

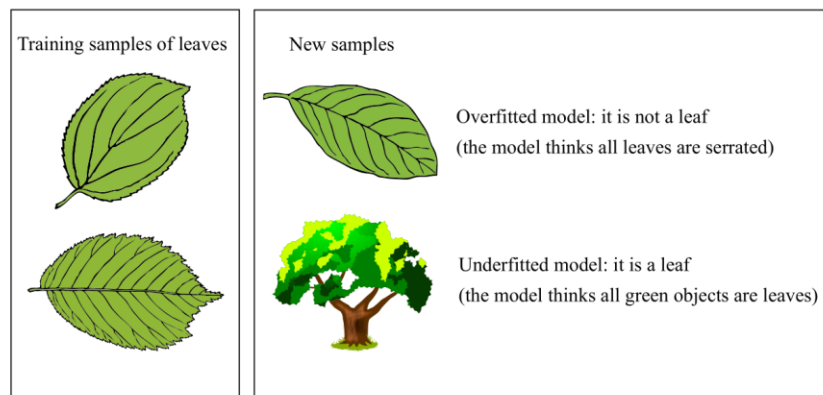
Lesson3: Linear Models

Instructor: Xiang Deng

. 2025.10.10 .

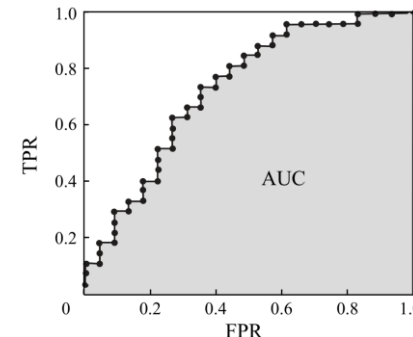
review

□ Empirical Error and Overfitting: overfitting, underfitting



The confusion matrix of binary classification

Ground-truth class	Predicted class	
	Positive	Negative
Positive	TP	FN
Negative	FP	TN



$$\text{Precision } P = \frac{TP}{TP + FP}$$

$$\text{Recall } R = \frac{TP}{TP + FN}$$

$$F1 = \frac{2 \times P \times R}{P + R}$$

- Evaluation methods: generalization error, Hold-out, Cross-Validation, Bootstrapping
- Performance Measure: Mean Squared Error, Error rate, accuracy, Precision, recall, f1, P-R curve, break-even points, ROC curve and AUC, cost-sensitive error rate
- Performance Comparisons: Hypothesis testing
- Bias and Variance: the generalization error can be decomposed into the sum of bias, variance, and noise.

$$E(f; D) = \text{bias}^2(\mathbf{x}) + \text{var}(\mathbf{x}) + \varepsilon^2$$

Outline

- Linear Regression
 - Least Squares Method

- Binary Classification Problem
 - Logistic Regression
 - Linear Discriminant Analysis

- Multiclass Classification Problem
 - One vs. One (OvO)
 - One vs. Rest (OvR)
 - Many vs. Many (MvM)

- Class Imbalance Problem

The Basic Form

- The Basic form of the linear model

$$f(\mathbf{x}) = w_1x_1 + w_2x_2 + \dots + w_dx_d + b$$

$\mathbf{x} = (x_1; x_2; \dots; x_d)$ is a sample where x takes the value x_i on the i -th variable.

- The vector form

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$$

where $\mathbf{w} = (w_1; w_2; \dots; w_d)$

Advantages of the Linear Model

- ❑ Simple form and ease of modeling
- ❑ Comprehensibility
- ❑ Nonlinear models can be derived from linear models
 - Introducing multi-layer structures or high-dimensional mapping.
- ❑ An example
 - Determine the ripeness of a watermelon by considering its **color**, **root** and **sound** information.
$$f_{\text{ripe}}(\mathbf{x}) = 0.2 \cdot x_{\text{color}} + 0.5 \cdot x_{\text{root}} + 0.3 \cdot x_{\text{sound}} + 1.$$
 - From the coefficients, we know that root is the most important variable, and sound is more important than color.

Linear Regression

- Given a data set $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$
where $\mathbf{x}_i = (x_{i1}; x_{i2}; \dots; x_{id})$ $y_i \in \mathbb{R}$
- The aim of linear regression
 - Learn a linear model that can accurately predict the real-valued output labels
- For discrete variables
 - An ordinal relationship exists between values
 - Convert the variables into real-valued variables
 - No ordinal relationship exists
 - Convert the discrete variable with k possible values into a k-dimensional vector

Linear Regression

- Linear regression aims to learn the function:

$$f(x) = wx + b, \text{ such that } f(x_i) \simeq y_i, i = 1, \dots, m$$

- Parameter/model estimation: least square method

$$\begin{aligned}(w^*, b^*) &= \arg \min_{(w, b)} \sum_{i=1}^m (f(x_i) - y_i)^2 \\ &= \arg \min_{(w, b)} \sum_{i=1}^m (y_i - wx_i - b)^2\end{aligned}$$

Linear Regression - Least Square Method

- Minimize mean-square error (MSE)

$$E_{(w,b)} = \sum_{i=1}^m (y_i - wx_i - b)^2$$

- Calculate the derivatives of $E_{(w,b)}$ with respect to w and b respectively:

$$\frac{\partial E_{(w,b)}}{\partial w} = 2 \left(w \sum_{i=1}^m x_i^2 - \sum_{i=1}^m (y_i - b) x_i \right)$$

$$\frac{\partial E_{(w,b)}}{\partial b} = 2 \left(mb - \sum_{i=1}^m (y_i - wx_i) \right)$$

$$\begin{aligned} \frac{\partial E_{(w,b)}}{\partial w} &= \frac{\partial}{\partial w} \left[\sum_{i=1}^m (y_i - wx_i - b)^2 \right] \\ &= \sum_{i=1}^m \frac{\partial}{\partial w} [(y_i - wx_i - b)^2] \\ &= \sum_{i=1}^m [2 \cdot (y_i - wx_i - b) \cdot (-x_i)] \\ &= \sum_{i=1}^m [2 \cdot (wx_i^2 - y_i x_i + bx_i)] \\ &= 2 \cdot \left(w \sum_{i=1}^m x_i^2 - \sum_{i=1}^m y_i x_i + b \sum_{i=1}^m x_i \right) \\ &= 2 \left(w \sum_{i=1}^m x_i^2 - \sum_{i=1}^m (y_i - b) x_i \right) \end{aligned}$$

$$\begin{aligned} \frac{\partial E_{(w,b)}}{\partial b} &= \frac{\partial}{\partial b} \left[\sum_{i=1}^m (y_i - wx_i - b)^2 \right] \\ &= \sum_{i=1}^m \frac{\partial}{\partial b} [(y_i - wx_i - b)^2] \\ &= \sum_{i=1}^m [2 \cdot (y_i - wx_i - b) \cdot (-1)] \\ &= \sum_{i=1}^m [2 \cdot (b - y_i + wx_i)] \\ &= 2 \cdot \left[\sum_{i=1}^m b - \sum_{i=1}^m y_i + \sum_{i=1}^m wx_i \right] \\ &= 2 \left(mb - \sum_{i=1}^m (y_i - wx_i) \right) \end{aligned}$$

Linear Regression - Least Square Method

- We have the closed-form solutions (derivatives w.r.t. w and b equal to 0) :

$$w = \frac{\sum_{i=1}^m y_i (x_i - \bar{x})}{\sum_{i=1}^m x_i^2 - \frac{1}{m} \left(\sum_{i=1}^m x_i \right)^2}$$

$$b = \frac{1}{m} \sum_{i=1}^m (y_i - wx_i)$$

where

$$\bar{x} = \frac{1}{m} \sum_{i=1}^m x_i$$

convex function

$$E_{(w,b)} = \sum_{i=1}^m (y_i - wx_i - b)^2$$

How are they obtained?

