

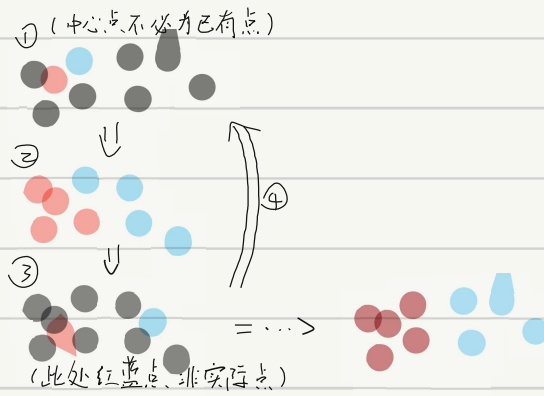
# 人工智能求解与实践 (注: 非期末小抄)

K-Means ① (随机) 选取  $K$  个中心点.

② 将每个点分配到最近的中心点.

③ 重新计算中心点.

④ 回到 ② 直到中心点收敛



## Spectral Clustering

无向图切图: for  $G(V, E)$ , 子图点集

$$\text{假设 } \bigcup A_i = V, A_i \cap A_j = \emptyset$$

$$\text{for set } A, B \subset V, \text{ def } W(A, B) = \sum_{i \in A, j \in B} w_{ij}$$

$$(Ncut) \text{ def } cut(A_1, A_2, \dots, A_k) = \sum_{i=1}^k \sum_{j=1}^k W(A_i, \bar{A}_j), \bar{A}_j = (V \setminus A_j) \text{ 补集}$$

$$\text{考虑平衡: } Ncut(A_1, A_2, \dots, A_k) = \frac{1}{2} \sum_{i=1}^k \frac{W(A_i, \bar{A}_i)}{vol(A_i)}$$

$$\text{其中 } vol(A) = \sum_{i \in A} d_i, d_i = \sum_{j=1}^n w_{ij} \text{ (与一点相连边权和)}$$

将  $vol(A)$  替换为  $|A|$ , 得到 Ratocut, 性能不如 Ncut

卢老师上课以  $N=2$  为例, 且边权都为 1,  $W(A_1, A_2) = |\{(i, j) \in E; i \in A_1, j \in A_2\}|$

故上述  $d_i$  为一点的边数

$vol(A)$  即为 2 倍  $A$  点间的边 + 经过  $A$  边界的边

(注, 虽然  $A$  只是点集, 但将  $A$  看作区域更为直观)

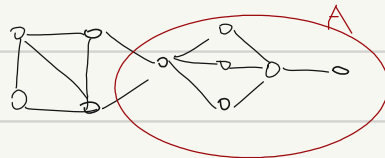
$$\text{代回 Ncut, } Ncut(A_1, A_2) = \frac{W(A_1, A_2)}{vol(A_1)} + \frac{W(A_1, A_2)}{vol(A_2)}$$

$$vol(A_1) + vol(A_2) = 2|E|,$$

于是卢老师实际给出的简化优化对象为

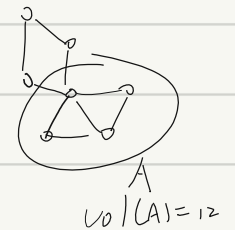
$$\text{argmin} \frac{\text{连结区域 } A, B \text{ 的边数}}{\text{较小区域的 } vol} = \frac{|\{(i, j) \in E, i \in A, j \notin A\}|}{\min(vol(A), 2m - vol(A))}$$

例



$$vol(A) = 16, 2m - vol(A) = 12$$

$$\text{结果 } \phi = \frac{6}{12} = \frac{1}{2}$$



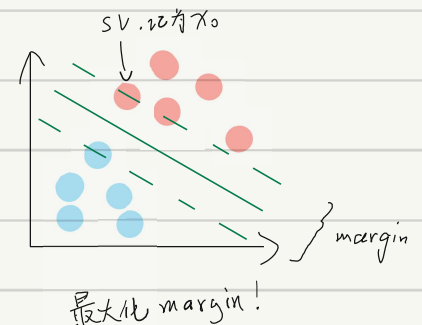
SVM 寻找  $M-1$  维超平面将  $n$  个  $M$  维样本分为 2 类

