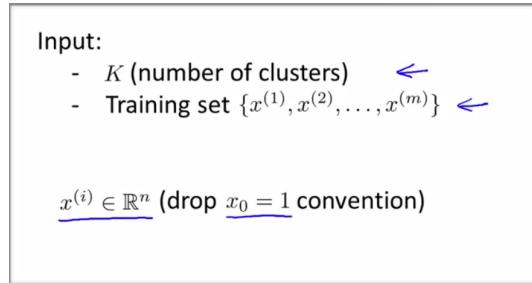
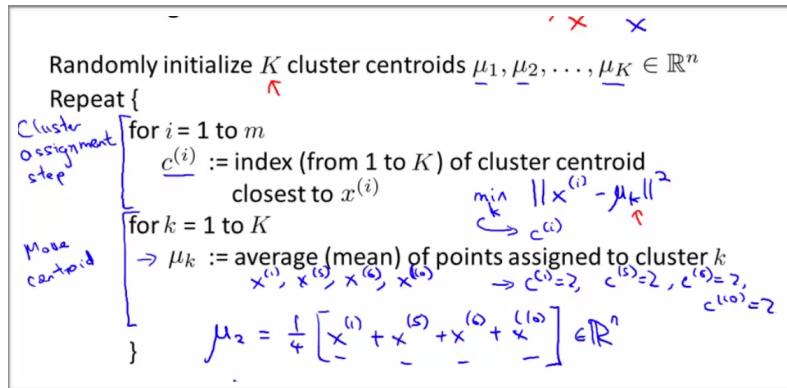


- **Unsupervised Learning Introduction**
- Supervised Learning
  - Find the decision boundary
- Unsupervised learning
  - no labels
  - give unlabeled training set to the algorithm and find some structure to the data
  - clustering
- **Clustering K-Means**
- Cluster centroids
  - the number of groups that are wanted
- Iterative process
  - compute the average of all the points and reassign the centroid centers
- K means



- K-Means algorithm



- **Optimization Objective or Cost Function**
- K-means optimization objective

$\rightarrow c^{(i)} =$  index of cluster ( $1, 2, \dots, K$ ) to which example  $x^{(i)}$  is currently assigned

$\rightarrow \mu_k =$  cluster centroid  $k$  ( $\mu_k \in \mathbb{R}^n$ )  $\leftarrow k \in \{1, 2, \dots, K\}$

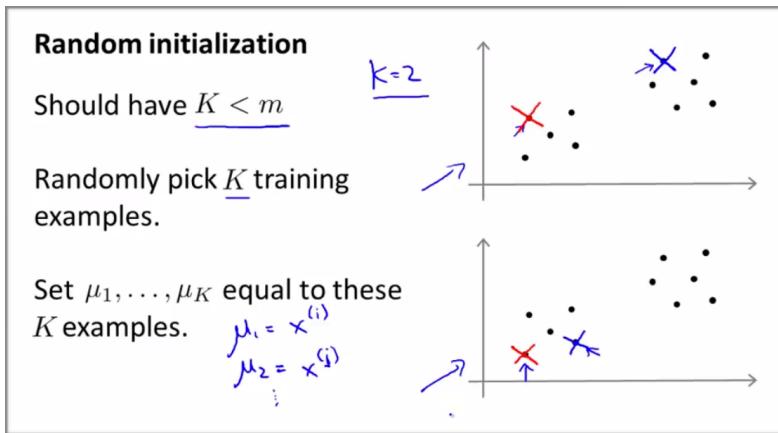
$\mu_{c^{(i)}} =$  cluster centroid of cluster to which example  $x^{(i)}$  has been assigned  $x^{(1)} \rightarrow 1 \quad c^{(1)}=1 \quad \mu_{c^{(1)}} = \mu_1$

Optimization objective:

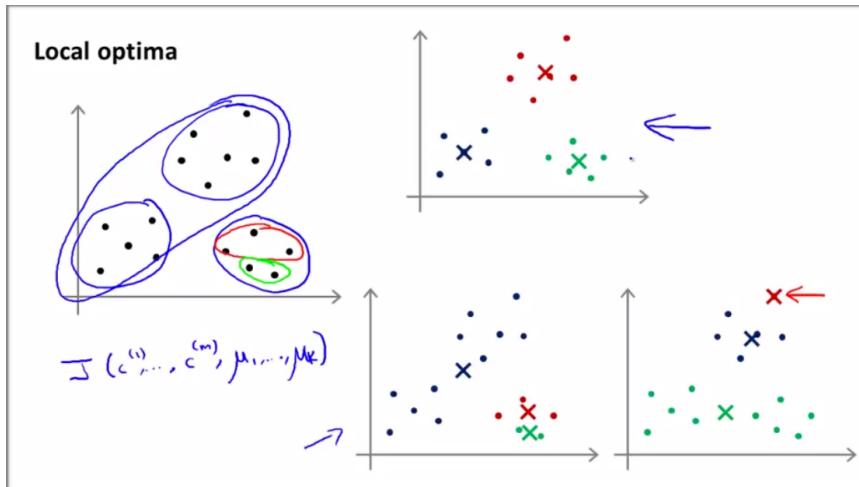
$$\rightarrow J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2$$

$$\min_{\substack{c^{(1)}, \dots, c^{(m)}, \\ \mu_1, \dots, \mu_K}} J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$$

- Random initialization
- Random initialization of cluster centers



- Local optima



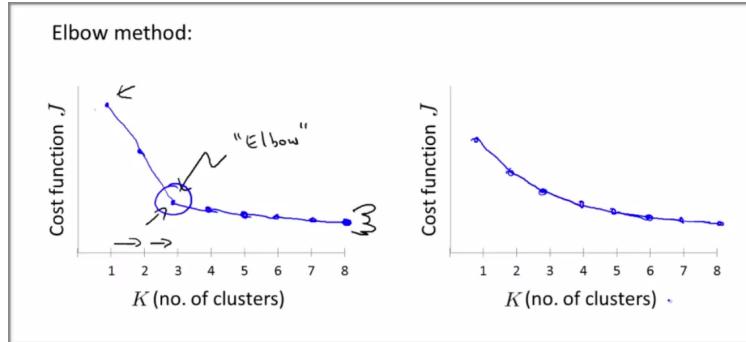
- To prevent this run k means multiple times
- Random initialization

```
For i = 1 to 100 { 50 - 1000
    → Randomly initialize K-means.
    Run K-means. Get  $c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K$ .
    Compute cost function (distortion)
    →  $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$ 
}
```

Pick clustering that gave lowest cost  $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$

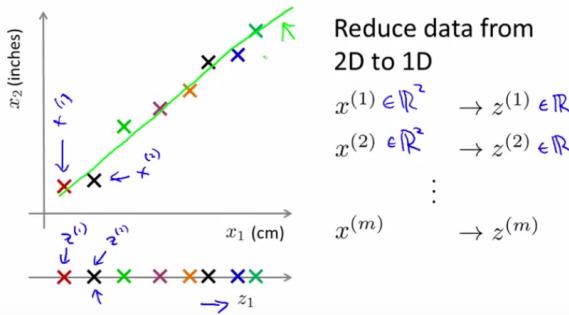
*K = 2 - 10*

- If key is very large run multiple random initialization will make not that big of a difference
- **Choosing the number of clusters**
- Choosing the value of K
  - Elbow method



- Choosing the value of K
  - Sometimes, you're running K-means to get cluster to use for some later/downstream purpose. Evaluate K-means based on a metric for how well it performs for that later purpose.
- **Dimensionality Reduction**
- Motivation 1: Data compression

