

$$\frac{\partial L}{\partial \omega} = \sum_{i=1}^n -\frac{1}{2\sigma^2} 2(\omega^T x_i - y_i) x_i - 2 \frac{\omega}{\sigma^2}$$

$$\Rightarrow -\frac{1}{\sigma^2} (\omega^T x - y) x^T - \frac{\omega}{\sigma^2} = 0$$

$$\Rightarrow -\frac{\omega^T x^T x}{\sigma^2} + \frac{y x^T}{\sigma^2} - \frac{\omega}{\sigma^2} = 0$$

$$\Rightarrow \omega \left[-\frac{x^T x}{\sigma^2} - \frac{1}{\sigma^2} \right] = -\frac{y x^T}{\sigma^2}$$

$$\Rightarrow \omega \left(\frac{\beta^2 x^T x + \sigma^2}{\sigma^2 \sigma^2} \right) = \frac{y x^T}{\sigma^2}$$

$$\Rightarrow \omega \left(x^T x + \frac{\sigma^2}{\sigma^2} \right) = x^T y$$

$$\Rightarrow \omega_{MAP} = (x^T x + \lambda I)^{-1} x^T y$$

$$\boxed{\lambda = \frac{\sigma^2}{\beta^2}}$$

Bias Variance Decomposition:

$\gamma(x|D)$ \Rightarrow function learnt from training data

$h(x)$ \Rightarrow actual regression function (we don't know this but we can assume this to be an optimal function)

$$E_D \left[\{y(x; D) - h(x)\}^2 \right]$$

$$\Rightarrow E_D \left[y(x; D) - E_D(y(x; D)) + E_D(y(x; D)) - h(x) \right]^2$$

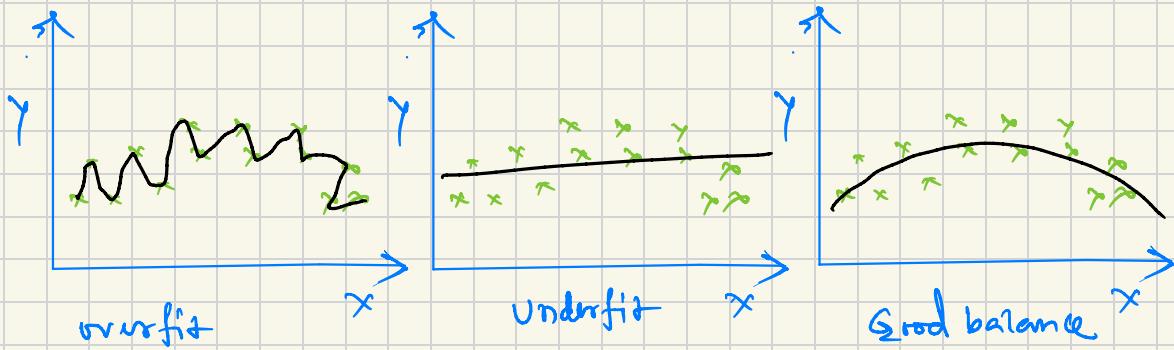
$$\Rightarrow E_D \left[\{E_D(y(x; D)) - h(x)\}^2 + \{y(x; D) - E_D(y(x; D))\}^2 \right. \\ \left. - 2 \{E_D(y(x; D)) - h(x)\} \{y(x; D) - E_D(y(x; D))\} \right]$$

$$\Rightarrow E_D \left[\{E_D(y(x; D)) - h(x)\}^2 \right] + E_D \left[\{y(x; D) - E_D(y(x; D))\}^2 \right] \\ - 2 E_D \left\{ E_D(y(x; D)) - h(x) \right\} E_D \left\{ y(x; D) - E_D(y(x; D)) \right\}$$

this becomes 0.

$$\Rightarrow E_D \left[\{E_D(y(x; D)) - h(x)\}^2 \right] + E_D \left[\{y(x; D) - E_D(y(x; D))\}^2 \right]$$

Bias² variance.



Linear ML algo often have a high bias & a low variance

Nonlinear ML algo often have a low bias but a high variance.

Bias variance trade off?

