

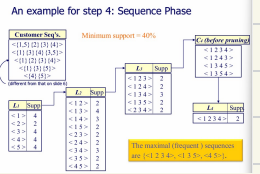
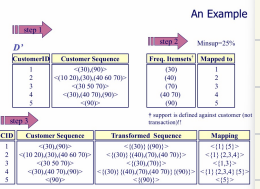
ARM1:
Rule measure support, confidence
for $X \Rightarrow Z$ $sup = P(X, Z)$ $con = P(Z|X)$
Apriori principle: 频繁项的子集也应该是频繁集
Apriori 算法: C_k : 大小为 k 的候选集 L_k 为频繁集
 $L_1 = \{frequent\ items\}$
for $k=1; k < k_{max}; k++$
 C_{k+1} : candidates from L_k
计算 C_{k+1} 支持度 (scan 数据库)
 $L_{k+1} = C_{k+1}$ with min. support
end
return $\cup_{k=1}^n L_k$

假设 L_{k-1} 是 sup 充分的频繁集
step 1: self-joining L_{k-1}
inserts into C_k
select p.item, ... p.item_n, q.item_1
from $L_{k-1} P L_{k-1} Q$
where p.item = q.item, ... p.item_n = q.item_n
step 2: pruning

for all itemsets c in C_k do
for all $(k-1)$ subsets s of c do
if s is not in L_{k-1} then delete c from C_k
Rule Generation: 对每个 frequent itemset L
生成所有非空子集, 对每个 L , 生成规则 $L \Rightarrow I$
如果 $con(L \Rightarrow I) \geq support(L) / support(I) > min-conf$
则 $L \Rightarrow I$ 是一个 strong association rule

The case of Apriori algorithm
Use frequent $(k-1)$ itemsets to generate C_k
Use database scanning and pattern match to collect counts for the candidate itemsets
Need $(k-1)$ scans, where n is the length of the longest pattern. (最大的频繁集)
Interest / life: $PLA(B)$ 1-1 关联
ARM2: sequence consists of a list of itemsets in temporal order and ordered c_1, s_2, \dots, s_n
sequential Association Rule Mining steps

1. sort phase: 转化为升序项集
2. frequent itemset phase: Find the set of all frequent large itemset L
3. transformation phase: 将 frequent 转化成 frequent item 的形式
4. sequence phase: 用 Apriori 发现频繁 sequence
5. 最大化: 找到所有不满足 sequence 的 sequence



Apriori All Algorithm

AprioriAll (Algorithm):
1. $L_1 = \{frequent\ 1\text{-itemsets}\}$
2. $k = 2$, "k" represents the pass number. * while $(L_{k-1} \neq \emptyset)$ do
 $P = \bigcup_{i=1}^k L_i$
 $C_k = \text{New candidates of size } k \text{ generated from } L_{k-1}$
 for each candidate-sequence $c \in C_k$ do
 increase the count of all candidates in C_k that are contained in c
 $L_k = \text{All candidates in } C_k \text{ with minimum support}$
 $k++$
return (P)

How to generate candidate sequences?

The candidate generation steps are similar to those of non-sequential association rule mining. However, be mindful of the order!

Step 1: self-joining L_{k-1}
insert into C_k
select p.itemset, p.itemset_n, q.itemset_1, q.itemset_n
from $L_{k-1} P L_{k-1} Q$
where p.itemset = q.itemset_1, ... p.itemset_n = q.itemset_n

Step 2: pruning
forall sequence c in C_k do
 forall $(k-1)$ -subsequence s of c do
 if s not in L_{k-1} then delete c from C_k

A sequence $\langle a_1, a_2, \dots, a_n \rangle$ contained in $\langle b_1, b_2, \dots, b_m \rangle$ if exist $i_1 < i_2 < \dots < i_n$
 $a_1 \leq b_{i_1}, \dots, a_n \leq b_{i_n} \Rightarrow$ sequence is maximal if no other s contains a
Sequence Association Rules: 对每个 frequent sequence into $S_1 \& S_2$ (不能破项)
如果 $conf(L_1 \Rightarrow S_2) = support(S_2) / support(S_1) > min-conf$, R 为 strong sequential association Rule

ARM2: 如果 rule 的 support 与另一个 rule 的 support 值相等, 则此 rule 是 redundant 的
在每个 item 旁加上大类名 (ancestor)

Algorithm Basic

- Form the extended transactions
Extended Transactions T
- Find frequent itemsets (minsup=30%, minconf=60%)
- Find the rules

Transaction	Item bought
200	Shirt (Clothes)
300	Jacket, Outerwear, Clothes, Hiking Boots, Footwear
400	Shirts, Footwear
500	Shoes, Footwear
600	Jacket, Outerwear, Clothes

