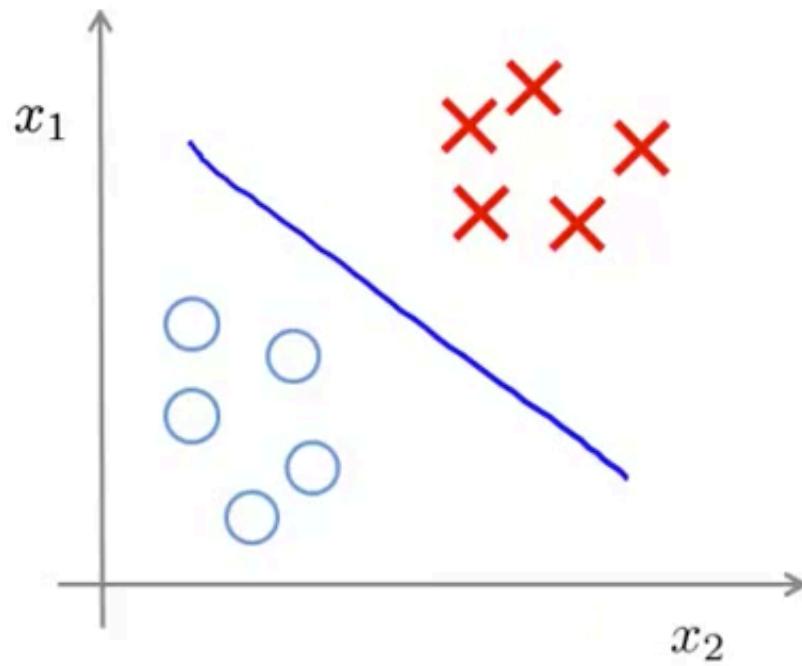


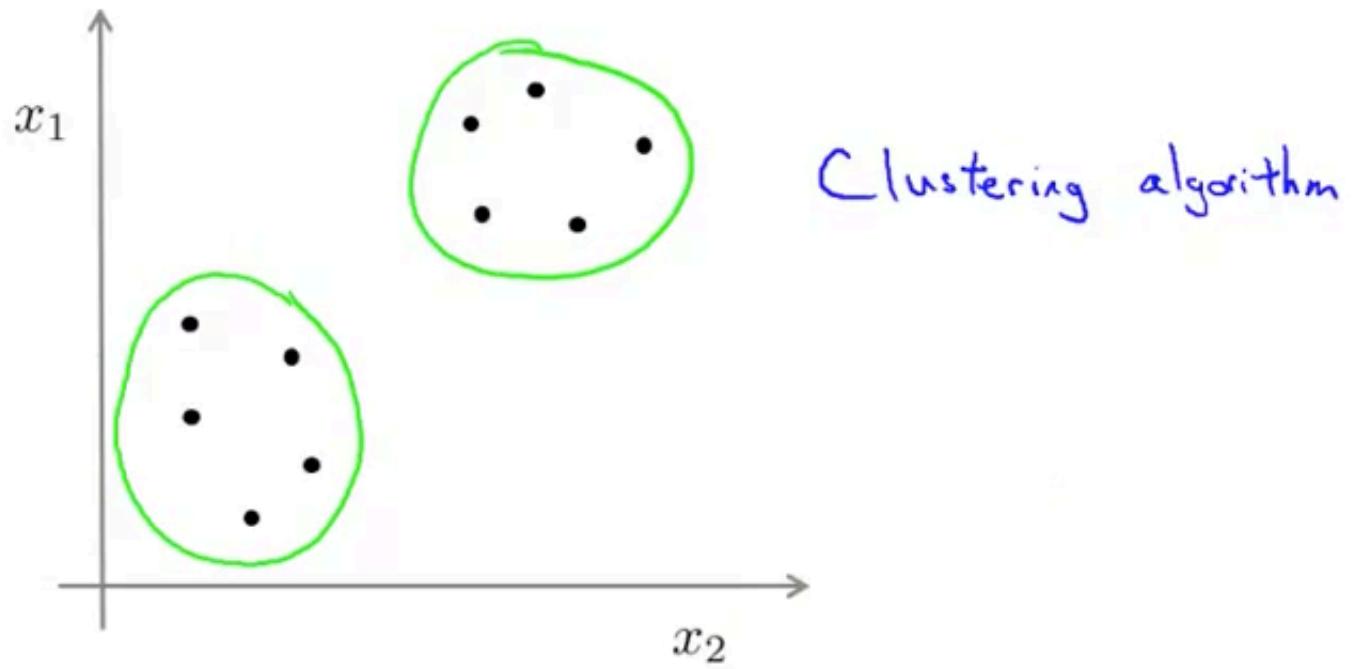
Supervised learning



Training set: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}), \dots, (x^{(m)}, y^{(m)})\}$.



Unsupervised learning

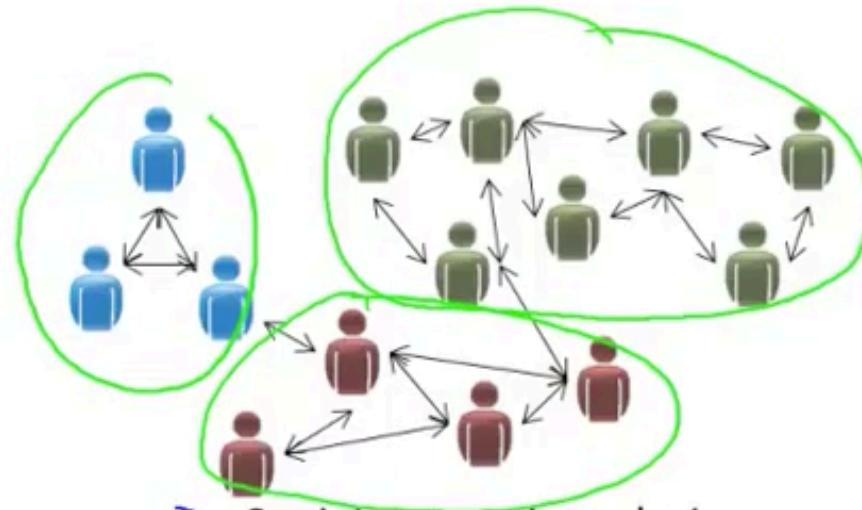


Training set: $\{\underline{x}^{(1)}, \underline{x}^{(2)}, \underline{x}^{(3)}, \dots, \underline{x}^{(m)}\}$ ←

