

Young People Survey

Explore the preferences, interests, habits, opinions, and fears
of young people



計財所 碩二 108071601 賴冠維



Data Introduction



In 2013, students of the Statistics class at FSEV UK were asked to invite their friends to participate in this survey.

- Music preferences (19 items)
- Movie preferences (12 items)
- Hobbies & interests (32 items)
- Phobias (10 items)
- Health habits (3 items)
- Personality traits (57 items)
- Spending habits (7 items)
- Demographics (10 items)

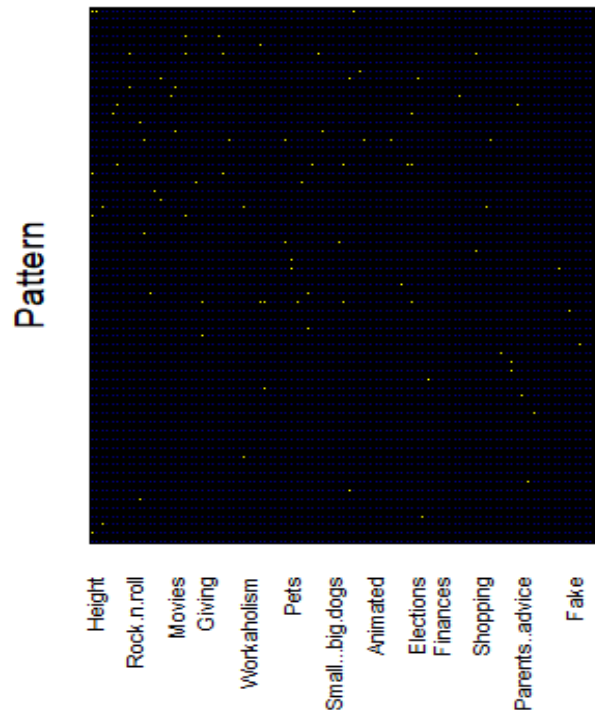
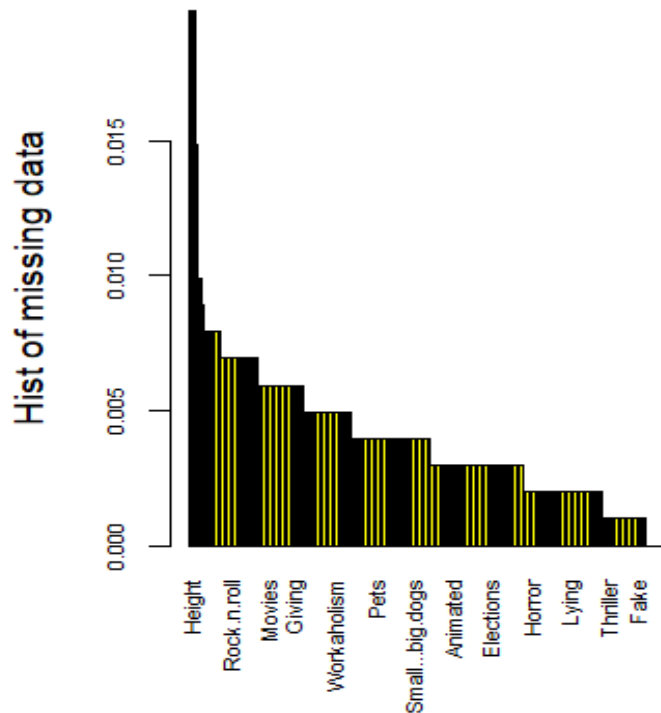
Outline



- ① **Clustering** Describe the composition of the participants
- ② **Classification** Use Xgboost to classify someone whether have alcohol addiction with other features.
- ③ **Relationship** Use factor analysis to find out the relationship among the features that be selected in former.
- ④ **GMM** Weighted the answer and classification by GMM

Missing Value Imputation

Use CART (Classification And Regression Tree) to impute the missing value



Hierarchical Clustering



Describe the composition of the participants

- Ward's minimum variance method :

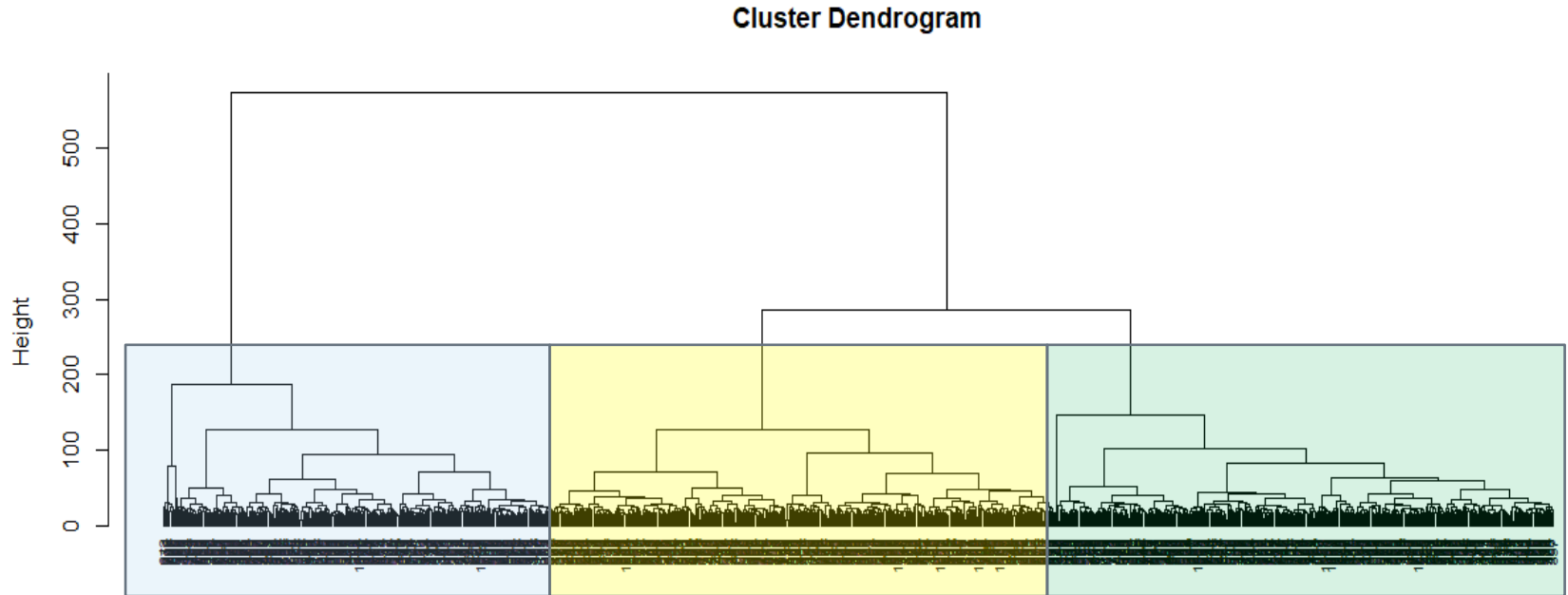
$$Total ESS = ESS_1 + ESS_2 + \cdots + ESS_k$$

$$ESS_k = \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^T (x_{ij} - \bar{x}_i)$$

- x_{ij} : j^{th} number of component in i^{th} cluster
- \bar{x}_i : Mean of the i^{th} cluster

