

# Case Study: Modeling **Self-Perceived Health** and **Longevity** with SHARE Dataset

Калоян Томев  
Михаил Ангелов  
Сергей Филипов

# Outline



Introduction



Dataset



Aim



Charts



Methodology



Results



Conclusions



Appendix

*a research infrastructure for studying the effects of health, social, economic and environmental policies over the life-course of European citizens and beyond.*

9 waves from 2004 to 2021/2022

In Bulgaria – from the 7<sup>th</sup> wave (2018–)



# Dataset



46733  
observations



3540  
variables



26 European  
countries + Israel

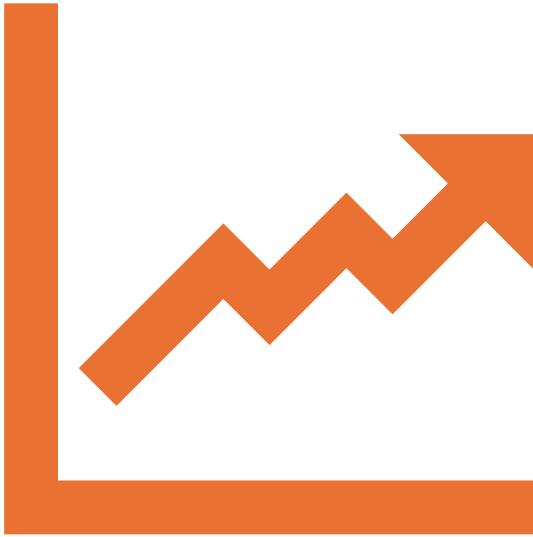
## Aim

**Self-Perceived Health**

*how individuals rate their health*

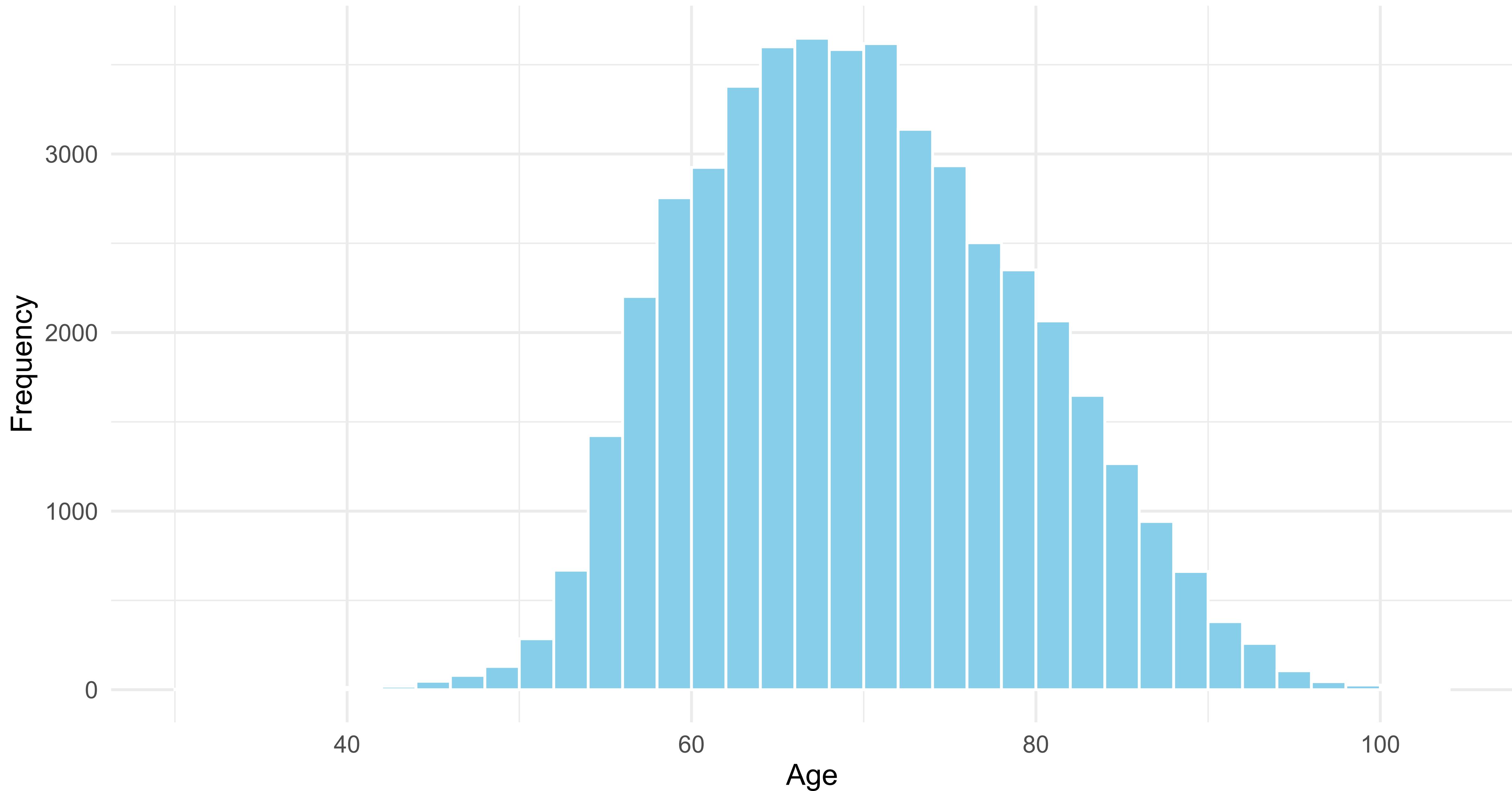
**Age**

*used as a proxy for longevity*

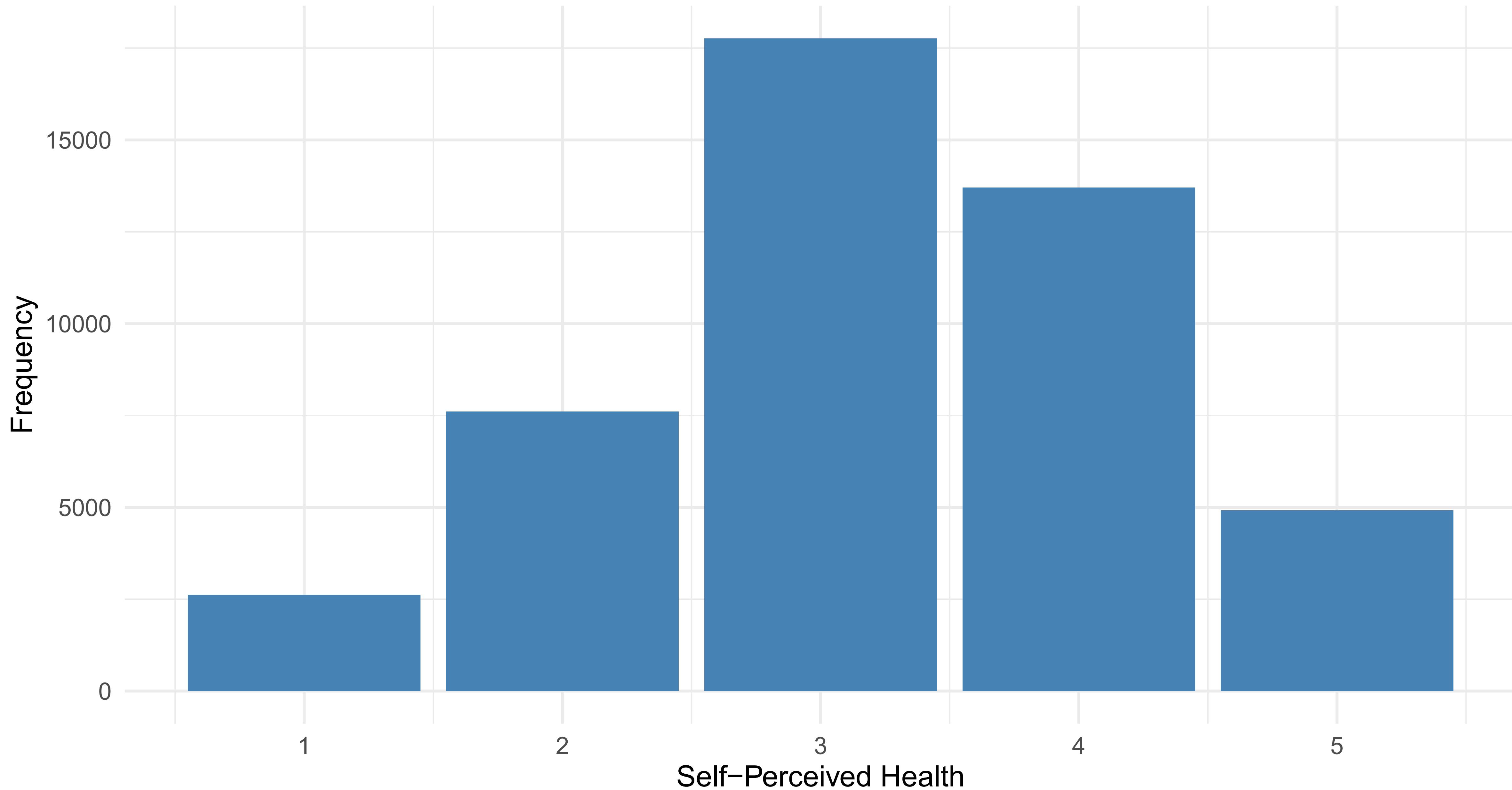


# Charts

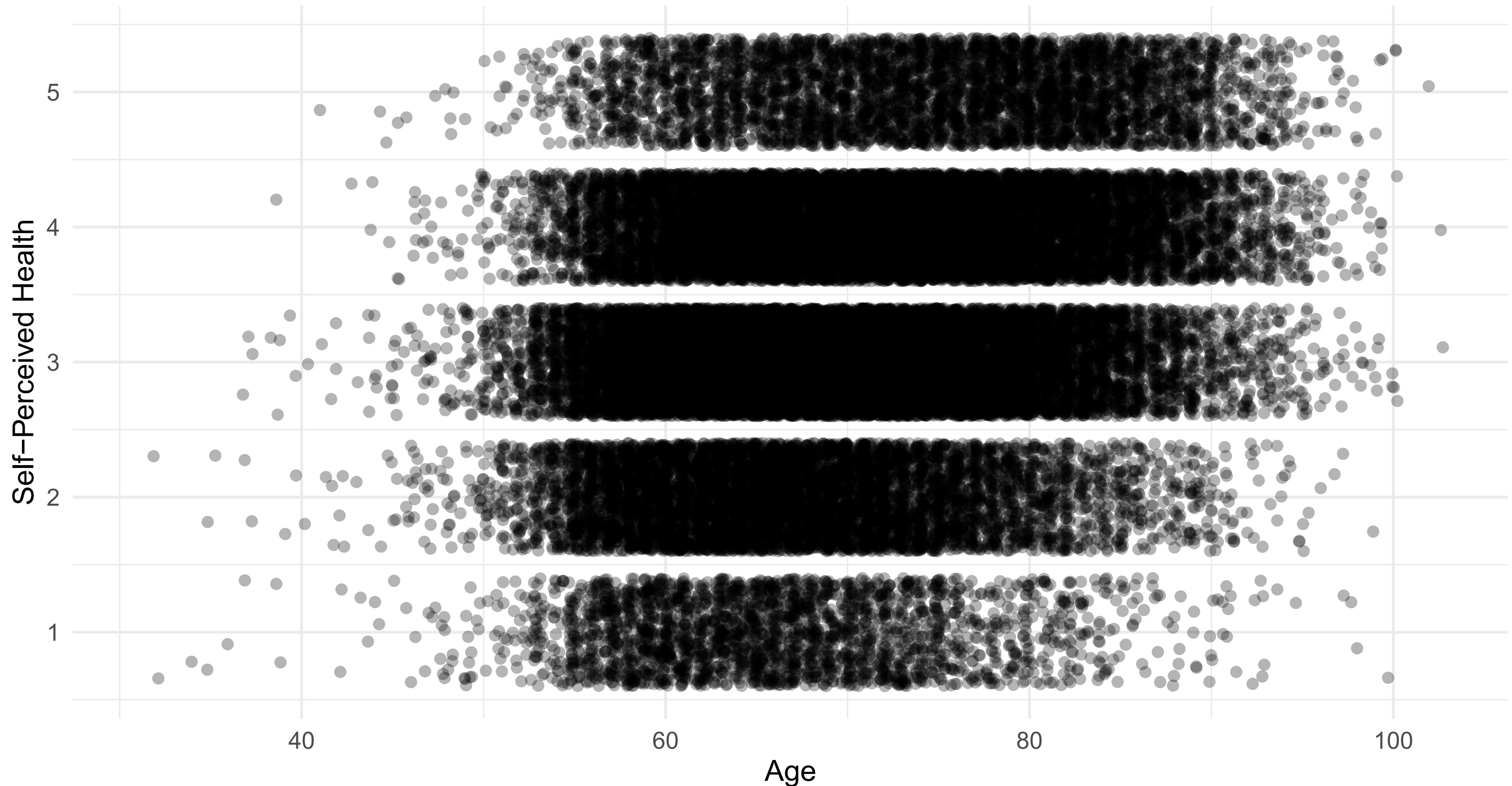
# Age Distribution



# Self-perceived health Distribution



# Age gradient in Self-perceived health

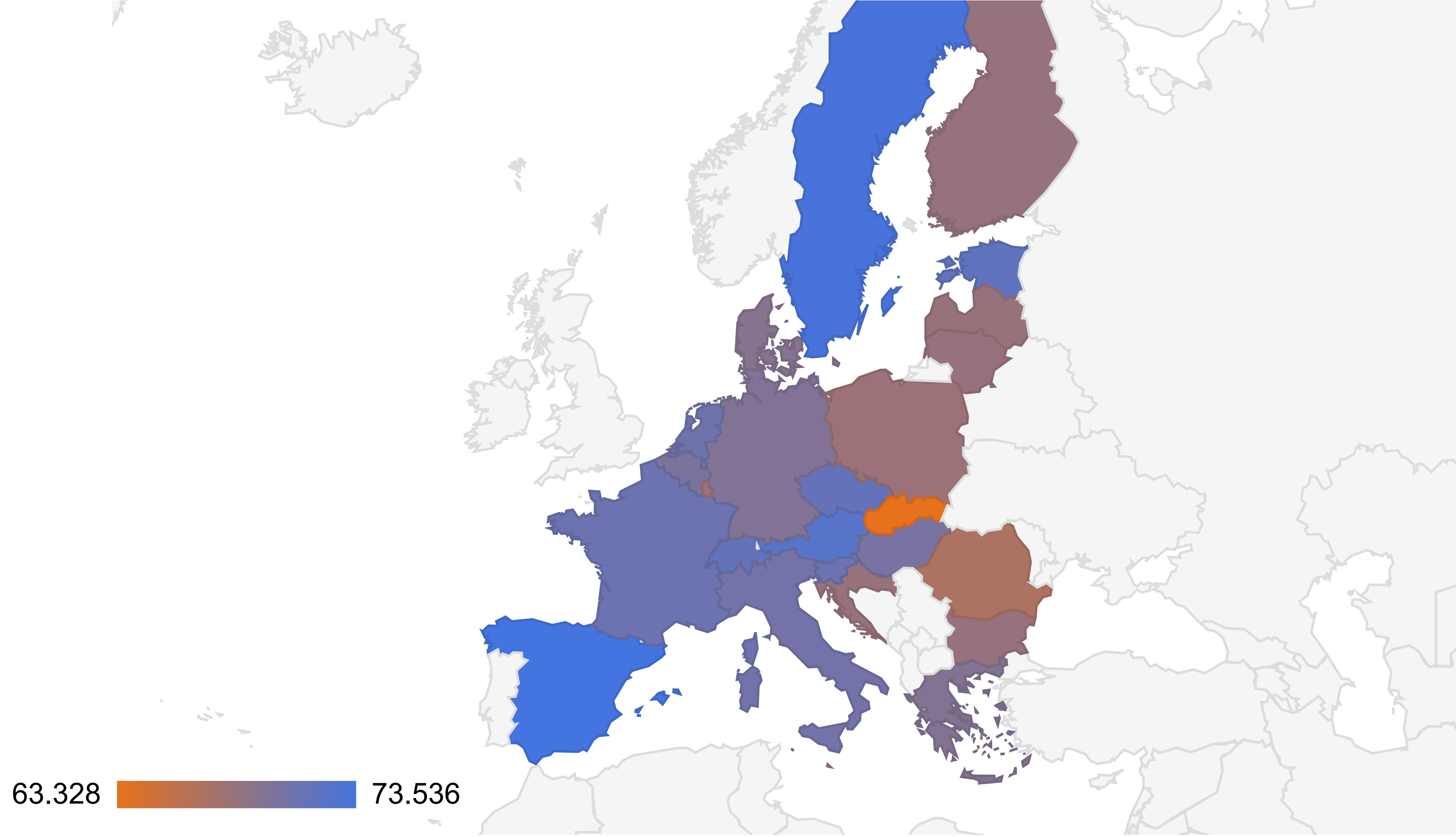


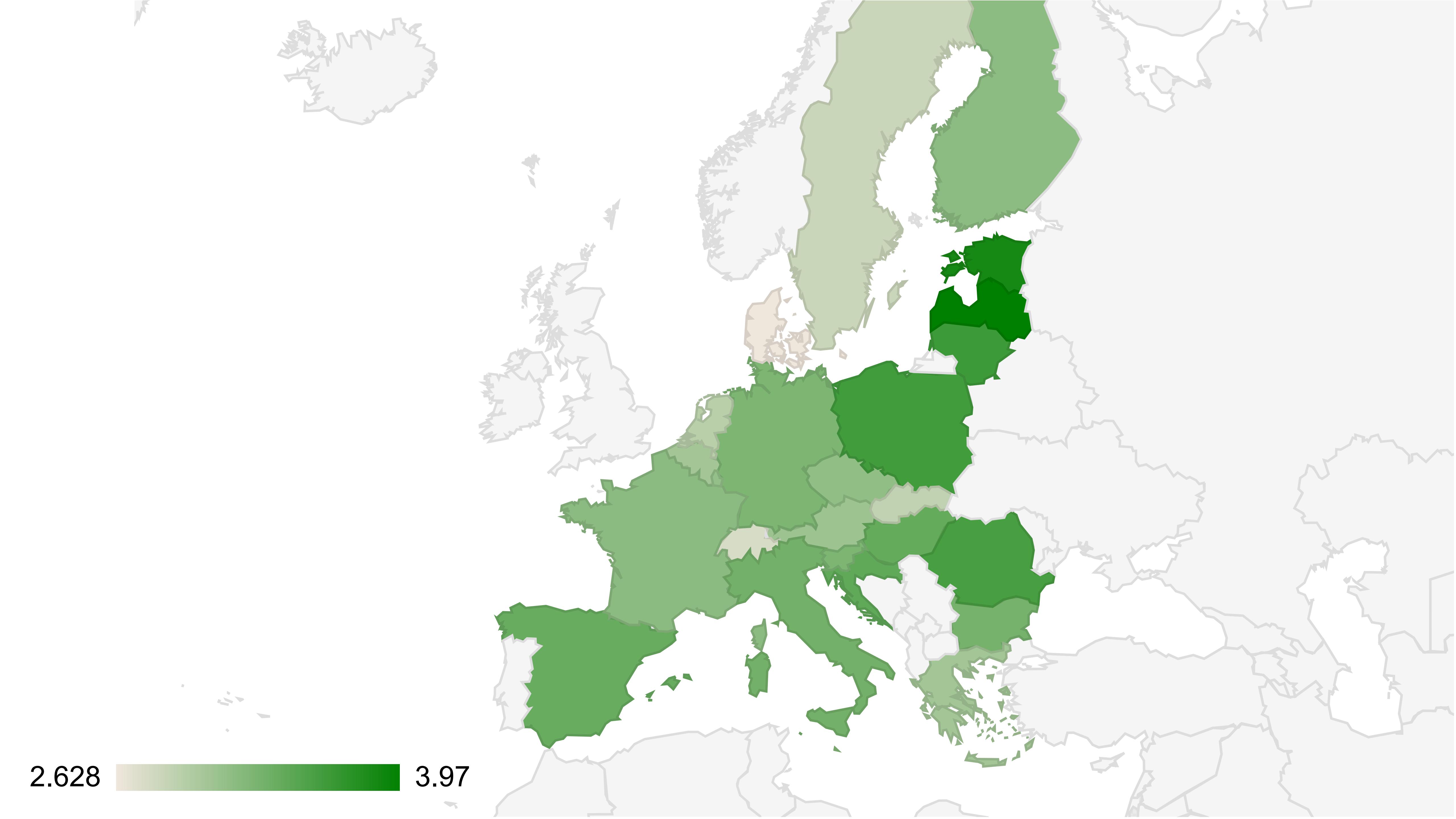
# Self-perceived health distribution by Country



