

Multivariate Linear Regression 多元线性回归

$$h_0(x) = \theta_0 + \theta_1 x$$

↓

$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

($x_0^{(i)} = 1$: 第 i 个样本的第 0 个特征固定为 1.) $x_j^{(i)}$: 第 i 个样本的第 j 个特征

↓

Hypothesis

$$h_0(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$= \theta^T x = [\theta_0 \ \theta_1 \ \dots \ \theta_n] \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

$1 \times (n+1) \quad (n+1) \times 1$

intuition: x_1 面积 x_2 卫生间个数 x_3 卧室
房价 \Leftarrow θ_0 : 每平方米 θ_1 : 每个卫生间 θ_2 : ...
单价 价格

Parameter: $\theta \in \mathbb{R}^{n+1} = \theta_0, \theta_1, \theta_2, \dots, \theta_n$

Cost Function $J(\theta) = \frac{1}{2m} \sum_{i=1}^m [h_0(x^{(i)}) - y^{(i)}]^2$

Gradient descent:

Repeat until convergence: $\{ \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \}$ for $j = 0, \dots$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{m} \sum_{i=1}^n [h_0(x^{(i)}) - y^{(i)}] \cdot x_j^{(i)}$$

eg: $\theta_0 = \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^n (h_0(x^{(i)}) - y^{(i)}) \underbrace{x_0^{(i)}}_{=1}$

$\theta_1 = \theta_1 - \dots \dots \dots x_1^{(i)}$

$\theta_2 = \theta_2 - \dots \dots \dots x_2^{(i)}$

Gradient Descent in practice

Way to speed up: Feature Scaling

因为如果 θ 的范围太大, 那么会收敛得很慢, 且有可能在最优值附近振荡

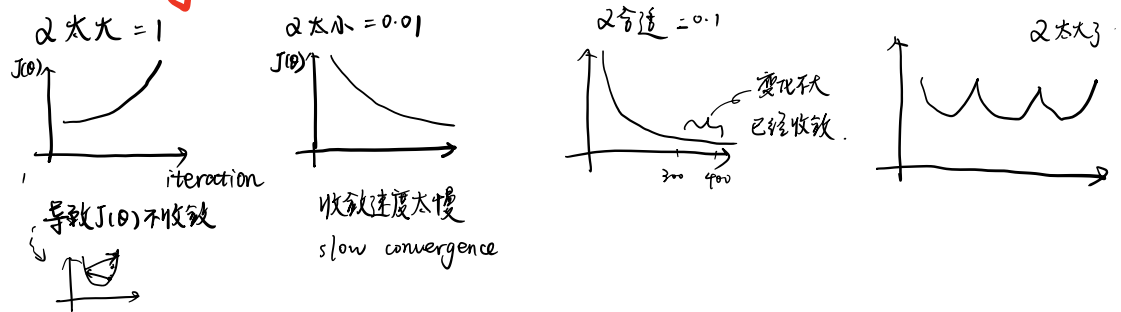


feature scaling

mean normalization: $\frac{x - \mu}{s}$ (scale of the data / standard deviation)

$X \in (-0.5, 0.5) / (-1, 1) / \dots$

Learning Rate α



Polynomial Regression

$$y = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

$$\text{or: } y = \theta_0 + \theta_1 x + \theta_2 \sqrt{x}$$

eg: size = 1 ~ 1000
size² = 1 ~ 10⁶
size³ = 1 ~ 10¹²
Feature scale 注意!

Normal Equation

$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 = (\theta_0 \ \theta_1 \ \theta_2 \ \theta_3) \begin{pmatrix} 1 \\ x_1 \\ x_2 \\ x_3 \end{pmatrix}$$

$$X = \begin{pmatrix} x_1 & x_2 & x_3 & y \\ x_1^{(1)} & x_2^{(1)} & x_3^{(1)} & y^{(1)} \\ x_1^{(2)} & x_2^{(2)} & x_3^{(2)} & y^{(2)} \\ x_1^{(3)} & x_2^{(3)} & x_3^{(3)} & y^{(3)} \\ x_1^{(4)} & x_2^{(4)} & x_3^{(4)} & y^{(4)} \end{pmatrix} = X$$

$X = \begin{pmatrix} | & | & | & | \\ 1 & x_1 & x_2 & x_3 \\ | & | & | & | \end{pmatrix} \quad y = \begin{pmatrix} | \\ y \\ | \end{pmatrix}$
 $m \times (n+1) \quad m \times 1$

$$y = \begin{pmatrix} 1 & x_1^{(1)} \\ 1 & x_1^{(2)} \\ 1 & x_1^{(3)} \\ 1 & x_1^{(4)} \end{pmatrix} \quad m \times 2$$

$$\theta = (X^T X)^{-1} X^T y \rightarrow \text{in octave: pinv}(X' * X) * X' * y$$

如果采用了这种形式表示, 就不用做 Feature Scaling; $x \in [0, 10^{-5}] \checkmark x \in [0, 10^4] \checkmark \text{OK}$.

m training samples, n features

Gradient Descent VS Normal Equation

1. need to choose α 1. 不需要选择 α

2. Need many iterations 2. 不需迭代

3. Work well even n is large 3. 需要计算 $(X^T X)^{-1}_{n \times n}$: $O(n^3)$

✓ 更常用. 4. 当 n 很大时, 会很慢.
eg n > 10000 时.

若 $X^T X$ 不可逆 (non-invertible) $\begin{cases} \text{① 有重复特征 (2列相同)} \\ \text{② 太多特征 (m \leq n)} (m \square n) \end{cases}$