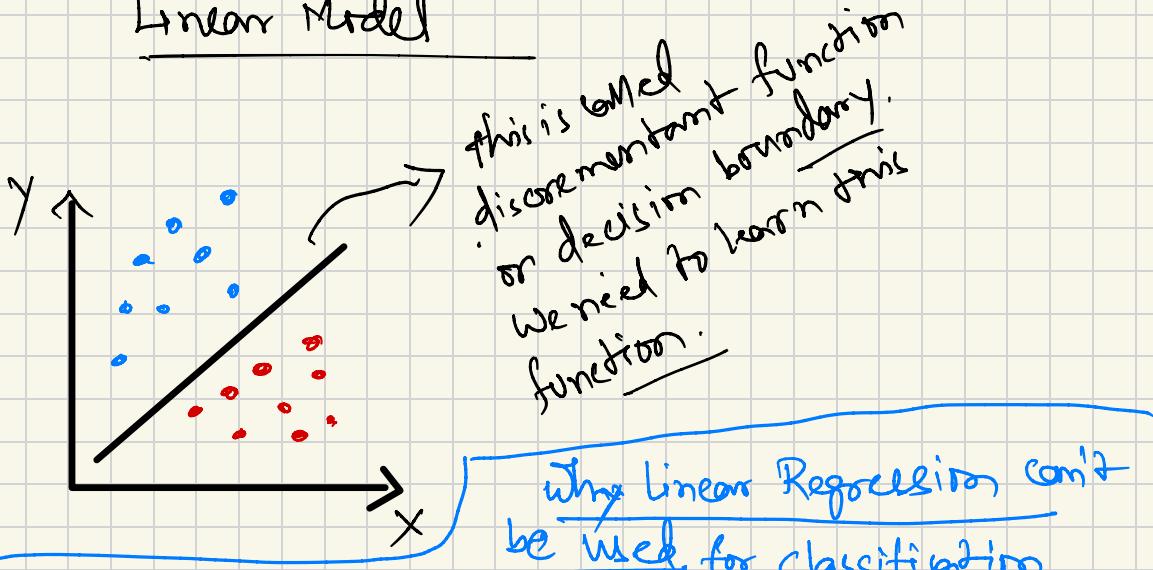
- ① Dimensionality reduction & feature selection can decrease variance by simplifying model.
- ② Similarly a larger training set tends to decrease variance.
- ③ Adding features tends to decrease bias, at the expense of introducing additional variance.

- ④ Linear & generalized linear models can be regularized to decrease their variance at the cost of increasing their bias.

Logistic Regression:

Classification Problem

Linear Model



→ For classification problem using linear regression, the probability is a Normal distribution. So, output $P(y=0|x) + P(y=1|x) \neq 1$.

→ even if you provide binary output to linear regression, it may provide output in real number like 1, 2, -2 etc which doesn't make sense in case of binary classification.

→ To make proper decision we need probabilistic output which sum upto 1 for all classes.

So, use of linear regression for classification problem is not a good option. We need good model which can distribute our probability value within 0 & 1 and for K-class classification within $\{1, K\}$.

our task is to,

$$P(y=1|x) \Rightarrow \mathbb{R}^D \xrightarrow{w^T x} \text{sigmoid function} \rightarrow \{0, 1\}$$

also called probit function.

$$\text{If } P(x) = N(x|0, 1)$$

$$\text{CDF}(x) = \int_{-\infty}^x p(x) dx$$

$$\text{CDF}(-\alpha) = 0$$

$$\begin{aligned} \text{CDF}(0) &= \int_{-\infty}^0 p(x) dx \\ &= 0.5 \end{aligned}$$

$$\text{CDF}(\alpha) = 1.$$

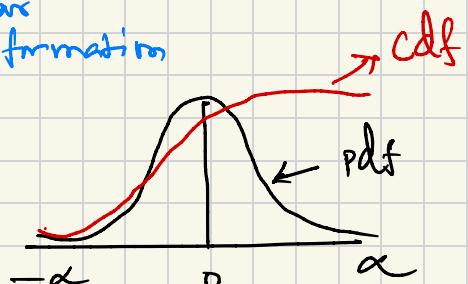
$$\sigma(-\alpha) = 0$$

$$\sigma(0) = 0.5$$

$$\sigma(\alpha) = 1$$

so, sigmoid function we can use to convert $w^T x$ to $\{0, 1\}$

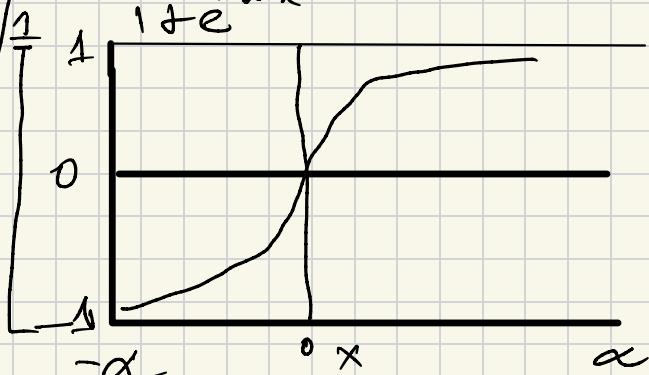
$w^T x$
Linear transformation



Logistic Sigmoid:

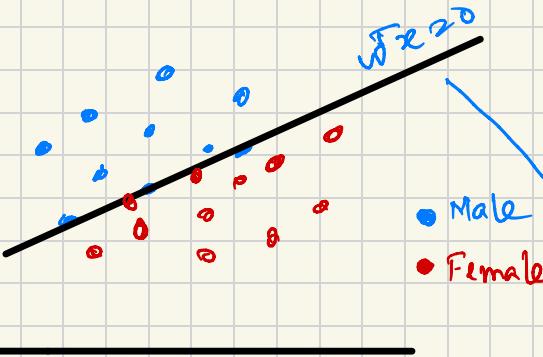
$$P(y=1|x) = \sigma(w^T x)$$

$$= \frac{1}{1 + e^{-w^T x}}$$



$$P(y=1|X) = \sigma(w^T x) = \frac{1}{1+e^{-w^T x}}$$

$$P(y=0|X) = 1 - \sigma(w^T x) = 1 - \frac{1}{1+e^{-w^T x}}$$



$$= \frac{e^{-w^T x}}{1+e^{-w^T x}}$$

$$= \frac{1}{1+e^{w^T x}}$$

→ this is decision boundary with equation
 $w^T x = 0$

$$P(y=1|X) = P(y=0|X) = 0.5$$

points on lies on the line

$w^T x = 0$ have same probability

value = 0.5.

points above $w^T x = 0$ line have high probability towards class 1.

For Male class & points lies below

$w^T x = 0$ have high probability towards class 0 or Female class.

This probability distribution over the output is good for decision making. But sometime we don't need probability to find confidence towards a positive class. In such

scenario we can go with approaches which just provide decision boundary. Ex - LDA, FLD, SVM..

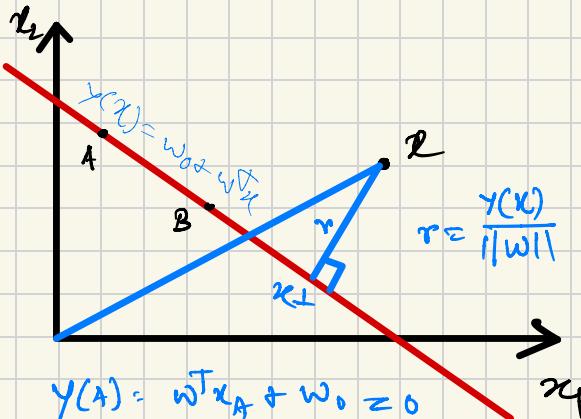
→ All they learn the decision boundary which is called Discriminant function.

Discriminant function:

$$y(x) = w_0 + w^T x$$

w_0 denotes how far your decision boundary from origin. It's a perpendicular direction towards decision boundary from origin.

w is orthogonal to every vector lying within the decision surface, so w determines the orientation of the decision surface.



$$y(A) = w^T x_A + w_0 = 0$$

$$y(B) = w^T x_B + w_0 = 0$$

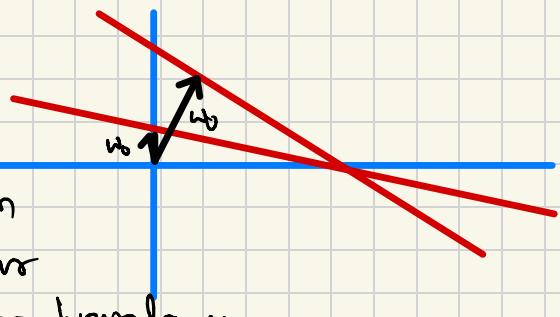
Both are on the line.

so,

$$w^T (x_A - x_B) = 0$$

product of two vectors is 0

so they are perpendicular to each other. so, w is orthogonal to every vector lying within the decision surface.



Notes:

Simplistic discriminant function $y \geq w_0 + w^T x$ represent a hyperplane in $(m-1)$ dimensional space. where m is # of features.

