

RESEARCH

Introduction of Profile Areas of Data Science :Project 9

Suresh Kumar Choudhary, Frenny Macwan and Aman Jain

Full list of author information is available at the end of the article

Abstract

Goal of the Project: To study the data given, about the students' personality traits, social behaviour and predict how co-related ((Or Different) are the individuals from others in the group along with predicting the individual's personality.

Main Result of the Project: We were able to fetch the information of the students who are highly co-related and highly not-co-related with the others members of the group, along with creating a classification and regression models to correctly classify individual traits. We found best classification with logistic regression model with auc 0.74 with target trait neurotic-ism and best regression model with mean squared error 0.42 and R^2 (coefficient of determination) 0.01 with trait conscientiousness.

Personal Key Leanings: Big five trait, 44 items, Prophet Library, Big Query, ICC (Inter-class Co-relation Coefficient)

Estimate working hours: 24

Project Evaluation: 2

Number of Words: 1746

Keywords: Time Series data trend; Inter-Class Corelation Coefficient; Classification; Regression; Personality Traits

Scientific Background

By usage of Mobile Sensing Methods in four studies, data has been collected from young adults across four different communication channels : Conversations, Phone Calls, Text Messages and Use of Messaging and Social Media, thorough which Behavioural Traits like Duration and Frequency of each channel were fetched along with, the Personality Traits through Surveys were fetched. Variances of individuals were calculated on their behaviour traits.

Goal

To analyse the social life of individuals using the mobile sensing (Majorly the Frequency and Duration of Social and Messaging Apps) and assessment techniques of personality traits .

Data

We used three Sample Data Sets [1] ; S1, S2 and S3 for the tasks.

- 1 **S1** - Data of 48 students were taken for a total of 66 days (10 Weeks). Data was acquired by the self tracked psychological experience by the participants and

also behaviours by the smart phone data. The data consist of the participant's age, id, their average call duration and average frequency of the calls of the entire 66 days period, along with the daily count of the duration and frequency of the calls in separate table as well. We create the tables(schema) in the GoogleBigQuery and imported them using gbq commands (gbq.read_gbq). As part of preprocessing, since most of the data which went missing were numeric values, we took the mean of the factor and filled it where the values weren't available.

- 2 **S2** - Data of 118 1st year Students of a UK College were taken in two different phases of two weeks each, in a gap of three months. Phase 1, where the mobile details were tracked by the app Easy-M faced some technical issues and hence only Phase 2 Data was considered where the behavioural factors like Frequency and Duration of Incoming-Outgoing calls, Incoming-Outgoing SMS were tracked for 28 Students through MyLife Logger App. All tables having _long as the suffix tracked the daily frequency and duration of calls and SMS for each participant, whereas the tables having _wide suffix (In this Sample only One), tracked the average behaviours of the participants.
- 3 **S3** - Data of 137 Students and Employees of Southern German University were taken. The tracking happened for 8 Weeks (60 Days). From Day 2 to Day 31 (Total of 30 Days) data was taken for this study, where Participants were rewarded with 30 Euros in exchange of their Behaviour Tracking on Phone and Personality feed-backs. Two factors (Duration and Frequency) were considered. For each user Behaviour factors (Duration and Frequency) were tracked for 31 days from apps like Messaging, Social, Calls In-Out, SMS In-Out. Along with daily tracking, an average duration and frequency was also given on daily average basis and on average time-of-day basis (Morning, Afternoon, Evening, and Night) for each participant along with their points on each personality traits after their inputs on the questionnaires. The five personality traits given in the data set are Openness, Conscientiousness, Extroversion, Agreeableness, and Neuroticism.

Result

We have followed the research paper [2] for Task1 and Task 2.

Task 1 : Are there individual differences in the daily social behavior?

