# Chapter 4 Linear Models

LECTURER:   ZHIHUA JIANG

# Content

3.1 Linear model

3.2 Linear regression

3.3 Logistic regression

3.4 Linear discriminant analysis (LDA)

# 3.1 Linear model

- Linear model: linear function of attributes

$$f(\boldsymbol{x}) = w_1 x_1 + w_2 x_2 + \ldots + w_d x_d + b$$

- Matrix form:

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

$$\boldsymbol{\beta} = (\boldsymbol{w}; b), \quad \hat{\boldsymbol{x}} = (\boldsymbol{x}; 1)$$

$$\boldsymbol{w}^{\mathrm{T}} \boldsymbol{x} + b \longrightarrow \boldsymbol{\beta}^{\mathrm{T}} \hat{\boldsymbol{x}}$$

- Advantages:
  simple model; basic model; good interpretability

# *3.2 Linear regression*

- **Univariate linear regression**

$$f(x_i) = wx_i + b$$ such that $f(x_i) \approx y_i$

*Where $x_i$ is a scalar*

- **Multivariate linear regression**

$$f(\boldsymbol{x}_i) = \boldsymbol{w}^{\mathrm{T}}\boldsymbol{x}_i + b$$ such that $f(x_i) \approx y_i$

*Where $x_i$ is a vector*

- **Generalized linear model**

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

*Where $g(\cdot)$ is a monotone differentiable function*

# *Univariate linear regression*

●How to determine *w* and *b*?

a) To minimize the *MSE*

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

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

# *Univariate linear regression*

- How to determine *w* and *b*? (cont.)
  b) Parameter estimation based on least square method

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

- Differentiate with *w* and *b*, respectively

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

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

$$\Longrightarrow$$

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

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

# Multivariate linear regression

- Given $D = \{(x_1, y_1), \ldots, (x_m, y_m)\}$,

$x_i = (x_{i1}; \ldots; x_{id}), y_i \in \mathbb{R} \rightarrow$ *m×(d+1) matrix X*

$$\hat{w} = (w; b)$$

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1d} & 1 \\ x_{21} & x_{22} & \cdots & x_{2d} & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_{m1} & x_{m2} & \cdots & x_{md} & 1 \end{pmatrix} = \begin{pmatrix} x_1^{\mathrm{T}} & 1 \\ x_2^{\mathrm{T}} & 1 \\ \vdots & \vdots \\ x_m^{\mathrm{T}} & 1 \end{pmatrix}$$

$$y = (y_1; y_2; \ldots; y_m)$$

# *Multivariate linear regression*

- Similarly, to minimize the MSE

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

$$E_{\hat{w}} = \left(y - X\hat{w}\right)^{\mathrm{T}} \left(y - X\hat{w}\right)$$

- Differentiate with $\hat{w}$

$$\frac{\partial E_{\hat{w}}}{\partial \hat{w}} = 2\,X^{\mathrm{T}}\left(X\hat{w} - y\right)$$

full-rank matrix, or

$$\Longrightarrow \qquad \hat{w}^* = \left(X^{\mathrm{T}}X\right)^{-1} X^{\mathrm{T}} y$$

positive definite matrix

# Generalized linear model

● Problem: how to let the linear prediction to approximate some function of real labels?

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

*Where g(·) is* the link function (*monotone & differentiable*),

*g⁻¹(·) is the inverse function*

● Example: lo                          ssion (对数线性回归)

$$\ln y = \boldsymbol{w}^{\mathrm{T}}\boldsymbol{x} + b \implies y = e^{W^{T}X+b}$$

