

Machine Learning 2014-2015—Extra Resit

22 April 2015

The exam is **open book**: you can use Flach's *Machine Learning* (as a print of the PDF or as a proper book) as well as the lecture slides and any notes you've taken. You can use a calculator as long as it doesn't provide (internet) connectivity.

The exam grade determines 50% of your final grade for the course, the other half is determined by your project grade. The exam grade must at least be 5.5 to pass the course.

Good luck!

Question	Points
1. Short answers	15
2. Problem Type	5
3. Decision Trees	15
4. Neural Networks and Regression	20
5. Support Vector Machines	15
6. Clustering	15
7. Performance Measures and Generalisation	15
Total	100

Questions

1. Short answers no justification required (3 Points for each question)

- (True or False): Decision trees are a typical example of unsupervised learning.
- (True or False): Neural networks without hidden nodes can learn anything that more general neural networks can, it just takes longer.
- (True or False): Kernel functions allow nearest neighbour algorithms to take all datapoints into account instead of only the k nearest.
- (True or False): If a model overfits, its performance on a validation set is typically worse than on the training set.
- (True or False): Regression models can only use numerical features.

2. Problem Type (5 points)

For each of the following learning problems, please indicate whether it is a prediction, regression or classification problem. (An explanation is not required.)

- A biologist has given different amounts of food to different rats in his laboratory. He has recorded the weight of each rat after two months. Now he wants to learn how the weight of the rats depends on the amount of food they get.
- Each spring a farmer counts the number of newborn sheep. Based on his counts of the previous years he wants to estimate the number of newborn sheep in the coming year.
- A computer program tries to determine whether a newspaper article is about politics based on the number of times the article contains the following words/phrases: 'law', 'sports', 'newspaper', 'hockey', 'elections', 'human rights' and 'party'.

1

3. Decision Trees

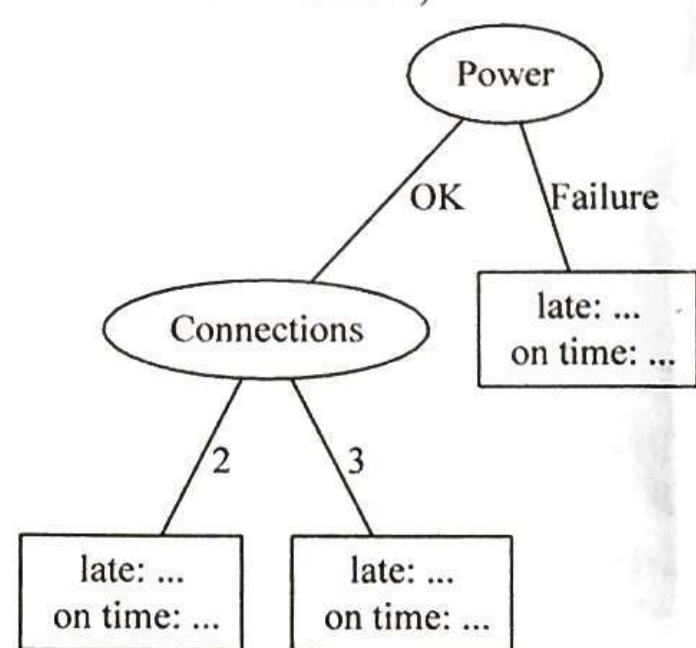
- (5 points) Based on recent experiences, you decide to build a decision tree to predict if you'll arrive on time for an exam based on the weather, if there is a power outage and the number of connections you have to make with public transport. You collect the data in table 1.

Weather	Power	Connections	On time
Fine	Failure	2	On time
Fine	OK	2	On time
Fine	OK	2	On time
Cloudy	OK	2	On time
Cloudy	OK	2	Late
Rain	OK	2	Late
Rain	OK	2	Late
Rain	OK	2	Late
Rain	Failure	3	Late

Table 1: On time or not?

Which attribute would the ID3 algorithm choose to use for the root of the tree (assume no pruning)? Show your calculations.

- (5 points) Using a different technique, someone developed the decision tree below. Draw the corresponding probability estimation tree (basing the probability estimates on the data in table 1).

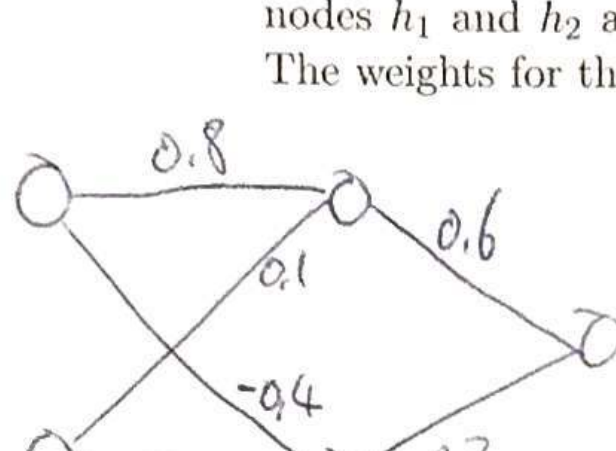


- (5 points) When pruning decision trees, the decision to prune a node is based on the accuracy of the resulting tree on a validation set. What would go wrong if we used the train set instead of this validation set?

