# Logistic regression (classification)

## Example:

Suppose you're an actuary and want to predict that a given patient may have some major health issue in the next 5 years.
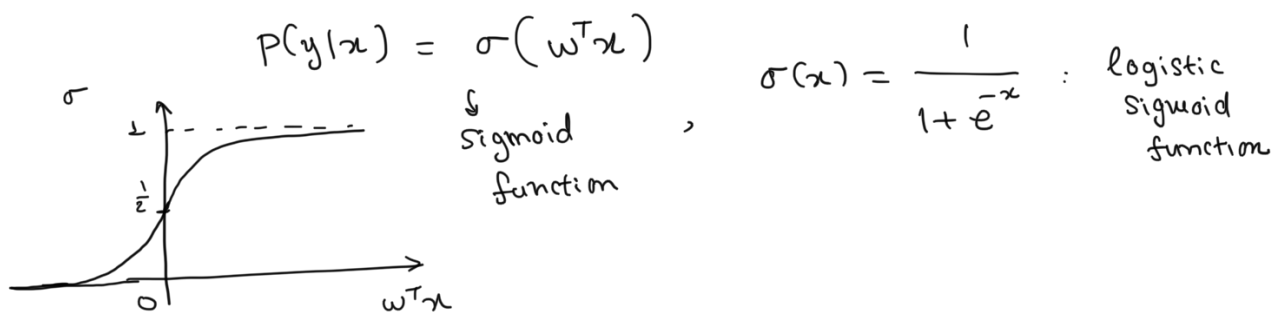
i.e. $P(\text{major health issue} \mid x)$, $x = (x_1, x_2, x_3)$

$$x_1 = \text{age}, \quad x_2 = M/F, \quad x_3 = \text{cholesterol levels}$$

The simplest model would be to consider a linear combination of the input variables:

$$\Big\{ \; w_0 + w_1 x_1 + w_2 x_2 + \dots + w_d x_d = w^T x, \quad x = (1, x_1, x_2, \dots$$

$\hookrightarrow$ this will not yield a probability

We can fix this by introducing a simple warping transformati

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

$\underset{\text{sigmoid function}}{\downarrow}$ , $\sigma(x) = \dfrac{1}{1 + e^{-x}}$ : logistic sigmoid function



## Workflow:

$\begin{cases} 1. \text{ Prior/model specification} \longrightarrow \text{Likelihood} \\ 2. \text{ Training} \\ 3. \text{ Prediction} \end{cases}$

## Setup:

Given $\mathcal{D} := \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$, $x_i \in \mathbb{R}^d$, $y_i \in \{$

$$i = 1, \dots, n$$

## Model:

$y_i \sim \text{Ber}\left(\sigma(w^T x_i)\right)$   $y_i$ are i.i.d.

**Pros :**
- <u>interpetable</u> : the model parameters $w$ have a meaning.

    e.g. $\underline{w_1 >> 0}$ then probability of disease increases with age

- it reveals which input features are most influention.

- small number of parameters ( $d+1$ )

- computationally efficient ways to estimate $w$.

- easy extension to multi-class classification

**Cons :**

Being a simple model, its performance is often inferior to more sophisticated models.

## Maximum Likehood Estimation :

$$W_{MLE} = \arg\max_{w} p(\mathcal{D}|w) \quad , \quad p(\mathcal{D}|w) =$$

$$= p(y_1, y_2, ..., y_n | x_1, x_2, ..., x_n, w) =$$

$$= \prod_{i=1}^{n} p(y_i | x_i, w)$$

Let $a_i = \sigma(w^T x_i)$  
predicted → class probability

then, $$p(\mathcal{D}|w) = \prod_{i=1}^{n} \underbrace{a_i^{y_i} (1-a_i)^{1-y_i}}_{\text{Bernoulli pmf}}$$

$$W_{MLE} = \arg\min_{w} - \log p(\mathcal{D}|w) := \mathcal{L}(w)$$

$$\mathcal{L}(w) = - \log p(\mathcal{D}|w) = - \sum_{i=1}^{n} y_i \log a_i + (1-y_i) \log (1-a_i)$$

Binary cross-entropy loss

Before we compute $\nabla_w \mathcal{L}(w)$, let's derive the following :

- $\log a = \log \sigma(w^T x) = \log \dfrac{1}{1+e^{-w^T x}} = - \log (1+e^{-w^T x})$

- $\log(1-a) = \log(1-\sigma(w^T x)) = \dots = -w^T x - \log(1+e^{-w^T x})$

- $\dfrac{\partial}{\partial w} \log a = -\dfrac{-x e^{-w^T x}}{1+e^{-w^T x}} = x(1-a)$

