

# Transunion-data-science-analytics-hiring-challenge-

## 2022-Approach

Machine learning model to classify the credit score based on people bank details and credit related information.

### 1.Exploratory Data Analysis

- Pandas, seaborn, matplotlib libraries are used in Exploratory data analysis.

### 2.Data Pre-Processing

- The following columns have non\_numeric values. So replaced all the numeric values and convert the column to numerical data type. Also, removed the unusual values.
  - Age
  - Annual income
  - Num\_of\_loan
  - Num\_of\_delayed\_payment
  - Changed\_credit\_limit
  - Amount\_invested\_monthly
  - Monthly\_balance
- The following column has special characters. So replaced all the special characters.
  - Occupation
  - Credit\_mix
  - Payment\_behaviour
- The credit history age column has text and numbers. converted into numeric column by extracting numbers.
- The type of loan column has comma-separated text. So using the TF-IDF vectorizer convert the columns values to numeric format.
- Missing values are replaced by mean and group mean.

- Create a missing row indicator for the columns which have missing values.

.

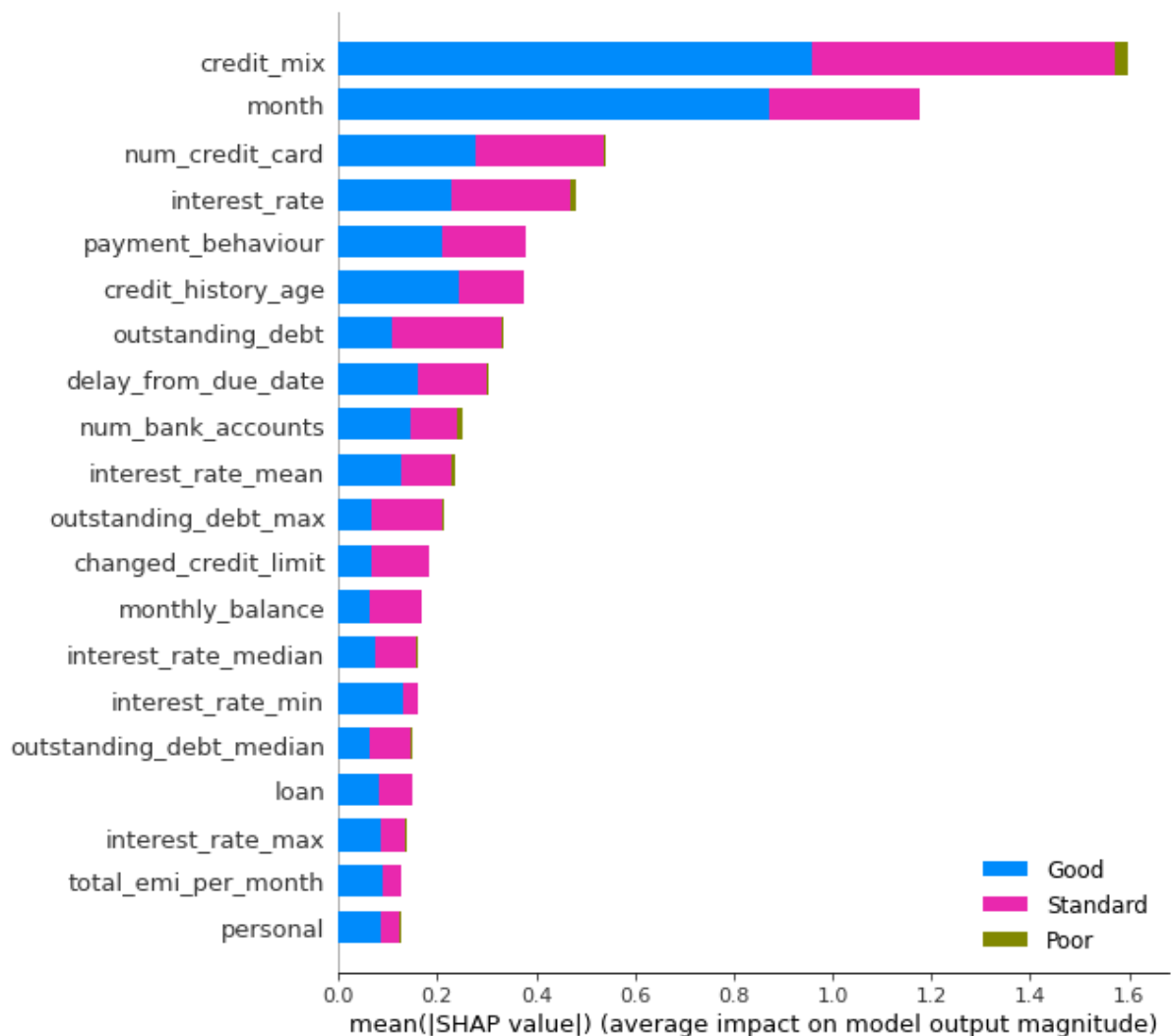
### 3.Model

- After pre-processing finally 113 columns are selected for the classification model.
- Selected columns are,
  - 0\_month
  - 1\_age
  - 2\_annual\_income
  - 3\_monthly\_inhand\_salary
  - 4\_num\_bank\_accounts
