

Written Assignment 2

Issued: Wednesday 16th October, 2019

Due: Wednesday 30th October, 2019

- 2.1. (2 points) Define the **design matrix** \mathbf{X} to be the m-by-n matrix that the training examples input values in its rows. Geometrically the solution of least-squares could be interpreted of as a vector in an M-dimensional space whose coordinates are $\mathbf{y} = [y^{(1)}, y^{(2)}, \dots, y^{(m)}]^T$. The least-squares regression function is obtained by finding the orthogonal projection of the target vector \mathbf{y} onto the subspace spanned by the column vectors of \mathbf{X} , in which the i-th column vector is denoted as \mathbf{X}_i .

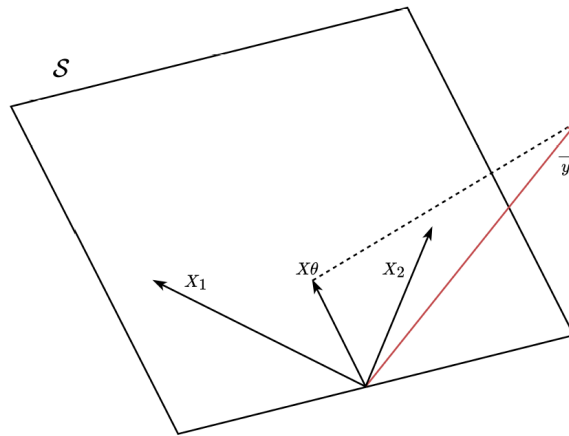


Figure 1: Projection of \mathbf{y} on column space of \mathbf{X}

As shown in Figure 1, please show that the matrix

$$\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

takes any vector \mathbf{v} and projects it onto the space spanned by the columns of \mathbf{X} . Use this result to show that the least-squares solution $\boldsymbol{\theta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ correspond to an orthogonal projection of the vector \mathbf{y} onto the column space of \mathbf{X} .

Solution: \mathbf{X} can be written as $[\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n]$.

Therefore,

$$\mathbf{X}\boldsymbol{\theta} = \sum_{i=1}^n \theta_i \mathbf{X}_i$$

It is on the column space of \mathbf{X}

Note that

$$\begin{aligned} \mathbf{X}^T \mathbf{X} \boldsymbol{\theta} &= \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ &= (\mathbf{X}^T \mathbf{X}) (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ &= \mathbf{X}^T \mathbf{y} \end{aligned}$$

which means

$$\langle \mathbf{X}_i, \mathbf{X}\boldsymbol{\theta} \rangle = \langle \mathbf{X}_i, \mathbf{y} \rangle$$

Therefore,

$$\mathbf{X}\boldsymbol{\theta} = \sum_{i=1}^n \langle \mathbf{X}_i, \mathbf{X}\boldsymbol{\theta} \rangle \mathbf{X}_i = \sum_{i=1}^n \langle \mathbf{X}_i, \mathbf{y} \rangle \mathbf{X}_i$$

$\mathbf{X}\boldsymbol{\theta}$ is the orthogonal projection of the vector \mathbf{y} onto the column space of \mathbf{X} .

- 2.2. (2 points) Suppose we are given a dataset $\{(\mathbf{x}^{(i)}, y^{(i)}): i = 1, 2, \dots, m\}$ consisting of m independent examples, where $\mathbf{x}^{(i)} \in \mathbb{R}^n$ are n -dimension vector, and $y^{(i)} \in \{1, 2, \dots, k\}$. We will model the joint distribution of (\mathbf{x}, y) according to:

$$y^{(i)} \sim \text{Multinomial}(\phi)$$

$$\mathbf{x}^{(i)} | y^{(i)} = j \sim \mathcal{N}(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$$

where the parameter ϕ_j gives $p(y^{(i)} = j)$ for each $j \in \{1, 2, \dots, k\}$.

In Gaussian Discriminant Analysis (GDA), Linear Discriminant Analysis (LDA) just assume that the classes have a common covariance matrix $\Sigma_j = \Sigma, \forall j$. If the Σ_j are not assumed to be equal, we get Quadratic Discriminant Analysis (QDA). The estimates for QDA are similar to those for LDA, except that separate covariance matrices must be estimated for each class. Give the maximum likelihood estimate of Σ_j in the case that $k = 2$.

