

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 performance on T , as measured by P , improves with experience E .

1.1

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 \rightarrow Regression: Predict continuous valued output.

Classification (分类): Discrete valued outputs (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”
 \downarrow
like 子是否只有
肿瘤恶性/良性?)

1-3 unsupervised learning

(相同或没有)

· 所用数据没有任何标签

· 可以将数据分成不同的簇 (聚类算法) \rightarrow UL中一种

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

· T 马鹿函数计算法上 找出数据结构. \rightarrow 处理/分离音频

$$[W, S, V] = \text{svd}(\text{repmat}(\text{sum}(X.^* X, 1), \text{size}(X, 1), 1).^* X)^* X';$$

\downarrow
奇异值分解

1.2 - Linear regression with one variable

2.1 模型描述

- Linear regression with one variable

Training set of housing prices	Size in feet ² (x)	Price in 1000's (y)
	2104	1600
	1616	132
	1534	315
	852	178
	:	:

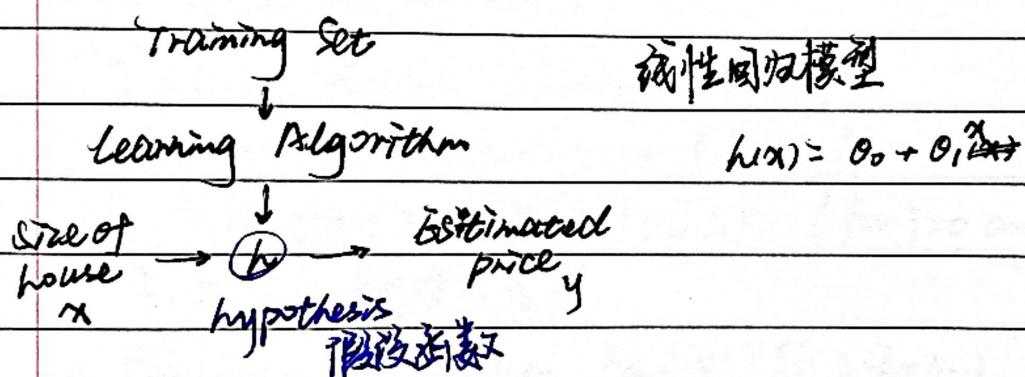
$\left. \begin{matrix} \\ \\ \\ \\ \end{matrix} \right\} m = 107$

m = number of training examples

x 's = "input" variable / feature

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

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



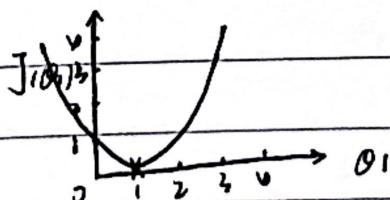
2.2 代价函数 为了让实际和预测更加拟合

Hypothesis: $h_0(x) = \theta_0 + \theta_1 x$ Parameters: θ_0, θ_1

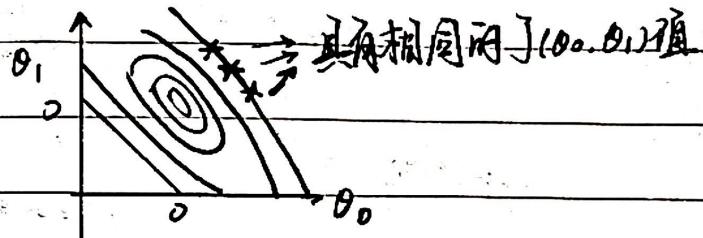
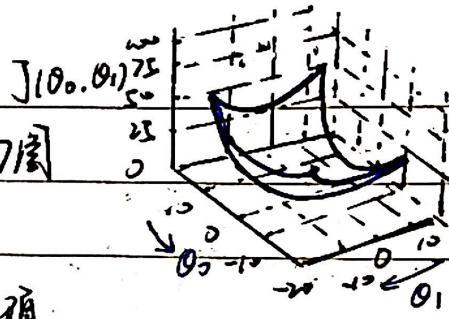
最小二乘法? Cost function: $J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$

Goal \rightarrow minimize $J(\theta_0, \theta_1)$
 (θ_0, θ_1)