- $\dfrac{\partial}{\partial w} \log(1-a) = -x + x(1-a) = -ax$

Therefore,

$$\frac{\partial}{\partial w_j} L(w) = -\sum_{i=1}^{n} y_i x_{ij}(1-a_i) - (1-y_i)x_{ij} a_i$$

$$= \dots = \sum_{i=1}^{n} (a_i - y_i) x_{ij}$$

Notice that: $\quad \nabla_w L(w) = X^T(a-y)$
$$\underset{(d+1)\times 1}{} \quad \underset{(d+1)\times n}{} \underset{n\times 1}{}$$

$$X = \begin{bmatrix} 1 & x_{11} & \cdots & \cdots & x_{1d} \\ 1 & & & & \\ \vdots & & & & \\ 1 & x_{n1} & \cdots & \cdots & x_{nd} \end{bmatrix} \underset{n\times(d+1)}{} \quad \overset{w_0}{\curvearrowleft}$$

$$a = \begin{bmatrix} \sigma(w^T x_1) \\ \vdots \\ \sigma(w^T x_n) \end{bmatrix}_{n\times 1} \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_r \end{bmatrix}_{n\times 1}$$

④ We cannot solve for $w_{MLE}$ analytically since $w$ shows up in a non-linear fashion in $\nabla_w L(w) = 0$.

<u>Hessian</u>: $\nabla_w^2 L(w)$

$$\left\{ \frac{\partial^2}{\partial w_j \partial w_k} L(w) = \sum_{i=1}^{n} x_{ij} \left( \frac{\partial}{\partial w_k} a_i \right) = \overbrace{\sum_{i=1}^{n} x_{ij} x_{ik} a_i (1-a_i)}^{\text{quadratic form}} \right.$$

$$= z_j^T A z_k \quad, \quad z_j := (x_{1j}, \dots, x_{nj})^T \to \begin{matrix} j\text{-th} \\ \text{of} \\ \text{(design} \\ \text{ma} \end{matrix}$$

where, $A := \begin{bmatrix} a_1(1-a_1) & & 0 \\ & \ddots & \\ 0 & & a_n(1-a_n) \end{bmatrix}$
$\underset{n\times n}{}$

$\nabla^2 \dots \quad \dots^T A \dots \quad$ one can show that this is

$$\implies \nabla_w L(w) = X A X$$
$$\underset{(d+1)\times(d+1)}{} \quad \underset{(d+1)\times n}{} \quad \underset{n\times n}{} \quad \underset{n\times(d+1)}{}$$

a positive-semidefinite matrix

$$\implies L(w) \text{ is convex in } w.$$

## Iterative re-weighted least squares :

Recall Newton : $\quad w_{t+1} = w_t - H_t^{-1} g_t \quad , \quad \boxed{\begin{array}{l} H_t = X^T A_t X \\ g_t = X^T(a_t - y) \end{array}}$

$$\implies w_{t+1} = w_t - (X^T A_t X)^{-1} X^T (a_t - y)$$

we re-write this as :

$$w_{t+1} = (X^T A_t X)^{-1} X^T A_t \overbrace{\left[ X w_t - A_t^{-1}(a-y) \right]}^{V_t}$$

$$= \underbrace{(X^T A_t X)^{-1} X^T A_t V_t}_{\text{Hessian}} \longrightarrow \begin{array}{l} \text{is the solution to a } \underline{\text{weighted}} \\ \text{least squares problems} \end{array}$$

Recall the MLE solution in linear regression :

$$\longrightarrow w_{MLE} = (X^T X)^{-1} X^T y = (X^T A X)^{-1} X^T A y$$

where $A$ is the identity matrix.

---

## Multi-class logistic regression :

soft-max function, a
generalization of 1
logistic sigmoic
the multi-class
setting

Model : $\quad p(y = c \mid x, w) = \dfrac{\exp(w_c^T x)}{\sum\limits_{c'=1}^{G} \exp(w_{c'}^T x)}$

where $w_c$ is the $c$-th column of $W$
$$(d+1)\times G$$

an $y$ is a <u>one-hot</u> encoding matrix : $\quad y = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$
$$n\times G$$

i.e. $y_{ic} = \mathbb{1}_{\{y_i = c\}}$

Likelihood :  $p(\mathcal{D}|w) = \prod\limits_{i=1}^{n} \prod\limits_{c=1}^{G} p(y_i = c | x_i, W)$

$\Rightarrow -\log p(\mathcal{D}|w) = -\sum\limits_{i=1}^{n} \left[ \left( \sum\limits_{c=1}^{G} y_{ic} W_c^T x_i \right) - \log \left( \sum\limits_{c'=1}^{G} \exp(w_{c'}^T x_i) \right. \right.$

$\hookrightarrow$ multi-class cross-entropy loss