**STAT 412 STATISTICAL DATA ANALYSIS**

**2023-2024 ACADEMIC YEAR SPRING SEMESTER**

**INTERIM REPORT**

**Eren Duralı 2361244**

**Introduction**

This study, Includes Credit score classification data taken from Kaggle
(https://www.kaggle.com/datasets/parisrohan/credit-score-classification/data)

Data includes 28 variables and 100000 observations.

**Aim**

Our target in this study, classify the credit score correctly by using some
statistical methods, and models with machine learning algorithms. Trying to find the best
classification model.

**Data Cleaning**

In the data cleaning part, we deleted some columns that are not our concern for analysis and model building. These columns are Id, Customer_id, Month, Name, SSN, and Occupation.



After that, we look at the summary statistics for detecting unusual moments in the variables.

Figure 1: Summary statistics

We detect some of them and deal with filtering the data. After this process, we were left with 74179 rows. After that, we look at the boxplots for detecting the outliers and distribution of the variables in our data.



Figure 2: Some examples of boxplots

For the categorical variables, we use bar plots to obtain their distributions.



Figure 3: Credit score bar plot

After these processes, we deal with missing values. Firstly, detect the missing values and we have 34113 missing values in our dataset. Monthly in hand salary column has 11150 missing

```
1] 34113
                   age           annual_income     monthly_inhand_salary       num_bank_accounts         num_credit_card
                     0                    5195                     11150                       0                       0
         interest_rate              num_of_loan         delay_from_due_date    num_of_delayed_payment     changed_credit_limit
                     0                       0                         0                       0                    1558
    num_credit_inquiries         outstanding_debt     credit_utilization_ratio      total_emi_per_month   amount_invested_monthly
                  1445                     754                         0                       0                    6473
        monthly_balance                   month                 occupation               type_of_loan                credit_mix
                   892                       0                         0                       0                       0
     credit_history_age      payment_of_min_amount        payment_behaviour             credit_score
                  6646                       0                         0                       0
```

values so we consider that column might be removed from data.

Figure 4: Total Missing Values

3

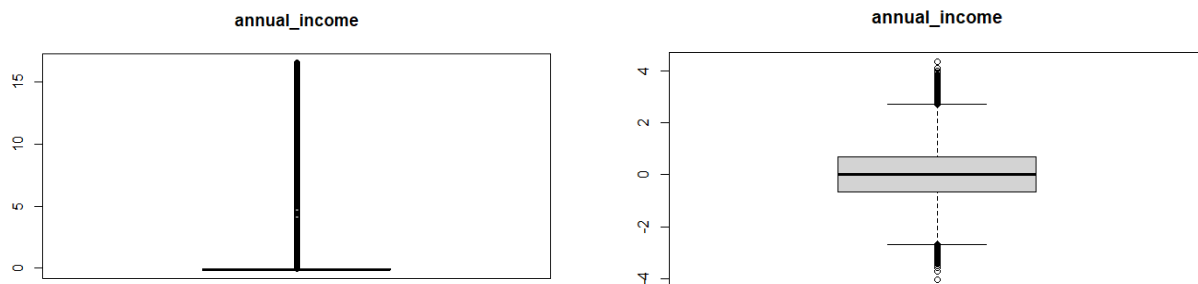| age | annual_income | monthly_inhand_salary | num_bank_accounts | num_credit_card |
|---|---|---|---|---|
| 0.00000000 | 0.07003330 | 0.15031208 | 0.00000000 | 0.00000000 |
| interest_rate | num_of_loan | delay_from_due_date | num_of_delayed_payment | changed_credit_limit |
| 0.00000000 | 0.00000000 | 0.00000000 | 0.00000000 | 0.02100325 |
| num_credit_inquiries | outstanding_debt | credit_utilization_ratio | total_emi_per_month | amount_invested_monthly |
| 0.01947991 | 0.01016460 | 0.00000000 | 0.00000000 | 0.08726189 |
| monthly_balance | month | occupation | type_of_loan | credit_mix |
| 0.01202497 | 0.00000000 | 0.00000000 | 0.00000000 | 0.00000000 |
| credit_history_age | payment_of_min_amount | payment_behaviour | credit_score | |
| 0.08959409 | 0.00000000 | 0.00000000 | 0.00000000 | |

Figure 5: Missing Values Percentage

Since, 60% our thresh hold for deleting the column we keep the monthly in hand salary. We use na.approx function from zoo package in R to deal with missing values in numerical columns. This package fills the na values by interpolated values. After dealing with numerical columns, we continue to deal with missing values in categorical columns. We have 6647 missing values in categorical columns and all of them in the credit history age column (Represents the age of credit history of the person). Since this column consist of years and months. We change this column to months and apply na.approx function. When missing



values are gone, we try to normalize our data because most of the statistical method assumptions are normality. We use bestNormalize package for that.