Applications of clustering



→ Market segmentation



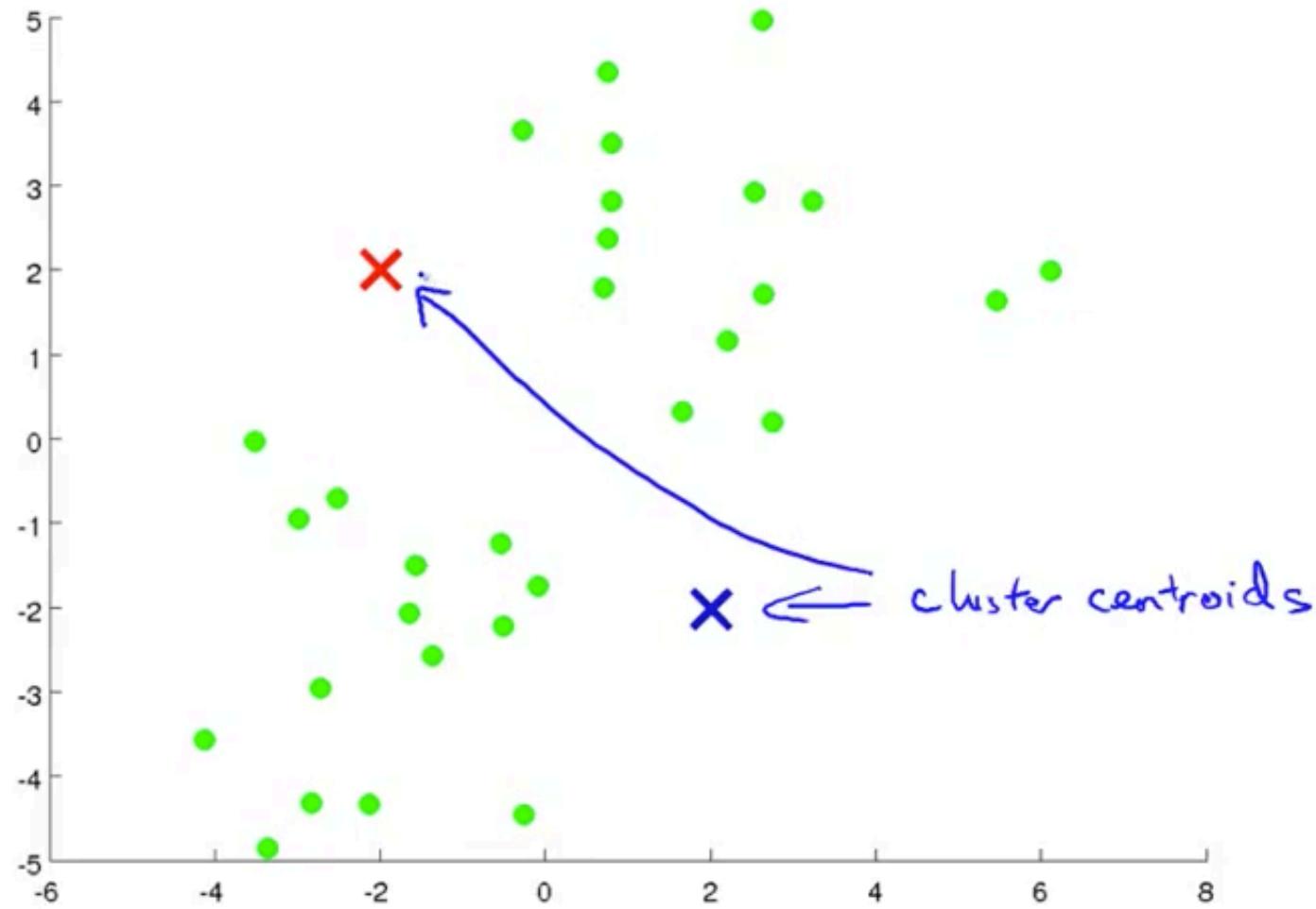
→ Social network analysis

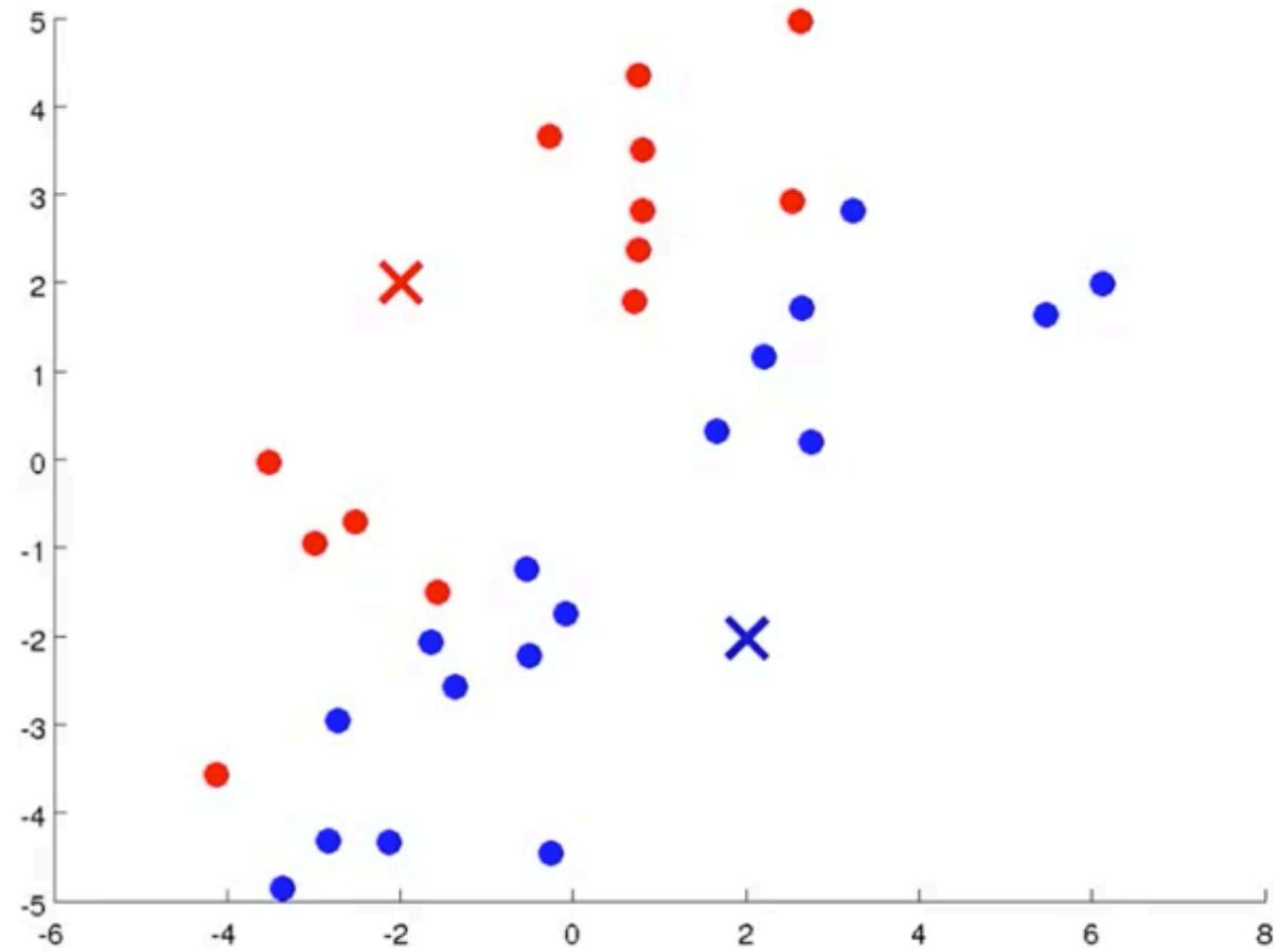


Organize computing clusters

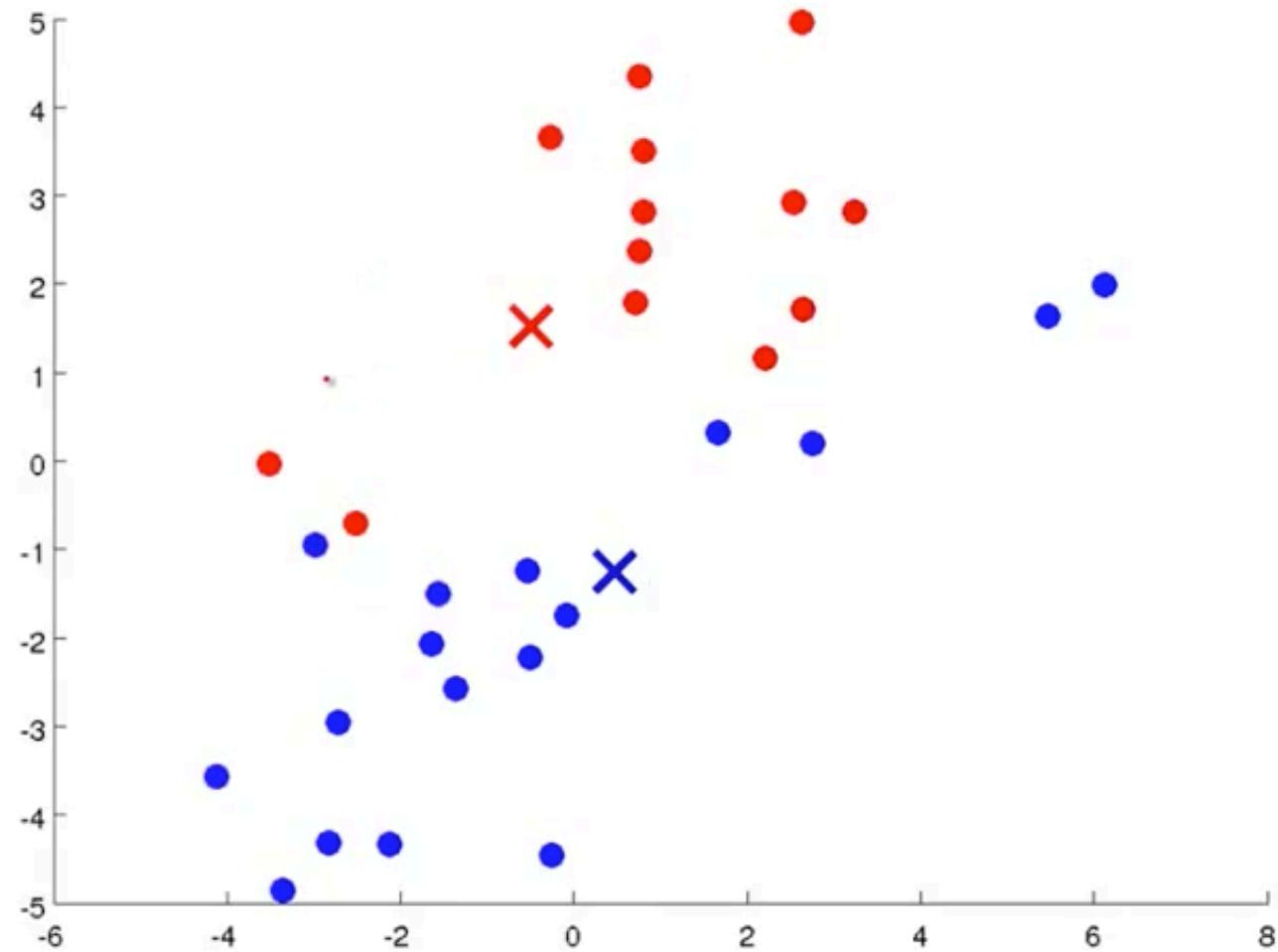


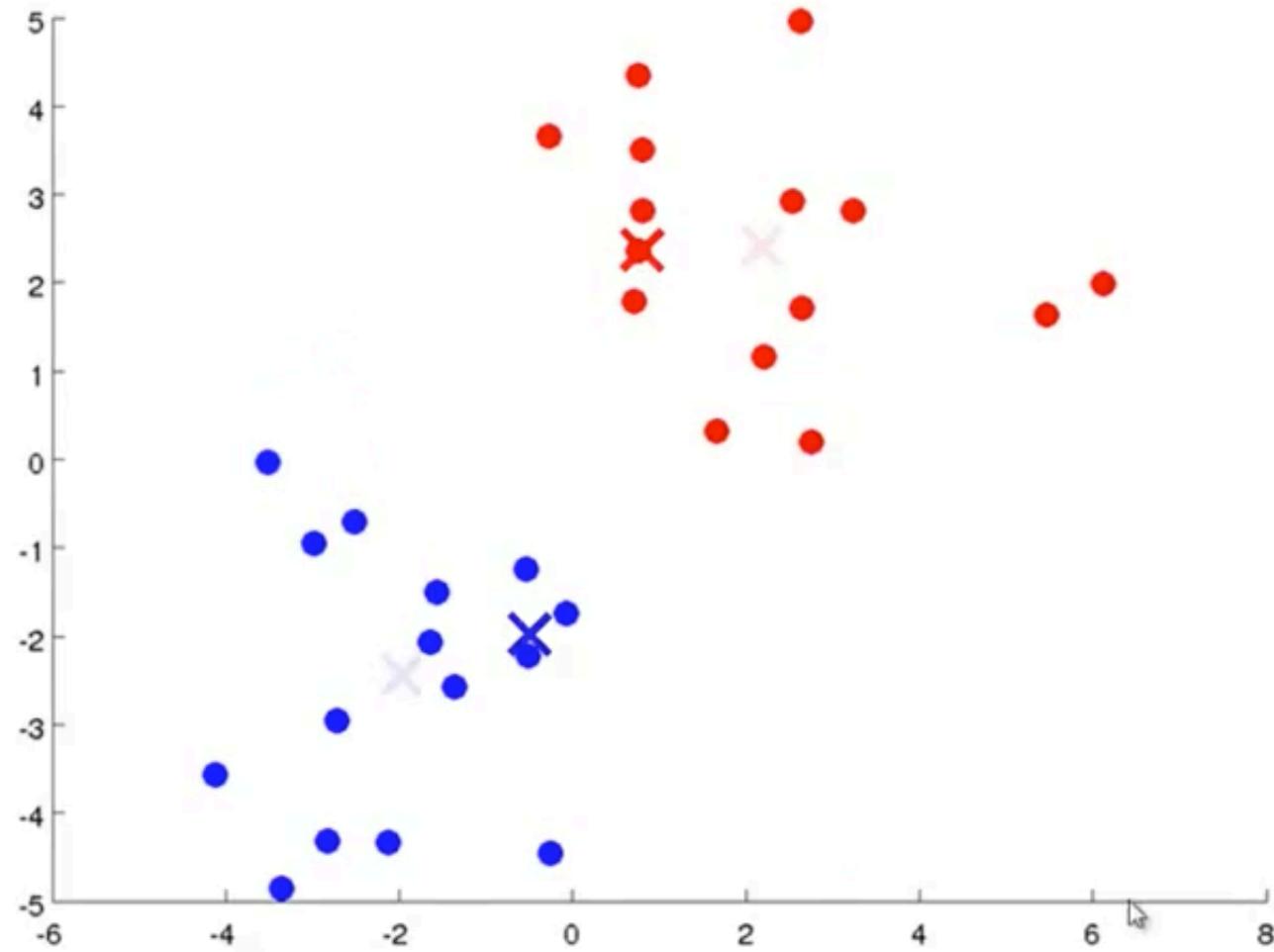
→ Astronomical data analysis

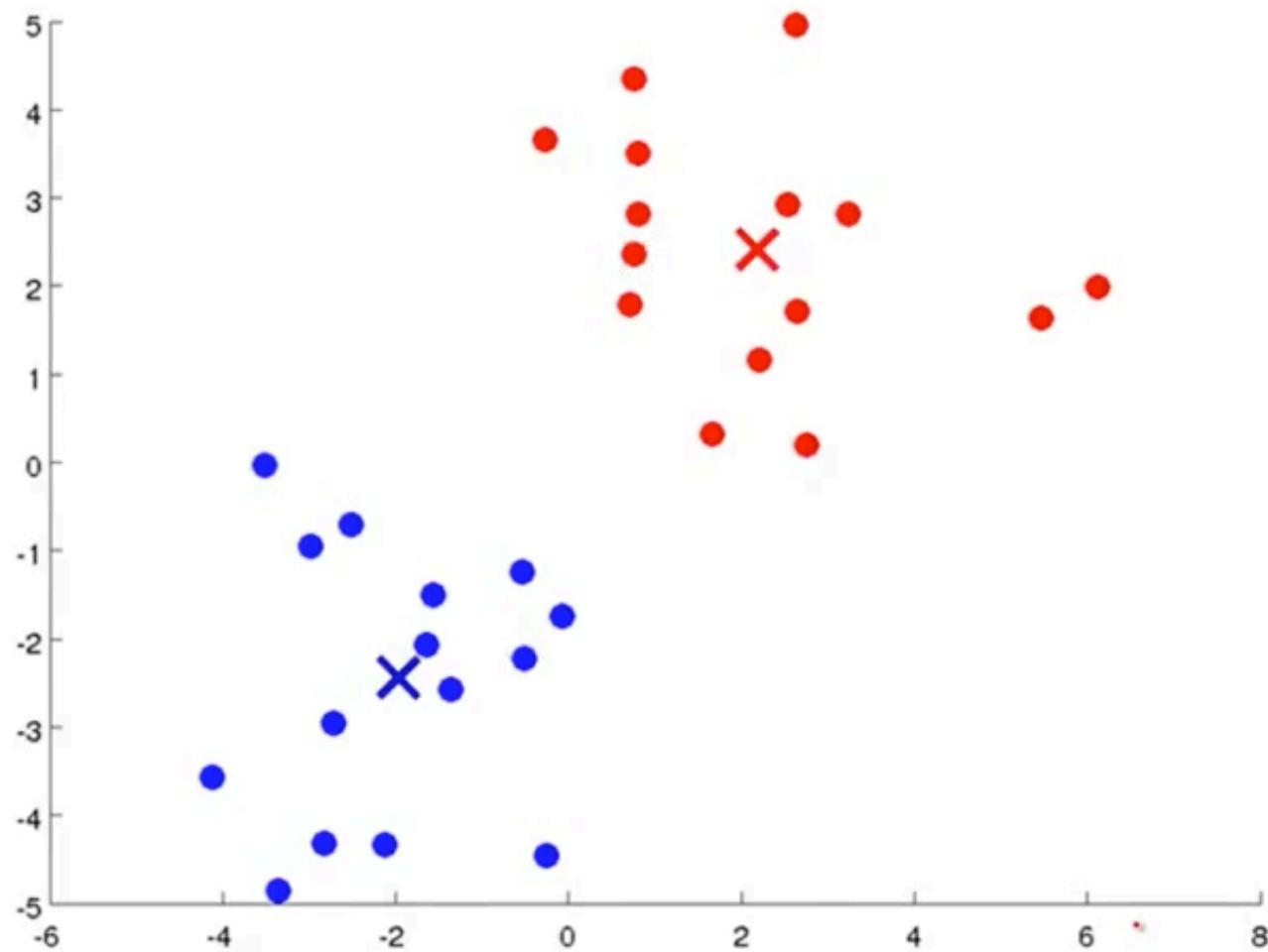




Andrew Ng







K-means algorithm

Input:

- K (number of clusters) 
- Training set $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ 

$x^{(i)} \in \mathbb{R}^n$ (drop $x_0 = 1$ convention)

K-means algorithm

$$\mu_1 \quad \mu_2$$


Randomly initialize K cluster centroids $\underline{\mu}_1, \underline{\mu}_2, \dots, \underline{\mu}_K \in \mathbb{R}^n$

Repeat {

Cluster
Assignment
step

for $i = 1$ to m

$c^{(i)}$:= index (from 1 to K) of cluster centroid
closest to $x^{(i)}$

$$\min_{c^{(i)}} \|x^{(i)} - \underline{\mu}_k\|^2$$

for $k = 1$ to K

$\rightarrow \underline{\mu}_k$:= average (mean) of points assigned to cluster k

$$x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}$$

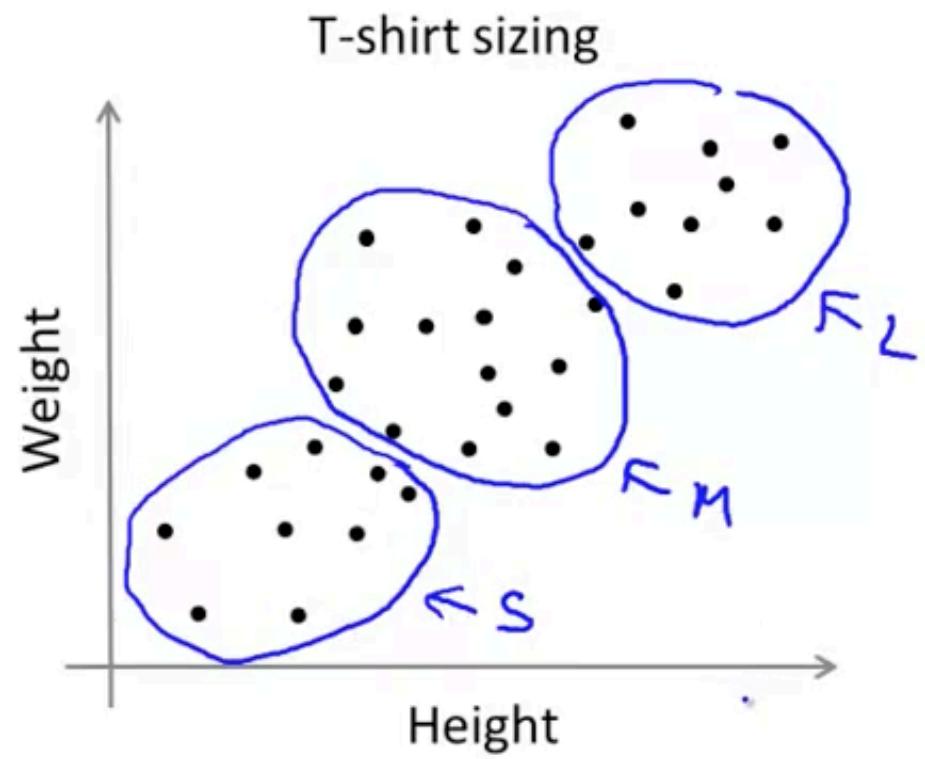
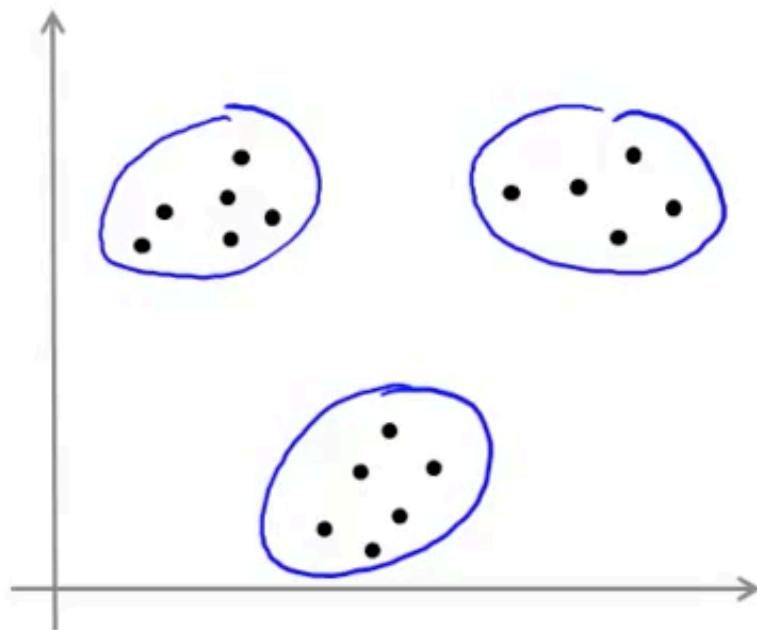
