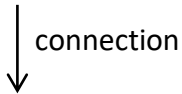


Linear classification (线性分类)

Linear regression



Linear classification

(Transform the linear function of \mathbf{w} using a nonlinear function $f(\cdot)$, $f(\cdot)$ is known as an activation function, whereas its inverse is called a link function)

$$Y = f(\mathbf{w}^T \mathbf{X} + \mathbf{b}) \longrightarrow \begin{cases} y \in \{0, 1\} & \left\{ \begin{array}{l} \text{Linear Discriminant Analysis (LDA 线性判别分析)} \\ \text{Perceptron (感知机)} \end{array} \right. \\ y \in \{0, 1\} & \left\{ \begin{array}{l} \text{Generative Model (生成式模型): GDA/Naïve Bayes ...} \\ \text{Discriminative Model (判别式模型): logistic regression...} \end{array} \right. \end{cases}$$

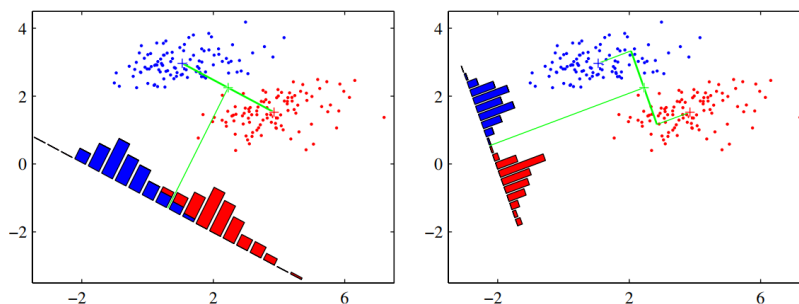
1. Linear Discriminant Analysis (LDA 线性判别分析)

Source(prml chapter4)

View a linear classification model is in terms of dimensionality reduction

Idea: give a large separation between the projected class means while also giving a small variance within each class

(类内距离小，类间距离大)



The second is better than the first one since the left one exists some overlaps between classes and the mean in the right one is also smaller.

$$X_{c1} = \{x_i | y_i = 1\} \quad X_{c2} = \{x_i | y_i = -1\}$$

$$|X_{c1}| = N_1, |X_{c2}| = N_2 \quad N_1 + N_2 = N$$

$$C_1: \bar{z}_1 = \frac{1}{N_1} \sum_{i=1}^{N_1} w^T x_i$$

$$S_1 = \frac{1}{N_1} \sum_{i=1}^{N_1} (w^T x_i - \bar{z}_1)(w^T x_i - \bar{z}_1)^T$$

$$C_2: \bar{z}_2 = \frac{1}{N_2} \sum_{i=1}^{N_2} w^T x_i$$

$$S_2 = \frac{1}{N_2} \sum_{i=1}^{N_2} (w^T x_i - \bar{z}_2)(w^T x_i - \bar{z}_2)^T$$

Model:

From the idea above, our goal function will be the ratio of the between-class variance to the within-class variance

$$J(w) = \frac{(\bar{z}_1 - \bar{z}_2)^2}{S_1 + S_2}$$

$$w^* = \underset{w}{\operatorname{argmax}} J(w)$$

$$\text{the between-class variance} = (\bar{z}_1 - \bar{z}_2)^2$$

$$\text{the within-class variance} = S_1 + S_2$$

$$(\bar{z}_1 - \bar{z}_2)^2$$

$$= \frac{1}{N_1} \sum_{i=1}^{N_1} w^T x_i - \frac{1}{N_2} \sum_{i=1}^{N_2} w^T x_i$$

$$= w^T \left(\frac{1}{N_1} \sum_{i=1}^{N_1} x_i - \frac{1}{N_2} \sum_{i=1}^{N_2} x_i \right)$$

