

RESEARCH

Introduction to the profile areas of data sciences: project 8

Rohan Chitte] Natalja Amiridze] Michael Onyebuchi Ohaya

Abstract

Goal of the project: The goal of this project was to perform data handling as well as data analysis on StudentLife study.

Main results of the project: The main results showed which features contributed in stress occurrence among students. Furthermore, Evaluation of student's current state based on other features.

Personal key learnings: We learnt how to create a schema/database like structure to handle large and complex datasets. Moreover, we performed data imputation dynamically and found correlation between features using Pearson's correlation coefficient.

Estimated working hours: 15

Project evaluation: 2

Number of words: 1268

1 Scientific Background

Many people experience stress, in one form or another, throughout their lives. Stress can be defined as “a state, which is accompanied by physical, psychological or social complaints or dysfunctions and which results from individuals feeling unable to bridge a gap with the requirements or expectations placed on them. Over the years, there has been a significant growth in the use of smartphones for health monitoring. Smartphone applications can be used to monitor physical activity, food intake, sleep quality, the menstrual cycle, and other issues related to health and well-being. Furthermore, monitoring can take place automatically through the sensors embedded within the smartphone, such as accelerometer and microphone, whereas others require users to interact with the application to record data. Subjective self-assessed stress can be measured using smart mobile devices via ecological momentary assessment (EMA). EMA is a collection of methods used to collect “assessments of subjects' current or recent states, sampled repeatedly over time, in their natural environment” Subjective stress can be assessed throughout the day using a time-based EMA where people are prompted to rate or answer questions about their “current stress level”

2 Goal

The aim of this project was to analyse the StudentLife study in particular the correlation stress with activity level and other features. This required importing the data into a database-like system as tables. Then it was necessary to create and

compare summary statistics for the data, pre-process the data with imputation of missing values, perform time series and correlation analyses. It was also to check other features for correlation with stress. Finally, it was to find out if the student's state can be predicted using other features

3 Data

The data used for the project came from the paper of Wang et al. The dataset was collected from all subjects, including automatic sensor data, behavioral interferences and self-reported EMA data with ground truth data including behavioral and mental health outcomes computed from survey instruments, as well as academic performance from spring and cumulative GPA scores provided by the registrar.

4 Results

4.1 Task 1: Is stress correlated to the activity-level?

Firstly, the data (.json-files containing EMA stress data for each student separately and also csv.files containing the activity data for each student separately) were loaded for students with available data as dictionary and the converted to pandas dataframe respectively. For EMA stress data some errors in the data could be identified (some rows of the column stress_level and location were in the column null_list to find. These incorrection could be solved by moving of respective rows to respective columns). As next the data were also checked for NaN-values. It could be 241 NaNs for stress_level identified and 0 NaNs for activity. After this preprocessing the data were imported to the database system Google Cloud Platform. As next, queries were generated to extract important figures for further statistical evaluations. First, statistics like average and standard deviation values for all data and for each student were extracted (see Figures 1, the output for activity data is for visibility not shown). In the next step, the average daily values as mean values

Index	all_student_id	total_stress_level_average	total_stress_level_std	total_stress_level_min	total_stress_level_max
0	2408	2.2691131490470916	1.350060561552417	1	5

Index	student_id	stress_level_average	stress_level_std	stress_level_min	stress_level_max
0	u00	2.263157894736842	1.310149315181356	1	5
1	u02	2.096774193646037	1.326514710334	1	5
2	u03	2.8709677419354044	1.7076677171680067	1	5
3	u05	3.3636363636363638	1.5666898036012806	1	5
4	u07	3.1754385964912273	1.4407687310393895	1	5
5	u08	1.813186813186813	0.8552074200916009	1	4
6	u09	2.428571428571429	1.3972762620115438	1	4
7	u10	2.4074074074074083	1.4599771924798823	1	5
8	u12	2.1875000000000004	1.3304740847993173	1	5
9	u13	3.0	NaN	3	3
10	u14	2.1749999999999999	0.7120753345337011	1	4
11	u15	2.0	1.2403473458920846	1	4

