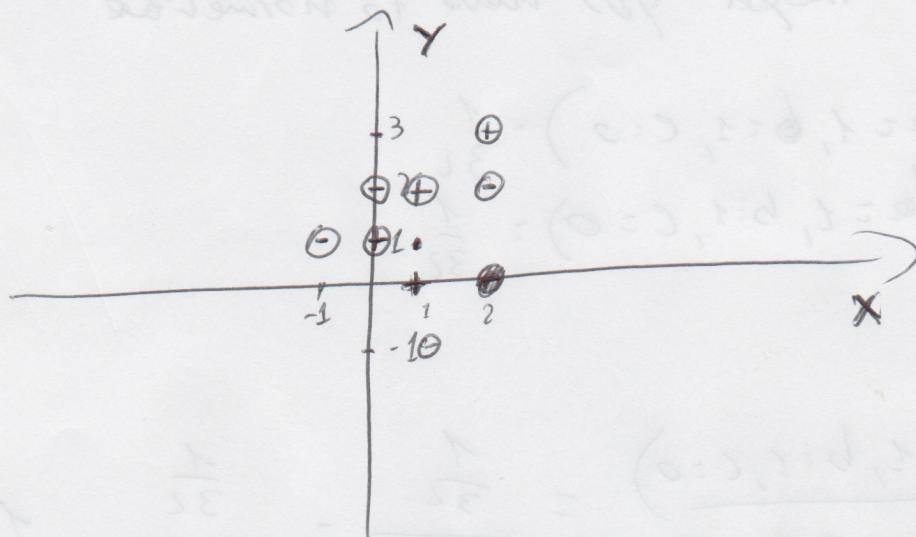


Problem 1



- a) $(1, 1)$ 3-nearest neighbor $\rightarrow +$ (3+)
 b) $(1, 1)$ 8-nearest neighbor $\rightarrow +$ (3+) (2-)
 c) $(1, 1)$ 7-nearest neighbor $\rightarrow -$ (3+) (4-)

Problem 2

perfect.

Naive Bayes classifier

$$\begin{aligned}
 P(K=1 | a=1, b=1, c=0) &= \frac{1}{32} && \text{premise} \\
 &= P(a=1 | K=1) P(b=1 | K=1) P(c=0 | K=1) \cdot P(K=1) \\
 &= \cancel{\frac{2}{4}} \cdot \cancel{\frac{1}{4}} \cdot \cancel{\frac{2}{2}} \cdot \cancel{\frac{4}{4}} \cdot \cancel{\frac{8}{2}} = \\
 &= \frac{2}{3} \cdot \frac{1}{3} \cdot \frac{2}{3} \cdot \frac{4}{8} = \frac{1}{32}
 \end{aligned}$$

$$\begin{aligned}
 P(K=1 | a=1, b=1, c=?) &= \\
 &= P(a=1 | K=1) P(b=1 | K=1) P(c=? | K=1) P(K=1) && \text{Unknown} \\
 &= \cancel{\frac{2}{3}} \cdot \cancel{\frac{1}{3}} \cdot \cancel{\frac{2}{3}} \cdot \cancel{\frac{4}{8}} = \left(\frac{1}{16}\right)
 \end{aligned}$$

Bayes yes, no and normalization.

For with bayes you have to normalize

$$P(K=1 | a=1, b=1, c=0) = \frac{1}{32}$$

$$P(K=0 | a=1, b=1, c=0) = \frac{1}{32}$$

normalize

$$\frac{P(K=1) P(K=1 | a=1, b=1, c=0)}{P(K=1) + P(K=0)} = \frac{\frac{1}{32}}{\frac{1}{32} + \frac{1}{32}} = \frac{\frac{1}{32}}{2 \cdot \frac{1}{32}} = \frac{1}{2} = 0,5\% \\ = 50\%$$

b)

$$P(K=1 | a=1, b=1) = P(a=1 | K=1) P(b=1 | K=1) = \frac{8}{3} \cdot \frac{1}{4} = \frac{2}{3}$$
$$P(K=0 | a=1, b=1) = P(a=1 | K=0) P(b=1 | K=0) = \frac{2}{3} \cdot \frac{1}{4} = \frac{1}{6}$$

