# Explainable AI for psychological profiling from behavioral data

**Yanou Ramon, Sandra Matz, Robert Farrokhnia, David Martens**
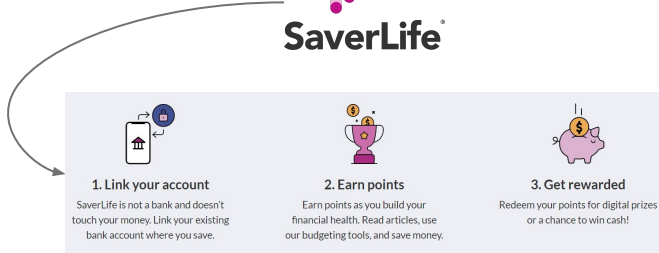
*Joint seminar adm+adrem, Dec 15, 12:30 pm*

# Context

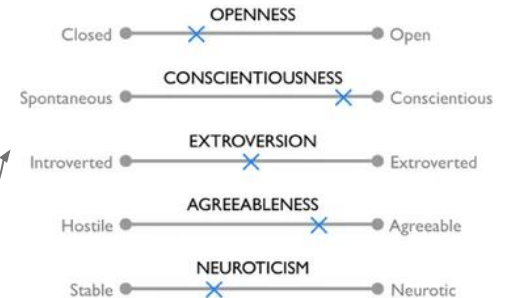Psychological profiling = "the automated assessment of psychological traits from digital footprints" (Matz, 2020)
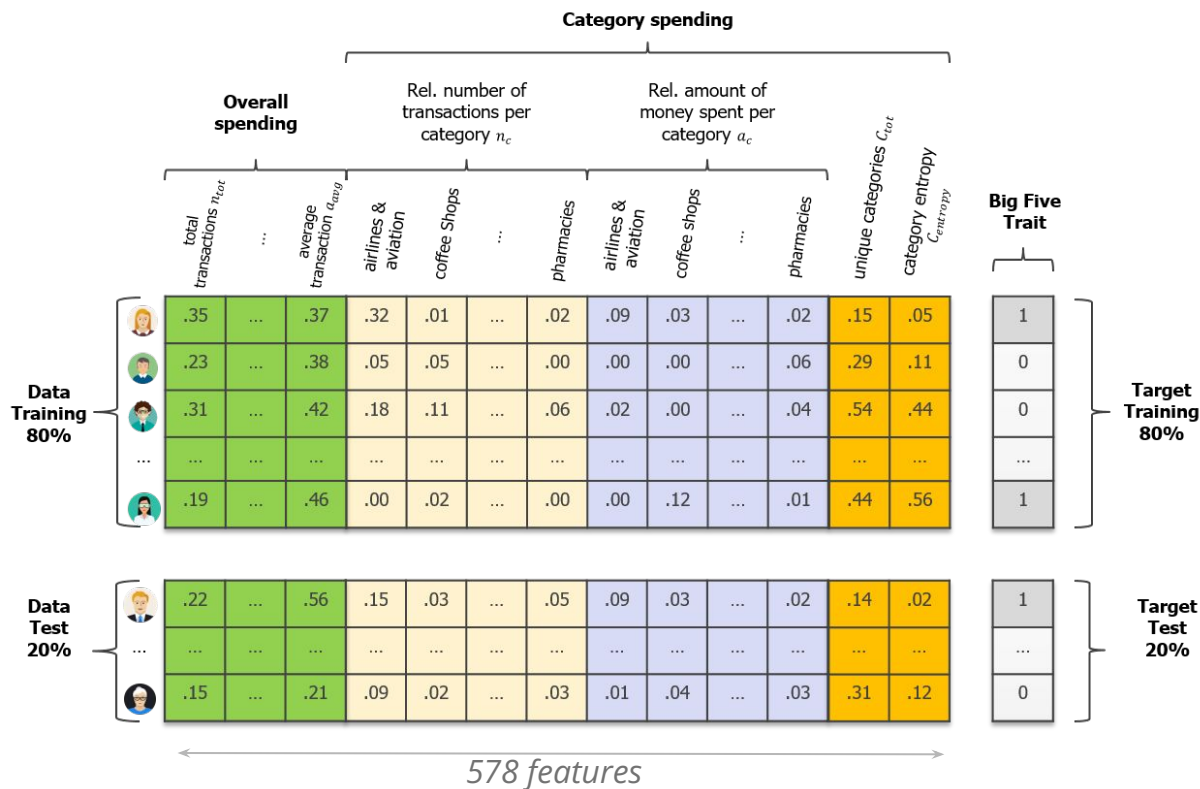
# Case study: personality profiling from consumer spending