**Self-Perceived Health**

- Excellent
- Very Good
- Good
- Fair
- Poor

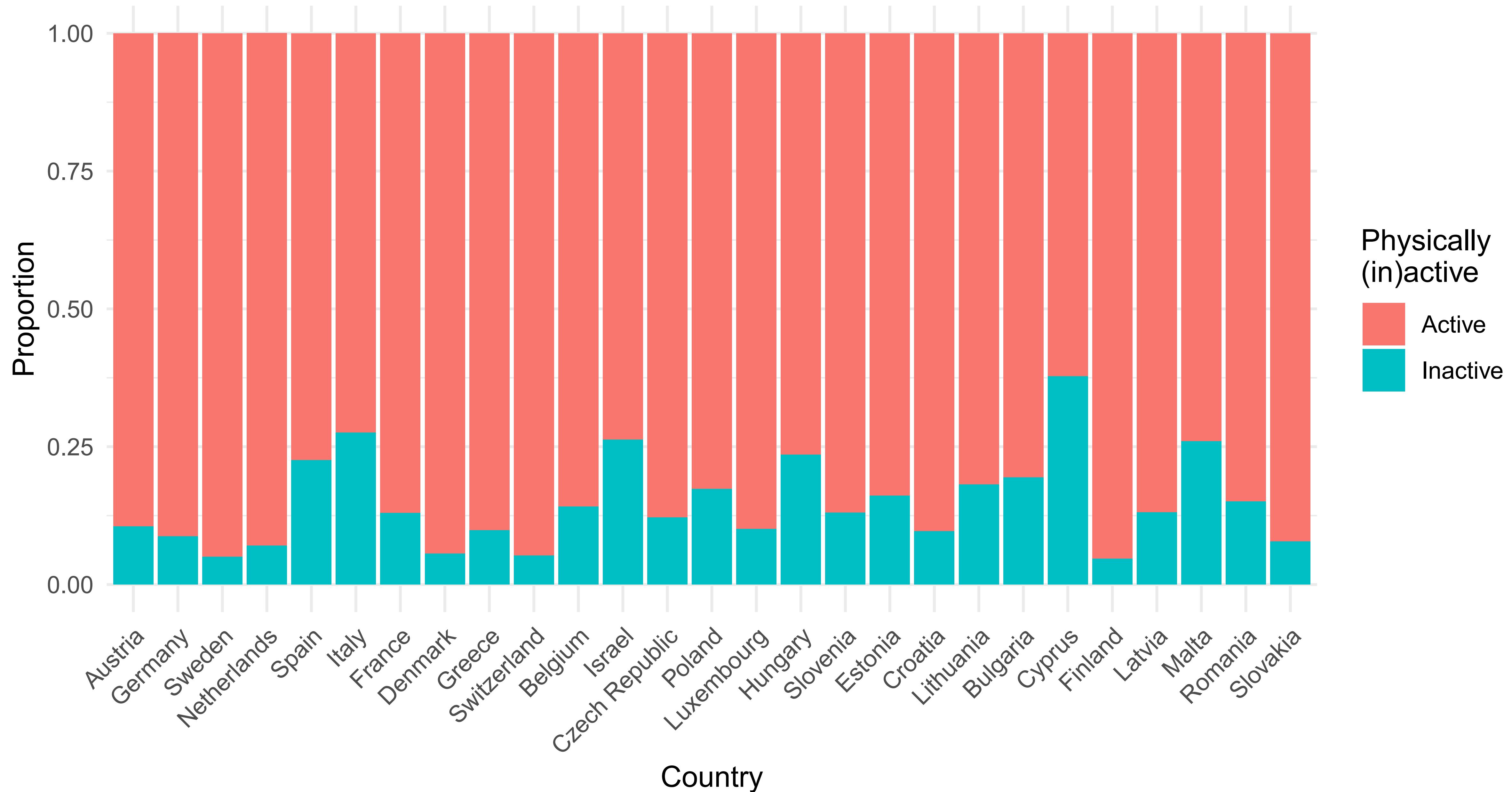




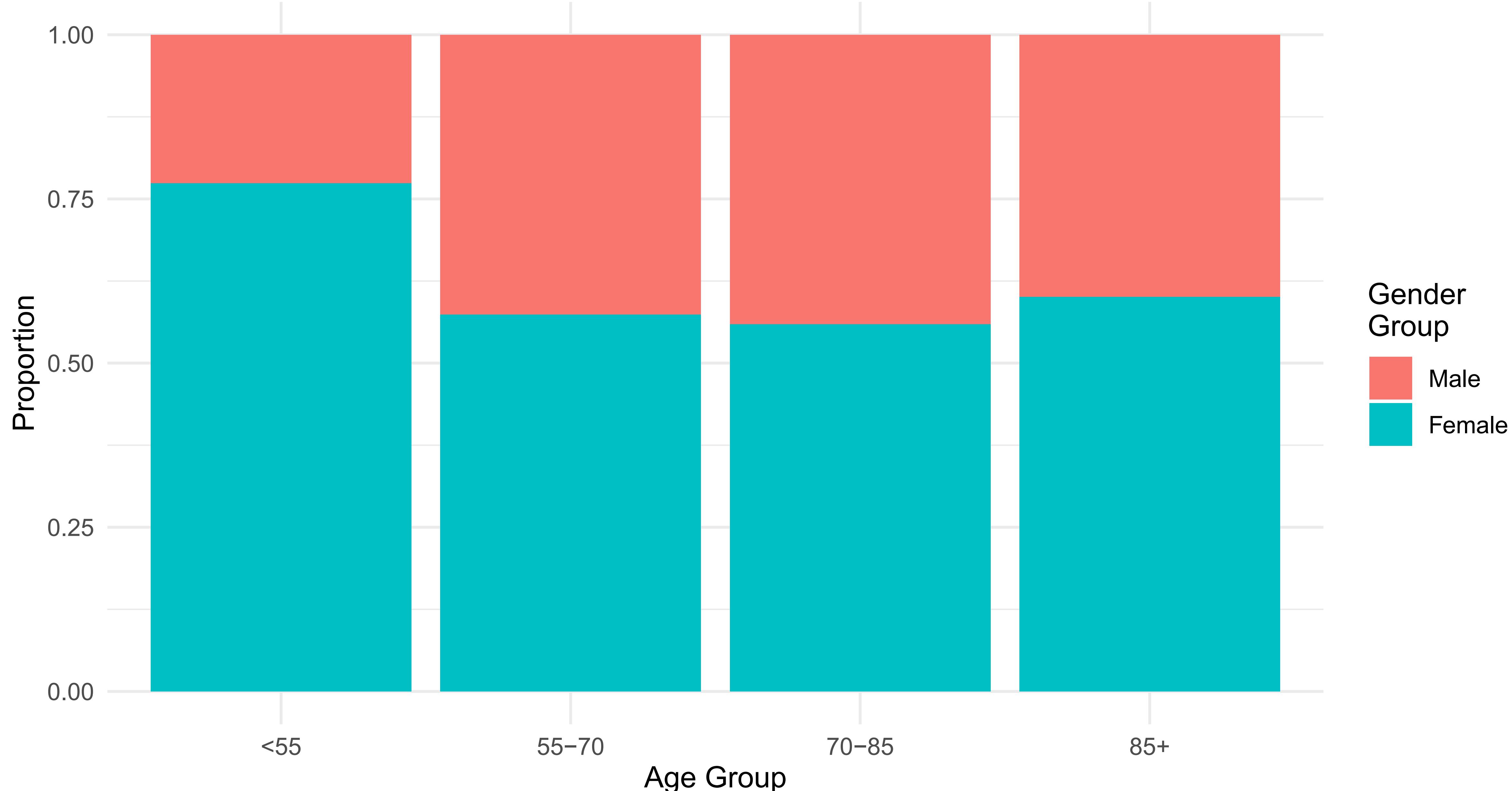
2.628

3.97

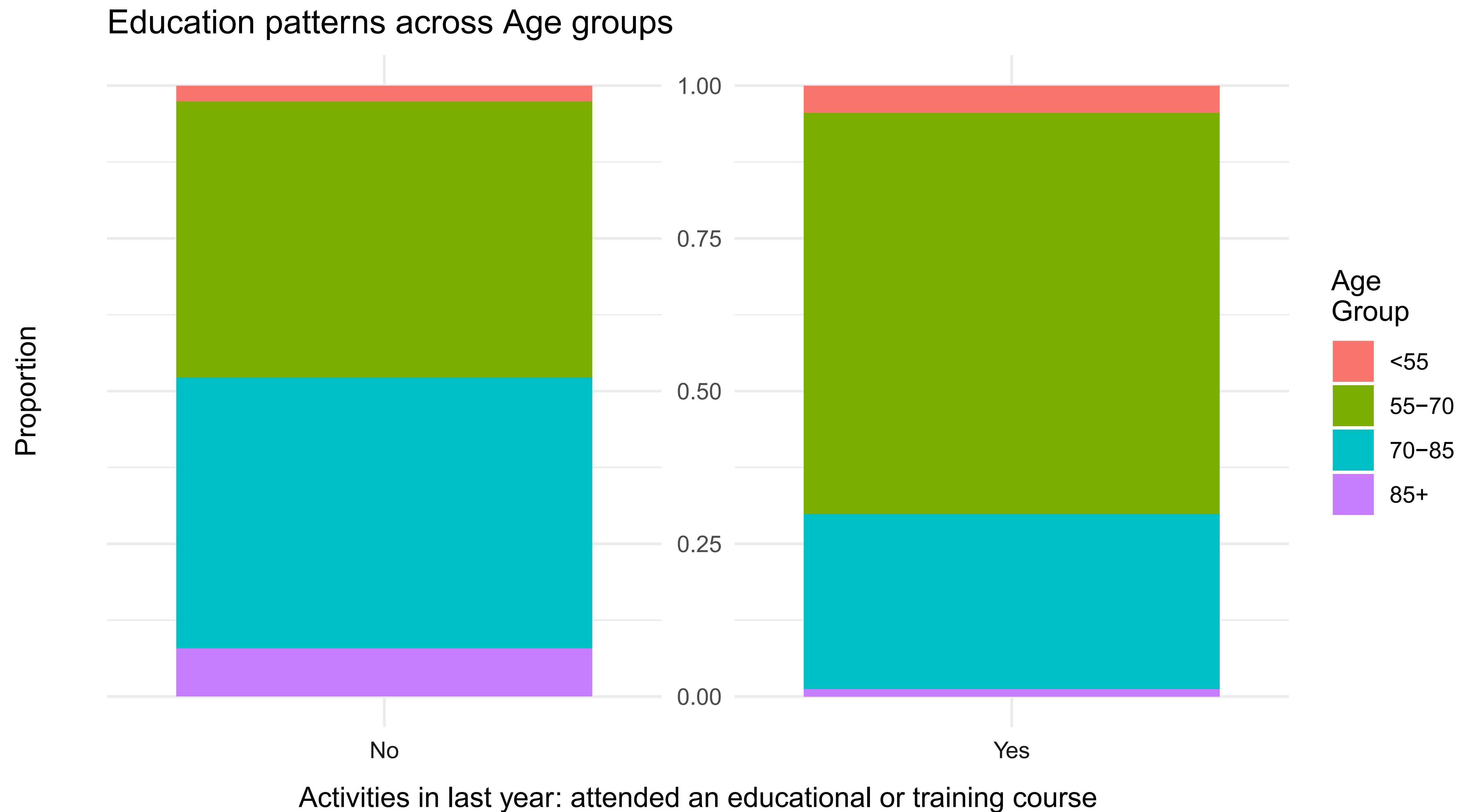
# Physical activity by Country



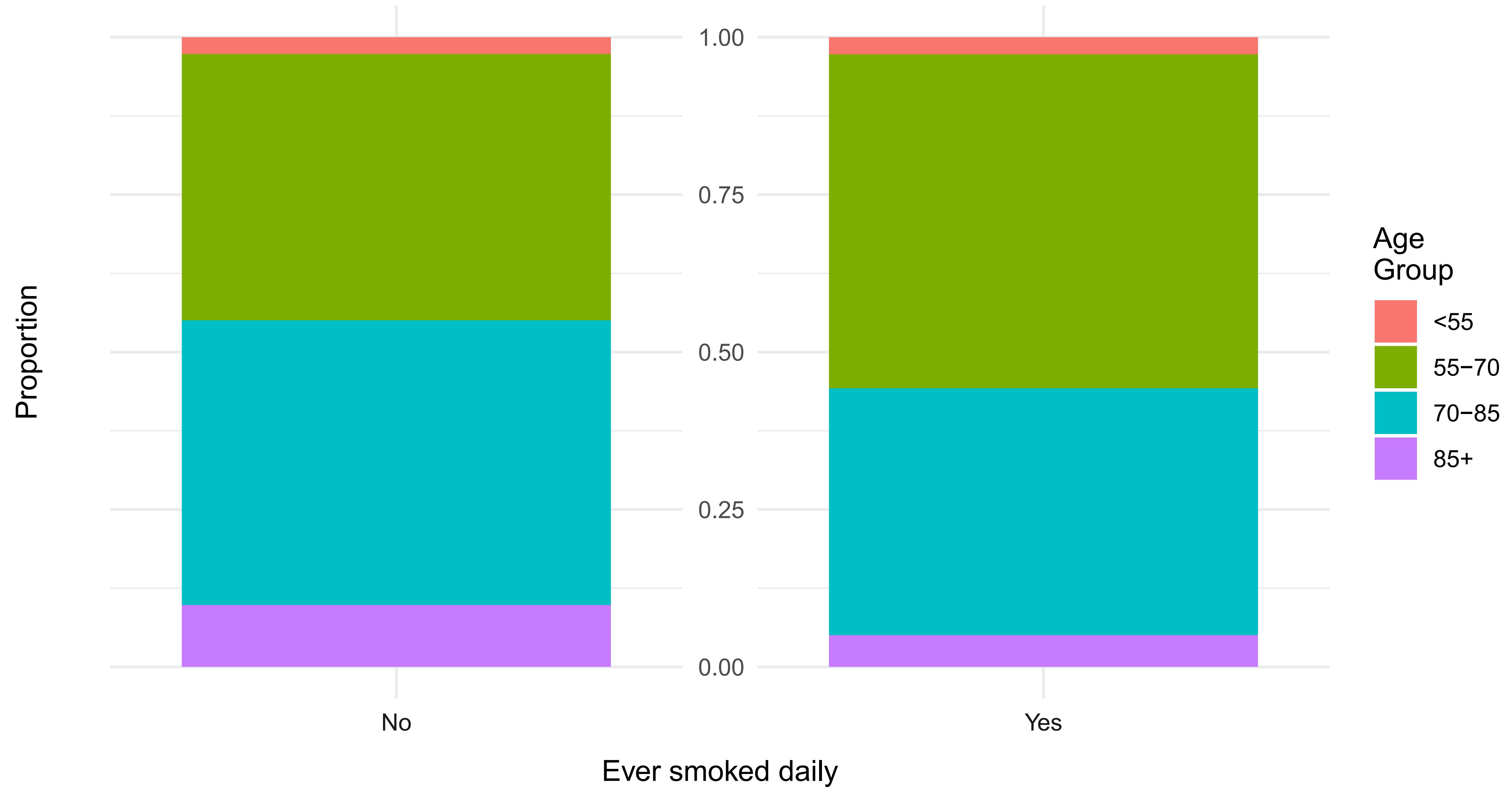
# Female survival advantage by Age group



# Education patterns across Age groups



## Smoking by Age group



# Feature Selection

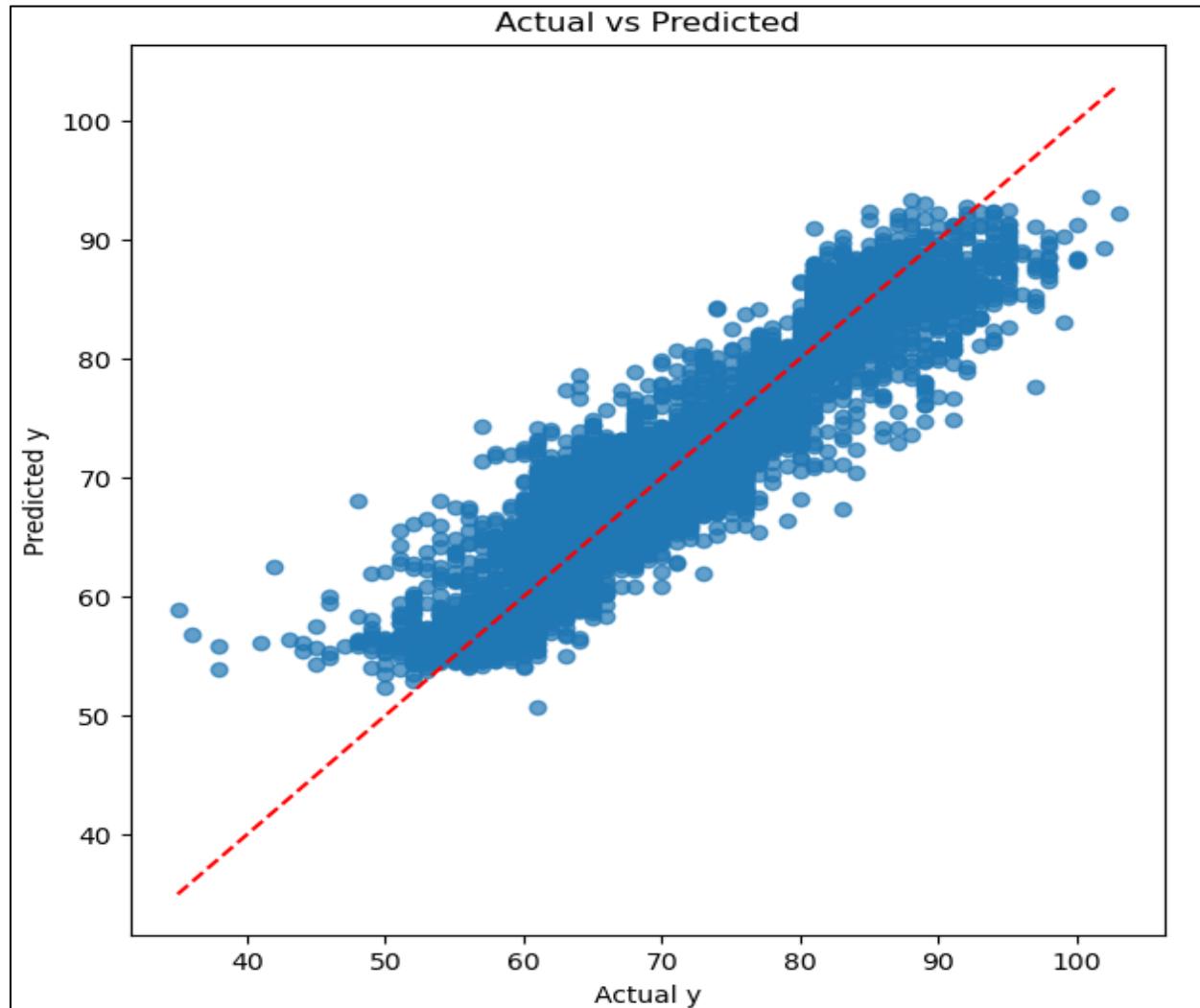
- Remove predefined redundant and multicollinear columns (from two exclusion lists)
- Replace all negative values with NaN (in numeric columns)
- Drop columns with  $\geq 70\%$  missing values
- Impute remaining missing numeric values with the median
- Split features into 10 chunks and compute feature importance using a Decision Tree for each chunk
- Select features