$$\rightarrow c^{(1)}=2, c^{(2)}=2, c^{(3)}=2, c^{(4)}=2$$

}

$$\underline{\mu}_2 = \frac{1}{4} \left[\underline{x}^{(1)} + \underline{x}^{(2)} + \underline{x}^{(3)} + \underline{x}^{(4)} \right] \in \mathbb{R}^n$$

K-means for non-separated clusters

S, M, L



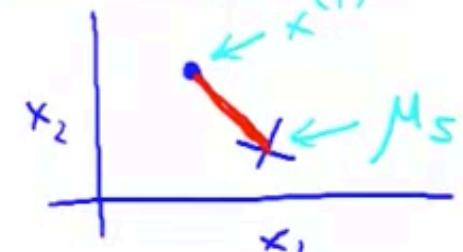
K-means optimization objective

- $c^{(i)}$ = index of cluster ($1, 2, \dots, K$) to which example $x^{(i)}$ is currently assigned
- μ_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 $x^{(i)} \rightarrow S$ $c^{(i)} = 5$ $\underline{\mu_{c^{(i)}}} = \underline{\mu_5}$

Optimization objective:

$$\rightarrow J(\underbrace{c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K}_{\text{Distortion}}) = \frac{1}{m} \sum_{i=1}^m \boxed{||x^{(i)} - \mu_{c^{(i)}}||^2} \leftarrow$$
$$\min_{\substack{c^{(1)}, \dots, c^{(m)}, \\ \mu_1, \dots, \mu_K}} J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$$

Distortion



K-means algorithm

Randomly initialize K cluster centroids $\mu_1, \mu_2, \dots, \mu_K \in \mathbb{R}^n$

Repeat {

Cluster assignment step

Minimize $J(\dots)$ w.r.t
(holding μ_1, \dots, μ_K fixed)

