

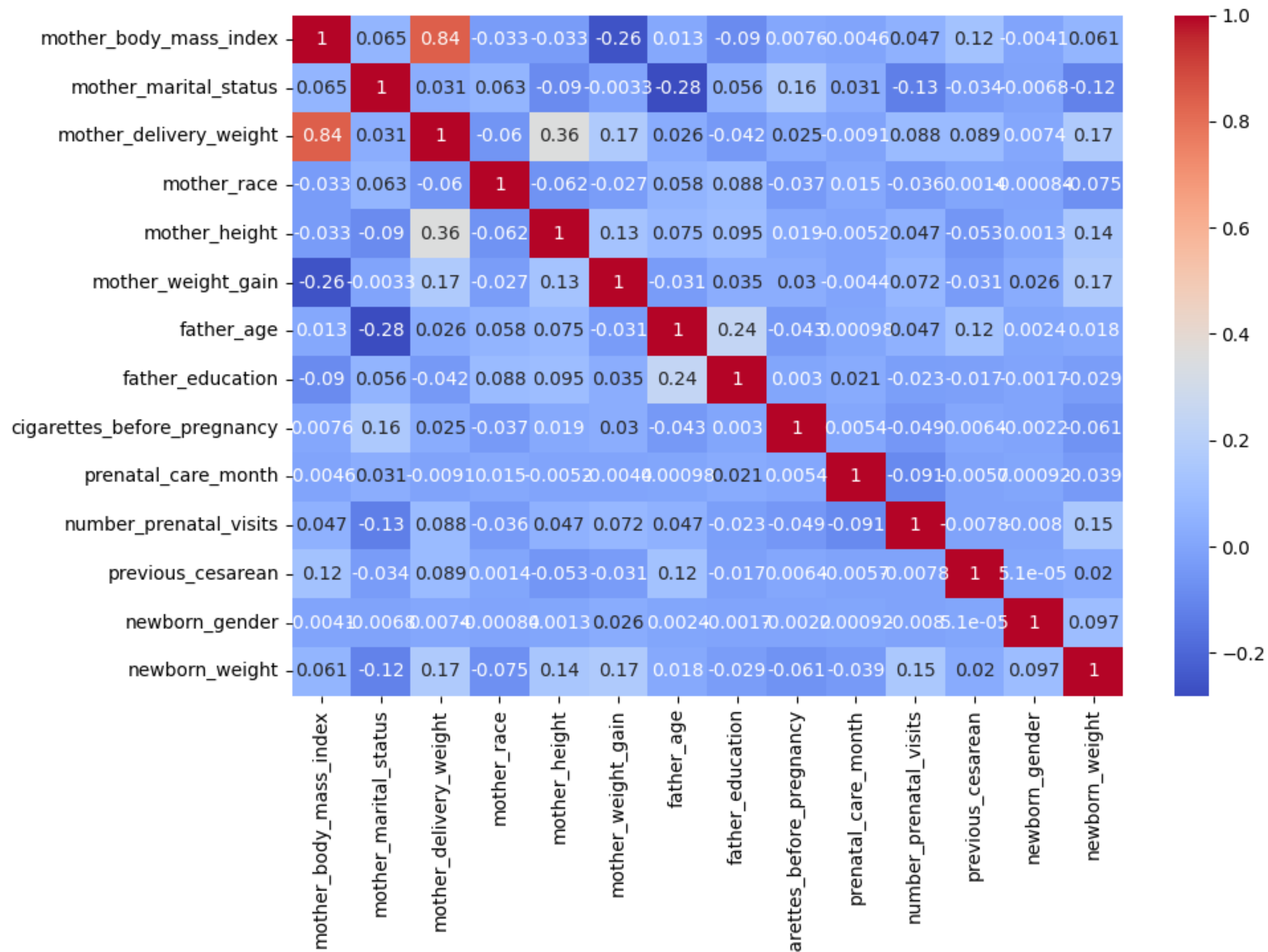
ML Project for Regression

Yiqing Hu 455858

Raffaele Ricci K-15305

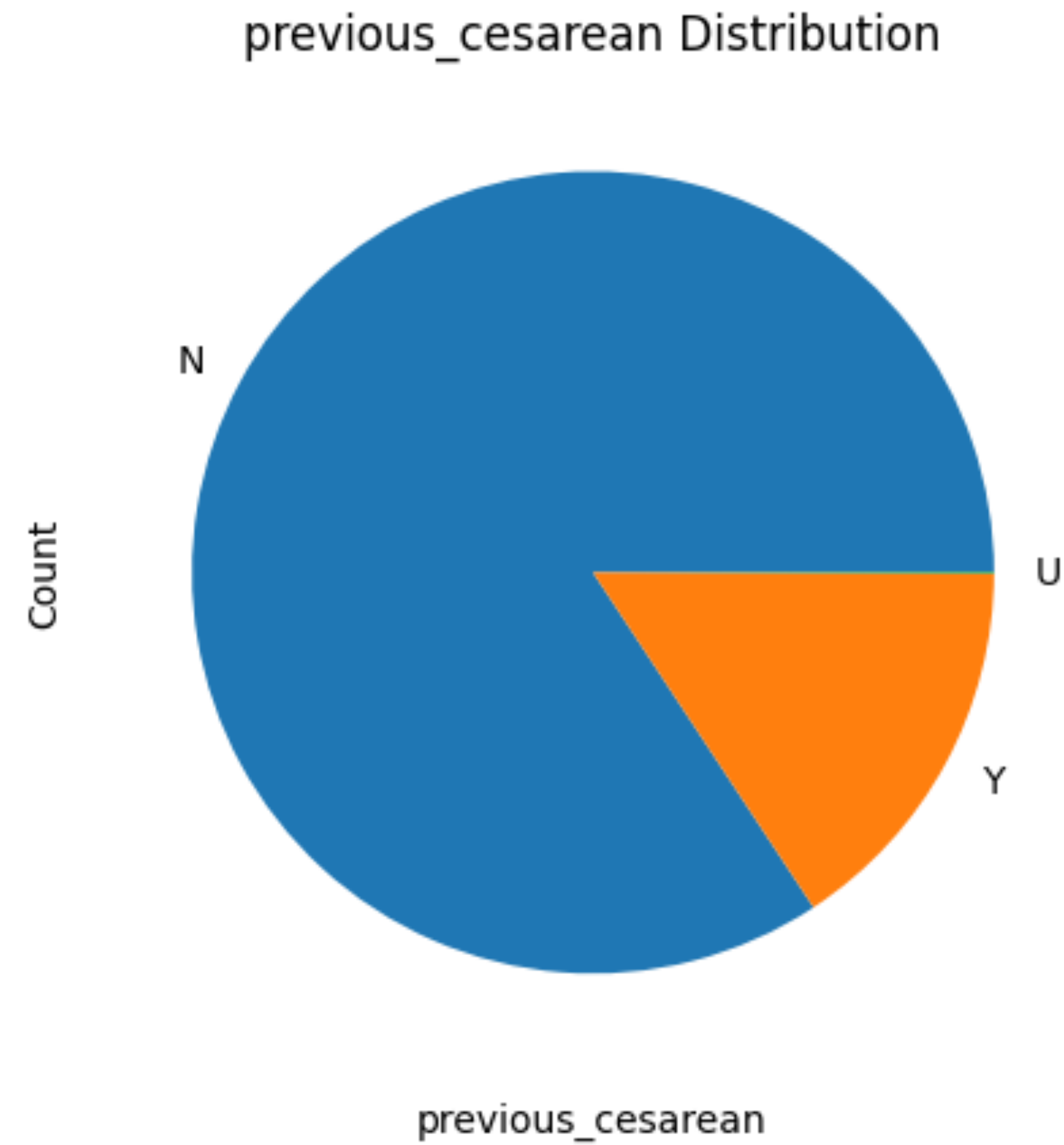
Data Checking

- Correlation
- Distribution
- NA



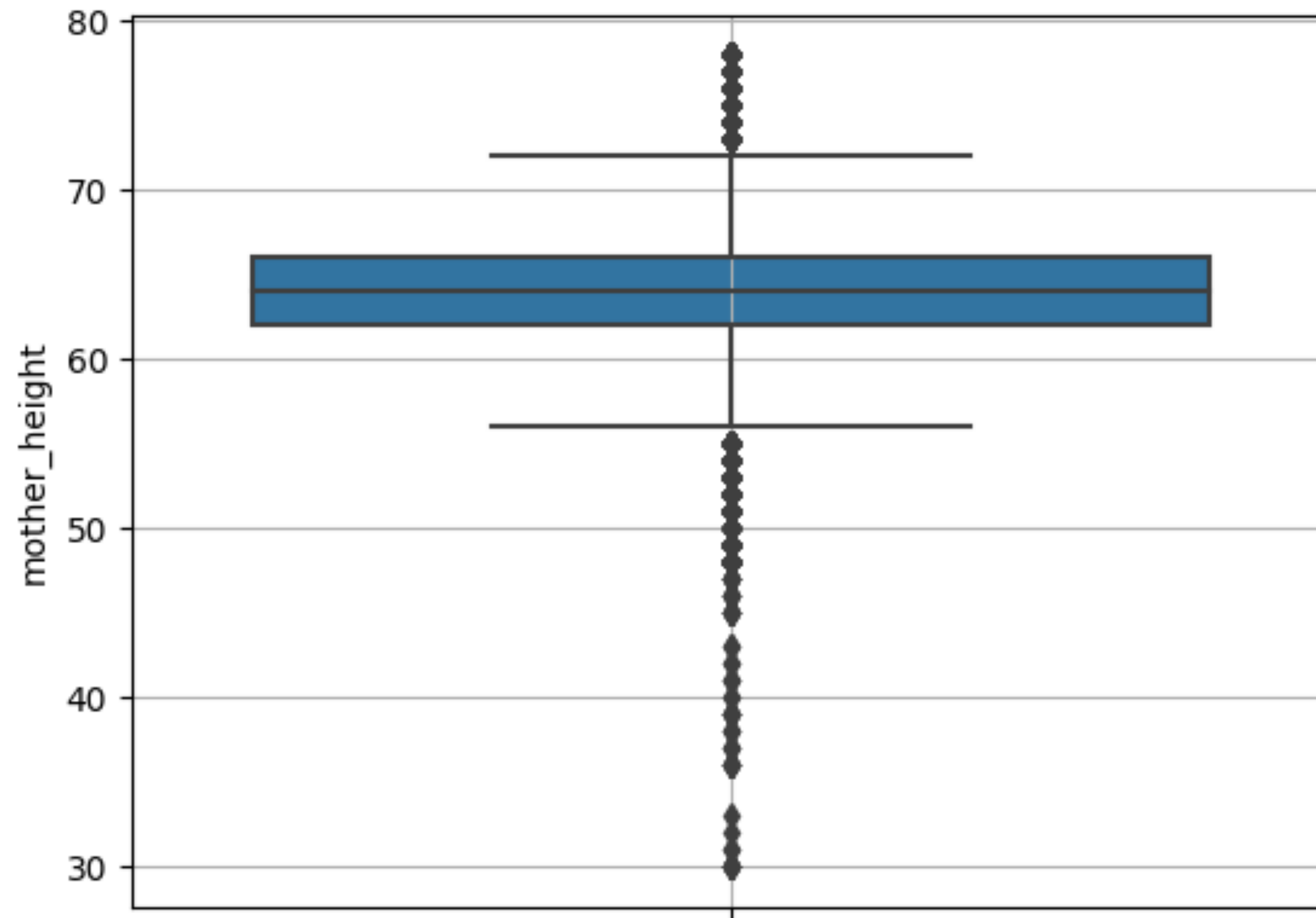
Visualization

Distribution of Discrete Data



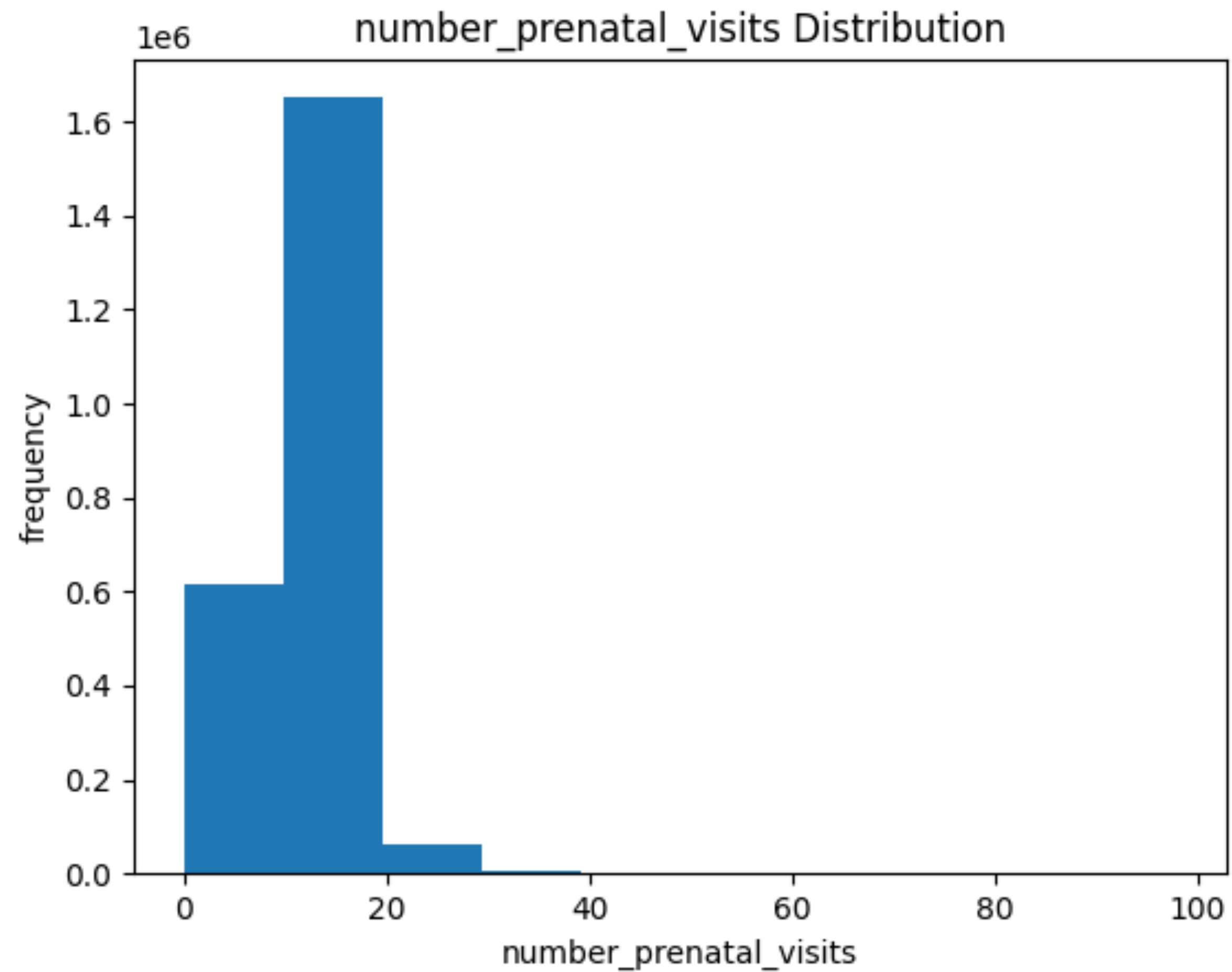
Visualization

Distribution of Float Data



Visualization

Distribution of Int Data



NA Checking

col_name	col_type	NA_ratio
