

PART I. The Decision Theory:

For the two classes case:

$$C_1, C_2, \{x_i\}$$

$\{x_i\}$ are divided into regions $R_k \Leftarrow$ decision region

$$\begin{aligned} P\{\text{mistake}\} &= P\{x \in R_1, C_2\} + P\{x \in R_2, C_1\} \\ &= \int_{R_1} P\{x, C_2\} dx + \int_{R_2} P\{x, C_1\} dx \end{aligned}$$

To minimize the mistake, we should rearrange the labels.

If $P\{x, C_1\} > P\{x, C_2\} \Rightarrow$ assign C_1 to x

$$\therefore P\{x, C_1\} = P\{C_1|x\} P\{C_1\}$$

$$P\{x, C_2\} = P\{C_2|x\} P\{C_2\}$$

\therefore we only need to care about

$$P\{C_1|x\}$$
 and $P\{C_2|x\}$.

Thus, we could broke down the classification problem into two stages:

I) Inference stage: use training data to learn a model for

$$P\{C_k|x\}$$

II) Decision stage: use the estimated $P\{C_k|x\}$ to assign labels.

Otherwise, we could map x directly to the decision stage by

a discriminant function.

ii) Complexity

2) Generative Models.

1. Assume or determine the class conditional density $p_j^g(x|c_k)$ for each class k conditionally.

2. Infer or assume $p_j^g(c_k)$ i.e. $p_j^g(c_k) = \frac{N_k}{\sum_i N_i}$

$$3. p_j^g(c_k|x) = \frac{p_j^g(x|c_k) p_j^g(c_k)}{p_j^g(x)}$$

4. Use the decision theory based on $p_j^g(c_k|x)$

Advantage: Could determine $p(x)$ useful for outlier detection.

Disadvantage: Requires $p_j^g(x|c_k)$

b) Discriminant models.

1. Directly model the posterior distribution $p_j^g(c_k|x)$

2. Use the decision theory to assign labels

c) Directly find a discriminant function $f(x)$
maps to the class labels.

PART II

Approach c: Discriminant functions.

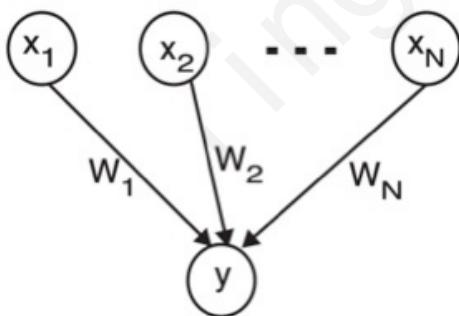
The perceptron algorithm:

Defn. $y = \text{sign}(w^T \phi(x))$

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

Error function:

$$E_p(w) = -\sum_{n \in M} w^T \phi_n t_n.$$



Since $\text{sign}(.)$ is a piece-wise function, we can use

Stochastic gradient descent to update w

$$w^{(t+1)} = w^{(t)} - \eta \nabla E_p(w) = w^{(t)} + \eta \phi_n t_n.$$

Substitute $w^{(T+1)}$ into $E(w)$:

$$E(w) = -w^{(T)} \phi_{n+1}$$

$$= -w^{(T)\top} \phi_{n+1} - \eta (\phi_{n+1})^\top \phi_{n+1}.$$

$$< -w^{(T)\top} \phi_{n+1}.$$

Perception Convergence Theorem

For any finite set of linearly separable labeled examples,
the Perception Learning Algorithm will halt after a finite
number of iterations.

Proof. Let w^* be the weight vector that
perfectly separates all samples.

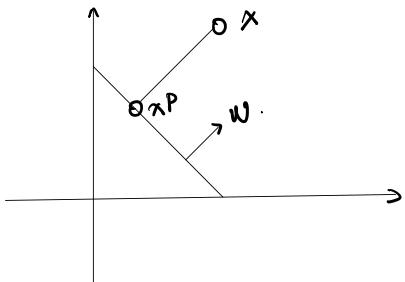
$$\Rightarrow \forall n. \quad w^{*\top} x_n t_n > 0$$

Let w^n denote the weight vector in the
 n th step.

Find: the angle between w^* and w^n

$$C(n) = \cos \theta(n) = \frac{\mathbf{w}^n \cdot \mathbf{w}^*}{\|\mathbf{w}^n\| \|\mathbf{w}^*\|}$$

Margin of a point. (The distance from the point x to w^*)



d : the vector from x to w^* with the minimum length

$$x^P = x - d$$

$\because d$ is parallel to w

$$\therefore d = \alpha \cdot w$$

$\therefore x^P \in \text{the separating plane } w^T x = 0$

$$w^T x^P = 0$$

$$w^T (x - d) = w^T (x - \alpha \cdot w) = 0$$

$$\Rightarrow \alpha = \frac{w^T \cdot x}{w^T \cdot w}$$

$$\Rightarrow \|d\|_2 = \frac{|w^T \cdot x|}{\|w\|_2}$$

$$g^n := \frac{t^n \cdot w^T \cdot x^n}{\|w^*\|_2} \Rightarrow \text{the signed margin.}$$

$\because w^T$ correctly separates all samples

$\therefore g^n$ is strictly positive.