  - 5\_num\_credit\_card
  - 6\_interest\_rate
  - 7\_num\_of\_loan
  - 8\_delay\_from\_due\_date
  - 9\_num\_of\_delayed\_payment
  - 10\_changed\_credit\_limit
  - 11\_num\_credit\_inquiries
  - 12\_outstanding\_debt
  - 13\_credit\_utilization\_ratio
  - 14\_credit\_history\_age
  - 15\_payment\_of\_min\_amount
  - 16\_total\_emi\_per\_month
  - 17\_amount\_invested\_monthly
  - 18\_payment\_behaviour
  - 19\_monthly\_balance
  - 20\_name\_isnull
  - 21\_age\_isnull
  - 22\_ssn\_isnull
  - 23\_occupation\_isnull
  - 24\_monthly\_inhand\_salary\_isnull
  - 25\_num\_bank\_accounts\_isnull
  - 26\_num\_credit\_card\_isnull
  - 27\_interest\_rate\_isnull

- 28\_num\_of\_loan\_isnull
- 29\_type\_of\_loan\_isnull
- 30\_num\_of\_delayed\_payment\_isnull
- 31\_changed\_credit\_limit\_isnull
- 32\_num\_credit\_inquiries\_isnull
- 33\_credit\_mix\_isnull
- 34\_credit\_history\_age\_isnull
- 35\_amount\_invested\_monthly\_isnull
- 36\_payment\_behaviour\_isnull
- 37\_monthly\_balance\_isnull
- 38\_credit\_score\_isnull
- 39\_ssn
- 40\_occupation
- 41\_type\_of\_loan
- 42\_credit\_mix
- 43\_auto
- 44\_builder
- 45\_consolidation
- 46\_credit
- 47\_debt
- 48\_equity
- 49\_home
- 50\_loan
- 51\_mortgage
- 52\_payday
- 53\_personal
- 54\_specified
- 55\_student
- Created catboost classifier with 5-fold stratified cross validation and tuned the hyperparameters with optuna framework.
- Fitted catboost classifier with tuned parameters