Linear Regression - Least Square Method

$$0 = w \sum_{i=1}^m x_i^2 - \sum_{i=1}^m (y_i - b)x_i \quad w \sum_{i=1}^m x_i^2 = \sum_{i=1}^m y_i x_i - \sum_{i=1}^m b x_i$$

$$0 = 2 \left(mb - \sum_{i=1}^m (y_i - w x_i) \right) \quad b = \frac{1}{m} \sum_{i=1}^m (y_i - w x_i)$$

$$w \sum_{i=1}^m x_i^2 = \sum_{i=1}^m y_i x_i - \sum_{i=1}^m (\bar{y} - w \bar{x}) x_i$$

$$\frac{1}{m} \sum_{i=1}^m y_i = \bar{y}, \quad \frac{1}{m} \sum_{i=1}^m x_i = \bar{x}$$

$$w \sum_{i=1}^m x_i^2 = \sum_{i=1}^m y_i x_i - \bar{y} \sum_{i=1}^m x_i + w \bar{x} \sum_{i=1}^m x_i$$

$$w \left(\sum_{i=1}^m x_i^2 - \bar{x} \sum_{i=1}^m x_i \right) = \sum_{i=1}^m y_i x_i - \bar{y} \sum_{i=1}^m x_i$$

$$\bar{y} \sum_{i=1}^m x_i = \frac{1}{m} \sum_{i=1}^m y_i \sum_{i=1}^m x_i = \bar{x} \sum_{i=1}^m y_i$$

$$w = \frac{\sum_{i=1}^m y_i x_i - \bar{y} \sum_{i=1}^m x_i}{\sum_{i=1}^m x_i^2 - \bar{x} \sum_{i=1}^m x_i}$$

$$\bar{x} \sum_{i=1}^m x_i = \frac{1}{m} \sum_{i=1}^m x_i \sum_{i=1}^m x_i = \frac{1}{m} (\sum_{i=1}^m x_i)^2$$

$$w = \frac{\sum_{i=1}^m y_i (x_i - \bar{x})}{\sum_{i=1}^m x_i^2 - \frac{1}{m} (\sum_{i=1}^m x_i)^2}$$

Multivariate Linear Regression

- Given a data set

$$D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$$

$$\mathbf{x}_i = (x_{i1}; x_{i2}; \dots; x_{id}) \quad y_i \in \mathbb{R}$$

- The objective function of multivariate linear regression

$$f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i + b \quad \text{such that} \quad f(\mathbf{x}_i) \simeq y_i$$

Multivariate Linear Regression

□ Rewrite

$$\hat{\mathbf{w}} = (\mathbf{w}; b) = (w_1; \dots; w_d; b) \in \mathbb{R}^{(d+1) \times 1}, \hat{\mathbf{x}}_i = (x_{i1}; \dots; x_{id}; 1) \in \mathbb{R}^{(d+1) \times 1},$$

□ Least square method

$$\begin{aligned}\hat{\mathbf{w}}^* &= \arg \min_{\hat{\mathbf{w}}} \sum_{i=1}^m \left(y_i - \hat{\mathbf{w}}^T \hat{\mathbf{x}}_i \right)^2 \\ &= \arg \min_{\hat{\mathbf{w}}} \sum_{i=1}^m \left(y_i - \hat{\mathbf{x}}_i^T \hat{\mathbf{w}} \right)^2\end{aligned}$$

$$\hat{\mathbf{w}}^* = \arg \min_{\hat{\mathbf{w}}} \begin{bmatrix} y_1 - \hat{\mathbf{x}}_1^T \hat{\mathbf{w}} & \cdots & y_m - \hat{\mathbf{x}}_m^T \hat{\mathbf{w}} \end{bmatrix} \begin{bmatrix} y_1 - \hat{\mathbf{x}}_1^T \hat{\mathbf{w}} \\ \vdots \\ y_m - \hat{\mathbf{x}}_m^T \hat{\mathbf{w}} \end{bmatrix}$$

Multivariate Linear Regression

$$\begin{aligned} \begin{bmatrix} y_1 - \hat{\mathbf{x}}_1^T \hat{\mathbf{w}} \\ \vdots \\ y_m - \hat{\mathbf{x}}_m^T \hat{\mathbf{w}} \end{bmatrix} &= \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} - \begin{bmatrix} \hat{\mathbf{x}}_1^T \hat{\mathbf{w}} \\ \vdots \\ \hat{\mathbf{x}}_m^T \hat{\mathbf{w}} \end{bmatrix} \\ &= \mathbf{y} - \begin{bmatrix} \hat{\mathbf{x}}_1^T \\ \vdots \\ \hat{\mathbf{x}}_m^T \end{bmatrix} \cdot \hat{\mathbf{w}} \\ &= \mathbf{y} - \mathbf{X}\hat{\mathbf{w}} \end{aligned}$$

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1d} & 1 \\ x_{21} & x_{22} & \cdots & x_{2d} & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_{m1} & x_{m2} & \cdots & x_{md} & 1 \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1^T & 1 \\ \mathbf{x}_2^T & 1 \\ \vdots & \vdots \\ \mathbf{x}_m^T & 1 \end{pmatrix}$$
$$\mathbf{y} = (y_1; y_2; \dots; y_m)$$

$$\hat{\mathbf{w}}^* = \arg \min_{\hat{\mathbf{w}}} (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})$$

Let $E_{\hat{\mathbf{w}}} = (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})$ and find the derivative with respect to $\hat{\mathbf{w}}$

Multivariate Linear Regression - Least Square Method

□ The derivative with respect to $\hat{\mathbf{w}}$

$$E_{\hat{\mathbf{w}}} = (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\hat{\mathbf{w}})$$

$$E_{\hat{\mathbf{w}}} = \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\hat{\mathbf{w}} - \hat{\mathbf{w}}^T \mathbf{X}^T \mathbf{y} + \hat{\mathbf{w}}^T \mathbf{X}^T \mathbf{X} \hat{\mathbf{w}}$$

$$\begin{aligned} \frac{\partial E_{\hat{\mathbf{w}}}}{\partial \hat{\mathbf{w}}} &= 0 - \mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{y} + (\mathbf{X}^T \mathbf{X} + \mathbf{X}^T \mathbf{X}) \hat{\mathbf{w}} \\ &= 2\mathbf{X}^T (\mathbf{X}\hat{\mathbf{w}} - \mathbf{y}) \end{aligned}$$

$$\frac{\partial E_{\hat{\mathbf{w}}}}{\partial \hat{\mathbf{w}}} = 2\mathbf{X}^T (\mathbf{X}\hat{\mathbf{w}} - \mathbf{y})$$

