

Introduction to Machine Learning - Summer 2019

Final Exam

Instructions

1. There are 10 questions (each question is 10% of the total grade).
2. Provide full solutions (explain your answers).
3. You can keep the questions form with you (so don't write your solution on it).
4. You can use a draft notebook (you don't need to submit it).
5. Write your student ID on the notebook you are submitting.
6. Good Luck!

1 Estimation

Let $\hat{\theta}$ be an estimator of θ .

- Bias:

$$b(\hat{\theta}) = \mathbb{E}[\hat{\theta}] - \theta$$

- Variance:

$$\text{Var}(\hat{\theta}) = \mathbb{E}\left[\left(\hat{\theta} - \mathbb{E}[\hat{\theta}]\right)^2\right]$$

- MSE:

$$\text{MSE}(\hat{\theta}) \triangleq \mathbb{E}\left[\left(\hat{\theta} - \theta\right)^2\right]$$

Prove that:

$$\text{MSE}(\hat{\theta}) = \text{Var}(\hat{\theta}) + b^2(\hat{\theta})$$

2 ML

The MLE $\hat{\theta}_{ML}$ of the parameter θ is defined by:

$$\hat{\theta}_{ML} \triangleq \arg \max_{\theta} \mathcal{L}(\theta) = \arg \max_{\theta} p(\{x_i\}; \theta)$$

where $\mathcal{L}(\theta) = p(\{x_i\}; \theta)$ is the likelihood function.

Consider the random variable X with the following probability density function:

$$f_X(x) = \begin{cases} \lambda x^{-2} \exp(-\lambda/x) & x > 0 \\ 0 & \text{otherwise.} \end{cases}$$

A set $\mathcal{D} = \{x_i\}_{i=1}^N$ of i.i.d samples from f_X is given.

1. Write the log likelihood function $\ell(\lambda)$ and compute the maximum likelihood estimator $\hat{\lambda}_{ML}$ of λ .
2. Write the value of your estimation given two observations $x_1 = 2, x_2 = \frac{1}{2}$.

3 MAP

The MAP estimator $\hat{\theta}_{MAP}$ of the random variable θ is defined by:

$$\hat{\theta}_{MAP} \triangleq \arg \max_{\theta} p(\theta | \{x_i\}) = \arg \max_{\theta} p(\{x_i\} | \theta) p(\theta)$$

Consider a random variable $K \sim \text{Beta}(\alpha, \beta)$ (beta distribution) whose probability density function with parameters $\alpha, \beta > 0$ is given by

$$f_K(k; \alpha \beta) = \begin{cases} C \cdot k^{\alpha-1} (1-k)^{\beta-1} & 0 \leq k \leq 1, \\ 0 & \text{else} \end{cases}$$

where C is some constant that depends on α, β .

1. Write an expression for C as a function of α, β . It may contain sums and integrals.
2. Consider the binary random variable X (a coin toss):

$$P_X(x) = \begin{cases} 1-k & x=0, \quad (\text{tails}) \\ k & x=1, \quad (\text{heads}) \end{cases}$$

Write the probability to get one *heads* and one *tails* out of two independent coin flips (as a function of k).

3. Assume we got one *heads* and one *tails* out of two coin flips and that $k \sim \text{Beta}(3, 1)$. What is the MAP estimator \hat{k}_{MAP} of k ?

4 Non-parametric estimation

The CDF of the random variable X is given by:

$$F_X(x) = \Pr\{X \leq x\}$$

Given $\{x_i\}_{i=1}^N$, N i.i.d realizations of X we define the following estimator for F_X :

$$\hat{F}_X(x_0) = \frac{1}{N} \sum_{i=1}^N \mathbf{I}\{x_i \leq x_0\}$$

1. Compute the bias of $\hat{F}_X(x_0)$ (for a fixed x_0).
2. Compute the MSE $\left(\hat{F}_X(x_0)\right)$ (for a fixed x_0).
3. Find the limit value of the MSE as $N \rightarrow \infty$:

$$\text{MSE} \left(\hat{F}_X(x_0) \right) \xrightarrow{N \rightarrow \infty} ?$$

5 PCA I

Let $\{\mathbf{x}_i\}_{i=1}^N$ be a set of vectors such that $\mathbf{x}_i \in \mathbb{R}^d$.

- Empirical mean and covariance:

$$\boldsymbol{\mu}_x = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \triangleq \bar{\mathbf{x}}, \quad \boldsymbol{\Sigma}_x = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}_x)(\mathbf{x}_i - \boldsymbol{\mu}_x)^T$$

- The eigen decomposition of $\boldsymbol{\Sigma}_x$ is given by:

$$\boldsymbol{\Sigma}_x = \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^T$$

where $\mathbf{U} \mathbf{U}^T = \mathbf{U}^T \mathbf{U} = \mathbf{I}$, and $\boldsymbol{\Lambda}$ is a diagonal matrix with non-negative elements.

