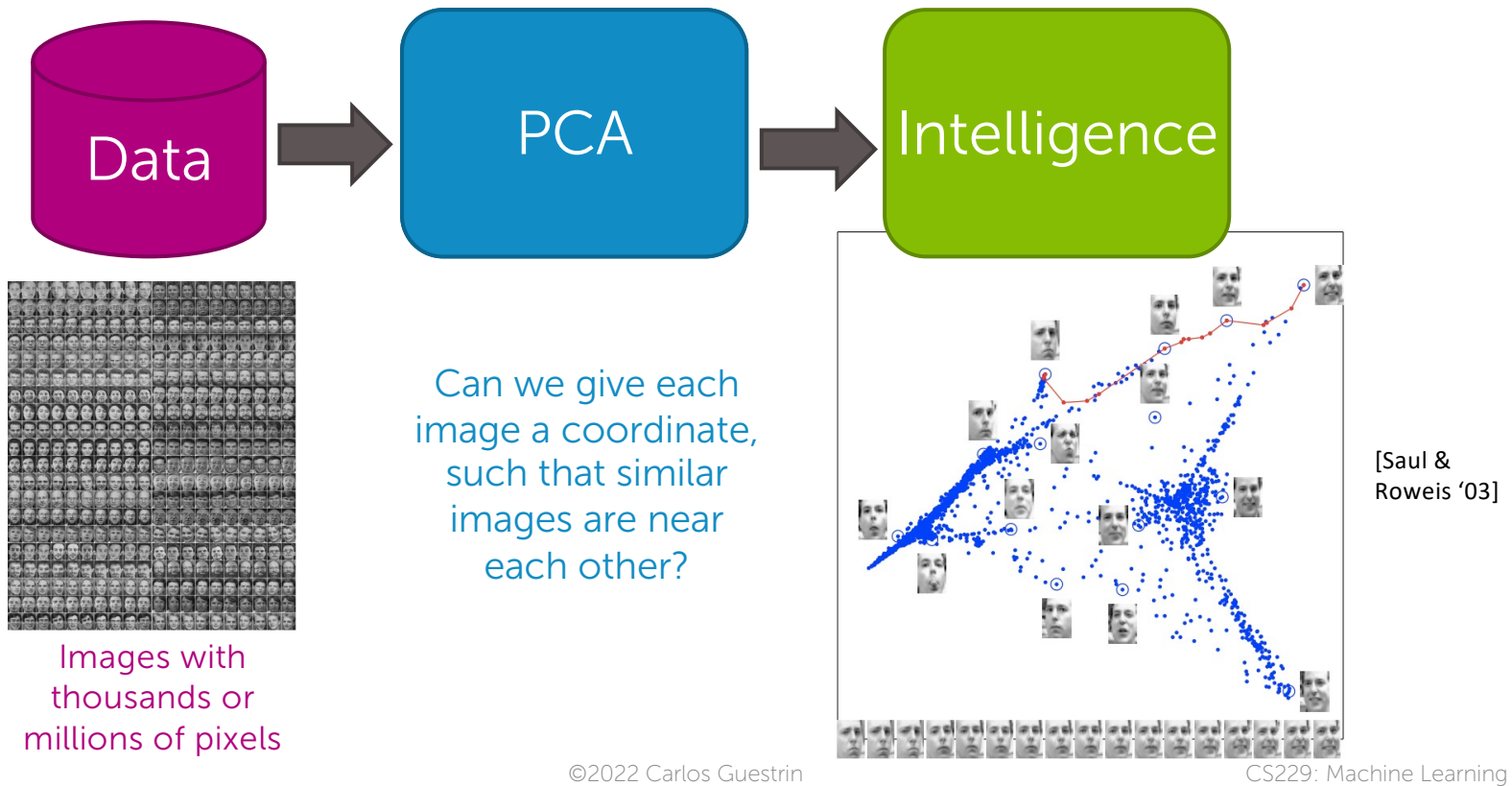# Dimensionality Reduction Principal Component Analysis (PCA)

