

Big Five Personality Test

CLUSTERING PROJECT



Uri Levy, Elie Cahen, Shlomi Cohen

Big Five Personality Test (or OCEAN model):

Openness (OPN)

- **Openness to experience, conservatism\liberalism, adventurous, in/conventional thinking**

Conscientiousness (CSN)

- **Self discipline, responsibility, self organization, carelessness**

Extroversion (EXT)

- **How Extroverted or Introverted an individual would be?**

Agreeableness (AGR)

- **Sociality, care\carelessness for others, aggression and selfishness**

Neuroticism (EST)

- **Mood stability, indifference, calmness, self-confidence**

DATA

- ▶ Retrieved from Kaggle
- ▶ +1,000,000 cases
- ▶ 109 columns:
 - ▶ 50 questions columns
 - ▶ Answer times columns
 - ▶ Date
 - ▶ No. of entries from the same IP (IPC)
 - ▶ Screen data
 - ▶ Location data

	EXT1	EXT2	EXT3	EXT4	EXT5	EXT6	EXT7	EXT8	EXT9	EXT10	EST1	EST2	EST3	EST4	EST5	EST6	EST7	EST8	EST9	EST10	AGR1	AGR2	AGR3	AGR4	AGR5	AGR6	AGR7	AC
1015336	4.0	2.0	4.0	3.0	4.0	3.0	3.0	3.0	3.0	3.0	4.0	3.0	3.0	3.0	4.0	3.0	4.0	3.0	3.0	3.0	5.0	4.0	2.0	5.0	2.0	4.0	2.0	4.0
1015337	4.0	3.0	4.0	3.0	3.0	3.0	4.0	4.0	3.0	3.0	4.0	3.0	5.0	1.0	5.0	5.0	4.0	4.0	4.0	5.0	2.0	4.0	1.0	4.0	3.0	5.0	3.0	4.0
1015338	4.0	2.0	4.0	3.0	5.0	1.0	4.0	2.0	4.0	4.0	3.0	2.0	4.0	3.0	2.0	2.0	4.0	2.0	4.0	1.0	3.0	5.0	5.0	3.0	2.0	3.0	2.0	4.0
1015339	2.0	4.0	3.0	4.0	2.0	2.0	1.0	4.0	2.0	4.0	4.0	3.0	4.0	2.0	4.0	4.0	2.0	2.0	4.0	4.0	2.0	3.0	2.0	4.0	3.0	4.0	2.0	4.0
1015340	4.0	2.0	4.0	2.0	4.0	1.0	4.0	2.0	4.0	4.0	4.0	3.0	4.0	3.0	2.0	3.0	3.0	1.0	4.0	2.0	1.0	5.0	2.0	4.0	3.0	5.0	2.0	3.0

```
| ▶ M↳ df.shape  
| df.shape  
| (1015341, 109)
```

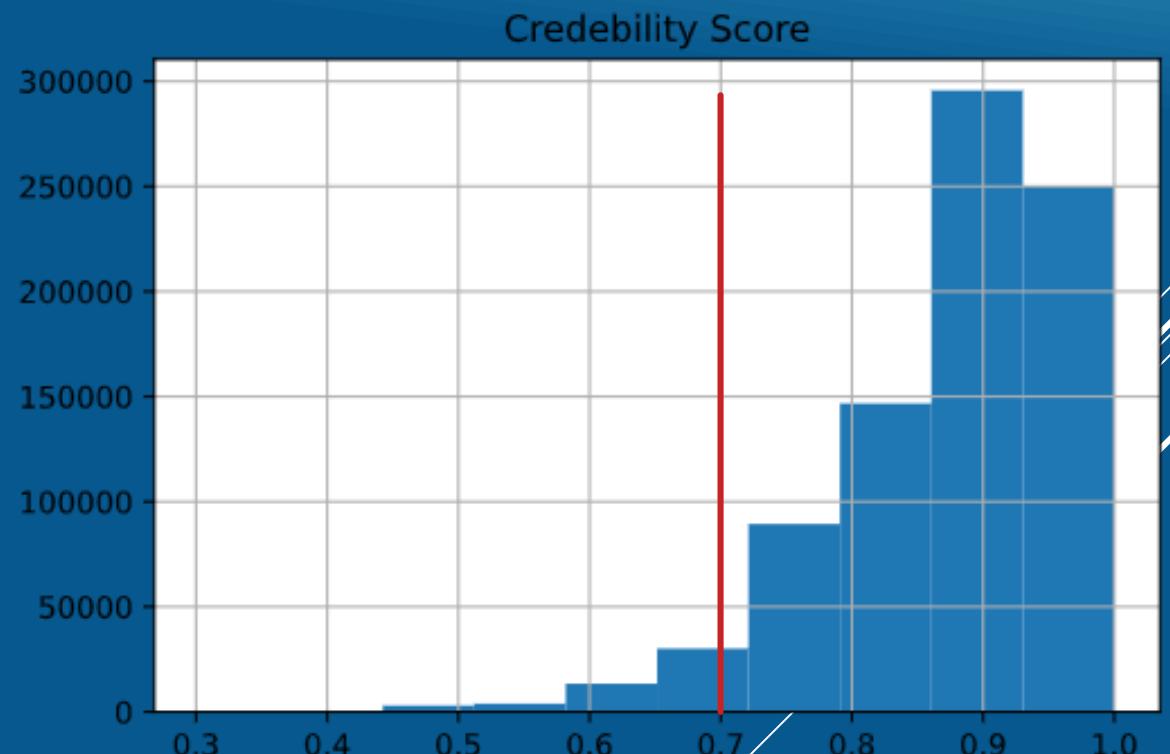
OUR OBJECTIVE:
IDENTIFYING TYPICAL PERSONALITIES AMONG SUBJECTS

EDA:

EDA:

Subjects authentication:

- ▶ Multiple subjects of same IP –
Same user or shared PC?
- ▶ Rapid repliers – less than 3 sec\questions
- ▶ Subject credibility - credibility scale (0-1):



R=0.76

EST7: I change my mood a lot.
EST8: I have frequent mood swings.

1 2 3 4 5
1 2 3 4 5

d=|5-2|=3

EDA:

Missing data:

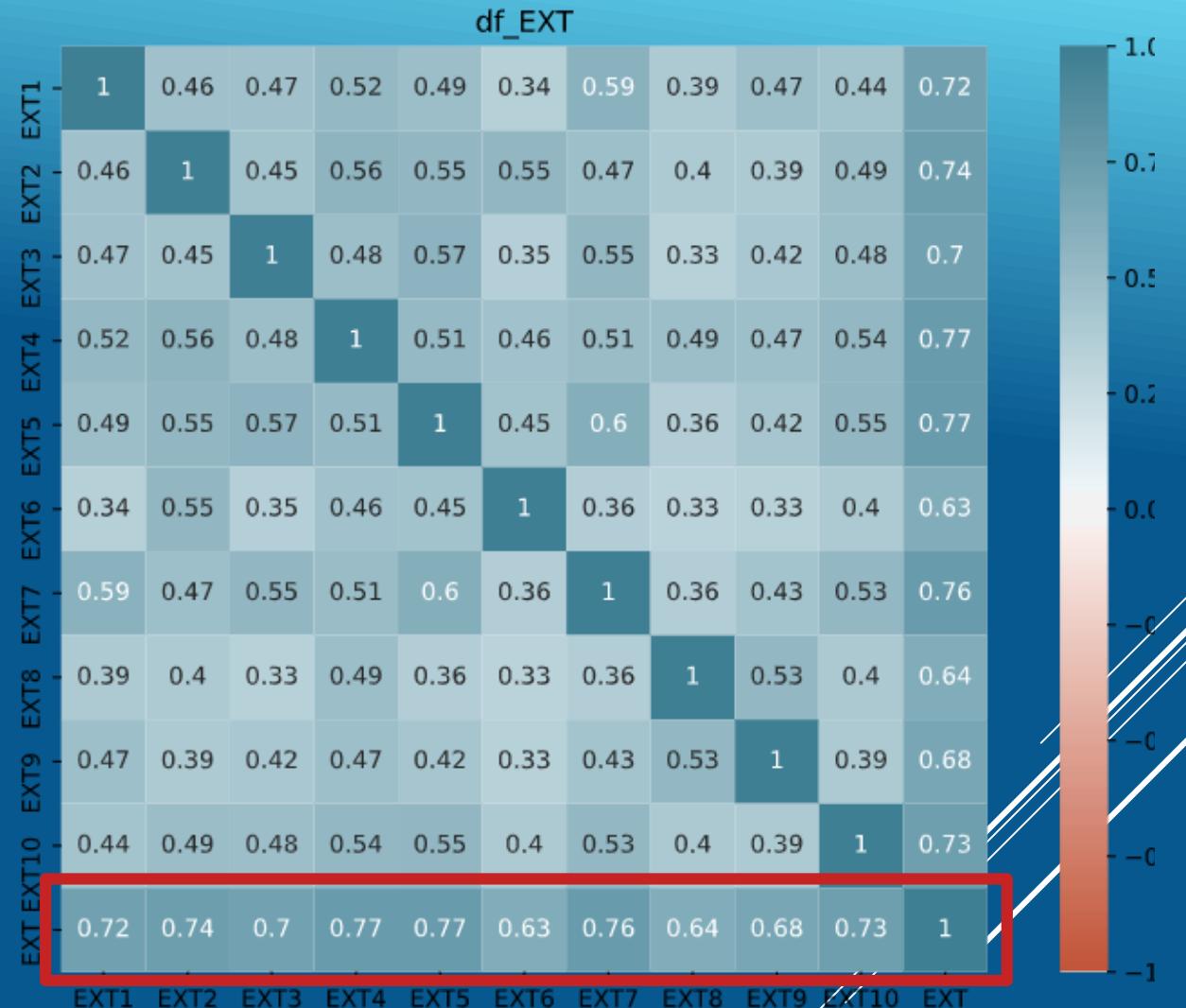
- ▶ Missing > 5 = drop
- ▶ Using inner scale correlations to replace missing data (linear model) -
- ▶ Iterative generation of regression model for each column:

Bayesian Ridge estimator

EXT_7 = ?



EXT_7: I talk to a lot of different people at parties
1 2 3 4 5



FEATURE ENGINEERING:

- ▶ Drop all irrelevant data (GIS, screen, times, etc.)

Questions data:

- ▶ Aligning negative questions:

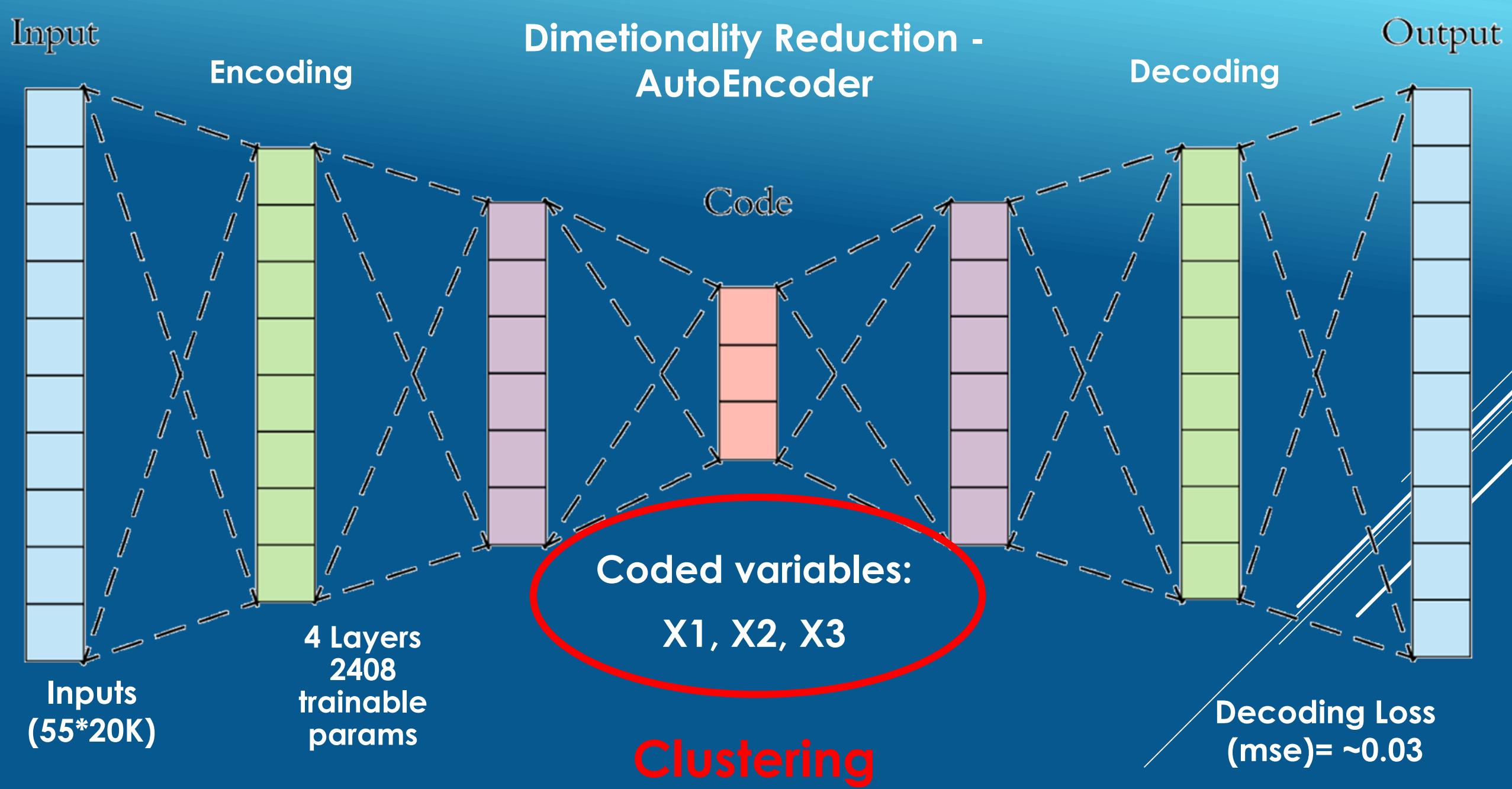
OPN_3: I have a vivid imagination (+)
OPN_6: I don't have a good imagination. (-)

Dimensionality Reduction:

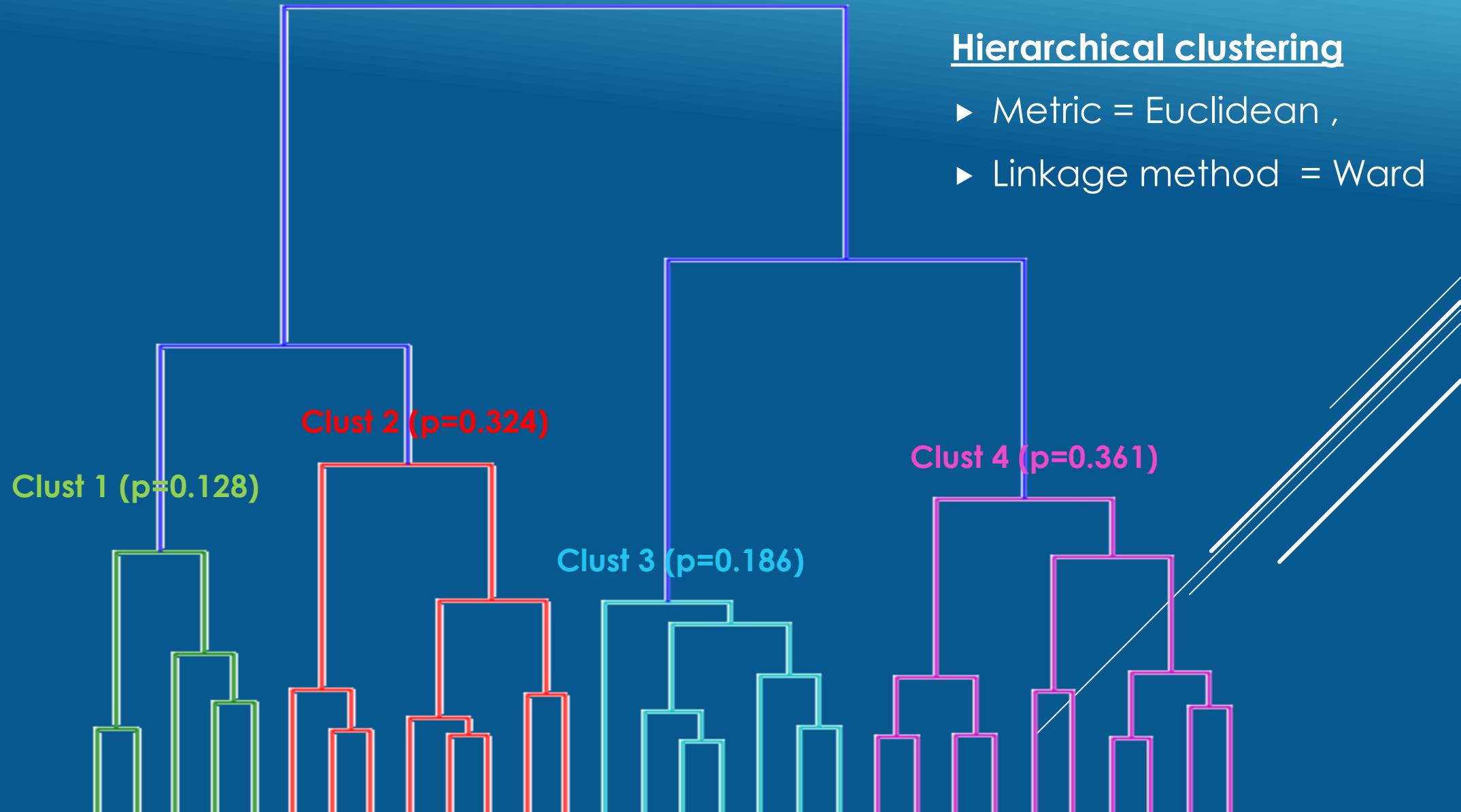
- ▶ For each trait:
 - ▶ Sums of all trait-related scores to one scale:



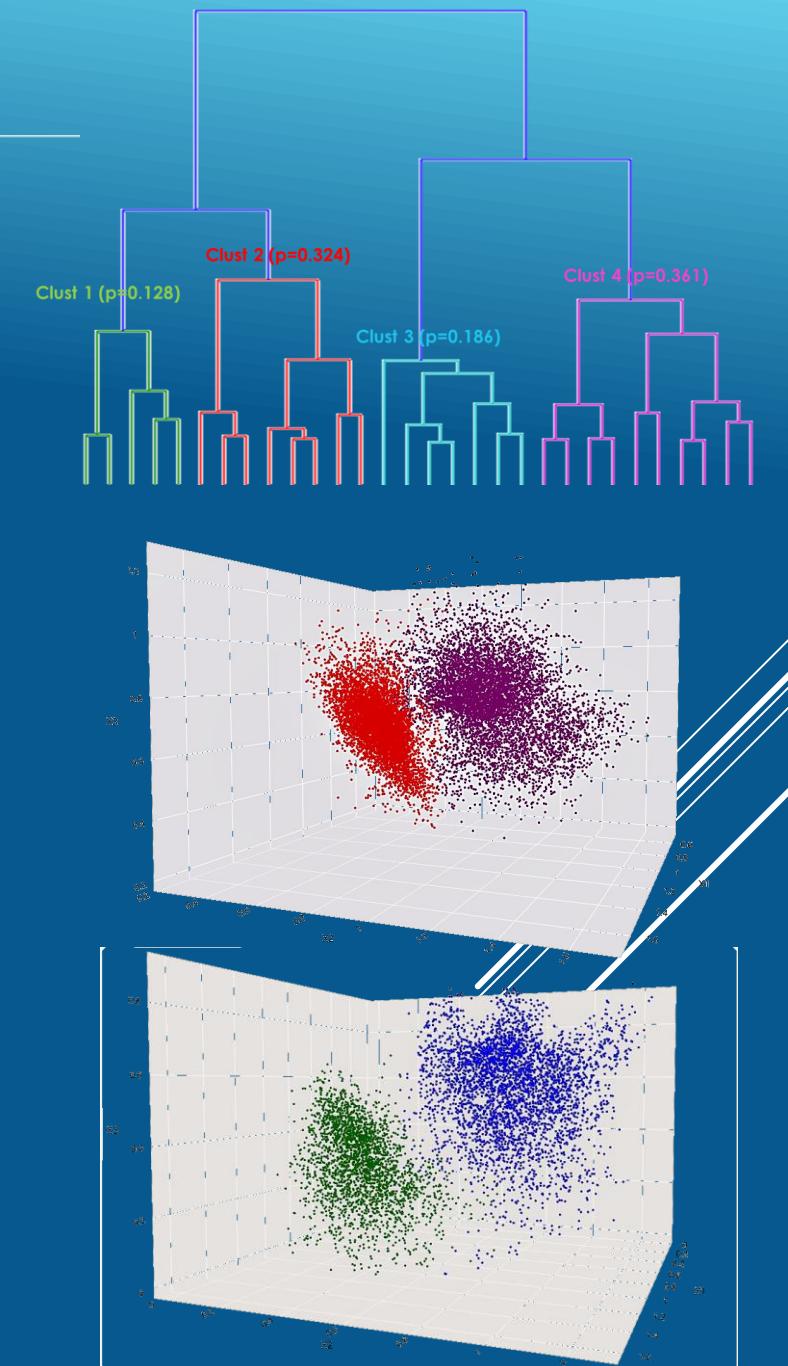
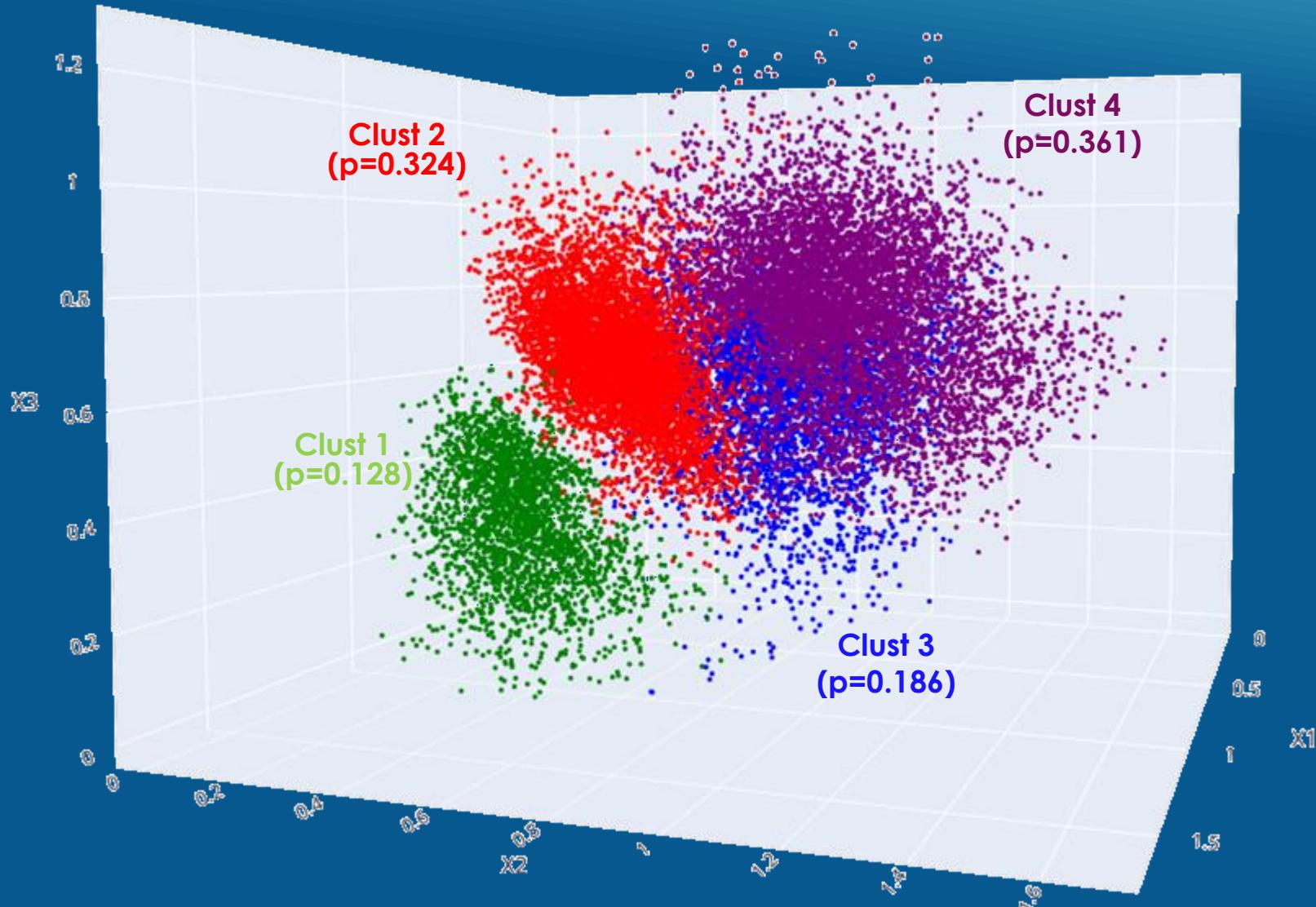
Increase of column variance
(range 1-5) → (range 10-50)



CLUSTERING (AE X's dimensions):



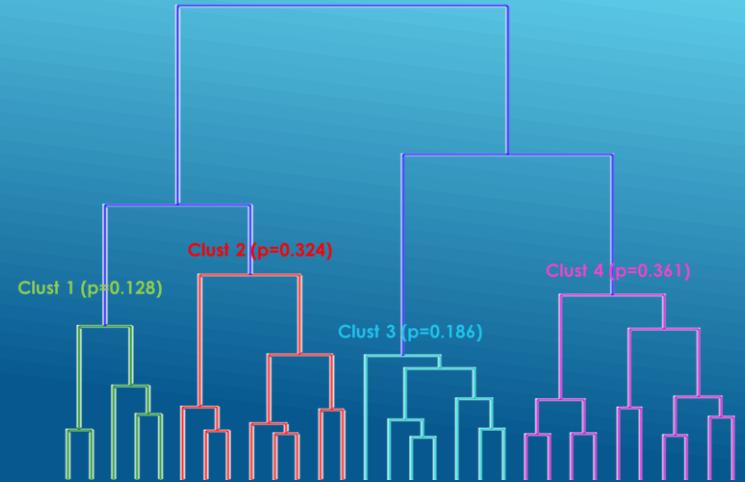
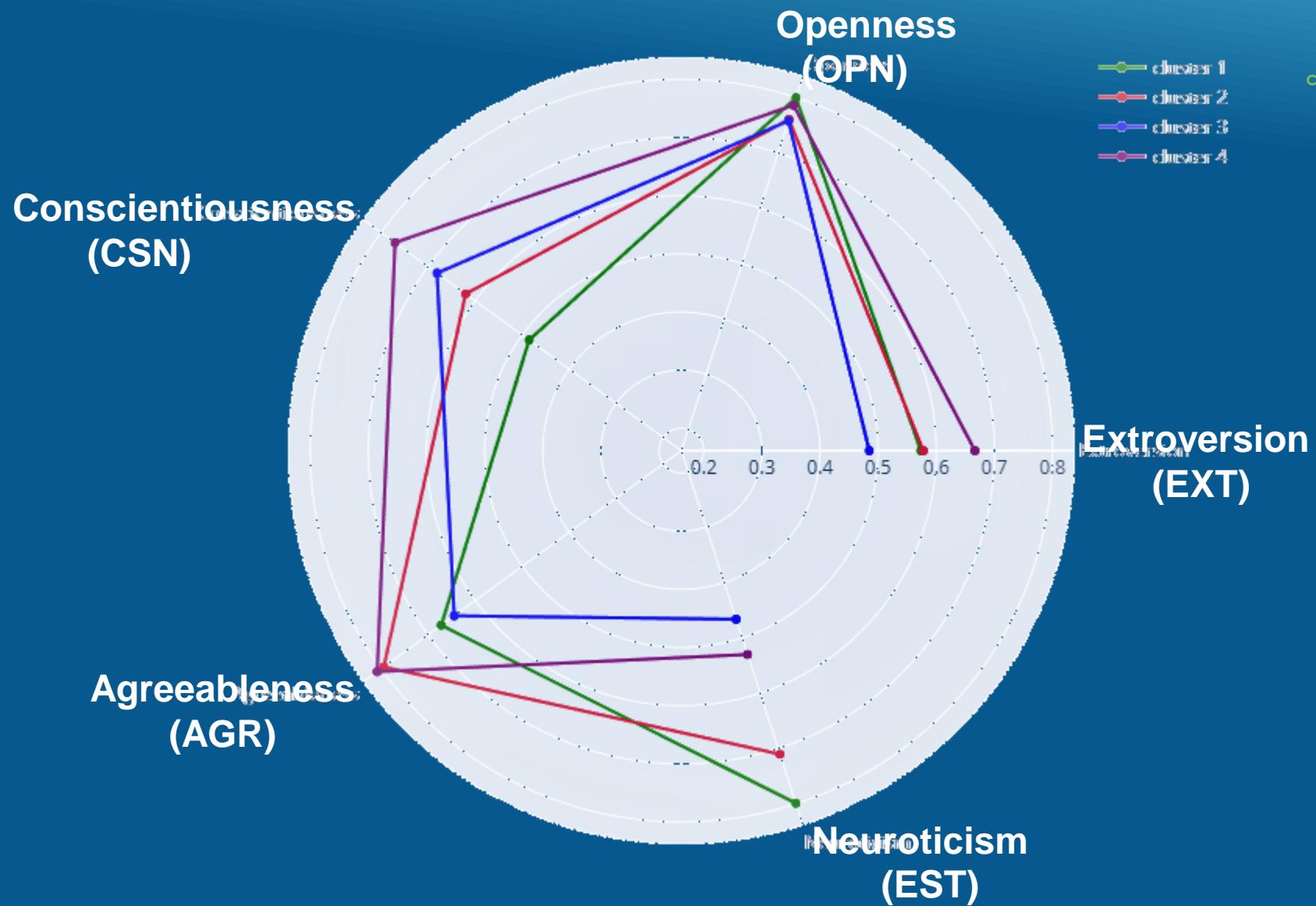
CLUSTERING (AE X's dimensions):



