# Bayesian Unsupervised Clustering Method For Uncovering Latent Personality Types

Boxuan Li, Nianli Peng, Danny Luo

4/24/2021

## Contents

## 1 Introduction

The Five Factor Model (FFM) of personality is a model for personality assessment that has been widely studied and applied in the field of Psychology. [1] It proposes 5 domains across which one's personality could be characterized. They are Openness to Experience, Conscientiousness, Extraversion, Agreeableness and Neuroticism (or in abbreviation, OCEAN) respectively.

While FFM presents a viable framework to evaluate individual personality's scores on those five traits, it does not identify any personality type by itself. To fully extract the value from FFM data usually means analyzing in depths the interaction between each dimensions or moving a step further in classifying individuals into homogenous personality profiles that could be interpretable under FFM.[2] Identifying those latent personality types will be of trememdous psychometric values. It will not only reveal correlations between each dimension of personality traits, but will also present us a fuller picture of compositions of human personalities. An ideal latent personality classification would also yield a simple and univariate measure of individual personality, that could be used in causal inference and prediction widely in the field of psychology and behavioral science.

Recent literatures have attempted with various techniques to approach this clustering tasks to identify personality types from FFM, including Latent Profile Analysis, Gaussian Mixture Model combined with Factor Analysis. [2][3]

We propose an Bayesian unsupervised clustering algorithm that leverages a two-fold modeling structure:

- A non-parametric Dirichlet process Gaussian mixture model to estimate size of clusters and their respective subpopulation parameters using a small portion of data
- Feed the above result as prior into Gaussian mixture model with fixed cluster size, utilizing the rest of our data

We adopt this two phase modeling due to expensive computational cost given the gigantic dataset. The final output will yield a clustering of all individuals into different latent personalities type that is highly interpretable using FFM framework.

## 2    Data

This dataset contains $1,015,342$ questionnaire answers collected through an interactive online personality test by Open Psychometrics from 2016 to 2018. The personality test was constructed with the "Big-Five Factor Markers" from the International Personality Item Pool, developed by Goldberg (1992). It consists of fifty items that the respondent must rate on how true they are about him/her on a five point scale from "Very Inaccurate," "Moderately Inaccurate," "Neither Inaccurate nor Accurate," "Moderately Accurate," and "Very Accurate." Responses to this test was recorded anonymously. More information about each question is included in the appendix.

In this study we will analyze the data set and use a Bayesian unsupervised learning algorithm for clustering the participants.

```
library(dplyr)
```

```
##
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:stats':
##
##     filter, lag
```

```
## The following objects are masked from 'package:base':
##
##     intersect, setdiff, setequal, union
```

```
library(dirichletprocess)
library(ggplot2)
library(patchwork)
```

```
data_raw = read.csv("data-final.csv", sep='\t', na.strings = "NULL")
```

```
data = data.frame(data_raw)
```

```
data <- data[ -c(51:107) ]
data <- data[ -c(52:53) ]
```

```
print(nrow(data))
```

```
## [1] 1015341
```

```
head(data)
```

```
##   EXT1 EXT2 EXT3 EXT4 EXT5 EXT6 EXT7 EXT8 EXT9 EXT10 EST1 EST2 EST3 EST4 EST5
## 1    4    1    5    2    5    1    5    2    4     1    1    4    4    2    2
## 2    3    5    3    4    3    3    2    5    1     5    2    3    4    1    3
## 3    2    3    4    4    3    2    1    3    2     5    4    4    4    2    2
## 4    2    2    2    3    4    2    2    4    1     4    3    3    3    2    3
## 5    3    3    3    3    5    3    3    5    3     4    1    5    5    3    1
## 6    3    3    4    2    4    2    2    3    3     4    3    4    3    2    2
##   EST6 EST7 EST8 EST9 EST10 AGR1 AGR2 AGR3 AGR4 AGR5 AGR6 AGR7 AGR8 AGR9 AGR10
## 1    2    2    2    3     2    2    5    2    4    2    3    2    4    3     4
## 2    1    2    1    3     1    1    4    1    5    1    5    3    4    5     3
## 3    2    2    2    1     3    1    4    1    4    2    4    1    4    4     3
```

```
## 4      2      2      2      4      3      2      4      3      4      2      4      2      4      3      4
## 5      1      1      1      3      2      1      5      1      5      1      3      1      5      5      3
## 6      1      2      1      2      2      2      3      1      4      2      3      2      3      4      4
##    CSN1 CSN2 CSN3 CSN4 CSN5 CSN6 CSN7 CSN8 CSN9 CSN10 OPN1 OPN2 OPN3 OPN4 OPN5
## 1    3    4    3    2    2    4    4    2    4     4    5    1    4    1    4
## 2    3    2    5    3    3    1    3    3    5     3    1    2    4    2    3
## 3    4    2    2    2    3    3    4    2    4     2    5    1    2    1    4
## 4    2    4    4    4    1    2    2    3    1     4    4    2    5    2    3
## 5    5    1    5    1    3    1    5    1    5     5    5    1    5    1    5
## 6    3    2    4    1    3    2    4    3    4     3    5    1    5    1    3
##    OPN6 OPN7 OPN8 OPN9 OPN10 country
## 1    1    5    3    4     5      GB
## 2    1    4    2    5     3      MY
## 3    2    5    3    4     4      GB
## 4    1    4    4    3     3      GB
## 5    1    5    3    5     5      KE
## 6    1    5    4    5     2      SE
```

```r
table(is.na(data))
```

```
##
##    FALSE     TRUE
## 51693241    89150
```