Solution:

$$\begin{aligned} \log L(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \phi) &= \log \prod_{i=1}^m P(\mathbf{x}_i | y_i; \boldsymbol{\mu}_{y_i}, \boldsymbol{\Sigma}_{y_i}) P(y_i; \phi_i) \\ &= \log \prod_{i=1}^m \prod_{j=1}^k \mathbf{1}_{\{y_i=j\}} \frac{1}{(2\pi)^{\frac{n}{2}} |\boldsymbol{\Sigma}_j|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu}_j)^T \boldsymbol{\Sigma}_j^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_j)\right) P(y_i=j; \phi_j) \\ &= \sum_{i=1}^m \sum_{j=1}^k \mathbf{1}_{\{y_i=j\}} \left(-\frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu}_j)^T \boldsymbol{\Sigma}_j^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_j) - \frac{n}{2} \log(2\pi) - \frac{1}{2} \log |\boldsymbol{\Sigma}_j| + \log P(y_i=j; \phi_j)\right) \end{aligned}$$

Note that

$$\begin{aligned} \frac{\partial \log |\boldsymbol{\Sigma}|}{\partial \boldsymbol{\Sigma}} &= \frac{1}{|\boldsymbol{\Sigma}|} \frac{\partial |\boldsymbol{\Sigma}|}{\partial \boldsymbol{\Sigma}} = \frac{|\boldsymbol{\Sigma}| \boldsymbol{\Sigma}^{-1}}{|\boldsymbol{\Sigma}|} = \boldsymbol{\Sigma}^{-1} \\ \frac{\partial \mathbf{v}^T \boldsymbol{\Sigma}^{-1} \mathbf{v}}{\partial \boldsymbol{\Sigma}} &= -\boldsymbol{\Sigma}^{-1} \mathbf{v} \mathbf{v}^T \boldsymbol{\Sigma}^{-1T} \end{aligned}$$

The equation above may need some procedures to be proved.

$$\begin{aligned} \left(\frac{\partial \mathbf{v}^T \boldsymbol{\Sigma}^{-1} \mathbf{v}}{\partial \boldsymbol{\Sigma}} \right)_{ij} &= \frac{\partial \mathbf{v}^T \boldsymbol{\Sigma}^{-1} \mathbf{v}}{\partial \Sigma_{ij}} = \mathbf{v}^T \frac{\partial \boldsymbol{\Sigma}^{-1}}{\partial \Sigma_{ij}} \mathbf{v} \\ &= -\mathbf{v}^T \boldsymbol{\Sigma}^{-1} \frac{\partial \boldsymbol{\Sigma}}{\partial \Sigma_{ij}} \boldsymbol{\Sigma}^{-1} \mathbf{v} \\ &= -\mathbf{v}^T \boldsymbol{\Sigma}^{-1T} \frac{\partial \boldsymbol{\Sigma}}{\partial \Sigma_{ij}} \boldsymbol{\Sigma}^{-1} \mathbf{v} \\ &= -\left(\boldsymbol{\Sigma}^{-1} \mathbf{v} \mathbf{v}^T \boldsymbol{\Sigma}^{-1T} \right)_{ij} \end{aligned}$$

Therefore, the derivative on likelihood function is

$$\begin{aligned}\frac{\partial \log L}{\partial \Sigma_j} &= \frac{1}{2} \sum_{i=1}^m \mathbf{1}\{y_i = j\} \left(-\Sigma_j^{-1} + \Sigma_j^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_j)(\mathbf{x}_i - \boldsymbol{\mu}_j)^T \Sigma_j^{-1} \right) \\ &= \mathbf{0}\end{aligned}$$

The result is

$$\Sigma_j = \frac{\sum_{i=1}^m \mathbf{1}\{y_i = j\} (\mathbf{x}_i - \boldsymbol{\mu}_j)(\mathbf{x}_i - \boldsymbol{\mu}_j)^T}{\sum_{i=1}^m \mathbf{1}\{y_i = j\}}$$

where $j = 1, 2$

2.3. Suppose the data are linearly separable. The optimization problem of SVM is

$$\begin{aligned}\underset{\mathbf{w}, b}{\text{minimize}} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{subject to} \quad & y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \quad i = 1, \dots, l,\end{aligned}\tag{P}$$

and let (\mathbf{w}^*, b^*) denote its optimal solution.

