

# Theoretical Homework 1

*For DATA MINING AND NEURAL NETWORKS*  
*MA4022*

Gleb Vorobchuk

Leicester

March 2018

gv@student.le.ac.uk

## Task N°1.

Give the Bayes formula for two events and for multivalued random variables.

The probability a woman in the general population has breast cancer is 0.003. The probability that a positive mammography result occurs given that a woman has cancer is 0.50. (This is the sensitivity of the test.) The probability that a negative result will occur given that a woman does not have breast cancer is 0.97. (This is the specificity of the test.) Suppose a woman has a mammography exam and gets a positive result. What is the probability that she actually has breast cancer?

Bayes Formula:

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

Given:

$$P(\text{Cancer}) = 0.003$$

$$P(\text{Positive} | \text{Cancer}) = 0.5$$

$$P(\text{Negative} | \text{Cancer}) = 0.97$$

$$P(\text{Positive} | \text{Non - Cancer}) = 1 - P(\text{Positive} | \text{Cancer}) = 1 - 0.97 = 0.03$$

$$\begin{aligned} P(\text{Cancer} | \text{Pos}) &= \frac{P(\text{Pos} | \text{Cancer}) \times P(\text{Cancer})}{P(\text{Pos} | \text{Cancer}) \times P(\text{Cancer}) + P(\text{Pos} | \text{Non - Cancer}) \times P(\text{Non - Cancer})} = \\ &= \frac{0.5 \times 0.003}{0.5 \times 0.003 + 0.03 \times 0.997} = 4.77 \% \end{aligned}$$

## Task №2.

Give a description of classification and clustering problems. What is the difference between them? Describe KNN approach and Hart's algorithm for data reduction. Describe the K-means algorithm and prove that it terminates after a finite number of steps. Construct an example of a dataset with the minimal number of data points for which the K-means algorithm (K=2) gives different clusters for different initial positions of centres.

Classification - Predictive task of identifying a set that new observation belongs to. Based on previous observations (training set). In Machine Learning defined as an one of the subtasks of supervised learning. Classification task approximate a mapping function  $f$  from inputs variables  $X$  to discrete output variables  $y$ . The output variables  $y$  are called labels or classes. Classification requires examples of training set to be classified into at least two categories.

Clustering - Descriptive task of grouping objects in way that objects in the same group are more similar. Clustering can be formulated as a multi-objective optimization problem. Require: scalability, ability to deal with many attributes, ability to deal with outliers, high-dimensionality, unlabelled data with no predefined classes. In Machine Learning defined as an one of the subtasks of unsupervised learning.

The difference between Classification and Clustering is that the training set for classification should contain predefined classes. Also Classification predict class for the new observation when Clustering just suggest groups based on patterns in data.

KNN approach.

KNN is the algorithm used for Classification and Regression. KNN makes prediction based on the outcome of the  $K$  closest neighbours to the examine point.

In order to do that we need to have a metric of distance between a neighbour and query point. The common metric is Euclidian distance -  $\sqrt{(x - p)^2}$  where  $x$  -query point,  $p$  - neighbour.

So after define the distance metric we can predict membership of the query point as averaged

sum of the KNN:  $y = \frac{1}{K} \sum_{i=1}^K y_i$

Hart's Algorithm:

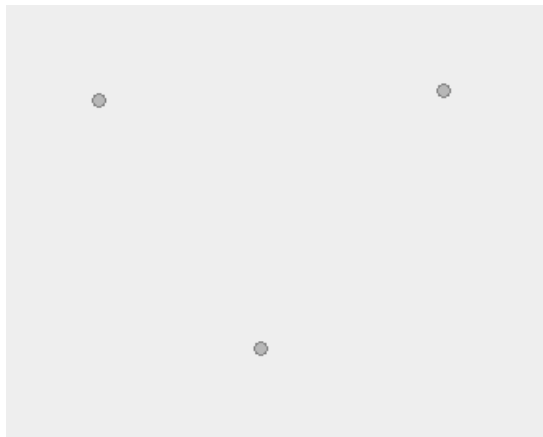
Also called Condensed nearest neighbour - is an algorithm based on KNN and was developed to reduce the data set for KNN. It selects the set of prototypes  $U$  from the training data, such that 1NN with  $U$  can classify the examples almost as accurately as 1NN does with the whole data set.