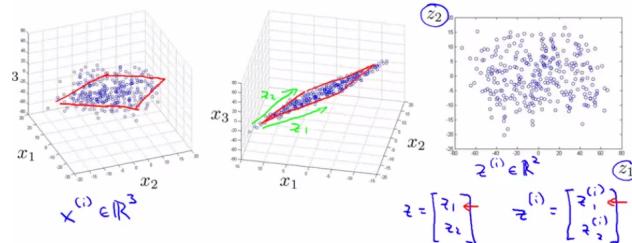
### Data Compression



- Data compression Reduce data from 3d to 2d

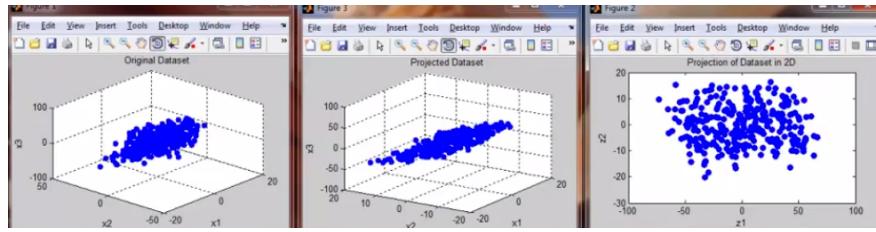
### Data Compression

$1000 \rightarrow 100$   
Reduce data from 3D to 2D



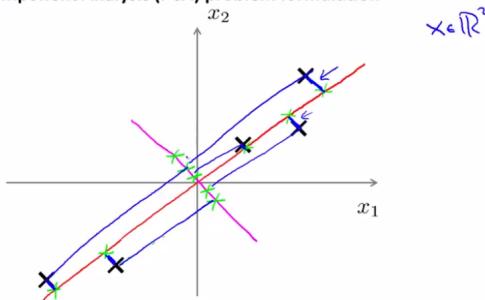
- Left is a 3d point cloud
- Middle is a projection of the data onto a 2d surface or plane

- Right now only need to numbers to represent the data on the plane



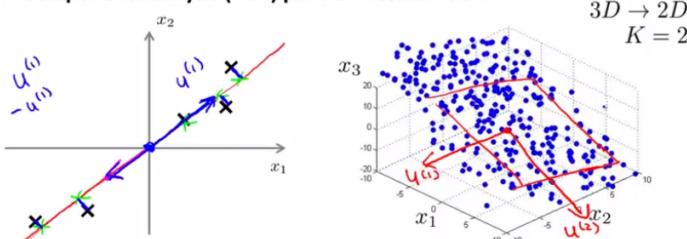
- Process to reduce the data from 3 dimensional to 2 dimensional
- **Motivation 2: Data Visualization**
- **Principal Component Analysis**
- Principal Component
- Analysis (PCA) problem formulation
  - tries to find a lower dimensional surface or a line in this case such that the sum of squares of the blue line segments is minimized. This is called the projection error.
  - should perform mean normalization and feature scaling before perform PCA

Principal Component Analysis (PCA) problem formulation



- Principal Component Analysis (PCA) problem formulation

Principal Component Analysis (PCA) problem formulation



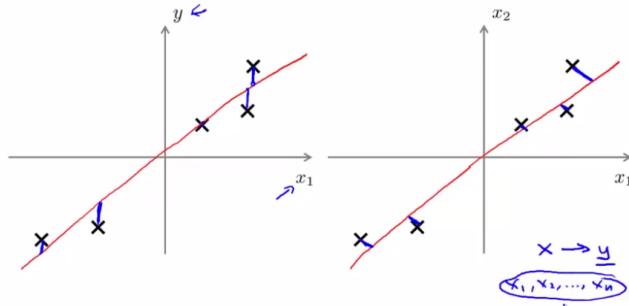
Reduce from 2-dimension to 1-dimension: Find a direction (a vector  $u^{(1)} \in \mathbb{R}^n$ ) onto which to project the data so as to minimize the projection error.

Reduce from  $n$ -dimension to  $k$ -dimension: Find  $k$  vectors  $u^{(1)}, u^{(2)}, \dots, u^{(k)}$  onto which to project the data, so as to minimize the projection error.

- PCA is not Linear regression
  - Linear regression we are minimizing the square magnitude of the blue lines
    - the blue lines are the value of the point and the value predicted by the hypothesis
  - Linear regression is taking the values of  $x$  and trying to predict  $y$

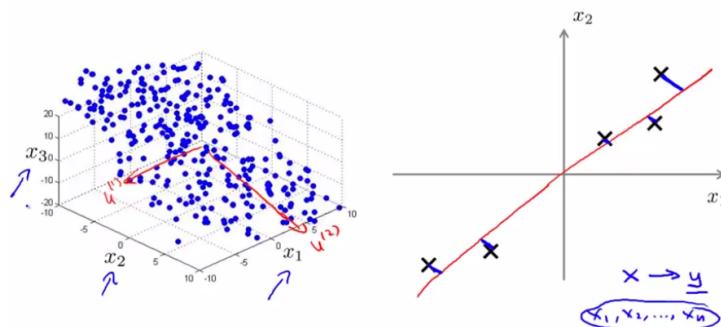
- PCA tries to minimize the magnitude of the blue lines which are drawn at an angle which are the shortest orthogonal or shortest distance between the point  $x$  and the red line
- PCA list of features where all the features are treated equally

### PCA is not linear regression



- 3 dimensional data and want to reduce to 2 dimensional

### PCA is not linear regression



- PCA trying to find a lower dimensional surface to project the data
- **Principal Component Analysis Algorithm**
- Data preprocessing