# Case study: personality profiling from consumer spending

# Case study: sample data (N=6,408)

# Predictability of personality

- Random Forest models work best
- Decent performance: min=53.4%, max=61.8%
- Best accuracy for Neuroticism
- Conscientiousness & Neuroticism easier to predict than Agreeableness & Openness

# Complication

*Black box models* → Why?

(1) High dimensionality
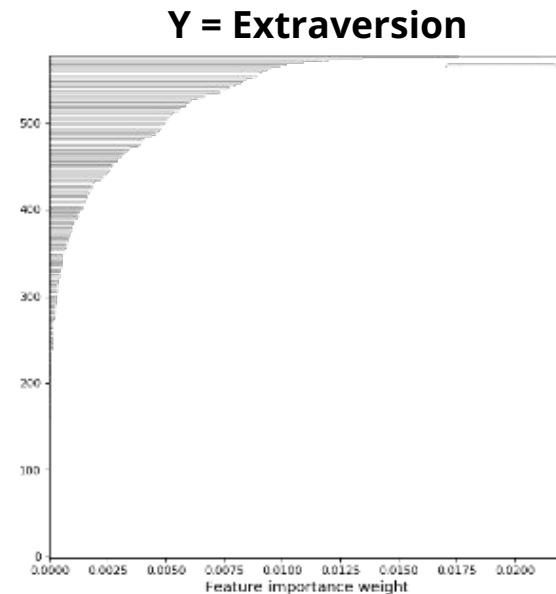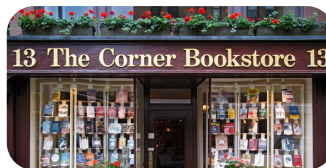(2) Sparsity
(3) Non-redundancy

(e.g., De Cnudde et al., 2019; Clark & Provost, 2019;
Junqué de Fortuny et al., 2013)

# Complication

*Black box models* → Why?

(1) High dimensionality
(2) Sparsity
(3) Non-redundancy

(e.g., De Cnudde et al., 2019; Clark & Provost, 2019; Junqué de Fortuny et al., 2013)



**Y = Extraversion**

# Main claim

Global and local rule-based Explainable AI (XAI) methods are important to gain insight into models for psychological profiling & particularly suitable for digital footprints data

# Global XAI: rule-extraction

Extract if-then-else rules using data and predictions of the model:

if **<condition1>** and **<condition2>** and ... ⇒ **class 1**

elif **<condition3>** ⇒ **class 1**

else **class 2**

*Fidelity (%)* → measures overlap between predictions of the model and predictions of the explanation rules

