

week 1. Machine Learning

1.3 supervised M.L. { regression
classification }

1.4 unsupervised M.L. { clustering }

待分类书本 book: hands-on-sklearn

2. Linear regression with one variable

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

to fit the line to data points θ_0, θ_1

modeling error:



该模型误差是指 θ_0 和 θ_1 的值。

即 cost func 最小

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

自动指出使 $J(\theta_0, \theta_1)$ 变小的策略

Gradient Descent

不断指明pk下降的方向。



Batch gradient descent:

$$\text{repeat until convergence} \left\{ \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1), (j=0, 1) \right.$$

α is learning rate.

correct: simultaneous update:

$$\text{temp } 0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

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

$$\begin{aligned} \theta_0 &:= \text{temp } 0 && \text{同时更新 } \theta_0, \theta_1, \text{之后也就更新了 } J(\theta_0, \theta_1) \\ \theta_1 &:= \text{temp } 1 \end{aligned}$$

2.7 Gradient Descent for Linear Regression.

Gradient Descent algorithm

repeat until converge

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1), \quad \left. \begin{array}{l} (j=0 \text{ or } 1) \end{array} \right\}$$

Linear Regression model

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

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

$$j=0, \quad \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0} = \frac{1}{m} \sum_{i=1}^m [h_{\theta}(x^{(i)}) - y^{(i)}]^2$$

$$j=1, \quad \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1} = \frac{1}{m} \sum_{i=1}^m [h_{\theta}(x^{(i)}) - y^{(i)}] \cdot x^{(i)}$$

3. Linear Algebra Review

矩阵乘法:

$$\text{性质: } A \times B \neq B \times A$$

$$A \times (B \times C) = (A \times B) \times C$$

identity matrix: $I = \begin{bmatrix} 1 & 0 & 0 & \dots \\ 0 & 1 & 0 & \dots \\ 0 & 0 & 1 & \dots \end{bmatrix}$

$$AA^{-1} = A^{-1}A = I$$

$$AI = IA = A$$

转置 transpose: $\begin{vmatrix} a & b \\ c & d \\ e & f \end{vmatrix}^T = \begin{vmatrix} a & c & e \\ b & d & f \end{vmatrix}$

$$(A \pm B)^T = A^T \pm B^T$$

$$(A \times B)^T = B^T \times A^T$$

$$(A^T)^T = A$$

$$(kA)^T = kA^T ? \text{不是?}$$

Week 2

4. Linear Regression with Multiple Variables

$x_j^{(i)}$ 第*i*条第*j*个特征.

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

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

cost function:

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

Gradient Descent:

Repeat $\left\{ \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1, \dots, \theta_n) \right\}$

步骤:

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m [h_{\theta}(x^{(i)}) - y^{(i)}] x_j^{(i)}$$

simultaneously update θ_j for $j=0, 1, \dots, n$.

Feature Scaling: Make sure features are on a similar scale.

mean normalization, $x_n = \frac{x_n - \mu_n}{s_n}$, μ_n 是第*n*个特征的均值, s_n 是第*n*个特征的标准差

Learning Rate: 太大: 收敛慢
太小: 收敛慢

4.5 Features and Polynomial Regression.

线性特征：

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

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^3$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 \sqrt{x_1}$$

若用 polynomial regression, feature scaling 很重要, eg.:

$$x_1 \sim 1 \sim 1000$$

$$x_2 \sim 1 \sim 10^6$$

4.6 Normal equations. — 适用于 linear model.

有些 θ 不适用 normal equ. 更好。

Normal equ. 是指 $\theta = (X^T X)^{-1} X^T y$ 求出的

$J(\theta)$ 最小的参数。

$$\begin{matrix} x_0, x_1, x_2, \dots, x_n \\ \vdots \end{matrix}, y$$

$$\underbrace{\quad}_{X} \quad \underbrace{\quad}_{y}$$

与 gradient descent 对比：见 P51.

$$\theta = (X^T X)^{-1} X^T y$$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m$$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m [h_{\theta}(x^{(i)}) - y^{(i)}]^2, h_{\theta}(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

$$J(\theta) = \frac{1}{2} (X\theta - y)^T (X\theta - y), X \text{ 为 } m \times n \text{ 矩阵}, m \neq n, n \neq 1, \theta \text{ 为 } n \times 1 \text{ 矩阵}, y \text{ 为 } m \times 1 \text{ 矩阵}.$$

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

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

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

对 $J(\theta)$ 求偏导，

$$\text{求偏导}: \frac{dAB}{dB} = A^T$$

$$\frac{\partial J(\theta)}{\partial \theta} = \frac{1}{2} [2X^T X\theta - X^T y - (y^T X)^T - 0], \frac{dX^T Ax}{dx} = 2AX$$

$$= \frac{1}{2} [2X^T X\theta - X^T y - X^T y - 0], \frac{dX^T A}{dx} = A$$

$$= X^T X\theta - X^T y$$

$$\text{let } \frac{\partial J(\theta)}{\partial \theta} = 0 \Rightarrow \theta = (X^T X)^{-1} X^T y$$

Week 3

6. Logistic Regression.

对分类， y 值是离散的
值位于 $0 \sim 1$, $0 \leq h_\theta(x) \leq 1$

Logistic Regression 模型:

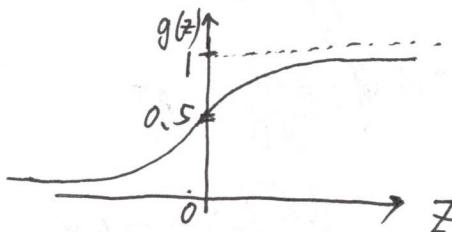
$$h_\theta(x) = g(\theta^T x)$$

x 是特征向量。

g 表示 logistic function

常用名 - ~~logistic~~ function 是 sigmoid function:

$$g(z) = \frac{1}{1 + e^{-z}}$$



$h(x)$ 的作用: 为了将输入, 计算该输入
入计算出输出 $\hat{y} = 1/h_\theta(x)$
概率, $h_\theta(x) = P(y=1|x; \theta)$

6.3 Precision boundary

Logistic Regression

$$h_\theta(x) = g(\theta^T x)$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

$\Rightarrow h_\theta(x) \geq 0.5$ 时, predict $y=1$
 < 0.5 , $y=0$

or:

$z=0$ 时, $g(z)=0.5$

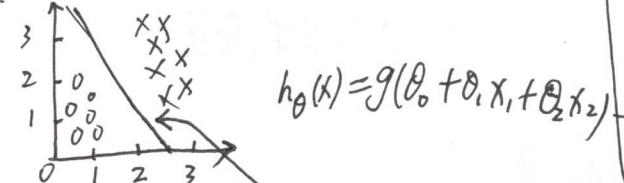
$z>0$ 时, $g(z) > 0.5$

$z<0$ 时, $g(z) < 0.5$

or $\theta^T x > 0$ 时, predict $y=1$

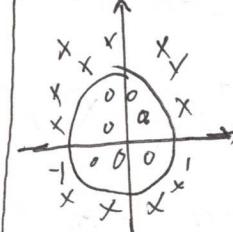
$\theta^T x < 0$, $y=0$

e.g.:



$$\theta = [-3, 1, 1]$$

当 $-3 + x_1 + x_2 > 0$ 即 $x_1 + x_2 > 3$ 时, predict $y=1$.



$$\theta = [-1, 0, 0, 1, 1]$$

6.4 Cost Function.

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}[h_\theta(x^{(i)}), y^{(i)}]$$

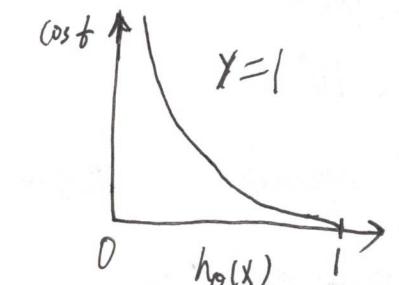
$$\text{Cost}[h_\theta(x^{(i)}), y^{(i)}]$$

$$= \begin{cases} -\log(h_\theta(x)) & , \text{if } y=1 \\ -\log(1-h_\theta(x)) & , \text{if } y=0 \end{cases}$$

