# BFS CAPSTONE PROJECT

## MID-SUBMISSION

# Objective

**CredX** is a leading credit card provider that gets thousands of credit card applications every year. But in the past few years, it has experienced an increase in credit loss. The CEO believes that the best strategy to mitigate credit risk is to acquire the right customers.

In this project, our task is to help CredX identify the right customers using predictive models. Using past data of the bank's applicants, we need to determine the factors affecting credit risk, create strategies to mitigate the acquisition risk and assess the financial benefit of your project.
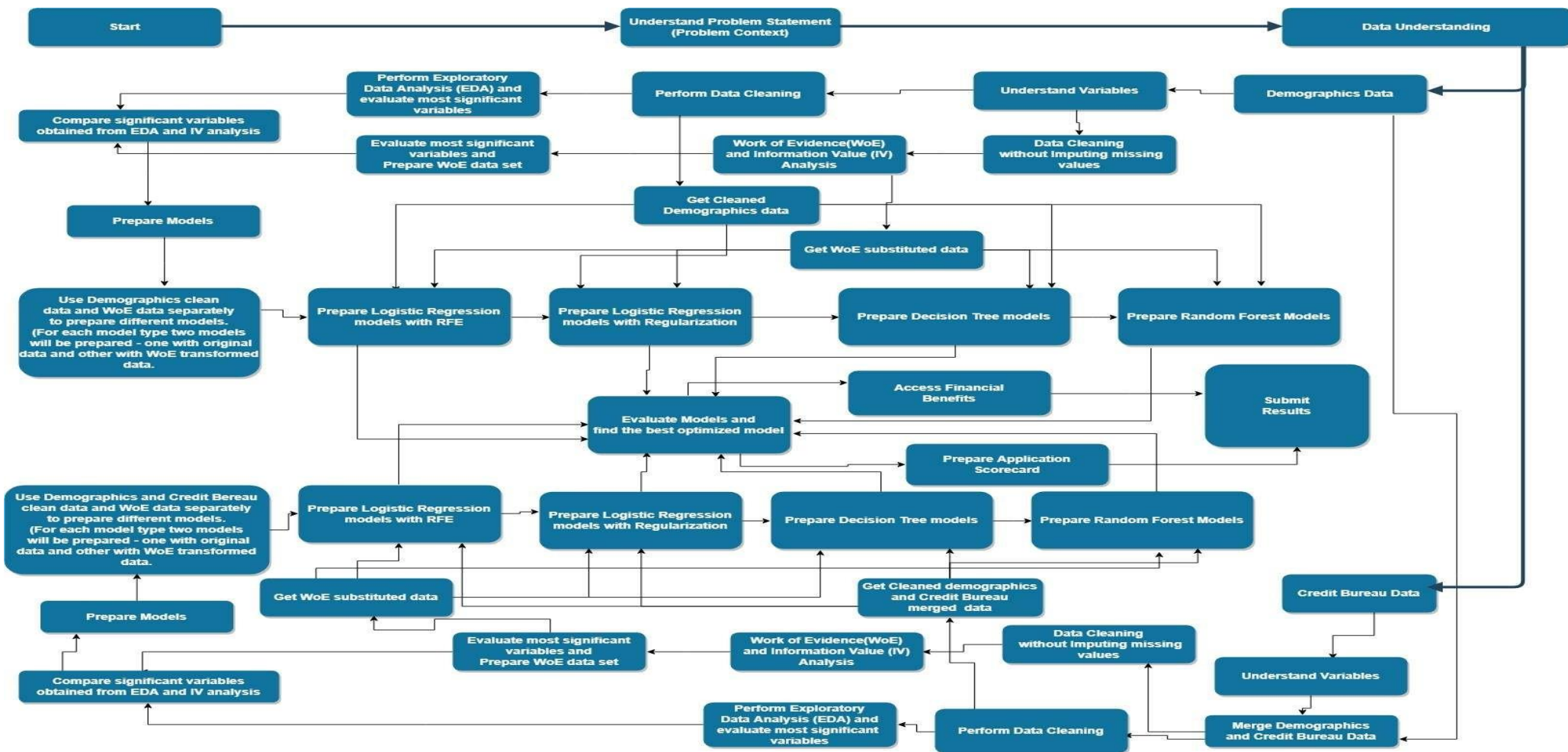
# Steps to Problem Solving

- Understand the underlying problem and the domain for the problem.

- Understand the dataset provided (both Demographic and Credit Bureau) and inspect each attribute of both.

- Perform data cleaning on both the data sets. For input dataset of IV analysis, imputing values needs to be ignored.

- Choose each data-set individually and perform Exploratory Data Analysis to anticipate the significant variables.

- Work of Evidence (WoE) and Information Value (IV) analysis and preparing WoE transformed dataset.
    - ☐ Take Demographic data-set and perform WoE transformation and also obtain significant variables based on IV.
    - ☐ Merge Demographic data-set with Credit Bureau dataset and do WoE transformation and get significant variables.

- Use both the original clean data-set and WoE transformed data set of demographics separately to prepare data models. For this bi-logit problem model preparation, begin with simple models like Logistic Regression model with RFE and step by step move on to relatively complex models like Logistic Regression with Regularization, Decision Tree, Random Forest etc. Following steps needs to be considered in each model building process:
    - ☐ Initially the dataset needs to be divided into Test and Train dataset.
    - ☐ There is a class imbalance in the dataset. This needs to be handled using balanced class during each model preparation.
    - ☐ Cross Validation needs to be done for each model.
    - ☐ Additional validation of data should be done on the dataset on the rejected applications (performance tag null) ignored for model building.
    - ☐ Hyperparameters for each type of models need to be optimized properly using GridSearch and model coming with optimized parameter should be chosen.

# Steps to Problem Solving (continued ...)

- The above step also needs to be carried out for merged data set of demographic and Credit Bureau.
- Evaluate all the models based on the following parameters :
  - ☐ Confusion matrix should be prepared for each model.
  - ☐ Sensitivity, specificity, accuracy curve for each model with different cut-offs.
  - ☐ AUC-ROC curve for the model using cut-off values for each model.
  - ☐ Precision and Recall curve for cut-off should be generated.
  - ☐ Gini-Index needs to be evaluated for Tree based models like decision tree and random forest.
  - ☐ Within each model type evaluation using GridSerach based on recall values should be done to get models with optimized hyperparameters.
  - ☐ For evaluation among models, the dataset for rejected applications (with performance tag missing), which were assumed as potentially defaulters should be considered for evaluations. Ideally, the output for all these applications should be defaulters.
  - ☐ Choose an easy model but not the easiest one.
- The apt two stable and optimized models (with stable characteristics) - one for demographics and second for combined data needs to be chosen.
- On the basis of the chosen model and significant variables in the model, two application scorecard should be prepared for the two models.
- Access the financial benefits of the project by checking the underlying matrices that get optimized.
- Present all the results obtained in all the above steps to the management.

# Problem Solving Methodology

# Data Understanding

There are two data sets in this project: Demographic and Credit bureau data.

- **Demographic/application data:** This is obtained from the information provided by the applicants at the time of credit card application. It contains customer-level information on age, gender, income, marital status, etc.

- **Credit bureau data:** This is taken from the credit bureau and contains variables such as 'number of times 30 DPD or worse in last 3/6/12 months', 'outstanding balance', 'number of trades', etc.

Both files contain a performance tag, which indicates whether the applicant has gone 90 days past due (DPD) or worse in the past 12 months (i.e. defaulted) after getting a credit card.

