

DDA3020 Programming Report 4

Yu Kangqi

May 10, 2023

1 PCA and K-means from scratch

1.1 PCA

Denote the i th data item as \mathbf{x}_i , the subspace that we would like to project \mathbf{x}_i to is \mathcal{S} . The projection of \mathbf{x}_i to \mathcal{S} is formed by $\{\mathbf{u}_k\}_{k=1}^K$.

The projection length from x_i to u_k is defined as:

$$z_{ik} = \mathbf{u}_k^\top (\mathbf{x}_i - \boldsymbol{\mu})$$

where $\boldsymbol{\mu} = \frac{1}{N} \sum_i^N \mathbf{x}_i$.

The representation of \mathbf{x}_i will be $\tilde{\mathbf{x}}_i$:

$$\mathbf{z}_i = \mathbf{U}^\top (\mathbf{x}_i - \boldsymbol{\mu})$$

Then the reconstruction of \mathbf{x}_i will be:

$$\tilde{\mathbf{x}}_i = \mathbf{U} \mathbf{z}_i + \boldsymbol{\mu}$$

The object of PCA is to maximize the variance of the representation in the subspace \mathcal{S} . The variance is defined as:

$$\sum_i^N \|\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}}\|^2$$

where $\tilde{\boldsymbol{\mu}} = \frac{1}{N} \sum_i^N \tilde{\mathbf{x}}_i$.

So, to reduce the dimension, we can solve the following optimization problem:

$$\begin{aligned} \max_{\mathbf{U}} \quad & \sum_i^N \|\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}}\|^2 \\ \text{s.t.} \quad & \mathbf{U}^\top \mathbf{U} = \mathbf{I} \end{aligned}$$

Claim: $\tilde{\boldsymbol{\mu}} = \boldsymbol{\mu}$

Proof: $\tilde{\boldsymbol{\mu}} = \sum_i^N \tilde{\mathbf{x}}_i = \sum_i^N \mathbf{U} \mathbf{z}_i + \boldsymbol{\mu} = \mathbf{U} \sum_i^N \mathbf{z}_i + \boldsymbol{\mu} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i - \boldsymbol{\mu} + \boldsymbol{\mu} = \boldsymbol{\mu}$

The optimization problem turns to:

$$\begin{aligned} \max_{\mathbf{U}} \quad & \sum_i^N \|\tilde{\mathbf{x}}_i - \boldsymbol{\mu}\|^2 \\ \text{s.t.} \quad & \mathbf{U}^\top \mathbf{U} = \mathbf{I} \end{aligned}$$

As $\tilde{\mathbf{x}}_i = \mathbf{U} \mathbf{z}_i + \boldsymbol{\mu}$, the objective function can be written as:

$$\sum_i^N \|\mathbf{U} \mathbf{z}_i\|^2$$

Then, for $\|\mathbf{U} \mathbf{z}_i\|^2 = (\mathbf{U} \mathbf{z}_i)^\top (\mathbf{U} \mathbf{z}_i) = \mathbf{z}_i^\top \mathbf{U}^\top \mathbf{U} \mathbf{z}_i = \mathbf{z}_i^\top \mathbf{z}_i = \|\mathbf{z}_i\|^2$, the objective function can be written as:

$$\sum_i^N \|\mathbf{z}_i\|^2$$

As $\mathbf{z}_i = \mathbf{U}^\top (\mathbf{x}_i - \boldsymbol{\mu})$, the objective function can be written as:

$$\sum_i^N \|\mathbf{U}^\top (\mathbf{x}_i - \boldsymbol{\mu})\|^2 = \sum_i^N \text{Trace}(\mathbf{U}^\top (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top \mathbf{U})$$

Note that $\sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top$ is N times the empirical covariance matrix $\boldsymbol{\Sigma}$, so the objective function can be written as:

$$\text{Trace}(\mathbf{U}^\top \boldsymbol{\Sigma} \mathbf{U})$$

Until now, the optimization problem turns to:

$$\begin{aligned} \max_{\mathbf{U}} \quad & \text{Trace}(\mathbf{U}^\top \boldsymbol{\Sigma} \mathbf{U}) \\ \text{s.t.} \quad & \mathbf{U}^\top \mathbf{U} = \mathbf{I} \end{aligned}$$

The Lagrangian of the optimization problem is:

$$L(\mathbf{U}, \mathbf{\Lambda}) = \text{Trace}(\mathbf{U}^\top \mathbf{\Sigma} \mathbf{U}) + \text{Trace}(\mathbf{\Lambda}^\top (\mathbf{U}^\top \mathbf{U} - \mathbf{I}))$$

where $\mathbf{\Lambda}$ is the Lagrange multiplier.