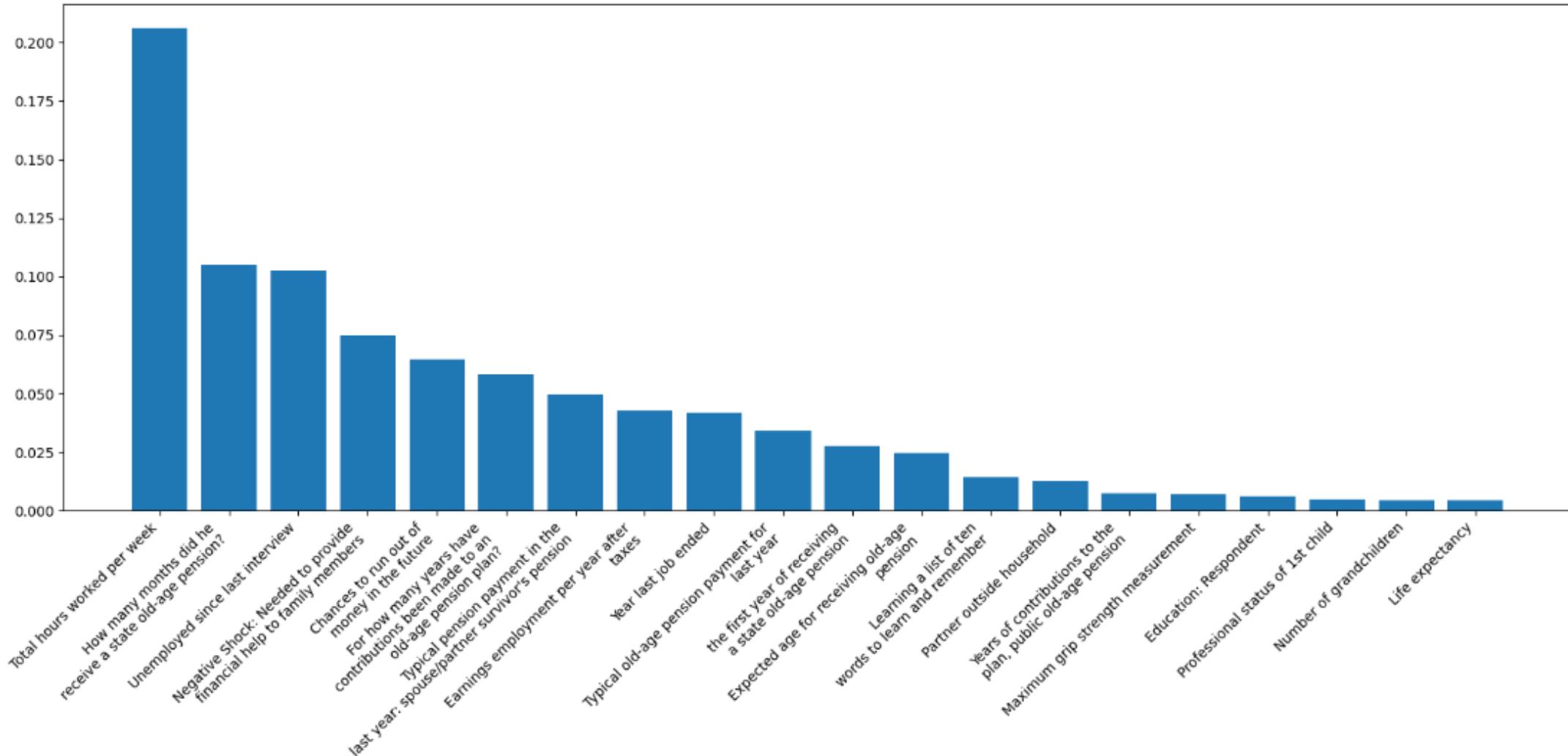
TEST set MAE: **2.534**  
TEST set MAPE: **3.62%**  
TEST set R<sup>2</sup>: **0.871**  
TEST set Adjusted R<sup>2</sup>: **0.868**

TRAIN set MAE: **2.489**  
TRAIN set MAPE: **3.56%**  
TRAIN set R<sup>2</sup>: **0.878**  
TRAIN set Adjusted R<sup>2</sup>: **0.878**

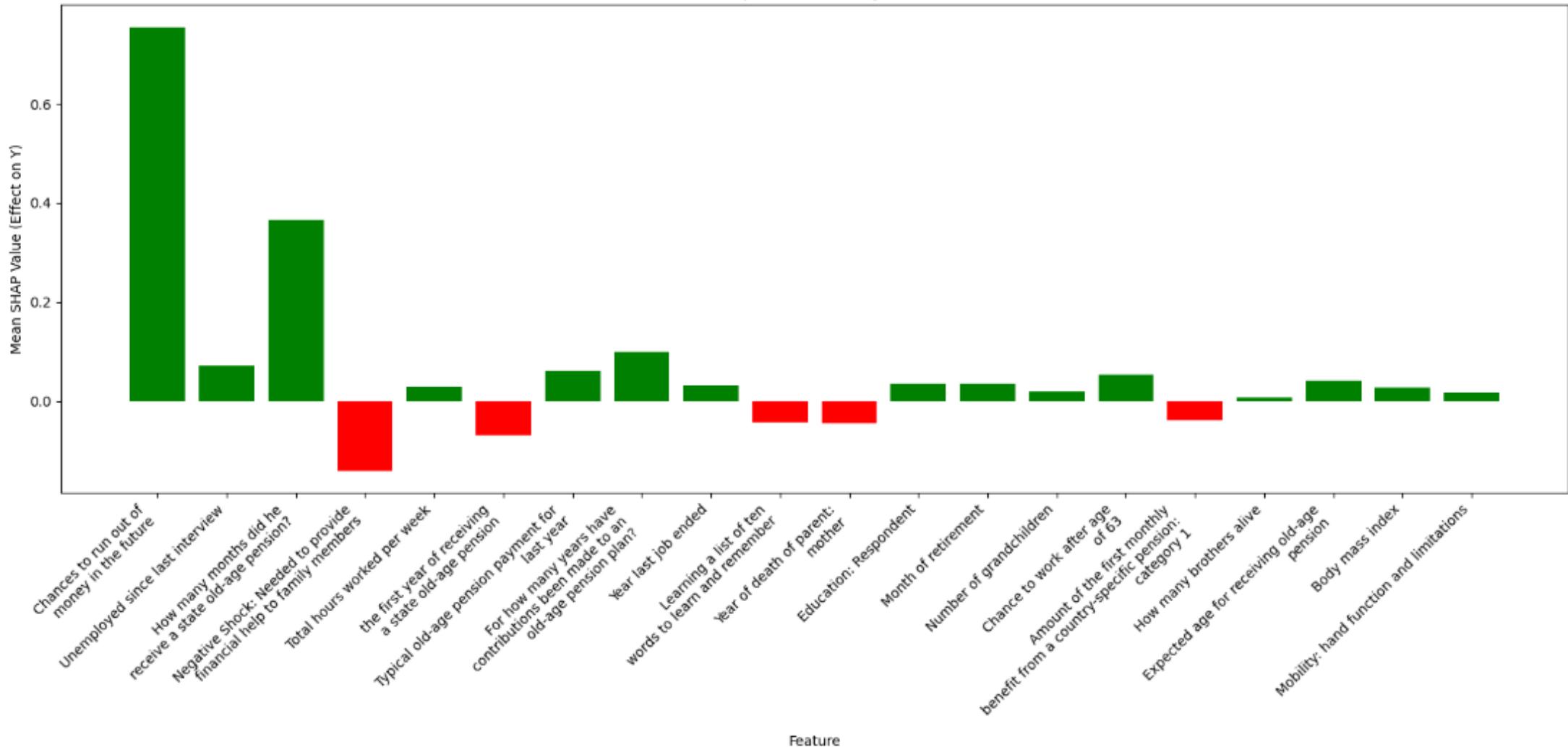
The data is split into 80% training and 20% test sets.