↓
 Simplified: $\theta_0 = 0 \Rightarrow J(\theta_1)$ & minimize $J(\theta_1)$



→ 保留 θ_0 和 θ_1 , 限制则为 → 三维曲面图
→ 维度为等高线图



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

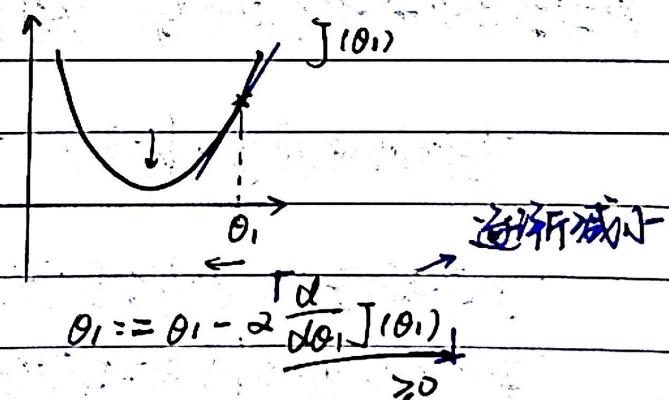
Simultaneous update (同时更新 θ_0 和 θ_1)

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

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

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

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



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

冲销性回归的梯度下降

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

→ 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.

3.3 Linear Regression with multiple variables

3.3.1 Multiple features

Size x_1	number of x_2 bedroom	number x_3 of floors	x_4	Age of home	Price y	$x_1 \sim x_4$	$n=4$ — number of features
210b	5	1	65	460			
161b	3	2	60	232			$m=10$
1534	3	2	30	315			
852	2	1	36	178			
...			

$x^{(i)}$ = input of i th training example

$$x^{(i)} = \begin{bmatrix} 1 \\ x_1^{(i)} \\ x_2^{(i)} \\ x_3^{(i)} \end{bmatrix}$$

$x_j^{(i)}$ = value of feature j in i th training example

$$x_3^{(2)} = 2$$

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

假定回归
假设假设
形式

define $x_0 = 1$ $\Rightarrow h(x) = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n$

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

多元线性回归

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

3-2 多元梯度下降法 New algorithm ($n \geq 1$):

Repeat ?

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

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

}

$$\rightarrow \theta_0 := \theta_0 - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)} \right] = 1 \cdot \dots$$

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

...

3-3 Gradient descent in Practice 1: feature Scaling 特征缩放

Idea: Make sure features are on a similar scale

→ 梯度下降法就稳定性快收敛

E.g. $x_1 = \text{size } (0 - 200 \text{ m}^2)$ $\Rightarrow x_1 = \frac{\text{size}}{200}$
 $x_2 = \text{number of bedrooms } (1-5)$ $\Rightarrow x_2 = \frac{\text{n of bedrooms}}{5}$

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

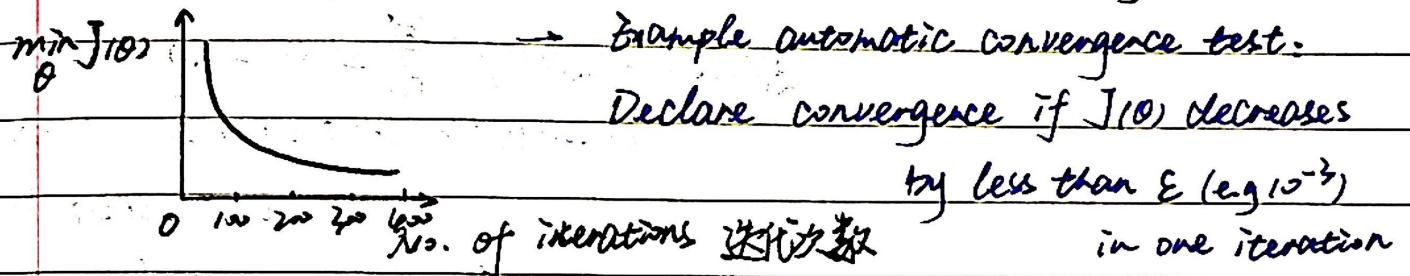
Replace x_i with $\tilde{x}_i = \frac{x_i - \mu_i}{\sigma_i}$ to make features have approximately zero mean (Do not apply to x_0)

E.g. $x_1 = \frac{\text{size} - 100}{200} \rightarrow -0.5 \leq x_1 \leq 0.5$
 $x_2 = \frac{\# \text{bedrooms} - 2}{5} \rightarrow -0.5 \leq x_2 \leq 0.5$

2.4 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 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? ⇒ 使用新的特征来获得更好的模型