# features	Discriminant function
1	point
2	line
3	plane
4	Hyperplane in 3D space
...	...
m	Hyperplane in $(m-1)$ D space

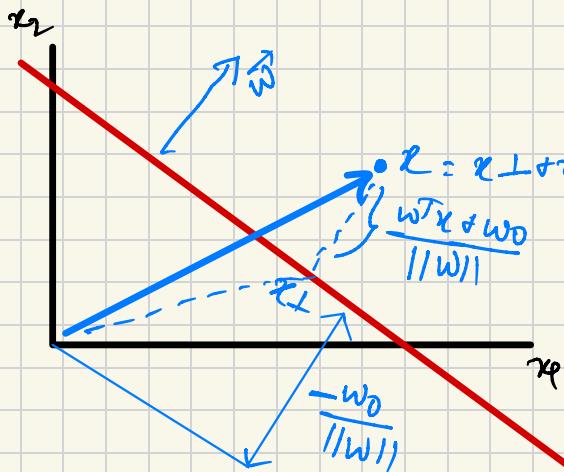
We have $w_0 + w^T x = 0$

$$w^T x = -w_0$$

Normalizing both sides with the length of the vector $\|w\|$, we get normal distance from the origin to the decision surface.

$$d = \frac{w^T x}{\|w\|} = -\frac{w_0}{\|w\|}$$

So, w_0 shifts the boundary away from origin.



this is unit vector for direction.

$$y(x) = w^T x + w_0$$

replacing x with $x + r \frac{w}{\|w\|}$

$$y(x) \geq \boxed{w^T x +} + r \frac{w^T w}{\|w\|} + w_0$$

both addition is 0 bcz they are on same line/plane.

$$y(x) = r \frac{w^T w}{\|w\|} = r \frac{\|w\|^2}{\|w\|}$$

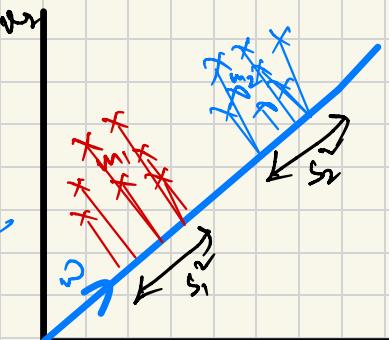
$$y(x) = r \|w\|$$

$$\text{So, } \boxed{r \geq \frac{y(x)}{\|w\|}}$$

$y(x)$ determines the distance of a point to surface.

Supervised Linear Discriminant Analysis (LDA)

- w as the direction to project x .
- Find w such that when x is projected, classes are well separated.
- Simplest measure of the separation of the classes, when projected onto w , is the separation of the projected class means.



$$m_1 = \frac{1}{N_1} \sum_{n \in C_1} x_n \quad m_2 = \frac{1}{N_2} \sum_{n \in C_2} x_n$$

we need to maximize $m_2 - m_1 = w^T(m_2 - m_1)$

this expression can be made arbitrarily large simply by increasing the magnitude of w . To solve this problem we could constrain w to have unit length, so that $\sum w_i^2 = 1$

$$L = \underset{w}{\operatorname{argmax}} w^T(m_2 - m_1) - \lambda \|w\|^2$$

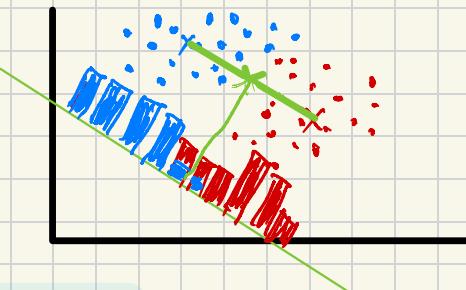
$$\nabla L = (m_2 - m_1) - 2\lambda w = 0$$

$$\Rightarrow w = \frac{1}{2\lambda} (m_1 - m_2)$$

$$\text{so, } w \propto (m_1 - m_2)$$

value if w is still depending on mean of two classes.

therefore, we can see there are still considerable overlap when projected onto the line joining their means. To eliminate this, Fisher proposes a function which also considers variance of both class points.



The within class variance of the transformed data from class C_k is given by

$$S_k^2 = \sum_{n \in C_k} (y_n - m_k)^2 \quad \text{where } y_n = w^T x_n$$

Fisher criterion is defined to be the ratio of the between class variance to the within class variance,

$$J(w) = \frac{(m_2 - m_1)^2}{S_1^2 + S_2^2}$$

$$\begin{aligned} & (m_2 - m_1)^2 \\ &= (w^T m_2 - w^T m_1)^2 \\ &= w^T [(m_2 - m_1)(m_1 - m_2)]^T w \\ &= w^T S_B w \end{aligned}$$

$$\begin{aligned} S_1^2 &= \sum_{n \in C_1} (w^T x_n - m_1)^2 \\ &= \sum_{n \in C_1} w^T [(x_n - m_1)(x_n - m_1)^T]^T w \\ &= w^T S_W w \\ S_2^2 &= w^T S_W w \quad \text{where } S_W = S_1 + S_2 \end{aligned}$$