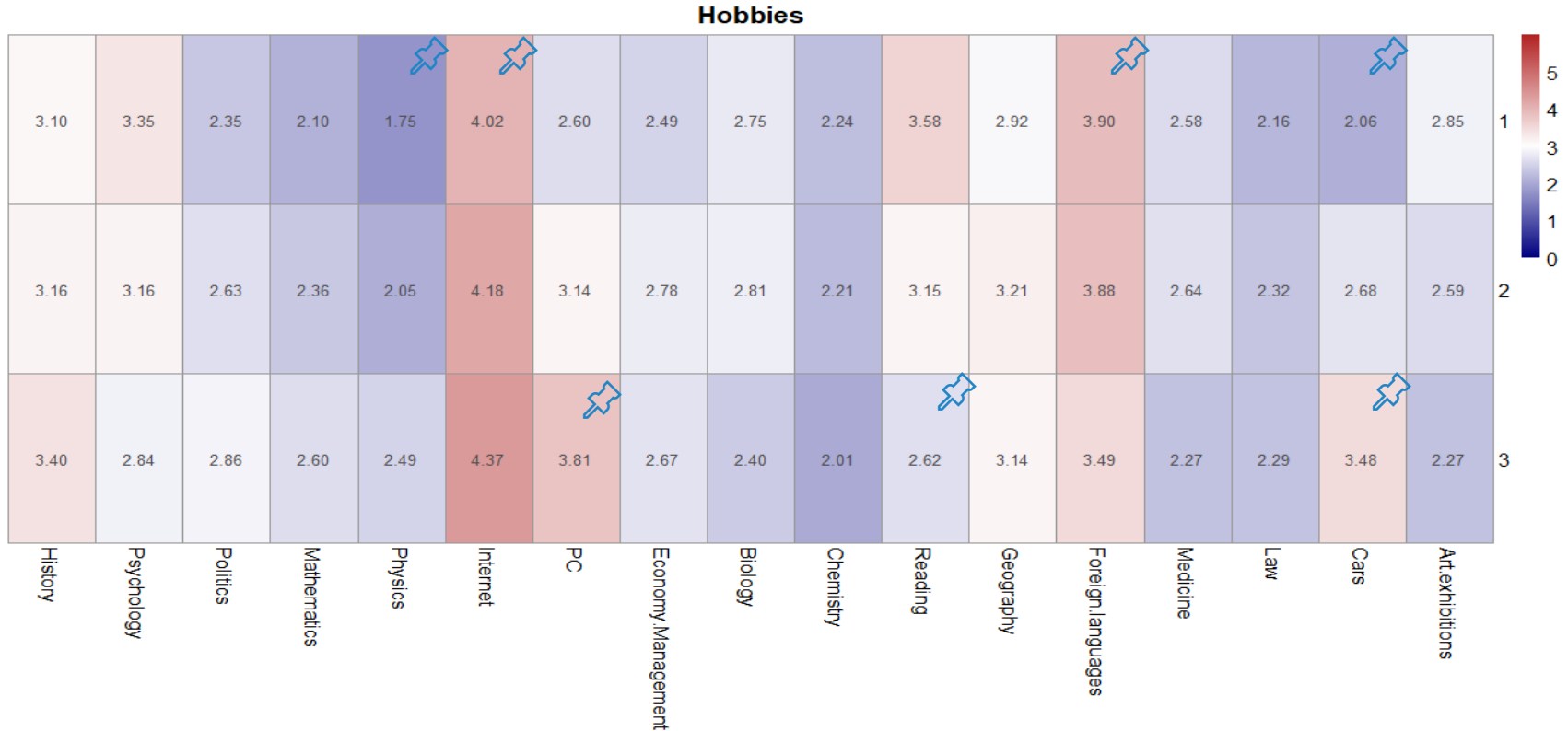
Clustering

Describe the composition of the participants



Clustering

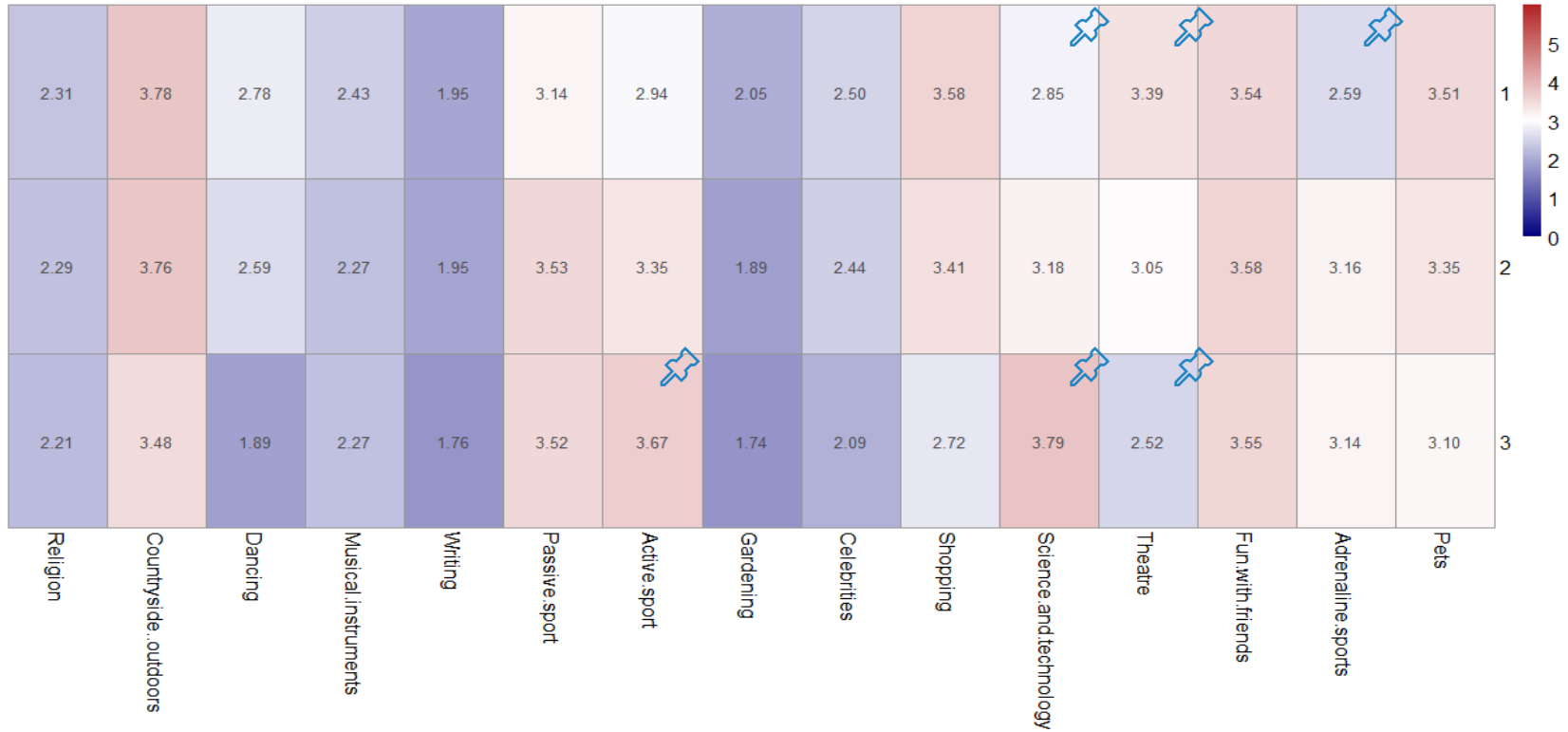
Describe the composition of the participants



Clustering

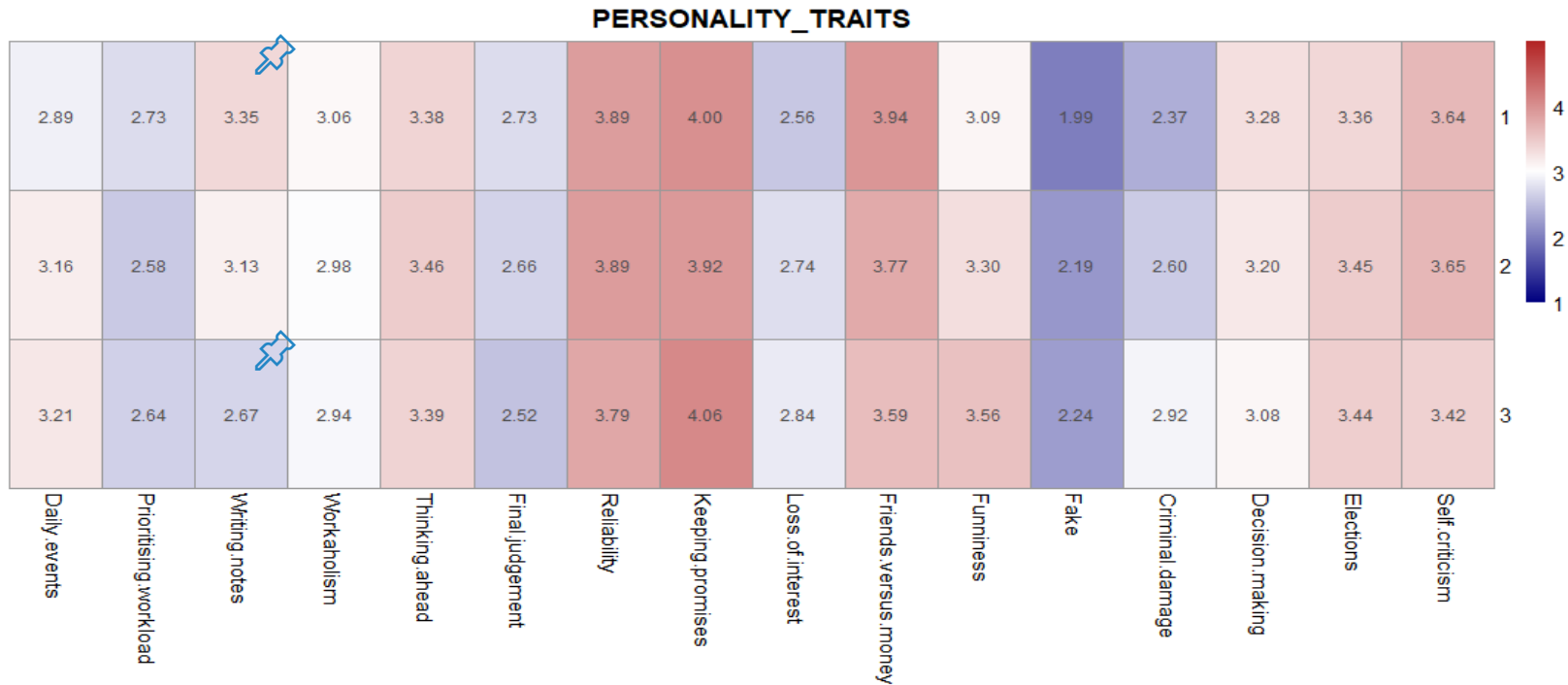
Describe the composition of the participants