father_age	float64	18.53563
mother_marital_status	float64	17.20142
mother_height	float64	10.19671
mother_body_mass_index	float64	6.11313
mother_weight_gain	float64	3.06378
number_prenatal_visits	float64	2.49784
mother_delivery_weight	float64	1.45773
cigarettes_before_pregnancy	float64	0.47124
mother_race	int64	0.0
father_education	int64	0.0
prenatal_care_month	int64	0.0
previous_cesarean	object	0.0
newborn_gender	object	0.0
newborn_weight	int64	0.0

Feature Engineering

int/float/discrete

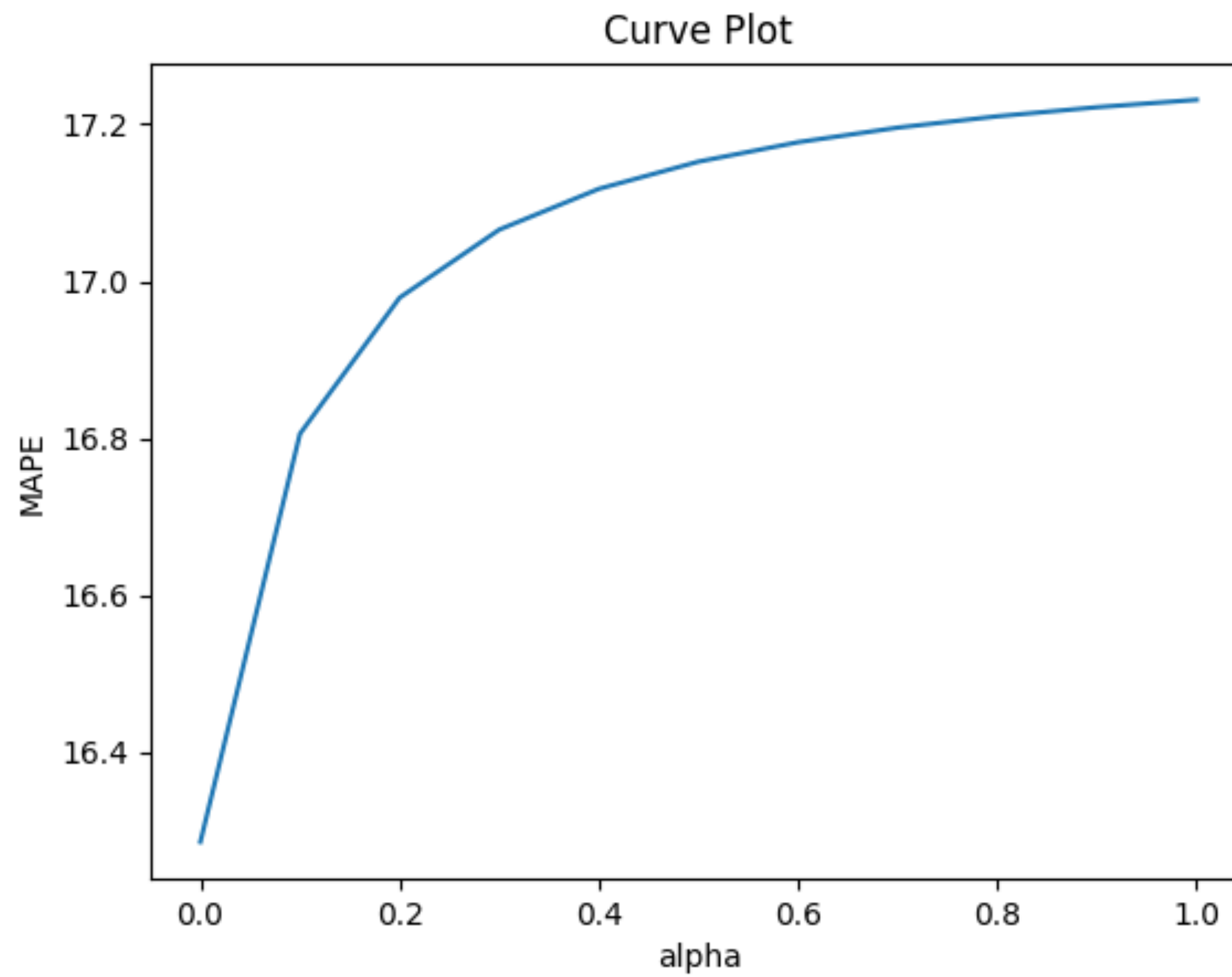
- Fill NA/outlier with median/mean to int and float features
- Make a new feature for old with which row has NA
- Truncate outliers
- Standardization and normalization for int and float features
- Make discrete features as dummies
- Power features which have high correlation with newborn weight

Algorithms

GLM/SVM/GBT

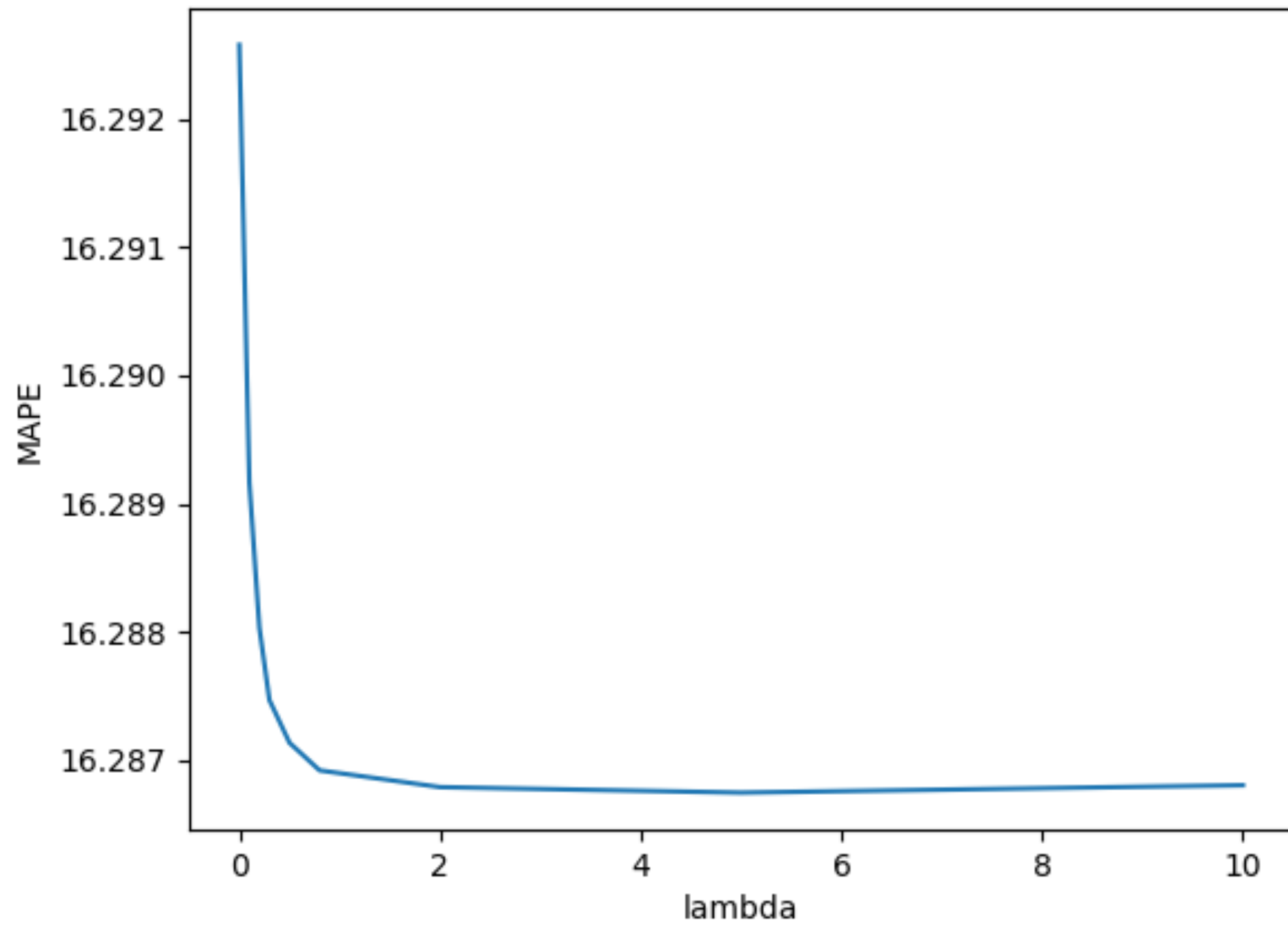
- LM, Ridge, Lasso, ElasticNet, FM
- SVR
- LGBMRegressor