$$= w^T (\bar{x}_{c1} - \bar{x}_{c2})^T$$

$$\begin{aligned}
& S_1 + S_2 \\
&= \frac{1}{N_1} \sum_{i=1}^{N_1} (w^T x_i - \bar{z}_1) (w^T x_i - \bar{z}_1)^T + \frac{1}{N_2} \sum_{i=1}^{N_2} (w^T x_i - \bar{z}_2) (w^T x_i - \bar{z}_2)^T \\
&= \frac{1}{N_1} \sum_{i=1}^{N_1} (w^T x_i - \frac{1}{N_1} \sum_{j=1}^{N_1} w^T x_j) (w^T x_i - \frac{1}{N_1} \sum_{j=1}^{N_1} w^T x_j)^T + \frac{1}{N_2} \sum_{i=1}^{N_2} (w^T x_i - \frac{1}{N_2} \sum_{j=1}^{N_2} w^T x_j) (w^T x_i - \frac{1}{N_2} \sum_{j=1}^{N_2} w^T x_j)^T \\
&= \frac{1}{N_1} \sum_{i=1}^{N_1} w^T (x_i - \frac{1}{N_1} \sum_{j=1}^{N_1} x_j) (x_i - \frac{1}{N_1} \sum_{j=1}^{N_1} x_j)^T w + \frac{1}{N_2} \sum_{i=1}^{N_2} w^T (x_i - \bar{x}_2) (x_i - \bar{x}_2)^T w \\
&= w^T \underbrace{\left\{ \frac{1}{N_1} \sum_{i=1}^{N_1} (x_i - \bar{x}_1) (x_i - \bar{x}_1)^T \right\}}_{\text{方差矩阵 } S_1} w + w^T \left\{ \frac{1}{N_2} \sum_{i=1}^{N_2} (x_i - \bar{x}_2) (x_i - \bar{x}_2)^T \right\} w
\end{aligned}$$

$$= w^T S_{c1} w + w^T S_{c2} w$$

$$= w^T (S_{c1} + S_{c2}) w$$

$$\therefore J(w) = \frac{w^T (\bar{x}_1 - \bar{x}_2) (\bar{x}_1 - \bar{x}_2)^T w}{w^T (S_{c1} + S_{c2}) w}$$

Let $S_b = (\bar{x}_1 - \bar{x}_2) (\bar{x}_1 - \bar{x}_2)^T$ (between class variance matrix) $p \times p$ matrix
 $S_w = (S_{c1} + S_{c2})$ (within class variance matrix) $p \times p$ matrix

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

$$J(w) = (w^T S_b w) (w^T S_w w)^{-1}$$

$$\frac{dJ(w)}{dw} = \frac{2S_b w (w^T S_w w)^{-1} + (w^T S_b w) \cdot (-1) (w^T S_w w)^{-2} (2S_w w)}{S_b w (w^T S_w w) - w^T S_b w (S_w \cdot w)} = 0$$

$$w^T S_b w S_w \cdot w = S_b w \cdot \underbrace{(w^T S_w \cdot w)}_{\text{实数} \in \mathbb{R}}$$

$$S_w \cdot w = \frac{w^T S_w \cdot w}{w^T S_b \cdot w} \cdot S_b \cdot w \quad \text{关心 } w \text{ 方向而非大小}$$

$$w = \frac{w^T S_w \cdot w}{w^T S_b \cdot w} S_w^{-1} S_b \cdot w \propto S_w^{-1} S_b \cdot w$$

$$= \frac{S_w^{-1} S_b \cdot w}{(\bar{x}_1 - \bar{x}_2) (\bar{x}_1 - \bar{x}_2)^T w}$$

(1xp) \cdot (px1) = 维

$$\propto S_w^{-1} (\bar{x}_1 - \bar{x}_2)$$

↓
 对称矩阵, 各向同性 $S_w^{-1} \propto I$
 进一步简化为 $w \propto (\bar{x}_1 - \bar{x}_2)$

2. Perceptron (感知机)

Idea: determine the parameters \mathbf{w} of the perceptron can most be motivated by error function minimization.(错误驱动)

Model:

$$f(x) = \text{sign}(w^T x), x \in \mathbb{R}^p, w \in \mathbb{R}^p$$

$$\text{sign}(a) = \begin{cases} +1 & a \geq 0 \\ -1 & a < 0 \end{cases}$$

Similar to the linear regression we want to minimize the loss function, use the idea above we know that if (x_i, y_i) is classified correctly:

$$\begin{matrix} w^T x_i > 0 & y_i = 1 \\ w^T x_i < 0 & y_i = -1 \end{matrix} \Rightarrow y_i w^T x_i > 0$$

So for the misclassified samples, we want them to be as less as possible.

$$L(w) = \sum_{i=1}^N \mathbb{I} \{ y_i w^T x_i < 0 \}$$

(the error is a piecewise constant function of \mathbf{w} with discontinuities, thus hard to find gradient)

Therefore consider an alternative error function known as the *perceptron criterion*.