Hobbies



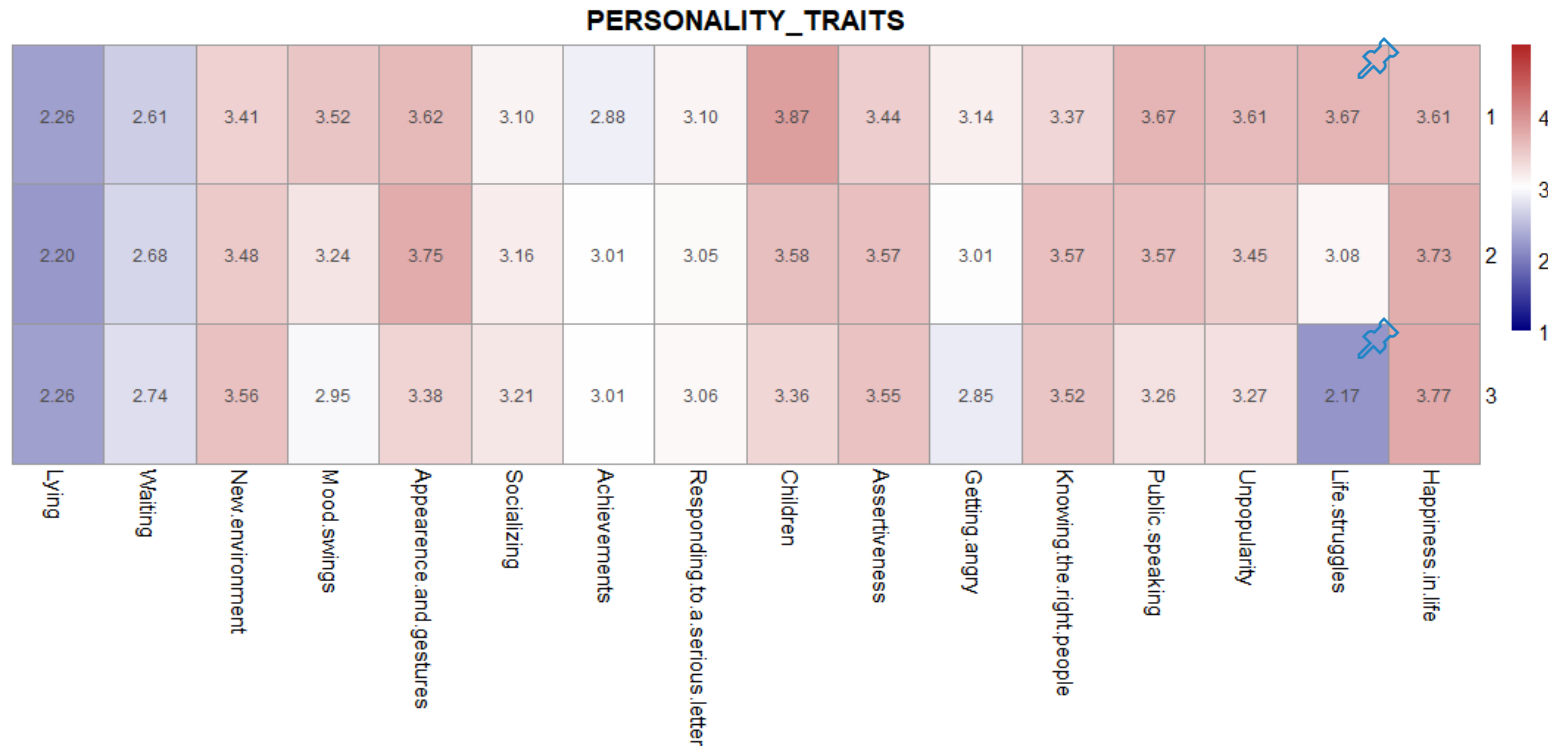
Clustering

Describe the composition of the participants



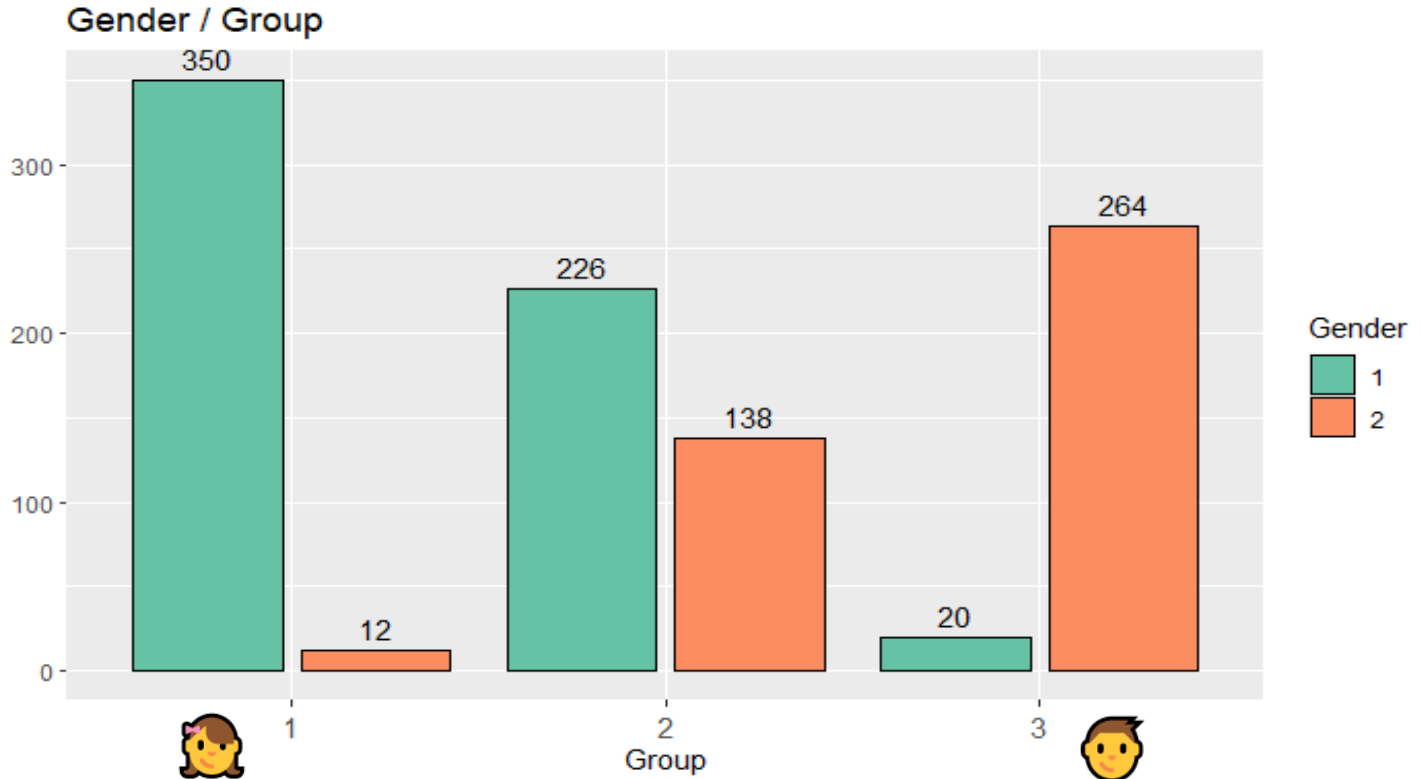
Clustering

Describe the composition of the participants



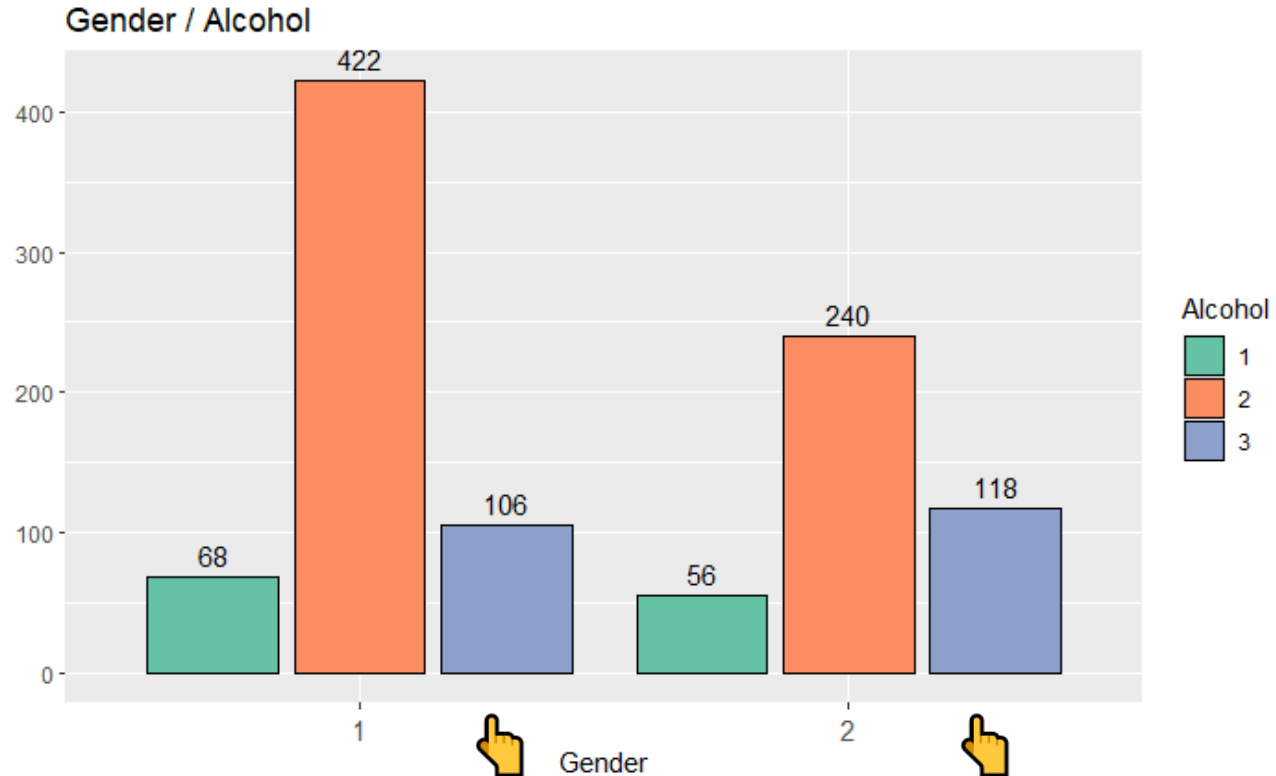
Clustering

Describe the composition of the participants



Predict