Index	student_id	stress_level	student_id_count	student_id_ratio
0	u00	1	3	0.0379746835443038
1	u00	2	32	0.4058632911392405
2	u00	3	14	0.1772151898734772
3	u00	4	11	0.13924050632911392
4	u00	5	16	0.20253164556862025
5	u00	6	3	0.0379746835443038
6	u01	1	16	0.4444444444444444
7	u01	2	6	0.16666666666666666
8	u01	3	3	0.08333333333333333
9	u01	4	9	0.25
10	u01	5	2	0.05555555555555555
11	u02	1	4	0.14285714285714285
12	u02	2	15	0.42857142857142855
13	u02	3	7	0.2
14	u02	4	1	0.02857142857142857
15	u02	5	7	0.2
16	u02	6	1	0.02857142857142857

Figure 1 Statistics: average and standard deviation values for EMA stress data

for all students were imputed and extracted (see Figures 2, the output for activity data is for visibility not shown) In the next step, time series for the students u10 and u16 were analysed with the help of Prophet library by extracting components

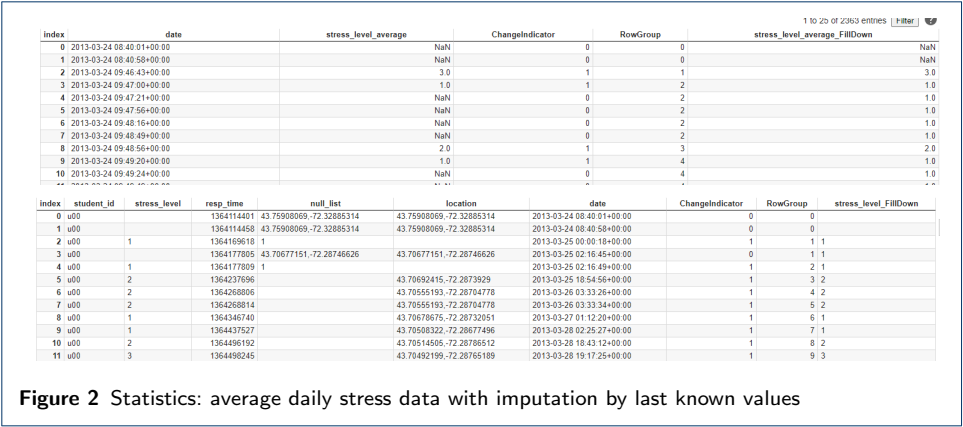


Figure 2 Statistics: average daily stress data with imputation by last known values

such as trend, weekly and daily (see Figures ?? and 4). For both students is trend of stress development growing with time. Looking at weekly component, it looks similar for both students, that the stress is high in the beginning and goes down until Friday. By the student u10 the stress grows on weekend again, by student u16 the stress grows on Saturday and goes down on Sunday. Looking at the daily component, it can be noticed that for the student u10 the stress is in the beginning of the day high and goes gradually down. For the student u16 the stress increases for period for 10 till 17 and is low for the rest of time. For both students is trend of