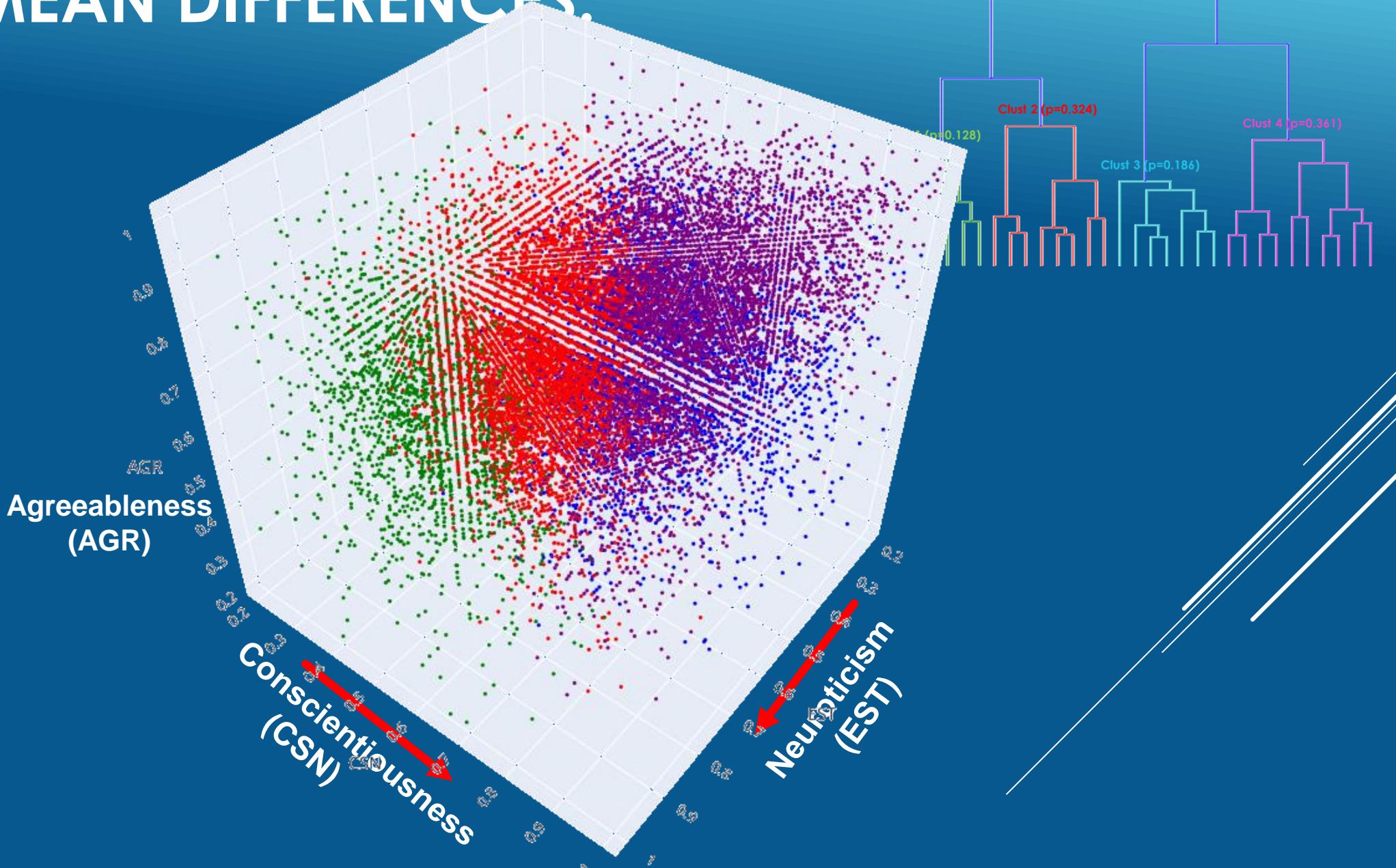


SO WHERE ARE THE DIFFERENCES?

TRAIT'S MEAN DIFFERENCES:



TRAIT'S MEAN DIFFERENCES:



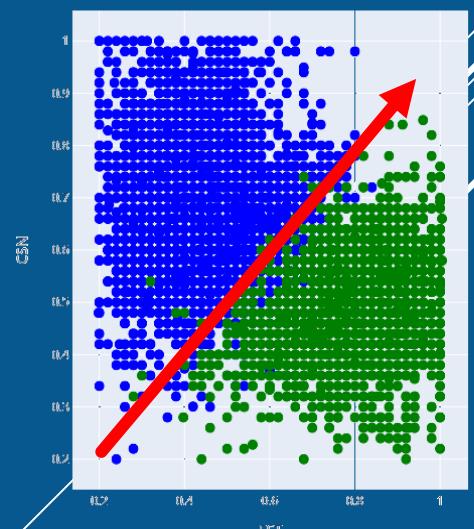
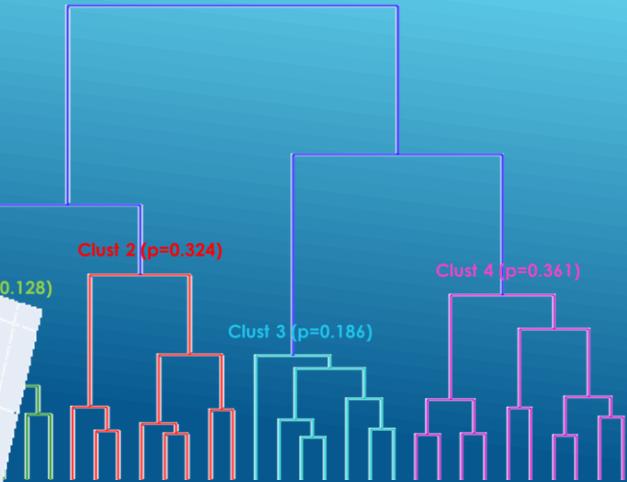
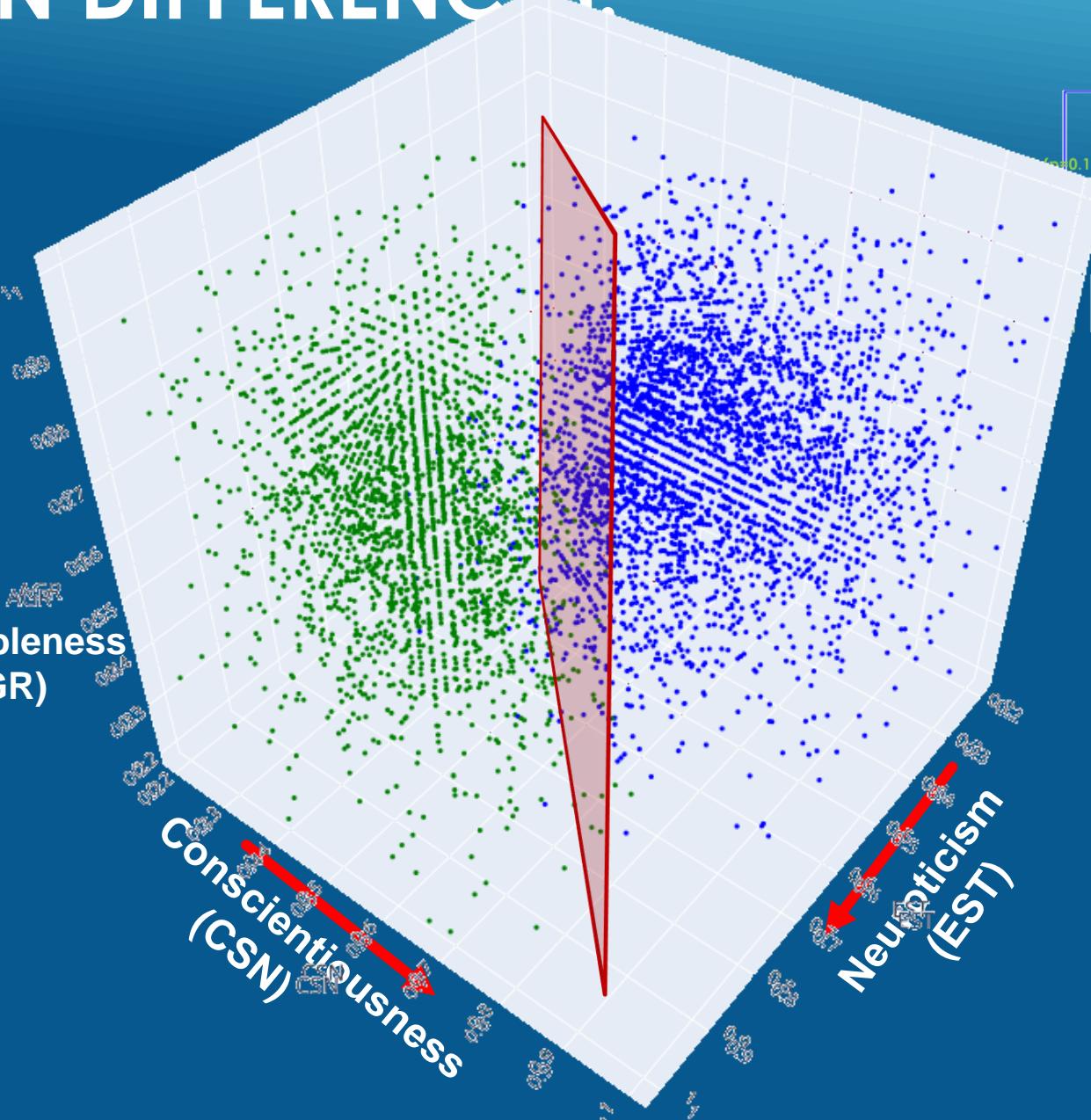
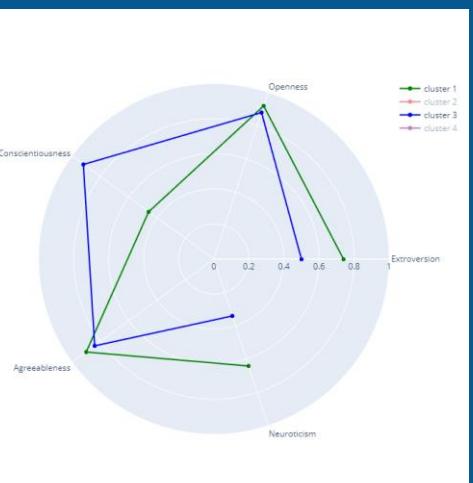
TRAIT'S MEAN DIFFERENCES:

Agreeableness
(AGR)

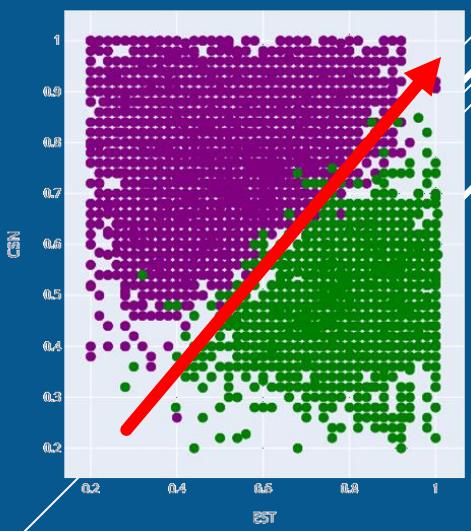
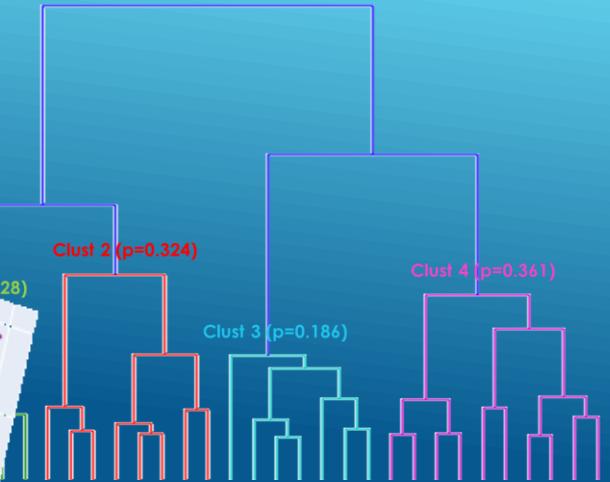
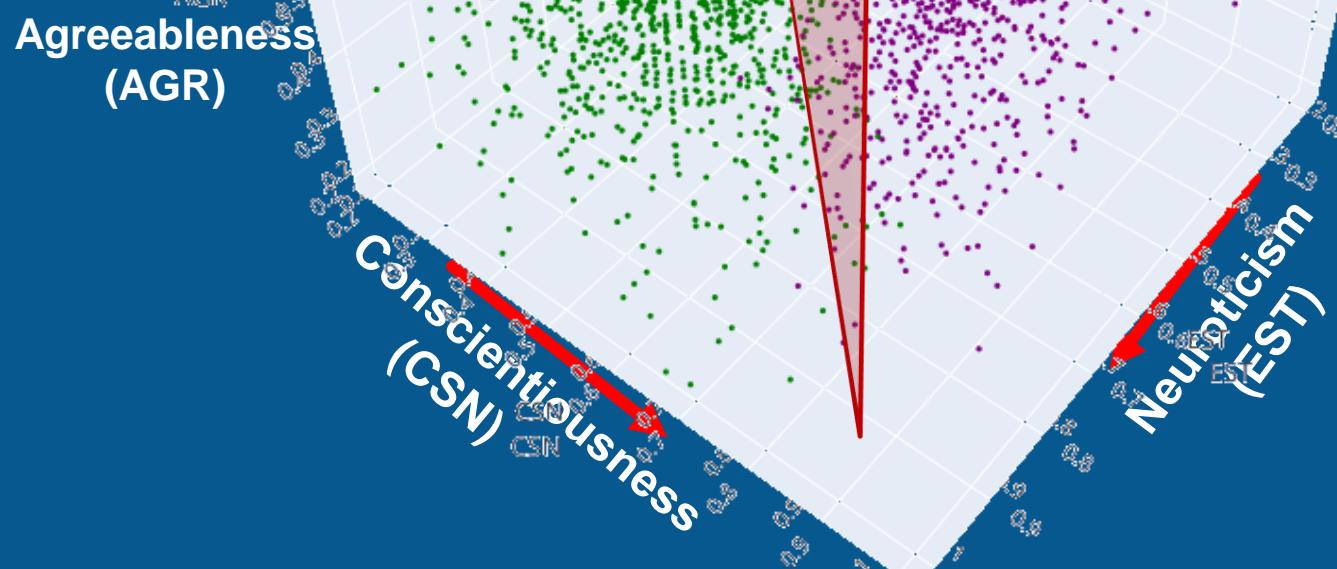
CSN

