

WEEK - 8

⇒ In unsupervised learning, the training set is ~~is~~ of the form $\{x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(m)}\}$ without labels $y^{(i)}$.

⇒ In unsupervised learning, you are given an unlabeled dataset and are asked to find 'structure' in the data.

⇒ K-means is an iterative algorithm and it does two things —
1) = Cluster assignment step
2) = Move centroid step

⇒ K-means Algo

Input → 1) = K (number of clusters)
2) = Training set (m)

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

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

→ Repeat {

cluster assignment { for $i = 1$ to m {
 $c^{(i)} = \text{index (from 1 to } K) \text{ of cluster centroid closest to } x^{(i)}$ }
centroid move { for $k = 1$ to K {
 $\mu_k = \text{average (mean) of points assigned to cluster } k$ }
}

$$\Rightarrow c^{(i)} = \min_k \|x^{(i)} - \mu_k\|^2 \quad (\text{cluster assignment})$$

$$\Rightarrow \mu_2 = \frac{1}{Q_4} [x^{(i)}, x^{(j)}, x^{(k)}, x^{(n)}] \in R \rightarrow \text{centroid movement}$$

the points whose $c = 2$ (belong to cluster 2)

\Rightarrow Optimization Objective of K-Means algorithm

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

aka Distortion f^n

\Rightarrow K should be less than $\underbrace{m}_{\text{no. of examples}}$ ($K < m$)

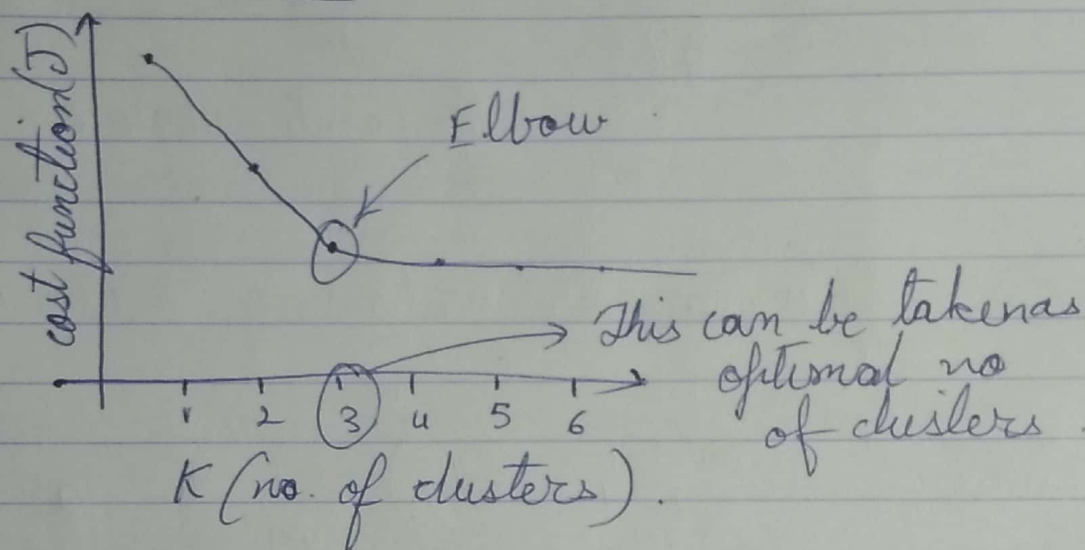
\Rightarrow If $K = 2 - 10$, then having multiple random initialization can help to find min cost f^n .

\Rightarrow If $K > 100$, then having multiple random initialization would not help much.

P.T.O

⇒ For choosing the value of K .

→ Elbow Method



⇒ If the above graph rises at any point, that means K -means got stuck in a bad local minimum. We should try ~~rerunning~~ rerunning K -means with multiple random initializations.

⇒ By hand → manually → by visualizing for each value of K and deciding ~~the~~ which no. of clusters fit correctly.

⇒ Based on a metric to use for later/downstream purpose.

⇒ Dimensionality Reduction is the second type of unsupervised machine ₂ algo.