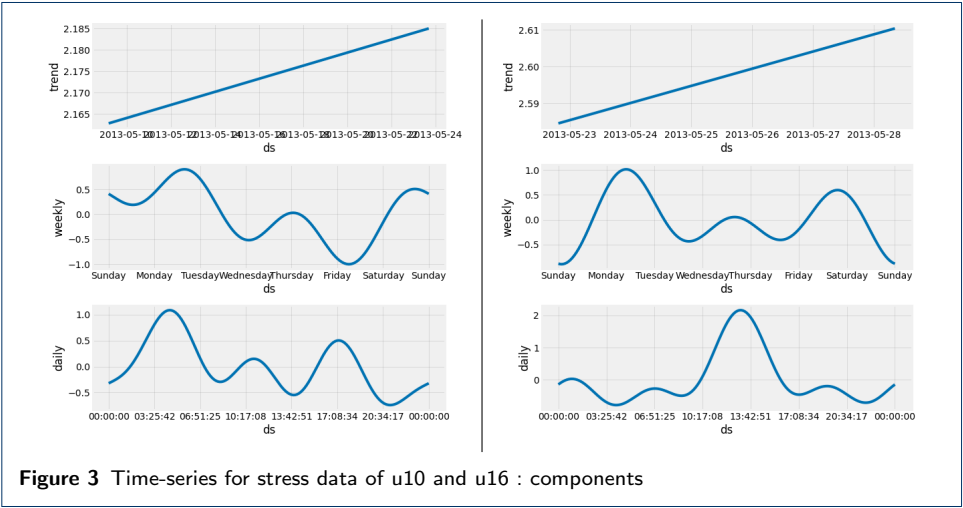


Figure 3 Time-series for stress data of u10 and u16 : components

activity development growing with time. By weekly component it can be seen the both students have high activity on Tuesday and Wednesday and on Sunday and low activity on Monday and Wednesday. Looking at daily component it is to see that student u10 has low activity from 3am till 10am and high activity from 10 am till 12pm. Student u16 has low activity from 6 am till 10am and high activity from 10am till 6am. To estimate influence of activity on stress the Pearson correlation was estimated for student u10 and u16 for all the time (see Figure 6 and ??). As can be seen, no significant correlation could be estimated. Performing correlation analysis on components show that activity has influence on the stress (see Figure 6 and 7) As it can be seen stress correlates negatively with activity both for weekly

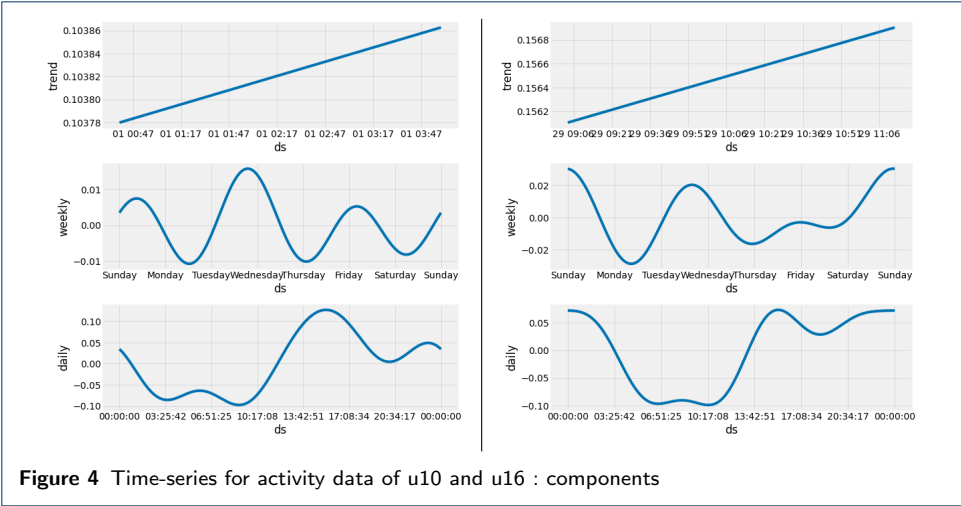


Figure 4 Time-series for activity data of u10 and u16 : components

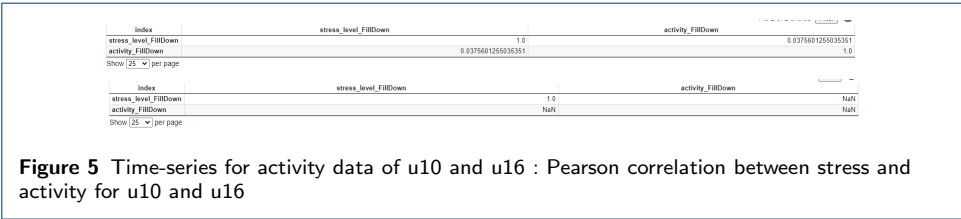


Figure 5 Time-series for activity data of u10 and u16 : Pearson correlation between stress and activity for u10 and u16

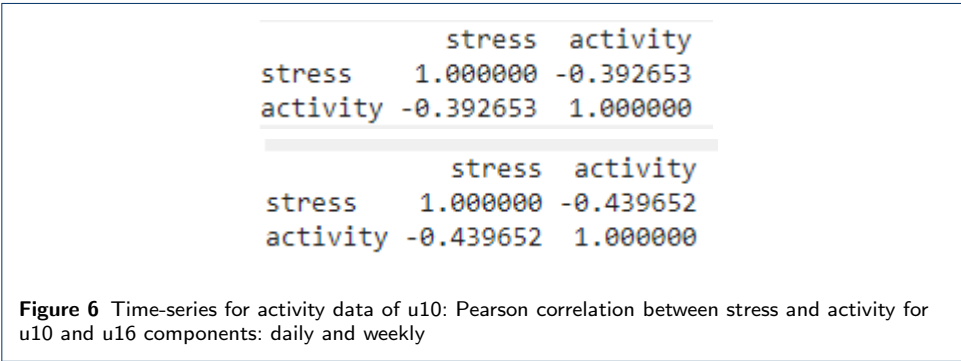


Figure 6 Time-series for activity data of u10: Pearson correlation between stress and activity for u10 and u16 components: daily and weekly

