

Machine Learning Part IV

Unsupervised Learning
K-means, PCA

K-means Algorithm

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

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

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

\swarrow
 $c^{(i)}$

Move
centroid

for $k = 1$ to K

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

$x^{(1)}, x^{(5)}, x^{(6)}, x^{(10)}$

$\rightarrow c^{(1)}=2, c^{(5)}=2, c^{(6)}=2, c^{(10)}=2$

$$\underline{\mu}_2 = \frac{1}{4} \begin{bmatrix} x^{(1)} + x^{(5)} + x^{(6)} + x^{(10)} \\ - \quad - \quad - \quad - \end{bmatrix} \in \mathbb{R}^n$$

}

Optimization Objective

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

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

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

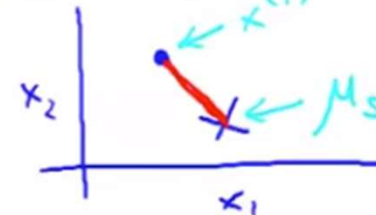
$x^{(i)} \rightarrow \underline{5}$ $\underline{c^{(i)} = 5}$ $\underline{\mu_{c^{(i)}} = \mu_5}$

Optimization objective:

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

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

Distortion



Optimization Objective

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 $c^{(1)}, c^{(2)}, \dots, c^{(n)}$ ←
(holding μ_1, \dots, μ_K fixed)

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

Random initialization

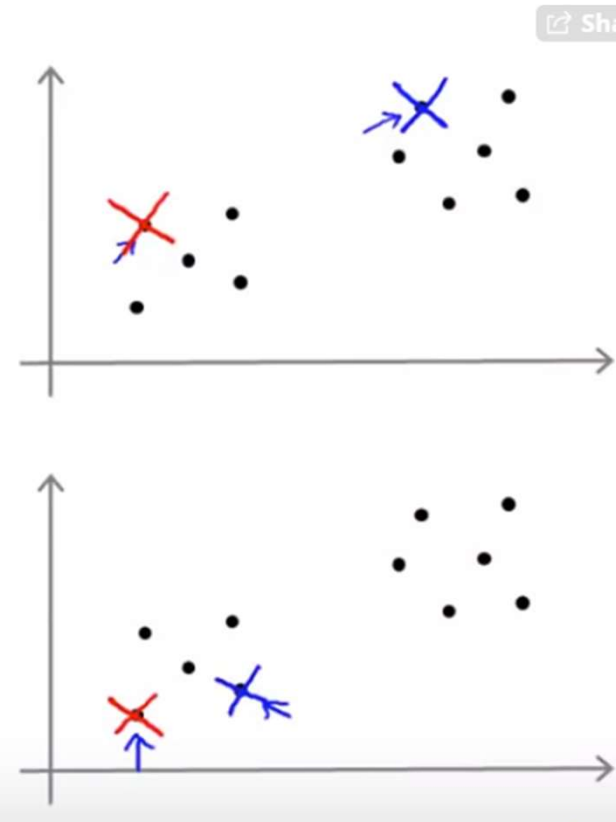
Should have $K < m$

$K=2$

Randomly pick K training examples.

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

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



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

$K = 2 - 10$

Choosing number of clusters

Suppose you run k-means using $k = 3$ and $k = 5$. You find that the cost function J is much higher for $k = 5$ than for $k = 3$. What can you conclude?

- ☐ This is mathematically impossible. There must be a bug in the code.
- ☐ The correct number of clusters is $k = 3$.
- ☒ In the run with $k = 5$, k-means got stuck in a bad local minimum. You should try re-running k-means with multiple random initializations.