- **Goal:** To find out the Inter-class Correlation Coefficient (ICC) between the students of a particular group (i.e. Sample S2)
- **Method:** Inter-class Correlation Coefficient describes how strongly units in the same group resemble each other. ICC values less than 0.5 are indicative of poor reliability, values between 0.5 and 0.75 indicate moderate reliability, values between 0.75 and 0.9 indicate good reliability, and values greater than 0.90 indicate excellent reliability.
- **Result:** To find individual differences in daily social behaviours, we are using the ICC method. ICC describes how strongly the data in the same group are related. To implement the ICC we have used the R package rpy2. We have created tables in the GoogleBigQuery and using the gbq.read_gbq method we

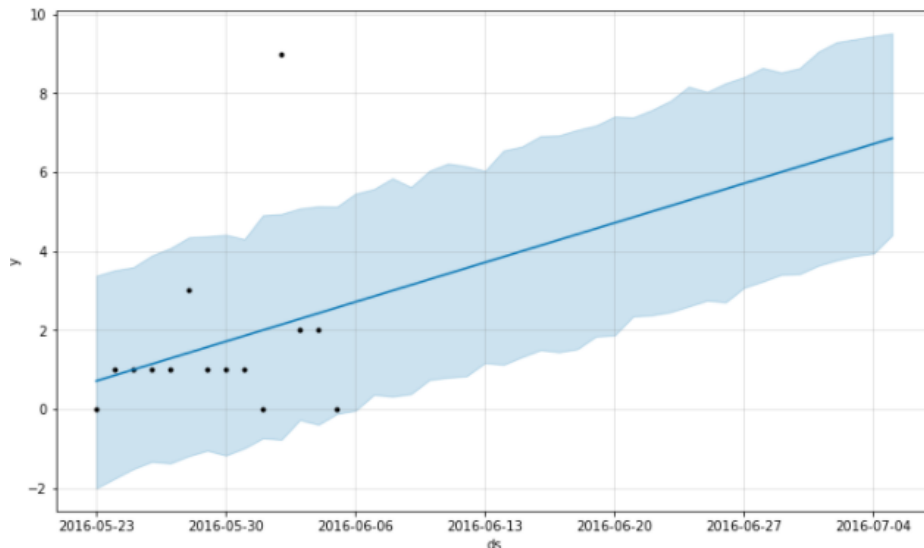
imported the data into data frame from GoogleBigQuery to Google Colab. We have considered Sample 2 data-set for analysis. For the Sample S2 we have different features like CALL IN FREQ, CALL OUT FREQ, CALL IN DURATION, CALL OUT DURATION, SMS IN FREQ, SMS OUT FREQ, SMS IN DURATION, SMS OUT DURATION. For four of the features, the data associated to one feature are comparatively strongly related to each other, which are SMS In Freq, SMS out Freq, SMS in Dur, SMS Out Dur, each having ICC of little more than 0.3, whereas Call In Freq, Call out Freq, Call in Dur and Call our Dur features have ICC of less than 0.2 and hence data points aren't strongly related to each other. For all the features and the exact ICC, please refer to the Fig 1

Figure 1 ICC calculation on sample 2 data-set

	CALL IN FREQ	CALL OUT FREQ	CALL IN DUR	CALL OUT DUR	SMS IN FREQ	SMS OUT FREQ	SMS IN DUR	SMS OUT DUR
ICC	0.1115708	0.1944612	0.1060321	0.1464291	0.39286	0.3549673	0.3001231	0.3347373
Lower	0.04356998	0.1025994	0.03976695	0.06790283	0.2598036	0.227266	0.1821257	0.2103536
Higher	0.2397746	0.3518347	0.2317776	0.2885876	0.5768427	0.5391687	0.4811719	0.5182721