Recall Matrix transpose:

$$(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T$$

$$(\mathbf{A} + \mathbf{B})^T = \mathbf{A}^T + \mathbf{B}^T$$

Recall Matrix Differentiation:

$$\frac{\partial \mathbf{a}^T \mathbf{x}}{\partial \mathbf{x}} = \frac{\partial \mathbf{x}^T \mathbf{a}}{\partial \mathbf{x}} = \mathbf{a}, \quad \frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{A} + \mathbf{A}^T) \mathbf{x}$$

The closed-form solution of $\hat{\mathbf{w}}$ can be obtained by making the equation equal to 0.

Multivariate Linear Regression - Discussion of full-rank

- $\mathbf{X}^T \mathbf{X}$ is a full-rank matrix or a positive definite matrix, then by making the above equation equal to 0:

$$\hat{\mathbf{w}}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

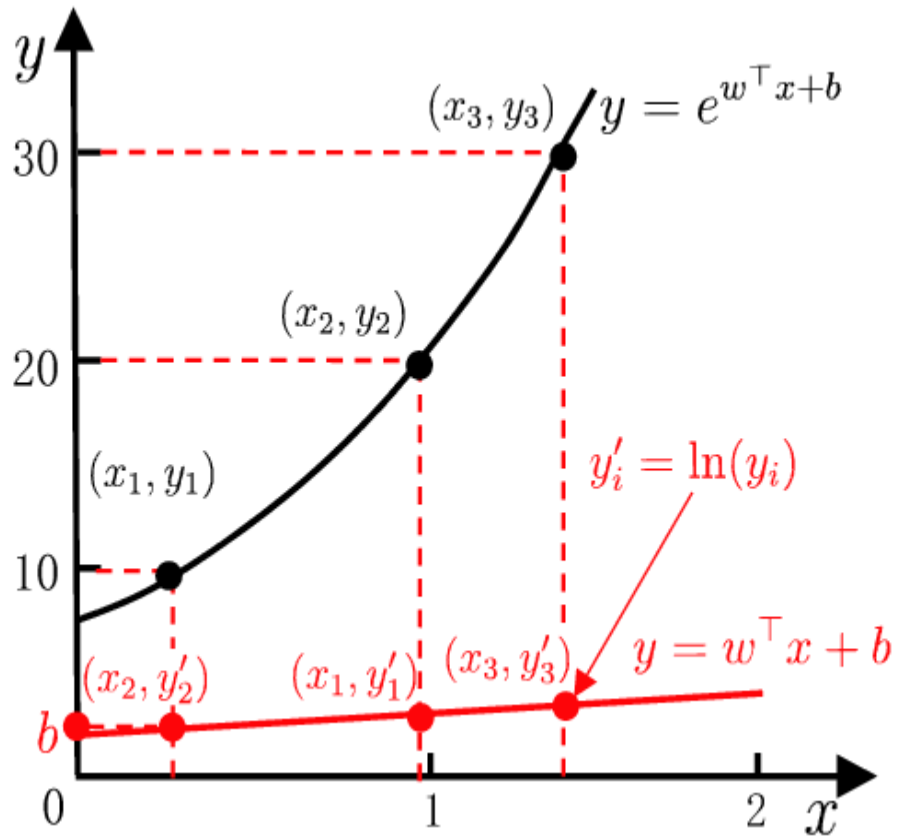
where $(\mathbf{X}^T \mathbf{X})^{-1}$ is the inverse of $\mathbf{X}^T \mathbf{X}$, the learned multivariate linear regression model is

$$f(\hat{\mathbf{x}}_i) = \hat{\mathbf{x}}_i^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- $\mathbf{X}^T \mathbf{X}$ is often not full-rank
 - introducing regularization.

Log-Linear Regression

- The logarithm of the output label can be used for approximation



$$\ln y = w^T x + b$$



$$y = w^T x + b$$

Linear Regression - Generalized Linear Model

- The general form

$$y = g^{-1}(\boldsymbol{w}^T \boldsymbol{x} + b)$$

- Where the function $g(\cdot)$ is the link function.
 - a monotonic differentiable function
- Log-linear regression is a special case of generalized linear models when $g(\cdot) = \ln(\cdot)$

Binary Classification

- The predictions and the output labels

$$z = \mathbf{w}^T \mathbf{x} + b \quad y \in \{0, 1\}$$

- The real-valued predictions of the linear regression model need to be converted into 0/1.

- Ideally, the unit-step function is desired

$$y = \begin{cases} 0, & z < 0; \\ 0.5, & z = 0; \\ 1, & z > 0, \end{cases}$$

- which predicts positive for z greater than 0, negative for z smaller than 0, and an arbitrary output when z equals to 0.

Binary Classification

❑ Disadvantages of unit-step function

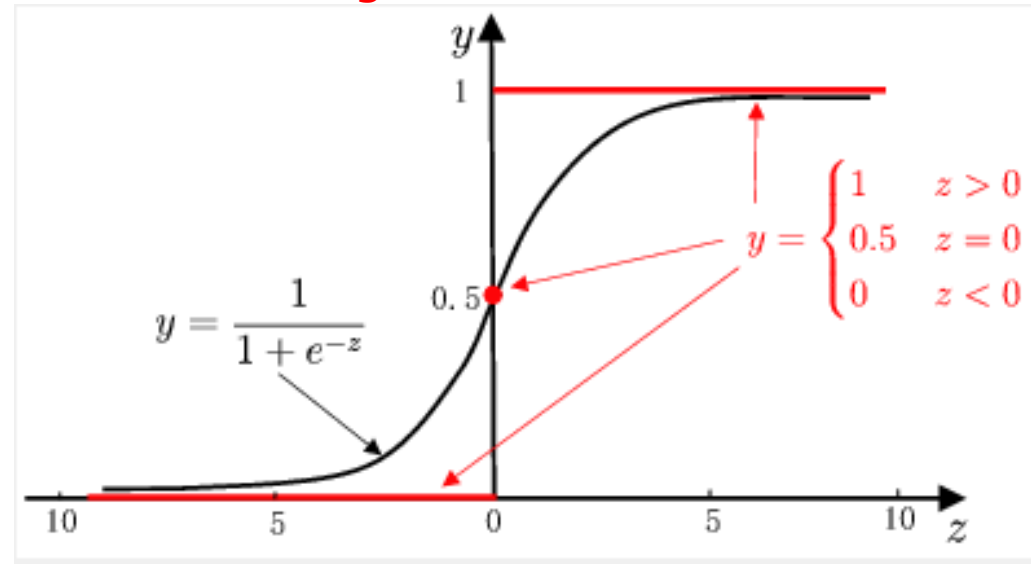
- not continuous

❑ Logistic function: a surrogate function to approximate the unit-step function

- monotonic differentiable

Comparison between unit-step function and logistic function

$$y = \frac{1}{1 + e^{-z}}$$



Logistic regression

- Apply logistic function

$$y = \frac{1}{1 + e^{-z}} \quad \text{transform into} \quad y = \frac{1}{1 + e^{-(w^T x + b)}}$$

- Log odds

- the logarithm of the relative likelihood of a sample being a positive sample

$$\ln \frac{y}{1 - y}$$

- Logistic regression has several nice properties
 - without requiring any prior assumptions on the data distribution
 - it predicts labels together with associated probabilities
 - it is solvable with numerical optimization methods.

