

# Credit EDA & Credit Score Calculation with Python

## Problem statement:

- To conduct a thorough exploratory data analysis (EDA) and deep analysis of a comprehensive dataset containing basic customer details and extensive credit-related information. The aim is to create new, informative features, calculate a hypothetical credit score, and uncover meaningful patterns, anomalies, and insights within the data. ### Suggestions for learners:

## Exploratory Data Analysis (EDA):

- Perform a comprehensive EDA to understand the data's structure, characteristics, distributions, and relationships. Identify and address any missing values, mismatch data types, inconsistencies, or outliers.
- Utilize appropriate visualizations (e.g., histograms, scatter plots, box plots, correlation matrices) to uncover patterns and insights.

## Feature Engineering:

- Create new features that can be leveraged for the calculation of credit scores based on domain knowledge and insights from EDA. Aggregate the data on the customer level if required

## Hypothetical Credit Score Calculation:

- Develop a methodology to calculate a hypothetical credit score using relevant features(use a minimum of 5 maximum of 10 features).
- Clearly outline the developed methodology in the notebook, providing a detailed explanation of the reasoning behind it. (use inspiration from FICO scores and try to use relevant features you created)
- Explore various weighting schemes to assign scores.
- Provide a score for each individual customer

## Analysis and Insights

- Add valuable insights from EDA and credit score calculation Can credit score and aggregated features be calculated at different time frames like the last 3 months/last 6 months (recency based metrics)

## Data Description

Column Name	Description
ID	Represents a unique identification of an entry
Customer_ID	Represents a unique identification of a person
Month	Represents the month of the year
Name	Represents the name of a person
Age	Represents the age of the person
SSN	Represents the social security number of a person
Occupation	Represents the occupation of the person
Annual_Income	Represents the annual income of the person
Monthly_Inhand_Salary	Represents the monthly base salary of a person
Num_Bank_Accounts	Represents the number of bank accounts a person holds
Num_Credit_Card	Represents the number of other credit cards held by a person
Interest_Rate	Represents the interest rate on credit card
Num_of_Loan	Represents the number of loans taken from the bank
Type_of_Loan	Represents the types of loan taken by a person
Delay_from_due_date	Represents the average number of days delayed from the payment date
Num_of_Delayed_Payment	Represents the average number of payments delayed by a person
Changed_Credit_Limit	Represents the percentage change in credit card limit
Num_Credit_Inquiries	Represents the number of credit card inquiries
Credit_Mix	Represents the classification of the mix of credits
Outstanding_Debt	Represents the remaining debt to be paid (in USD)
Credit_Utilization_Ratio	Represents the utilization ratio of credit card
Credit_History_Age	Represents the age of credit history of the person
Payment_of_Min_Amount	Represents whether only the minimum amount was paid by the person
Total_EMI_per_month	Represents the monthly EMI payments (in USD)
Amount_invested_monthly	Represents the monthly amount invested by the customer (in USD)
Payment_Behaviour	Represents the payment behavior of the customer (in USD)
Monthly_Balance	Represents the monthly balance amount of the customer (in USD)

**Collab Link :** [Click Here to explore the collab file \(https://colab.research.google.com/drive/1wZ7mo7g2\\_WG8UyYjCk1usp=sharing\)](https://colab.research.google.com/drive/1wZ7mo7g2_WG8UyYjCk1usp=sharing)



```
In [3]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import scipy
import scipy.stats as stats
import warnings
warnings.filterwarnings("ignore")
```

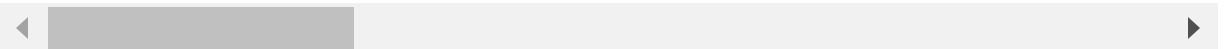
```
In [4]: df = pd.read_csv('Credit_score.csv')
```

```
In [5]: df.head()
```

```
Out[5]:
```