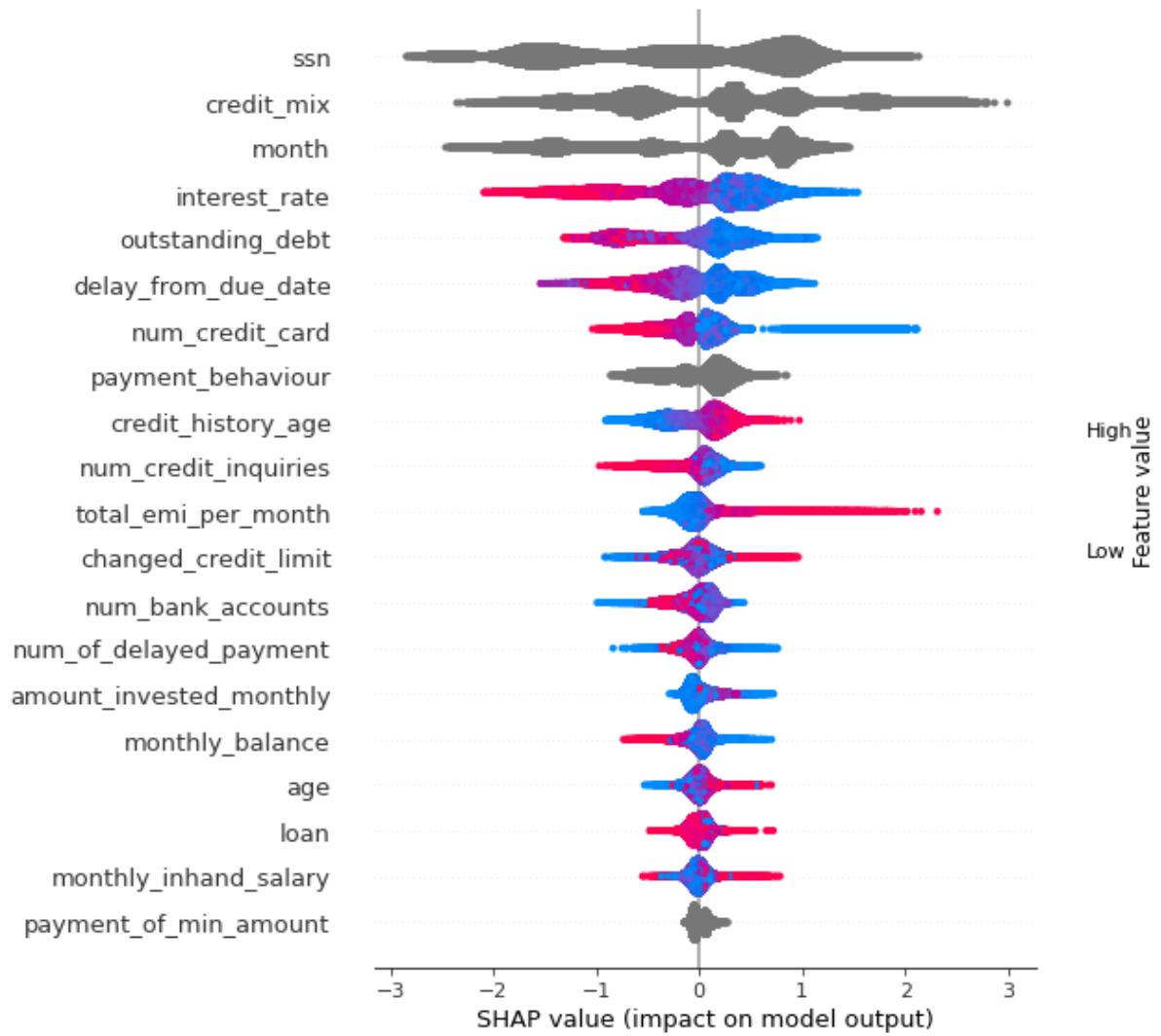
- Parameters are,  
`{'reg_lambda': 0.9138676838297956,`  
`'learning_rate': 0.08135287540629096,`  
`'n_estimators': 887,`  
`'max_depth': 10,`  
`'random_state': 2020,`  
`'boosting_type': 'Plain',`  
`'bootstrap_type': 'Bernoulli',`  
`'subsample': 0.9414567528861059}`

