

Coursera – Machine Learning (taught Andrew Ng) – Notes by Xiaofan Wang

Name: Xiaofan Wang
Chapter 1.

- What's Machine Learning.

1. Two definitions

Old: The field of study that gives computers the ability to learn without being explicitly programmed.

"计算机无须进行明确编程即可学习" 的领域

Modern: A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P , if it improves with experience E .

即某计算机程序在 T 任务中的性能 (由 P 测量)

随着经验 E 的提高而提高，则该它是从经验 E 中学习有关某类任务 T 和性能度量 P 的。

Example: Playing Checkers. 下国际象棋

E = experience of playing many games

P = probability that program will win next game.

T = task of playing checkers.

2. ML Algorithm

Supervised learning & Unsupervised Learning.

二. Supervised Learning.

已知: Dataset & correct Output

探求: Relationship between "input" and "output"

Regression: predict continuous "output" (values)

Classification: predict discrete "output" (categories)

三. Unsupervised learning.

Derive structure from data, don't necessarily know the effect of variables. No feedback on prediction results.

Derive structure by clustering data based on relationships among variables in the data.

基于变量之间的关系对数据进行聚类来获得此结构
不必知道变量的影响。没有基于预测结果的反馈。

Example:

Clustering: Take a collection of 1000+ different genes, and find a way to automatically group these genes into groups that are somehow similar or related to by different variables (e.g., lifespan, location, roles...)

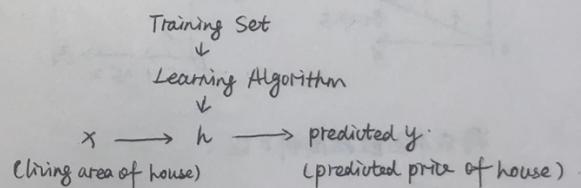
收集许多不同基因，并找到一种方法将基因组自动分组 (相似，相关)

Non-clustering: find structure in a chaotic environment.

(e.g. identify individual voices and music from a mesh of sounds at cocktail party).

四. Model and Cost Function.

1. Model Representation.

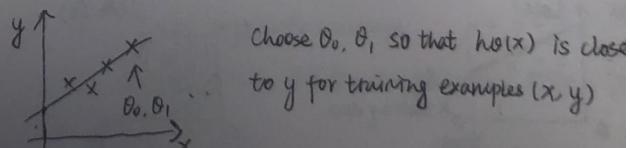


2. Cost Function

Measure the accuracy of our hypothesis function.

This takes an average difference of all the results of the hypothesis with inputs from x 's and the actual output y :

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}_i - y_i)^2 = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x_i) - y_i)^2$$



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

Parameters: θ_0, θ_1 / θ_1

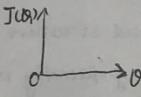
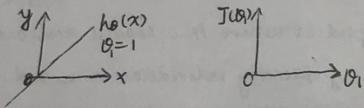
Cost Function: $J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

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

1. 训练数据集散落在 xy 平面上，尝试画条直线 ($h_{\theta}(x)$) 穿过这些分散的数据点。

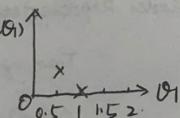
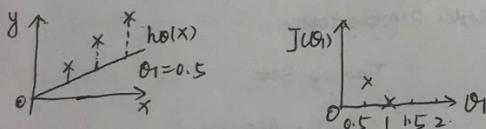
2. Best possible: Average squared vertical distances of the scatter points from the line will be the least.

最佳：散射点与该线的平均垂直距离最小。 $J(\theta_0, \theta_1) = 0$

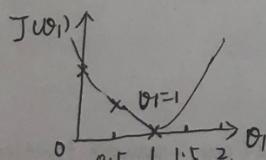


3. $\theta_1 = 1$ 时，每个数据点的斜率均为 1。

$\theta_1 = 0$ 时，从拟合到数据点的垂直距离增加了。



将成本函数提高到 0.58：



应尽量减小成本函数。

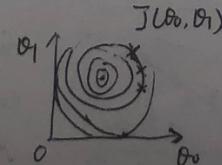
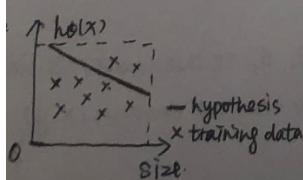
$\theta_1 = 1$ 是 global minimum.

4. Contour plot.

等高线：

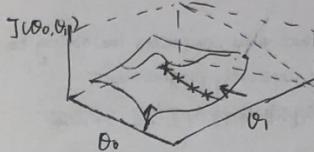
同一条线上的所有点具有恒定值！

A contour line of a two variable function has a constant value at all points of the same line.



五. Parameter Learning

1. 测量假设函数与数据的拟合程度：估计假设函数中的参数（梯度下降）



成本函数位于凹坑最底
部时，其值最小。

采用成本函数的导数（函数的切线），切线斜率是该点的导数，它将为我们提供一个方向，沿下降最快的方向逐步降低成本函数。（由 para or 学习率决定）

para: 步长。

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

$$\text{temp}_0 = \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\text{temp}_1 = \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 = \text{temp}_0$$

$$\theta_1 = \text{temp}_1$$

① If α is too small, gradient descent can be slow.

② If α is too large, gradient descent can overshoot the minimum. It may fail to converge. or even diverge. 收敛，发散。

2. 线性回归的梯度下降

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

$$= \frac{1}{2} \cdot 2 (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_j} (h_{\theta}(x) - y)$$

$$= (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_j} \left(\sum_{i=1}^m \theta_j x_j^{(i)} - y \right)$$

$$= (h_{\theta}(x) - y) x_j$$



椭圆是二次函数的轮廓。

图中的 x 标记了梯度下降

收敛到最小值时经历的
的连读值。

Name: Xiaofan Wang

Chapter 2. Linear Regression with Multiple Variables

1. Multiple Features

Hypothesis Function: $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$

$$h_{\theta}(x) = [\theta_0, \theta_1, \dots, \theta_n] \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} = \theta^T x$$

2. Gradient Descent for Multiple Variables

1. The gradient descent equation itself is generally the same form.

repeat it for "n" features.

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

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

$$\theta_2 = \theta_2 - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_2^{(i)}$$

:

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

2. hypothesis: $h_{\theta}(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n$

parameters: $\theta_0, \theta_1, \dots, \theta_n$ ($\vec{\theta}$)

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

3. Gradient Descent (Feature Scaling) I.

1. Feature Scale: Make sure features are on a similar scale.

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

2. Mean Normalization

Replace x_i with $\frac{x_i - \mu_i}{s_i}$ to make features have approximately zero mean.

$$x_i = \frac{x_i - \mu_i}{s_i}$$

4. Gradient Descent II (Learning Rate)

1. Debugging gradient descent

Make a plot with number of iterations on the x-axis.

Now plot the cost function, $J(\theta)$ over the number of iterations of gradient descent. If $J(\theta)$ ever increases, then you probably need to decrease α .

调试梯度下降：绘制一个在x轴上有迭代次数的图。

现在在梯度下降的迭代次数上绘制成本函数 $J(\theta)$.

如果 $J(\theta) \uparrow$. 则 $\downarrow \alpha$.

2. Automatic convergence test

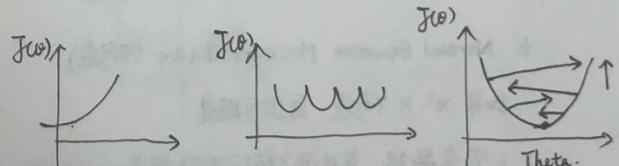
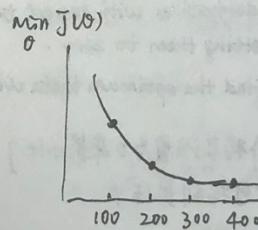
Declare convergence if $J(\theta)$ decreases by less than E in one iteration, where E is small value such as 10^{-3} .

However, in practice it's difficult to choose this threshold value.

自动收敛测试：如果在一次迭代中 $J(\theta) \downarrow < E$.

则声明收敛. 其中 E 是一些小值. 例如 10^{-3} .

但实际上，很难选择此阈值.



- For sufficiently small α , $J(\theta)$ should decrease on every iteration. 如果 α 太小，收敛缓慢.

- But if α is too small, gradient descent can be slow to converge.

如果 α 太大，不会每次迭代都减小. 可能不会收敛

2. Features & Polynomial Regression

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

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

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

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

3. Normal Equation

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

Gradient Descent. VS. Normal Equation.

Need to choose alpha. No Need to choose alpha.

Need many iterations. No need to iterate.

$O(kn^2)$. $O(n^3)$, calculate $X^T X$

Work well when n is large. Slow if n is very large.