It looks like there are 89150 missing values.

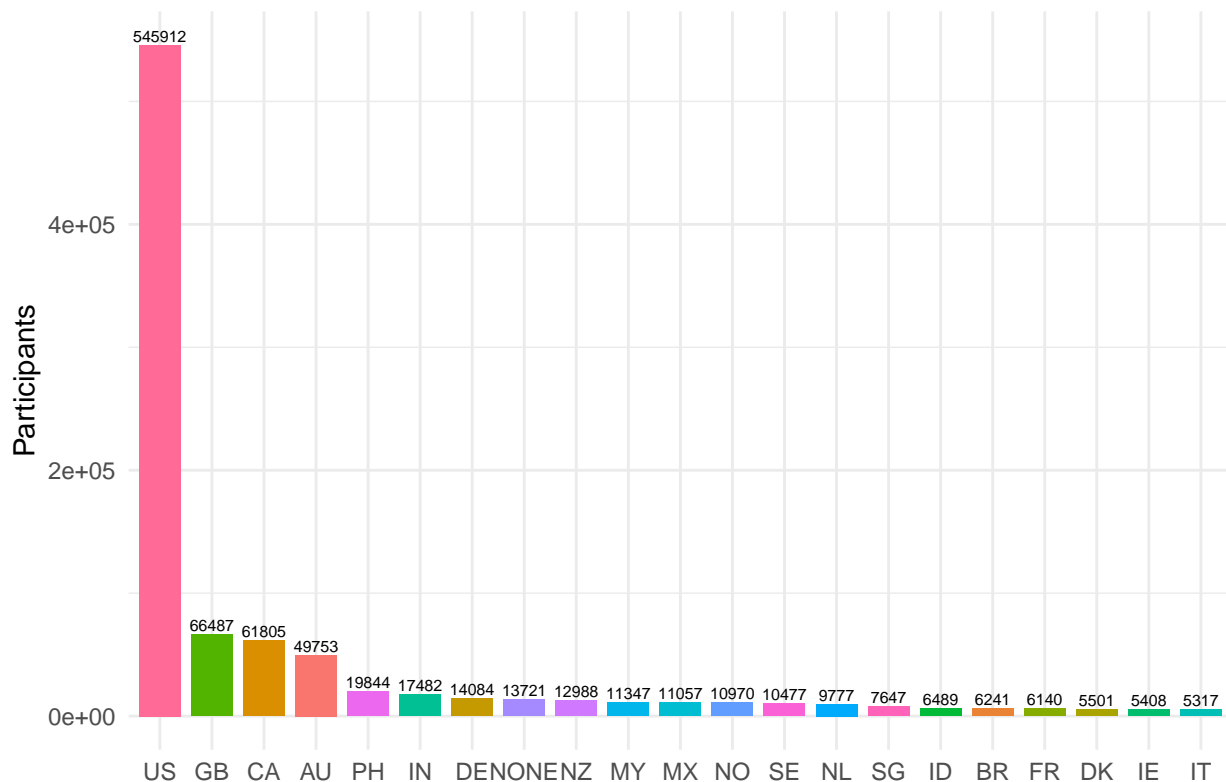```r
data <- na.omit(data)
nrow(data)
```

```
## [1] 1013558
```

After eliminating missing values, we have 1013558 valid observations.

