

# Machine Learning final Project

ME



Eyal Michaeli  
Economics & Data analytics graduate

BIG DATA

# What are we trying to predict?

We are trying to predict the  
**Default of Credit Card Clients -**  
**Classification**



# Our Data set

**Demographic Data**  
Education, age, sex, marriage status



Add Text  
Simple  
PowerPoint  
Presentation



**bill statement – 6 months back**  
Amount of bill statement

**Limit Balance**  
Amount of given credit in NT dollars



**Previous payment – 6 months back**  
Amount of previous payment in each of the last 6 months

**Repayment status – 6 months back**  
Whether a client has payed duly, and if not, by how many months back?



**Default**  
If a client has defaulted

# Motivation and applications for solving the problem

As a basis for credit card Rating



Give Higher interest rate for 'dangerous' clients



Don't give or give less credit



**630-689**

Fair

**690-719**

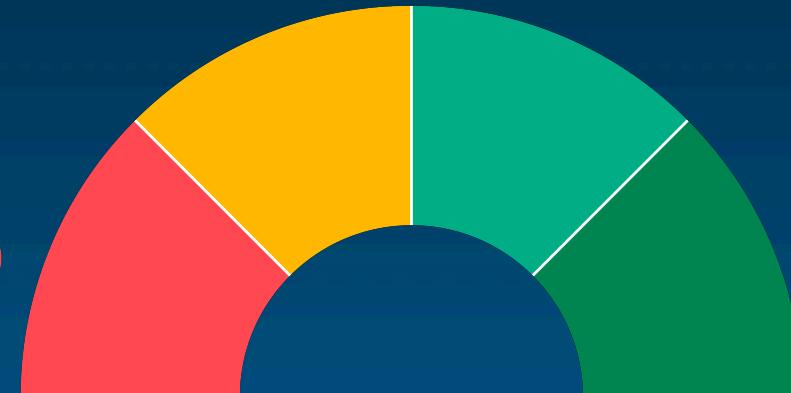
Good

**300-629**

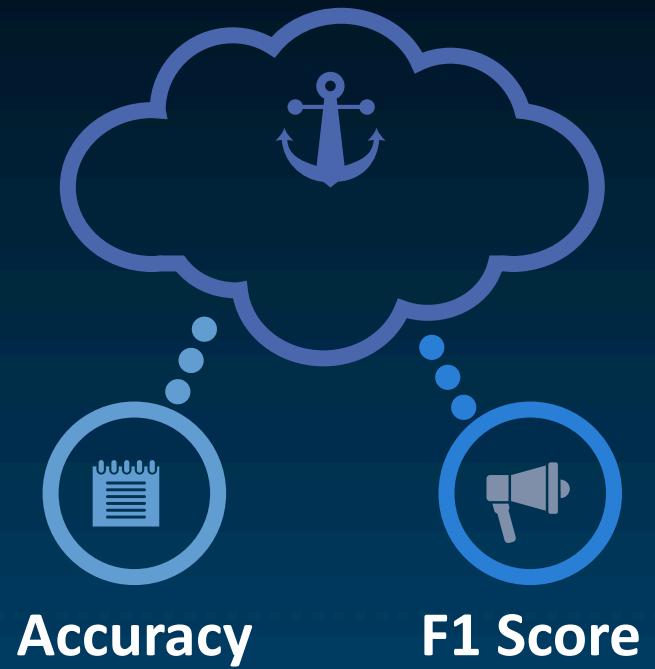
Bad

**720-850**

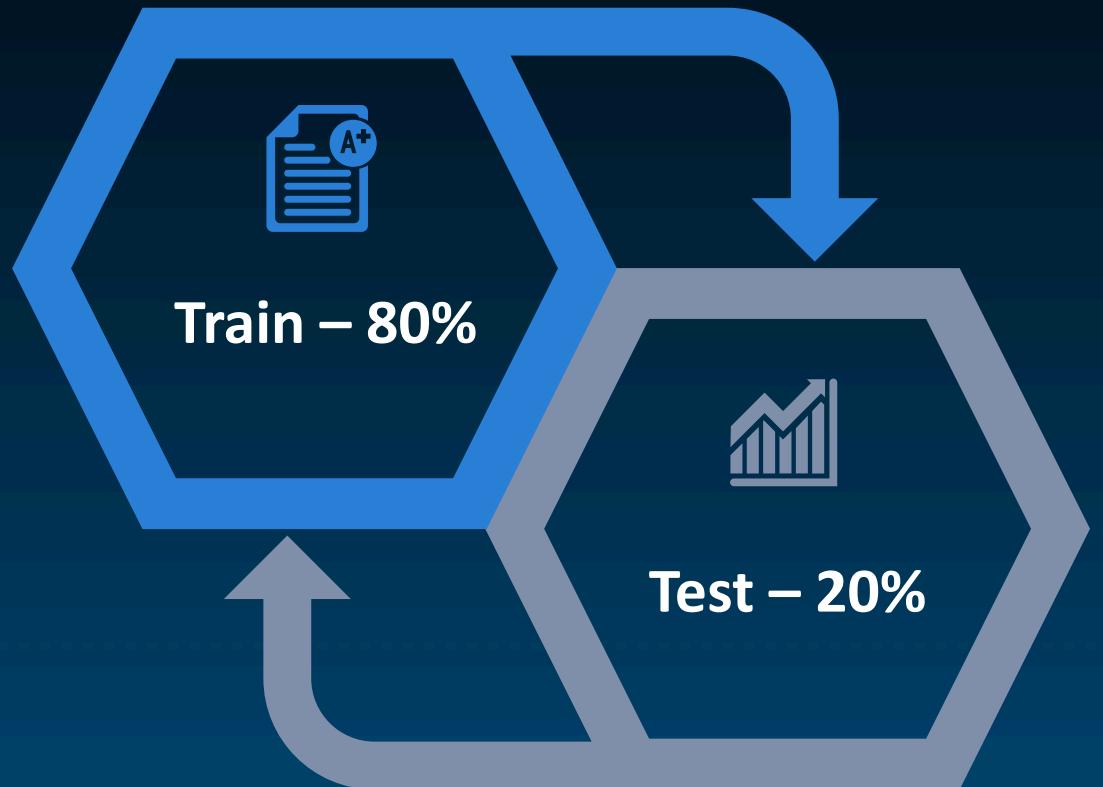
Excellent



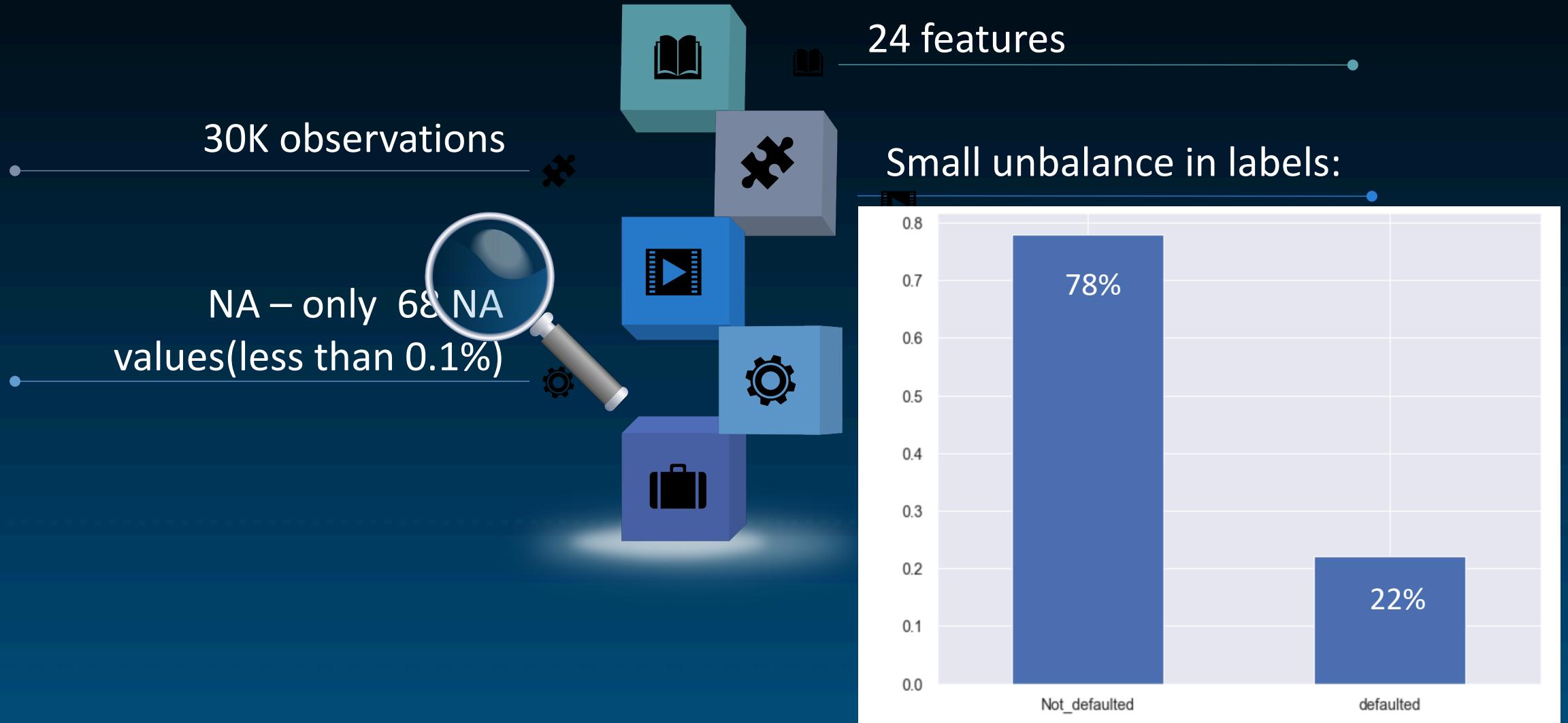