Predict someone have alcohol addiction with their other features



XGBoost

假設其分配為 Bernouli 分配 如下，其 Logistic 的損失函數：

$$\prod_1^n \pi(x_i)^{y_i} [1 - \pi(x_i)]^{1-y_i}, \pi(x) = \frac{\exp(\beta_i X_i)}{1 + \exp(\beta_i X_i)}$$

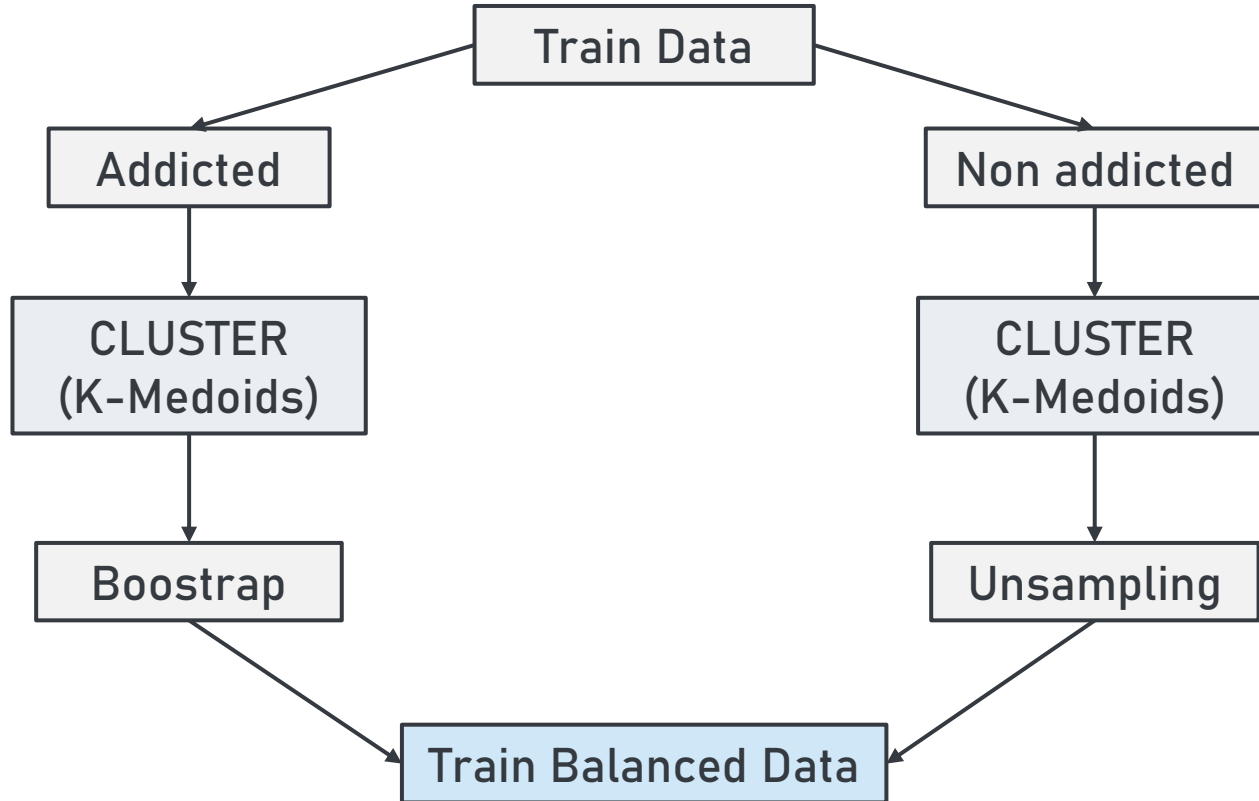
Boosting 中， $h_m(x_i)$ 代表經過 m 棵樹迭代後的估計值，同 Logistic 裡的 $\sum_{i=1}^n \beta_i X_i$