# Example: log-linear regression

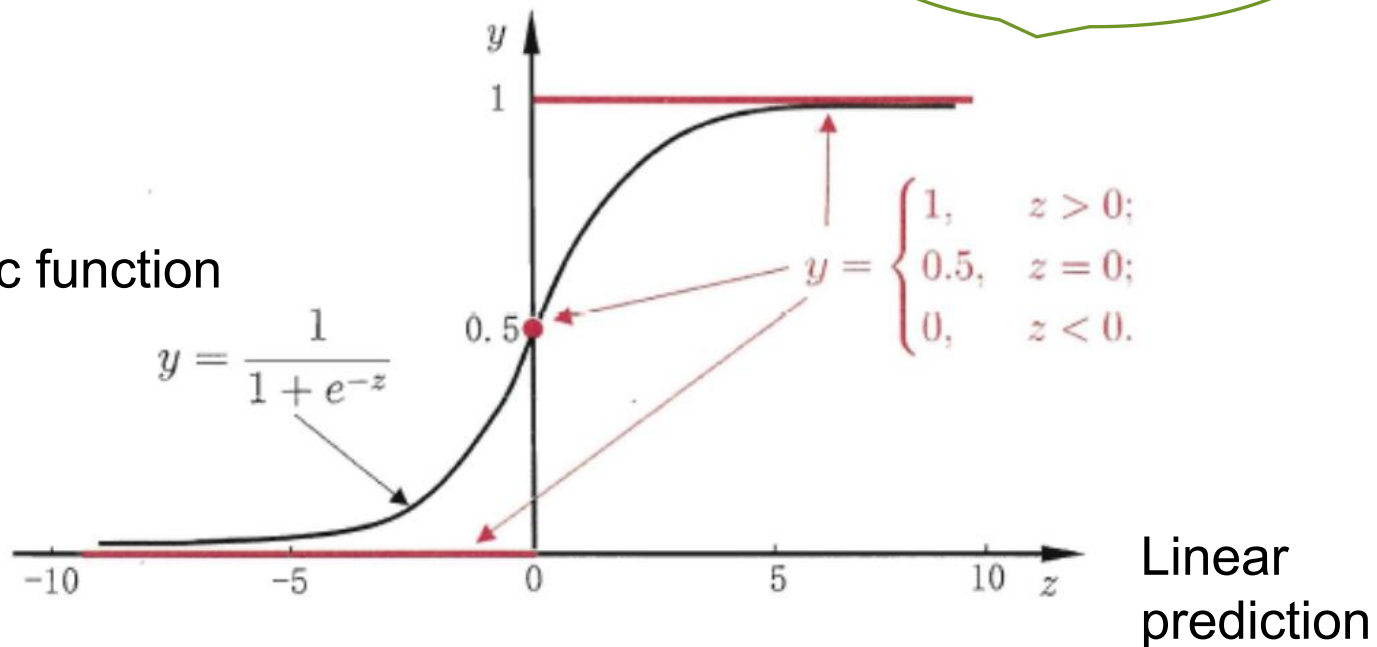# *3.3 logistic regression*

- Problem: use the linear regression for the classification problem?

- Solution: *Generalized linear model*

- Candidates of *g*(·):
  a) Unit-step function (单位阶跃函数)
  b) Logistic function (对数几率函数)

Unit-step function

Class label

Logistic function

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

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

Linear prediction

# Logistic function

● Odds(几率)

y: probability of x being a positive example

1-y: probability of x being a negative example

y/(1-y): relative probability of x being a positive example

● Logistic function: any order differentiable convex function has good mathematic properties

$$y = \frac{1}{1 + e^{-z}} \Longrightarrow y = \frac{1}{1 + e^{-(\boldsymbol{w}^{\mathrm{T}}\boldsymbol{x}+b)}} \Longrightarrow \ln\frac{y}{1-y} = \boldsymbol{w}^{\mathrm{T}}\boldsymbol{x} + b$$

# Logistic function

● Posterior probability estimation

$$\ln \frac{y}{1-y} = \boldsymbol{w}^{\mathrm{T}}\boldsymbol{x} + b \implies \ln \frac{p(y=1 \mid \boldsymbol{x})}{p(y=0 \mid \boldsymbol{x})} = \boldsymbol{w}^{\mathrm{T}}\boldsymbol{x} + b .$$

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

maximum likelihood method

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

$$\prod_{i=1}^{m} p(y_i \mid x_i; w, b)$$

?

# Logistic function

- Log likelihood function:

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

$$\boldsymbol{\beta} = (\boldsymbol{w}; b), \quad \hat{\boldsymbol{x}} = (\boldsymbol{x}; 1)$$

$$\boldsymbol{w}^{\mathrm{T}} \boldsymbol{x} + b \longrightarrow \boldsymbol{\beta}^{\mathrm{T}} \hat{\boldsymbol{x}}$$

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

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

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

# Logistic function

● Likelihood function (cont.)