Figure 6: Shows the boxplot before & after normalization

**Research Questions**

**Question1) How does annual income influence the number of bank accounts and credit cards an individual holds?**

We want to investigate is annual income has an influence on the number of bank accounts and credit cards and if it has an influence how is it?
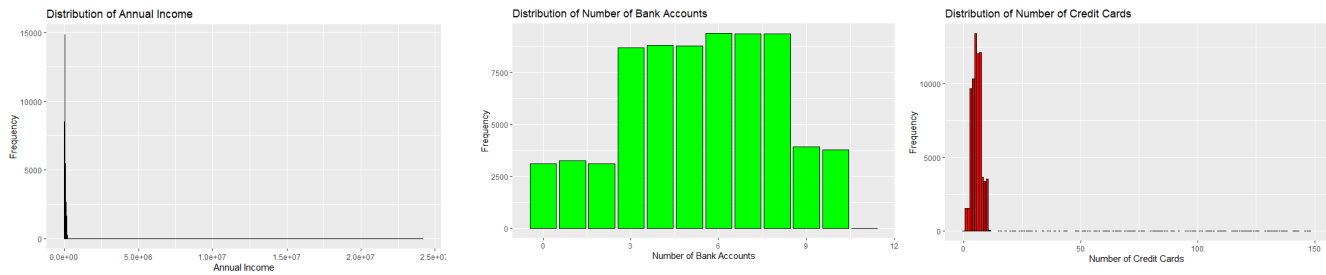
Figure 7: Distribution of Variables

Annual income and number of credit cards has a left skewed distribution which means that people with lower annual income holds most of the population.
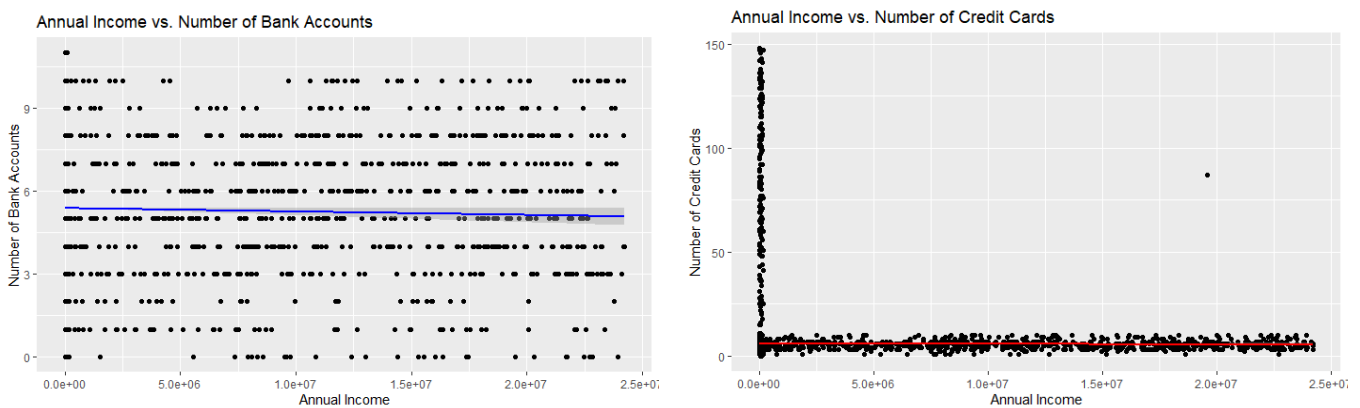


Figure 8: Scatter plot between independent and dependent variables

This scatter plot shows that we don't have a strong linear relationship between Annual income between number of bank accounts and credit cards. Since it is reasonable because of



the number of bank accounts and credit cards values are just countable numbers.

