# Expectation Maximization

Machine Learning – CSE546

Carlos Guestrin

University of Washington

November 13, 2014

1

---

# 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 $z$ is missing (could be many variables)
  - □ Compute, for each data point $\mathbf{x}^j$, for each value $i$ of $z$:
    - $P(z=i|\mathbf{x}^j, \theta^{(t)})$

- 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

2

---

1

# The general learning problem with missing data

■ Marginal likelihood – **x** is observed, **z** is missing:

$$\ell(\theta : \mathcal{D}) = \log \prod_{j=1}^{m} P(\mathbf{x}_j \mid \theta)$$

*max θ* · *Observed part*

$$= \sum_{j=1}^{m} \log P(\mathbf{x}_j \mid \theta)$$

$$= \sum_{j=1}^{m} \log \sum_{\mathbf{z}} P(\mathbf{x}_j, \mathbf{z} \mid \theta)$$

*k* · *estimate z's* · *max θ ℓ*

*latent*

3

---

# E-step
*weight*

■ **x** is observed, **z** is missing
■ Compute probability of missing data given current choice of $\theta^{(t)}$
  □ $Q(\mathbf{z}|\mathbf{x}_j)$ for each $\mathbf{x}_j$
    ■ e.g., probability computed during classification step
    ■ corresponds to "classification step" in K-means

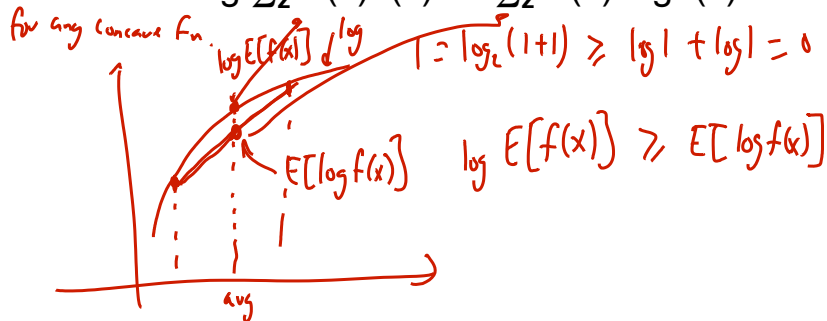$$Q^{(t+1)}(\mathbf{z}^i \mid \mathbf{x}_j) = P(\mathbf{z}^i \mid \mathbf{x}_j, \theta^{(t)})$$

*anything*     *Got previous iteration*

4

# Jensen's inequality

$$\ell(\theta : \mathcal{D}) \;=\; \sum_{j=1}^{m} \log \sum_{\mathbf{z}} P(\mathbf{z} \mid \mathbf{x}_j) P(\mathbf{x}_j \mid \theta)$$

*log*    *avg*    *f function*    *avg*   *logs*

- **Theorem**: $\log \sum_{\mathbf{z}} P(\mathbf{z})\, f(\mathbf{z}) \;\geq\; \sum_{\mathbf{z}} P(\mathbf{z}) \log f(\mathbf{z})$

*for any concave fn.*

$\log E[f(x)]$   *log*    $1 = \log_2(1+1) \geq |g| + \log|1 = 0$

$E[\log f(x)]$    $\log E[f(x)] \geq E[\log f(x)]$

*avg*

©Carlos Guestrin 2005-2014

# Applying Jensen's inequality

$\log \frac{a}{b} = \log a - \log b$

- Use: $\log \sum_{\mathbf{z}} P(\mathbf{z})\, f(\mathbf{z}) \geq \sum_{\mathbf{z}} P(\mathbf{z}) \log f(\mathbf{z})$

*f complex, instead max a lower bound*

$$\boxed{\max_{\theta}}\;\ell(\theta^{(t)} : \mathcal{D}) \;=\; \sum_{j=1}^{m} \log \sum_{\mathbf{z}} \underbrace{Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j)}_{P} \underbrace{\frac{P(\mathbf{z}, \mathbf{x}_j \mid \theta^{(t)})}{Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j)}}_{f(\mathbf{z})}$$

$$\geq \sum_{j=1}^{m} \sum_{z=1}^{k} Q^{(t+1)}(z \mid x_j) \log \frac{P(z, x_j \mid \theta^{(t)})}{Q^{(t+1)}(z \mid x_j)}$$

$$= \underbrace{\sum_{j=1}^{m} \sum_{z=1}^{k} Q^{(t+1)}(z \mid x_j) \log P(z, x_j \mid \theta^{(t)})}_{\text{weighted max likelihood}} - \underbrace{\sum_{j=1}^{n} \sum_{z=1}^{k} Q^{(t+1)}(z \mid x_j) \log Q^{(t+1)}(z \mid x_j)}_{\substack{m\,H(Q^{(t+1)}) \\ \text{doesn't depend on } \theta!!}}$$