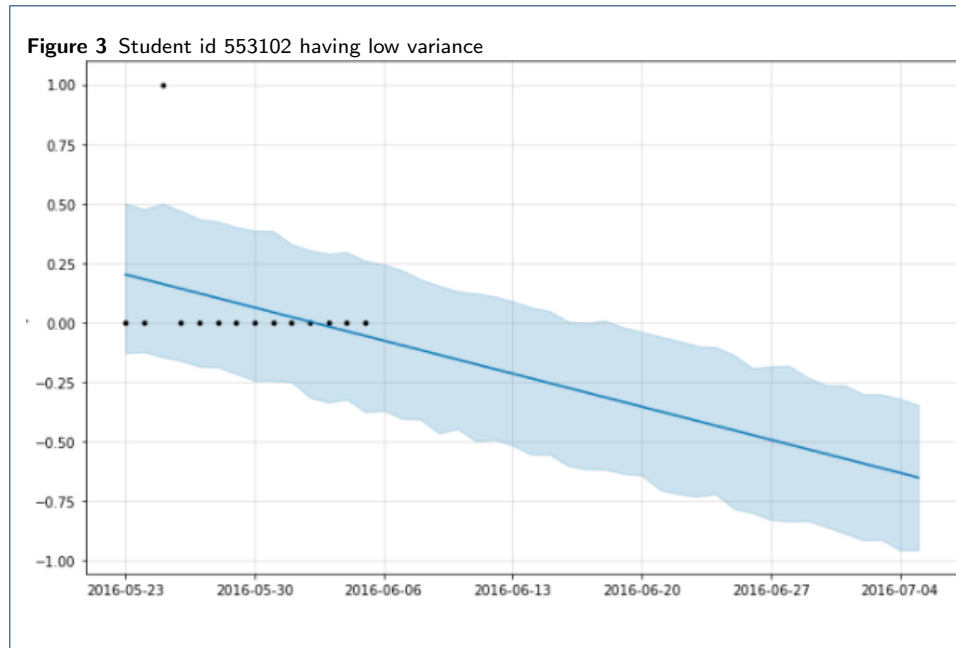
When we tried to find variances of an individual, for incoming call frequency (CALL IN FREQ), we found out it varies among individual students. Student ID 173512 shows the highest variance, where Student ID 553102 shows the least variance. That can be seen in Fig 2 and Fig 3

Figure 2 Student id 173512 having high variance



Task 2 : Which behavioral dispositions are related to personality traits?

- **Goal:** To find out the behavioural traits which are highly linked with the personality traits.



- **Method:** Using the Inter-class Correlation Coefficient, we took Daily Call Out, Daily Call In and Daily SMS In behavioural traits to co relate them with personality traits (Extra-Version).
- **Result:** To evaluate the extent to which behavioral sociability tendencies map on to standard self-reported measures of personality traits, we selected Sample no. 3 as our data-set which had 10 daily social behavior. The Spearman's Rank Correlation Coefficient is used to discover the strength of a link between two sets of data. With calculation of Spearman correlation it was possible to estimate how social behavior activities relate with personality traits (see Figure 5). Most BIG FIVE TRAITS showed low values of correlation with social behavior dispositions. Out of the three Behaviours we Call Out matches the most with the Extra-version Personality and are positively co-related amongst the three. Daily Call out (Figure : 4), Daily Call In (Figure : 5) and Daily SMS(Figure : 6). It can be seen in the Figure 4: Call out Behaviour is the most co-related with the Extra-version Personality. Most BIG FIVE TRAITS showed low values of correlation with social behavior dispositions. From the table, it can be observed that extroverts tend to engage in a longer calling behaviour [duration] during the night time ($r = 0.10$ to $r = 0.33$). A strong correlation is found during night, evening and weekend for calling frequency ($r = 0.10$ to $r = 0.28$, $r = 0.10$ to $r = 0.23$ $r = 0.12$ to $r = 0.21$) Fig 7.

Task 3 : Can we predict an individual's personality? In this task, our goal was to build a classifier to predict an individual's personality traits (extroversion, agreeableness, conscientiousness, neuroticism, and openness). So we divided this task into two sub tasks which as follows: 1. Classification task: build a classifier to predict an individual traits as categories (high,low) 2. Regression task: build a regressor to predict an individual traits as numbers. Sample 3 data set [1] was used to train and test to address both subtasks. We splitted the dataset into train and test set as 85% and 15% respectively to train and test our models. As a part of data

Figure 4 Correlation between Call Out and Personality

Daily CALL OUT DUR Extraversion ICC:[1] -0.05082625
 Lower ICC :[1] -0.4449291
 Upper ICC :[1] 0.2587485

Daily CALL OUT DUR openness ICC:[1] -0.367005
 Lower ICC :[1] -0.7814345
 Upper ICC :[1] -0.007952622

Daily CALL OUT DUR conscientiousness ICC:[1] 0.148483
 Lower ICC :[1] -0.2130452
 Upper ICC :[1] 0.4153091

Daily CALL OUT DUR agreeableness ICC:[1] 0.1350457
 Lower ICC :[1] -0.2292025
 Upper ICC :[1] 0.4050161

Daily CALL OUT DUR neuroticism ICC:[1] -0.2569622
 Lower ICC :[1] -0.6683824
 Upper ICC :[1] 0.08756804