Conscientiousness
(CSN)

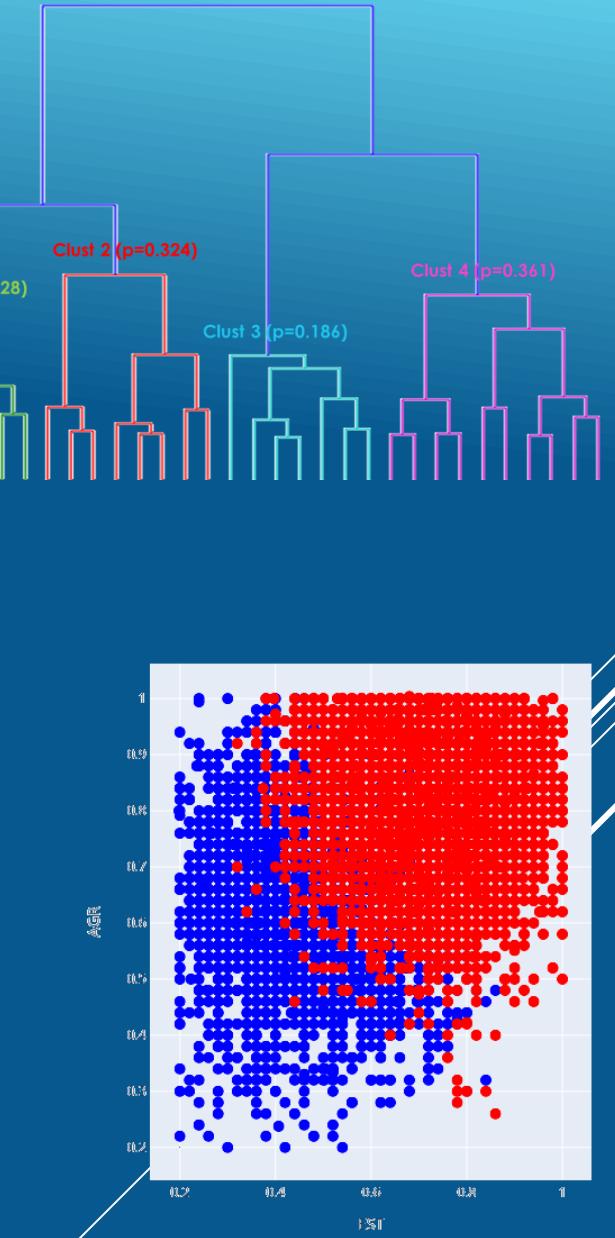
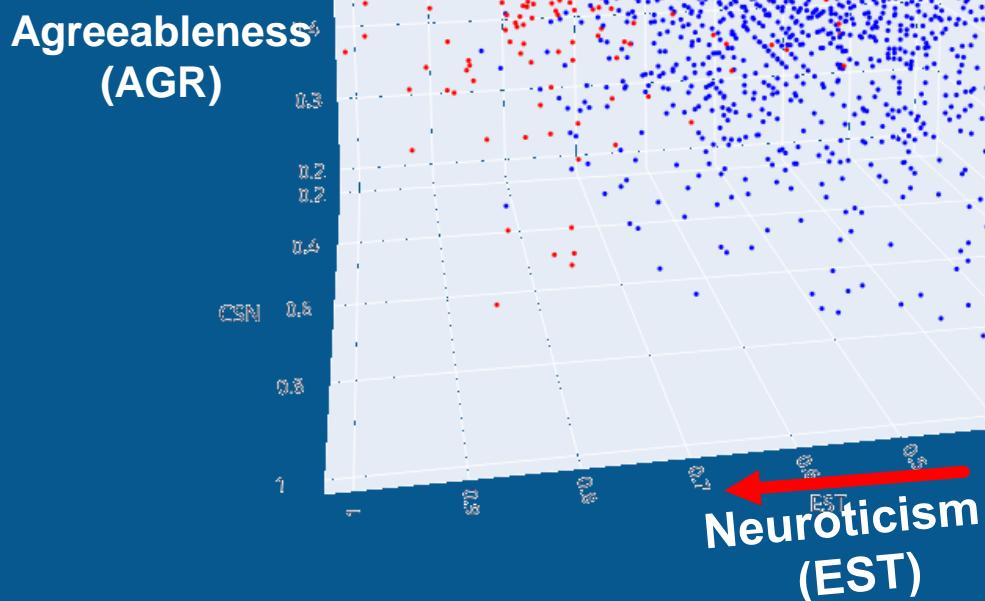
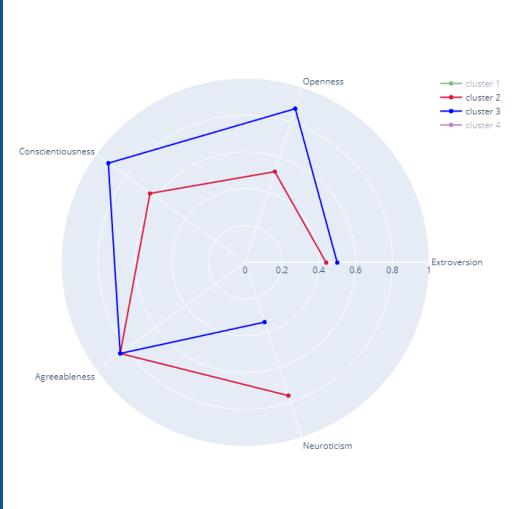
Neuroticism
(EST)