normalize

$$P(\text{dark prob}) P(K=0) = \frac{\frac{1}{8}}{\frac{1}{8} + \frac{1}{16}} = \frac{\frac{1}{8}}{\frac{2+1}{16}} = \frac{1}{8} \cdot \frac{16}{3} = \frac{16}{24} = 0,66$$

$$P(K=1) = \frac{\frac{1}{16}}{\frac{1}{8} + \frac{1}{16}} = \frac{\frac{1}{16}}{\frac{2+1}{16}} = \frac{1}{16} \cdot \frac{16}{3} = \frac{1}{3} = 0,33 \\ \rightarrow 33\%$$

66

Association rules

Association rules are ~~rules~~

- * Given a set of transactions, an association rule is a rule that predict the occurrences of related items in a dataset.

Association rules are built by inspecting the dataset.

The goal of Association rules ^{may} is to find all rules having (check, page 8)
Support > minsup threshold
Confidence > conf threshold etc.

What is a decision tree? Is it true or false that decision tree mining can be applied to any type of data? how? if false, why?

A decision tree is a tree which is used to classify data. Decision trees have nodes, branches and leaves a node is a test on an attribute, a branch is an outcome of a test, Decision trees do not work very well with large datasets with a lot of attributes, In those cases the purity functions ~~does not work~~ do not work as expected.

Information gain: biased towards multivalued attributes

Gini Ratio: tends to prefer unbalanced splits (page 81, peck 1)

Question 5

Explain the difference between apriori and fp-growth

Apriori follows the generate-and-test paradigm while fp-growth follows the divide-and-conquer paradigm.

Apriori is more resource consumptive because it computes the minsup for each level of the tree while fp-growth calculates the minsup while building the tree itself.

Fp-growth wins in term of computation complexity against apriori algorithm.

Question 6 (Ensembles)

What is Bagging? Is there any relation between Bagging and Bootstrap? If yes, which one? If no, why?

Bagging is an ensemble method to aggregate classifiers. The training procedure is based on Bootstrap (sampled with replacement).

The Bagging method gives the same weight to each classifier and classify new data based on the majority vote of all the classifiers.

It's the easiest way to perform ensembles, the others ways are Boosting and Random Forests.

(slides 41, page 15)

Question 7

Your company has around $2 \cdot 10^6$ customers (too long to write all the question)

Decision trees should be avoided because they don't perform well on large datasets with a lot of attributes like in this case.

K-nearest-neighbor can't be used neither due to time constraints (low CPU/Ram power) because this method needs to scan all the dataset to give a prediction. So the best algorithm to chose is the Naive Bayes classifier which performs very well on a large amount of data and is lighter than K-nearest-neighbor.

2007 07 20

Problem 1

Consider the following training date set with one real attribute X and the class attribute Y .

Using information gain, draw the first node of the decision tree for this dataset.

| | s_1 | s_2 |
|-------|-------|-------|
| 0 0 0 | ++- | 22 |
| - + - | - - + | ++ |

$$\text{info}(p_1, p_2) = \text{entropy } p_1 \cdot \text{entropy}_{p_1}$$

$$s_1 \swarrow \quad \searrow s_2 \quad \text{info}(D) = -2 \cdot \frac{3}{8} \log_2 \frac{3}{8} = 1$$

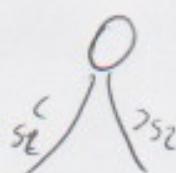
$$\text{info}_{CS_1} = -\frac{1}{3} \log_{10} \frac{1}{3} - \frac{2}{3} \log_{10} \frac{2}{3} = 0,276$$

$$\text{info}_{S_1} = -\frac{1}{5} \log_{10} \frac{1}{5} - \frac{2}{5} \log_{10} \frac{2}{5}$$

$$\text{change base } \frac{0,276}{\log_{10} 2} = 0,917$$

$$\frac{0,276}{\log_{10} 2} = 0,917$$

$$\text{info}_{S_1} = \frac{3}{8} \cdot 0,917 + \frac{5}{8} \cdot 0,917 = 0,85 \quad \text{info gain}_1 = 1 - 0,85 = 0,15$$