- Model Explanation with shap library
- Feature importance plot

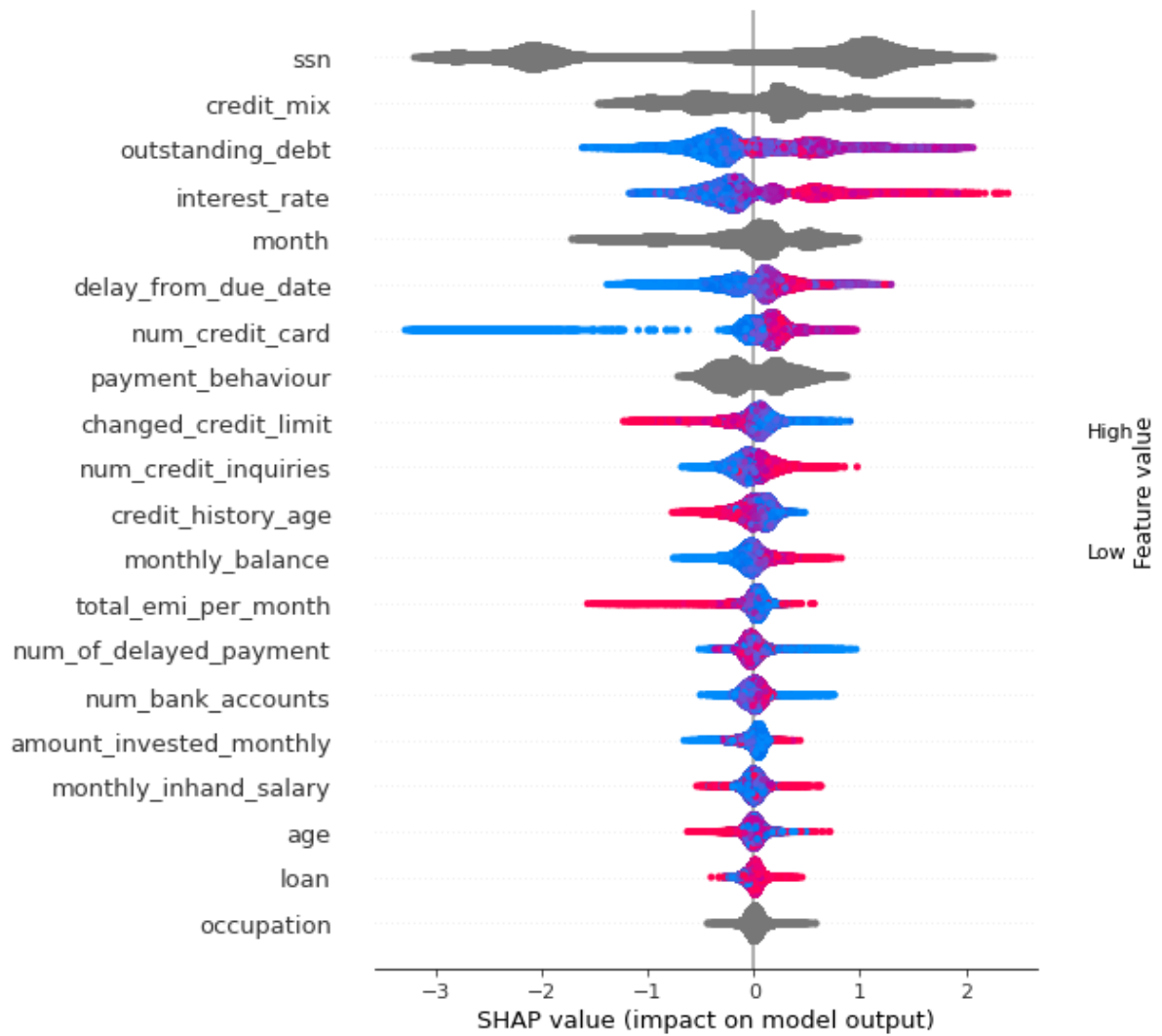


## Top features impact the model

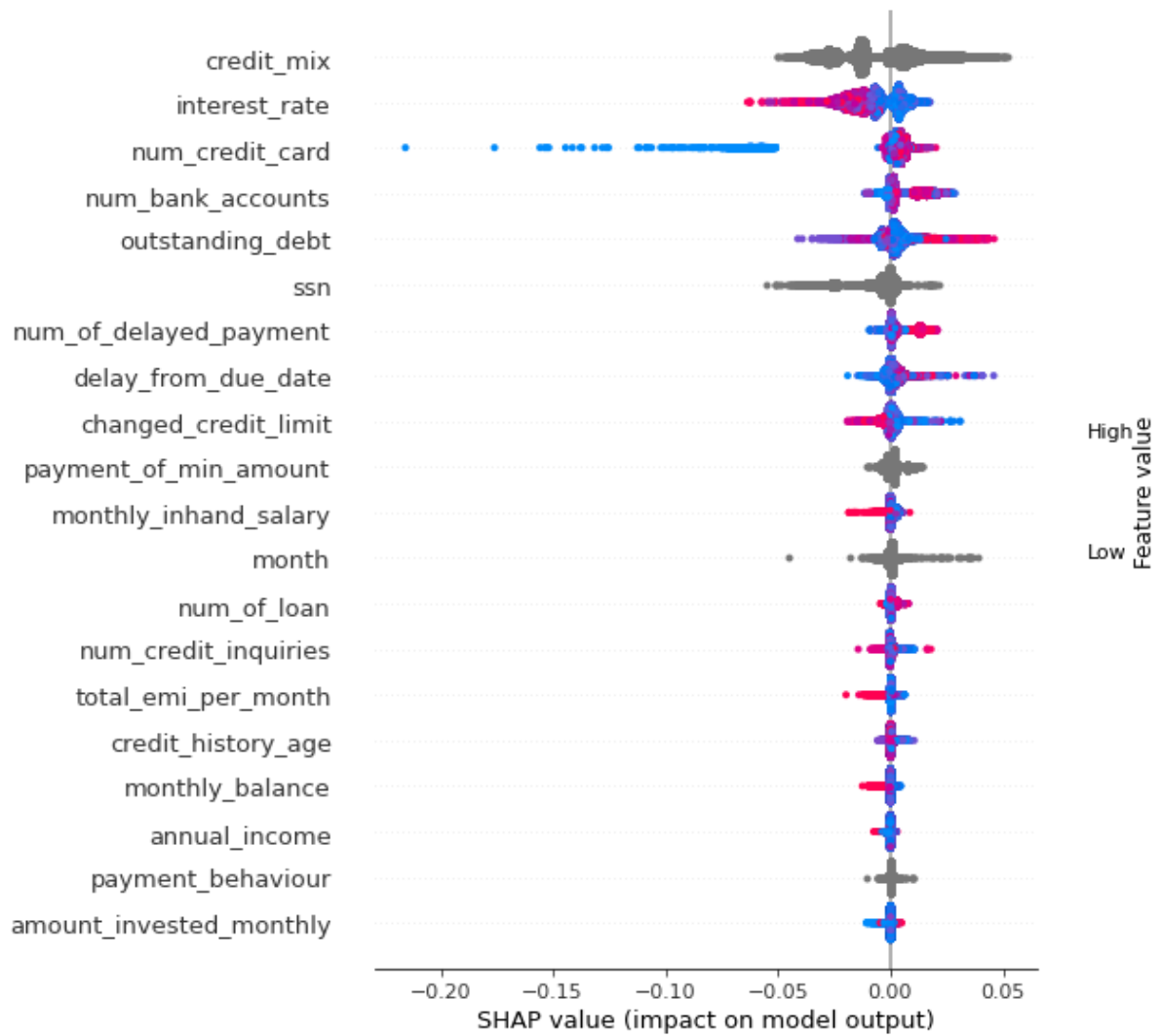
- For Good Class



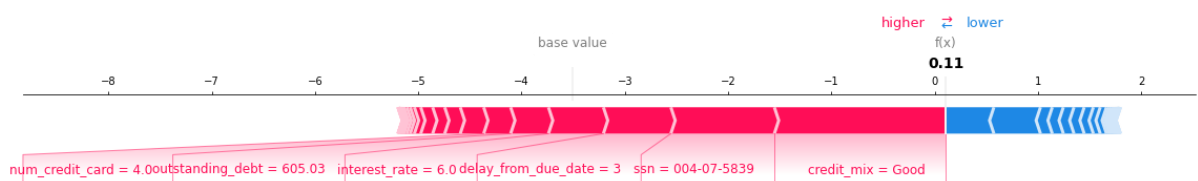
○ For Standard Class



