

Dimensionality Reduction

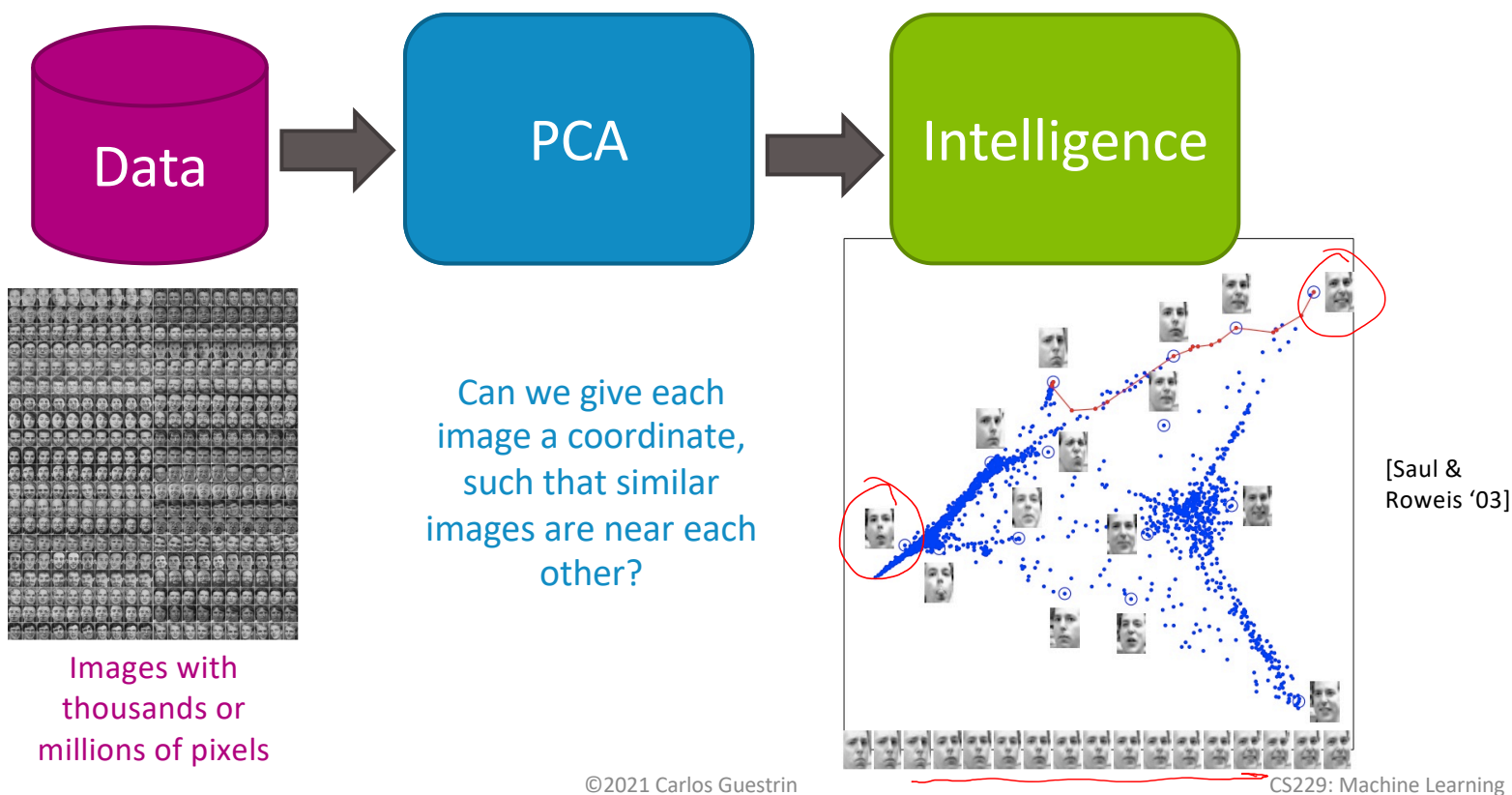
Principal Component Analysis (PCA)