$$\hat{f}(x_i) = \sum_{m=1}^M h_m(x_i) = \sum_{i=1}^n \beta_i X_i$$

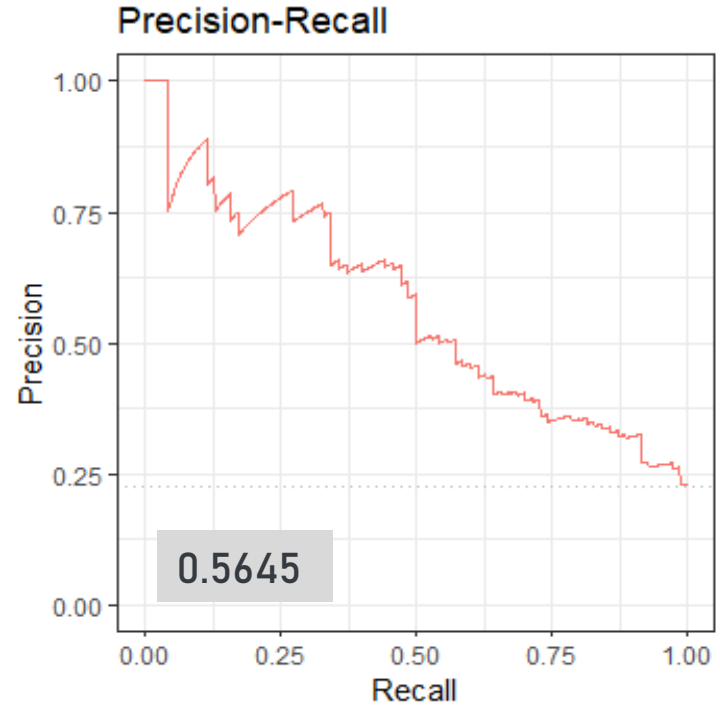
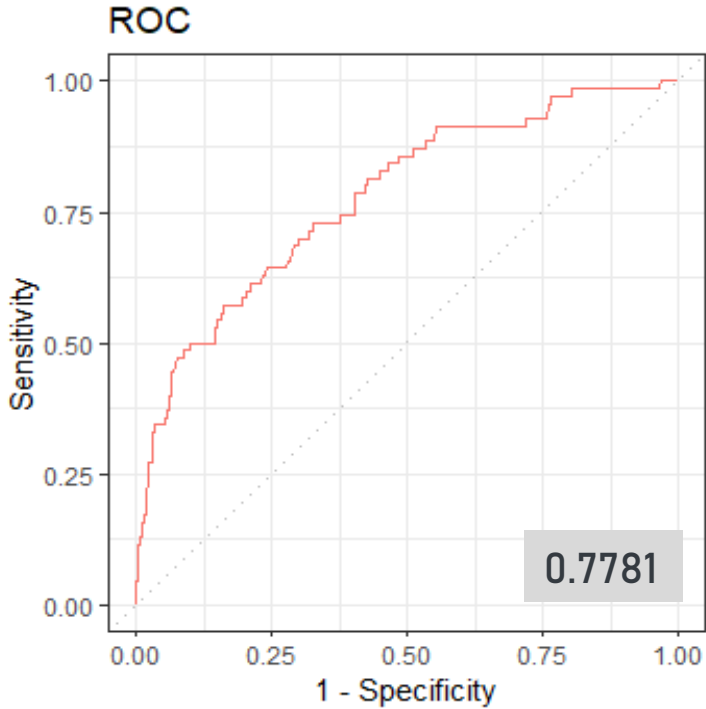
將其帶入上述損失函數後，可得下式：

$$L(y_i, f(x)) = y \ln(1 + e^{-f(x)}) + (1 - y) \ln(1 + e^{f(x)})$$

Rebalance

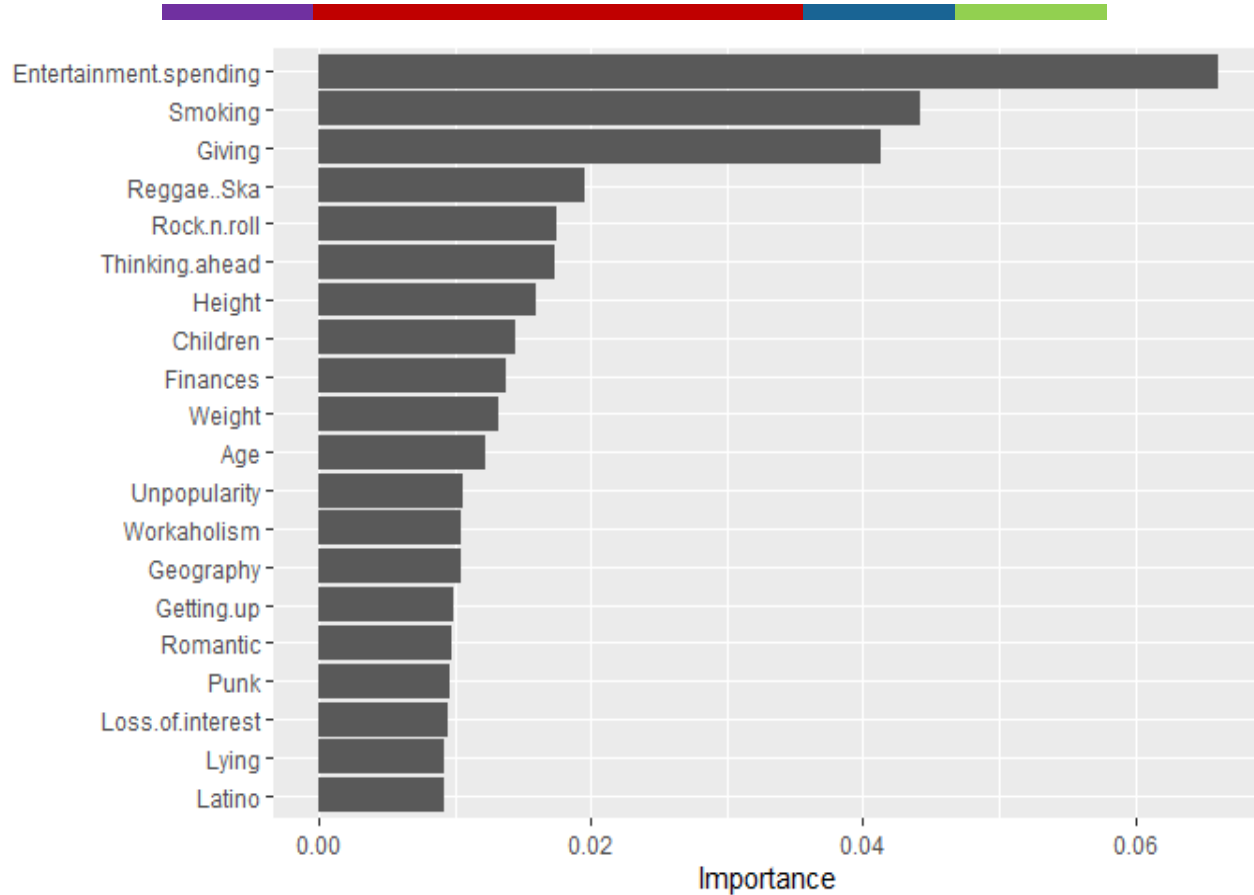


ROC / Recall



Test Set

IMPORTANCE PLOT



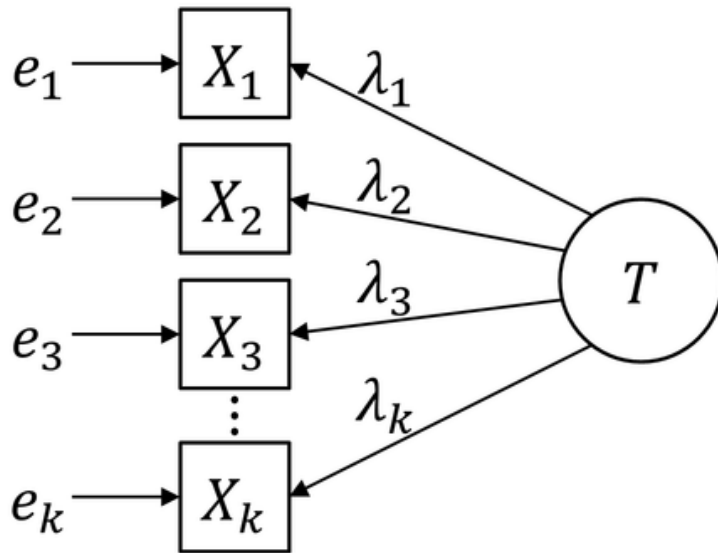
Relationship

Find out the relationship among the features that be selected.

