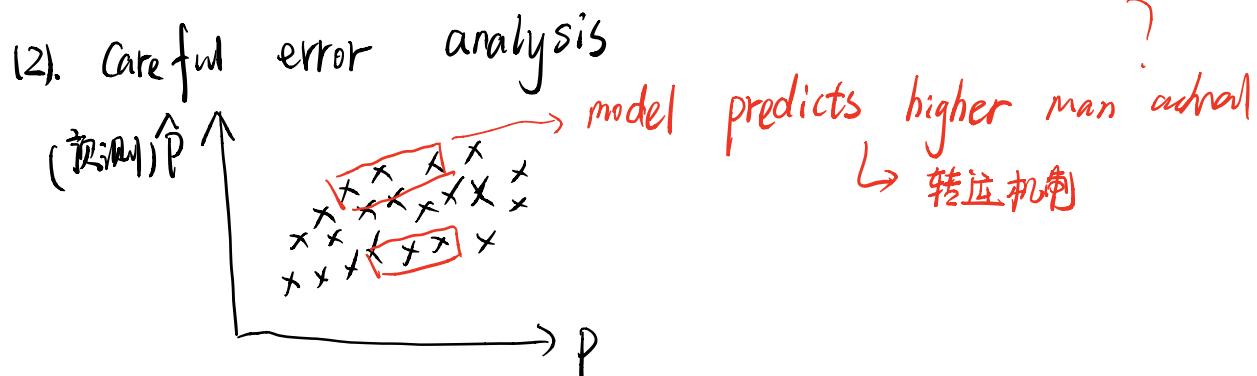
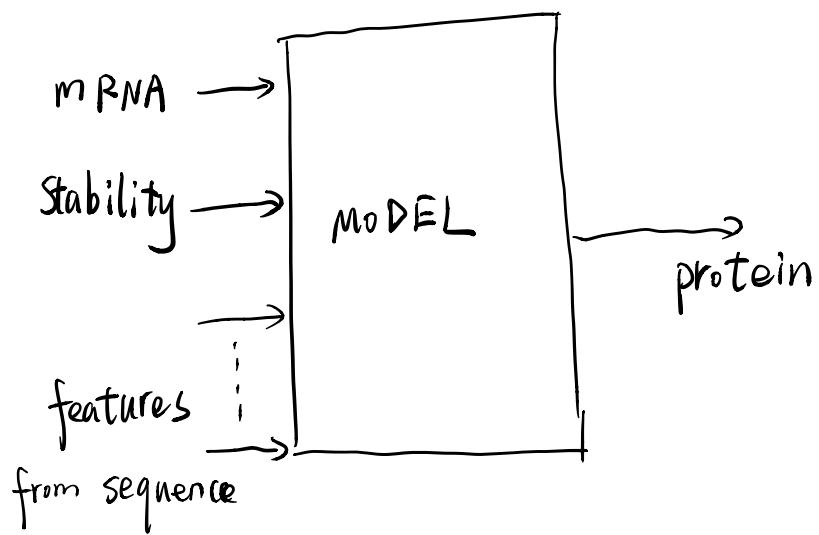


RIS HOR

PR ML

1. Genomics 基因组学



- 1) build model \rightarrow predict
measurements + derive features
 \Rightarrow regression model (回归模型)
- 2) Error analysis (Statement about problem)

2. Finance

time series analysis

预测: $\hat{s}(n) = s(n-1), s(n-2)$

	tomorrow	today	yesterday
residual (\Leftrightarrow) error	$r(n)$	$s(n) - \hat{s}(n)$	
剩余	true	predict	

unexplained information (\Leftrightarrow 例如 macro economics)

