

Optimization Problems: Linear Classification

Chih-Jen Lin
National Taiwan University

Last updated: April 23, 2021

Minimizing Training Errors

- Basically a classification method starts with **minimizing the training errors**

$$\min_{\text{model}} \quad (\text{training errors})$$

- That is, all or most training data with labels should be correctly classified by our model
- A model can be a decision tree, a neural network, or other types

Minimizing Training Errors (Cont'd)

- For simplicity, let's consider the model to be a vector \mathbf{w}
- That is, the decision function is

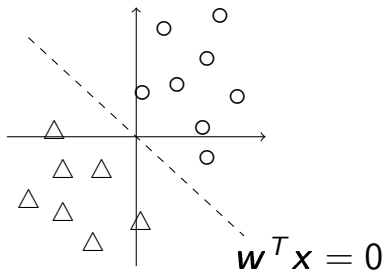
$$\text{sgn}(\mathbf{w}^T \mathbf{x})$$

- For any data, \mathbf{x} , the predicted label is

$$\begin{cases} 1 & \text{if } \mathbf{w}^T \mathbf{x} \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

Minimizing Training Errors (Cont'd)

- The two-dimensional situation



- This seems to be quite restricted, but practically x is in a much **higher dimensional space**

Minimizing Training Errors (Cont'd)

- To characterize the training error, we need a **loss function** $\xi(\mathbf{w}; y, \mathbf{x})$ for each instance (y, \mathbf{x}) , where

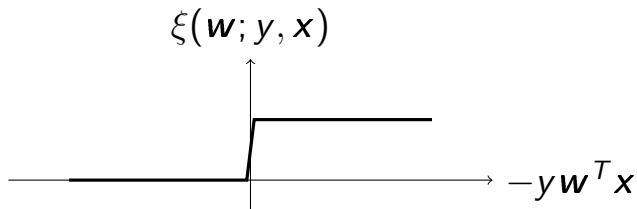
$y = \pm 1$ is the label and \mathbf{x} is the feature vector

- Ideally we should use 0–1 training loss:

$$\xi(\mathbf{w}; y, \mathbf{x}) = \begin{cases} 1 & \text{if } y\mathbf{w}^T\mathbf{x} < 0, \\ 0 & \text{otherwise} \end{cases}$$

Minimizing Training Errors (Cont'd)

- However, this function is **discontinuous**. The optimization problem becomes difficult



- We need **continuous approximations**

Common Loss Functions

- Hinge loss (l1 loss)

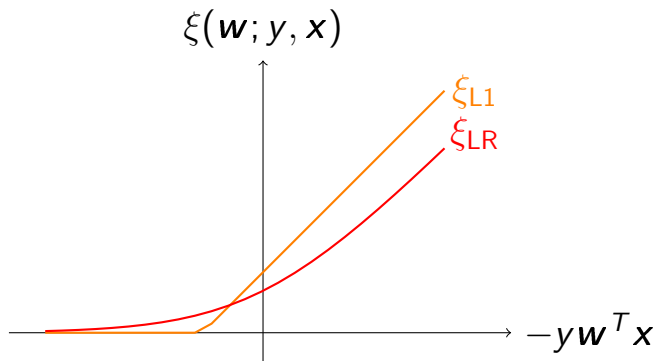
$$\xi_{L1}(\mathbf{w}; y, \mathbf{x}) \equiv \max(0, 1 - y\mathbf{w}^T \mathbf{x}) \quad (1)$$

- Logistic loss

$$\xi_{LR}(\mathbf{w}; y, \mathbf{x}) \equiv \log(1 + e^{-y\mathbf{w}^T \mathbf{x}}) \quad (2)$$

- Support vector machines (SVM): Eq. (1). Logistic regression (LR): (2)
- SVM and LR are two very fundamental classification methods

Common Loss Functions (Cont'd)



- Logistic regression is very related to SVM
- Their performance is usually similar

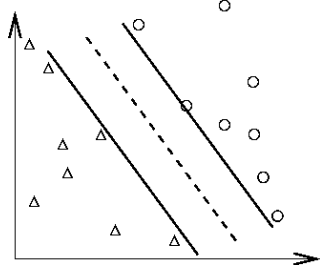
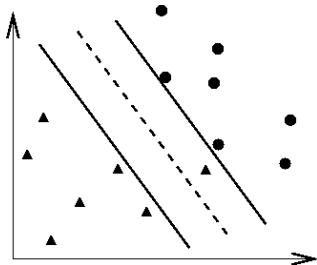
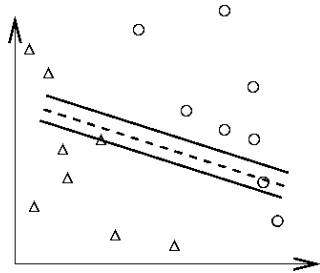
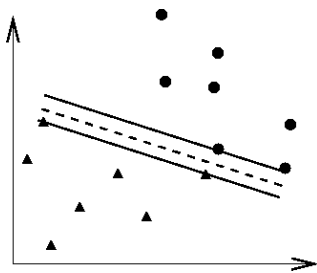
Common Loss Functions (Cont'd)

- However, minimizing training losses may not give a good model for future prediction
- **Overfitting occurs**

Overfitting

- See the illustration in the next slide
- For classification, you can easily achieve 100% training accuracy
- This is useless
- When training a data set, we should
Avoid **underfitting**: small training error
Avoid **overfitting**: small testing error

● and ▲: training; ○ and △: testing



Regularization

- To minimize the training error we manipulate the \mathbf{w} vector so that it fits the data
- To avoid overfitting we need a way to make \mathbf{w} 's values **less extreme**.
- One idea is to make **\mathbf{w} values closer to zero**
- We can add, for example,

$$\frac{\mathbf{w}^T \mathbf{w}}{2} \quad \text{or} \quad \|\mathbf{w}\|_1$$

to the function that is minimized

General Form of Linear Classification

- Training data $\{y_i, \mathbf{x}_i\}$, $\mathbf{x}_i \in R^n, i = 1, \dots, l, y_i = \pm 1$
- l : # of data, n : # of features

$$\min_{\mathbf{w}} f(\mathbf{w}), \quad f(\mathbf{w}) \equiv \frac{\mathbf{w}^T \mathbf{w}}{2} + C \sum_{i=1}^l \xi(\mathbf{w}; y_i, \mathbf{x}_i)$$

- $\mathbf{w}^T \mathbf{w}/2$: regularization term
- $\xi(\mathbf{w}; y, \mathbf{x})$: loss function
- C : regularization parameter (chosen by users)