

Machine learning

1-1 概念

A computer program is said to learn from experience E with respect to some task T and some performance measure P , if its introduction its performance on T , as measured by P , improves with experience E .

最常使用) $\begin{cases} \text{supervised learning} \\ \text{-unsupervised learning} \\ \text{others. Reinforcement learning, recommender systems.} \end{cases}$

1-2

Supervised learning

"right answers" given → Regression: Predict continuous valued output.

Classification (分类): Discrete valued output (0 or 1)

在监督学习中，对于数据集中的每个样本

我们想要算法预测，并给出“正确答案”

like 良性?

不是只有
两类

1-3 unsupervised learning

• 所用数据没有任何标签 (相同或不同)

• 可以将数据分成不同的簇 (聚类算法) → UL中一种

应用: Organizing computing clusters, Social network analysis, market segmentation, Astronomical data analysis, etc.

• T 鸡尾酒会算法 找出数据结构. → 处理/分离音频

$$[w, s, v] = \frac{\text{sum}(\text{repmax}(\text{sum}(x.^* x.^* 1), \text{size}(x.^* 1).1).^* x)^* x'}{\text{size}(x.^* 1).1};$$

奇异值分解

1.2 - Linear regression with one variable

1.1 模型描述

• Linear regression with one variable

Training set of
housing prices

Size in feet² (x)

214
1616
1534
852
⋮

Price in 1000's (y)

460
232
315
178
⋮

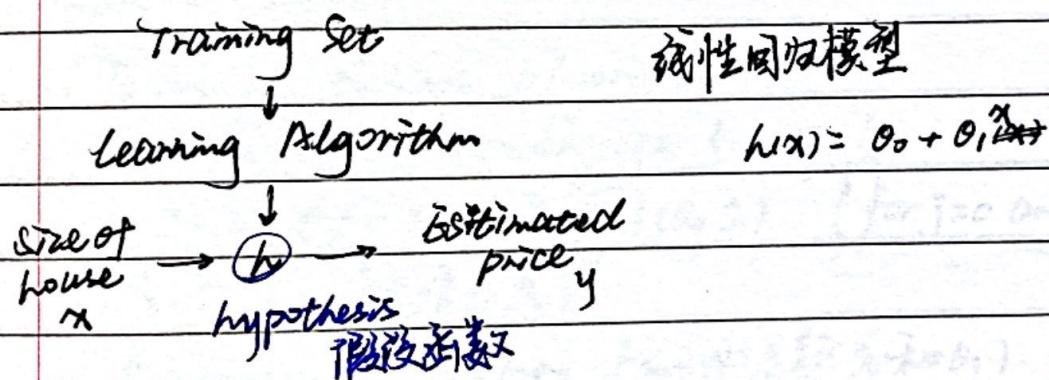
} $m = 10$

m = number of training examples

$(x^{(i)}, y^{(i)}) \rightarrow$ i th training example

x 's = "input" variable / feature

y 's = "output" variable / "target" variable.



2.2 代价函数 为了让实际和预测更加接近

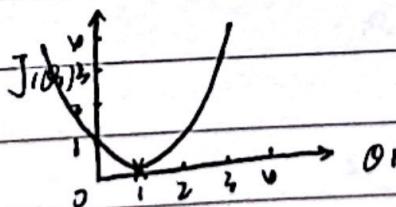
Hypothesis: $h_{\theta}(x) = \theta_0 + \theta_1 x$ Parameters: θ_0, θ_1

最小二乘法 Cost function: $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (\lambda_{\theta}(x^{(i)}) - y^{(i)})^2$

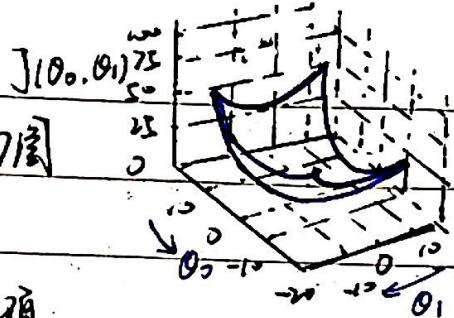
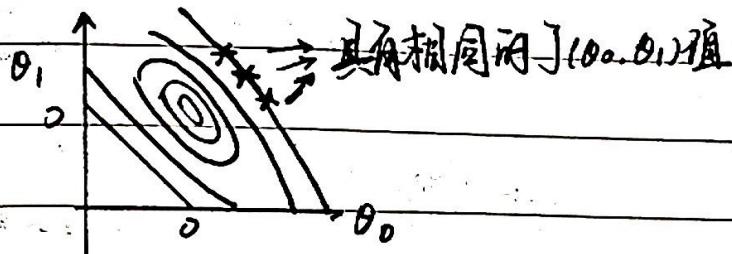
Goal \rightarrow minimize $J(\theta_0, \theta_1)$

↓

Simplified: $\theta_0 = 0 \Rightarrow J(\theta_1) \& \text{minimize } J(\theta_1)$



→ 保留 θ_0 和 θ_1 : 该制则为 → 三维曲面图
→ 继续为等高线图



2-3 梯度下降

Outline:

- Start with some θ_0, θ_1

- Keeping changing θ_0, θ_1 to reduce $J(\theta_0, \theta_1)$ until we hopefully end up at a minimum.