CS229: Machine Learning
Carlos Guestrin
Stanford University

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

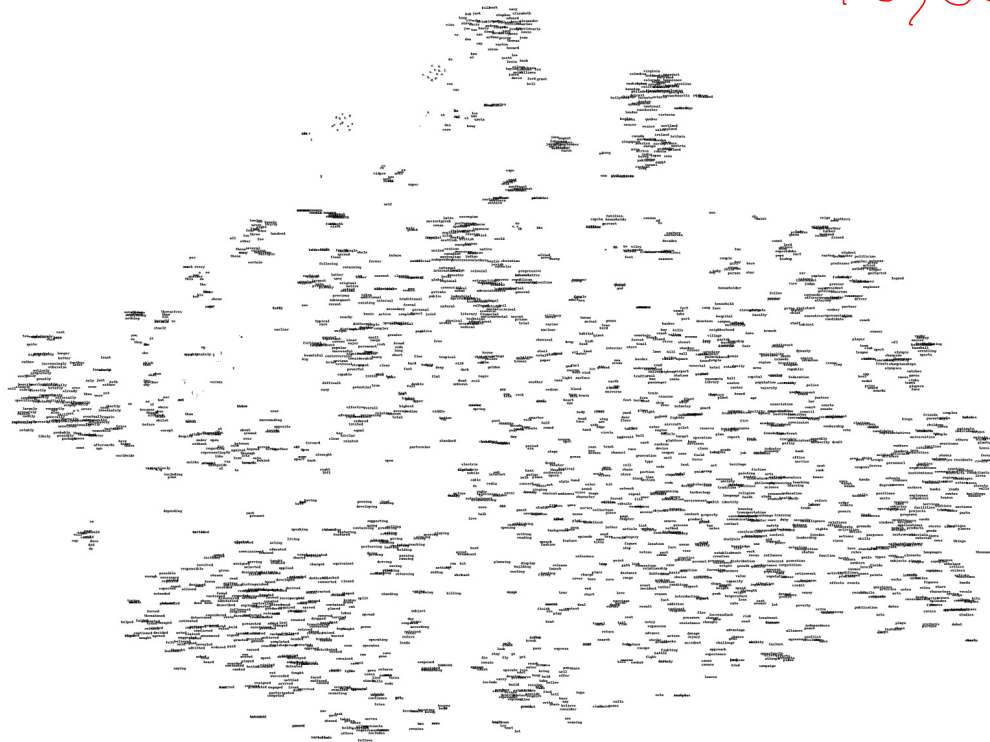
Embedding

Example: Embedding images to visualize data



Embedding words

10,000 words



[Joseph Turian]