CS229: Machine Learning
Carlos Guestrin
Stanford University

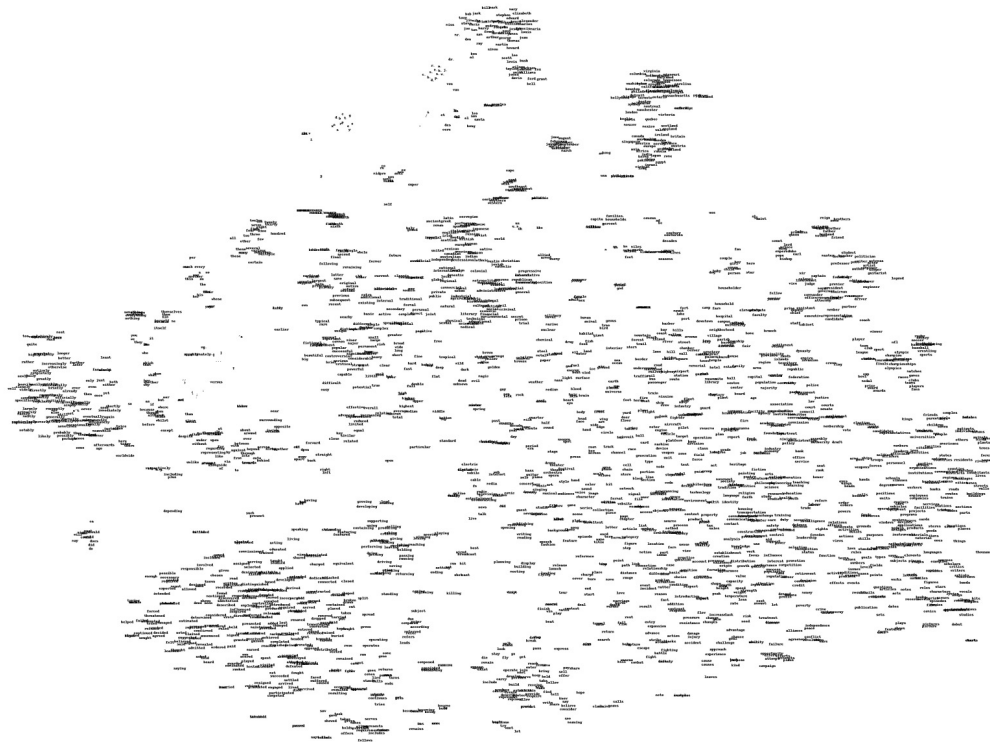Slides include content developed by and co-developed with Emily Fox

# Embedding
## Example: Embedding images to visualize data



Data → PCA → Intelligence
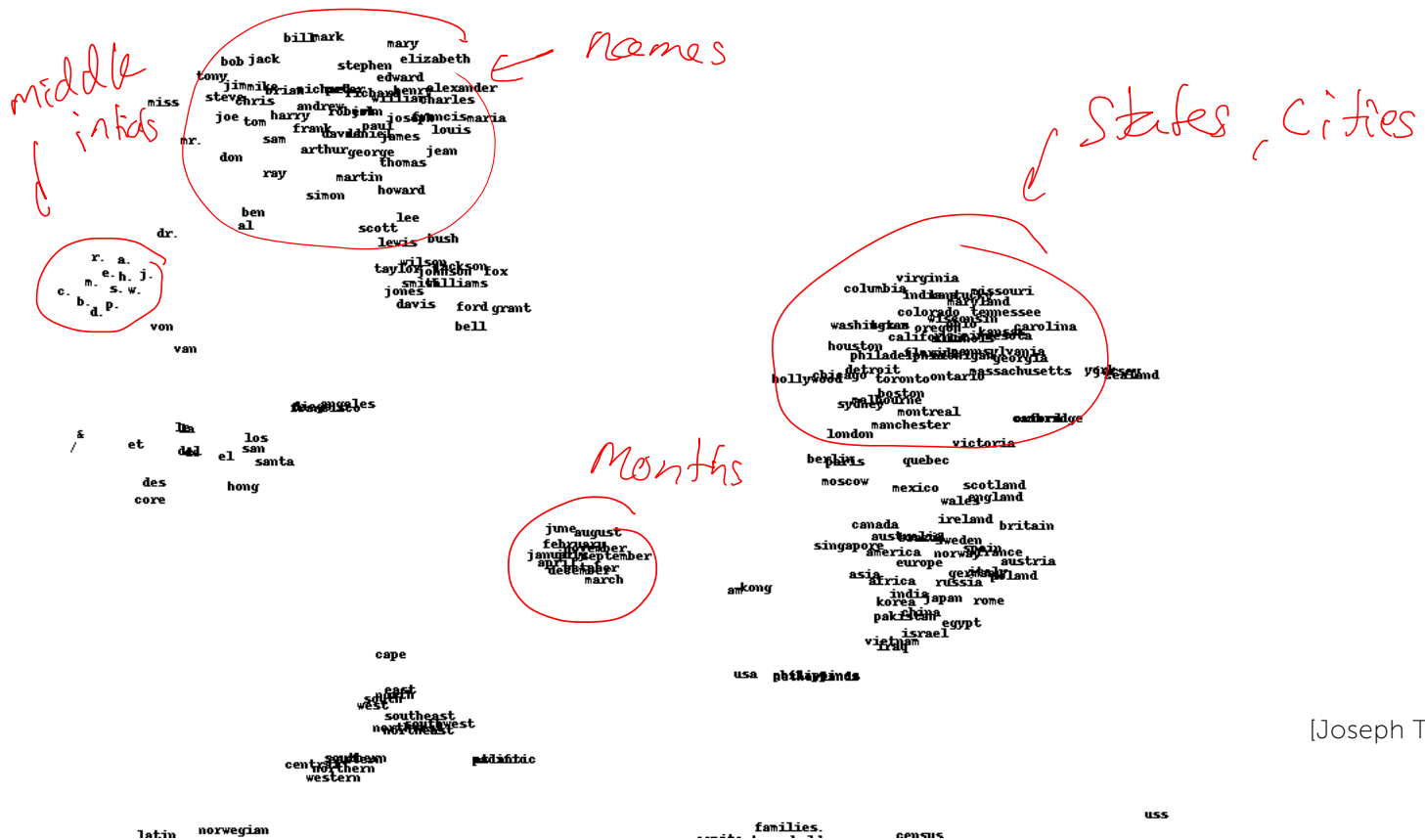
Images with thousands or millions of pixels