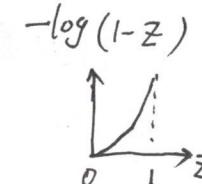
if $y=1$

cost = 0 if $y=1, h_\theta(x)=1$

as $h_\theta(x) \rightarrow 0$, cost $\rightarrow \infty$



if $y=0$



$$\text{cost}[h_{\theta}(x), y] = -y \log(h_{\theta}(x)) - (1-y) \log(1-h_{\theta}(x))$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m [-y^{(i)} \log(h_{\theta}(x^{(i)})) - (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)}))]$$

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(h_{\theta}(x^{(i)})) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)}))]$$

↳ Gradient Descent:

$$\text{Repeat } \left\{ \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \right. \\ \left. (\text{simultaneously update all } \theta_j) \right\}$$

解説 18E:

$$\text{Repeat } \left\{ \theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \right. \\ \left. (\text{simultaneously update all } \theta_j) \right\}$$

解説:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(h_{\theta}(x^{(i)})) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)}))]$$

$$h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}}$$

$$y^{(i)} \log(h_{\theta}(x^{(i)})) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)}))$$

$$= y^{(i)} \log\left(\frac{1}{1+e^{-\theta^T x}}\right) + (1-y^{(i)}) \log\left(1-\frac{1}{1+e^{-\theta^T x}}\right)$$

$$= -y^{(i)} \log(1+e^{-\theta^T x^{(i)}}) - (1-y^{(i)}) \log(1+e^{\theta^T x^{(i)}})$$

$$\text{So } \frac{\partial J(\theta)}{\partial \theta_j}$$

$$= \frac{\partial}{\partial \theta_j} \left[-\frac{1}{m} \sum_{i=1}^m [-y^{(i)} \log(1+e^{-\theta^T x^{(i)}}) - (1-y^{(i)}) \log(1+e^{\theta^T x^{(i)}})] \right]$$

$$= -\frac{1}{m} \sum_{i=1}^m \left[-y^{(i)} \frac{-x_j^{(i)} e^{-\theta^T x^{(i)}}}{1+e^{-\theta^T x^{(i)}}} - (1-y^{(i)}) \frac{x_j^{(i)} e^{\theta^T x^{(i)}}}{1+e^{\theta^T x^{(i)}}} \right]$$

$$= -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \frac{x_j^{(i)}}{1+e^{\theta^T x^{(i)}}} - (1-y^{(i)}) \frac{x_j^{(i)} e^{\theta^T x^{(i)}}}{1+e^{\theta^T x^{(i)}}} \right]$$

$$= -\frac{1}{m} \sum_{i=1}^m \frac{y^{(i)} x_j^{(i)} - x_j^{(i)} e^{\theta^T x^{(i)}} + y^{(i)} x_j^{(i)} e^{\theta^T x^{(i)}}}{1+e^{\theta^T x^{(i)}}}$$

$$= -\frac{1}{m} \sum_{i=1}^m \frac{y^{(i)} (1+e^{\theta^T x^{(i)}}) - e^{\theta^T x^{(i)}}}{1+e^{\theta^T x^{(i)}}} x_j^{(i)}$$

$$= -\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} - \frac{e^{\theta^T x^{(i)}}}{1+e^{\theta^T x^{(i)}}} \right) x_j^{(i)}$$

$$= -\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} - \frac{1}{1+e^{-\theta^T x^{(i)}}} \right) x_j^{(i)}$$

$$= -\frac{1}{m} \sum_{i=1}^m [y^{(i)} - h_{\theta}(x^{(i)})] x_j^{(i)}$$

$$= \frac{1}{m} \sum_{i=1}^m [h_{\theta}(x^{(i)}) - y^{(i)}] x_j^{(i)}$$

∴ Gradient descent と feature scaling

还有其他方法：conjugate Gradient, ..., BFGS, LBFGS

6.5 Simplified Cost Function and Gradient Descent.

Logistic regression cost function:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(h_\theta(x^{(i)}), y^{(i)})$$

$$\text{cost}(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)), & \text{if } y=1 \\ -\log(1-h_\theta(x)), & \text{if } y=0 \end{cases}$$

合并为：

$$\text{cost}(h_\theta(x), y) = -y \log(h_\theta(x)) - (1-y) \log(1-h_\theta(x))$$

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)}))]$$

找出使 $J(\theta)$ 的参数 θ .

通过 gradient descent:

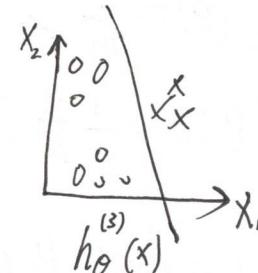
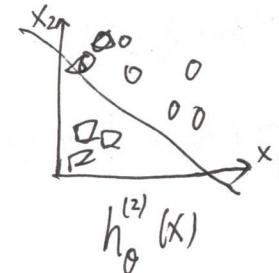
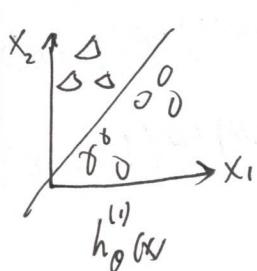
$$\text{Repeat } \left\{ \begin{array}{l} \theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \\ \text{(对所有 } \theta_j \text{)} \end{array} \right\}$$

A vectorized implementation is:

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

6.7 Multi-class Classification: One-vs-all

将多分类中的一类标记为正向类 ($y=1$), 其他所有类标记为负向类:



7. Regularization.

7.1 The problem of overfitting.

过拟合问题:

① 是一些不能帮助理解/理解模型的特征

② 正则化，保留所有特征，但减少其大小 (magnitude)

7.2 Cost Function.

$$\text{模型 } h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2^2 + \theta_3 x_3^3 + \theta_4 x_4^4$$

修改为带惩罚项的：

θ_3, θ_4 可以忽略。

$$\min_{\theta} \frac{1}{2m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + 1000 \theta_3^2 + 10000 \theta_4^2 \right]$$

若特征很多，不知选哪些。

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

λ is regularization parameter.

(梯度不对 θ_0 有影响)

7.3 Regularized Linear Regression.

Repeat until convergence {

$$\begin{aligned}\theta_0 &:= \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)} \\ \theta_j &:= \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \\ &= \theta_j (1 - \alpha \frac{\lambda}{m}) - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}\end{aligned}$$

regularized linear regression by gradient descent 在原有

梯度上加 λ 倍 θ 值减去一个额外的值。

normal eqn.

$$\theta = \left(\underbrace{x^T x + \lambda I}_{(n+1) \times (n+1) 矩阵} \right)^{-1} x^T y$$

$(n+1) \times (n+1)$ 矩阵

7.4 Regularized logistic Regression.

$J(\theta)$ 为 regularization 项

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)}))] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

Gradient Descent:

repeat until converge. {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\begin{aligned}\theta_j &:= \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \\ \text{for } j &= 1, 2, 3, \dots, n\end{aligned}\}$$

Week 4

8. Neural Networks: Representation.

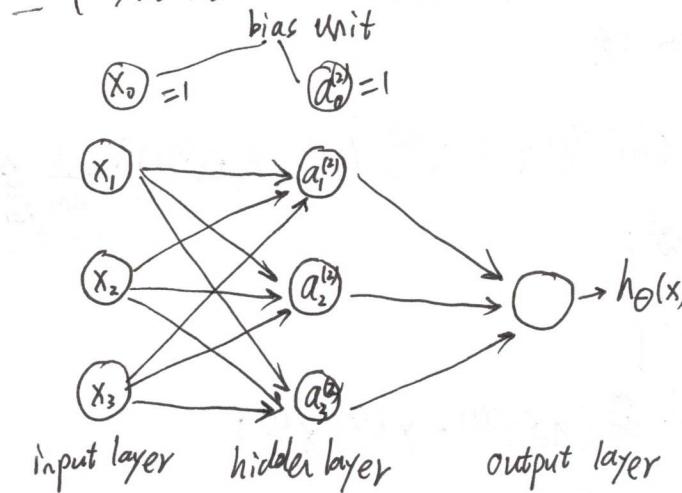
8.1 Non-linear Hypotheses.