Tuning For Lasso

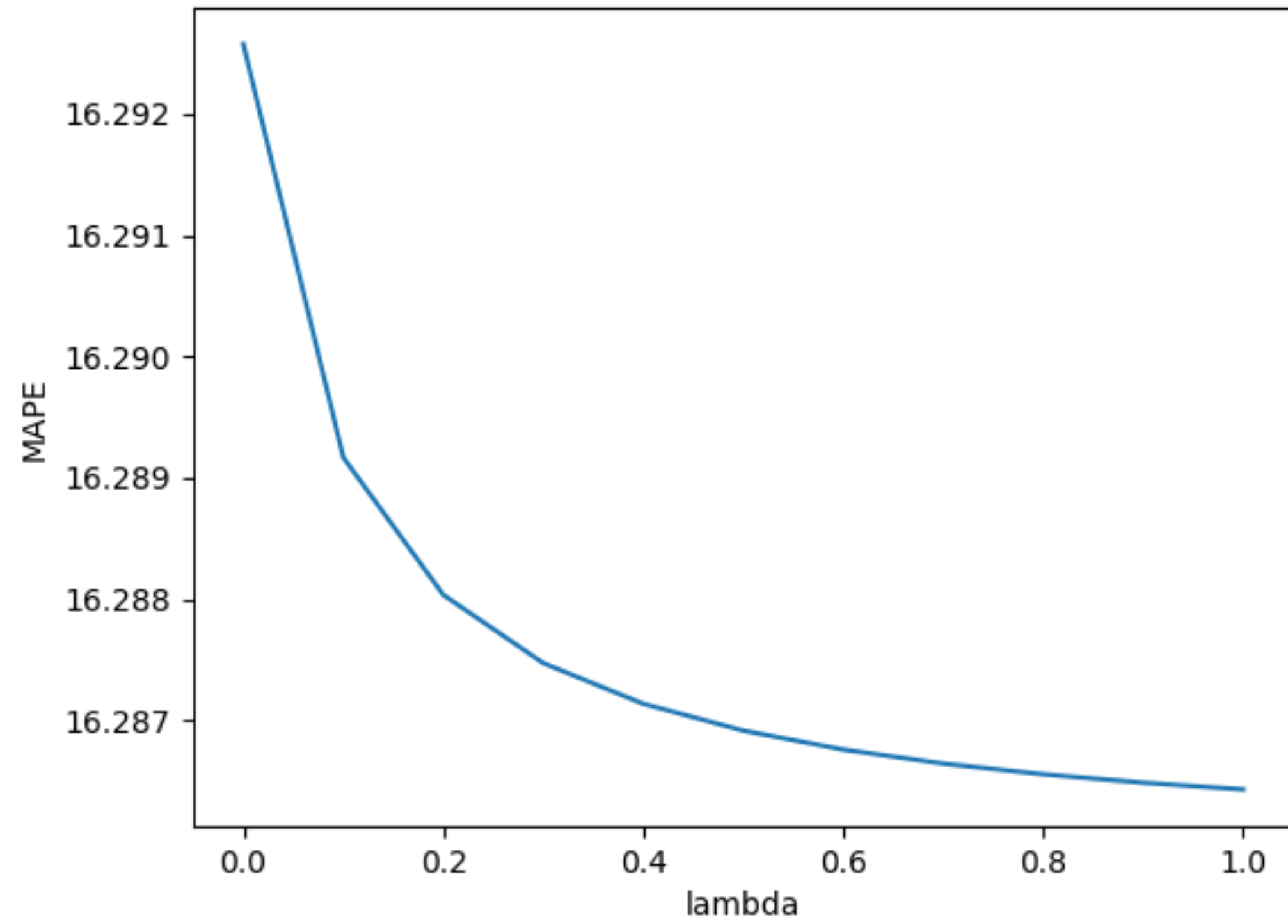


Tuning For Ridge

Curve Plot



Curve Plot

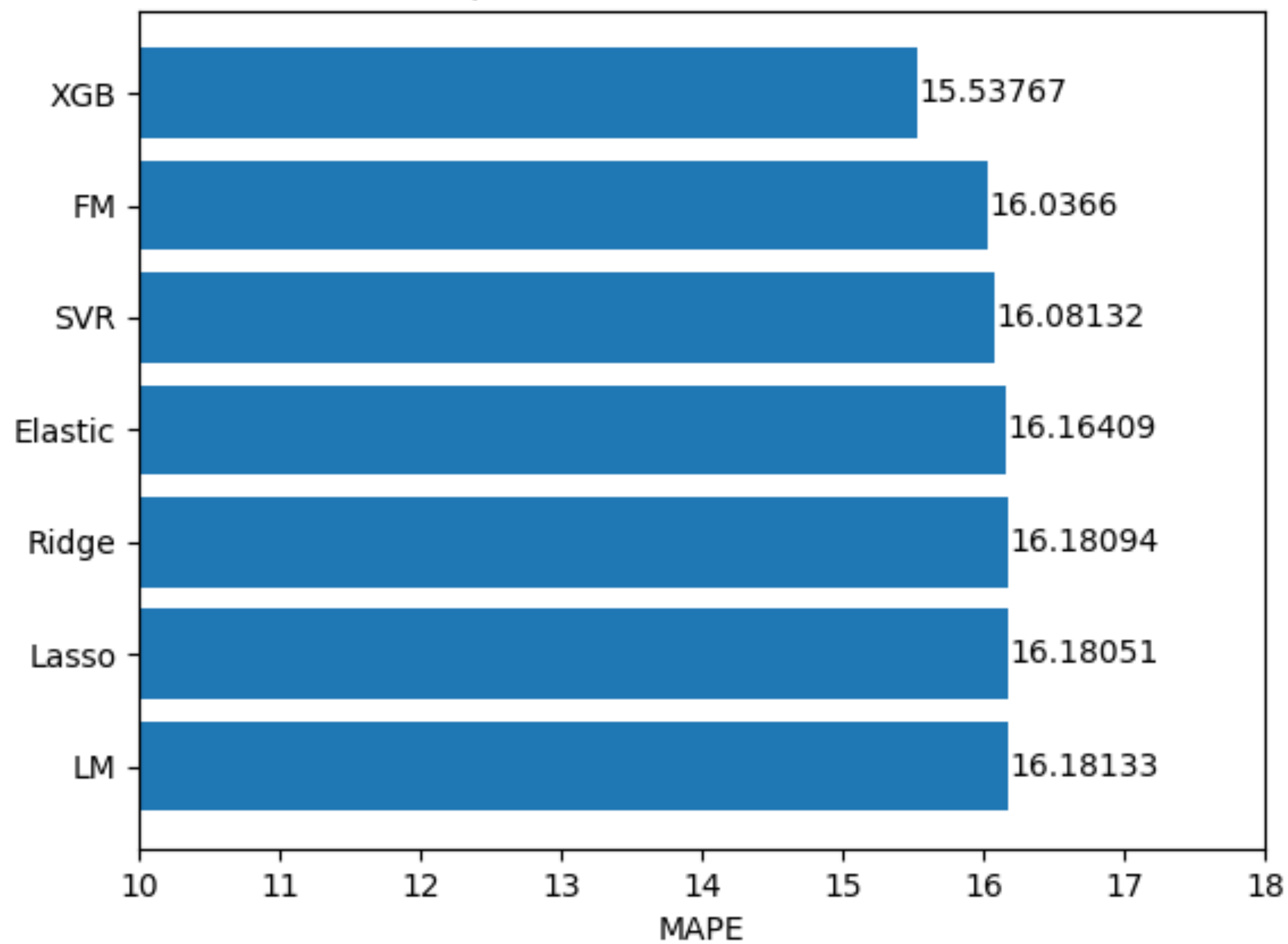


Tuning

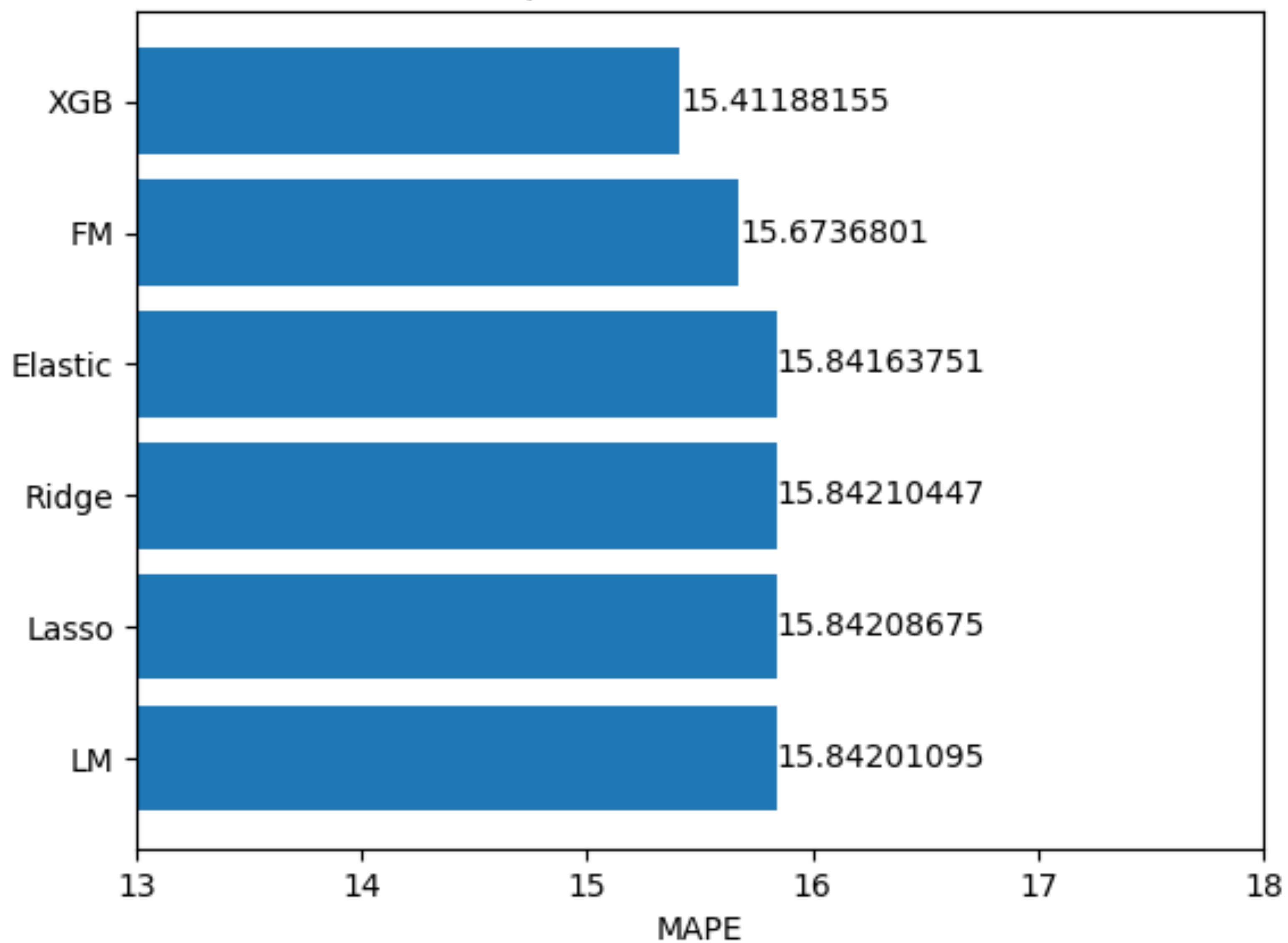
For ElasticNet

```
grid_search = GridSearchCV(ElasticNet(),  
                           param_grid,  
                           cv=5,  
                           scoring='mean_absolute_percentage_error')  
grid_search.fit(X_train, y_train)
```