$$\delta^* := \min_n \delta^n > 0.$$

$$w^n - w^{n-1} = \Delta w^n = \gamma t^n x^n.$$

\Leftrightarrow The update will happen if and only if

$$\text{there is an error. : } t^n \cdot w^{n-1} T \cdot x^n < 0.$$

$$\Rightarrow w^{(n-1)T} \cdot \Delta w^n$$

$$= w^{(n-1)T} \cdot \gamma t^n x^n < 0.$$

$$w^n T \cdot w^* = w^{(n-1)T} \cdot w^* + \Delta w^n T \cdot w^*$$

$$= w^{(n-1)T} w^* + \gamma t^n w^* T \cdot x^n$$

$$= w^{(n-1)T} w^* + \eta \delta^n \|w^*\|$$

$$\geq w^{(n-1)T} w^* + \eta \delta^* \|w^*\|$$

$$\geq w^{(n-1)T} w^* + 2\eta \delta^* \|w^*\|$$

$$\geq \eta n \delta^* \|w^*\|$$

$$\begin{aligned}
\|w^*\|^2 &= \|w^{n+1} + \Delta w^n\|_2 \\
&= \|w^{n+1}\|_2 + \|\Delta w^n\|_2 + 2w^{n+1}^\top \Delta w^n \\
&\quad \Leftarrow t^n = \underbrace{\gamma}_{-1}^{n+1} \\
&= \|w^{n+1}\|_2 + \gamma^2 \|x^{(n)}\|^2 + 2w^{n+1}^\top \Delta w^n \\
&< \|w^{n+1}\|_2 + \gamma^2 D^2
\end{aligned}$$

$$D := \max_n \|x^n\|$$

By applying the inequality n times:

$$\|w^n\|^2 \leq \|w^0\|_2 + n\gamma^2 D^2$$

\Rightarrow for sufficiently large n :

$$c_{1n} \geq \frac{\eta n \delta^* \|w^*\|}{\sqrt{n} \eta D \|w^*\|} = \frac{\delta^*}{D} \sqrt{n}$$

Assume, for contradiction, if the algorithm will go on forever,

c_{1n} will be greater than 1 for large n .

\Rightarrow After a finite step, the algorithm will halt.

Convergence Time of the Perceptron Algorithm

Define:

$$\delta_{\text{new}} = \max_{w^*} \delta^*$$

$$\therefore \frac{\delta^*}{D} \sqrt{n} \leq 1$$

$$\therefore n \leq \frac{D^2}{\delta^{*2}} \leq \frac{D^2}{\delta_{\text{new}}^2}$$

The convergence time depends on the margin.

If the two points with different labels are close to each other, the perceptron algorithm will take a long time

Definition. Augmente Margin: $\nu := \min_{i \in T(w)} |x_i \cdot w^*| = \delta^* \|w\|$

PART III: Probabilistic Generative Models: Bayes classifier

and Naive Bayes.

For the two classes cases:

$$\begin{aligned} P(c_1|x) &= \frac{P(x|c_1) P(c_1)}{P(x|c_1) P(c_1) + P(x|c_2) P(c_2)} \\ &= \frac{1}{1 + \frac{P(x|c_1) P(c_1)}{P(x|c_2) P(c_2)}} \end{aligned}$$

$$\text{let } a := \ln \frac{P\{x|c_1\} P\{c_1\}}{P\{x|c_2\} P\{c_2\}}$$

$$* = \frac{1}{1 + \exp(-a)}$$

Defn: logistic sigmoid function:

$$\sigma(a) := \frac{1}{1 + \exp(-a)}$$

Properties: \Rightarrow symmetric: $\sigma(a) = 1 - \sigma(-a)$

$$\text{II) Inverse: } a = \ln \frac{\sigma}{1-\sigma} \quad (\text{logistic function})$$

$$= \ln \frac{P\{c_1|x\}}{P\{c_2|x\}}$$

For $k=2$ classes:

$$P\{c_k|x\} = \frac{P\{x|c_k\} P\{c_k\}}{\sum_j P\{x|c_j\} P\{c_j\}}$$

$$= \frac{\exp\{a_k\}}{\sum_j \exp\{a_j\}} \quad \text{"softmax function"}$$

$$a_i := \ln P\{x|c_i\} P\{c_i\}$$

softmax function is a smoothed version of the max function. When $a_k \gg a_j$, $k \neq j$: $P\{c_k|x\} \approx 1$

1° If the inputs were all continuous variables

Assume, the class conditional densities were Gaussian.

$$P\{x|c_k\} = \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \exp\left\{-\frac{1}{2}(x-\mu_k)^T \Sigma^{-1} (x-\mu_k)\right\}.$$

For the two-class cases.

$$P\{c_1|x\} = \frac{P\{x|c_1\} P\{c_1\}}{P\{x|c_1\} P\{c_1\} + P\{x|c_2\} P\{c_2\}}$$

$$\begin{aligned} \ln \frac{P\{x|c_1\}}{P\{x|c_2\}} &= \ln \frac{\exp\left\{-\frac{1}{2}(x-\mu_1)^T \Sigma^{-1} (x-\mu_1)\right\}}{\exp\left\{-\frac{1}{2}(x-\mu_2)^T \Sigma^{-1} (x-\mu_2)\right\}} \\ &= \Sigma^{-1} \left[(\mu_1 - \mu_2)^T x - \frac{1}{2} \mu_1^T \mu_1 + \frac{1}{2} \mu_2^T \mu_2 \right] \end{aligned}$$

$$\Rightarrow P\{c_1|x\} = \epsilon(w^T x + w_0)$$

$$w := \Sigma^{-1}(\mu_1 - \mu_2)$$

$$w_0 := -\frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2 + \ln \frac{P(c_1)}{P(c_2)}$$

For the $k \geq 2$ classes cases:

$$a_{k|x} = \ln P\{x|c_k\} P\{c_k\}$$

$$\Rightarrow P\{c_k|x\} = \epsilon(w_k^T x + w_0)$$

$$w_k := \bar{\Sigma}^{-1} u_k$$

$$w_0 := -\frac{1}{2} u_k \bar{\Sigma}^{-1} u_k + \ln p_{C_k}^x \}$$

Once we have determined the distribution of $p_{C_k}^x$, together with p_{C_k} , we can get the maximum likelihood selection of the parameters.

