# Big 5 personality traits

Alexia

# Data Exploration

```
> summary(data)
   X.AUTHID              STATUS                 X.1              X.2              X.3
 Length:9916         Length:9916         Mode:logical     Mode:logical     Mode:logical     Mode:logical
 Class :character    Class :character    NA's:9916        NA's:9916        NA's:9916        NA's:9916
 Mode  :character    Mode  :character


   X.4                 DATE              NETWORKSIZE        BETWEENNESS         NBETWEENNESS       DENSITY
 Mode:logical        Length:9916        Min.   :  24.0     Min.   :    185.7   Min.   :31.21     Min.   :0.00000
 NA's:9916           Class :character   1st Qu.: 196.0     1st Qu.:  16902.2   1st Qu.:93.77     1st Qu.:0.01000
                     Mode  :character   Median : 317.0     Median :  47166.9   Median :96.44     Median :0.02000
                                        Mean   : 426.4     Mean   : 135439.0   Mean   :94.67     Mean   :0.03029
                                        3rd Qu.: 633.0     3rd Qu.: 196606.0   3rd Qu.:97.88     3rd Qu.:0.03000
                                        Max.   :1596.0     Max.   :1251780.0   Max.   :99.82     Max.   :0.40000

   BROKERAGE            NBROKERAGE         TRANSITIVITY           cEXT                cNEU
 Min.   :    241     Min.   :0.32       Min.   :0.0000      Length:9916         Length:9916
 1st Qu.:  17982     1st Qu.:0.49       1st Qu.:0.0600      Class :character    Class :character
 Median :  48683     Median :0.49       Median :0.0900      Mode  :character    Mode  :character
 Mean   : 137656     Mean   :0.49       Mean   :0.1288
 3rd Qu.: 198186     3rd Qu.:0.50       3rd Qu.:0.1700
 Max.   :1263790     Max.   :0.50       Max.   :0.6300
   cAGR                cCON               cOPN
 Length:9916         Length:9916        Length:9916
 Class :character    Class :character   Class :character
 Mode  :character    Mode  :character   Mode  :character
```

| STATUS |
| --- |
| likes the sound of thunder. |
| is so sleepy it's not even funny that's she can't get to … |
| is sore and wants the knot of muscles at the base of … |
| likes how the day sounds in this new song. |
| is home. <3 |
| www.thejokerblogs.com |

**Task: to cluster the 250 users and explore the relationship between Big 5 personality and each cluster**

# Data Wrangling

| STATUS |
| --- |
| likes the sound of thunder. |
| is so sleepy it's not even funny that's she can't get to … |
| is sore and wants the knot of muscles at the base of … |
| likes how the day sounds in this new song. |
| is home. <3 |
| www.thejokerblogs.com |

- num_characters: The number of characters in each post.
- num_punc: The number of punctuation marks used in each post (including !, ~, and #).
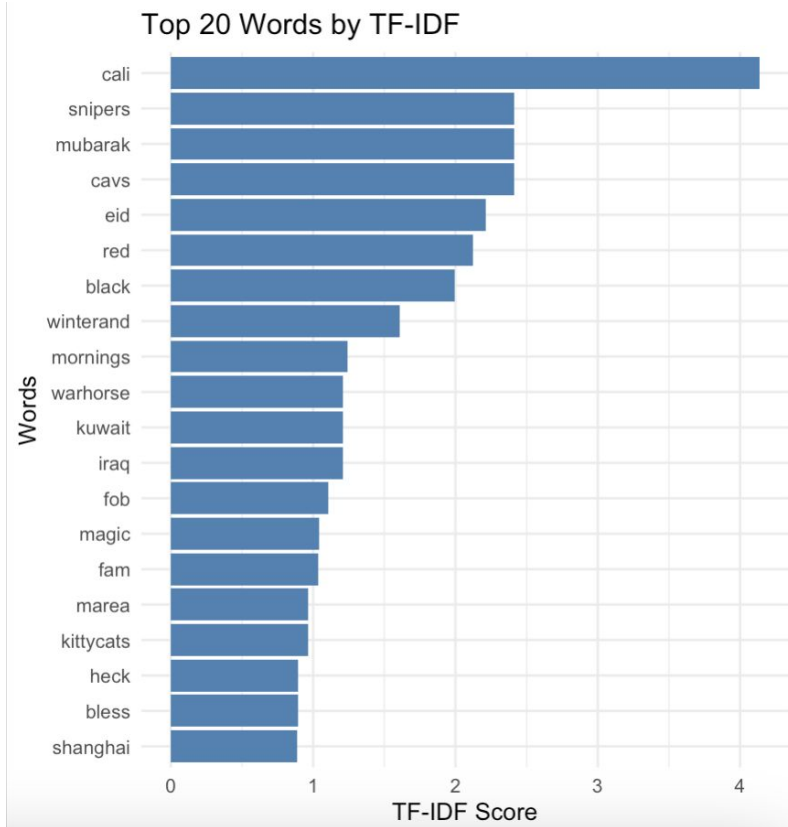
# Data Wrangling