$$S_1, J(w) = \frac{w^T S_B w}{w^T S_W w}$$

$$\frac{\partial J(w)}{\partial w} = \frac{(w^T S_W w) 2S_B w - (w^T S_B w) 2S_W w}{(w^T S_W w)^2} \geq 0$$

$$\Rightarrow S_B w - J(w) S_W w = 0 \quad \Rightarrow S_B^{-1} S_B w - S_W^{-1} S_W w = 0$$

scalar

so we can drop this

$$\Rightarrow w = S_W^{-1} S_B w$$

$$S_W^{-1} w \propto S_B^{-1} (m_1 - m_2)$$

Multiclass classification:

given $\{(x_i, y_i)\}_{i=1}^N$ $x_i \in \mathbb{R}^n$ $y \in \{0, 1\}^K$

So, Probability that $y \in K$ given $x \in \mathbb{R}^n$

$$\frac{P(y=c_k | x, w)}{\text{Softmax function}} = \sigma \left(\begin{bmatrix} w_1^T x \\ w_2^T x \\ \vdots \\ w_K^T x \end{bmatrix} \right)_k = \frac{e^{w_K^T x}}{\sum_i e^{w_i^T x}}$$

\downarrow

$$\begin{bmatrix} -w_1- \\ -w_2- \\ \vdots \\ -w_K- \end{bmatrix} \begin{bmatrix} x \end{bmatrix}$$

Now we need to find w_1, w_2, \dots, w_K given x, y

Likelihood

$$P(y | x, w) = \prod_{n=1}^N \prod_{k=1}^K \left(\frac{e^{w_k^T x_n}}{\sum_i e^{w_i^T x_n}} \right)^{y_{nk}} \begin{array}{l} \text{if } y_n \in K \\ \text{otherwise} \\ \text{it will be 0} \end{array}$$

Now we need to minimize negative log likelihood

$$L(w) = - \sum_{n=1}^N \sum_{k=1}^K y_{nk} \ln \left[e^{w_k^T x_n} - \ln \left(\sum_{i=1}^K e^{w_i^T x_n} \right) \right]$$

$$L(w) = - \sum_{n=1}^N \sum_{k=1}^K y_{nk} w_k^T x_n + \sum_{n=1}^N \boxed{\sum_{k=1}^K y_{nk}} \ln \left(\sum_{i=1}^K e^{w_i^T x_n} \right)$$

Now take derivative wrt w

→ this will be 1. ✓

$$\frac{\partial L(w)}{\partial w_j} = -\sum_{n=1}^N y_{nj} x_n + \sum_{n=1}^N \frac{1}{\sum_i e^{w_i^T x_n}} e^{w_j^T x_n} \cdot x_n$$

$$= \sum_{n=1}^N \left(\boxed{\frac{e^{w_j^T x_n}}{\sum_i e^{w_i^T x_n}} - y_{nj}} \right) x_n$$

$\hookrightarrow P(y \in c_j | w, x)$

$$= \sum_{n=1}^N (\overbrace{y_{nj} - \hat{y}_{nj}}^{\text{diff between predicted value}}) x_n$$

& true value.

Now we can apply gradient descent to get optimal value of w .

GD:

$$\text{init } \{w\}$$

for itr in iterations:

for every class k :

$$w^{k+1} = w^k - \eta \nabla_{w_k} l(w)$$

After training

$$\begin{aligned} \text{label}(x_{new}) &= \underset{k}{\operatorname{argmax}} \left[\sigma(w^{\text{final}} \cdot x_{new}) \right] \\ &= \underset{k}{\operatorname{argmax}} \left(\begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_K \end{bmatrix} \begin{bmatrix} x_{new} \end{bmatrix}_{K \times M} \right)_{M \times 1} \\ &= \underset{k}{\operatorname{argmax}} (\text{K} \times 1 \text{ matrix}) \end{aligned}$$