卢老师上课时仍先以 2 分类为例

推导大致过程: 设平面  $w_1 x_1 + w_2 x_2 + b = 0$

经过两个支持向量的平面为  $w_1 x_1 + w_2 x_2 + b = \pm c$

归一化: 三式均除以  $c$ , 直线不变, 参数变化



本笔记为复习期间回顾，独立撰写。如与上课不同，以上课为准

得表达式  $w_1'x_1 + w_2'x_2 + b' = 0, \pm 1$  三个平面

简化表达，以向量形式记为  $w^T x + b = 0$  (注，此处  $w = \begin{pmatrix} \frac{w_1}{c} \\ \frac{w_2}{c} \end{pmatrix}$ )

距离公式以  $\frac{w^T x + b}{\|w\|}$ ，于是分类  $y$  为  $\begin{cases} 1, & \frac{w^T x + b}{\|w\|} \geq d \\ -1, & \frac{w^T x + b}{\|w\|} \leq -d \end{cases}$ ， $d$  为  $\frac{1}{2}$  margin

注意到  $d = \frac{w^T x_0 + b}{\|w\|}$  未知，而  $\|w\|$  可调

$w^T x_0 + b = \|w\| d = \pm 1$ ，故  $\begin{cases} w^T x + b \geq 1, & y = 1 \\ w^T x + b \leq -1, & y = -1 \end{cases}$

即  $y(w^T x_0 + b) \geq 1$  confidence

SVM 欲使 margin 最大，而  $\text{margin} = 2d = 2 \frac{|w^T x + b|}{\|w\|} = 2 \frac{y(w^T x + b)}{\|w\|} = \frac{2}{\|w\|}$

故优化目标  $\arg\max \frac{1}{\|w\|}$  s.t.  $y(w^T x + b) \geq 1, \forall x \in D$

或  $\arg\min \|w\|$  或  $\arg\min \|w\|^2$

此处进入 SVM 精华部分，然卢老师止步于此 (即考试不考)

忆 (我考砸了的) 数分 Lagrange 乘子法

对优化问题  $\min f(x)$ , s.t.  $g_n(x) = 0, n = 1, 2, \dots, L$  (此  $f, g$  为向量值函数)

令  $L(x, \lambda) = f(x) + \sum_{n=1}^L \lambda_n g_n(x)$

可能极值点为  $\frac{\partial L}{\partial x_i}, \frac{\partial L}{\partial \lambda_i}$  的零点

于是 SVM 优化目标  $\min f(w) = \min \|w\|^2$  s.t.  $g_i(w) = 1 - y(x_i + b) \leq 0$

记  $L(w, \lambda) = \frac{1}{2} \|w\|^2 + \sum_{i=1}^l \lambda_i g_i(w)$

要求  $\min \frac{1}{2} \|w\|^2$ ，由  $g_i(w) \leq 0$ ，即求  $\min_w \max_{\lambda} L(w, \lambda)$  s.t.  $\lambda_i \geq 0$

由强对偶性，化为  $\max_{\lambda} \min_w L(w, \lambda)$

偏导可得  $w = \sum_{i=1}^l \lambda_i x_i y_i$ ,  $\sum_{i=1}^l \lambda_i y_i = 0$

代回，SMD 求解，此处略去

以上称为 Hard Margin SVM，不可解决 inseparable 问题

引入 slack variable  $\xi_i$ ，目标改为  $\arg\min \frac{1}{2} \|w\|^2 + C \sum_{i=1}^l \xi_i$  s.t.  $g_i(w) - \xi_i \leq 0, \xi_i \geq 0$   
 $C$  为错误分类惩罚因子， $\xi_i$  为错误程度

要求  $\max_{\lambda, \mu} \min_{w, b, \xi} L(w, b, \xi, \lambda, \mu)$ ，求法同上

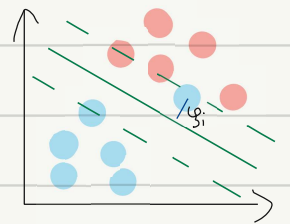
当维度相当低，不足以寻找平面 (如  $\begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$ ) 时，