Idea I

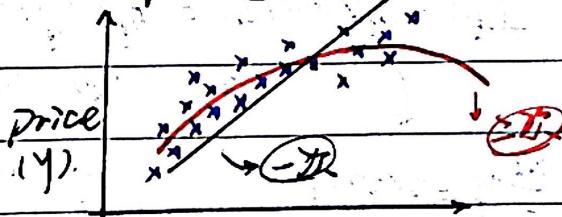
E.g. Housing prices prediction

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

房屋宽度

但可以選擇将两个特征合并，即 $\text{Area} = \text{frontage} \times \text{depth}$

$$\Rightarrow h_\theta(x) = \theta_0 + \theta_1 x$$



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 (会出现下降趋势)

⇒ 考虑三次函数

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

or

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

(要使用特征缩放来下序)

3-6 Normal equation 正规方程(区别于迭代方法的直接解法)

- 未用解析解，一步得到
- 原理：

$$\text{OG} \in \mathbb{R}^{m \times m} \quad J(\theta_0, \dots, \theta_m) = \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

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

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

Example

$m \times (n+1)$ 个数

	x_0	x_1	x_2	x_3	x_4	y	
对数	1	2.4	5	-1	4.5	6.5	
统计	1	1.6	3	2	1.0	2.3	
特征	1	1.5	3	2	2.0	3.2	
将征	1	8.5	2	1	3.0	3.5	
变量	1	8.5	2	1	3.0	1.78	

(矩阵分析中讲的小二乘法)

求解公式)

$$y = \begin{bmatrix} -4.6 \\ 2.32 \\ 3.15 \\ 1.78 \end{bmatrix}$$

$m \times 1$

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

• for m example, n features

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

(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

4-1 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. $m \geq n$)

- Delete some features, or use regularization

4-2 Logistic Regression

4-1 Classification 分类 预测性问 y ∈ {0, 1} 二分类

使用 → threshold classifier output $h_{\theta}(x)$ at 0.5:

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

(逻辑回归) If $h_{\theta}(x) > 0.5$, predict "y=1"

→ 新算法 logistic regression

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

4-2 Hypothesis Representation 假设表达式

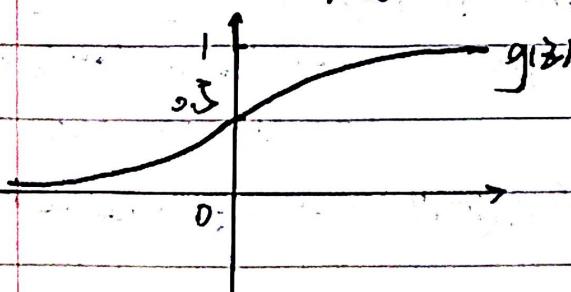
logistic Regression model

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

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

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

sigmoid function
logistic



• Interpretation of Hypothesis Output

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

Example

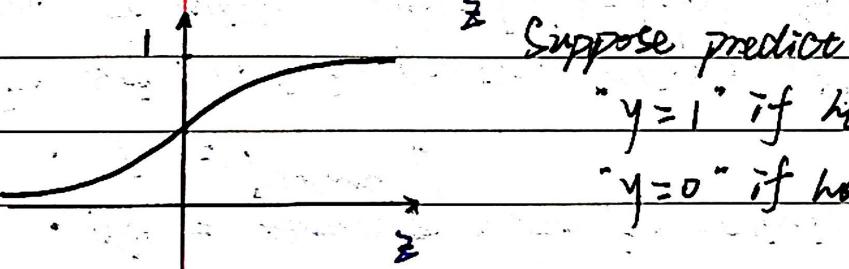
$$\text{if } x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumor-size} \end{bmatrix}$$

$h_{\theta}(x) = 0.7 \rightarrow$ 特征为 x , $y=1$ 的概率.

→ Tell patient that 70% chance of tumor being malignant

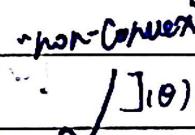
4-3 Decision boundary 決策边界 \rightarrow 从 $\hat{y} = \theta^T x + b$ 得到 $\hat{y} = 0$

$$h_{\theta}(x) = g(\theta^T x) \rightarrow P(y=1/x, \theta) \quad y=1 \text{ 概率估计}$$



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

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



4-4 Cost Function

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

$$\rightarrow \text{Cost}(h_{\theta}(x) - y) = \frac{1}{2} (h_{\theta}(x) - y)^2$$
 范例：最小二乘法