In some cases, it is observed that all the variables in the credit bureau data are zero and credit card utilization is missing. These represent cases in which there is a no-hit in the credit bureau. The cases with missing credit card utilization are also observed. These are the cases in which the applicant does not have any other credit card.

# Data Understanding and Handling Data Issues for Demographic Data

| Variable Name | Description | Data Issues and their handling |
|---|---|---|
| Application Id | Unique Ids of the customer | 3 non-unique Application Id values (6 records) were found. On manual check, it was found that the underlying customers with both records were different. All the 6 records, being less in number were removed considering them as junk data. |
| Age | Age of customer | • One customer had negative age (age was -3).<br>• 19 customers had age as 0.<br>• 45 customers had age less than 18 (assuming a person with age less than 18 is not eligible for applying the card).<br>The age for all the above record was set as 18. |
| Gender | Gender of Customer | Gender was missing for 2 customers. These were imputed as 'M' based on maximum value counts (mode). |
| Marital Status | Marital status of customer (at the time of application) | Marital Status was missing for 6 customers. Missing values were imputed as 'Married' based on maximum value counts (mode). |
| Income | Income of customers | Income was negative for 81 customers. This negative count was imputed with median value. |

# Data Understanding and Handling Data Issues for Demographic Data (continued ...)

| Variable Name | Description | Data Issues and their handling |
|---|---|---|
| No. of Dependents | No. of children of customers | No. of children were missing for 3 customers. All these 3 customers had Marital status 'Married'. So, the values were imputed with value 3 based on maximum occurring value. |
| Education | Education of customers | Education was missing for 119 customers. Missing values were imputed with value 'UNKNOWN'. |
| Profession | Profession of customers | Profession was missing for 14 customers. Missing values were imputed with 'SAL' based on maximum occurring values. |
| Type of Residence | Type of residence of customers | Residence type was missing for 8 customers. Missing values were imputed with 'Rented' based on maximum occurring values. |
| No. of months in current residence | No of months in current residence of customers | |
| No of months in current company | No of months in current company of customers | |
| Performance Tag | Status of customer performance (" 1 represents "Default") | Performance Tag was missing for 1425 customers. All these customers were removed from the data set as Credit Card was never issued to these customers. |

# Data Understanding and Handling Data Issues for Credit Bureau Data

| Variable Name | Description | Data Issues and their handling |
|---|---|---|
| Application Id | Customer application ID | 3 non-unique Application Id values (6 records) were found. On manual check, it was found that the underlying customers with both records were different. All the 6 records, being less in number were removed considering them as junk data. |
| No of times 90 DPD or worse in last 6 months | Number of times customer has not payed dues since 90 days in last 6 months | |
| No of times 60 DPD or worse in last 6 months | Number of times customer has not payed dues since 60 days in last 6 months | |
| No of times 30 DPD or worse in last 6 months | Number of times customer has not payed dues since 30 days in last 6 months | |
| No of times 90 DPD or worse in last 12 months | Number of times customer has not payed dues since 90 days in last 12 months | |
| No of times 60 DPD or worse in last 12 months | Number of times customer has not payed dues since 60 days in last 12 months | |

# Data Understanding and Handling Data Issues for Credit Bureau Data (continued ...)

| Variable Name | Description | Data Issues and their handling |
|---|---|---|
| No of times 30 DPD or worse in last 12 months | Number of times customer has not payed dues since 30 days in last 12 months | |
| Avgas CC Utilization in last 12 months | Average utilization of credit card by customer | This was missing for 1058 customers. This missing values were imputed by median value except for the values where Outstanding Balance was missing (in that scenario it was imputed with 0). |
| No of trades opened in last 6 months | Number of times the customer has done the trades in last 6 months | This was missing for 1 customer. Missing value was imputed with the maximum occurring value. |
| No of trades opened in last 12 months | Number of times the customer has done the trades in last 12 months | |
| No of PL trades opened in last 6 months | No of PL trades in last 6 month of customer | |
| No of PL trades opened in last 12 months | No of PL trades in last 12 month of customer | |
| No of Inquiries in last 6 months (excluding home & auto loans) | Number of times the customers has inquired in last 6 months | |

# Data Understanding and Manipulation for Credit Bureau Data (continued ...)

| Variable Name | Description | Data Issues and their handling |
|---|---|---|
| No of Inquiries in last 12 months (excluding home & auto loans) | Number of times the customers has inquired in last 12 months | |
| Presence of open home loan | If the customer has home loan (1 represents "Yes") | This was missing for 272 customers. This was assumed that these haven't taken home loan so imputed with 0. |
| Outstanding Balance | Outstanding balance of customer | This was missing for 272 customers. This was assumed that these haven't taken/applied any credit card yet so imputed with 0. |
| Total No of Trades | Number of times the customer has done total trades | |
| Presence of open auto loan | If the customer has auto loan (1 represents "Yes") | |
| Performance Tag | Status of customer performance (" 1 represents "Default") | Performance Tag was missing for 1425 customers. All these customers were removed from the data set as Credit Card was never issued to them. |

# EDA and WoE / IV Analysis for Demographic Data

```
iv_demographics.sort_values(by='IV', ascending = False)
```

Following predictor variables were obtained as part of EDA and WoE Analysis for Demographic data

(top five based on IV values).

- Months_Current_Residence
- Income
- Months_Current_Company
- Age
- Dependents_No

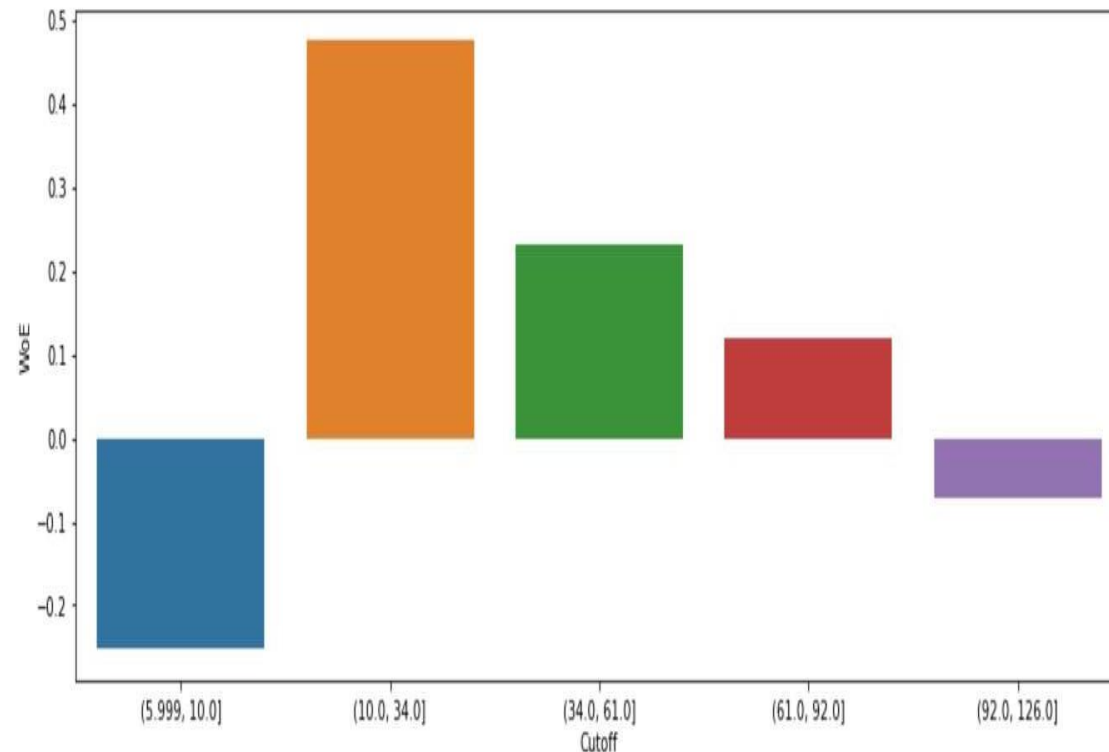The overall information value of the data set was **0.1467**.