- Gradient descent algorithm

repeat until convergence } → 等价于

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad (\text{for } j=0 \text{ and } j=1)$$

} 梯度运算符

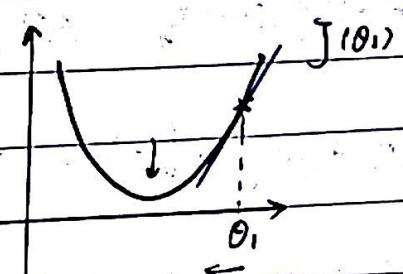
Simultaneous update (同时对更新 θ_0 和 θ_1)

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

$$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{tempo}$$

$$\theta_1 := \text{temp1}$$



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

$$\theta_1 := \theta_1 - \alpha \cdot (\text{positive number})$$

冲线性回归的梯度下降

$$\begin{aligned}\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) &= \frac{\partial}{\partial \theta_j} \cdot \frac{1}{m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2 \\ &= \frac{\partial}{\partial \theta_j} \cdot \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2\end{aligned}$$

\Rightarrow repeat until convergence ^{收敛}

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

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

}

"Batch" Gradient Descent

"Batch": Each step of gradient descent uses all the training examples.

3. Linear Regression with multiple variables

3.1 Multiple features

Size x_1	number of bedroom x_2	number of floors x_3	x_4 Age of home	Price y	$n=4$ number of features
2104	5	1	65	460	
1616	3	2	30	232	$m=10$
1534	3	2	30	315	
852	2	1	36	178	
...	

$x^{(i)} = \text{input of } i^{\text{th}} \text{ training example} \rightarrow x^{(1)} = \begin{bmatrix} 2104 \\ 5 \\ 1 \\ 65 \end{bmatrix}$

$x_j^{(i)} = \text{value of feature } j \text{ in } i^{\text{th}} \text{ training example} \rightarrow x_3^{(1)} = 1$

$\star \Rightarrow$ hypothesis: $h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$

线性回归
假设函数
形式

$$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}, \quad \theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix} \Rightarrow h(x) = \theta^T x$$

"multivariate
linear regression"
多变量线性回归

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

3.2 多元梯度下降法 New algorithm ($n \geq 1$):

Repeat ?

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

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

(Simultaneously update θ_j for $j = 0, \dots, n$)

}

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

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

...

3.3 Gradient descent in practice 1: feature Scaling 特征缩放

Idea: Make sure features are on a similar scale

→ 梯度下降法能更快收敛

E.g. $x_1 = \text{size } (0 - 2000 \text{ m}^2) \Rightarrow x_1 = \frac{\text{size}}{2000}$

$x_2 = \text{number of bedrooms } (1-5) \Rightarrow x_2 = \frac{n \text{ of bedrooms}}{5}$

so. $0 \leq x_i \leq 1$

• Feature Scaling

Get every feature into approximately a $-1 \leq x_i \leq 1$ range.

• Mean normalization 均值归一化

Replace x_i with $\tilde{x}_i = \frac{x_i - \mu_i}{\sigma_i}$ to make features have approximately zero mean (Do not apply to x_0)

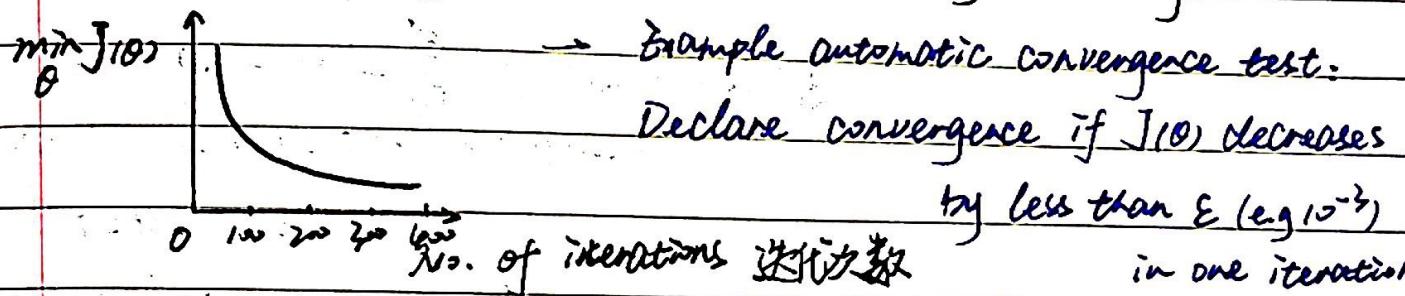
E.g. $x_1 = \frac{\text{size} - 1200}{2000} \rightarrow -0.5 \leq x_1 \leq 0.5$

$x_2 = \frac{\# \text{bedrooms} - 2}{5} \rightarrow -0.5 \leq x_2 \leq 0.5$

3.4 Gradient descent in practice II: learning rate

$$\text{Gradient descent } \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

- Making sure gradient descent is working correctly



- Gradient descent not working → Use smaller α .
 - for sufficiently small α , $J(\theta)$ should ~~decrease~~ increase on every iteration
 - But if α is too small, gradient descent can be slow to converge

3.5 Features and polynomial regression

- 选择 features 并如何使用? ⇒ 引入新的特征来获得更好的模型

Idea I

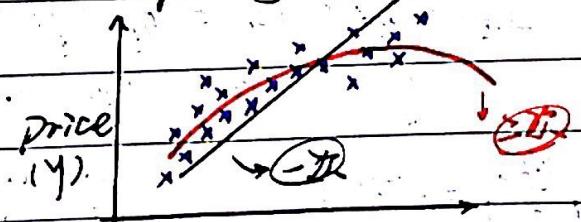
E.g. Housing prices prediction

$$h_0(x) = \theta_0 + \theta_1 \times \text{frontage} + \theta_2 \times \text{depth}$$

frontage
深度
前方宽度

还可以选择将两个特征合并, 即 $\text{Area} = \text{frontage} \times \text{depth}$

$$\Rightarrow h_0(x) = \theta_0 + \theta_1 \times \text{Area}$$



Idea II

线性回归 ⇔ 非线性回归

E.g. 对于 size ↔ price. 如果直线不能很好地拟合

→ 考虑选择二次函数 (会产生凹形趋势)