# Evaluation



# Train – test split

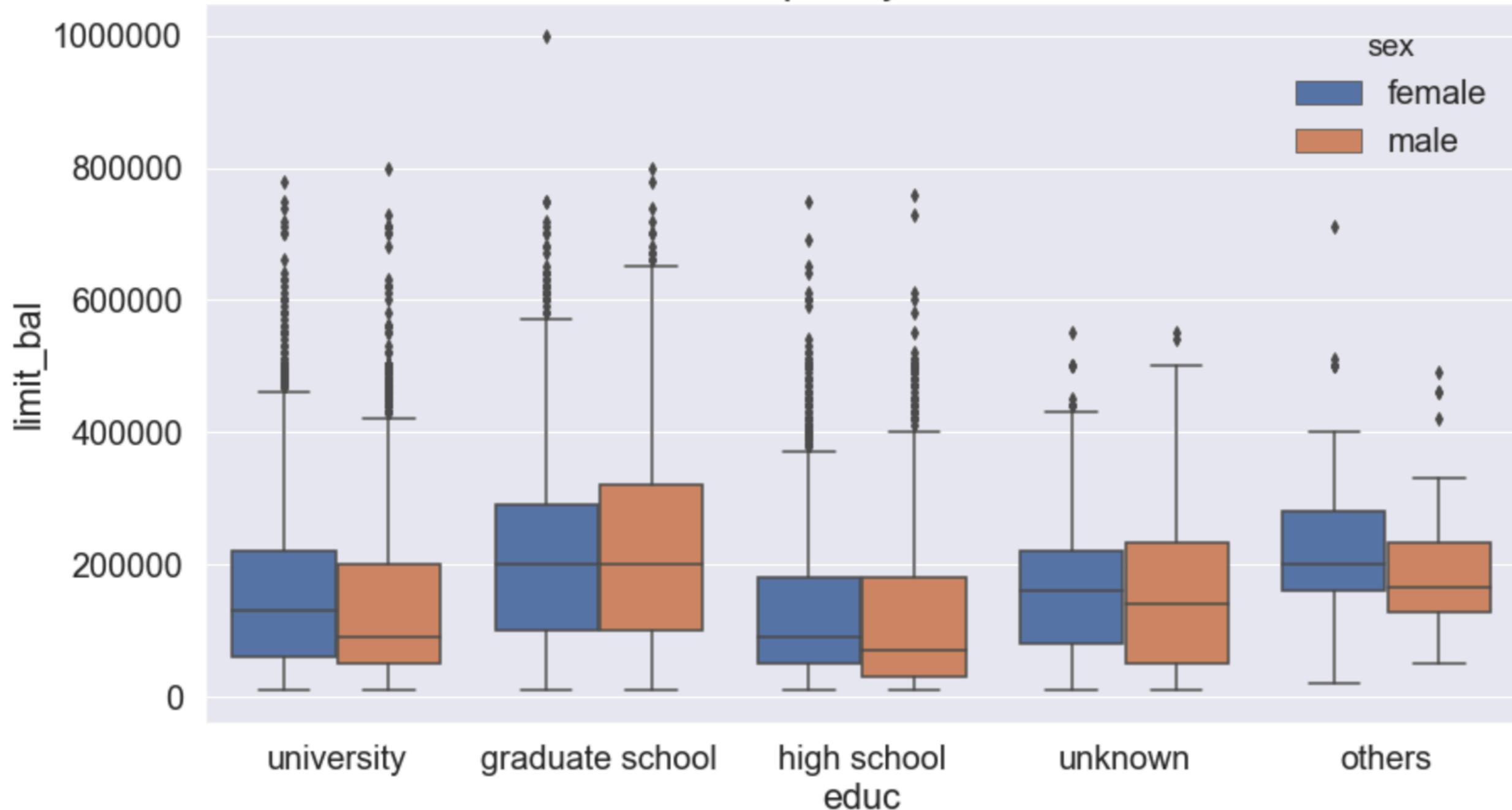


# Data description

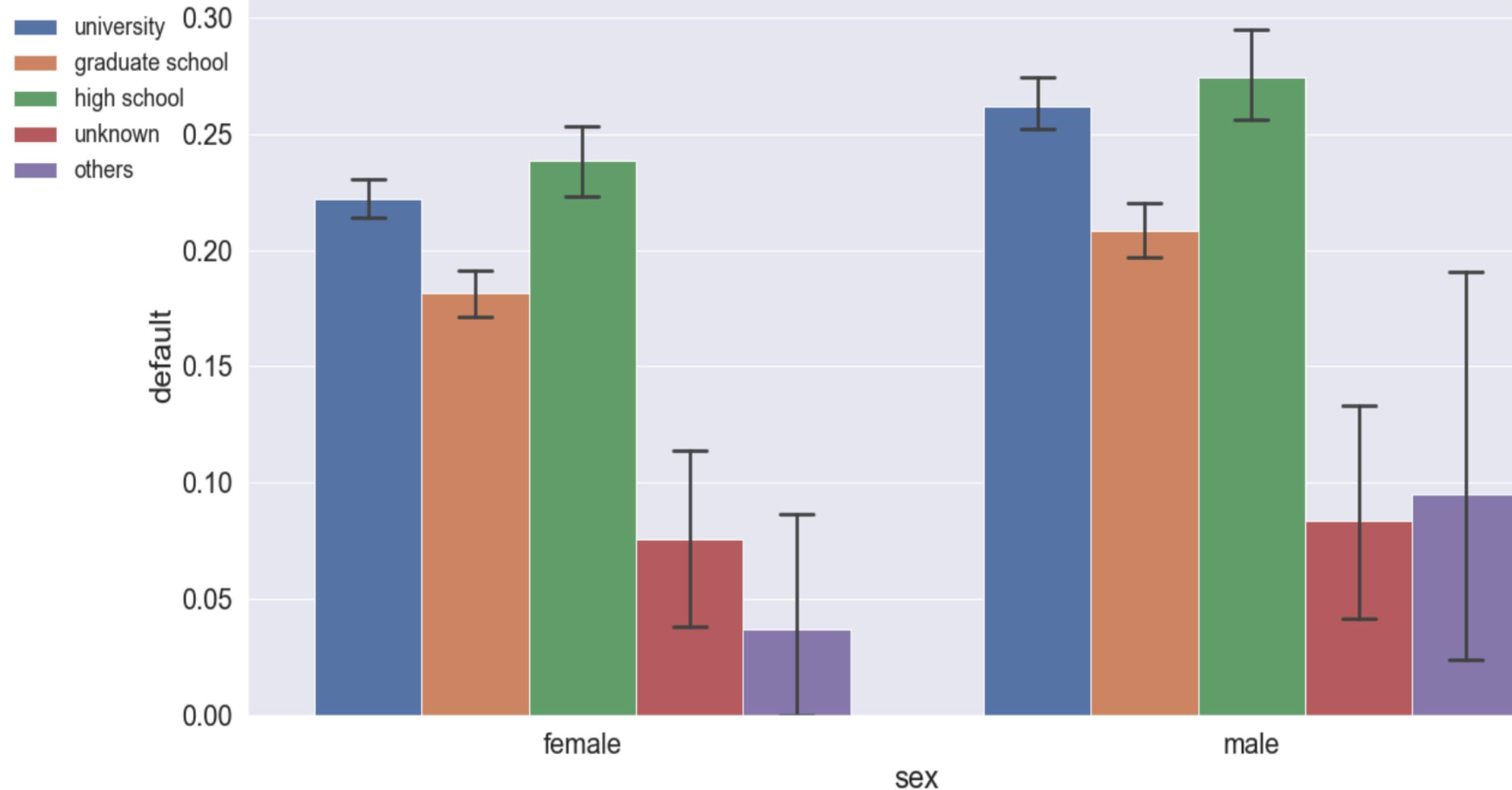


- Graphs describing various aspects of the data

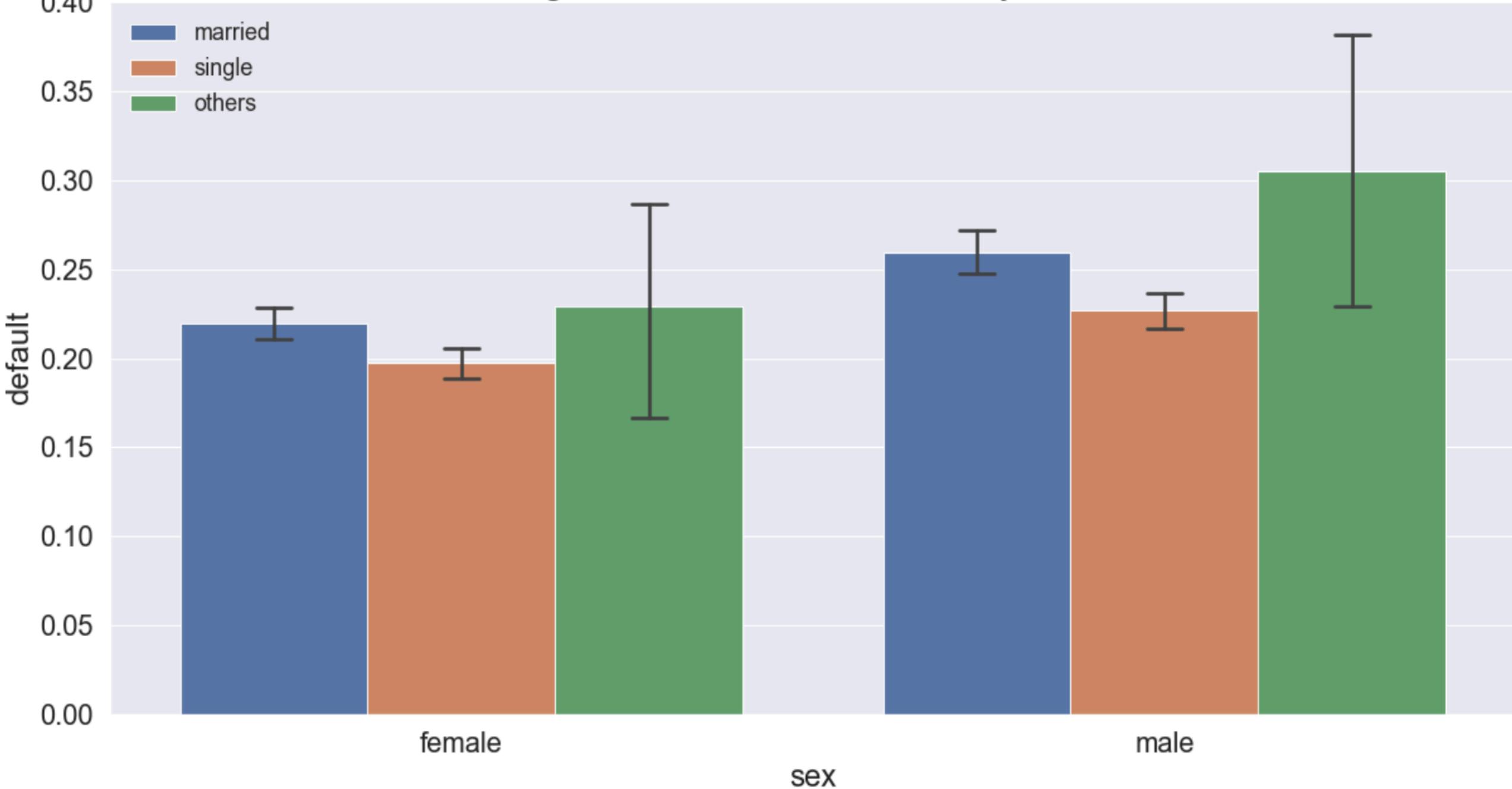
# Limit Balance Box plot by Gender & Education



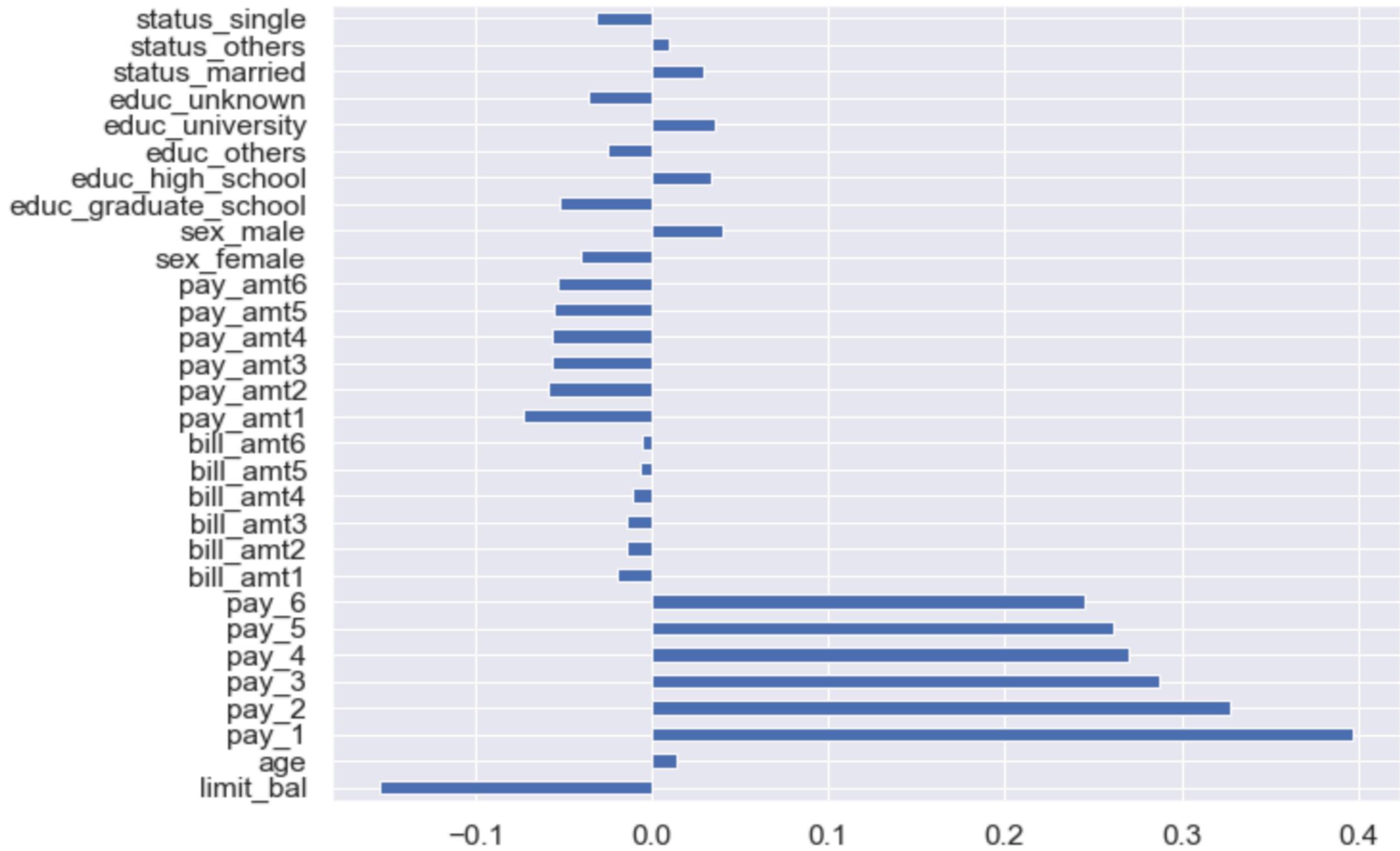
# Default average & confidence intervals by Education level & Gender



# Default average & confidence intervals by Status & Gender



### Default correlation with features



# Data engineering



Did you add any features?