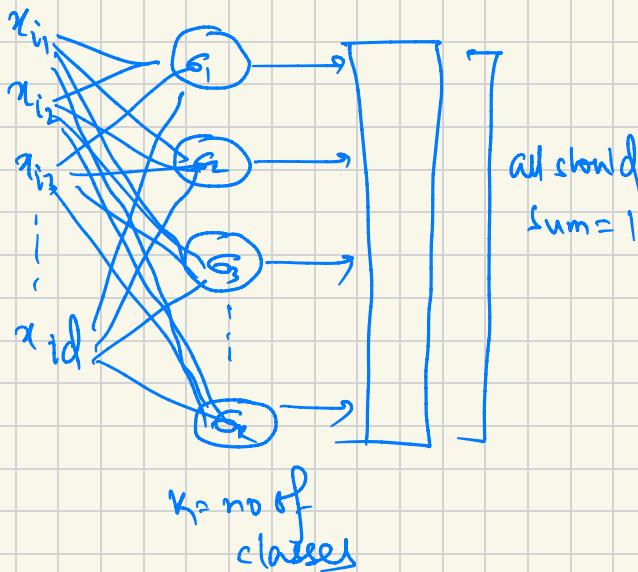
Multi Class Classification

$$P(y=c | x, w) = \frac{\exp(w_c^T x)}{\sum_{c=1}^C \exp(w_c^T x)}$$

log likelihood

$$L(w) = \log \prod_{i=1}^N \prod_{c=1}^C y_{ic} \log p_{ic}$$

this is cross entropy error function for multi class classification.



$$\sigma_i(z_i) = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$$

where $z_i = \sum_{j=1}^d w_{ij} x_{ij}$

probabilistic Generative Models

$$P(G_1 | x) = \frac{P(x|G_1) P(G_1)}{P(x|G_1) P(G_1) + P(x|G_2) P(G_2)}$$

$$= \frac{1}{1 + \exp(-\alpha)} \quad \text{where } \alpha = \ln \frac{P(x|G_1) P(G_1)}{P(x|G_2) P(G_2)}$$

$$= \sigma(\alpha)$$

$$P(G_k | x) = \frac{P(x|G_k) P(G_k)}{\sum_{k=1}^K P(x|G_k) P(G_k)}$$

This is for binary classification.
for multiclass classification
we have below form
of the logit equation.

$$P(a_k | x) \Rightarrow \frac{\exp(\alpha_k)}{\sum_j \exp(\alpha_j)}$$

$$\alpha_k = \ln P(x|G_k) P(G_k)$$

→ similar to softmax function we have seen in multiclass classification where we have taken $w^T x_i$ inside exp function.

$P(x|y)$ is class conditioned distribution.

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

Property:

$$\sigma(-a) = 1 - \sigma(a)$$

$$\text{Inverse } a = \ln\left(\frac{\sigma}{1-\sigma}\right)$$

if $\sigma(a) = p(G|x)$

$$a = \ln\left(\frac{p(G|x)}{1 - p(G|x)}\right)$$

$$a = \ln\left(\frac{p(G|x)}{p(G|x)}\right)$$

This is also known as logit function.

It represents the log of the ratio of probability of two classes. It also

known as log odds.

*#

let's say given logistic regression classifier

$$p = \sigma(x, y) = \frac{1}{1 + \exp(-w_0 - 7.5x - 7.5y)}$$

$$\text{then } -w_0 - 7.5x - 7.5y = \log\left(\frac{p}{1-p}\right)$$

this is inverse of logit.

$$p(x|c_k) \sim N(x; \mu_k, \Sigma_k)$$

$x \in \mathbb{R}^d$, $\mu = d \times 1$ mean vector

$\Sigma = \text{covariance matrix}$
($d \times d$) $|\Sigma| = \text{determinant}$

$$\frac{1}{\sqrt{(2\pi)^d |\Sigma|}} e^{-\frac{(x-\mu)^T \Sigma^{-1} (x-\mu)}{2}}$$

$$p(x|\mu, \Sigma_k) = \frac{1}{\sqrt{2\pi)^d |\Sigma|}} \exp \left(-\frac{1}{2} (x-\mu)^T \Sigma_k^{-1} (x-\mu) \right)$$

$$P(G|x) = \sigma(w^T x + w_0)$$

$$w_2 \in \Sigma^{-1}(\mu_1 - \mu_2)$$

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

$$P(G|x) = \sigma \left(\ln \frac{P(x|G) P(G)}{P(x|\bar{G}) P(\bar{G})} \right) =$$

$$\ln \frac{\frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp \left\{ -\frac{1}{2} (x-\mu_1)^T \Sigma^{-1} (x-\mu_1) \right\}}{\frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp \left\{ -\frac{1}{2} (x-\mu_2)^T \Sigma^{-1} (x-\mu_2) \right\}} + \ln \frac{P(G)}{P(\bar{G})}$$

<https://regenerativetoday.com/univariate-and-bivariate-gaussian-distribution-clear-explanation-with-visuals/>

https://mmuratarat.github.io/2019-10-05/univariate-multivariate_gaussian#:~:text=Like%20the%20univariate%20normal%20distribution,they%20change%20together.