|   | Variable | IV |
|---|---|---|
| 0 | Months_Current_Residence | 0.073546 |
| 0 | Income | 0.039054 |
| 0 | Months_Current_Company | 0.021682 |
| 0 | Age | 0.004896 |
| 0 | Dependents_No | 0.002818 |
| 0 | Profession | 0.002221 |
| 0 | Residence_Type | 0.000942 |
| 0 | Education | 0.000782 |
| 0 | Gender | 0.000568 |
| 0 | Marital_Status | 0.000147 |

# Understanding Months_Current_Residence as predictor demographic variable

The WoE values across rising in bins show monotonic decrease in WoE as months of current residence increase across bins (except the lowest value bin).

Similar trend of monotonic decrease is also observed in the bar plot for the bins created for Number of months in current residence.

```
plt.figure(figsize=(15,5))
sns.barplot(y='WoE', x="Cutoff", data=woe_demographics[woe_demographics.Variable=='Months_Current_Residence'])
plt.show()
```



Percent Stack Plot for Months_Current_Residence_bins

Default_rate values by bin:
- (5.999, 10.0]: 3.31 %
- (10.0, 34.0]: 6.61 %
- (34.0, 61.0]: 5.25 %
- (61.0, 92.0]: 4.72 %
- (92.0, 126.0]: 3.94 %

# Understanding Income as predictor demographic variable

The WoE and bins plot chart shows monotonic decrease in default rate as income increase across bins.

Similar trend is observed across Income bins plotted along the bar plot which shows decrease in percentage of defaulters across income bins.

```
plt.figure(figsize=(15,5))
sns.barplot(y='WoE', x="Cutoff", data=woe_demographics[woe_demographics.Variable=='Income'])
plt.show()
```
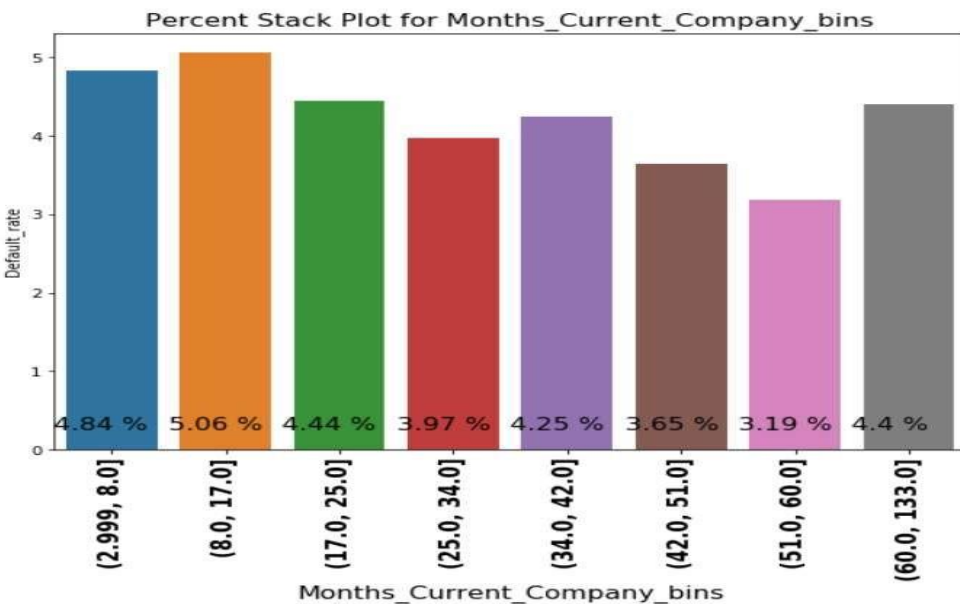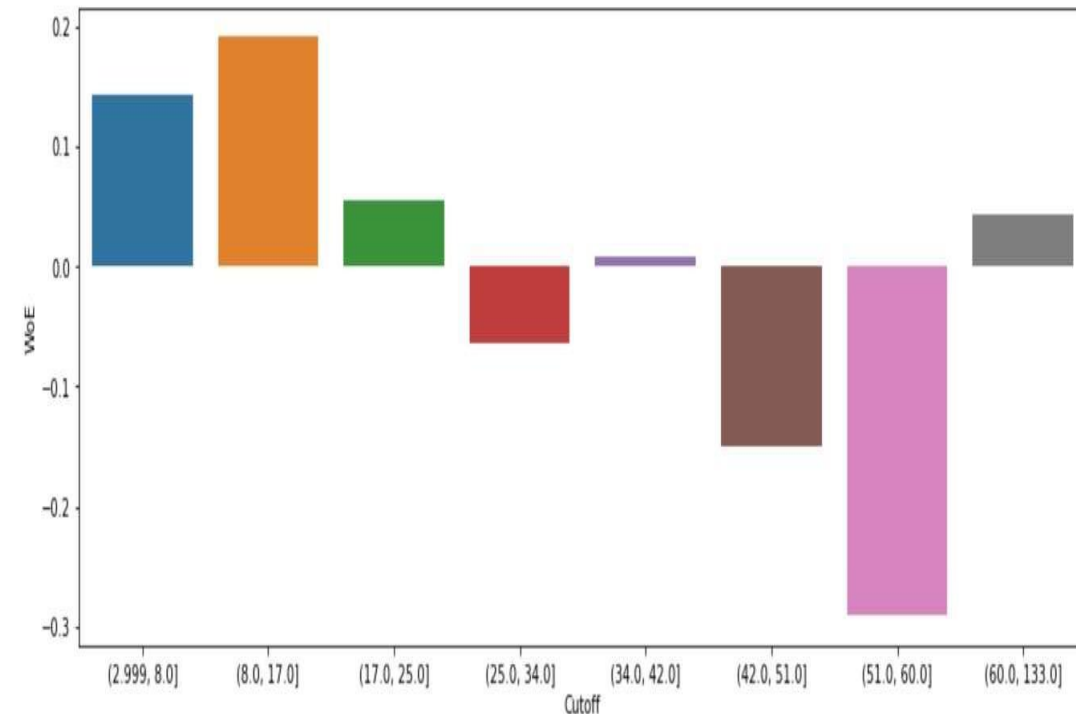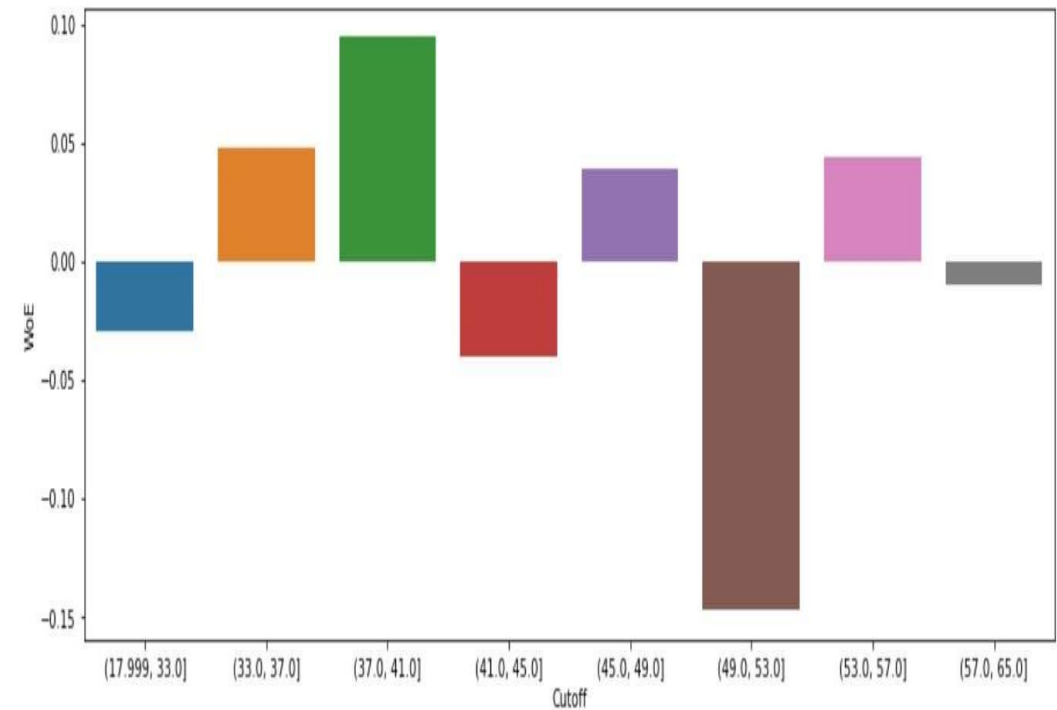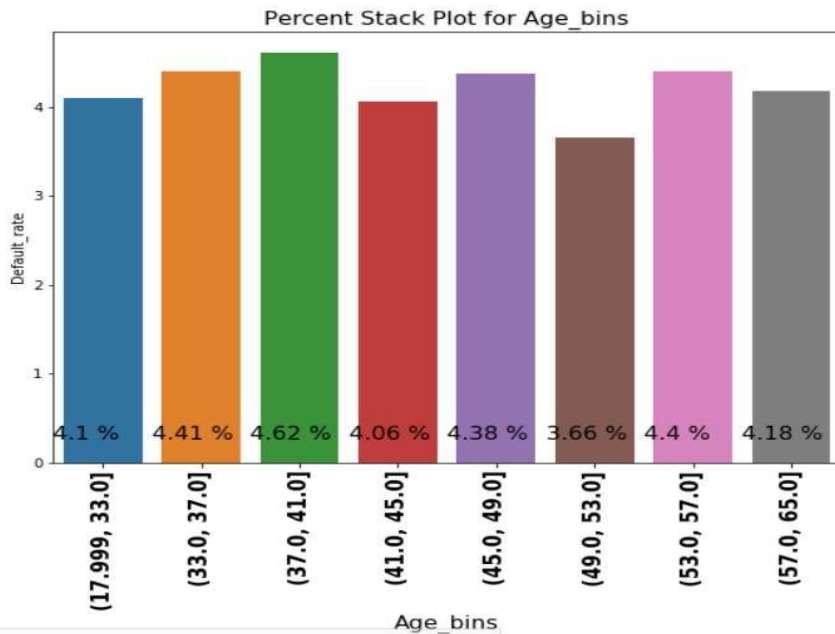