For two - classes case:

$$\text{let } p_{C_1} = \pi, p_{C_2} = 1-\pi$$

$$\Rightarrow p_{X, C_1} = p_{C_1} p_{X|C_1} = \pi N(x_n | \mu_1, \Sigma)$$

$$p_{X, C_2} = p_{C_2} p_{X|C_2} = (1-\pi) N(x_n | \mu_2, \Sigma)$$

$$\Rightarrow p_{\text{tot}}(\pi, \mu_1, \mu_2, \Sigma) = \prod_{i=1}^N [\pi N(x_n | \mu_1, \Sigma)]^{t_i} [(1-\pi) N(x_n | \mu_2, \Sigma)]^{1-t_i}$$

$$\log p_{\text{tot}}(\pi, \mu_1, \mu_2, \Sigma) = \sum_{i=1}^N t_i \log \pi N(x_n | \mu_1, \Sigma) + (1-t_i) \log (1-\pi) N(x_n | \mu_2, \Sigma)$$

w.r.t. π :

The terms containing π are:

$$\sum_{i=1}^N \{ t_i \ln \pi + (1-t_i) \ln (1-\pi) \}$$

Take derivative, get:

$$\frac{1}{T_0} = \frac{1}{N} \sum_{n=1}^N t_n = \frac{N_1}{N} = \frac{N_1}{N_1 + N_2} \quad (\text{Note that } t_i = \begin{cases} 1 & \text{if } i \in C \\ 0 & \text{otherwise} \end{cases})$$

w.r.t. μ_1 :

$$\sum_{n=1}^N t_n \ln N \{x_n | \mu_1, \Sigma\} = -\frac{1}{2} \sum_{n=1}^N t_n (x_n - \mu_1)^T \Sigma^{-1} (x_n - \mu_1) + \text{const.}$$

Taking derivative:

$$\sum_{n=1}^N t_n (x_n - \mu_1)^T \Sigma^{-1} = 0.$$

$$\Rightarrow \hat{\mu}_1 = \frac{1}{N_1} \sum_{n=1}^N t_n x_n.$$

Similarly for μ_2 :

$$\hat{\mu}_2 = \frac{1}{N_2} \sum_{n=1}^N (1 - t_n) x_n.$$

w.r.t. Σ :

$$\begin{aligned} -\frac{1}{2} \sum_{n=1}^N t_n \ln |\Sigma| - \frac{1}{2} \sum_{n=1}^N t_n (x_n - \mu_1)^T \Sigma^{-1} (x_n - \mu_1) \\ -\frac{1}{2} \sum_{n=1}^N (1 - t_n) \ln |\Sigma| - \frac{1}{2} \sum_{n=1}^N (1 - t_n) (x_n - \mu_2)^T \Sigma^{-1} (x_n - \mu_2) \\ = -\frac{N}{2} \ln |\Sigma| - \frac{N}{2} \text{Tr} \{ \Sigma^{-1} S \}. \end{aligned}$$

$$S := \frac{N_1}{N} S_1 + \frac{N_2}{N} S_2$$

$$S_1 = N_1^{-1} \sum_{n \in C_1} (x_n - \mu_1) (x_n - \mu_1)^T$$

$$S_2 = N_2^{-1} \sum_{n \in C_2} (x_n - \mu_2) (x_n - \mu_2)^T$$

For multiple classes

$$P\{C_1, \dots, C_k\} = \pi_1 \cdot \dots \cdot \pi_k = 1 - \bar{\pi}_1 - \dots - \bar{\pi}_{k-1}$$

$$P\{x_n, c_j\} = P\{c_j\} P\{x_n | c_j\} = \pi_j N(x_n | \mu_j, \Sigma)$$

Similarly, the likelihood function is given by:

$$P\{x_n, t_u\} = \prod_{n=1}^N \prod_{k=1}^K \{P(\phi_n | c_k) \pi_k\}^{t_{nk}}$$

$$\ln P\{x_n, t_u\} = \sum_{n=1}^N \sum_{k=1}^K t_{nk} \{ \ln P(\phi_n | c_k) + \ln \pi_k \}$$

The problem is now a convex optimization with constraints.

$$\max \ln P\{x_n, t_u\} \quad \text{s.t. } \sum_k \pi_k = 1$$

⇒ By Lagrange:

$$\ln P\{x_n, t_u\} + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right)$$

Set the derivative to zero:

$$\sum_{n=1}^N \frac{t_{nk}}{\pi_k} + \lambda = 0$$

$$\Rightarrow -\pi_k \cdot \lambda = \sum_{n=1}^N t_{nk} = N k.$$

$$\sum_k -\bar{t}_{lk} \lambda = -\lambda = \sum_k N_k = N$$

$$\Rightarrow \lambda = -N$$

$$\Rightarrow \hat{\bar{t}}_{lk} = \frac{N_k}{N}$$

$$\ln p(\{x_n, t_n\} | \{\bar{t}_{lk}\}) = -\frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K t_{nk} \left\{ \ln |\Sigma| + (x_n - u_k)^T \Sigma^{-1} (x_n - u_k) \right\}$$

w.r.t. u_k :

$$\sum_{n=1}^N \sum_{k=1}^K t_{nk} \Sigma^{-1} (x_n - u_k) = 0$$

for $k=j$:

$$\sum_{n=1}^N t_{nj} \Sigma^{-1} (x_n - u_j) = 0$$

$$\Rightarrow \hat{u}_j = \frac{1}{N_j} \sum_{n=1}^N t_{nj} x_n.$$

For Σ :

Rewriting the log likelihood function as:

$$-\frac{1}{2} b \sum_{n=1}^N \sum_{k=1}^K t_{nk} \left\{ \ln |\Sigma| + \text{Tr} [\Sigma^{-1} (x_n - u_k)(x_n - u_k)^T] \right\}$$

Taking derivative:

$$\frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K t_{nk} \left\{ \Sigma - (x_n - u_k)(x_n - u_k)^T \right\} = 0.$$

$$\Rightarrow \hat{\mu}_k = \sum_{k=1}^K \frac{N_k}{N} \bar{x}_k$$

$$S_k = \frac{1}{N_k} \sum_{n=1}^N I(x_n \in C_k) (x_n - \bar{x}_k)^T$$

If we have multiple features, $x \in \mathbb{R}^{N \times D}$ is a matrix.

We may use the Naive Bayes Assumption.

$$P\{x|y\} = \prod_{\alpha=1}^D P\{x_\alpha | y\} = \prod_{i=1}^N \prod_{\alpha=1}^D P\{x_{i\alpha} | y\},$$

which means: the feature values are independent of the given label.