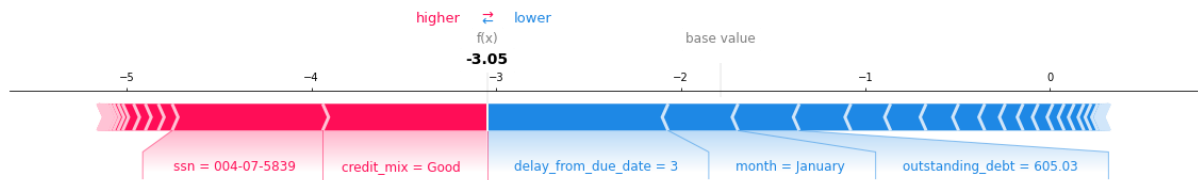
- For Poor Class



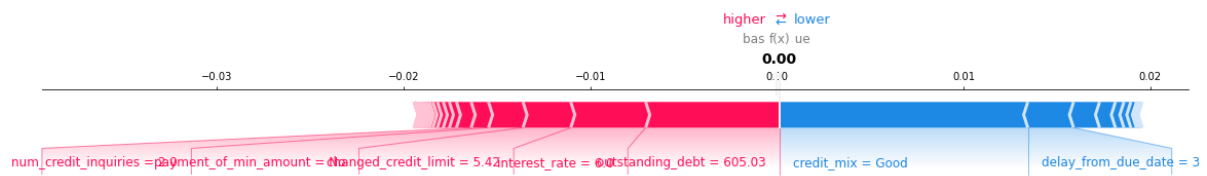
- Feature influences the model prediction for a single observation in each class
  - For Good Class



## ○ For Standard Class



## ○ For Poor Class



- Final score is 78.03