# Global XAI: results (max. features per rule = 3)

| Trait | Explanation rules |
|---|---|
| Neurotic | *if (Square cash($) ≤ 0.3%) and (Average transaction ≤ $57.08) and (Clothing & Accessories ≤ 0.7%) → Model predicts High Neuroticism*<br>*if (Square cash($) > 0.3%) and (Subscription($) > 0.5%) and (Loans & Mortgages($) ≤ 3.9%) → Model predicts High Neuroticism*<br>*else: Model predicts Default* |
| Conscientious | *if (Square cash > 0.4%) and (Beauty Products > 0.3%) → Model predicts High Conscientiousness*<br>*if (Square cash > 0.4%) and (Beauty Products ≤ 0.3%) and (Clothing & Accessories($) > 0.8%) → Model predicts High Conscientiousness*<br>*if (Square cash ≤ 0.4%) and (Discount Stores > 0.8%) and (Shops > 0.5%) → Model predicts High Conscientiousness*<br>*else: Model predicts Default* |
| Extroverted | *if (Square cash ≤ 0.7%) and (Clothing & Accessories ($) > 0.7%) and (Hotels & Motels > 0.1%) → Model predicts High Extraversion*<br>*if (Square cash > 0.7%) and (Variability transaction amount ≤ 0.31) → Model predicts High Extraversion*<br>*if (Square cash > 0.7%) and (Variability transaction amount > 0.31) and (Service > 0.3%) → Model predicts High Extraversion*<br>*else: Model predicts Default* |
| Agreeable | *if (Square cash ≤ 0.5%) and (Discount Stores($) > 0.1%) and (Shops ≤ 0.6%) → Model predicts High Agreeableness*<br>*if (Square cash > 0.5%) and (Discount Stores > 0.7%) → Model predicts High Agreeableness*<br>*if (Square cash > 0.5%) and (Discount Stores ≤ 0.7%) and (ATM > 5.7%) → Model predicts High Agreeableness*<br>*else: Model predicts Default* |
| Open | *if (Venmo($) > 0.1%) → Model predicts High Openness*<br>*if (Venmo($) ≤ 0.1%) and (Square cash($) > 0.5%) and (Digital purchase > 2.5%) → Model predicts High Openness*<br>*if (Venmo($) ≤ 0.1%) and (Square cash($) ≤ 0.5%) and (Taxi($) > 0.4%) → Model predicts High Openness*<br>*else: Model predicts Default* |

# Global XAI: results for Conscientiousness

if (Square Cash > 0.4%) and (Beauty Products > 0.3%) ⇒ Model predicts **High C**

elif (Square Cash > 0.4%) and (Beauty Products <= 0.3%) and (Clothing & Accessories($) > 0.8%)
⇒ Model predicts **High C**

elif (Square Cash <= 0.4%) and (Discount Stores > 0.8%) and (Shops > 0.5%) ⇒ Model predicts **High C**

else: Model predicts **Default**

*Fidelity: 75.8%*

*Joint Seminar, Dec 15, Rule-based XAI to gain insight into models for psychological profiling from behavioral data*
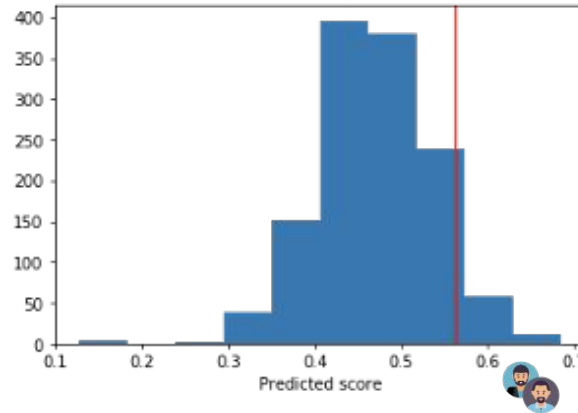
# Local XAI: counterfactual explanations

Extract if-then-else rules using instance **x** and scoring function:

if **<condition1>** and **<condition2>** and ...