1. Gradient Descent : Give one way of minimizing.

2. Normal Equation : minimize J by explicitly taking its
~~the~~ derivatives with respect to θ_j 's
and setting them to zero.
⇒ Find the optimum theta without iteration.

针对 θ_j 取导数并将其设置为 0 来最小化 J.

使得无梯迭代即可找到最佳 θ .

4. Normal Equation Noninvertibility (不可逆)

如果 $X^T X$ 不可逆. 原因可能是:

1. 冗余特征: 其中两个特征线性相关.

2. 特征过多 (删除某些功能用正则化)

Name: Xiaofan Wang

Chapter 3. Classification.

- Classification and Representation.

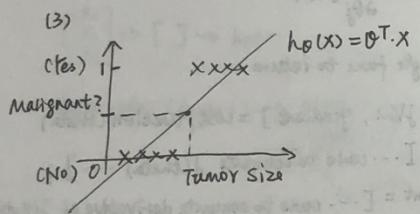
1. Classification.

- (1) $y \in \{0, 1\}$. 0: "Negative Class" (e.g. benign tumor)
- 1: "Positive Class" (e.g. malignant tumor)

(2) Example: Email (Spam / Not Spam)

Online Transactions: Fraudulent (Yes / No)

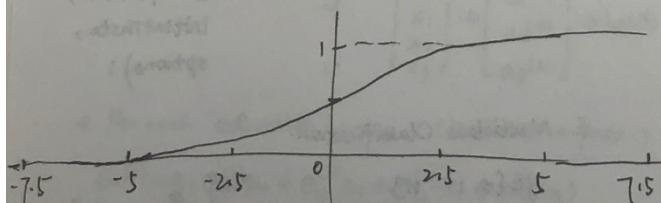
Tumor: Malignant / Benign.



2. Hypothesis Representation. 假设表示.

Logistic Function (Sigmoid Function)

$$\begin{cases} h_\theta(x) = g(\theta^T x) & 0 \leq h_\theta(x) \leq 1 \\ z = \theta^T x \\ g(z) = \frac{1}{1+e^{-z}} \end{cases}$$



$$h_\theta(x) = P(y=1|x; \theta) = 1 - P(y=0|x; \theta)$$

$$P(y=1|x; \theta) + P(y=0|x; \theta) = 1$$

3. Decision Boundary. 决策边界.

(1) Logistic Regression: $h_\theta(x) = g(\theta^T x)$

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

↓

$$h_\theta(x) \geq 0.5 \rightarrow y=1$$

$$h_\theta(x) < 0.5 \rightarrow y=0$$

$$\begin{cases} z=0, g(z) = \frac{1}{1+e^{-z}} = \frac{1}{1+e^0} = \frac{1}{2} = 0.5 \\ z=\infty, g(z) = \frac{1}{1+\frac{1}{e^\infty}} = 1 \\ z=-\infty, g(z) = \frac{1}{1+e^{-\infty}} = 0 \end{cases}$$

$g(z) \geq 0.5$. When $z \geq 0$

$$h_\theta(x) = g(\theta^T x) \geq 0.5 \text{ when } \theta^T x \geq 0$$

↓

$$\theta^T x \geq 0 \rightarrow y=1$$

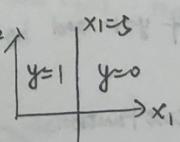
$$\theta^T x < 0 \rightarrow y=0$$

(2) Decision Boundary.

$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

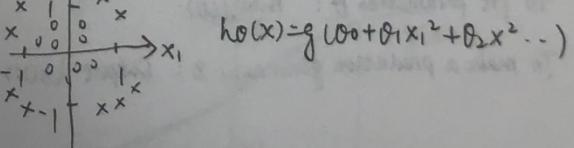
$$\text{By: } \theta = \begin{bmatrix} 5 \\ -1 \\ 0 \end{bmatrix}, h_\theta(x) = g(5 - x_1)$$

$$5 - x_1 \geq 0, x_1 \leq 5$$



(3) Non-linear Decision Boundary.

$$z = \theta_0 + \theta_1 x_1^2 + \theta_2 x_2^2$$



| |
|--------|
| $0: +$ |
| $x: -$ |

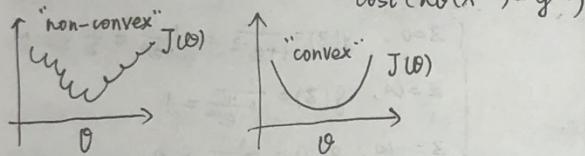
2. Logistic Regression Model.

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

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

2. Cost Function

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

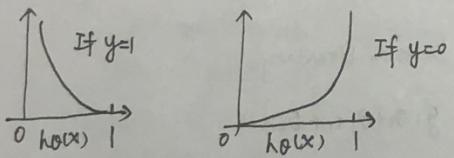


(2) cost func for LR: ($y=0/1$ always)

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

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

$$\text{cost}(h_{\theta}(x), y) = -\log(1-h_{\theta}(x)) \quad \text{if } y=0$$



$$\text{cost}(h_{\theta}(x), y) = 0. \quad \text{if } h_{\theta}(x) = y.$$

$$\text{cost}(h_{\theta}(x), y) \rightarrow \infty \quad \text{if } y=0 \text{ and } h_{\theta}(x) \rightarrow 1$$

$$\text{cost}(h_{\theta}(x), y) \rightarrow \infty \quad \text{if } y=1 \text{ and } h_{\theta}(x) \rightarrow 0.$$

3. Logistic Regression Cost Function.

$$\begin{aligned} J(\theta) &= \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)}) \\ &= \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)})) \right]. \end{aligned}$$

To fit parameters θ : $\min_{\theta} J(\theta)$.

To make a prediction given new x : Output $h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}}$

4. Gradient Descent.

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

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

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

5. Advanced Algorithm 高级算法

"共轭梯度" 是更复杂、更快的优化目的方法。
"BFGS" \Rightarrow 可用于梯度下降
"L-BFGS" (不需自己满维矩阵算法，使用一些库)

(1) Func: $J(\theta) \cdot \frac{\partial}{\partial \theta_j} J(\theta)$

write a single func to return:

```
>> function [JVal, gradient] = costFunction(theta)
>> JVal = [... code to compute J(theta) ...];
>> gradient = [... code to compute derivative of J(theta)...];
>> end.
```

(2) optimization Algorithm: fminunc(), optimset()

```
options = optimset('GradObj', 'on', 'MaxIter', 100);
initialTheta = zeros(2, 1);
```

```
[optTheta, functionVal, exitFlag] = fminunc
```

```
(@costFunction,
initialTheta,
options);
```

6. Multiclass Classification.

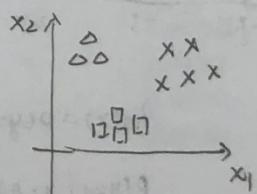
$$y \in \{0, 1, \dots, n\}$$

$$h_{\theta}^{(0)}(x) = P(y=0 | x; \theta)$$

$$h_{\theta}^{(1)}(x) = P(y=1 | x; \theta)$$

$$h_{\theta}^{(i)}(x) = P(y=i | x; \theta)$$

$$\text{prediction} = \max_i h_{\theta}^{(i)}(x)$$



Name: Xiaofan Wang

Chapter 4. Neural Networks

- Model Representation. (I)



1. Neurons are basically computational units that take inputs as electrical inputs that are channeled to outputs.
(dendrites ~~神经元~~ 树突) (axons ~~轴突~~)

神经元是计算单元，将输入(树突)作电输入。
并传输到输出(轴突)的计算单元。

2. In the model, Input features: x_1, x_2, \dots, x_n
output is the result of hypothesis function: $\frac{1}{1+e^{-\theta^T x}}$
(activation function) (θ : weights)

$$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} \rightarrow [] \rightarrow h_{\theta}(x)$$

Input nodes (layer 1) \rightarrow another node \rightarrow hypothesis function
(layer 2) (output layer)

Intermediate layers : hidden layers. ($a_0^2 \dots a_n^2$ activation units)

(Between input & output layer)

$\{a_i^{(j)}\}$ = "activation" of unit i in layer j
 $\{\theta^{(j)}\}$ = matrix of weights controlling function
mapping from layer j to layer $j+1$

3. If having one hidden layer:

$$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \rightarrow \begin{bmatrix} a_1^{(2)} \\ a_2^{(2)} \\ a_3^{(2)} \end{bmatrix} \rightarrow h_{\theta}(x)$$