Itemset	Support
{Jacket}	4
{Outerwear}	4
{Clothes}	4
{Hiking Boots}	2
{Footwear}	4
{Shoes, Hiking Boots}	2
{Clothes, Hiking Boots}	2
{Outerwear, Outerwear}	2
{Clothes, Footwear}	2

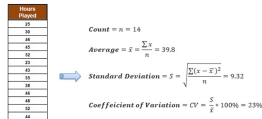
Rule	Support	Conf.
{Outerwear} \Rightarrow {Hiking Boots}	2/4	50%
{Outerwear} \Rightarrow {Footwear}	2/4	66%
{Hiking Boots} \Rightarrow {Outerwear}	2/2	100%
{Hiking Boots} \Rightarrow {Shoes}	2/2	100%

设置有效的 min support
 \Rightarrow uniform / reduced 支持与关联度一样 强化关联
Classification: seen Data
Training Data Test Data (模型训练的是 unseen)
对决策树来说, 每个特征都是 test, 似可作为 \dots then

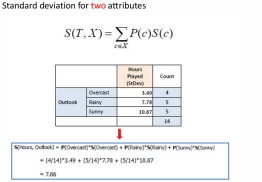
Information Gain (I) (分支)
 $I(p, n) = -\frac{p}{n} \log_2 \frac{p}{n} - \frac{n-p}{n} \log_2 \frac{n-p}{n}$
A 对 p 的信息增益: $E(A) = \sum \frac{p_i \cdot n_i}{n} I(p_i, n_i)$
A 对 n 的信息增益为 $Gain(A) = I(p, n) - E(A)$
Example: $p = 9, n = 5, I(p, n) = 0.99$
 $a_0 \quad p: n$
 $\leq 3 \quad 2 \quad 3 \quad 0.971$
 $3 > 4 \quad 0 \quad 0$
 $> 4 \quad 3 \quad 2 \quad 0.971$
 $E(A) = \frac{1}{5} I(2, 3) + \frac{1}{5} I(4, 0)$
 $\frac{1}{5} I(3, 2) = 0.694$

过拟合: ① 太多分支 ② 对未知数据效果不好
① 预剪枝 ② 后剪枝
大数据降下的分类任务, 为 million 级别的
hundreds 个特征下有 2 率运行速度

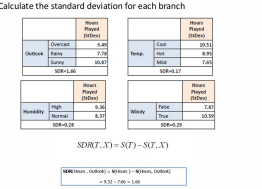
Classification: 回归树



Standard deviation (S) is for tree building (branching)
Coefficient of variation (CV) is used to determine branching termination.
Average (Avg) is the value in the leaf nodes



Standard Deviation Reduction (SDR)



Bayesian Classification: Why?

- Probabilistic learning: Calculate explicit probabilities for hypothesis, among the most practical approaches to certain types of learning problems
- Incremental: Each training example can incrementally increase/decrease the probability that a hypothesis is correct. Prior knowledge can be combined with observed data.
- Probabilistic prediction: Predict multiple hypotheses, weighted by their probabilities
- Standard: Even when Bayesian methods are computationally intractable, they can provide a standard of optimal decision making against which other methods can be measured

$maximize PchID) = \frac{PchID) Pch)}{P(cD)}$
 $\Rightarrow maximize PchID) Pch)$

$PCC_j | X) \propto P(C_j) \prod_{i=1}^n P(x_i | C_j)$
1-i 的类别, 用 C_j 的类别 $P(x_i | C_j)$
如 i 的连续 用高斯分布去算
eager learning: 提前构建假设, 必须新样本输入即可得出结论 (新样本不改变学习规律) (决策树, 神经网络, 朴素贝叶斯)
Lazy learning: 预测时才构建假设, 预测时用遍历史大量样本, 速度慢

Instance-Based Classifiers

- store the training records / case
- Use training records to predict of unseen case 找最相似的训练样本类别
- 不可解释性
- 用 NN 类优势: 能拟合高维, 输出结果最复杂, 训练时间大
- shallow: perception, MLP, SVM Deep: CNN, GNN 不能拟合高维知识
- RNN LSTM Transformer

停止标准
① 该特征下所有样本一致
② 没有其他的特征
③ 多数类才保留
④ 没有多余的特征了
选 SDR 最大的作为根节点

ification 11

classification II:

Measuring Error

True Class	Predicted class	
	Yes	No
Yes	TP: True Positive	FN: False Negative
No	FP: False Positive	TN: True Negative

- Error rate** = # of errors / # of instances = (FN+FP) / N
- Recall** = # of found positives / # of positives = TP / (TP+FN) = sensitivity = hit rate
- Precision** = # of found positives / # of found = TP / (TP+FP)
- Specificity** = TN / (TN+FP)
- False alarm rate** = FP / (FP+TN) = 1 - Specificity