```r
countries <- data %>% count(country)
countries <- countries[countries$n>=5000,]
```

```r
ggplot(countries, aes(reorder(country, -n, sum), n, fill = country)) +
  geom_bar(stat="identity", width = 0.8)+
  geom_text(aes(label=n), vjust=-0.3, size=2)+
  theme_minimal()+
  theme(legend.position="none")+
  labs(title= "Countries With More Than 5000 Participants",
                   y="Participants", x = element_blank())
```

## Countries With More Than 5000 Participants



We see that the vast majority of the participants are from the U.S. We might be exposed to selection bias.

Among the 50 items in the survey, some are positive (e.g. "I am the life of the party") while some are negative (e.g. "I don't talk a lot").

For + keyed items, the response "Very Inaccurate" is assigned a value of 1, "Moderately Inaccurate" a value of 2, "Neither Inaccurate nor Accurate" a 3, "Moderately Accurate" a 4, and "Very Accurate" a value of 5.

For − keyed items, the response "Very Inaccurate" is assigned a value of 5, "Moderately Inaccurate" a value of 4, "Neither Inaccurate nor Accurate" a 3, "Moderately Accurate" a 2, and "Very Accurate" a value of 1.

Once numbers are assigned for all of the items in the scale, we will sum all the values to obtain a total scale score for each of the five personality traits.

```
pos_keyed_vars <-  c('EXT1', 'EXT3', 'EXT5', 'EXT7', 'EXT9',
                     'EST1', 'EST3', 'EST5', 'EST6', 'EST7',
                     'EST8', 'EST9', 'EST10',
                     'AGR2', 'AGR4', 'AGR6', 'AGR8', 'AGR9', 'AGR10',
                     'CSN1', 'CSN3', 'CSN5', 'CSN7', 'CSN9', 'CSN10',
                     'OPN1', 'OPN3', 'OPN5', 'OPN7', 'OPN8', 'OPN9',
                     'OPN10')
neg_keyed_vars <-  c('EXT2', 'EXT4', 'EXT6', 'EXT8', 'EXT10',
                     'EST2', 'EST4',
                     'AGR1', 'AGR3', 'AGR5', 'AGR7',
                     'CSN2', 'CSN4', 'CSN6', 'CSN8',
                     'OPN2', 'OPN4', 'OPN6')
```