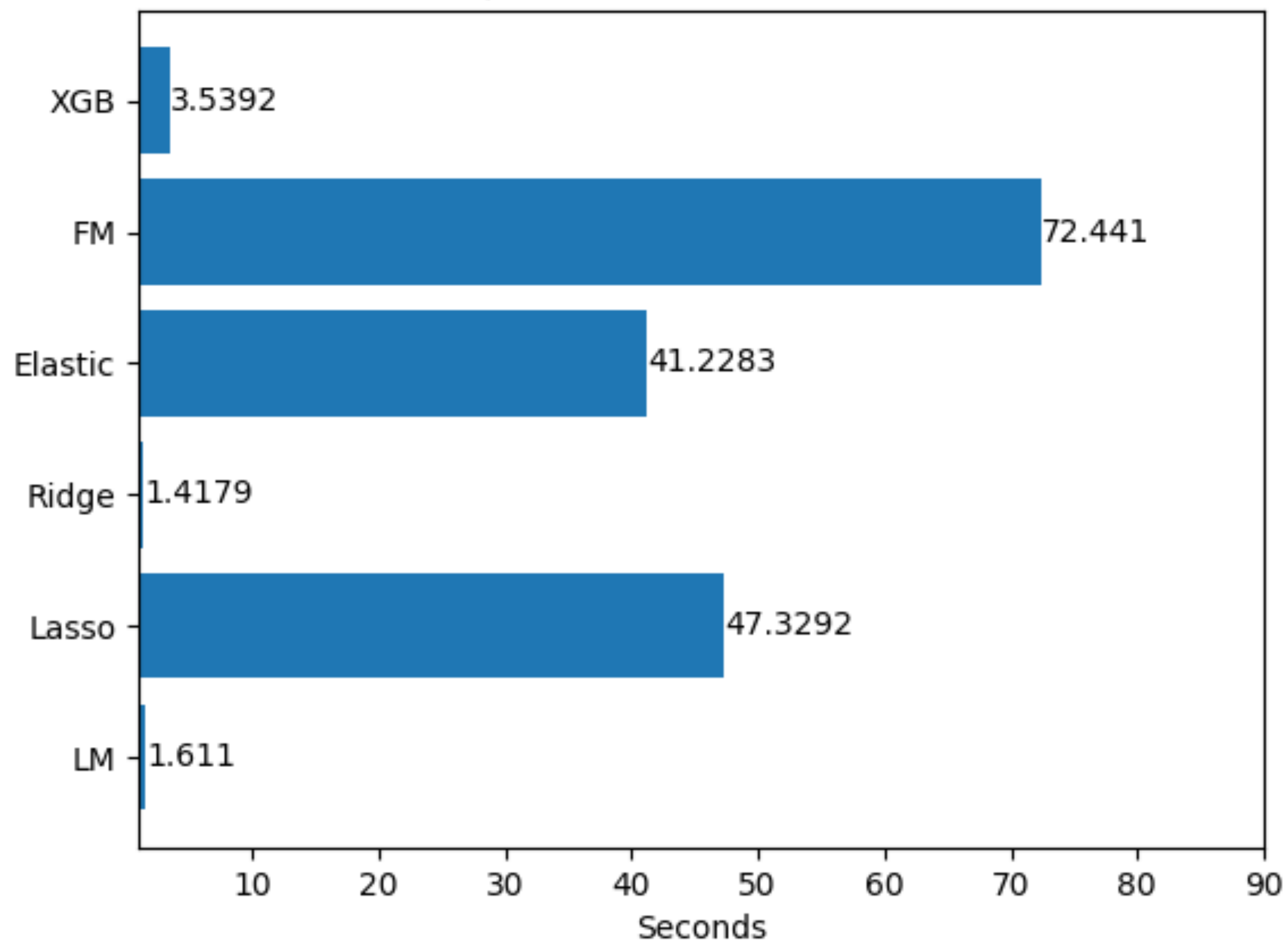
Comparison of MAPE in Validation



Comparison of MAPE in Test



Comparison of Time Cost in Test



Result

LGT

- Train time cost: within 3s for full train data
- MAPE: 15.411% for 10% testset

ML Project for Classification

Yiqing Hu 455858

Raffaele Ricci K-15305

Dimension of dataset = 10127, 21

Na Values

customer sex = 1018
total transaction amount 407
customer age = 624
customer salary range = 681

Unknown Values

customer education = 1519
customer civil status = 749

Feature Selection

#customer_sex

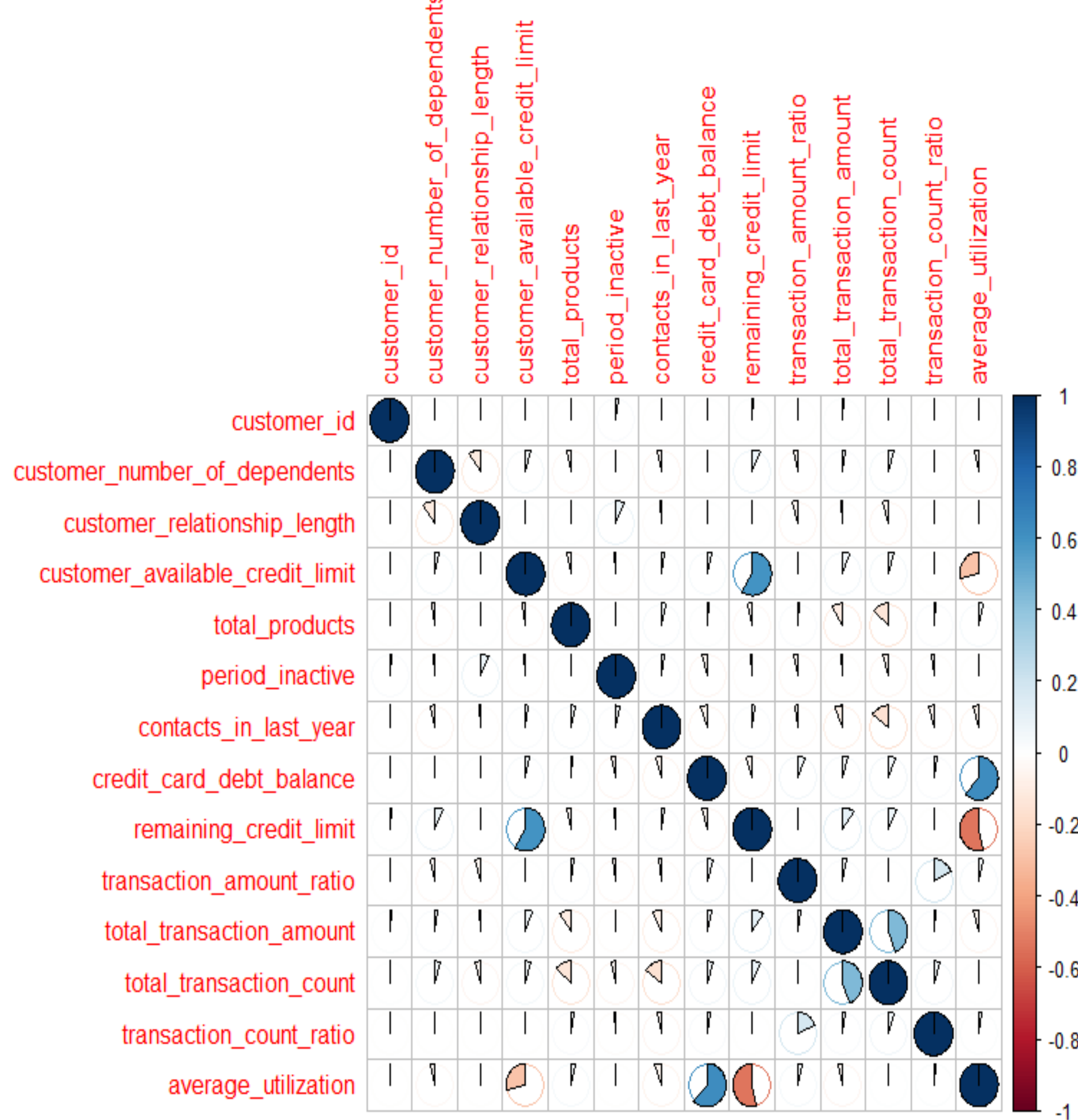
#customer_age

0.003424409 corr. account_status

mean in group closed	mean in group open
46.62196	46.25996
Test t: p-value = 0.09492	

- we replace the NA value in total transaction amount with the mean
- Replace NA values with "Unknown"
- actually in the variable of salary range we believe that the people with lower or higher salary are unlikely to say own salary