Percent Stack Plot for Income_bins

# Understanding Months_Current_Company as predictor demographic variable

The WoE and bins plot chart shows monotonic decrease in default rate as Months_Current_Company increase across bins (with some exceptions).

Similar trend is observed across Months_Current_Company bins plotted along the bar plot which shows decrease in percentage of defaulters across bins.

```
plt.figure(figsize=(15,5))
sns.barplot(y='WoE', x="Cutoff", data=woe_demographics[woe_demographics.Variable=='Months_Current_Company'])
plt.show()
```



Percent Stack Plot for Months_Current_Company_bins

4.84 % 5.06 % 4.44 % 3.97 % 4.25 % 3.65 % 3.19 % 4.4 %

# Understanding Age as predictor demographic variable

The WoE and bins plot chart shows variations across Age group bins.

Similar kind of variations are also observed in bar plot for Age bins created.

```
plt.figure(figsize=(15,5))
sns.barplot(y='WoE', x="Cutoff", data=woe_demographics[woe_demographics.Variable=='Age'])
plt.show()
```
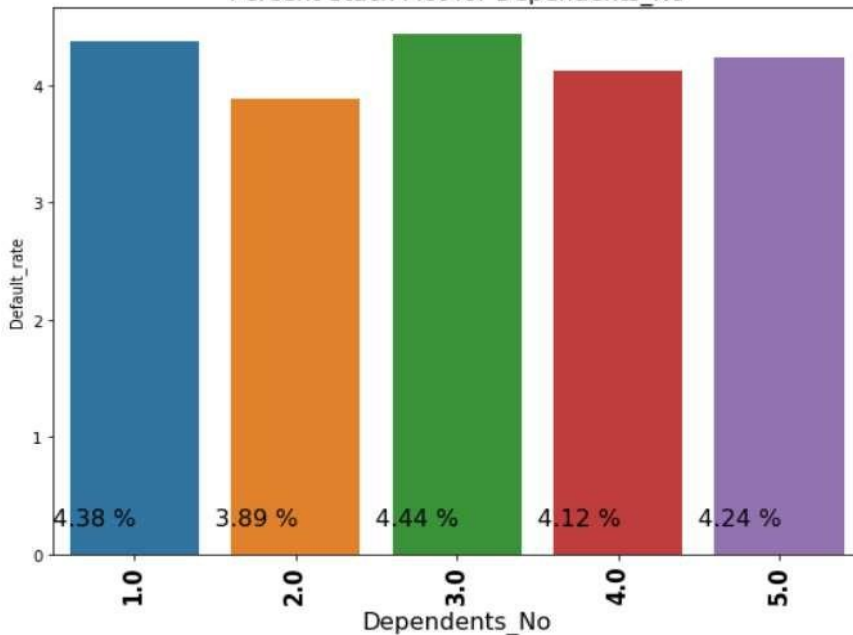