→ 考虑三次函数

$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 = \theta_0 + \theta_1 (\text{size}) + \theta_2 (\text{size})^2 + \theta_3 (\text{size})^3$$

x_1 x_2 x_3

$$h_0(x) = \theta_0 + \theta_1 (\text{size}) + \theta_2 \sqrt{\text{size}}$$

(尝试用特征缩放来降低梯度)

3.6 Normal equation 正规方程 (区别于迭代方法的直接解法)

- 未归一化解，一步得到
- 原理：

$$\text{目标函数 } J(\theta_0, \dots, \theta_m) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \dots = 0 \quad (\text{for every } i)$$

Solve for $\theta_0, \dots, \theta_m$.

Example

	x_0	x_1	x_2	x_3	x_4	y
对数	1	2.04	5	-1	4.5	6.66
辐射	1	1.616	3	2	1.0	2.32
特征	1	1.534	3	2	3.0	3.15
质量	1	8.52	2	1	3.6	1.78

$$X = \begin{bmatrix} 1 & 2.04 & 5 & -1 & 4.5 \\ 1 & 1.616 & 3 & 2 & 1.0 \\ 1 & 1.534 & 3 & 2 & 3.0 \\ 1 & 8.52 & 2 & 1 & 3.6 \end{bmatrix}$$

(矩阵运算中涉及小二乘法)

求解公式

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

$$y = \begin{bmatrix} -1.66 \\ 2.32 \\ 3.15 \\ 1.78 \end{bmatrix}$$

$m \times (n+1)$ 矩阵

for m example, n features

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \rightarrow X = \begin{bmatrix} \cdots (X^{(1)})^T \cdots \\ \vdots \\ \cdots (X^{(m)})^T \cdots \end{bmatrix}_{m \times (n+1)}$$

(design matrix)

Gradient Descent

- Need to choose α
- Needs many iterations
- Works well even when n is large

Normal Equation

- No need to choose α
- No need to iterate
- Need to compute $(X^T X)^{-1}$
- Slow if n is very large

3.7 Normal equation and non-invertibility

What if $X^T X$ is non-invertible?

- Redundant features (linearly dependent)

E.g. $x_1 = \text{Size in feet}^2$

$x_2 = \text{Size in m}^2$

- Too many features (e.g. $m \leq n$)

- Delete some features, or use regularization

4. Logistic Regression

4.1 Classification 分类

動取性의 $y \in \{0, 1\}$ 二分类

使用 \rightarrow Threshold classifier output $h_{\theta}(x)$ at 0.5.

线性回归 If $h_{\theta}(x) \leq 0.5$, predict "y=0"

(λ -推荐) If $h_{\theta}(x) > 0.5$, predict "y=1"

\rightarrow 新算法 logistic regression

特点: 预测值一直在 $[0, 1]$ 之间 ($0 \leq h_{\theta}(x) \leq 1$)

4.2 Hypothesis Representation 假设表示

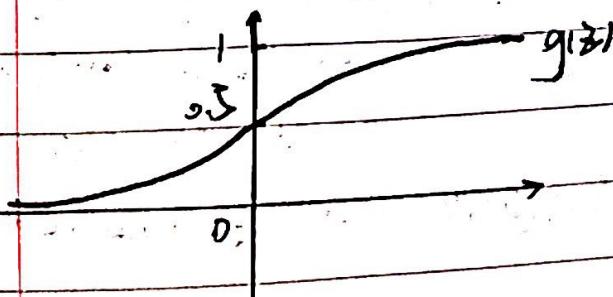
logistic Regression model

sigmoid function
logistic

want $0 \leq h_{\theta}(x) \leq 1$

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

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



Interpretation of hypothesis output

$h_{\theta}(x)$ = estimated probability that $y=1$ on input x 痘症

Example

If $x = [x_0 \ x_1] = [\text{tumor size}]$

$h_{\theta}(x) = 0.7 \rightarrow$ 痘症为 x , $y=1$ 时的概率.

→ Tell patient that 70% chance of tumor being malignant

4-3 Decision boundary 決策邊界 \rightarrow 假設函數的 T 屬性

$h_{\theta}(x) = g(\theta^T x) \rightarrow P(y=1 | x, \theta)$ \rightarrow 痘症為 x 的 $y=1$ 機率估計

Suppose predict

" $y=1$ " if $h_{\theta}(x) \geq 0.5 \rightarrow \theta^T x \geq 0$

" $y=0$ " if $h_{\theta}(x) < 0.5 \rightarrow \theta^T x < 0$ non-convex

4-4 Cost function

对于 linear regression $J(\theta) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$ \uparrow 損失

$\rightarrow \text{Cost}(h_{\theta}(x) - y) = \frac{1}{2} (h_{\theta}(x) - y)^2$ 范圍 $\frac{1}{2}$ 乘以

对于 logistic regression

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

if $y=1$

$\text{Cost} = 0$ if $y=1, h_{\theta}(x)=1$

But as $h_{\theta}(x) \rightarrow 0$

$\text{Cost} \rightarrow \infty$

\uparrow convex $J(\theta)$

$h_{\theta}(x)$

Captures intuition that if $h_{\theta}(x)=0$,
but $y=1$, we'll penalize learning algorithm
by a very large cost.