(a) (2 points) Show that

$$b^* = -\frac{1}{2} \left(\max_{i: y_i = -1} \mathbf{w}^{*T} \mathbf{x}_i + \min_{i: y_i = 1} \mathbf{w}^{*T} \mathbf{x}_i \right).$$

The corresponding Lagrange dual problem is given by

$$\begin{aligned}\underset{\boldsymbol{\alpha}}{\text{maximize}} \quad & \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \\ \text{subject to} \quad & \alpha_i \geq 0, \quad i = 1, \dots, l, \\ & \sum_{i=1}^l \alpha_i y_i = 0.\end{aligned}\tag{D}$$

Suppose the optimal solution of (D) is $\boldsymbol{\alpha}^* = (\alpha_1^*, \dots, \alpha_l^*)^T$, from the KKT conditions we know that

$$\begin{aligned}\mathbf{w}^* &= \sum_{i=1}^l \alpha_i^* y_i \mathbf{x}_i, \\ \sum_{i=1}^l \alpha_i^* [y_i(\mathbf{w}^{*T} \mathbf{x}_i + b^*) - 1] &= 0.\end{aligned}\tag{1}$$

(b) (1 point) Based on (1), verify that

$$\frac{1}{2} \|\mathbf{w}^*\|_2^2 = \sum_{i=1}^l \alpha_i^* - \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i^* \alpha_j^* y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle = \frac{1}{2} \sum_{i=1}^l \alpha_i^*.$$

Solution:

- (a) *The key point is to show that $\gamma_1 = \gamma_2 = 1$ (γ is defined below). Geometrically, it means that we can always find samples lying on the two boundaries $\mathbf{w}^T \mathbf{x} + b = \pm 1$ when the parameters are optimal.*

Since the data are linearly separable, we have $\|\mathbf{w}^*\| < \infty$. Define γ_1, γ_2 as

$$\gamma_1 \triangleq \min_{i: y_i=1} \mathbf{w}^{*T} \mathbf{x}_i + b^*,$$

$$\gamma_2 \triangleq - \left(\max_{i: y_i=-1} \mathbf{w}^{*T} \mathbf{x}_i + b^* \right),$$

then it suffices to prove $\gamma_1 = \gamma_2 = 1$. Note that the constraints of (P) imply $\gamma_1 \geq 1, \gamma_2 \geq 1$. If $\gamma_1 + \gamma_2 > 2$, we can define $\mathbf{w}^\dagger, b^\dagger$ as

$$\hat{\mathbf{w}}^\dagger \triangleq \frac{2}{\gamma_1 + \gamma_2} \mathbf{w}^*, \quad b^\dagger \triangleq \frac{2}{\gamma_1 + \gamma_2} \left(b^* - \frac{\gamma_1 - \gamma_2}{2} \right),$$

such that $(\mathbf{w}^\dagger, b^\dagger)$ is a feasible solution of (P) and $\|\mathbf{w}^\dagger\|_2^2 \leq \|\mathbf{w}^*\|_2^2$, which contradicts to the optimality of (\mathbf{w}^*, b^*) .

- (b) *For convex optimization problems, the KKT conditions are sufficient conditions for strong duality.*

From the first equality of (1), we have

$$\mathbf{w}^{*T} \mathbf{w}^* = \left(\sum_{i=1}^l \alpha_i^* y_i \mathbf{x}_i \right)^T \left(\sum_{j=1}^l \alpha_j^* y_j \mathbf{x}_j \right) = \sum_{i=1}^l \sum_{j=1}^l \alpha_i^* \alpha_j^* y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

Moreover,

$$\mathbf{w}^{*T} \mathbf{w}^* = \mathbf{w}^{*T} \left(\sum_{i=1}^l \alpha_i^* y_i \mathbf{x}_i \right) = \sum_{i=1}^l \alpha_i^* y_i \mathbf{w}^{*T} \mathbf{x}_i = \sum_{i=1}^l \alpha_i^*,$$

where the last equality has used the second equation of (1). As a result,

$$\|\mathbf{w}^*\|_2^2 = \sum_{i=1}^l \sum_{j=1}^l \alpha_i^* \alpha_j^* y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle = \sum_{i=1}^l \alpha_i^*,$$

which immediately yields the result.

2.4. When the data are not linearly separable, consider the soft-margin SVM given by

$$\begin{aligned} & \underset{\mathbf{w}, b, \xi}{\text{minimize}} && \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^l \xi_i \\ & \text{subject to} && \xi_i \geq 0, \quad i = 1, \dots, l, \\ & && y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i, \quad i = 1, \dots, l, \end{aligned} \tag{2}$$

where $C > 0$ is a fixed parameter.

- (a) (1 point) Show that (2) is equivalent¹ to

$$\underset{\mathbf{w}, b}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^l \ell(y_i, \mathbf{w}^T \mathbf{x}_i + b), \quad (3)$$

where $\ell(\cdot, \cdot)$ is the hinge loss defined by $\ell(y, z) \triangleq \max\{1 - yz, 0\}$.

