

1.8 Advice for applying machine learning

8-1 Deciding what to try next

- How to debugging a learning algorithm
 - Get more training examples
 - Try smaller sets of features
 - Try getting additional features
 - Try adding polynomial features (x_1, x_2, \dots etc.)
 - Try decreasing or increasing λ

(但如何选择?)

8-2 Evaluating a hypothesis 评估假设

- 评估假设的一种常用方式
将 "Dataset" 分割为 {
 - Training set (70%)
 - Test set (30%)
- [note: 随机分组]

• Training / testing procedure for linear regression

- Learn parameter θ from training set
- Compute test set error:

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

• ... for logistic regression

- learn parameter θ from training set
- Compute test set error:

$$J_{\text{test}}(\theta) = -\frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} y^{(i)}_{\text{test}} \log(\theta_0(x^{(i)}_{\text{test}})) + (1 - y^{(i)}_{\text{test}}) \log(1 - \theta_0(x^{(i)}_{\text{test}}))$$

(0/1 misclassification error):

$$\text{error}(\theta_0(x), y) = \begin{cases} 1 & \text{if } \theta_0(x) \geq 0.5, y=0 \\ 0 & \text{or if } \theta_0(x) < 0.5, y=1 \\ \text{0.5} & \text{otherwise} \end{cases}$$

$$\text{Test error} = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \text{error}(\theta_0(x^{(i)}_{\text{test}}), y^{(i)}_{\text{test}})$$

8.3 Model selection and training / validation / test set

为了检验模型在新样本上的泛化能力。

将“Dataset”分为三部分，
 Training set (60%)
 Cross validation (20%) 交叉验证
 Test set (20%) Set

训练误差 Training error 即代价函数

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

【选择模型】· Cross validation error

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

【测试泛化能力】· Test error

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

e.g.

Model selecting

$$\theta=1 \quad 1. h_{\theta}(x) = \theta_0 + \theta_1 x \xrightarrow{\min J_{\text{train}}(\theta)} \theta^{(1)} \rightarrow J_{\text{CV}}(\theta^{(1)})$$

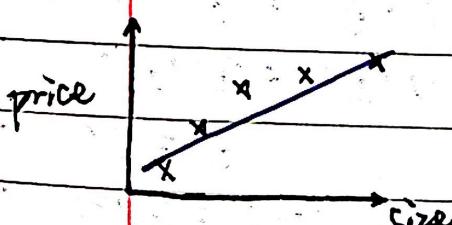
$$\theta=2 \quad 2. h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \xrightarrow{\theta^{(2)}} J_{\text{CV}}(\theta^{(2)}) \quad \left. \begin{array}{l} \text{选择 } \min J_{\text{CV}}(\theta^{(i)}) \\ \text{作为最佳模型} \end{array} \right.$$

$$\theta=10 \quad 10. h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_{10} x^{10} \xrightarrow{\theta^{(10)}} J_{\text{CV}}(\theta^{(10)})$$

⇒ Estimate generalization error for test set $J_{\text{test}}(\theta^{(i)})$

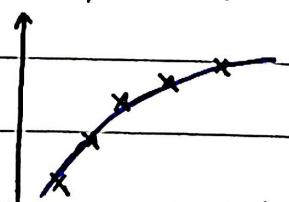
8.4 Diagnosing bias vs. variance 诊断偏差和方差

最佳模型

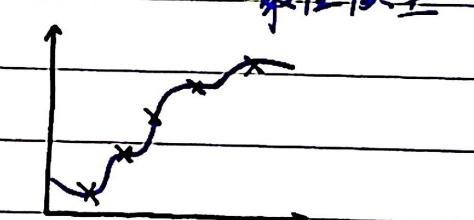


High bias
(underfit)

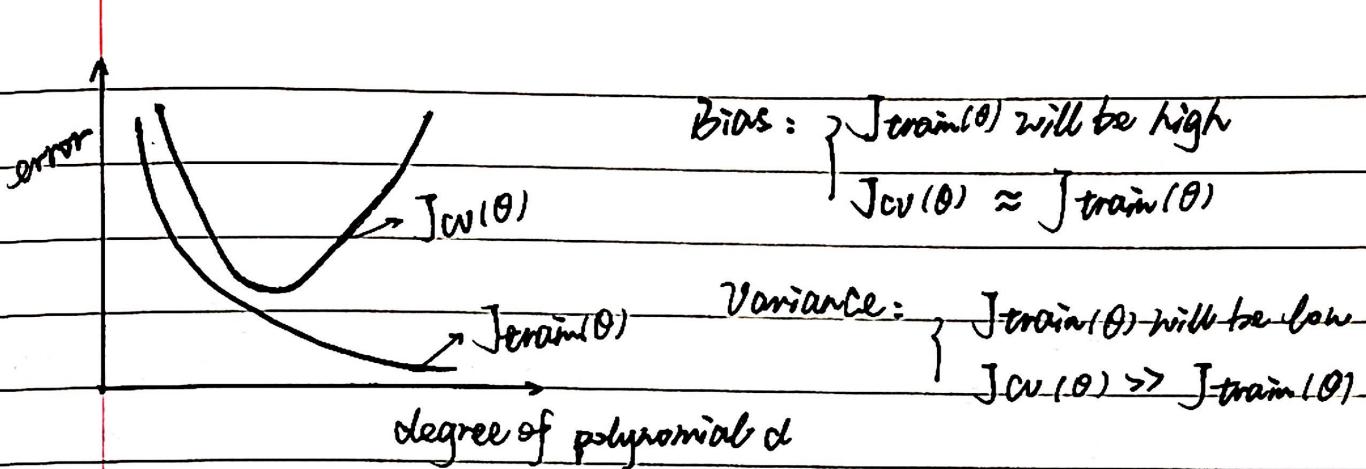
欠拟合 (拟合很差)