⇒ class *changes* from **class 1** to **class 2**

# Local XAI: results for Neuroticism
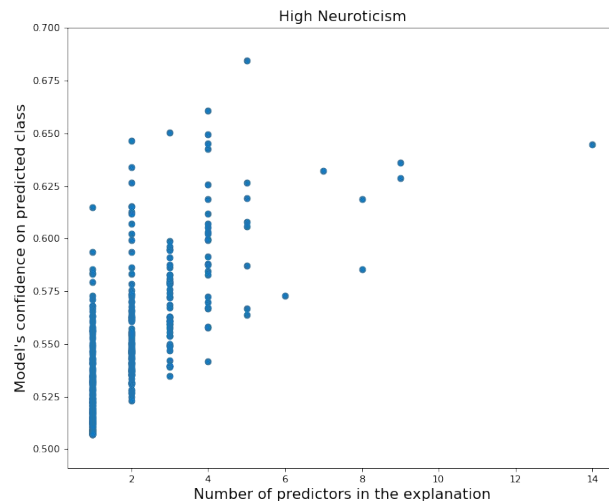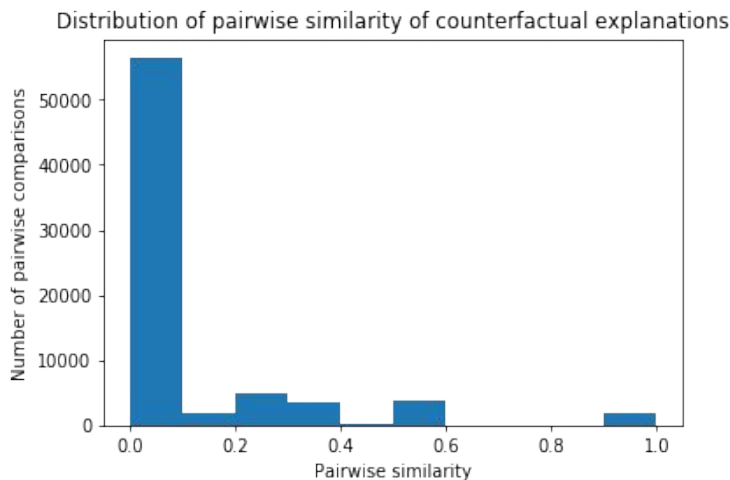


**Person A:**
IF person A had spent *less frequently* in Computers & Electronics, Insurance and Shops, and *more frequently* in Clothing & Accessories and Restaurants ⇒ THEN he would not have been predicted as Neurotic

**Person B:**
IF person B had spent *less frequently* in Shops and Tobacco, and *less money* on Subscription and Tobacco ⇒ THEN he would not have been predicted as Neurotic
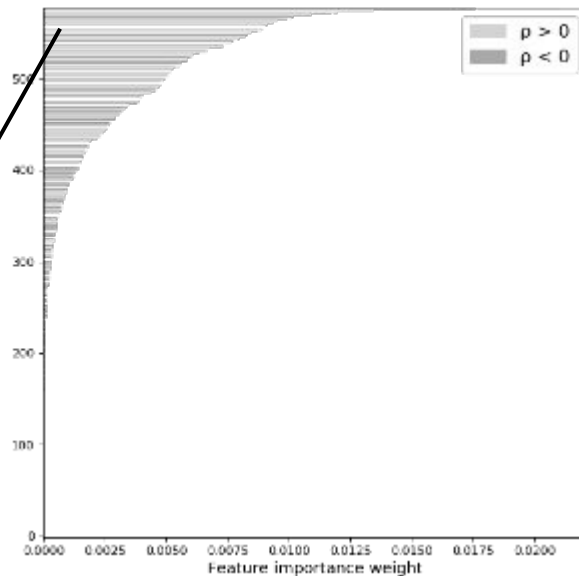
# Local XAI: results for Neuroticism

(1)  **Uniqueness**: variety of features in explanations
(2)  **Concise**: on average, 0.3% of the features in the explanation
(3)  **Comply** with regulatory requirements (e.g., GDPR)



Distribution of pairwise similarity of counterfactual explanations



High Neuroticism

*Joint Seminar, Dec 15, Rule-based XAI to gain insight into models for psychological profiling from behavioral data*

# Local XAI (vs. global)



*'Tobacco'* ranked 73rd out of 578 features in feature relevance list of *'Neuroticism'* model

**Person B:**
IF person B had spent *less frequently* in Shops and Tobacco, and *less money* on Subscription and Tobacco ⇒ THEN he would not have been predicted as Neurotic

*Joint Seminar, Dec 15, Rule-based XAI to gain insight into models for psychological profiling from behavioral data*

# Conclusions

- Both global & local XAI methods are important to open black box, especially when modeling digital footprints data
- Different use cases:

→ **Global:** (i) trust & validation, (ii) audit functionality, (iii) insights, (iv) improve

→ **Local:** (i) provide unique & personalized insight into how data is used, (ii) validate individual predictions

# Thank you!