

Writing Assignment 3

Issued: Saturday 14th November, 2020

Due: Friday 27th November, 2020

3.1. (K-means) Given input data $\mathcal{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$, $\mathbf{x}^{(i)} \in \mathbb{R}^d$, the k -means clustering partitions the input into k sets C_1, \dots, C_k to minimize the within-cluster sum of squares:

$$\arg \min_C \sum_{j=1}^k \sum_{\mathbf{x} \in C_j} \|\mathbf{x} - \boldsymbol{\mu}_j\|^2,$$

where $\boldsymbol{\mu}_j$ is the center of the j -th cluster:

$$\boldsymbol{\mu}_j \stackrel{\text{def}}{=} \frac{1}{|C_j|} \sum_{\mathbf{x} \in C_j} \mathbf{x}, \quad j = 1, \dots, k.$$

- (a) i. (2 points) Show that the k -means clustering problem is equivalent to minimizing the pairwise squared deviation between points in the same cluster:

$$\sum_{j=1}^k \frac{1}{2|C_j|} \sum_{\mathbf{x}, \mathbf{x}' \in C_j} \|\mathbf{x} - \mathbf{x}'\|^2.$$

- ii. (2 points) Show that the k -means clustering problem is equivalent to maximizing the between-cluster sum of squares:

$$\sum_{i=1}^k \sum_{j=1}^k |C_i| |C_j| \|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|^2.$$

- (b) Define the distortion of k -means clustering as

$$J(\{c^{(i)}\}_{i=1}^m, \{\boldsymbol{\mu}_j\}_{j=1}^k) = \sum_{i=1}^m \|\mathbf{x}^{(i)} - \boldsymbol{\mu}_{c^{(i)}}\|^2.$$

- i. (0.5 points) Show that the distortion J does not increase in each step of Lloyd's algorithm (refer to the lecture slides).
ii. (0.5 points) Does this algorithm always converge? Prove it or give a counterexample.

Solution:

$$\begin{aligned} \sum_{j=1}^k \sum_{\mathbf{x} \in C_j} \|\mathbf{x} - \boldsymbol{\mu}_j\|^2 &= \sum_{j=1}^k \sum_{\mathbf{x} \in C_j} (\|\mathbf{x}\|^2 + \|\boldsymbol{\mu}_j\|^2 - 2\langle \mathbf{x}, \boldsymbol{\mu}_j \rangle) \\ &= \sum_{\mathbf{x} \in \mathcal{X}} \|\mathbf{x}\|^2 + \sum_{j=1}^k \left(|C_j| \|\boldsymbol{\mu}_j\|^2 - 2 \left\langle \sum_{\mathbf{x} \in C_j} \mathbf{x}, \boldsymbol{\mu}_j \right\rangle \right) \\ &= \sum_{\mathbf{x} \in \mathcal{X}} \|\mathbf{x}\|^2 + \sum_{j=1}^k (|C_j| \|\boldsymbol{\mu}_j\|^2 - 2|C_j| \langle \boldsymbol{\mu}_j, \boldsymbol{\mu}_j \rangle) \\ &= \sum_{\mathbf{x} \in \mathcal{X}} \|\mathbf{x}\|^2 - \sum_{j=1}^k |C_j| \|\boldsymbol{\mu}_j\|^2. \end{aligned}$$

Note that the first term is independent of the clustering results. Hence, the k-means algorithm is equivalent to maximizing

$$\sum_{j=1}^k |C_j| \|\boldsymbol{\mu}_j\|^2.$$

(a) i. Since

$$\begin{aligned} \sum_{\mathbf{x}, \mathbf{x}' \in C_j} \|\mathbf{x} - \mathbf{x}'\|^2 &= \sum_{\mathbf{x} \in C_j} \sum_{\mathbf{x}' \in C_j} (\|\mathbf{x}\|^2 + \|\mathbf{x}'\|^2 - 2\langle \mathbf{x}, \mathbf{x}' \rangle) \\ &= 2|C_j| \sum_{\mathbf{x} \in C_j} \|\mathbf{x}\|^2 - 2 \left\langle \sum_{\mathbf{x} \in C_j} \mathbf{x}, \sum_{\mathbf{x} \in C_j} \mathbf{x}' \right\rangle \\ &= 2|C_j| \sum_{\mathbf{x} \in C_j} \|\mathbf{x}\|^2 - 2|C_j|^2 \|\boldsymbol{\mu}_j\|^2, \end{aligned}$$

