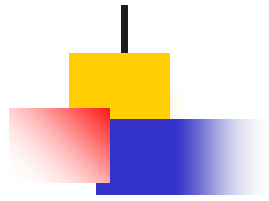


Fundamentals of Machine Learning

LOGISTIC REGRESSION

Amit K Roy-Chowdhury

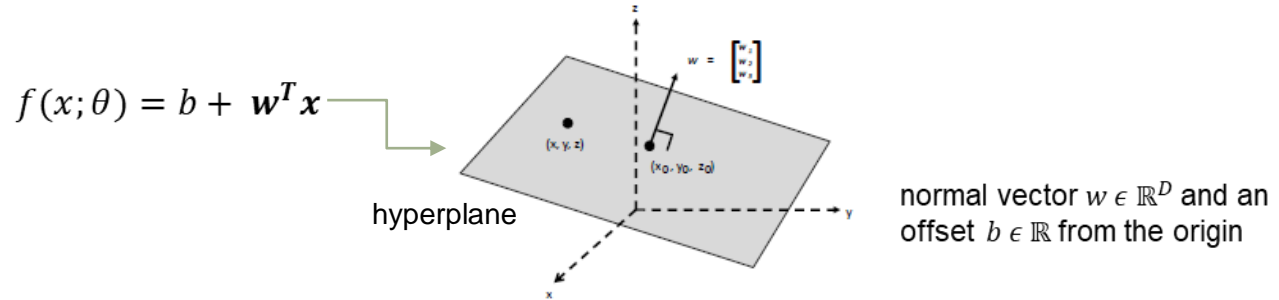
Acknowledgments: Adapted from slides at <https://probml.github.io/pml-book/teaching1.html> by Prof. Saw Shier Nee



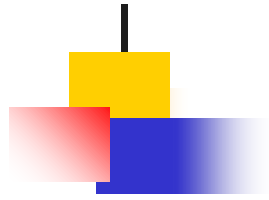
Linear Classifier

The prediction can be written as

$$f(x) = \mathbb{I}(p(y = 1|x) > p(y = 0|x)) = \mathbb{I}\left(\log \frac{p(y = 1|x)}{p(y = 0|x)} > 0\right) = \mathbb{I}(a > 0)$$



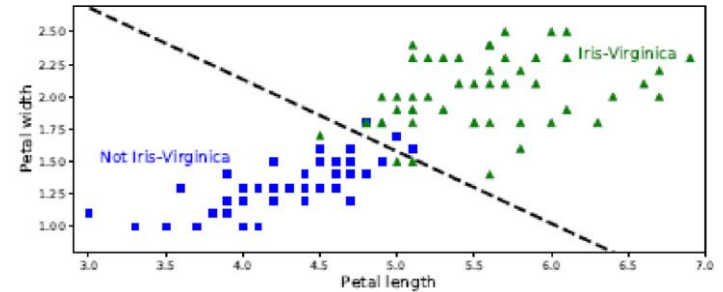
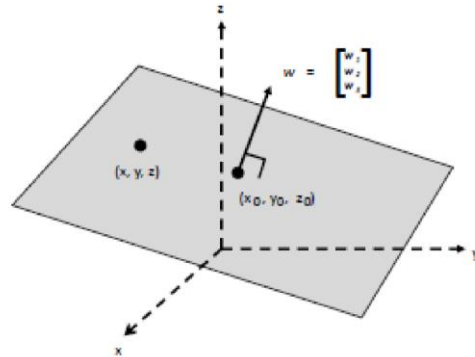
This linear hyperplane (decision boundary) separates 3d space into half spaces.

