

COMP2261 ARTIFICIAL INTELLIGENCE / MACHINE LEARNING

# Multivariate Linear Regression

-- Cost Function & Gradient Descent

Dr SHI Lei

# Supervised Learning

- To build a model represented as a hypothesis function  $h_{\theta}(x)$ .



data



model

# Supervised Learning

income (input  $x$ , dependent variable)



model (**hypothesis function**, mapping  $x \rightarrow y$ )



happiness (output  $y$ , independent variable)



**cost function**  
(loss function)

# Univariate Linear Regression

Hypothesis Function

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

Independent Variable

$$x$$

Parameters

$$\theta_0, \theta_1$$

Cost Function

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x)^{(i)} - y^{(i)})^2$$

# Univariate Linear Regression

## Hypothesis Function

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

## Independent Variable

$x$

## Model Parameters

$\theta_0, \theta_1$

## Cost Function

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x)^{(i)} - y^{(i)})^2$$

# Multivariate Linear Regression

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

$$\mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad (x_0 = 1, \text{constant})$$

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad h_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x}$$

vectors

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

# Univariate Linear Regression

( $n = 1$ )

Cost Function

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x)^{(i)} - y^{(i)})^2$$

Partial derivative of  $J$  with respect to parameters

$$\frac{\partial J}{\partial \theta_0}(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x)^{(i)} - y^{(i)}) \boxed{x_0^{(i)}}_{x_0^{(i)} = 1}$$

$$\frac{\partial J}{\partial \theta_1}(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x)^{(i)} - y^{(i)}) \boxed{x_1^{(i)}}$$

# Multivariate Linear Regression

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

$$\frac{\partial J}{\partial \theta_0}(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(\mathbf{x})^{(i)} - y^{(i)}) \boxed{x_0^{(i)}}_{x_0^{(i)} = 1}$$

$$\frac{\partial J}{\partial \theta_1}(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(\mathbf{x})^{(i)} - y^{(i)}) \boxed{x_1^{(i)}}$$

...

$$\frac{\partial J}{\partial \theta_n}(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(\mathbf{x})^{(i)} - y^{(i)}) x_n^{(i)}$$

# Univariate Linear Regression

## Gradient Descent

Repeat until convergence {  $\frac{\partial J}{\partial \theta_0}(\theta_0, \theta_1)$

$$\theta_0 := \theta_0 - \underbrace{\alpha}_{\text{learning rate}} \cdot \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x)^{(i)} - y^{(i)})$$

$$\theta_1 := \theta_1 - \underbrace{\alpha}_{\text{learning rate}} \cdot \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x)^{(i)} - y^{(i)})x^{(i)}$$

} simultaneously update  $\theta_0, \theta_1$   $\frac{\partial J}{\partial \theta_1}(\theta_0, \theta_1)$

# Multivariate Linear Regression

Repeat until convergence {  $\frac{\partial J}{\partial \theta_0}(\theta)$

$$\theta_0 := \theta_0 - \underbrace{\alpha}_{\text{learning rate}} \cdot \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x)^{(i)} - y^{(i)})x_0^{(i)}$$

$$\theta_1 := \theta_1 - \underbrace{\alpha}_{\text{learning rate}} \cdot \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x)^{(i)} - y^{(i)})x_1^{(i)}$$

...

$$\theta_n := \theta_n - \underbrace{\alpha}_{\text{learning rate}} \cdot \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x)^{(i)} - y^{(i)})x_n^{(i)}$$

} simultaneously update  $\theta_0, \dots, \theta_n$   $\frac{\partial J}{\partial \theta_n}(\theta)$

## ✓ Takeaway Points

- Univariate linear regression is a special case of multivariate linear regression when the number of features  $n=1$ .
- $n+1$  dimensional column vectors to denote features and model parameters.
- Feature vectors and parameters vectors to express hypothesis function and cost function.
- Each iteration in gradient descent, all the  $n$  parameters ( $\theta$ s) need to be updated simultaneously.