$$\text{info}_{CS_2} = -\frac{2}{6} \log_{10} \frac{2}{6} - \frac{4}{6} \log_{10} \frac{4}{6} =$$

$$\text{info}_{S_2} = -\frac{2}{8} \log_{10} \frac{2}{8} =$$

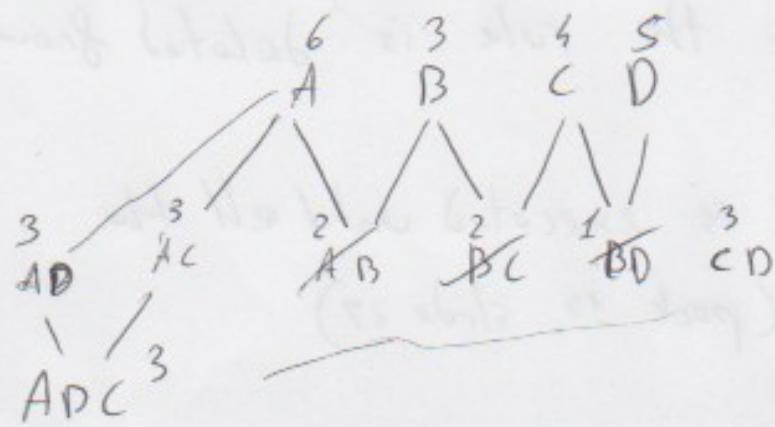
$$\text{info gain}_2 = \frac{6}{8} \cdot \text{info}_{CS_2} + \frac{2}{8} \cdot \text{info}_{S_2} = \quad \text{info gain} = 1 - \text{info}_{S_2} =$$

Problem 2 (Association rules)

Consider the following dataset. Extract all the frequent itemsets with a minimum support of 3 and a minimum confidence equal to 0.8

- Confidence is defined on rules, not on itemsets

Using the Apriori algorithm



Freq itemsets: A, B, C, D, AD, AC, CD, ADC

Problem 3 ~~not done~~

Illustrate the typical steps of a KDD process

KDD is a multi-step process involving data preparation, pattern searching, knowledge evaluation and refinement with iteration after modification.

The steps are: selection, preprocessing, transformation, data mining, evaluation and at the end what it's built is knowledge.

Problem 4

Briefly illustrate sequential covering algorithms

The sequential covering algorithm iterates over a dataset trying to find a rule which classifies part of the dataset, when a rule is found, all the itemsets identified by the rule is deleted from the dataset.

The rule creation stage is executed until all the dataset is covered (week 12, slide 27)

Problem 5

Error in a 1-nearest-neighbor

| | | error | |
|----|---|-------|--|
| 1 | - | 1 | |
| 2 | + | 0 | Cross validation error |
| 3 | + | 0 | = $\log_{10} \text{error} = \frac{1}{10} = 0,4 \quad 40\%$ |
| 4 | - | 1 | |
| 5 | + | 1 | Accuracy: $1 - 0,4 = 0,6$ |
| 6 | + | 0 | |
| 7 | - | 0 | |
| 8 | - | 0 | |
| 9 | + | 0 | |
| 10 | + | 0 | |

Question 5

What is data streaming mining? What are synopses?

Data stream mining is data mining applied to streams of data and not datasets. ~~is~~ Streams are not finite and usually contain unordered data. Data stream mining is an actual research field in data mining.

Synopses are tradeoff between storage needed resources and performance between while analyzing and researching on the streams data streams.

Question 6

Too long to copy

We have ~~as~~ a dataset without any information about it. so our goal is to understand what the data ~~does~~ is.

Consultant A is completely wrong. A decision tree in a big amount of data can introduce problems.

Consultant B suggests a not practical way to inspect the datasets coz decision trees are more efficient way to classify data but in this case we just want to understand what the data is.

Consultant C ^{says} is the best way to approach this problem. Hierarchical cluster will form initial groups of data, then K-means can show effectively the differences on each cluster.