We define the following map:

$$\phi(\mathbf{x}) = \mathbf{U}^T (\mathbf{x} - \boldsymbol{\mu}_x)$$

and denote:

$$\mathbf{y}_i = \phi(\mathbf{x}_i) = \mathbf{U}^T (\mathbf{x}_i - \boldsymbol{\mu}_x), \quad \forall i \in \{1, 2, \dots, N\}$$

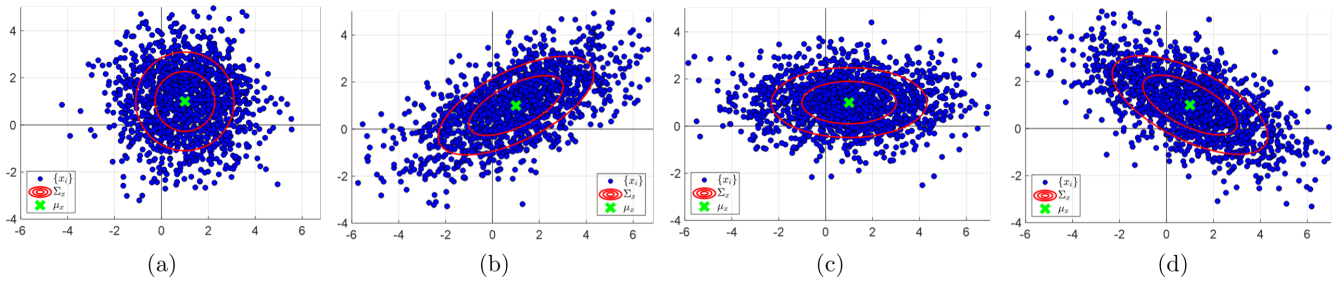
Prove that:

1. $\boldsymbol{\mu}_y = \mathbf{0}$.
2. $\boldsymbol{\Sigma}_y$ is a diagonal matrix.
3. $\|\mathbf{x}_i - \mathbf{x}_j\|_2 = \|\mathbf{y}_i - \mathbf{y}_j\|_2$ for all i, j .

6 PCA II + K-means

6.1

Match between each of the four data sets and their corresponding covariance matrix:



$$(1) \boldsymbol{\Sigma}_x = \begin{bmatrix} 5 & 0 \\ 0 & 2 \end{bmatrix}, \quad (2) \boldsymbol{\Sigma}_x = \begin{bmatrix} 5 & 2 \\ 2 & 2 \end{bmatrix}, \quad (3) \boldsymbol{\Sigma}_x = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}, \quad (4) \boldsymbol{\Sigma}_x = \begin{bmatrix} 5 & -2 \\ -2 & 2 \end{bmatrix}$$

6.2

In K-means, we seek to minimize the following objective function:

$$\min_{\{\mathcal{C}_k\}, \{\boldsymbol{\mu}_k\}} \sum_{k=1}^K \sum_{\mathbf{x}_i \in \mathcal{C}_k} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|_2^2$$

Where $K \in \mathbb{N}$ is the desired number of clusters.

Consider the set $\mathcal{D} = \{\mathbf{x}_i \in \mathbb{R}^d\}_{i=1}^N$.

Given fixed and known clusters \mathcal{C}_k such that $\bigsqcup_{k=1}^K \mathcal{C}_k = \mathcal{D}$, find the optimal centroids $\{\boldsymbol{\mu}_k \in \mathbb{R}^d\}$ which minimize the objective function.

7 MAP Classifier

- The MAP classifier is given by:

$$f_{\text{MAP}}(\mathbf{x}) = \arg \max_{C_k \in \mathcal{Y}} p(\mathbf{x}|C_k) P_{\mathcal{Y}}(C_k)$$

In the binary case ($\mathcal{Y} = \{C_1, C_2\}$), the decision rule is given by:

$$p(\mathbf{x}|C_1) P_{\mathcal{Y}}(C_1) \underset{C_2}{\overset{C_1}{\gtrless}} p(\mathbf{x}|C_2) P_{\mathcal{Y}}(C_2)$$

and the decision boundary is given by:

$$p(\mathbf{x}|C_1) P_{\mathcal{Y}}(C_1) = p(\mathbf{x}|C_2) P_{\mathcal{Y}}(C_2)$$

- The multivariate Gaussian distribution is given by ($X \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$):

$$P_X(\mathbf{x}) = |2\pi\boldsymbol{\Sigma}|^{-1} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right), \quad \boldsymbol{\mu} \in \mathbb{R}^d, \boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$$

Consider the random vector $X \in \mathbb{R}^d$ such that:

$$X|C_1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \mathbf{I})$$

$$X|C_2 \sim \mathcal{N}(\boldsymbol{\mu}_2, \mathbf{I})$$

with the prior:

$$P_{\mathcal{Y}}(C_1) = p_1, \quad P_{\mathcal{Y}}(C_2) = 1 - p_1$$

Show that the decision boundary is linear, that is, it can be written as

$$\mathbf{w}^T \mathbf{x} - b = 0$$

Write \mathbf{w} and b explicitly as a function of $p_1, \boldsymbol{\mu}_1$ and $\boldsymbol{\mu}_2$.