```
for(key in neg_keyed_vars){
  data[key]=6-data[key]
}
```
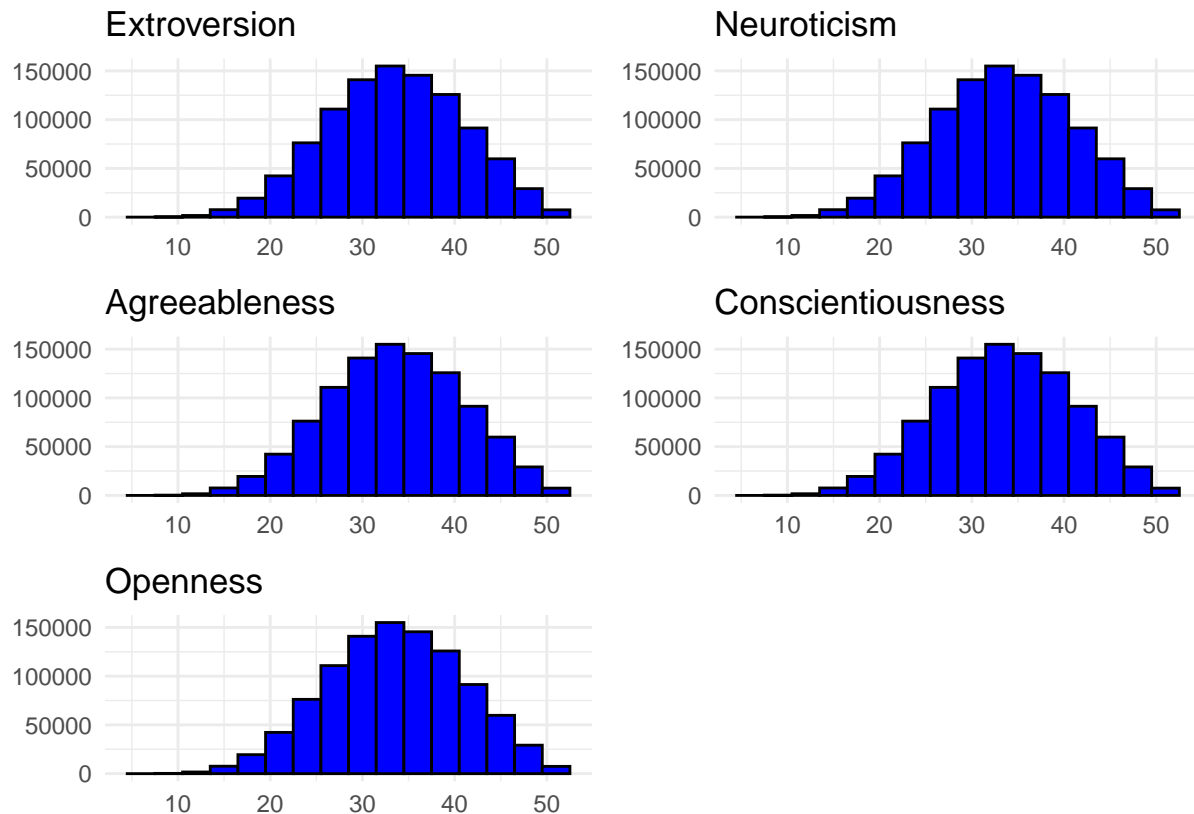
```
data <- data %>% mutate(
  EST=rowSums(data[11:20]),
  EXT=rowSums(data[1:10]),
  OPN=rowSums(data[41:50]),
  AGR=rowSums(data[21:30]),
  CSN=rowSums(data[31:40]))
score_data_final <- data[,52:56]

traits = c('EXT', 'EST', 'AGR', 'CSN', 'OPN')
trait_labels = c('Extroversion', 'Neuroticism', 'Agreeableness', 'Conscientiousness', 'Openness')
myplots <- list()
for(i in 1:5){
  p1 <- ggplot(score_data_final, aes(x= score_data_final[,i] )) +
    geom_histogram(colour="black", fill="blue",binwidth=3)+
    theme_minimal()+
    theme(axis.title.x=element_blank(),
        axis.ticks.x=element_blank(),
        axis.title.y=element_blank())+
    labs(title= trait_labels[i])
  myplots[[i]] <- p1
}