Logistic regression - maximum likelihood

- Log odds can be rewritten as

$$\ln \frac{p(y = 1 \mid \mathbf{x})}{p(y = 0 \mid \mathbf{x})} = \mathbf{w}^T \mathbf{x} + b$$

$$p(y = 1 \mid \mathbf{x}) = \frac{e^{\mathbf{w}^T \mathbf{x} + b}}{1 + e^{\mathbf{w}^T \mathbf{x} + b}}$$

$$p(y = 0 \mid \mathbf{x}) = \frac{1}{1 + e^{\mathbf{w}^T \mathbf{x} + b}}$$

Logistic regression - maximum likelihood

□ Maximum likelihood

- Given a data set

$$\{(\mathbf{x}_i, y_i)\}_{i=1}^m$$

- Maximizing the probability of each sample being predicted as the ground-truth label
 - the log-likelihood to be maximized is

$$\ell(\mathbf{w}, b) = \sum_{i=1}^m \ln p(y_i \mid \mathbf{x}_i; \mathbf{w}_i, b)$$

Logistic regression - maximum likelihood

□ Transform into minimize negative log-likelihood

- Let $\beta = (\mathbf{w}; b)$, $\hat{\mathbf{x}} = (\mathbf{x}; 1)$, $\mathbf{w}^T \mathbf{x} + b$ can be rewritten as $\beta^T \hat{\mathbf{x}}$

- Let

$$p_1(\hat{\mathbf{x}}_i; \beta) = p(y = 1 \mid \hat{\mathbf{x}}_i; \beta)$$

$$p_0(\hat{\mathbf{x}}_i; \beta) = p(y = 0 \mid \hat{\mathbf{x}}_i; \beta) = 1 - p_1(\hat{\mathbf{x}}_i; \beta)$$

the likelihood term in can be rewritten as

$$p(y_i \mid \mathbf{x}_i; \mathbf{w}_i, b) = y_i p_1(\hat{\mathbf{x}}_i; \beta) + (1 - y_i) p_0(\hat{\mathbf{x}}_i; \beta)$$

- the log-likelihood can be written:

$$\ell(\beta) = \sum_{i=1}^m \ln(y_i p_1(\hat{\mathbf{x}}_i; \beta) + (1 - y_i) p_0(\hat{\mathbf{x}}_i; \beta))$$

Logistic regression - maximum likelihood

- since $p_1(\hat{\mathbf{x}}_i; \boldsymbol{\beta}) = \frac{e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}}{1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}}, p_0(\hat{\mathbf{x}}_i; \boldsymbol{\beta}) = \frac{1}{1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}},$

$$\begin{aligned}\ell(\boldsymbol{\beta}) &= \sum_{i=1}^m \ln \left(\frac{y_i e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i} + 1 - y_i}{1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}} \right) \\ &= \sum_{i=1}^m \left(\ln(y_i e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i} + 1 - y_i) - \ln(1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}) \right)\end{aligned}$$

- y_i can be 0 or 1:

$$\ell(\boldsymbol{\beta}) = \begin{cases} \sum_{i=1}^m (-\ln(1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i})), & y_i = 0 \\ \sum_{i=1}^m (\boldsymbol{\beta}^T \hat{\mathbf{x}}_i - \ln(1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i})), & y_i = 1 \end{cases}$$

$$\ell(\boldsymbol{\beta}) = \sum_{i=1}^m \left(-y_i \boldsymbol{\beta}^T \hat{\mathbf{x}}_i + \ln(1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}) \right)$$

Logistic regression

□ We have

$$\boldsymbol{\beta}^* = \arg \min_{\boldsymbol{\beta}} \ell(\boldsymbol{\beta})$$

□ Taking Newton's method as an example, the update rule at the (t+1)th iteration is

$$\boldsymbol{\beta}^{t+1} = \boldsymbol{\beta}^t - \left(\frac{\partial^2 \ell(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T} \right)^{-1} \frac{\partial \ell(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}}$$

Need to calculate the first- and second-order derivatives with respect to $\boldsymbol{\beta}$

higher order differentiable convex function, gradient
descent method / Newton's method [Boyd and Vandenberghe, 2004]

Logistic regression

the first-order derivatives with respect to β are

$$\begin{aligned}\frac{\partial \ell(\beta)}{\partial \beta} &= \frac{\partial \sum_{i=1}^m \left(-y_i \beta^T \hat{\mathbf{x}}_i + \ln \left(1 + e^{\beta^T \hat{\mathbf{x}}_i} \right) \right)}{\partial \beta} \\&= \sum_{i=1}^m \left(\frac{\partial \left(-y_i \beta^T \hat{\mathbf{x}}_i \right)}{\partial \beta} + \frac{\partial \ln \left(1 + e^{\beta^T \hat{\mathbf{x}}_i} \right)}{\partial \beta} \right) \\&= \sum_{i=1}^m \left(-y_i \hat{\mathbf{x}}_i + \frac{1}{1 + e^{\beta^T \hat{\mathbf{x}}_i}} \cdot \hat{\mathbf{x}}_i e^{\beta^T \hat{\mathbf{x}}_i} \right) \\&= - \sum_{i=1}^m \hat{\mathbf{x}}_i \left(y_i - \frac{e^{\beta^T \hat{\mathbf{x}}_i}}{1 + e^{\beta^T \hat{\mathbf{x}}_i}} \right) \\&= - \sum_{i=1}^m \hat{\mathbf{x}}_i (y_i - p_1(\hat{\mathbf{x}}_i; \beta))\end{aligned}$$

