

Model & Cost Function

Training Dataset



Learning Algorithm



input $x^{(i)} \rightarrow h \rightarrow \text{output } y^{(i)}$
(hypothesis)

how to measure hypothesis?

$$h(x) = \theta_0 + \theta_1 x \quad (\text{linear function})$$

linear regression

how to choose the best parameters?

so that $h(x)$ is close to y

$$J(\theta_0, \theta_1) = \underset{\theta_0, \theta_1}{\text{minimize}} \left(\frac{1}{2m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2 \right)$$

training example

$$h(x^{(i)}) = \theta_0 + \theta_1 x^{(i)}$$

均方误差
MSE (mean square error)

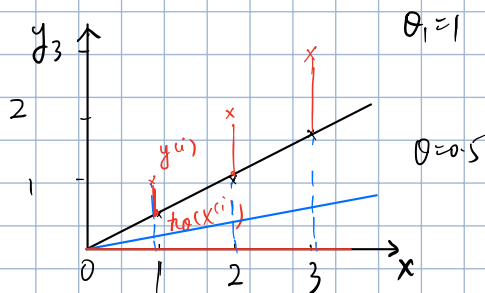
$$= \frac{1}{2} \bar{x} \rightarrow \bar{x} \text{ 指平方的均值}$$

$\frac{1}{2}$ 是为了计算梯度时的方便, 因为在求导时会和平方项的 $()^2$ 抵消.

minimize $J(\theta_0, \theta_1)$
 θ_0, θ_1 Cost function

Cost Function

hypothesis function
 $h(x)$



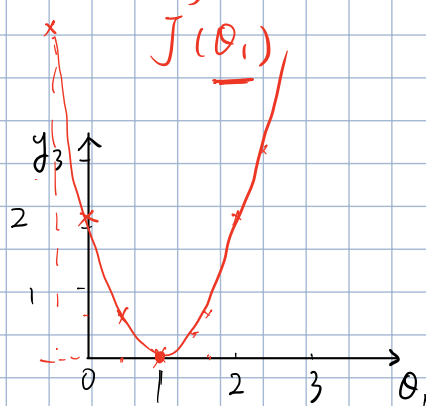
$$J(1) = \frac{1}{2m} \sum_{i=1}^m [h_0(x^{(i)}) - y^{(i)}]^2$$

$$= 0$$

$$J(0.5) = \frac{1}{2 \times 3} (0.5^2 + 1^2 + 1.5^2)$$

$$= \frac{3.5}{6} \approx 0.58$$

Loss function



$$J(0) = \frac{1}{2 \times 3} (1 + 4 + 9) = \frac{14}{6} \approx 2.3$$

cost function

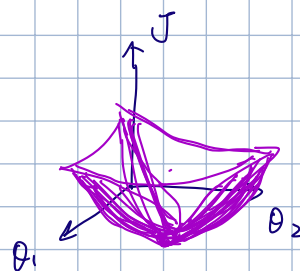
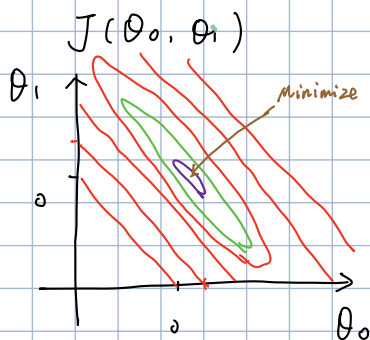
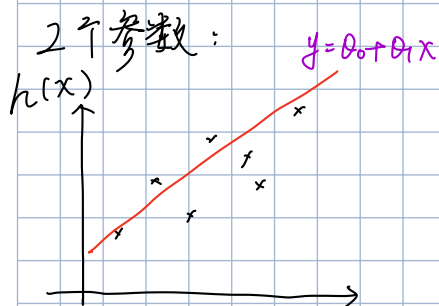
minimize $J(\theta_1)$, will find the best parameter.

Hypothesis: $H(x) = \theta_0 + \theta_1 x$

Parameter: θ_0, θ_1

Loss Function: $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m [h(x^{(i)}) - y^{(i)}]^2$

Goal: Minimize $J(\theta_0, \theta_1)$
 θ_0, θ_1



用等高线 (contour plot)

在2D上表示3D图象

同一条线上, $J(\theta_0, \theta_1)$ 即 Loss 相同

Parameter Learning

gradient descent Algorithm:

use to estimate the parameters in the hypothesis function

Any function $\rightarrow J(\theta_0, \theta_1) \dots$

$$\begin{cases} \theta_0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) \\ \theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) \end{cases}$$

simultaneously update:

1. 参数同步更新, 而非更新一个后, 用该值去更新另一个.

2. 对初始位置敏感, 可能会收敛到局部最小

(步长会自动越来越短, 越来越精细). local minimum

no need to decrease α over time.

