

Assignment 1 - Advanced *Data Mining Techniques* A Special Dataset

Deadline: 16/04/2021, 23:59

INTRODUCTION

This document introduces you to the first assignment of the Data Mining Techniques course at the VU. This is a group task (3 members), please make sure all team members contribute to the work as expected. The assignment essentially consists of creating predictive models for a challenging dataset, not from the perspective of the size of the dataset, but more the nature of the data. You can earn a total of 100 points.

DATASET AND PROBLEM

As a first step, let us look at the dataset we are faced with. The domain from which the dataset originates is the domain of mental health. More and more smartphone applications are becoming available to support people suffering a depression. These applications record all kinds of sensory data about the behavior of the user and in addition frequently ask the user for a rating of the mood. A snapshot of the resulting dataset is shown in Table 1. The dataset contains ID's, reflecting the user the measurement originated from. Furthermore, it contains time-stamped pairs of variables and values. The variables and their interpretation are shown in Table 2.

Using this dataset, we would like to build a predictive model that is able to *predict* the *average mood* of the user *on the next day* based on the data we obtained from the user *on the days*

Table 1: Snapshot of the data

ID	Timestamp	Variable	Value
AS14.01	26-02-2014 15:00.00	mood	6
AS14.01	26-02-2014 15:21.00	activity	0.031
AS14.01	26-02-2014 15:55.00	screen	103.1
AS14.01	27-02-2014 16:00.00	mood	6
AS14.01	27-02-2014 12:00.00	appCat.builtin	0.052

Table 2: Variables in the dataset

Variable	Explanation
mood	The mood scored by the user on a scale of 1-10
circumplex.arousal	The arousal scored by the user, on a scale between -2 to 2
circumplex.valence	The valence scored by the user, on a scale between -2 to 2
activity	Activity score of the user (number between 0 and 1)
screen	Duration of screen activity (time)
call	Call made (indicated by a 1)
sms	SMS sent (indicated by a 1)
appCat.builtin	Duration of usage of builtin apps (time)
appCat.communication	Duration of usage of communication apps (time)
appCat.entertainment	Duration of usage of entertainment apps (time)
appCat.finance	Duration of usage of finance apps (time)
appCat.game	Duration of usage of game apps (time)
appCat.office	Duration of usage of office apps (time)
appCat.other	Duration of usage of other apps (time)
appCat.social	Duration of usage of social apps (time)
appCat.travel	Duration of usage of travel apps (time)
appCat.unknown	Duration of usage of unknown apps (time)
appCat.utilities	Duration of usage of utilities apps (time)
appCat.weather	Duration of usage of weather apps (time)

before. This is illustrated graphically in Figure 1 below. In order to create such a predictive model, we need to perform some transformations and we need to decide on what features we want to use for these predictions.

TASK 1: PRE-PROCESS THE DATASET (40 POINTS)

Essentially there are two approaches you can consider to create a predictive model using this dataset: (1) use an machine learning approach that can deal with temporal data (e.g.

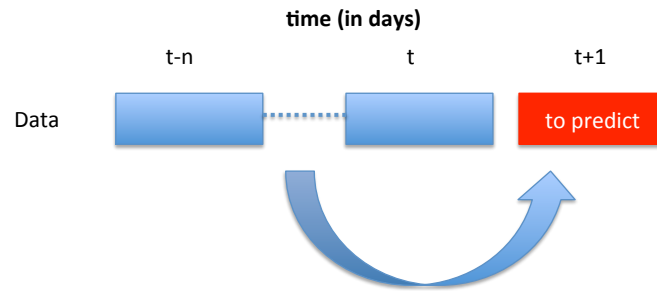


Figure 1: Predictive model

ARIMA, recurrent neural networks) or you can try to aggregate the history somehow to create attributes that can be used in a more common machine learning approach (e.g. SVM, decision tree). For instance, you use the average mood during the last five days as a predictor. Ample literature is present in the area of temporal data mining that describes how such a transformation can be made. We are going to focus on such a transformation in this part of the assignment. What we are trying to do is illustrated in Figure 2.