The new loss function becomes

$$\begin{aligned} L(w) &= \sum_{x_i \in D} -y_i w^T x_i \\ \nabla_w L &= -y_i x_i \end{aligned}$$

Where D is the set that been separated wrong.

Then using the stochastic gradient descent, we got

$$\begin{aligned} w^{(t+1)} &\leftarrow w^{(t)} - \lambda \nabla_w L \quad (\lambda: \text{learning rate}) \\ w^{(t)} &\leftarrow w^{(t)} - \lambda_i y_i x_i \end{aligned}$$

Note: for data sets that are not linearly separable, the perceptron learning algorithm will never converge.

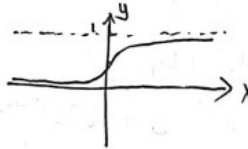
3. Logistic regression (逻辑回归)

(Source: cs 229 note1 & PRML chapter 4)

(1) binary classification

logistic Regression

Sigmoid function $\sigma(z) = \frac{1}{1+e^{-z}}$
 $\sigma: \mathbb{R} \mapsto (0, 1)$
 $w^T x \mapsto p$



$$\begin{cases} P(y=1|x) = \sigma(w^T x) = \frac{1}{1+\exp(-w^T x)}, y=1 & \text{let } P(y=1|x) = \psi(x_i, w) \\ P(y=0|x) = 1 - P(y=1|x) = \frac{\exp(-w^T x)}{1+\exp(-w^T x)}, y=0 & \text{let } P(y=0|x) = \log(1 - \psi(x_i, w)) \end{cases}$$

$$\rightarrow P(y|x) = p_1^y p_0^{1-y}$$

$$\begin{aligned} \text{MLE } \hat{w} &= \arg \max \log P(Y|X) \\ &= \arg \max \log \prod_{i=1}^N P(Y_i | X_i) \\ &= \arg \max \sum_{i=1}^N \log P(Y_i | X_i) \\ &= \arg \max \sum_{i=1}^N y_i \log \psi(x_i, w) + (1-y_i) \log (1 - \psi(x_i, w)) \end{aligned}$$

★ (max) MLE \Rightarrow loss function (min cross entropy)

Maximize the likelihood:

$$\begin{aligned} p_i' &= \left(\frac{1}{1+\exp(-w^T x)} \right)' = p_i(1-p_i) \\ \frac{\partial}{\partial w} \log P(Y|X) &= \sum_{i=1}^N (y_i - p_i) x_i \end{aligned}$$

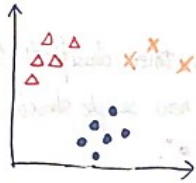
Using the stochastic gradient descent rule:

$$w^{(t+1)} \leftarrow w^{(t)} + \lambda \frac{\partial}{\partial w} \log P(Y|X)$$

λ : learning rate

(2) multi-class classification

(1) One-vs-all

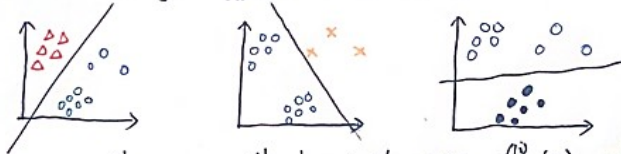


思想 idea: trans one multi-classification problem to many binary-classification problem.

Choose one of the classes and labeled as " $y=1$ ", the others are all labeled as " $y=0$ "
then similarly, choose one of the classes and labeled as " $y=2$ ", the others are all labeled as " $y=0$ "

\vdots

we got $f_w^i(x) = P(y=i|x, w)$ ($i=1, 2, \dots, k$) k : # of classes



then we will choose the max $f_w^i(x)$ of sample x that represent x belongs to the i th class.