参数: α 过小: 下降速度太慢

α 过大: 可能导致不收敛 (发散)
fail to converge



\rightarrow because gradient descent will automatically take smaller steps as the derivative get

$\frac{d}{d\theta_1} J(\theta_1)$ approaches 0 as we approach the bottom of our convex function. smaller.

so that $\theta_1 := \theta_1 - \alpha * 0$

Gradient descent & Cost Function

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

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

$$= \frac{1}{2m} \sum_{i=1}^m [\theta_0 + \theta_1 x^{(i)} - y^{(i)}]^2$$

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

$$\begin{aligned} \text{① } J=0: & \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) \\ &= \frac{1}{m} \sum_{i=1}^m [h(x^{(i)}) - y^{(i)}] \end{aligned}$$

$$\text{② } J=1: \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$= \frac{1}{m} \sum_{i=1}^m [h(x^{(i)}) - y^{(i)}] \cdot x^{(i)}$$

→ Gradient Descent Algorithm:

repeat until convergence {

$$\theta_0 := \theta_0 - \alpha \cdot \frac{1}{m} \sum_{i=1}^m [h_0(x^{(i)}) - y^{(i)}] \rightarrow \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) \quad \text{update } \theta_0 \text{ and } \theta_1,$$

$$\theta_1 := \theta_1 - \alpha \cdot \frac{1}{m} \sum_{i=1}^m [h_0(x^{(i)}) - y^{(i)}] \cdot x^{(i)} \rightarrow \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) \quad \text{simultaneously.}$$

即计算完一个参数后，不要代回原来式子，更新参数都是相同输入了

in linear regression:

Cost function are convex function ~~so~~ so that it doesn't have any local optimal except for the one global optimum.

实际中，有时找到局部最优也够用了。

"Batch" Gradient Descent:

each step of gradient descent using all the training samples.

Linear Algebra Review :

对于 $h(x) = \theta_0 + \theta_1 x$

Prediction = Data matrix * Parameters.

$$x \text{ (房间面积)} = \begin{pmatrix} 100 \\ 121 \\ 195 \\ 130 \end{pmatrix}, \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} \begin{pmatrix} 1 & 100 \\ 1 & 121 \\ 1 & 195 \\ 1 & 130 \end{pmatrix} \begin{pmatrix} 50 \\ 0.25 \end{pmatrix} = \begin{pmatrix} 1 \times 50 + 100 \times 0.25 \\ 1 \times 50 + 121 \times 0.25 \\ \vdots \end{pmatrix}$$

(补充) $4 \times 2 \quad 2 \times 1 \quad 4 \times 1$

矩阵向量乘法: $\square \times \square = \square$

$5 \times 2 \quad 2 \times 1 \quad 5 \times 1$
内层 决定能否相乘; 外层定型。

矩阵矩阵乘法，就像矩阵向量乘法的并行式操作:

$$\square \times \square = \square$$

$5 \times 2 \quad 2 \times 3 \quad 5 \times 3$

例如: 对于 $\begin{cases} h_1(x) = 50 + 0.25x \\ h_2(x) = 70 + 0.3x \\ h_3(x) = 20 + 0.7x \end{cases}$

the result of the prediction of the first hypothesis
the third hypothesis.

$$\begin{pmatrix} 1 & 100 \\ 1 & 121 \\ 1 & 195 \\ 1 & 130 \end{pmatrix} \begin{pmatrix} 50 & 70 & 20 \\ 0.25 & 0.3 & 0.7 \end{pmatrix} = \begin{pmatrix} \text{ } & \text{ } & \text{ } \\ \text{ } & \text{ } & \text{ } \\ \text{ } & \text{ } & \text{ } \\ \text{ } & \text{ } & \text{ } \end{pmatrix}$$

(补充) $4 \times 3 \quad 2 \times 3$

4x2

↑ 4x3
the second hypothesis

identity matrix:

$$A \cdot I = I \cdot A = A \rightarrow \text{其实是2个不一样的矩阵.}$$

\uparrow \uparrow \uparrow \uparrow \uparrow
 $m \times n$ $n \times n$ $m \times m$ $n \times n$ $m \times n$

$AB \neq BA$ in general
except A/B is I .

$$(A \cdot B) \cdot C = A \cdot (B \cdot C) \checkmark$$

Matrix Inverse.

eg: $12 \times \left(\frac{1}{12}\right) = 1$; $0 \times \underline{X} = 1?$

$\rightarrow = 12^{-1}$

$A \cdot A^{-1} = A^{-1} \cdot A = I$, only if A is a square matrix. $\rightarrow m \times m$ 零矩阵不存在逆矩阵.

Matrix Transpose

$$A = \begin{pmatrix} 1 & 2 & 0 \\ 3 & 5 & 9 \end{pmatrix} \quad A^T = \begin{pmatrix} 1 & 3 \\ 2 & 5 \\ 0 & 9 \end{pmatrix}$$

\uparrow \uparrow \uparrow \uparrow \uparrow \uparrow
 $= A$ 2×3 $= B$ 3×2

$$A_{ij} = B_{ji}$$