"Just right"



"High variance"
(overfit)
过拟合 (泛化误差)

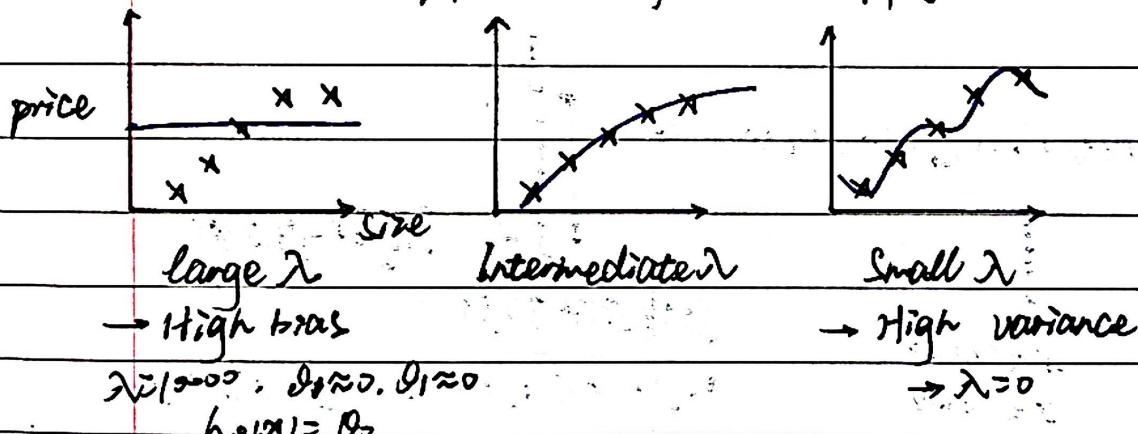


8.5 Regularization and bias/variance

- linear regression with regularization

Model: $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

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



那么该不该 λ 大小？ \rightarrow 训练时才正则化

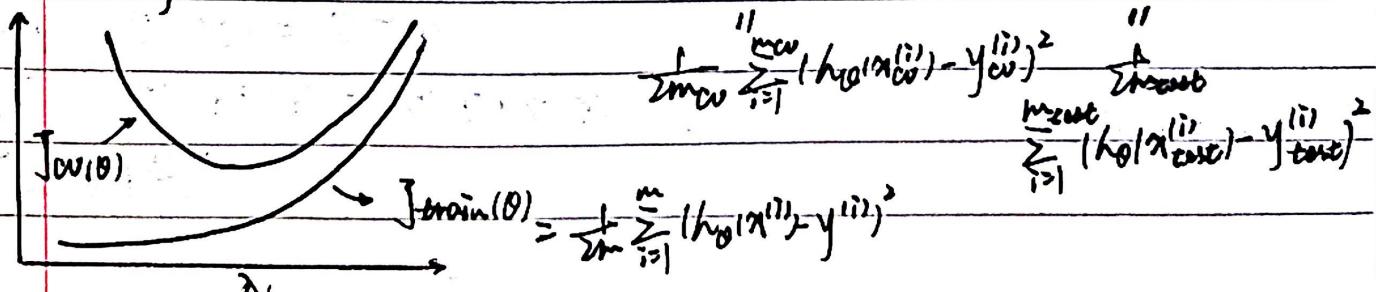
1. Try $\lambda = 0 \rightarrow \min J(\theta) \rightarrow \theta^{(1)} \rightarrow J_{cv}(\theta^{(1)})$