Given a training set  $X$ , CNN works iteratively:

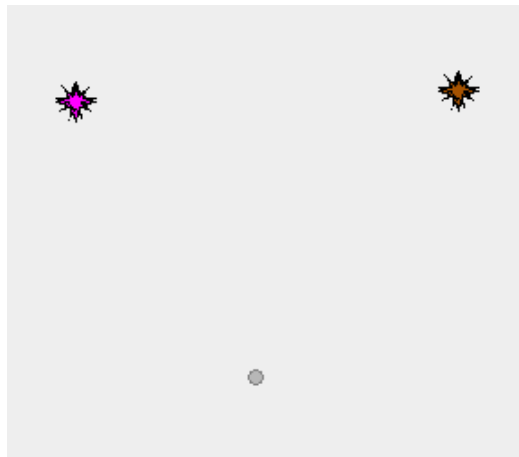
- 1      Scan all elements of  $X$ , looking for an element  $x$  whose nearest prototype from  $U$  has a different label than  $x$ .
- 2      Remove  $x$  from  $X$  and add it to  $U$
- 3      Repeat the scan until no more prototypes are added to  $U$ .

Use  $U$  instead of  $X$  for classification. The examples that are not prototypes are called "absorbed" points.<sup>[1]</sup>

Since we got  $k=2$  we need at least 3 data points. Screenshots made in Java Applet<sup>[2]</sup>



Step 1:

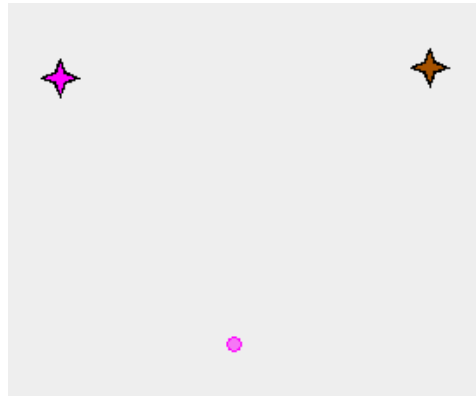


---

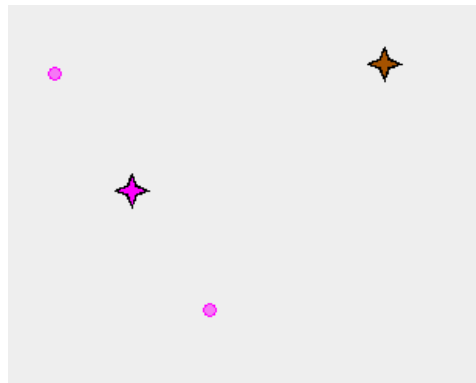
<sup>1</sup> [https://en.wikipedia.org/wiki/K-nearest\\_neighbors\\_algorithm](https://en.wikipedia.org/wiki/K-nearest_neighbors_algorithm)

<sup>2</sup> E.M. Mirkes, [K-means and K-medoids applet](#). University of Leicester, 2011.

Step 2:



Step 3



As the result we observe 2 clusters.

### Task N°3.

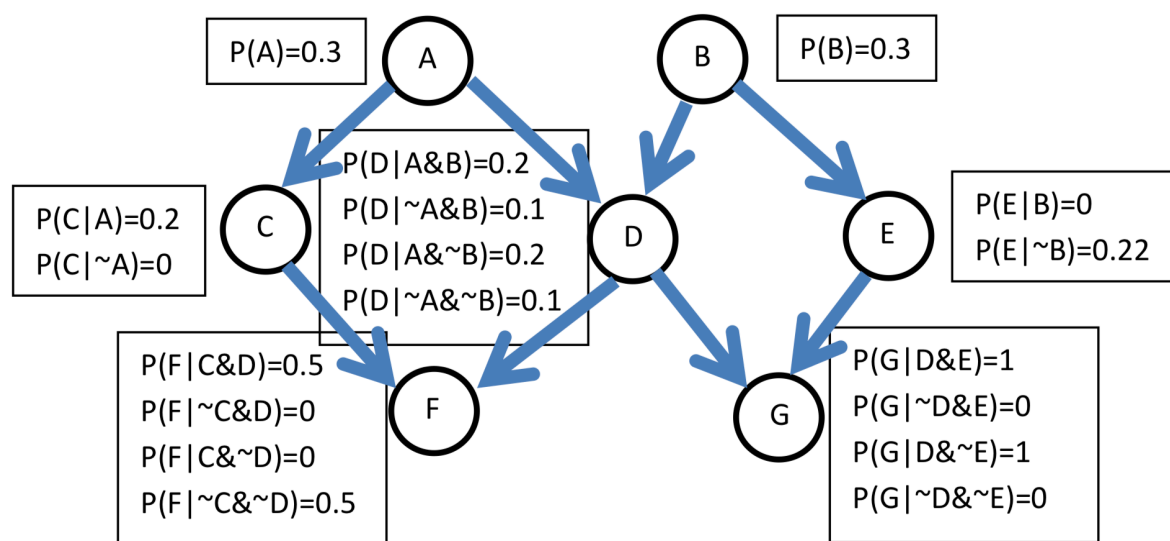
Describe the structure of Bayes net: what are vertices, directed edges between vertices, and what information do we keep with vertices? Give the definition of conditional independence. Compute the probabilities  $P(F \& G)$  and  $P(B | G)$  from the given Bayes net:

Vertex of Bayes Network - is the event that can occurs.