4. For each "activation" nodes: (3x4 matrix para-)

$$a_1^{(2)} = g(\theta_{01}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{21}^{(1)} x_2 + \theta_{31}^{(1)} x_3)$$

$$a_2^{(2)} = g(\theta_{02}^{(1)} x_0 + \theta_{12}^{(1)} x_1 + \theta_{22}^{(1)} x_2 + \theta_{32}^{(1)} x_3)$$

$$a_3^{(2)} = g(\theta_{03}^{(1)} x_0 + \theta_{13}^{(1)} x_1 + \theta_{23}^{(1)} x_2 + \theta_{33}^{(1)} x_3)$$

$$h_{\theta}(x) = a_1^{(3)} = g(\theta_{00}^{(2)} a_0^{(2)} + \theta_{10}^{(2)} a_1^{(2)} + \theta_{20}^{(2)} a_2^{(2)} + \theta_{30}^{(2)} a_3^{(2)})$$

$$\theta^{(j)} \text{ dimension: } S_{j+1} \times (S_j + 1)$$

= Model Representation (II). Replace parameter with "z")

$$a_1^{(2)} = g(z_1^{(2)})$$

$$a_2^{(2)} = g(z_2^{(2)})$$

$$a_3^{(2)} = g(z_3^{(2)})$$

1. For layer $j=2$. Node k :

$$z_k^{(2)} = \theta_{k,0}^{(1)} x_0 + \theta_{k,1}^{(1)} x_1 + \dots + \theta_{k,n}^{(1)} x_n$$

2. The vector representation of x , $z^{(j)}$:

$$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \quad z^{(j)} = \begin{bmatrix} z_1^{(j)} \\ z_2^{(j)} \\ \vdots \\ z_n^{(j)} \end{bmatrix} \Rightarrow z^{(j)} = \theta^{(j-1)} a^{(j-1)}$$
$$z^{(j+1)} = \theta^{(j)} a^{(j)}$$

3. Final result: $h_{\theta}(x) = a^{(j+1)} = g(z^{(j+1)})$

III. Examples and Intuitions (I)

1. AND.

$$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} \rightarrow [g(z^{(2)})] \rightarrow h_{\theta}(x)$$

$$\theta^{(1)} = [-30 \ 20 \ 20]$$

$$h_{\theta}(x) = g(-30 + 20x_1 + 20x_2)$$

$$x_1=0 \text{ and } x_2=0 \cdot g(-30) \approx 0$$

$$x_1=0 \text{ and } x_2=1 \cdot \text{then } g(-10) \approx 0$$

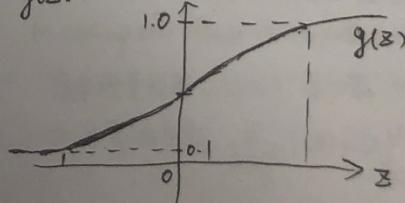
$$x_1=1 \text{ and } x_2=0 \cdot \text{then } g(10) \approx 0$$

$$x_1=1 \text{ and } x_2=1 \cdot \text{then } g(30) \approx 1$$

2. OR.

| | | | x_1 | x_2 | $h_{\theta}(x)$ |
|------|-----|---|-------|-------|--------------------|
| (+) | -10 | | 0 | 0 | $g(-10) \approx 0$ |
| (x1) | 20 | → | 0 | 1 | $g(10) \approx 1$ |
| (x2) | 20 | → | 1 | 0 | $g(10) \approx 1$ |
| | | | 1 | 1 | $g(30) \approx 1$ |

$$g(z) :$$



III. Examples and Intuitions (II)

$$1. \begin{cases} \text{AND: } \theta^{(1)} = [-30, 20, 20] \\ \text{NOR: } \theta^{(1)} = [10, -20, -20] \\ \text{OR: } \theta^{(1)} = [-10, 20, 20] \end{cases}$$

2. combine these to "XNOR":

$$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} \rightarrow \begin{bmatrix} \alpha_0^{(2)} \\ \alpha_1^{(2)} \\ \alpha_2^{(2)} \end{bmatrix} \rightarrow \begin{bmatrix} \alpha^{(3)} \end{bmatrix} \rightarrow h_{\theta}(x)$$

$$3. \theta^{(1)} = \begin{bmatrix} -30 & 20 & 20 \\ 10 & -20 & -20 \end{bmatrix} \quad \theta^{(2)} = \begin{bmatrix} -10 & 20 & 20 \end{bmatrix}$$

$$4. \alpha^{(2)} = g(\theta^{(1)} x)$$

$$\alpha^{(3)} = g(\theta^{(2)} \alpha^{(2)})$$

$$h_{\theta}(x) = \alpha^{(3)}$$

5. put it together: $x_1 \text{ XNOR } x_2$

$$\begin{array}{c} \oplus \\ \otimes \\ \otimes \end{array} \begin{array}{c} -30 \\ 20 \\ 20 \end{array} \rightarrow h_{\theta}(x)$$

$$\begin{array}{c} \oplus \\ \otimes \\ \otimes \end{array} \begin{array}{c} 10 \\ -20 \\ -20 \end{array} \rightarrow h_{\theta}(x)$$

$$\begin{array}{c} \oplus \\ \otimes \\ \otimes \end{array} \begin{array}{c} -10 \\ 20 \\ 20 \end{array} \rightarrow h_{\theta}(x)$$

$x_1 \text{ AND } x_2$

$(\text{NOT } x_1) \text{ AND } (\text{NOT } x_2)$

$x_1 \text{ OR } x_2$

IV. Multiclass Classification

$$y^{(i)} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

Final hypothesis function:

$$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \rightarrow \begin{bmatrix} \alpha_0^{(2)} \\ \alpha_1^{(2)} \\ \alpha_2^{(2)} \\ \vdots \end{bmatrix} \rightarrow \begin{bmatrix} \alpha_0^{(3)} \\ \alpha_1^{(3)} \\ \alpha_2^{(3)} \\ \vdots \end{bmatrix} \rightarrow \dots \rightarrow \begin{bmatrix} h_{\theta}(x)_1 \\ h_{\theta}(x)_2 \\ h_{\theta}(x)_3 \\ h_{\theta}(x)_4 \end{bmatrix}$$

Name: Xiaofan Wang

Chapter 5. Neural Networks : Learning.

1. Cost Function. (generalization of logistic regression)

1. define variables :

L : total number of layers in the network 层数

S_l : number of units in layer l .

k : number of output units

2. For regularized logistic regression:

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

For Neural Networks, it's more complicated :

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

① 添加了一些嵌套求和，该求和遍历输出节点的数量。

② 考虑多个 θ matrix. { 列数 = 当前层的节点数 . }

{ 行数 = 下一层的节点数 . }

3. Notes.

① Double sum adds up the logistic regression costs calculated for each cell in the output layer.

② Triple sum adds up the squares of all the individual θ s in the entire network.

③ the i in the triple sum does not refer to training example.

2. Backpropagation Algorithm. (minimize cost function)

$$\text{min}_{\theta} J(\theta) \rightarrow \frac{\partial}{\partial \theta_{i,j}} J(\theta)$$

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

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

For training example $t=1$ to m : Set $a^{(1)} = x^{(t)}$

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

Gradient computation :

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

$$z^{(2)} = \theta^{(1)} a^{(1)}$$

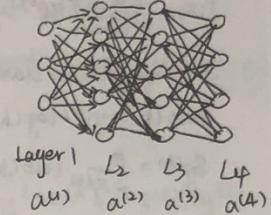
$$a^{(2)} = g(z^{(2)}) \text{ (add } a_0^{(2)})$$

$$z^{(3)} = \theta^{(2)} a^{(2)}$$

$$a^{(3)} = g(z^{(3)}) \text{ (add } a_0^{(3)})$$

$$z^{(4)} = \theta^{(3)} a^{(3)}$$

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



3. Using $y^{(t)}$, compute $\delta^{(L)} = a^{(L)} - y^{(t)}$

$a^{(l)}$: vector output of activation units for last year.

To get the $\delta^{(L)}$ before last layer, back from right to left

4. Compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$

$$\text{using } \delta^{(l)} = ((\theta^{(l)})^T \delta^{(l+1)}) * a^{(l)} * (1-a^{(l)})$$

下一层增量 * θ Matrix

第一层增量值

$$g'(z^{(l)}) = a^{(l)} * (1-a^{(l)})$$

activation function

$$\begin{cases} \Delta_{i,j}^{(l)} = \Delta_{i,j}^{(l)} + a_j^{(l)} s_i^{(l+1)} \\ \Delta^{(l)} = \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T \end{cases}$$

Update New Matrix :

$$\begin{cases} D_{i,j}^{(l)} = \frac{1}{m} (\Delta_{i,j}^{(l)} + \lambda \theta_{i,j}^{(l)}), \text{ if } j \neq 0 \\ D_{i,j}^{(l)} = \frac{1}{m} \Delta_{i,j}^{(l)}, \text{ if } j = 0. \end{cases}$$

Capital-delta matrix D is used as an "accumulator" to add up our values → compute partial derivative.