Figure 5 Correlation between Call In and Behaviour

Daily CALL IN DUR Extraversion ICC:[1] -0.3457488
 Lower ICC :[1] -0.7677827
 Upper ICC :[1] 0.01497609

Daily CALL IN DUR openness ICC:[1] -0.2941983
 Lower ICC :[1] -0.714321
 Upper ICC :[1] 0.05939391

Daily CALL IN DUR conscientiousness ICC:[1] 0.06655132
 Lower ICC :[1] -0.3135421
 Upper ICC :[1] 0.3535099

Daily CALL IN DUR agreeableness ICC:[1] 0.1568647
 Lower ICC :[1] -0.2052632
 Upper ICC :[1] 0.4228592

Daily CALL IN DUR neuroticism ICC:[1] -0.1955795
 Lower ICC :[1] -0.6095069
 Upper ICC :[1] 0.1426439

Figure 6 Correlation between SMS In and Behaviour

```

Daily SMS IN NUM Extraversion ICC:[1] -0.02712702
Lower ICC :[1] -0.07489065
Upper ICC :[1] 0.1544323

Daily SMS IN NUM openness ICC:[1] 0.06718044
Lower ICC :[1] -0.02747302
Upper ICC :[1] 0.3545156

Daily SMS IN NUM conscientiousness ICC:[1] 0.09117272
Lower ICC :[1] -0.01460414
Upper ICC :[1] 0.3956336

Daily SMS IN NUM agreeableness ICC:[1] 0.050627
Lower ICC :[1] -0.0361532
Upper ICC :[1] 0.3241456

Daily SMS IN NUM neuroticism ICC:[1] -0.03156072
Lower ICC :[1] -0.07700371
Upper ICC :[1] 0.1431406

```

cleaning and feature preprocessing, we found one feature 'demogsex' as categorical, so we converted it into numerical feature (m:1 , f:0). We used all features(200) to train our models.

1. Classification Task:: As a part of this task, we used logistic regression scikit learn library model to learn the data. As sample 3 dataset had continuous values for personality traits. So, we converted those values into high and low categories using the threshold as zero, means values greater than zero were considered high and less than or equal to zero were considered as low. Model parameters were taken as follows: learning rate as 0.01, max iteration as 100 and error threshold as 0.001. As a part of result of this subtask, we found the best model with target trait as Neuroticism, which was providing the promising result with area under curve (auc) value as 0.74 and worst model with trait openness with auc 0.39 on test dataset. Other models are slightly better than the random model. The results can be found in table 1 and implementation screenshot can be found in figure 8. Surely a question will raise in mind why we have used the logistic regression model, why not other model. The very simple reason is that it gives the probability of each class that can be used as the score and that can be used to further interpretation of individual trait.

2. Regression Task:: In this task, we tried several models like linear regression, lasso, decision tree regressor, random forest regressor, support vector regressor, gradient boosting regressor, adaboost regressor using scikit learn library. We tried to reduce the number of features to model input with auto feature library but that

Figure 7 Correlation between Time of the week day/week social behaviour tendencies and Extraversion