Logistic regression

the second-order derivatives with respect β are

$$\begin{aligned}\frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta^T} &= - \frac{\partial \sum_{i=1}^m \hat{\mathbf{x}}_i \left(y_i - \frac{e^{\beta^T \hat{\mathbf{x}}_i}}{1 + e^{\beta^T \hat{\mathbf{x}}_i}} \right)}{\partial \beta^T} \\&= - \sum_{i=1}^m \hat{\mathbf{x}}_i \frac{\partial \left(y_i - \frac{e^{\beta^T \hat{\mathbf{x}}_i}}{1 + e^{\beta^T \hat{\mathbf{x}}_i}} \right)}{\partial \beta^T} \\&= - \sum_{i=1}^m \hat{\mathbf{x}}_i \left(\frac{\partial y_i}{\partial \beta^T} - \frac{\partial \left(\frac{e^{\beta^T \hat{\mathbf{x}}_i}}{1 + e^{\beta^T \hat{\mathbf{x}}_i}} \right)}{\partial \beta^T} \right) \\&= \sum_{i=1}^m \hat{\mathbf{x}}_i \cdot \frac{\partial \left(\frac{e^{\beta^T \hat{\mathbf{x}}_i}}{1 + e^{\beta^T \hat{\mathbf{x}}_i}} \right)}{\partial \beta^T}\end{aligned}$$

Logistic regression

according to matrix differential formulas $\frac{\partial \mathbf{a}^T \mathbf{x}}{\partial \mathbf{x}^T} = \frac{\partial \mathbf{x}^T \mathbf{a}}{\partial \mathbf{x}^T} = \mathbf{a}^T$,

$$\begin{aligned}
 \frac{\partial \left(\frac{e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}}{1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}} \right)}{\partial \boldsymbol{\beta}^T} &= \frac{\text{numerator} \quad \text{denominator} \quad \text{numerator} \quad \text{denominator}}{\frac{\partial e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}}{\partial \boldsymbol{\beta}^T} \cdot (1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}) - e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i} \cdot \frac{\partial (1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i})}{\partial \boldsymbol{\beta}^T}} \\
 &= \frac{\text{natural constant} \quad \hat{\mathbf{x}}_i^T e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i} \cdot (1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}) - e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i} \cdot \hat{\mathbf{x}}_i^T e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}}{(1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i})^2} \\
 &= \hat{\mathbf{x}}_i^T e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i} \cdot \frac{(1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}) - e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}}{(1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i})^2} \\
 &= \hat{\mathbf{x}}_i^T e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i} \cdot \frac{1}{(1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i})^2} \\
 &= \hat{\mathbf{x}}_i^T \cdot \frac{e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}}{1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}} \cdot \frac{1}{1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}}
 \end{aligned}$$

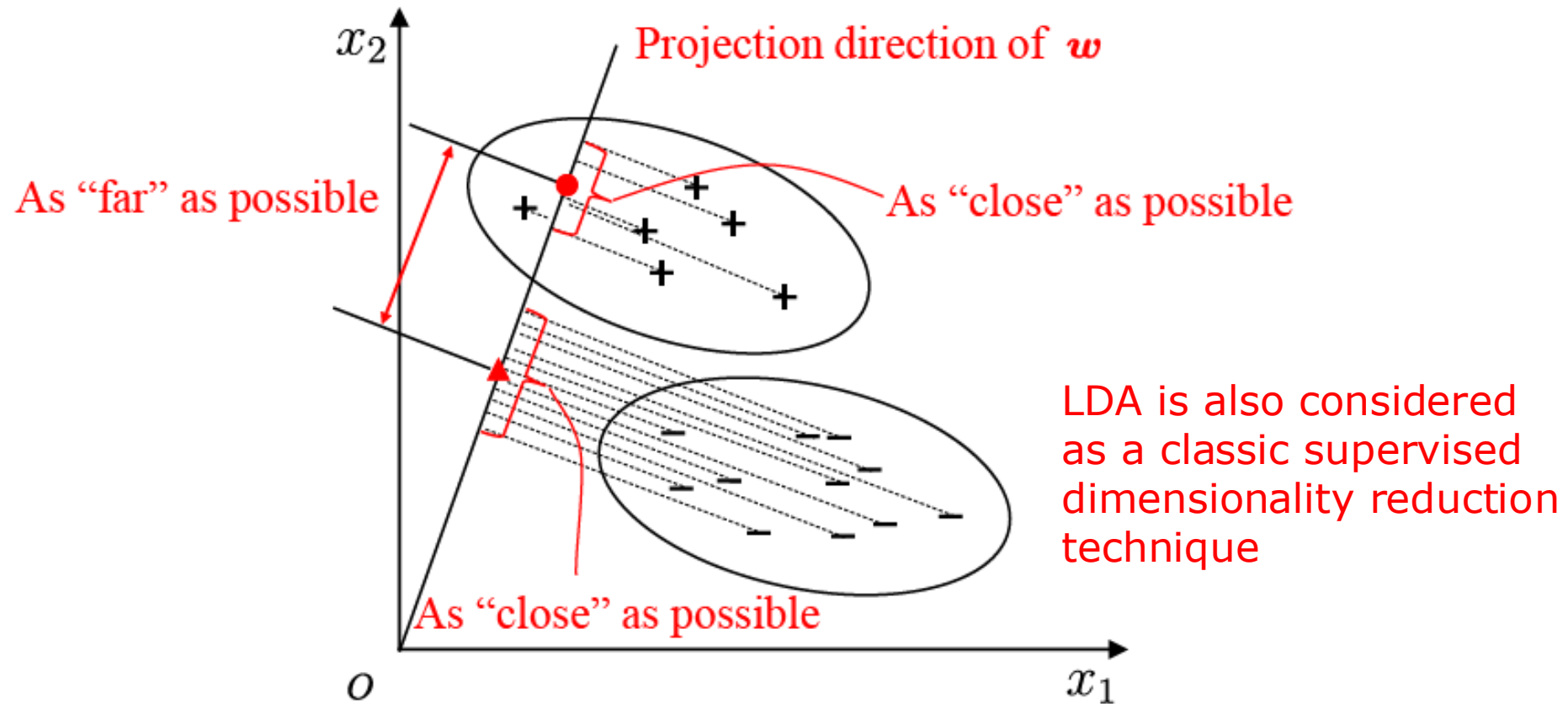
from the previous slide:

$$\frac{\partial^2 \ell(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T} = \sum_{i=1}^m \hat{\mathbf{x}}_i \cdot \frac{\partial \left(\frac{e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}}{1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}} \right)}{\partial \boldsymbol{\beta}^T}$$

$$\begin{aligned}
 \frac{\partial^2 \ell(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T} &= \sum_{i=1}^m \hat{\mathbf{x}}_i \cdot \hat{\mathbf{x}}_i^T \cdot \frac{e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}}{1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}} \cdot \frac{1}{1 + e^{\boldsymbol{\beta}^T \hat{\mathbf{x}}_i}} \\
 &= \sum_{i=1}^m \hat{\mathbf{x}}_i \hat{\mathbf{x}}_i^T p_1(\hat{\mathbf{x}}_i; \boldsymbol{\beta}) (1 - p_1(\hat{\mathbf{x}}_i; \boldsymbol{\beta}))
 \end{aligned}$$

Binary Classification - Linear Discriminant Analysis

□ Linear Discriminant Analysis [Fisher, 1936]



Binary Classification - Linear Discriminant Analysis

- The idea of LDA: project the same class samples onto a line, while samples of different classes are far away from each other.
 - To make the projection points of similar samples as close as possible, we can make the covariance of the projection points of similar samples as small as possible
 - To make the projection points of examples from different classes as far away as possible, we can make the distance between the class centers as large as possible