For the continuous features, we just need to repeat the above procedures:

$$P\{x_\alpha | c_k\} = N\{\mu_{\alpha k}, \sigma_{\alpha k}^2\}.$$

Taking derivative to the log likelihood w.r.t. $\mu_{\alpha k}$ and $\sigma_{\alpha k}^2$.

$$\Rightarrow \mu_{\alpha k} = \frac{1}{N_k} \sum_{i=1}^N I(t_i = c_k) x_{i\alpha}.$$

$$\sigma_{\alpha k}^2 = \frac{1}{N_k} \sum_{i=1}^N I(t_i = c_k) (x_{i\alpha} - \mu_{\alpha k})^2$$

2^o Dealing with Discrete Features.

If $\{x\}_a \in \{f_1, f_2, \dots, f_K\}$.

i.e. male / female single / married / widowed.

Under the Naive Bayes assumption, the log-likelihood function could be written as:

$$\begin{aligned}
 L(\theta) &= \sum_{i=1}^n \log P\{x^{(i)}, y^{(i)}\} \\
 &= \sum_{i=1}^n \log \left[\pi_{y_i} \prod_{j=1}^D P\{x_j^{(i)} | y^{(i)}\} \right] \\
 &= \sum_{i=1}^n \log \pi_{y_i} + \sum_{i=1}^n \sum_{j=1}^D \log P\{x_j^{(i)} | y^{(i)}\} \\
 &= \sum_{c_i \in \{1, \dots, k\}} \text{count}(c_i) \cdot \log \pi_{c_i} + \\
 &\quad \sum_{j=1}^D \sum_{c_i \in \{1, \dots, j\}} \sum_{x_j \in \{f_1, \dots, f_k\}} \text{count}_{j,y}(x_j | y) \log P\{x_j | c_i\}.
 \end{aligned}$$

By the previous derivation, we have verified that

the MLE of π_{c_i} is $\frac{\text{count}(c_i)}{N}$

Thus, to maximize the log Likelihood, we need to

maximize.

$$\sum_{x_j \in \{f_1, \dots, f_k\}} \text{const}_j \log P\{x_j | c_i\}.$$

Note that, x_j is a categorical variable, thus,

$P\{x_j | c_i\}$ is a multinomial distribution.

We simplify the above notations:

we want to find a distribution q , s.t.

$$q^* = \underset{q \in P_Y}{\operatorname{argmax}} \sum_{y \in Y} c_y \log q_y$$

$$P_Y := \left\{ q \in \mathbb{R}^{|\mathcal{Y}|} : \forall y \in \mathcal{Y}, q_y \geq 0, \sum_{y \in \mathcal{Y}} q_y = 1 \right\}.$$

↑:

$$\max \sum_{y \in Y} c_y \log q_y \quad \text{s.t. } \sum_{y \in Y} q_y = 1$$

By Lagrange:

$$g(u, q) = \sum_{y \in Y} c_y \log q_y - \alpha \left[\sum_{y \in Y} q_y - 1 \right]$$

$$\frac{d}{dq_y} g(u, q) = \frac{c_y}{q_y} - \alpha = 0$$

$$\Rightarrow q_y = \frac{c_y}{\lambda}$$

$$\text{By } \sum_y q_y = 1 : \quad q_y = \frac{c_y}{\sum_y c_y}$$

Substituting $c_y = \text{count}\{x_j | c_i\}$.

$$\Rightarrow P_j(x|y) = \frac{\text{count}(x_j | c_i)}{\sum_{x \in \{x_1, \dots, x_b\}} \text{count}(x_j | c_i)}$$

$$\Leftrightarrow P\{x_\alpha=j | y=c\} = \frac{\sum_{i=1}^n I[y_i=c \wedge x_{i\alpha}=j]}{\sum_{i=1}^n I[y_i=c]}$$

3^o Dealing with multinomial inputs.

$$x_\alpha \in \{0, 1, 2, \dots, m\} \quad \text{and} \quad m = \sum_{\alpha=1}^D x_\alpha.$$

By multinomial distribution:

$$P\{x | m, y=c\} = \frac{m!}{x_1! x_2! \cdots x_d!} \prod_{\alpha=1}^D [\theta_{\alpha c}]^{x_\alpha}$$

$\theta_{\alpha c}$: the parameter of the multinomial distribution

the probability of selecting x_α .

→ The likelihood function can be written as:

$$P\{X | y=c\} = \prod_{i=1}^N \prod_{\alpha=1}^D \frac{M_i!}{x_{i1}! \cdots x_{iD}!} [\theta_{i\alpha}]^{x_{i\alpha}}$$

By the above derivation, to obtain the MLE of $\theta_{i\alpha}$,

we just need to maximize:

$$\log P\{X | y=c\} = \sum_{i=1}^N \sum_{\alpha=1}^D x_{i\alpha} \log(\theta_{i\alpha}) + \log(M_i) - \log(x_{i1}! \cdots x_{iD}!)$$

$$\text{s.t. } \sum_{\alpha=1}^D \theta_{i\alpha} = 1$$

By Lagrange:

$$f(\theta, \lambda) = \sum_{i=1}^N \sum_{\alpha=1}^D x_{i\alpha} \log(\theta_{i\alpha}) + \lambda \left(1 - \sum_{j=1}^D \theta_{j\alpha}\right)$$

$$\frac{\partial}{\partial \theta_k} f(\theta, \lambda) = \sum_{i=1}^N \frac{x_{ik}}{\theta_k} - \lambda = 0$$

$$\Rightarrow \sum_{i=1}^N x_{ik} = \lambda \theta_{k\alpha}$$

$$\sum_{\alpha=1}^D \sum_{i=1}^N x_{i\alpha} = \sum_{\alpha=1}^D \lambda \theta_{k\alpha} = \lambda$$

$$\Rightarrow \hat{\theta}_{k\alpha} = \frac{\sum_{i=1}^N x_{ik}}{\sum_{\alpha=1}^D \sum_{i=1}^N x_{i\alpha}}$$

Note that the above derivation under the known of label.

$$\hat{\theta}_{kc} = \frac{\sum_{i=1}^n I[y_i = +1] x_{ik}}{\sum_{i=1}^n \sum_{\alpha=1}^D I[y_i = +1] x_{i\alpha}}$$

Theorem: Under some certain cases, the naive bayes classifier leads to a linear combination.