#### Data preprocessing

Training set:  $x^{(1)}, x^{(2)}, \dots, x^{(m)}$

Preprocessing (feature scaling/mean normalization):

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

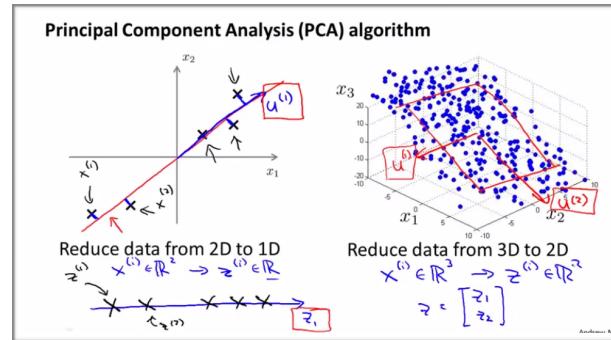
Replace each  $x_j^{(i)}$  with  $x_j^{(i)} - \mu_j$ .

If different features on different scales (e.g.,  $x_1$  = size of house,  $x_2$  = number of bedrooms), scale features to have comparable range of values.

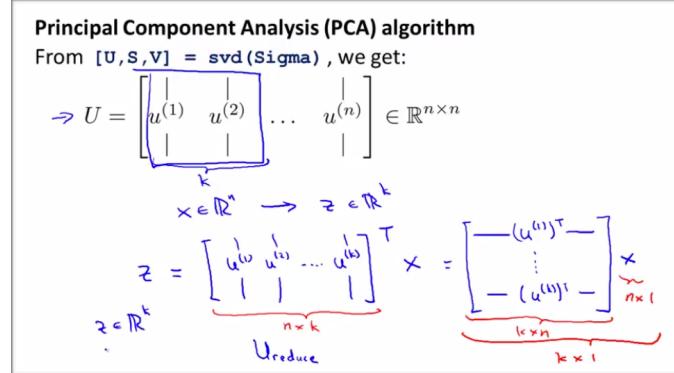
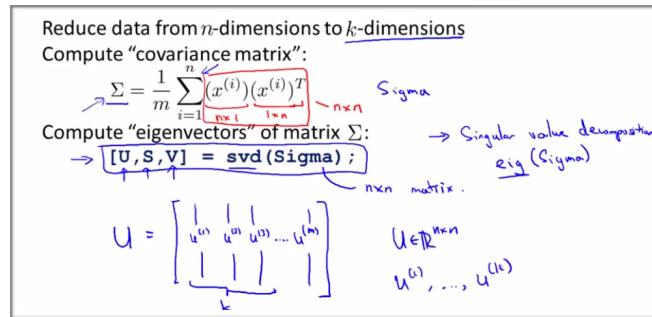
$$x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \mu_j}{s_j}$$

- $s_j$  can be the max - min value or more common the std

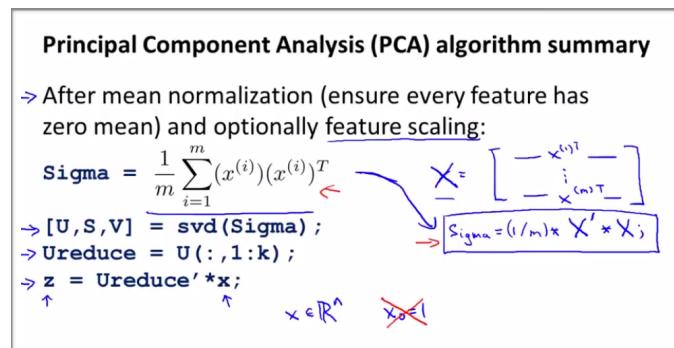
- Principal Component Analysis (PCA) algorithm



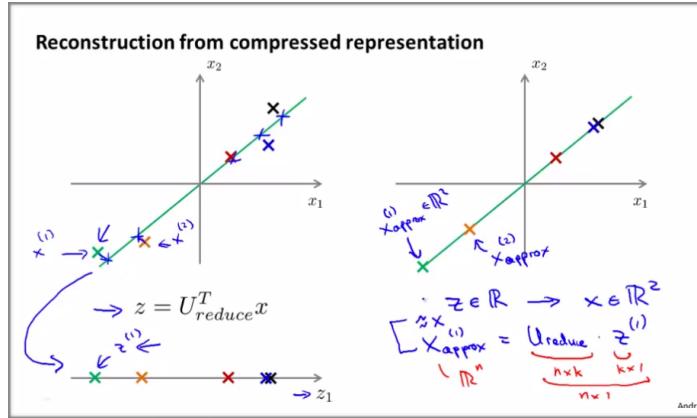
- Principal Component Analysis (PCA) algorithm
  - reduce data from n-dimensions to k-dimensions
  - Compute “covariance matrix”:
  - SVD singular value decomposition