8 Regression

Given a set of points $\{x_i, y_i\}_{i=1}^N$, the regression L_2 error of the function \hat{f} is given by:

$$L_2(\hat{f}) = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{f}(x_i))^2$$

The M order polynomial function is given by:

$$\hat{f}(x) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{m=0}^M w_mx^m$$

where $\{w_m\}_{m=0}^M$ are the polynomial coefficients.

For a given set of points $\{x_i, y_i\}_{i=1}^N$, find the optimal $\{w_m\}_{m=0}^M$ which minimize the L_2 error.

The solution can be written in a vector form

(make sure to explain the dimensions and content of each matrix or vector).

Hint: define $\phi(x) = [1 \ x \ x^2 \ \dots \ x^M]$.

9 Linear SVM

Given a training set $\{\mathbf{x}_i, y_i\}_{i=1}^N$, the support vector machine (SVM) classification problem is given by:

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1, \quad i = 1, 2, \dots, N. \end{aligned}$$

Consider a training set $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N$ where $y_i \in \{-1, 1\}$ and \mathcal{D} is linearly separable.

Denote the optimal SVM solution by \mathbf{w}^* and b^* .

We define a new training set $\tilde{\mathcal{D}} = \{\tilde{x}_i, y_i\}_{i=1}^N$ using the following transformation:

$$\tilde{x}_i = \langle \mathbf{w}^*, \mathbf{x}_i \rangle + b^*, \quad i = 1, 2, \dots, N.$$

Is $\tilde{\mathcal{D}}$ linearly separable?

If so, write an expression for the **optimal** SVM solution (\tilde{w}, \tilde{b}) which perfectly separates $\tilde{\mathcal{D}}$.

If $\tilde{\mathcal{D}}$ is linearly non-separable, provide an example.

10 Kernels

- A **kernel** function $k : \mathbb{R}^d \times \mathbb{R}^d \longrightarrow \mathbb{R}$ satisfies:
 1. Symmetry: $k(\mathbf{x}_i, \mathbf{x}_j) = k(\mathbf{x}_j, \mathbf{x}_i)$
 2. k is positive definite, that is, the matrix $\mathbf{K}[i, j] = k(\mathbf{x}_i, \mathbf{x}_j)$ is PSD for any set $\{\mathbf{x}_i\}_i$

- A kernel function k can be written as (for some $M \in \{\mathbb{N} \cup \infty\}$):

$$k(\mathbf{x}_i, \mathbf{x}_j) = \sum_{m=1}^M \phi_m(\mathbf{x}_i) \phi_m(\mathbf{x}_j) = \langle \boldsymbol{\phi}(\mathbf{x}_i), \boldsymbol{\phi}(\mathbf{x}_j) \rangle$$

- Given a training set $\{\mathbf{x}_i, y_i\}_{i=1}^N$, the optimal weights for the SVM task are given by:

$$\mathbf{w}^* = \sum_{i=1}^N \alpha_i^* y_i \mathbf{x}_i$$

where $\{\alpha_i^*\}_i$ are the solution of the dual problem.

Hence, $\langle \mathbf{w}^*, \mathbf{x}_0 \rangle = \sum_{i=1}^N \alpha_i^* y_i \langle \mathbf{x}_i, \mathbf{x}_0 \rangle$ depends only on the inner product between the training vectors. This allows to extend SVM to the non-linear case in which:

$$\langle \mathbf{w}^*, \mathbf{x}_0 \rangle = \sum_{i=1}^N \alpha_i^* y_i k(\mathbf{x}_i, \mathbf{x}_0)$$

10.1

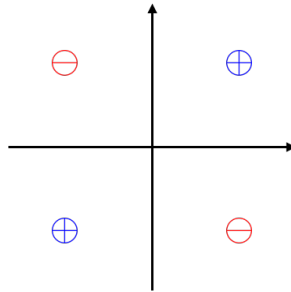
Show that the following kernel can be written as an inner product, namely:

$$k(x, y) = (1 + x \cdot y)^2 = \langle \phi(x), \phi(y) \rangle, \quad x, y \in \mathbb{R}$$

Write the transformation $\phi(x)$ explicitly.

10.2

Consider the following training set:



For each kernel, write if the classification task is linear separable or linear non-separable.

- | | |
|--|--|
| 1. $k(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y}$ | 3. $k(\mathbf{x}, \mathbf{y}) = (1 + \mathbf{x}^T \mathbf{y})^2$ |
| 2. $k(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^T \mathbf{y})^2$ | 4. $k(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{\ \mathbf{x} - \mathbf{y}\ _2^2}{2\sigma^2}\right)$
(you may chose the value of σ) |