# Lesson B01 (b):
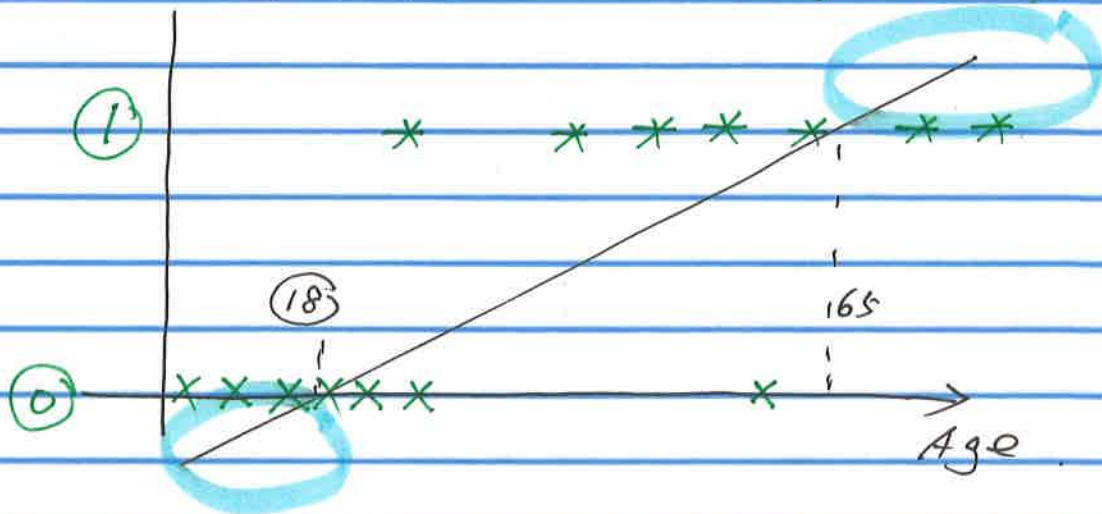## Logistic Regression.

Logistic Regression: one of the most ML algorithms for binary classification.

Main Implementation Steps:

① How to compute logistic function.

② How to learn the coefficient for a logistic regression, using stochastic gradient descent

③ How to predict?

Action (Y/N) (eg: insurance of health)



① 

⑱

165

⓪ 

Age.

why Linear regression is
not applicable on above
problem? → Range of DV
may lie outside of [0, 1].

Solution:

$$y = \beta_0 + \beta_1 x.$$
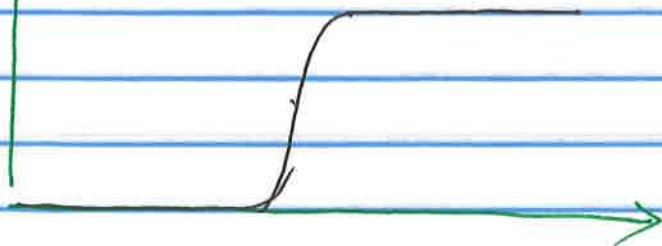(Linear regression)

$$P = \frac{1}{1 + e^{-y}}$$
(sigmoid function)

$$\boxed{\ln\left(\frac{P}{1-P}\right) = \beta_0 + \beta_1 x}$$

$$= e^y$$

$\hat{P}$ (probability)



②

# Principle of Logistic Regression.

→ linear regression of LOGIT.

( or $\ln(odds)$ )

① Define $P$ as the probability of an Event.

$$odds = \frac{P}{1-P} \iff P = \frac{odds}{1 + odds}.$$

② $\underline{\ln(Odds)} = \beta_0 + \beta_1 x_1 + \beta_2 x_2$
$$+ \dots + \beta_p \cdot x_p + e$$
$$\downarrow$$
LOGIT.

$$P \in [0,1] \longrightarrow \ln(Odds) \in (-\infty, +\infty)$$

e.g : $P = 0 \Rightarrow odds = 0 \Rightarrow \ln(odds) = -\infty$

e.g : $P = 1 \Rightarrow odds = +\infty \Rightarrow \ln(odds) = +\infty$

③ $\ln(Odds) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$ ⟶ $\boxed{sigmoid}$

$\Rightarrow \quad odds = e^{\beta_0 + \beta_1 x_1 + \dots \beta_p x_p}$

$$\Rightarrow P = \frac{e^{\beta_0 + \beta_1 x_1 + \dots}}{1 + e^{\beta_0 + \beta_1 x_1 + \dots}} = 1 - \frac{1}{1 + e^{\beta_0 + \beta_1 \dots}}$$

③

$\hat{P}$ (expected porbability)

$P = 91\%$

$P = 85\%$

$P = 27\%$

$P = 17\%$

20  30  40  50

$x$.

$\hat{y}$ (predicted DV)

1.0

0.5 — — — — — — threshold.

$x$

$\hat{y} = 0$.    $\hat{y} = 1$.

Logistic Regression.

① Supervised Learning for binary
classification.

② logit = log odds.
odds = $\frac{P(event)}{1 - P(event)}$ $(\in (-\infty, +\infty))$

$P(x) = Pr(y=1/x)$

③ Sigmoid func. $P(x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$

$\Rightarrow \ln(\frac{P}{1-P}) = \beta_0 + \beta_1 x$.

④

# Parameter Estimation

The goal of Learning is to estimate parameter vector $\hat{\vec{\beta}} = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}$.

Logistic Regression use Maximum Likelihood to estimate $\hat{\beta}$.