One-Hot-Encoding

Dropped rows with NA values



Removed ID column



# Benchmark algorithm

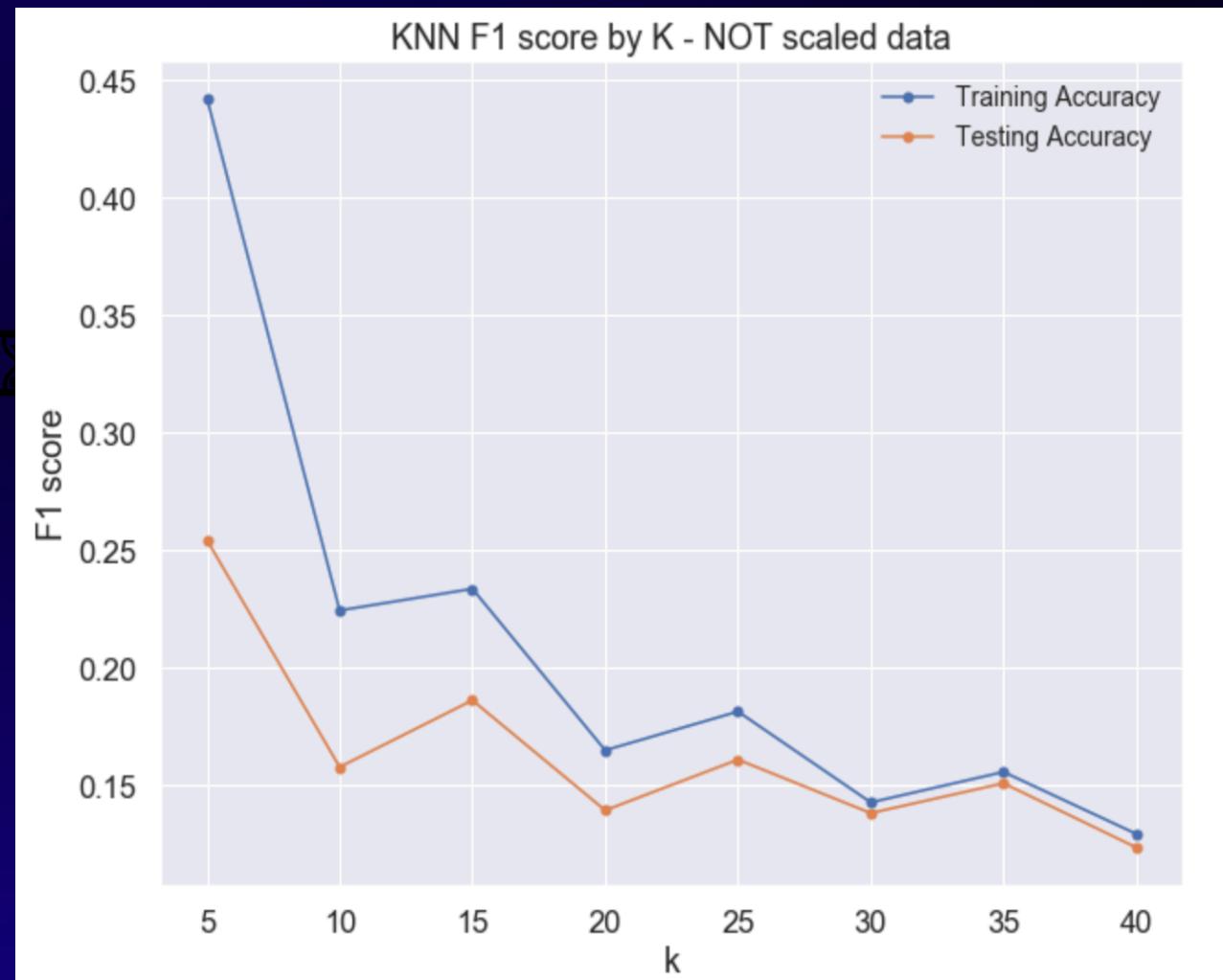
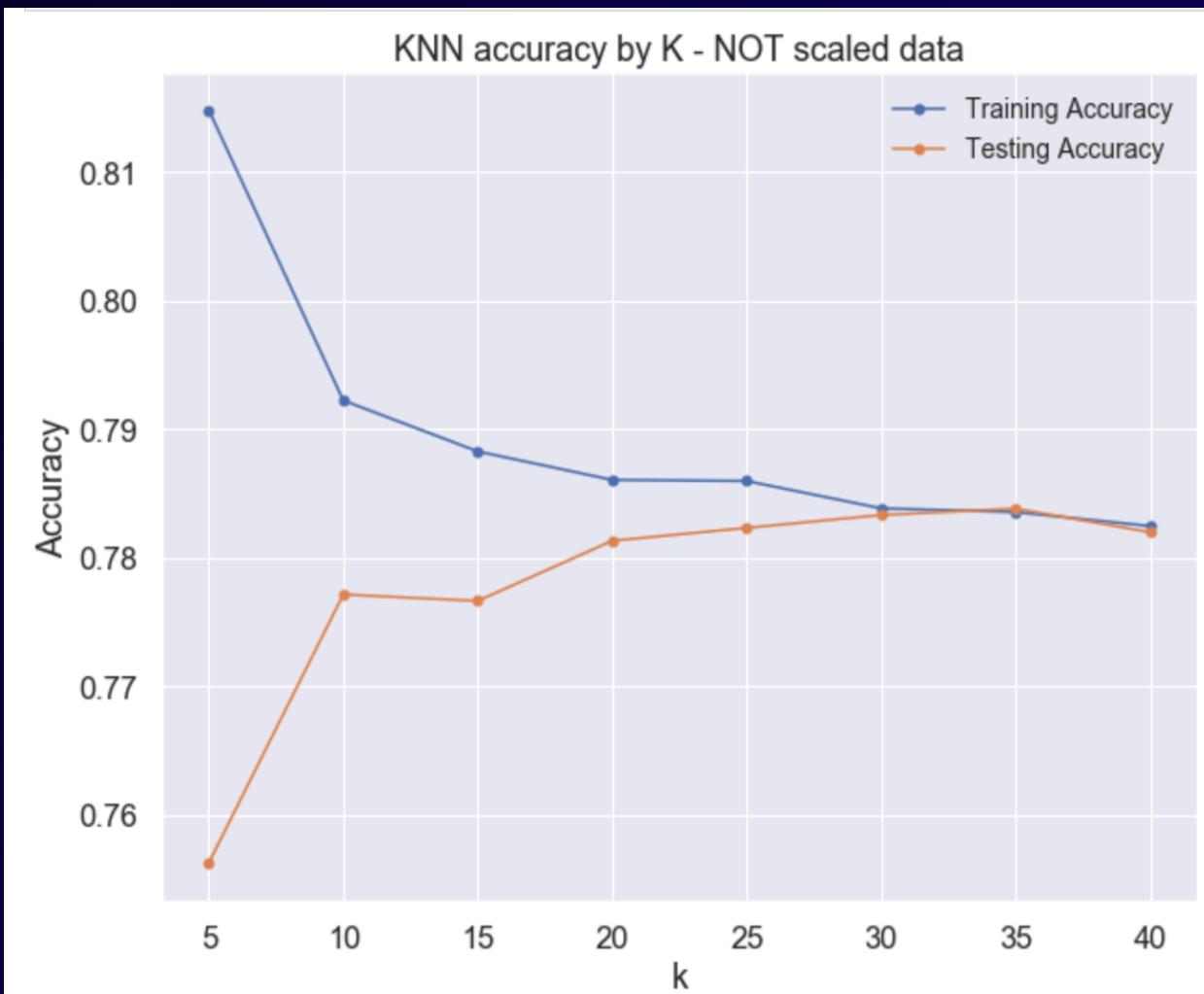
Prediction – always not default



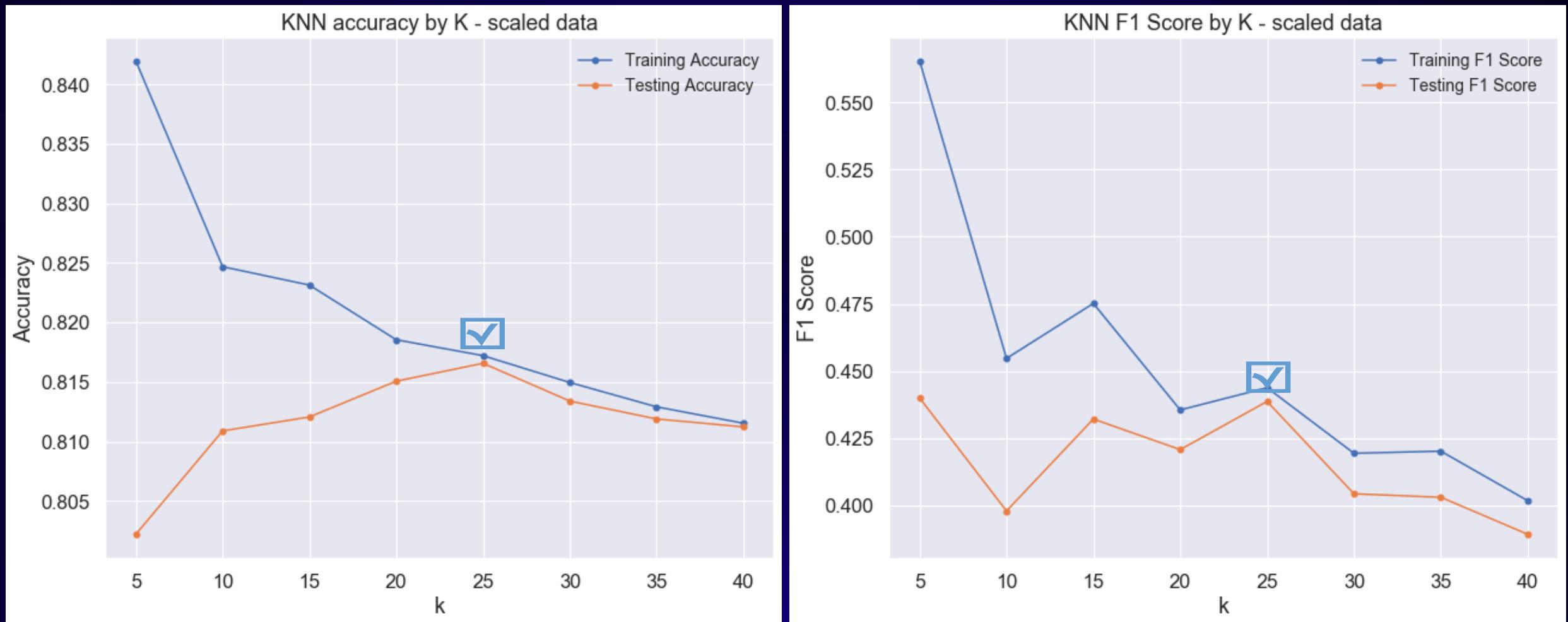
Accuracy of 78%  
F1 score of 0(of course).

We are going to try to beat that!