Figure 9: Correlation values and Linear regression models

Annual income has a 25% negative correlation between number of bank accounts and 19% negative correlation between number of credit cards. In the linear regression models, Annual income was significant for both models but R-squared values were low. This shows that Annual income explains the variance of the number of bank accounts at 6.2% and the variance of the number of credit cards 3.9%. Residuals were looks good. And it has a negative coefficient as expected.

**Q2) What is the impact of interest rate and outstanding debt on the credit score of individuals?**

We want to investigate the impact of interest rate and outstanding debt on the our target variable credit score.

Figure 10: Distributions of Variables

Since our credit score is a categorical variable, we apply a multinominal logistic regression using the "multinom" function in the "nnet" package. Before that we encode the credit score

```
# weights:  12 (6 variable)
initial  value 81492.862349
iter  10 value 64567.931419
final  value 64567.727685
converged
Call:
multinom(formula = credit_score ~ interest_rate + outstanding_debt,
    data = credit5_norm)

Coefficients:
  (Intercept) interest_rate outstanding_debt
2   0.8372937    -0.3621638       -0.6468804
3  -0.6739207    -1.1864338       -0.8583294

Std. Errors:
  (Intercept) interest_rate outstanding_debt
2  0.00985607    0.01139506       0.01152196
3  0.01502935    0.01609618       0.01595673

Residual Deviance: 129135.5
AIC: 129147.5
  (Intercept) interest_rate outstanding_debt
2   0.8372937    -0.3621638       -0.6468804
3  -0.6739207    -1.1864338       -0.8583294
  (Intercept) interest_rate outstanding_debt
2   2.3101066     0.6961683        0.5236769
3   0.5097063     0.3053081        0.4238696
```

by label encoding.

Figure 11: Logistic regression model

In this logistic regression model (Intercept) represents "Poor", 2 represents "Standard" and 3 represents "Good" credit score. Interest rate and outstanding debt is significant since their p-value is $< 0.05$ (it showed in the code). For every one unit increase in interest rate; The odds of having a "Standard" credit score versus a "Poor" credit score decrease by approximately 30.38%.The odds of having a "Good" credit score versus a "Poor" credit score decrease by approximately 69.47%.For every one-unit increase in outstanding debt; The odds of having a

"Standard" credit score versus a "Poor" credit score decrease by approximately 47.63%.The odds of having a "Good" credit score versus a "Poor" credit score decrease by approximately 57.61%.

**Question 3) Is there a significant association between the type of loan (e.g., mortgage, car loan, personal loan) and the frequency of late payments?**

We try to find an association between the type of loan and the frequency of late payments. We conduct a Pearson's Chi-squared test.



```
          Pearson's Chi-squared test

data:  contingency_table
X-squared = 586953, df = 231583, p-value < 2.2e-16

X-squared
 17.11059
```

Figure 12: Chi-squared test

Since Cramer's V = 17.11 type of loan and the frequency of delayed payments has a strong association. This shows that type of loan has a significant impact on the frequency of delayed payments.

**Question 4) Does the age of individuals vary significantly across different levels of payment behavior?**

We conduct ANOVA to analyze this.

H0: The mean age is the same across all levels of payment behavior.