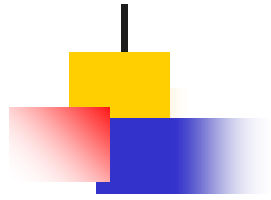


Linear Classifier

$$f(x) = \mathbb{I}(p(y = 1|x) > p(y = 0|x)) = \mathbb{I}\left(\log \frac{p(y = 1|x)}{p(y = 0|x)} > 0\right) = \mathbb{I}(a > 0)$$

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





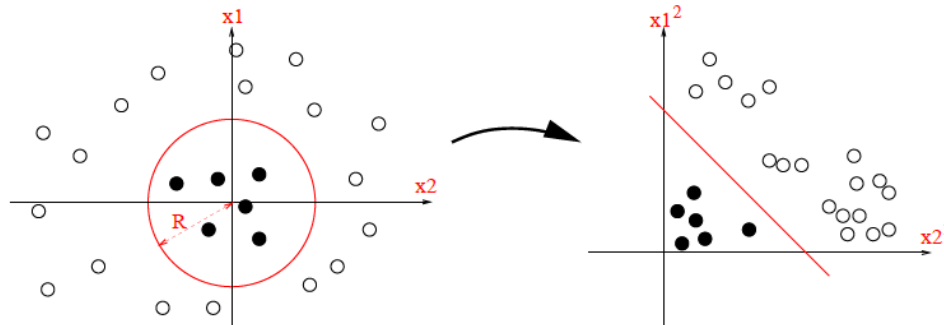
Non Linear Classifier

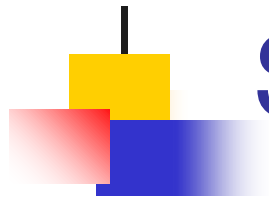
Transform input features in suitable way

$$\phi(x_1, x_2) = [1, x_1^2, x_2^2]$$

$$w = [-R^2, 1, 1]. \text{ Then } w^T \phi(x) = x_1^2 + x_2^2 - R^2$$

Decision boundary (where $f(x) = 0$) defines a circle with radius R





Sigmoid (Logistic) Function

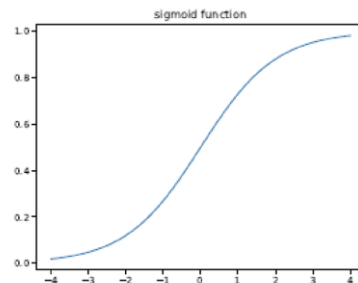
Predict binary random variable y given inputs x .

$$p(y|x, \theta) = \text{Ber}(y|f(x; \theta))$$

$$p(y|x, \theta) = \text{Ber}(y|\sigma(f(x; \theta)))$$

$$p(y = 1|x, \theta) = \frac{1}{1 + e^{-a}} = \frac{e^a}{1 + e^a} = \sigma(a)$$

$$p(y = 0|x, \theta) = 1 - \frac{1}{1 + e^{-a}} = \frac{e^{-a}}{1 + e^{-a}} = \frac{1}{1 + e^a} = \sigma(-a)$$

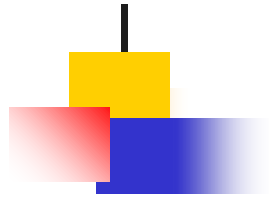


Logistic/Sigmoid function

$$\sigma(a) \triangleq \frac{1}{1 + e^{-a}} = \frac{e^a}{1 + e^a}$$

Logit function

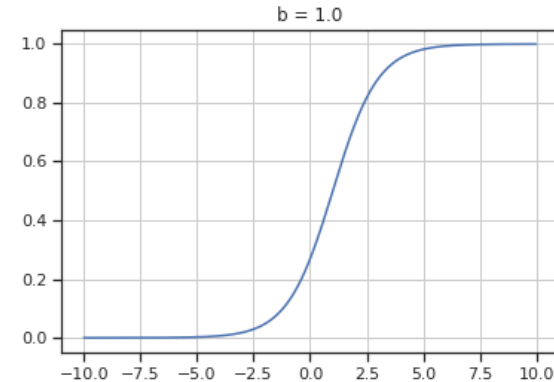
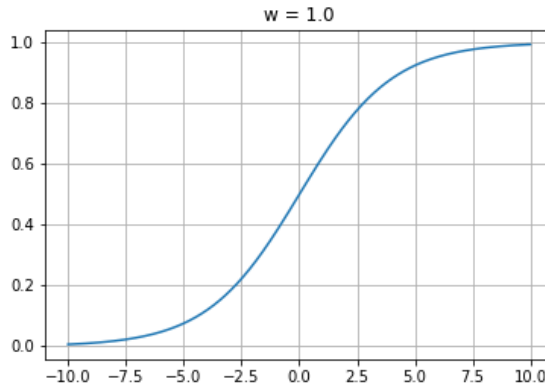
$$a = \text{logit}(p) = \sigma^{-1}(p) \triangleq \log\left(\frac{p}{1-p}\right)$$