| STATUS |
|---|
| likes the sound of thunder. |
| is so sleepy it's not even funny that's she can't get to … |
| is sore and wants the knot of muscles at the base of … |
| likes how the day sounds in this new song. |
| is home. <3 |
| www.thejokerblogs.com |

lowercase,
remove punc,
remove stop words,
tokenise

| word |
|---|
| likes |
| sound |
| thunder |
| sleepy |
| funny |
| sleep |
| sore |

# Data Wrangling– TF-IDF


Top 20 Words by TF-IDF

Term Frequency Inverse Document Frequency:

How important a word is to a document in a collection (or corpus) of documents

# Data Wrangling – N-grams

# Data Wrangling – Sentiment Analysis

TheMysMan_bh@d4jk:

"My friend Kanye **punched** me because I cheated on him with his girlfriend Kim, but I still feel **happy** today for Kim's **kiss** ! xoxo!!!"

Positive ratio: 3/(1+3)

# Data Wrangling

💡

```
$ cOPN           : chr  "y" "y" "y"
$ num_punc       : int  0 0 1 0 0 0
$ num_characters : int  26 59 116 4
$ mean_tf_idf    : num  0.00634 0.0(
$ top_word       : chr  "3" "3" "3"
$ positive_ratio : num  0.637 0.637
```

# Data Wrangling

RF for feature importance (eg. **cEXT**)

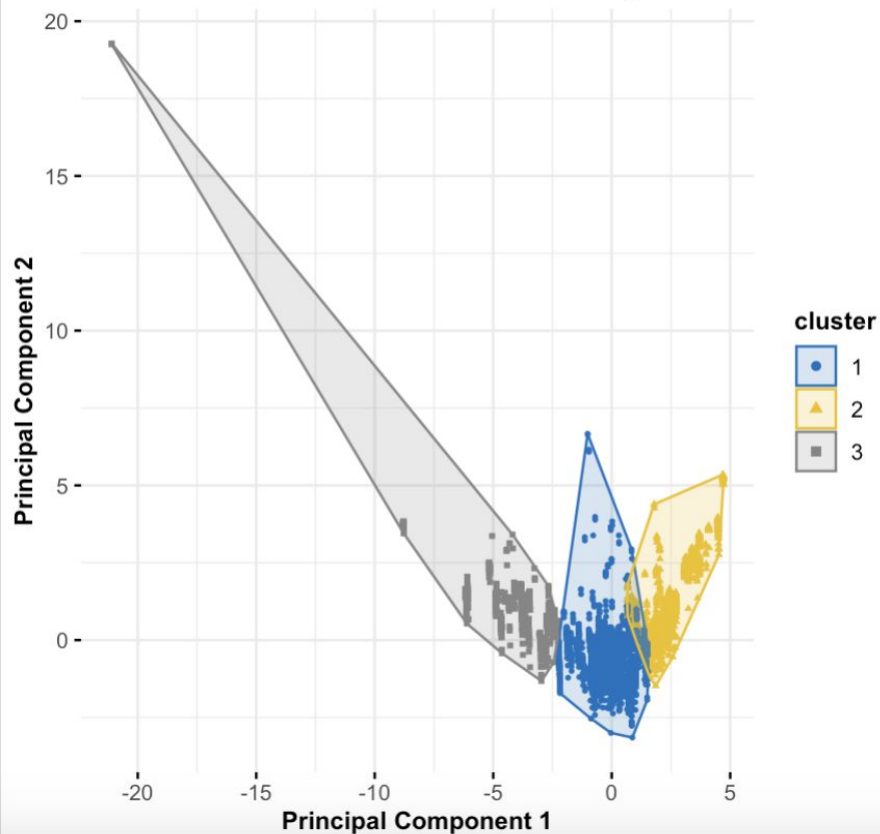| | n | y | MeanDecreaseAccuracy | MeanDecreaseGini |
|---|---|---|---|---|
| NETWORKSIZE | 11.023787 | 10.437391 | 12.831892 | 185.0561380 |
| BETWEENNESS | 11.447184 | 11.205105 | 13.358052 | 185.7450730 |
| NBETWEENNESS | 12.154502 | 10.844694 | 13.555254 | 171.8532026 |
| NBROKERAGE | 6.382967 | 7.028083 | 8.207925 | 70.0791159 |
| DENSITY | 8.588331 | 5.830041 | 9.444128 | 45.7282323 |
| TRANSITIVITY | 12.670610 | 11.238558 | 13.815618 | 249.3083881 |
| num_punc | 2.837492 | 1.181884 | 2.923152 | 0.4568379 |
| num_characters | 2.248371 | 1.396597 | 2.541610 | 0.6145732 |
| mean_tf_idf | 14.767135 | 10.831288 | 14.233036 | 159.0392216 |
| positive_ratio | 11.631051 | 11.407710 | 12.691495 | 88.4420601 |