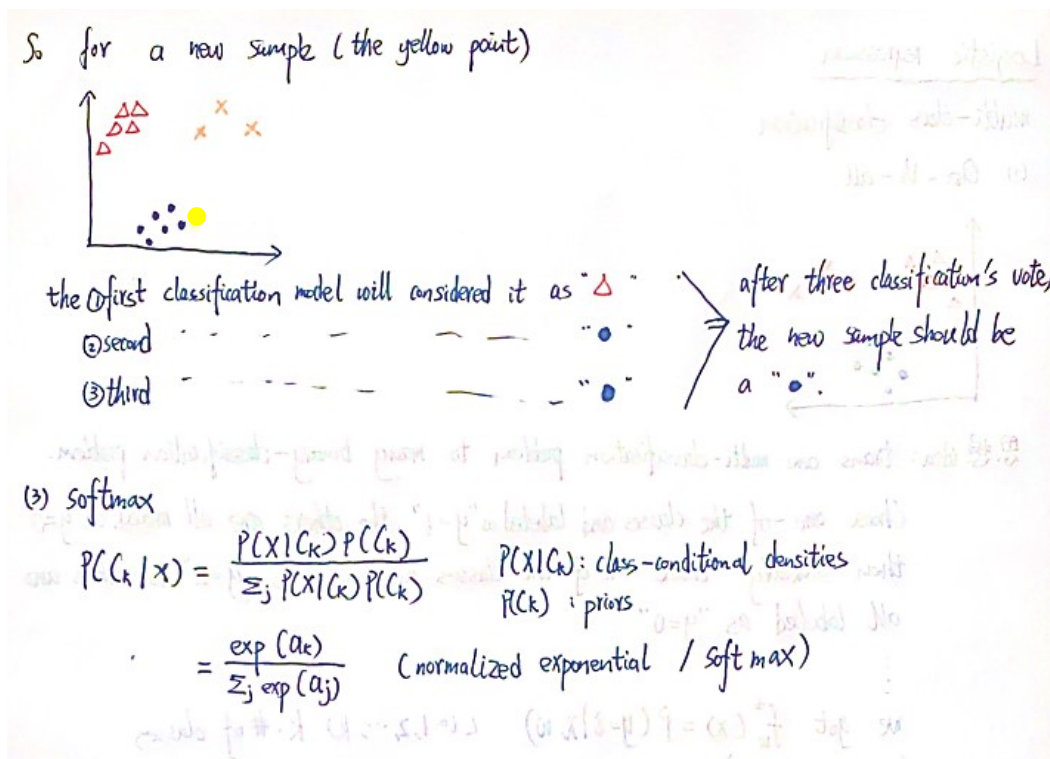
advantages: big range to use

disadvantages: the binary classification is unbalanced.

(2) One-vs-One

still the same sample





4. Gaussian Discriminant Analysis(GDA 高斯判别分析)

(source: <http://cs229.stanford.edu/notes/cs229-notes2.pdf>)

Idea: judge $p(y=0|x)$ and $p(y=1|x)$ which is bigger.

For binary classification

$$y \sim \text{Bernoulli}(\phi)$$

$$x|y=0 \sim \mathcal{N}(\mu_0, \Sigma)$$

$$x|y=1 \sim \mathcal{N}(\mu_1, \Sigma)$$

$$p(y) = \phi^y (1 - \phi)^{1-y}$$

$$p(x|y=0) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_0)^T \Sigma^{-1} (x - \mu_0)\right)$$

$$p(x|y=1) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1} (x - \mu_1)\right)$$

$$\begin{aligned}
 \text{log-likelihood: } \log(\theta) &= \log \prod_{i=1}^N P(X_i, y_i) \leftarrow \text{联合概率} \\
 \theta = (\mu, \Sigma, \phi) &= \sum_{i=1}^N \log [P(X_i | y_i) P(y_i)] \\
 \hat{\theta} = \arg \max_{\theta} (\log \theta) &= \sum_{i=1}^N [\log P(X_i | y_i) + \log P(y_i)] \\
 &= \sum_{i=1}^N [\log N(\mu, \Sigma)^{y_i} \cdot N(\mu_b, \Sigma)^{1-y_i} + \log \phi^{y_i} (1-\phi)^{1-y_i}] \\
 &= \sum_{i=1}^N [\underbrace{\log N(\mu, \Sigma)^{y_i}}_{\textcircled{1}} + \underbrace{\log N(\mu_b, \Sigma)^{1-y_i}}_{\textcircled{2}} + \underbrace{\log \phi^{y_i} (1-\phi)^{1-y_i}}_{\textcircled{3}}]
 \end{aligned}$$

求 ϕ : $\sum_{i=1}^N \log \phi^{y_i} (1-\phi)^{1-y_i}$

$$\frac{\partial \log \theta}{\partial \phi} = \sum_{i=1}^N \left[\frac{y_i}{\phi} - \frac{(1-y_i)}{1-\phi} \right] = 0$$

$$\Rightarrow \sum_{i=1}^N y_i (1-\phi) - \phi (1-y_i) = 0$$

$$\sum_{i=1}^N y_i - \phi y_i - \phi + \phi y_i = 0$$

$$\sum_{i=1}^N y_i - \phi = 0$$

$$\therefore \boxed{\hat{\phi} = \frac{1}{N} \sum_{i=1}^N y_i = \frac{1}{N_b} \sum_{i=1}^N y_i + \frac{1}{N_b} \sum_{i=1}^N y_i = \frac{1}{N_b} \sum_{i=1}^N y_i = \frac{M}{N}}$$

求 μ :

$$\textcircled{1} \sum_{i=1}^N \log N(\mu, \Sigma)^{y_i} \quad y_i = 0$$

$$= \sum_{i=1}^N y_i \log \left(\frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right\} \right) \quad \text{可去掉: 与 } \Sigma \text{ 无关只与 } \mu \text{ 有关}$$

$$\mu_1 = \arg \max_{\mu} \textcircled{1} = \arg \max_{\mu} \sum_{i=1}^N y_i \left(-\frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right)$$

$$\sum_{i=1}^N y_i \left(-\frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right)$$

$$= -\frac{1}{2} \sum_{i=1}^N y_i (x_i^T \Sigma^{-1} - \mu^T \Sigma^{-1}) (x_i - \mu)$$

$$= -\frac{1}{2} \sum_{i=1}^N y_i \underbrace{(x_i^T \Sigma^{-1} x_i)}_{\text{constant}} - \underbrace{x_i^T \Sigma^{-1} \mu}_{(1 \times p)(p \times p)(p \times 1)} - \underbrace{\mu^T \Sigma^{-1} x_i}_{(1 \times p)(p \times p)(p \times 1)} + \underbrace{\mu^T \Sigma^{-1} \mu}_{(1 \times p)(p \times p)(p \times 1)}$$

$\in \mathbb{R} \quad \quad \quad \in \mathbb{R}$

$$= -\frac{1}{2} \sum_{i=1}^N y_i (x_i^T \Sigma^{-1} x_i - 2 x_i^T \Sigma^{-1} \mu + \mu^T \Sigma^{-1} \mu)$$

$$\frac{\partial \log \theta}{\partial \mu} = -\frac{1}{2} \sum_{i=1}^N y_i (-2 \Sigma^{-1} x_i + 2 \Sigma^{-1} \mu) = 0$$

$$= \sum_{i=1}^N y_i (\Sigma^{-1} \mu - \Sigma^{-1} x_i) = 0$$

$$= \sum_{i=1}^N y_i (\mu - x_i) = 0$$

$$= \sum_{i=1}^N (y_i \mu - y_i x_i) = 0$$

$$\boxed{\hat{\mu}_1 = \frac{\sum_{i=1}^N y_i x_i}{\sum_{i=1}^N y_i} = \frac{1}{N_1} \sum_{i=1}^N y_i x_i}$$

$$\begin{aligned} \text{求 } \mu_2 : \textcircled{2} & \frac{N}{i=1} \log N(\mu_2, \Sigma)^{-1} y_i \\ & = \frac{N}{i=1} (1-y_i) \log \left(\frac{1}{(2\pi)^{p/2} |\Sigma|^{p/2}} \exp \left\{ -\frac{1}{2} (x_i - \mu_2)^T \Sigma^{-1} (x_i - \mu_2) \right\} \right) \end{aligned}$$

$$\mu_2 = \arg \max_{\mu_2} \textcircled{2} = \arg \max_{\mu_2} \frac{N}{i=1} (1-y_i) \left(-\frac{1}{2} (x_i - \mu_2)^T \Sigma^{-1} (x_i - \mu_2) \right)$$

$$\text{Similarly } \frac{\partial \textcircled{2}}{\partial \mu_2} = -\frac{1}{2} \frac{N}{i=1} (1-y_i) (-2 \Sigma^{-1} x_i + 2 \Sigma^{-1} \mu_2) = 0$$

$$\frac{N}{i=1} (1-y_i) (\mu_2 - x_i) = 0$$

$$\frac{N}{i=1} (1-y_i) \mu_2 - \frac{N}{i=1} (1-y_i) x_i = 0$$

$$\boxed{\mu_2 = \frac{\sum_{i=1}^N (1-y_i) x_i}{\sum_{i=1}^N (1-y_i)} = \frac{1}{N_2} \sum_{i=1}^N (1-y_i) x_i}$$

$$\text{求 } \Sigma : \textcircled{2} = \arg \max_{\Sigma} \textcircled{1} + \textcircled{2}$$

$$\text{let } C_1 = \{x_i | y_i = 1, i=1, 2, \dots, N\} \quad |C_1| = N_1$$

$$C_2 = \{x_i | y_i = 0, i=1, 2, \dots, N\} \quad |C_2| = N_2$$

$$\text{thus } \textcircled{1} + \textcircled{2} = \sum_{x_i \in C_1} \log N(\mu_1, \Sigma) + \sum_{x_i \in C_2} \log N(\mu_2, \Sigma)$$

$$\begin{aligned} \frac{N}{i=1} \log N(\mu, \Sigma) &= \frac{N}{i=1} \left(\log \frac{1}{(2\pi)^{p/2} |\Sigma|^{p/2}} \exp \left\{ -\frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right\} \right) \\ &= \frac{N}{i=1} \left(\log \frac{1}{(2\pi)^{p/2}} + \log |\Sigma|^{-\frac{p}{2}} - \frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right) \end{aligned}$$

$$\frac{N}{i=1} \left[-\frac{1}{2} \log |\Sigma| - \frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right]$$

$$= \tilde{C} - \frac{1}{2} \frac{N}{i=1} \log |\Sigma| - \frac{1}{2} \frac{N}{i=1} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$$

$$= \tilde{C} - \frac{1}{2} N \log |\Sigma| - \frac{1}{2} \frac{N}{i=1} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$$

$$x_i \in \mathbb{R}^p$$

$$(x_i - \mu)^T \in 1 \times p$$

$$(1 \times p)(p \times p)(p \times 1) \in \mathbb{R}$$

for real number this equals to its trace

$$\text{for } -\frac{1}{2} \frac{N}{i=1} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$$

$$= -\frac{1}{2} \frac{N}{i=1} \text{tr} \left[(x_i - \mu) (x_i - \mu)^T \Sigma^{-1} \right]$$

$$= -\frac{1}{2} \text{tr} \left\{ \sum_{i=1}^N (x_i - \mu) (x_i - \mu)^T \Sigma^{-1} \right\}$$

$$= -\frac{1}{2} N \text{tr}(\Sigma^{-1})$$

$$= \tilde{C} - \frac{1}{2} N \log |\Sigma| - \frac{1}{2} N \text{tr}(\Sigma^{-1})$$

$$\text{thus } \textcircled{1} + \textcircled{2} = -\frac{1}{2} N_1 \log |\Sigma| - \frac{1}{2} N_1 \text{tr}(\Sigma^{-1}) - \left(\frac{1}{2} N_2 \log |\Sigma| + \frac{1}{2} N_2 \text{tr}(\Sigma^{-1}) \right) + \tilde{C}$$

$$= -\frac{1}{2} N \log |\Sigma| - \frac{1}{2} N_1 \text{tr}(\Sigma^{-1}) - \frac{1}{2} N_2 \text{tr}(\Sigma^{-1}) + \tilde{C}$$

$$\frac{\partial (\textcircled{1} + \textcircled{2})}{\partial \Sigma} = -\frac{1}{2} N \frac{1}{|\Sigma|} |\Sigma| \cdot \Sigma^{-1} - \frac{1}{2} N_1 S_1^T \cdot (-1) \Sigma^{-2} - \frac{1}{2} N_2 S_2^T \cdot (-1) \Sigma^{-2}$$