Correct

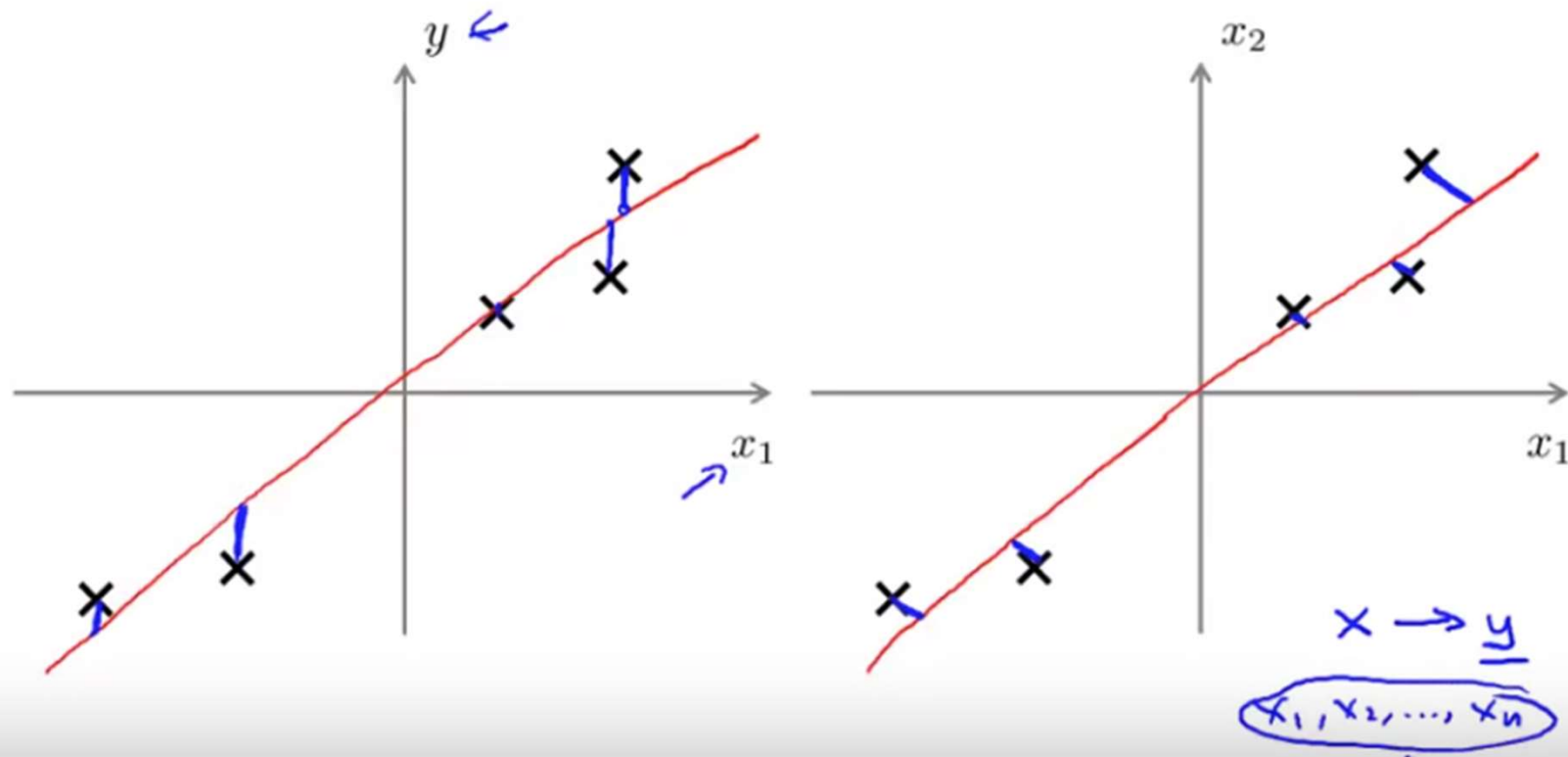
- ☐ In the run with $k = 3$, k-means got lucky. You should try re-running k-means with $k = 3$ and different random initializations until it performs no better than with $k = 5$.