升维或利用核函数转化 kernel trick，代  $x = \varphi(x)$

因升维并不能加秩，故用核函数可简化大量计算

如升维  $x' = (x_1^2, \dots, x_n^2, \dots, \sqrt{2}x_1, \dots, \sqrt{2}x_n, 1)$

相当  $\varphi(x) = (x \cdot x + 1)^2$



## Spectral Clustering

补充求解过程，为方便，权记为1

先得到邻接(adjacency)矩阵  $W$  和度矩阵  $D$ , Laplacian 矩阵  $L = D - W$

引入归一化向量矩阵,  $h_{ij} = \begin{cases} 0, & i \notin A_j \\ \frac{1}{\sqrt{\text{vol}(A_j)}}, & i \in A_j \end{cases}$

列向量  $h_i$ ,  $h_i^T L h_i = \frac{\text{cut}(A_i, \bar{A}_i)}{\text{vol}(A_i)}$ ,  $N\text{cut}(A_1, A_2, \dots, A_k) = \text{tr}(H^T L H)$

卢老师的二分类情况, 用了更为简单的表达

令  $x_i = \begin{cases} 1, & i \in A \\ -1, & i \notin A \end{cases}$ ,  $x$  为一些由  $\pm 1$  组成的向量  
但二元化使同一类所有点贡献相同, 不能很好区分中间点和边界点.

于是不妨  $x_i > 0$ ,  $i \in A$

优化  $\arg\min_x f(x) = \sum_{(i,j) \in E} (x_i - x_j)^2$  s.t.  $\sum x_i^2 = 0$ ,  $\sum x_i = 0$  (归一化条件)

化简  $f(x) = \sum_{(i,j) \in E} (x_i^2 + x_j^2 - 2x_i x_j)$ , 结合度的定义

$= \sum_{i=1}^n D_{ii} x_i^2 - \sum_{(i,j) \in E} 2x_i x_j$ , 再由  $W$  定义

$= \sum_{i=1}^n (D_{ii} x_i^2 - \sum_j W_{ij} x_i x_j) = \sum_{i=1}^n L_{ij} x_i x_j$

$= x^T L x$ , s.t.  $x^T x = 1$ ,  $\sum x_i = 0$

即优化  $\arg\min_x \frac{x^T L x}{x^T x}$  s.t.  $\sum x_i = 0$

之后用线性代知求解, 略

## Expectation - Maximization

发现时间不足, 开始简略

EM 在 K-Means 加入 covariance  $\sigma^2$

求解 Gaussian( $\mu, \sigma^2$ ), 流程与 K-Means 一致

补充知识点: 最大似然估计

## Neural Networks

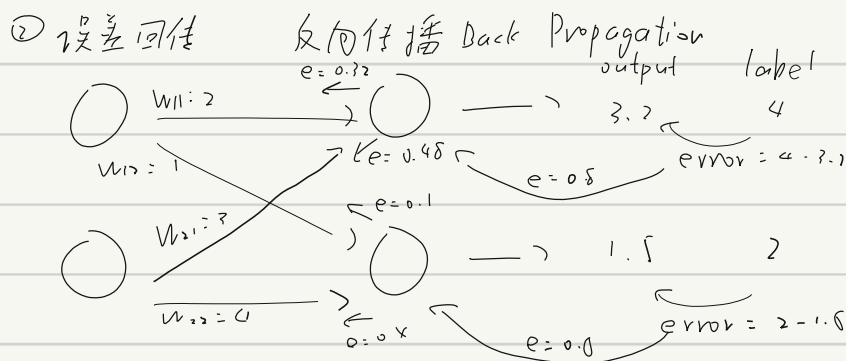
还不会的, 埋了吧

此处用 sigmoid  $y = \frac{1}{1+e^{-x}}$

FC 计算过程  $a^{(n+1)} = \sigma(W a^{(n)} + b)$  由上-层算下-层 前向传播

input 1  $\rightarrow$   $\frac{w_{11}}{0.9} \rightarrow 1.01 \xrightarrow{\sigma} 0.7408$  output

0.1  $\rightarrow$   $\frac{w_{21}}{0.8} \rightarrow 0.6 \xrightarrow{\sigma} 0.6457$



### ③ 更新权重

误差对权求导:  $\frac{\partial E}{\partial w_{ij}} = \frac{\partial (l_j - o_j)}{\partial w_{ij}} = -2(l_j - o_j) \frac{\partial o_j}{\partial w_{ij}}$

$$= -2(l_j - o_j) \cdot \frac{\partial \text{sigmoid}(\sum_i w_{ij} \cdot o_i)}{\partial w_{ij}}$$