对于 logistic regression

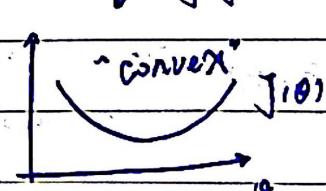
$$\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)), & \text{if } y=1 \\ -\log(1-h_{\theta}(x)), & \text{if } y=0 \end{cases}$$
 ↑ 损失函数

if $y=1$

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

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

Cost $\rightarrow \infty$



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 $\rightarrow \infty$

$h_\theta(x)$

4-5 Simplified Cost function and gradient descent

logistic regression cost function [交叉熵损失函数 Cross-entropy]

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

本质上也是一种对数似然函数

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

note: $y=0$ or 1 always ★ loss(error) function

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

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

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 {

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

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

($\ln \rightarrow \log$)

Logistic regression 逻辑回归
"j=2" 实际上差异不同

4-6 Advanced optimization

Given, we have code that can compute:

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

Optimization algorithms

- Conjugate gradient
- BFGS
- L-BFGS } 共轭梯度法

Advantages:

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

Disadvantages:

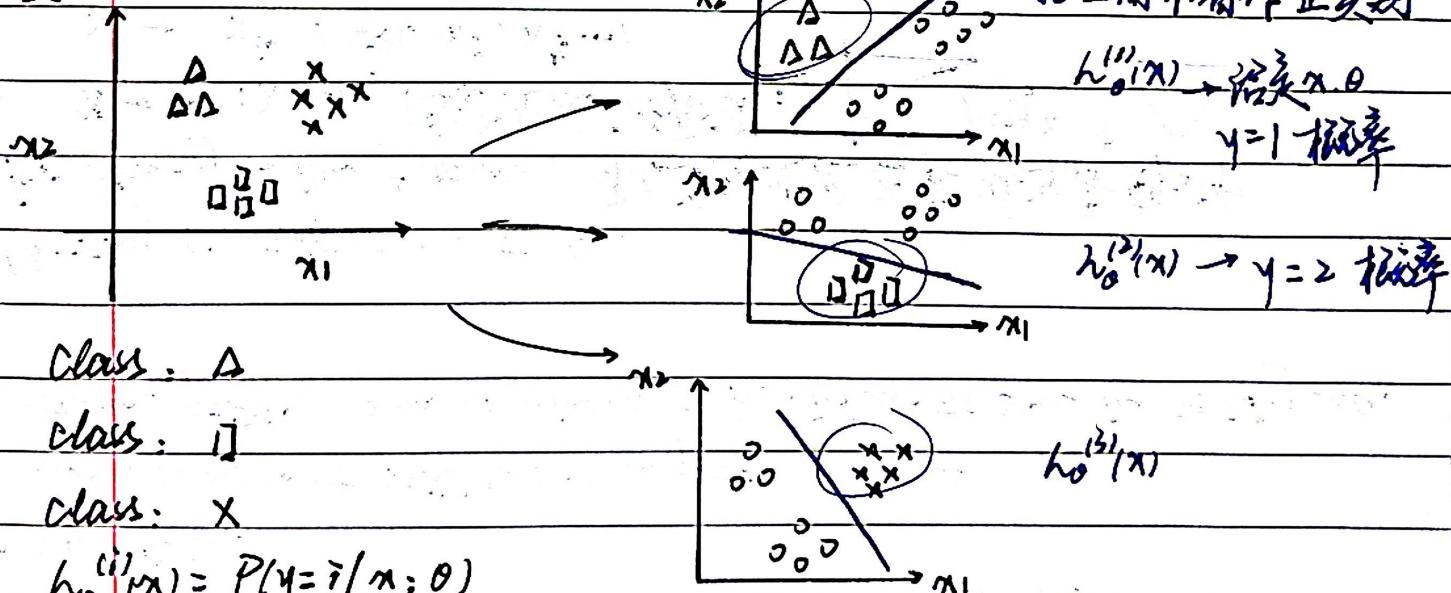
- More complex

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

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

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

原理



Class: A

Class: I

Class: X

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

总结:

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

In a new input x , to make a prediction, pick the class i that maximizes $\max h_0^{(i)}(x)$.