	ID	Customer_ID	Month	Name	Age	SSN	Occupation	Annual_Income	Monthly_Inl
0	0x1602	CUS_0xd40	January	Aaron Maashoh	23	821-00-0265	Scientist	19114.12	
1	0x1603	CUS_0xd40	February	Aaron Maashoh	23	821-00-0265	Scientist	19114.12	
2	0x1604	CUS_0xd40	March	Aaron Maashoh	-500	821-00-0265	Scientist	19114.12	
3	0x1605	CUS_0xd40	April	Aaron Maashoh	23	821-00-0265	Scientist	19114.12	
4	0x1606	CUS_0xd40	May	Aaron Maashoh	23	821-00-0265	Scientist	19114.12	

5 rows × 27 columns



```
In [6]: df.describe()
```

```
Out[6]:
```

	Monthly_Inhand_Salary	Num_Bank_Accounts	Num_Credit_Card	Interest_Rate	Delay_fron
count	84998.000000	100000.000000	100000.000000	100000.000000	100
mean	4194.170850	17.091280	22.47443	72.466040	
std	3183.686167	117.404834	129.05741	466.422621	
min	303.645417	-1.000000	0.00000	1.000000	
25%	1625.568229	3.000000	4.00000	8.000000	
50%	3093.745000	6.000000	5.00000	13.000000	
75%	5957.448333	7.000000	7.00000	20.000000	
max	15204.633330	1798.000000	1499.00000	5797.000000	



```
In [7]: df1 = df.copy()
```

```
In [8]: df1.shape
```

```
Out[8]: (100000, 27)
```

```
In [9]: df1.isna().sum()
```

```
Out[9]: ID                                0
        Customer_ID                       0
        Month                             0
        Name                             9985
        Age                               0
        SSN                               0
        Occupation                        0
        Annual_Income                     0
        Monthly_Inhand_Salary             15002
        Num_Bank_Accounts                 0
        Num_Credit_Card                   0
        Interest_Rate                     0
        Num_of_Loan                       0
        Type_of_Loan                      11408
        Delay_from_due_date               0
        Num_of_Delayed_Payment            7002
        Changed_Credit_Limit              0
        Num_Credit_Inquiries              1965
        Credit_Mix                        0
        Outstanding_Debt                  0
        Credit_Utilization_Ratio          0
        Credit_History_Age                9030
        Payment_of_Min_Amount             0
        Total_EMI_per_month               0
        Amount_invested_monthly           4479
        Payment_Behaviour                 0
        Monthly_Balance                   1200
        dtype: int64
```

## Treating outliers and missing columns

### 1) Name Column :

- Handled missing data by forward-filling within customer segments.

```
In [10]: df1['Name'] = df1.groupby('Customer_ID')['Name'].fillna(method='ffill')
```

### 2) Age

- Eliminated redundant symbols within the column data.
- Addressed missing data points through forward and backward imputation for each customer group.
- Imposed an upper age limit of 95.

```
In [11]: df1['Age'].nunique()
```

```
Out[11]: 1788
```

```
In [12]: df1['Age']=df1['Age'].str.replace('-', '')
df1['Age']=df1['Age'].str.replace('_', '')
df1['Age'] = df1['Age'].astype('int')

df1['Age'].dtypes
```

```
Out[12]: dtype('int64')
```

```
In [13]: df1['Age'].unique()
```

```
Out[13]: array([ 23, 500, 28, ..., 4808, 2263, 1342])
```

```
In [14]: def Age(group):
    mode_age = group['Age'].mode()[0]
    group['Age'] = group['Age'].apply(lambda y: mode_age if y >= 95 else y)
    return group

# Apply the function to each group
df2 = df1.groupby('Customer_ID').apply(Age)
df2 = df2.reset_index(drop=True)

df1 = df2
df1['Age'].nunique()
```

```
Out[14]: 43
```

