

Introduction to Machine Learning - Summer 2019

Final Exam – Solution

1 Estimation

$$\begin{aligned}\text{MSE}(\hat{\theta}) &\triangleq \mathbb{E} \left[(\hat{\theta} - \theta)^2 \right] \\&= \mathbb{E} \left[(\hat{\theta} - \mathbb{E}[\hat{\theta}] + \mathbb{E}[\hat{\theta}] - \theta)^2 \right] \\&= \mathbb{E} \left[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2 + 2(\hat{\theta} - \mathbb{E}[\hat{\theta}])(\mathbb{E}[\hat{\theta}] - \theta) + (\mathbb{E}[\hat{\theta}] - \theta)^2 \right] \\&= \underbrace{\mathbb{E} \left[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2 \right]}_{=V(\hat{\theta})} + \underbrace{2\mathbb{E} \left[(\hat{\theta} - \mathbb{E}[\hat{\theta}]) \right]}_{=0} (\mathbb{E}[\hat{\theta}] - \theta) + \underbrace{(\mathbb{E}[\hat{\theta}] - \theta)^2}_{b^2(\hat{\theta})} \\&= V(\hat{\theta}) + b^2(\hat{\theta})\end{aligned}$$

2 ML

2.1

$$\begin{aligned}\ell(\lambda) &= \log p(\mathcal{D}; \lambda) \\&= \log \left(\prod_{i=1}^N \lambda x_i^{-2} \exp\left(-\frac{\lambda}{x_i}\right) \right) \\&= N \log(\lambda) - 2 \sum_{i=1}^N \log(x_i) - \lambda \sum_{i=1}^N \frac{1}{x_i} \\&\Rightarrow \hat{\lambda}_{ML} = \arg \max_{\lambda} \ell(\lambda)\end{aligned}$$

$$\begin{aligned}\frac{\partial}{\partial \lambda} \ell(\lambda) &= 0 \\&\Rightarrow \frac{N}{\lambda} - \sum_{i=1}^N \frac{1}{x_i} = 0 \\&\Rightarrow \hat{\lambda}_{ML} = \frac{N}{\sum_{i=1}^N \frac{1}{x_i}}\end{aligned}$$

2.2

$$\hat{\lambda}_{ML} = \frac{2}{2 + \frac{1}{2}} = \frac{4}{5}$$

3 MAP

3.1

$$\begin{aligned} \int_{-\infty}^{\infty} f_K(k; \alpha \beta) dk &= 1 \\ \Rightarrow \int_0^1 C \cdot k^{\alpha-1} (1-k)^{\beta-1} dk &= 1 \\ \Rightarrow C &= \frac{1}{\int_0^1 k^{\alpha-1} (1-k)^{\beta-1} dk} \end{aligned}$$

3.2

$$\begin{aligned} P(\text{Heads and Tails}) &= P(\text{Heads} \cap \text{Tails}) + P(\text{Tails} \cap \text{Heads}) \\ &= 2P(\text{Heads})P(\text{Tails}) \\ &= 2k(1-k). \end{aligned}$$

3.3

$$\begin{aligned} \hat{k}_{MAP} &= \arg \max_{0 \leq k \leq 1} P(\text{Heads} \cap \text{Tails}) \cdot f_K(k; \alpha \beta) \\ &= \arg \max_{0 \leq k \leq 1} 2k(1-k)C \cdot k^{\alpha-1} (1-k)^{\beta-1} \\ &= \arg \max_{0 \leq k \leq 1} k(1-k)k^{\alpha-1} (1-k)^{\beta-1} \\ &= \arg \max_{0 \leq k \leq 1} k^{\alpha} (1-k)^{\beta} \\ &= \arg \max_{0 \leq k \leq 1} k^3 (1-k) \end{aligned}$$

$$\begin{aligned} \frac{d}{dk} (k^3(1-k)) &= 0 \\ 3k^2(1-k) - k^3 &= 0 \\ k^2(3-3k-k) &= 0 \\ \Rightarrow \hat{k}_{MAP} &= \frac{3}{4} \end{aligned}$$

4 Non-parametric estimation

1.

$$\begin{aligned} \mathbb{E} [\hat{F}_X(x_0)] &= \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N \mathbf{I} \{x_i \leq x_0\} \right] = \frac{1}{N} \sum_{i=1}^N \mathbb{E} [\mathbf{I} \{x_i \leq x_0\}] \\ &= \frac{1}{N} \sum_{i=1}^N \Pr \{x_i \leq x_0\} = \Pr \{x_1 \leq x_0\} \\ &= F_X(x_0) \\ \Rightarrow b(\hat{F}_X(x_0)) &= \mathbb{E} [\hat{F}_X(x_0)] - F_X(x_0) = 0 \end{aligned}$$

Hence, \hat{F}_X is unbiased.