输出三个分类器

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** 具有高方差
- 假设函数几乎能拟合所有数据 → 函数太庞大，变量太多
无法约束 → 拟合训练集，导致系统泛化到新样本中 → **过拟合**
通常在变量太多时

• 解决方式

- 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：密布变量但无密布信息

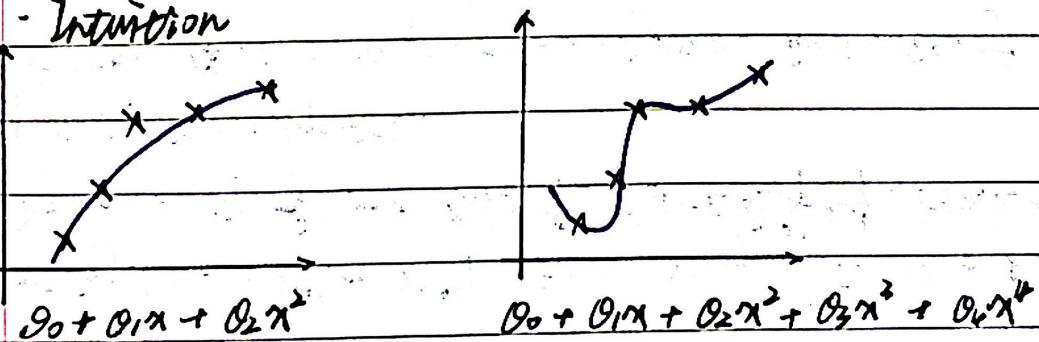
• Regularization 正则化

- 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.

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (\theta_0 x^{(i)} + \theta_1 x_1^{(i)} + \dots + \theta_n x_n^{(i)} - y^{(i)})^2 + 1000 \theta_3^2 + 1000 \theta_4^2$$

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

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

- 'Simpler' hypothesis \rightarrow 函数更平滑.
 - Less prone to overfitting \rightarrow 如果给 $\theta_0 \sim \theta_n$ 都加上惩罚项 \rightarrow 简化函数
- $J(\theta) = \frac{1}{2m} \sum_{i=1}^m (\theta_0 x^{(i)} + \theta_1 x_1^{(i)} + \dots + \theta_n x_n^{(i)} - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2$
- 因为不知道哪些系数，添加一个额外的正则化项以减小参数值

- 将 $J(\theta)$ 分为二项

- 第一项：为了更好地拟合训练集
- 第二项：保持参数尽可能小

三：正则化参数，控制以上两项间平衡关系

λ过大，会导致 $\theta_1 \sim \theta_n \rightarrow 0$ ，只剩 θ_0 (直线)

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

- Gradient descent

Repeat ?

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (\theta_0 x^{(i)} + \theta_1 x_1^{(i)} + \dots + \theta_n x_n^{(i)} - y^{(i)}) x_0^{(i)}$$

$$\theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (\theta_0 x^{(i)} + \theta_1 x_1^{(i)} + \dots + \theta_n x_n^{(i)} - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \quad j = 1, 2, 3, \dots, n$$

$$\rightarrow \theta_j := \theta_j (1 - \alpha \frac{\lambda}{m}) - \alpha \frac{1}{m} \sum_{i=1}^m (\theta_0 x^{(i)} + \theta_1 x_1^{(i)} + \dots + \theta_n x_n^{(i)} - y^{(i)}) x_j^{(i)}$$

由于 α 偏小 $\rightarrow \alpha \frac{\lambda}{m}$ 是一个很小的数

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

也就是说每次迭代时都将 θ_j 乘以一个比 1 小的数
然后进行和之前一样更新操作

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

- Normal equation

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

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

Regularized

5-4 Logistic regression logistic 回归的正则化

- Cost function

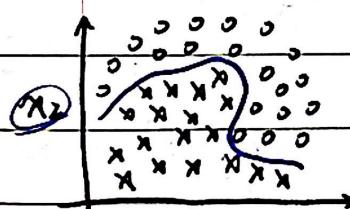
$$J(\theta) = - \left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h(\mathbf{x}^{(i)}) + (1-y^{(i)}) \log (1-h(\mathbf{x}^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