©Carlos Guestrin 2005-2014

3

# The M-step maximizes lower bound on weighted data

- Lower bound from Jensen's:

*want max $\theta$*

$$\ell(\theta^{(t)} : \mathcal{D}) \geq \sum_{j=1}^{m} \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j \mid \theta^{(t)}) + m.H(Q^{(t+1)})$$

*doesn't depend on $\theta$*

*focus on this first term in M-step*

- Corresponds to weighted dataset:
  - $\langle \mathbf{x}_1, \mathbf{z}=1 \rangle$ with weight $Q^{(t+1)}(\mathbf{z}=1|\mathbf{x}_1)$ *0.3*
  - $\langle \mathbf{x}_1, \mathbf{z}=2 \rangle$ with weight $Q^{(t+1)}(\mathbf{z}=2|\mathbf{x}_1)$ *0.6*
  - $\langle \mathbf{x}_1, \mathbf{z}=3 \rangle$ with weight $Q^{(t+1)}(\mathbf{z}=3|\mathbf{x}_1)$ *0.1*
  - $\langle \mathbf{x}_2, \mathbf{z}=1 \rangle$ with weight $Q^{(t+1)}(\mathbf{z}=1|\mathbf{x}_2)$ ⋮
  - $\langle \mathbf{x}_2, \mathbf{z}=2 \rangle$ with weight $Q^{(t+1)}(\mathbf{z}=2|\mathbf{x}_2)$ ⋮
  - $\langle \mathbf{x}_2, \mathbf{z}=3 \rangle$ with weight $Q^{(t+1)}(\mathbf{z}=3|\mathbf{x}_2)$ ⋮
  - …

7

# The M-step

$$\ell(\theta^{(t)} : \mathcal{D}) \geq \sum_{j=1}^{m} \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j \mid \theta^{(t)}) + m.H(Q^{(t+1)})$$

*fixed choice of Q*

- Maximization step:

$$\theta^{(t+1)} \leftarrow \arg\max_{\theta} \sum_{j=1}^{m} \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j \mid \theta)$$
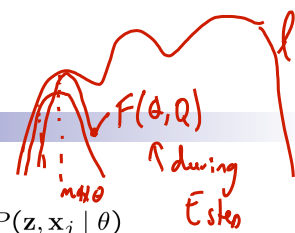
- Use expected counts instead of counts:
  - If learning requires Count($\mathbf{x}$,$\mathbf{z}$)
  - Use $E_{Q(t+1)}[\text{Count}(\mathbf{x},\mathbf{z})]$

8

4

# Convergence of EM

*[handwritten: $F(\theta,Q)$, $m4/\theta$, ↑ during E step]*

- Define potential function F(θ,Q):

$$\ell(\theta : \mathcal{D}) \;\geq\; F(\theta, Q) = \sum_{j=1}^{m} \sum_{\mathbf{z}} Q(\mathbf{z} \mid \mathbf{x}_j) \log \frac{P(\mathbf{z}, \mathbf{x}_j \mid \theta)}{Q(\mathbf{z} \mid \mathbf{x}_j)}$$

- EM corresponds to coordinate ascent on F
  - □ Thus, maximizes lower bound on marginal log likelihood
  - □ We saw that M-step corresponds to fixing Q, max θ  *[handwritten: coord ascent]*
  - □ E-step fix θ and max Q

*[handwritten: ↑ Same kind of analysis]*

---

# M-step is easy

$$\theta^{(t+1)} \leftarrow \arg\max_{\theta} \sum_{j=1}^{m} \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j \mid \theta)$$

- Using potential function

$$F(\theta, Q^{(t+1)}) \;=\; \sum_{j=1}^{m} \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j \mid \theta) + m.H(Q^{(t+1)})$$

# E-step also doesn't decrease potential function 1

- Fixing θ to θ(t):

$$\ell(\theta^{(t)} : \mathcal{D}) \geq F(\theta^{(t)}, Q) = \sum_{j=1}^{m} \sum_{\mathbf{z}} Q(\mathbf{z} \mid \mathbf{x}_j) \log \frac{P(\mathbf{z}, \mathbf{x}_j \mid \theta^{(t)})}{Q(\mathbf{z} \mid \mathbf{x}_j)}$$