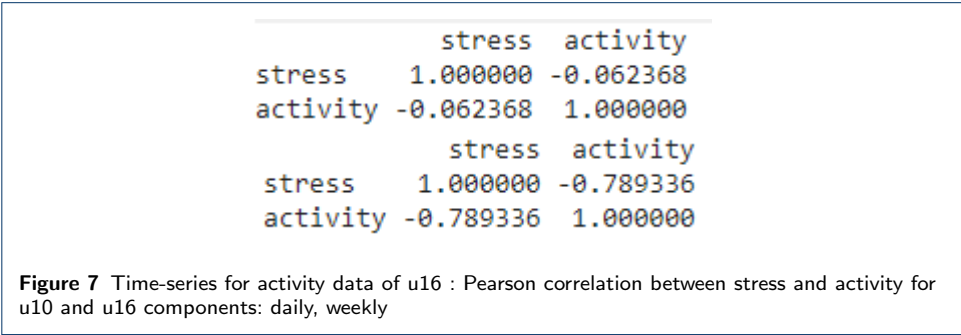


Figure 7 Time-series for activity data of u16 : Pearson correlation between stress and activity for u10 and u16 components: daily, weekly

and daily components for u10 and correlates negatively only for weekly component for u16 and has no correlation for daily component of u16.

4.2 Task 2: What else is stress correlated to?

To identify other features that might be correlated to stress, we first chose student no 10 and 16 as our data. To discover correlations we computed Pearson correlation coefficient also referred to as Pearson's r . The Pearson correlation coefficient, r , can take a range of values from +1 to -1. The stronger the association of the two variables, the closer the Pearson correlation coefficient, r , will be to either +1 or -1 depending on whether the relationship is positive or negative, respectively. P-value evaluates how well your data rejects the null hypothesis, which states that there is no relationship between two compared groups. Higher value of p means how likely it is that you would receive these results assuming that there is no correlation or relationship among the subjects. We used scipy stat's `pearsonr` class to compute r and p values. Below table shows the correlation of stress with sleep, exercise and events. As one would naturally think, stress does relate to sleep. The correlation value r between stress and sleep came out to be 0.040(Shown in table) which suggests a strong positive correlation. Further, we analysed if events or exercise are correlated to stress. Although, computed Pearson's r showed positive correlation of 0.011 and 0.012 between stress and exercise and stress and events respectively, it wasn't as significant when compared to sleep.

Correlations of stress with other features:			
	Other Features	r	p -value
0	Stress and Sleep	0.040342	0.672792
1	Stress and Exercise	0.011554	0.950811
2	Stress and Events	0.012278	0.952529

Figure 8 Correlation of Stress with other features

4.3 Task 3: Can we predict a student's state?

The classification using knn Classifier and logistic regression with `multi_class` set to multinomial was done using GPS, activity, SMS and CALL as features to predict Perceived Stress. As outcome was estimated the accuracy of 0.26, which is too low with knn classifier. Using logistic regression as classifier accuracy 0 was estimated. So it seems that the combination of selected features is not good for predicting StressPerceived values

5 Discussion

Such a study can prove to be beneficial in multiple ways on the social system withing a university. Prevalence of smart phones and web provides ease in conducting such a study. One could monitor specific as well as average student's health condition. Based upon on the identified results, the university could draw further analyses on what affected student's health. For instance, if the average students' health is being

impacted by the study or work load then it suggests that revisions can be made in the study program itself. On the other hand, if a student is performing well then his life pattern can be observed and adapted by other students to improve their performance.

6 Appendix

Name	Work Description
Michael	Abstract, Scientific Background, Results 8.2 and Code, Discussion Goal, Data, Results 8.1, 8.3 and Code
Rohan	
Natalja	

Table 1 Task responsibilities

References