Direct edges of Bayes Network is describe connections between parent and child vertices.

Parent event have affection on child event.

Conditional Independence - In [probability theory](https://en.wikipedia.org/wiki/Conditional_independence), two events  $R$  and  $B$  are **conditionally independent** given a third event  $T$  precisely if the occurrence of  $R$  and the occurrence of  $B$  are [independent](https://en.wikipedia.org/wiki/Conditional_independence) events in their [conditional probability distribution](https://en.wikipedia.org/wiki/Conditional_independence) given  $T$ . [3]



<sup>3</sup> [https://en.wikipedia.org/wiki/Conditional\\_independence](https://en.wikipedia.org/wiki/Conditional_independence)

$P(F \& G)$  computed as follows:

Since some of the probabilities are equal to zero we have only 3 terms:

$$\begin{aligned} P(G \& F) &= (P(G|D \& E) \times P(F|C \& D) \times P(E|\neg B) \times P(C|A) \times P(D|A \& \neg B) \times P(A) \times P(\neg B)) + \\ &+ (P(G|D \& \neg E) \times P(F|C \& D) \times P(\neg E|B) \times P(C|A) \times P(D|A \& B) \times P(A) \times P(B)) + \\ &+ (P(G|D \& \neg E) \times P(F|C \& D) \times P(\neg E|\neg B) \times P(C|A) \times P(D|A \& \neg B) \times P(A) \times P(\neg B)) \\ &= 0.000924 + 0.0018 + 0.003696 = 0.00642 \end{aligned}$$

$P(B|G)$  computed as follows:

$$P(B|G) = \frac{P(B \& G)}{P(G)}$$

$$P(G) = P(B \& G) + P(\neg B \& G)$$

$$\begin{aligned} P(B \& G) &= (P(G|D \& \neg E) \times P(\neg E|B) \times P(D|A \& B) \times P(A) \times P(B)) + \\ &+ (P(G|D \& \neg E) \times P(\neg E|B) \times P(D|\neg A \& B) \times P(\neg A) \times P(B)) = 0.039 \end{aligned}$$

$$\begin{aligned} P(\neg B \& G) &= (P(G|D \& E) \times P(E|\neg B) \times P(D|A \& \neg B) \times P(A) \times P(\neg B)) + \\ &+ (P(G|D \& \neg E) \times P(\neg E|\neg B) \times P(D|A \& \neg B) \times P(A) \times P(\neg B)) + \\ &+ (P(G|D \& E) \times P(E|\neg B) \times P(D|\neg A \& \neg B) \times P(\neg A) \times P(\neg B)) + \\ &+ (P(G|D \& \neg E) \times P(\neg E|\neg B) \times P(D|\neg A \& \neg B) \times P(\neg A) \times P(\neg B)) = 0.091 \end{aligned}$$

So:

$$P(G) = 0.039 + 0.091 = 0.13$$

$$P(B|G) = \frac{0.039}{0.13} = 0.3$$

You can find table used for calculations attached to report.

## Task N°4.

A Paradox. A PhD student analysed data about efficiency of video lectures. He has analysed performance of 300 students in Cambridge and 300 students from another university. The data are summarised in the table below. Do the video courses increase performance? According to the table for Cambridge, the answer is 'No'. Similarly, the answer is negative for another university. However, ignoring the information about the university, and collating the data into one combined table, we find that video helps! Should we recommend the video courses? Create a Bayes net for this problem. How is it possible to avoid the 'paradox' by sampling?

| Cambridge | 1 <sup>st</sup> class | Not 1 <sup>st</sup> class | % of 1 <sup>st</sup> |
|-----------|-----------------------|---------------------------|----------------------|
| Video     | 120                   | 80                        | 60%                  |
| Standard  | 75                    | 25                        | 75%                  |