① Suppose that $y_i \in \{-1, +1\}$, and features are multinomial.

To prove the linearity, we need to show that:

$$h(x) = \arg \max_y \prod_{\alpha=1}^D P[x_\alpha | y] = \text{sign}(w^\top x + b)$$

which means:

$$w^\top x + b > 0 \Leftrightarrow h(x) = +1$$

Define:

$$P[x_\alpha | y = +1] \propto \theta_{\alpha+}^{x_\alpha}$$

$$P[Y = +1] = \pi_+$$

$$\text{let } [w]_\alpha = \log(\theta_{\alpha+}) - \log(\theta_{\alpha-})$$

$$b = \log(\pi_+) - \log(\pi_-)$$

Then:

$$w^T x + b > 0$$

$$\Leftrightarrow \sum_{\alpha=1}^D [x]_\alpha (\log(\theta_{\alpha+}) - \log(\theta_{\alpha-}))$$

$$+ \log(\pi_+) - \log(\pi_-)$$

$$\Leftrightarrow \exp \left\{ \sum_{\alpha=1}^D [x]_\alpha (\log(\theta_{\alpha+}) - \log(\theta_{\alpha-})) + \log(\pi_+) - \log(\pi_-) \right\} > 1$$

$$\Leftrightarrow \prod_{\alpha=1}^D \frac{\exp \left\{ \log \theta_{\alpha+}^{[x]_\alpha} + \log(\pi_+) \right\}}{\exp \left\{ \log \theta_{\alpha-}^{[x]_\alpha} + \log(\pi_-) \right\}} > 1$$

$$\Leftrightarrow \prod_{\alpha=1}^D \frac{\theta_{\alpha+}^{[x]_\alpha} \pi_+}{\theta_{\alpha-}^{[x]_\alpha} \pi_-} > 1$$

$$\Leftrightarrow \frac{\prod_{\alpha=1}^D P \{ x | \alpha | Y=+1 \} \pi_+}{\prod_{\alpha=1}^D P \{ x | \alpha | Y=-1 \} \pi_-} > 1$$

By NB
assumption

$$\frac{P \{ x | Y=+1 \} \pi_+}{P \{ x | Y=-1 \} \pi_-} > 1$$

$$\Leftrightarrow P \{ Y=+1 | x \} > P \{ Y=-1 | x \}$$

$$\Leftrightarrow \underset{y}{\operatorname{argmax}} P\{Y=y|x\} = +1$$

$$\text{Thus: } w^T x + b > 0 \Leftrightarrow +1 = \underset{y}{\operatorname{argmax}} P\{y\} \prod_{\alpha=1}^D P\{x_\alpha | y\}$$

② Gaussian Naive Bayes

If the covariance is shared between classes

$$h(x) = \underset{y}{\operatorname{argmax}} P\{y\} \prod_{\alpha=1}^D P\{x_\alpha | y\} = 1$$

$$\Leftrightarrow P\{x|y=1\} > P\{x|y=2\}$$

$$\Leftrightarrow \log \pi_1 - \frac{1}{2} (x - \mu_1)^T \Sigma^{-1} (x - \mu_1) > \log \pi_0 - \frac{1}{2} (x - \mu_0)^T \Sigma^{-1} (x - \mu_0)$$

$$\Leftrightarrow x^T \Sigma^{-1} x - 2\mu_1^T \Sigma^{-1} x + \mu_1^T \Sigma^{-1} \mu_1 + c > x^T \Sigma^{-1} x - 2\mu_0^T \Sigma^{-1} x + \mu_0^T \Sigma^{-1} \mu_0$$

$$\Leftrightarrow [2(\mu_0 - \mu_1)^T \Sigma^{-1}] x - (\mu_0 - \mu_1)^T \Sigma^{-1} (\mu_0 - \mu_1) > c$$

$$\Rightarrow w^T x + b > 0$$

$$w := 2(\mu_0 - \mu_1)^T \Sigma^{-1} \bar{x}$$

$$b := -(\mu_0 - \mu_1)^T \bar{\Sigma}^{-1} (\mu_0 - \mu_1) - c$$

PART IV. Probabilistic Discriminative Models.: Logistic Regression.

By the discussion above, we see that

$$P\{c_1 | x\} = \gamma(\phi) = \sigma(w^T x)$$

In this section, we directly model the conditional distribution $P\{c_1 | x\}$ rather than assuming the class conditional function $P\{x | c_i\}$ first.

Given a dataset: $\{\phi_i, t_i\} \quad i \in \{0, 1\}$.

$$\text{let } y_n = P\{c_1 | \phi_n\}$$

$$\Rightarrow P\{t | w\} = \prod_{n=1}^N y_n^{t_n} (1 - y_n)^{1-t_n}$$

Define the error function as:

$$E(w) = -\ln P\{t | w\} = -\sum_{n=1}^N \{t_n \ln y_n + (1-t_n) \ln (1-y_n)\} \quad \text{"cross entropy"}$$

$$\text{where } y_n = \sigma(w^T \phi_n)$$

For $\sigma(\cdot)$:

$$\begin{aligned} \frac{\partial \sigma}{\partial w} &= -\frac{-\exp(-w)}{(1+\exp(-w))^2} = \frac{\exp(-w)}{\exp(-w) + 2\exp(-w)} \\ &= \frac{\exp(-w)}{1+\exp(-w)} \cdot \frac{1}{1+\exp(-w)} = \sigma(1-\sigma) \end{aligned}$$

$$\begin{aligned}
\Rightarrow D E(\omega) &= \sum_{n=1}^N \left[-t_n \frac{1}{y_n} \cdot (1-y_n) \cdot y_n + \frac{1-t_n}{1-y_n} \cdot y_n (1-y_n) \right] \phi_n \\
&= \sum_{n=1}^N \left[-t_n (1-y_n) + y_n (1-t_n) \right] \phi_n \\
&= \sum_{n=1}^N \left[-t_n + y_n + t_n y_n - y_n t_n \right] \phi_n \\
&= \sum_{n=1}^N (y_n - t_n) \phi_n.
\end{aligned}$$

For $k \geq 2$ classes.

