

① 英文试题, 中文答题, 填空题写过程

② 可带计算器

③ 范围: 作业 + 课件

第一章: 基础 (强化学习2考)

第三章: 公式 (条件概率, 方差, 距离)

第三章: 逻辑 (通过逻辑连接词计算真假, 量词, 范围不同导致真值不同.)

有监督学习

① KNN

② NB (习题, 计算方法)

③ 决策树 (记得1-2号, 条件熵的权重 取平均离散化, 习题2)

④ PLA (公式, 习题)

⑤ LR, NN, SVM (了解) 适用范围, 特点

无监督学习

① Apriori

② k-means (k值多少, 初始中心点多少)

① 联合概率 $P(A, B) = P(A)P(B|A) = P(B)P(A|B)$

$$P(A) = P(A, B) + P(A, \bar{B})$$

$$P(A) = \sum_{i=1}^n P(A, B_i) = \sum_{i=1}^n P(A|B_i)P(B_i)$$

② 期望: $E[X] = \sum_{i=1}^n x_i P(X=x_i)$ (离散型)

$$E[X] = \int_{-\infty}^{+\infty} x f(x) dx \quad (\text{连续型})$$

③ 性质: $E[aX + bY] = aE[X] + bE[Y]$ (不需要独立)

$$E[aX] = aE[X]$$

$$E[XY] = E[X] \cdot E[Y] \quad (\text{独立})$$

④ 方差: $\text{Var}(X) = E[(X - E[X])^2] = E[X^2] - (E[X])^2$

⑤ 协方差: $\text{Cov}(X, Y) = E[XY] - E[X]E[Y]$

⑥ 相关系数 (Correlation): $\text{corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}}$

⑦ 正态分布 (Normal Random Variables)

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

$$E[X] = \mu \quad \text{Var}(X) = \sigma^2$$

⑧ 距离:

Minkowski distance: $d(x, y) = \left(\sum_{k=1}^n |x_k - y_k|^r \right)^{1/r}$

$r=1 \rightarrow$ City block distance (曼哈顿距离)

$r=2 \rightarrow$ Euclidean distance

$r=\infty \rightarrow$ Supremum distance

⑨ Logic:

Propositional logic (命题逻辑) $\wedge \vee \Rightarrow \Leftrightarrow \neg$ (真假值一定)
First order predicate logic (一阶谓词逻辑) $\forall \exists, A(\dots), B(\dots)$

⑩ "⇒": 条件为真, 后边也为真才为真, 条件为假, 说什么都是对的.
 等价关系 $(P \Rightarrow Q) \Leftrightarrow (\neg P \vee Q)$

⑪ 一阶谓词逻辑:

object (个体词): a, b, ...

predicate (谓词): A(...), B(...), Z(...)

Quantifier (量词) $\left\{ \begin{array}{l} \text{Universal } \sim (\text{全称量词}) \forall \\ \text{Existential } \sim (\text{存在量词}) \exists \end{array} \right.$

Morgan's Law: $\begin{cases} \neg \forall x L \Leftrightarrow \exists x \neg L \\ \neg \exists x L \Leftrightarrow \forall x \neg L \end{cases}$

⑫ Quantifier Scope (量词作用域)

1. $\forall x (\forall y (x \Rightarrow F(y)))$ false \rightarrow 会飞的都是人, 比如鸟会飞, 但不是人

2. $\forall x F(x) \Leftrightarrow F(h)$ true \rightarrow 是人就会飞, 条件错误导致结果说错

$F: \dots$ can fly
 $h: \dots$ is human being

1. $\forall x (F(x) \Rightarrow F(h))$ 可以写成 $F(x) \Leftrightarrow F(h) \wedge F(x) \Leftrightarrow F(h) \wedge \dots F(x) \Leftrightarrow F(h)$
 只要其中一个错误即为错误

