



# Workshop 3

COMP90051 Statistical Machine Learning  
Semester 2, 2024

# Learning outcomes

At the end of this workshop you should:

- be able to implement **linear regression** and **logistic regression**
- be able to explain how the **optimisation problems** for linear regression and logistic regression differ
- be able to implement **gradient descent**

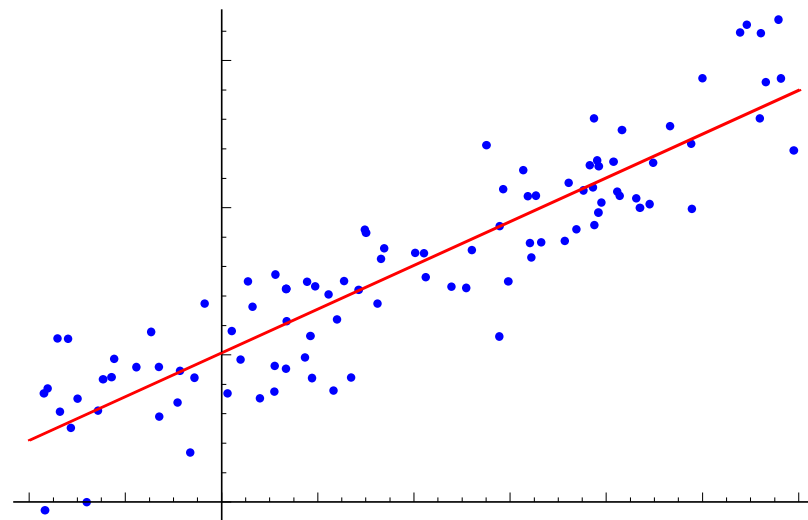
# Linear regression

Assume the response  $y$  is a *linear* function of the features  $\mathbf{x} = [x_1, \dots, x_m]^T$ :

$$y = w_0 + \sum_{i=1}^m w_i \cdot x_i$$

Write this more compactly as  $y = \mathbf{x}^T \mathbf{w}$  by redefining  $\mathbf{x} = [x_0, x_1, \dots, x_m]^T$  with  $x_0 = 1$  and defining  $\mathbf{w} = [w_0, \dots, w_m]^T$

**If we encode noise:**  $y = \mathbf{x}^T \mathbf{w} + \varepsilon_{\text{noise}}$



**Question:** How do we choose the weights?

# Solving linear regression

## Decision theoretic view

Make decision that minimises the empirical risk

$$\hat{R} = \frac{1}{n} \sum_{i=1}^n L(y_i, \hat{y}_i)$$

and choose the square loss

$$L(y, \hat{y}) = (\hat{y} - y)^2. \text{ take derivative}$$

Optimal decision for  $\mathbf{w}$

minimises the sum-squared error.