myplots[[1]]+myplots[[2]]+myplots[[3]]+myplots[[4]]+
  myplots[[5]]+plot_layout(nrow = 3)
```



**MAY wanna put this in Appendix**

The following items were presented on one page and each was rated on a five point scale using radio buttons. The order on page was was EXT1, AGR1, CSN1, EST1, OPN1, EXT2, etc. The scale was labeled 1 =

Disagree, 3 = Neutral, 5 =Agree

- EXT1 I am the life of the party.
- EXT2 I don't talk a lot.
- EXT3 I feel comfortable around people.
- EXT4 I keep in the background.
- EXT5 I start conversations.
- EXT6 I have little to say.
- EXT7 I talk to a lot of different people at parties.
- EXT8 I don't like to draw attention to myself.
- EXT9 I don't mind being the center of attention.
- EXT10 I am quiet around strangers.
- EST1 I get stressed out easily.
- EST2 I am relaxed most of the time.
- EST3 I worry about things.
- EST4 I seldom feel blue.
- EST5 I am easily disturbed.
- EST6 I get upset easily.
- EST7 I change my mood a lot.
- EST8 I have frequent mood swings.
- EST9 I get irritated easily.
- EST10 I often feel blue.
- AGR1 I feel little concern for others.
- AGR2 I am interested in people.
- AGR3 I insult people.
- AGR4 I sympathize with others' feelings.
- AGR5 I am not interested in other people's problems.
- AGR6 I have a soft heart.
- AGR7 I am not really interested in others.
- AGR8 I take time out for others.
- AGR9 I feel others' emotions.
- AGR10 I make people feel at ease.
- CSN1 I am always prepared.
- CSN2 I leave my belongings around.
- CSN3 I pay attention to details.
- CSN4 I make a mess of things.
- CSN5 I get chores done right away.
- CSN6 I often forget to put things back in their proper place.
- CSN7 I like order.
- CSN8 I shirk my duties.
- CSN9 I follow a schedule.
- CSN10 I am exacting in my work.
- OPN1 I have a rich vocabulary.
- OPN2 I have difficulty understanding abstract ideas.
- OPN3 I have a vivid imagination.
- OPN4 I am not interested in abstract ideas.
- OPN5 I have excellent ideas.
- OPN6 I do not have a good imagination.
- OPN7 I am quick to understand things.
- OPN8 I use difficult words.
- OPN9 I spend time reflecting on things.
- OPN10 I am full of ideas.

The time spent on each question is also recorded in milliseconds. These are the variables ending in _E. This was calculated by taking the time when the button for the question was clicked minus the time of the most

recent other button click.

- dateload The timestamp when the survey was started.
- screenw The width the of user's screen in pixels
- screenh The height of the user's screen in pixels
- introelapse The time in seconds spent on the landing / intro page
- testelapse The time in seconds spent on the page with the survey questions
- endelapse The time in seconds spent on the finalization page (where the user was asked to indicate if they has answered accurately and their answers could be stored and used for research. Again: this dataset only includes users who answered "Yes" to this question, users were free to answer no and could still view their results either way)