⇒ DIMENSIONALITY REDUCTION (DR)

→ DATA COMPRESSION

⇒ Reduce data from 2D to 1D ^{or 3D to 2D} → we can compress two similar features (features who have ~~the~~ the same meaning) into a single feature which holds or represents both their data.

→ This saves memory and also helps the algo run faster.

⇒ Most Common (DR) algo is PCA
Principal Components Analysis.

→ PCA tries to find a lower dimensional surface onto which to project the data so that the distance b/w the points & the surface is minimum.

⇒ The distance b/w the points & the surface onto which they are projected is called projection error.

⇒ Formula →

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

→ For n -directions, we define k vectors $u^{(1)}, u^{(2)} \dots u^{(k)}$

PCA Algo

$$\Rightarrow \|u^{(1)}\| = \sqrt{(u_1^{(1)})^2 + (u_2^{(1)})^2}$$

⇒ Before performing PCA, we have to preprocess the data i.e. perform feature scaling or mean normalization.

Mean normalization

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

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

⇒ To reduce data from n -dimensions to k -dimensions

→ Compute → "covariance matrix"

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

Summation symbol

→ Compute "eigen vectors" of matrix Σ

$$[U, S, V] = \text{svd}(\Sigma) \rightarrow \begin{matrix} n \times n \\ \text{matrix} \end{matrix}$$

svd → Singular Value Decomposition

or

$$\Rightarrow \text{eig}(\Sigma);$$

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

and we are compressing from n dimensions
to k dimensions columns
→ we ~~do~~ take the first k ~~elements~~ of U

we get

get

$$Z = \begin{bmatrix} | & | & & | \\ u^{(1)} & u^{(2)} & \dots & u^{(k)} \\ | & | & & | \end{bmatrix}^T \times X \Rightarrow \begin{bmatrix} - & u^{(1)} & - \\ & \vdots & \\ - & u^{(k)} & - \end{bmatrix} \times \begin{matrix} \downarrow \\ X \end{matrix}$$

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

(Ureduce)

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

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

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

$\rightarrow U_{\text{reduce}} = U(:, 1:k);$

→ $Z = U_{\text{reduce}} * x;$

⇒ Reconstruction from compressed representation

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

$$\therefore x_{\text{approx}} = \underbrace{U_{\text{reduce}}}_{n \times k} * \underbrace{Z^{(i)}}_{k \times 1}$$

$n \times k$

which is very near to the original data.

⇒ If we run PCA with $k=n$ (no reduction)

$\therefore x_{\text{approx}} = x$ for every example x

\therefore The percentage of variance retained will be 100%.

$\therefore U_{\text{reduce}}$ will be a $n \times n$ matrix.

⇒ Percent/Fraction of variance ^{retained} is given by

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

⇒ k = no of principal components.

⇒ Choosing k (number of principal components)

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

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

→ Typically choose k to be smallest value so that \rightarrow

$$\frac{\frac{1}{m} \sum_{i=1}^m (\|x^{(i)} - x_{\text{approx}}^{(i)}\|)^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq \begin{matrix} 0.01 & (1\%) \\ \text{or} & \\ 0.05 & (5\%) \\ \text{or} & \\ 0.10 & (10\%) \\ \text{or} & \\ 0.15 & (85\%) \end{matrix}$$

→ "99% of variance is retained"

⇒ Algo for choosing k (very inefficient)

→ Try PCA with $k=1$

→ Compute $U_{\text{reduce}}, z^{(1)}, z^{(2)} \dots z^{(m)}, x^{(1)}_{\text{approx}}, \dots, x^{(m)}_{\text{approx}}$

→ Check if 99% of variance is retained
i.e. \rightarrow by using the big formula above

$m \Rightarrow$ no. of training examples.

\Rightarrow Easier & faster method (using $\text{svd}()$)

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

so

$$S = \begin{bmatrix} s_{11} & & & \\ & s_{22} & & \\ & & s_{33} & \\ & & & \ddots \\ & & & & s_{nn} \end{bmatrix}$$

For given k .

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

OR

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

\Rightarrow PCA tries to minimize $\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2$

\Rightarrow You should run the PCA only on the training set.