2. In Lecture 1, we proved that:

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

$$\begin{aligned} \Rightarrow \text{MSE}(\hat{F}_X(x_0)) &= \underbrace{b^2(\hat{F}_X(x_0))}_{=0} + \text{Var}(\hat{F}_X(x_0)) = \text{Var}\left(\frac{1}{N} \sum_{i=1}^N \underbrace{\mathbf{I}\{x_i \leq x_0\}}_{\triangleq Y_i}\right) \\ &= \frac{1}{N^2} \text{Var}\left(\sum_{i=1}^N Y_i\right), \quad \{Y_i\}_i \text{ are i.i.d} \\ &= \frac{1}{N^2} N \text{Var}(Y_1) = \frac{\mathbb{E}[Y_1^2] - \mathbb{E}^2[Y_1]}{N} \\ &= \frac{F_X(x_0) - F_X^2(x_0)}{N} \end{aligned}$$

3. Using the previous result, we have:

$$\text{MSE}(\hat{F}_X(x_0)) \xrightarrow{N \rightarrow \infty} 0$$

5 PCA I

1. By definition:

$$\boldsymbol{\mu}_y = \frac{1}{N} \sum_{i=1}^N \mathbf{y}_i = \frac{1}{N} \sum_{i=1}^N \mathbf{U}^T (\mathbf{x}_i - \boldsymbol{\mu}_x) = \mathbf{U}^T \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}_x) = \mathbf{U}^T \left(\underbrace{\frac{1}{N} \sum_{i=1}^N \mathbf{x}_i}_{\boldsymbol{\mu}_x} - \boldsymbol{\mu}_x \right) = 0$$

2. Consider the eigen-value decomposition: $\boldsymbol{\Sigma}_x = \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^T$ (\mathbf{U} is unitary) Hence,

$$\begin{aligned} \boldsymbol{\Sigma}_y &= \frac{1}{N} \sum_{i=1}^N (\mathbf{y}_i - \boldsymbol{\mu}_y) (\mathbf{y}_i - \boldsymbol{\mu}_y)^T = \frac{1}{N} \sum_{i=1}^N \mathbf{y}_i \mathbf{y}_i^T \\ &= \frac{1}{N} \sum_{i=1}^N \mathbf{U}^T (\mathbf{x}_i - \boldsymbol{\mu}_x) (\mathbf{x}_i - \boldsymbol{\mu}_x)^T \mathbf{U} = \mathbf{U}^T \underbrace{\frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}_x) (\mathbf{x}_i - \boldsymbol{\mu}_x)^T}_{=\boldsymbol{\Sigma}_x} \mathbf{U} \\ &= \mathbf{U}^T \boldsymbol{\Sigma}_x \mathbf{U} = \mathbf{U}^T \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^T \mathbf{U} = \boldsymbol{\Lambda} \end{aligned}$$

3. Since \mathbf{U}^T is unitary, we have $\|\mathbf{U}^T \mathbf{v}\|_2 = \|\mathbf{v}\|_2$, thus:

$$\|\mathbf{y}_i - \mathbf{y}_j\|_2 = \|\mathbf{U}^T (\mathbf{x}_i - \boldsymbol{\mu}_x) - \mathbf{U}^T (\mathbf{x}_j - \boldsymbol{\mu}_x)\|_2 = \|\mathbf{U}^T (\mathbf{x}_i - \mathbf{x}_j)\|_2 = \|\mathbf{x}_i - \mathbf{x}_j\|_2$$

6

6.1 PCA II

- (a) – (3)
- (b) – (2)
- (c) – (1)
- (d) – (4)