- IPC The number of records from the user's IP address in the dataset. For max cleanliness, only use records where this value is 1. High values can be because of shared networks (e.g. entire universities) or multiple submissions
- country The country, determined by technical information (NOT ASKED AS A QUESTION)
- lat_appx_lots_of_err approximate latitude of user. determined by technical information, THIS IS NOT VERY ACCURATE.
- long_appx_lots_of_err approximate longitude of user
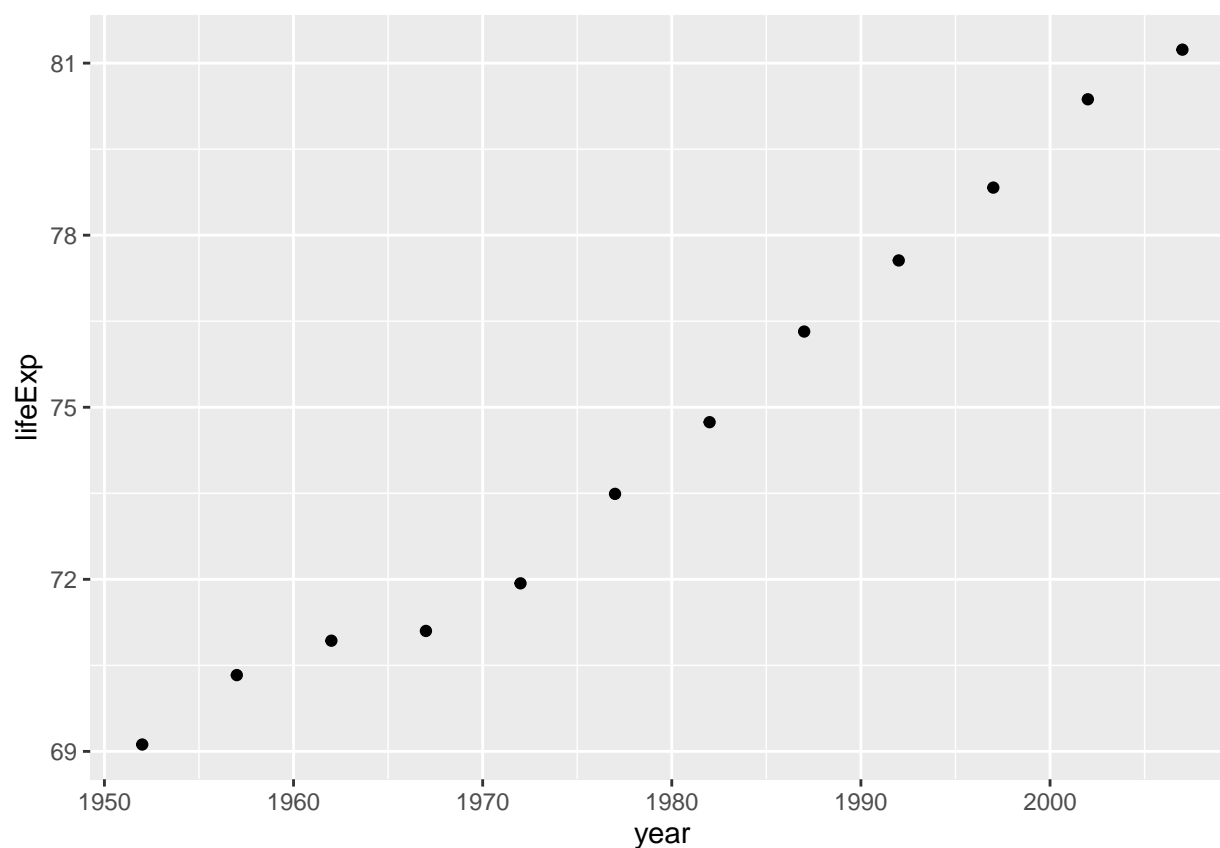


Figure 1: Life expectancy from 1952 - 2007 for Australia. Life expentancy increases steadily except from 1962 to 1969. We can safely say that our life expectancy is higher than it has ever been!

# 3 Model

$$f(k) = \binom{n}{k} p^k (1-p)^{n-k}$$

$$f(k) = \binom{n}{k} p^k (1-p)^{n-k} \tag{1}$$

: Discussion & justification of the proposed Bayesian model framework (prior and sampling model). This discussion should elicit prior information on the problem, the data sources available, and relevant project goals. Any "downstream" uses of the model (e.g., for prediction, optimization, ranking) should be discussed in detail here. See project rubric for details.

This is test for referencing equation. See equation (1)

## 4   Results

Posterior analyses from the fitted Bayesian model, and a translation of such findings into meaningful & understandable conclusions for the target audience (e.g., engineers, business managers, policy-makers, etc). See project rubric for details.

## 5   Conclusion

A summary of key findings and potential impacts of your project.

## References

[1]     N. Taylor and G. P. de Bruin, "The basic traits inventory," in *Psychological assessment in south africa: Research and applications*, Wits University Press, 2013, pp. 232–243.

[2]     E. L. Merz and S. C. Roesch, "A latent profile analysis of the five factor model of personality: Modeling trait interactions," *Pers Individ Dif*, vol. 51, no. 8, pp. 915–919, Dec. 2011, doi: 10.1016/j.paid.2011.07.022.

[3]     M. Gerlach, B. Farb, W. Revelle, and L. A. Nunes Amaral, "A robust data-driven approach identifies four personality types across four large data sets," *Nat Hum Behav*, vol. 2, no. 10, pp. 735–742, Oct. 2018, doi: 10.1038/s41562-018-0419-z.

[4]     M. A. Hayat, J. Wu, and Y. Cao, "Unsupervised bayesian learning for rice panicle segmentation with UAV images," *Plant Methods*, vol. 16, no. 1, p. 18, Feb. 2020, doi: 10.1186/s13007-020-00567-8.

[5]     "Big five personality test  kaggle." https://www.kaggle.com/tunguz/big-five-personality-test (accessed Apr. 24, 2021).

[6]     "Big five personality test. Open-source psychometrics project," Aug. 02, 2019. https://openpsychome trics.org/tests/IPIP-BFFM/.

[7]     G. J. Ross and D. Markwick, "Dirichletprocess: An r package for fitting complex bayesian nonparametric models," p. 42.

[8]     D. Görür and C. Edward Rasmussen, "Dirichlet process gaussian mixture models: Choice of the base distribution," *J. Comput. Sci. Technol.*, vol. 25, no. 4, pp. 653–664, Jul. 2010, doi: 10.1007/s11390-010-9355-8.

## Appendix