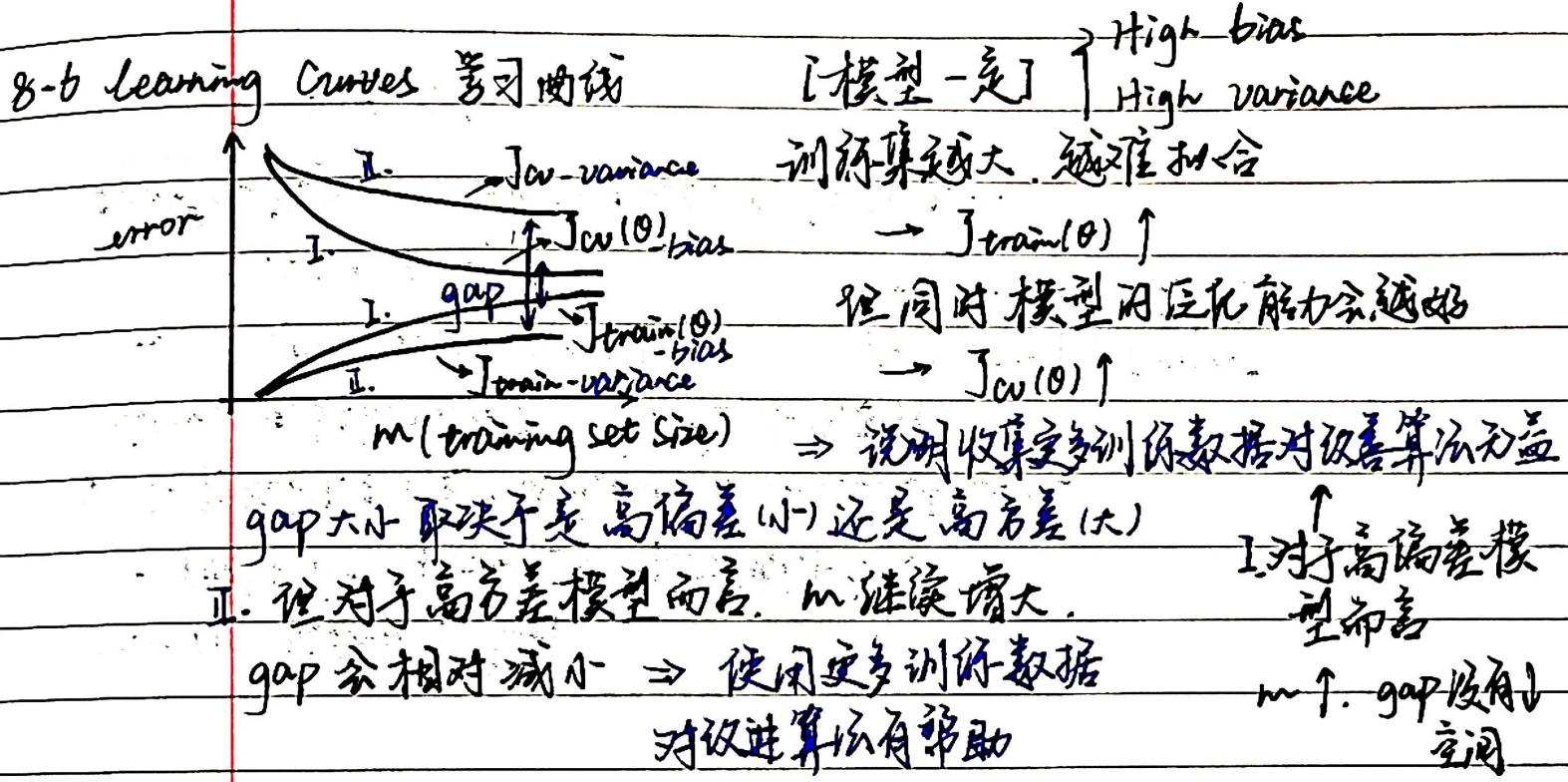
2. Try $\lambda = 100 \rightarrow \theta^{(2)} \rightarrow J_{cv}(\theta^{(2)})$

3. Try $\lambda = 0.02 \rightarrow \theta^{(3)} \rightarrow J_{cv}(\theta^{(3)})$ \rightarrow Pick $\theta^{(3)}$ (最佳模型)

Test error:

12. Try $\lambda = 10.24 \rightarrow \theta^{(12)} \rightarrow J_{cv}(\theta^{(12)}) \rightarrow J_{test}(\theta)$





8-7 神经

- Get more training example \rightarrow 适用于 High variance

通过绘制学习曲线来判断, i.e. $J_{cv}(\theta) \gg J_{train}(\theta)$

- Try ~~more~~ smaller sets of features \rightarrow 适用于 High variance
- Try getting additional features \rightarrow 适用于 High bias
- Try adding polynomial features
- Try decreasing λ \rightarrow 适用于 High bias
- Try increasing λ \rightarrow 适用于 High variance

对于神经网络来说

- "small" Neural Network $\begin{cases} \text{神经元个数较少, 易出现欠拟合, } \\ \text{hidden layer}\end{cases}$ 计算量小
- "large" Neural Network, \sim 较多, 易出现过拟合, 计算量大, 但上限更高
Use Regularization (λ)

29 Machine learning system design

9-1 Prioritizing what to work on 优先级排序

E.g. Building a spam classifier

Supervised learning: $x = \text{features of email}$, $y = \text{spam}(1) \text{ or not spam}(0)$.

Feature x : Choose 100 words indicative of spam/not spam.

like A, B, C, D ... $\xrightarrow{\text{假设}} x \text{ input} = [0 \ 1 \ 1 \ 0 \ 0 \dots 1 \dots]^T$

→ how to make it have low error?

(这件垃圾邮件是1. 非垃圾邮件0)

- Collect lots of data?
- Develop sophisticated features based on email routing information
- Develop sophisticated features for message body. (from email header)
e.g. Should 'discount' and 'discounts' be treated as the same word?
- Develop sophisticated algorithm to detect misspellings
e.g. medicine, w4tch's

$$x_j = \begin{cases} 1 & \text{if word } j \text{ appears in email} \\ 0 & \text{otherwise.} \end{cases}$$

9-2 Error analysis

• Recommended approach

- Start with a simple algorithm that you can implement quickly.

Implement it and test it on your cross-validation data.

- Plot learning curves to decide if more data, more features, etc. are likely to help.

- Error analysis: Manually examine the examples (in cross-validation set) that your algorithm made errors on.

⇒ 从而启发如何改进算法 (新特征?)

错误率

5.4%

3.7%

• The importance of numerical evaluation e.g. 正确大小写

数值评估

不正确

⇒ 错误不区分

〔统计分类〕

9.3 Error metrics for skewed classes 不对称性分类的误差评估

• Cancer classification example

Train logistic algorithm model $h_0(x)$ to find that you got 1% error on the test, but only 0.5% of patients have cancer.