Top 20 XGBoost feature importances



XGBoost: Top 20 Features by SHAP Effect



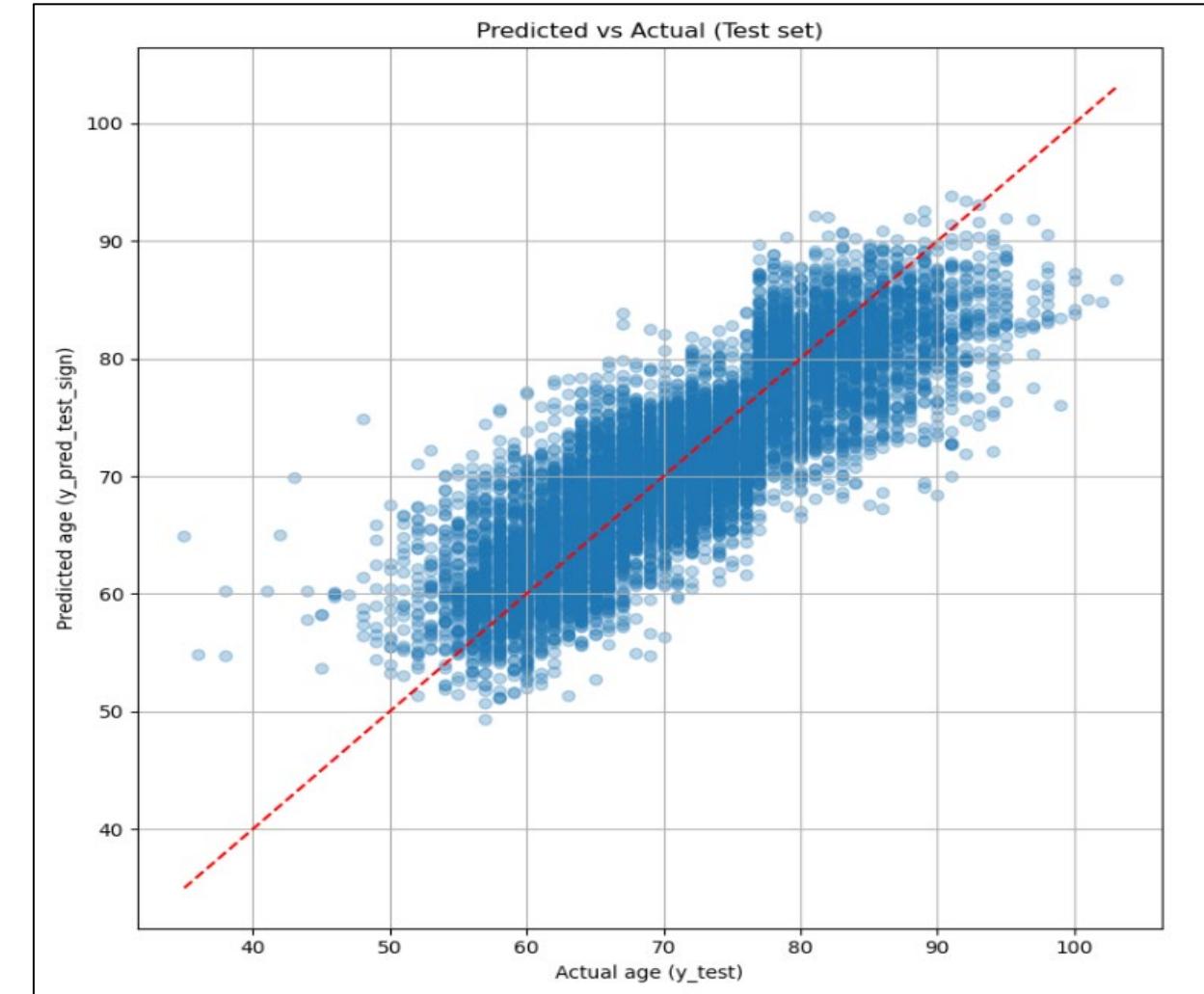
## OLS Regression in R



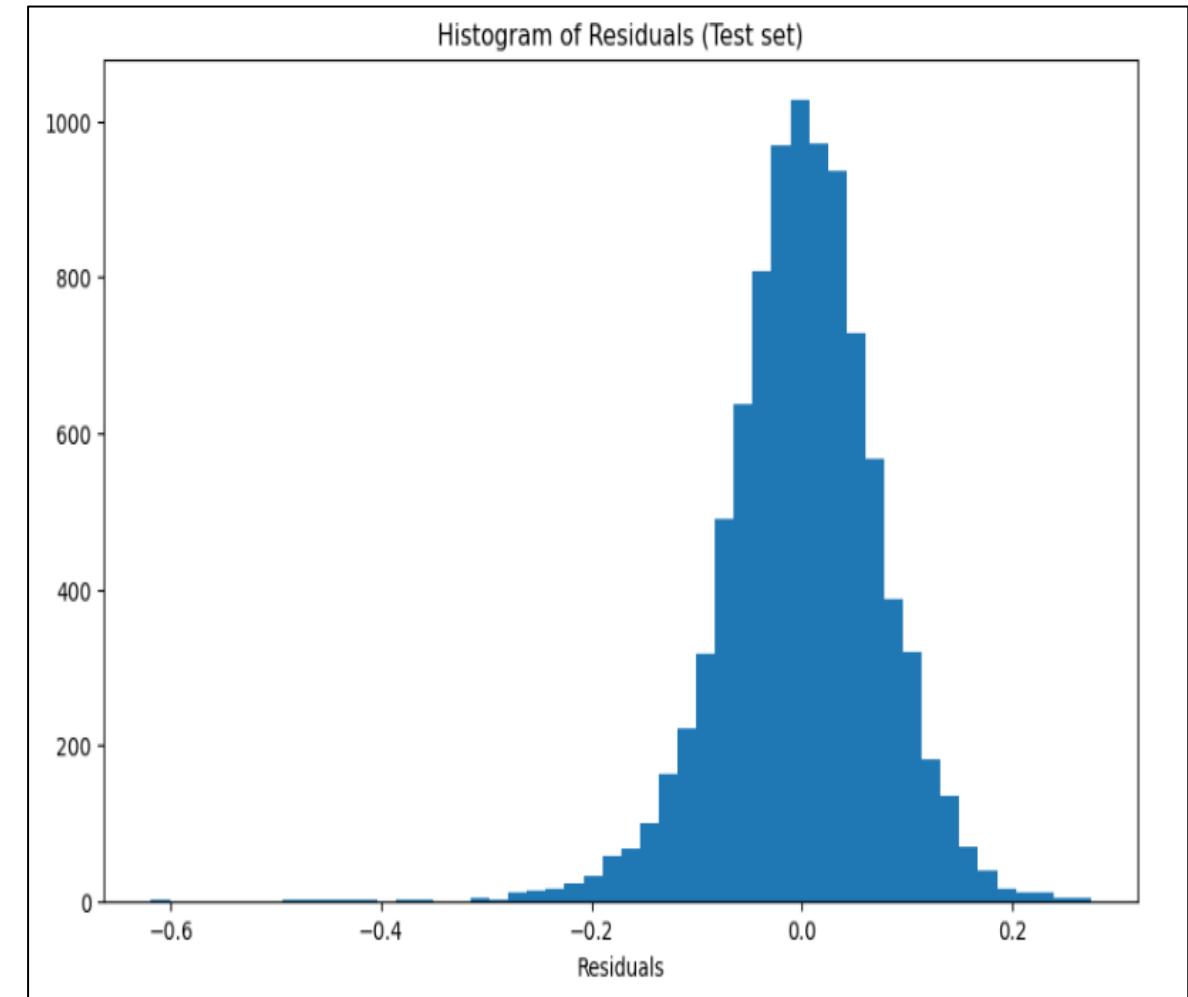
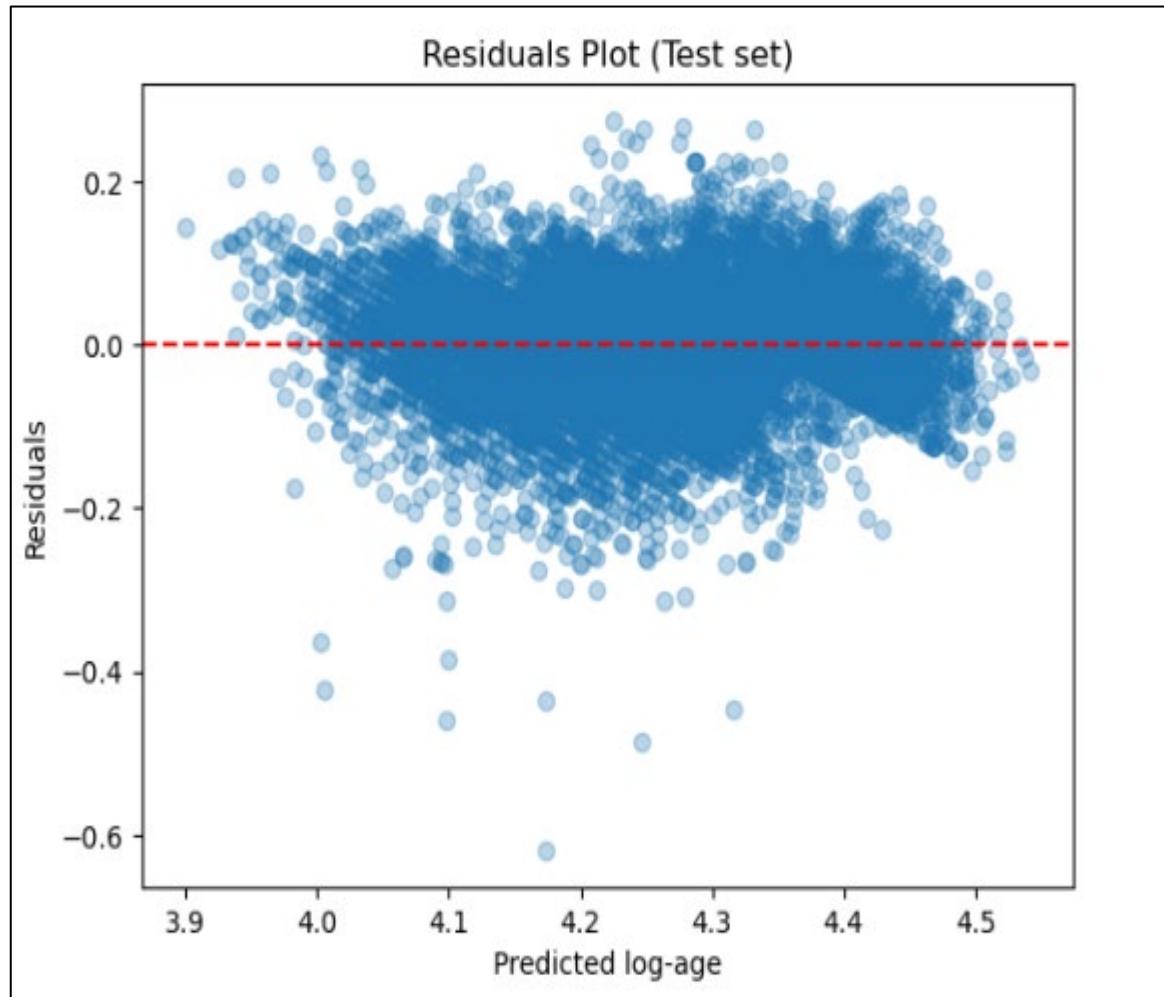
TEST set MAE: **3.916**  
 TEST set MAPE: **13.25%**  
 TEST set MSE: **25.989**  
 TEST set R<sup>2</sup>: **0.708**  
 TEST set Adjusted R<sup>2</sup>: **0.706**

