

第一周

监督学习：回归，分类

无监督：聚类

单变量线性回归：

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

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

$$\text{Hypothesis: } h(x) = \theta_0 + \theta_1 x$$

Parameters: θ_0, θ_1

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

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

梯度下降：

batch gradient descent:

repeat until convergence

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

}

注：需同时更新 θ_0, θ_1

多变量线性回归：

$$\text{Hypothesis: } h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$\text{or } h(x) = \theta^T X \quad \text{其中 } x_0 \equiv 1$$

$$\text{Loss: } J(\theta_0, \theta_1, \theta_2, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2$$

1



Parameters: $\theta_0, \theta_1, \dots, \theta_n$

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

Goal: minimize $J(\theta_0, \theta_1, \dots, \theta_n)$
 $\theta_0, \theta_1, \dots, \theta_n$

梯度下降法 tips:

① Feature scaling: $x_{in} = \frac{x_n - \mu_n}{s_n}$ 其中 μ_n 为平均值
 s_n 为标准差.

② Learning rate: 如何调整

特征与多项式回归:

可以添加多项式项整合特征联系

注: 采用多项式回归, 特征缩放非常有必要.

Normal Equation (正规方程):

用于求解最优解: $\theta = (X^T X)^{-1} X^T y$

注: 仅适用于可逆矩阵

梯度下降

VS

正规方程

需要选择 α

不需要

多次迭代

一次运算

适于 n 大

需计算 $(X^T X)^{-1}$, $O(n^3)$, $n < 10000$ 可接受

各种模型

只适用于线性模型.



逻辑回归 (Logistic Regression)

① 分类问题

△ 为什么是线性分类器?

① $\hat{y} = \hat{\theta} \cdot x$ 为 x 的线性组合

② 决策边界是线性的.

$$\text{Hypothesis: } h_{\theta}(x) = g(\theta^T x)$$

$$\hat{\theta} \cdot x = 0$$

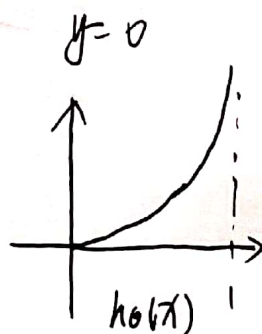
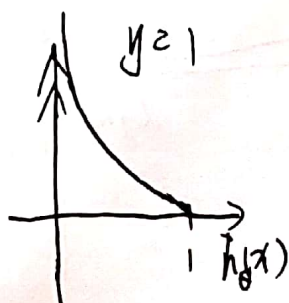
$$g(z) = \frac{1}{1 + e^{-z}} \quad (\text{Sigmoid Function})$$

决策边界 (decision boundary)

Loss: 如果用平方误差将 $g(z)$ 代入, 得到一个非凸函数

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

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



$$\text{缩写为: } \text{Cost}(h_{\theta}(x), y) = -y \log(h_{\theta}(x)) - (1-y) \log(1-h_{\theta}(x))$$

注: Feature scaling is important.

其它求解法: 共轭梯度, 局部优化法, 有限内存局部优化法



多类别分类: 训练多个分类器.

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

神经网络:

Loss Function:

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

反向传播算法:

梯度检验:

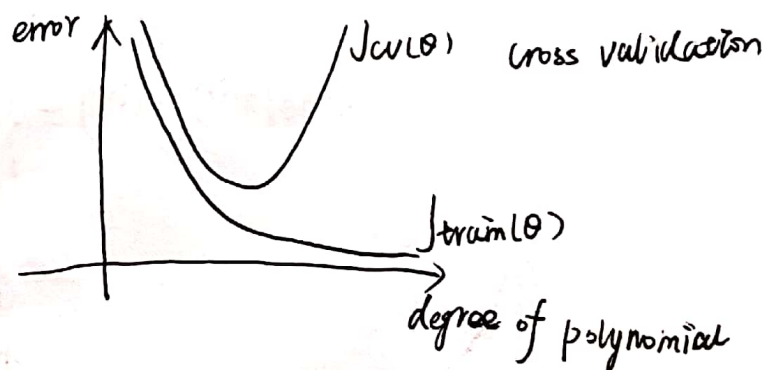
检验 θ 的斜率与反向传播计算差距.

注意随机初始化.

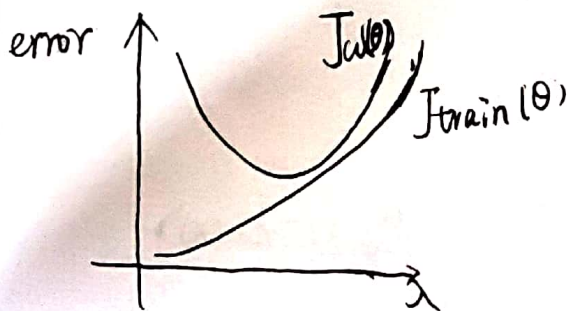
偏差与方差的诊断:

高偏差: 欠拟合

高方差: 过拟合.