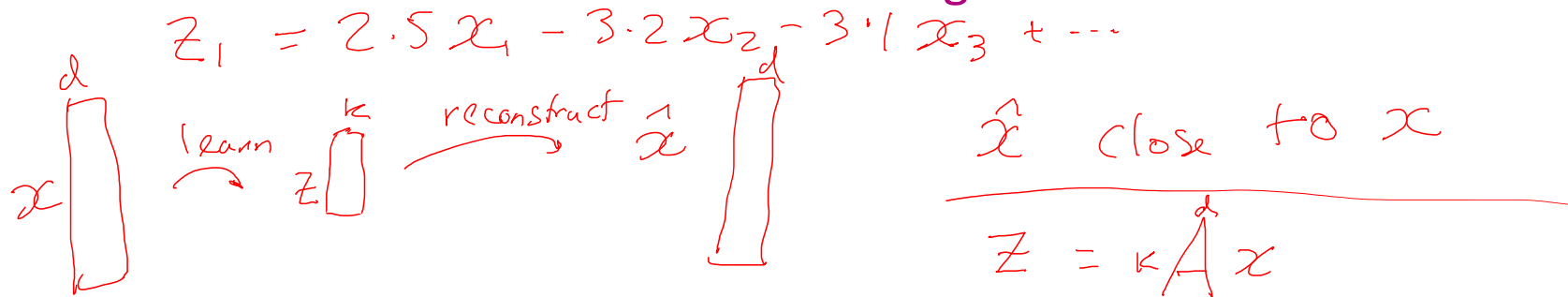
Dimensionality reduction

- 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

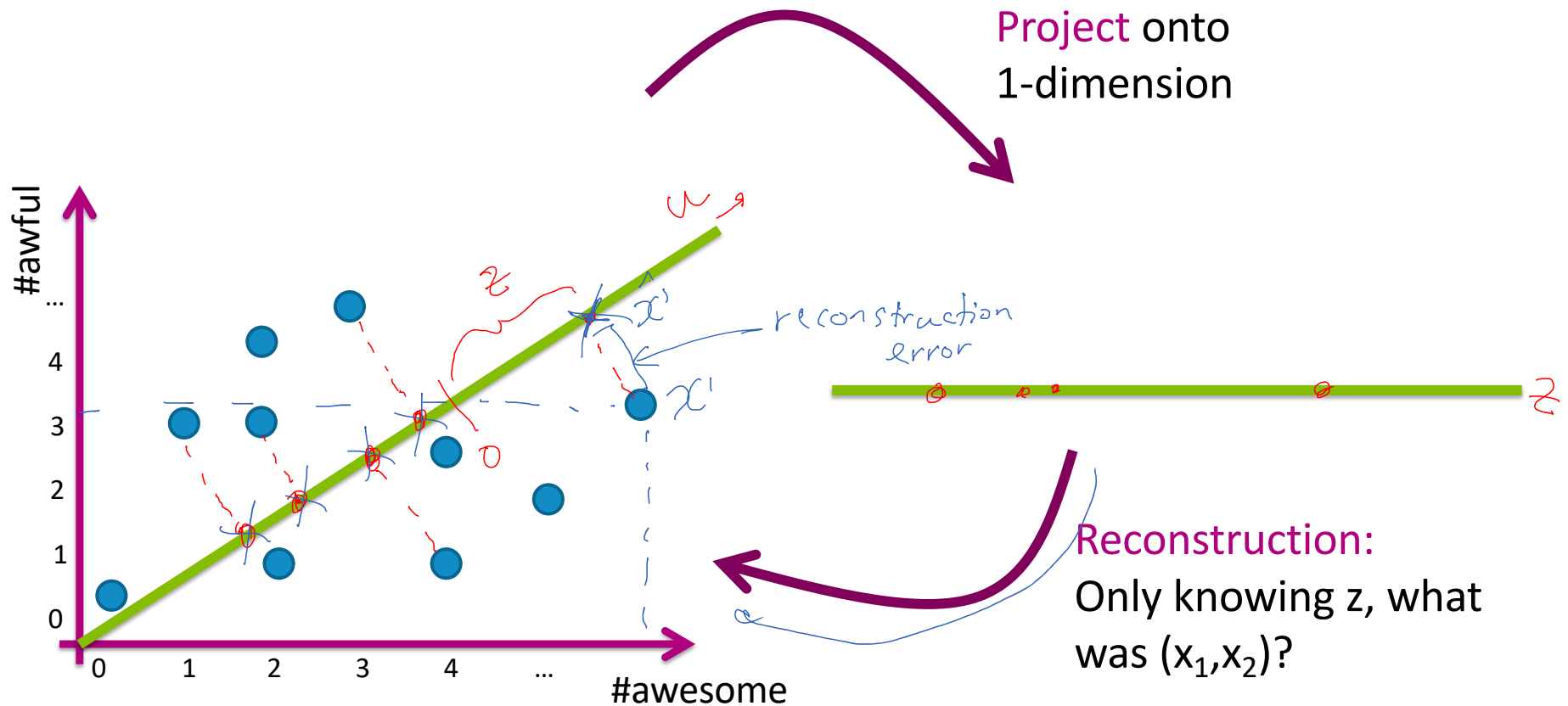
$$d \gg k$$

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



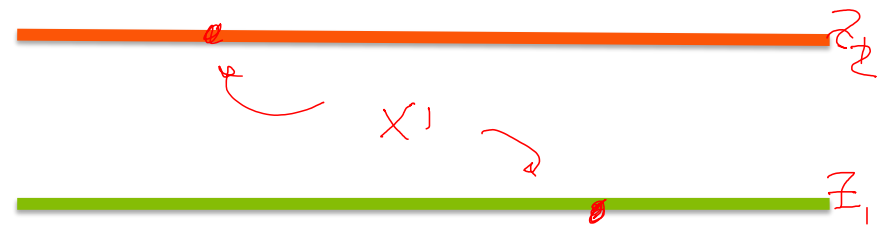
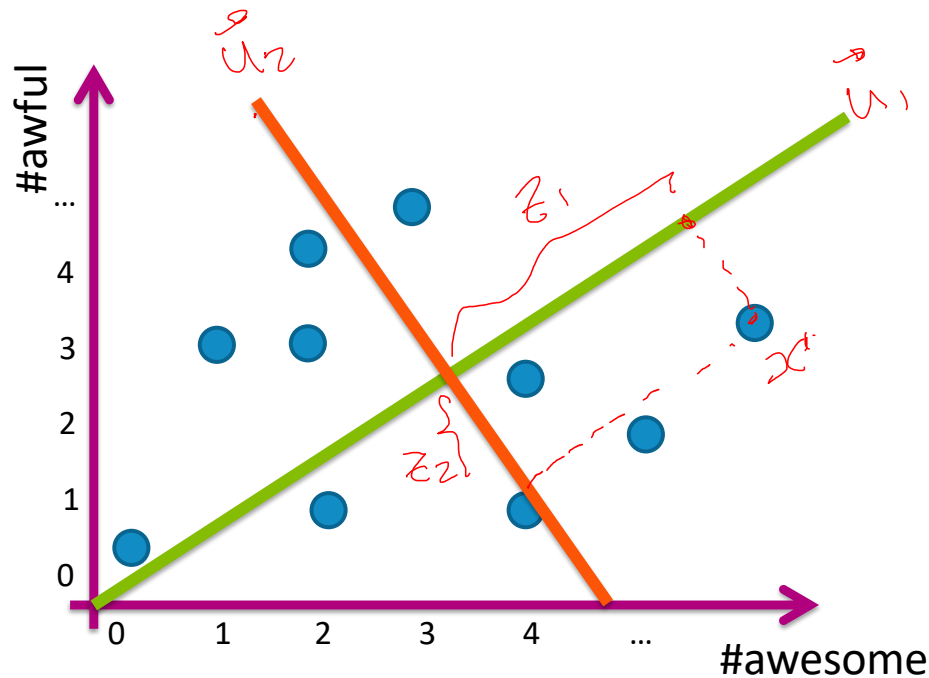
- Let's see this in the unsupervised setting
 - just x , but no y

Linear projection and reconstruction



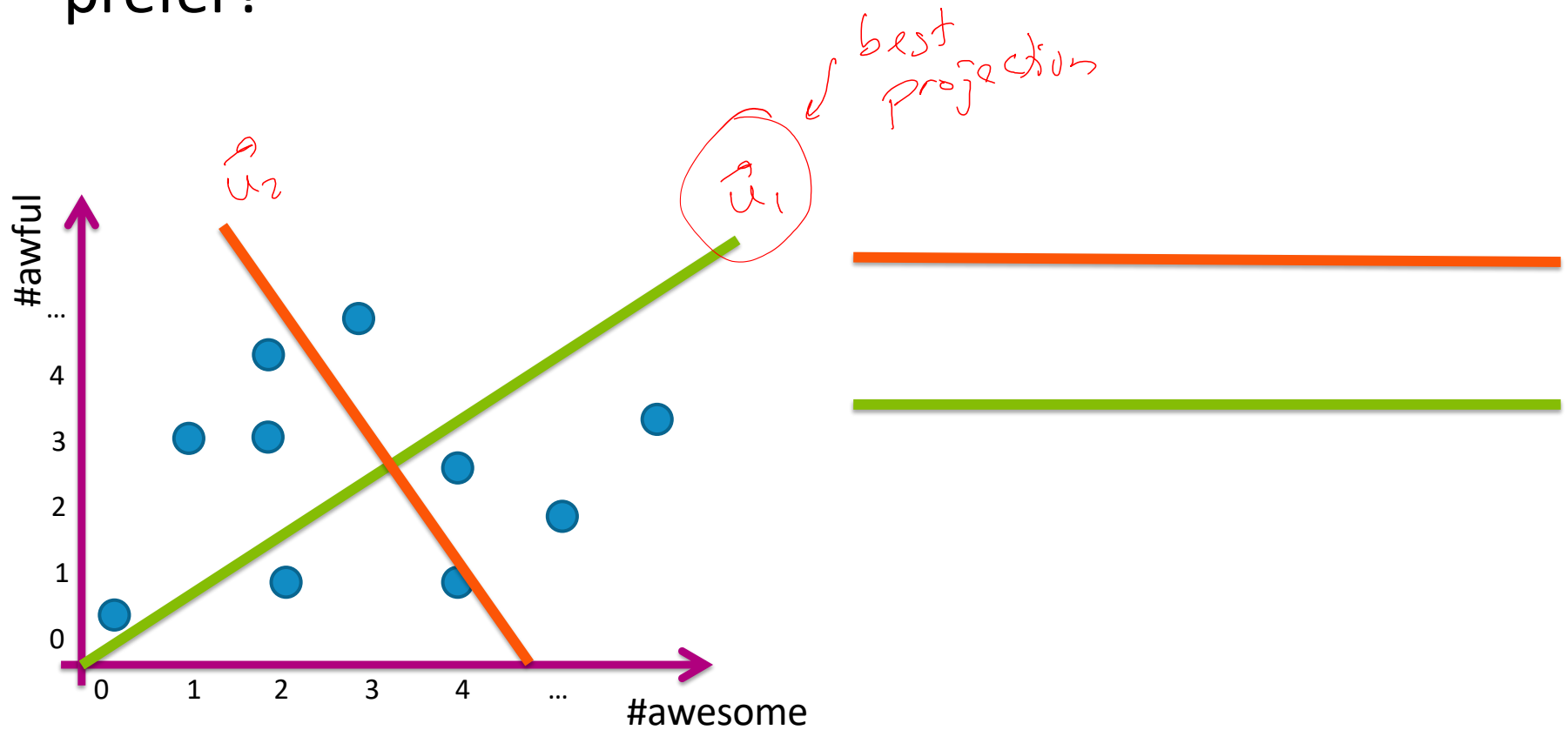
What if we project onto d vectors?

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



Perfect reconstruction!

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



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/>

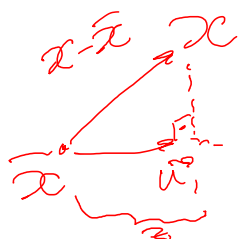
Linear projections, a review

- Project a point into a (lower dimensional) space:
 - **point:** $\mathbf{x} = (x_1, \dots, x_d)$ *(0,0,0,0)*
 - **select a basis** – set of basis vectors – $(\mathbf{u}_1, \dots, \mathbf{u}_k)$ *(0,0)*
 - we consider orthonormal basis: *\vec{u}_1*
 - $\mathbf{u}_i \bullet \mathbf{u}_i = 1$, and $\mathbf{u}_i \bullet \mathbf{u}_j = 0$ for $i \neq j$ *\vec{u}_2*
 - **select a center** – $\bar{\mathbf{x}}$, defines offset of space *mean of \mathcal{X}*
 - **best coordinates** in lower dimensional space defined by dot-products:
 - (z_1, \dots, z_k) , $z_i = (\mathbf{x} - \bar{\mathbf{x}}) \bullet \mathbf{u}_i$
 - minimum squared error

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

if $k = d$:

$$\hat{\mathbf{x}} = \mathbf{x}$$



$$z_1 = (\mathbf{x} - \bar{\mathbf{x}}) \bullet \vec{u}_1 \quad \left. \vphantom{z_1} \right\} z_1 = \underset{z}{\operatorname{arg\,min}} ((\mathbf{x} - \bar{\mathbf{x}}) - z \vec{u}_1)^2$$

PCA finds projection that minimizes reconstruction error

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

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

index of data point (pointing to i in z_j^i)

projection direction (under \mathbf{u}_j)

avg over data (under $\bar{\mathbf{x}}$)

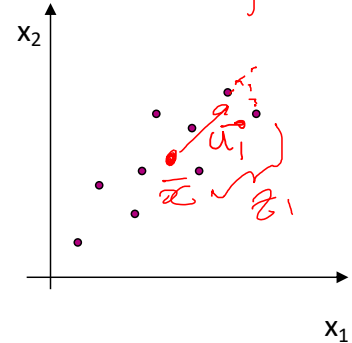
$$z_j^i = (\mathbf{x}^i - \bar{\mathbf{x}}) \cdot \mathbf{u}_j$$

projection from previous slide

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

want to min $\mathbf{u}_1, \dots, \mathbf{u}_k$

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



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:

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

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

Given $k \ll d$, find $(\mathbf{u}_1, \dots, \mathbf{u}_k)$ minimizing reconstruction error:

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

$$\begin{aligned}
 \min_{\mathbf{u}} \text{error}_k &= \sum_{i=1}^N (\mathbf{x}^i - \hat{\mathbf{x}}^i)^2 = \sum_{i=1}^N \left[\bar{\mathbf{x}} + \sum_{j=1}^d z_j^i \mathbf{u}_j - \left(\bar{\mathbf{x}} + \sum_{j=1}^k z_j^i \mathbf{u}_j \right) \right]^2 = \sum_{i=1}^N \left[\sum_{j=k+1}^d z_j^i \mathbf{u}_j \right]^2 \\
 &= \sum_{i=1}^N \left[\sum_{j=k+1}^d z_j^i \mathbf{u}_j \cdot \mathbf{u}_j z_j^i + \sum_{j=k+1}^d \sum_{j' \neq j} z_j^i \mathbf{u}_j \cdot \mathbf{u}_{j'} z_{j'}^i \right] \\
 &= \sum_{i=1}^N \sum_{j=k+1}^d (z_j^i)^2 \leftarrow \text{minimizing reconstruction error} \equiv \text{min square of throwa out coefficients } \begin{matrix} | \\ \vdots \\ | \end{matrix}
 \end{aligned}$$

Reconstruction error and covariance matrix

$$\begin{aligned}
 \text{error}_k &= \sum_{i=1}^N \sum_{j=k+1}^d \overbrace{[u_j \cdot (x^i - \bar{x})]^2}^{z_j} \\
 &= \sum_{i=1}^N \sum_{j=k+1}^d u_j^T (x^i - \bar{x}) (x^i - \bar{x})^T u_j \\
 &\quad \text{push sum over } i \text{ in, because } u_j \text{ doesn't depend on } i \\
 &= \sum_{j=k+1}^d u_j^T \left[\sum_{i=1}^N (x^i - \bar{x}) (x^i - \bar{x})^T \right] u_j \\
 &\quad \underbrace{\sum_{i=1}^N (x^i - \bar{x}) (x^i - \bar{x})^T}_{N \Sigma} \\
 &= N \sum_{j=k+1}^d u_j^T \Sigma u_j
 \end{aligned}$$

choose u_j that minimizes this error

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

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

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

matrix notation

Minimizing reconstruction error and eigen vectors

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

$$\text{error}_k = \frac{1}{N} \sum_{j=k+1}^d \mathbf{u}_j^T \Sigma \mathbf{u}_j = N \sum_{j=k+1}^d \lambda_j$$