6.2 K-means

For fixed clusters $\{\mathcal{C}_k\}$,

we can compare the gradient with respect to $\boldsymbol{\mu}_s$ to zero:

$$\begin{aligned} \nabla_{\boldsymbol{\mu}_s} \left(\sum_{k=1}^K \sum_{\mathbf{x}_i \in \mathcal{C}_k} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|_2^2 \right) &= \mathbf{0} \\ -2 \sum_{\mathbf{x}_i \in \mathcal{C}_s} (\mathbf{x}_i - \boldsymbol{\mu}_s) &= \mathbf{0} \\ \Rightarrow \boxed{\boldsymbol{\mu}_s = \frac{1}{|\mathcal{C}_s|} \sum_{\mathbf{x}_i \in \mathcal{C}_s} \mathbf{x}_i} \end{aligned}$$

In words, the optimal centroid $\boldsymbol{\mu}_k$ of the k th cluster \mathcal{C}_k is the mean of the cluster.

7 MAP classifier

$$\begin{aligned} p(\mathbf{x}|C_1) p_Y(C_1) &= p(\mathbf{x}|C_2) p_Y(C_2) \\ p(\mathbf{x}|C_1) p_1 &= p(\mathbf{x}|C_2) (1 - p_1) \\ \frac{p_1}{1 - p_1} e^{-\frac{1}{2} \|\mathbf{x} - \boldsymbol{\mu}_1\|_2^2} &= e^{-\frac{1}{2} \|\mathbf{x} - \boldsymbol{\mu}_2\|_2^2} \\ \log \left(\frac{p_1}{1 - p_1} \right) - \frac{1}{2} \|\mathbf{x} - \boldsymbol{\mu}_1\|_2^2 &= -\frac{1}{2} \|\mathbf{x} - \boldsymbol{\mu}_2\|_2^2 \\ 2 \log \left(\frac{p_1}{1 - p_1} \right) - \|\mathbf{x}\|^2 + 2\boldsymbol{\mu}_1^T \mathbf{x} - \|\boldsymbol{\mu}_1\|^2 &= -\|\mathbf{x}\|^2 + 2\boldsymbol{\mu}_2^T \mathbf{x} - \|\boldsymbol{\mu}_2\|^2 \\ 2 \log \left(\frac{p_1}{1 - p_1} \right) + \|\boldsymbol{\mu}_2\|^2 - \|\boldsymbol{\mu}_1\|^2 + 2\boldsymbol{\mu}_1^T \mathbf{x} - 2\boldsymbol{\mu}_2^T \mathbf{x} &= 0 \\ \underbrace{2 \log \left(\frac{p_1}{1 - p_1} \right) + \|\boldsymbol{\mu}_2\|^2 - \|\boldsymbol{\mu}_1\|^2}_{=-b} + \underbrace{2(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \mathbf{x}}_{=\mathbf{w}^T} &= 0 \\ \mathbf{w}^T \mathbf{x} - b &= 0 \end{aligned}$$

8 Regression

$$\begin{aligned} \mathbf{w} &\triangleq \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_M \end{bmatrix}, \quad \boldsymbol{\phi}(x) \triangleq \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^M \end{bmatrix} \\ \Rightarrow \hat{f}(x) &= \boldsymbol{\phi}^T(x) \cdot \mathbf{w} \end{aligned}$$

$$\Rightarrow L(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \left(y_i - \hat{f}(x_i) \right)^2 = \frac{1}{N} \sum_{i=1}^N \left(y_i - \boldsymbol{\phi}^T(x_i) \mathbf{w} \right)^2 = \frac{1}{N} \|\mathbf{y} - \boldsymbol{\Phi} \mathbf{w}\|_2^2$$

where:

$$\boldsymbol{\Phi} \triangleq \begin{bmatrix} | & & | \\ \boldsymbol{\phi}(x_1) & \cdots & \boldsymbol{\phi}(x_N) \\ | & & | \end{bmatrix}^T = \begin{bmatrix} 1 & x_1 & x_1^2 & \cdots & x_1^M \\ \vdots & \vdots & \vdots & \cdots & \vdots \\ 1 & x_N & x_N^2 & \cdots & x_N^M \end{bmatrix}$$

