COMP2261 ARTIFICIAL INTELLIGENCE / MACHINE LEARNING

# Multivariate Linear Regression

-- Vectorisation

Dr SHI Lei





- Vectorisation of Hypothesis Function
- Vectorisation of Cost Function









$$h_{\theta}(\mathbf{x}) = \theta_0 + \theta_1 \cdot x_1 + \theta_2 \cdot x_2 + \dots + \theta_n \cdot x_n$$

$$h_{\theta}(\mathbf{x}) = \theta_0 \cdot 1 + \theta_1 \cdot x_1 + \theta_2 \cdot x_2 + \dots + \theta_n \cdot x_n$$

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \dots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1} \qquad \boldsymbol{x} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \dots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

parameter vector

variable vector





$$h_{\theta}(x) = \theta_0 + \theta_1 \cdot x_1 + \theta_2 \cdot x_2 + \dots + \theta_n \cdot x_n$$

$$h_{\theta}(\mathbf{x}) = \theta_0 \cdot x_0 + \theta_1 \cdot x_1 + \theta_2 \cdot x_2 + \dots + \theta_n \cdot x_n \qquad (x_0 = 1, constant)$$

$$= [\theta_0, \theta_1, \theta_2, \dots, \theta_n] \cdot \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \dots \\ x_n \end{bmatrix}$$

 $= \theta^T x$  for single instance





# Vectorisation of Hypothesis Function for the whole Training Set





features label

income (k)	age	# of children	tea (cups/week)	happiness	
38.63	46	1	7	2.314	
49.79	37	0	15	3.433	
40.34	52	3	20	4.03	
21.18	25	0	3	1.45	





#### EXAMPLE.

#### Happiness predictor

$x_1$	$\underline{x_2}$		$x_3$		$\mathcal{X}_4$	${\mathcal Y}$		
income (k)	age		# of children	te	a (cups/weel	<b>(</b> )	happiness	
38.63	46		1		7		2.314	
49.79	37		0		15		3.433	
40.34	52		3		20		4.03	
21.18	25		0		3		1.45	





#### EXAMPLE. Happiness predictor

$x_1$	$x_2$	$x_3$	$x_4$		$\boldsymbol{\mathcal{X}}_n$	y
income (k)	age	# of children te	a (cups/week)	th	ne n <sup>th</sup> feature	happiness
38.63	46	1	7		₹	2.314
49.79	37	0	15		₹	3.433
40.34	52	3	20		₹	4.03
21.18	25	0	3		₹	1.45





#### EXAMPLE. Happiness predictor

		income (k)	age	# of children	tea (cups/week)	 the n <sup>th</sup> feature	happiness	
$(x^{(1)})^T$	1	38.63	46	1	7	 水	2.314	$y^{(1)}$
$(x^{(2)})^T$	1	49.79	37	0	15	 た	3.433	$y^{(2)}$
$(x^{(3)})^T$	1	40.34	52	3	20	 た	4.03	$y^{(3)}$
$(x^{(4)})^T$	1	21.18	25	0	3	 ネ	1.45	$y^{(4)}$
$(x^{(m)})^T$	1	₹	☆	₹	τ̈	 ☆	1.45	$y^{(m)}$

EXAMPLE.

Happiness predictor

$$x_0^{(2)}$$
  $x_1^{(2)}$   $x_2^{(2)}$   $x_3^{(2)}$   $x_4^{(2)}$  ...  $x_n^{(2)}$   $(x^{(2)})^T = \begin{bmatrix} 1 & 49.79 & 37 & 0 & 15 & ... & * \end{bmatrix}$ 





#### Happiness predictor

$$\begin{bmatrix} 1 \\ 49.79 \\ x_1^{(2)} \\ x_1^{(2)} \\ x_1^{(2)} \\ x_2^{(2)} \\ 0 \\ x_3^{(2)} \\ 15 \\ x_4^{(2)} \\ \dots \\ x_n^{(2)} \end{bmatrix}$$

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ x_2^{(i)} \\ x_3^{(i)} \\ x_4^{(i)} \\ \dots \\ x_n^{(i)} \end{bmatrix}$$

For the i-th instance



$$X = \begin{bmatrix} (x^{(1)})^T \\ (x^{(2)})^T \\ (x^{(3)})^T \\ (x^{(4)})^T \\ \dots \\ (x^{(m)})^T \end{bmatrix}$$





$$\boldsymbol{X} = \begin{bmatrix} (\boldsymbol{x^{(1)}})^T \\ (\boldsymbol{x^{(2)}})^T \\ (\boldsymbol{x^{(3)}})^T \\ (\boldsymbol{x^{(4)}})^T \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(i)} & x_2^{(i)} & \cdots & x_n^{(i)} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 1 & x_1^{(i)} & x_2^{(i)} & \cdots & x_n^{(i)} & \cdots & x_n^{(i)} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 1 & x_1^{(i)} & x_2^{(i)} & \cdots & x_n^{(i)} & \cdots & x_n^{(i)} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 1 & x_1^{(m)} & x_2^{(m)} & \cdots & x_n^{(m)} \end{bmatrix} \qquad \boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_j \\ \vdots \\ \theta_n \end{bmatrix}$$