2. $\forall x F(x) \Rightarrow F(h)$ 因为当然所有x都会飞, 条件当然错误, 即该式为True, 如果为 \Leftrightarrow , 则则需要两边等价. 两边为True都为True

⑬ KNN, ... 算欧氏距离

⑭ NB:

求 max $P(C_i | X)$

$P(C_1 | X) = ?$

$P(C_2 | X) = ?$

$P(C_m | X) = ?$

$\left. \begin{array}{l} P(C_1 | X) = ? \\ P(C_2 | X) = ? \\ \vdots \\ P(C_m | X) = ? \end{array} \right\} \Rightarrow \text{out put } \hat{p}_{\max}$

$P(C_i | X) = \frac{P(X, C_i)}{P(X)} \rightarrow$ 比较大小用比较分子

$\Rightarrow P(X, C_i) = P(X | C_i) P(C_i)$

\downarrow
 条件概率
 \downarrow
 统计可得

$= P(X | C_1) \dots P(X | C_i)$

\downarrow
 考试用 C_i 对应的文本表示, 比如出现词为 1

⑮ PLA:

$h(X) = \text{sign}(w^T X)$

w 更新 $w_{t+1} \leftarrow w_t + y_{\text{new}} x_{\text{new}}$

⑯ 决策树

1. ID3: 信息增益 $\max \leftarrow$ 熵 min

熵: $H(D) = -\sum_{i=1}^n P(x_i) \log_2 P(x_i)$
 (Entropy)

条件熵 (conditional Entropy) $H(Y|D) = \sum_x P(D=x) H(Y|D=x)$
 联合熵 $H(X, Y) = -\sum_{x,y} p(x,y) \log_2 p(x,y)$
 $= H(X) + H(Y|X)$

互信息 (Mutual Information) $I(X; Y) = H(Y) - H(Y|X)$
 (等度量离散变量 X 和 Y 的相关度)
 $= H(X) - H(X|Y)$

$g(D, A) = H(D) - H(D|A)$

2. C4.5: 信息增益率 max

$GainRatio(D) = Gain_A(D) / SplitInfo(D)$
 其中 $SplitInfo(D) = -\sum_{j=1}^J \frac{|D_j|}{|D|} \log_2 \frac{|D_j|}{|D|}$

3. CART (Gain Index) 基尼指数 max

$gini(D) = \sum_{j=1}^J P_j(1-P_j) = 1 - \sum_{j=1}^J P_j^2$
 $gini_{split}(D) = \sum_{n \in D} \frac{1}{N} gini(D_n)$

4. 决策树剪枝: 惩罚因子, 人为规定.

Penalty term: $ec = \sum_{i=1}^L (r(n_i) + s(n_i))$
 计算错误率: $ec = \sum_{i=1}^L m(n_i)$

分母: 样本数.

分子: 错误数 + 节点数乘惩罚因子

如果 ec 个则剪枝.

(17) K-means

1. 先给出质心 (prototype): 计算其他点到质心的距离, 划为 cluster 重新计算中心点, 取平均值.
 2. 先给出 cluster 计算质心, 然后与 1 相同
- } \rightarrow 直至收敛

(18) DBSCAN

1. 核心点 (core points): 满足 EPS (半径) 和 MinPts (簇数)
 2. 边界点 (border points): 图中有核心点.
 3. 噪声点 (noise points): 除上面之外点.
- 算法: 以每个点为圆心画圆. 然后以 1, 2, 3 判断.

(19) 关联规则挖掘 Apriori

1. 支持度 (support): $S(A \rightarrow B) = S(A, B) \Rightarrow$ 联合概率
 2. 置信度 (confidence): $C(A \rightarrow B) = P(A, B) / P(A) \Rightarrow$ 条件概率
- 算法:
- 候选集 (candidate itemset)
- 通过最小支持度

频繁项集 (frequent)

最大频繁项集 (Maximal): 叶子节点, 没有直接超集

闭 \sim (closed): 其直接超集没有和它支持度相同的.