资本增量矩阵 D 被用于“累加器”来求和。

$$\text{以计算偏导数. } \frac{\partial}{\partial \theta_{i,j}} J(\theta) = D_{i,j}^{(l)}$$

3. Backpropagation Intuition.

1. Cost Function for NN:

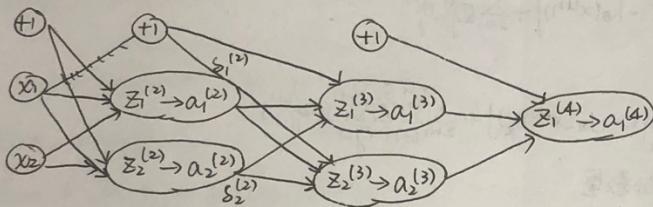
$$J(\theta) = \frac{-1}{m} \sum_{t=1}^m \sum_{k=1}^K [y_k^{(t)} \log(h_\theta(x^{(t)}))_k + (1-y_k^{(t)}) \log(1-h_\theta(x^{(t)}))_k] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} (\theta_{j,i}^{(l)})^2$$

If non-multiclass classification ($K=1$) and disregard regularization:

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

$$S_j^{(l)} = \frac{\partial}{\partial \theta_j^{(l)}} \text{cost}(t)$$

2. Forward propagation:



$$S_2^{(2)} = \theta_{12}^{(2)} * S_1^{(1)} + \theta_{22}^{(2)} * S_2^{(1)}$$

$$S_2^{(3)} = \theta_{12}^{(3)} * S_1^{(2)} + \theta_{22}^{(3)} * S_2^{(2)}$$

12. Backpropagation Implementation: Unrolling Parameters.

1. Use Optimizing function "fminunc()":

fminunc(@costFunction, initialTheta, options)

function [fval, gradientVec] = costFunction(thetaVec)

2. From thetaVec : $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$

Use forward prop/back prop : $D^{(1)}, D^{(2)}, D^{(3)}, J(\theta)$

\downarrow
gradientVec

3.

Matrix. $\begin{bmatrix} \theta^{(1)}, \theta^{(2)}, \theta^{(3)}, \dots \\ D^{(1)}, D^{(2)}, D^{(3)}, \dots \end{bmatrix}$

put into long Vector.

$\begin{bmatrix} \text{thetaVector} = [\text{Theta1}(:); \text{Theta2}(:); \text{Theta3}(:)]; \\ \text{deltaVector} = [D1(:); D2(:); D3(:)]; \end{bmatrix}$

4. If the dimension of Theta1 (10x11), Theta2 (10x11), Theta3 (1x11).

Get back original matrices

$\begin{cases} \text{Theta1} = \text{reshape}(\text{thetaVector}(1:110), 10, 11) \\ \text{Theta2} = \text{reshape}(\text{thetaVector}(111:220), 10, 11) \\ \text{Theta3} = \text{reshape}(\text{thetaVector}(221:231), 1, 11) \end{cases}$

4. Gradient Checking.

1. Approximate the derivative of Cost Function :

$$\frac{\partial}{\partial \theta} J(\theta) \approx \frac{J(\theta + \varepsilon) - J(\theta - \varepsilon)}{2\varepsilon}$$

2. With multiple theta matrix :

$$\frac{\partial}{\partial \theta} J(\theta) \approx \frac{J(\theta_1 \dots \theta_l + \varepsilon \dots \theta_n) - J(\theta_1 \dots \theta_l - \varepsilon \dots \theta_n)}{2\varepsilon}$$

5. Random Initialization. (symmetry breaking)

Initial each $\theta_{ij}^{(l)}$ to a random value in $[-\varepsilon, \varepsilon]$
 $(-\varepsilon \leq \theta_{ij}^{(l)} \leq \varepsilon)$.

Theta1 = rand(10, 11) * (2 * INIT_EPSILON) - INIT_EPSILON
Theta2 = rand(10, 11) * (2 * INIT_EPSILON) - INIT_EPSILON

t. Pick a network architecture.

Number of Input units = dimension of features $x^{(1)}$
Number of output units = Number of classes.
Number of hidden units per layer = usually more the better.
Defaults : 1 hidden layer. (If > 1 , same number of units in every hidden layer).

6. Training a Neural Network.

1. Randomly initialize the weights.

2. Implement forward propagation $\rightarrow h_\theta(x^{(i)})$

3. Implement cost Function.

4. Implement backpropagation \rightarrow partial derivatives

5. Gradient Checking. \rightarrow Then disable it.

6. Gradient descent / optimization function to minimize the cost function with the weights in theta.

② Solving the Problem of Overfitting

1. The problem of overfitting.



(1) Underfitting: hypothesis func h maps poorly to the trend of the data. (Too simple func. too few features)

(2) over fitting: not generalize well to predict new data. (Too complicated func. create unnecessary curves and angles unrelated to data)

(3) To address overfitting:

① reduce the number of features:

- Manually select which features to keep
- Use a model selection algorithm

② Regularization:

- Keep all features, but reduce the magnitude of θ_j
- Regularization works well when we have lots of slightly useful features.

2. cost Function.

$$\theta_0 + \theta_1 x_1 + \theta_2 x_2^2 + \theta_3 x_3^3 + \theta_4 x_4^4$$

$$\text{min}_{\theta} \frac{1}{2m} \sum_{i=1}^m (\theta_0 h(\mathbf{x}^{(i)}) - y^{(i)})^2 + \lambda \underbrace{\theta_0^2 + \theta_1^2 + \dots + \theta_n^2}_{\text{regularize}}$$

(When λ too large \rightarrow underfitting)

3. Regularized Linear Regression.

(1) Gradient Descent.

$$\theta_0 := \theta_0 - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)}) \mathbf{x}_0^{(i)}$$

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

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

(2) Normal Equation.

$$\theta = (\mathbf{x}^T \mathbf{x} + \lambda \mathbf{I})^{-1} \cdot \mathbf{x}^T \cdot \mathbf{y}$$

$$\mathbf{L} = \begin{bmatrix} 0 & & & \\ & 1 & & \\ & & 1 & \\ & & & \ddots & 1 \end{bmatrix}$$

如果 $m < n$, $\mathbf{x}^T \mathbf{x}$ 是不可逆的。

但当添加入 λ 时, $\mathbf{x}^T \mathbf{x} + \lambda \mathbf{I}$ 变为可逆的。

4. Regularized Logistic Regression.

$$\mathbf{x}^T \begin{bmatrix} x_0 & x_1 & x_2 \end{bmatrix} \quad h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2^2 + \theta_3 x_3^3 + \theta_4 x_4^4 + \theta_5 x_5^5 + \dots)$$

Cost Function:

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

Regularize this equation by adding a term:

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

Name: Xiaofan Wang

Chapter 6. Apply Machine Learning.

一些人可以高效运用这些学习算法。

一些人将时间浪费在无效尝试上。

本课：对于机器学习的漫谈，如何选择正确的道路。

1. Evaluating a Learning Algorithm.

1. If hypothesis makes unacceptable large errors in its prediction.

- Get more training examples

- Try smaller sets of features.

- Try getting additional features.

- Try adding polynomial features ($x_1^2, x_2^2, x_1x_2, \dots$)

- Try decreasing λ .

- Try increasing λ .

↓

Machine Learning Diagnose.

2. Hypothesis evaluation.

A hypothesis may have a low error for training examples but still be inaccurate (because of overfitting).

Typically { 70% training set.

{ 30% test set.

↓

New procedure using two sets:

① Learn θ , minimize $J_{\text{train}}(\theta)$ using training set.

② Compute the test set error $J_{\text{test}}(\theta)$

3. Test Error.

① For linear regression: $J_{\text{test}}(\theta) = \frac{1}{2m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} (\theta_0 x_i^{(i)} - y_i^{(i)})^2$

② For classification: (Misclassification Error).

$$\text{err}(\theta_0(x), y) = \begin{cases} 1 & \text{if } \theta_0(x) > 0.5 \text{ and } y=0 \\ 0 & \text{or } \theta_0(x) < 0.5 \text{ and } y=1 \end{cases}$$

$$\text{Test Error} = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \text{err}(\theta_0(x_i^{(i)}), y_i^{(i)})$$

This gives the proportion of test data that was misclassified.

4. Model Selection. & Train. Validation. Test Sets.

① Overfitting problem.

② Model Selection.

$d = \text{degree of polynomial}$.

先选择第一个模型，然后求训练误差的最小值。

可得一个参数向量 θ .

再选第二个模型，过程同理，得到另一个参数向量 θ .