- Principal Component Analysis (PCA) algorithm summary



- Applying PCA
- Reconstruction from Compressed Representation
- Reconstruction from compressed representation



- Choosing the Number of Principal Components
- Choosing k (number of principal components)

**Choosing  $k$  (number of principal components)**

Average squared projection error:  $\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{approx}^{(i)}\|^2$

Total variation in the data:  $\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2$

Typically, choose  $k$  to be smallest value so that

$$\rightarrow \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{approx}^{(i)}\|^2 \leq 0.01 \quad (1\%)$$

$$\rightarrow \frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2$$

⇒ “99% of variance is retained”

- Choosing K (number of principal components)

### Choosing $k$ (number of principal components)

Algorithm:

Try PCA with  $k=1$   ~~$k=2$~~   ~~$k=3$~~   ~~$k=4$~~

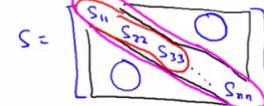
Compute  $U_{reduce}, z^{(1)}, z^{(2)}, \dots, z^{(m)}$ ,  $x_{approx}^{(1)}, \dots, x_{approx}^{(m)}$

Check if

$$\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{approx}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01?$$

$k=1$

$$\rightarrow [U, S, V] = svd(\text{Sigma})$$



For given  $k$

$$1 - \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \leq 0.01$$

$$\frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \geq 0.99$$

- allows to select the value of K without having to run PCA over and over
- Choosing k (number of principal components)

### Choosing $k$ (number of principal components)

$\Rightarrow [U, S, V] = \text{svd}(\Sigma)$

Pick smallest value of  $k$  for which

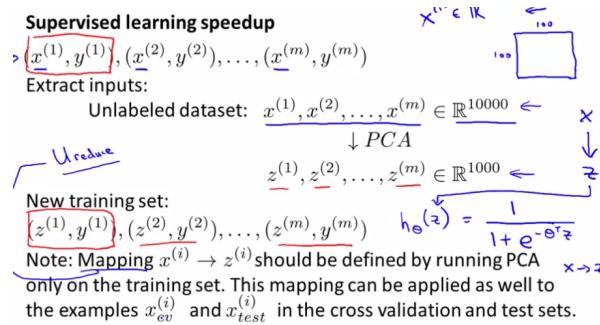
$$\frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^m S_{ii}} \geq 0.99$$

$k \leq 100$

(99% of variance retained)

### - Advice for applying PCA

#### - Supervised Learning speedup



#### - Application of PCA

- Compression
  - Reduce memory/disk needed to store data
  - Speed up learning algorithm
  - Choose  $k$  by percentage of variance retained
- Visualization
  - $k = 2$  or  $k = 3$

#### - Bad use of PCA: To prevent overfitting

##### Bad use of PCA: To prevent overfitting

$\Rightarrow$  Use  $z^{(i)}$  instead of  $x^{(i)}$  to reduce the number of features to  $k < n$ .

Thus, fewer features, less likely to overfit.

Bad!

This might work OK, but isn't a good way to address overfitting. Use regularization instead.

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

- PCA reduces the dimension of the data without knowing what the value of  $y$  is, but this will throw away some of the valuable information

- PCA is sometimes used where it shouldn't be

### PCA is sometimes used where it shouldn't be

Design of ML system:

- - Get training set  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$
- - ~~Run PCA to reduce  $x^{(i)}$  in dimension to get  $z^{(i)}$~~
- - Train logistic regression on  $\{(\cancel{x^{(1)}}, y^{(1)}), \dots, (\cancel{x^{(m)}}, y^{(m)})\}$
- - Test on test set: Map  $x_{test}^{(i)}$  to  $z_{test}^{(i)}$ . Run  $h_\theta(z)$  on  $\{(z_{test}^{(1)}, y_{test}^{(1)}), \dots, (z_{test}^{(m)}, y_{test}^{(m)})\}$

→ How about doing the whole thing without using PCA?

→ Before implementing PCA, first try running whatever you want to do with the original/raw data  $x^{(i)}$ . Only if that doesn't do what you want, then implement PCA and consider using  $z^{(i)}$ .