Having obtained above,

$$p_{\text{class } k}(\phi) = y_k(\phi) = \frac{\exp\{a_k\}}{\sum_j \exp\{a_j\}}$$

$$a_k := \omega_k^\top \phi$$

$$\frac{\partial y_k}{\partial a_j}$$

$$\text{If } j=k: \quad \frac{\exp\{a_k\} \sum_j \exp\{a_j\} - \exp\{a_k\} \cdot \exp\{a_k\}}{(\sum_j \exp\{a_j\})^2}$$

$$= \frac{\exp\{a_k\}}{\sum_j \exp\{a_j\}} - \frac{\exp\{a_k\} \cdot \exp\{a_k\}}{\sum_j \exp\{a_j\}} \frac{\exp\{a_k\}}{\sum_j \exp\{a_j\}}$$

$$= y_k - \underline{y_k^2} = y_k(1 - y_k)$$

$\# j \neq k:$

$$= - \frac{\exp(a_k) \cdot \exp(a_j)}{\left(\sum_j \exp(a_j)\right)^2}$$

$$= - y_k \cdot y_j$$

$$\Rightarrow \frac{\partial y_k}{\partial a_j} = y_k(I_{kj} - y_j)$$

The likelihood function. By one-hot encoding,

with N samples and k classes, our target label

matrix is $T \in \mathbb{R}^{N \times k}$

$$P\{T | w_1, \dots, w_k\} = \prod_{n=1}^N \prod_{k=1}^K p\{c_k | \phi_n\}^{t_{nk}}$$

$$= \prod_{n=1}^N \prod_{k=1}^K y_{nk}^{t_{nk}}$$

$$\Rightarrow E(w_1, \dots, w_k) = -\ln P\{T | w_1, \dots, w_k\}$$

$$= - \sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk}$$

$$\begin{aligned}
 \nabla_{w_j} E(w_1, \dots, w_k) &= \sum_{n=1}^N \sum_{k=1}^K \frac{\partial E}{\partial y_{nk}} \cdot \frac{\partial y_{nk}}{\partial w_j} \cdot \frac{\partial w_j}{\partial w_j} \\
 &= \sum_{n=1}^N - \sum_{b=1}^B \frac{t_{nb}}{y_{nk}} y_{nk} (I_{kj} - y_{nj}) \cdot \phi_n \\
 &= \sum_{n=1}^N (y_{nj} - t_{nj}) \cdot \phi_n.
 \end{aligned}$$

\Rightarrow Logistic Regression does not have a closed-form

solution, we need a sequential algorithm to approximate it.

Iterative reweighted least squares.

The Newton-Raphson iteration.

$f(x)$ takes root at r , we want to find r .

given a good estimate of r : x_0 , $r = x_0 + h$, h is "small"

$$0 = f(r) = f(x_0 + h) \approx f(x_0) + h f'(x_0)$$

$$\Rightarrow h \approx -\frac{f(x_0)}{f'(x_0)}$$

$$\Rightarrow r = x_0 + h \approx x_0 - \frac{f(x_0)}{f'(x_0)}$$

$$\text{let } x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}$$

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

To our problem, since $E(w)$ is clearly concave, we want to find the root of $\triangledown E(w)$. thus:

$$w^{(n+1)} = w^{(n)} - H^{-1} \triangledown E(w)$$

For $k=2$:

$$\triangledown E(w) = \sum_{n=1}^N (y_n - t_n) \phi_n = \Phi^T (y - t)$$

$$H = \triangledown \triangledown E(w) = \sum_{n=1}^N y_n (1-y_n) \phi_n \phi_n^T = \Phi^T R \Phi$$

R is a diagonal matrix: $R_{nn} = y_n (1-y_n)$

$$w^{(n+1)} = w^{(n)} - (\Phi^T R \Phi)^{-1} \Phi^T (y - t)$$

$$= (\Phi^T R \Phi)^{-1} \{ \Phi^T R \Phi w^{(n)} - \Phi^T (y - t) \}$$

$$= (\Phi^T R \Phi)^{-1} \Phi^T R z \quad \leftarrow \begin{array}{l} \text{the normal equation} \\ \text{in the least squares regression.} \end{array}$$

$$\text{where } z := \Phi^T w^{(n)} - R^{-1} (y - t)$$

$$\Rightarrow w^{(n+1)} = \underset{w}{\operatorname{arg\min}} (z - w^T \Phi)^T R (z - w^T \Phi)$$

For $k \geq 2$:

$$H_{ij} = \nabla w_j \nabla w_i^T E_{w0} = -\sum_{n=1}^N y_n (I_{nj} - y_j) \phi_n \phi_n^T$$

Ridge Logistic Regression.

Why we need Ridge estimation?

1° n : number of samples ; p : number of features.

If $p > n$: our discriminative boundary:

$$w^T \phi = 0.$$

$\because p > n \therefore$ the null space of w is non-trivial

\therefore for some $r > 0$, there exists:

$$w^T \phi = 0$$

$$w^T \phi + r \phi = 0$$

\Rightarrow an infinite number of estimators of the logistic regression exists.

2° When the dataset is perfectly linearly separable.

\therefore there exists an optimal hyperplane that separates correctly all the samples

\therefore For the optimal fit: we have $P\{y_i = 1 | \phi_i\} \in \{0, 1\}$.

⇒ The Loss function becomes:

$$L = - \sum_{i=1}^n t_i \log [P\{Y_i=1 | x_i\}] + (1-t_i) \log [P\{Y_i=0 | x_i\}] = 0$$

which does not contain any logistic parameters.

$$\Rightarrow P\{Y_i=1 | x_i\} = \frac{1}{1 + \exp[-w^T x_i]} \in [0, 1].$$

This will only happen when w approaches negative infinity or positive infinity.

⇒ The logistic regression cannot be learned from the separable data.

Thus, we may want to add a ridge penalty to the loss function.

$$E(w, \lambda) = - \sum_{i=1}^n [t_i \log(y_i) + (1-t_i) \log(1-y_i)] + \frac{\lambda}{2} \|w\|_2^2$$

$$\nabla_w E(w, \lambda) = \nabla_w E(w) + \lambda w$$

$$\nabla_{w_i} \nabla_{w_j} E(w, \lambda) = \nabla_{w_i} \nabla_{w_j} E(w) + \lambda I_{ij}$$

For $k=2$ case, we reformulate the Newton-Raphson algorithm.

$$\Delta E(w, \lambda) = \underline{\Phi}^T (y - \underline{\Phi}) + \lambda w$$

$$\nabla \Delta E(w, \lambda) = \underline{\Phi}^T R \underline{\Phi} + \lambda I_{pp}, \quad R_{nn} = y_n(1-y_n)$$

$$\Rightarrow \hat{w}^{(W)} = w^N - (\Phi^T R \Phi + \lambda I_{PP})^{-1} [\Phi^T (y - t) + \lambda v]$$

$$\text{let } V = \Phi^T R \Phi + \lambda I_{PP}$$

$$\begin{aligned} * &= \hat{w}^{(W)} - V^{-1} [\Phi^T (y - t) + \lambda v] \\ &= V^{-1} V \hat{w}^{(W)} - \lambda V^{-1} w^{(W)} - V^{-1} \Phi^T R V^{-1} (y - t) \\ &= V^{-1} \Phi^T R [\hat{w}^{(W)} - V^{-1} (y - t)] \end{aligned}$$

$$\text{let } z := \Phi \hat{w}^{(W)} - V^{-1} (y - t)$$

$$* = (\Phi^T R \Phi + \lambda I_{PP})^{-1} \Phi^T R \cdot z \quad \text{the Ridge estimation of the weighted least square}$$

Bayesian Logistic Regression.

The bayesian treatment of the logistic regression is more complex than the linear regression, since the posterior is no longer Gaussian.

Laplace Approximation.

Laplace Approximation is used to find a Gaussian approximation to a probability density over a set of continuous variables.

Suppose:

$$P(z) = \frac{1}{Z} f(z)$$

$$Z := \int f(z) dz \quad \text{"normalization coefficient"}$$

The Laplace method wants to find a Gaussian approximation whose mean is the mode of $P(z)$

To find the mode of $P(z)$:

We want to find a z_0 st:

$$\frac{d}{dz} P(z) \Big|_{z=z_0} = 0$$

$$\Leftrightarrow \frac{d}{dz} f(z) \Big|_{z=z_0} = 0$$

By Taylor expansion:

$$\ln f(z) \approx \ln f(z_0) - \frac{1}{2} A(z-z_0)^2$$

where: $A = - \left. \frac{d}{dz^2} \ln f(z) \right|_{z=z_0}$

$$\Rightarrow f(z) \approx f(z_0) \exp \left\{ - \frac{A}{2} (z-z_0)^2 \right\}$$

$$\begin{aligned}
 Z &= \int f(z) dz \\
 &= f(z_0) \int \exp\left\{-\frac{A}{2}(z-z_0)^2\right\} dz \\
 &= f(z_0) \cdot \sqrt{\frac{2\pi}{A}} \quad \Leftarrow \text{Gaussian Integral}
 \end{aligned}$$

\Rightarrow The approximation $f(z)$:

$$f(z) = \frac{1}{Z} f(z_0) = \left(\frac{A}{2\pi}\right)^{\frac{1}{2}} \exp\left\{-\frac{A}{2}(z-z_0)^2\right\}$$

For the multi-dimensional case:

The Taylor expansion step:

$$\ln f(z) \approx \ln f(z_0) - \frac{1}{2} (z-z_0)^T A (z-z_0)$$

$$A = -\nabla \ln f(z_0) \Big|_{z=z_0}$$

$$\Rightarrow f(z) \approx f(z_0) \exp\left\{-\frac{1}{2} (z-z_0)^T A (z-z_0)\right\}$$

\Rightarrow The approximation $f(z)$:

$$f(z) = \frac{|A|^{\frac{1}{2}}}{(2\pi)^{\frac{m}{2}}} \exp\left\{-\frac{1}{2} (z-z_0)^T A (z-z_0)\right\}$$

For the logistic regression, we firstly assume the prior is Gaussian:

$$p(w) = N\{w|w_0, S_0\}$$

⇒ The posterior is:

$$p\{w|t\} \propto p\{w\} p\{t|w\}$$

$$\Rightarrow \ln p\{w|t\} = -\frac{1}{2}(w-w_0)^T S_0^{-1}(w-w_0)$$

$$+ \sum_{n=1}^N \left\{ t_n \ln y_n + (1-t_n) \ln (1-y_n) \right\} + \text{const.}$$

$$\text{and } y_n = \sigma(w^T \phi)$$

To perform the Laplace approximation, we need firstly to

compute the w_{MAP} by maximizing the posterior distribution.

By Laplace Approximation, the value is:

$$S_N = -\nabla \nabla \ln p(w|t) = S_0^{-1} + \sum_{n=1}^N y_n(1-y_n) \phi_n \phi_n^T$$

⇒ The Gaussian Approximation:

$$q(w) = N\{w|w_{MAP}, S_N\}$$

The predictive distribution is then given by:

$$P\{c_1 | \phi, f\} = \int P\{c_1 | \phi, \omega\} P\{\omega | f\} d\omega$$

$$\simeq \int \epsilon(\omega^T \phi) f(\omega) d\omega.$$

$$\text{let } a = \omega^T \phi$$

$$\Rightarrow \epsilon(\omega^T \phi) = \int \delta(a - \omega^T \phi) \epsilon(\omega) da.$$

$\delta(\omega)$ is a Dirac delta function

$$\delta(x) = \begin{cases} 1 & , x=0 \\ 0 & , \text{otherwise.} \end{cases}$$

Thus,

$$\int \epsilon(\omega^T \phi) f(\omega) d\omega = \int \epsilon(a) p(a) da$$

$$\text{where } p(a) = \int \delta(a - \omega^T \phi) f(\omega) d\omega.$$

$$m_a = E[a] = \int p(\omega) \cdot a d\omega = \int f(\omega) \cdot \omega^T \phi d\omega = \omega_{MAP}^T \phi.$$

$$\sigma_a = \text{Var}[a] = \int p(\omega) \{a^2 - E[a]^2\} d\omega$$

$$= \int f(\omega) \{(a^T \phi)^2 - (m_a^T \phi)^2\} d\omega$$

$$= \phi^T S_N \phi$$

Thus:

$$P\{c_1 | \phi, +\} = \int \sigma(\omega) P(\omega) d\omega$$

$$= \int \sigma(\omega) N\{\alpha(\omega), \sigma^2\} d\omega. \quad "A convolution of the sigmoid and Gaussian"$$

The above does not have an analytical solution.