we have

$$\begin{aligned} \sum_{j=1}^k \frac{1}{2|C_j|} \sum_{\mathbf{x}, \mathbf{x}' \in C_j} \|\mathbf{x} - \mathbf{x}'\|^2 &= \sum_{j=1}^k \frac{1}{2|C_j|} \left(2|C_j| \sum_{\mathbf{x} \in C_j} \|\mathbf{x}\|^2 - 2|C_j|^2 \|\boldsymbol{\mu}_j\|^2 \right) \\ &= \sum_{j=1}^k \left(\sum_{\mathbf{x} \in C_j} \|\mathbf{x}\|^2 - |C_j| \|\boldsymbol{\mu}_j\|^2 \right) \\ &= \sum_{\mathbf{x} \in \mathcal{X}} \|\mathbf{x}\|^2 - \sum_{j=1}^k |C_j| \|\boldsymbol{\mu}_j\|^2, \end{aligned}$$

which is the same as the original objective function.

ii. Note that

$$\begin{aligned} \sum_{i=1}^k \sum_{j=1}^k |C_i| |C_j| \|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|^2 &= \sum_{i=1}^k \sum_{j=1}^k |C_i| |C_j| (\|\boldsymbol{\mu}_i\|^2 + \|\boldsymbol{\mu}_j\|^2 - 2\langle \boldsymbol{\mu}_i, \boldsymbol{\mu}_j \rangle) \\ &= \sum_{i=1}^k |C_i| \|\boldsymbol{\mu}_i\|^2 \sum_{j=1}^k |C_j| + \sum_{i=1}^k |C_i| \sum_{j=1}^k |C_j| \|\boldsymbol{\mu}_j\|^2 \\ &\quad - 2 \left\langle \sum_{i=1}^k |C_i| \boldsymbol{\mu}_i, \sum_{i=1}^k |C_i| \boldsymbol{\mu}_i \right\rangle \\ &= 2|\mathcal{X}| \sum_{i=1}^k |C_i| \|\boldsymbol{\mu}_i\|^2 - 2|\mathcal{X}|^2 \|\boldsymbol{\mu}\|^2 \end{aligned}$$

where $\boldsymbol{\mu}$ is defined as the mean of all samples, i.e.,

$$\boldsymbol{\mu} = \frac{1}{|\mathcal{X}|} \sum_{\mathbf{x} \in \mathcal{X}} \mathbf{x},$$

and we have used the fact that

$$\sum_{i=1}^k |C_i| \boldsymbol{\mu}_i = \sum_{i=1}^k \sum_{\mathbf{x} \in C_i} \mathbf{x} = \sum_{\mathbf{x} \in \mathcal{X}} \mathbf{x} = |\mathcal{X}| \boldsymbol{\mu}.$$

Again, since the second term is independent of the choice of clustering, maximizing this object is equivalent to maximizing the first term, which is equivalent to the original optimization problem as illustrated at the beginning.

- (b) i. The Lloyd's algorithm contains two steps:

$$c^{(i)} \leftarrow \arg \min_j \|\mathbf{x}^{(i)} - \boldsymbol{\mu}_j\|^2, i = 1, \dots, m,$$

$$\boldsymbol{\mu}_j \leftarrow \frac{\sum_{i=1}^m \mathbf{1}\{c^{(i)} = j\} \mathbf{x}^{(i)}}{\sum_{i=1}^m \mathbf{1}\{c^{(i)} = j\}}, j = 1, \dots, k.$$

In the first step, denote the updated c as \hat{c} , then we have

$$\|\mathbf{x}^{(i)} - \boldsymbol{\mu}_{c^{(i)}}\| \geq \|\mathbf{x}^{(i)} - \boldsymbol{\mu}_{\hat{c}^{(i)}}\|, \forall i = 1, \dots, m.$$

Taking sum over i yields

$$J(\{c^{(i)}\}_{i=1}^m, \{\boldsymbol{\mu}_j\}_{j=1}^k) \geq J(\{\hat{c}^{(i)}\}_{i=1}^m, \{\boldsymbol{\mu}_j\}_{j=1}^k).$$

Similarly, in the second step, we have

$$\begin{aligned} J(\{c^{(i)}\}_{i=1}^m, \{\boldsymbol{\mu}_j\}_{j=1}^k) &= \sum_{i=1}^m \|\mathbf{x}^{(i)} - \boldsymbol{\mu}_{c^{(i)}}\|^2 \\ &= \sum_{j=1}^k \sum_{\mathbf{x} \in C_j} \|\mathbf{x} - \boldsymbol{\mu}_j\|^2 \\ &\geq \sum_{j=1}^k \sum_{\mathbf{x} \in C_j} \|\mathbf{x} - \hat{\boldsymbol{\mu}}_j\|^2 \\ &= J(\{c^{(i)}\}_{i=1}^m, \{\hat{\boldsymbol{\mu}}_j\}_{j=1}^k), \end{aligned}$$

where $\hat{\boldsymbol{\mu}}$ is the updated $\boldsymbol{\mu}$.