不适用

8.2 Neurons and Brain

8.3 Model Representation. I.

一个简单的 Neural Network



$a_i^{(j)}$ 代表第 j 层的第 i + 1 个单元.

$\theta^{(j)}$ 代表从 j 层映射到 j+1 层的权重矩阵

↑ size: $v_h \times t+1$ 行数

v_h j 层为 ~~行数~~ 单元数 + 1 为列数

$$a_1^{(2)} = g(\theta_{10}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2 + \theta_{13}^{(1)} x_3)$$

$$a_2^{(2)} = g(\theta_{20}^{(1)} x_0 + \theta_{21}^{(1)} x_1 + \theta_{22}^{(1)} x_2 + \theta_{23}^{(1)} x_3)$$

$$a_3^{(2)} = g(\theta_{30}^{(1)} x_0 + \theta_{31}^{(1)} x_1 + \theta_{32}^{(1)} x_2 + \theta_{33}^{(1)} x_3)$$

$$h_{\theta}(x) = a_1^{(3)} = g(\theta_{10}^{(2)} a_0^{(2)} + \theta_{11}^{(2)} a_1^{(2)} + \theta_{12}^{(2)} a_2^{(2)} + \theta_{13}^{(2)} a_3^{(2)})$$

从左向右称为 forward propagation.

矩阵乘子 $\theta \cdot X = a$

$$X = \begin{matrix} x_0 \\ x_1 \\ x_2 \\ \vdots \end{matrix}, \quad \theta = \begin{matrix} \theta_{10} & \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots & \theta_{33} \end{matrix}, \quad a = \begin{matrix} a_1 \\ a_2 \\ a_3 \end{matrix}$$

8.4 Model Representation II

用向量计算第二层的值:

$$X = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}, \quad Z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix}, \quad z^{(2)} = \theta^{(1)} X, \quad a^{(2)} = g(z^{(2)})$$

即 $\theta^{(1)} \cdot X = a^{(2)}$

$$g\left(\left[\theta_{10}^{(2)}, \theta_{11}^{(2)}, \theta_{12}^{(2)}, \theta_{13}^{(2)}\right] \times \begin{bmatrix} a_0^{(2)} \\ a_1^{(2)} \\ a_2^{(2)} \\ a_3^{(2)} \end{bmatrix}\right) = \dots = h_{\theta}(x)$$

$$\text{令 } z^{(3)} = \theta^{(2)} \cdot a^{(2)}, \text{ 则 } h_{\theta}(x) = a^{(3)} = g(z^{(3)})$$

这是针对 training set + 1 training unit 的计算, 要算整个 training set,
需将训练集特征矩阵转置, 使同一实例的特征在同一列.

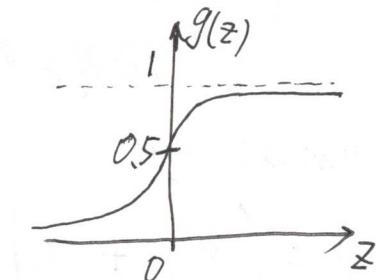
$$Z^{(2)} = \theta^{(1)} \cdot X^T$$

$$a^{(2)} = g(Z^{(2)})$$

8.5 Examples and Intuitions I

神经网络中, 单元神经元的计算可

用来表示逻辑运算, And, Or, ...



8.6 Examples and Intuitions II

Binary logical operations 二进制输入 bool 值 0, 1,
我们只使用一个单一的激活函数作为二元逻辑
运算符，为逻辑网加运符 (and, or, not)，是否
选择不同加权。

看图，理解。

8.7 Multiclass Classification.

Week 5.

9. Neural Networks: Learning.

9.1 Cost Function.

假设训练样本 m 个

每个样本包含一组输入 X , 一组输出 y ,

L 表示神经网络层数,

S_L 表示每层 neuron 个数, S_i 表示 # of units (not contain bias unit) in layer i .

神经网络分为 $\begin{cases} \text{二分类: } S_L = 1, Y = 0 \text{ or } 1 \text{ 表示那一类.} \\ \text{k 类分类: } S_L = k, Y_i = 1 \text{ 表示分到第 } i \text{ 类, } (k > 2) \end{cases}$

cost function in logistic regression.

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m Y^{(i)} \log(h_\theta(x^{(i)})) + (1-Y^{(i)}) \log(1-h_\theta(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

在 logistic regression 中, 只有一个输出变量, 又一个因变量 y .

在 neural network 中, ~~有多个输出~~, $h_\theta(x)$ 是维度为 k 的向量.

$$h_\theta(x) \in \mathbb{R}^k, (h_\theta(x))_i = i^{\text{th}} \text{ output.}$$

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K Y_k^{(i)} \log(h_\theta(x^{(i)}))_k + (1-Y_k^{(i)}) \cdot \log(1-(h_\theta(x^{(i)}))_k) \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} (\theta_{ji}^{(l)})^2$$

对于每一特征, 那么输出 K 个预测。

regularized term 帮助了正则化 θ , 防止 θ 短阵的 go.

随机化时, 随机范围: $[-\epsilon_{int}, \epsilon_{int}]$, $\epsilon_{int} = \frac{\sqrt{6}}{\sqrt{L_{in} + L_{out}}}$, $L_h = S_h$, $L_{out} = S_{out} + 1$

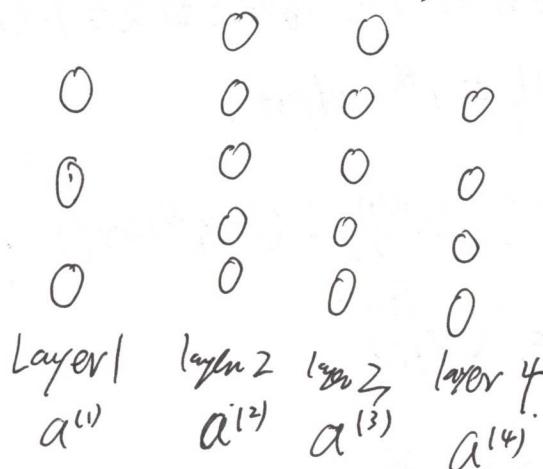
5

9.2 Backward propagation Algorithm

为了计算 $\frac{\partial J(\theta)}{\partial \theta_{ij}^{(l)}}$, 需用反向传播算法.

例子：

只有一个实例 $(x^{(1)}, y^{(1)})$, neural net with \rightarrow 4 层, $K=4$, $S_L=4$, $L=4$



最后一层误差: $\delta^{(4)} = a^{(4)} - y$ 为什么有这一项?

$$\delta^{(3)} = (\theta^{(3)})^T \delta^{(4)} \cdot g'(z^{(3)})$$
 ? 不理解.

其中: $g'(z^{(3)}) = a^{(3)} \cdot (1 - a^{(3)})$, $a^{(3)} = g(z^{(3)})$
 $(\theta^{(3)})^T \delta^{(4)}$ 是权重导致的误差 bag

$$\delta^{(2)} = (\theta^{(2)})^T \delta^{(3)} \cdot g'(z^{(2)})$$

第一层是输入, 无误差

$$g(z) = \frac{1}{1+e^{-z}}, g'(z) = \frac{1}{(1+e^{-z})^2} \cdot (-1) \cdot e^{-z} = \frac{e^{-z}}{(1+e^{-z})^2}$$

$$= \frac{1}{(1+e^{-z})} \cdot \frac{1+e^{-z}-1}{(1+e^{-z})} =$$

$$= g(z)(1-g(z))$$

若 $\lambda=0$, $\frac{\partial J(\theta)}{\partial \theta_{ij}^{(l)}}$ $= a_j^{(l)} \otimes \delta_i^{(l+1)}$

l : 代表当前所在层

j : 当前层 activation unit 下标

i : $l+1$ 层中误差单元的下标.

有很多 training unit.