Feature Selection



customer_number_of_dependents	-0.0189905963
customer_relationship_length	-0.0136868512
customer_available_credit_limit	0.0174582135
total_products	0.0795483222
period_inactive	-0.1524488063
contacts_in_last_year	-0.2044905100
credit_card_debt_balance	0.2630528831
remaining_credit_limit	0.0002850775
transaction_amount_ratio	0.1310628478
total_transaction_amount	0.0959313407
total_transaction_count	0.3714027012
transaction_count_ratio	0.1208377823
average_utilization	0.1784103316

Feature Selection

- Recode ordinal variables with levels, including Unknown as a level

Correlation with account_status

customer_education = 0.035

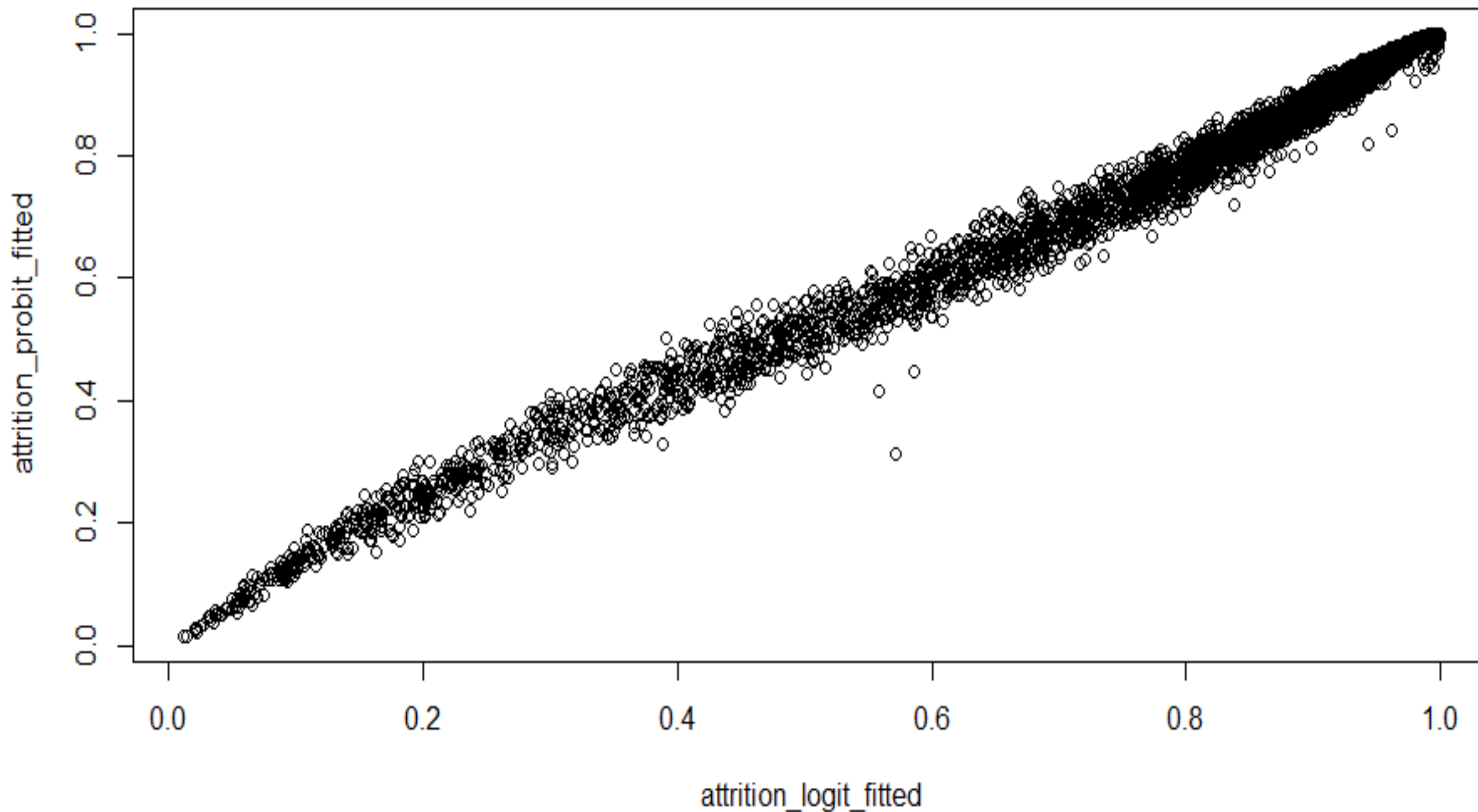
customer_civil_status = 0.024

customer_salary_range = 0.037

credit_card-classification = 0.015

- we don't remove anything from the earlier apart from the variable sex and age

Logistic Regression



- Division of the dataset in train and test sample
- No big difference in predictions using probit and logit
- Balanced accuracy = 0.735

Feature Selection

In logistic regression we have some insignificant variables

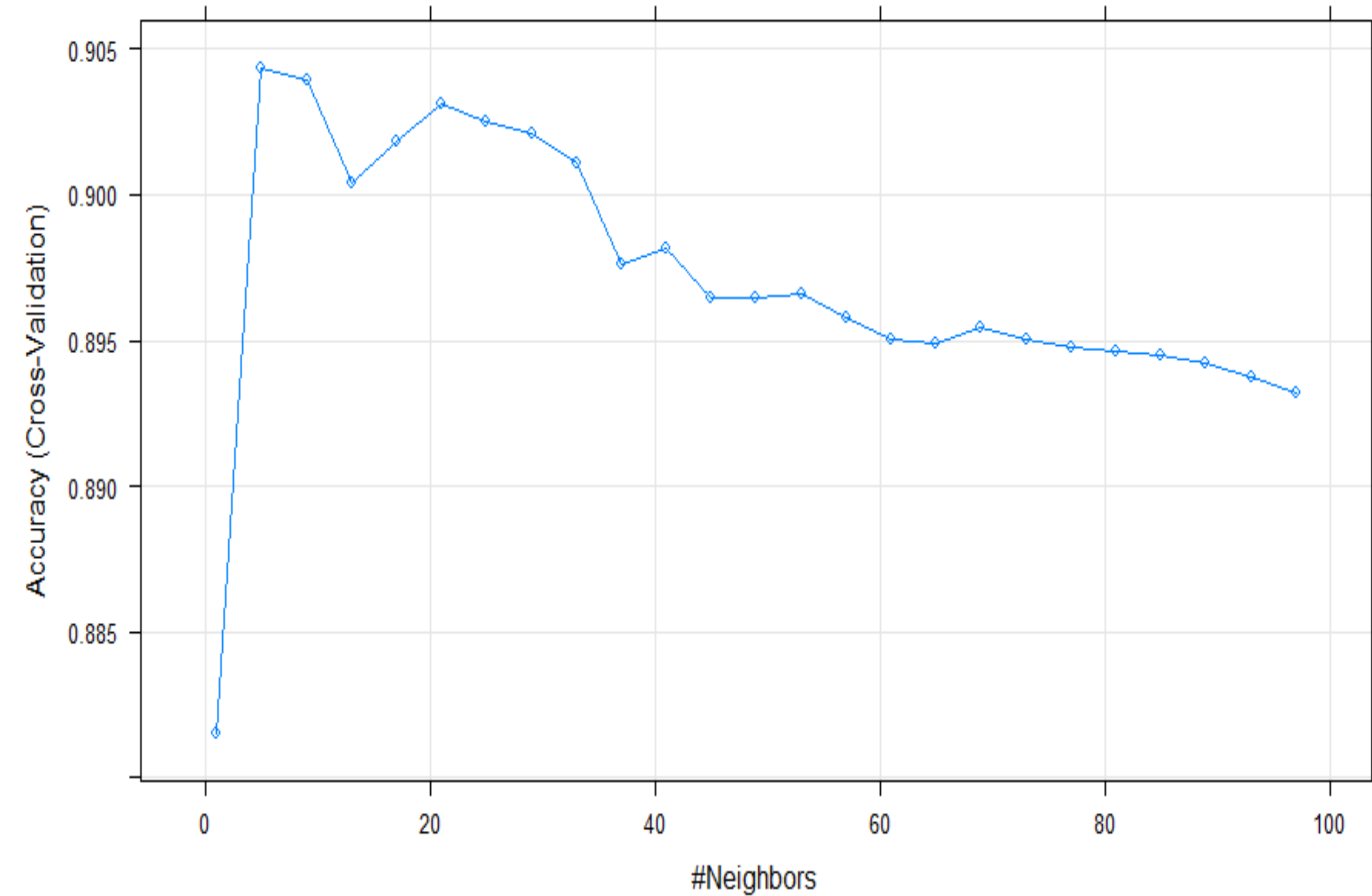
IN LOGIT :