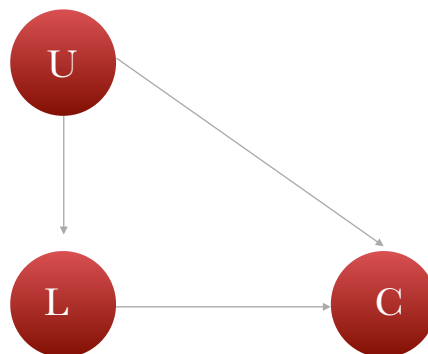
| Another Uni | 1 <sup>st</sup> class | Not 1 <sup>st</sup> class | % of 1 <sup>st</sup> |
|-------------|-----------------------|---------------------------|----------------------|
| Video       | 10                    | 90                        | 10%                  |
| Standard    | 30                    | 170                       | 15%                  |

| Combined | 1 <sup>st</sup> class | Not 1 <sup>st</sup> class | % of 1 <sup>st</sup> |
|----------|-----------------------|---------------------------|----------------------|
| Video    | 130                   | 170                       | 43.3%                |
| Standard | 105                   | 195                       | 35%                  |

This paradox called Simpson's Paradox. It appears in several different groups of data but disappears or reverses when these groups are combined. It is sometimes given the descriptive title reversal paradox or amalgamation paradox.<sup>[4]</sup>

In order to solve this task I used example<sup>[5]</sup>

We can write distribution as: (U-Universiy, L-lecture type, C -class)



<sup>4</sup> [https://en.wikipedia.org/wiki/Simpson%27s\\_paradox](https://en.wikipedia.org/wiki/Simpson%27s_paradox)

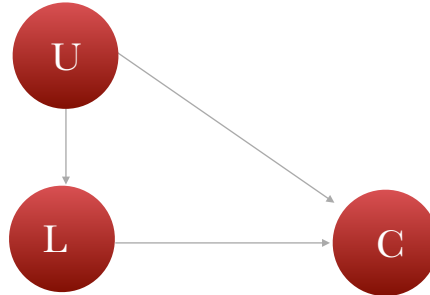
<sup>5</sup> <http://www.cimat.mx/~src/slides4.pdf>



Observable Calculation computed as:

$$P(U \& L \& C) = P(C | U \& L) \times P(U | L) \times P(U)$$

In order to resolve the paradox we need to consider new distribution conditioned on lecture type, not on university. Term  $P(U | L)$  play no role :



$$\tilde{P}(U \& C | L) = P(C | U \& L) \times P(U)$$

$$P(C || L) = \sum_U \tilde{P}(U \& C | L) = P(C | U \& L) \times P(U)$$

Applying new distribution we can get away from paradox:

$$P(1_{st}Class | Video) = 0.5 \times 0.6 + 0.5 \times 0.1 = 0.35 = 35 \%$$

$$P(1_{st}Class | Standard) = 0.5 \times 0.75 + 0.5 \times 0.15 = 0.45 = 45 \%$$

As we can see paradox resolved. Video doesn't increase performance and we should not recommend them. Standard courses are better.

## Task N°5.

Give the definitions of entropy, information gain and relative information gain. Calculate the entropies  $H(C)$ ,  $H(C|A)$ ,  $H(C|B)$ , information gain and relative information gain  $IG(C|A)$ ,  $RIG(C|A)$ ,  $IG(C|B)$ , and  $RIG(C|B)$  for the following data table and create the decision tree for the target attribute C.

**Information entropy** is defined as the [average](#) amount of [information](#) produced by a [stochastic](#) source of data.<sup>[6]</sup> In other words entropy is the measure of uncertainty.

Information gain is the change in Information Entropy.

It can be written as:

$$H(X) = \sum_{i=1}^n P(x_i) I(x_i) = - \sum_{i=1}^n P(x_i) \log_b P(x_i),$$

Where X is system with positive values  $x_1, \dots, x_n$ , and b is the base of logarithm. Often b is equal to 2.

Here is the entropy for the system:

$$H(C) = -\frac{1}{4} \log \frac{1}{4} - \frac{1}{4} \log \frac{1}{4} - \frac{1}{2} \log \frac{1}{2} = \frac{3}{2} = 1.5$$