又  $\sigma'(x) = \sigma(x)(1 - \sigma(x))$

于是  $\frac{\partial E}{\partial w_{ij}} = -2(l_j - o_j) \cdot \sigma(\sum_i w_{ij} \cdot o_i) (1 - \sigma(\sum_i w_{ij} \cdot o_i)) o_i$

$$\sim e_j \sigma(\sum_i w_{ij} \cdot o_i) (1 - \sigma(\sum_i w_{ij} \cdot o_i)) o_i$$

更新即  $\text{new } w = \text{old } w - \alpha \frac{\partial E}{\partial w}$  其中  $\alpha$  称为学习率

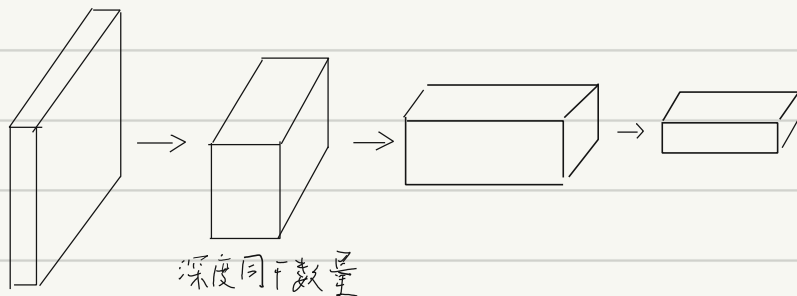
(卷积) CNN

CONV

对应块与Filter值相乘

再套ReLU

池化层: 降采样 (2x) 取 max



常用W初始化:  $W = \text{np.random.randn}(\text{in\_dim}, \text{out\_dim}) / \text{np.sqrt}(\text{in\_dim})$

应用: LeNet, AlexNet, VGG net,

DenseNet 比 ResNet 节约显存

Naive Bayes

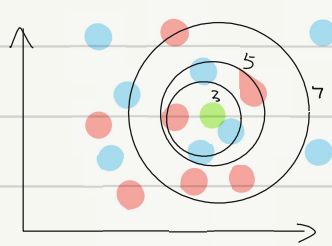
Theorem:  $P(B|A) = \frac{P(A|B)P(B)}{P(A)}$

So  $P(C|x_1, x_2, \dots, x_n) = \frac{P(C) \cdot P(x_1, x_2, \dots, x_n|C)}{P(x_1, x_2, \dots, x_n)}$

$\arg \max_c P(x_1, x_2, \dots, x_n|c) P(c) = \arg \max_c P(c) \prod_i P(x_i|c)$

if features are strongly independent

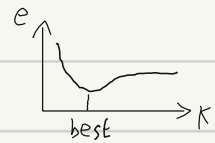
## K-Nearest-Neighbours



Green dot is class Blue for K 3.5

Red for K 7

Best K means min error.



Search with K-d Tree

## Decision Tree (C4.5)

$$\text{Entropy} = - \sum P(i) \log_2 P(i) \quad (\log_2 \text{ 或换为 } \log_{10})$$

设 data: label = 0 5个, label = 1 5个,

$$H(P) = - (\frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2}) = 1, \text{ 熵大, 不确定程度高. 信息量大}$$

$$1 \text{ 个 } 9 \text{ 个}, H(P) = - (0.1 \log_2 0.1 + 0.9 \log_2 0.9) \approx 0.47$$