$$h_{\theta}\left(x_1^{(1)}, x_2^{(1)}, \dots, x_n^{(1)}\right) = h_{\theta}(x^{(1)})$$

$$\boldsymbol{X}\boldsymbol{\theta} = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \cdots & x_j^{(1)} & \cdots & x_n^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \cdots & x_j^{(2)} & \cdots & x_n^{(2)} \\ 1 & x_1^{(3)} & x_2^{(3)} & \cdots & x_j^{(3)} & \cdots & x_n^{(3)} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 1 & x_1^{(i)} & x_2^{(i)} & \cdots & x_j^{(i)} & \cdots & x_n^{(i)} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 1 & x_1^{(m)} & x_2^{(m)} & \cdots & x_n^{(m)} \end{bmatrix} \cdot \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} = \begin{bmatrix} \theta_0 + \theta_1 \cdot x_1^{(1)} + \theta_2 \cdot x_2^{(1)} + \cdots + \theta_n \cdot x_n^{(1)} \\ \theta_0 + \theta_1 \cdot x_1^{(2)} + \theta_2 \cdot x_2^{(2)} + \cdots + \theta_n \cdot x_n^{(2)} \\ \theta_0 + \theta_1 \cdot x_1^{(3)} + \theta_2 \cdot x_2^{(3)} + \cdots + \theta_n \cdot x_n^{(3)} \\ \vdots \\ \theta_0 + \theta_1 \cdot x_1^{(i)} + \theta_2 \cdot x_2^{(i)} + \cdots + \theta_n \cdot x_n^{(i)} \\ \vdots \\ \theta_0 + \theta_1 \cdot x_1^{(i)} + \theta_2 \cdot x_2^{(i)} + \cdots + \theta_n \cdot x_n^{(i)} \\ \vdots \\ \theta_0 + \theta_1 \cdot x_1^{(m)} + \theta_2 \cdot x_2^{(m)} + \cdots + \theta_n \cdot x_n^{(m)} \end{bmatrix}$$

 $\mathbb{R}^{m \times (n+1)}$ 

 $\mathbb{R}^{(n+1)\times 1}$ 

 $\mathbb{R}^{m \times 1}$ 





 $\mathbb{R}^{m \times (n+1)}$ 

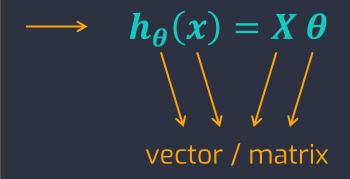
$$\boldsymbol{X}\boldsymbol{\theta} = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \cdots & x_j^{(1)} & \cdots & x_n^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \cdots & x_j^{(2)} & \cdots & x_n^{(2)} \\ 1 & x_1^{(3)} & x_2^{(3)} & \cdots & x_j^{(3)} & \cdots & x_n^{(3)} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 1 & x_1^{(i)} & x_2^{(i)} & \cdots & x_j^{(i)} & \cdots & x_n^{(i)} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 1 & x_1^{(m)} & x_2^{(m)} & \cdots & x_n^{(m)} & \cdots & x_n^{(m)} \end{bmatrix} \cdot \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_j \\ \vdots \\ \theta_n \end{bmatrix} = \begin{bmatrix} h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(1)}) \\ h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(2)}) \\ h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(3)}) \\ \vdots \\ h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}) \\ \vdots \\ \theta_n \end{bmatrix}$$

 $\mathbb{R}^{(n+1)\times 1}$ 





 $\mathbb{R}^{m \times 1}$ 



A vector containing the results of the hypothesis function, i.e., the predicted value, for all the instances from the training set.



$$h_{\theta}(\mathbf{x}) = \mathbf{\theta}^T \mathbf{x}$$

VS

$$h_{\theta}(x) = X \theta$$

takes single instance outputs predicted label

takes all instances

outputs a vector of predicted labels









$$\boldsymbol{X}\boldsymbol{\theta} - \boldsymbol{y} = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \cdots & x_n^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \cdots & x_n^{(2)} \\ 1 & x_1^{(3)} & x_2^{(3)} & \cdots & x_n^{(3)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(m)} & x_2^{(m)} & \cdots & x_n^{(m)} \end{bmatrix} \cdot \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} - \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ y^{(3)} \\ \vdots \\ y^{(m)} \end{bmatrix} = \begin{bmatrix} h_{\theta}(x^{(1)}) - y^{(1)} \\ h_{\theta}(x^{(2)}) - y^{(2)} \\ h_{\theta}(x^{(3)}) - y^{(3)} \\ \vdots \\ h_{\theta}(x^{(m)}) - y^{(m)} \end{bmatrix}$$





$$\mathbf{X}\boldsymbol{\theta} - \mathbf{y} = \begin{bmatrix} h_{\theta}(x^{(1)}) - y^{(1)} \\ h_{\theta}(x^{(2)}) - y^{(2)} \\ h_{\theta}(x^{(3)}) - y^{(3)} \\ \vdots \\ h_{\theta}(x^{(m)}) - y^{(m)} \end{bmatrix}$$