TRAIT'S MEAN DIFFERENCES



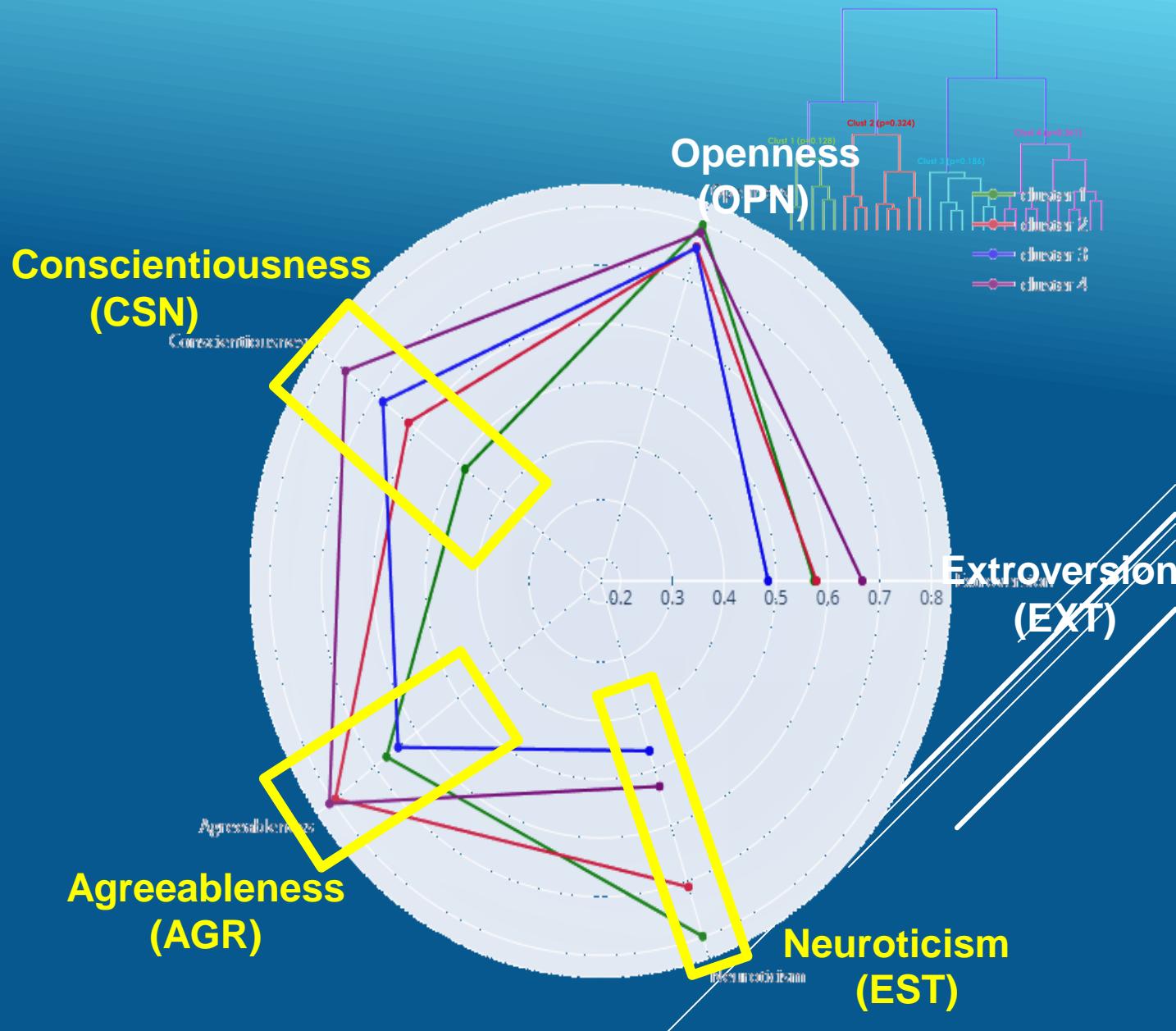
TRAIT'S MEAN DIFFERENCES:



INSIGHTS:

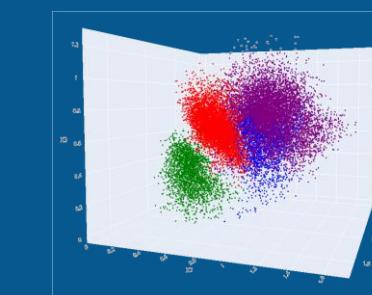
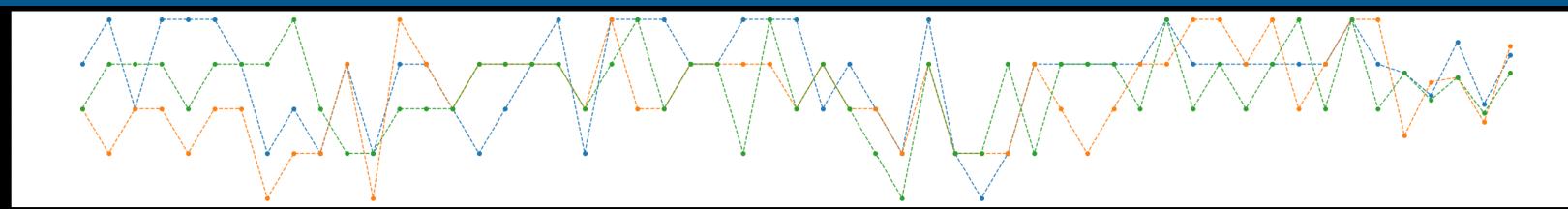
Most robust predictors for OCEAN clustering are:

- A. Neuroticism (i.e. mood stability, calmness, self-confidence, etc.)
- B. Conscientiousness (i.e. self discipline, responsibility, self organization, etc.)
- C. Agreeableness (Sociality, care for others, selfishness, etc.)

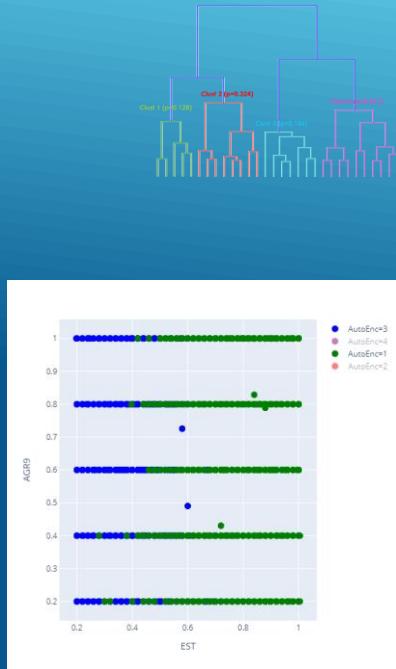
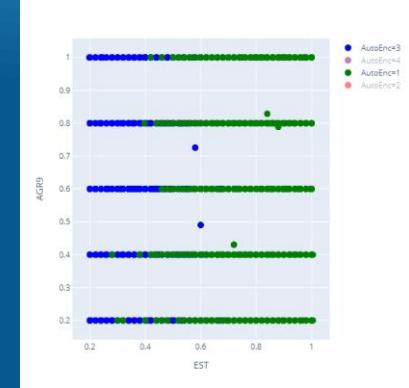


INSIGHTS:

- Clustering is only a half job done
- Subjective reporting is not a continuous measure
 - Designed metrics (cosine?)
 - Very low variance
- Distinctive features aren't Encoded features
- Possible analyses directions:
CNN (image-like), NLP (text-like)



≠



Special Thanks:

Amit Rappel

NAYA-team

THANK YOU FOR LISTENING