- 算法与 5-3 中 Gradient descent 一致. 但是用 $h(\mathbf{x})$ 表达式不同

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

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

⑨ - 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

b-3 模型表示 / Model representation

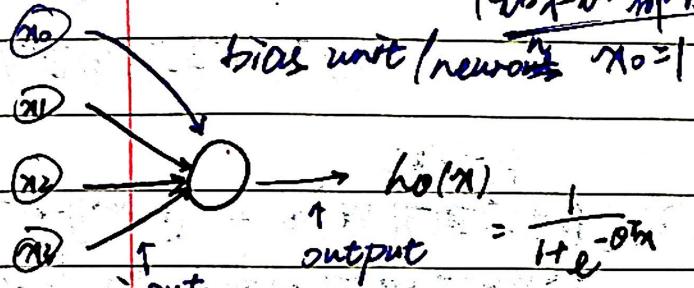
输出: 输出

树突: 输入通道
动作电位 \leftarrow 通道

• logistic unit

$1 + \text{exp}(-\theta x)$ 表示

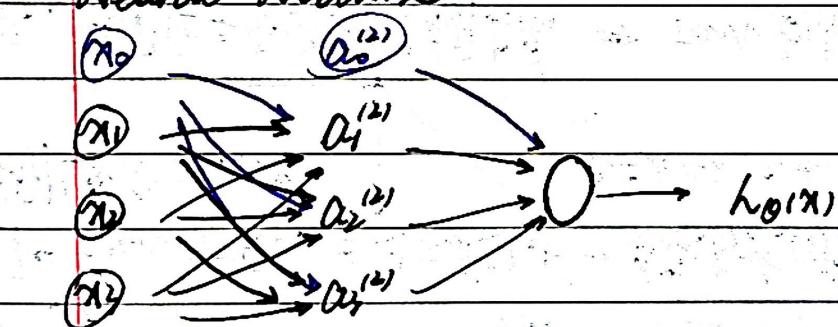
具体情况具体分析 双重 weights
或参数



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

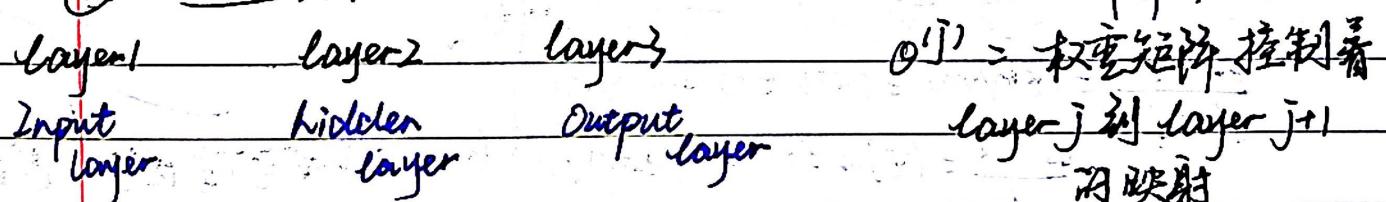
input units = 带有 sigmoid (logistic) 激活函数的人工神经元

• Neural Network



$a_i^{(j)}$ = activation of unit i in layer j

激活函数: 由一个具体神经元计算并输出数值



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

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

⋮

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

$\rightarrow s_j$ units in layer j

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

s_{j+1} units in layer $j+1$

6-4 模型展示 II

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

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

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

- 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

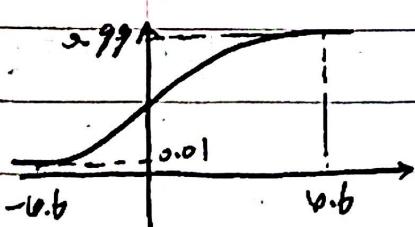
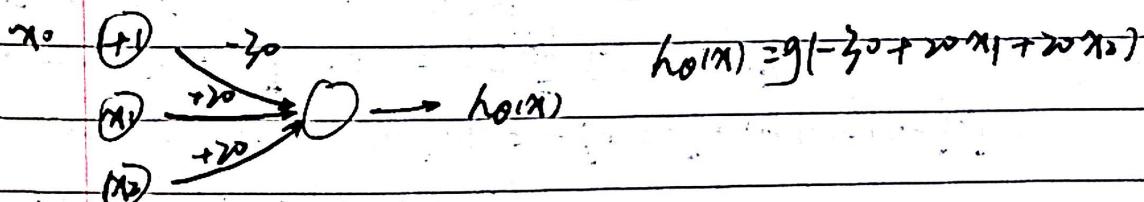
\Rightarrow 为 $\theta^{(1)}$ 选择不同参数，可以得到一个更好的假设函数
相较于直接使用原始特征 x_1, x_2, x_3 。