Percent Stack Plot for Age_bins

# Understanding Dependent No. as predictor demographic variable
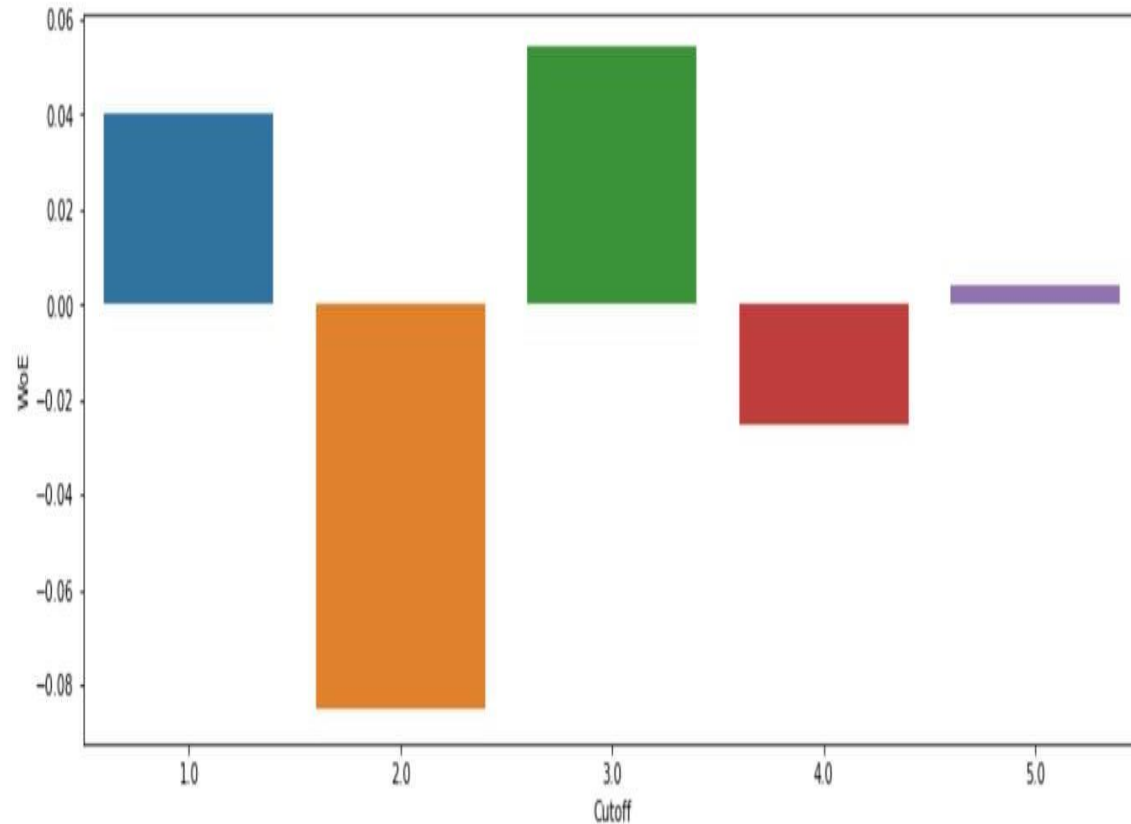
The Dependent No. is the fifth significant  variable as per IV analysis and there is a big  fall in the Information value as comparison  to fourth significant variable.

Although  this  variable  is  significant  as per IV   value.  the  bar  plot  is  unable  to

```
plt.figure(figsize=(15,5))
sns.barplot(y='WoE', x="Cutoff", data=woe_demographics1[woe_demographics1.Variable=='Dependents_No'])
plt.show()
```



Percent Stack Plot for Dependents_No

# EDA and WoE / IV Analysis for Credit Bureau Data

Following predictor variables were obtained as part of EDA and WoE Analysis for Credit Bureau data

(top five based on IV values).

- Avgas_CC_Utilization_12_months
- Trades_opened_last_12_months
- PL_Trades_opened_last_12_months
- Outstanding_Balance
- Inquiries_last_12_months

The overall information value of the data set was **3.287**.

```
iv_creditbureau.sort_values(by='IV', ascending = False)
```
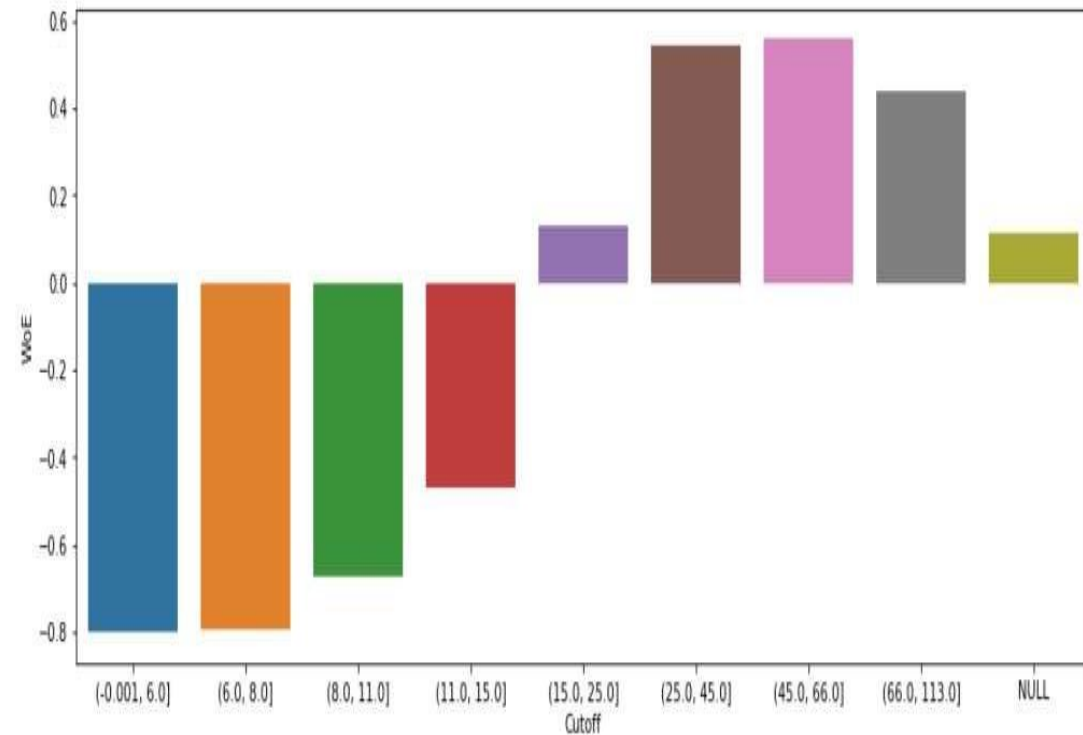
|   | Variable | IV |
|---|---|---|
| 0 | Avgas_CC_Utilization_12_months | 0.308660 |
| 0 | Trades_opened_last_12_months | 0.291635 |
| 0 | PL_Trades_opened_last_12_months | 0.255968 |
| 0 | Outstanding_Balance | 0.253333 |
| 0 | No_30_DPD_last_6_months | 0.244250 |
| 0 | Total_Trades | 0.238455 |
| 0 | PL_Trades_opened_last_6_months | 0.224219 |
| 0 | No_30_DPD_last_12_months | 0.218609 |
| 0 | No_90_DPD_last_12_months | 0.215653 |
| 0 | No_60_DPD_last_6_months | 0.211274 |
| 0 | No_60_DPD_last_12_months | 0.188230 |
| 0 | Trades_opened_last_6_months | 0.186148 |
| 0 | Inquiries_last_12_months | 0.172713 |
| 0 | No_90_DPD_last_6_months | 0.162659 |
| 0 | Inquiries_last_6_months | 0.113141 |
| 0 | Presence_open_auto_loan | 0.001655 |
| 0 | Presence_open_home_loan | 0.000462 |