## SSN Number

- Eliminated redundant symbols within the column data.
- Addressed missing data points through forward and backward imputation for each customer group.

```
In [15]: # df1.groupby('Customer_ID')['SSN'].apply(lambda y : None if y == '#F%D@*&8' else y['SSN'] )
df1['SSN'] = df1['SSN'].replace("#F%D@*&8", np.nan)
df1['SSN'].isna().sum()
def SSN(group):
    mode_ssn = group['SSN'].mode()[0] # Calculate the mode of the 'SSN' column
    group['SSN'] = group['SSN'].apply(lambda y: mode_ssn if pd.isna(y) else y)
    return group

# Apply the function to each group
df2 = df1.groupby('Customer_ID').apply(SSN)
df2 = df2.reset_index(drop=True)
df1 = df2
```

## Occupation

- Eliminated redundant symbols within the column data.
- Addressed missing data points through forward and backward imputation for each customer group.

```
In [16]: df1['Occupation'].replace('_____', "").unique()
```

```
Out[16]: array(['Lawyer', 'Mechanic', '', 'Media_Manager', 'Doctor', 'Journalist',
               'Accountant', 'Manager', 'Entrepreneur', 'Scientist', 'Architect',
               'Teacher', 'Engineer', 'Writer', 'Developer', 'Musician'],
              dtype=object)
```

```
In [17]: df1['Occupation'] = df1['Occupation'].replace('_____', "")
df1['Occupation'] = df1['Occupation'].replace("", np.nan)
df1['Occupation'] = df1.groupby('Customer_ID')['Occupation'].fillna(method='ffill')
df1['Occupation'] = df1.groupby('Customer_ID')['Occupation'].fillna(method='bfill')
```