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

- Eigen vector:

memory lane:

$$\Sigma \mathbf{u} = \lambda \mathbf{u}$$

↖ eigen vector ↘ eigen value

if \mathbf{u}_j is an eigen vector

$$\mathbf{u}_j^T \Sigma \mathbf{u}_j = \lambda_j \mathbf{u}_j^T \mathbf{u}_j = \lambda_j$$

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

ignored dimensions

Min error_k \equiv throwing out $\mathbf{u}_{k+1}, \dots, \mathbf{u}_d$ with smallest eigen values of Σ
 $\mathbf{u}_1, \dots, \mathbf{u}_k$ \equiv keeping $\mathbf{u}_1, \dots, \mathbf{u}_k$ with largest eigen values of Σ

Basic PCA algorithm

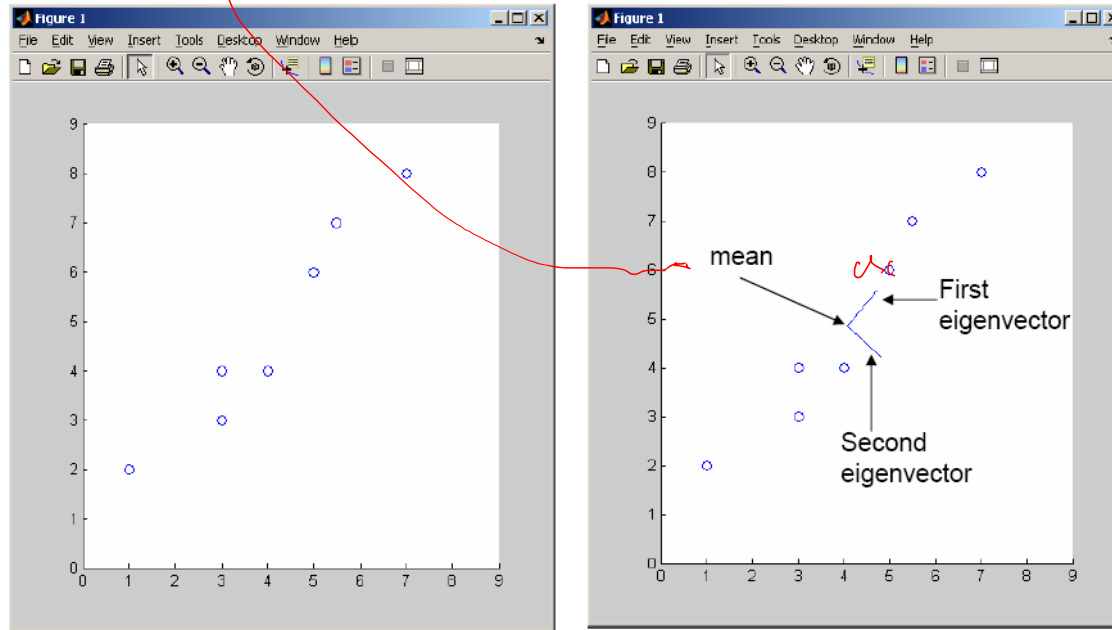
- Start from N by d data matrix \mathbf{X}
- **Recenter:** subtract mean from each row of \mathbf{X}
 - $\mathbf{X}_c \leftarrow \mathbf{X} - \bar{\mathbf{X}}$
- **Compute covariance matrix:**
 - $\Sigma \leftarrow 1/N \mathbf{X}_c^T \mathbf{X}_c$
- Find **eigen vectors and values** of Σ
- **Principal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

$$\mathbf{X}_c = N \left(\begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} \right) x^i - \bar{x}$$

The image shows a handwritten red equation: $\mathbf{X}_c = N \left(\begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix} \right) x^i - \bar{x}$. The number 'k' is written above the opening parenthesis. The three horizontal lines inside the parentheses represent rows of a matrix. A red arrow points from the opening parenthesis to the first line, and another red arrow points from the closing parenthesis to the second line. To the right of the parentheses is the expression $x^i - \bar{x}$.