Parcel ID	Origin	Destination	Type	Weight
1	HK	HK	Parcel	Light
2	Kln	Kln	Letter	Light
3	NT	Kln	Letter	Light
4	HK	HK	Parcel	Light
5	Kln	Kln	Parcel	Light
6	NT	NT	Letter	Light
7	HK	HK	Letter	Light
8	Kln	Kln	Parcel	Heavy
9	Kln	Kln	Letter	Light
10	HK	HK	Letter	Light
11	HK	HK	Parcel	Heavy
12	Kln	Kln	Letter	Light
13	HK	HK	Letter	Light
14	Kln	Kln	Parcel	Light
15	HK	NT	Parcel	Heavy
16	NT	HK	Letter	Light
17	HK	NT	Letter	Light
18	Kln	HK	Parcel	Light
19	HK	HK	Parcel	Heavy
20	HK	HK	Parcel	Light
21	Kln	Kln	Letter	Light
22	Kln	HK	Parcel	Heavy
23	Kln	Kln	Letter	Light
24	Kln	Kln	Letter	Light
25	HK	HK	Parcel	Light

- By setting the support to 25% and confidence to 50%, use the Apriori algorithm to find all interesting association rules in the above Table.
- Assume that a user sets the lift ratio to 1.75, which rules you discovered for a) above are still interesting?

Suggested Answer: (Please try it first!)

a) min_sup=25% (i.e., ≥ 7 records) and min_conf=50%.
The items in (black) with count ≥ 7 records below are frequent items.

1-itemset	Count	2-itemset	Count	3-itemset	Count
OHK	11	OHK, DHK	8	OHK, DHK, P	5
OKLN	11	OHK, DKLN	0	OHK, DHK, Li	6
ONT	3	OHK, P	7	OHK, P, Li	6
DHK	10	OHK, L	4		
DKLN	11	OHK, Li	7	OKLN, DKLN, Li	8
DNT	4	OKLN, DHK	2	DHK, P, Li	6
P	12	OKLN, DKLN	9	DKLN, Li, Li	8
L	13	OKLN, P	5		
Li	19	OKLN, Li	6		
H	6	OKLN, Li	9		
		DHK, P	7		
		DHK, L	3		
		DHK, Li	7		
		DKLN, P	3		
		DKLN, Li	8		
		DKLN, Li	10		
		P, Li	6		
		L, Li	13		

Estimate Error Rate

- partition: 折分 training a test set (大型数据集)
- 欠一折法: 把数据折成k个, 用k-1个, 单独k个(中小)

Ensemble Methods

多个模型一起把错的纠正率很小

The rules NOT highlighted in red are interesting ones w.r.t min_sup=25% and min_conf=50%. The same procedure can be repeated for frequent 2-itemsets but is omitted here.

b) Recall from lecture notes that

$$\text{Interest}(A \rightarrow B) = \frac{P(A \cap B)}{P(A)P(B)} \times \frac{1}{P(A \rightarrow B)} \times \frac{1}{P(B)} = \text{conf}(A \rightarrow B) \times \frac{1}{P(B)}$$

OKLN, DKLN, Li	DKLN, L, Li
OKLN, DKLN → Li	(8/9)(19/25)
OKLN, Li → DKLN	(8/9)(11/25)
DKLN, Li → OKLN	(8/10)(11/25)
OKLN → DKLN, Li	(8/11)(10/25)
DKLN → OKLN, Li	(8/11)(9/25)
Li → OKLN, DKLN	(8/19)(9/25)
DKLN, L → Li	(8/9)(19/25)
DKLN, Li → L	(8/10)(13/25)
L, Li → DKLN	(8/13)(11/25)
DKLN → L, Li	(8/13)(13/25)
L → L, DKLN	(8/13)(10/25)
Li → DKLN, L	(8/19)(9/25)

The rules above with interest > 1.75 are considered interesting.

Customer	Transaction Time	Items
David	2 Feb 2002	30,50
David	10 Feb 2002	50
David	21 Feb 2002	70
John	5 Feb 2002	10,30
John	7 Feb 2002	50
John	16 Feb 2002	40,60,70
John	27 Feb 2002	50,90
Peter	16 Feb 2002	30,50,70
Aaron	5 Feb 2002	30,30
Aaron	7 Feb 2002	30,50
Aaron	17 Feb 2002	70
Aaron	27 Feb 2002	50
Leon	13 Feb 2002	30

Build different "experts", and let them vote

- Advantages:**
 - Improve predictive performance
 - Easy to implement
 - No too much parameter tuning
- Disadvantages:**
 - The combined classifier is not so transparent (black box; low interpretability)
 - Not a compact representation
- Two approaches:** Bagging (Bootstrap Aggregating) and Boosting