## Treating numeric columns

- Eliminated redundant symbols within the column data.
- Addressed missing data points through forward and backward imputation for each customer group.
- Performed mathematical corrections for necessary columns
- Imposed upper limit to necessary columns.

```
In [18]: pd.set_option('display.float_format', lambda x: '%.6f' % x)
df1['Annual_Income'] = df1['Annual_Income'].str.replace('_', "").astype('float')
```

```
In [19]: df1['Monthly_Inhand_Salary'] = df1.groupby('Customer_ID')['Monthly_Inhand_Salary'].fillna(method='ffill')
df1['Monthly_Inhand_Salary'] = df1.groupby('Customer_ID')['Monthly_Inhand_Salary'].fillna(method='bfill')
```

```
In [20]: columns = ['Num_Bank_Accounts', 'Num_Credit_Card']
for column in columns:

    df1[column] = df1[column].apply(lambda x: np.nan if x>10 else x)
    df1[column] = df1.groupby('Customer_ID')[column].fillna(method='ffill')

    df1[column] = df1.groupby('Customer_ID')[column].fillna(method='bfill')
    df1[column] = df1[column].astype('int')
```

```
In [21]: df1['Num_Bank_Accounts'].unique()
```

```
Out[21]: array([ 6,  1,  3,  7,  2,  5,  8,  4,  0, 10,  9, -1])
```

```
In [22]: df1['Num_Bank_Accounts'] = df1['Num_Bank_Accounts'].apply(lambda x: -1*x if x== -1 else x)
df1['Num_Credit_Card'] = df1['Num_Credit_Card'].apply(lambda x: -1*x if x== -1 else x)
```

```
In [23]: df1['Interest_Rate'].unique()
```

```
Out[23]: array([ 27, 17, 1, ..., 4892, 4378, 3808])
```

```
In [24]: df1['Interest_Rate'] = df1['Interest_Rate'].apply(lambda x : np.nan if x>45 else x)
df1['Interest_Rate'] = df1.groupby('Customer_ID')['Interest_Rate'].fillna(method='ffill')
df1['Interest_Rate'] = df1.groupby('Customer_ID')['Interest_Rate'].fillna(method='bfill')
df1['Interest_Rate'].unique()
```

```
Out[24]: array([27., 17., 1., 6., 16., 23., 9., 11., 2., 10., 30., 26., 5.,
18., 14., 4., 24., 8., 15., 21., 7., 19., 31., 33., 34., 13.,
20., 28., 32., 29., 12., 25., 3., 22.])
```

```
In [25]: df1['Num_of_Loan'] = df1['Num_of_Loan'].str.replace('_', '')
df1['Num_of_Loan'] = df1['Num_of_Loan'].astype('int')
df1['Num_of_Loan'] = df1['Num_of_Loan'].apply(lambda x: np.nan if x>=23 or x== -100 else x)
df1['Num_of_Loan'] = df1.groupby('Customer_ID')['Num_of_Loan'].fillna(method='ffill')
df1['Num_of_Loan'] = df1.groupby('Customer_ID')['Num_of_Loan'].fillna(method='bfill')
df1['Num_of_Loan'].unique()
```

```
Out[25]: array([ 2., 4., 0., 3., 8., 1., 9., 7., 6., 5., 17., 18., 19.])
```

```
In [26]: df1['Num_of_Loan'] = df1['Num_of_Loan'].astype('int')
```

```
In [27]: df1['Type_of_Loan'] = df1['Type_of_Loan'].fillna('Not Specified')
```

```
In [28]: df1['Delay_from_due_date'].unique()
```

```
Out[28]: array([62, 64, 67, 57, 8, 10, 3, 5, 14, 19, 9, 27, 29, 12, 16, 6, 24,
0, -5, -4, 1, 15, 28, 23, 18, 13, 11, 25, 48, 47, 50, 46, 7, 2,
-3, 4, 30, 17, 21, 20, 22, 35, 40, 31, 26, 59, 63, 58, 37, 42, 43,
38, 55, 41, 36, 52, 53, 54, 49, -2, 44, 39, 61, 34, 33, -1, 45, 60,
51, 66, 56, 32, 65])
```

```
In [29]: df1['Delay_from_due_date'] = df1['Delay_from_due_date'].apply(lambda x: -1*x if x<0 else x)
df1['Delay_from_due_date'].unique()
```

```
Out[29]: array([62, 64, 67, 57, 8, 10, 3, 5, 14, 19, 9, 27, 29, 12, 16, 6, 24,
0, 4, 1, 15, 28, 23, 18, 13, 11, 25, 48, 47, 50, 46, 7, 2, 30,
17, 21, 20, 22, 35, 40, 31, 26, 59, 63, 58, 37, 42, 43, 38, 55, 41,
36, 52, 53, 54, 49, 44, 39, 61, 34, 33, 45, 60, 51, 66, 56, 32, 65])
```





```
In [30]: df1['Num_of_Delayed_Payment'].unique()
```

```

Out[30]: array(['25', '23', '28', '26', '1749', '16', '18', '19', '7', '9', '8',
               '12', nan, '17_', '15', '13', '10', '22', '20', '2', '1', '5',
               '11', '17', '15_', '14', '4', '3', '6', '21', '8_', '11_', '2230',
               '0', '24', '18_', '-2', '19_', '1636', '20_', '-1', '16_', '921',
               '9_', '1766', '21_', '12_', '6_', '1_', '25_', '0_', '-3', '1572',
               '5_', '14_', '3_', '3162', '27', '1034', '4211', '4_', '2712',