$c^{(1)}, c^{(2)}, \dots, c^{(n)} \leftarrow$

for $i = 1$ to m

$c^{(i)} :=$ index (from 1 to K) of cluster centroid
closest to $x^{(i)}$

move
centroid

for $k = 1$ to K

$\mu_k :=$ average (mean) of points assigned to cluster k

}

minimize $J(\dots)$ w.r.t

μ_1, \dots, μ_K

Random initialization

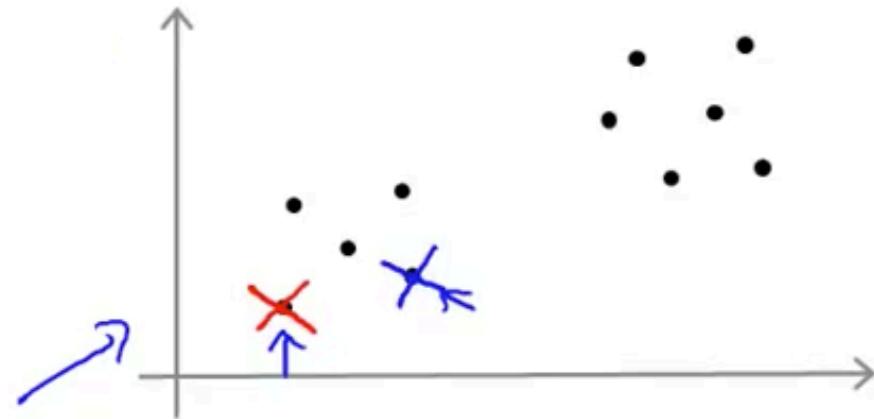
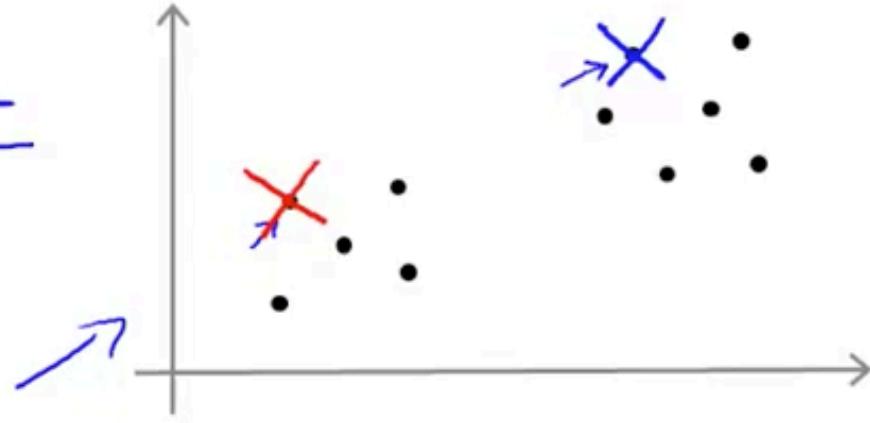
Should have $K < m$

Randomly pick K training examples.

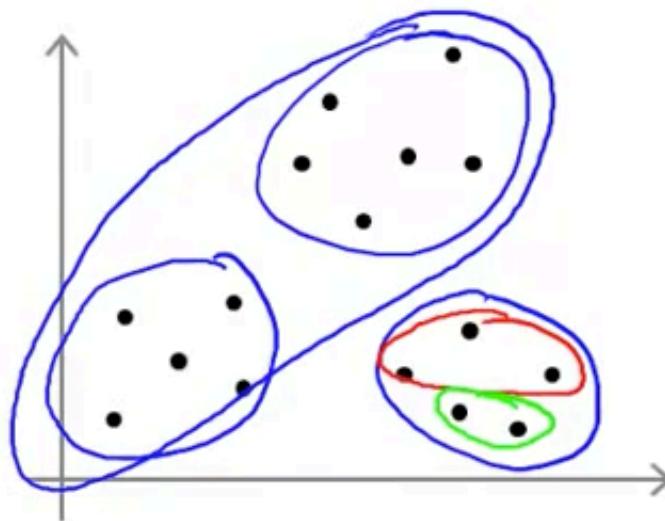
Set μ_1, \dots, μ_K equal to these K examples.

$$\begin{aligned}\mu_1 &= x^{(i)} \\ \mu_2 &= x^{(j)} \\ &\vdots\end{aligned}$$

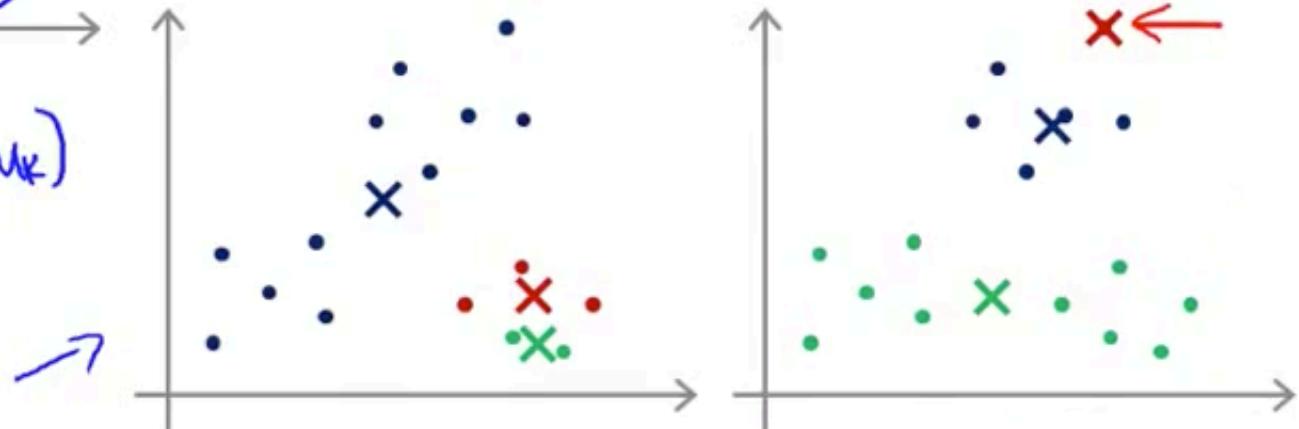
$$K=2$$



Local optima



$$\mathcal{J}(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k)$$



Andrew Ng

Random initialization

For $i = 1$ to 100 { $50 - 1000$

 → Randomly initialize K-means.

 Run K-means. Get $c^{(1)}, \dots, c^{(m)}$, μ_1, \dots, μ_K .

 Compute cost function (distortion)

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

}

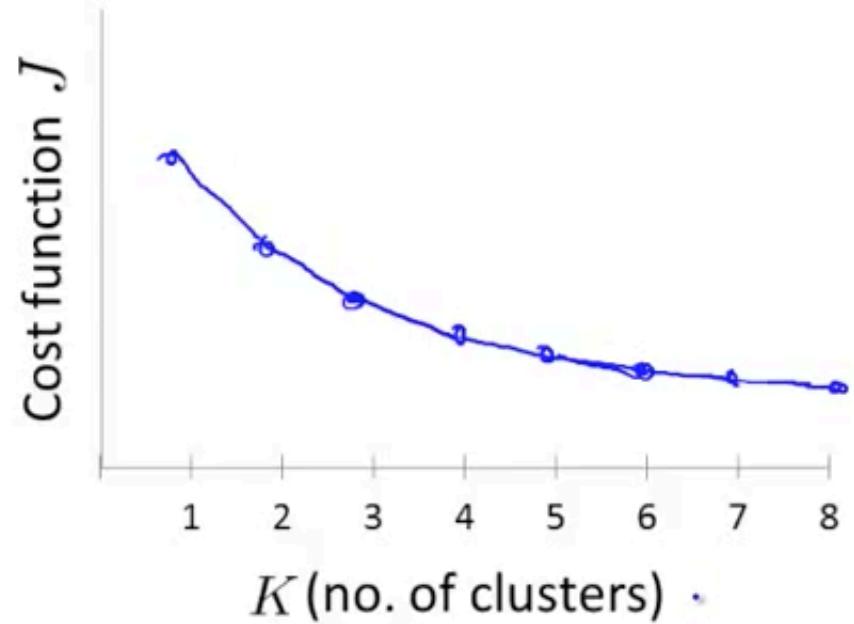
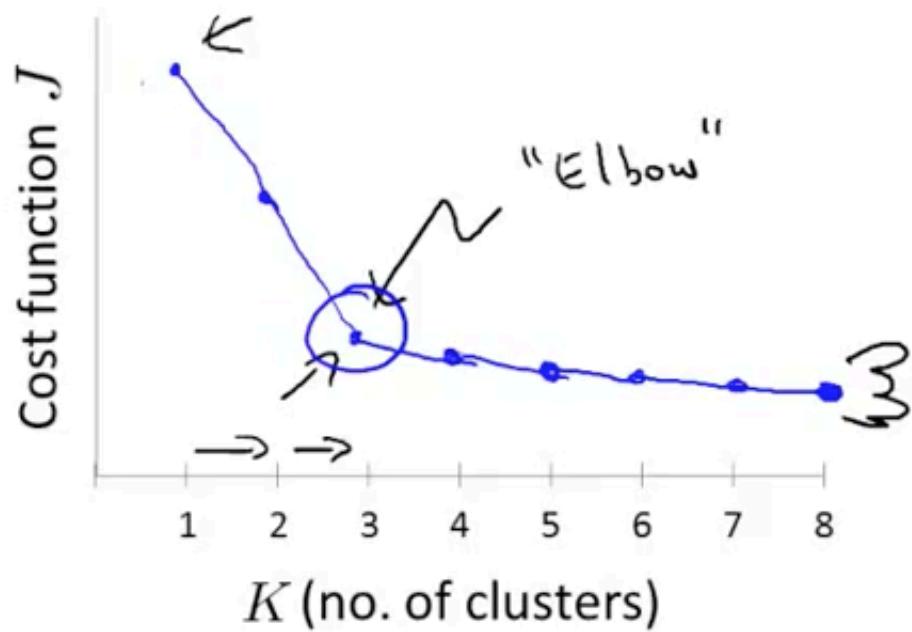
Pick clustering that gave lowest cost $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$

$k = 2 - 10$

\hat{i}

Choosing the value of K

Elbow method:

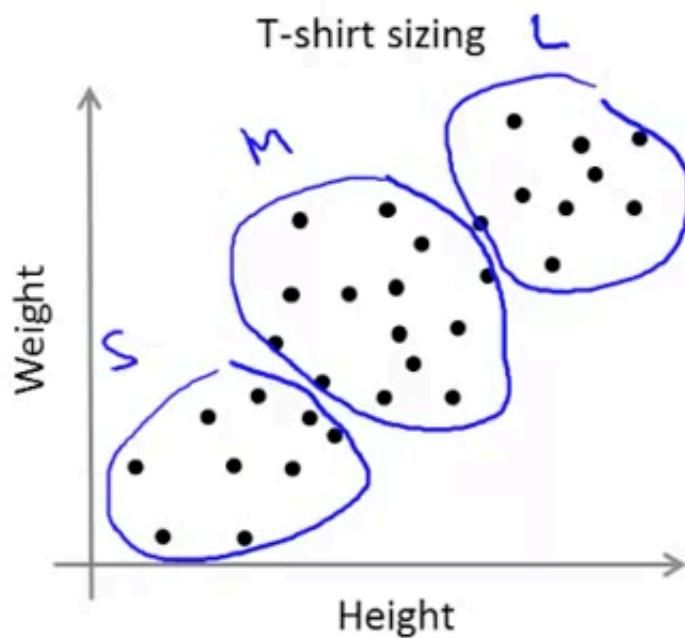


Choosing the value of K

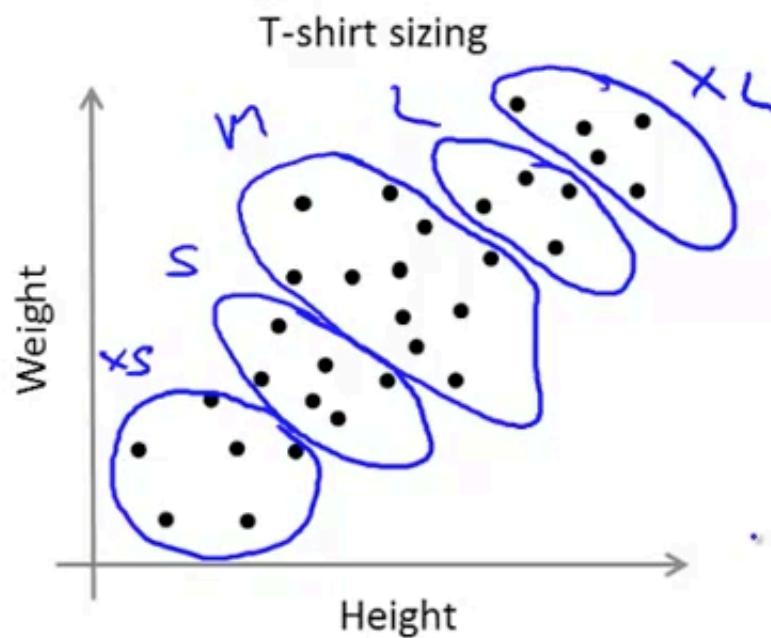
Sometimes, you're running K-means to get clusters to use for some later/downstream purpose. Evaluate K-means based on a metric for how well it performs for that later purpose.

$$K=3 \quad S, M, L$$

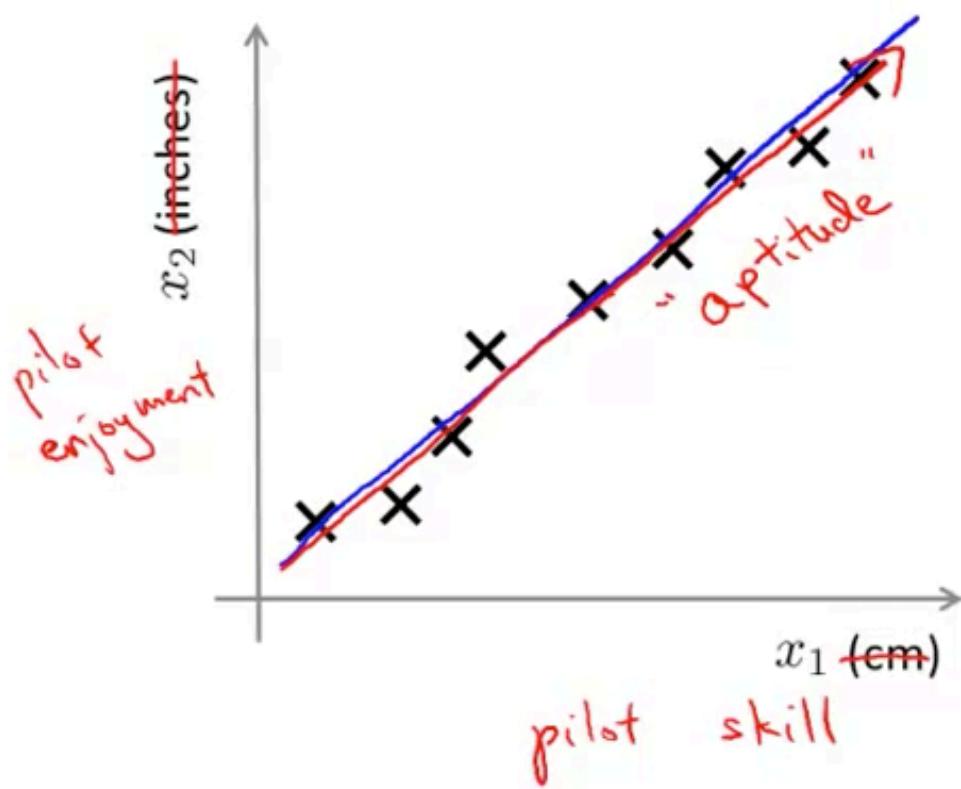
E.g.



$$K=5 \quad XS, S, M, L, XL$$

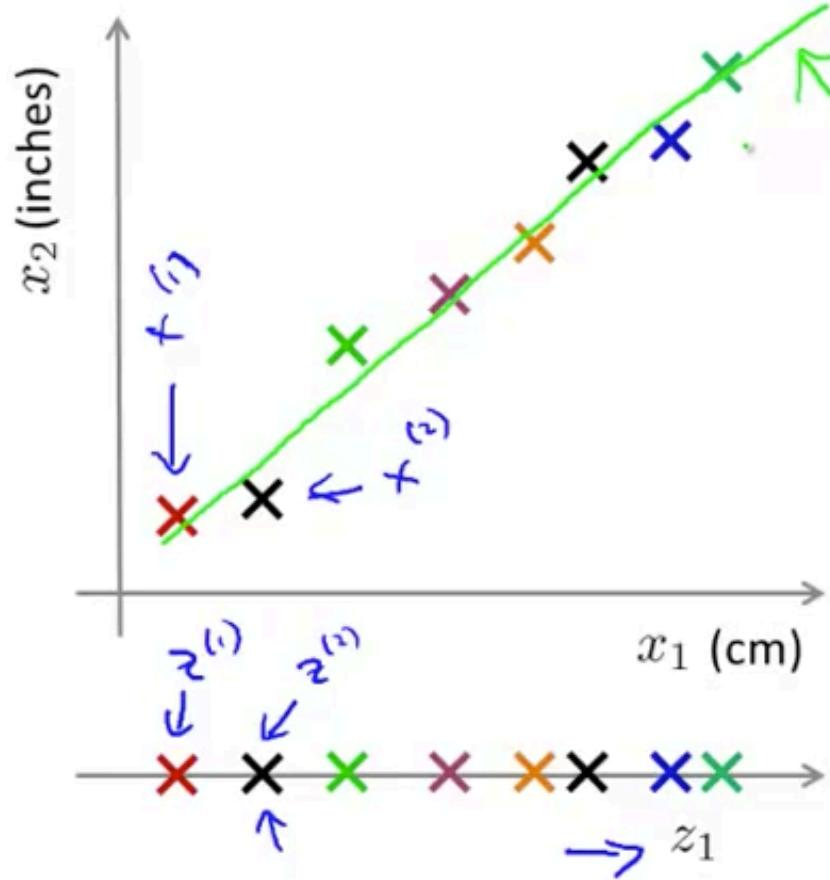


Data Compression



Reduce data from
2D to 1D

Data Compression



Reduce data from
2D to 1D

$$x^{(1)} \in \mathbb{R}^2 \rightarrow z^{(1)} \in \mathbb{R}$$

$$x^{(2)} \in \mathbb{R}^2 \rightarrow z^{(2)} \in \mathbb{R}$$

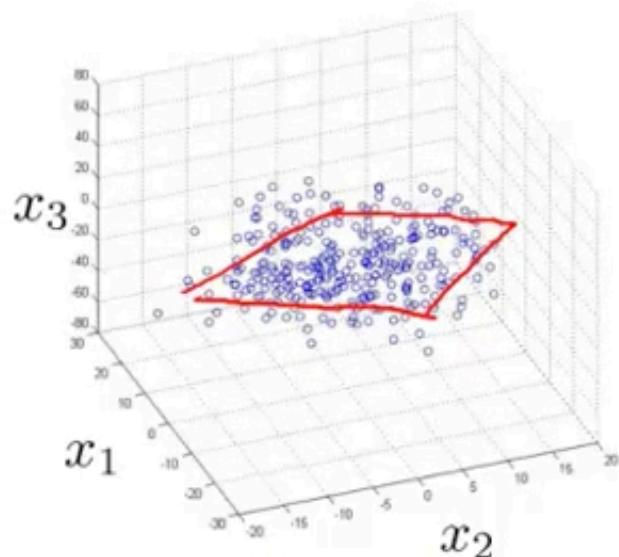
⋮

$$x^{(m)} \rightarrow z^{(m)}$$

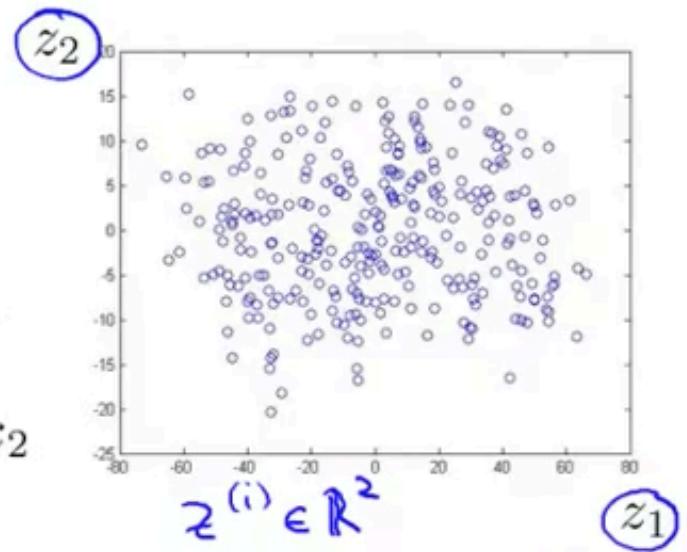
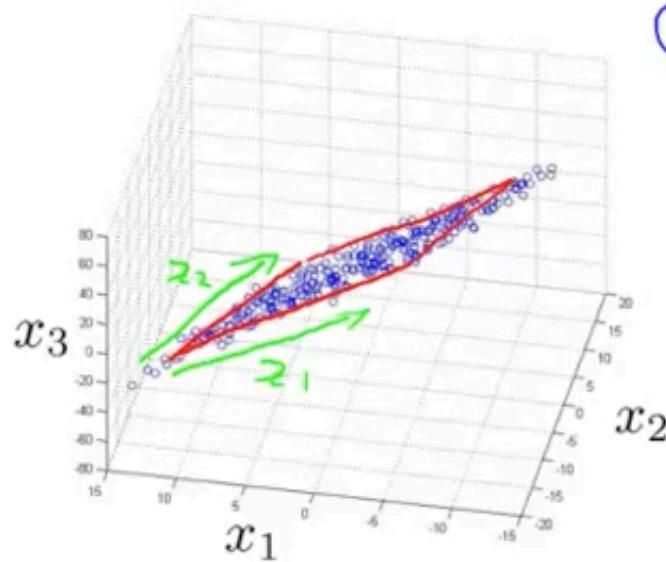
Data Compression

10000 \rightarrow 1000

Reduce data from 3D to 2D



$$x^{(i)} \in \mathbb{R}^3$$



$$z = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \quad z^{(i)} = \begin{bmatrix} z_1^{(i)} \\ z_2^{(i)} \end{bmatrix}$$

Data Visualization

$$x \in \mathbb{R}^{50}$$

$$x^{(i)} \in \mathbb{R}^{50}$$

x_6

Country	x_1 GDP (trillions of US\$)	x_2 Per capita GDP (thousands of intl. \$)	x_3 Human Development Index	x_4 Life expectancy	x_5 Poverty Index (Gini as percentage)	x_6 Mean household income (thousands of US\$)	...
Canada	1.577	39.17	0.908	80.7	32.6	67.293	...
China	5.878	7.54	0.687	73	46.9	10.22	...
India	1.632	3.41	0.547	64.7	36.8	0.735	...
Russia	1.48	19.84	0.755	65.5	39.9	0.72	...
Singapore	0.223	56.69	0.866	80	42.5	67.1	...
USA	14.527	46.86	0.91	78.3	40.8	84.3	...
...