- (b) (2 points) Show that the objective function of (3), denoted by $f(\mathbf{w}, b)$, is convex, i.e.,

$$f(\theta \mathbf{w}_1 + (1 - \theta) \mathbf{w}_2, \theta b_1 + (1 - \theta) b_2) \leq \theta f(\mathbf{w}_1, b_1) + (1 - \theta) f(\mathbf{w}_2, b_2)$$

for all $\mathbf{w}_1, \mathbf{w}_2 \in \mathbb{R}^n, b_1, b_2 \in \mathbb{R}$, and $\theta \in [0, 1]$.

Solution:

- (a) Consider the optimal value of ξ for given (\mathbf{w}, b) in (2). The constraints are equivalent to

$$\xi_i \geq \max\{0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b)\} = \ell(y_i, \mathbf{w}^T \mathbf{x}_i + b), \quad \forall i = 1, \dots, l.$$

Hence, we have

$$\frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^l \xi_i \geq \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^l \ell(y_i, \mathbf{w}^T \mathbf{x}_i + b),$$

where the equality holds if and only if $\xi_i = \ell(y_i, \mathbf{w}^T \mathbf{x}_i + b)$ for all $i = 1, \dots, l$. As a consequence, to minimize the objective function in (2), the optimal ξ_i shall be chosen as $\xi_i = \ell(y_i, \mathbf{w}^T \mathbf{x}_i + b)$, then the optimization problem (2) becomes the optimization problem (3).

- (b) *It can be shown that both $\|\mathbf{w}\|_2^2$ and $\mathbf{w}^T \mathbf{x}_i + b$ are convex functions of (\mathbf{w}, b) . Then the conclusion is obvious by noting that the nonnegative weighted sum and pointwise maximum are operations that preserve convexity. See Chapter 3.2 of [1] for more details.*

Since (3) is a convex optimization problem without constraints, it can be solved efficiently.

The function $\|\mathbf{w}\|_2^2$ is convex since

$$\theta \|\mathbf{w}_1\|_2^2 + (1 - \theta) \|\mathbf{w}_2\|_2^2 - \|\theta \mathbf{w}_1 + (1 - \theta) \mathbf{w}_2\|_2^2 = \theta(1 - \theta) \|\mathbf{w}_1 - \mathbf{w}_2\|_2^2 \geq 0.$$

Let $t^+ \triangleq \max\{0, t\}$, then the hinge loss $\ell(y, \mathbf{w}^T \mathbf{x} + b)$ as a function of (\mathbf{w}, b) is convex, since

¹Two optimization problems are called equivalent if from a solution of one, a solution of the other is readily found, and vice versa.

$$\begin{aligned}
& \ell(y, (\theta \mathbf{w}_1 + (1 - \theta) \mathbf{w}_2)^T \mathbf{x} + \theta b_1 + (1 - \theta) b_2) \\
&= \max\{0, 1 - y [(\theta \mathbf{w}_1 + (1 - \theta) \mathbf{w}_2)^T \mathbf{x} + \theta b_1 + (1 - \theta) b_2]\} \\
&= \max\{0, \theta[1 - y(\mathbf{w}_1^T \mathbf{x} + b_1)] + (1 - \theta)[1 - y(\mathbf{w}_2^T \mathbf{x} + b_2)]\} \\
&\leq \max\left\{0, (\theta[1 - y(\mathbf{w}_1^T \mathbf{x} + b_1)])^+ + ((1 - \theta)[1 - y(\mathbf{w}_2^T \mathbf{x} + b_2)])^+\right\} \\
&= (\theta[1 - y(\mathbf{w}_1^T \mathbf{x} + b_1)])^+ + ((1 - \theta)[1 - y(\mathbf{w}_2^T \mathbf{x} + b_2)])^+ \\
&= \theta (1 - y(\mathbf{w}_1^T \mathbf{x} + b_1))^+ + (1 - \theta) (1 - y(\mathbf{w}_2^T \mathbf{x} + b_2))^+ \\
&= \theta \cdot \ell(y, \mathbf{w}_1^T \mathbf{x} + b_1) + (1 - \theta) \cdot \ell(y, \mathbf{w}_2^T \mathbf{x} + b_2),
\end{aligned}$$

where the inequality follows from $t \leq t^+$, and the penultimate equality follows from the fact that $(\theta t)^+ = \theta \cdot t^+$ for $\theta \geq 0$.

Finally, the convexity of the objective function $f(\mathbf{w}, b)$ follows from the fact that $f(\mathbf{w}, b)$ is a nonnegative weighted sum of convex functions $\|\mathbf{w}\|_2^2$ and $\ell(y_i, \mathbf{w}^T \mathbf{x}_i + b), i = 1, \dots, l$. (*Easily verified by definition*)

References

- [1] Stephen Boyd and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.