③ Break down the dataset into three sets:

Training set: 60%. $(x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})$

Cross Validation set: 20%. $(x_{cv}^{(1)}, y_{cv}^{(1)}) \dots (x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$

Test Set: 20%. $(x_{\text{test}}^{(1)}, y_{\text{test}}^{(1)}) \dots (x_{\text{test}}^{(m_{\text{test}})}, y_{\text{test}}^{(m_{\text{test}})})$

④ Train/Validation/Test error

$$J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^{m_{\text{train}}} (\theta_0(x_i^{(i)}) - y_i^{(i)})^2$$

$$J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} (\theta_0(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

$$J_{\text{test}}(\theta) = \frac{1}{2m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} (\theta_0(x_{\text{test}}^{(i)}) - y_{\text{test}}^{(i)})^2$$

- Optimize the parameters in θ using the training set for each polynomial degree. 对每个多项式的训练集优化参数 θ .
- Find the polynomial degree d with the least error using the cross validation set. ~~用交叉验证集找到误差最小的模型~~. 通过用交叉验证集找到误差最小的模型
- Estimate the generalization error using the test set with $J_{\text{test}}(\theta)$ ($d = \theta$ from polynomial with lower error). 使用带有 $J_{\text{test}}(\theta)$ 的测试集估计泛化误差。

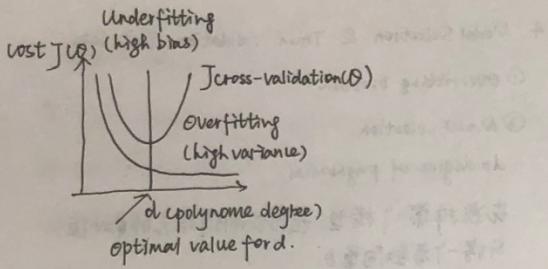
2. Bias & Variance.

1. Diagnosing Bias & Variance

当学习算法不理解根时，偏差较大。（欠拟合）

| 偏差较小。（过拟合）

| 过拟合。欠拟合。



High Bias (underfitting) : $J_{train}(\theta), J_{cv}(\theta)$ will be high.

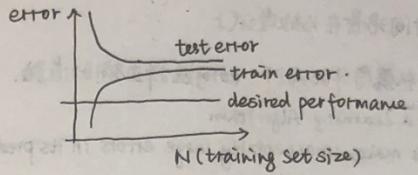
$$J_{train}(\theta) \approx J_{cv}(\theta)$$

High Variance (overfitting) : $J_{train}(\theta)$ will be low.

$$J_{cv}(\theta) \text{ will be much greater than } J_{train}(\theta)$$

3. Learning Curves.

① Experiencing High Bias :



Low training Set Size : $J_{train}(\theta)$ to be low.

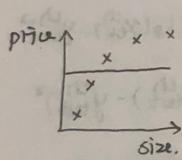
$$J_{cv}(\theta) \text{ to be high.}$$

$$\text{Large training Set Size: } J_{train}(\theta) \approx J_{cv}(\theta)$$

2. Regularization and Bias/Variance.

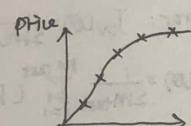
① Linear Regression : $h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_n x^n$

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

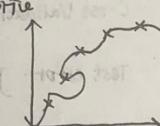


Large λ .

High Bias (underfit)
 $\lambda = 10000, \theta_0 = 0, \theta_1 \approx 0, h_\theta(x) \approx \theta_0$



Intermediate λ .
Just Right.



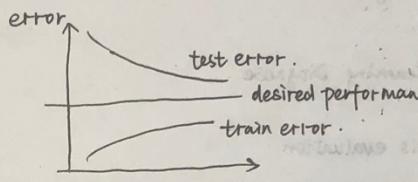
Small λ .

High Variance (overfit)
 $\lambda = 0, \theta_0 = 0, \theta_1 \approx 0, \dots, \theta_n \approx 0$

② Choose the model and the regularization term λ .

- Create a list of lambdas (λ)
- Create a set of models with different degrees or any other variants
- Iterate through λ s and for each λ go through all models to learn some θ .
- Compute the cross validation error using the learned θ (with λ) on the $J_{cv}(\theta)$ without regularization or $\lambda = 0$.
- Select best combo that produces the lowest error on CV test.
- Using the best combo θ and λ , apply it on $J_{test}(\theta)$ to see if it has a good generalization of the problem.

② Experiencing High Variance :



Low training Set Size : $J_{train}(\theta)$ will be low.

$$J_{cv}(\theta) \text{ will be high.}$$

$$\text{Large training Set Size: } J_{train}(\theta) < J_{cv}(\theta)$$

4. Diagnose.

① Getting more training examples: Fix high variance.

② Trying smaller sets of features: Fix high variance.

③ Adding features: Fix high Bias.

④ Adding Polynomial features: Fix high Bias.

⑤ Decrease λ : Fix high Bias.

⑥ Increase λ : Fix high Variance.

5. Diagnose Neural Network.

- ① Fewer parameters \rightarrow underfitting (computationally cheaper)
- ② More parameters \rightarrow overfitting (computationally expensive)

(Hidden Layers \sim Cross Validation Set)

6. Model Complexity Effects

- ① Low-order polynomials (low model complexity)
have high bias and low variance. The model fits poorly consistently.
- ② Higher-order polynomials (high model complexity)
fit training data extremely well and the test data
extremely poorly. These have low bias on the training data,
but very high variance.
- ③ In reality, choose a model somewhere in between.

Name: Xiaofan Wang

Chapter 7 Support Vector Machines.

1. Large Margin Classification.

1. Optimization Objective.

① SVM: 在学习非线性的方程时, 提供了更清晰强大的方法

② Alternative view of logistic regression:

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

If $y=1$, we want $h_{\theta}(x) \approx 1$. $\theta^T x \geq 0$

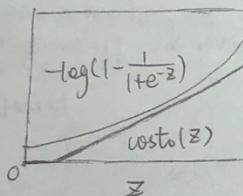
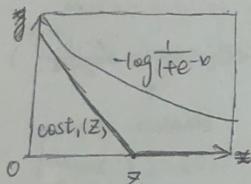
If $y=0$, we want $h_{\theta}(x) \approx 0$. $\theta^T x \leq 0$

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

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

If $y=1$, ($\theta^T x \geq 0$)

If $y=0$ ($\theta^T x \leq 0$)



③ Support Vector Machine

$$\text{logistic regression: } \min_{\theta} \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \underbrace{(-\log h_{\theta}(x^{(i)}) + (1-y^{(i)}) \underbrace{(-\log(1-h_{\theta}(x^{(i)})))}_{\text{cost. } (\theta^T x^{(i)})})}_{\text{cost. } (\theta^T x^{(i)})} \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

$$\text{support vector machine: } \min_{\theta} \left(\frac{1}{m} \sum_{i=1}^m y^{(i)} \underbrace{\text{cost. } (\theta^T x^{(i)})}_{\text{常量}} + (1-y^{(i)}) \underbrace{\text{cost. } (\theta^T x^{(i)})}_{\text{常量}} + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2 \right)$$

可以去掉

$$\text{By } \min_u (u-5)^2 + 1 \Rightarrow u=5$$

$$\min_u 10(u-5)^2 + 10 \Rightarrow u=5$$

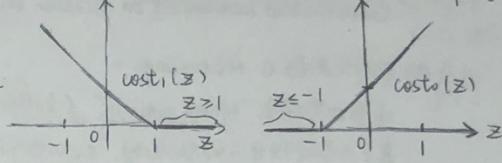
hypothesis:

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

$$\Rightarrow \text{SVM: } \min_{\theta} C \sum_{i=1}^m [y^{(i)} \underbrace{\text{cost. } (\theta^T x^{(i)})}_{\text{常量}} + (1-y^{(i)}) \underbrace{\text{cost. } (\theta^T x^{(i)})}_{\text{常量}}] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

2. Large Margin Intuition

$$\text{SVM: } \min_{\theta} C \sum_{i=1}^m [y^{(i)} \underbrace{\text{cost. } (\theta^T x^{(i)})}_{\text{常量}} + (1-y^{(i)}) \underbrace{\text{cost. } (\theta^T x^{(i)})}_{\text{常量}}] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$



If $y=1$, we want $\theta^T x \geq 1$ (not just ≥ 0)

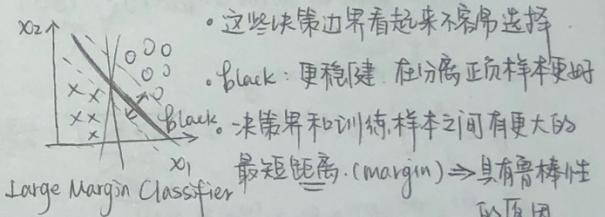
If $y=0$, we want $\theta^T x \leq -1$ (not just < 0)

① SVM Decision Boundary:

$$\begin{cases} \text{Whenever } y^{(i)}=1: \theta^T x^{(i)} \geq 1 \\ \text{Whenever } y^{(i)}=0: \theta^T x^{(i)} \leq -1 \end{cases}$$

$$\min_{\theta} C \times 0 + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

② Linearly Separable case



Large Margin Classifier 最短距离 (margin) \Rightarrow 具有鲁棒性

的原因

用一个最大间距来分离样本

(大间距分类器)

logistic regression: $A + \lambda B$

SVM: $CA + CB$ (C=1时, 两者意义相同. Θ同)

(用参数决定, 是侧重 A Feature, or B Feature)

③ Large Margin Classifier in presence of outliers.