# KL-divergence

- Measures distance between distributions

$$KL(Q||P) = \sum_{z} Q(z) \log \frac{Q(z)}{P(z)}$$

- KL=zero if and only if Q=P

# E-step also doesn't decrease potential function 2

- Fixing θ to θ<sup>(t)</sup>:

$$
\begin{aligned}
\ell(\theta^{(t)} : \mathcal{D}) \geq F(\theta^{(t)}, Q) &= \ell(\theta^{(t)} : \mathcal{D}) + \sum_{j=1}^{m} \sum_{\mathbf{z}} Q(\mathbf{z} \mid \mathbf{x}_j) \log \frac{P(\mathbf{z} \mid \mathbf{x}_j, \theta^{(t)})}{Q(\mathbf{z} \mid \mathbf{x}_j)} \\
&= \ell(\theta^{(t)} : \mathcal{D}) - \sum_{j=1}^{m} KL\left(Q(\mathbf{z} \mid \mathbf{x}_j) || P(\mathbf{z} \mid \mathbf{x}_j, \theta^{(t)})\right)
\end{aligned}
$$

©Carlos Guestrin 2005-2014

13

# E-step also doesn't decrease potential function 3

$$
\ell(\theta^{(t)} : \mathcal{D}) \geq F(\theta^{(t)}, Q) = \ell(\theta^{(t)} : \mathcal{D}) - m \sum_{j=1}^{m} KL\left(Q(\mathbf{z} \mid \mathbf{x}_j) || P(\mathbf{z} \mid \mathbf{x}_j, \theta^{(t)})\right)
$$

- Fixing θ to θ<sup>(t)</sup>
- Maximizing F(θ<sup>(t)</sup>,Q) over Q → set Q to posterior probability:

$$
Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j) \leftarrow P(\mathbf{z} \mid \mathbf{x}_j, \theta^{(t)})
$$

- Note that

$$
F(\theta^{(t)}, Q^{(t+1)}) = \ell(\theta^{(t)} : \mathcal{D})
$$

14

# EM is coordinate ascent

$$\ell(\theta : \mathcal{D}) \;\geq\; F(\theta, Q) = \sum_{j=1}^{m} \sum_{\mathbf{z}} Q(\mathbf{z} \mid \mathbf{x}_j) \log \frac{P(\mathbf{z}, \mathbf{x}_j \mid \theta)}{Q(\mathbf{z} \mid \mathbf{x}_j)}$$

- **M-step**: Fix Q, maximize F over θ (a lower bound on $\ell(\theta : \mathcal{D})$ ):

$$\ell(\theta : \mathcal{D}) \;\geq\; F(\theta, Q^{(t)}) = \sum_{j=1}^{m} \sum_{\mathbf{z}} Q^{(t)}(\mathbf{z} \mid \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j \mid \theta) + m.H(Q^{(t)})$$

- **E-step**: Fix θ, maximize F over Q:

$$\ell(\theta^{(t)} : \mathcal{D}) \geq F(\theta^{(t)}, Q) \;=\; \ell(\theta^{(t)} : \mathcal{D}) - m \sum_{j=1}^{m} KL\left(Q(\mathbf{z} \mid \mathbf{x}_j) || P(\mathbf{z} \mid \mathbf{x}_j, \theta^{(t)})\right)$$

- □ "Realigns" F with likelihood:

$$F(\theta^{(t)}, Q^{(t+1)}) \;=\; \ell(\theta^{(t)} : \mathcal{D})$$

15

---

# 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
- General case for EM

16

# Dimensionality Reduction PCA

Machine Learning – CSE4546

Carlos Guestrin

University of Washington

November 13, 2014
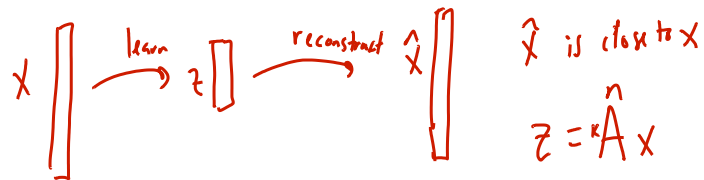
17

---

# Dimensionality reduction

- Input data may have thousands or millions of dimensions! $x$
  - e.g., text data has 10,000 words — 10 000 000 000
- **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