(为了解决斜偏类问题) → Precision / Recall [精确度] 和 [召回率] 评估指标

E.g. Actual class

		1	0
		True positive	False positive
Predicted class	1	True positive	
	0	False negative	True negative

• Precision

$$= \frac{\text{True Positive}}{\# \text{predicted positive}} = \frac{TP}{TP + FP} \quad (\text{只算正例})$$

• Recall

$$= \frac{\text{True positive}}{\# \text{Actual positive}} = \frac{TP}{TP + FN}$$

9.4 Trading off precision and recall

• Cancer classification example

- logistic regression: $0 \leq h_0(x) \leq 1$ { predict 1 if $h_0(x) \geq 0.5$
0 if $h_0(x) < 0.5$

Suppose we want to predict $y=1$ (cancer only if very confident)

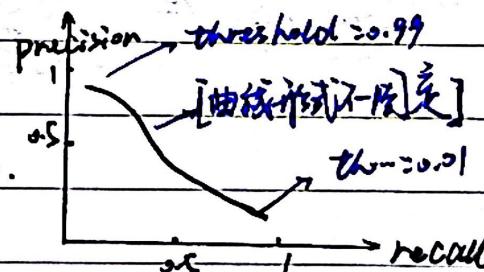
→ 改进阈值 like { 1 if $h_0(x) \geq 0.9$ ⇒ higher precision, lower recall.

Suppose we don't want to avoid missing too many cases of cancer

→ like { 1 if $h_0(x) \geq 0.3$ ⇒ higher recall, lower precision

More generally, predict $y=1$ if $h_0(x) \geq \text{threshold}$

对于大多数分类器，需要权衡 precision vs recall.



How to compare?

$$\rightarrow F_1 \text{ Score} = 2 \frac{PR}{P+R}$$

Algorithm	Precision	Recall	F_1 Score
1	0.5	0.16	0.464
2	0.7	0.1	0.175
3	0.02	1.0	0.0292

9-5 Data for machine learning

- When feature $x \in \mathbb{R}^{n+1}$ has sufficient information to predict y accurately → 使用大量数据才有意义
- Large data rationale
 - Use a learning algorithm with many parameters (e.g. logistic regression / linear regression with many features; neural network with many hidden units). → 防止过拟合
 - Use a large training set (足够大, \Rightarrow parameters \Rightarrow unlikely to overfit) → 防止方差

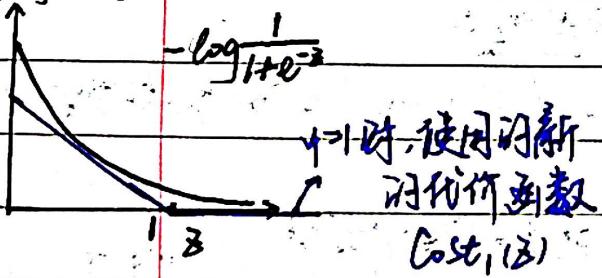
10 Support Vector Machines 支持向量机

10-1 Optimization objective 线性目标

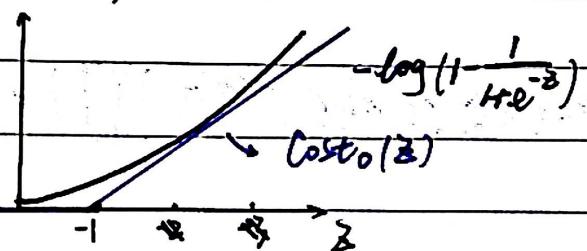
- Alternative view of logistic regression

$$\begin{aligned} \text{Cost of example: } & -y \log h_{\theta}(x) + (1-y) \log(1-h_{\theta}(x)) \\ & = -y \log \frac{1}{1+e^{-\theta^T x}} - (1-y) \log \left(1 - \frac{1}{1+e^{-\theta^T x}}\right) \end{aligned}$$

- If $y=1$ (want $\theta^T x \gg 0$)



- If $y=0$ (want $\theta^T x \ll 0$)



Support vector machine 逻辑回归的对称

$$\frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2$$

- Logistic regression: $\min_{\theta} \frac{1}{m} \left[\sum_{i=1}^m [y^{(i)}(-\log h_{\theta}(x^{(i)}) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)}))] \right] + \text{constant}$

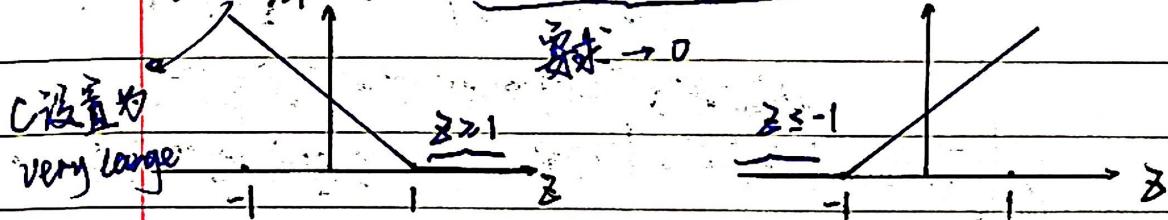
- SVM: $\min_{\theta} C \sum_{i=1}^m [y^{(i)} \text{Cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{Cost}_0(\theta^T x^{(i)})] + \frac{1}{2} \sum_{j=1}^m \theta_j^2$

其中该类数为 1
 $h_{\theta}(x) \begin{cases} 1, & \text{if } \theta^T x \geq 0 \\ 0, & \text{otherwise} \end{cases}$ 并不会输出概率

对大间隔的直观理解

10-2 Large Margin Intuition (SVM有时被称为大间隔分类器)

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



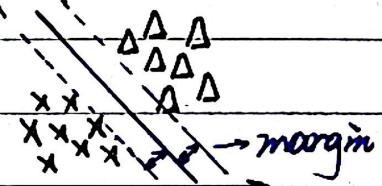
if $y=1$, want $\theta^T x \geq 1$ (not just > 0) → 在“安全间隔”

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

(1)
②

SVM的决策边界会使得margin变大.