Take first order derivative of $L(\mathbf{U}, \mathbf{\Lambda})$ with respect to \mathbf{U} and set it to zero, we have:

$$\mathbf{\Sigma} \mathbf{U} + \mathbf{U} \mathbf{\Sigma}^\top - \mathbf{U} \mathbf{\Lambda} - \mathbf{\Lambda}^\top \mathbf{U} = 0$$

$$\mathbf{\Sigma} \mathbf{U} = \mathbf{U} \mathbf{\Lambda}$$

$$\mathbf{\Sigma} \mathbf{u}_k = \lambda_k \mathbf{u}_k$$

As we can see, \mathbf{u}_k is the eigenvector of $\mathbf{\Sigma}$, and λ_k is the corresponding eigenvalue.

Further, we can turn the optimization problem to:

$$\begin{aligned} \max_{\mathbf{U}} \quad & \sum_{k=1}^K \mathbf{u}_k^\top \mathbf{\Sigma} \mathbf{u}_k \\ \text{s.t.} \quad & \mathbf{U}^\top \mathbf{U} = \mathbf{I} \end{aligned}$$

From simplicity, we can use SVD to decompose $\mathbf{\Sigma}$ to $\mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^\top$, where \mathbf{U} is the eigenvector matrix and $\mathbf{\Lambda}$ is the eigenvalue matrix.

Then, the optimization problem turns to:

$$\begin{aligned} \max_{\mathbf{U}} \quad & \sum_{k=1}^K \sum_{d=1}^D \lambda_d \mathbf{u}_k^\top \mathbf{q}_d \mathbf{q}_d^\top \mathbf{u}_k \\ \text{s.t.} \quad & \mathbf{U}^\top \mathbf{U} = \mathbf{I} \end{aligned}$$

To maximize the objective function, we need to pick K eigenvectors with top-K eigenvalues of $\mathbf{\Sigma}$.

The way that I implement the PCA algorithm is as following:

1. Calculate the empirical covariance matrix $\mathbf{\Sigma} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}^{(n)} - \boldsymbol{\mu})(\mathbf{x}^{(n)} - \boldsymbol{\mu})^\top$, where $\boldsymbol{\mu}$ is a vector containing the mean of each attribution.
2. Do SVD decomposition of $\mathbf{\Sigma}$ to obtain its D eigenvalues $\{\lambda_i\}_{i=1}^D$ and eigenvectors $\{\mathbf{q}_i\}_{i=1}^D$, and rank them from large to small according to the eigenvalues.
3. Pick the top-K eigenvectors to form the matrix $\mathbf{U} = [\mathbf{q}_1, \dots, \mathbf{q}_K] \in \mathbb{R}^{D \times K}$.
4. The new representation of $\mathbf{x}^{(n)}$ is $\mathbf{U}^\top (\mathbf{x}^{(n)} - \boldsymbol{\mu})$.

1.2 K-means

Denote the i th data item to x_i , the k th cluster center to c_k . If x_i belongs to c_k , $r_{ik} = 1$, otherwise 0.

The object of K-means is to minimize the variance between the data items and the cluster center. The variance is defined as:

$$\sum_{i=0}^n \sum_{k=0}^K r_{ik} (x_i - c_k)^2$$

The optimization perspective of K-means is solve the following objective function:

$$\begin{aligned} \min_{\mathbf{c}, \mathbf{r}} \quad & \sum_{i=0}^n \sum_{k=0}^K r_{ik} (x_i - c_k)^2 \\ \text{s.t.} \quad & \begin{cases} \mathbf{r} \in \{0, 1\}^{n \times K} \\ \sum_{k=1}^K r_{ik} = 1, \forall i \end{cases} \end{aligned}$$

We can solve this problem by using the following iterative algorithm:

1. Initialize the cluster centers $\mathbf{c} = \{c_1, c_2, \dots, c_K\}$.
2. Repeat the following steps until convergence:
 - (a) Given the cluster centers \mathbf{c} , update \mathbf{r} . So, the optimization problem in this part is as follows:

$$\begin{aligned} \min_{\mathbf{r}} \quad & \sum_{i=0}^n \sum_{k=0}^K r_{ik} (x_i - c_k)^2 \\ \text{s.t.} \quad & \begin{cases} \mathbf{r} \in \{0, 1\}^{n \times K} \\ \sum_{k=1}^K r_{ik} = 1, \forall i \end{cases} \end{aligned}$$

r_{ik} can be solved independently in the above problem, so we can solve it by:

$$r_{ik^*} = 1$$

$$\text{where } k^* = \underset{k=1}{\operatorname{argmin}} \sum_{i=1}^n \{(x_i - c_k)^2\}^K$$

- (b) Given the cluster centers \mathbf{r} , update \mathbf{c} . So, the optimization problem in this part is as follows:

$$\min_{\mathbf{c}} \sum_{i=0}^n \sum_{k=0}^K r_{ik} (x_i - c_k)^2$$

c_k can be solved independently in the above problem, so we can solve it by taking derivative of c_k and set it to 0, then we get:

$$c_k = \frac{\sum_{i=1}^n r_{ik} x_i}{\sum_{i=1}^n r_{ik}}$$

The way that I implement the K-means algorithm is as following:

1. Randomly pick K (K is the cluster number) points as the initial centroids.
2. Go into a loop, in each iteration, do the following:
 - (a) Assign each point to the nearest centroid.
 - (b) Update the centroids by calculating the mean of all points assigned to it.
3. Stop the loop when the centroids do not change or the iteration items goes to the maximum iteration.

