# Stress Detection

## Group 15:

Ben Butler, Ian Franda, Jack Grosskreuz, Annika Kennerhed, Sara Bay

# Introduction to dataset

3000 rows of daily data (100 participants, 30 days)

Perceived stress scores: range 0-40, (the dataset only ranges from 10-40)

- low: 0-13
- moderate: 13-26
- high: 27-40

Key variables include personality traits, phone usage, sleep metrics, skin conductance and physical activity.

# Proposal questions

Our analysis will focus on the following questions:

- Which factors are most correlated with stress levels?
    - Which personality traits are most correlated with high stress?
    - How can stress levels vary throughout the week?
    - Is phone usage correlated with stress level?
    - How is sleep/physical activity correlated with stress levels?

We will use the results from our data to address these questions:

- Are these determining factors easy to control or manage?
- What would be potential ways to mitigate high stress levels?

# Methods

- No missing values
- Split data 80/10/10 by participant IDs into train/val/test
- Perceived Stress Score (PSS) column used as y values

# Attempt #1: Decision Tree Regressor

- Reasons for using DecisionTreeRegressor:
  - All variables are numeric
  - Decision tree will choose most relevant variables and sort out the least useful variables
  - Many variables are on a scale, which means linearity might not make sense
- Attempted DecisionTreeRegressor
  - Used grid search to find best parameters
    - max_depth = 1
    - min_impurity_decrease = 0.45
    - Finds that the best tree has no splits
  - Guesses 24.82 (mean) for every point

Our Decision "Tree"

squared_error = 73.91
samples = 2400
value = 24.82

# Attempt #2: Decision Tree Classifier

- Bin stress levels into high stress level (27-40) and low/moderate stress (10-26)
- Attempted DecisionTreeClassifier
  - Score of 0.53 on validation data
  - 57% of validation data is low/moderate stress
  - Guessing low/moderate every time makes a better model than the Decision Tree Classifier

# Attempt #3: Linear models

- Linear model using "Openness", "Neuroticism" and "mobility_distance"
- Lasso chooses "Openness", "call_duration" and "num_sms"

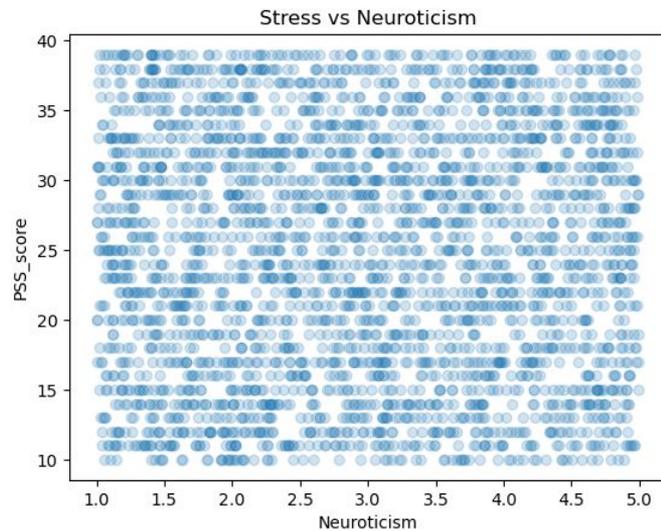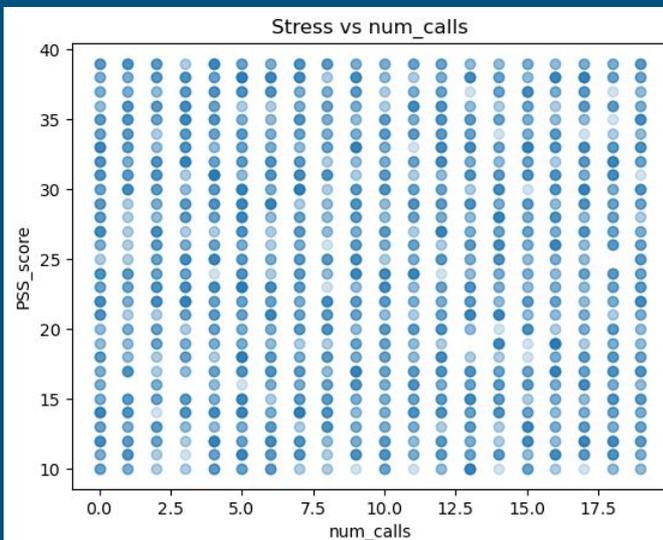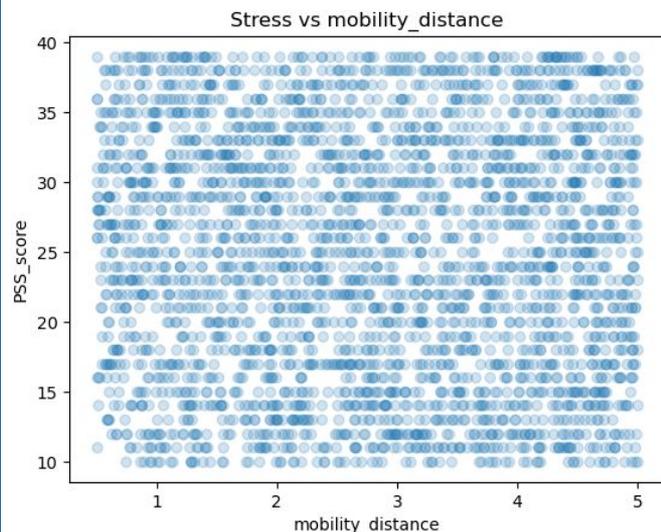# Performance of regression models on validation data