H1: At least one level of payment behavior has a different mean age compared to others.



```
                    Df Sum Sq Mean Sq F value Pr(>F)
payment_behaviour    6    118  19.624   19.92  <2e-16 ***
Residuals        74171  73073   0.985
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
  Tukey multiple comparisons of means
    95% family-wise confidence level

Fit: aov(formula = age ~ payment_behaviour, data = credit5_norm)

$payment_behaviour
                                                                         diff          lwr          upr      p adj
High_spent_Large_value_payments-!@9#%8                             0.0683536537  0.01959456  1.171127e-01 0.0007087
High_spent_Medium_value_payments-!@9#%8                            0.0488922186  0.00214186  9.564258e-02 0.0334964
High_spent_Small_value_payments-!@9#%8                             0.0102354990 -0.04023577  6.070677e-02 0.9968931
Low_spent_Large_value_payments-!@9#%8                              0.0220475592 -0.02930193  7.339705e-02 0.8675553
Low_spent_Medium_value_payments-!@9#%8                             0.0104056948 -0.03820409  5.901548e-02 0.9958110
Low_spent_Small_value_payments-!@9#%8                             -0.0472100226 -0.09176053 -2.659520e-03 0.0295330
High_spent_Medium_value_payments-High_spent_Large_value_payments  -0.0194614351 -0.05819285  1.926998e-02 0.7560740
High_spent_Small_value_payments-High_spent_Large_value_payments   -0.0581181547 -0.10126760 -1.496870e-02 0.0013942
Low_spent_Large_value_payments-High_spent_Large_value_payments    -0.0463060945 -0.09047957 -2.132616e-03 0.0327612
Low_spent_Medium_value_payments-High_spent_Large_value_payments   -0.0579479589 -0.09890449 -1.699143e-02 0.0006014
Low_spent_Small_value_payments-High_spent_Large_value_payments    -0.1155636763 -0.15160911 -7.951824e-02 0.0000000
High_spent_Small_value_payments-High_spent_Medium_value_payments  -0.0386567196 -0.07952263  2.209186e-03 0.0778212
Low_spent_Large_value_payments-High_spent_Medium_value_payments   -0.0268446594 -0.06879038  1.510106e-02 0.4890324
Low_spent_Medium_value_payments-High_spent_Medium_value_payments  -0.0384865238 -0.07702980  5.675586e-05 0.0506283
Low_spent_Small_value_payments-High_spent_Medium_value_payments   -0.0961022412 -0.12938017 -6.282431e-02 0.0000000
Low_spent_Large_value_payments-High_spent_Small_value_payments     0.0118120602 -0.03424439  5.786851e-02 0.9888779
Low_spent_Medium_value_payments-High_spent_Small_value_payments    0.0001701958 -0.04281047  4.315086e-02 1.0000000
Low_spent_Small_value_payments-High_spent_Small_value_payments    -0.0574455216 -0.09577532 -1.911572e-02 0.0002009
Low_spent_Medium_value_payments-Low_spent_Large_value_payments    -0.0116418644 -0.05565048  3.236675e-02 0.9869055
Low_spent_Small_value_payments-Low_spent_Large_value_payments     -0.0692575817 -0.10873662 -2.977854e-02 0.0000048
Low_spent_Small_value_payments-Low_spent_Medium_value_payments    -0.0576157174 -0.09345893 -2.177251e-02 0.0000441
```

Figure 13: ANOVA between age and payment behavior

Since P score <2e-16 lower than 0.05 payment behavior is associated with variations on age. Individuals with "High_spent_Large_value_payments" tend to be older compared to those with "High_spent_Medium_value_payments" and "High_spent_Small_value_payments" (p < 0.05).There are significant age differences between various combinations of payment behavior categories, such as "Low_spent_Small_value_payments" versus "High_spent_Small_value_payments" (p < 0.05).

**Question 5) How accurate a multinominal regression model can classify credit score correctly?**

Since credit score data aim to build a model which classifies credit scores. We investigate that question. Firstly, we encode all categorical variables for putting in the model. Secondly, we split the data as a train and test data. Thirdly, we use 5-fold cross validation. Lastly, We

```
              Reference
Prediction    1    2     3
         1 2073 1216    53
         2 1901 5759 1314
         3  341  925 1252

Overall Statistics

               Accuracy : 0.6124
                 95% CI : (0.6045, 0.6202)
    No Information Rate : 0.5326
    P-Value [Acc > NIR] : < 2.2e-16

                  Kappa : 0.3343

 Mcnemar's Test P-Value : < 2.2e-16

Statistics by Class:

                     Class: 1 Class: 2 Class: 3
Sensitivity            0.4804   0.7290   0.4780
Specificity            0.8794   0.5363   0.8964
Pos Pred Value         0.6203   0.6417   0.4972
Neg Pred Value         0.8049   0.6346   0.8890
Prevalence             0.2909   0.5326   0.1766
Detection Rate         0.1397   0.3882   0.0844
Detection Prevalence   0.2253   0.6050   0.1697
Balanced Accuracy      0.6799   0.6327   0.6872
```

build our model by using preprocess("center", "scale", "pca") in the multinom function.

Figure 14: Prediction of model in the test set and confusion matrix

The confusion matrix shows the counts of correct and incorrect predictions for each class. It looks good but it can be better also the model accuracy is 61.24%. Model accuracy achieved

by always predicting the most frequent class around 53.26%. A measure of agreement between actual and predicted classes, adjusted for chance agreement. Here, it's around 33.43%, this means fair agreement. Sensitivity is higher for class 2 as 72.9% but sensitivity in the other class below 50% so we should be deal with that in the following process.