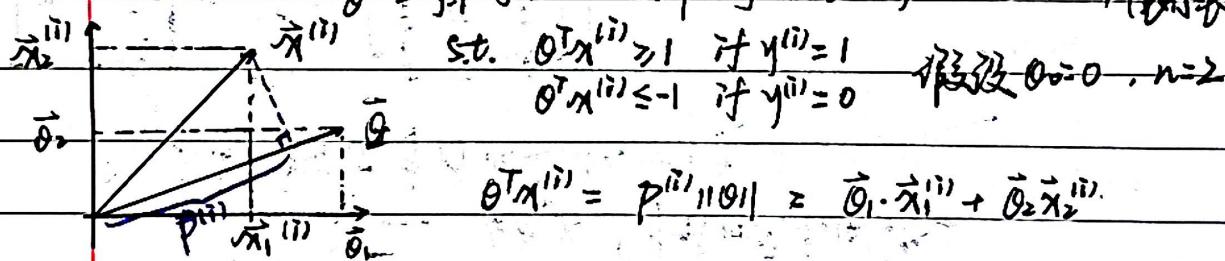
会让SVM具有鲁棒性.



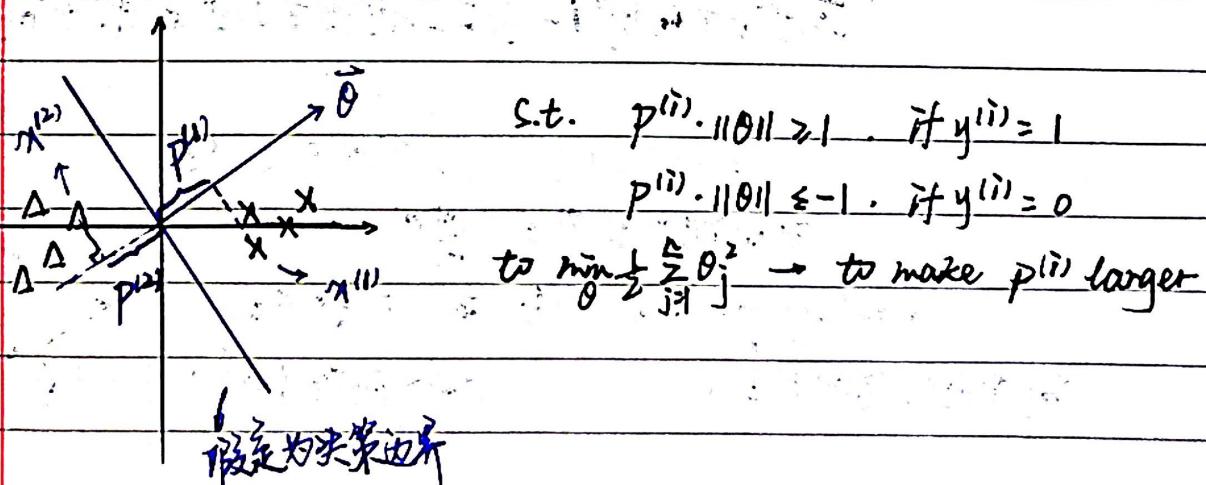
- C值设置越大, 对新数据点(异常点)越敏感. → 所以让C小一点)

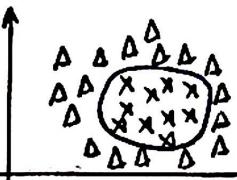
10-3 大间隔分类器的数学原理

- 向量化. $\min \frac{1}{2} \sum_j \theta_j^2 \rightarrow$ 看作 $\vec{\theta} = \vec{\theta}_1 + \vec{\theta}_2$ → $\frac{1}{2} \| \vec{\theta} \|_2^2$ (表示长度/范数)



- 系数向量 θ 事实上是与决策边界垂直的 $\Rightarrow \theta_1 x_1 + \theta_2 x_2 = 0$





对 Non-linear 问题

$$h_0(x) = \begin{cases} 1 & \text{if } \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 + \theta_5 x_2^2 \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

10-4. Kernels I

\Rightarrow 相比于高阶项，是否有更好的特征选择？

- 建立新特征 (Given x , compute new feature depending on proximity)
- 先标记: on marks $l^{(1)}, l^{(2)}, l^{(3)}$) 实际可选择更多
(and)

$$f_1 = \text{similarity}(x, l^{(1)}) = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right) = \exp\left(-\frac{\sum_j (x_j - l_j^{(1)})^2}{2\sigma^2}\right)$$

$$f_2 = \text{similarity}(x, l^{(2)}) = \exp\left(-\frac{\|x - l^{(2)}\|^2}{2\sigma^2}\right)$$

$$f_3 = \dots \rightarrow k(x, l^{(i)}) \quad \text{Gaussian kernels}$$

if $x \approx l^{(1)}$: $f_1 \approx \exp(-\frac{\sigma^2}{2\sigma^2}) \approx 1$ - σ^2 决定了特征变量减小速度

if x is far from $l^{(1)}$: $f_1 \approx 0$, 值越小，速度越快。

- 通过标准化和相似性函数定义新的特征变量 → 训练复杂的非线性边界

10-5 kernels II • Where / How to get $l^{(1)}, l^{(2)}, l^{(3)}, \dots$? (选择样本点)

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

\rightarrow Choose $x^{(1)} = l^{(1)}, x^{(2)} = l^{(2)}, \dots, x^{(m)} = l^{(m)}$

- Given example x : $f_1 = \text{similarity}(x, l^{(1)}), f_2 = \text{similarity}(x, l^{(2)}), \dots$

$$\rightarrow f = [f_0, f_1, \dots, f_m]^T, f_0 = 1 \quad f_m = \text{sim}(x, l^{(m)})$$