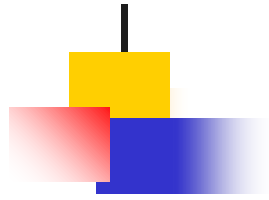


Sigmoid Function

Sigmoid function

$$\sigma(a) = \frac{1}{1 + e^{-a}}, \text{ where } a = \log \frac{p}{1-p} = b + \mathbf{w}^T \mathbf{x}$$





Logistic Regression

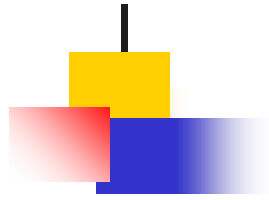
Binary logistic regression often follows the following model

$$p(y|x; \theta) = \text{Ber}(y | \sigma(w^T x + b))$$

Bernoulli sigmoid weight bias

$$p(y = 1|x; \theta) = \sigma(a) = \frac{1}{1 + e^{-a}}, \text{ where } a = \log \frac{p}{1 - p}$$

$$f(k; p) = \begin{cases} p & \text{if } k = 1, \\ q = 1 - p & \text{if } k = 0. \end{cases}$$



Logistic Regression – Cost Function

Maximize Maximum Likelihood Estimation / Minimize Negative Log Likelihood

$$\text{NLL}(w) = -\frac{1}{N} \log p(\mathcal{D}|w) = -\frac{1}{N} \log \prod_{n=1}^N \text{Ber}(y_n|\mu_n)$$

$$= -\frac{1}{N} \sum_{n=1}^N \log[\mu_n^{y_n} \times (1 - \mu_n)^{1-y_n}]$$

$$= -\frac{1}{N} \sum_{n=1}^N [y_n \log \mu_n + (1 - y_n) \log(1 - \mu_n)]$$

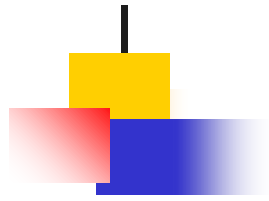
No of Samples \rightarrow $= \frac{1}{N} \sum_{n=1}^N \mathbb{H}(y_n, \mu_n)$ \leftarrow probability

$$\mathbb{H}(p, q) = -[p \log q + (1 - p) \log(1 - q)]$$

Binary cross-entropy

$\mu_n = \sigma(a_n)$ is the probability of class 1

$a_n = w^T x_n$ is the **logit**



Logistic Regression – Cost Function

To find the MLE, we must solve

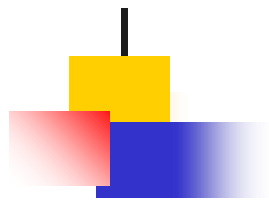
$$\longrightarrow \nabla_{\mathbf{w}} \text{NLL}(\mathbf{w}) = \mathbf{g}(\mathbf{w}) = \mathbf{0}$$

Use chain rule to work out derivatives

$$\begin{aligned}\nabla_{\mathbf{w}} \text{NLL}(\mathbf{w}) &= -\frac{1}{N} \sum_{n=1}^N [y_n(1 - \mu_n)x_n - (1 - y_n)\mu_n x_n] \\ &= -\frac{1}{N} \sum_{n=1}^N [y_n x_n - y_n x_n \mu_n - x_n \mu_n + y_n x_n \mu_n] \\ &= \frac{1}{N} \sum_{n=1}^N (\mu_n - y_n) x_n\end{aligned}$$