用 $\Delta_{ij}^{(l)}$ 表示误差矩阵,

第 l 层的第 i 个训练单元对第 j 个输出单位的误差。

算式为:

for $i=1:m$ {

set $a^{(i)} = x^{(i)}$

perform "forward propagation" to compute $a^{(l)}$ for $l=1, 2, \dots, L$

using $\delta^{(l)} = a^{(l)} - y^{(l)}$

perform "backward prop" to compute previous layer

error vector $\Delta_{ij}^{(l)} := \Delta_{ij}^{(l+1)} + a_j^{(l)} \delta_i^{(l+1)}$

找出 $\Delta_{ij}^{(l)}$ 及 $\frac{\partial J}{\partial \theta}$ will be:

$$P_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \theta_{ij}^{(l)} \text{ if } j \neq 0$$

$$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} \text{ if } j = 0$$

9.3 Backpropagation intuition.

看 pdf, p₁₄₆

9.4: Implementation Note - Unshuffling parameters

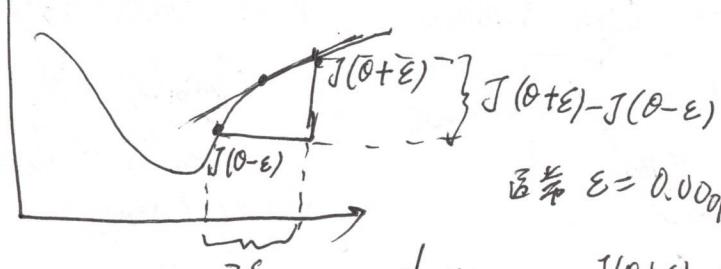
不能选择把参数从矩阵层形式而是

看 pdf, p₁₄₆

9.5 Gradient checking

虽然 $J(\theta)$ 在 \downarrow , 但结果可能不是最优的

用 gradient numerical gradient checking 来做这个梯度值来验证。



$$\frac{d}{d\theta} J(\theta) \approx \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$$

parameter vector θ

$$\theta \in \mathbb{R}^n, \quad \theta = [\theta_1, \theta_2, \dots, \theta_n]$$

$$\frac{\partial}{\partial \theta_i} J(\theta) \approx \frac{J(\theta_1 + \epsilon, \theta_2, \dots, \theta_n) - J(\theta_1 - \epsilon, \theta_2, \dots, \theta_n)}{2\epsilon}$$

~~200~~:

$$\frac{\partial}{\partial \theta_n} J(\theta) \approx \frac{J(\theta_1, \theta_2, \dots, \theta_n + \epsilon) - J(\theta_1, \theta_2, \dots, \theta_n - \epsilon)}{2\epsilon}$$

9.6 Random Initialization.

虽然 logistic Regression 可以直接初始化，但 neural network

不能用 0 初始化，否则第二层所有单元值相同。

故用随机数初始化，范围 $-\epsilon \sim \epsilon$, $-\epsilon \leq \theta_{ij}^{(l)} \leq \epsilon$

9.7 Put it together.

Training a network:

- ① Randomly initialize the weights,
- ② Implement forward propa. to get $h_\theta(x^{(i)})$ for any $x^{(i)}$
- ③ Implement the cost func.
- ④ Implement backpropa. To compute partial derivative
- ⑤ Use gradient checking to confirm that your backpropa. works, then disable gradient checking
- ⑥ Use gradient descent or built-in optimization func to minimize the cost func.

week 6

10 Advice for Applying Machine Learning

10.1 模型误差大，如何修改：

获得更多样本，但代价大，不推荐

- ① 增加特征加权重
- ② 获得更多的特征
- ③ 增加多项式特征
- ④ 增加正则化比例系数
- ⑤ ↑ _____.

10.2 Evaluating a Hypothesis

分成两部分 training set (70%)， test set (30%)
然后重 shuffle。

对线性回归模型，用 test set 算了。

对 logistic regression，用 test set 算了

还能算：

misclassification error:

$$\text{err}(h_\theta(x), y) = \begin{cases} 1, & \text{if } h_\theta(x) \geq 0.5, y=0 \\ & \text{or } h_\theta(x) < 0.5, y=1 \\ 0, & \text{otherwise.} \end{cases}$$

$$\text{Test error} = \frac{1}{M_{\text{test}}} \sum_{i=1}^m \text{err}(h_\theta(x_{\text{test}}^{(i)}), y^{(i)})$$

10.3 Model Selection and Train-Validation-Test sets

60% 为 training set, 20% 为 validation set, 20% 为 test set

- ① 用 training set 训练 10 个 model
- ② 用 10 个 model 和 validation set 算 CV error, (cost function)
- ③ 选出 CV error 最小的模型
- ④ 用 ③ 中选出的 model 对 test set 算测试误差 (cost function)

Training error:

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

Cross validation error:

$$J_{\text{cv}}(\theta) = \frac{1}{2M_{\text{cv}}} \sum_{i=1}^m (h_\theta(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)})^2$$

Test error:

$$J_{\text{test}}(\theta) = \frac{1}{2M_{\text{test}}} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

10.4 Diagnose Bias vs. Variance

under fit ~ high bias

over fit ~ high variance



Bias (under fit)

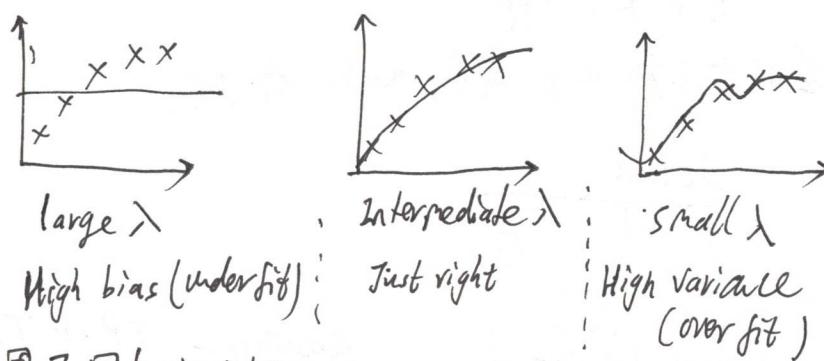
Jtrain(θ) high,
Jtrain(θ) ≈ Jcv(θ)

Variance (overfit)

Jtrain(θ) low

10.5 Regularization and Bias-variance.

要选择合适的 λ .

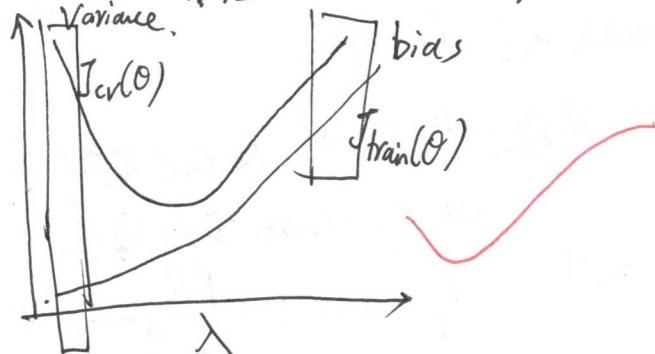


- ① 用不同的 λ 训练出 12 个 model
- ② 12 个 model 算 $J_{cv}(\theta)$
- ③ 选 $J_{cv}(\theta)$ 较低的 model
- ④ 用③中的 model 计算 generalized error $J_{test}(\theta)$.

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2$$

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

$$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_\theta(x_{cv}^{(i)}) - y^{(i)})^2$$

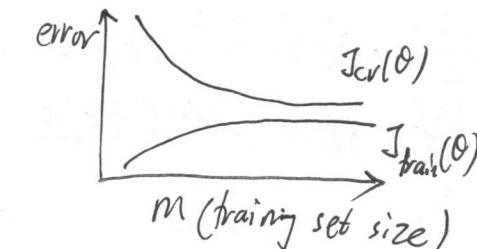


10.6 Learning Curves

learning curve 是很好的合理性检查 (sanity check).

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

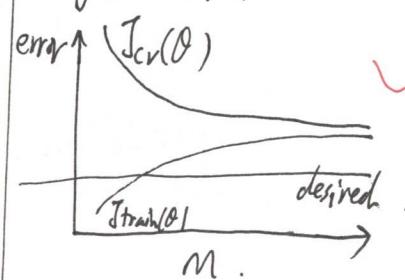
$$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$



在高 high bias (underfit) 的情况下, ↑训练集 无助.