- For training example $(x^{(i)}, y^{(i)})$ $x^{(i)} \in \mathbb{R}^{m+1}$

$$x^{(i)} \Rightarrow f^{(i)} = \text{sim}(x^{(i)}, l^{(i)}) \Rightarrow f^{(i)} = \begin{bmatrix} f_0^{(i)} \\ f_1^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix}, f_0^{(i)} = 1$$

Hypothesis: Given x , compute feature $f \in \mathbb{R}^{m+1}$

Predict " $y=1$ " if $\theta^T f \geq 0$ ($\theta_0 f_0 + \theta_1 f_1 + \theta_2 f_2 + \dots + \theta_m f_m$)

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

参数选择

$C (=\frac{1}{\lambda}) \rightarrow$ large C : lower bias, higher variance

small C : higher bias, lower variance

$\sigma^2 \rightarrow$ large σ^2 : Feature f_i vary more smoothly, higher bias, lower variance

small σ^2 : f_i vary less smoothly, lower bias, higher variance.

10-6 Using an SVM

• Use SVM software package (e.g. Sklearn) 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 $\theta^T x \geq 0 \rightarrow \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n \geq 0, x \in \mathbb{R}^n$

• Gaussian kernel

$$f_i = \exp\left(-\frac{\|x - b^{(i)}\|^2}{2\sigma^2}\right), \text{ where } b^{(i)} = x^{(i)}, \text{ Need to choose } \sigma^2$$

Note: Do perform [feature scaling] before using the Gaussian kernel.

因为 $\|x - b^{(i)}\|^2 \rightarrow (x_1 - b_1)^2 + (x_2 - b_2)^2 + \dots + (x_n - b_n)^2$

教室

不同特征维度的数据值差异可能会很大，例如 x_1 表示房子面积， x_2 表示房间数。

• 对于 Multi-class classification (i.e. kSVM)

原理和 4-7-一致，one to distinguish $y=i$ from the rest, for $i=1, 2, \dots, k$,
get $\theta^{(1)} (y=1), \theta^{(2)} (y=2), \dots, \theta^{(k)}$: pick class i with largest $(\theta^{(i)})^T x$

• 选择区间 vs. SVM: (n 表示特征数, m 表示训练数据数)

- If n is large (relative to m), e.g. $n=10000, m=10 \sim 100$

use a logistic regression, or SVM without a kernel ("linear kernel")

- If n is small, m is intermediate, e.g. $n=1 \sim 100, m=10 \sim 10000$

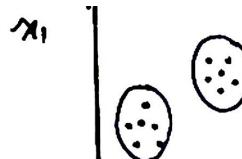
use SVM with Gaussian kernel

- If n is small, m is large, e.g. $n=1 \sim 100, m=50000 +$

add/create more feature, then use logistic regression or SVM
without a kernel.

△ Neural network likely to work well for most settings,
but may be slower to train.

4.1 Clustering 聚类 → 划值



Clustering algorithm

4.1-1 Intro- · 未标签, train set: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\} \rightarrow$ 找隐藏在数据中的结构

Unsupervised learning · 应用: ①-3

随机选择

4.1-2 K-means Algorithm

- K-means 是一个迭代算法 \rightarrow 移动聚类中心

每次内循环第一步: 进行簇分配, 样本与哪个中心

哪个更近来将样本点分配给哪个聚类中心之一 (遍历每一个点)

第二步: 移动聚类中心, 将两个中心分别移动到其分配点的均值处
然后进行迭代

Algorithm

Set number of clusters

(drop $x_0 = 1$ convention)

Inputs: Training set $\{x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(m)}\}, x^{(i)} \in \mathbb{R}^n$.

Randomly initialize K cluster centroid $\mu_1, \mu_2, \dots, \mu_K \in \mathbb{R}^n$

Repeat { for $i = 1$ to m

样本点所属 (C⁽ⁱ⁾) = index (from 1 to K) of cluster centroid closest to $x^{(i)}$
归簇类 for $k = 1$ to K

$$\min_k \|x^{(i)} - \mu_k\|^2$$

将 x_i 分配给 C⁽ⁱ⁾

$$C^{(i)} = k$$

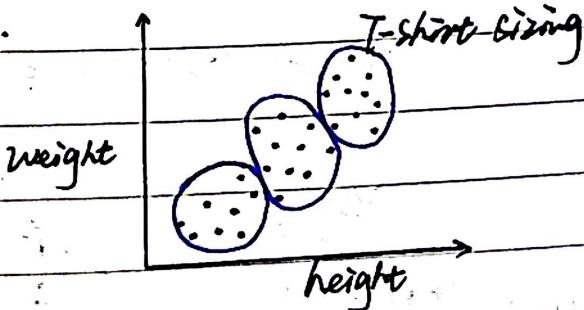
μ_k := average (mean) of points assigned to cluster k

$$\text{e.g. } x^{(1)}, x^{(5)}, x^{(6)}, x^{(10)} \rightarrow C^{(1)} = 2, C^{(5)} = 2, C^{(6)} = 2, C^{(10)} = 2$$

$$\Rightarrow \mu_2 = \frac{1}{4} [x^{(1)} + x^{(5)} + x^{(6)} + x^{(10)}] \in \mathbb{R}^n$$

• K-means 对于 non-separated clusters 也有效

E.g.