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

$$\boldsymbol{w}^{\mathrm{T}}\boldsymbol{x} + b \longrightarrow \boldsymbol{\beta}^{\mathrm{T}}\hat{\boldsymbol{x}}$$

Max $l(w;b) \leftrightarrow$ Min $l(\beta)$

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

Newton method (update for the $t+1$ iteration)

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

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

$$l(w,b) = \sum_{i=1}^{m} \ln p(y_i \mid x_i; w, b)$$

$$\text{Max } l(w,b) \leftrightarrow \text{Min } l(\beta)$$

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

$$= \sum_{i=1}^{m} \ln[y_i \frac{e^{\beta^T \hat{x}_i}}{1 + e^{\beta^T \hat{x}_i}} + (1 - y_i) \frac{1}{1 + e^{\beta^T \hat{x}_i}}]$$

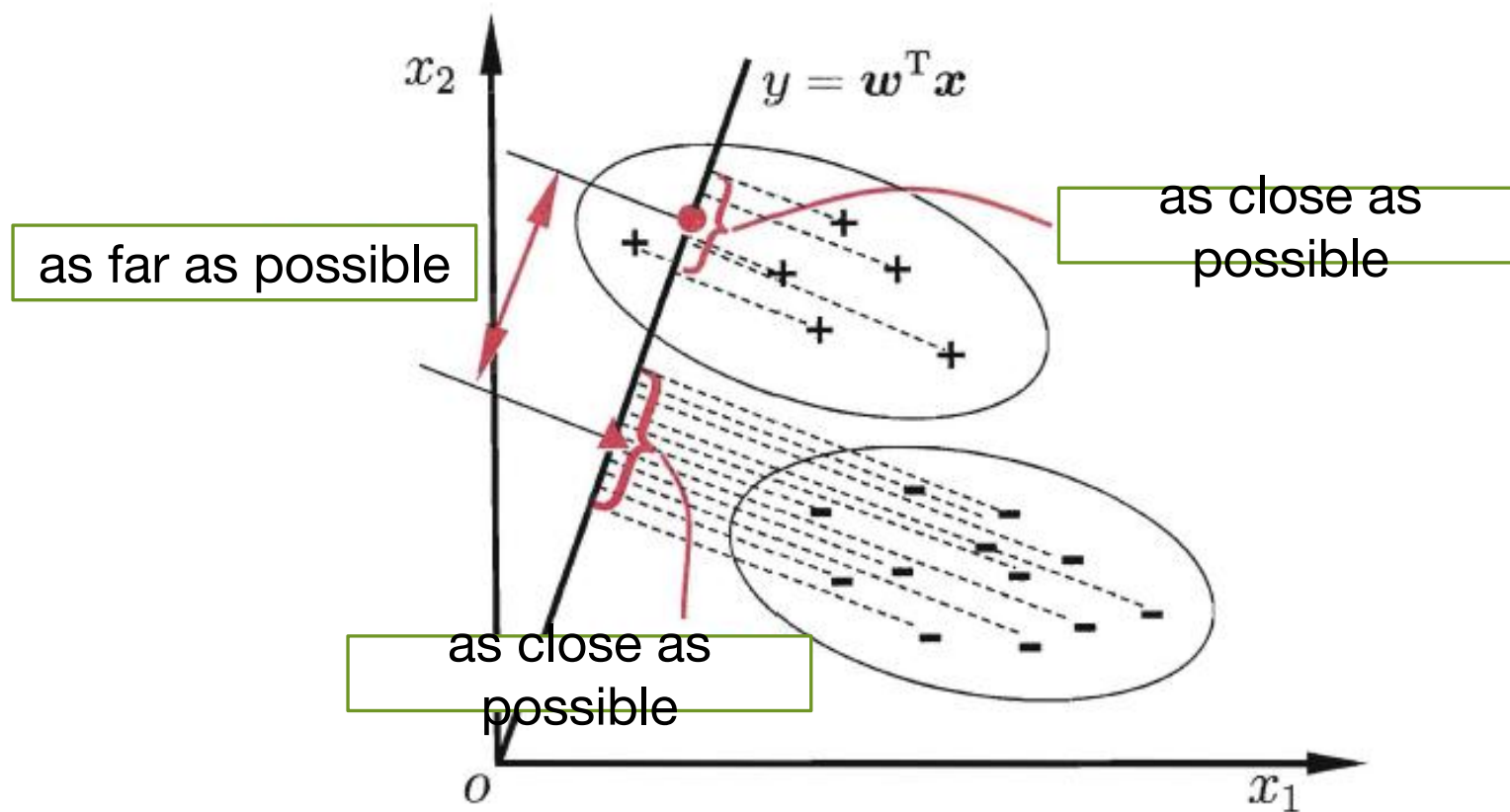$$= \sum_{i=1}^{m} \ln \frac{y_i e^{\beta^T \hat{x}_i} + (1 - y_i)}{1 + e^{\beta^T \hat{x}_i}}$$