不可预测的原因 \Rightarrow 会导致 $r(n)$ 的存在

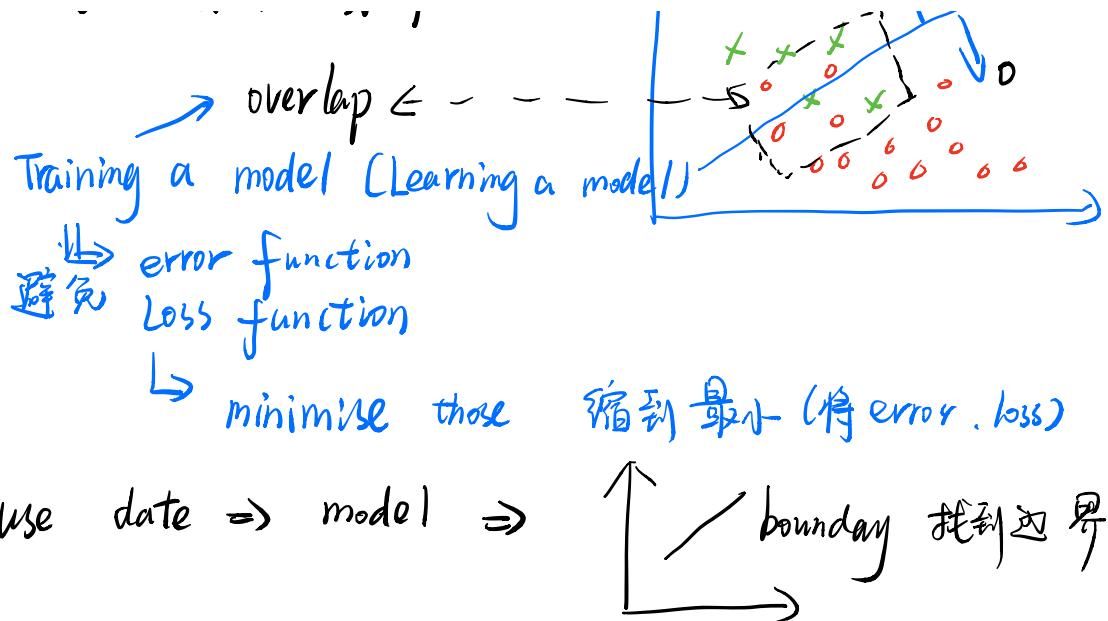
3. classification 分类

① Features proxies ml knowledge
代理

代替考试的形式

② Distribution 分布 $x \circ$





总结：

- ① Regression 回归
- ② Time Series 时间序列
- ③ classification 分类

Part 2. Chapter 1 of the book

Data \rightarrow Regulation/ties \rightarrow predicts

Representation Learning

$$\begin{array}{l} \text{Training} \\ \{x_n, t_n\}_{n=1}^N \end{array} \Rightarrow \begin{array}{l} \text{Test} \\ \hat{t}_{N+1} = f(x_{N+1}) \end{array}$$

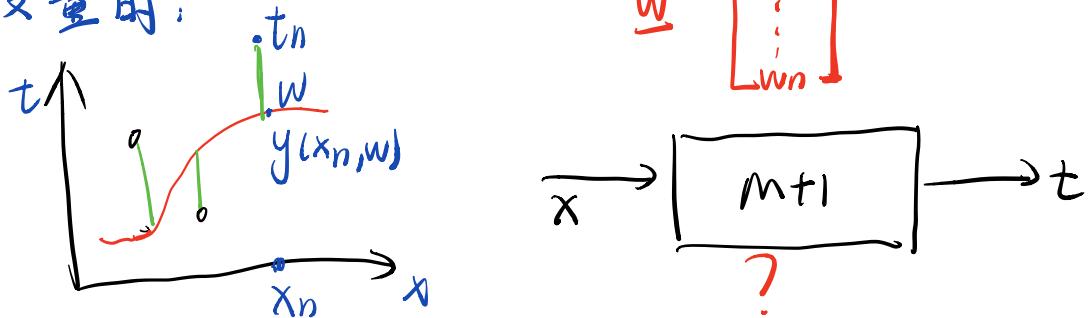
1.1 Polynomial Curve Fitting

Model: Polynomial

$$y(x, w) = w_0 + w_1 x + w_2 x^2 + \dots + w_m x^m$$

多个变量时的模型表达式.

单变量时:



$$\text{? } E(w) = \frac{1}{2} \sum_{n=1}^N \{ y(x_n, w) - t_n \}^2$$

$E(\approx \frac{1}{2} x^2)$? 误差

问题: over fitting

原因: 跨度太大 | 数据量小

解决: constraining the weight magnitude \downarrow hyper parameter
(约束) (量级) tune in (调整)

$$\hat{E}(w) = \frac{1}{2} \sum_{n=1}^N \{ y(x_n, w) - t_n \}^2 + \frac{\lambda}{2} \|w\|^2$$

