# Dimensionality Reduction
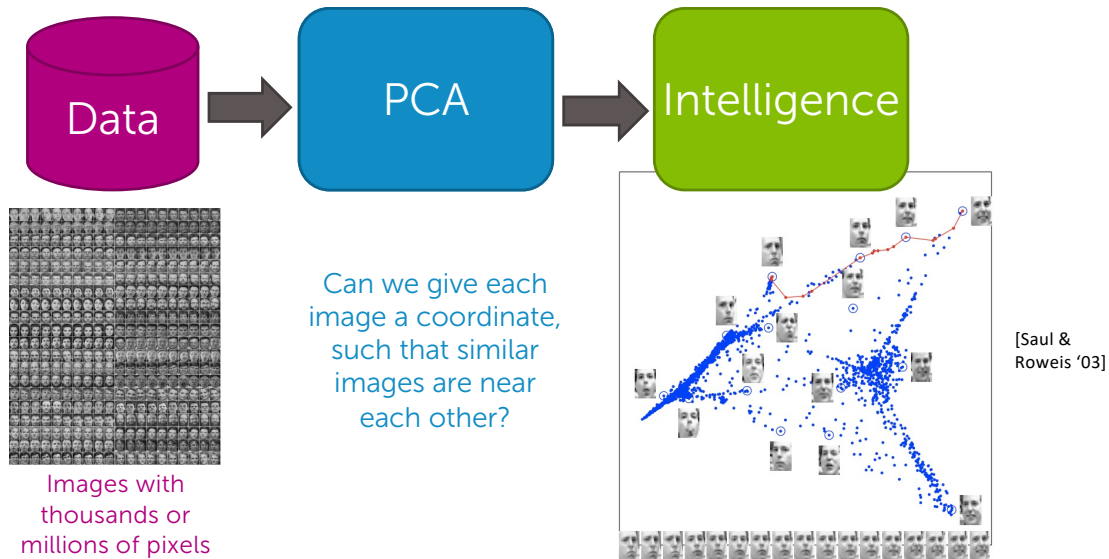# Principal Component Analysis (PCA)

CS229: Machine Learning
Carlos Guestrin
Stanford University
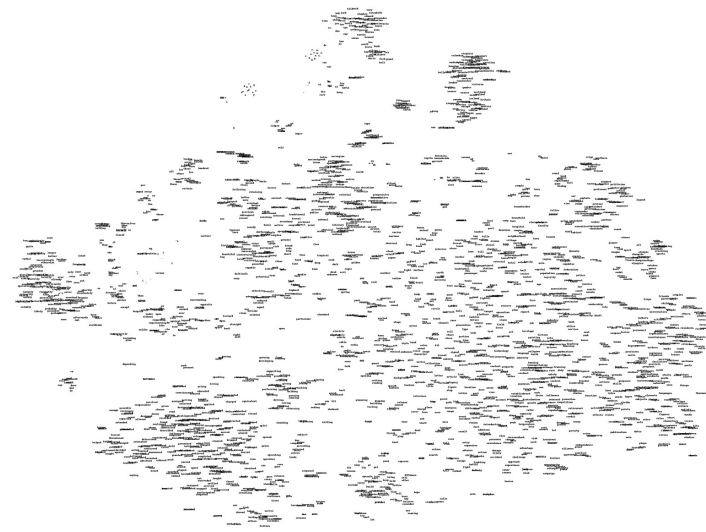
Slides include content developed by and co-developed with Emily Fox

# Embedding
## Example: Embedding images to visualize data
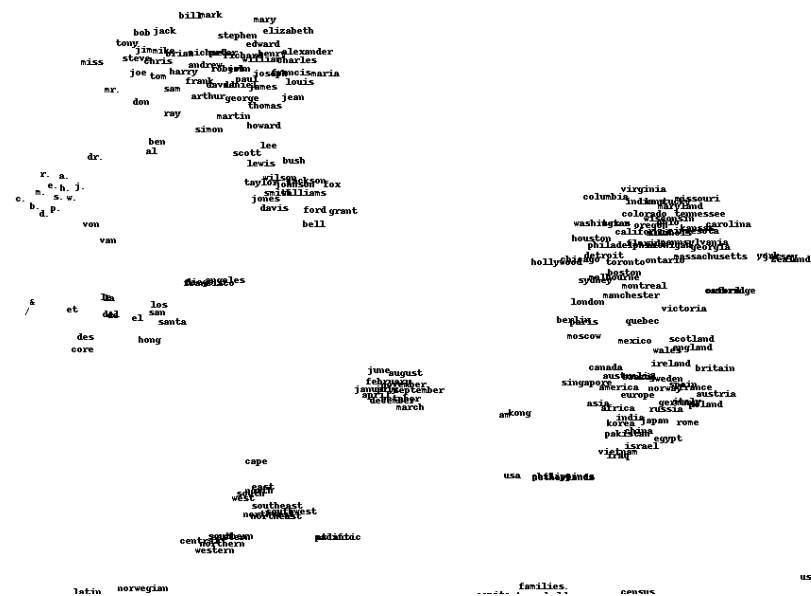


Data → PCA → Intelligence

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]