$$\ln \frac{\frac{1}{\sqrt{2\pi}\sigma} \cdot \frac{1}{\sqrt{|\Sigma|}} \exp \left\{ -\frac{1}{2} (x - \mu_1)^T \Sigma^{-1} (x - \mu_1) \right\}}{\frac{1}{\sqrt{2\pi}\sigma} \cdot \frac{1}{\sqrt{|\Sigma|}} \exp \left\{ -\frac{1}{2} (x - \mu_2)^T \Sigma^{-1} (x - \mu_2) \right\}} + \ln \frac{p(y)}{p(x)}$$

$$\frac{x^T \Sigma^{-1} x}{2} = q.$$

$$\Rightarrow \ln \frac{\exp \left\{ -\frac{1}{2} (x^T - \mu_1^T) \Sigma^{-1} (x - \mu_1) \right\}}{\exp \left\{ -\frac{1}{2} (x^T - \mu_2^T) \Sigma^{-1} (x - \mu_2) \right\}}$$

$$\Rightarrow \ln \frac{\exp \left\{ -\frac{1}{2} [x^T \Sigma^{-1} x - \mu_1^T \Sigma^{-1} x - x^T \Sigma^{-1} \mu_1 + \mu_1^T \Sigma^{-1} \mu_1] \right\}}{\exp \left\{ -\frac{1}{2} [x^T \Sigma^{-1} x - \mu_2^T \Sigma^{-1} x - x^T \Sigma^{-1} \mu_2 + \mu_2^T \Sigma^{-1} \mu_2] \right\}}$$

$$\Rightarrow \ln \exp \left\{ \frac{1}{2} \left[\cancel{x^T \Sigma^{-1} x} + \cancel{\mu_1^T \Sigma^{-1} x} + x^T \Sigma^{-1} \mu_1 - \mu_1^T \Sigma^{-1} \mu_1 \right. \right. \\ \left. \left. - \cancel{x^T \Sigma^{-1} x} - \cancel{\mu_2^T \Sigma^{-1} x} - x^T \Sigma^{-1} \mu_2 + \mu_2^T \Sigma^{-1} \mu_2 \right] \right\}$$

$$\Rightarrow \ln \exp \left\{ \frac{1}{2} \left[(\mu_1^T - \mu_2^T) \Sigma^{-1} x + (\mu_1 - \mu_2) \Sigma^{-1} x^T \right. \right. \\ \left. \left. - (\mu_1^T \Sigma^{-1} \mu_1 - \mu_2^T \Sigma^{-1} \mu_2) \right] \right\}$$

$$\Rightarrow \frac{1}{2} \left(\underset{\partial x^T}{\cancel{\Sigma^{-1}(\mu_1 - \mu_2)^T x}} + \underset{\partial x^T}{\cancel{\Sigma^{-1}(\mu_1 - \mu_2) x^T}} \right. \\ \left. - \mu_1^T \Sigma^{-1} \mu_1 + \mu_2^T \Sigma^{-1} \mu_2 \right)$$

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

$$w_0 = -\frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2$$

For multiclass classification

$$P(c_k|x) \Rightarrow \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

$$a_k = \ln P(x|c_k) P(c_k)$$

$$a_k(x) = w_k^T x + w_{k0}$$

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

$$w_{k0} = -\frac{1}{2} u_k^T \Sigma^{-1} u_k + \ln P(c_k)$$

If covariance is not shared

Quadratic terms of x will be present.

It is also called quadratic discriminant analysis

$$P(y=c|x, w) = \frac{\frac{f_c}{\sqrt{2\pi\Sigma_c}} \exp\left[-\frac{1}{2}(x-u_c)^T \Sigma_c^{-1} (x-u_c)\right]}{\sum_c \frac{f_c}{\sqrt{2\pi\Sigma_c}} \exp\left[-\frac{1}{2}(x-u_c')^T \Sigma_c'^{-1} (x-u_c')\right]}$$

parameter estimation technique

$$P(x_n, q) = P(q) P(x_n|q) = \pi N(x_n | \mu_1, \Sigma)$$

$$P(x_n, q) = P(q) P(x_n|q) = (1-q) N(x_n | \mu_2, \Sigma)$$

$$P(y|q, \mu_1, \mu_2, \Sigma) = \prod_{n=1}^N \left[P N(x_n | \mu_1, \Sigma) \right]^{y_n} \left[(1-p) N(x_n | \mu_2, \Sigma) \right]^{(1-y_n)}$$

It is convenient to maximize the log of the likelihood function.