Given $N$ sample with Label $0/1$.

$$\begin{cases} \longrightarrow \text{sample labelled "1"}: & \hat{P}(x) \longrightarrow 1. \\ \longrightarrow \text{sample labelled "0"}: & 1 - \hat{P}(x) \longrightarrow 1. \end{cases}$$

$$L(\vec{\beta}) = \prod_{s \in y_i = 1} P(x_i) \cdot \prod_{s \in y_i = 0} (1 - P(x_i))$$

(likilyhood function)

$$= \prod_{s} P(x_i)^{y_i} \cdot (1 - P(x_i))^{(1 - y_i)}$$

$$\Rightarrow \log L(\vec{\beta}) = \sum_{i=1}^{N} \left[ y_i \log(P(x_i)) + (1 - y_i) \log(1 - P(x_i)) \right]$$

$$= \sum_{i=1}^{N} \left[ y_i \log\left( \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_i)}} \right) + (1 - y_i) \log\left( \frac{e^{-\beta_0 - \beta_1 x_i}}{1 + e^{-\beta_0 - \beta_1 x_i}} \right) \right]$$

$$\Rightarrow \log L(\vec{\beta}) = \sum_{i=1}^{N} \left[ y_i (\beta_0 + \beta_1 x_i) x_i - \log(1 + e^{\beta_0 + \beta_1 x}) \right]$$

⑤

$$\vec{\beta} = \underset{\vec{\beta}}{argmax} \; log(L(\vec{\beta}))$$

$$log(L(\vec{\beta})) = \sum_{i=1}^{N} \left[ y_i (\beta_0 + \beta_1 x_i) x_i \right.$$

$$\left. - log(1 + e^{\beta_0 + \beta_1 x_i}) \right]$$

— transcendental Equation

Define $\ell(\vec{\beta}) = log(L(\vec{\beta}))$

$$\vec{\beta} = \underset{\vec{\beta}}{argmax} \; \ell(\vec{\beta})$$

2  Methods : (Iterative)

① Gradient

$$\vec{\beta}^{t+1} = \vec{\beta}^{t} + \alpha \cdot \nabla_{\beta} \ell(\vec{\beta}^{t})$$

② Newton Raphson.

$$\vec{\beta}^{t+1} = \vec{\beta}^{t} - \frac{\nabla_{\beta} \ell(\vec{\beta}^{t})}{\nabla_{\beta}^{2} \ell(\vec{\beta}^{t})}$$

Hessian

⑥

$$\nabla_{\vec{\beta}} \, \ell(\vec{\beta}) = \nabla_{\vec{\beta}} \sum_{i=1}^{N} \left[ y_i (\beta_0 + \beta_1 x_i) \right.$$

$$\left. - \log \left( 1 + e^{\beta_0 + \beta_1 x_i} \right) \right]$$

$$= \sum_{i=1}^{N} \left[ y_i - P(x_i) \right] \cdot x_i$$

$$\nabla_{\vec{\beta}}^{2} \, \ell(\vec{\beta}) = - \sum_{i=1}^{N} P(x_i) \left( 1 - P(x_i) \right) x_i^{\top} x_i$$

Hands Exercise of Logistic Regression
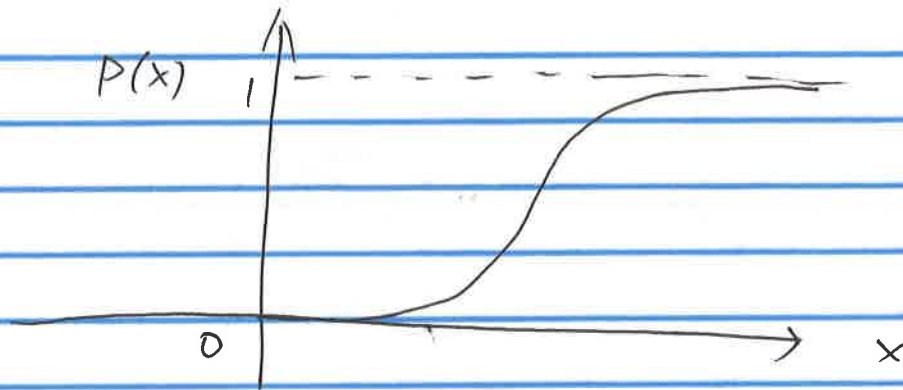(Training based on Gradient Ascent)

⓪  Given a training set

| $X_1$ | $X_2$ | $y$ |
|------|------|-----|
| 2.7 | 2.5 | 0 |
| 1.4 | 2.3 | 0 |
| 3.3 | 4.4 | 0 |
| 3.06 | 3.05 | 0 |
| 5.3 | 2.75 | 1 |



Observations:

① The data is linearly seperable

② We need to transform data points using LOGIT or sigmoid functions.

step ②    Find the parameter of model.



$$P(\vec{x}) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \beta_2 x_2)}}$$

$$(\text{threshold} = 0.5)$$

✓ Initially, we assume $\beta_0 = 0$. $\beta_1 = 0$.
$$\beta_2 = 0.$$

✓ we will calculate the prediction using the Above model parameter

◁— e.g. for 1$^{st}$ observation:

$$\boxed{x_1 = 2.7. \qquad x_2 = 2.5. \qquad y = 0}$$

$$\text{prediction} = \frac{1}{1 + e^{-(0 + 0.0 * 2.7 + 0.0 * 2.5)}} = 0.5$$

Next we will update parameter via:

(learn-rate)

$$\beta = \beta + \alpha * (y - \text{prediction}) * \text{prediction} *$$
$$(1 - \text{prediction}) * x$$

⑨ (β₀, β₁, or β₂) ⟹ $\beta_0 = -0.0375$, $\beta_1 = -0.1$; $\beta_2 = -0.095$

## Step 3

Use the newly parameter
compute the prediction again.

★ One round of calculating the
predicted value for all training
set using newly updated
parameter is "Epoch".