               '1832', '22_', '3251', '7_', '867', '13_', '4106', '3951', '2216',
               '24_', '10_', '2_', '1640', '2142_', '754', '974', '1180', '1359',
               '320', '2250', '3621', '2438', '531', '3738', '2566', '719',
               '4326', '223', '1833', '3881', '23_', '439', '1614', '3495', '960',
               '4075', '3119', '4302', '121', '2081', '3894', '3484', '2594',
               '4126', '3944', '2553', '1820', '819', '27_', '3629', '2080',
               '1480', '2801', '359', '94', '473', '2072', '2604', '306', '1633',
               '4262', '2488', '2008', '2955', '1647', '1691', '468', '1150',
               '3491', '4178', '1215', '3793', '3623', '2672', '2508', '1867',
               '4340', '1862', '1282', '1422', '441', '1204', '519', '2938',
               '371', '594', '663_', '46', '3458', '2658', '4134', '2907', '2047',
               '4011', '2991', '4319', '674', '4216', '2671', '-2_', '2323',
               '271', '2184', '2628', '2381', '3191', '2376', '2260', '4005',
               '426', '399', '337', '3069', '3156', '4231', '1750', '372', '2378',
               '876', '2279', '3545', '1222', '3764', '1663', '3200', '1890',
               '2728', '4069', '559', '1598', '3316', '2753', '1687', '281', '84',
               '4047', '1354', '4135', '2533', '2018', '708', '1509', '4360',
               '3726', '1825', '1864', '3112', '1329', '-3_', '733', '1765',
               '775', '3684', '3212', '3478', '2400', '4278', '3636', '871',
               '3946', '3900', '2534', '49', '26_', '197', '1295_', '1841',
               '1478', '4172', '2638', '3972', '1211', '905', '1699', '2324',
               '1325', '1706', '2056', '2903', '2569', '4293', '2621', '2924',
               '1792', '1338', '3107', '430', '714', '2015', '2879', '4024',
               '1673', '415', '2569_', '-1_', '1900', '1852', '2945', '4249',
               '195', '2280', '132', '384', '3148', '642', '3539', '3905', '3171',
               '3050', '1911', '804', '2493', '85', '1463', '3208', '3031',
               '2560', '1795', '1664', '3739', '1481', '3861_', '1172', '1014',
               '1106', '4219', '3751', '3051', '1989', '2149', '1323_', '739',
               '47', '1735', '2255', '1263', '1718', '2566_', '4002', '4295',
               '1402', '3329', '1086', '2873', '4113', '3037', '848_', '813',
               '2413', '2142', '2521', '926', '3707', '210', '2348', '3216',
               '1450', '2021', '2766', '3340', '3447', '1328', '2913', '615',
               '4241', '3313', '1994', '2420', '532', '538', '1411', '2511',
               '3529', '4169', '107', '1191', '2823', '283', '3580', '2354',
               '3765', '1332', '1530', '3926', '3706', '3099', '3790', '1850',
               '2131', '2697', '2239', '162', '2590', '904', '1370', '847',
               '3103', '3661', '1216', '544', '1985', '3502', '4185', '3533',
               '3368', '106', '1301', '853', '3840_', '4191', '523', '3318',
               '2128', '1015', '4022', '4280', '585', '2578', '3819', '972',
               '602', '2060', '2278', '264', '3845', '1502', '3688', '221',
               '1154', '1473', '666', '3920_', '2237_', '1243', '1976', '1192',
               '450', '1552', '1278', '3097_', '851', '3040', '2127', '1685',
               '4096', '4042', '1511', '1523', '3815', '3855', '4161', '133',
               '3750', '252', '2397', '217', '88', '2529', '309', '2286', '273',
               '1079', '2694', '166', '3632', '1443', '1199', '4107', '2875',
               '834', '808', '2429', '3457', '2219_', '577', '3721', '3011',
               '2492', '2729', '4282', '182', '3858', '1743', '2615', '3092',
               '2950', '3536', '3355', '1823', '238', '4077', '2943', '4095',
               '3865', '1861', '3708', '183_', '1184', '846', '709', '4239',
               '2926', '1087_', '2707', '4159', '1371', '3142', '2882', '787',
               '3392', '2793', '3568', '845', '1073', '1975', '3919', '3909',