$\Rightarrow C$ very large:

Outlier would have impact on decision boundary

\Rightarrow 正则化参数 C not too large:

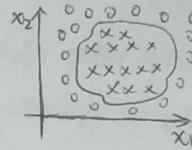
决策边界将会 "ignore outliers". (better decision boundary)

甚至不是线性可分的情况下，SVM也可将它们恰当分开。

= kernels. (核函数)

1. kernels I.

① Non-linear Decision Boundary



Predict $y=1$.

\Rightarrow If $\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 + \theta_5 x_2^2 > 0$.

$$h_0(x) = \begin{cases} 1 & (\text{if } \theta_0 + \theta_1 x_1 + \dots \geq 0) \\ 0 & \text{otherwise.} \end{cases}$$

3. The mathematics behind large margin classification. (Optional)

① vector inner product.

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

$$\vec{u}^T \cdot \vec{v} = [u_1 \ u_2] \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = u_1 v_1 + u_2 v_2.$$

$$(\vec{u}^T \cdot \vec{v} = p \cdot \|u\|, \text{ PGR.})$$

$p = \text{length of projection of } v \text{ onto } u.$

$$\|u\| = \text{length of vector } u = \sqrt{u_1^2 + u_2^2}. \text{ (GR.)}$$

② SVM Decision Boundary.

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

$$\theta^T \cdot x^{(i)} = p^{(i)} \|\theta\| = \theta_1 \cdot x_1^{(i)} + \theta_2 \cdot x_2^{(i)}$$

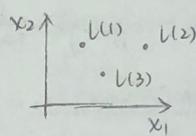
$$(\vec{u}^T \cdot \vec{v})$$

$$\text{s.t. } \begin{cases} p^{(i)} \cdot \|\theta\| \geq 1 & \text{if } y^{(i)} = 1 \\ p^{(i)} \cdot \|\theta\| \leq -1 & \text{if } y^{(i)} = -1 \end{cases}$$

where $p^{(i)}$ is the projection of $x^{(i)}$ onto the vector θ .

Simplification: $\theta_0 = 0$

② Kernel.



Given x , compute new feature depending on proximity to landmarks $L(1), L(2), L(3)$.

$$\text{Given } x: f_1 = \text{similarity}(x, L^{(1)}) = \exp\left(-\frac{\|x - L^{(1)}\|^2}{2\sigma^2}\right)$$

Kernel

Gaussian kernels
 $k(x, L^{(1)})$

$$= \exp\left(-\frac{\sum_{i=1}^n (x_i - L^{(1)})^2}{2\sigma^2}\right)$$

$$\begin{cases} \text{If } x = L^{(1)}: f_1 \approx \exp\left(-\frac{0^2}{2\sigma^2}\right) \approx 1 \\ \text{If } x \text{ far from } L^{(1)}: f_1 = \exp\left(-\frac{(\text{large number})^2}{2\sigma^2}\right) \approx 0 \end{cases}$$

Example:

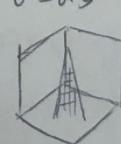
$$L^{(1)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}, \quad f_1 = \exp\left(-\frac{\|x - L^{(1)}\|^2}{2\sigma^2}\right)$$

σ^2 大时，特征变量的值减小速度会变慢。

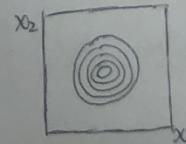
$$\sigma^2 = 1$$

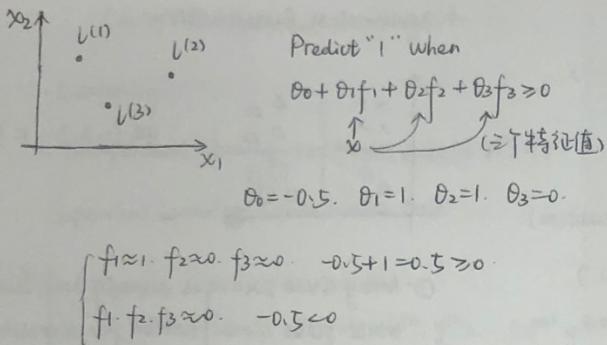


$$\sigma^2 = 0.5$$



$$\sigma^2 = 3$$





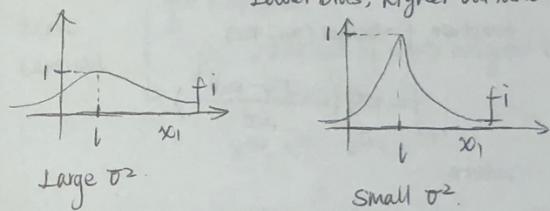
Hypothesis: Given x , compute features $f \in \mathbb{R}^{m+1}$
predict " $y=1$ " if $\theta^T f \geq 0$.
Training: $\min_{\theta} \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T f^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T f^{(i)}) + \frac{1}{2} \sum_{j=1}^m \theta_j^2$.

③ SVM parameters:

$C (= \frac{1}{\lambda}) \Rightarrow$ Large C : Lower bias, high variance. ($s: \lambda$)
Small C : Higher bias, low variance ($L: \lambda$)

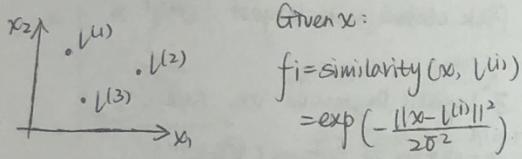
$\sigma^2 \Rightarrow$ Large σ^2 : Features f_i vary more smoothly.
Higher bias, lower variance.

Small σ^2 : Features f_i vary less smoothly.
Lower bias, higher variance.



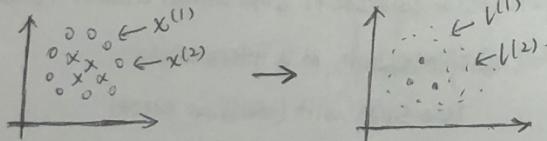
2. kernels II.

① Choose the landmarks.



Predict $y=1$ if $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$.

Where to get $v^{(1)}$, $v^{(2)}$, $v^{(3)}$... ?



将训练样本 $x^{(1)} \dots$ 作为标记点 $v^{(1)}, v^{(2)} \dots v^{(m)}$.

② SVM with Kernels

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

choose $v^{(1)} = x^{(1)}, v^{(2)} = x^{(2)} \dots v^{(m)} = x^{(m)}$

Given example x :

$$f_1 = \text{similarity}(x, v^{(1)}) \quad f = \begin{bmatrix} f_0 \\ f_1 \\ \vdots \end{bmatrix} \quad f_0 = 1$$

$$f_2 = \text{similarity}(x, v^{(2)})$$

For training example $(x^{(i)}, y^{(i)})$:

$$x^{(i)} \rightarrow \begin{cases} f_1^{(i)} = \text{similarity}(x^{(i)}, v^{(1)}) \\ f_2^{(i)} = \text{similarity}(x^{(i)}, v^{(2)}) \end{cases}$$

$$f^{(i)} = \begin{bmatrix} f_0^{(i)} \\ f_1^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix} \quad x^{(i)} \in \mathbb{R}^{m+1}$$

描述样本的特征向量.

≡ Using an SVM.

1. Use SVM software package (e.g. liblinear, libsvm --) to solve for parameters θ .
- Need to specify:
 - Choice of parameter C .
 - Choice of kernel (similarity function).

E.g. No kernel ("linear kernel")
 Predict " $y=1$ " if $\theta^T x \geq 0$

- Gaussian kernel: $f_i = \exp\left(-\frac{\|x - v^{(i)}\|^2}{2\sigma^2}\right)$, where $v^{(i)} = x^{(i)}$
 Need to choose σ^2 .

2. Kernel (similarity) functions:

function $f = \text{kernel}(x_1, x_2)$

$$f = \exp\left(-\frac{\|x_1 - x_2\|^2}{2\sigma^2}\right)$$

Return.

(Note: Do perform feature scaling before using the Gaussian Kernel.)

$$\|x - v\|^2 \Rightarrow v = x - l$$

$$\begin{aligned} \|v\|^2 &= v_1^2 + v_2^2 + \dots + v_n^2 \\ &= (x_1 - l_1)^2 + (x_2 - l_2)^2 + \dots + (x_n - l_n)^2 \end{aligned}$$

3. Other choices of kernel.

Note: Not all similarity functions "similarity (x, v) " make valid kernels.