1° We could use MCMC to get an approximation of the analytical solution

2° We could approximate the sigmoid function with probit function to get an analytical expression.

$$\text{probit function: } \Phi(a) = \int_{-\infty}^a N(\theta | 0, 1) d\theta = \frac{1}{2} \left\{ 1 + \frac{1}{\sqrt{2}} \operatorname{erf}\left(\frac{a}{\sqrt{2}}\right) \right\}$$

For the method 2,

To obtain the best approximation, we need firstly to rescale the horizontal axis.

$$\begin{aligned} \left. \frac{d\sigma}{d\alpha} \right|_{\alpha=0} &= \sigma(0)(1 - \sigma(0)) \\ &= \frac{1}{2}(1 - \frac{1}{2}) = \frac{1}{4} \end{aligned}$$

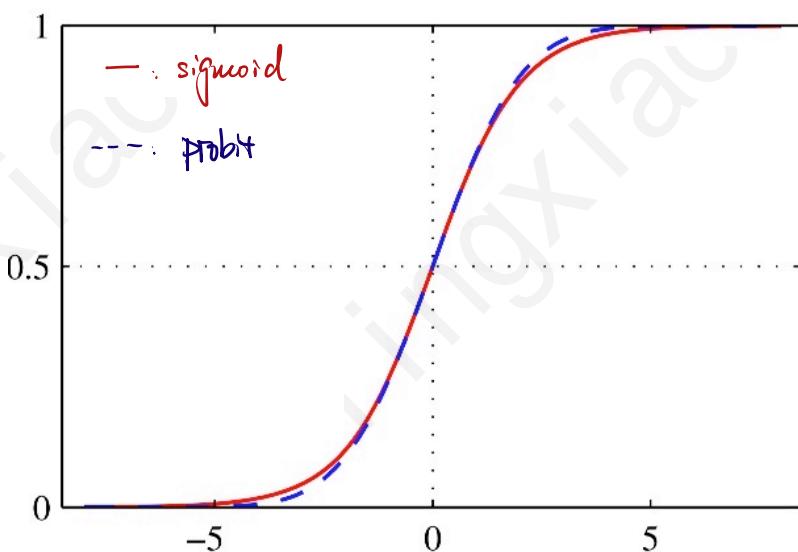
$$\frac{d\Phi}{da} \Big|_{a=0} = \lambda N(0|0,1)$$

$$= \frac{\lambda}{\sqrt{2\pi}}$$

$$\text{let } \frac{\lambda}{\sqrt{2\pi}} = \frac{1}{4}$$

$$\lambda = \frac{\sqrt{2\pi}}{4} \quad \lambda^2 = \frac{\pi}{8}$$

we approximate $\sigma(a)$ by $\Phi(ua)$.



With the probit function, we can show that:

$$\int \Phi(ua) N(u|0, \sigma^2) da = \Phi\left(\frac{u}{\sigma^2 + u^2}\right)$$

Proof:

Statement: Suppose we have a random variable satisfied normal distribution denoted as $X \sim N(\mu, \sigma^2)$.

The probability of $X \leq x$ is $P\{X \leq x\} = \Phi\left(\frac{x-\mu}{\sigma}\right)$

$$\begin{aligned} P\{X \leq x\} &= \int_{-\infty}^x \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\} dx \\ &= \int_{-\infty}^{\frac{x-\mu}{\sigma}} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2}\nu^2\right\} d\nu \\ &= \int_{-\infty}^{\frac{x-\mu}{\sigma}} \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}\nu^2\right\} d\nu \\ &= \Phi\left(\frac{x-\mu}{\sigma}\right) \end{aligned}$$

For two random variables $X \sim N(10, \lambda^2)$ and $Y \sim N(\mu, \sigma^2)$,

$$P\{X \leq Y | Y=a\} = P\{X \leq a\} = \Phi\left(\frac{a-10}{\lambda}\right) = \Phi(a)$$

$$P\{X \leq Y\} = \int_{-\infty}^{+\infty} P\{X \leq Y | Y=a\} P\{Y=a\} da$$

$$= \int_{-\infty}^{+\infty} \Phi(\mu_a) N(a|\mu, \sigma^2) da.$$

$\because x - Y$ also follows a Gaussian

$$\therefore E[x - Y] = E[x] - E[Y] = \theta - \mu = -\mu$$

$$\text{Var}(x - Y) = \text{Var}(x) + \text{Var}(Y) = \lambda^{-2} + \sigma^2$$

$$\therefore P\{x - Y \leq 0\} = \Phi\left\{\frac{\theta - (-\mu)}{\sqrt{\lambda^{-2} + \sigma^2}}\right\} = \Phi\left(\frac{\mu}{\sqrt{\lambda^{-2} + \sigma^2}}\right)$$

Thus, $\int \Phi(\mu_a) N(a|\mu, \sigma^2) da = \Phi\left(\frac{\mu}{\sqrt{\lambda^{-2} + \sigma^2}}\right)$

Hawing approximated the convolution by the probit function, we can approximate the probit again by sigmoid:

$$\int \sigma(u) N(a|\mu, \sigma^2) da \approx \sigma(k(\sigma^2) \mu)$$

$$k(\sigma^2) = \sqrt{1 + \frac{\pi \sigma^2}{8}}$$

Thus, the predictive distribution now becomes:

$$P\{c_i | \phi, t\} = \sigma(k(\sigma^2) \mu)$$