singular vectors with highest singular values (rows of $V^T$)
  - **Coefficients** become:

# What you need to know

- Dimensionality reduction
  - why and when it's important
- Simple feature selection
- Principal component analysis
  - minimizing reconstruction error
  - relationship to covariance matrix and eigenvectors
  - using SVD