$$w_0^2 + w_1^2 + \dots + w_m^2 \quad \text{和}$$

- 1. $\lambda = 0$, 不影响
- 2. λ 很大, $w_i^2 = 0$. 也不影响

1.2 Probability Theory

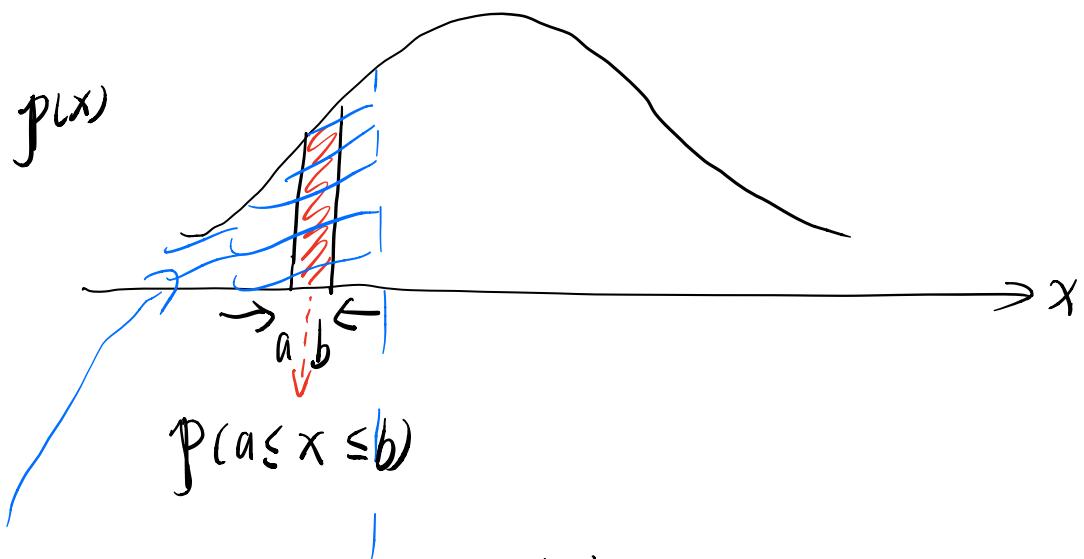
$$p(x) \quad p(y) \quad p(x,y)$$

$$p(x) = \sum_y p(x,y)$$

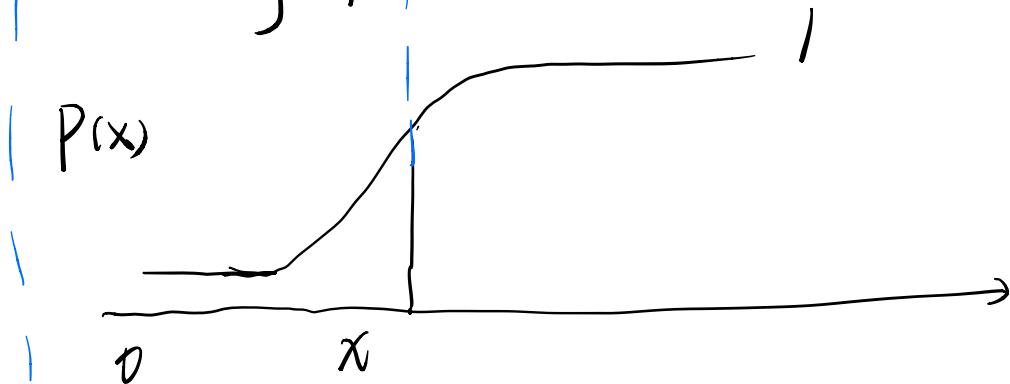
$$\begin{aligned} p(x,y) &= p(x|y) \cdot p(y) \\ &= p(y|x) \cdot p(x) \end{aligned}$$

$$\Rightarrow p(x|y) = \frac{p(y|x) \cdot p(x)}{p(y)}$$

概率密度



cumility function 45'26''



$$P[x] = P[X < x] = \int_{-\infty}^x P(z) dz$$

Expectation \bar{x} $f(x)$

$$E[f] = \sum_x p(x) \cdot f(x) \quad \rightarrow (1.33)$$

probability * integrate over the whole space
 $= \int p(x) f(x) dx \quad \rightarrow (1.34)$

协方差 Covariance

Frequentist

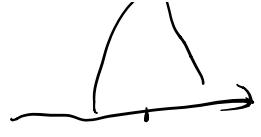
frequency of occurrence
时间点

Bayesian \Rightarrow quantifying

贝叶斯

量化

uncertainty

(估算) 

0,

likelihood of Date

$P(D|\underline{w}) P(\underline{w})$

1.43. $P(\underline{w} | D) = \frac{P(D|\underline{w}) P(\underline{w})}{P(D)}$

Posterior Bayes inference on \underline{w}

conditional of Data
(数据后验条件)

D : Data $\rightarrow \{x_n, t_n\}_{n=1}^N$

coefficients of the polynomial
系数 多项式

1.44. Posterior \propto likelihood \times prior

$$P(D) = \int P(D|\underline{w}) P(\underline{w}) d\underline{w}$$

- 高斯

§ 1: - Decisions

- Model selection
- Curse of dimensionality
维度灾难?
- Information Theory (Entropy)

火鸡

Part 3

独立同分布 (1.54)

IID : independently identically distributed?

$$\text{likelihood} = P(X | \mu, \sigma^2) = \prod_{n=1}^N N(x_n | \mu, \sigma^2)$$

关于参数的函数 (输入 X) 已知参数

$\frac{1}{\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$...

需要给出: $X = [x_1, x_2, \dots, x_n]$

$$\ln P(X | \mu, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln 2\pi$$

$$\left. \begin{aligned} \frac{\partial \ln P(X | \mu, \sigma^2)}{\partial \mu} &= 0 \\ \frac{\partial}{\partial \sigma^2} &= 0 \end{aligned} \right\}$$

作业

max likelihood ((μ, σ^2))

得到 \downarrow

$$\mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n \quad \text{平均值}$$

$$\sigma_{ML}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML})^2$$

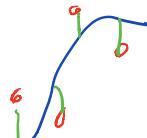
1.5b

自己学(作业)