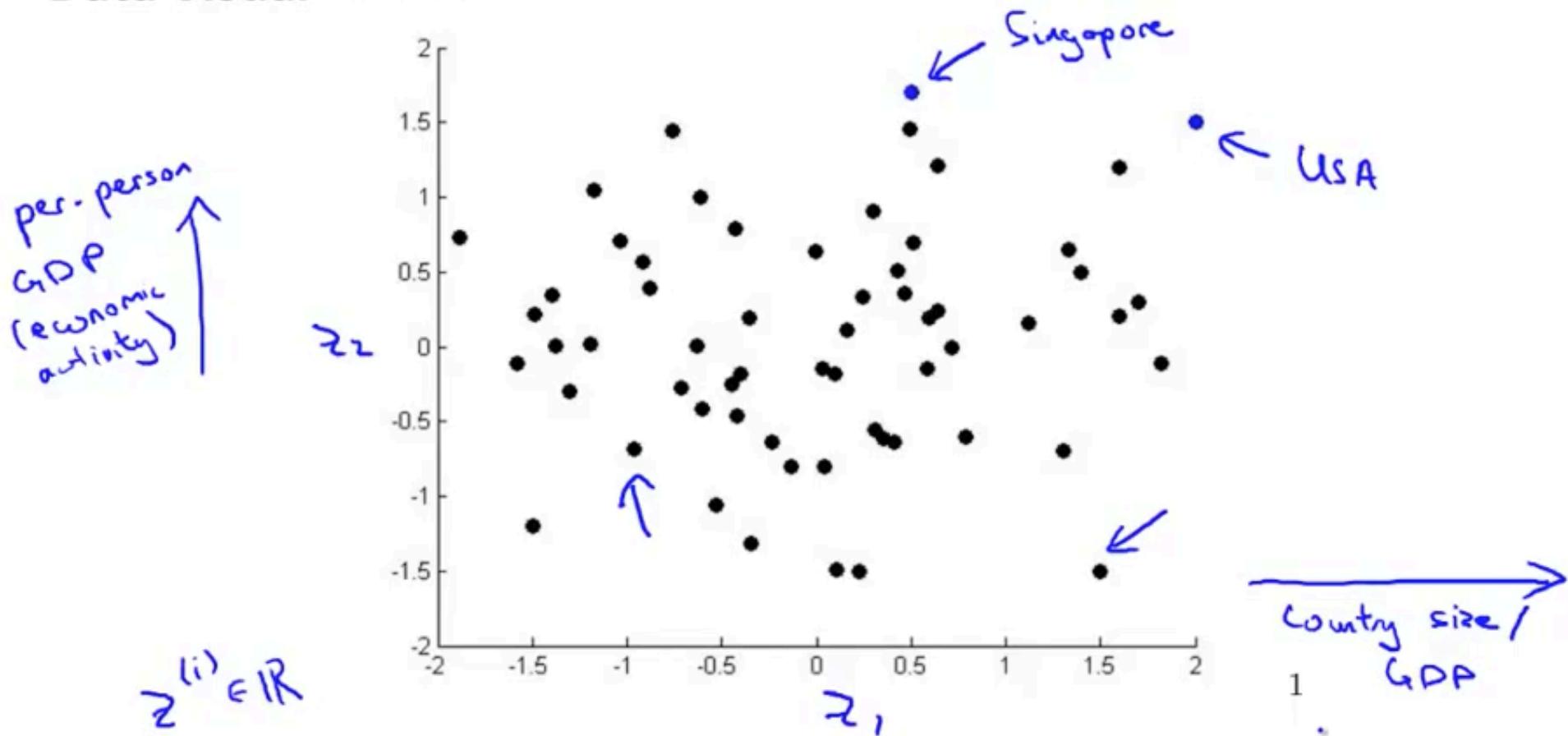
[resources from en.wikipedia.org]

Andrew Ng

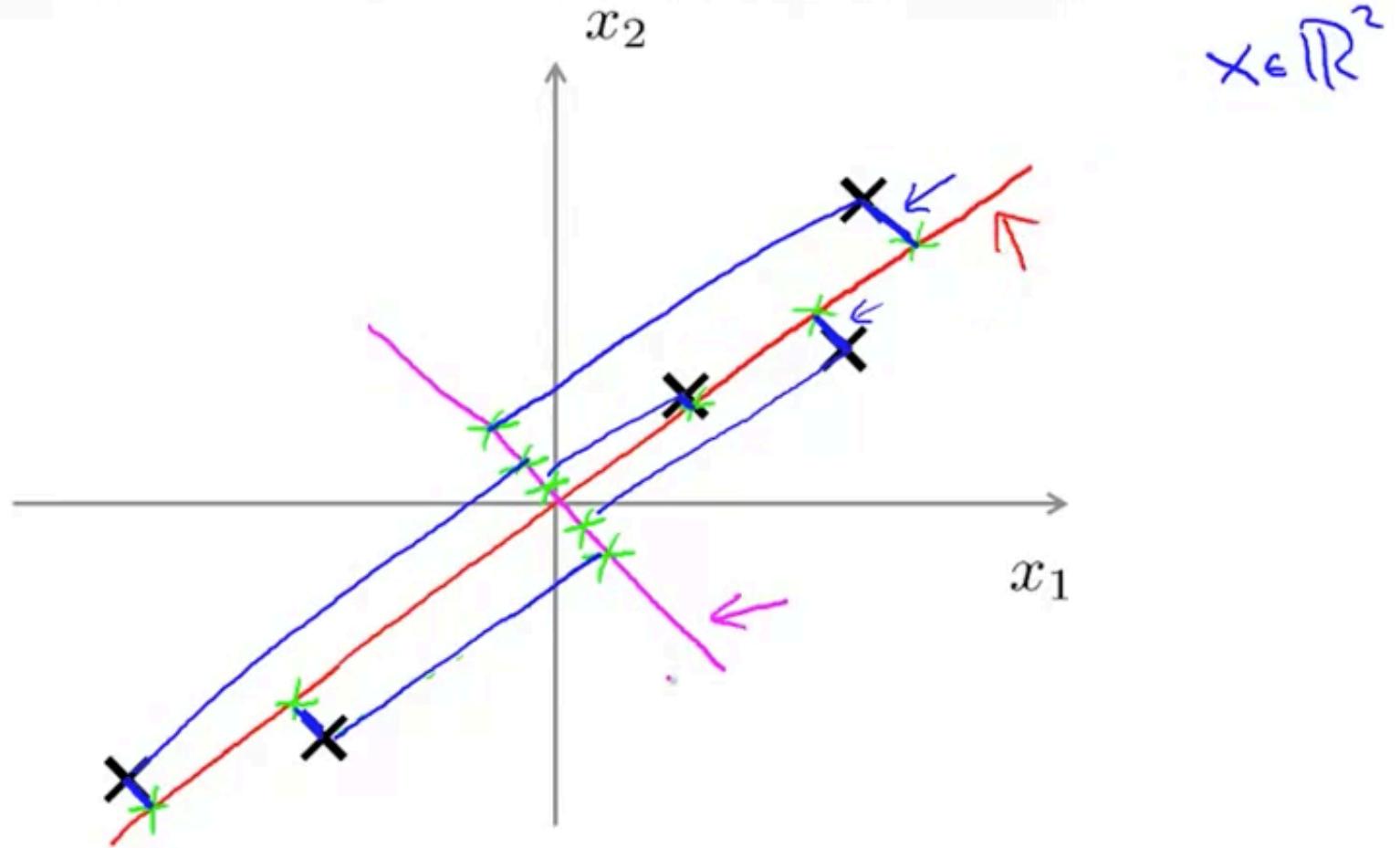
Data Visualization

Country	z_1	z_2	$z^{(i)} \in \mathbb{R}^2$
Canada	1.6	1.2	
China	1.7	0.3	Reduce data
India	1.6	0.2	from 50D
Russia	1.4	0.5	to 2D
Singapore	0.5	1.7	
USA	2	1.5	
...	