20070803

Problem 2

Explain the similarities and the differences between naive Bayes classifiers and Bayesian networks.

Naive Bayes classifier uses probability to compute a class label and works on all the attributes.

It assumes that attributes are independent

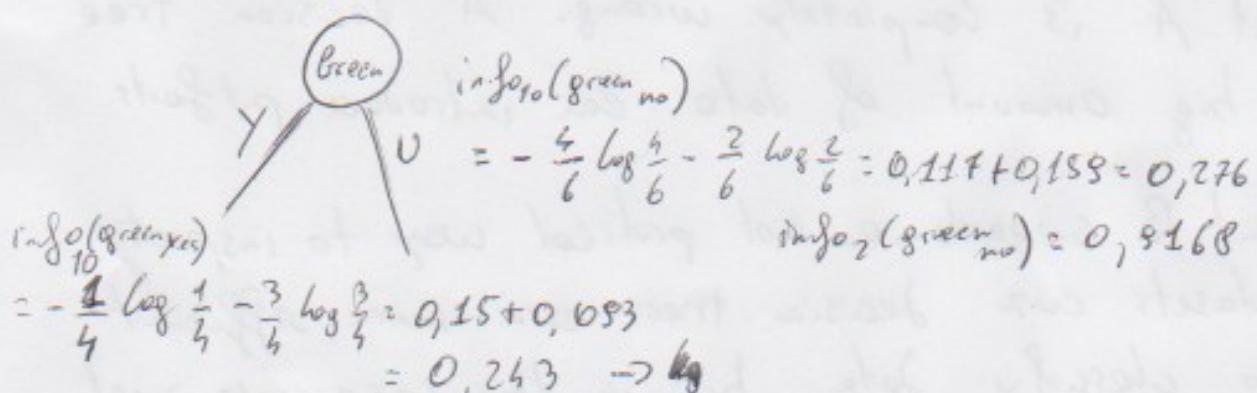
But it's often not so: class properties can be related.

Bayesian networks show in a quick and graphical way what are the related attributes of a class.

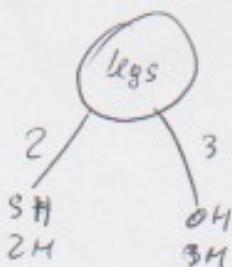
Bayesian networks can be used to improve performances of Bayes classifiers.

Problem 2

Decision tree



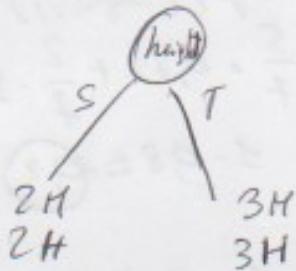
$$\text{info}_{\text{green}}(\text{green}) = \frac{4}{10} \cdot 0,81 + \frac{6}{10} \cdot 0,8168 = 0,324 + 0,55 = 0,87408$$
$$\text{info}_{\text{green}}(\text{green}) = 1 - 0,87408 = 0,12592$$



$$\text{info}_{10}(\text{legs}(2)) = -\frac{5}{7} \log \frac{5}{7} - \frac{2}{7} \log \frac{2}{2} = 0,104 + 0,1555 = 0,2593$$

$$\text{info}_2(\text{legs}(2)) = 0,8617 \quad \Rightarrow \text{info}_{10}(\text{legs}) = \frac{7}{10} \cdot 0,8617 = 0,60$$

$$\text{info}_{10}(\text{legs}(3)) = -\frac{3}{3} \log \frac{3}{3} = -1 \log 1 = 0 \quad \downarrow \quad \text{info}_{10}(\text{legs}) = 1 - 0,60 = 0,4$$