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

- 异或 XOR: 同为0(假), 异为1(真) 计算逻辑函数

- 同或 NOR: 同真异假

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



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

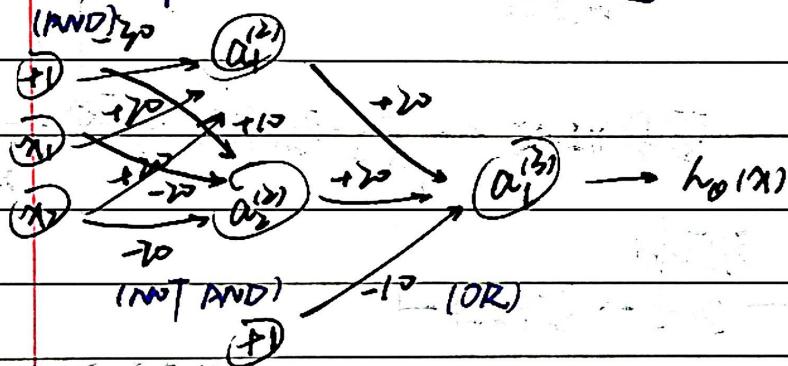
	x_1	x_2	$h_0(x)$
- OR	0	0	$g(-1) \approx 0$
	1	0	$g(+1) \approx 1$
	0	1	$g(1) \approx 1$
	1	1	$g(3) \approx 1$

	x_1	$h_0(x)$
- NOT	0	$g(1) \approx 1$
	1	$g(-1) \approx 0$

- $(\text{NOT } x_1) \text{ AND } (\text{NOT } x_2) \rightarrow g(1 - 2x_1 - 2x_2)$

6-6 例子和直觉理解 II

- 构造 x_1, x_2 XNOR x_2 (putting AND, OR, (NOT AND) together)



x_1	x_2	$a_1^{(2)}$	$a_2^{(2)}$	$a_1^{(3)}$	$a_2^{(3)}$	$h_0(x)$
0	0	0	1	1	0	1
0	1	0	0	0	0	0
1	0	0	0	0	0	0
1	1	1	0	0	1	1

6-7 多输出类

- Multiple output units : one-vs-all

$$[N \times N_o] \rightarrow h_0(x) \in \mathbb{R}^N$$

want $h_0(x) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$, $h_0(x) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$, $h_0(x) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$, etc.

when pedestrian when car when motorcycle

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

$y^{(i)}$ one of $\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$ $\Rightarrow h_0(x^{(i)}) \approx y^{(i)}$

pedestrian car motor truck

L7: Neural Networks: Learning

7-1 Cost function

- Neural Network (Classification)

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

- Binary classification

$y = 0 \text{ or } 1$, 1 output unit, $h_\theta(x) \text{ GR}$, $S_L = 1$ ($k=1$)

- Multi-class classification (K class)

$y \in R^k$ (多维向量), K output unit, $h_\theta(x) \text{ GR}^k$, $S_L = k$.

- 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:

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

$\frac{\lambda}{2m} \sum_{l=1}^L \sum_{i=1}^m \sum_{j=1}^{S_{l+1}} (\theta_{j+1}^{(l)})^2$

• 不对 bias, 即 $i=0$ 未加

• [前行-后行- \rightarrow 前行和后行所有单元连接]

• 第二-第 K 个输出单元

7-2 Backpropagation algorithm 反向传播算法

- Need to compute $\frac{\partial}{\partial \theta^{(l)}} J(\theta)$

$$\rightarrow \frac{\partial}{\partial \theta^{(l)}} J(\theta)$$

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

$$a^{(1)} = x, z^{(1)} = \theta^{(1)} \# a^{(1)}$$

$$a^{(2)} = g(z^{(2)}) \quad (\text{add } a^{(1)}), \quad z^{(2)} = \theta^{(2)} a^{(1)}$$

$$a^{(3)} = g(z^{(3)}) \quad (\text{add } a^{(2)}), \quad z^{(3)} = \theta^{(3)} a^{(2)}$$

$$a^{(4)} = h_\theta(x) = g(z^{(4)})$$

引入“误差”的概念，然后通过“误差”来得训练度

- Backpropagation

Intuition: $s_j^{(l)} = \text{"error" of node } j \text{ in layer } l$

for each output unit (layer $L=4$)

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

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

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

no $s^{(1)}$

激活函数偏导数

→ 第一层 (layer 1) 不存在误差

$$\rightarrow \text{可简单证明 } \frac{\partial J(\theta)}{\partial \theta_{ij}^{(l)}} \Rightarrow a_j^{(l)} s_i^{(l+1)}$$