CALL OUT DURATION				
Variable	r		ci	p-value
0 Morning	0.067001	(-0.1018591675536811, 0.23210818774972677)		0.436619
1 Afternoon	0.115695	(-0.05304943295583833, 0.2780163047484722)		0.178206
2 Evening	0.134760	(-0.03371785017953301, 0.2957899577752689)		0.116406
3 Night	0.330748	(0.17260760758212476, 0.4722663857180523)		0.000079
4 Weekday	0.092950	(-0.0759499255672151, 0.25666388570837767)		0.280000
5 Weekend	0.143516	(-0.024796522911973838, 0.3039160935415045)		0.094309
CALL IN DURATION				
Variable	r		ci	p-value
0 Morning	0.074287	(-0.09460727811732919, 0.23902439469617903)		0.388285
1 Afternoon	0.061641	(-0.10718255105843766, 0.22700963396192522)		0.474263
2 Evening	0.047848	(-0.12083737173275416, 0.21384721659867822)		0.578737
3 Night	0.107541	(-0.06127967677260702, 0.2703799037277814)		0.210995
4 Weekday	0.041250	(-0.12734684590027479, 0.20752941016818247)		0.632227
5 Weekend	0.080939	(-0.08797071210987562, 0.24532423406519038)		0.347100
CALL OUT FREQ				
Variable	r		ci	p-value
0 Morning	0.135354	(-0.03311322983150563, 0.29634219667753)		0.114789
1 Afternoon	0.190428	(0.02346180299179199, 0.3470591828217841)		0.025819
2 Evening	0.238631	(0.07387240662476555, 0.3907097921373484)		0.004982
3 Night	0.284499	(0.12263476184928929, 0.43161957657572037)		0.000753
4 Weekday	0.182117	(0.014854865071530972, 0.3394636151615405)		0.033177
5 Weekend	0.215116	(0.049175026521131496, 0.3695008685309908)		0.011589
CALL IN FREQ				
Variable	r		ci	p-value
0 Morning	0.107604	(-0.06121577106062771, 0.27043935960246895)		0.210725
1 Afternoon	0.114453	(-0.054304655725312234, 0.2768543491064921)		0.182943
2 Evening	0.100131	(-0.06873931285769022, 0.26342226559253556)		0.244348
3 Night	0.107776	(-0.06104239736316285, 0.27060064832475994)		0.209992
4 Weekday	0.125804	(-0.04281474387167832, 0.2874545047476547)		0.142966
5 Weekend	0.122840	(-0.045819265953581295, 0.2846904369958763)		0.152701
SMS IN FREQ				
Variable	r		ci	p-value
0 Morning	0.167129	(-0.000603155016352276, 0.32571488607706656)		0.050935
1 Afternoon	0.158854	(-0.009104145791840648, 0.318094625424333)		0.063726
2 Evening	0.218851	(0.05308409528880508, 0.37288016894311765)		0.010189
3 Night	0.088319	(-0.08059065178285614, 0.2522969426818806)		0.304761
4 Weekday	0.186112	(0.01898943662819745, 0.34311763929797506)		0.029444
5 Weekend	0.200332	(0.0337500617853631, 0.3560831442972565)		0.018916

Table 1 Classification Model Result: Logistic regression model to predict the personality traits

Target Personality Trait	AUC
neuroticism	0.74
extroversion	0.55
conscientiousness	0.53
agreeableness	0.52
openness	0.39

didn't work for us. However we found good results with random forest regressor among all these models, so here we are listing only the results of it. The results of other models can be seen in code file. We used the parameter of random forest regressor as follows: number of estimators=200, max depth=5, max features='log2'. We used the Multiple Output Regressor class of scikit learn library to train all model together with their respective targets. We evaluated the model with mean squared error and R^2 (coefficient of determination) regression score function. Best possible

Figure 8 Classification model(logistic regression) result screenshot

```

λ
> [376] ytrain1=np.where(ytrain>0,1,0)
> ytest1=np.where(ytest>0,1,0)
⌵
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import classification_report, confusion_matrix
from sklearn import metrics
for i, val in enumerate(model_list):
    model5 = LogisticRegression(C=0.01)
    model5.fit(xtrain, ytrain1[:,i])
    y_pred=model5.predict_proba(xtest)
    #cm = confusion_matrix(ytest1[:,2], y_pred)
    pred=np.where(y_pred[:,1]>0.5,1,0)
    #print(classification_report(ytest1[:,i],pred))
    fpr, tpr, thresholds = metrics.roc_curve(ytest1[:,i],z)
    print(val,'auc:',metrics.auc(fpr, tpr))

BFSI_0 auc: 0.3942307692307692
BFSI_C auc: 0.5333333333333333
BFSI_E auc: 0.5555555555555556
BFSI_A auc: 0.5181818181818182
BFSI_N auc: 0.7361111111111112

```

score of R^2 is 1.0 and it can be negative (because the model can be arbitrarily worse). A constant model that always predicts the expected value of y , disregarding the input features, would get a R^2 score of 0.0. Further information about the R^2 can be seen at Wikipedia [3]. We found the regressor model slightly better than the constant model with target traits neuroticism and conscientiousness. The results of regression models can be found in table 2 and implementation screenshot can be found in Figure 9.