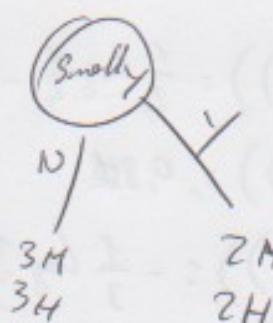
$$\text{info}_{10}(\text{height}(S)) = -\frac{2}{5} \log \frac{2}{5} - \frac{2}{5} \log \frac{2}{5} = -2 \log \frac{2}{5} = 0,30$$

$$\text{info}_{10}(\text{height}(T)) = -\frac{3}{6} \log \frac{3}{6} - \frac{3}{6} \log \frac{3}{6} = -2 \frac{3}{3} \log \frac{3}{6} = 0,30$$

$$\text{info}_2(\text{height}) = \frac{0,30}{\log 2} \approx 1$$

$$\text{info}(\text{height}) = \frac{8}{10} \cdot 0,1 + \frac{2}{10} \cdot 1 = \frac{8}{10} + \frac{2}{10} = 1$$

$$\text{info}_{10}(\text{height}) = 1 - 1 = 0$$



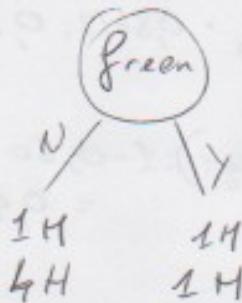
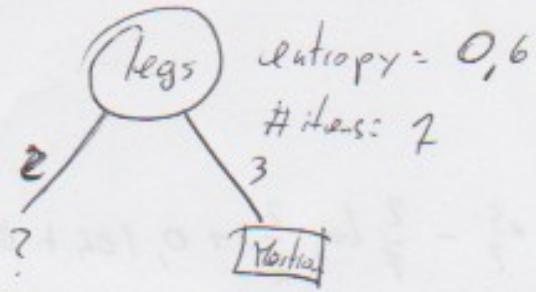
* Same as height

$$\text{info}_{10}(\text{smelly}) = 0,30$$

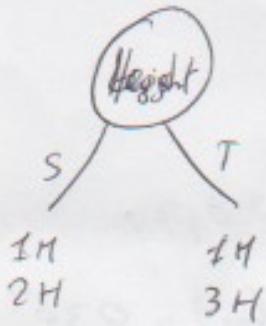
$$\text{info}_{10}(\text{smelly}) = \frac{5}{10} \cdot 1 + \frac{5}{10} \cdot 1 = 1$$

$$\text{info}_{10}(\text{smelly}) = 1 - 1 = 0$$

First split is legs.

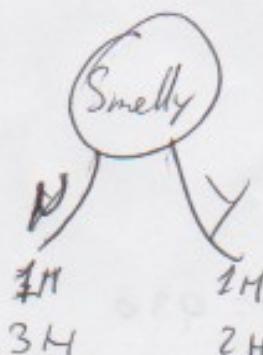


$$\text{info}_{10}(\text{legs}(n)) = -\frac{1}{5} \log \frac{1}{3} - \frac{4}{5} \log \frac{4}{5} = 0,130 + 0,0775 = 0,2075$$
$$\text{info}_2(\text{legs}(n)) = 0,72262$$
$$\text{info}_{10}(\text{legs}(y)) = 0,30 \Rightarrow \text{info}_2(\text{legs}(y)) = 1$$
$$\text{info}_{10}(\text{legs}) = \frac{5}{7} \cdot 0,72 + \frac{2}{7} \cdot 1 = 0,8$$
$$\text{info}_{10}(\text{green}) = \cancel{0,2075} - 1 - 0,8 = \textcircled{0,2}$$



$$\text{info}_{10}(\text{height}(S)) = -\frac{1}{3} \log \frac{1}{3} - \frac{2}{3} \log \frac{2}{3} = 0,159 + 0,117 = 0,276$$
$$\text{info}_2(\text{height}(S)) = 0,9181$$
$$\text{info}_{10}(\text{height}(T)) = -\frac{1}{3} \log \frac{1}{3} - \frac{3}{3} \log \frac{3}{3} = 0,15 + 0,0837 = 0,2337$$
$$\text{info}_2(\text{height}(T)) = 0,8085$$
$$\text{info}_{10}(\cancel{\text{height}}) = \frac{3}{7} 0,9181 + \frac{4}{7} 0,8085 = 0,3833 + 0,4628 = \textcircled{0,856}$$