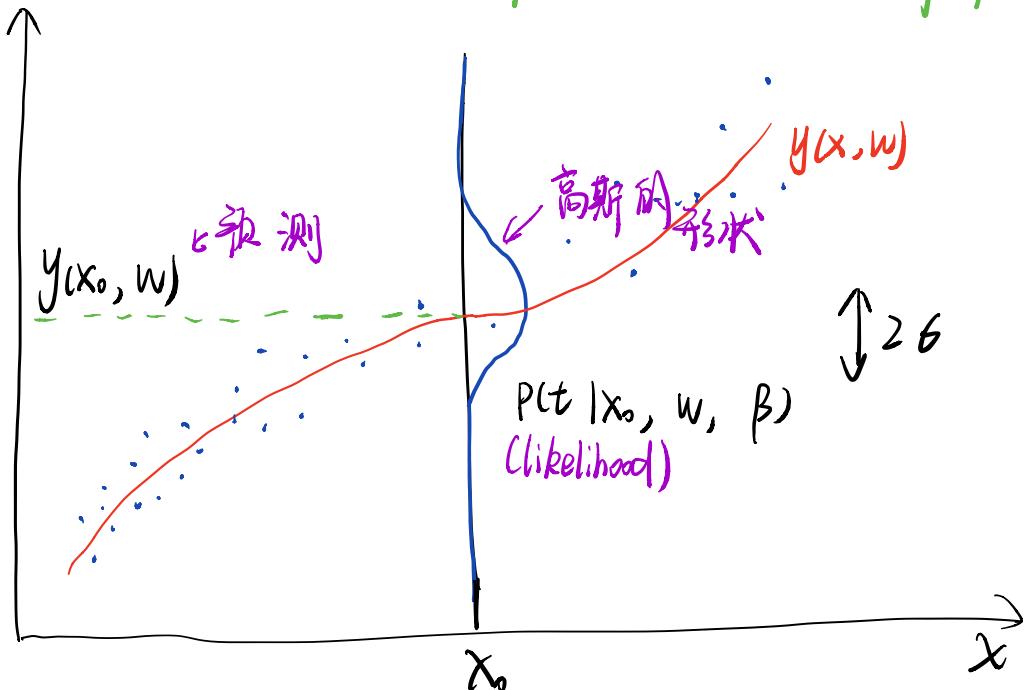
1.59

$\min \sum \text{error}^2$ 尽量减小误差

从条件 condition minor input



1.40 $p(t | \underline{x}, \underline{w}, \beta)$
 probability over targets
 noise term $\beta = \frac{1}{\sigma^2} \rightarrow$ 精确度参数
 precision parameter variance
 weights (coefficients of polynomial)



预测的 $y(x_0, w)$ 与实际有偏差，是因为有噪音的存在 (26)

$$P(t | \underline{x}, \underline{w}, \beta) = \prod_{n=1}^N N(t_n | y(x_n, \underline{w}), \beta^{-1})$$

原始高斯: $N(x | M, \sigma^2)$

$$\ln P(t | \underline{x}, \underline{w}, \beta) = -\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \underline{w}) - t_n\}^2 + \frac{N}{2} \ln \beta$$

$\approx \ln 2\pi$

可能对且不

JMLR 12 481~500

- 2 -

$$\frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{n=1}^N \{y(x_n, \underline{w}_{ML}) - t_n\}^2$$

贝叶斯公式：

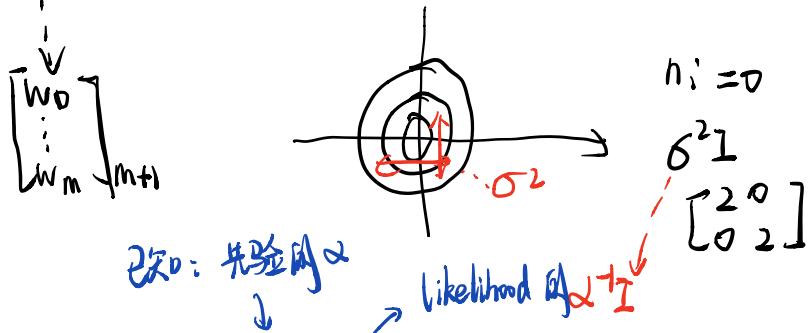
w prior distribution (先验分布)

$$\begin{aligned} & \{x_n, t_n\} \\ & y(x, w) \end{aligned} \quad \left\{ \begin{array}{l} \rightarrow \text{得到 likelihood} \end{array} \right.$$

计算 = posterior distribution

并且最大化 maximize posterior

$$p(w|\alpha) = N(w | \underline{\alpha}, \alpha^{-1} I) \quad \dots \dots (1.65)$$



改写：先验则 α

likelihood $\alpha^{-1} I$

$$P(w|x, t, \underline{\alpha}, \beta) \propto P(t|x, w, \beta) p(w|\alpha)$$

posterior

若要取最大值，则



$$\text{Max}_{\underline{w}}: \frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \underline{w}) - t_n\}^2 + \frac{\alpha}{2} \underline{w}^T \underline{w}$$

和 1.4 提到的公式 : $\sum_{n=1}^N \{ \text{data}_n - y \}^2 + \lambda \|w\|^2$ 一样

即: maximizing posterior
= minimizing regularized ?
square

Week 4

overfitting:

1. posterior \propto likelihood \times prior

$$p(w|x, t, \alpha, \beta) \propto p(t|x, w, \beta) \times p(w|\alpha)$$

$\uparrow_{\text{max}} \quad \hat{w}$

Bayesian

new point (input) $\leftarrow x$

predictive distribution on the targets $\rightarrow t$

$$p(t|x, \underline{x}, \underline{t}) = \int p(t|x, w) p(w|x, \underline{t}) dw$$