remaining_credit_limit
average_utilization
customer_available_credit_limit
customer_civil_status

IN PROBIT :

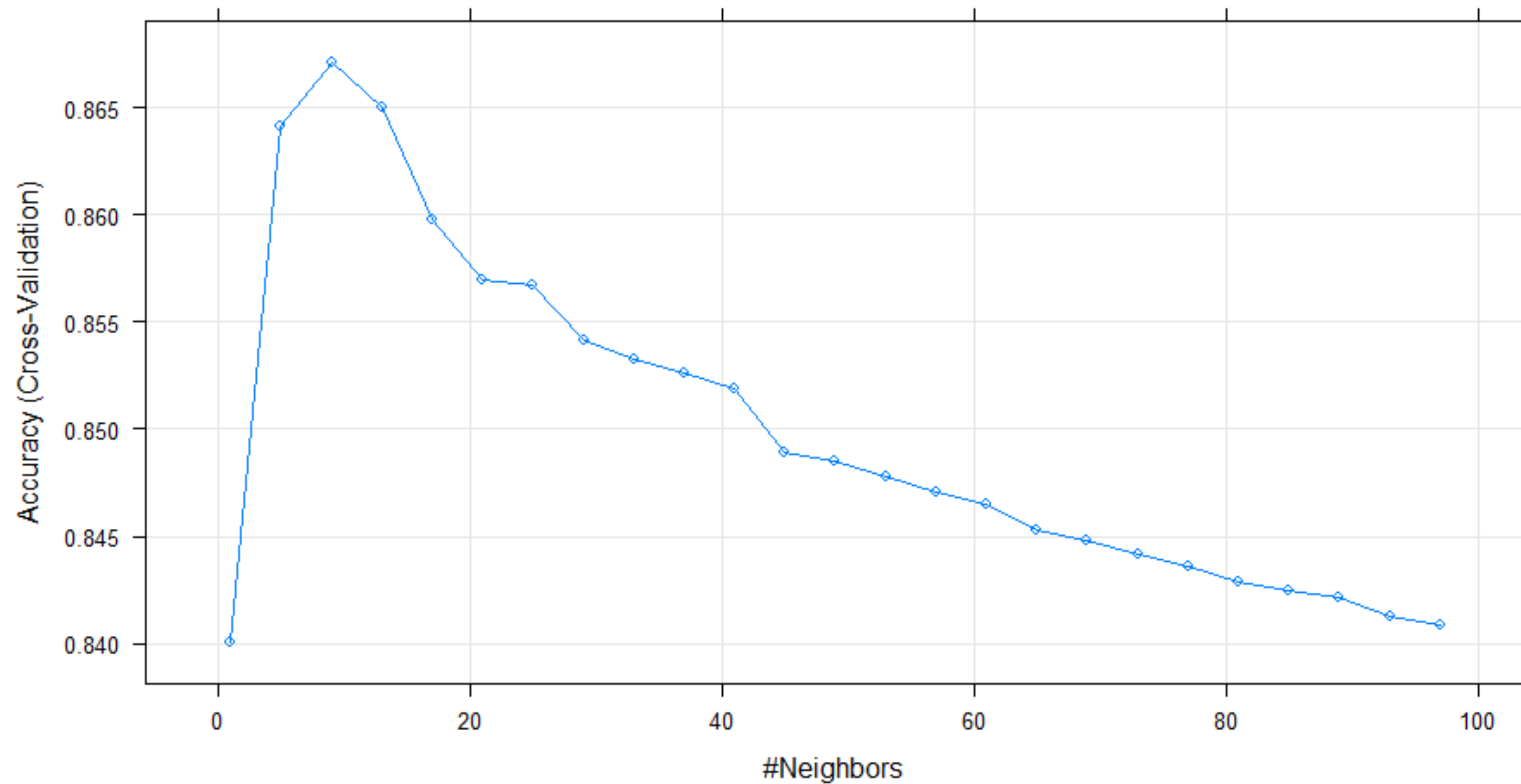
customer_civil_status 4
customer_available_credit_limit 7
remaining_credit_limit 13

KNN Model



Using cross-validation :
Accuracy was used to select the optimal model using the largest value.
The final value used for the model was $k = 5$.
Balanced Accuracy = 0.7995

KNN Model



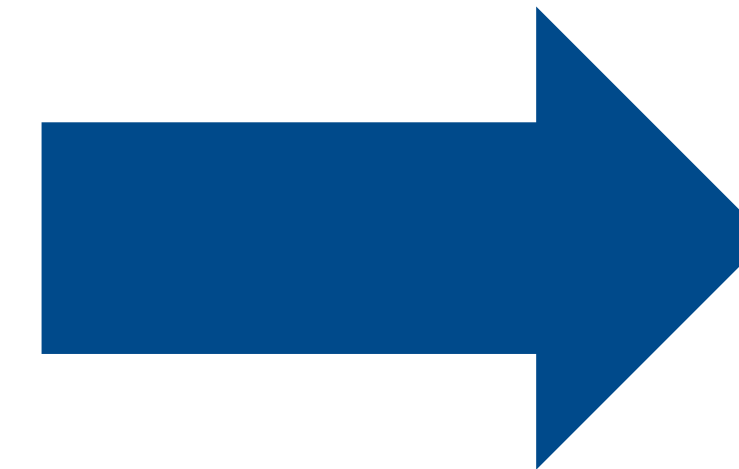
We compared all the models with scaling the variables but still the best balanced Accuracy remains 0.7995 of the model with $K = 5$

SVM Model

the value of 0.5 is indicated as optimal,

but the accuracy is identical for all C apart from 0.001

# C	Accuracy	Kappa
#0.001	0.8393286	0.00000000
#0.010	0.8918042	0.5004691
0.020	0.8931677	0.5161437
0.050	0.8936379	0.5267727
0.100	0.8938730	0.5299907
0.200	0.8940611	0.5328512
0.500	0.8941081	0.5334885
1.000	0.8940141	0.5330844
2.000	0.8939670	0.5330912
5.000	0.8940611	0.5334981



**Balanced
accuracy : 0.73**

Accuracy was used to select the optimal model using the largest value.

The final value used for the model was $C = 0.5$.

Thanks
guys