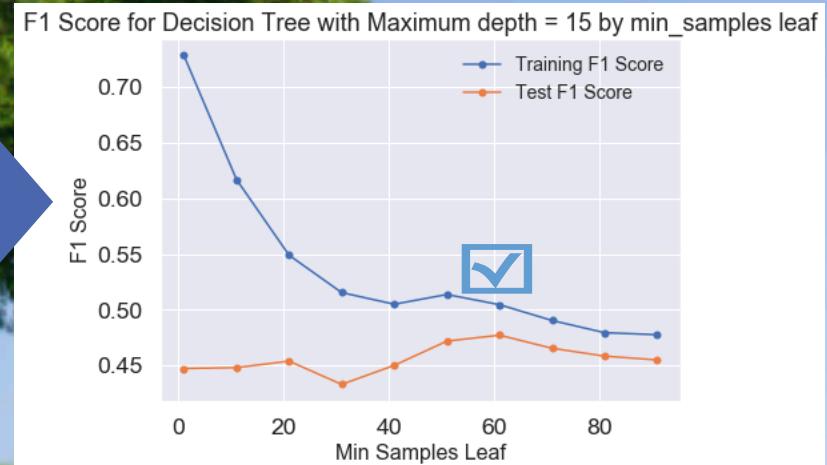
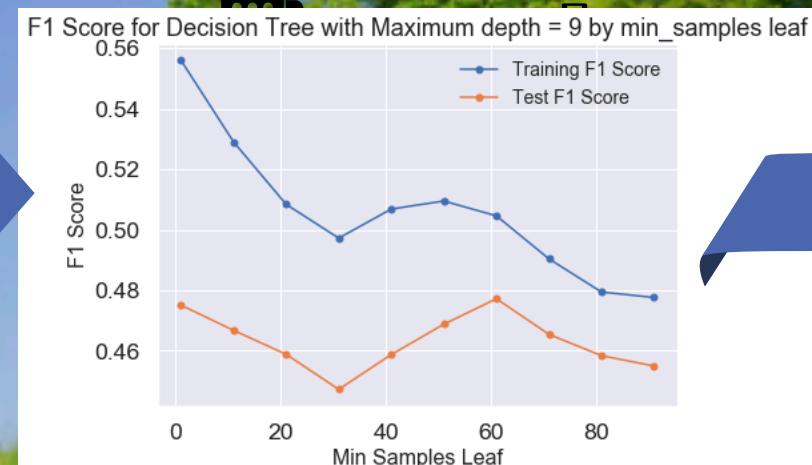
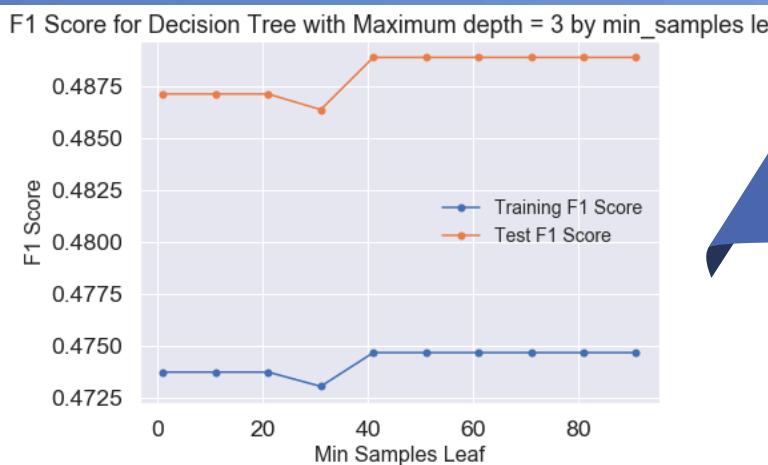
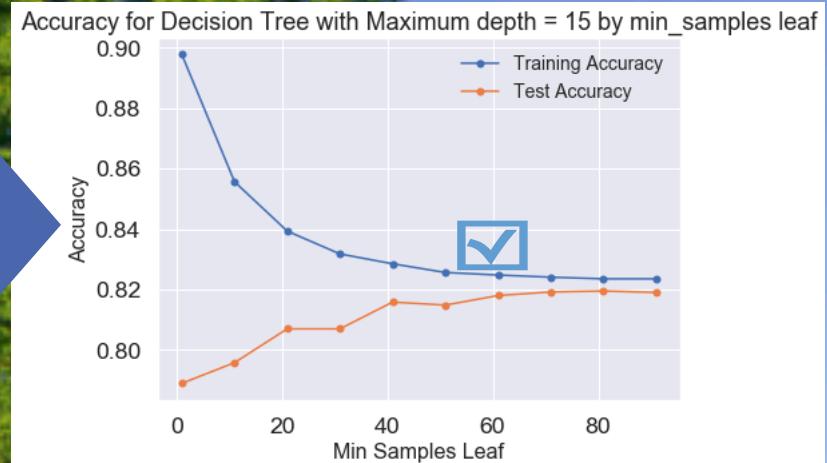
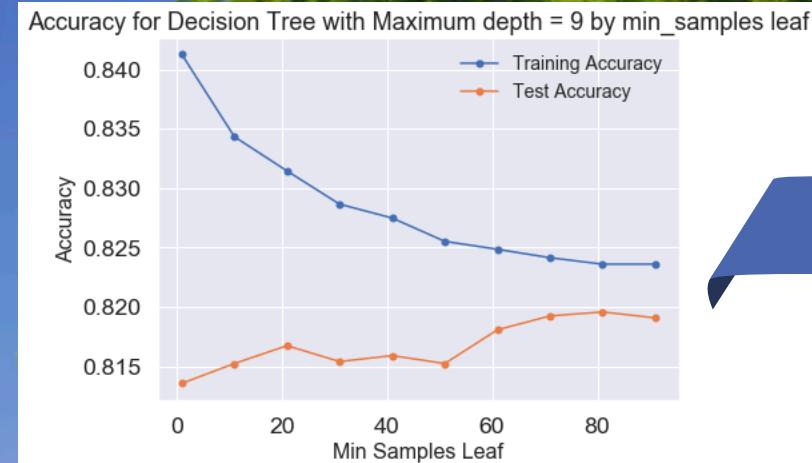
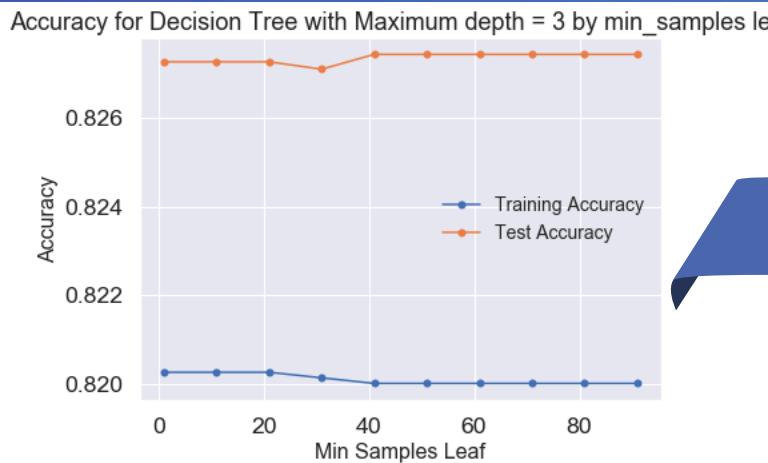
# KNN – NOT scaled data