If $y=0$

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

But as $h_\theta(x) \rightarrow 1$

Cost $\rightarrow \infty$

4-5 Simplified Cost function and gradient descent

logistic regression cost function (交叉熵损失函数 Cross-entropy)

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

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

Note: $y=0$ or 1 always

$$\text{损失函数} \Rightarrow -y \log(h_\theta(x)) - (1-y) \log(1-h_\theta(x)) \Rightarrow \text{Cost}(h_\theta(x), y)$$

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

To fit parameter θ :

$$\min_{\theta} J(\theta) \rightarrow \text{Cost}(\theta)$$

To make a prediction given new x :

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

• Gradient Descent

want $\min_{\theta} J(\theta)$:

Repeat {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad | \text{ Simultaneously update all } \theta_j$$

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

($\ln \leftrightarrow \log$)

Logistic regression 二分类
"3-2" 离线上标注不同

11-6 Advanced optimization 高級最優化

Given, we have code that can compute:

- $J(\theta)$
- $\frac{\partial}{\partial \theta_j} J(\theta)$

Optimization algorithms

- Conjugate gradient
- BFGS 共變梯度法
- LBFGS

Advantages:

- no need to manually pick α
- often faster than gradient descent

Disadvantages:

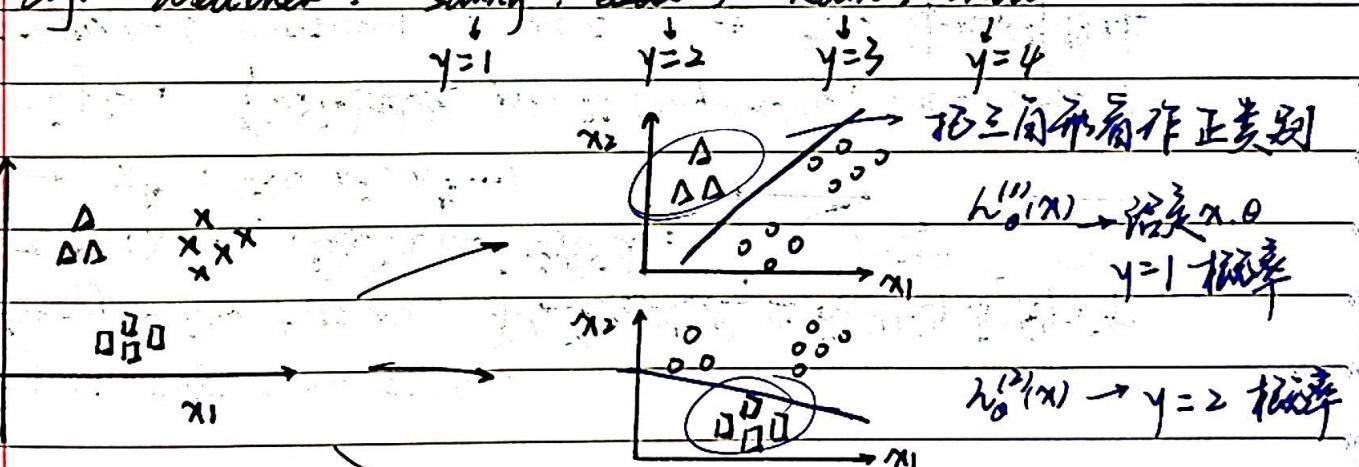
- More complex

4-7 Multi-class classification: One-vs-all. 多類別

E.g. Weather: Sunny, Cloudy, Rain, Snow

$$y=1 \quad y=2 \quad y=3 \quad y=4$$

原理

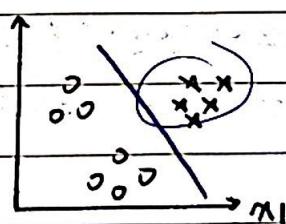


Class: Δ

Class: □

Class: X

$$h_0^{(i)}(x) = P(y=i/x; \theta)$$



$$h_0^{(3)}(x)$$

输出三个分类器

Train a logistic regression classifier $h_0^{(i)}(x)$ for each class i to predict the probability that $y=i$.

On a new input x , to make a prediction, pick the class i that maximizes $\max_{\theta_0}^{(i)}(x)$.

4.5 Regularization

5-1 The problem of overfitting 过拟合问题

- 没有很好拟合，具有高偏差 → 欠拟合 具有高偏差
- 假设函数几乎能拟合所有数据 → 系数过大，变量太多无法约束 → 拟合训练集，导致无法泛化到新样本中 → 过拟合 通常在变量太多时

解决方法

- Reduce the number of feature

• 人工检查变量清单 → 重要 → 保留

• 模型选择算法 → 缺点：舍弃变量但未舍弃信息

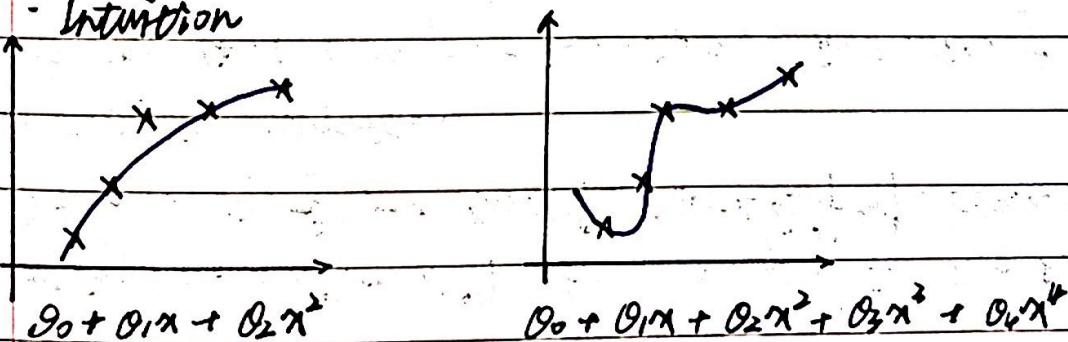
- Regularization 正则化

• Keep all the features, but reduce magnitude / values of parameters θ_j

• Works well when we have a lot of features, each of which contributes a bit to predicting y .