# Embedding words



[Joseph Turian 2008]

CS229: Machine Learning

# Embedding words (zoom in)
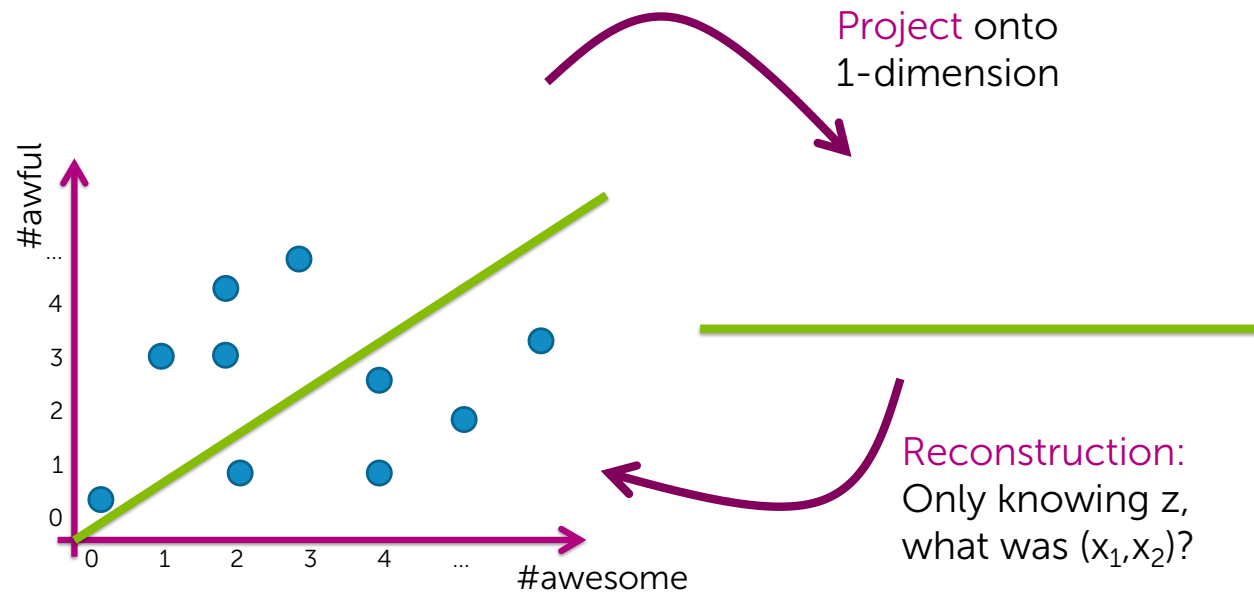


[Joseph Turian 2008]

# 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
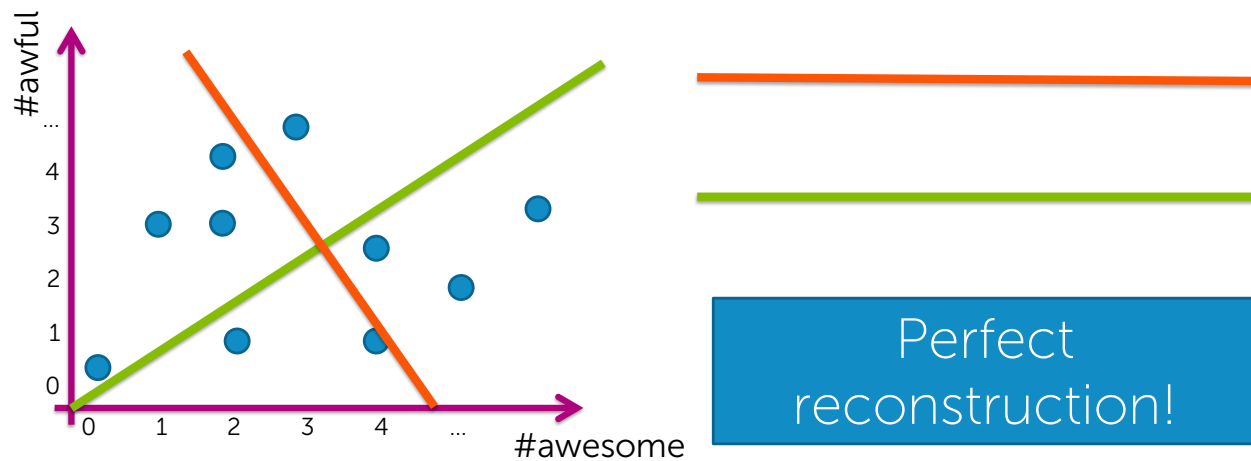    - high dimensional data that is truly lower dimensional

# Lower dimensional projections

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

- Let's see this in the unsupervised setting
  - just **x**, but no y

©2022 Carlos Guestrin                    CS229: Machine Learning

# Linear projection and reconstruction

Project onto
1-dimension

#awful

#awesome

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

# What if we project onto d vectors?



Perfect reconstruction!