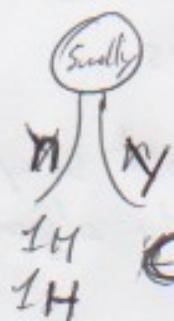
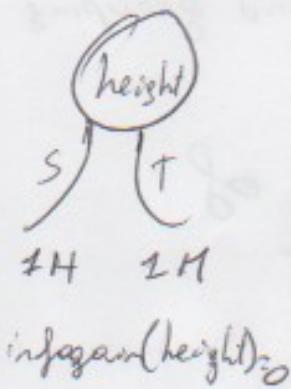
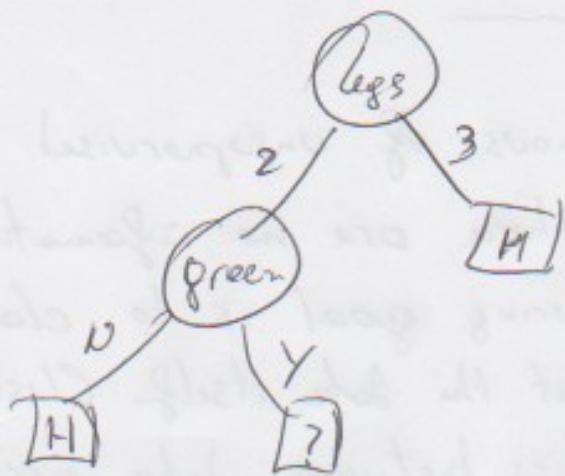
$$\text{info}_{10}(\text{height}) = 1 - 0,856 = \textcircled{0,144}$$



$$\text{info}_{10}(\text{smelly}(n)) = -\frac{1}{5} \log \frac{1}{4} - \frac{3}{5} \log \frac{3}{4} = 0,15 + 0,053$$
$$\text{info}_2(\text{smelly}(n)) = \cancel{0,857}$$
$$\text{info}_{10}(\text{smelly}(y)) = -\frac{1}{3} \log \frac{1}{3} - \frac{2}{3} \log \frac{2}{3} = 0,160 + 0,117 = 0,277$$
$$\text{info}_2(\text{smelly}(y)) = 0,82$$
$$\text{info}_6(\text{smelly}) = \frac{4}{7} \cdot 0,81 + \frac{3}{7} 0,82 = 0,4628 + 0,3342 = 0,857$$