Error

Here, we can see that the gradient is weighed by the error for each input



Calculating Derivatives

$$a_n = w^T x_n \text{ and } \mu_n = \sigma(a_n)$$

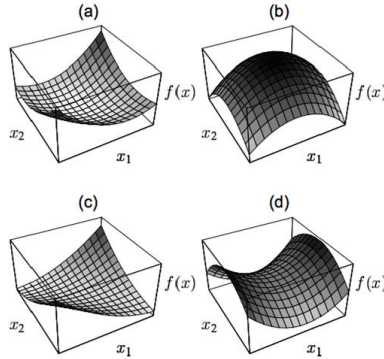
$$\frac{d\mu_n}{da_n} = \sigma(a_n)(1 - \sigma(a_n)) \quad (\text{try! nice exercise in taking derivatives})$$

$$\frac{\partial}{\partial w_d} \mu_n = \frac{\partial}{\partial w_d} \sigma(w^T x_n) = \frac{\partial}{\partial a_n} \sigma(a_n) \frac{\partial a_n}{\partial w_d} = \mu_n (1 - \mu_n) x_{nd}$$

$$\nabla_w \log(\mu_n) = \frac{1}{\mu_n} \nabla_w \mu_n = (1 - \mu_n) x_n$$

$$\nabla_w \log(1 - \mu_n) = \frac{-\mu_n (1 - \mu_n) x_n}{1 - \mu_n} = -\mu_n x_n$$

Logistic Regression – Cost Function



Ensure NLL has bowl shape (global minimum)
check Hessian matrix

$$\mathbf{H}(w) = \nabla_w \nabla_w^T \text{NLL}(w) = \frac{1}{N} \sum_{n=1}^N (\mu_n(1 - \mu_n) x_n) x_n^T = \frac{1}{N} \mathbf{X}^T \mathbf{S} \mathbf{X}$$

$$\mathbf{S} \triangleq \text{diag}(\mu_1(1 - \mu_1), \dots, \mu_N(1 - \mu_N))$$

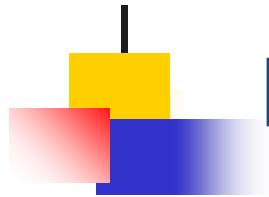
Check convexity

We see that \mathbf{H} is positive definite, since for any nonzero vector v , we have

$$v^T \mathbf{X}^T \mathbf{S} \mathbf{X} v = (v^T \mathbf{X}^T \mathbf{S}^{\frac{1}{2}})(\mathbf{S}^{\frac{1}{2}} \mathbf{X} v) = \|v^T \mathbf{X}^T \mathbf{S}^{\frac{1}{2}}\|_2^2 > 0$$

$$\mathbf{H}_f = \frac{\partial^2 f}{\partial x^2} = \nabla^2 f = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & & \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{pmatrix}$$

This follows since $\mu_n > 0$ for all n , because of the use of the sigmoid function. Consequently the NLL is strictly convex. However, in practice, values of μ_n which are close to 0 or 1 might cause the Hessian to be close to singular. We can avoid this by using ℓ_2 regularization, as we discuss in Section 10.2.7.



Logistic Regression – Cost Function

Our goal is to solve the following optimization problem

$$\hat{w} \triangleq \underset{w}{\operatorname{argmin}} \mathcal{L}(w)$$

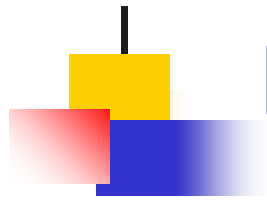
where $\mathcal{L}(w)$ is the loss function, in this case the negative log likelihood:

$$\text{NLL}(w) = -\frac{1}{N} \sum_{n=1}^N [y_n \log \mu_n + (1 - y_n) \log(1 - \mu_n)]$$

where $\mu_n = \sigma(a_n)$ is the probability of class 1, and $a_n = w^T x_n$ is the log odds.