$$= \begin{cases} y_i = 1, & \sum_{i=1}^{m} \ln \dfrac{e^{\beta^T \hat{x}_i}}{1 + e^{\beta^T \hat{x}_i}} = \sum_{i=1}^{m} (\beta^T \hat{x}_i - \ln(1 + e^{\beta^T \hat{x}_i})) \\[3em] y_i = 0, & \sum_{i=1}^{m} \ln \dfrac{1}{1 + e^{\beta^T \hat{x}_i}} = \sum_{i=1}^{m} (-\ln(1 + e^{\beta^T \hat{x}_i})) \end{cases}$$

# 3.4 Linear discriminant analysis (LDA)

- Idea:

  ○ Cast the samples onto a straight line

  ○ Project the similar samples as close as possible

  ○ Project the dissimilar samples as far as possible

  ○ For a new sample, determine the class according to the relative position of its projection point.

# 3.4 Linear discriminant analysis (LDA)



as far as possible

as close as possible

as close as possible

$x_2$

$y = \boldsymbol{w}^{\mathrm{T}}\boldsymbol{x}$

$x_1$

$o$

# Goal of LDA

- Given dataset $D = \{(x_i, y_i)\}_{i=1}^{m}, y_i \in \{0,1\}$
  - $X_i$: sample set of the $i$th class
  - $\mu_i$: mean vector of the $i$th class
  - $\Sigma_i$: covariance matrix of the $i$th class
  - Projection points of the two centers in the line: $w^{\mathsf{T}}\mu_0$ and $w^{\mathsf{T}}\mu_1$
  - Covariance: $w^{\mathsf{T}}\Sigma_0 w$ and $w^{\mathsf{T}}\Sigma_1 w$

- Similar samples as close as possible → $w^{\mathsf{T}}\Sigma_0 w + w^{\mathsf{T}}\Sigma_1 w$ as small as possible

$$\|\boldsymbol{w}^{\mathrm{T}}\boldsymbol{\mu}_0 - \boldsymbol{w}^{\mathrm{T}}\boldsymbol{\mu}_1\|_2^2$$

- Dissimilar samples as far as possible → as large as possible

$$\|\boldsymbol{x}\|_2 := \sqrt{x_1^2 + \cdots + x_n^2}.$$

If the entries in the column vector

$$\mathbf{X} = \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix}$$

are random variables, each with finite variance, then the covariance matrix Σ is the matrix whose (i, j) entry is the covariance

$$\Sigma_{ij} = \text{cov}(X_i, X_j) = \mathrm{E}[(X_i - \mu_i)(X_j - \mu_j)] = \mathrm{E}[X_i X_j] - \mu_i \mu_j,$$

where the operator E denotes the expected (mean) value of its argument, and

$$\mu_i = \mathrm{E}(X_i)$$

is the expected value of the i-th entry in the vector **X**. In other words,

$$\Sigma = \begin{bmatrix} \mathrm{E}[(X_1 - \mu_1)(X_1 - \mu_1)] & \mathrm{E}[(X_1 - \mu_1)(X_2 - \mu_2)] & \cdots & \mathrm{E}[(X_1 - \mu_1)(X_n - \mu_n)] \\ \mathrm{E}[(X_2 - \mu_2)(X_1 - \mu_1)] & \mathrm{E}[(X_2 - \mu_2)(X_2 - \mu_2)] & \cdots & \mathrm{E}[(X_2 - \mu_2)(X_n - \mu_n)] \\ \vdots & \vdots & \ddots & \vdots \\ \mathrm{E}[(X_n - \mu_n)(X_1 - \mu_1)] & \mathrm{E}[(X_n - \mu_n)(X_2 - \mu_2)] & \cdots & \mathrm{E}[(X_n - \mu_n)(X_n - \mu_n)] \end{bmatrix}.$$