# K-Means



Elbow Method for Optimal Clusters



**K-Means Clustering Visualization**

Clusters visualized after PCA dimensionality reduction

# K-means — evaluation

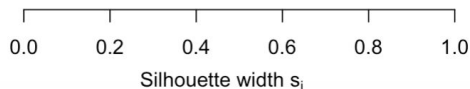**Silhouette Plot for K-Means Clustering**

n = 9886

3 clusters $C_j$

$j : n_j | ave_{i \in Cj} \, s_i$

1 : 6128 | 0.32

2 : 2577 | 0.28

3 : 1181 | 0.29
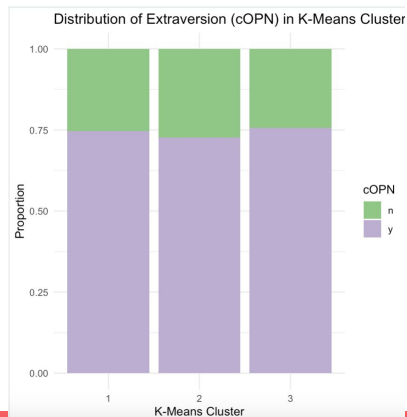
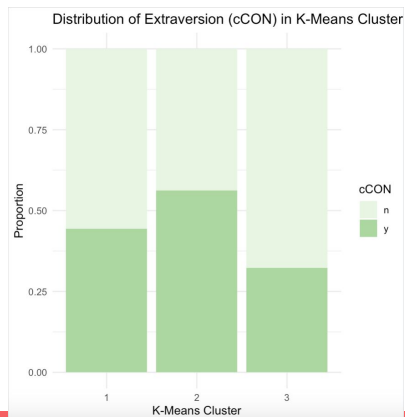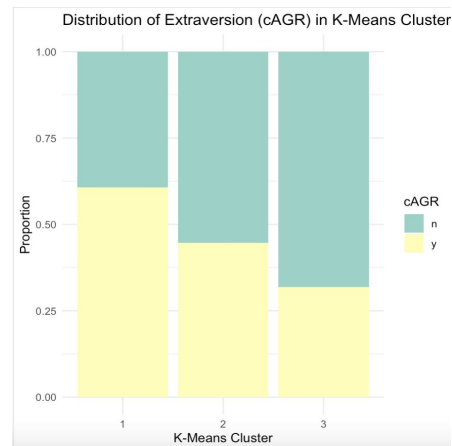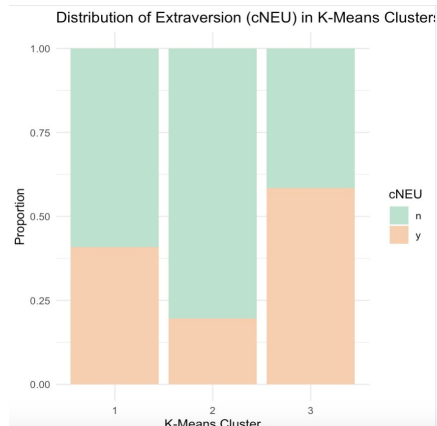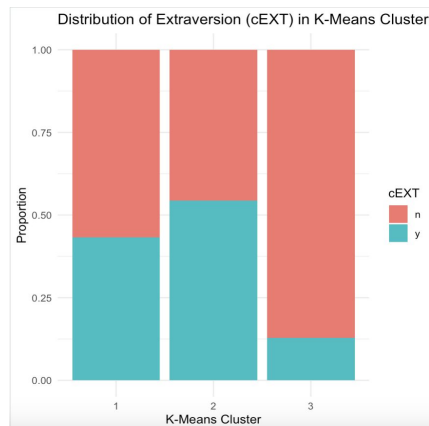0.0    0.2    0.4    0.6    0.8    1.0