Table 2 Regression Model Result: Random forest regressor model to predict the personality traits

Target Personality Trait	MSE	R^2 score
neuroticism	0.6137	0.08
extroversion	0.5476	-0.13
conscientiousness	0.4157	0.01
agreeableness	0.7901	-0.13
openness	0.5158	-0.04

Figure 9 Regression model(Random Forest Regressor)result screenshot

```

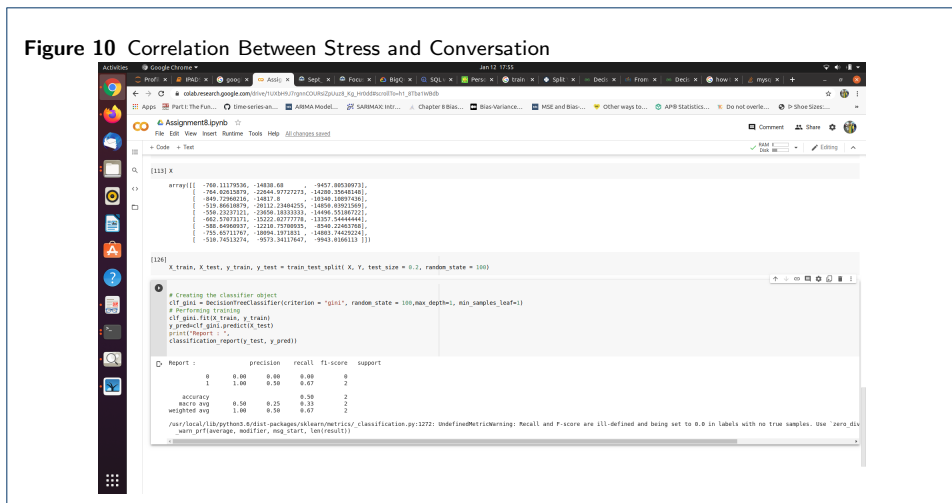
[544] ypred = model2.predict(xtest)
model_list=['BFSI_0 ','BFSI_C ','BFSI_E ','BFSI_A ','BFSI_N']
print("y1 BFSI_0 MSE:%.4f" % mean_squared_error(ytest[:,0], ypred[:,0]))
print("y2 BFSI_C MSE:%.4f" % mean_squared_error(ytest[:,1], ypred[:,1]))
print("y3 BFSI_E MSE:%.4f" % mean_squared_error(ytest[:,2], ypred[:,2]))
print("y4 BFSI_A MSE:%.4f" % mean_squared_error(ytest[:,3], ypred[:,3]))
print("y5 BFSI_N MSE:%.4f" % mean_squared_error(ytest[:,4], ypred[:,4]))

y1 BFSI_0 MSE:0.5158
y2 BFSI_C MSE:0.4157
y3 BFSI_E MSE:0.5476
y4 BFSI_A MSE:0.7901
y5 BFSI_N MSE:0.6137

from sklearn.metrics import mean_squared_error, r2_score
for i,model in enumerate(model_list):
    print(model,"R2 score : %.2f" % r2_score(ytest[:,i],ypred[:,i]))

BFSI_0 R2 score : -0.04
BFSI_C R2 score : 0.01
BFSI_E R2 score : -0.13
BFSI_A R2 score : -0.13
BFSI_N R2 score : 0.08

```

Discussion

As per these kind of analysis, we can understand how individual's daily behavior plays a role in their daily lives. We can obtain basic descriptive details on how much people tend to socialize and when the pattern which is followed. We can map daily life behavior to psychological characteristics (e.g., personality traits, attitudes, values) and life outcomes (e.g., mental health, physical health) through these kind of analysis.

Appendix

Code -

- Aman - Task1
- Suresh - Task3
- Frenny - Task2

Report -

- Aman - Scientific background, Data, Task1
- Suresh - Discussion, Task 3
- Frenny - Abstract, Goal and Task 2.

Author details

References

1. Sample 3 Dataset. <https://osf.io/8ztd6/>
2. Harari, G.M., Müller, S.R., Stachl, C., Wang, R., Wang, W., Bühner, M., Rentfrow, P.J., Campbell, A.T., Gosling, S.D.: Sensing sociability: Individual differences in young adults' conversation, calling, texting, and app use behaviors in daily life. *Journal of personality and social psychology* (2019)
3. R^2 Score Function. https://en.wikipedia.org/wiki/Coefficient_of_determination