$$\ln P(y|q, \mu_1, \mu_2, \Sigma) = \sum_{n=1}^N \ln \left[P N(x_n | \mu_1, \Sigma) \right]^{y_n} + \ln \left[(1-p) N(x_n | \mu_2, \Sigma) \right]^{(1-y_n)}$$

$$\Rightarrow \sum_{n=1}^N y_n \ln p \left[-\frac{1}{2} (x_n - \mu_1)^T \Sigma^{-1} (x_n - \mu_1) \right]$$

$$(1-y_n) \ln(1-p) \left[-\frac{1}{2} (x_n - \mu_2)^T \Sigma^{-1} (x_n - \mu_2) \right]$$

$$L \rightarrow \sum_{n=1}^N y_n \ln p + (1-y_n) \ln(1-p)$$

$$\frac{\partial L}{\partial p} = \frac{y_n}{p} - \frac{(1-y_n)}{1-p} = 0$$

$$y_n - \cancel{y_n/p} - p + \cancel{y_n/p} = 0$$

$$p = \sum_{n=1}^N y_n$$

$$\frac{\partial L}{\partial \mu_i} \Rightarrow \sum_{n=1}^N y_n \left[-\frac{1}{2} (x_n - \mu_i)^T \Sigma^{-1} (x_n - \mu_i) \right]$$

$$= \sum_{n=1}^N y_n \left(-\frac{1}{2} \Sigma^{-1} (x_n - \mu_i) \Sigma^{-1} \right)$$

$$= \sum_{n=1}^N y_n \left[+ (x_n - \mu_i) \Sigma^{-1} \right]$$

$$\Rightarrow \sum_{n=1}^N y_n (x_n - \mu_i) \Sigma^{-1} = 0$$

$$\sum_{n=1}^N y_n (x_n - \mu_i) = 0$$

$$\boxed{\mu_i = \frac{1}{N} \left(\sum_{n=1}^N x_n y_n \right)}$$

$$\begin{aligned} \frac{\partial L}{\partial \mu} &= -\frac{1}{2} \left(\frac{\partial (y - \mu)^T \Sigma^{-1} (y - \mu)}{\partial \mu} \right) \\ &= -\frac{1}{2} (-2 \Sigma^{-1} (y - \mu)) \\ &= \Sigma^{-1} (y - \mu) \end{aligned}$$

PUBLIC
Questions
Tags
Users
Unanswered

TEAMS

Stack Overflow for Teams – Start collaborating and sharing organizational knowledge.



Create a free Team
[Why Teams?](#)

83

In chapter 2 of the [Matrix Cookbook](#) there is a nice review of matrix calculus stuff that gives a lot of useful identities that help with problems one would encounter doing probability and statistics, including rules to help differentiate the multivariate Gaussian likelihood.

If you have a random vector \mathbf{y} that is multivariate normal with mean vector $\boldsymbol{\mu}$ and covariance matrix Σ , then use equation (86) in the matrix cookbook to find that the gradient of the log likelihood L with respect to $\boldsymbol{\mu}$ is

$$\begin{aligned} \frac{\partial L}{\partial \boldsymbol{\mu}} &= -\frac{1}{2} \left(\frac{\partial (\mathbf{y} - \boldsymbol{\mu})' \Sigma^{-1} (\mathbf{y} - \boldsymbol{\mu})}{\partial \boldsymbol{\mu}} \right) \\ &= -\frac{1}{2} (-2 \Sigma^{-1} (\mathbf{y} - \boldsymbol{\mu})) \\ &= \Sigma^{-1} (\mathbf{y} - \boldsymbol{\mu}) \end{aligned}$$

I'll leave it to you to differentiate this again and find the answer to be $-\Sigma^{-1}$.

As "extra credit", use equations (57) and (61) to find that the gradient with respect to Σ is

$$\begin{aligned} \frac{\partial L}{\partial \Sigma} &= -\frac{1}{2} \left(\frac{\partial \log(|\Sigma|)}{\partial \Sigma} + \frac{\partial (\mathbf{y} - \boldsymbol{\mu})' \Sigma^{-1} (\mathbf{y} - \boldsymbol{\mu})}{\partial \Sigma} \right) \\ &= -\frac{1}{2} (\Sigma^{-1} - \Sigma^{-1} (\mathbf{y} - \boldsymbol{\mu}) (\mathbf{y} - \boldsymbol{\mu})' \Sigma^{-1}) \end{aligned}$$

I've left out a lot of the steps, but I made this derivation using only the identities found in the matrix cookbook, so I'll leave it to you to fill in the gaps.

I've used these score equations for maximum likelihood estimation, so I know they are correct :)