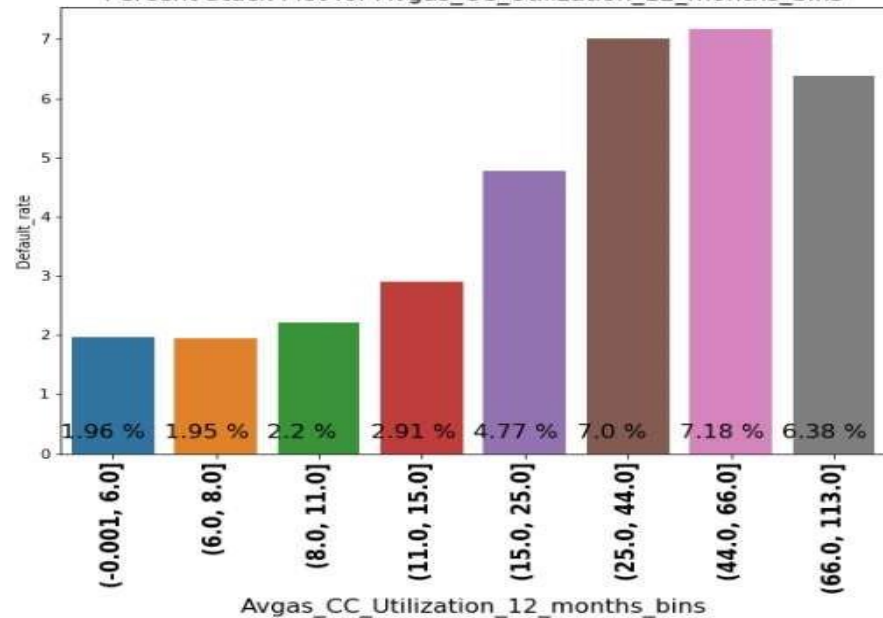# Understanding Avgas_CC_Utilization_12_months as predictor Credit Bureau variable

The WoE values across rising in bins show monotonic increase in WoE as Avgas_CC_Utilization_12_months increase across bins (except highest values bin).

Similar trend of monotonic increase in defaulter's percent is also observed in the bar plot for the bins created for Avgas_CC_Utilization_12_months.

```
plt.figure(figsize=(15,5))
sns.barplot(y='WoE', x="Cutoff", data=woe_creditbureau[woe_creditbureau.Variable=='Avgas_CC_Utilization_12_months'])
plt.show()
```


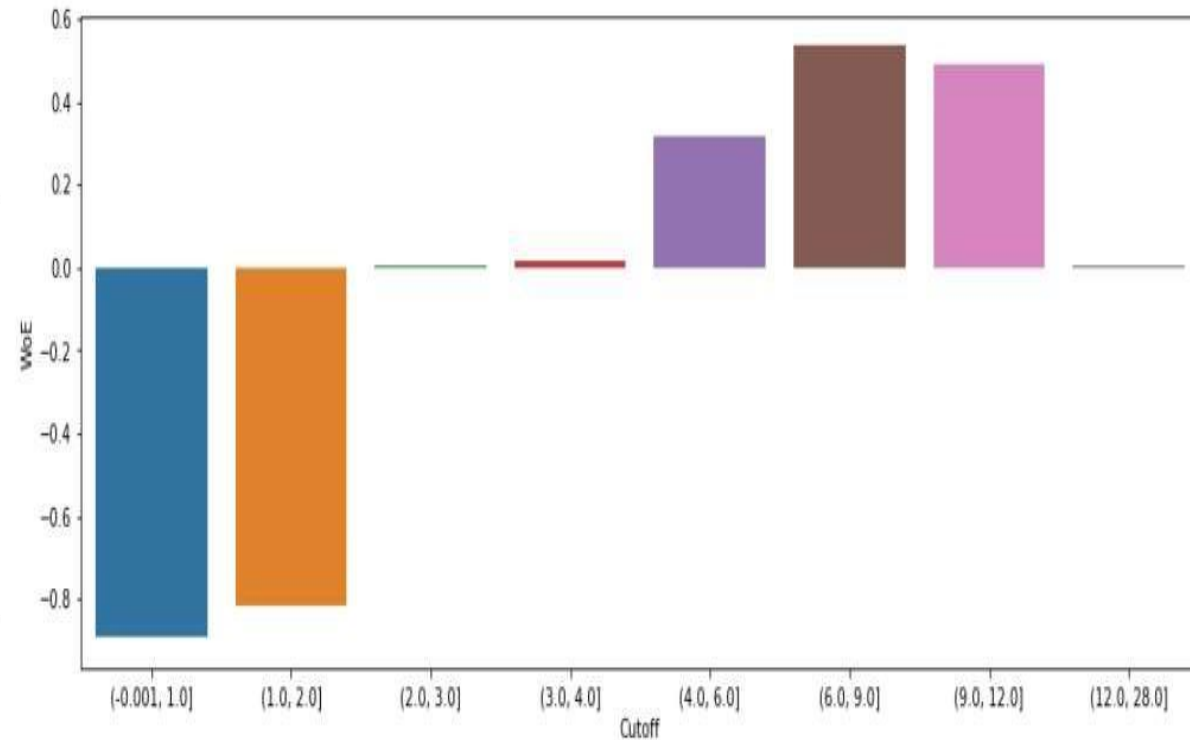Percent Stack Plot for Avgas_CC_Utilization_12_months_bins

# Understanding Trades_opened_last_12_months as predictor Credit Bureau variable

The WoE values across rising in bins show monotonic increase in WoE as Trades_opened_last_12_months increase across bins (except the highest value bins).

Similar trend of monotonic increase in default percent is also visible in the bar plot for the bins created for Trades_opened_last_12_months.
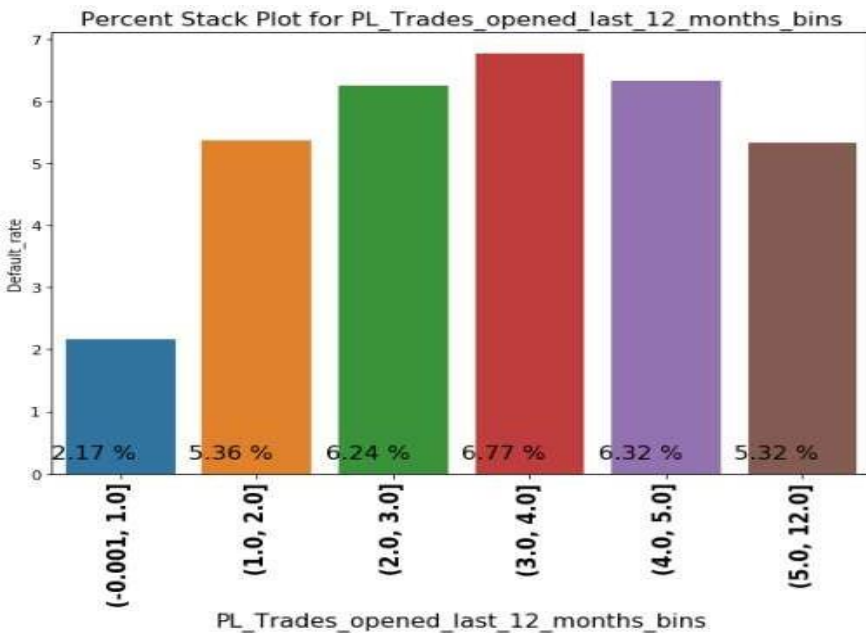
```
plt.figure(figsize=(15,5))
sns.barplot(y='WoE', x="Cutoff", data=woe_creditbureau[woe_creditbureau.Variable=='Trades_opened_last_12_months'])
plt.show()
```