# KNN –scaled data

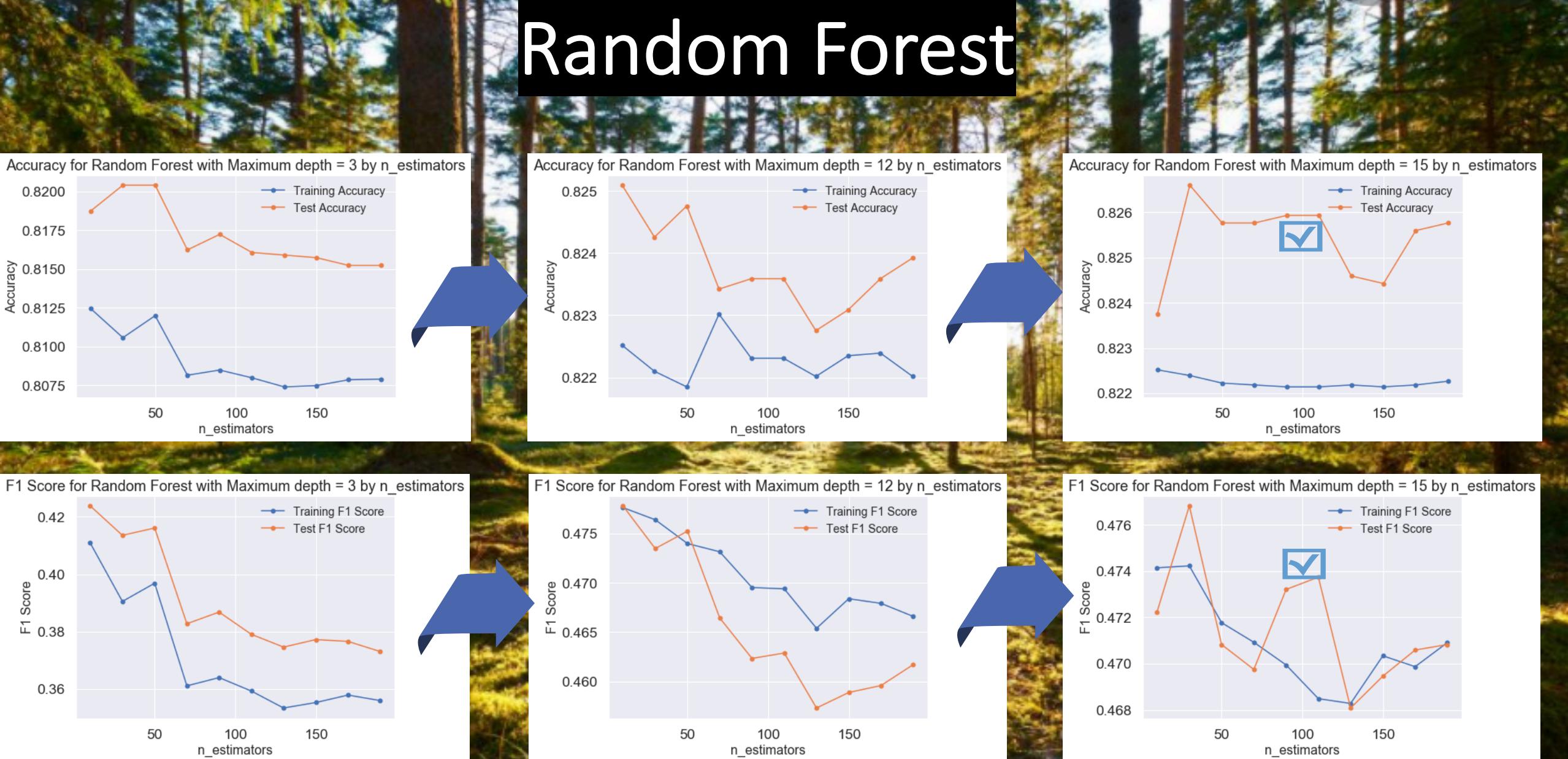


# Decision tree



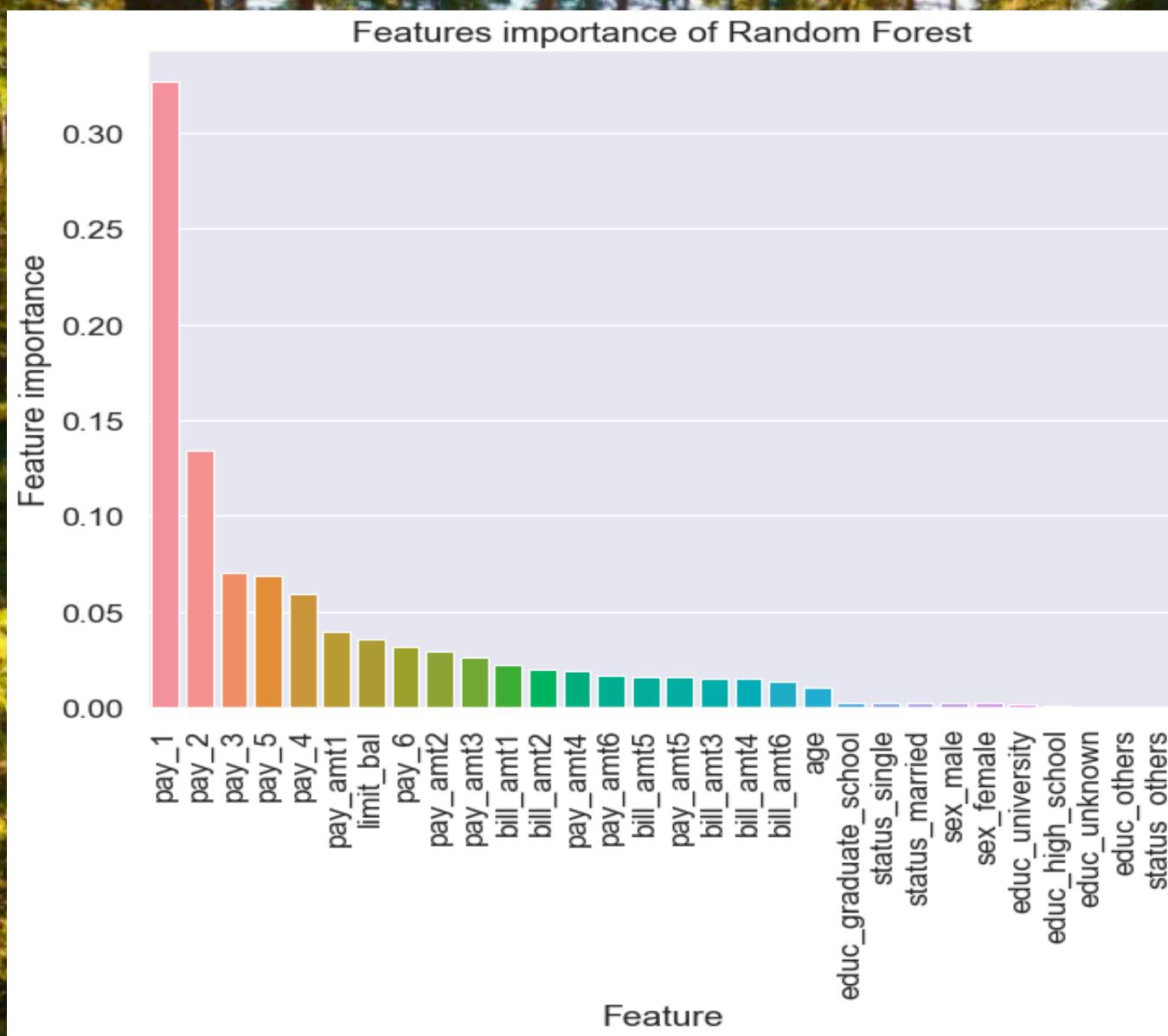
We chose max depth of 15 and Min samples leaf of 60

# Random Forest

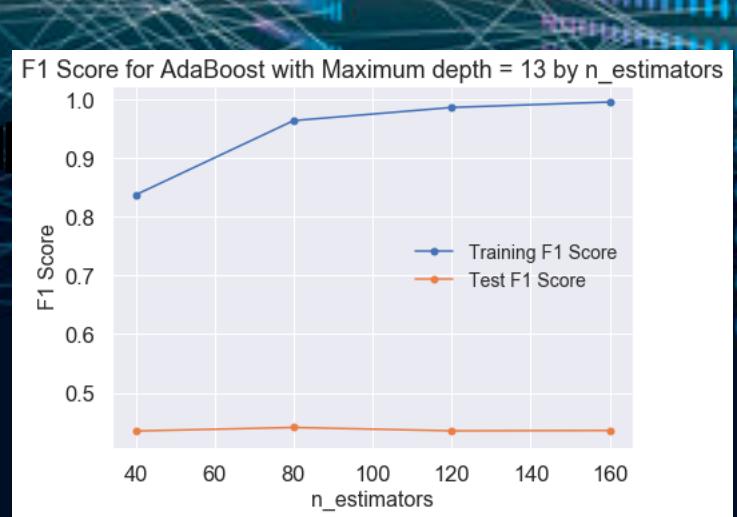
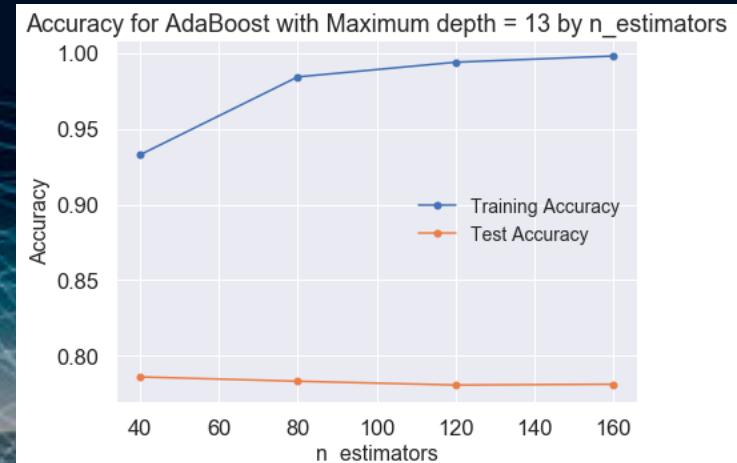
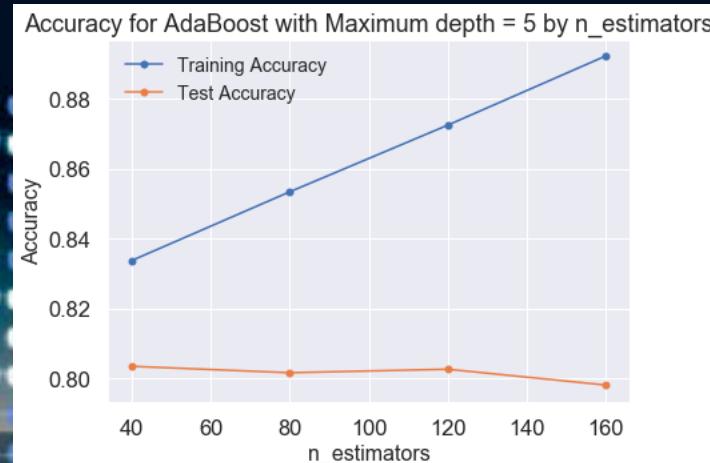