Factor Analysis :

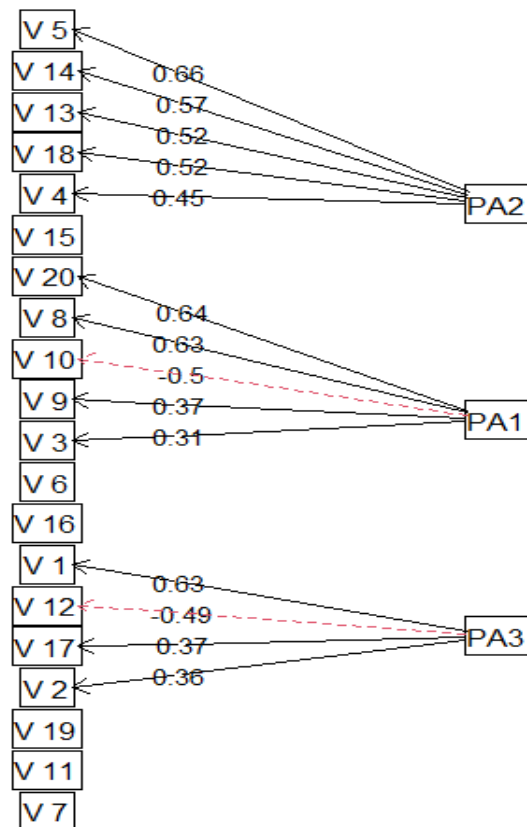
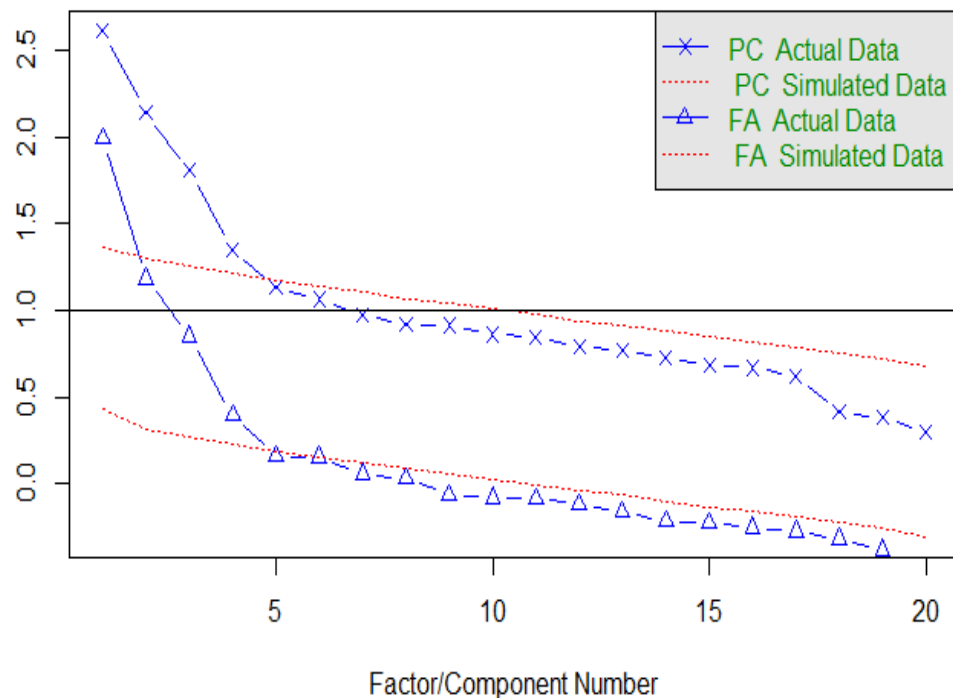
$$\mathbf{x} = \boldsymbol{\mu} + \mathbf{F}\mathbf{z} + \boldsymbol{\epsilon}$$

- x : Random Variables
- μ : Expection of x
- F : Factor Loading
- z : Hidden Facotr
- ϵ : Idiosyncratic factor