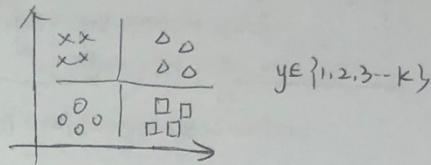
(Need to satisfy technical condition called "Mercer's Theorem")

to make sure SVM packages "optimizations run correctly, and do not diverge".

- ① Many off-the-shelf kernels available:

- Polynomial kernel: $k(x, v) = (x^T \cdot v)^2, (x^T \cdot v)^3, (x^T \cdot v + 1)^3, (x^T \cdot v + 1)^4 \Rightarrow (x^T \cdot v + \text{constant})^{\text{degree}}$
- More esoteric: String kernel, Chi-square kernel, Histogram intersection kernel.

4. Multi-class Classification.



- ① Many SVM packages already have built-in multi-class classification functionality.

- ② Otherwise, use one-vs-all method.

(Train K SVMs, one to distinguish $y=i$ from the rest, for $i=1, 2, \dots, K$)

Pick class i with largest $(\theta^{(i)})^T \cdot x$

5. Logistic Regression vs. SVMs.

n = number of features ($x \in \mathbb{R}^{n+1}$)

m = number of training examples

- ① If n is large, (relative to m):

Use LR or SVM without a kernel (Linear).

- ② If n is small, m is intermediate:

Use SVM with Gaussian kernel.

- ③ If n is small, m is large:

Create/add more features, then use logistic regression or SVM without a kernel.

[Neural Network likely to work well for most of these ~~these~~ settings, but maybe slower to train].

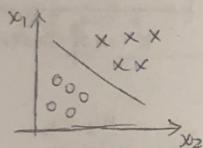
Name: Xiaofan Wang

Ch 8 Unsupervised Learning

- Clustering

1. Unsupervised learning. (从未标记的数据中学习)

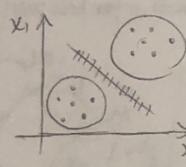
Supervised Learning (with labeled Data)



Training Set:

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

Unsupervised Learning (with unlabeled Data)



Training Set:

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

\Rightarrow Clustering Algorithm

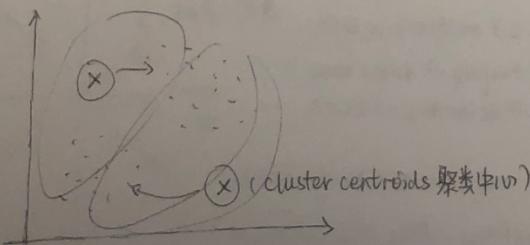
Applications of Clustering:

- Market Segmentation.
- Social Network Analysis.
- Organize Computing Clusters
- Astronomical Data Analysis

2. k-means Algorithm (most popular Clustering Algorithm)



帮助找到紧密的子集或簇。



(1) 簇分配: 将所有点分配到两个不同的聚类中心中。

(2) 移动聚类中心: 迭代下去, 聚类中心不再变, 聚集点也不变. \Rightarrow (k-means 收敛)

k-means Algorithm:

- ① Input:
 - K (number of clusters)
 - Training Set $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$
 - $x^{(i)} \in \mathbb{R}^n$. (drop $x_0=1$ convention).

- ② Randomly Initialize K cluster centroids
 $\mu_1, \mu_2, \mu_3, \dots, \mu_K \in \mathbb{R}^n$.

~~Repeat~~

Cluster Assignment Step

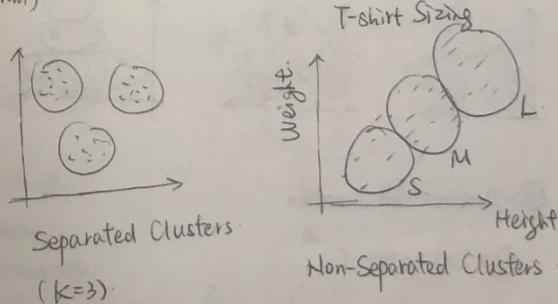
[
for $i=1$ to m
 $c^{(i)}$: index (from 1 to K) of cluster centroid closest to $x^{(i)}$
]

Move Centroid

[
for $k=1$ to K
 μ_k : average of points assigned to cluster k .
]

$$\begin{aligned} & x^{(1)}, x^{(5)}, x^{(6)}, x^{(10)} \\ & \downarrow \\ & c^{(1)}=2, c^{(5)}=2, c^{(6)}=2, c^{(10)}=2 \\ & \mu_2 = \frac{1}{4}[x^{(1)} + x^{(5)} + x^{(6)} + x^{(10)}] \\ & \in \mathbb{R}^n. \end{aligned}$$

③ k-means for non-separated clusters



By: 市场细分
 \downarrow
不同的顾客群体
不同的需求。

3. Clustering: Optimization Objective.

$c^{(i)}$ = index of cluster ($1, 2, \dots, k$) to which example $x^{(i)}$ is currently assigned.

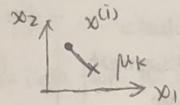
μ_k = cluster centroid k ($\mu_k \in \mathbb{R}^n$)

$\mu_c^{(i)}$ = cluster centroid of cluster to which example $x^{(i)}$ has been assigned.

Optimization Objective:

$$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2$$

$$\begin{aligned} \min \\ c^{(1)}, \dots, c^{(m)} \\ \mu_1, \dots, \mu_k \end{aligned}$$



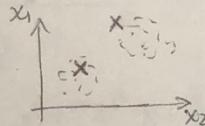
4. Clustering: ~~Random~~ Random Initialization.

① Random Initialization.

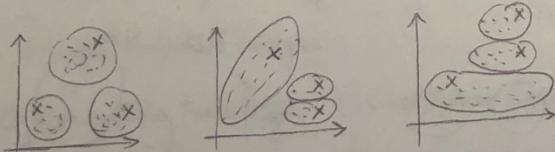
Should have $K < m$

Randomly pick K training examples.

Set μ_1, \dots, μ_K equal to these K examples



② Local Optima.



③ Random Initialization:

For $i = 1$ to 100 {

Randomly initialize K-means.

Run K-means. Get $c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k$.

Compute cost function $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k)$

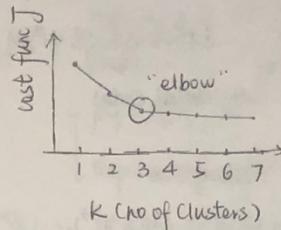
}

Pick Clustering that gave lowest cost $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k)$

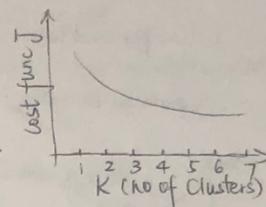
$c_k = 2-10$, 尝试多次随机初始化. 较大影响 $\rightarrow \min J$.

5. Choosing the Number of Clusters

① Elbow Method. (肘部法则)

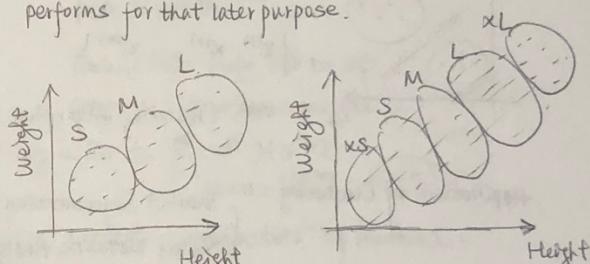


$K=3$. 之后下降非常慢.



② For later / downstream purpose.

Evaluate K-means based on a metric for how well it performs for that later purpose.

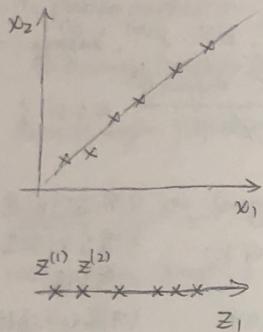


T-shirt Sizing, $K=3$ or $K=5$. \Rightarrow later purpose

二. Dimensionality Reduction. (减少内存，算法提速)

1. Data Compression.

Reduce memory/disk needed to store data.
Speed up learning algorithm.



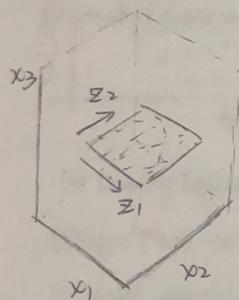
Reduce data from 2D to 1D

$$x^{(1)} \in \mathbb{R}^2 \rightarrow z^{(1)} \in \mathbb{R}$$

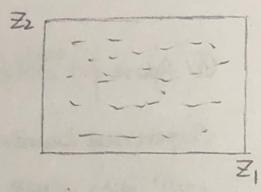
$$x^{(2)} \in \mathbb{R}^2 \rightarrow z^{(2)} \in \mathbb{R}$$

$$\vdots$$

$$x^{(m)} \rightarrow z^{(m)}$$



Reduce Data from 3D to 2D.