- Some variables
 - the sample set of the i -th class X_i
 - the mean vector of the i -th class μ_i
 - the covariance matrix of the i -th class Σ_i
 - the centers of those two classes samples $w^T \mu_0$ and $w^T \mu_1$
 - the covariances of the two classes samples $w^T \Sigma_0 w$ and $w^T \Sigma_1 w$

Binary Classification - Linear Discriminant Analysis

- We have the objective to be maximized

$$\begin{aligned} J &= \frac{\|w^T \mu_0 - w^T \mu_1\|_2^2}{w^T \Sigma_0 w + w^T \Sigma_1 w} \\ &= \frac{w^T (\mu_0 - \mu_1) (\mu_0 - \mu_1)^T w}{w^T (\Sigma_0 + \Sigma_1) w} \end{aligned}$$

To make the projection points of similar samples as close as possible, we can make the covariance of the projection points of similar samples as small as possible.

To make the projection points of examples from different classes as far away as possible, we can make the distance between the class centers as large as possible

- The within-class scatter matrix

$$\begin{aligned} S_w &= \Sigma_0 + \Sigma_1 \\ &= \sum_{x \in X_0} (x - \mu_0) (x - \mu_0)^T + \sum_{x \in X_1} (x - \mu_1) (x - \mu_1)^T \end{aligned}$$

- The between-class scatter matrix

$$S_b = (\mu_0 - \mu_1) (\mu_0 - \mu_1)^T$$

Binary Classification - Linear Discriminant Analysis

- Generalized Rayleigh quotient

$$J = \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}}$$

- Let $\mathbf{w}^T \mathbf{S}_w \mathbf{w} = 1$, maximizing generalized Rayleigh quotient is equivalent to

$$\begin{aligned} \min_{\mathbf{w}} \quad & -\mathbf{w}^T \mathbf{S}_b \mathbf{w} \\ \text{s.t.} \quad & \mathbf{w}^T \mathbf{S}_w \mathbf{w} = 1 \end{aligned}$$

- Using the method of Lagrange multipliers

$$L(\mathbf{w}, \lambda) = -\mathbf{w}^T \mathbf{S}_b \mathbf{w} + \lambda(\mathbf{w}^T \mathbf{S}_w \mathbf{w} - 1)$$

Binary Classification - Linear Discriminant Analysis

the derivatives are:

$$\begin{aligned}\frac{\partial L(\mathbf{w}, \lambda)}{\partial \mathbf{w}} &= -\frac{\partial(\mathbf{w}^T \mathbf{S}_b \mathbf{w})}{\partial \mathbf{w}} + \lambda \frac{\partial(\mathbf{w}^T \mathbf{S}_w \mathbf{w} - 1)}{\partial \mathbf{w}} \\ &= -(\mathbf{S}_b + \mathbf{S}_b^T) \mathbf{w} + \lambda(\mathbf{S}_w + \mathbf{S}_w^T) \mathbf{w}\end{aligned}$$

since $\mathbf{S}_b = \mathbf{S}_b^T, \mathbf{S}_w = \mathbf{S}_w^T$,

$$\frac{\partial L(\mathbf{w}, \lambda)}{\partial \mathbf{w}} = -2\mathbf{S}_b \mathbf{w} + 2\lambda \mathbf{S}_w \mathbf{w}$$

Let it be 0: $-2\mathbf{S}_b \mathbf{w} + 2\lambda \mathbf{S}_w \mathbf{w} = 0$

$$\mathbf{S}_b \mathbf{w} = \lambda \mathbf{S}_w \mathbf{w}$$

$$(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^T \mathbf{w} = \lambda \mathbf{S}_w \mathbf{w}$$

Binary Classification - Linear Discriminant Analysis

$$(\mu_0 - \mu_1)(\mu_0 - \mu_1)^T \mathbf{w} = \lambda \mathbf{S}_w \mathbf{w}$$

Let $(\mu_0 - \mu_1)^T \mathbf{w} = \gamma$:

$$\gamma(\mu_0 - \mu_1) = \lambda \mathbf{S}_w \mathbf{w}$$

$$\mathbf{w} = \frac{\gamma}{\lambda} \mathbf{S}_w^{-1} (\mu_0 - \mu_1)$$

Since the final solution of \mathbf{w} only concerns its direction rather than its magnitude, its magnitude can be assigned arbitrarily. Moreover, as μ_0 and μ_1 have fixed magnitudes, γ is affected only by the magnitude of \mathbf{w} . Therefore, by adjusting the magnitude of \mathbf{w} , we can make $\gamma = \lambda$.

It finally gives: $\mathbf{w} = \mathbf{S}_w^{-1} (\mu_0 - \mu_1)$

Extend LDA to multiclass classification problems

- The global scatter matrix

$$\begin{aligned} \mathbf{S}_t &= \mathbf{S}_b + \mathbf{S}_w \\ &= \sum_{i=1}^m (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T \end{aligned}$$

- The within-class scatter matrix

$$\mathbf{S}_w = \sum_{i=1}^N \mathbf{S}_{w_i}$$

where

$$\mathbf{S}_{w_i} = \sum_{\mathbf{x} \in X_i} (\mathbf{x} - \boldsymbol{\mu}_i)(\mathbf{x} - \boldsymbol{\mu}_i)^T$$