We chose max depth of 15 and Number of trees = 100

# What is the random forest feature importance?

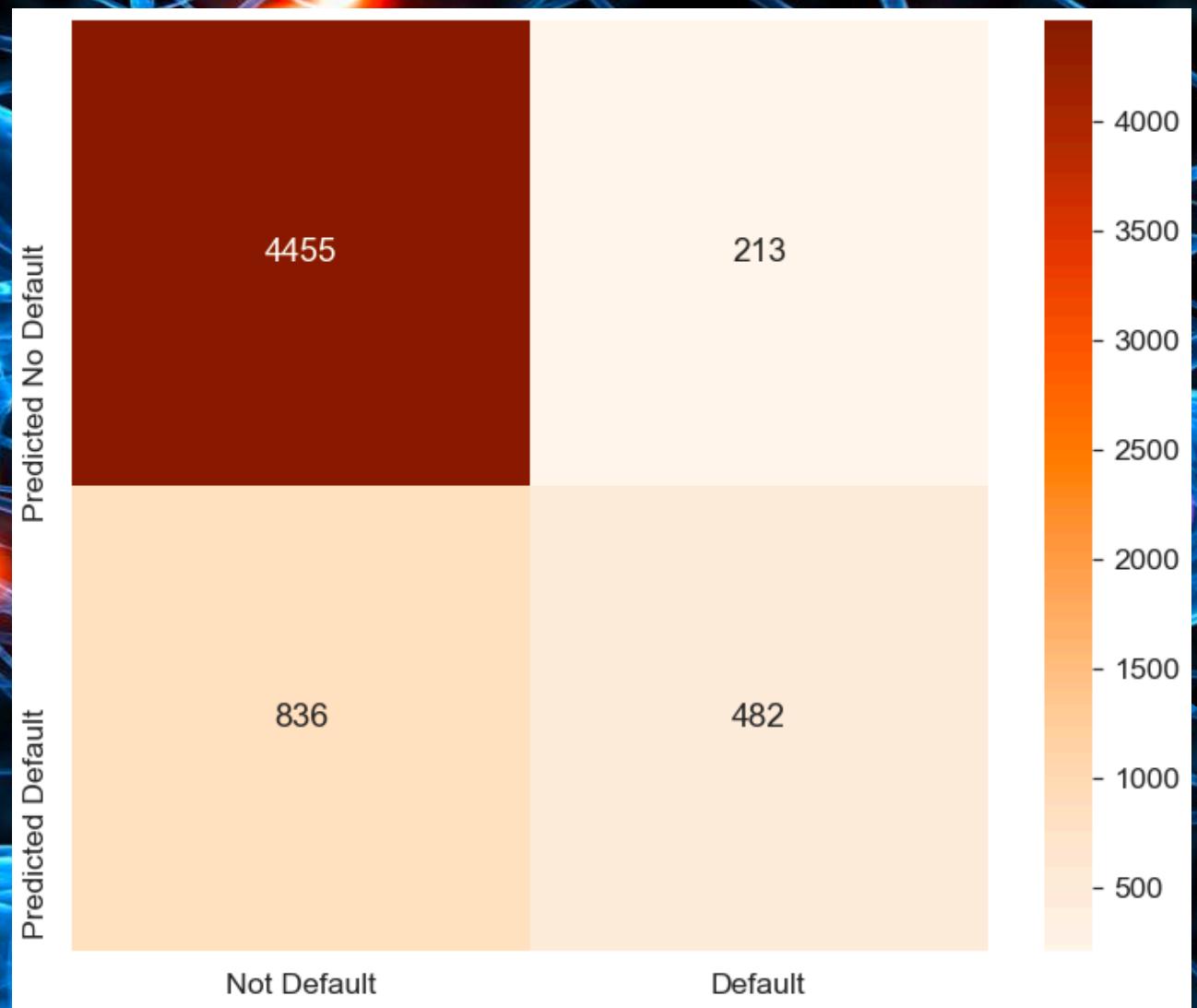


# Ada Boost



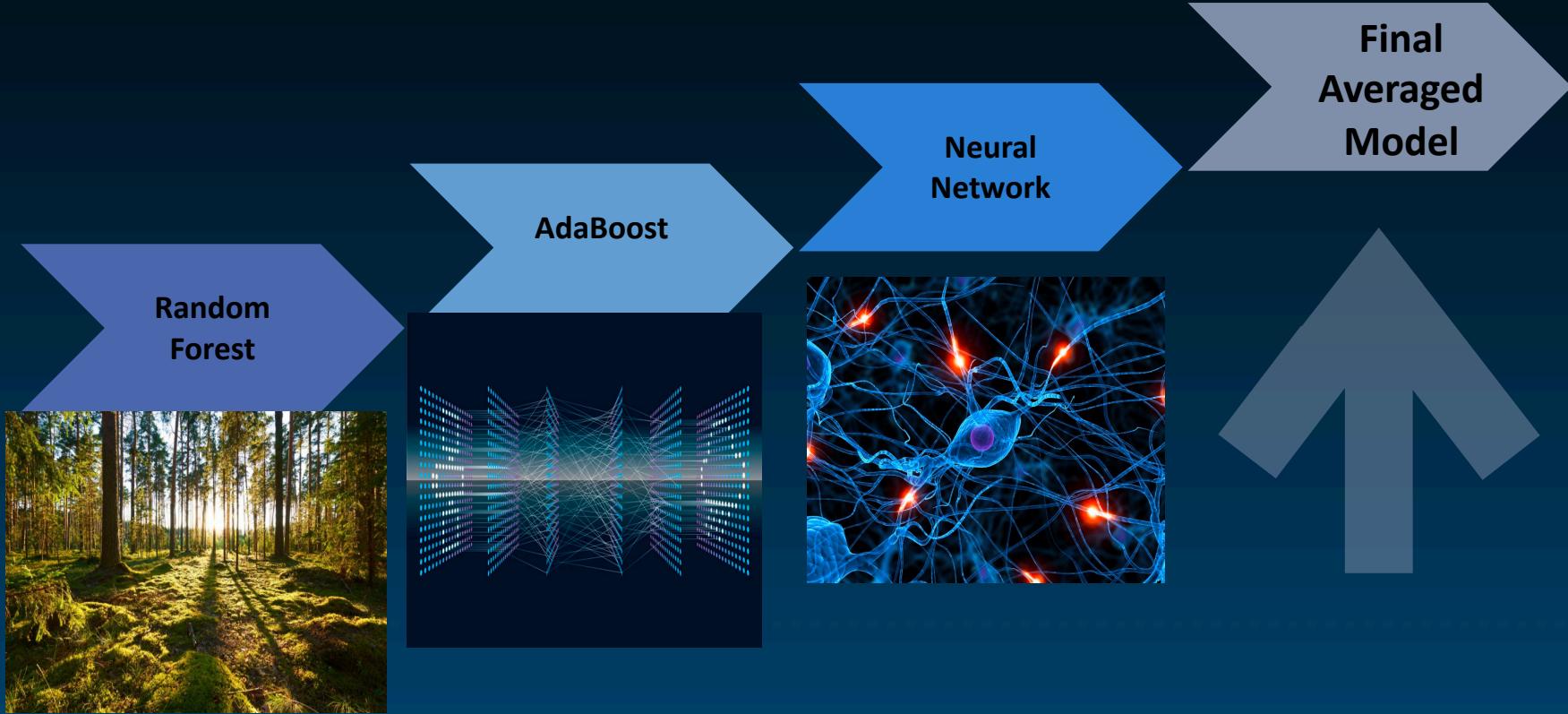
We chose max depth of 5 and Number of trees = 40

# Neural network

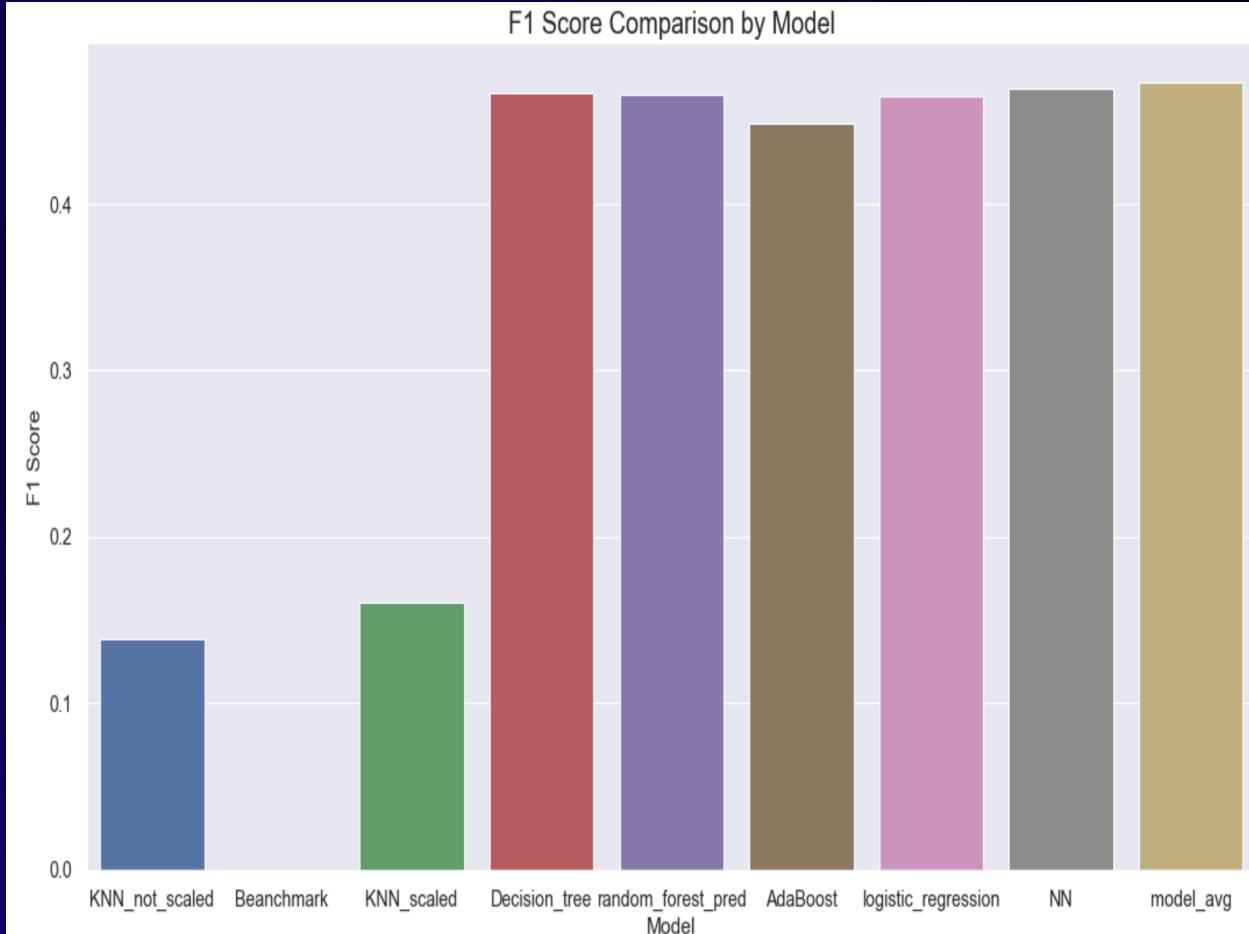
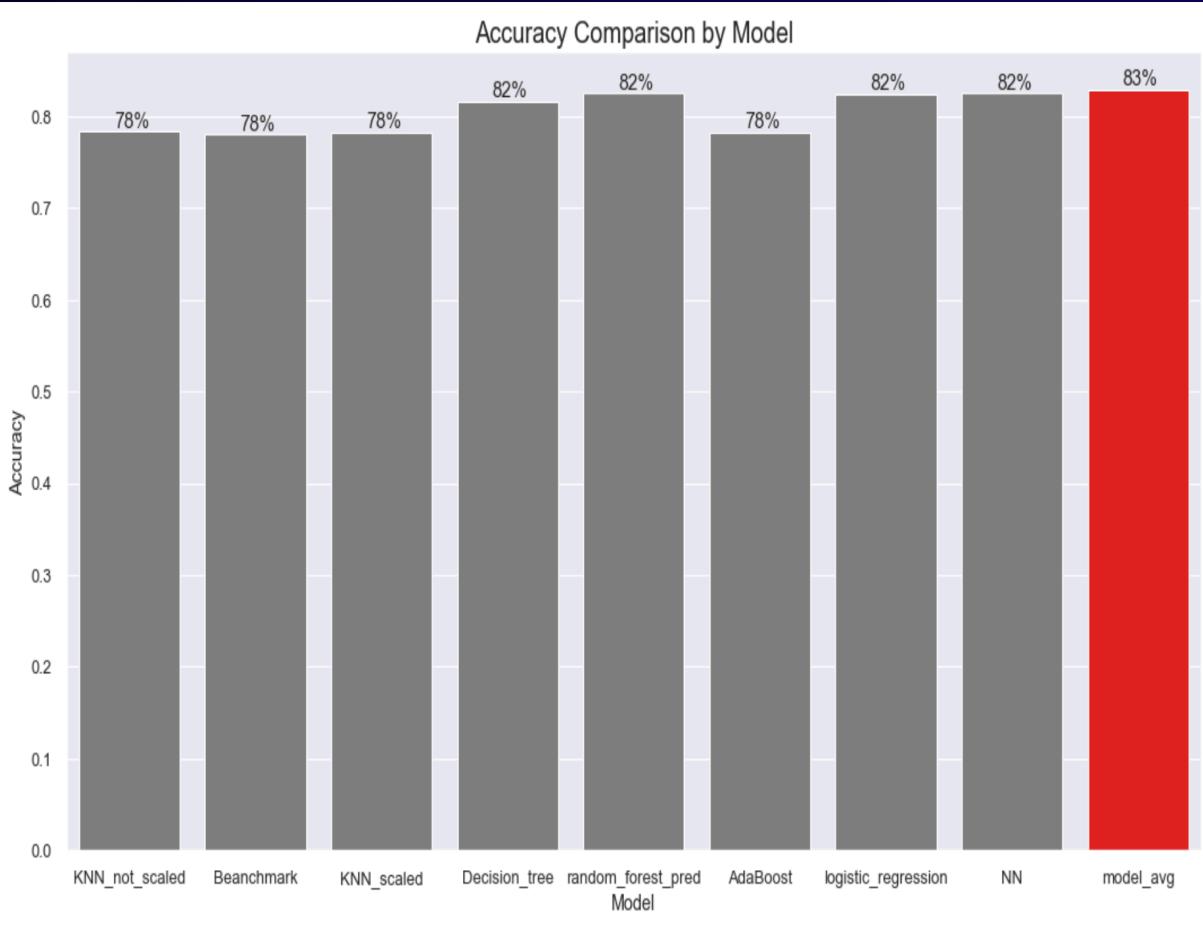