Let  $\text{Info}(Y) = H(P)$ , Y is a result and A is an attribute

$$\text{Info}_A(Y) = \sum P(i) \cdot \text{Info}(Y)$$

$$\text{Gain}(A) = \text{Info}(Y) - \text{Info}_A(Y)$$

$$\text{SplitInfo}(A) = - \sum P(A) \log P(A)$$

$$\text{GainRatio}(A) = \frac{\text{Gain}(A)}{\text{SplitInfo}(A)}$$

PPT例) + 挺明白的

建树: ① max info gain attribute

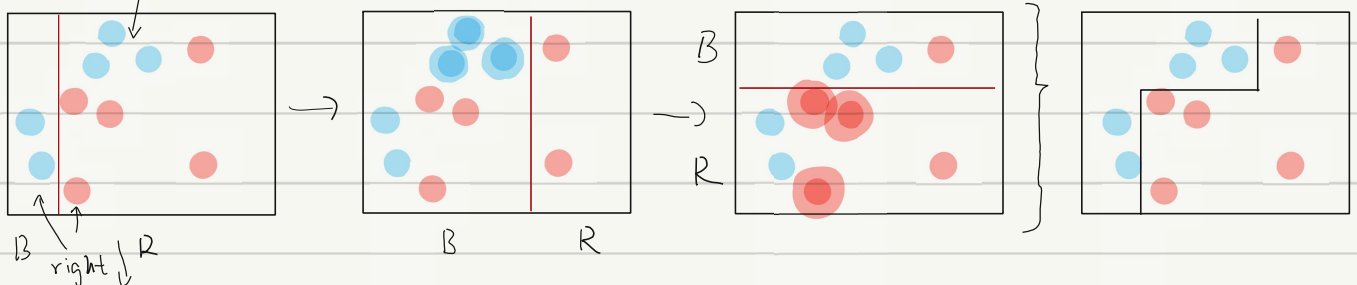
② delete selected new node

③ recalculate

(实际上, 先挑出 ratio 高于平均值的, 再算最高的)

Adaboost: 与 bagging (RF 所属算法) 并列, 增强 (boosting) 2-

思想: 相信弱分类器堆料可以成为强分类器 (AMD 战胜 Intel)

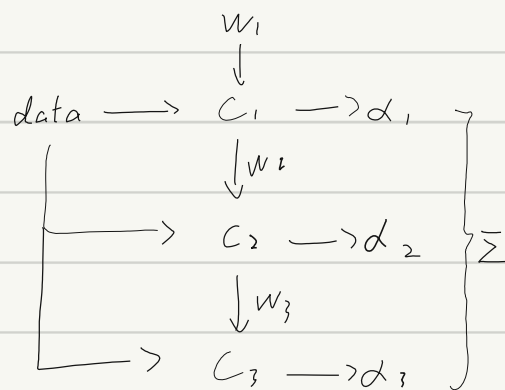


① initialize  $w_1 = \frac{1}{N}$ , equally contribute

② for range (Num\_of\_c)

fit c to all data with w

$$\text{error} = \frac{\sum_{i=1}^N w_i I(\text{wrong ans of } c)}{\sum_{i=1}^N w_i}$$



if error > 0.5 abort it

else compute  $\alpha = \log((1-\text{err})/\text{err})$

if wrong prediction for dot  $x_i$ ,  $w_i^* = e^{\alpha_m}$  ↑

if prediction correct for dot  $x_j$ ,  $w_j^* = e^{-\alpha_m}$  ↓

③ output =  $\sum \alpha(C(x))$

## PageRank

PR(A) 指, 每个人的平均分. A 得票数, L(A) 为 A 投给的人数

$$PR(A) = \sum \frac{PR(x)}{L(x)}$$

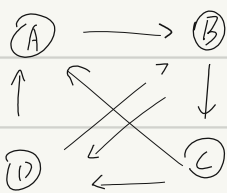
修正系数 d, 以免全变成 0 等情况

$$PR(A) = \left( \sum \frac{PR(x)}{L(x)} \right) \alpha + \frac{1-\alpha}{N}$$

化为矩阵形式  $P_{i+1} = \alpha S P_i + (1-\alpha) P_0$

S 为边权矩阵, 如例

$$S = \begin{pmatrix} 0 & 0 & \frac{1}{2} & \frac{1}{2} \\ 1 & 0 & 0 & \frac{1}{2} \\ 0 & \frac{1}{2} & 0 & \frac{1}{2} \\ 0 & \frac{1}{2} & \frac{1}{2} & 0 \end{pmatrix}$$



$$P_0 = \frac{\vec{e}}{N} = \begin{pmatrix} \frac{1}{4} \\ \frac{1}{4} \\ \frac{1}{4} \\ \frac{1}{4} \end{pmatrix}$$

S 列之和为 1,  $P_0$  和为 1

令  $A = \alpha S + (1-\alpha) \frac{e^T}{N}$  ( $P_i$  列之和也为 1)

于是  $P = \lim_{n \rightarrow \infty} A^n P_0$ , 由线代知识,  $P = A^n P$ , 解特征向量即可

## 基础知识 list

① 监督与非监督

② 分类与回归

③ Train: samples used to train      validation: calculate error (k-fold validation)

test: see final performance

④ distance