Linear regression does best on validation data, but difference is very small

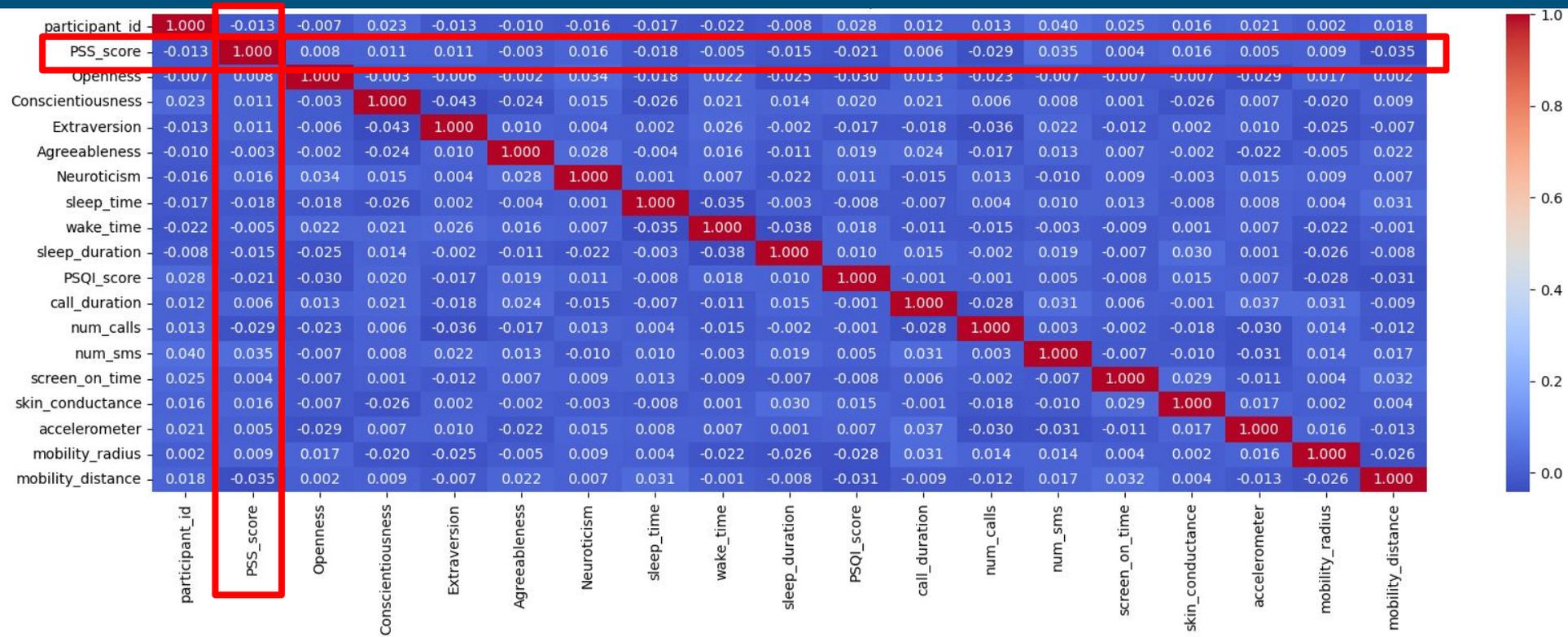Since the performance measures are so close, we looked more into correlation between features

| Method | MSE on validation data |
|---|---|
| Decision tree regressor (guess mean every time) | 75.535 |
| LASSO | 75.532 |
| Linear Regression | 75.498 |

# Correlation check

Plots show no correlation between PSS_score and mobility_distance, num_calls and neuroticism.
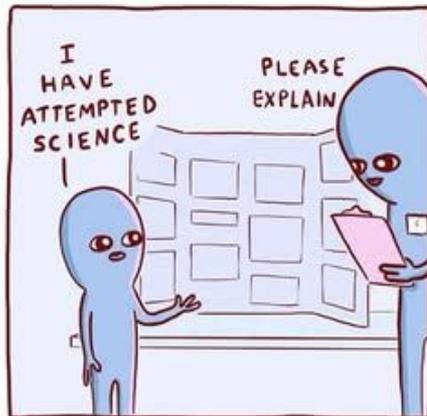
# Correlation Heatmap

# Outputs/Analysis

- The dataset shows that stress levels cannot necessarily be predicted from:
  - Personality traits
  - Sleep metrics
  - Skin conductance
  - Phone usage
  - Physical activity
- Possible explanations:
  - PSS_scores are conducted from a questionnaire
    - Not precise
    - Different interpretations of questions
  - Different people have different causes/symptoms for stress

# Conclusion