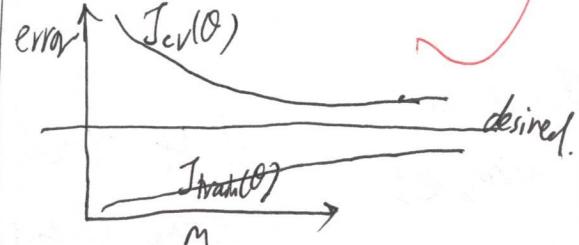
在 High Variance (Overfit) : ~, ↑ ~ 有帮助. ✓

High bias.



Low training set size: $J_{train}(\theta)$ low, $J_{cv}(\theta)$ high
large ~: $J_{train}(\theta)$, $J_{cv}(\theta)$ both are high

High variance.



low training set size: $J_{train}(\theta)$ low, $J_{cv}(\theta)$ high
large ~: $J_{train}(\theta) \uparrow$ as $m \uparrow$
 $J_{cv}(\theta) \downarrow$

10.7 Recapping what to do next

- ① 过拟合案例, $\lambda =$ 高方差 high variance
 ② 偏向性过拟合量 — Variance
 ③ 基于多项式特征 bias
 ④ 增加多项式特征 bias
 ⑤ $\downarrow \lambda$ bias
 ⑥ $\uparrow \lambda$ Variance.

若 neural network: high bias, (underfit)
 太大 neural network: high variance (overfit).

11. Machine Learning System Design.

11.1 prioritizing What to work on.

11.2 · Error Analysis.

- ① 通过简单的能快速发现的算法方法,
 实现之并用 CV, learning curve, 检查
 特征, 更多数据 或者 什么。
 ② 进行 error analysis: · 拉直 CV set 中产生
 误差的原因, L2 正则

11.3 · Error Metrics for Skewed classes.

skew classes: training set 中有非常多的同一类, 很少 or 没有
 其他类的数据。

rest 误差过大不能用单类别 model 来拟合。

		Actual	
		1	0
predict	1	True positive (TP)	False positive (FP)
	0	False negative (FN)	True negative (TN)

$$\text{accuracy} = \frac{TP + TN}{\text{total examples}}$$

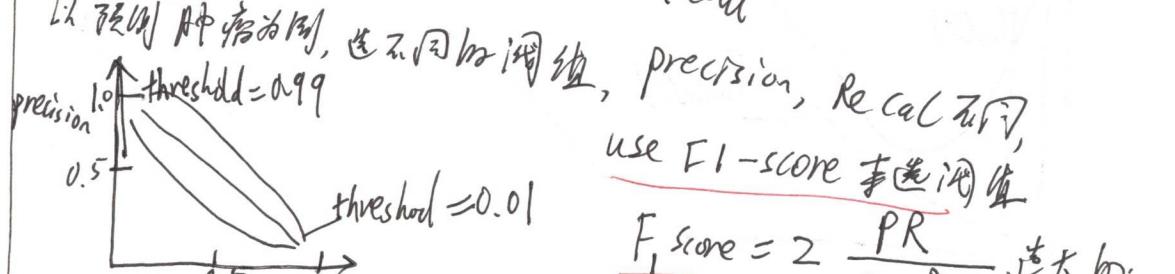
precision: $\frac{TP}{TP+FP}$

of all patients we predict $y=1$,
 what fraction actually has cancer?

Recall: $\frac{TP}{TP+FN}$

of all patients that actually have cancer,
 what fraction did we detect as having cancer?

11.4 Trade off precision and recall



$$F_1 \text{ score} = 2 \cdot \frac{PR}{P+R}$$

11.5 Data for Machine Learning

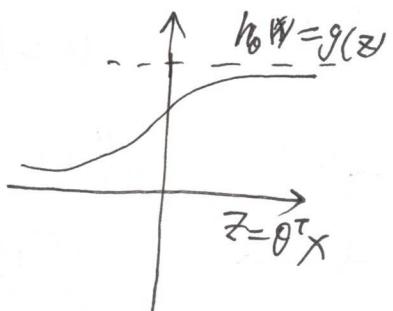
Week 7

12. Support Vector Machine

12.1 Optimization Objective

logistic regression.

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$



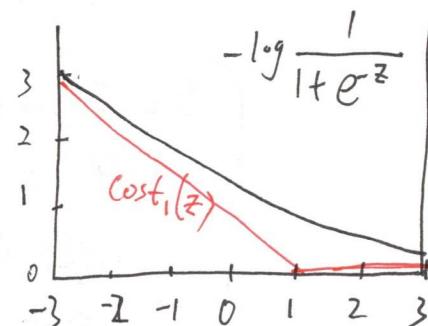
if $y=1$, we want $h_{\theta}(x) \approx 1, \theta^T x \gg 0$

$y=0$ $h_{\theta}(x) \approx 0, \theta^T x \ll 0$

$$\text{cost} : -(y \log h_{\theta}(x)) + (1-y) \log(1-h_{\theta}(x))$$

$$= -y \log \frac{1}{1 + e^{-\theta^T x}} - (1-y) \log \left(1 - \frac{1}{1 + e^{-\theta^T x}}\right)$$

if $y=1, (\theta^T x > 0)$



用线性/凸形代替

logistic regression ..

$$\min_{\theta} \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} (-\log h_{\theta}(x^{(i)})) + (1-y^{(i)}) (-\log (1-h_{\theta}(x^{(i)}))) \right] + \lambda \sum_{j=1}^n \theta_j^2$$

忽略掉常数M.

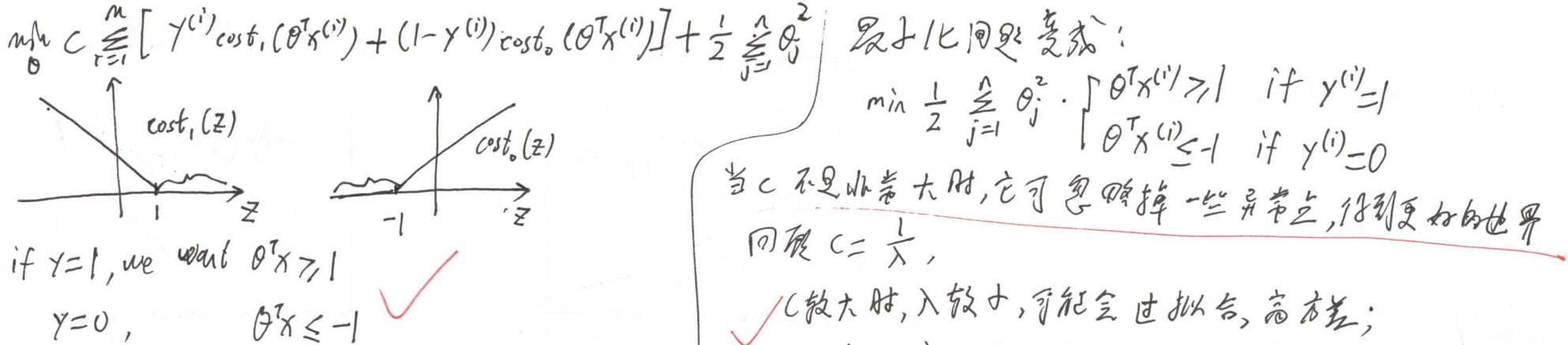
support vector machine:

$$\min_{\theta} C \sum_{i=1}^m [y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)})] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

$$h_{\theta}(x) \begin{cases} 1, & \text{if } \theta^T x \geq 0 \\ 0, & \text{otherwise.} \end{cases}$$

12.2 Large Margin Intuition

$$\min_{\theta} C \sum_{i=1}^m [y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)})] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$



若 C 很大, 在 $\min J(\theta)$ 时, 第一项为 0.