$$(X\theta - y)^{T} \cdot (X\theta - y)$$

$$= [h_{\theta}(x^{(1)}) - y^{(1)}, h_{\theta}(x^{(2)}) - y^{(2)}, h_{\theta}(x^{(3)}) - y^{(3)}, ..., h_{\theta}(x^{(m)}) - y^{(m)}] \cdot \begin{bmatrix} h_{\theta}(x^{(1)}) - y^{(1)} \\ h_{\theta}(x^{(2)}) - y^{(2)} \\ h_{\theta}(x^{(3)}) - y^{(3)} \\ \vdots \\ h_{\theta}(x^{(m)}) - y^{(m)} \end{bmatrix}$$

$$=(h_{\theta}(x^{(1)})-y^{(1)})^2 + (h_{\theta}(x^{(2)})-y^{(2)})^2 + (h_{\theta}(x^{(3)})-y^{(3)})^2 + ... + (h_{\theta}(x^{(m)})-y^{(m)})^2$$

$$= \sum_{i=1}^{m} (h_{\theta}(x^{(i)}) - y^{(i)})^2$$





$$(X\theta - y)^{T} \cdot (X\theta - y) = \sum_{i=1}^{m} (h_{\theta}(x^{(i)}) - y^{(i)})^{2}$$





$$\sum_{i=1}^{m} (h_{\theta}(x^{(i)}) - y^{(i)})^2 = (X\theta - y)^T (X\theta - y)$$





$$\frac{1}{2m} \sum_{i=1}^{m} (h_{\theta}(x^{(i)}) - y^{(i)})^{2} = \frac{1}{2m} (X\theta - y)^{T} (X\theta - y)$$





$$J(\boldsymbol{\theta}) = \frac{1}{2m} \sum_{i=1}^{m} (h_{\theta}(x^{(i)}) - y^{(i)})^2 = \frac{1}{2m} (X\boldsymbol{\theta} - y)^T (X\boldsymbol{\theta} - y)$$





$$J(\boldsymbol{\theta}) = \frac{1}{2m} (\boldsymbol{X}\boldsymbol{\theta} - \boldsymbol{y})^T (\boldsymbol{X}\boldsymbol{\theta} - \boldsymbol{y})$$









#### Repeat until convergence {

$$\theta_0 := \theta_0 - \alpha \, \frac{\partial J}{\partial \theta_0}(\boldsymbol{\theta})$$

$$\theta_1 := \theta_1 - \alpha \, \frac{\partial J}{\partial \theta_1}(\boldsymbol{\theta})$$

• • •

$$\theta_n := \theta_1 - \alpha \, \frac{\partial J}{\partial \theta_n}(\boldsymbol{\theta})$$

$$\frac{\partial J}{\partial \theta_k}(\theta) = \frac{1}{m} \sum_{i=1}^{m} (h_{\theta}(x)^{(i)} - y^{(i)}) x_k^{(i)} 
= \frac{1}{m} \sum_{i=1}^{m} x_k^{(i)} (h_{\theta}(x)^{(i)} - y^{(i)})$$

$$= \frac{1}{m} (X\theta - y)$$







#### Repeat until convergence {

$$\theta_0 := \theta_0 - \alpha \left[ \frac{\partial J}{\partial \theta_0}(\boldsymbol{\theta}) \right] = \left[ \frac{1}{m} \boldsymbol{x}_k^T (\boldsymbol{X} \boldsymbol{\theta} - \boldsymbol{\theta}) \right]$$

$$\theta_1 := \theta_1 - \alpha \left[ \frac{\partial J}{\partial \theta_1} (\boldsymbol{\theta}) \right]$$

. . .

$$\theta_n := \theta_1 - \alpha \left| \frac{\partial J}{\partial \theta_n}(\boldsymbol{\theta}) \right|$$

}





#### Repeat until convergence {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} x_0^T (X\theta - y)$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} x_1^T (X\theta - y)$$

$$\vdots$$

$$\theta_n := \theta_1 - \alpha \frac{1}{m} x_n^T (X\theta - y)$$





#### Repeat until convergence {

$$\boldsymbol{\theta} := \boldsymbol{\theta} - \alpha \ \frac{1}{m} \boldsymbol{X}^T (\boldsymbol{X} \boldsymbol{\theta} - \boldsymbol{y})$$

}





$$h_{\theta}(x) = \theta^T x$$
  $h_{\theta}(x) = X \theta$ 

Vectorisation of Cost Function

$$J(\boldsymbol{\theta}) = \frac{1}{2m} (\boldsymbol{X}\boldsymbol{\theta} - \boldsymbol{y})^T (\boldsymbol{X}\boldsymbol{\theta} - \boldsymbol{y})$$

Vectorisation of Gradient Descent

Repeat until convergence {

$$\boldsymbol{\theta} := \boldsymbol{\theta} - \alpha \ \frac{1}{m} \boldsymbol{X}^T (\boldsymbol{X} \boldsymbol{\theta} - \boldsymbol{y})$$





#### Comparing with un-vectorised

- ✓ Easier to implement
- ✓ Computationally more efficient