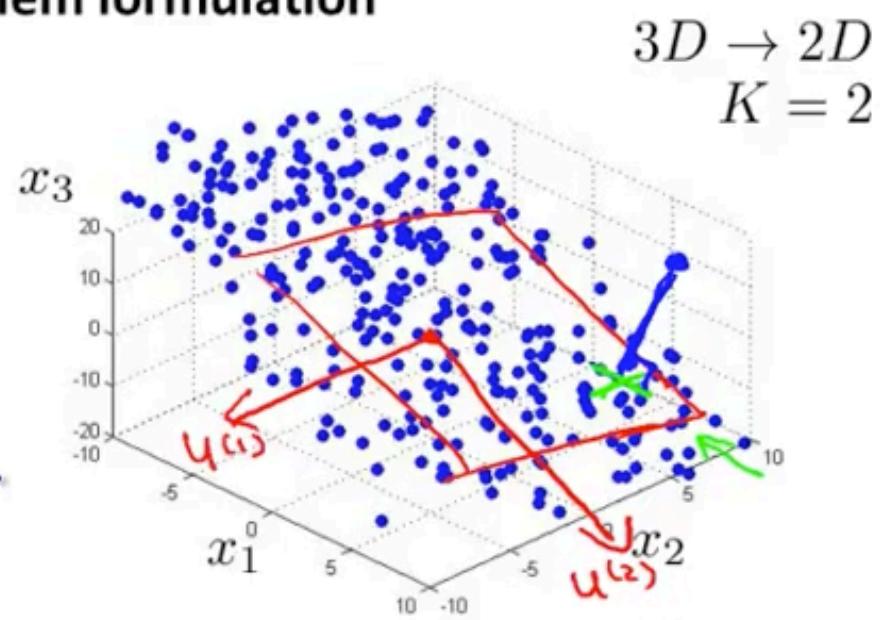
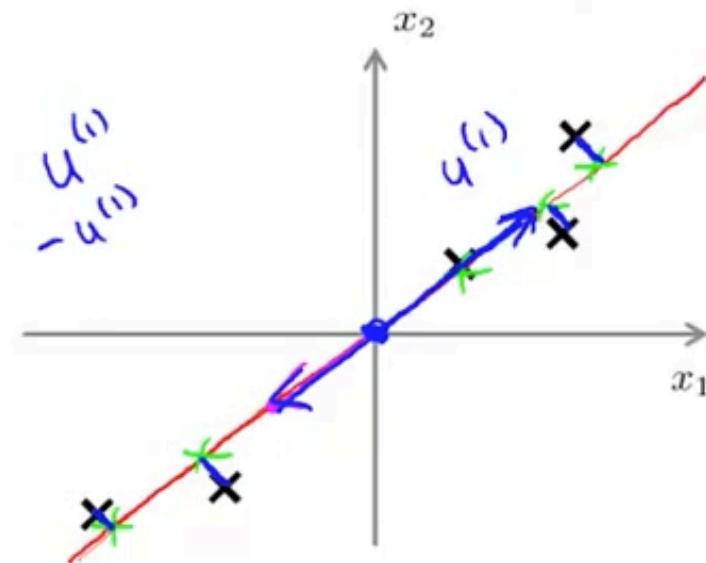
Data Visualization



Principal Component Analysis (PCA) problem formulation



Principal Component Analysis (PCA) problem formulation

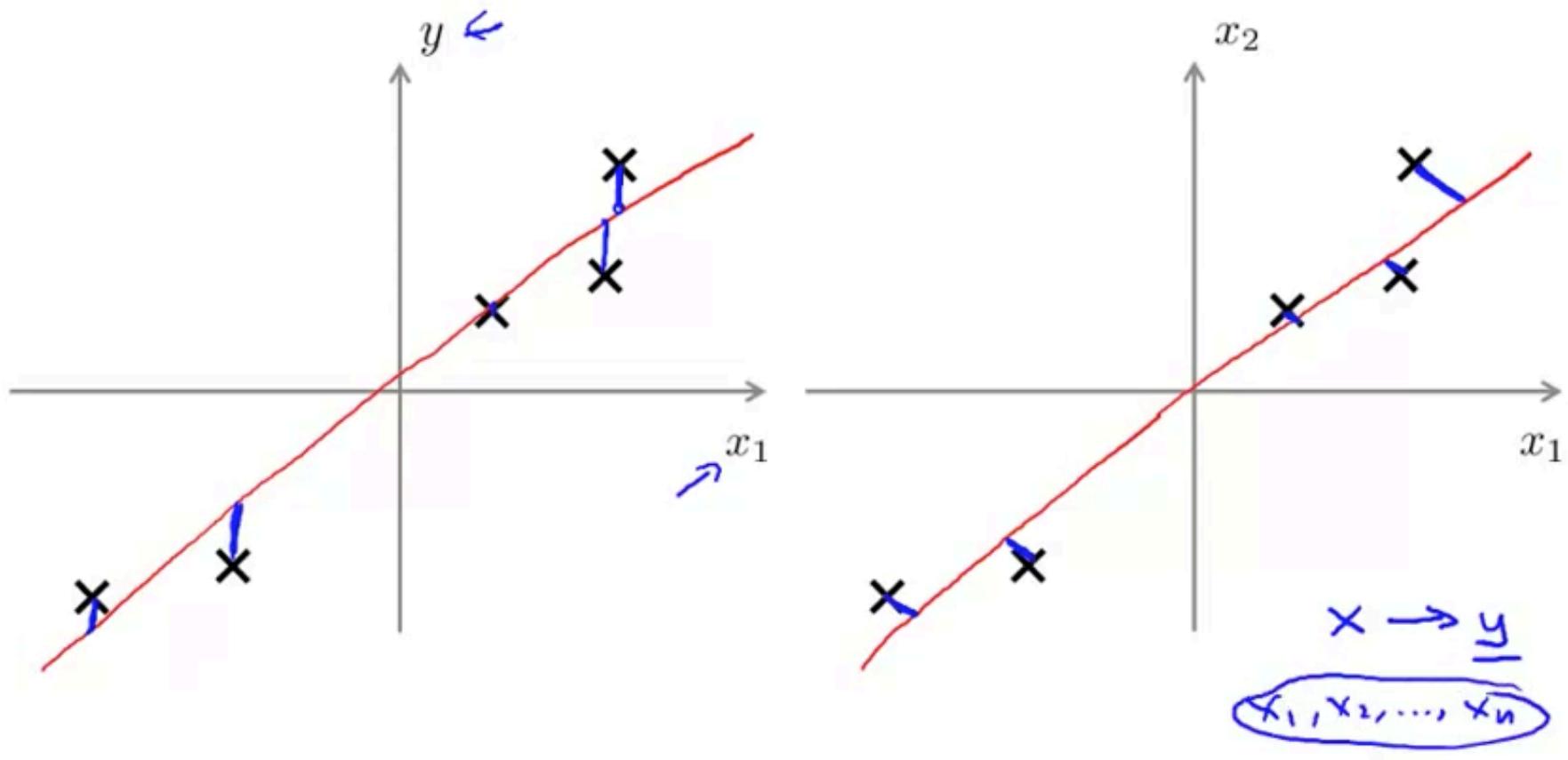


$3D \rightarrow 2D$
 $K = 2$

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

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

PCA is not linear regression



Data preprocessing

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

Preprocessing (feature scaling/mean normalization):

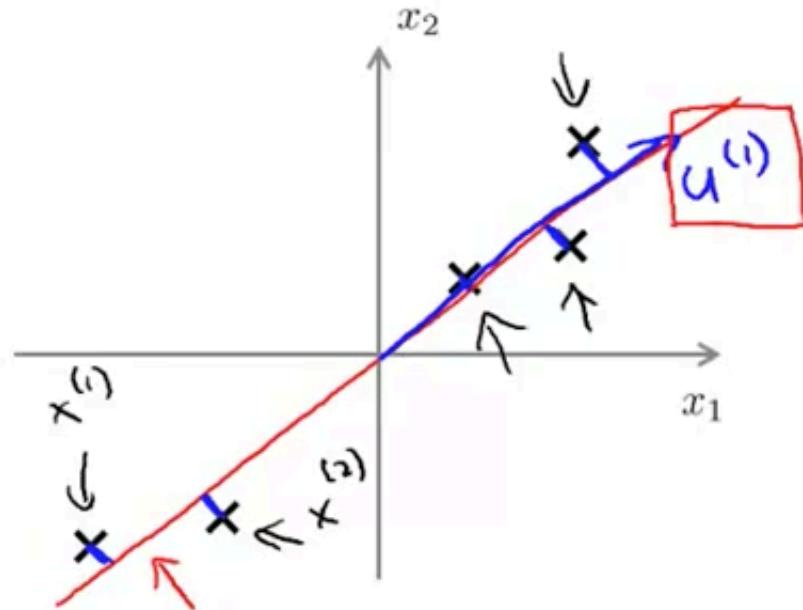
$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

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

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.

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

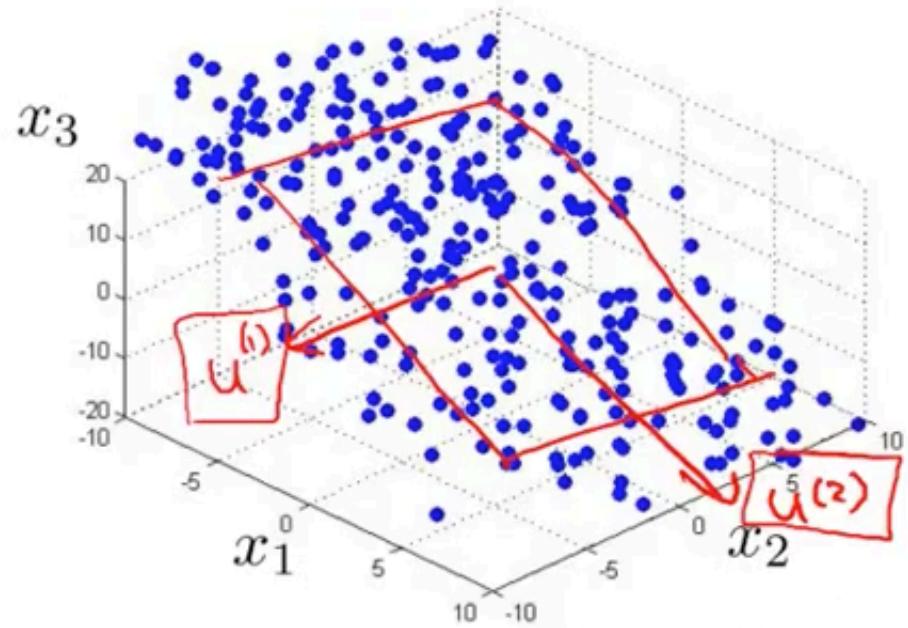
Principal Component Analysis (PCA) algorithm



Reduce data from 2D to 1D

$$z^{(i)} \quad x^{(i)} \in \mathbb{R}^2 \rightarrow z^{(i)} \in \mathbb{R}$$

A horizontal line representing the projection of 2D data points onto a 1D axis, labeled $z^{(i)}$.



Reduce data from 3D to 2D

$$x^{(i)} \in \mathbb{R}^3 \rightarrow z^{(i)} \in \mathbb{R}^2$$

$$z = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} .$$

Principal Component Analysis (PCA) algorithm

Reduce data from n -dimensions to $\underline{k\text{-dimensions}}$

Compute "covariance matrix":

$$\Sigma = \frac{1}{m} \sum_{i=1}^n (x^{(i)}) (x^{(i)})^T$$

$n \times 1$ $1 \times n$

Sigma

Compute "eigenvectors" of matrix Σ :

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

n × n matrix.

