

Written Assignment 3

Issued: Sunday 10th November, 2019

Due: Sunday 24th November, 2019

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:

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

- (a) i. (1 point) 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. (1 point) 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.

3.2. (PCA) The covariance matrix of a random vector $\mathbf{x} \in \mathbb{R}^d$ is defined as

$$\text{Cov}(\mathbf{x}) \triangleq \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T] \in \mathbb{R}^{d \times d},$$

where $\mathbb{E}[\cdot]$ is the mathematical expectation.

- (a) (1 point) Show that
i. $\mathbf{u}^T \text{Cov}(\mathbf{x}) \mathbf{u} \geq 0, \forall \mathbf{u} \in \mathbb{R}^d$.
ii. $\text{tr}(\text{Cov}(\mathbf{x})) = \mathbb{E}[\|\mathbf{x} - \mathbb{E}[\mathbf{x}]\|^2]$.

- (b) (1 point) Suppose we want to estimate the covariance matrix $\hat{\mathbf{C}}$ of the dataset \mathcal{X} in 3.1 using the following formula:

$$\hat{\mathbf{C}} = \frac{1}{m-1} \sum_{i=1}^m (\mathbf{x}^{(i)} - \hat{\boldsymbol{\mu}})(\mathbf{x}^{(i)} - \hat{\boldsymbol{\mu}})^T,$$

p+1

where

$$\hat{\boldsymbol{\mu}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}^{(i)}.$$

Give the minimum value of m required such that $\hat{\mathbf{C}}$ is non-singular.

3.3. (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) (0.5 points) 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 $(\mathbf{x}^{(i)T} \mathbf{u}) \mathbf{u}$. Therefore the residual of a projection will be

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

Please show that

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

- (b) (0.5 points) 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 .

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

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

- (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 $\mathbf{u}_i^T \mathbf{u}_j = \delta_{ij}$, 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.

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

Please show that

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

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.

- 3.4. (Kernel PCA 1 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 $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}') = \mathbf{x}^T \mathbf{x}'$.
- 3.5. (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: $A = U\Sigma V^T$, where $U = [u_1, \dots, u_r] \in \mathbb{R}^{m \times r}$, $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_r)$, $V = [v_1, \dots, v_r] \in \mathbb{R}^{n \times r}$, $U^T U = V^T V = I_r$, $\sigma_1 \geq \dots \geq \sigma_r > 0$.

- (a) (1 point) Show that $Av_i = \sigma_i u_i$, $A^T u_i = \sigma_i v_i$, $i = 1, \dots, r$.
- (b) (1 point) The 2-norm of A is defined as

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

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

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