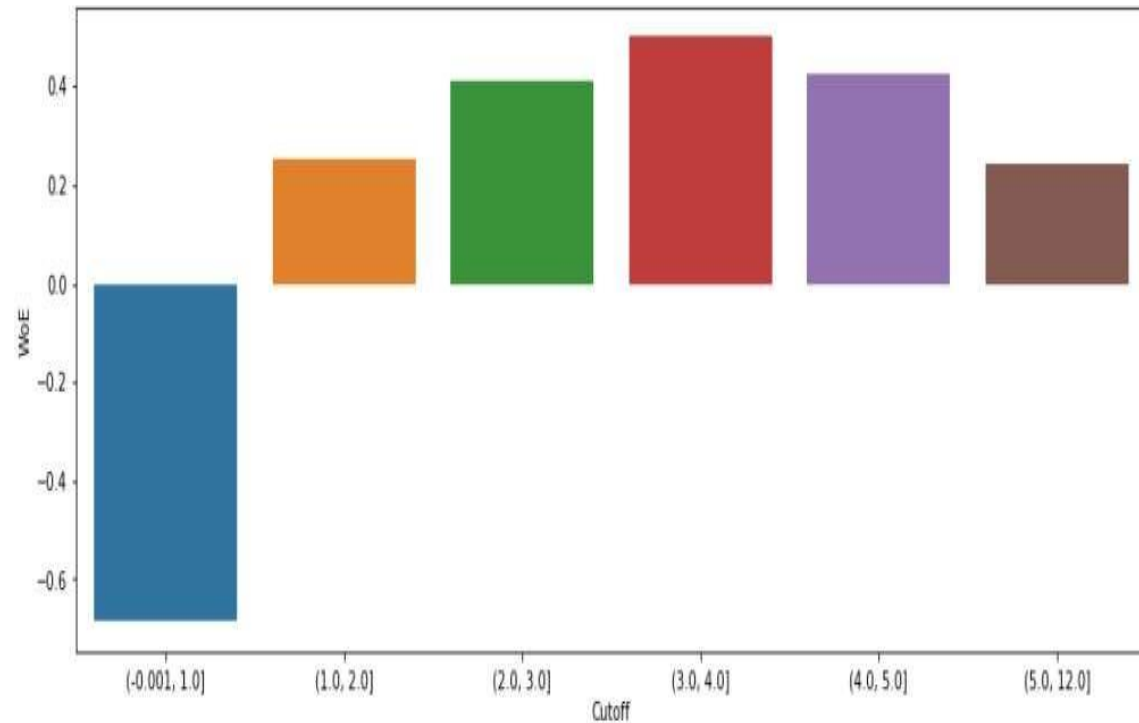Percent Stack Plot for Trades_opened_last_12_months_bins

# Understanding PL_Trades_opened_last_12_months as predictor Credit Bureau variable

The WoE values across rising in bins show monotonic increase and then monotonous decrease as PL_Trades_opened_last_12_months increase across bins.

Similar trend of monotonic increase and decrease in default percent is also observed in the bar plot for the bins created for PL_Trades_opened_last_12_months.
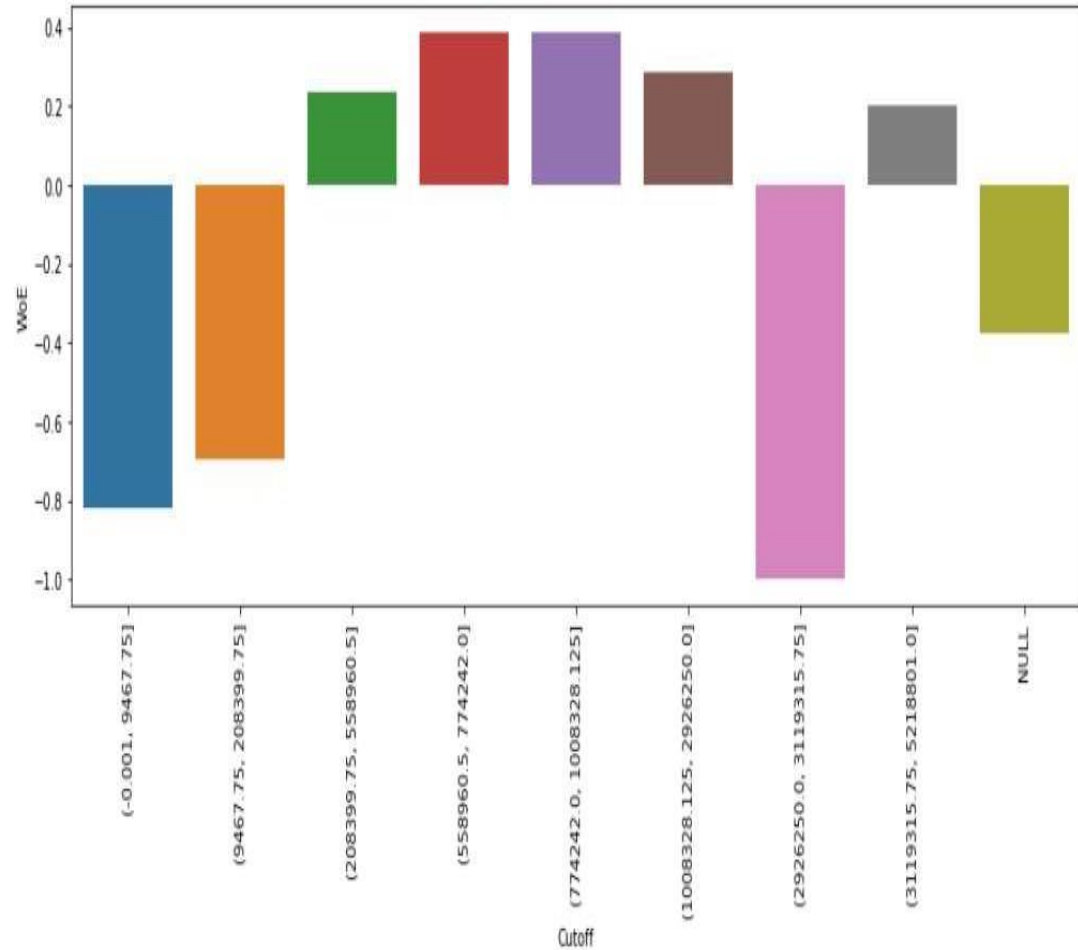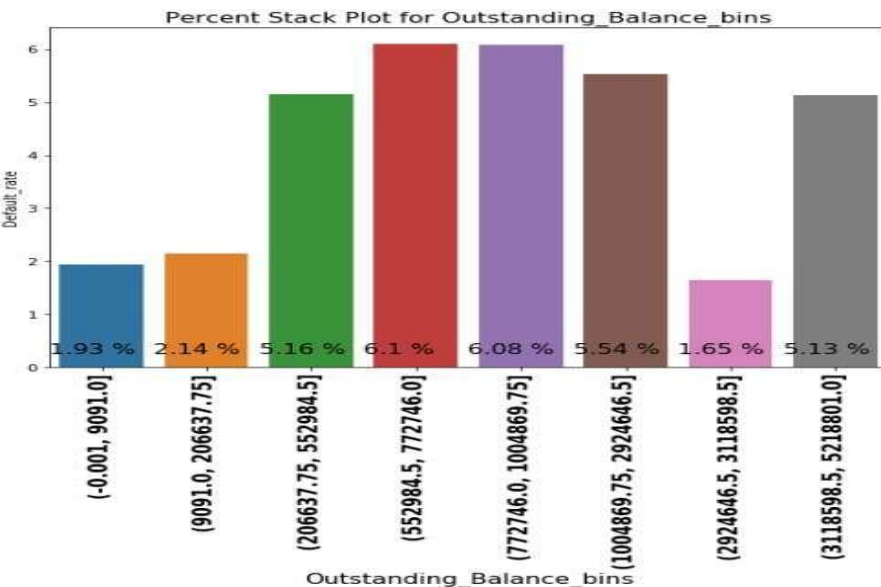
```python
plt.figure(figsize=(15,5))
sns.barplot(y='WoE', x="Cutoff", data=woe_creditbureau[woe_creditbureau.Variable=='PL_Trades_opened_last_12_months'])
plt.show()
```



Percent Stack Plot for PL_Trades_opened_last_12_months_bins

| Bin | Default rate |
|-----|--------------|
| (-0.001, 1.0] | 2.17 % |
| (1.0, 2.0] | 5.36 % |
| (2.0, 3.0] | 6.24 % |
| (3.0, 4.0] | 6.77 % |
| (4.0, 5.0] | 6.32 % |
| (5.0, 12.0] | 5.32 % |

# Understanding Outstanding_Balance as predictor Credit Bureau variable

The WoE values across rising in bins show monotonic increase in WoE and then monotonous decrease as Outstanding_Balance increase across bins.