As a result, the distortion J does not increase in each step.

- ii. We have shown that the distortion $J > 0$ does not grow in each iteration; thus J must converge. To make the algorithm converge, we need one further assumption.

Note that the value of $\arg \min$ in the first step of the algorithm can be non-unique, which can eventually give different clustering results. Suppose that in this situation, we will always choose the smallest j that minimizes $\|\mathbf{x}^{(i)} - \boldsymbol{\mu}_j\|^2$. With this assumption, the convergence of J will imply the convergence of c and $\boldsymbol{\mu}$, since $(c, \boldsymbol{\mu})$ will be updated, if and only if the updated J is strictly less than the previous one.

3.2. (PCA) We will talk about a natural way to define PCA called Projection Residual Minimization. Suppose we have m samples $\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)} \in \mathbb{R}^n\}$, then we try

to use the projections or image vectors to represent the original data. There will be some errors (projection residuals) and naturally we hope to minimize such errors.

- (a) (1 point) First consider the case with one-dimensional projections. Let \mathbf{u} be a non-zero unit vector. The projection of sample $\mathbf{x}^{(i)}$ on vector \mathbf{u} is represented by $(\mathbf{x}^{(i)\top}\mathbf{u})\mathbf{u}$. Therefore the residual of a projection will be

$$\|\mathbf{x}^{(i)} - (\mathbf{x}^{(i)\top}\mathbf{u})\mathbf{u}\|$$

Please show that

$$\arg \min_{\mathbf{u}: \mathbf{u}^\top \mathbf{u} = 1} \|\mathbf{x}^{(i)} - (\mathbf{x}^{(i)\top}\mathbf{u})\mathbf{u}\|^2 = \arg \max_{\mathbf{u}: \mathbf{u}^\top \mathbf{u} = 1} (\mathbf{x}^{(i)\top}\mathbf{u})^2$$

- (b) (1 point) Follow the proof above and the discussion of the variance of projections in the lecture. Please show that minimizing the residual of projections is equivalent to finding the largest eigenvector of covariance matrix Σ .

$$\mathbf{u}^* = \arg \min_{\mathbf{u}: \mathbf{u}^\top \mathbf{u} = 1} \frac{1}{m} \sum_{i=1}^m \|\mathbf{x}^{(i)} - (\mathbf{x}^{(i)\top}\mathbf{u})\mathbf{u}\|^2$$

then \mathbf{u}^* is the largest eigenvector of $\Sigma = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$

- (c) (1 point) Now for a n-dimensional projection where the basis is a complete orthonormal set $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n\}$ that satisfies $\mathbf{u}_i^\top \mathbf{u}_j = \delta_{ij}$, where

$$\delta_{ij} = \begin{cases} 0 & i \neq j \\ 1 & i = j \end{cases}$$

If we pick up a k-dimension projection, the residual will be the linear combination of the remaining bases.

$$\mathbf{x}^{(i)} - \sum_{j=1}^k (\mathbf{x}^{(i)\top}\mathbf{u}_j)\mathbf{u}_j = \sum_{j=k+1}^n (\mathbf{x}^{(i)\top}\mathbf{u}_j)\mathbf{u}_j$$

Please show that

$$\min_{\mathbf{u}_1, \dots, \mathbf{u}_k: \mathbf{u}_i^\top \mathbf{u}_j = \delta_{ij}} \frac{1}{m} \sum_{i=1}^m \left\| \mathbf{x}^{(i)} - \sum_{j=1}^k (\mathbf{x}^{(i)\top}\mathbf{u}_j)\mathbf{u}_j \right\|^2 = \sum_{i=k+1}^n \lambda_i,$$

where $\lambda_1 > \lambda_2 > \dots > \lambda_n$ is the eigenvalues of Σ . It leads to the conclusion that the minimum average projection error is the sum of the eigenvalues of those eigenvectors that are orthogonal to the principal subspace.