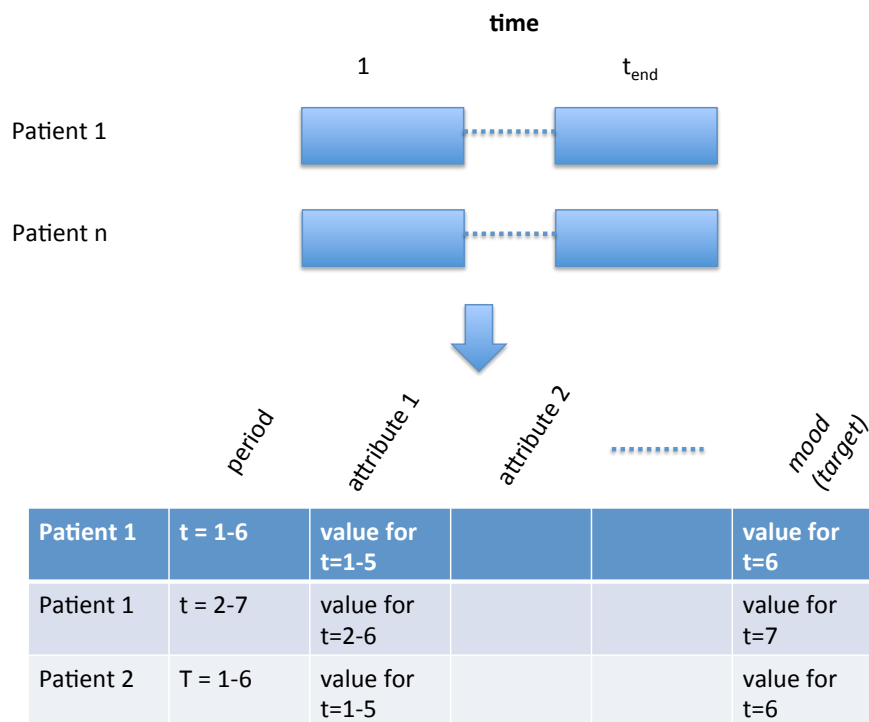


Figure 2: Predictive model

In the end, we end up with a dataset with a number of training instances per patient (as you have a number of time points for which you can train), i.e. an instance that concerns

the mood at $t=1$, $t=2$, etc. Of course it depends on your choice of the history you consider relevant from what time point you can start predicting (if you use a windows of 5 days of history to create attributes you cannot create training instances before the 6th day). To come to this dataset, you need to:

1. Define attributes that aggregate the history, draw inspiration from the field of temporal data mining.
2. Define the target by averaging the mood over the entire day.
3. Create an instance-based dataset as described in Figure 2.

Talk to your TA to discuss your ideas before finalizing them. In the end, you should describe and argue your choices clearly and back them up with scientific literature. Next, perform a preliminary analysis (e.g. look at correlations) on the usefulness of the attributes you have defined.

TASK 2: LEARN USING THE DATASET (40 POINTS)

In the next step, we are going to use our dataset to create a predictive model. You can make your own choice whether you want to create individual models per patient or a single model for all patients. You will need to study three variants of predictive models:

1. A variant where you use the pre-processed dataset you identified in Task 1 in combination with a machine learning technique you consider appropriate.
2. A variant where you apply a learning algorithm that is able to cope with this temporal data (e.g. ARIMA, recurrent neural networks, etc.).
3. Implement a benchmark: predict the mood on the next day by just saying it is the same as the previous day.

Define a proper performance metric and create a solid evaluation setup. Describe and argue your choices again and show the results you have obtained. Create graphs to illustrate the performance in an insightful way.

TASK 3: EVALUATE AND REFLECT ON YOUR RESULTS (20 POINTS)

Finally, analyse the results in detail both using a more statistical view and by means of your interpretation. Argue what the pros and cons of the different approaches are.

REPORT

We would like you as a group of 3 to prepare a report with the following in mind:

- The report should be submitted via Canvas by 16/04/2021 23:59. This is a strict deadline, please try to respect that, otherwise points will be deducted (1 full point per day)
- Please format the document according to the LNCS guidelines. Templates are available on Canvas for both LaTeX and Microsoft Word, do not deviate from these templates. Note that you do not need to include an abstract in your report. The paper **should not exceed 10 pages including all figures and tables, but excluding references** (references do not count for the number of pages to encourage you to cite all relevant work). With the page limit, my aim is to challenge you to report only what is necessary.
- Make sure we can identify your report, i.e., your group number, names and student numbers should be in the document's header.
- Make an attempt to make the report look professional. Have a short introduction of your document, use appropriate language, etc. Let's say, if you gave your report to the manager of your DM project at a company, they would need to be able to understand it and conclude that it's a good project start.

GRADING

Marking will be based on the tasks as reflected by quality of the report (so content, style, etc. all matter). You can get maximum 100 marks for this assignment. You will need at least 55 to pass. Also, 100 points are only given to students whose reports are of exceptional quality, and they also should report something we did not specifically ask for (in other words, we value proactivity and creativity).

The grading scheme we apply for this assignment is shown in Table 3 below.

Table 3: Grading scheme

Task	Grading Component	Weight
1	Exploratory data analysis	10
	Explain general setup of feature engineering	10
	Use of scientific literature supporting the setup	5
	Rationale for choice of final attributes	15
2	Feature Engineered model	8
	Temporal Model	8
	Benchmark Model	4
	Evaluation (validation / test)	10
	Illustrate performance with graphs	10
3	Analyse results using statistics	8
	Analyse results by interpretation	8
	Pros and cons of different approaches	4
Deductions	Extra page	-10
	Wrong formatting	-10
Total		100