# 3.4 Goal of LDA

- Goal: maximize

$$J = \frac{\|\boldsymbol{w}^{\mathrm{T}}\boldsymbol{\mu}_0 - \boldsymbol{w}^{\mathrm{T}}\boldsymbol{\mu}_1\|_2^2}{\boldsymbol{w}^{\mathrm{T}}\boldsymbol{\Sigma}_0\boldsymbol{w} + \boldsymbol{w}^{\mathrm{T}}\boldsymbol{\Sigma}_1\boldsymbol{w}}$$

$$= \frac{\boldsymbol{w}^{\mathrm{T}}(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^{\mathrm{T}}\boldsymbol{w}}{\boldsymbol{w}^{\mathrm{T}}(\boldsymbol{\Sigma}_0 + \boldsymbol{\Sigma}_1)\boldsymbol{w}}$$

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

- Within-class scatter matrix

$$\mathbf{S}_w = \boldsymbol{\Sigma}_0 + \boldsymbol{\Sigma}_1$$
$$= \sum_{\boldsymbol{x} \in X_0}(\boldsymbol{x} - \boldsymbol{\mu}_0)(\boldsymbol{x} - \boldsymbol{\mu}_0)^{\mathrm{T}} + \sum_{\boldsymbol{x} \in X_1}(\boldsymbol{x} - \boldsymbol{\mu}_1)(\boldsymbol{x} - \boldsymbol{\mu}_1)^{\mathrm{T}}$$

- Between-class scatter matrix

$$\mathbf{S}_b = (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^{\mathrm{T}}$$

# 3.4 Goal of LDA

- How to solve *w*?

$$L(w, \lambda) = -w^T S_b w + \lambda(w^T S_w w - 1)$$

$$\frac{\partial L}{\partial w} = -2S_b w + 2\lambda S_w w = 0 \Rightarrow S_b w = \lambda S_w w$$

Lagrange multipliers

$$J = \frac{w^T S_b w}{w^T S_w w}$$

set $w^T S_w w = 1$

$$\min_{w} \quad -w^T S_b w$$
$$\text{s.t.} \quad w^T S_w w = 1$$

$$S_b w = \lambda S_w w$$

[推导]：由公式 (3.36) 可得拉格朗日函数为

$$L(w, \lambda) = -w^{\mathrm{T}} S_b w + \lambda(w^{\mathrm{T}} S_w w - 1)$$

对 $w$ 求偏导可得

$$\frac{\partial L(w, \lambda)}{\partial w} = -\frac{\partial(w^{\mathrm{T}} S_b w)}{\partial w} + \lambda \frac{\partial(w^{\mathrm{T}} S_w w - 1)}{\partial w}$$
$$= -(S_b + S_b^{\mathrm{T}})w + \lambda(S_w + S_w^{\mathrm{T}})w$$

由于 $S_b = S_b^{\mathrm{T}}, S_w = S_w^{\mathrm{T}}$，所以

$$\frac{\partial L(w, \lambda)}{\partial w} = -2 S_b w + 2\lambda S_w w$$

令上式等于 0 即可得

$$-2 S_b w + 2\lambda S_w w = 0$$

$$S_b w = \lambda S_w w$$

由于我们想要求解的只有 $w$，而 $\lambda$ 这个拉格朗乘子具体取值多少都无所谓，因此我们可以任意设定 $\lambda$ 来配合我们求解 $w$。我们注意到

$$S_b w = (\mu_0 - \mu_1)(\mu_0 - \mu_1)^{\mathrm{T}} w$$

如果我们令 $\lambda$ 恒等于 $(\mu_0 - \mu_1)^{\mathrm{T}} w$，那么上式即可改写为

$$S_b w = \lambda(\mu_0 - \mu_1)$$

将其代入 $S_b w = \lambda S_w w$ 即可解得

$$w = S_w^{-1}(\mu_0 - \mu_1)$$

# Lagrange multipliers (拉格朗日乘子法)

- Idea:

  *d* variables, *k* constraints → *d+k* variables, 0 constraint

- Goal:

  *x* is *d*-vector, minimize *f*(*x*) s.t. *g*(*x*)=0

- ⬤        n:

$$L(\boldsymbol{x}, \lambda) = f(\boldsymbol{x}) + \lambda g(\boldsymbol{x})$$,

$$\nabla f(\boldsymbol{x}^*) + \lambda \nabla g(\boldsymbol{x}^*) = 0$$
$$g(\boldsymbol{x}) = 0.$$

Set partial
derivative to be 0