$y=1$ 时, 找一个 θ 使 $\theta^T x \geq 1$

$y=0$ 时, $\theta^T x \leq -1$

这样会有边界:



这个距离叫做 margin (margin)

这个距离叫做 SVM 的 margin, 也叫 SVM 是有 robust 性的, 因其用的是大间隔分离样本。

故 SVM 又称为 大间隔分类器
 (large margin classifier)

最小化问题变成:
 $\min \frac{1}{2} \sum_{j=1}^n \theta_j^2 \cdot \begin{cases} \theta^T x^{(i)} \geq 1 & \text{if } y^{(i)} = 1 \\ \theta^T x^{(i)} \leq -1 & \text{if } y^{(i)} = 0 \end{cases}$

当 C 不是很大时, 它可以忽略掉一些异常点, 得到更好的结果

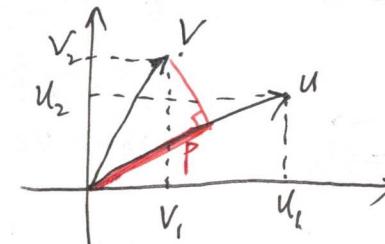
$$\text{因此 } C = \frac{1}{\lambda},$$

C 放大时, 入放大, 可能会过拟合, 高方差;
 C 缩小, 入缩小, \sim 欠拟合, 高 bias.

12.3. Mathematics behind Large margin classifier.

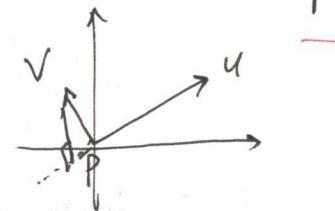
Vector inner product

$$u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}, v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$



$$\|u\| = \text{length of vector } u$$

$$= \sqrt{v_1^2 + u_1^2 + u_2^2}$$



$p = \text{length of the projection of } v \text{ on to } u$

$$u^T v = p \cdot \|u\|$$

$$p \perp u$$

$$= u_1 v_1 + u_2 v_2$$

SVM decision bdry.

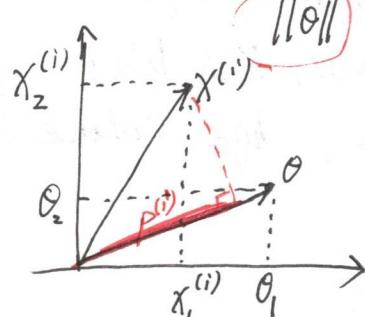
$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

$$\text{s.t. } \theta^T x^{(i)} \geq 1 \quad \text{if } y^{(i)} = 1 \\ \theta^T x^{(i)} \leq -1 \quad \text{if } y^{(i)} = 0$$

e.g. let $\theta_0 = 0, n=2$

$$\frac{1}{2} \sum_{j=1}^2 \theta_j^2 = \frac{1}{2} (\theta_1^2 + \theta_2^2) = \frac{1}{2} (\sqrt{\theta_1^2 + \theta_2^2})^2 = \frac{1}{2} \|\theta\|^2$$

$\oplus \theta^T x^{(i)} = ?$



$$\begin{aligned} \theta^T x^{(i)} &= P^{(i)} \cdot \|\theta\| \\ &= \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} \end{aligned}$$

Now, SVM decision bdry:

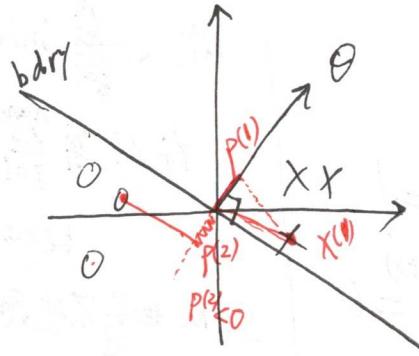
$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

$$\text{s.t. } P^{(i)} \cdot \|\theta\| \geq 1 \quad \text{if } y^{(i)} = 1$$

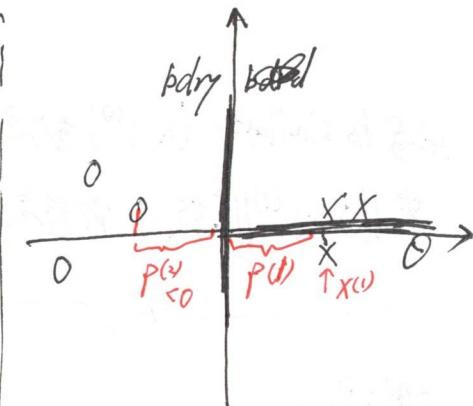
$$P^{(i)} \cdot \|\theta\| \leq -1 \quad \text{if } y^{(i)} = 0$$

$P^{(i)}$ is the project of $x^{(i)}$ onto the vector θ .

simplifying: $\theta_0 = 0 \rightarrow \theta_0 = 0$ decision bdry is θ



$P^{(1)} \cdot \|\theta\| \geq 1$, $P^{(1)}$ is small, so $\|\theta\|$ should be large
 $P^{(2)} \cdot \|\theta\| \leq -1$, $P^{(2)}$ small, so $\|\theta\|$ large
 not good, we want small θ .



$P^{(1)}, P^{(2)}$ larger, θ small

12.4 kernels I

一个多边形边界的模型可能是 $\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 + \theta_5 x_2^2 + \dots$

除了对原有特征组合，有无更好的产生构造 f_1, f_2, f_3, \dots

用 kernels 计算出新的特征。

给定训练集 X, Y 用 X 的各个特征与 Y 的 n 个 landmark $l^{(1)}, l^{(2)}, l^{(3)}, \dots$ 的相似度来构造新特征 f_1, f_2, f_3, \dots

$$\begin{aligned} f_1 &= \text{similarity}(x, l^{(1)}) = e^{-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}} \\ \text{其中 } \|x - l^{(1)}\|^2 &= \sum_{j=1}^n (x_j - l_j^{(1)})^2 \end{aligned}$$

这个是 similarity $(x, l^{(i)})$ 就是 kernels.

若 $\|x - l^{(i)}\|^2$ 较大, PP 距离较大, 则 $e^{-\theta} = 1 - \dots - \dots - \dots - \dots - \dots = 0$

参考: P199

12.5 kernels II.

如何选取地 f?

Given $(x^{(1)}, y^{(1)})$, $(x^{(2)}, y^{(2)})$, ..., $(x^{(m)}, y^{(m)})$

choose $f^{(1)} = x^{(1)}$, $f^{(2)} = x^{(2)}$, ..., ~~$f^{(m)} = x^{(m)}$~~

$$f^{(i)} = \begin{cases} f_0^{(i)} = 1 \\ f_1^{(i)} = \text{sim}(x^{(i)}, l^{(1)}) \\ f_2^{(i)} = \text{sim}(x^{(i)}, l^{(2)}) \\ \vdots \\ f_i^{(i)} = \text{sim}(x^{(i)}, l^{(i)}) = e^{\theta} = 1 \\ \vdots \\ f_m^{(i)} = \text{sim}(x^{(i)}, l^{(m)}) \end{cases}$$

$$\min_{\theta} \sum_{i=1}^m [y^{(i)} \text{cost}_1(\theta^T f^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T f^{(i)})] + \frac{1}{2} \sum_{j=1}^m \theta_j^2$$

在计算 $\sum_{j=1}^m \theta_j^2 = \theta^T \theta$ 时, 用 $\theta^T M \theta$ 替代 $\theta^T \theta$, M 是根据选择的 kernels 不同而不用的矩阵, 为了简化计算.

SVM 也使用 kernel, 称为 linear kernel.

$$C = \frac{1}{\lambda}$$

C 太大, λ 小, overfit, high variance

C 太小, λ 大, underfit, high bias

? C 太大, low variance, high bias

? C 太小, low bias, high variance

12.6 Using SVM.

修改 SVM 假设为:

给定 x , 计算拟合值 f , 当 $\theta^T f \geq 0$ 时, $y=1$, 否则 $y=0$

修改 cost func 为 $\sum_{j=1}^m \theta_j^2 = \theta^T \theta$

Week 8.

13 Clustering

13.1 Unsupervised Learning - Introduction.

13.2 k-Means Algorithm.