5-2 Cost function

- Intuition



• Suppose we generalize and make θ_3, θ_4 very small.

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (\hat{y}_i - y_i)^2 + 1000\theta_3^2 + 1000\theta_4^2$$

为了让上式尽可能小 $\Rightarrow \theta_3, \theta_4$ 接近 0.

正则化思想 - Small values for parameters $\theta_0, \theta_1, \dots, \theta_n$

• Simpler hypothesis

• Less prone to overfitting

函数更平滑

如果令 $\theta_0 \sim \theta_n$ 都加上惩罚项 \rightarrow 满足斜率

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

因为不加惩罚项

加 θ_0 与 θ_n 影响不大

系数，添加一个额外

一般不对 θ_0 正则化

正则化项以 \downarrow 每个参数值

• 将 $J(\theta)$ 分为二项式

• 第一项：为了更好地拟合训练集

• 第二项：将参数值变小

入：正则化参数，控制以上两项间平衡关系

入过大，会导致 $\theta_1 \sim \theta_n \rightarrow 0$ ，只剩 θ_0 (僵化)

5-3 Regularized linear regression 线性回归的正则化

• Gradient descent

Repeat ?

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

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

↓

$$\rightarrow \theta_j := \theta_j \left(1 - \alpha \frac{\lambda}{m} \right) - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

由于又偏 λ $\Rightarrow \alpha \frac{\lambda}{m}$ 是一个很小的数

$\Rightarrow (1 - \alpha \frac{\lambda}{m})$ 接近于 1

也就是说每次迭代时都将 θ_j 乘以一个较小的数
然后进行和之前一样更新操作

$$x^{(m)} = [x_0^{(m)} \ x_1^{(m)} \ \dots \ x_n^{(m)}]$$

- Normal equation

$$X = \begin{bmatrix} (x^{(1)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix}_{m \times n+1} \quad y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix} \in \mathbb{R}^m$$

$$\text{求 } \min_{\theta} J(\theta) \rightarrow \theta = (X^T X + \lambda \begin{bmatrix} 0 & 1 & \dots & 1 \end{bmatrix})^{-1} X^T y$$

Regularized

5-4 logistic regression

logistic 回归的正规化

- Cost function

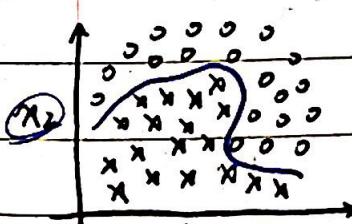
$$J(\theta) = - \left[\frac{1}{m} \sum_{j=1}^m y^{(j)} \log h_{\theta}(x^{(j)}) + (1-y^{(j)}) \log (1-h_{\theta}(x^{(j)})) \right] + \frac{\lambda}{m} \sum_{j=1}^n \theta_j^2$$

与梯度下降法 5-3 中 Gradient descent 一致。但各向 $h_{\theta}(x)$ 表达式不同

6.6 Neural Networks: Representation (表示)

6-1 Non-linear hypotheses 非线性假设

Non-linear classification



可构造包含很多非线性项的 logistic 回归函数

$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 x_2 \\ + \theta_5 x_1^3 x_2 + \theta_6 x_1 x_2^2 + \dots)$$

④ 但当初始特征个数 n 很大时，将高阶多项式项数
包括到特征里，会使特征空间急剧膨胀

6-2 Neurons and brain 神经元和大脑

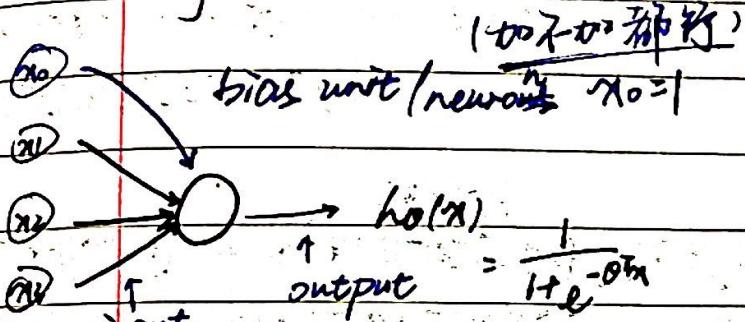
Origins: Algorithms that try to mimic the brain

6-3 模型表示 [Model representation]

输出: 编码

动作定位 \leftarrow 通道
树突: 强入通道

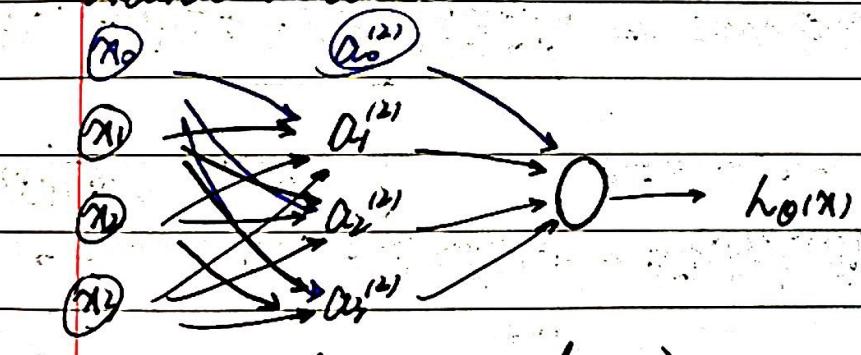
• Logistic unit



$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}, \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