PCA example

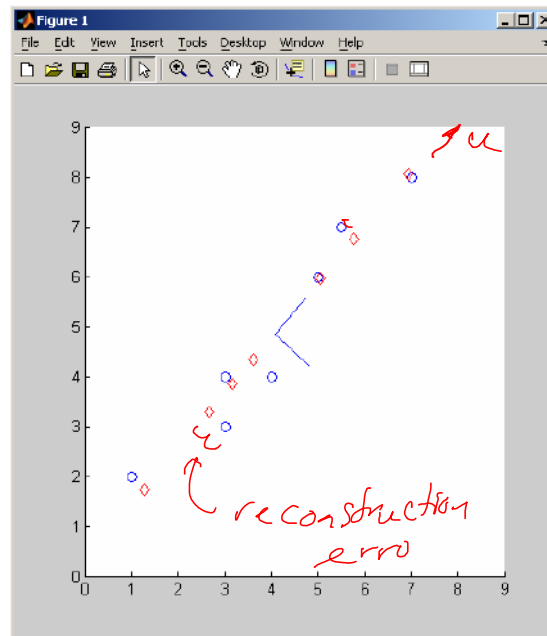
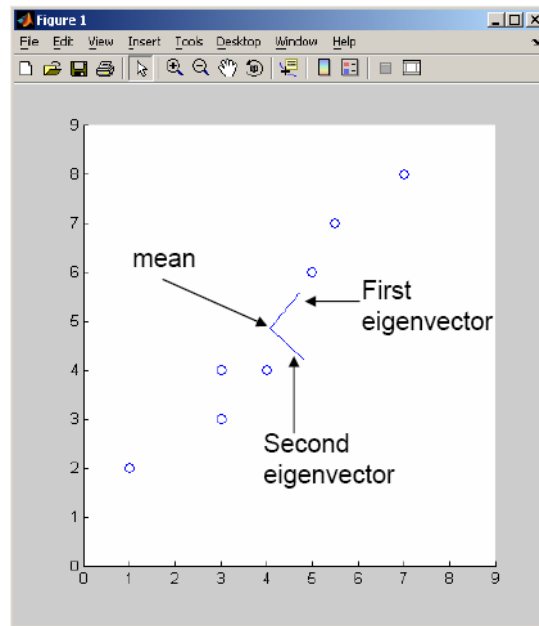
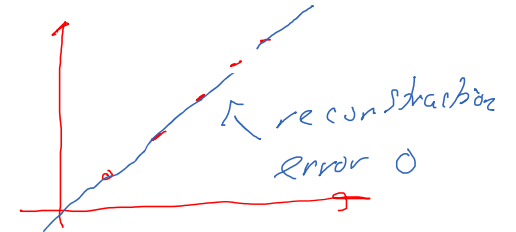
$$\hat{\mathbf{x}}^i = \bar{\mathbf{x}} + \sum_{j=1}^k z_j^i \mathbf{u}_j$$



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:



λ_1
short
or
longer

glasses

Eigenfaces reconstruction

- Each image corresponds to adding 8 principal components:



target face x

Scaling up

- Covariance matrix can be really big!
 - Σ is d by d *← 100 000 000 entries*
 - Say, only 10000 features
 - finding eigenvectors is very slow...
- Use singular value decomposition (SVD)
 - finds to k eigenvectors, *without forming Σ explicitly*
 - 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 λ_j
 - $\mathbf{V}^T \leftarrow$ singular vector matrix
 - in our setting each row is eigenvector \mathbf{v}_j

PCA using SVD algorithm

- Start from m by n data matrix \mathbf{X}
- **Recenter:** subtract mean from each row of \mathbf{X}
 - $\mathbf{x}_c \leftarrow \mathbf{x} - \bar{\mathbf{x}}$
- 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