Application

- Customer segmentation
- Topic group
- Image compression (as in ex7)

Principal Component Analysis

PCA is not linear regression



Principal Component Analysis (PCA) algorithm

Reduce data from n -dimensions to k -dimensions

Compute "covariance matrix":

$$\Sigma = \frac{1}{m} \sum_{i=1}^m \underbrace{(x^{(i)})}_{n \times 1} \underbrace{(x^{(i)})^T}_{1 \times n} \quad \text{Sigma}$$

$n \times n$

Compute "eigenvectors" of matrix Σ :

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

\rightarrow Singular value decomposition
 $\text{eig}(\text{Sigma})$

$n \times n$ matrix.

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

k

$$U \in \mathbb{R}^{n \times n}$$

$$u^{(1)}, \dots, u^{(k)}$$

Principal Component Analysis (PCA) algorithm

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

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

$\underbrace{\hspace{10em}}_k$

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

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

$\underbrace{\hspace{10em}}_{n \times k}$
 U_{reduce}

$z \in \mathbb{R}^k$

$$X^{(i)} = \begin{bmatrix} \text{---} (u^{(1)})^T \text{---} \\ \vdots \\ \text{---} (u^{(k)})^T \text{---} \end{bmatrix}$$

$\underbrace{\hspace{10em}}_{k \times n}$
 $\underbrace{\hspace{10em}}_{k \times 1}$