Random Forests:

- Basic Bagging 对数据集进行树
- Feature Bagging 不同特征用不同树

Ada boost:

- 弱分类模型提高精度(样本权重再判)投票

- Given: Training data $(x_1, y_1), \dots, (x_N, y_N)$ with $y_i \in \{-1, +1\}$, $\forall n$
- Initialize weight of each example (x_n, y_n) : $D_1(n) = 1/N, \forall n$
- For round $t = 1, T$
 - Learn a weak $h_t(x) \rightarrow \{-1, +1\}$ using training data weighted as per D_t
 - Compute the **weighted** fraction of errors of h_t on this training data

$$\epsilon_t = \sum_n D_t(n) |h_t(x_n) \neq y_n|$$

Set "importance" of h_t : $\alpha_t = \frac{1}{2} \log \left(\frac{1-\epsilon_t}{\epsilon_t} \right)$ (gets larger as ϵ_t gets smaller)

- Update the weight of each example
 - $D_{t+1}(n) \propto \begin{cases} D_t(n) \times \exp(-\alpha_t) & \text{if } h_t(x_n) = y_n \\ D_t(n) \times \exp(\alpha_t) & \text{if } h_t(x_n) \neq y_n \end{cases}$ (correct prediction: decrease weight)
 - $D_{t+1}(n) \propto \begin{cases} D_t(n) \times \exp(\alpha_t) & \text{if } h_t(x_n) = y_n \\ D_t(n) \times \exp(-\alpha_t) & \text{if } h_t(x_n) \neq y_n \end{cases}$ (incorrect prediction: increase weight)

- Normalize D_{t+1} so that it sums to 1: $D_{t+1}(n) = \frac{D_t(n) \exp(\pm \alpha_t)}{\sum_{i=1}^N D_t(i) \exp(\pm \alpha_t)}$
- Output the "boosted" final hypothesis $H(x) = \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(x) \right)$

Table 1 Social network user popularity data

User ID	A1: Number of followers	A2: Number of following	A3: Privacy level	Popularity Score
B	Many	Low	High	75
B	Few	Many	High	72
C	Few	High	Low	30
D	Many	Many	Low	90
E	Few	Few	Low	37
F	Few	Many	High	42

Suppose you are asked to construct a regression tree based on the six records in Table 1.

- Given the following statistics about the popularity score, determine the root node of the regression tree to be constructed based on standard deviation reduction. Show the involved computation:

AVERAGE \bar{x}	6
STDEV s	66.00
Std.Deviation σ	$\frac{SD \times \sqrt{2}}{\sqrt{n}}$
Coef. of Variation CV	$\frac{s}{\bar{x}} \times 100\%$
	34%

The following pre-computed statistics may help to speed up your calculation of the SDR values:
 Average of (75, 90, 87)=84
 Std. Deviation of (75, 90, 87)=6.48
 Average of (30, 90, 42)=54
 Std. Deviation of (30, 90, 42)=37.66
 Average of (30, 90, 42)=54
 Std. Deviation of (30, 90, 42)=37.66
 Average of (75, 87)=81
 Std. Deviation of (75, 87)=6

(18 marks)

Ans. Standard Deviation for two attributes (target and predictor).

Consider A1

	SD	COUNT
A1	Many	6
	Few	17,66

S(TP,A1)	12.07
SDR(TP,A1)	10.31

SDR: Standard Deviation Reduction

Consider A2

	SD	COUNT
A2	Many	23,76
	Few	6

Consider A3

	SD	COUNT
A3	High	17,66
	Low	6,48

S(TP,A3)	12.07
SDR(TP,A3)	10.31

So, either A1 or A3 should be selected as the root node.

- Based on your regression tree with only the root node in part (a), predict the popularity score of users X and Y below.

User ID	A1: Number of followers	A2: Number of following	A3: Privacy level	Popularity Score
X	Many	Many	High	?
Y	Few	Few	Low	?

Ans. Based on A1 as the root node, the popularity score of X is 84
 Based on A1 as the root node, the popularity score of Y is 48

- WITHOUT guessing/inducing the missing values (denoted as "?"), show how the popularity score of the following data records should be predicted by your regression in part (a).

User ID	A1: Number of followers	A2: Number of following	A3: Privacy level	Popularity Score
W	Many	---	High	?
V	---	Many	Low	?

(6 marks)

Ans. Based on A1 as the root node, the popularity score of V is 84
 Based on A1 as the root node, the popularity score of W is undetermined. If we consider also A3 as a complementary solution, the popularity score of W is 48. Both solutions are acceptable.