$$\mathbf{S}_b = \mathbf{S}_t - \mathbf{S}_w$$

- We have

$$= \sum_{i=1}^N m_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T$$

Extend LDA to multiclass classification problems

$$\mathbf{S}_b = \mathbf{S}_t - \mathbf{S}_w$$

$$\begin{aligned} &= \sum_{i=1}^m (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T - \sum_{i=1}^N \sum_{\mathbf{x} \in X_i} (\mathbf{x} - \boldsymbol{\mu}_i)(\mathbf{x} - \boldsymbol{\mu}_i)^T \\ &= \sum_{i=1}^N \left(\sum_{\mathbf{x} \in X_i} ((\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T - (\mathbf{x} - \boldsymbol{\mu}_i)(\mathbf{x} - \boldsymbol{\mu}_i)^T) \right) \\ &= \sum_{i=1}^N \left(\sum_{\mathbf{x} \in X_i} ((\mathbf{x} - \boldsymbol{\mu})(\mathbf{x}^T - \boldsymbol{\mu}^T) - (\mathbf{x} - \boldsymbol{\mu}_i)(\mathbf{x}^T - \boldsymbol{\mu}_i^T)) \right) \\ &= \sum_{i=1}^N \left(\sum_{\mathbf{x} \in X_i} (\mathbf{x}\mathbf{x}^T - \mathbf{x}\boldsymbol{\mu}^T - \boldsymbol{\mu}\mathbf{x}^T + \boldsymbol{\mu}\boldsymbol{\mu}^T - \mathbf{x}\mathbf{x}^T + \mathbf{x}\boldsymbol{\mu}_i^T + \boldsymbol{\mu}_i\mathbf{x}^T - \boldsymbol{\mu}_i\boldsymbol{\mu}_i^T) \right) \\ &= \sum_{i=1}^N \left(\sum_{\mathbf{x} \in X_i} (-\mathbf{x}\boldsymbol{\mu}^T - \boldsymbol{\mu}\mathbf{x}^T + \boldsymbol{\mu}\boldsymbol{\mu}^T + \mathbf{x}\boldsymbol{\mu}_i^T + \boldsymbol{\mu}_i\mathbf{x}^T - \boldsymbol{\mu}_i\boldsymbol{\mu}_i^T) \right) \\ &= \sum_{i=1}^N \left(-\sum_{\mathbf{x} \in X_i} \mathbf{x}\boldsymbol{\mu}^T - \sum_{\mathbf{x} \in X_i} \boldsymbol{\mu}\mathbf{x}^T + \sum_{\mathbf{x} \in X_i} \boldsymbol{\mu}\boldsymbol{\mu}^T + \sum_{\mathbf{x} \in X_i} \mathbf{x}\boldsymbol{\mu}_i^T + \sum_{\mathbf{x} \in X_i} \boldsymbol{\mu}_i\mathbf{x}^T - \sum_{\mathbf{x} \in X_i} \boldsymbol{\mu}_i\boldsymbol{\mu}_i^T \right) \\ &= \sum_{i=1}^N (-m_i\boldsymbol{\mu}_i\boldsymbol{\mu}^T - m_i\boldsymbol{\mu}\boldsymbol{\mu}_i^T + m_i\boldsymbol{\mu}\boldsymbol{\mu}^T + m_i\boldsymbol{\mu}_i\boldsymbol{\mu}_i^T + m_i\boldsymbol{\mu}_i\boldsymbol{\mu}_i^T - m_i\boldsymbol{\mu}_i\boldsymbol{\mu}_i^T) \\ &= \sum_{i=1}^N (-m_i\boldsymbol{\mu}_i\boldsymbol{\mu}^T - m_i\boldsymbol{\mu}\boldsymbol{\mu}_i^T + m_i\boldsymbol{\mu}\boldsymbol{\mu}^T + m_i\boldsymbol{\mu}_i\boldsymbol{\mu}_i^T) \\ &= \sum_{i=1}^N m_i (-\boldsymbol{\mu}_i\boldsymbol{\mu}^T - \boldsymbol{\mu}\boldsymbol{\mu}_i^T + \boldsymbol{\mu}\boldsymbol{\mu}^T + \boldsymbol{\mu}_i\boldsymbol{\mu}_i^T) \\ &= \sum_{i=1}^N m_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T \end{aligned}$$

Extend LDA to multiclass classification problems

□ Objective

$$\max_{\mathbf{W}} \frac{\text{tr}(\mathbf{W}^T \mathbf{S}_b \mathbf{W})}{\text{tr}(\mathbf{W}^T \mathbf{S}_w \mathbf{W})}$$

$$\mathbf{W} = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_i, \dots, \mathbf{w}_{N-1}) \in \mathbb{R}^{d \times (N-1)}$$

where $\mathbf{W} \in \mathbb{R}^{d \times (N-1)}$

$$\begin{cases} \text{tr}(\mathbf{W}^T \mathbf{S}_b \mathbf{W}) = \sum_{i=1}^{N-1} \mathbf{w}_i^T \mathbf{S}_b \mathbf{w}_i \\ \text{tr}(\mathbf{W}^T \mathbf{S}_w \mathbf{W}) = \sum_{i=1}^{N-1} \mathbf{w}_i^T \mathbf{S}_w \mathbf{w}_i \end{cases}$$



$$\mathbf{S}_b \mathbf{W} = \lambda \mathbf{S}_w \mathbf{W}$$

Concatenating the eigenvectors corresponding to the d' largest non-zero eigenvalues of $\mathbf{S}_w^{-1} \mathbf{S}_b$ leads to the closed-form solution of \mathbf{W} , where $d' \leq N - 1$

- Since the projection reduces the data while considering the class information, LDA is also considered as a classic supervised dimensionality reduction technique

Multiclass Classification

❑ Multiclass Classification learning methods

- Some binary classification methods can be directly extended to accommodate multiclass cases
- Apply some strategies to solve multiclass classification problems with any existing binary classification methods
(more general)
 - Decompose the problem and then train a binary classifier for each divided binary classification problem
 - Ensemble the outputs collected from all binary classifiers into the final multiclass predictions

❑ Dividing strategies

- One vs. One (OvO)
- One vs. Rest (OvR)
- Many vs. Many (MvM)

Multiclass Classification - OvO

□ In the decomposing phase

- puts the N classes into pairs
 - $N(N-1)/2$ binary classification tasks
- trains a classifier for each task
 - $N(N-1)/2$ classifiers

□ In the testing phase

- a new sample is classified by all classifiers
 - $N(N-1)/2$ classification outputs
- the final prediction can be made via voting
 - the predicted class is the one received the most votes