$$\text{info}_{10}(\text{smelly}) = 1 - 0,857 = \textcircled{0,143}$$

Second split is green

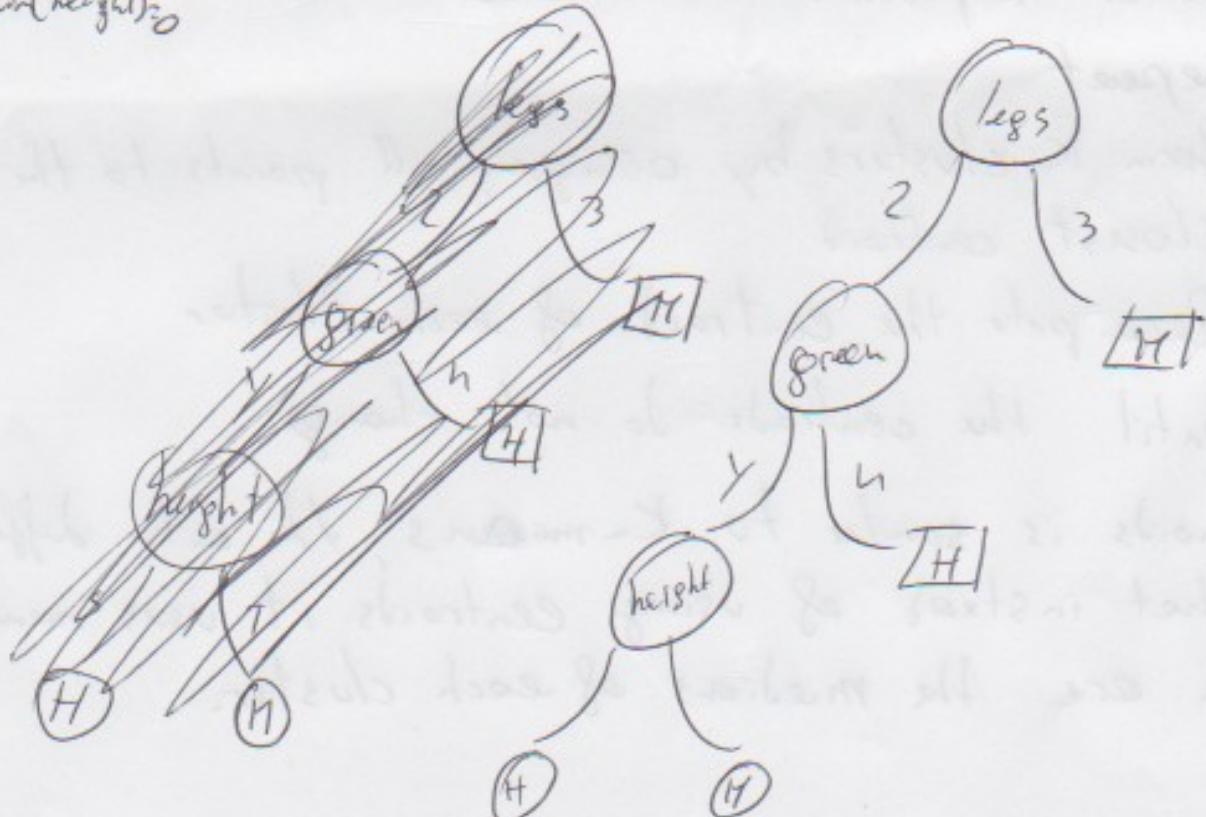


$$\text{info}_{\text{H}}(\text{sneaky}) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2}$$

$$\text{info}_{\text{H}}(\text{sneaky}) = 1$$

$$\text{info}_{\text{H}}(\text{sneaky}) = \frac{2}{2} \cdot 1 = 1$$

$$\text{info}_{\text{green}}(\text{sneaky}) = 0$$



Problem 3

Explain what is clustering —

Clustering is synonymous of unsupervised learning and it's applied when there are no informations about the data. Clustering goal is to classify similar data looking at the data itself. Clustering is done finding similarities between data according to the characteristics found in the data and grouping similar data objects into clusters.

The K-means algorithm is the most simple of the clustering algorithms. This is the algorithm:

Select k points in the dataset

Repeat

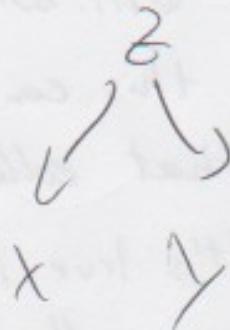
Form K clusters by assigning all points to the closest centroid

Recompute the centroid of each cluster until the centroids do not change

Kmedoids is similar to K-means, the only difference is that instead of using centroids it uses medoids which are the medians of each cluster.

Problem 4

e) $X \in \{0, 1\}$ $Y \in \{0, 1\}$



b) $10 \rightarrow P(X=1|Z=1), P(X=1|Z=0), P(X=0|Z=0)$
 $P(X=0|Z=1), P(Y=0|Z=1), P(Y=1|Z=1)$
 $P(Y=0|Z=0), P(Y=1|Z=0), P(Z=0)$
 $P(Z=1)$

Problem 5

Bootstrap, oversampling with replacement to form the training set and it's ~~not~~ suitable for small datasets coz it produces new items in the testset.

Cross-validation is sampling without replacement and in it leave one out version can reach an accuracy of 80%

With bootstrap you can achieve ~60% accuracy but the dataset must be uniformly distributed.

Problem 6

(1LY)

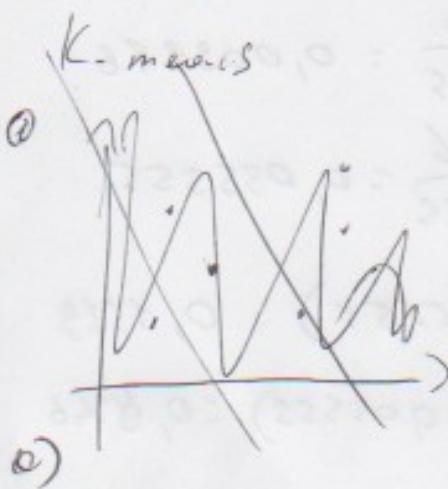
Overfitting is when our model of the class to recognize is full of useless attributes. ~~This also~~ ~~says that the class is~~ this can reduce the power of mining algorithms that collapses a simple search into a database. It's true the statement because in small datasets some attributes could look ~~selected~~ relevant when they aren't, thus this can lead to an overfitted model.