\downarrow
 $x^{(i)}$
 \sim
 $n \times 1$

$$Z = X * U[:, 1:k]$$

X is m by n,

U is n by k,

Z is m by k

Principal Component Analysis (PCA) algorithm summary

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

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

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

→ $\text{Ureduce} = U(:, 1:k) ;$

→ $\mathbf{z} = \text{Ureduce}' * \mathbf{x} ;$

↑

↑

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

$$\mathbf{x}_0 = 1$$

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

→ $\text{Sigma} = (1/m) * X' * X ;$

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

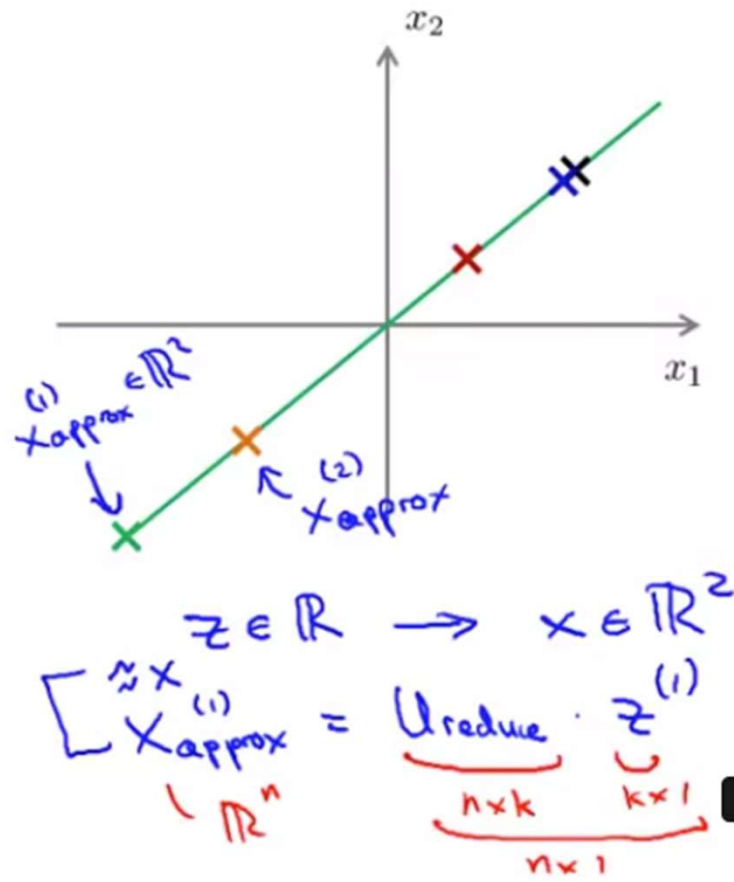
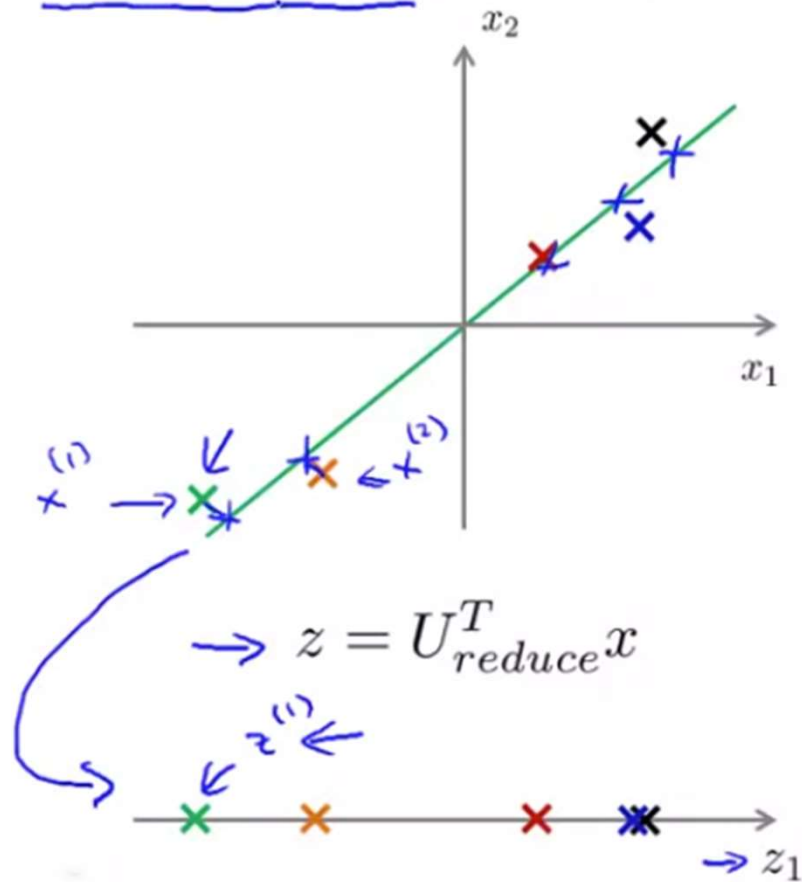
$$z = \begin{bmatrix} | & | & \dots & | \\ u^{(1)} & u^{(2)} & & u^{(k)} \\ | & | & & | \end{bmatrix}^T x = \begin{bmatrix} \text{---} & (u^{(1)})^T & \text{---} \\ \text{---} & (u^{(2)})^T & \text{---} \\ & \vdots & \\ \text{---} & (u^{(k)})^T & \text{---} \end{bmatrix} x$$

Which of the following is a correct expression for z_j ?

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

Correct

Reconstruction from compressed representation



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

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

\Rightarrow “99% of variance is retained”

Choosing k (number of principal components)

Algorithm:

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

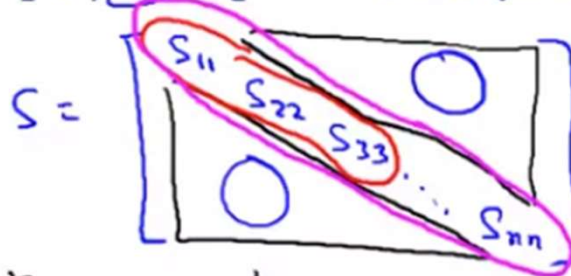
Compute $U_{reduce}, z^{(1)}, z^{(2)}, \dots, z^{(m)}, x_{approx}^{(1)}, \dots, x_{approx}^{(m)}$

Check if

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

$k=17$

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



For given k

$k=3$

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

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

Supervised learning speedup

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

Extract inputs:

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

$\downarrow \text{PCA}$

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

New training set:

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

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.



Share

$x \downarrow z$

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

$x \rightarrow z$