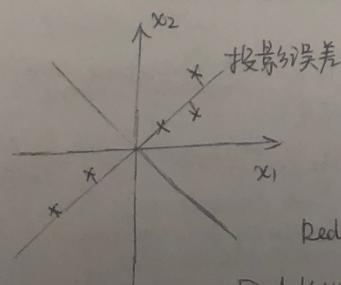


$$z = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \quad z^{(i)} = \begin{bmatrix} z_1^{(i)} \\ z_2^{(i)} \end{bmatrix}$$

2. Data Visualization. (将高维数据可视化为二维或三维)

3. Principal Component Analysis.

① PCA problem Formulation. (不是线性回归)



Reduce from 2D to 1D:

Find a direction (a vector \$u \in \mathbb{R}^n\$)

onto which to project the data

so as to minimize the projection error.

Reduce from nD to kD:

Find \$k\$ vectors \$u^{(1)}, u^{(2)}, \dots, u^{(k)}\$ onto which to project the data, so as to minimize the projection error. (寻找一组向量 \$k\$ 个方向，\$k\$ 维的平面。)

K个方向来定义平面中点的位置。)

② PCA Algorithm.

Data preprocessing:

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

Preprocessing: (feature scaling / mean normalization)

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

Replace each \$x_j^{(i)}\$ with \$x_j^{(i)} - \mu_j

If different features on different scales

(\$x_1 = \text{size of house}, x_2 = \text{number of bedrooms}\$),

Scale features to have comparable range of values.

$$x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \mu_j}{s_j}$$

PCA Algorithm:

Reduce data from nD to kD

Compute "covariance matrix":

$$\Rightarrow \Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)}) (x^{(i)})^T \quad (n \times n) \quad (n \times n)$$

compute "eigenvectors" of matrix \$\Sigma\$:

$$\Rightarrow [U, S, V] = \text{Svd}(\Sigma); \rightarrow n \times n \text{ matrix.}$$

$$\Rightarrow U_{\text{reduce}} = U(:, 1:k);$$

$$\Rightarrow z = U_{\text{reduce}} * x; \quad \begin{cases} \text{Svd: singular value decomposition} \\ \text{eig}(\Sigma): 奇偶值分解. \end{cases}$$

$$U = \underbrace{\begin{bmatrix} & & & \\ & & & \\ U^{(1)} & U^{(2)} & \dots & U^{(m)} \\ & & & \\ & & & \end{bmatrix}}_{K} \in \mathbb{R}^{n \times n}$$

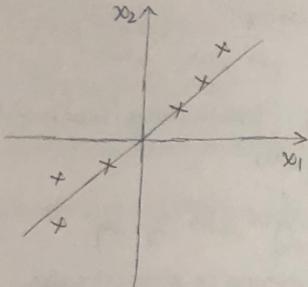
$$K \cdot x \in \mathbb{R}^K \rightarrow z \in \mathbb{R}^K$$

$$\Downarrow$$

$$z = \underbrace{\begin{bmatrix} & & & \\ & & & \\ U^{(1)} & U^{(2)} & \dots & U^{(k)} \\ & & & \\ & & & \end{bmatrix}^T}_{n \times K} \cdot x \quad (n \times 1)$$

$$U_{\text{reduce}}$$

③ Reconstruction from compressed representation



$$z = U_{\text{reduce}}^T x$$

$$z \in \mathbb{R}^k \rightarrow x \in \mathbb{R}^n$$

$$\begin{aligned} x_{\text{approx}}^{(i)} &= U_{\text{reduce}} \cdot z^{(i)} && (\text{z 映射回原分布}) \\ (\in \mathbb{R}^n) &\quad \underbrace{n \times k}_{\text{ }} \quad \underbrace{k \times 1}_{\text{ }} && (\text{高维数据的近似值}) \\ &\quad \underbrace{\text{ }}_{n \times 1} \end{aligned}$$

④ Choosing the number of Principal Components (k)

(主成分数量)

$$\text{Average squared projection error: } \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2$$

$$\text{Total variation in the data: } \frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2$$

Typically, choose k to be smallest value so that:

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

| | |
|------|-----|
| 0.05 | 5% |
| 0.10 | 10% |

99% of variance is retained.

95%

90%

许多特征高度相关

Choosing k (number of Principal components):

Algorithm:

Try PCA with k=1, 2, ...

$$\text{Compute } U_{\text{reduce}}, z^{(1)}, z^{(2)}, \dots, z^{(m)}, x_{\text{approx}}^{(1)}, \dots, x_{\text{approx}}^{(m)}$$

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

$[U, S, V] = \text{svd}(\Sigma)$ \Rightarrow 只需运行此步

Pick smallest value of k: 可得 S 对角线上的值

$$\frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^m S_{ii}} \geq 0.99 \quad \begin{array}{l} \text{改变 k 值} \\ \text{(不需要从头运行)} \end{array}$$

$(99\% \text{ of variance retained})$

⑤ Advice for applying PCA

Supervised Learning Speedup:

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

Extract inputs:

$$\begin{aligned} \text{Unlabeled dataset: } x^{(1)}, x^{(2)}, \dots, x^{(m)} &\in \mathbb{R}^{10000} \\ x \rightarrow z: h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}} \rightarrow z^{(1)}, z^{(2)}, \dots, z^{(m)} &\in \mathbb{R}^{10000} \end{aligned}$$

$$\text{New training set: } (z^{(1)}, y^{(1)}), (z^{(2)}, y^{(2)}), \dots$$

Note: Mapping $x^{(i)} \rightarrow z^{(i)}$ should be defined by running PCA only on the training set. This mapping can be applied as well to the examples $x_{\text{cv}}^{(i)}$ and $x_{\text{test}}^{(i)}$ in the cross validation and test sets.

⑥ Bad Use of PCA: To prevent overfitting. (x)

(更容易舍弃有价值的信息)

$$\Rightarrow \text{Use regularization: } \min_{\theta} \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

⑦ PCA shouldn't be: design of ML system.

⑧ Application of PCA

Data Compression: reduce memory & speedup algorithm

Data Visualization: 2D 3D.

Name: Xiaofan Wang

Chapter 9. Problem Motivation. (Anomaly Detection)

- Density Estimation.

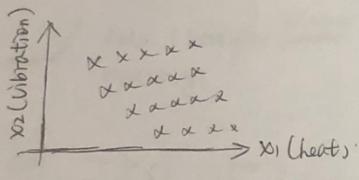
1. Problem Motivation.

① Anomaly Detection Example.

Aircraft engine features: $x_1 = \text{heat generated}$
 $x_2 = \text{vibration intensity}$

Dataset: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

New Engine: x_{test}



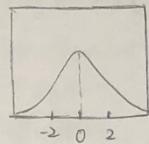
② Density Estimation

Dataset $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

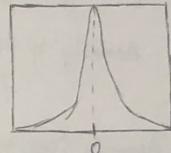
Is x_{test} anomalous? Model $p(x)$

2. Gaussian Distribution Example.

$\mu=0, \sigma=1$



$\mu=0, \sigma=0.5$



③ Example.

Fraud Detection: $x^{(i)} = \text{features of user } i \text{ activities}$

Model $p(x)$ from data.

Identify unusual users by checking $p(x) < \epsilon$.

3. Parameter Estimation.

Dataset $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}, x^{(i)} \in \mathbb{R}$

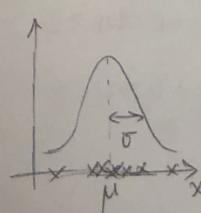
假设每个样本服从正态/高斯分布 $x^{(i)} \sim N(\mu, \sigma^2)$

Manufacturing.

Monitoring computers in a data center:

$x^{(i)} = \text{features of machine } i$
 $x_1 = \text{memory use},$
 $x_2 = \text{number of disk} \sim$
 $x_3 = \text{CPU load},$
 $x_4 = \text{CPU load / network traffic}$

参数估计：给定数据集，希望可以找到能估算 μ 和 σ 的值。



其中 μ 对应分布函数的中心。

σ 控制高斯分布的宽度。

很好地拟合了数据。
分布在中心区域的概率较大。

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

三. Algorithm. (应用高斯分布开发异常检测算法)

1. Anomaly Detection Algorithm.

① Choose features x_i that you think might be
indicative of anomalous examples. (描述数据相关属性特征)

② Fit parameters $\mu_1 \dots \mu_n, \sigma_1^2 \dots \sigma_n^2$.

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

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

③ Given new example x , compute $p(x)$:

$$p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi} \sigma_j} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)$$

Anomaly if $p(x) < \epsilon$.

2. Anomaly Detection Example.