Silhouette width $s_i$

Average silhouette width : 0.31

```
> purity_cEXT
> print(purity.
[1] 0.6613892
> purity_cNEU
> print(purity.
[1] 0.6598889
> purity_cAGR
> print(purity.
[1] 0.6143171
> purity_cCON
> print(purity.
[1] 0.59881
> purity_cOPN
> print(purity.
[1] 0.74366
```
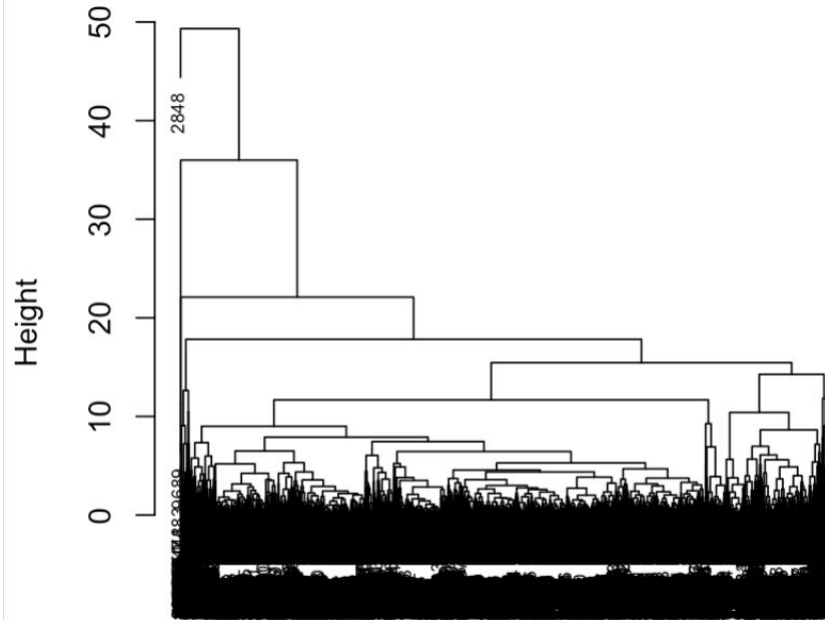
a measure of cluster quality in relation to each personality