©2022 Carlos Guestrin    CS229: Machine Learning

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

©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/

CS229: Machine Learning

# Linear projections, a review

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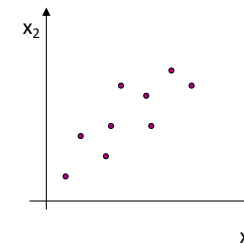    - minimum squared error

   CS229: Machine Learning

# 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:

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

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

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

©2022 Carlos Guestrin                    CS229: Machine Learning

# Understanding the reconstruction error

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

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

- Rewriting error:

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

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

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

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

CS229: Machine Learning

# Reconstruction error and covariance matrix

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

CS229: Machine Learning

# Minimizing reconstruction error and eigen vectors

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

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

- Eigen vector:

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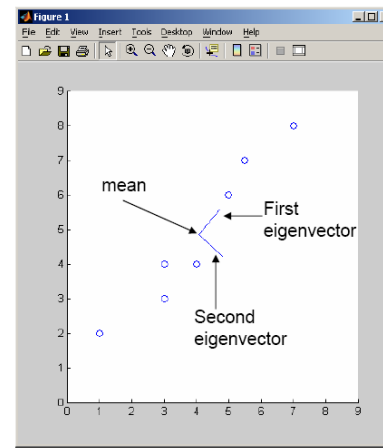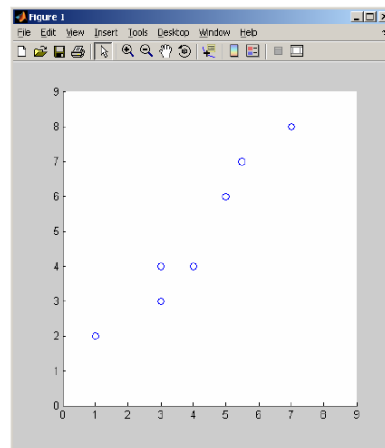
# Basic PCA algoritm

- Start from 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
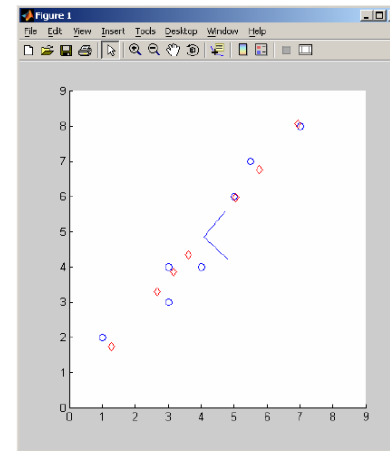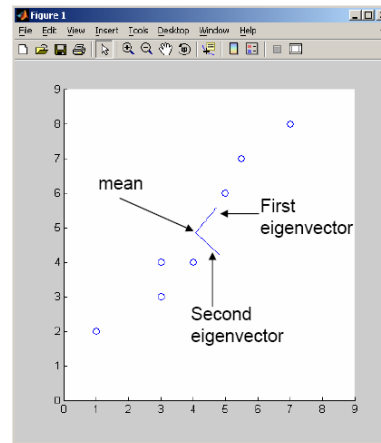
# PCA example

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

©2022 Carlos Guestrin · : 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

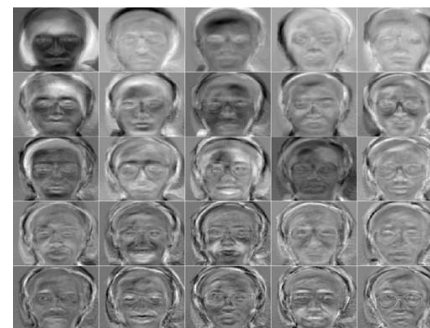©2022 Carlos Guestrin                    : Machine Learning

# Eigenfaces [Turk, Pentland '91]

- Input



- Principal components:

# Eigenfaces reconstruction

- Each image corresponds to adding 8 principal
  components:

©2022 Carlos Guestrin

CS229: Machine Learning

# 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 $\mathbf{X}$
- **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