Then we have table containing probabilities of C given A and B:

| A\C | -1  | 0   | 1   |
|-----|-----|-----|-----|
| 0   | 2/4 | 2/4 | 0   |
| 1   | 0   | 2/4 | 2/4 |

| B\C | -1  | 0   | 1   |
|-----|-----|-----|-----|
| 0   | 2/4 | 1/4 | 1/4 |
| 1   | 2/4 | 1/4 | 1/4 |

<sup>6</sup> [https://en.wikipedia.org/wiki/Entropy\\_\(information\\_theory\)](https://en.wikipedia.org/wiki/Entropy_(information_theory))

$$H(C|A) = \frac{1}{2}(-\frac{2}{4} \log \frac{2}{4} - \frac{2}{4} \log \frac{2}{4}) + \frac{1}{2}(-\frac{2}{4} \log \frac{2}{4} - \frac{2}{4} \log \frac{2}{4}) = 1$$

$$H(C|B) = \frac{1}{2}(-\frac{2}{4} \log \frac{2}{4} - \frac{1}{4} \log \frac{1}{4} - \frac{1}{4} \log \frac{1}{4}) + \frac{1}{2}(-\frac{2}{4} \log \frac{2}{4} - \frac{1}{4} \log \frac{1}{4} - \frac{1}{4} \log \frac{1}{4}) = \frac{3}{2}$$

Information Gains:

$$IG(C|A) = H(C) - H(C|A) = 1.5 - 1 = 1$$

$$IG(C|B) = H(C) - H(C|B) = 1.5 - 1.5 = 0$$

Relative Information Gains:

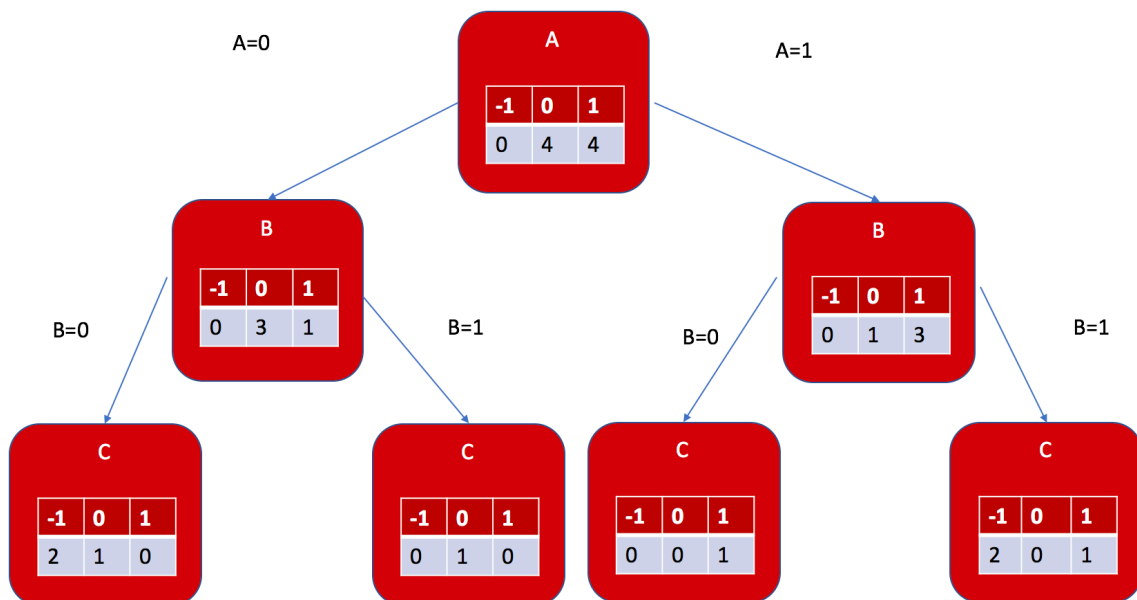
$$RIG(C|B) = \frac{H(C) - H(C|A)}{H(C)} = \frac{1.5 - 1}{1.5} = \frac{1}{3}$$

$$RIG(C|B) = \frac{H(C) - H(C|B)}{H(C)} = \frac{1.5 - 1.5}{1.5} = 0$$

Since Information Gain for B equal to 0 we will split node at A.

Now we have only one feature to work with.

Here is the final tree. Second row of the tables is the number of samples.



## References:

1. [https://en.wikipedia.org/wiki/K-nearest\\_neighbors\\_algorithm](https://en.wikipedia.org/wiki/K-nearest_neighbors_algorithm)
2. E.M. Mirkes, [K-means and K-medoids applet](#). University of Leicester, 2011
3. [https://en.wikipedia.org/wiki/Conditional\\_independence](https://en.wikipedia.org/wiki/Conditional_independence)
4. [https://en.wikipedia.org/wiki/Simpson%27s\\_paradox](https://en.wikipedia.org/wiki/Simpson%27s_paradox)
5. <http://www.cimat.mx/~src/slides4.pdf>
6. [https://en.wikipedia.org/wiki/Entropy\\_\(information\\_theory\)](https://en.wikipedia.org/wiki/Entropy_(information_theory))