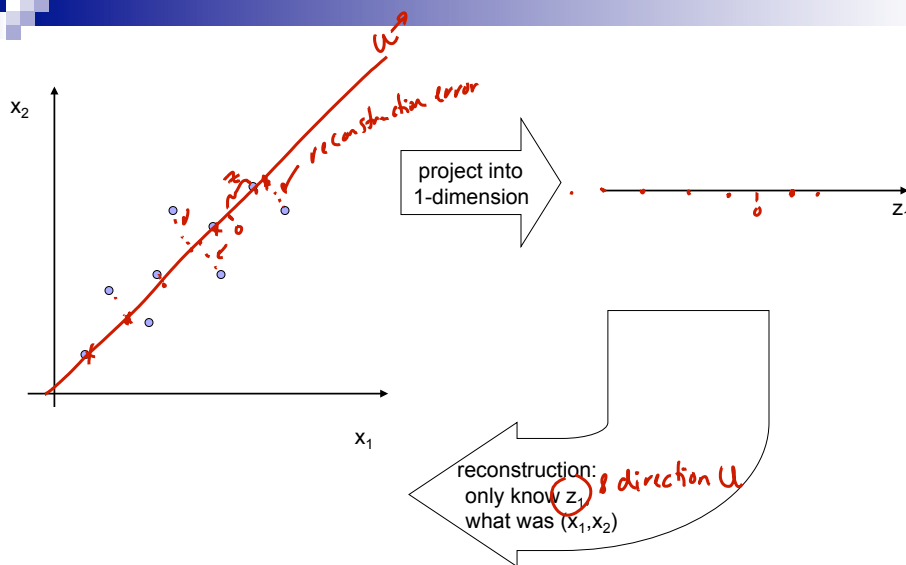- Rather than picking a subset of the features, we can new features that are combinations of existing features

$$z_1 = 2.5 X_1 - 3.2 X_2 - 3.1 X_3 + 70 X_4$$

$X \xrightarrow{\text{learn}} z \xrightarrow{\text{reconstruct}} \hat{X}$

$\hat{X}$ is close to $X$

$z = \overset{n}{\overset{k}{A}} X$

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

---

# Linear projection and reconstruction



$x_2$

$u$

reconstruction error

project into 1-dimension

$z_1$

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

direction $u$

$x_1$

# 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:
      - $\mathbf{u}_i \cdot \mathbf{u}_i = 1$, and $\mathbf{u}_i \cdot \mathbf{u}_j = 0$ for $i \neq j$
  - **select a center** – $\bar{\mathbf{x}}$, defines offset of space
  - **best coordinates** in lower dimensional space defined by dot-products: $(z_1,\ldots,z_k)$, $z_i = (\mathbf{x}-\bar{\mathbf{x}}) \cdot \mathbf{u}_i$
    - minimum squared error

$$z_1 = (x-\bar{x}) \cdot u_1 \quad\Big\} \quad z_1 = \underset{z}{\arg\min}\left((x-\bar{x}) - z\, u_1\right)^2$$

# PCA finds projection that minimizes reconstruction error

- Given N data points: $\mathbf{x}^i = (x_1^i,\ldots,x_d^i)$, i=1…N
- Will represent each point as a projection:

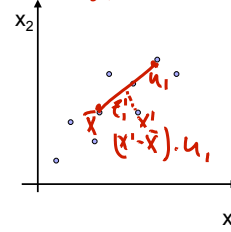  □ $\hat{\mathbf{x}}^i = \bar{\mathbf{x}} + \sum_{j=1}^{k} z_j^i \mathbf{u}_j$ where: $\bar{\mathbf{x}} = \dfrac{1}{N}\sum_{i=1}^{N}\mathbf{x}^i$ and $z_j^i = (\mathbf{x}^i - \bar{\mathbf{x}}) \cdot \mathbf{u}_j$

  *avg of data*    *projection from last slide*

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

    *for a choice of basis $u_1 \ldots u_k$*

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

---

# Understanding the reconstruction error

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

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

□ 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:

$$error_k = \sum_{i=1}^{N}(\mathbf{x}^i - \hat{\mathbf{x}}^i)^2 = \sum_{i=1}^{N}\left[\bar{x} + \sum_{j=1}^{d}z_j^i u_j - \left(\bar{x} + \sum_{j=1}^{k}z_j^i u_j\right)\right]^2 = \sum_{i=1}^{N}\left[\sum_{j=k+1}^{d}z_j^i u_j\right]^2$$

$$= \sum_{i=1}^{N}\left[\sum_{j=k+1}^{d} z_j^i u_j \cdot u_j z_j^i + \sum_{j=k+1, j'\neq j}^{d} z_j^i u_j \cdot u_{j'} z_{j'}^i\right]$$

$$= \sum_{i=1}^{N}\sum_{j=k+1}^{d}(z_j^i)^2 \quad \Leftarrow \quad minimizing\ projection\ error \equiv min\ square\ of\ thrown\ out\ coeffs.$$