权重 "weights" or 系数
parameters

• Neural Network



layer 1 layer 2 layer 3
Input layer hidden layer Output layer

$\theta^{(j)}$ = 权重矩阵 控制着
layer-j 到 layer-j+1
的映射

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

$$a_1^{(3)} = g(\theta_{10}^{(2)} a_0^{(2)} + \theta_{11}^{(2)} a_1^{(2)} + \theta_{12}^{(2)} a_2^{(2)} + \theta_{13}^{(2)} a_3^{(2)}) = h_0(x)$$

$\rightarrow s_j$ units in layer j
 s_{j+1} units in layer j+1

$$\theta^{(j)} \rightarrow s_{j+1} \times (s_j + 1)$$

6-4 模型展示 II

def $a_1^{(2)} = g(z_1^{(2)})$, $a_2^{(2)} = g(z_2^{(2)})$, $a_3^{(2)} = g(z_3^{(2)})$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}, z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix}, z^{(2)} = \theta^{(1)}x$$

$$\rightarrow a^{(3)} = g(z^{(3)}) \rightarrow z^{(3)} = \theta^{(2)} a^{(2)} \rightarrow h_{\theta}(x) = a^{(3)} = g(z^{(3)})$$

- 上述过程称为前向传播

$$\# h_{\theta}(x) = g(\theta_{00}^{(2)} a_0^{(2)} + \theta_{01}^{(2)} a_1^{(2)} + \theta_{02}^{(2)} a_2^{(2)} + \theta_{03}^{(2)} a_3^{(2)})$$

和 logistic regression 类似，但 neural network 没有用输入特征 x_1, x_2, x_3 来训练逻辑回归。

而是自己训练逻辑回归的输入 a_1, a_2, a_3

\Rightarrow 为 $\theta^{(1)}$ 选择不同参数，可以得到一个更好的假设函数

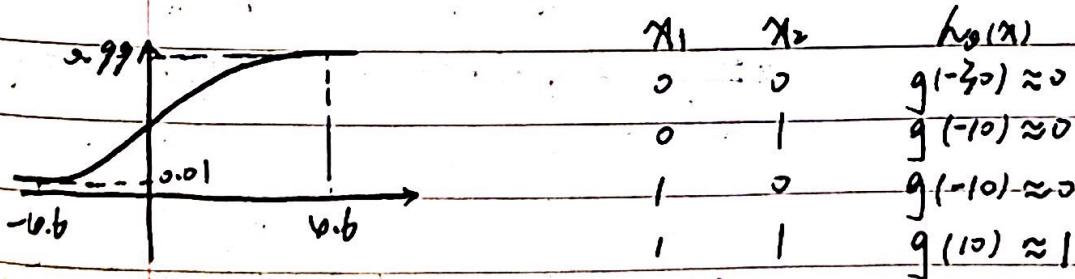
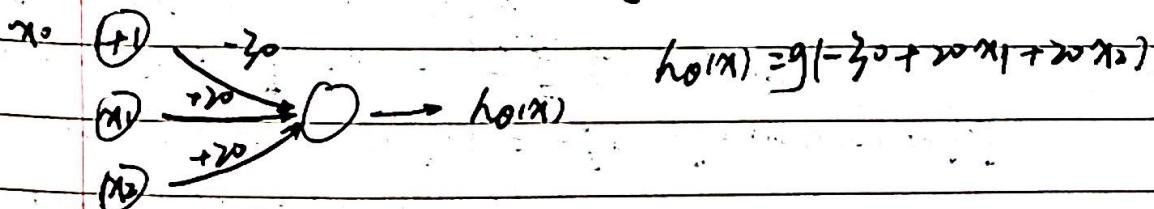
相对于直接使用原始特征 x_1, x_2, x_3 。

6-5 Examples and intuitions 例子和直觉理解 I (单个神经元如何被用来)

· 异或 XOR: 同为0(假), 异为1(真) 计算逻辑函数

· 同或 NOTOR: 同真异假

- AND $x_1, x_2 \in \{0, 1\}$, $y = x_1 \text{ AND } x_2$



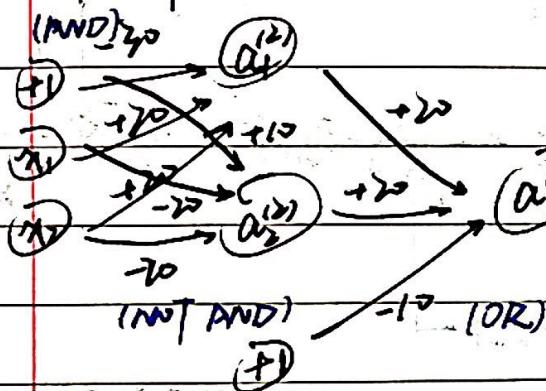
	x_1	x_2	$h_\theta(x)$
- OR	0	0	$g(-10) \approx 0$
	1	0	$g(+10) \approx 1$
	0	1	$g(10) \approx 1$
	1	1	$g(30) \approx 1$

	x_1	$h_\theta(x)$
- NOT	0	$g(10) \approx 1$
	1	$g(-10) \approx 0$

- $(\text{NOT } x_1) \text{ AND } (\text{NOT } x_2) \rightarrow g(10 - 20x_1 - 20x_2)$

6-6 例子和直覺理解 II

- 构造 $x_1 \text{ XNOR } x_2$ (putting AND, OR, (NOT AND) together)



x_1	x_2	$\alpha_1^{(12)}$	$\alpha_2^{(12)}$	$h_\theta(x)$
0	0	0	1	1
0	1	0	0	0
1	0	0	0	0
1	1	1	0	1