(假设 $x^{(i)}$ 被划为第 5 个簇)

$$\text{e.g. } x^{(i)} \rightarrow 5, \text{ 即 } c^{(i)} = 5, \Rightarrow \mu_{c^{(i)}} = \mu_5$$

11-3 Optimization objective

(1) 11-2) $c^{(i)} = \dots, \mu_1 = \dots, \mu_k = \dots$, cluster centroid of cluster to which example $x^{(i)}$ has been assigned.

优化目标函数: (Distortion, 失真代价函数)

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

$$\min_{c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k} J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k)$$

- 实际上, 11-2 中 (第一步(分配)) 就是在 $\min J$ w.r.t. $c^{(1)}, \dots, c^{(m)}$.
(第二步(移动)) 就是在 $\min J$ w.r.t. μ_1, \dots, μ_k .

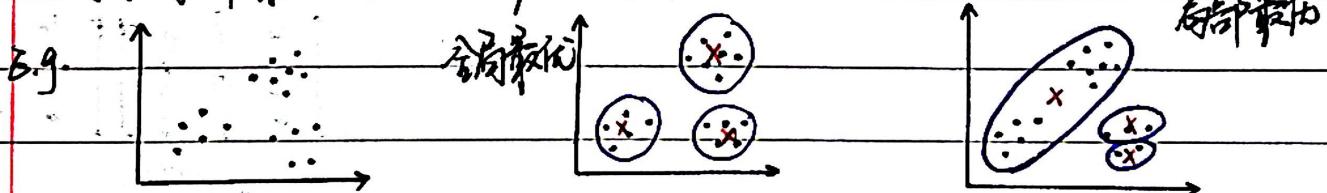
11-4 Random initialization 随机初始化

- 选择样本点作为聚类中心

Should have $K < m$, randomly pick k training examples.

Set μ_1, \dots, μ_k equal to these k examples. ($\mu_1 = x^{(i)}, \mu_2 = x^{(j)}, \dots$)

- 对于局部最优 (local optima)



\Rightarrow 对 k-means Algorithm 进行多次初始化 (对于 k 值相对较小的情况)
例如 $k=2 \sim 10$

for $i=1$ to 100 {

理论上，因为无标签 \Rightarrow 无法设计 Algorithm 因为选择大：

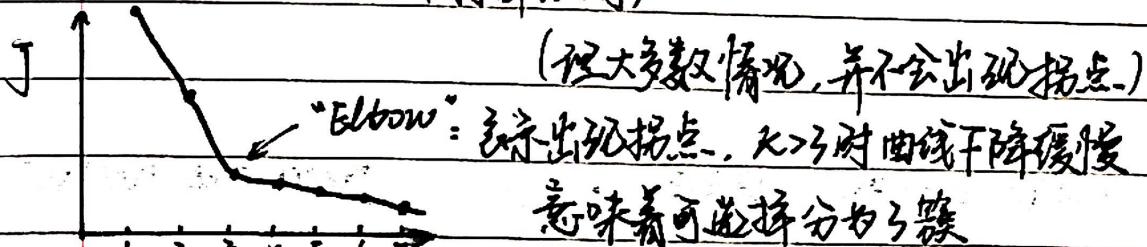
答案可能不止一个。

一般高类运动选择

观察可视化的内容：

11-5 Choosing the number of cluster 观察聚类算法的输出, etc.

• Elbow method (肘部法则)



(e.g. 市场分割)

• 下游目的, i.e. 看哪个聚类数量能更好地应用于后续目的

12 Dimensionality Reduction 降维

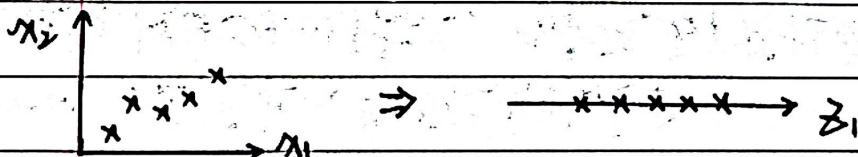
12-1 Motivation I: Data Compression 目标 I: 数据压缩

让数据占用较少的内存或硬盘空间, [如果特征高度相关
对学习算法进行加速] \rightarrow 高度冗余

E.g. Reduce data from 2D to 1D

度冗

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



Reduce data from 3D to 2D, i.e. 立体转化为平面

12-2 Motivation II: Data Visualization (可视化)

E.g. 各个国家特征 $x_1 = \text{GDP}$, $x_2 = \text{Poverty index}$, $x_3 = \text{Life expectancy} \dots$

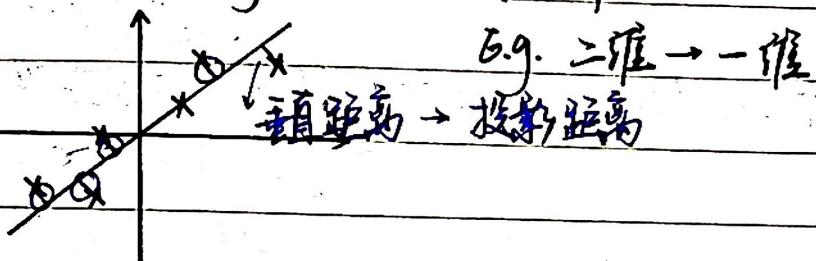
$x \in \mathbb{R}^{50}$ \Rightarrow 将 x 降低到 2 维, 2.1 表示-GDP 总值

2.2 表示人均 GDP

(PCA: 处理降维问题的算法)

12-3 Principal Component Analysis problem formulation 1 主成分分析问题规划 1

- PCA → find k vectors $u^{(1)}, \dots, u^{(k)}$, 将数据投影到这 k 个向量展开的 [Reduced from n -dimension to k -dimension] 线性空间上，并最小化投影误差
- PCA for linear regression 是不同的算法。 数据点到空间的 距离