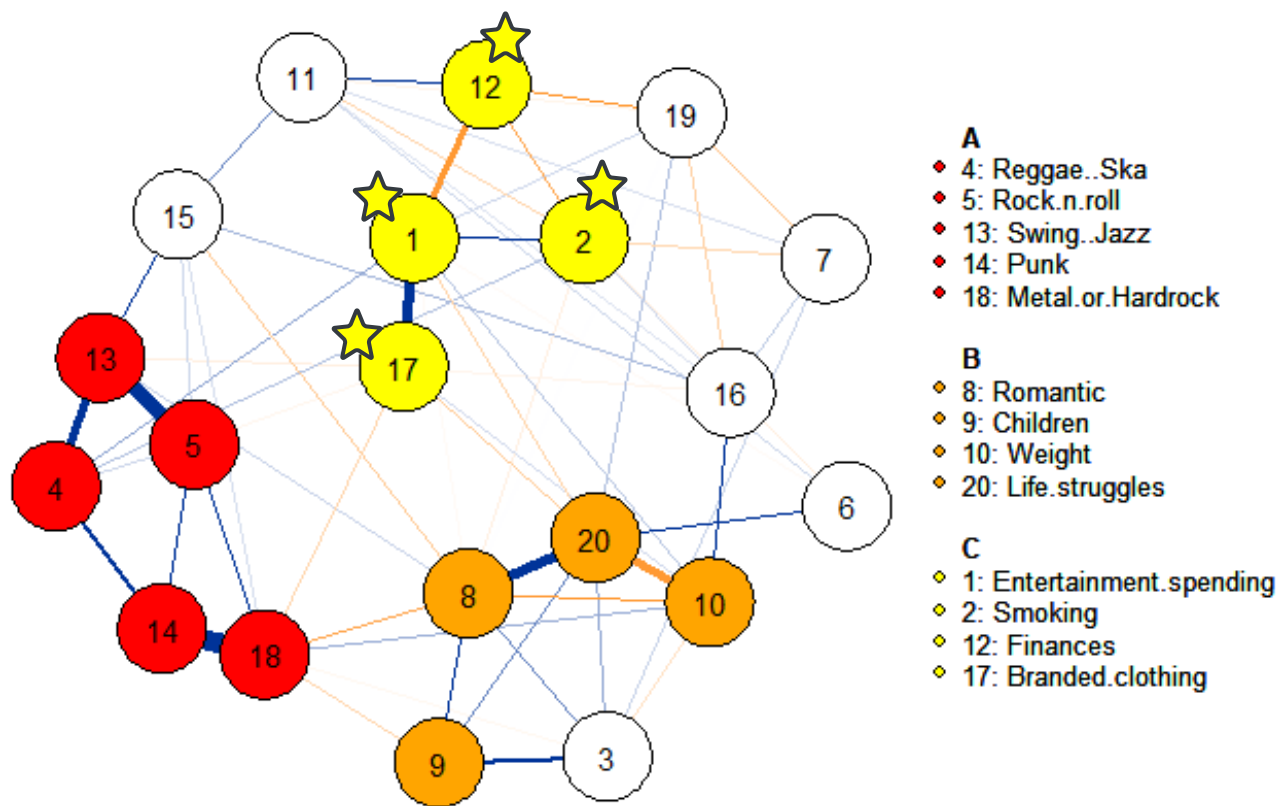


Factor Analysis

Parallel Analysis Scree Plots

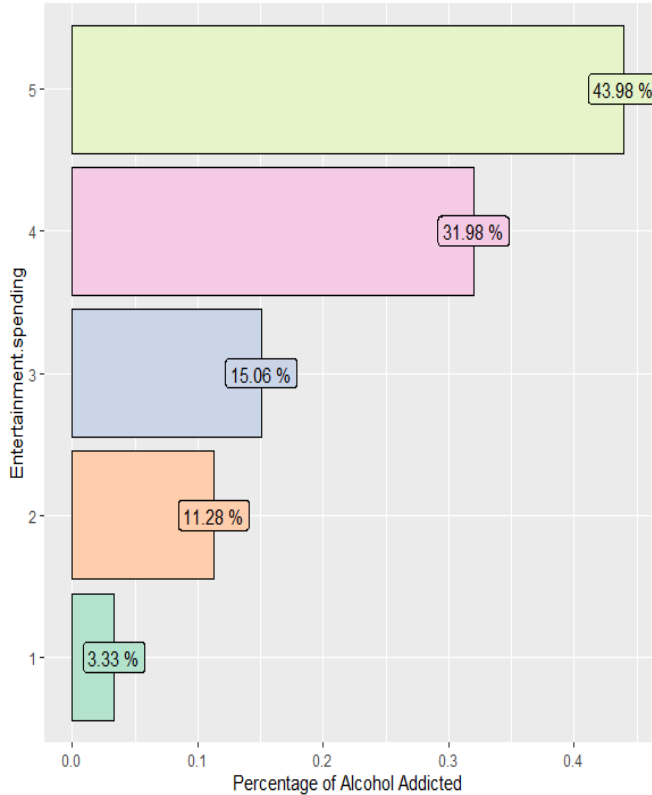


Factor Analysis

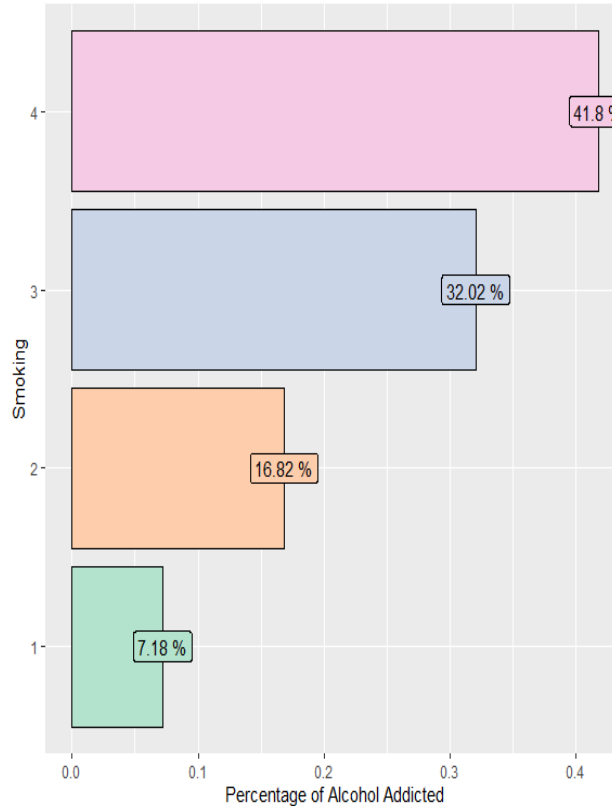


Visualization

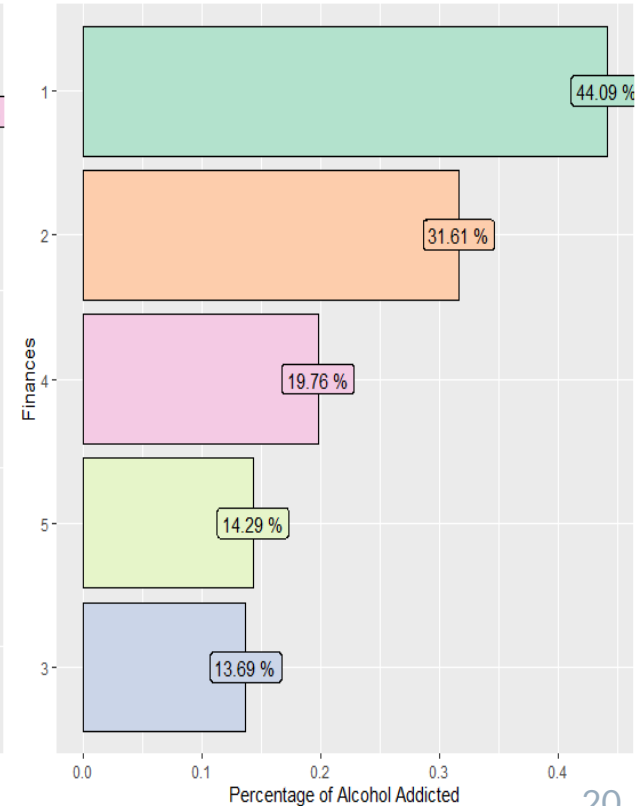
Entertainment.spending Levels v.s. Alcohol Addicted