6-7 多输出层

- Multiple output units: one-vs-all

$$[N_o \cdots N_{oo}] \rightarrow h_\theta(x) \in \mathbb{R}^n$$

want $h_\theta(x) \approx \begin{bmatrix} 1 \\ 0 \end{bmatrix}$, $h_\theta(x) \approx \begin{bmatrix} 0 \\ 1 \end{bmatrix}$, $h_\theta(x) \approx \begin{bmatrix} 0 \\ 0 \end{bmatrix}$, etc.

when pedestrian when car when motorcycle

Training Set $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}) \dots (x^{(m)}, y^{(m)})$

$$y^{(i)} \text{ one of } \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix} \Rightarrow h_\theta(x^{(i)}) \approx y^{(i)}$$

pedestrian car motor truck

L7: Neural Networks: learning

7-1 Cost function

- ## • Neural Network (Classification)

$\{(x^{(1)}, y^{(1)}), \dots, (x^{(L)}, y^{(L)})\}$
 } $L = \text{total no. of layers}$
 } $S_l = \text{no. of units (no bias)}$

- ## • Binary classification

~~$y=0 \text{ or } 1$~~ . 1 output unit, $h_0(x)GR$, $S_k > 1$ ($k=1$)

- Multi-class classification (K class)

$y \in R^k$ (预测向量), k output unit, $\log(R^k)$, $S_L = k$.

- ## - Cost function

- ## • logistic regression:

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

- ## Neural network:

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_\theta(x^{(i)})_k + (1-y_k^{(i)}) \log(1-h_\theta(x^{(i)})_k) \right]$$

$$+ \frac{\lambda}{2m} \sum_{l=1}^{L-1} \left| \sum_{i=1}^m \sum_{j=1}^{J_l} (\theta^{(l)}_{ij})^2 \right|$$

. 不对 bias, $\hat{\theta}^{(l)} = 0$ 未加

· 最后一篇大千输出单元

~~银行一年和至所有年数
干 连接】~~

7-2 Backpropagation algorithm 反向传播算法

- Need to compute $\rightarrow J(\theta)$

$$\rightarrow \frac{\partial}{\partial \theta^{(i)}_{(t)}} J(\theta)$$

- forward propagation ij (e.g. layer = 4)

$$\alpha^{(1)} = x \quad , \quad \beta^{(2)} = \theta_1^{(1)} \neq \alpha^{(1)}$$

$$\alpha^{(2)} = g(\beta^{(2)}) \text{ (add } \alpha_0^{(2)}\text{)}, \quad \beta^{(3)} = \theta^{(2)} \alpha^{(2)}$$

$$\alpha^{(3)} = q(\beta^{(2)}) \quad (\text{add } \alpha_0^{(3)}), \quad \beta^{(4)} = \theta^{(3)} \alpha^{(3)}$$

$$\alpha^{(4)} = h_0(x) = g(z^{(4)})$$

- Backpropagation

Intuition: $s_j^{(l)} = \text{"error" of node } j \text{ in layer } l$
 for each output unit (layer $L=4$)

$$s_j^{(4)} = a_j^{(4)} - y_j \rightarrow \text{向量表示} s^{(4)} = a^{(4)} - y$$

$$\rightarrow s^{(3)} = (\theta^{(3)})^T s^{(4)} * g'(z^{(3)})$$

$$\rightarrow s^{(2)} = (\theta^{(2)})^T s^{(3)} * g'(z^{(2)}) \rightarrow a^{(2)} * (1 - a^{(2)})$$

no $s^{(1)}$

激活函数偏导数

→ 第二层 (layer 1) 不包含误差

in layer 1
~~in layer 2~~

$$\Rightarrow \text{可简单证明 } \frac{\partial}{\partial \theta_{ij}^{(4)}} J(\theta) \Rightarrow a_j^{(4)} s_i^{(4)}$$

(ignore λ_j)

- Backpropagation algorithm

Training Set $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$

Set $\Delta_{ij}^{(l)} = 0$ (for all l, i, j)

for $i = 1$ to m

Set $a^{(1)} = x^{(i)}$

perform forward propagation to compute $a^{(l)}$ for $l=2, 3, \dots, L$

using $y^{(i)}$, compute $s^{(L)} = a^{(L)} - y^{(i)}$

Compute $s^{(L-1)}, s^{(L-2)}, \dots, s^{(2)}$

$$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} s_i^{(l+1)} \quad \text{? 误差项} \quad \Delta^{(l)} := \Delta^{(l)} + s^{(l+1)} (a^{(l)})^T$$

$$\text{计算梯度矩阵} \quad D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \theta_{ij}^{(l)} \quad \text{if } j \neq 0 \quad \rightarrow \frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = D_{ij}^{(l)}$$

