



Machine Learning F2024.

Sila. March 20th, 2024.

The Mandatory Assignment. Spring 2024.

You must hand-in individual reports, group hand-ins are not allowed in this assignment. Your hand-in can be written in danish or in english. Your solution consist of your Python code, plus a small document, 1 page, that tells what you have done.

*Your report must be handed in on **Canvas on Tuesday the 2nd of April 2024 at 20.00 O'clock at the latest.***

If/when the report consists of multiple files (report, scripts) then it must be zipped into one file. Make sure to make a zip file correctly. Expected size is 1-2 MB.

Generally, there will not be given individual feedback, for this exercise, but we will talk about the assignment next time in class.

Mandatory Exercise. Titanic.

In the class you should have worked with a limited version of real data from Titanic detailing whether passengers survived or died and other data. We want to predict who survives and who dies based on other passenger data.

In this exercise you will look at a larger dataset from the Titanic and use machine learning on this dataset.

The dataset can be downloaded in .csv format from Canvas. It consists of a header and 800 passenger records (And yes...Multiple versions of the Titanic passenger exist, and an alternative version with 887 records, Titanic_alternative.csv, is also uploaded to Canvas. **But: Don't** use this alternative file until you have completed this entire exercise with the 800 passenger excel file. The alternative file is only meant for those who want to dig a little deeper here, after having answered all of the questions in this exercise, or run into problems with the 800 passenger excel file).

The dataset can be imported into a text program or a .csv reader (I recommend using openoffice for this which is really good at displaying the data from .csv files).

Here is a small explanation of the data (the name and the passengerId should be self-explanatory):

Variable in data	Definition	Key
Survival	Survived or not	1 = survived, 0 = died
pclass	Passenger class	1 = 1 st class, 2 = second class, 3 = third class
Sex	Sex	Male or female

age	Age in years	
sibsp	# of siblings / spouses aboard the Titanic	
parch	# of parents / children aboard the Titanic	
Ticket	A ticket number	
fare	The price paid for the ticket	
cabin	The cabin number	
embarked	Port city where the passenger embarked	C = Cherbourg, Q = Queenstown, S = Southampton

There are several sections below marked with the capital letters. In each section there are several questions. Your job is to answer as many questions correctly as you can.

You should also hand in a .py file or files with your code included in the .zip file

A) Getting to know the data and the problem.

What kind of machine learning problem are we talking about in this exercise? Be as precise as possible.

How many different features (not counting labels (y-values) as a feature here) are there in this dataset?

How many of these features would the Y-data (the labels) consist of?

B) Cleaning the data

You have to decide on which of the features for the X training dataset you want to use and then later what to do with missing data for the features you keep.

Think about if there are some features that will most likely have no impact on the survival and discuss why this could be the case. Of course it can be “dangerous” to just remove features, as they might be of even some limited value to a training model, but in this dataset there are features which can be safely removed.

You can back up claims also by looking to see if there is some relationship or not between the feature and the survival rate. An easy way to see if there is a relationship between two variables is to check for correlation or doing a graph like a scatterplot to see it visually (this was also part of a previous exercise from class). Although, no correlation is not an absolute guarantee that a feature can be removed – remember there could be multi-dimensional relationships between 3 or 4 features that you cannot see.

The feature(s) removed will then not be used in the training and also not in the test set.

Some of the features also have missing data. For instance the age field has missing data. You need to discuss how you will handle missing data? (Think back to earlier discussions we have had in class about similar problems).

You also need to consider that some the data is not numerical data, so that data would need to be converted into numbers (in both training and test set). This is easy to do with pandas – here is an example: `xtrain['Sex'] = xtrain['Sex'].replace(['female'], 1.0)`

This changes all the 'female' labels into 1.0 instead of test for the sex (of course the male would also need a value).

After cleaning (and scaling your data – again see slides or previous exercise on this) then you should split your data into two sets. One for training and then one for evaluating performance.

There are 800 samples in this dataset.

How many would you put into the training set and how many would you put into the test set? Give reasons for your decision.

C) Choosing a model and doing training

We have worked with a few models, such as randomforests, decision trees, neural networks etc.

Choose at least 1 model for doing training and explain your reason for this choice. Train the model using the training data.

D) Evaluating performance on the test set.

Previously, in the course, in the exercises and on the slides (and in the book also) we have discussed how to evaluate performance.

Give the precision and recall rates for your model.

Also give the confusion matrix and use the values in the confusion matrix to calculate the percentage of how many of the samples in the set you got right?

Look at the numbers in the confusion matrix also – what is your model best at and what is it worst at?

E) Experiments

Try to do a few experiments – either with the parameters of your chosen model or by choosing another model and comparing prediction performance with your old model. Document your experiments and results.