2 Results

3 Evaluation to the results

3.1 Silhouette Coefficient

The silhouette coefficient is a measure of cluster cohesion and separation. It quantifies how well a data point fits into its assigned cluster based on two factors:

$$a(i) = \frac{\sum_{j \in C_i, j \neq i} d(i, j)}{|C_i| - 1}, \forall i \in C_i$$

$$b(i) = \min_{k \neq i} \frac{\sum_{j \in C_k} d(i, j)}{|C_k|}$$

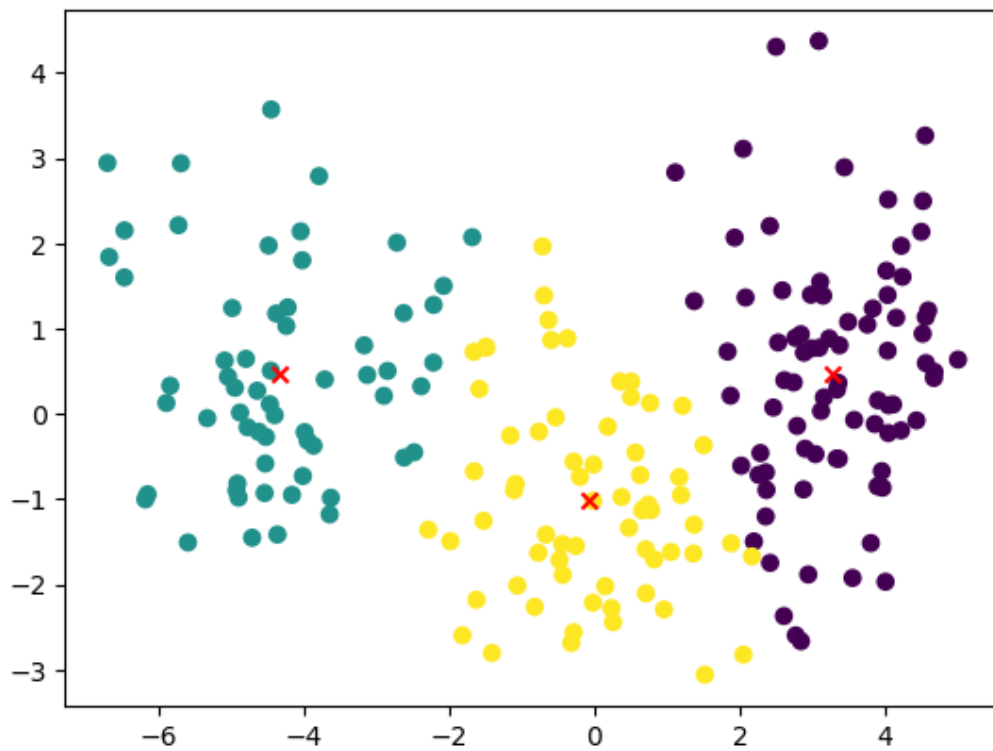


Figure 1: The result of K-means clustering for the data after PCA

where $d(i, j)$ is the distance between point i and j , C_i is the cluster that point i belongs to.

The $s(i)$ for a single point is then given as:

$$s(i) = \begin{cases} \frac{b(i) - a(i)}{\max(a(i), b(i))}, & \text{if } |C_i| > 1 \\ 0, & \text{if } |C_i| = 1 \end{cases}$$

For cluster k , $\bar{s}(k)$ is the mean of all the points in the cluster. The silhouette coefficient for the clustering is the maximum of $\bar{s}(k)$ for all clusters.

$$SC = \max \bar{s}(k)$$

3.2 Rand Index

The Rand index is a measure of the similarity between two data clusterings. Given a set of elements S and two partitions of S to compare, X and Y , define the following:

In this example, X is the observation and Y is the predicted result.

a , the number of pairs of elements in S that are in the same subset in X and in the same subset in Y .

b , the number of pairs of elements in S that are in different subsets in X and in different subsets in Y .

c , the number of pairs of elements in S that are in the same subset in X and in different subsets in Y .

d , the number of pairs of elements in S that are in different subsets in X and in the same subset in Y .

$$IR = \frac{a + b}{a + b + c + d}$$

3.3 Performance of the clustering

Here is the performance of the above clustering:

- Silhouette Coefficient: 0.5463
- Rand Index: 0.8744