正则参数的选择:



类偏斜:

实际 \ 预测	P	N
	TP	FN
实际 \ 预测	FP	TN
	FN	TP

准确率: $P = \frac{TP}{TP + FP}$

召回率: $R = \frac{TP}{TP + FN}$

F_1 score: $\frac{2PR}{R+P}$



(监督学习)

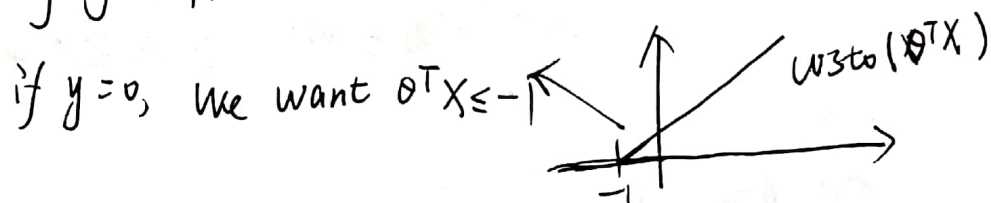
Loss Function: $CA + B$

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

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

(大间隔分类器)

if $y=1$, we want $\theta^T x \geq 1$



核函数: 与标记点

线性核、高斯核、多项式核、字符串核、卡方核。

Logestic regression vs ~~SVM~~ SVM

n : number of feature m : number of examples

- ① n large (relative to m) Logistic or SVM without kernel.
- ② n is small, m intermediate. SVM with Gaussian kernel
- ③ n is small, m is large. Create more features, then use Logistic or SVM without a kernel.



扫描全能王 创建

K-均值: ① 随机选取K个聚类中心.

② 代分类别

③ 更新聚类中心

④ 直到中心点不变。

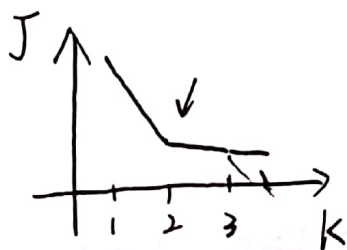
Loss Function: $J(c^{(1)}, c^{(2)}, \dots, c^{(m)}, u_1, \dots, u_K)$

$$= \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - u_{c(i)}\|^2$$

$u_{c(i)}$ 代表与 $x^{(i)}$ 最近的聚类中心。

选择K值:

使用“肘部法则”



距离/相似度度量:

① 闵可夫斯基距离 $\text{dist}(x, y) = \left(\sum_{i=1}^n |x_i - y_i|^p \right)^{\frac{1}{p}}$

② 杰卡德相似系数:

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

③ 余弦相似度:

$$\cos(\theta) = \frac{x^T y}{\|x\| \|y\|} = \frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum_{i=1}^n x_i^2} \sqrt{\sum_{i=1}^n y_i^2}}$$

④ Pearson 相似系数:

$$\rho_{xy} = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{E[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y} = \frac{\sum_{i=1}^n (x_i - \mu_X)(y_i - \mu_Y)}{\sqrt{\sum_{i=1}^n (x_i - \mu_X)^2} \sqrt{\sum_{i=1}^n (y_i - \mu_Y)^2}}$$



均一性 p , 完整性 r , V-measure: $V = \frac{(1+\beta) \times p \times r}{\beta^2 p + r}$

轮廓系数:

降维:

主成分分析 (PCA):

找一方向向量, 希望投影平面的方差尽可能小。

① 均值归一化 $x_j = \frac{x_j - \mu_j}{\sigma_j^2}$

② 协方差矩阵 $\Sigma = \frac{1}{m} \sum_{i=1}^n (x^{(i)})(x^{(i)})^T$

③ 计算 Σ 的特征向量。

$$[u, s, v] = \text{svd}(\Sigma)$$

选取 u 中前 k 个向量, 计算新特征向量 $z^{(i)} = U_{\text{reduce}}^T * x^{(i)}$

主成分分析 主成分数量选择:

$$\rho = \frac{\text{投影的平方误差}}{\text{训练集方差}} = \frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 1\%$$

则保留 99% 的偏差。

$$x_{\text{approx}}^{(i)} = U_{\text{reduce}} * z^{(i)}$$

异常检测:

通过高斯分布来解决异常。

高斯:

构建特征的协方差矩阵, 用所有特征一起来计算 $\rho(x)$ 。



原高斯

不能捕捉特征之间相关性,
可以通过组合方法解决
计算成本低.

多元高斯

自动捕捉特征之间相关性.

计算代价大.

必须有 $m > n$, 保证协方差矩阵不可逆.

推荐系统:

① 基于内容的推荐: 已知物品特征.

② 协同过滤: 拥有用户参数. 亦可同时学习物品和用户参数.

均值归一化解决新用户; 预测时: $(\theta^{(i)})^T x^{(i)} + \mu_i$, 但困难于新用户
评分为均值 (大众口味)

随机梯度下降: 小批量梯度下降

上限分析:

