

Pattern Recognition and Machine Learning: Homework 3

Qingru Hu 2020012996

March 13, 2023

Problem 1

The loss function for LDA (using Fisher's condition) is:

$$l_{\text{LDA}} = \frac{1}{2}(\mathbf{w}^\top \mathbf{x}_i + b - y_i)^2 \quad (1)$$

The loss function for logistic regression is:

$$l_{\text{LG}}(x) = \begin{cases} -\log(\theta(x)) & y_i = 1 \\ -\log(1 - \theta(x)) & y_i = 0 \end{cases}$$

where $\theta(x) = \frac{e^{\mathbf{w}^\top \mathbf{x}_i + b}}{1 + e^{\mathbf{w}^\top \mathbf{x}_i + b}}$. The loss function of soft SVM is:

$$l_{\text{hinge}}(x) = \max(1 - y_i \cdot (\mathbf{w}^\top \mathbf{x}_i + b), 0) \quad (2)$$

For a concrete class $y = 1$, plot curves of the three loss functions in Fig.1.

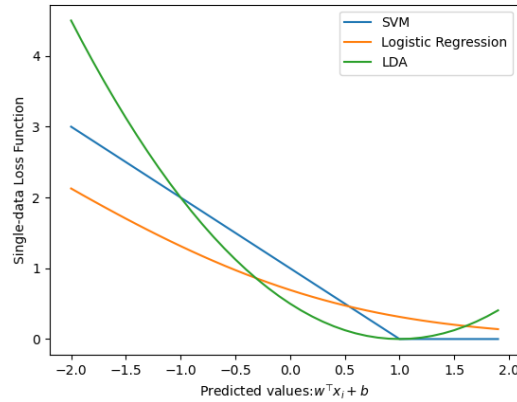


Figure 1: The single-data loss function for the three methods

How does that affect these methods' behavior and make them different?

Problem 2

The **hard-margin** problem is:

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & y_i \cdot (\mathbf{w}^\top \mathbf{x}_i + b) \geq 1, \quad 0 \leq i \leq n. \end{aligned}$$

The Lagrangian function is:

$$\begin{aligned} L(\mathbf{w}, b, \alpha, \xi, \mu) &= \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \alpha_i [1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b)] \\ \alpha_i &\geq 0, i = 1, \dots, n \end{aligned}$$

Take the partial derivatives of Lagrangian w.r.t \mathbf{w}, b and set to zero:

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{w}} &= 0 \quad \mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \\ \frac{\partial L}{\partial b} &= 0 \quad \sum_{i=1}^n \alpha_i y_i = 0 \end{aligned}$$

Pluge the above relations into the original Lagrangian function and we can get:

$$\begin{aligned} L(\mathbf{w}, b, \alpha, \xi, \mu) &= \frac{1}{2} \mathbf{w}^\top \mathbf{w} + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \alpha_i y_i (\mathbf{w}^\top \mathbf{x}_i + b) \\ &= \frac{1}{2} \left(\sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \right)^\top \left(\sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \right) + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \left(\alpha_i y_i \left(\sum_{j=1}^n \alpha_j y_j \mathbf{x}_j \right)^\top \mathbf{x}_i \right) \\ &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j \end{aligned}$$

Therefore, the original optimal problem is equivalent to the dual problem:

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j \\ \text{s.t.} \quad & \sum_{i=1}^n \alpha_i y_i = 0, \alpha_i \geq 0, i = 1, \dots, n \end{aligned}$$

Denote the solution in the first problem as \mathbf{w}^* and solution in the second problem as α^* , and the solutions satisfy:

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

From the complementary slackness of KKT conditions, we have:

$$\alpha_i[1 - y_i(\mathbf{w}^\top x_i + b)] = 0$$

If $\alpha_i > 0$, then $y_i(\mathbf{w}^\top x_i + b) = 1$. Notice that $y_i = \pm 1$, so $\mathbf{w}^\top x_i + b = y_i$.

Problem 3

Search the paper and conduct a literature review for how to add different regularizations to SVM, and more sophisticated forms of kernel function.

Problem 4

(a) Use SVM to classify the train set and predict on the test set

Use Linear SVM to classify the samples on the training set, and the accuracy on the test set is 97%.

(b) Different kernel functions and penalty strength C

Kernel	Penalty Strength	Accuracy(%)
Linear	0.1	96
Poly	0.1	79
RBF	0.1	84
Sigmoid	0.1	86
Linear	1	96
Poly	1	83
RBF	1	97
Sigmoid	1	81
Linear	10	96
Poly	10	90
RBF	10	97
Sigmoid	10	76

Table 1: Different kernel functions and penalty strength for SVM

The best composition is the RBF kernel with a penalty strength of 1 or 10, which both have an accuracy of 97%.

(c) The supporting vectors