Similar trend of monotonic increase and decrease in default percent is also observed in the bar plot for the bins created for Outstanding_Balance.



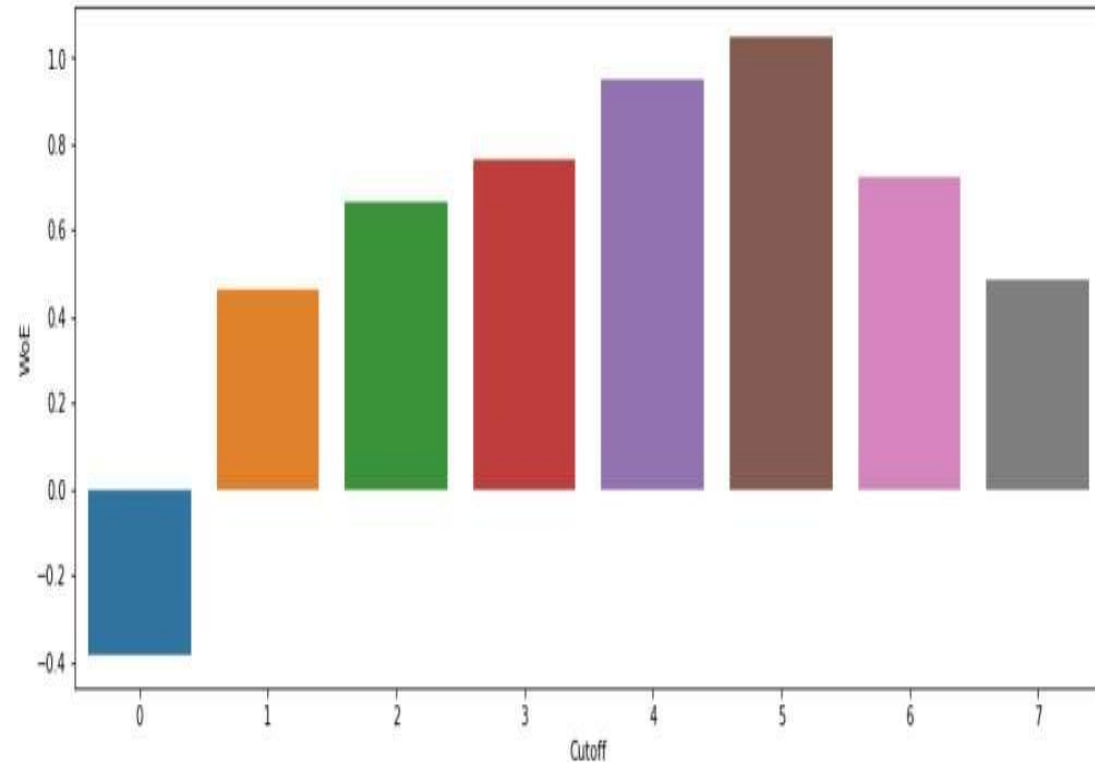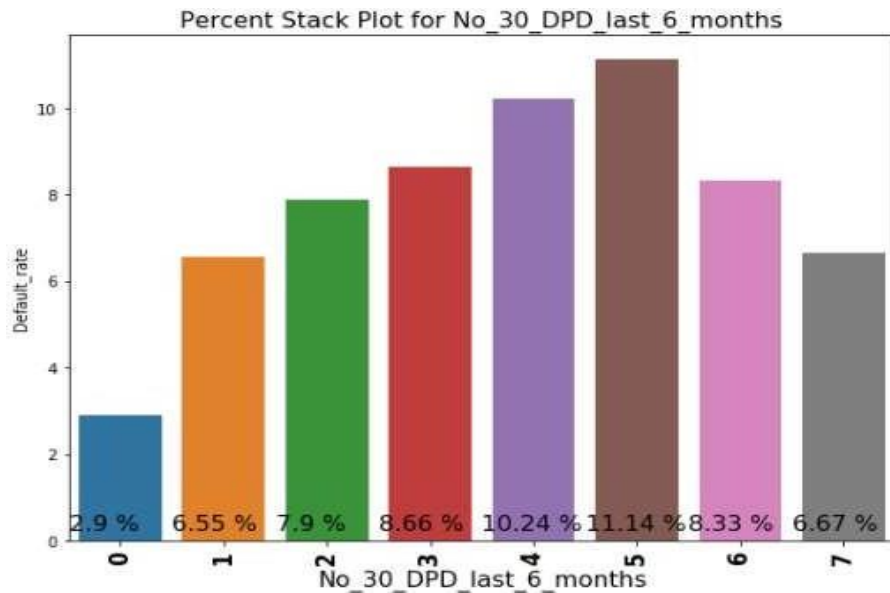Percent Stack Plot for Outstanding_Balance_bins

# Understanding No_30_DPD_last_6_months as predictor Credit Bureau variable

The WoE values across rising in bins show monotonic increase and then monotonous decrease as No_30_DPD_last_6_months increase across bins.

Similar trend of monotonic increase and decrease in default percent is also observed in the bar plot for the bins created for No_30_DPD_last_6_months.

```python
plt.figure(figsize=(15,5))
sns.barplot(y='WoE', x="Cutoff", data=woe_creditbureau[woe_creditbureau.Variable=='No_30_DPD_last_6_months'])
plt.show()
```



Percent Stack Plot for No_30_DPD_last_6_months

# Model Evaluation Techniques

For each model obtained, following needs to be carried out:

- Confusion matrix should be prepared for each model.
- Sensitivity, specificity, accuracy curve for each model with different cut-offs.
- AUC-ROC curve for the model using cut-off values for each model.
- Precision and Recall curve for cut-off should be generated.
- Gini-Index needs to be evaluated for Tree based models like decision tree and random forest.
- Within each model type evaluation using GridSerach based on recall values should be done to get models with optimized hyperparameters.
- For evaluation among models, the dataset for rejected applications (with performance tag missing), which were assumed as potentially defaulters should be considered for evaluations. Ideally, the output for all these applications should be defaulters.
- Choose an easy model but not the easiest one.

# Preparing Application Scorecards

- Two Applications Scorecard should be prepared – One based on "Demographic Data Model" and another based on "Combined data for Demographics and Credit Bureau".

- Among these scorecards, Demographic data scorecard can be used for screening purpose (to save cost on getting data from credit Bureau).

- Application scorecard needs be generated based on the mathematical function of the output value (probability) obtained from the model.

- The mathematical function for application scorecard generation may use some fixed coefficient and may take the values obtained for demographics and credit bureau variables.

- Python library for "scorecard" can be used to generate scorecards (as shared in the link https://pypi.org/project/scorecard/).

- More knowledge concerning scorecard generation still awaited from Project Mentor.

- More details would be more relevant when performed practically.

# Thank You