12

# Reconstruction error and covariance matrix

*(handwritten top right)* memory low $\Sigma u = \lambda u$ ← eigen vector

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

*(handwritten: $z_j^i$)*

$$= \sum_{i=1}^{N} \sum_{j=k+1}^{d} u_j^T (x^i - \bar{x})(x^i - \bar{x})^T u_j$$

flip sums, $u_j$ don't depend on $i$

$$= \sum_{j=k+1}^{d} u_j^T \left[ \sum_{i=1}^{N} (x^i - \bar{x})(x^i - \bar{x})^T \right] u_j$$

$$\underbrace{\qquad}_{N\Sigma}$$

$$= N \sum_{j=k+1}^{d} u_j^T \Sigma u_j \leftarrow \text{choose to ignore } u_j \text{ that minimize this error}$$

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

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \\ & \sigma_2^2 & \\ & & \ddots \end{pmatrix}$$

$$\sigma_{rs} \overset{MLE}{=} \frac{1}{N} \sum_{i=1}^{N} (x_r^i - \bar{x}_r)(x_s^i - \bar{x}_s)$$

in vector form

---

# 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:

$$u^T \Sigma u = \lambda u^T u = \lambda$$

$$\Sigma u = \lambda u$$
*(handwritten: ↑ eigenvector ← eigen value)*

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

Ignored dims

$$\min_{u_1 \ldots u_n} error_k \equiv \text{throwing out } u_{k+1} \ldots u_d$$

directions with smallest eigen values

$$\equiv \text{keep } u_1 \ldots u_k \text{ to be eigen vectors of } \Sigma$$
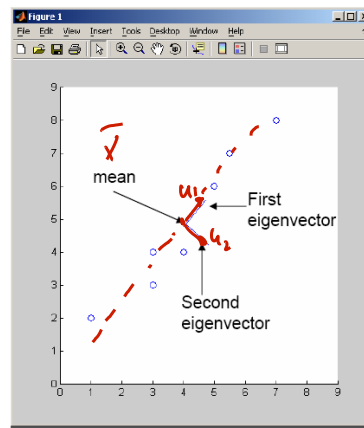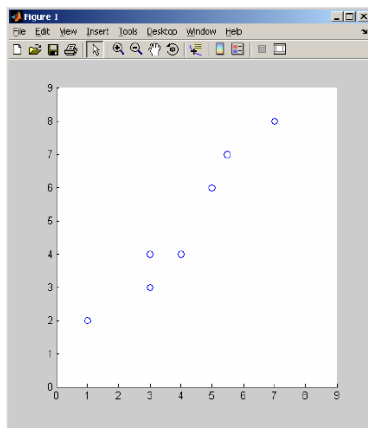
with largest eigen values

13

# 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

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

14

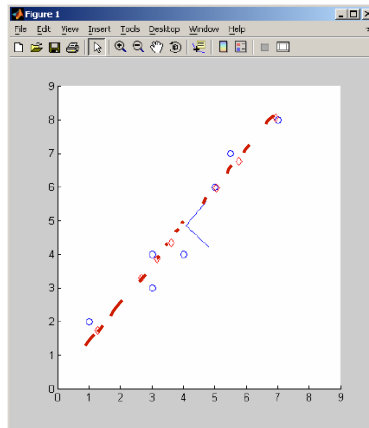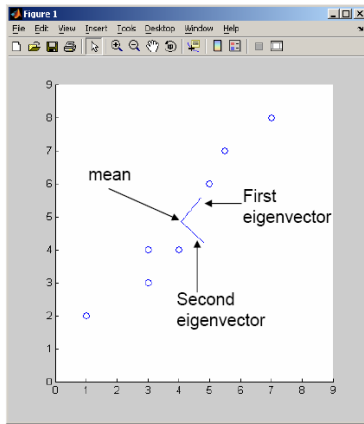# PCA example – reconstruction

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

only used first principal component

# Eigenfaces [Turk, Pentland '91]

- Input images:

- Principal components:

# Eigenfaces reconstruction

- Each image corresponds to adding 8 principal components:

$k=0$

# 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 of $\Sigma$ by just looking at $X_c$
  - great implementations available, e.g., python, R, Matlab svd

# SVD

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

---

# PCA using SVD algoritm

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