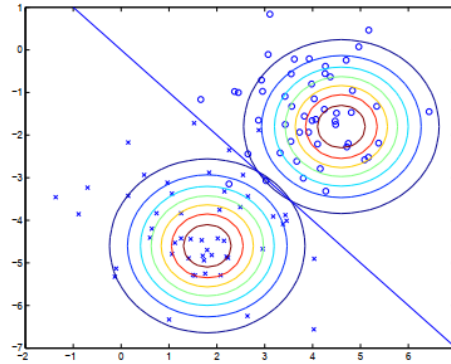
$$= -\frac{1}{2} N \Sigma^{-1} - \frac{1}{2} N_1 S_1 \Sigma^{-2} - \frac{1}{2} N_2 S_2 \Sigma^{-2} = 0$$

$$= 0 \cdot N \Sigma - N_1 S_1 - N_2 S_2 = 0$$

$$\boxed{\Sigma = \frac{1}{N} (N_1 S_1 + N_2 S_2)}$$

$\frac{\partial \text{tr}(AB)}{\partial A} = B^T$
$\frac{\partial A }{\partial A} = A \cdot A^{-1}$
$\text{tr}(AB) = \text{tr}(BA)$
$\text{tr}(ABC) = \text{tr}(CAB)$
$= \text{tr}(BCA)$

Pictorially, what the algorithm is doing can be seen in as follows:



In the figure is the straight line giving the decision boundary at which $p(y = 1 | x) = 0.5$.

Since we already calculate all parameters we need, we can calculate $p(y = 1 | x)$ and $p(y = 0 | x)$.

If $p(y = 1 | x) > p(y = 0 | x)$, we'll predict $y = 1$ to be the most likely outcome, otherwise, we'll predict $y = 0$.

Discussion: GDA and logistic regression

(1)

If $p(x | y)$ is multivariate Gaussian (with shared Σ), then $p(y | x)$ necessarily follows a logistic function. The converse, however, is not true. GDA makes stronger modeling assumptions, and is more data efficient.

(2)

Logistic regression is also more robust and less sensitive to incorrect modeling assumptions.

5. Naïve Bayes (朴素贝叶斯)

Source (cs229 note2)

(1) Naive Bayes (NB) assumption: x_i 's are conditionally independent given y .

$$P(x|y) = \prod_{j=1}^n p(x_j|y)$$

(2) Procedure:

The idea is to find the argmax posterior probability

$$\text{Data: } \{(x_i, y_i)\}_{i=1}^N$$
$$x_i \in \mathbb{R}^p, y_i \in \{1, 2, \dots, k\}$$

$$\hat{y} = \arg \max_y P(y|x) \quad P(y|x) = \frac{P(x,y)}{P(x)} = \frac{P(x|y) \cdot P(y)}{P(x)} \propto P(y) P(x|y)$$
$$= \arg \max_y P(y) P(x|y) \rightarrow \text{MLE 求解}$$

Step1: calculate the prior probability and condition probability

$$P(Y = c_k) = \frac{\sum_{i=1}^N I(y_i = c_k)}{N}, \quad k = 1, 2, \dots, K$$

$$P(X^{(j)} = a_{jl} | Y = c_k) = \frac{\sum_{i=1}^N I(x_i^{(j)} = a_{jl}, y_i = c_k)}{\sum_{i=1}^N I(y_i = c_k)}$$
$$j = 1, 2, \dots, n; \quad l = 1, 2, \dots, S_j; \quad k = 1, 2, \dots, K$$

Step2: calculate and decide the category of $x(j)$

$$P(Y = c_k) \prod_{j=1}^n P(X^{(j)} = x^{(j)} | Y = c_k), \quad k = 1, 2, \dots, K$$

$$y = \arg \max_{c_k} P(Y = c_k) \prod_j P(X^{(j)} = x^{(j)} | Y = c_k)$$

(3) Bayes estimation

Since that the probability might be 0, we need to use the Bayes estimation. The formula looks like this:

$$P_{\lambda}(X^{(j)} = a_{jl} | Y = c_k) = \frac{\sum_{i=1}^N I(x_i^{(j)} = a_{jl}, y_i = c_k) + \lambda}{\sum_{i=1}^N I(y_i = c_k) + S_j \lambda} \quad \lambda \geq 0$$

$$P_{\lambda}(Y = c_k) = \frac{\sum_{i=1}^N I(y_i = c_k) + \lambda}{N + K \lambda}$$

When $\lambda = 1$, we have Laplace smoothing

When $\lambda = 0$, we have MLE