假设分成 k 组

- ① 先 k 个随机点为 cluster centroids.
- ② 根据每个 point 与最近的中心关联，
或一类。
- ③ 计算每组的均值，将中心移到平均位置。

用 $\mu^1, \mu^2, \dots, \mu^k$ 表示聚类中心

用 $C^{(1)}, C^{(2)}, \dots, C^{(m)}$ 存储点；
其类别索引

13.3 Optimization Objective

k-means 最优化：最大化 point 与聚类中心距离之和最小。

k-means loss function : (Distortion func.)

$$J(C^{(1)}, C^{(2)}, \dots, C^{(m)}, \mu_1, \dots, \mu_k) = \frac{1}{m} \sum_{i=1}^m \|X^{(i)} - \mu_c^{(i)}\|^2$$

同时是指出除了 $(1, \mu^1, \mu^2, \dots, \mu^k, m, m^2, \dots, m^k)$

13.4 Random Initialization.

$k < m$, 随机选 k 个点作为聚类中心

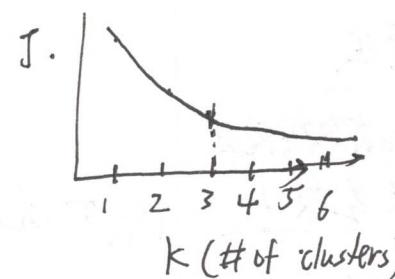
问题：可能停在 local minimum.

解决：多次运行 k-means, 每次重新 random initialization.

+ $O(m \cdot k + t(2 \sim 10))$ 可行, 大了不实用

13.5 Choosing the # of clusters.

选择 k , 用 elbow method.



参考资料 P213.

10

14. Dimensionality Reduction.

14.1 Motivation I — Data compression

14.2 Motivation II — Visualization.

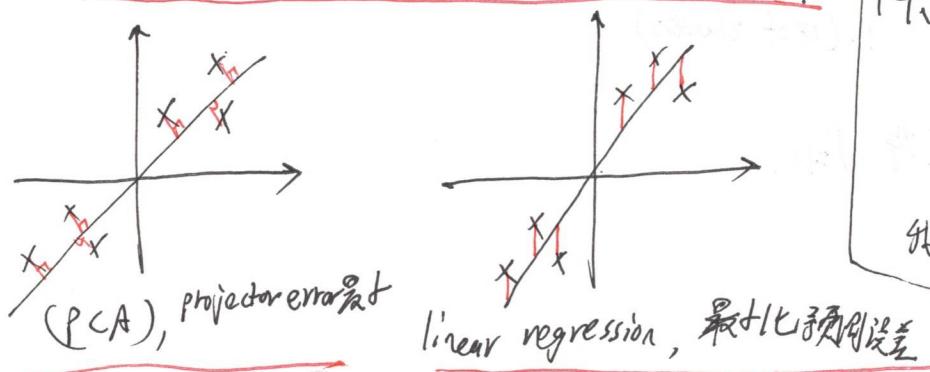
将高维数据降低才方便可视化。

降低维数的数据特征含义由自己定义。

14.3 Principle component analysis problem formulation.

(PCA) PCA是常见降维算法。

在PCA中，找到 vector direction，使所有数据投影到该 vector direction，使平均均方误差最小。



14.4 Principal Component Analysis Algorithm

PCA 算法 n 维到 k 维：

① 均值归一化：减均值 μ_j , $x_j = x_j - \mu_j$, 不同版块用不同 μ_j

② Covariance matrix Σ : $\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)}) (x^{(i)})^T$

③ 计算 covariance matrix 的 eigenvectors

通过 SVD, $[U, S, V] = SVD(\Sigma)$

$$U = \begin{bmatrix} | & | & | & | \\ u_1^{(1)} & u_1^{(2)} & u_1^{(3)} & \dots & u_1^{(n)} \\ | & | & | & & | \end{bmatrix}$$

$n \times n \rightarrow k \times n$, 取前 k 列的 U 中的 k 列

$$\text{向量 } z^{(i)} = U_{\text{reduce}}^T x^{(i)}$$

14.5 choosing the # of principle components

PCA 是通过投影的方式——平均 TS 方误差 (Avg. squared project error)

$$\text{训练集的方差} (\text{Total variation}) : \frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2$$

希望在平均均方误差与训练集方差的比值尽可能小的情况下选择尽可能大的 k 值。

choose k to be smallest. value so that

$$\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01$$

99% of variance is retained

14.6 Reconstruction from compressed representation

$$Z = U_{\text{reduce}}^T X$$

$$X_{\text{approx}} = U_{\text{reduce}} \cdot Z$$

$$X_{\text{approx}} \approx X$$

14.5 (continued)

Algorithm:

(1) Try PCA with $k=1$

(2) compute U_{reduce} , $X^{(1)}, Z^{(1)}, Z^{(2)}, \dots Z^{(m)}$

$$X_{\text{approx}}^{(1)}, X_{\text{approx}}^{(2)}, \dots X_{\text{approx}}^{(m)}$$

(3) check if

$$\frac{\frac{1}{m} \sum_{i=1}^m \|X^{(i)} - X_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|X^{(i)}\|^2} \leq 0.01$$

(4) 若(3)不满足, try $k=2, 3, \dots$ repeat ②~③.

$$[U, S, V] = \text{SVD}(X)$$

$$S = \begin{bmatrix} S_{11} & & & \\ & S_{22} & & \\ & & S_{33} & \\ & & & \ddots & S_{nn} \end{bmatrix}$$

Given k , is $1 - \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \leq 0.01$? or $\frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \geq 0.99$

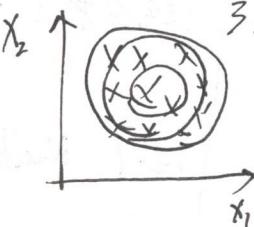
14.7 Advice for Applying PCA.

week 9.

15. Anomaly Detection

15.1 problem Motivation.

What's anomaly detection: 有一组数据 $X^{(1)}, (X^{(2)}, \dots, X^{(n)})$, 将其数据与已有 m 组数据比较, 是否异常?



若在圆内, 则属于该组的纯度较高.

作为密度估计, $P(x) \begin{cases} < \epsilon, \text{ anomaly} \\ \geq \epsilon, \text{ normal} \end{cases}$

$P(x)$ 为 x 属于另一组数据的频率

15.2 Gaussian Distribution.

若 X 遵循正态分布, $X \sim N(\mu, \sigma^2)$

probability density func. $P(X, \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

$$\mu = \frac{1}{m} \sum_{i=1}^m X^{(i)}, \quad \sigma^2 = \frac{1}{m} \sum_{i=1}^m (X^{(i)} - \mu)^2$$

15.3 Algorithm.

异常检测算法:

对给定数据集 $X^{(1)}, X^{(2)}, \dots, X^{(m)}$, 对每一个特征
计算 μ_j, σ_j^2 .

$$\mu_j = \frac{1}{m} \sum_{i=1}^m X_j^{(i)}$$

$$\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (X_j^{(i)} - \mu_j)^2$$

给定训练集，计算 $P(x)$

$$P(x) = \prod_{j=1}^n P(x_j; \mu_j; \sigma_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi\sigma_j^2}} e^{-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}}$$

$P(x) < \epsilon$ 时为异常

15.4 Developing and Evaluating an Anomaly Detection System

unsupervised

从带标记的 (正常, 异常) 数据着手, 选一部分正常数据
作 training set, 用剩下的正常和异常数据混合构成
交叉 CV set, Test set.

- ① 根据 training set, 估计 μ_j, σ_j^2 , 计算 $P(x)$ 及其 precision, recall 来选择 ϵ
- ② 对 CV set, 尝试不同 ϵ 作阈值, 根据 F1 score, recall 来选择 ϵ

③ 选出 ϵ , 对 Test set, 计算 F1 score, precision, recall.

15.5 Anomaly Detection vs. Supervised Learning

见 P232

15.6 Choosing what features to use

Anomaly detection assume features 服从 Gauss distribution

根据 不服从 Gauss dist. 及 logit Gauss dist.

如 $x = \log(x + c)$, $x = x^c$ 等方法

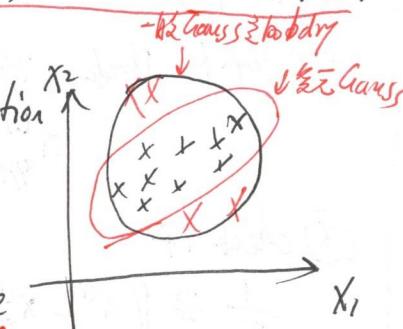
误差分析: 有些异常值可能也有高 $P(x)$ 值, 可以通过特征筛选

15.7 Multivariate Gaussian Distribution

一维 Gauss, 即 $P(x) = \text{计算各特征 } p(x) \text{ 并累乘}$

multivariate Gauss: 指由特征 covariance

matrix, 用所有特征一起算 $P(x)$.