**Accuracy of 82%  
F1 Score of 45%**

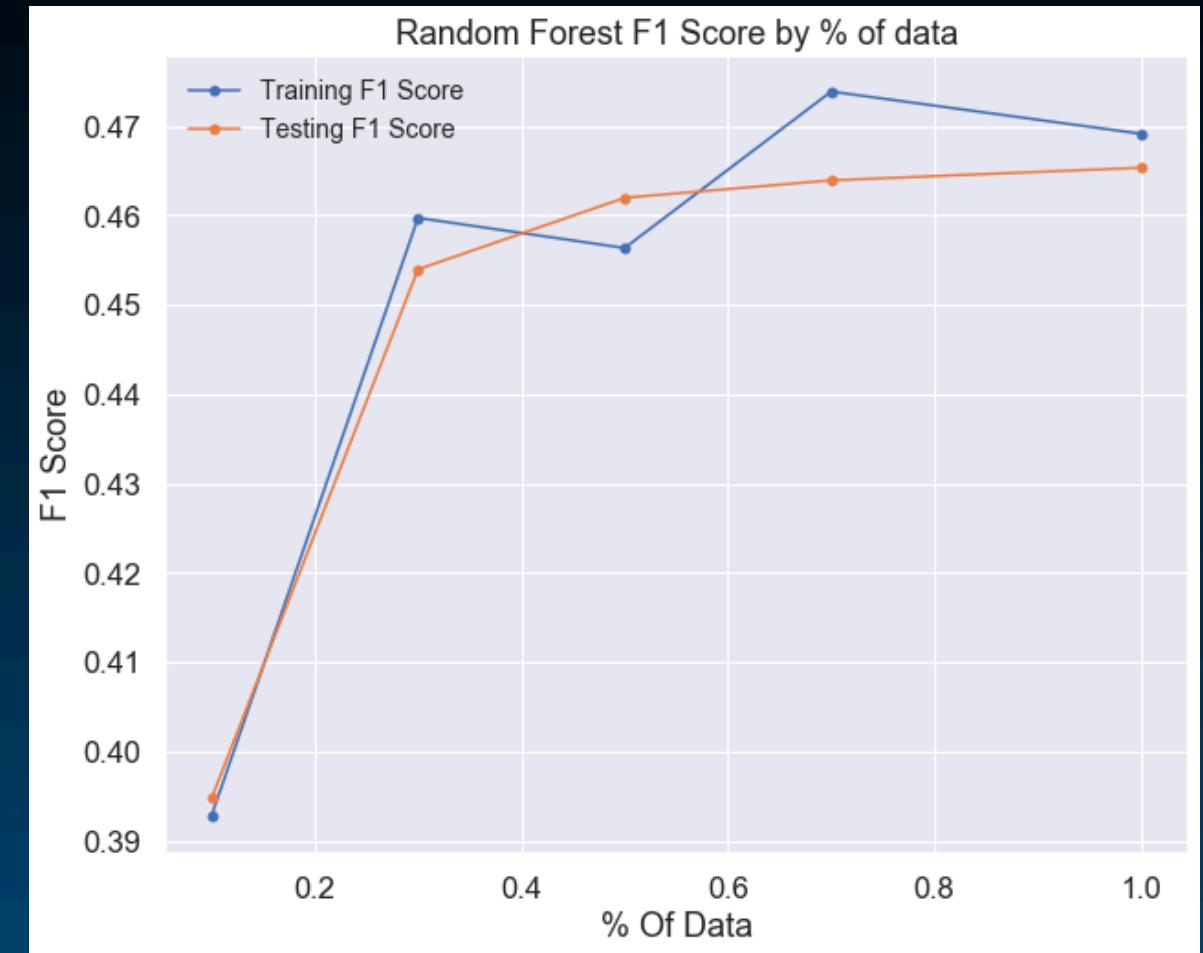
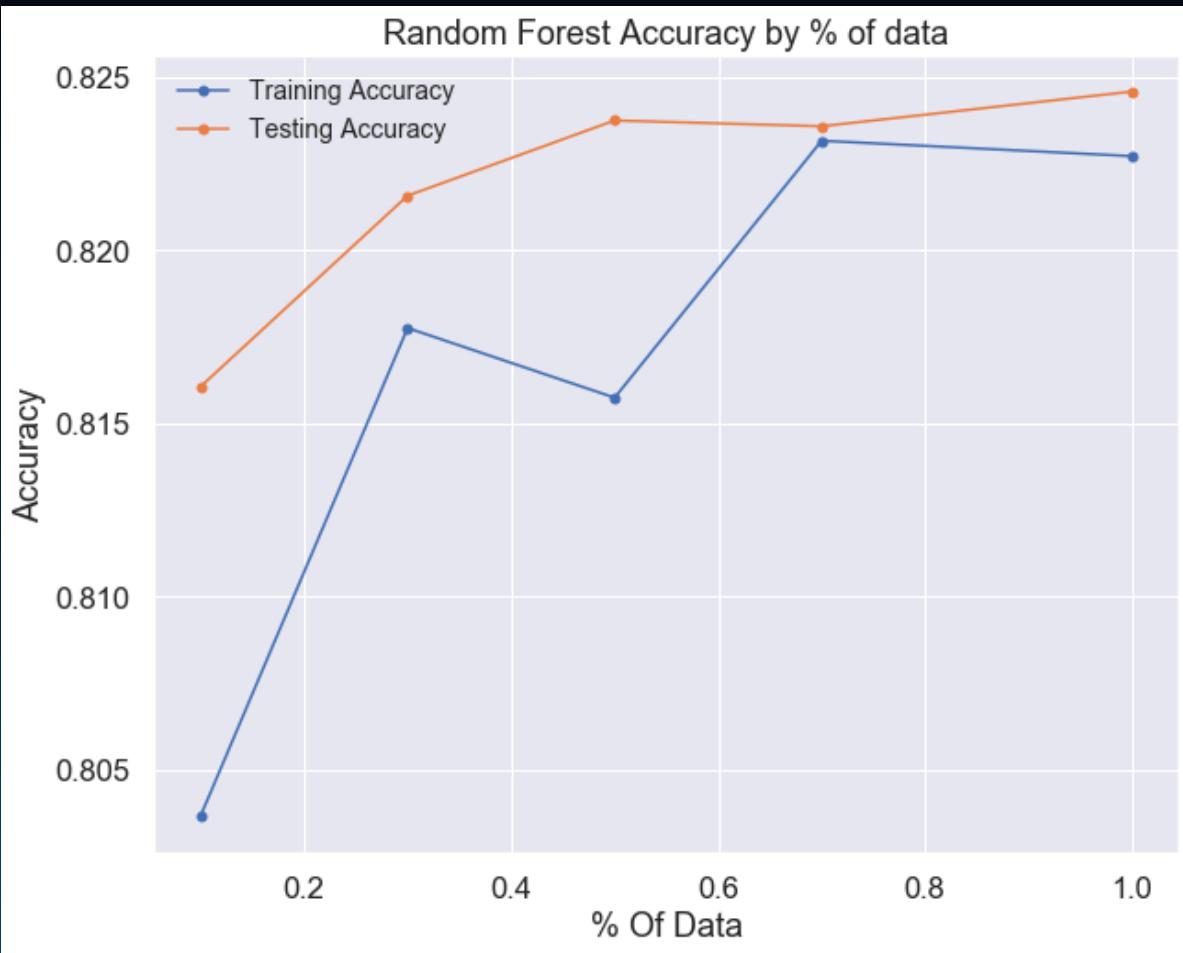
# Averaging models



# Averaging models



# Performance vs. amount of data



It seems that there is no need for more data



# THANK YOU