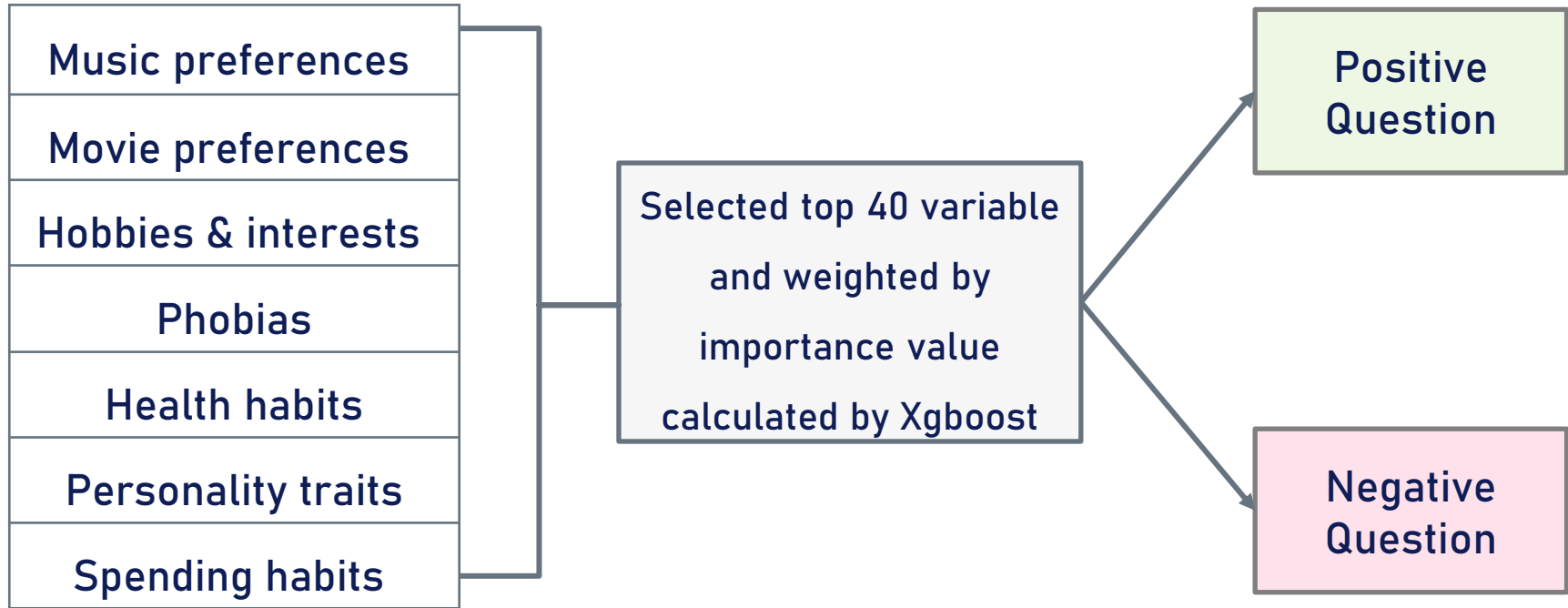
Smoking Levels v.s. Alcohol Addicted



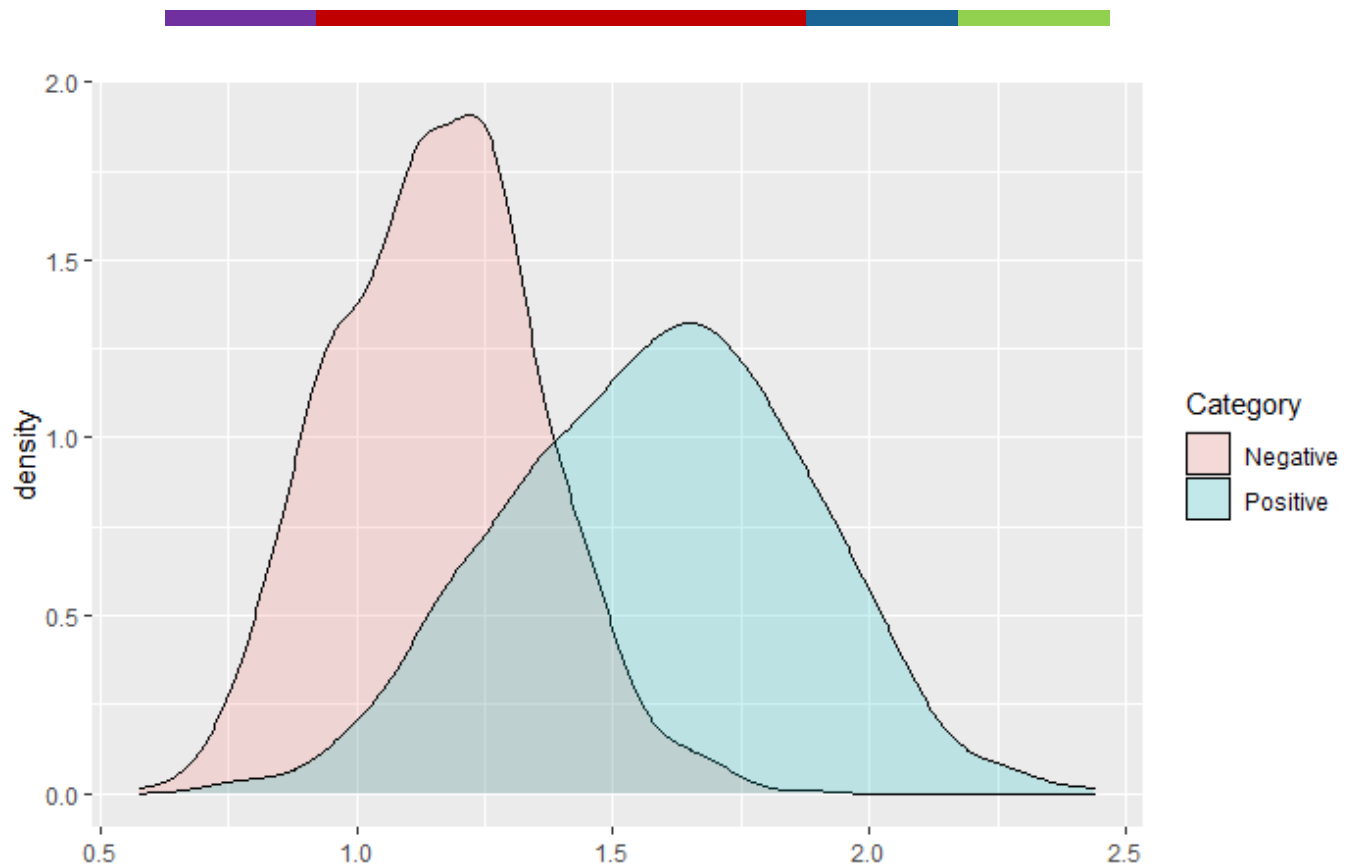
Finances Levels v.s. Alcohol Addicted



GMM



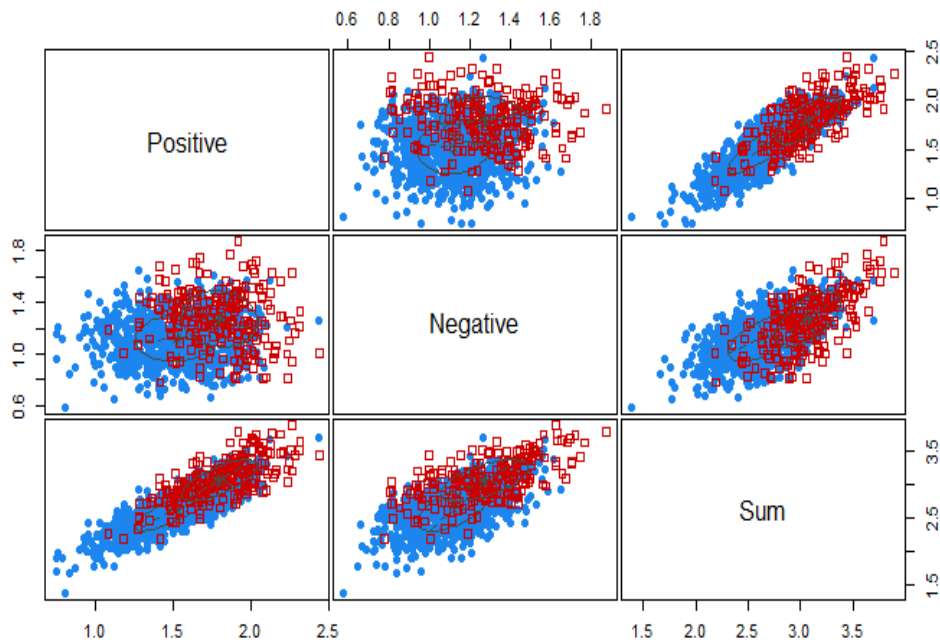
GMM



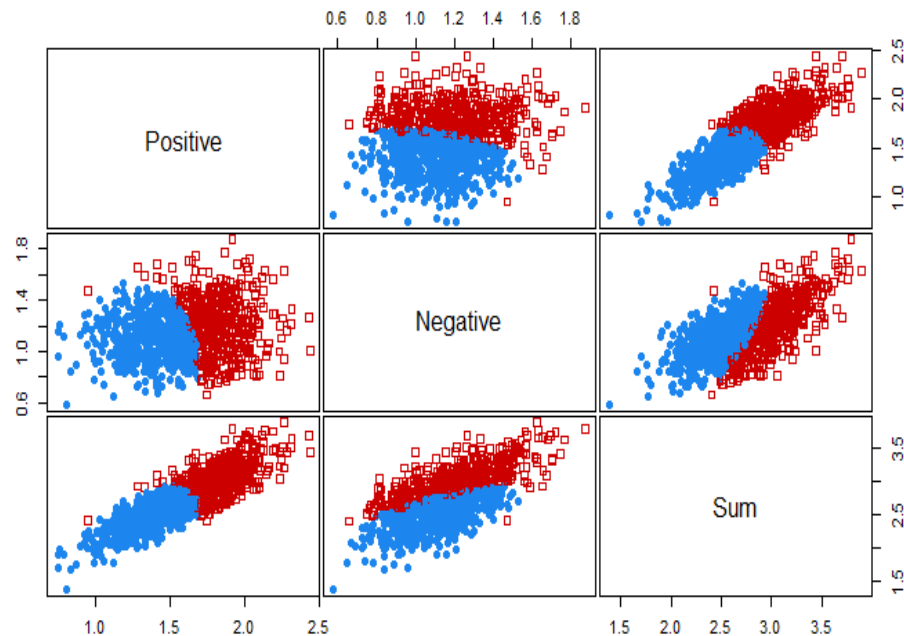
GMM



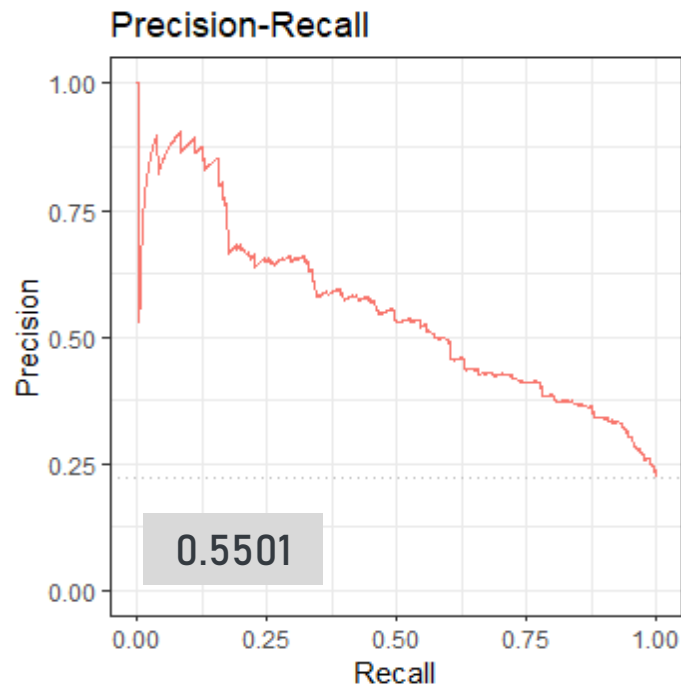
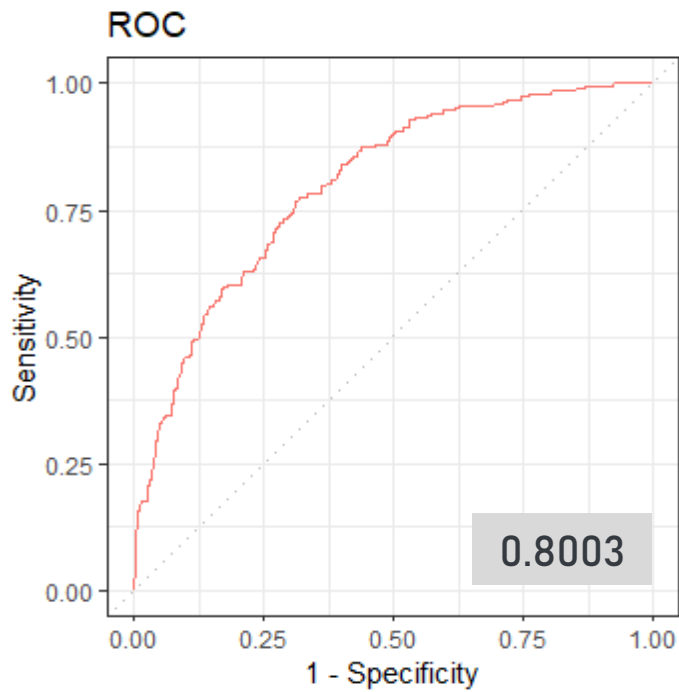
True Value Scatter Plot



Classification



GMM



Train Set