→ Singular value decomposition
eig(Sigma)

$$U = \begin{bmatrix} | & | & | & | \\ u^{(1)} & u^{(2)} & u^{(3)} & \dots u^{(m)} \\ | & | & | & | \end{bmatrix}$$

k

$U \in \mathbb{R}^{n \times n}$
 $u^{(1)}, \dots, u^{(k)}$

Principal Component Analysis (PCA) algorithm

From $[U, S, V] = \text{svd}(\Sigma)$, we get:

$$\Rightarrow U = \begin{bmatrix} u^{(1)} & u^{(2)} & \dots & u^{(n)} \end{bmatrix} \in \mathbb{R}^{n \times n}$$

$$x \in \mathbb{R}^n \rightarrow z \in \mathbb{R}^k$$

$$\cdot z^{(i)} = \underbrace{\begin{bmatrix} u^{(1)} & u^{(2)} & \dots & u^{(k)} \end{bmatrix}^T}_{n \times k} x^{(i)} = \underbrace{\begin{bmatrix} (u^{(1)})^T \\ \vdots \\ (u^{(k)})^T \end{bmatrix}}_{k \times n} \underbrace{\begin{bmatrix} z^{(i)} \\ \vdots \\ z^{(i)} \end{bmatrix}}_{k \times 1}$$

U_{reduce}

Principal Component Analysis (PCA) algorithm summary

- After mean normalization (ensure every feature has zero mean) and optionally feature scaling:

$$\text{Sigma} = \frac{1}{m} \sum_{i=1}^m (x^{(i)})(x^{(i)})^T$$

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

→ $U_{\text{reduce}} = U(:, 1:k)$;

→ $z = U_{\text{reduce}}' * x$;

↑

↑

$$x \in \mathbb{R}^n$$

$$x_0 \neq 1$$

$$X = \begin{bmatrix} x^{(1)\top} \\ \vdots \\ x^{(m)\top} \end{bmatrix}$$

$$\text{Sigma} = (1/m) * X' * X$$

In PCA, we obtain $z \in \mathbb{R}^k$ from $x \in \mathbb{R}^n$ as follows:

$$z = \begin{bmatrix} | & | & & | \\ u^{(1)} & u^{(2)} & \dots & u^{(k)} \\ | & | & & | \end{bmatrix}^T x = \begin{bmatrix} \dots & (u^{(1)})^T & \dots \\ \dots & (u^{(2)})^T & \dots \\ & \vdots & \\ \dots & (u^{(k)})^T & \dots \end{bmatrix} x$$

Which of the following is a correct expression for z_j ?

- $z_j = (u^{(k)})^T x$
- $z_j = (u^{(j)})^T x_j$
- $z_j = (u^{(j)})^T x_k$
- $z_j = (u^{(j)})^T x$

Continue

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

$$\rightarrow \frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{approx}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq \frac{0.01}{0.05} \quad \frac{(1\%)}{5\%}$$

\rightarrow “~~99%~~ of variance is retained”
~~95%~~ 90%.

Choosing k (number of principal components)

Algorithm:

Try PCA with $\underline{k=1}$ ~~$k=2$~~ ~~$k=3$~~ $\underline{k=4}$

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

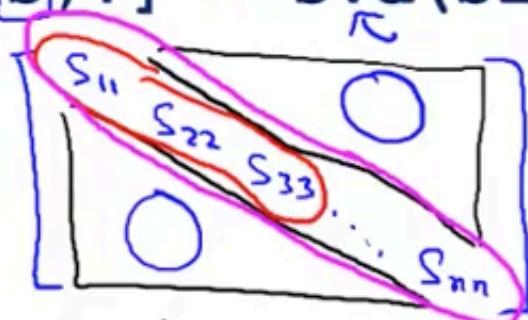
Check if

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

$$\underline{k=17}$$

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

$$\rightarrow \Sigma =$$



For given k $\underline{k=3}$

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

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

Choosing k (number of principal components)

→ $[U, S, V] = \text{svd}(\Sigma)$

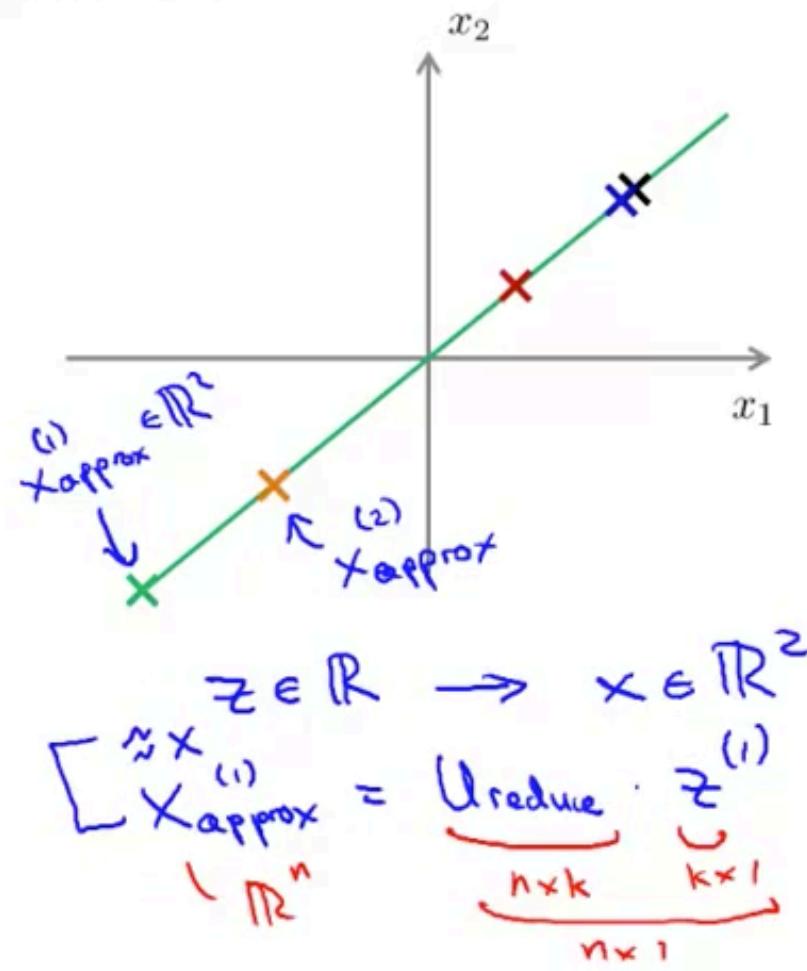
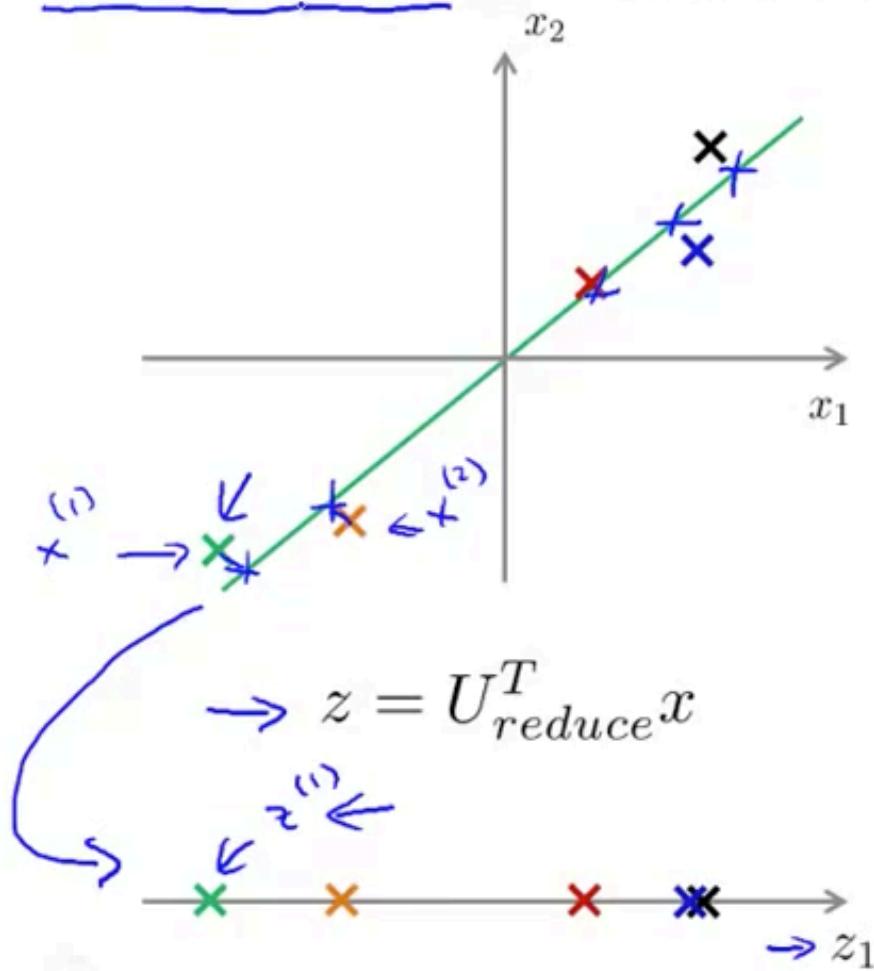
Pick smallest value of k for which

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

$k \leq 100$

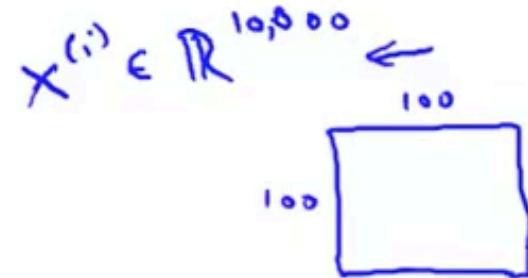
(99% of variance retained)

Reconstruction from compressed representation



Supervised learning speedup

→ $(\underline{x}^{(1)}, y^{(1)}), (\underline{x}^{(2)}, y^{(2)}), \dots, (\underline{x}^{(m)}, y^{(m)})$



Extract inputs:

Unlabeled dataset: $\underline{x}^{(1)}, \underline{x}^{(2)}, \dots, \underline{x}^{(m)} \in \mathbb{R}^{10000}$

↓ PCA

$\underline{z}^{(1)}, \underline{z}^{(2)}, \dots, \underline{z}^{(m)} \in \mathbb{R}^{1000}$

x

↓

z

New training set:

$(\underline{z}^{(1)}, y^{(1)}), (\underline{z}^{(2)}, y^{(2)}), \dots, (\underline{z}^{(m)}, y^{(m)})$

$$h_{\theta}(z) = \frac{1}{1 + e^{-\theta^T z}}$$

Note: Mapping $x^{(i)} \rightarrow z^{(i)}$ should be defined by running PCA only on the training set. This mapping can be applied as well to the examples $x_{cv}^{(i)}$ and $x_{test}^{(i)}$ in the cross validation and test sets.

Application of PCA

- Compression
 - Reduce memory/disk needed to store data
 - Speed up learning algorithm ←

Choose k by % of variance retain

- Visualization
 - $k=2$ or $k=3$

Bad use of PCA: To prevent overfitting

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

Thus, fewer features, less likely to overfit.

Bad!

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

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

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 $\{(\cancel{z^{(1)}}, y^{(1)}), \dots, (\cancel{z^{(m)}} y^{(m)})\}$
- - 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 $\underline{z^{(i)}}$.