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

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

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

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

$$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} s_i^{(l+1)}$$

误差表示

$$\Delta^{(l)} := \Delta^{(l)} + s^{(l+1)} (a^{(l)})^T$$

$$\begin{aligned} \text{计算梯度} \quad D_{ij}^{(l)} &:= \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \theta_{ij}^{(l)} \quad \text{if } j \neq 0 \\ \text{反向传播} \quad D_{ij}^{(l)} &:= \frac{1}{m} \Delta_{ij}^{(l)} \quad \text{if } j = 0 \end{aligned} \rightarrow \frac{\partial J(\theta)}{\partial \theta_{ij}^{(l)}} = D_{ij}^{(l)}$$

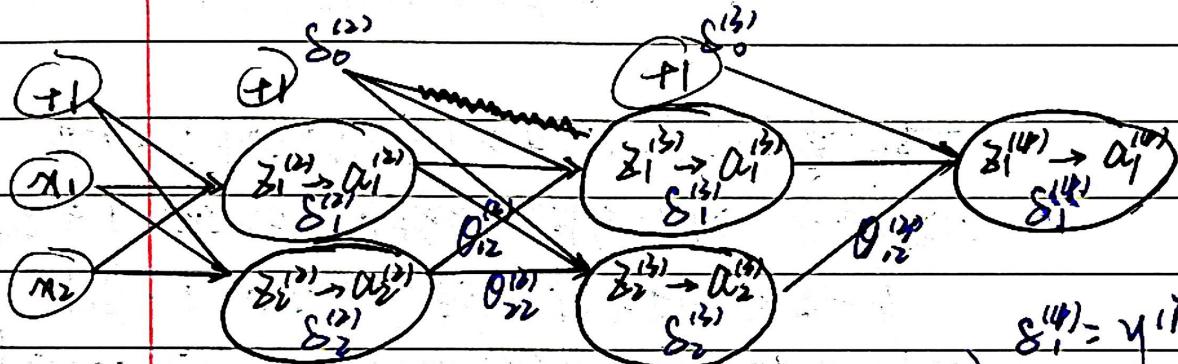
然后可以使用梯度下降或其变体。

7-3 Backpropagation intuition 反向传播直觉

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

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

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



$$\delta_1^{(4)} = y^{(i)} - \alpha_1^{(4)}$$

$$\delta_2^{(4)} \Rightarrow \theta_{12}^{(3)} \delta_1^{(4)}$$

$$\delta_2^{(4)} \Rightarrow \theta_{22}^{(3)} \delta_1^{(4)} + \theta_{23}^{(3)} \delta_2^{(4)}$$

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

7-4 使用注意：展开参数 Implementation note: Unrolling parameters

- Advanced optimisation 中输入的 theta / initial theta 为 n/m 维向量 gradient 也为 n/m 的 vector.

- Neural Network ($L=4$) \rightarrow 展开成什么?

参数矩阵 $\alpha_m, \theta^{(2)}, \theta^{(3)}$ - matrices (Theta1, Theta2, Theta3)

(返回值) 梯度矩阵 $D^{(1)}, D^{(2)}, D^{(3)}$ - matrices (D1, D2, D3)

如何展开为向量?

(unroll)

7-5 Gradient checking 梯度检测 → 避免算法产生 bug

- 原理 (双侧差分)

$$\frac{\partial}{\partial \theta} J(\theta) \approx \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon} \quad \epsilon \text{尽量小 e.g. } \epsilon = 10^{-4}$$

上述 θ 为参数

- Parameter vector θ , $\theta \in \mathbb{R}^n$ b.g. θ is "unrolled" version of $\theta = [\theta_1, \dots, \theta_n]$

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

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

并在训练中利用 for 循环得到上述所有参数 ~~依次~~ gradApprox 然后与反向传播中得到的梯度进行比较

一旦确定 gradApprox 和 Drec 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
- "Drec" 高效

7-6 Random initialization 随机初始化 T 可考 Deep learning 21 Note 1.1 节

- 如何对 θ 设初始值?

假设 zero 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
仅限于单层
学习

- Random initialization : Symmetry breaking 对称性破坏
Initialize each $\theta_{ij}^{(l)}$ 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 $-E$ 和 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）

7-7 Putting it together

• Training a neural network

- I. - pick a network architecture (connectivity pattern between neurons)
 - No. of input units: Dimension of features $x^{(i)}$
 - No. of output unit: Number of classes \rightarrow like $y = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ or $\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$ or $\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$
 - Reasonable default:
 - 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. \geq No. of input units

II. 步骤

注:

1. randomly initialize weights (已 ~~上~~ Summarize)
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 (见 7-5)
6. Using gradient descent or advanced optimization method, with backpropagation to try to minimize $J(\theta)$ as a function of parameters θ .