Thus, the optimal \mathbf{w} is given by:

$$\Rightarrow \mathbf{w} = \left(\Phi^T \Phi \right)^{-1} \Phi^T \mathbf{y} = \Phi^\dagger \mathbf{y}$$

9 Linear SVM

The SVM solution satisfies:

$$y_i (\langle \mathbf{w}^*, \mathbf{x}_i \rangle + b) \geq 1$$

Therefore,

$$\Rightarrow \begin{cases} \langle \mathbf{w}^*, \mathbf{x}_i \rangle + b \geq 1 & y_i = 1 \\ \langle \mathbf{w}^*, \mathbf{x}_i \rangle + b \leq -1 & y_i = -1 \end{cases}$$

Hence:

$$\Rightarrow \begin{cases} \tilde{x}_i \geq 1 & y_i = 1 \\ \tilde{x}_i \leq -1 & y_i = -1 \end{cases}$$

and specifically, the support vectors satisfy:

$$\Rightarrow \begin{cases} \tilde{x}_i = 1 & y_i = 1 \\ \tilde{x}_i = -1 & y_i = -1 \end{cases}, \quad \text{for all support vector } \tilde{x}_i$$

In words, all the positive samples are on the positive side (of the real line) and all the negative samples are on the negative side.

Thus, $\tilde{\mathcal{D}}$ is indeed linear separable and the SVM solution is given by:

$$\tilde{w}^* = 1, \quad \tilde{b} = 0$$

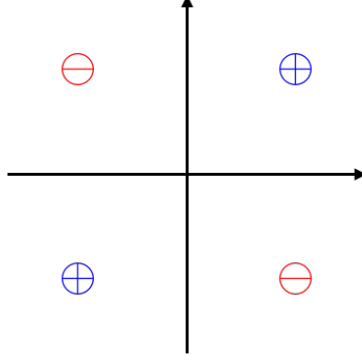
since the problem is centered around the origin.

10 Kernels

10.1

$$\begin{aligned} k(x, y) &= (1 + xy)^2 \\ &= 1 + 2xy + x^2 y^2 \\ &= \left\langle \begin{bmatrix} 1 \\ \sqrt{2}x \\ x^2 \end{bmatrix}, \begin{bmatrix} 1 \\ \sqrt{2}y \\ y^2 \end{bmatrix} \right\rangle \\ \Rightarrow \phi(x) &= \begin{bmatrix} 1 \\ \sqrt{2}x \\ x^2 \end{bmatrix} \end{aligned}$$

10.2



The original data set is:

$$\mathbf{X} = \begin{bmatrix} 1 & -1 & -1 & 1 \\ 1 & 1 & -1 & -1 \end{bmatrix}, \quad \mathbf{y} = [+1 \quad -1 \quad +1 \quad -1]$$

1. This problem is not linear separable with the standard kernel $k(\mathbf{x}, \mathbf{z}) = \mathbf{x}^T \mathbf{z}$.
- 2.

$$\begin{aligned} k(\mathbf{x}, \mathbf{z}) &= (\mathbf{x}^T \mathbf{z})^2 \\ &= (x_1 z_1 + x_2 z_2)^2 \\ &= x_1^2 z_1^2 + 2x_1 x_2 z_1 z_2 + x_2^2 z_2^2 \\ &= \left\langle \begin{bmatrix} x_1^2 \\ \sqrt{2}x_1 x_2 \\ x_2^2 \end{bmatrix}, \begin{bmatrix} z_1^2 \\ \sqrt{2}z_1 z_2 \\ z_2^2 \end{bmatrix} \right\rangle \end{aligned}$$

Hence, the features in the kernel space are:

$$\Rightarrow \Phi = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & -1 \\ 1 & 1 & 1 & 1 \end{bmatrix}, \quad \mathbf{y} = [+1 \quad -1 \quad +1 \quad -1]$$

This can be linearly separated using the second row.

3. The kernel $k(\mathbf{x}, \mathbf{z}) = (1 + \mathbf{x}^T \mathbf{z})^2$ contains all the features of the kernel $(\mathbf{x}^T \mathbf{z})^2$ and thus, can also be linearly separated.
4. Any data set can be linearly separated with a Gaussian kernel (and some small σ).