所有特征的均值: $\mu = \frac{1}{m} \sum_{i=1}^m X^{(i)} \rightarrow \mu \in \mathbb{R}^n, n \times 1$

covariance matrix: $\Sigma = \frac{1}{m} \sum_{i=1}^m (X^{(i)} - \mu)(X^{(i)} - \mu)^T = \frac{1}{m} (X - \mu)^T (X - \mu)$

$$P(x) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left[-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right]$$

↑ determinant of Σ

15.8 Anomaly Detection using the Multivariate Gaussian dist.

$$P(x; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left[-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) \right]$$

Fit model ($P(x)$) by setting

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)(x^{(i)} - \mu)^T$$

Given a new example x , compute

$$P(x) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left[-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) \right]$$

flag an anomaly if $P(x) < \epsilon$

16. Recommender System

16.1 Problem Formulation

- 一个系统对用户推荐商品，

n_u : 用户数

n_m : 商品数

$r(i, j)$: 用户 j 给商品 i 的评分, 则 $r(i, j) = 1$

$y^{(i, j)}$: 用户 j 对商品 i 的评分

m_j : 用户 j 评过分的商品数

16.2 Content-based Recommendations

分析物品特征, $x_1 \sim$ 浏览程度, $x_2 \sim$ action 程度 / 用 linear regression

$\theta^{(j)}$: 用户 j 特征向量

$x^{(i)}$: 商品 i 的特征向量

对用户 j 喜欢 i , 预测评分 $(\theta^{(j)})^T x^{(i)}$

针对用户 j 的 cost function:

$$\min_{\theta^{(j)}} \frac{1}{2} \sum_{i: r(i, j)=1} [(\theta^{(j)})^T x^{(i)} - y^{(i, j)}]^2 + \frac{\lambda}{2} (\theta_k^{(j)})^2$$

所有用户:

$$\min_{\theta^{(1)}, \dots, \theta^{(n_u)}} \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i: r(i, j)=1} [(\theta^{(j)})^T x^{(i)} - y^{(i, j)}]^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2$$

用 Gradient Descent,

$$\check{\theta}_k^{(j)} := \theta_k^{(j)} - \alpha \sum_{i: r(i, j)=1} [(\theta^{(j)})^T x^{(i)} - y^{(i, j)}] x_k^{(i)}, \text{ for } k=0$$

$$\check{\theta}_k^{(j)} := \theta_k^{(j)} - \alpha \left(\sum_{i: r(i, j)=1} [(\theta^{(j)})^T x^{(i)} - y^{(i, j)}] x_k^{(i)} + \lambda \theta_k^{(j)} \right), \text{ for } k \neq 0$$

16.3 Collaborative Filtering

反过来, 若拥有用户的历史数据, 可以预测商品的特征。

$$\min_{x^{(1)}, \dots, x^{(n_m)}} \frac{1}{2} \sum_{i=1}^{n_m} \sum_{j: r(i, j)=1} [(\theta^{(j)})^T x^{(i)} - y^{(i, j)}]^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2$$

collaborative filtering 通过用户对电影的评分，预测评分。

$$J(x^{(i)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}) \\ = \frac{1}{2} \sum_{(i,j) : r(i,j)=1} [(\theta^{(i)})^T x^{(i)} - y^{(i,j)}]^2 + \frac{\lambda}{2} \sum_{j=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2 \\ + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2$$

梯度：

$$x_k^{(i)} := x_k^{(i)} - \alpha \left(\sum_{j: r(i,j)=1} ((\theta^{(i)})^T x^{(i)} - y^{(i,j)}) \theta_k^{(i)} + \lambda x_k^{(i)} \right)$$

$$\theta_k^{(i)} := \theta_k^{(i)} - \alpha \left(\sum_{j: r(i,j)=1} ((\theta^{(i)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} + \lambda \theta_k^{(i)} \right)$$

Collaborative filtering algorithm

- ① Initialize $x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}$ to small random values.
- ② Minimize $J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)})$ use gradient descent
- ③ For user θ , movie x , predict a star rating $\theta^T x$.

16.5 Vectorization — Low Rank Matrix Factorization

5 movies, 4 users, 5星评价。

$$\begin{matrix} (\theta^{(1)})^T x^{(1)} & (\theta^{(2)})^T x^{(1)} & \dots & (\theta^{(n_u)})^T x^{(1)} \\ (\theta^{(1)})^T x^{(2)} & (\theta^{(2)})^T x^{(2)} & \dots & (\theta^{(n_u)})^T x^{(2)} \\ \vdots & \vdots & & \vdots \\ (\theta^{(1)})^T x^{(n_m)} & (\theta^{(2)})^T x^{(n_m)} & \dots & (\theta^{(n_u)})^T x^{(n_m)} \end{matrix}$$

For each product $x^{(i)}$, learn a feature vector $x^{(i)} \in \mathbb{R}^n$.
How to find movie j related to movie i ?

small $\|x^{(i)} - x^{(j)}\| \rightarrow$ movie i, j are similar.

16.6 Implementational Detail — Mean Normalization

首先对矩阵 mean normalization; ~~每行去均值~~ 每一用户去均值

用 collaborative filtering model,

使用 $\theta^T x + b$ @ $\theta^T x$, $(\theta^{(i)})^T x^{(i)} + m_i$.

Week 10

17. Large Scale Machine Learning

17.1 Learn with Large Datasets.

17.2 Stochastic Gradient Descent

要大规模训练模型，用 SGD stochastic gradient descent
代替批量梯度下降法.

In SGD: cost func:

$$\text{cost}(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2} (h_\theta(x^{(i)}) - y^{(i)})^2$$

SGD: 先对训练集随机打乱顺序，再

Repeat (usually anywhere between 1~10)

```
{ for i=1:m  
  { θ := θ_j - α(h_θ(x^{(i)}) - y^{(i)})x_j^{(i)}  
    (for j=0:n)  
  }  
}
```

缺点：可能无法到最优点，在其附近徘徊。

17.3 Mini-Batch Gradient Descent

Mini-batch GD 是介于 Batch GD 与 SGD 之间的算法，每训练 b 次训练实例，更新一次 θ 。 $b \in [2, 200]$

```
repeat {  
  for i=1:m  
  { θ := θ_j - α  $\frac{1}{b} \sum_{k=i}^{i+b-1} (h_θ(x^{(k)}) - y^{(k)})x_j^{(k)}$   
    (for j=0:n)  
  }  
  i += 10  
}
```

17.4 Stochastic GD convergence 在 SGD 中，每次只计算代价函数值，

收敛不平且不下降：↑α 使之平缓。

仍振荡且不下降：model 有问题

or 逐渐上升， $\alpha \downarrow \alpha$

也可令 α 随 iteration ↑ 而 ↓

$$\alpha = \frac{\text{const 1}}{\# \text{ of iteration} + \text{const 2}}$$

17.5 Online Learning

对网站中一个个用户，给出喜好预测 $P(Y=1)$ 。

对单一查询进行学习：

Repeat forever (as long as the website is running)

```
{ get (x, y) corresponding to the current user  
  θ := θ_j - α(h_θ(x) - y)x_j  
  (for j=0:n)  
}
```

17.6 Mapping reduce and Data parallelism.

18. Application Example: photo OCR

18.1 problem description and pipeline.

image → text detection → character segmentation

→ character recognition.

18.2 Sliding windows — 滑动窗口

18.3 Getting lots of data and artificial data.

获得更多数据的方法：

- ① 人工标注
- ② 手动收集，标记数据
- ③ 红包

8.4 Ceiling Analysis - what part of the pipeline
to work on Next

前一步输入 100% 正确的话，下一步是否 Accuracy 提升，说明这一步值得花精力提升

复习记录：(复习 notes, assignment)

2020.8.25 复习 Notes.