Multiclass Classification - OvR

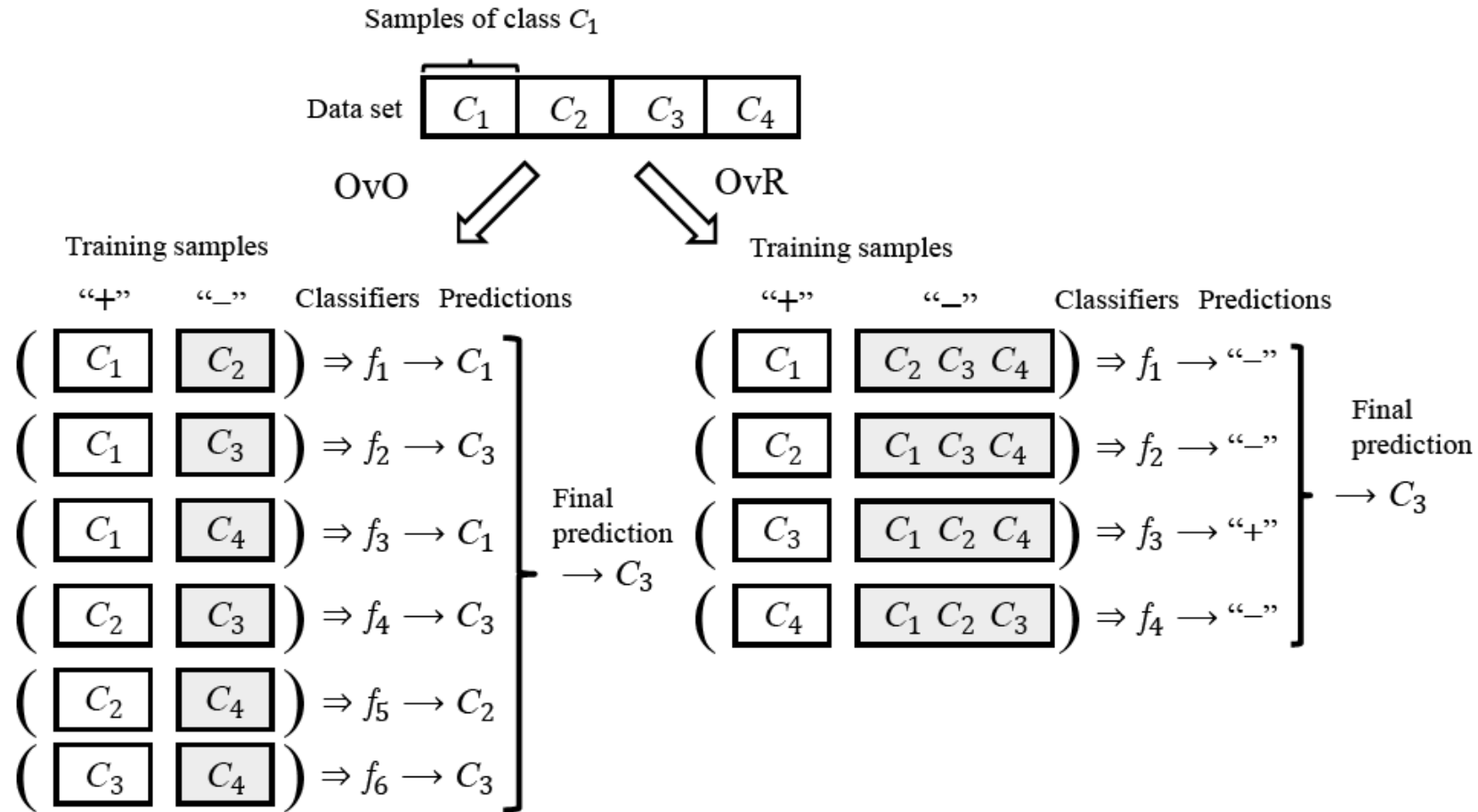
□ In the decomposing phase

- consider each class as positive in turn, and the rest classes are considered as negative
 - N binary classification tasks
- trains a classifier for each task
 - N classifiers

□ In the testing phase

- a new sample is classified by all classifiers
 - N classification outputs
- the prediction confidences are usually assessed
 - the class with the highest confidence is used as the classification result

Multiclass Classification – A comparison between OvO and OvR



Multiclass Classification – A comparison between OvO and OvR

OvO

- ❑ Train $N(N-1)/2$ classifiers, the memory and testing time costs are often higher
- ❑ Each classifier uses only samples of two classes. Hence, the computational cost of training OvO is lower

OvR

- ❑ Train N classifiers, the memory and testing time costs are often lower
- ❑ Each classifier uses all training samples. Hence, the computational cost of training OvO is higher

As for the prediction performance, it depends on the specific data distribution, and in most cases, the two methods have similar performance.

Multiclass Classification - MvM

- ❑ Many vs Many(MvM)
 - MvM conducts multiple trials, and each trial puts several classes as positive and several classes as negative.
- ❑ Error Correcting Output Code, ECOC

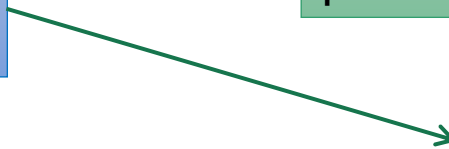
Encoding: split the N classes M times, where each time splits some classes as positive and some classes as negative.

A total of M training sets are Generated. The base codeword of each class are calculated.

Decoding: use the M classifiers to predict a testing sample

The class with the shortest distance is returned as the final prediction.

Combine the predicted labels into a codeword



Multiclass Classification - MvM

❑ Error Correcting Output Code, ECOC

	f_1	f_2	f_3	f_4	f_5	Hamming distance	Euclidean distance		f_1	f_2	f_3	f_4	f_5	f_6	f_7	Hamming distance	Euclidean distance
$C_1 \rightarrow$	-1	+1	-1	+1	+1	$\rightarrow 3$	$2\sqrt{3}$	$C_1 \rightarrow$	-1	-1	+1	+1	-1	+1	+1	$\rightarrow 4$	4
$C_2 \rightarrow$	+1	-1	-1	+1	-1	$\rightarrow 4$	4	$C_2 \rightarrow$	-1	0	0	0	+1	-1	0	$\rightarrow 2$	2
$C_3 \rightarrow$	-1	+1	+1	-1	+1	$\rightarrow 1$	2	$C_3 \rightarrow$	+1	+1	-1	-1	-1	+1	-1	$\rightarrow 5$	$2\sqrt{5}$
$C_4 \rightarrow$	-1	-1	+1	+1	-1	$\rightarrow 2$	$2\sqrt{2}$	$C_4 \rightarrow$	-1	+1	0	+1	-1	0	+1	$\rightarrow 3$	$\sqrt{10}$
Testing instance \rightarrow	-1	-1	+1	-1	+1	\uparrow	\uparrow	Testing instance \rightarrow	-1	+1	+1	-1	+1	-1	+1	\uparrow	\uparrow

Deactivated class,
Hamming distance is
0.5 not 1

(a) Binary ECOC coding.

[Dietterich and Bakiri, 1995]

(b) Ternary ECOC coding.

[Allwein et al. 2000]

- The ECOC codeword has the error tolerance and correction ability. a longer ECOC codeword produces better correction ability
- In theory, the correction ability of a fixed length codeword increases as the distances between classes increase.

Class Imbalance Problem

□ Class imbalance

- a significantly different number of samples for each class.
(the positive class is the minority)

the classes are balanced $\frac{y}{1-y} > 1$  $\frac{y}{1-y} > \frac{m^+}{m^-}$ the observed class ratio

□ Rescaling

- undersampling
 - some negative samples are selectively dropped so that the classes are balanced (EasyEnsemble [Liu et al.,2009])
- oversampling
 - increase the number of positive samples so that the classes are balanced (SMOTE [Chawla et al.2002])
- threshold-moving

$$\frac{y}{1-y} > \frac{m^+}{m^-}$$

Summary

- ❑ The objective of optimization of each model
 - Least squares method: Minimize the mean-square error
 - Logistic regression: Maximizing the likelihood of sample distribution
 - Linear discriminant analysis: minimize the within-class scatter matrix and maximize the between-classes scatter matrix

- ❑ Optimization method of parameters
 - Least squares method: Linear Algebra
 - Logistic regression: convex optimization, Newton's method
 - Linear discriminant analysis: matrix theory, generalized Rayleigh quotient

Summary

- Linear Regression
 - Least Squares Method (minimize the mean-square error)

- Binary Classification Problem
 - Logistic Regression
 - Unit-step function, logistic function, maximum likelihood method
 - Linear Discriminant Analysis
 - Maximizing generalized Rayleigh quotient

- Multiclass Classification Problem
 - One vs. One (OvO)
 - One vs. Rest (OvR)
 - Many vs. Many (MvM)
 - Error Correcting Output Code

- Class Imbalance Problem
 - Strategy: rescaling

Thanks!