## Probabilistic view *likelihood*

Assume

$$y|\mathbf{x}, \mathbf{w} \sim \mathcal{N}(\mathbf{x}^T \mathbf{w}; \sigma^2)$$

Can write down the likelihood for the observations

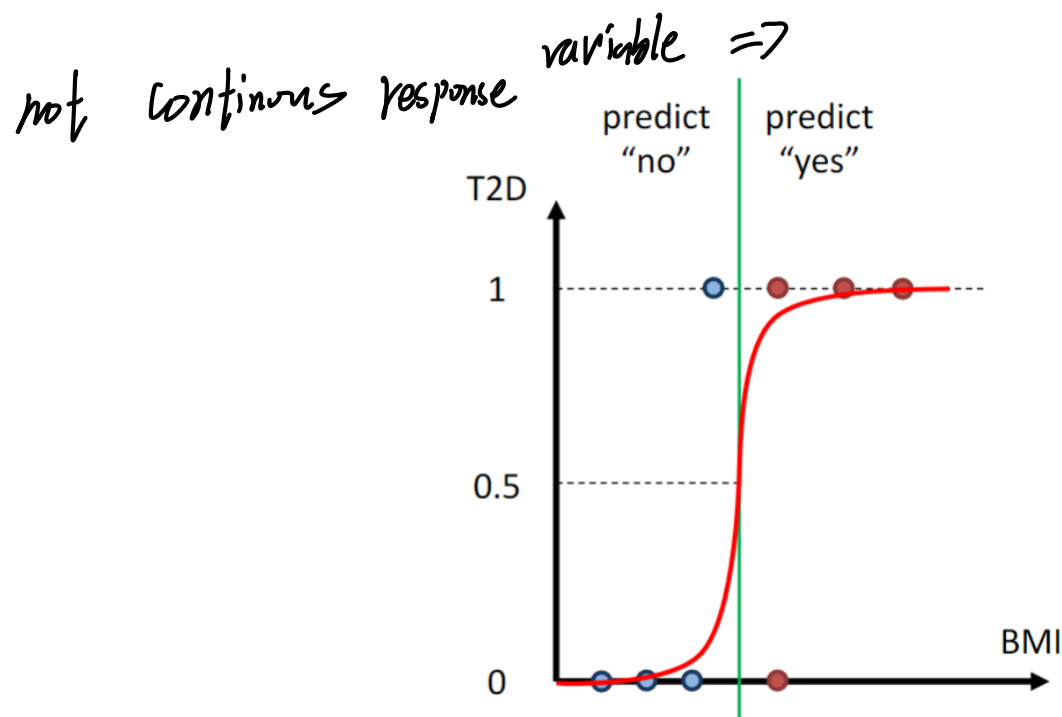
$$\begin{aligned} L(\mathbf{w}|\mathbf{X}, \mathbf{Y}) \\ = \prod_{i=1}^n p(y_i|\mathbf{x}_i, \mathbf{w}, \sigma) \end{aligned}$$

*MLE*

MLE for  $\mathbf{w}$  minimises the sum-squared error.

# Logistic regression

- Logistic regression is a linear binary (could be extend to multi-class) classifier for classification task
- Linear regression: gives a continuous value of **output  $y$**  for a given input  $X$ .
- Logistic regression: gives a continuous value of  $P(Y=1)$  for a given input  $X$ , which is later converted to  $Y=0$  or  $Y=1$  based on a threshold value.



# Solving logistic regression

Logistic regression optimisation problem:

$$\mathbf{w}^* \in \arg \min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \ell(y_i, \mu_i)$$

where  $\mu_i = \frac{1}{1 + e^{-\mathbf{x}_i^\top \mathbf{w}}}$  and  $\ell(y, \mu) = -y \log \mu - (1 - y) \log(1 - \mu)$   
*loss entropy*

Unfortunately, no closed form solution, need to use optimization techniques

*Gradient Descent: easy to compute slow to converge*

IRLS: *quick to convergence*

Least square 的 optimal weights:

$$(y - Xw)^T (y - Xw)$$

$$= y^T y - y^T Xw -$$

$$\begin{cases} \nabla_w W^T X = \nabla_w X^T W = X^T \\ \nabla_w A W = A \\ \nabla_w W^T A W = W^T (A^T + A) \end{cases}$$

$$\frac{\partial L}{\partial w} =$$

$$(Xw - y)^T (Xw - y)$$

$$X: n \times p$$

$$w: p \times 1$$

$$\nabla_w (R(w)) = w^T X^T X w - w^T X^T y - y^T X w - y^T y$$

$$= \underbrace{W^T (X^T X + X^T X)}_{=: 2X^T X} - y^T X - y^T X \Rightarrow$$

$$2 W^T X^T X - 2 X^T y$$

$$W^T = y^T X (X^T X)^{-1} \text{ invertible.}$$

$$W^* = (X^T X)^{-1} X^T y$$

$\frac{1}{n} X^T (M - y) \rightarrow$  logistic 的梯度  $\rightarrow$  update the weight at the opposite direction

$$W_t = W_{t-1} - \underbrace{\eta}_{\downarrow \text{学习率}} \nabla_w R(W_{t-1})$$

# Worksheet 3