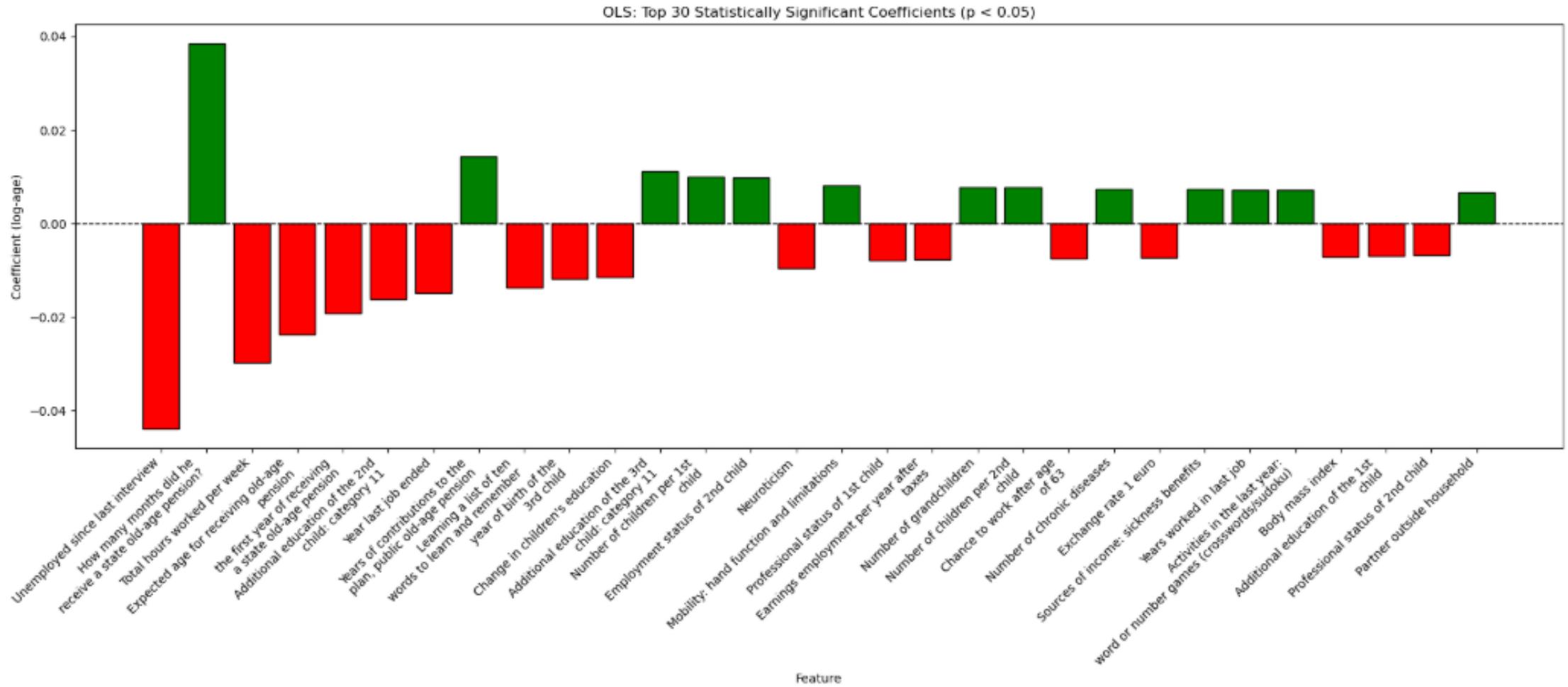
TRAIN set MAE: **3.921**  
 TRAIN set MAPE: **5.60%**  
 TRAIN set MSE: **26.027**  
 TRAIN set R<sup>2</sup>: **0.711**  
 TRAIN set Adjusted R<sup>2</sup>: **0.710**