Can we give each image a coordinate, such that similar images are near each other?

[Saul & Roweis '03]

CS229: Machine Learning

# Embedding words

10,000 word dictionary



[Joseph Turian 2008]

# Embedding words (zoom in)



bill mark          mary
bob jack       stephen  elizabeth
tony                    edward
jimmie                  henry alexander
steve chris  andrew  richard charles
joe  tom  harry  robert  joseph francis maria
     frank     paul  david daniel james  louis
        sam
     don    arthur george       jean
        ray          thomas
              martin  howard
        simon

ben
al
        lee
     scott  bush
     lewis

middle intls

names

States, Cities

miss

dr.

r.    a.
  e. h. j.
m.   s. w.
c.  b.  p.
    d.
        von
    van

wilson jackson
taylor johnson  fox
       smith williams
       jones
       davis    ford grant
              bell

virginia
columbia indiana missouri
          kentucky maryland
        colorado tennessee
washington oregon   carolina
    california    alaska
houston   pennsylvania
philadelphia   georgia
     detroit
hollywood chicago toronto ontario massachusetts ygukean zealand
        boston
     melbourne
  sydney  montreal
    manchester       oxford bridge
  london        victoria

los angeles

&
  et          los
/      da del     san
      el      santa
  des
  core      hong

berlin         quebec

moscow    mexico   scotland
              wales england

Months

june  august
  february
january september
april october
     december
       march

am kong

canada    ireland  britain
     australia
        sweden
singapore america  norway france
    europe         austria
asia             germany holland
africa   russia italy
    india  japan  rome
  korea china
pakistan    egypt
      israel
   vietnam iraq

cape

east
south
west
  southeast
 northwest southwest
     northeast

central southern atlantic
  northern
   western

usa  philippines

[Joseph Turian 2008]

latin   norwegian              families.                    uss

census
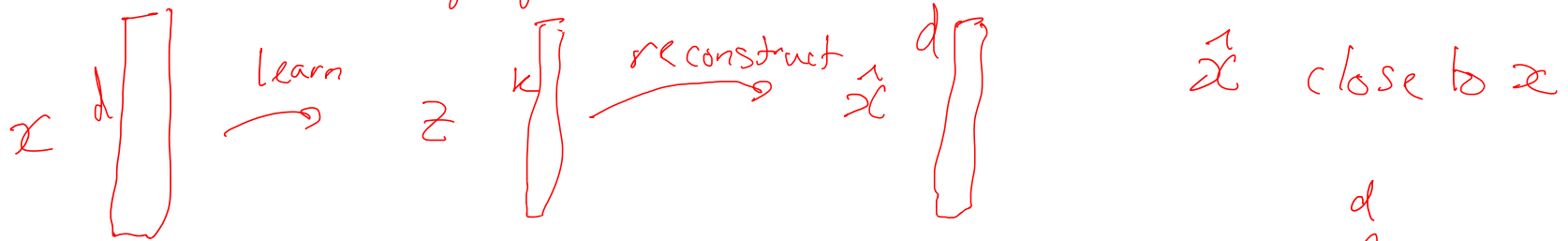
CS229: Machine Learning

# Dimensionality reduction

- Input data may have thousands or millions of dimensions!
  - e.g., text data
- **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