Solution:

(a)

$$\|\mathbf{x}^{(i)} - (\mathbf{x}^{(i)\top}\mathbf{u})\mathbf{u}\|^2 = \mathbf{x}^{(i)\top}\mathbf{x}^{(i)} - (\mathbf{x}^{(i)\top}\mathbf{u})^2$$

Since $\mathbf{x}^{(i)\top}\mathbf{x}^{(i)}$ is a constant,

$$\arg \min_{\mathbf{u}: \mathbf{u}^\top \mathbf{u} = 1} \|\mathbf{x}^{(i)} - (\mathbf{x}^{(i)\top}\mathbf{u})\mathbf{u}\|^2 = \arg \max_{\mathbf{u}: \mathbf{u}^\top \mathbf{u} = 1} (\mathbf{x}^{(i)\top}\mathbf{u})^2$$

(b) Therefore

$$\begin{aligned} \arg \min_{\mathbf{u}: \mathbf{u}^T \mathbf{u} = 1} \frac{1}{n} \sum_{i=1}^n \left\| \mathbf{x}^{(i)} - (\mathbf{x}^{(i)T} \mathbf{u}) \mathbf{u} \right\|^2 &= \arg \min_{\mathbf{u}: \mathbf{u}^T \mathbf{u} = 1} \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)T} \mathbf{x}^{(i)} - (\mathbf{x}^{(i)T} \mathbf{u})^2 \\ &= \arg \max_{\mathbf{u}: \mathbf{u}^T \mathbf{u} = 1} \frac{1}{n} \sum_{i=1}^n (\mathbf{x}^{(i)T} \mathbf{u})^2 \\ &= \arg \max_{\mathbf{u}: \mathbf{u}^T \mathbf{u} = 1} \mathbf{u}^T \Sigma \mathbf{u} \end{aligned}$$

From Generalized Lagrange function that has been introduced in the class, we can find \mathbf{u}^* is the largest eigenvector.

(c)

$$\begin{aligned} &\min_{\mathbf{u}_1, \dots, \mathbf{u}_k: \mathbf{u}_i^T \mathbf{u}_j = \delta_{ij}} \frac{1}{n} \sum_{i=1}^n \left\| \mathbf{x}^{(i)} - \sum_{j=1}^k (\mathbf{x}^{(i)T} \mathbf{u}_j) \mathbf{u}_j \right\|^2 \\ &= \min_{\mathbf{u}_{k+1}, \dots, \mathbf{u}_n: \mathbf{u}_i^T \mathbf{u}_j = \delta_{ij}} \frac{1}{n} \sum_{i=1}^n \left(\sum_{j=k+1}^n (\mathbf{x}^{(i)T} \mathbf{u}_j) \mathbf{u}_j \right)^2 \\ &= \min_{\mathbf{u}_{k+1}, \dots, \mathbf{u}_n: \mathbf{u}_i^T \mathbf{u}_j = \delta_{ij}} \frac{1}{n} \sum_{i=1}^n \sum_{j=k+1}^n (\mathbf{x}^{(i)T} \mathbf{u}_j)^2 \\ &= \min_{\mathbf{u}_{k+1}, \dots, \mathbf{u}_n: \mathbf{u}_i^T \mathbf{u}_j = \delta_{ij}} \sum_{j=k+1}^n \mathbf{u}_j^T \Sigma \mathbf{u}_j \end{aligned}$$

Follow how we find the j-th principal component in the class, you will realize it is the summation of $(n - k)$ least eigenvalues. However the last step above is not that strict. The strict proof is a little bit complex with Lagrangian Multiplier Method. The good news is the conclusion is obvious. Let us omit the proof.

3.3. (Kernel PCA 2 point) Show that the conventional linear PCA algorithm is recovered as a special case of kernel PCA if we choose the linear kernel function given by $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}') = \mathbf{x}^T \mathbf{x}'$.

Solution: The proof is also not that obvious. Thus, we want to introduce an intuitionistic explanation. It would be easy to understand. At least it is a common case that satisfies $\phi(\mathbf{x})^T \phi(\mathbf{x}') = \mathbf{x}^T \mathbf{x}'$.

$$\mathbf{x}^{(i)} \mathbf{x}^{(i)T} = \phi(\mathbf{x}^{(i)}) \phi(\mathbf{x}^{(i)})^T$$

A bit more strict proof is to verify that Σ_1 and Σ_2 have the same eigenvalue.