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

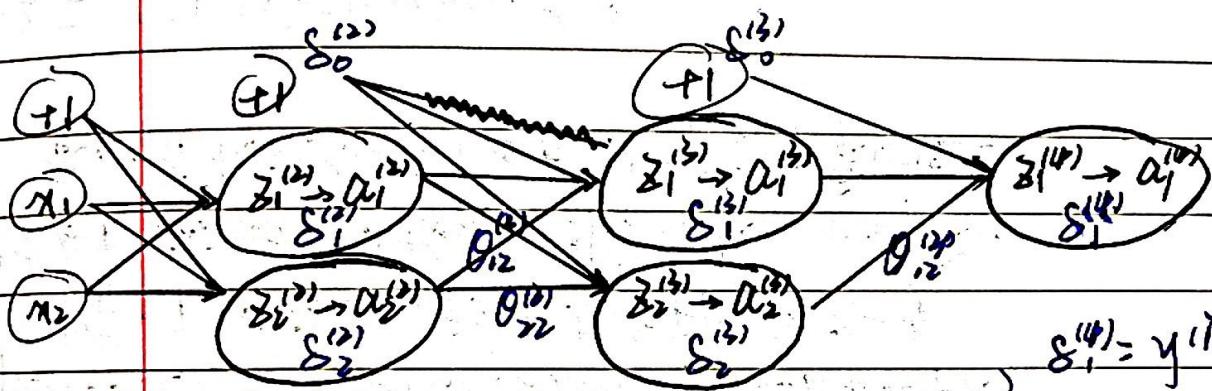
然后可以使用梯度下降或其变体。

7.3 Backpropagation intuition 逐层反向传播

- focusing on a single example $x^{(i)}, y^{(i)}$, the case of 1 output unit and ignoring regularization ($\lambda=0$)

$$\text{Cost}(i) = y^{(i)} \log \text{hol}(x^{(i)}) + (1-y^{(i)}) \log \text{hol}(1-x^{(i)})$$

(Think of $\text{Cost}(i) \approx (\text{hol}(x^{(i)}) - y^{(i)})^2$)



Formally, $\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{Cost}(i)$ (for $j \geq 0$)

$$\left. \begin{aligned} \delta_1^{(4)} &= y^{(i)} - \alpha_1^{(4)} \\ \delta_2^{(3)} &\rightarrow \theta_{12}^{(3)} \delta_1^{(4)} \\ \delta_2^{(2)} &\rightarrow \theta_{22}^{(2)} \delta_1^{(3)} + \theta_{12}^{(2)} \delta_2^{(3)} \end{aligned} \right\}$$

7.4 使用注意：展开参数 Implementation note: Unrolling parameters

- Advanced optimization of θ into $\theta_1, \theta_2, \theta_3$ to n/m gradient to n/m vector.

- Neural Network ($L=4$) → 展开为矩阵

参数矩阵 $\theta_1^{(1)}, \theta_2^{(2)}, \theta_3^{(3)}$ - matrices ($\Theta_1, \Theta_2, \Theta_3$)

(返回值) 梯度矩阵 $D_1^{(1)}, D_2^{(2)}, D_3^{(3)}$ - matrices (D_1, D_2, D_3)

如何展开为向量?
(unroll)

7.5 Gradient Checking 梯度检测 → 确保反向传播
→ 避免算法产生 bug

· 原理 (双侧差分)

$$\frac{d}{d\theta} J(\theta) \approx \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon} \quad \epsilon \text{ 很小} \rightarrow \text{e.g. } \epsilon = 10^{-4}$$

上述 θ 为实数

· Parameter vector θ , $\theta \in \mathbb{R}^n$ (e.g. θ is "unrolled" version of $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$)

$$\theta = [\theta_1, \dots, \theta_n]$$

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

$$\frac{\partial}{\partial \theta_n} J(\theta) \approx \frac{J(\theta_1, \theta_2, \dots, \theta_{n-1} + \epsilon) - J(\theta_1, \theta_2, \dots, \theta_{n-1} - \epsilon)}{2\epsilon}$$

在算法中可利用 for 循环得到上述所有参数偏导数 gradApprox
然后与反向传播中得到的梯度进行比较

一旦确定 gradApprox 和 Dvec 近似相等，在开始学习/训练
网络前关闭“梯度检测”。

→ · "gradApprox" 计算量大，运算速度慢

· "Dvec" 高效

7.6 Random initialization 随机初始化 T 可考 Deep learning 21 Note
1.1 节上

- 如何对 θ 设初值？

· 做假零 initialization, 即 $\theta_{ij}^{(l)} = 0$ for all l, i, j

→ After each update, parameters corresponding to inputs
going into each of all hidden units are identical.

→ 对称权重问题，即所有权重都一样

Solution

random initialization

只有一个特征向量
随机初始化
署名

- random initialization $\xrightarrow{10}$ symmetry breaking 对称性破坏
Initialize each θ_{ij} to a random value in $[-E, E]$

Summarize, 为了训练神经网络，应首先将权重随机初始化为一个 接近 0 的 范围在 $-E$ 到 E 之间的数，然后进行反向传播。
用梯度检测，最后使用梯度下降（或其高级优化算法）

7.1 Putting it together

• Training a neural network

I. - pick a network architecture (connectivity pattern between neurons)

\rightarrow No. of input units: Dimension of features $x^{(i)}$

\rightarrow No. of output unit: Number of classes \rightarrow like

- Reasonable default:

$$y = \begin{bmatrix} 1 \\ \vdots \\ 0 \end{bmatrix} \text{ or } \begin{bmatrix} 1 \\ \vdots \\ 0 \end{bmatrix} \text{ or } \begin{bmatrix} ? \\ \vdots \\ ? \end{bmatrix}$$

1 hidden layer, or if > 1 hidden layer, have same no. of hidden units in every layer (usually the more the better)

And hidden units no. \gg No. of input units

II. 步骤

1. randomly initialize weights (从上至下 summarize)

由于 $J(\theta)$ 是非凸函数 \rightarrow 2. Implement forward propagation to get $h_\theta(x^{(i)})$ for any $x^{(i)}$

3. Implement code to compute cost function $J(\theta)$

4. Implement backprop to compute partial derivatives $\frac{\partial}{\partial \theta^{(k)}} J(\theta)$

5. Using gradient checking (W 7-5)

6. Using gradient descent or advanced optimization method with backpropagation to try to minimize $J(\theta)$ as a function of parameters θ .