Stochastic gradient descent with minibatch of size 1

$$w_{t+1} = w_t - \eta_t \nabla_w \text{NLL}(w_t) = w_t - \eta_t (\mu_n - y_n) x_n$$



Logistic Regression – Optimizer

1. First order method

- Stochastic Gradient Descent

Slow convergence, when gradient is small

$$\begin{aligned}w_{t+1} &= w_t - \mathbf{H}^{-1}g_t \\&= w_t + (\mathbf{X}^T \mathbf{S}_t \mathbf{X})^{-1} \mathbf{X}^T (y - \mu_t) \\&= (\mathbf{X}^T \mathbf{S}_t \mathbf{X})^{-1} [(\mathbf{X}^T \mathbf{S}_t \mathbf{X})w_t + \mathbf{X}^T (y - \mu_t)] \\&= (\mathbf{X}^T \mathbf{S}_t \mathbf{X})^{-1} \mathbf{X}^T [\mathbf{S}_t \mathbf{X} w_t + y - \mu_t] \\&= (\mathbf{X}^T \mathbf{S}_t \mathbf{X})^{-1} \mathbf{X}^T \mathbf{S}_t z_t\end{aligned}$$

2. Second order method

- Newton Method (Iteratively reweighted least squares)

$$w_{t+1} = w_t - \eta_t \mathbf{H}_t^{-1} g_t$$

where $\mathbf{H}_t \triangleq \nabla^2 \mathcal{L}(w)|_{w_t} = \nabla^2 \mathcal{L}(w_t) = \mathbf{H}(w_t)$

$$z_t \triangleq \mathbf{X} w_t + \mathbf{S}_t^{-1} (y - \mu_t)$$



Logistic Regression – Optimizer

$$w_{t+1} = w_t - \eta_t H_t^{-1} g_t$$

where $H_t \triangleq \nabla^2 \mathcal{L}(w)|_{w_t} = \nabla^2 \mathcal{L}(w_t) = H(w_t)$

$$\begin{aligned} w_{t+1} &= w_t - H^{-1} g_t \\ &= w_t + (X^T S_t X)^{-1} X^T (y - \mu_t) \\ &= (X^T S_t X)^{-1} [(X^T S_t X) w_t + X^T (y - \mu_t)] \\ &= (X^T S_t X)^{-1} X^T [S_t X w_t + y - \mu_t] \\ &= (X^T S_t X)^{-1} X^T S_t z_t \end{aligned}$$

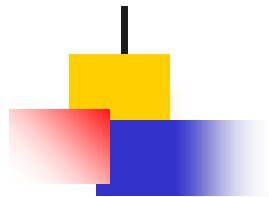
where $z_t \triangleq X w_t + S_t^{-1} (y - \mu_t)$

This is an iteratively weighted least squares problem:
$$\sum_{n=1}^N S_{t,n} (z_{t,n} - w_t^T x_n)^2$$

Least Squares

$$\begin{aligned} J(\theta) &= \sum_{n=0}^{N-1} (x[n] - s[n; \theta])^2 \\ &= (x - H\theta)^T (x - H\theta) \end{aligned}$$

$$\hat{\theta}_{LS} = (H^T H)^{-1} H^T x$$



Multinomial Logistic Regression

Binary Logistic Regression

Probability

$$p(y|x; \theta) = \text{Ber}(y|\sigma(w^\top x + b))$$

Activation function

σ = sigmoid activation

$$p(y = 1|x; \theta) = \sigma(a) = \frac{1}{1 + e^{-a}}$$

Cost function

$$\frac{1}{N} \sum_{n=1}^N \mathbb{H}_{ce}(y_n, \mu_n)$$

$$\mathbb{H}_{ce}(p, q) = -[p \log q + (1 - p) \log(1 - q)]$$

Multiple Logistic Regression

$$p(y|x; \theta) = \prod_{c=1}^C \text{Ber}(y_c|\sigma(w_c^\top x))$$

σ = softmax activation

$$\text{softmax}(a) \triangleq \left[\frac{e^{a_1}}{\sum_{c'=1}^C e^{a_{c'}}}, \dots, \frac{e^{a_C}}{\sum_{c'=1}^C e^{a_{c'}}} \right]$$

$$\frac{1}{N} \sum_{n=1}^N \mathbb{H}_{ce}(y_n, \mu_n)$$

$$\mathbb{H}_{ce}(p, q) = - \sum_{c=1}^C p_c \log q_c$$