$$\Sigma_1 = \frac{1}{m} \sum_{i=1}^m \phi(\mathbf{x}^{(i)}) \phi(\mathbf{x}^{(i)})^T$$

$$\Sigma_2 = \frac{1}{m} \sum_{i=1}^m \mathbf{x}^{(i)} \mathbf{x}^{(i)\top}$$

Note that the eigenvectors of Σ_1 can be written as a linear combination of $\phi(\mathbf{x}^{(i)})$ and the eigenvectors of Σ_2 can be written as a linear combination of $\mathbf{x}^{(i)}$

$$\begin{aligned} \phi(\mathbf{x}^{(j)\top}) \Sigma_1 \phi(\mathbf{x}^{(j)}) &= \frac{1}{m} \sum_{i=1}^m \phi(\mathbf{x}^{(j)\top}) (\phi(\mathbf{x}^{(i)}) \phi(\mathbf{x}^{(i)\top})) \phi(\mathbf{x}^{(j)}) \\ &= \frac{1}{m} \sum_{i=1}^m (\phi(\mathbf{x}^{(j)\top}) \phi(\mathbf{x}^{(i)})) (\phi(\mathbf{x}^{(i)\top}) \phi(\mathbf{x}^{(j)})) \\ &= \frac{1}{m} \sum_{i=1}^m (\mathbf{x}^{(j)\top} \mathbf{x}^{(i)}) (\mathbf{x}^{(i)\top} \mathbf{x}^{(j)}) \\ &= \frac{1}{m} \sum_{i=1}^m \mathbf{x}^{(j)\top} (\mathbf{x}^{(i)} \mathbf{x}^{(i)\top}) \mathbf{x}^{(j)} \\ &= \mathbf{x}^{(j)\top} \Sigma_2 \mathbf{x}^{(j)} \end{aligned}$$

Therefore the sample covariance matrix is still the linear PCA covariance matrix.

3.4. (Bonus Question) (SVD) In the CCA and maximal correlation lecture, we used singular value decomposition (SVD)¹ to extract important features from data. The following exercise explores several properties of SVD in details.

Suppose a rank- r matrix $A \in \mathbb{R}^{m \times n}$ has the singular value decomposition: $A = U \Sigma V^\top$, where $U = [u_1, \dots, u_r] \in \mathbb{R}^{m \times r}$, $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_r)$, $V = [v_1, \dots, v_r] \in \mathbb{R}^{n \times r}$, $U^\top U = V^\top V = I_r$, $\sigma_1 \geq \dots \geq \sigma_r > 0$.

(a) (1 point) Show that $Av_i = \sigma_i u_i$, $A^\top u_i = \sigma_i v_i$, $i = 1, \dots, r$.

(b) (1 point) The 2-norm of A is defined as

$$\|A\|_2 \stackrel{\text{def}}{=} \max_{x \in \mathbb{R}^n: \|x\| > 0} \frac{\|Ax\|}{\|x\|}.$$

Prove that $\|A\|_2 = \sigma_1$. (Hint: If $U^\top U = I$, then $\|Ux\| = \|x\|$.)

Solution:

(a) Let e_i denote the i -th standard basis in \mathbb{R}^r , then

$$Av_i = U \Sigma V^\top v_i = U \Sigma e_i = \sigma_i U e_i = \sigma_i u_i.$$

Similarly, we have $A^\top u_i = \sigma_i v_i$.

¹See https://en.wikipedia.org/wiki/Singular_value_decomposition for reference on SVD.

(b) $\forall x \in \mathbb{R}^n \setminus 0$,

$$\|Ax\| = \|U\Sigma V^T x\| = \|\Sigma V^T x\|$$

Let $\alpha = V^T x \in \mathbb{R}^r$, then we can show that

$$\|\Sigma\alpha\| \leq \sigma_1 \|\alpha\| \leq \sigma_1 \|x\|.$$

To prove the first inequality, assume $\alpha = (\alpha_1, \dots, \alpha_r)^T$, then

$$\|\Sigma\alpha\| = \sqrt{\sum_{i=1}^r \sigma_i^2 \alpha_i^2} \leq \sqrt{\sum_{i=1}^r \sigma_1^2 \alpha_i^2} = \sigma_1 \|\alpha\|.$$

To prove the second inequality, we can expand $\{v_1, \dots, v_r\}$ to an orthonormal basis $\{v_1, \dots, v_r, \dots, v_n\}$, then

$$\|\alpha\|^2 = \|V^T x\|^2 = \sum_{i=1}^r \langle v_i, x \rangle^2 \leq \sum_{i=1}^n \langle v_i, x \rangle^2 = \|x\|^2.$$

As a consequence, we have

$$\max_{x \in \mathbb{R}^n: \|x\| > 0} \frac{\|Ax\|}{\|x\|} \leq \sigma_1 = \frac{\|Av_1\|}{\|v_1\|} \leq \max_{x \in \mathbb{R}^n: \|x\| > 0} \frac{\|Ax\|}{\|x\|},$$

which means

$$\max_{x \in \mathbb{R}^n: \|x\| > 0} \frac{\|Ax\|}{\|x\|} = \sigma_1.$$