2

4. Neural Networks and Regression

In this question, we consider a neural network with two input nodes x_1 and x_2 , two hidden nodes h_1 and h_2 and one output node o . The nodes have sigmoid activation functions. The weights for the connections are as follows:



Connection	Weight
$x_1 \rightarrow h_1$	0.8
$x_1 \rightarrow h_2$	-0.4
$x_2 \rightarrow h_1$	0.1
$x_2 \rightarrow h_2$	-0.5
$h_1 \rightarrow o$	0.6
$h_2 \rightarrow o$	0.2

Table 2: Neural Network Weights

- (5 points) What is the output of the network if the input is (0.3, 0.3)? Show your calculations.
- (5 points) Let's assume that for this input, the target value $t = 0.0$ and the learning rate $\eta = 0.1$. Compute the new weight for the connection $h_1 \rightarrow o$ using backpropagation. Show your calculations.
- (5 points) When training a neural network, what can happen if the learning rate η is set too high?
 - The network becomes overfitted to the solution it is trained with.
 - The network does not stabilise and the algorithm loops indefinitely
 - The network's output values are no longer guaranteed to be between 0 and 1
- (5 points) Describe in your own words what reaching a local optimum (e.g. in gradient descent) means.

5. Support Vector Machines

In this question, we'll consider a dataset for two classes A and B, each with three points. Class A: $\{A_1 = (1, 1); A_2 = (-1, 3); A_3 = (2, 6)\}$ and Class B: $\{B_1 = (-1, -2); B_2 = (1, -3); B_3 = (-5, -7)\}$

- (5 points) Identify the support vectors for a linear Support Vector Machine on this data (Hint: draw a plot of the data).
- (5 points) What is the optimal separating line for this data? You may show your answer as an expression or draw a diagram.
- (5 points) The kernel trick allows SVMs to define a non-linear separation between classes. Which of the following statements is true (more than one may apply)?
 - The kernel function replaces the dot-product to calculate the distance between points.
 - The kernel trick requires an explicit transformation to a higher-dimensional feature space.
 - A kernel can be a polynomial function.

3

6. Clustering

- (5 points) Name two essential differences between hierarchical and k-means clustering.
- (5 points) Describe, in your own words, the initialisation procedure of the k-means algorithm.
- (5 points) Consider the data set with 5 cases labelled A,B,C,D and E in figure 1. Draw the dendrogram that would result from a single-link agglomerative hierarchical clustering of that data (assuming Euclidean distance).

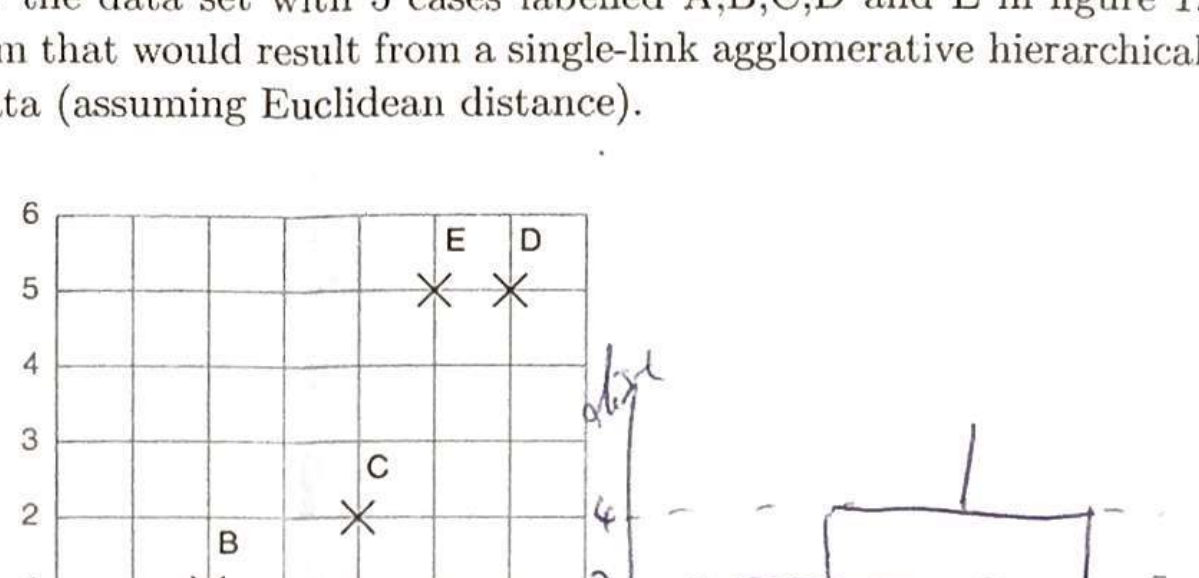


Figure 1: A clustering data set

7. Performance Measures and Generalisation

- (5 points) Given a binary classification problem with dataset D and two learning algorithms L_1 and L_2 , describe how you would decide which algorithm to use.
- (5 points) Fitting a 14th degree polynomial to 15 data points using least squares regression will result in overfitting. For each of the following two cases, argue whether they would help against overfitting:
 - Suppose that instead of minimising the sum of the squares of the errors, we would minimise the sum of the absolute values of the errors.
 - Suppose that instead of 15 data points we had 100 000.
- (5 points) For a binary classification problem, a classifier is evaluated on a testset of 1000 cases. There are 700 positive and 300 negative examples in the dataset. It correctly identifies 586 positive cases and 234 negative cases.

Give an estimate of the true accuracy of this classifier, and a 95% confidence interval around that estimate (you may give your confidence interval in the form of an expression). Note: $Z_{95} = 1.96$.

4