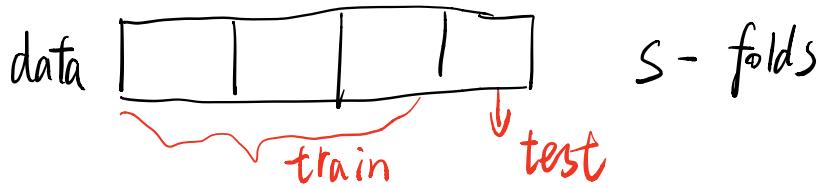
↑ 根據
訓練
數據
已知

test input ↗ posterior distribution
 計算出 w

2. 模型选择

问题 {
high over polynomial \rightarrow overfitting
low - - - - - \rightarrow not modelling data well

1) Cross validation



fit model m_2 m_3 .

2) Akaike information criterion AIC

$$\ln p(D | W_{m_L}) \sim M$$

↑
 best value of W
 I model (polynomial order)

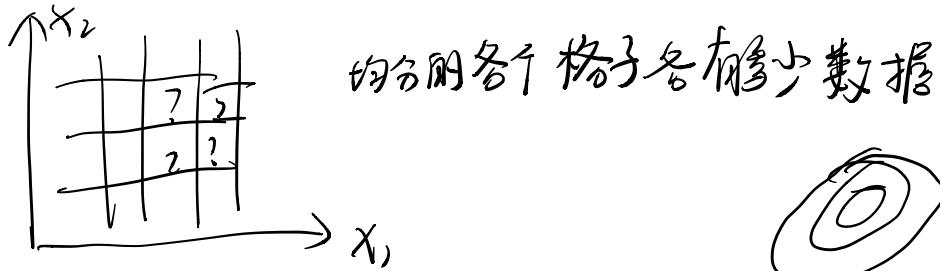
(重要概念)

3. Curse of dimensionality (多维)

① 一维:

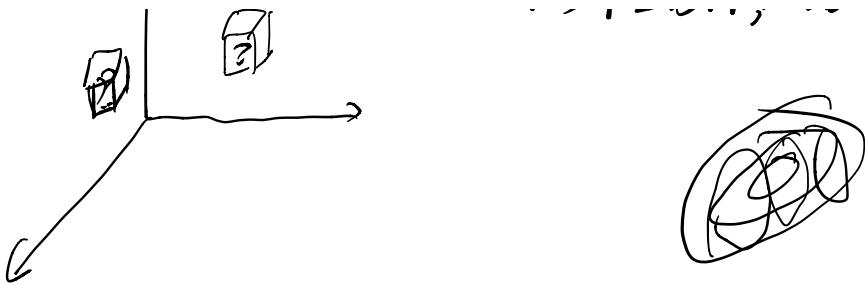


② 二维:



③ 三维:

↑
均分为多个立方体，看有多少个data

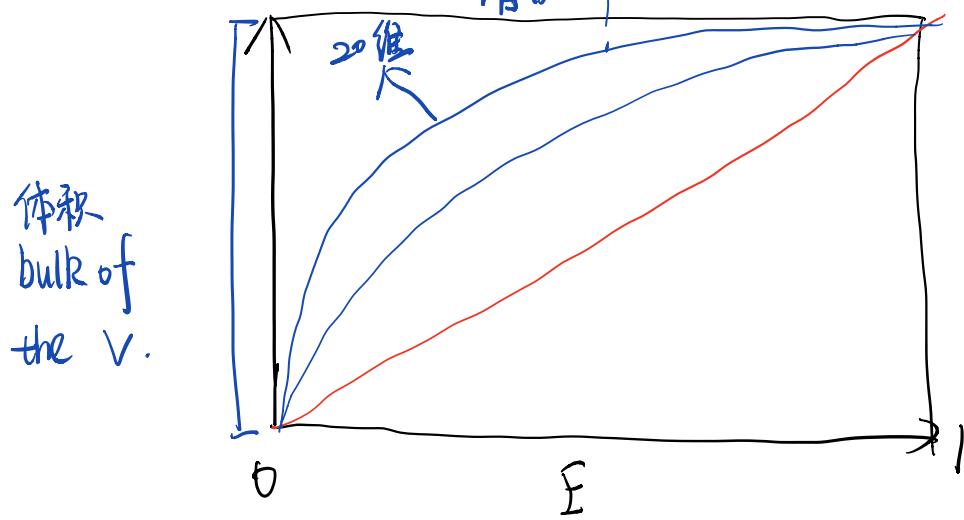


④ 多维: Volume contain in high dimensions

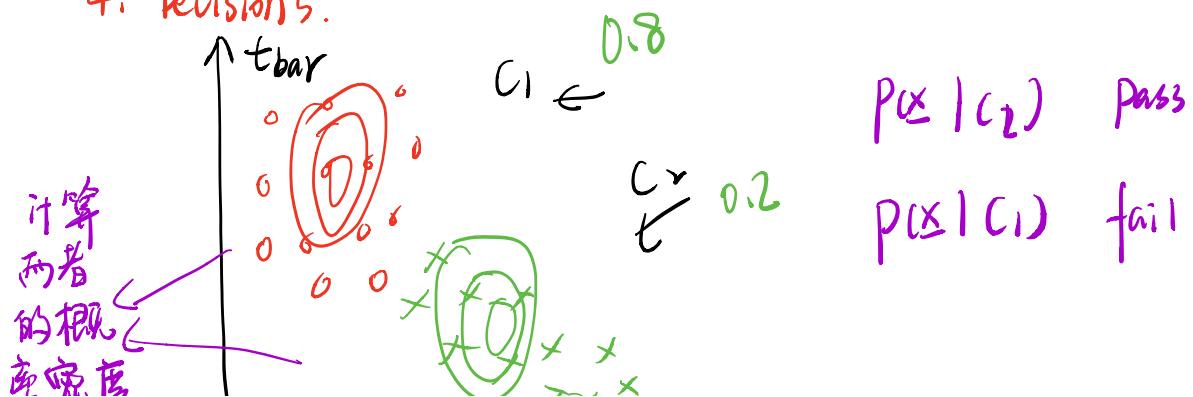
\downarrow
二维 \Rightarrow

计算阴影的 volume

随着维度的增加
 $r=1$
 $r=1-E$



4. Decisions.



t_{\max}

$$t_{\max} \xrightarrow{X \sim \cdot} t_{\text{lib}}$$

$$P(C_k | \underline{x}) = \frac{P(\underline{x} | C_k) \cdot P(C_k)}{P(\underline{x})} \rightarrow \text{prior}$$

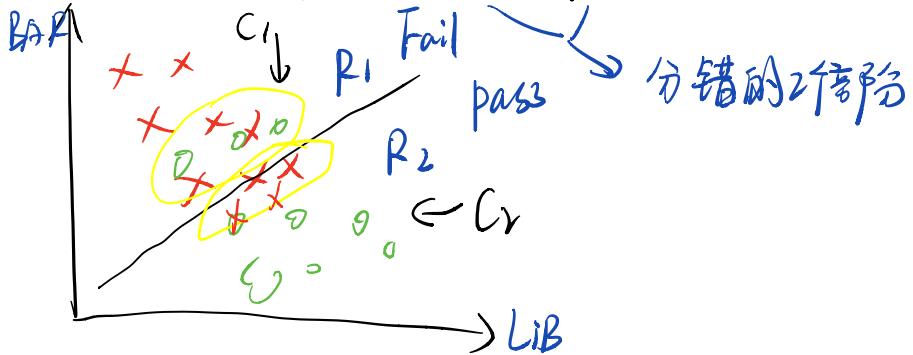
likelihood

所以，当给出 $\underline{x} = \begin{pmatrix} 3 \\ 0.1 \end{pmatrix}$ 可计算出 C_k (pass or fail)

$$\therefore P(C_k | \underline{x}) \geq 0$$

$$\sum_{k=1}^K P(C_k | \underline{x}) = 1$$

$$P(\text{error}) = P(\underline{x} \in R_1, C_2) + P(\underline{x} \in R_2, C_1)$$



$$P(\underline{x}, C_1) > P(\underline{x}, C_2) \Rightarrow \text{class } C_1 \vee$$

最小化 $P(\text{error})$

L_{kj} : how much it cost to misclassify
(Loss of k_j) class i as class k

期望：

$$\sum_k \sum_j \int L_{kj} P(x, c_k) dx$$

\uparrow
 \uparrow
mispredicted true class
as class j

如果不能计算，则可统计出错的次数

Reject

$$\text{if } P(c_1 | x) > P(c_2 | x)$$

在上述例子中，概率可分2种情况：

① 一种特别明显的 Pass 0.8 > fail 0.2

② 一种是边界处 Pass 0.501 ≈ fail 4.99.

2. 可 reject 第2种情况

if $P(c_1 | x) > P(c_2 | x) \rightarrow \text{class 1}$
else reject

week 5

Probability distribution

Bernoulli distribution \in binary

$$P(X=1) = \mu \quad (0 \leq \mu \leq 1)$$

$$P(X=0) = 1 - \mu$$

$$\text{Bern}(x|\mu) = \mu^x (1-\mu)^{1-x}$$

当 $x=1$ 时 $\Rightarrow \mu$

当 $x=0$ 时 $\Rightarrow (1-\mu)$

$$P(D|\mu) = \prod_{n=1}^N P(X_n|\mu)$$

Data

$$= \prod_{n=1}^N \mu^{x_n} (1-\mu)^{1-x_n}$$

$$\ln P(D|\mu) = \sum_{n=1}^N \ln p(X_n|\mu)$$

$$= \sum_{n=1}^N x_n \ln \mu + (1-x_n) \ln (1-\mu)$$

$$\text{最大化: 则 } \frac{\partial \ln p(D|\mu)}{\partial \mu} = 0$$

$$\begin{aligned}\text{得: } \mu_{ML} &= \frac{1}{N} \sum_{n=1}^N x_n \quad (x_n \in \mathbb{R}) \\ &= \frac{m}{N} \rightarrow \text{多个}\end{aligned}$$