Problem 7

No, it's not a good offer. Since nominal attributes can have many values an automatic tool, if not well designed, can lead to overfitted rules. I would like to really test that software and see how it handles cases with many values for nominal attributes.

2007 09 17

Problem 1



~~Expectation points~~

(page 32, rock Bayes)

a)

$P(\text{play?} = \text{yes} | X = \{\text{sunny, hot, high, N}\})$, $P(\text{play?} = \text{yes} | X = \{\text{sunny, hot, high, W}\})$.

$P(\text{play} = \text{yes} | X = \{\text{sunny, hot, high, N}\})$

$$= P(X = \{\text{sunny, hot, high, N}\} | \text{play} = \text{yes}) P(\text{play} = \text{yes}) \\ = \frac{2}{9} \cdot \frac{2}{9} \cdot \frac{3}{9} \cdot \frac{6}{9} \cdot \frac{9}{15} = 0,00658$$

$P(\text{play} = \text{no} | X = \{\text{sunny, hot, high, N}\})$

$$= P(X = \{\text{sunny, hot, high, N}\} | \text{play} = \text{no}) = \frac{3}{6} \cdot \frac{2}{6} \cdot \frac{5}{6} \cdot \frac{2}{6} \cdot \frac{6}{15} \\ = 0,01851$$

Normalize

$$P(\text{yes}) = 0,00658 / (0,00658 + 0,01851) = 0,263 \Rightarrow 26\%$$

$$P(\text{no}) = 0,01851 / (0,00658 + 0,01851) = 0,736 \Rightarrow 74\%$$

So the result is No

b)

like before

$$P(\text{sunny, hot, high} \mid \text{play=yes}) = \frac{2}{8} \cdot \frac{2}{3} \cdot \frac{3}{9} \cdot \frac{9}{15} = 0,008876$$

$$P(\text{sunny, hot, high} \mid \text{play=no}) = \frac{3}{8} \cdot \frac{2}{6} \cdot \frac{5}{6} \cdot \frac{6}{15} = 0,055555$$

normalization

$$P(\text{yes}) = 0,008876 / (0,008876 + 0,055555) = 0,1519$$

$$P(\text{no}) = 0,055555 / (0,008876 + 0,055555) = 0,8486$$

$$P(\text{yes}) = 15\%$$

$$P(\text{no}) = 85\%$$

c) as seen in point a) $P(\text{no}) = 75\%$ so the classifier classifies $x = (\text{sunny, hot, high, N}) \rightarrow \text{play=no!}$

Problem 2

Plant

$$\text{info}_{10}(\text{Plant(usa)}) = -\frac{2}{4} \log \frac{2}{4} - \frac{2}{4} \log \frac{2}{4} = 0,30$$

ignore ~~info~~ plant

$$\text{info}_{10}(\text{Plant(cu)}) = 0 - 1 \log 1$$

$$\text{info}_{10}(\text{Plant(cin)}) = -1 \log 1$$

$$\text{info}_{10}(\text{Plant}) = 1 \rightarrow \text{ignore} = 0$$

Type

$$\text{info}_{10}(\text{Type}(A)) = \cancel{\frac{2}{5} \log \frac{2}{5}} - \frac{2}{5} \log \frac{2}{5} - \frac{3}{5} \log \frac{3}{5} \rightarrow \text{info}_{10}(\text{Type}(A)) =$$

$$\text{info}_{10}(\text{Type}(B)) = -\frac{5}{5} \log \frac{5}{5} + \cancel{\frac{2}{5} \log \frac{2}{5}} = 0$$

$$\text{info}(\text{Type}) = \frac{5}{10}.$$