## OLS Regression: Feature Selection and Collinearity Control

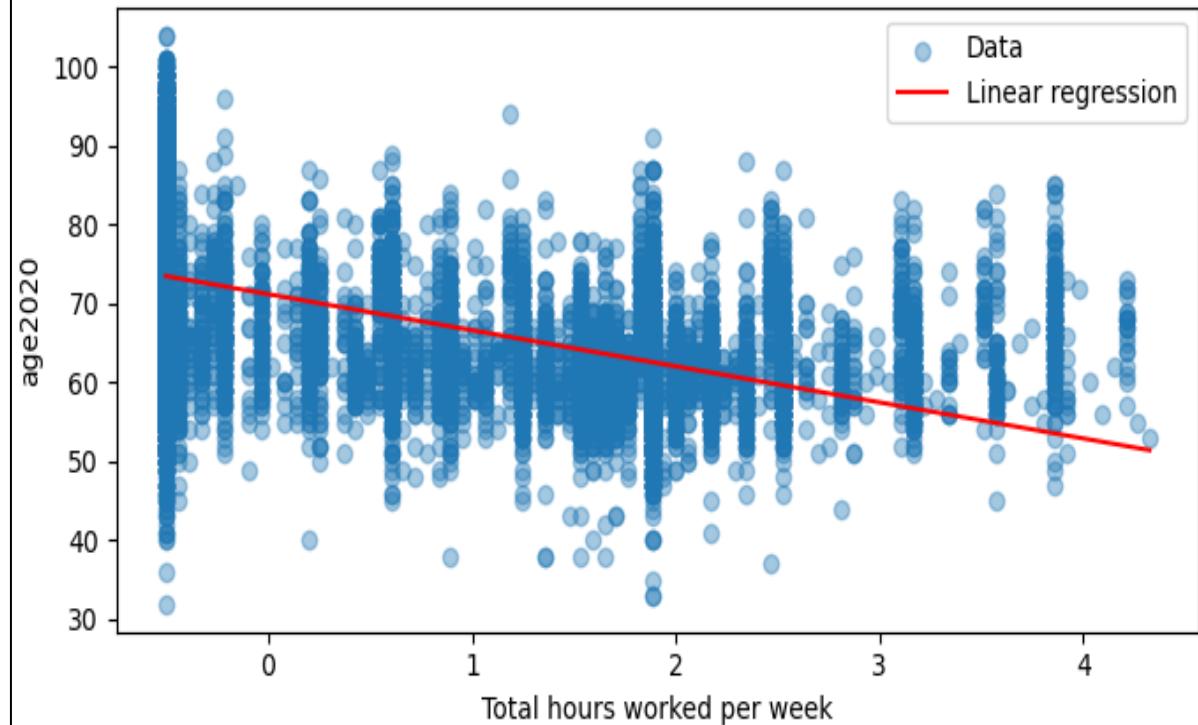


# Residue diagnostics

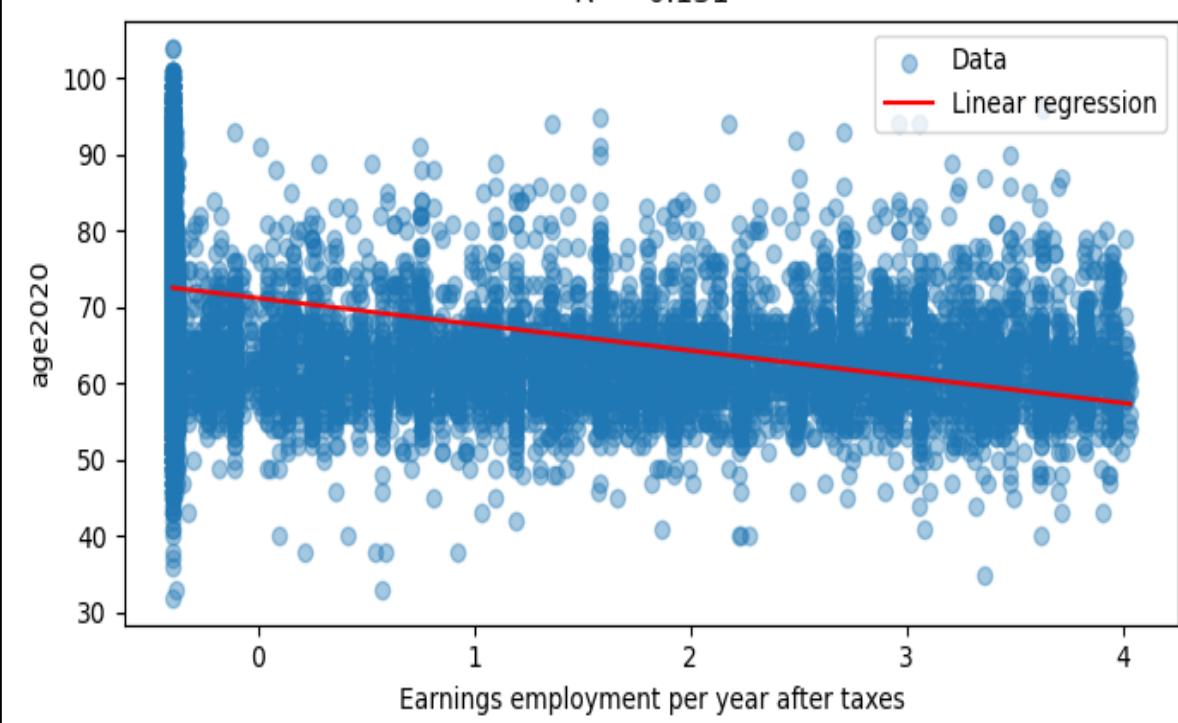




Total hours worked per week  
 $R^2 = 0.232$



Earnings employment per year after taxes  
 $R^2 = 0.131$

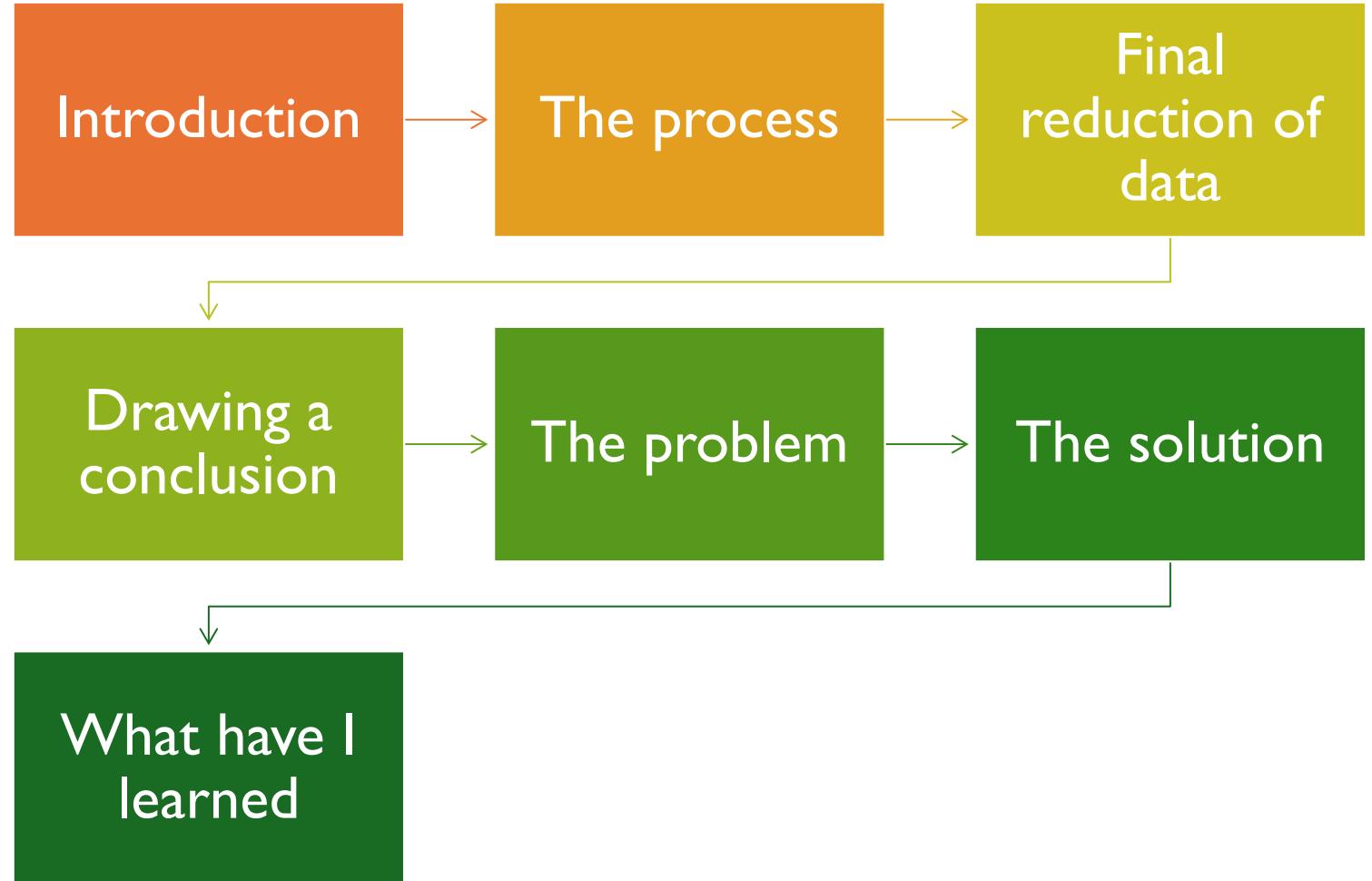


**0** = 20 hours per week  
**1** = 30 hours per week  
**2** = 40 hours per week (full-time)  
**3** = 50 hours per week (overtime)  
**4** = 60 hours per week (heavy workers)

**0** = 20 000  
**1** = 40 000  
**2** = 40 000  
**3** = 50 000  
**4** = 80 000

Variable_down	Effect_down	Label by SHARE	Variable_up	Effect_up	Label by SHARE3
ep325_	down	Unemployed since last interview	ep208_1	up	How many months did he receive a state old-age pension?
ep013_	down	Total hours worked per week	ch016_1	up	Years of contributions to the plan, public old-age pension
cf016tot	down	Learning a list of ten words to learn and remember	ch019_1	up	Number of children per 1st child
bfi10_neuro	down	Neuroticism	chronicw8c	up	Number of chronic diseases
bfi10_extra	down	Extraversion	co002ub	up	Amount spent on food at home
ep205e	down	Earnings employment per year after taxes	ep328_	up	Month of retirement
te035_	down	Hours spent on paid work	ac035d9	up	Activities in last year: did word or number games (crossword puzzles/Sudoku...)
ac035d5	down	Hours spent on paid work	ch019_2	up	Number of children per 2nd child
ac036_9	down	How often did word or number games the last 12 months	ch021_	up	Number of grandchildren
bmi	down	Body mass index	ac036_5	up	How often gone to a sport/social/other kind of club the last 12 months
exrate	down	Exchange rate 1 euro	sn007_2	up	Contact on social network
dn036_	down	How many brothers alive	ep671d7	up	Sources of income: sickness benefits
te005_	down	Hours spent on chores	ex801_	up	Chances to run out of money in the future
ex007_	down	Chance government reduces pension	ch_occupation_3	up	Professional status of the 3rd child
ep671d12	down	Income sources: public long-term care insurance	ep671d6	up	Sources of income: disability pension
maxgrip	down	Maximum grip strength measurement	ep649_	up	Years worked in last job
ex009_	down	Life expectancy	ch_outhh_receive_care_2	up	Outside the household: assistance received from 2nd child
ep152isco	down	ISCO code: respondent's last job	ho034_	up	Years in accommodation
dn037_	down	How many sisters alive	ch_school_education_2	up	School degree of child 2, based on ch017_2 & ch510_2
areabldgi	down	Area of building	ch_closeness_1	up	Emotional closeness with 1st child
as051e	down	Amount selling cars	ep078e_1	up	Typical old-age pension payment for last year
br003_	down	How many years smoked	ep678e	up	Amount received from all occupational pensions last year
br006_	down	Average amount of cigarettes per day	ft009_	up	Received financial gift of 250 or more
hc877_	down	Contacts specialist	ch022_	up	Has great-grandchildren
eurod	down	Depression on the EURO-D scale	isced2011_f	up	Parent's education: Father
numeracy2	down	Math result: subtraction	te047_	up	Hours spent on napping and resting during daytime
hc125_	down	Satisfaction with own coverage in basic health insurance	as070e	up	Interest or dividend income
ac036_10	down	How often played cards or games such as chess the last 12 months	ch_closeness_3	up	Emotional closeness with the 3rd child
ep616isco	down	Respondent's current main job	ep078e_2	up	Typical sickness pension payment for last year

# Outline



# Introduction

# The process



## **First contact with Share**

- Reading the data
- Sampling the data
- Dividing the dataset by feature sets
- Running different models
- Unsatisfactory results



## **New approach**

- Dropping features with more than 50% missing values
- Reducing the data by using feature importance
- Reducing the data by correlation
- Excluding leaky features
- Training logistic regression, random forest and decision tree
- Results: maximum accuracy on random forest: ~ 0.6 accuracy

```
print(f"original df shape: {df.shape}")

df2 = drop_admin(df)
print(f"drop_admin df2 shape: {df2.shape}")

df2 = filter_missing(df2)
print(f"filter_missing df2 shape: {df2.shape}")

df2 = filter_variance(df2)
print(f"filter_variance df2 shape: {df2.shape}")

df2 = drop_high_card(df2)
print(f"drop_high_card df2 shape: {df2.shape}")

df2 = corr_filter(df2)
print(f"corr_filter df2 shape: {df2.shape}")
```

```
model_df shape: (46614, 4)
```

```
Accuracy: 0.42797382816689905
```

	precision	recall	f1-score	support
1	0.00	0.00	0.00	2046
2	0.40	0.49	0.44	3553
3	0.46	0.60	0.52	3724
accuracy			0.43	9323
macro avg	0.28	0.36	0.32	9323
weighted avg	0.33	0.43	0.37	9323

```
y = df_clean['sphus']
X_feats = df_clean.drop(columns=['sphus'])
print(f"X shape before vif_filter: {X_feats.shape}")
X_feats = vif_filter(X_feats)
print(f"X afteer vif_filter: {X_feats.shape}")
print(f"df_clean after vif_filter: {df_clean.shape}")
df_clean = pd.concat([X_feats, y], axis=1)
# 3) Build X/y and train OVR
X = df_clean.select_dtypes(include='number').drop(columns=['sphus'])
y = df_clean['sphus']
print(f"final X shape: {X.shape}")
models = train_ovr(X, y, classes=[1,2,3])
```

**Model: rf**

**Accuracy: 0.9494357918221994**

	<b>precision</b>	<b>recall</b>	<b>f1-score</b>	<b>support</b>
<b>not_1</b>	<b>0.95</b>	<b>0.99</b>	<b>0.97</b>	<b>36384</b>
<b>is_1</b>	<b>0.95</b>	<b>0.82</b>	<b>0.88</b>	<b>10230</b>
<b>accuracy</b>			<b>0.95</b>	<b>46614</b>
<b>macro avg</b>	<b>0.95</b>	<b>0.90</b>	<b>0.92</b>	<b>46614</b>
<b>weighted avg</b>	<b>0.95</b>	<b>0.95</b>	<b>0.95</b>	<b>46614</b>

**Model: xgb**

**Accuracy: 0.8150341099240571**

	<b>precision</b>	<b>recall</b>	<b>f1-score</b>	<b>support</b>
<b>not_1</b>	0.82	0.98	0.89	36384
<b>is_1</b>	0.77	0.23	0.35	10230
<b>accuracy</b>			0.82	46614
<b>macro avg</b>	0.79	0.60	0.62	46614
<b>weighted avg</b>	0.81	0.82	0.77	46614

--- Class 2 vs Rest ---

**Model: log**

**Accuracy: 0.5754708885742481**

	<b>precision</b>	<b>recall</b>	<b>f1-score</b>	<b>support</b>
<b>not_1</b>	0.85	0.56	0.67	36384
<b>is_1</b>	0.29	0.64	0.40	10230
<b>accuracy</b>			0.58	46614
<b>macro avg</b>	0.57	0.60	0.53	46614
<b>weighted avg</b>	0.72	0.58	0.61	46614



# Reading the cookbook

## Final reduction of data

- Choosing about 1000 features, carefully evading leaky ones
- Processing these features through a pipeline
- Extracting the feature importance of the ~400 features left
- Grouping the features in 11 categories
- Training the model using different combinations of the groups

## Drawing conclusions

1. Mobility limitation ( mobility : 0.0227)
2. Chronic conditions ( chronicw8c : 0.0227)
3. Maximum grip strength ( maxgrip : 0.0183)
4. Numeracy-2 test ( numeracy2 : 0.0082)
5. Transfer variable TE-026 ( te026\_ : 0.0091)



# Health and physical

Functional limitation and chronic diseases are the prime factor of change in self perceived health

## Psychological & Cognitive

- Broader well-being (CASP) matters more than raw memory scores, though numeracy and loneliness also play roles.



# Less influational domains

These domains had very small individual importances

# The problem

```
== All features model ==
Accuracy: 0.6433551431942508
      precision    recall  f1-score   support
          1       0.66     0.47     0.55     2046
          2       0.55     0.65     0.59     3553
          3       0.75     0.73     0.74     3724

      accuracy                           0.64     9323
    macro avg       0.65     0.62     0.63     9323
weighted avg       0.65     0.64     0.64     9323
```



The solution

# What have I learned

Without the knowledge, but with a direction

Keeping a journal/log is more important than I initially thought

Avoiding leaky feature

Feature engineering

Dimensionality reduction



# Q & A

**Thank you for  
your attention!**



# Appendix

# List of Countries

- Austria
- Germany
- Sweden
- Netherlands
- Spain
- Italy
- France
- Denmark
- Greece
- Switzerland
- Belgium
- Israel
- Czech Republic
- Poland
- Luxembourg
- Hungary
- Slovenia
- Estonia
- Croatia
- Lithuania
- Bulgaria
- Cyprus
- Finland
- Latvia
- Malta
- Romania
- Slovakia