$$2. \text{Beta}(\mu | a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1}$$

$$E[\mu] = \frac{a}{a+b}$$

$$\text{Var}[\mu] = \frac{ab}{(a+b)^2(a+b+1)}$$

$$p(\mu | m, l, a, b) = \frac{\Gamma(m+a+l+b)}{\Gamma(m+a)\Gamma(l+b)} \mu^{m+a-1} (1-\mu)^{l+b-1}$$

predictions:

$$p(x=1 | D) = \frac{m+a}{m+a+l+b}$$

若 data only (只有数据) $\Rightarrow \frac{m}{N}$

prior $\Rightarrow \frac{a}{a+b}$

posterior $\Rightarrow \frac{a+m}{(a+b)+N}$

3. likelihood 最大值.

$$\mu \in [0, 1]$$

Week 6 . chapter 4

备注: $\exp(-a)$ 表示 e^{-a}

4.1

\approx

$\rightarrow \pi \rightarrow$

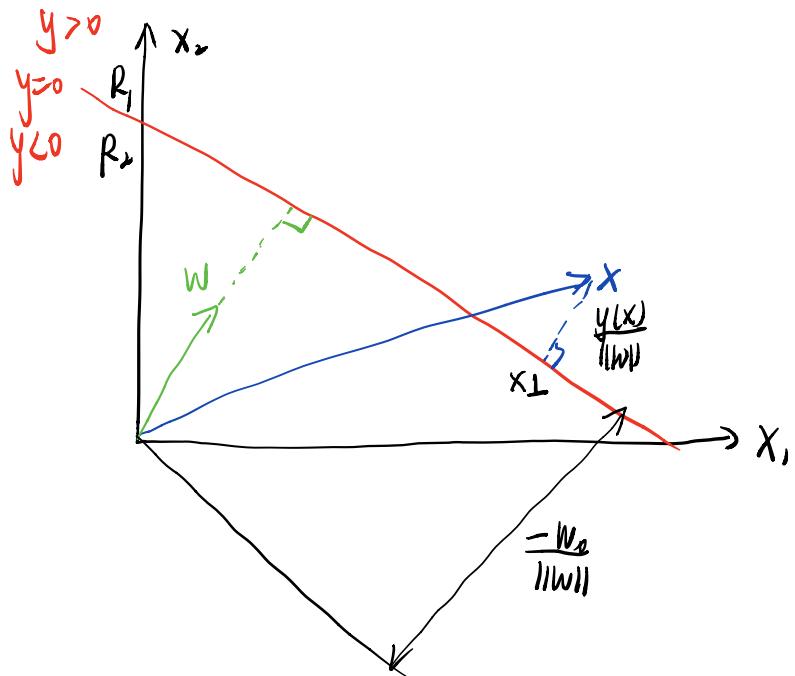
决策面
Decision surfaces are hyperplanes :

$$y(x) = w^T x + w_0$$

↑ weight ↓ 变量

① $y(x) \geq 0 \Rightarrow C_1$

②



有： $\frac{w^T x}{\|w\|} = -\frac{w_0}{\|w\|}$

③

设 k classes C_1, C_2, \dots, C_k

$(k-1)$ one ^{相对的} versus rest // $K \frac{(k-1)}{2}$ one versus
one classifiers.

若 $k=10$, 则可 build 9 ↑ classifiers,
KB 36 页

one versus

↓

one classifier.

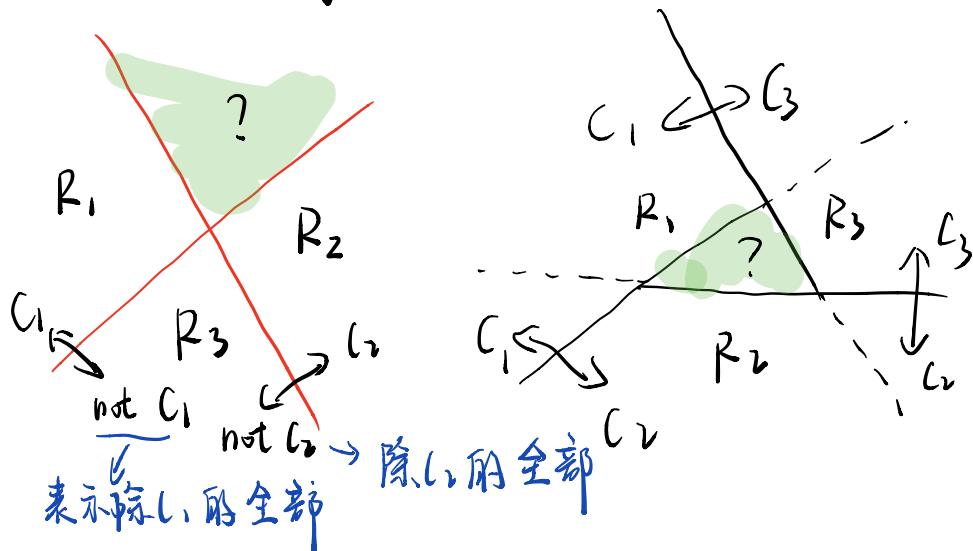
$C_2 \Rightarrow C_3$

$C_3 \Rightarrow C_4$

$C_4 \Rightarrow C_5$

↑

⊕ Ambiguous regions



定义了关于类的函数：

$$y_k(x) = w_k^T x + w_{k0}$$

指定 $x \rightarrow c_k$ if $y_k(x) > y_j(x), \forall j \neq k$

则有：Boundary between c_k and c_j :