12-4 主成分分析问题规划 2

- Data preprocessing: Training set: $x^{(1)}, x^{(2)}, \dots, x^{(m)}$

Preprocessing (feature scale / [mean normalization])

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)} \quad \text{Replace each } x_j^{(i)} \text{ with } x_j - \mu_j$$

If different features on different scales (e.g. x_1 = size of house, x_2 = number of bedrooms), scale features to have comparable range of values

- Algorithm: (Reduce data from n -d to k -d) n 维 → k 维

Compute "Covariance matrix" (协方差) $\rightarrow \Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)}) (x^{(i)})^T$

Compute "eigenvectors" of matrix Σ : $[U, S, V] = \text{svd}(\Sigma)$

To get: 特征向量

$$U = [u^{(1)} \dots u^{(n)}] \in \mathbb{R}^{n \times n} \quad \text{选取前 } k \text{ 列 } u^{(1)} \rightarrow u^{(k)}$$

$$Z = \underbrace{[u^{(1)} \dots u^{(k)}]}_{n \times k}^T X = \underbrace{\begin{bmatrix} -u^{(1)T} \\ \vdots \\ -u^{(k)T} \end{bmatrix}}_{k \times n} X_{n \times 1}$$

U reduce

i.e. Z 是 k 维向量

主成分数量选择

12-5 Choosing the number of principal components (k)

- Algorithm: Try PCA with $k=1, 2, \dots$ (迭代直到满足条件)

Compute $\mathbf{U}_{\text{reduce}}, \mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(m)}, \mathbf{x}^{(1)}_{\text{approx}}(\text{投影}), \dots, \mathbf{x}^{(m)}_{\text{approx}}$

Check if

$$\frac{\sum_{i=1}^m \|\mathbf{x}^{(i)} - \mathbf{x}^{(i)}_{\text{approx}}\|^2}{\sum_{i=1}^m \|\mathbf{x}^{(i)}\|^2}$$

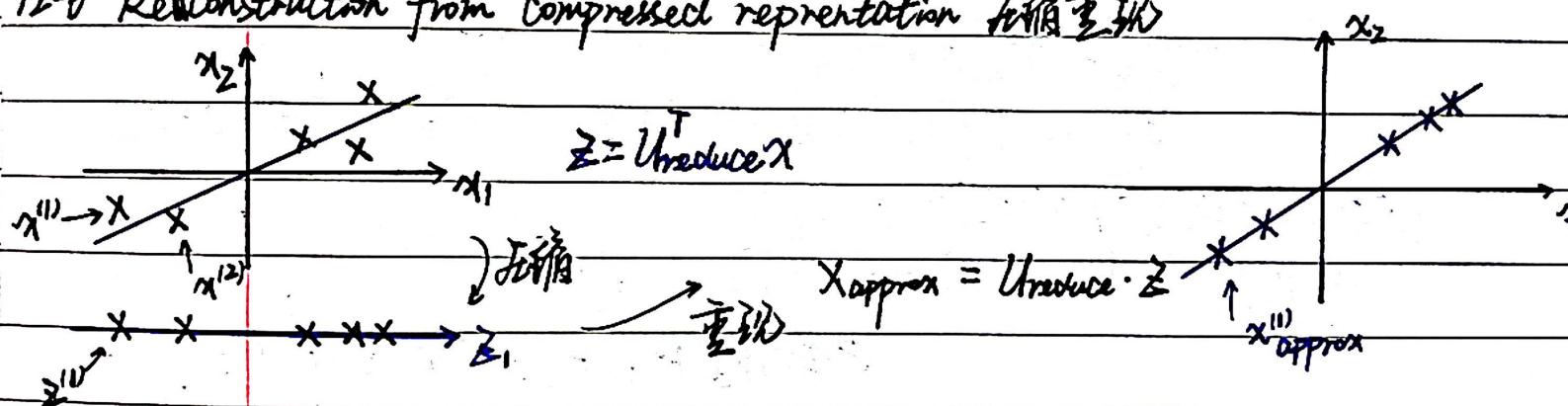
≤ 0.01 ($\rightarrow 99\%$ of variance is retained)

- $[\mathbf{U}, \mathbf{S}, \mathbf{V}] = \text{svd}(\mathbf{\Sigma})$ to get $\mathbf{S} = \begin{bmatrix} S_{11} & & \\ & S_{22} & 0 \\ & & \ddots & \end{bmatrix}$
(for given k) pick smallest of

value k for which

$$\frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \geq 0.99$$

12-6 Reconstruction from compressed representation (压缩重建)



12-7 Advice for applying PCA.

- Design
 - Get training set $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}) \dots, (x^{(m)}, y^{(m)})\}$
 - Run PCA to reduce $x^{(i)}$ in dimension to get $\mathbf{z}^{(i)}$
 - Train logistic regression on $\{(\mathbf{z}^{(1)}, y^{(1)}), \dots, (\mathbf{z}^{(m)}, y^{(m)})\}$
 - Test on test set: [Map $x^{(i)}_{\text{test}}$ to $\mathbf{z}^{(i)}_{\text{test}}$]. Run $h_{\theta}(z)$ on $\{(\mathbf{z}^{(1)}_{\text{test}}, y^{(1)}_{\text{test}}), \dots, (\mathbf{z}^{(m)}_{\text{test}}, y^{(m)}_{\text{test}})\}$ \rightarrow 通过映射关系利用同类
PCA在训练集中获得的参数
- Bad use of PCA: To prevent overfitting

[PCA正则化]