$$d >> k$$

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

PCA: linear projections: $z_1 = 2.5 x_1 - 3.2 x_2 + 3.7 x_3 \cdots$

$x$ [$d$] $\xrightarrow{\text{learn}}$ $z$ [$k$] $\xrightarrow{\text{reconstruct}}$ $\hat{x}$ [$d$]   $\hat{x}$ close to $x$

- Let's see this in the unsupervised setting   $z = \begin{bmatrix} k \\ A \end{bmatrix}^{d} x$
  - just **x**, but no y

learn $\longrightarrow$ projection matrix

# Linear projection and reconstruction

**Project** onto 1-dimension

Reconstruction: Only knowing z, what was $(x_1, x_2)$?

*#awful*

*#awesome*

©2022 Carlos Guestrin    CS229: Machine Learning

# What if we project onto d vectors?

$$\vec{x}' = z_1' \vec{u}_1 + z_2' \vec{u}_2 \quad \text{(ignoring offset)}$$



Perfect reconstruction!

# If I had to choose one of these vectors, which do I prefer?



best projection direction
on avg, lower
reconstruction
error

$\vec{u}_2$

$\vec{u}_1$

$\vec{u}_1$ reconstruction error

$\vec{u}_2$ reconstruction error

#awful

#awesome

©2022 Carlos Guestrin

CS229: Machine Learning

# Principal component analysis (PCA) – Basic idea

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

- Choose projection with minimum reconstruction error

# "PCA explained visually"

http://setosa.io/ev/principal-component-analysis/

# Linear projections, a review

- Project a point into a (lower dimensional) space:
  - **point**: $x = (x_1,...,x_d)$ $\overbrace{10,000}$ dims
  - **select a basis** − set of basis vectors − $(u_1,...,u_k)$ $\overbrace{100\ dims}$
    - we consider orthonormal basis:
      - $u_i \bullet u_i = 1$, and $u_i \bullet u_j = 0$ for $i \neq j$    $u_j$
  - **select a center** − $\overline{x}$, defines offset of space   mean value of X
  - **best coordinates** in lower dimensional space defined by dot-products:
    $(z_1,...,z_k)$, $z_i = (\overline{x}-x) \bullet u_i$
    - minimum squared error

*Handwritten notes (red):*

Project X into U
coordinates $z_1 \cdots z_k$

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

reconstruction

if $k = d$

$$\hat{x} = x$$

dot product minimizes the squared error of projection

$$z_1 = (x - \overline{x}) \cdot \vec{u}_1$$

$$z_1 = \underset{z}{argmin} \left( (x-\overline{x}) - z\, \vec{u}_1 \right)^2$$

$\overline{w}$   $\overline{x}$   $x - \overline{x}$   $x$

# PCA finds projection that minimizes reconstruction error

*mean vector*

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

*avg x over data*

*from previous slide: projection coordinate is dot product*

*reconstructed $x_i$*

*Projection direction*

$$\hat{\mathbf{x}}^i = \bar{\mathbf{x}} + \sum_{j=1}^{k} z_j^i \mathbf{u}_j \quad \text{and} \quad \bar{\mathbf{x}} = \frac{1}{N}\sum_{i=1}^{N}\mathbf{x}^i \qquad z_j^i = (\mathbf{x}^i - \bar{\mathbf{x}}) \cdot \mathbf{u}_j$$

*coordinate*

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

*Want minimize*
*over choice of $\vec{u}_1, ..., \vec{u}_k$*

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

$x_2$

*learn $\vec{u}_1$ minimize*

$\bar{x}$

$x_1$

# Understanding the reconstruction error

- 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:

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

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

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

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

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

$$= \sum_{i=1}^{N} \left[ \sum_{j=k+1}^{d} z_j^i u_j \right]^2 = \sum_{i=1}^{N} \left[ \sum_{j=k+1}^{d} z_j^i \underset{orthonormal}{u_j \cdot u_j} z_j^i + \sum_{j=k+1}^{d} \sum_{j'\neq j} z_j^i \underset{orthonormal}{u_j \cdot u_{j'}} z_{j'}^i \right]$$

$$= \sum_{i=1}^{N} \sum_{j=k+1}^{d} (z_j^i)^2 \leftarrow \text{minimizing reconstruction error} \equiv \min \text{ squared thrown out coefficients!!}$$

# Reconstruction error and covariance matrix

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

$z_j$ (bracket over the inner sum)

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

push sum over data in

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

$$= N \sum_{j=k+1}^{d} u_j^T \Sigma u_j$$

$\leftarrow$ choose $u_j$ to minimize this error

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

matrix notation

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \ddots \\ \sigma_{21} & \sigma_2^2 & \ddots \\ & & \ddots \end{pmatrix}$$

# Minimizing reconstruction error and eigen vectors

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

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

- Eigen vector:

$u_j$ is an eigen vector

$$u_j^T \Sigma u_j = u_j^T \lambda_j u_j$$
$$= \lambda_j u_j^T u_j$$ — orthonormal
$$= \lambda_j$$

Memory lane:
$$\Sigma u = \lambda u \leftarrow \text{eigen vector}$$
$$\nwarrow \text{eigen value}$$

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

$$\min_{u_1...u_k} error_k$$

$\equiv$ throwing out $u_{d+1}...u_d$ with smallest eigen values of $\Sigma$

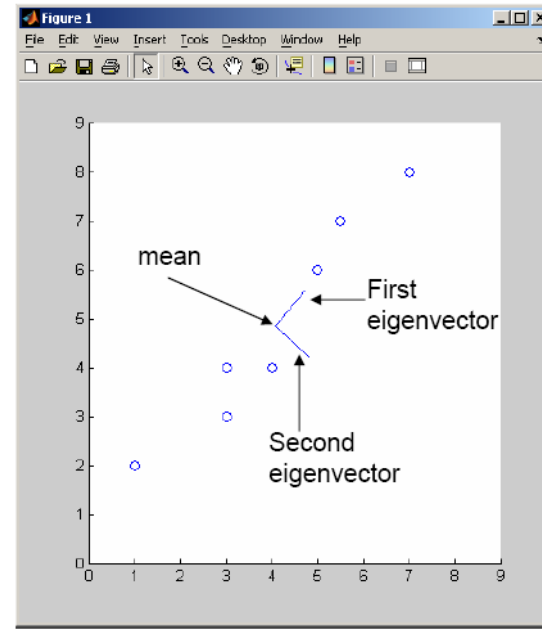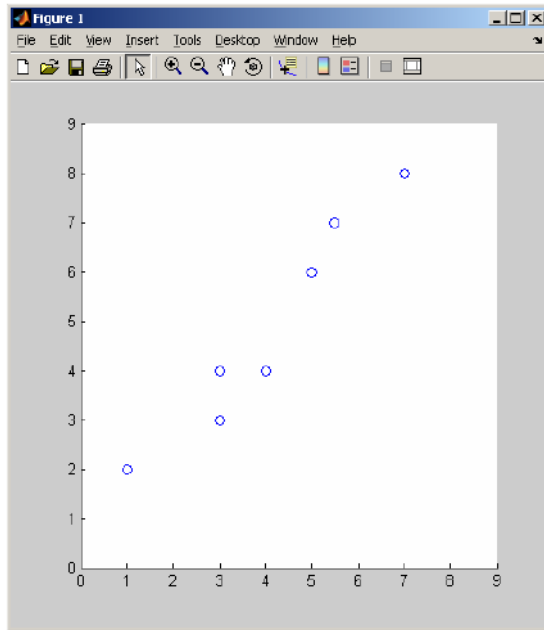$\equiv$ keeping $u_1,...u_k$ with highest eigen values

# Basic PCA algoritm

- Start from N by d 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

$$X_c = N \left( \underset{d}{\bigcirc} \right) x^j - \overline{x}$$

# PCA example

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

: Machine Learning

# PCA example – reconstruction

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

only used first principal component

: Machine Learning

# Eigenfaces [Turk, Pentland '91]

- Input

■ Principal components:



*Glasses*

# 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., python, R, 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