# K-means evaluation



Distribution of Extraversion (cEXT) in K-Means Clusters



Distribution of Extraversion (cNEU) in K-Means Clusters



Distribution of Extraversion (cAGR) in K-Means Cluster



Distribution of Extraversion (cCON) in K-Means Cluster



Distribution of Extraversion (cOPN) in K-Means Cluster

# Hierarchical Clustering



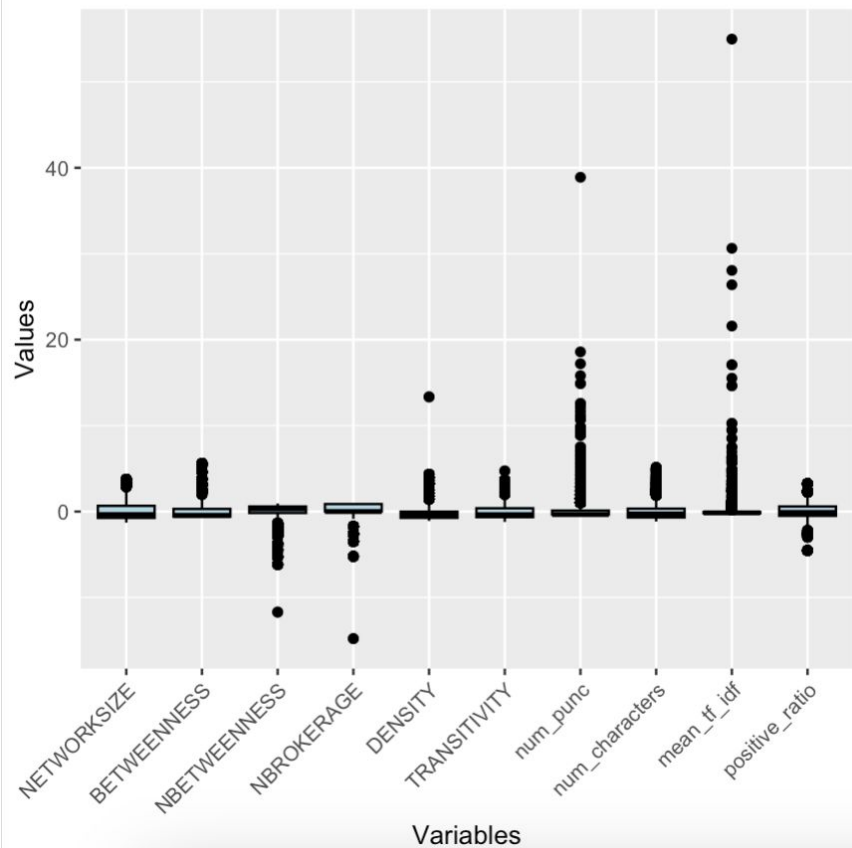Dendrogram for Hierarchical Clustering

# Hierarchical Clustering — evaluation

```
> print(Hmean_silhouette)
[1] 0.8726472
```

Well-separated pt;
Compact clusters —- minimal overlap between

```
> print(paste("Purity for cEXT:", purity_
[1] "Purity for cEXT: 0.858140240817566"
> purity_cNEU <- purity_score(dw_clean$he
> print(paste("Purity for cEXT:", purity_
[1] "Purity for cEXT: 0.875341495497319"
> purity_cAGR <- purity_score(dw_clean$he
> print(paste("Purity for cAGR:", purity_
[1] "Purity for cAGR: 0.843940773719181"
> purity_cCON <- purity_score(dw_clean$he
> print(paste("Purity for cCON:", purity_
[1] "Purity for cCON: 0.846639009747378"
> purity_cOPN <- purity_score(dw_clean$he
> print(paste("Purity for cOPN:", purity_
[1] "Purity for cOPN: 0.914938109211103"
```

Boxplots for Selected Variables

# Advantages & disadvantages

– Advantages

1. Integration of Structured and Unstructured Data
2. Model Evaluation
3. Real-World Application

– Disadvantages

1. Text Complexity
2. Moderate K-Means Performance
3. Feature Selection Limitations

# Potential Challenges and Considerations

1. Data Overlap

   overlapping nature of personality traits
   (e.g., Agreeableness and Conscientiousness)

2. Scalability


3. Text Data Variability

   Variations in language usage, slang, and punctuation

# Ethical considerations

1. User Privacy

   Maintaining anonymity and protecting users from re-identification is crucial

2. Data Sensitivity

   Social media data can be misused in harmful ways

3. Algorithmic Fairness

   Ensure that clustering methods do not reinforce biases or stereotypes