$$(w_k - w_j)^T x + (w_{k0} - w_{j0})$$

⑤ 指定 $t = [0, 0, 1, 0, 0]$ $\xrightarrow{x^0}$

$$\tilde{w}_k = [w_{k0}, w_k^T]^T, \tilde{x} = [1, x^T]^T$$

输入 class c_k if $y_k = \tilde{w}_k^T \tilde{x}$ is largest

Vector of outputs as weight matrix times input.

在以上条件的情况下，

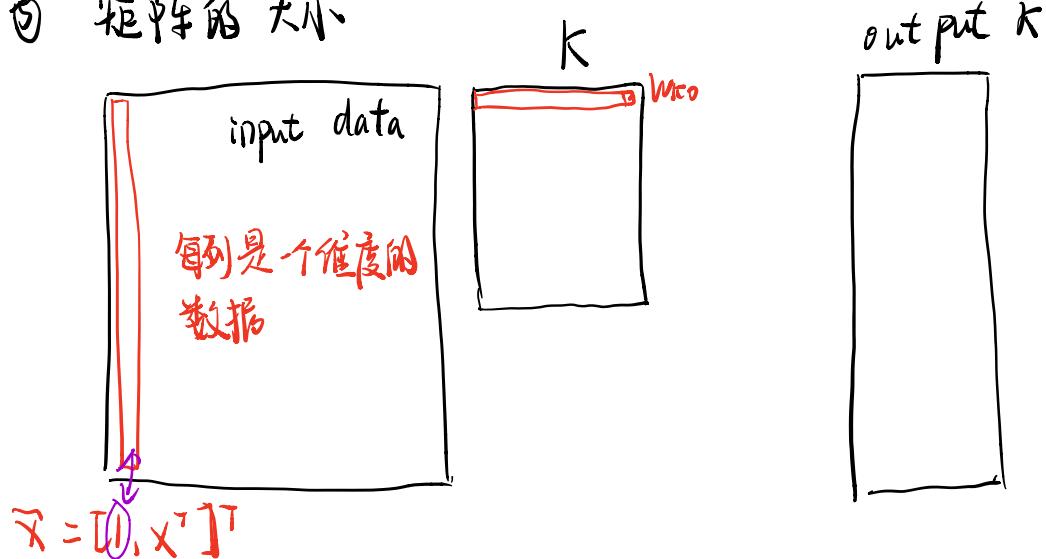
输入数据 Data $\{x_n, t_n\}$, $n = 1, 2, \dots, N$

最小化 error:

$$E_d(\tilde{w}) = \frac{1}{2} \text{Tr} \left\{ (\tilde{x}\tilde{w} - T)^T (\tilde{x}\tilde{w} - T) \right\}$$

↑
数据 × classes
↓
T的所有数据

⑤ 矩阵的大小



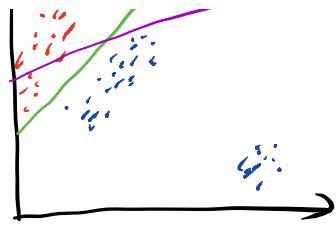
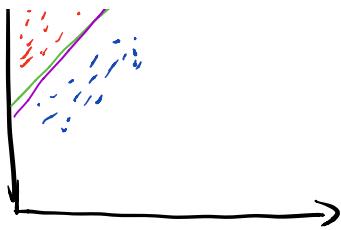
⑥ 结论：

$$\tilde{w} = (\tilde{x}^T \tilde{x})^{-1} \tilde{x}^T \cdot T$$

① 不能将分类器问题表述为插值(回归)问题：

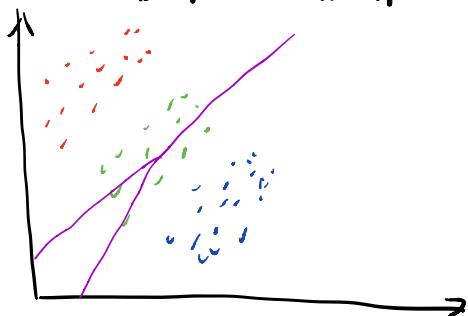
I. 原因 1

若有一小簇数据远离了，则两种方法结果不一

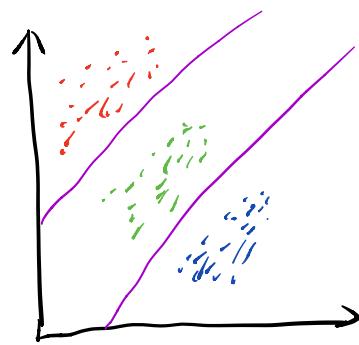


II. 原因 2.

若有三类数据



回归



分类