

Unsupervised learning.

Optimization objectives of k-means.

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

m_k = cluster centroid k ($m_k \in \mathbb{R}^n$)

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

Optimization objectives:

$$J(c^{(1)}, \dots, c^{(m)}, m_1, \dots, m_k) = \frac{1}{m} \sum \|x^{(i)} - m_{c^{(i)}}\|^2 \quad \leftarrow \text{Distortion Function}$$

(不变函数).

$$\min_{c^{(i)}, m_k} J(c^{(1)}, \dots, c^{(m)}, m_1, \dots, m_k)$$

最小化所有数据点与其所关联的数据中心点之间的距离之和.

Random initialization:

① $k < m$: the # of centroid n should be less than # of training sample m .

② Randomly pick k training examples.

③ Set m_1, \dots, m_k equal to these k examples.

Local Optimal:

Sometimes we might encounter the local optima, what we can do is try k-means multiple times with different random initialization.

Choose number of clusters:

Tibbons method

① Run k-means in a range of numbers.

② plot the cost corresponding to each k .

Dimensionality Reduction.

1. Data Compression:

Reduce data from 2D to 1D: use a new feature to represent those original features in a coherent relationship.

3D-2D: project data to new directions/features with lower dimensionality.

2. Data Visualization.

Visualize the data to help us analyse them.

Principal Component Analysis problem Formulation.

k -dimension
Find a vector/surface onto which to project the data so as to minimize the projection error.

找到一个向量，把所有的数据都投影到该向量上，希望平均物与误差 (distance) 小。

PCA is not linear regression

PCA

reduce dimensionality of data

minimize projected error

projection corresponding to direct vector.

Linear regression.

predict y according to feature x .

minimize prediction error.

projection corresponding to x -axis

principal component analysis algorithm

In PCA, we need to do :

- ① Compute direction vector u
- ② Compute the new features $z^{(i)}$

Algorithm = reduce from n to k .

Step 1 : 均值归一化, 计算所有特征的均值, ensure every feature has zero mean

Step 2 : Compute "Covariance matrix" Σ

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

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

$n \times n$

$$X = \begin{bmatrix} | & | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(n)} \\ | & | & | & | \end{bmatrix}$$

Step 3 : Compute "eigenvectors of matrix Σ " \Rightarrow 协方差矩阵的特征值.

$$[U \Sigma V^T] = \text{SVD}(\Sigma);$$

Σ is $n \times n$ matrix.

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

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

$u^{(1)} \sim u^{(n)}$

$$x \in \mathbb{R}^n \Rightarrow x \in \mathbb{R}^k.$$

$$Z^{(i)} = \begin{bmatrix} | & | & | & | \\ u_1^{(1)} & u_1^{(2)} & \dots & u_1^{(k)} \\ | & | & | & | \end{bmatrix}^T \quad X^{(i)} = \begin{bmatrix} -(u_1^{(1)})^T \\ -(u_1^{(2)})^T \\ \vdots \\ -(u_1^{(k)})^T \end{bmatrix} \cdot X^{(i)}$$

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

U_{reduce} .

$$\boxed{z_j = (u^{(j)})^T \cdot X} \Rightarrow \text{第 } j \text{ 个方向向量 } u^{(j)} \text{ 对应的新特征 } z_j.$$

Reconstruction from compressed representation

PCA is reducing the data from n -D to k -D ($n > k$), so that is should be reversible which means that $X_{\text{approx}} \approx X$.

$$X_{\text{approx}} = U_{\text{reduce}} \cdot Z$$

$$\text{so } X_{\text{approx}} \approx X.$$

choose the number of principal components (k):

PCA is minimizing the averaged projection error:

$$\frac{1}{m} \sum_{i=1}^m \|X^{(i)} - X_{\text{approx}}^{(i)}\|^2$$

Total variation of data, how much the data varies:

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

choose k to be smallest value so that

$$\frac{\frac{1}{m} \sum_{i=1}^m \|X^{(i)} - X_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum \|X^{(i)}\|^2} \leq 0.01 \quad \leftarrow \text{make sure the ratio is less than } 0.01 (\%).$$

99% of variance is retained.

many data are highly correlated.

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

for given k :

$$V = \begin{bmatrix} S_{11} & S_{21} & 0 \\ S_{12} & S_{22} & \dots \\ 0 & S_{13} & \dots \\ \vdots & \vdots & \ddots & S_{nn} \end{bmatrix}$$

$$\frac{\frac{1}{m} \sum \|X^{(i)} - X_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum \|X^{(i)}\|^2} = 1 - \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \leq 0.01$$

choose k for which

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

Advice:

Only run PCA on training set to learn Ureduce, then apply Ureduce in cross-validation or test sets in the same way.

$$\textcircled{1} \quad x^{(1)}, x^{(2)}, \dots, x^{(m)} \subset \mathbb{R}^{1000} \xrightarrow{\text{PCA}} z^{(1)}, z^{(2)}, \dots, z^{(m)} \subset \mathbb{R}^{100}$$

\textcircled{2} Training use new features.

\textcircled{3} In test/cross-validation, apply Ureduce to transfer X to z .

Bad use of PCA:

\textcircled{1} To prevent overfitting: PCA just throws some information of data/features, but sometimes this may cause the loss of valuable information.

\textcircled{2} To be a part of training:

Before running PLA, first try running whatever you want to do with the original/raw data, only if that doesn't do what you want, then implement PCA.