```

```
'2334', '640', '1541', '2759', '4023', '2751', '1471', '1256',
'2657', '2274', '1096', '3009', '1164', '3155', '2148', '2737',
'86', '3522', '2523', '4281', '3489', '3177', '3154', '3415',
'1606', '1967', '3864', '3300', '1392', '1869', '1177', '3407',
'887', '145', '4144', '4384', '3499', '969', '2854', '1538',
'3559', '3402', '2666', '1004', '2705', '2314', '2138', '3754',
'583', '98', '2044', '1697', '2959', '3722', '933', '4051', '2655',
'1849', '2689', '3222', '2552', '2794', '2006', '829', '1063',
'28_', '2162', '3105', '1045', '1859', '4397', '1337', '3060',
'3467', '683', '2677', '938', '2956', '1389', '1653', '351', '693',
'3243', '1941', '2165', '2070', '4270', '2141', '4019', '3260',
'2461', '3404', '2007', '2616', '482', '3268', '398', '1571',
'3488', '2617', '2810', '2311', '700', '2756', '1181', '2896',
'4128', '3083', '3078', '416', '2503', '1473_', '2506', '742',
'3229', '3253', '4053', '1553', '1236', '2591', '1732', '707',
'4164', '411', '4292', '3115', '749', '2185', '1946', '3584',
'1953', '3978', '541', '3827', '809', '142', '2276', '2317',
'3749', '2587', '2636', '3416', '3370', '3766', '2278_', '4311',
'1489', '130', '294', '827', '3796', '1801', '1218', '4059',
'2768', '4266', '1579', '1952', '2457', '3179', '290', '2589',
'1608', '2196', '2820', '2418', '3245', '2076', '2573', '1133',
'2812', '2498', '1668', '2777', '3870', '186', '2860', '2609',
'3955', '2300', '2570', '508', '793', '1954', '211', '80', '1775',
'676', '1049', '2384', '102', '1891', '4344', '1061', '1879',
'3574', '662', '529', '3043', '2834', '3104', '1060', '929',
'2297', '659', '2262', '3878', '4324', '3336', '4388', '2450',
'3511', '3763', '4251', '192', '3960', '4043', '1996', '1178',
'2660', '3776', '3660', '1874', '1534', '3175', '2645', '4139',
'996', '2351', '2352', '2001', '3880', '1018', '758_', '4337',
'3869', '823', '2544', '2585', '497', '3274', '3456', '2385',
'196', '923', '2431', '3010', '2243', '1884', '778', '175', '2167',
'2222', '1531', '72', '265', '2954', '800', '3847', '779', '4037',
'3391', '4298', '2919', '3492', '52', '1498', '328', '1536',
'2204', '1087'], dtype=object)
```

```
In [31]: df1['Num_of_Delayed_Payment'] = df1['Num_of_Delayed_Payment'].str.replace
('_', '')
df1['Num_of_Delayed_Payment'] = df1['Num_of_Delayed_Payment'].str.replace('-',
',')
df1['Num_of_Delayed_Payment'] = df1['Num_of_Delayed_Payment'].astype('float64')
df1['Num_of_Delayed_Payment'] = df1['Num_of_Delayed_Payment'].apply(lambda x:
np.nan if x>50 else x)
```

```
In [32]: df1['Num_of_Delayed_Payment'] = df1.groupby('Customer_ID')['Num_of_Delayed_Pay
ment'].fillna(method='ffill')
df1['Num_of_Delayed_Payment'] = df1.groupby('Customer_ID')['Num_of_Delayed_Pay
ment'].fillna(method='bfill')
# pd.set_option('display.float_format', lambda x: '%.3f' % x)
df1['Num_of_Delayed_Payment'] = df1['Num_of_Delayed_Payment'].astype('int')
df1['Num_of_Delayed_Payment'].unique()
```

```
Out[32]: array([25, 23, 28, 26, 16, 18, 19, 7, 9, 8, 12, 17, 15, 13, 10, 22, 20,
2, 1, 5, 11, 14, 4, 3, 6, 21, 0, 24, 27, 46, 49, 47])
```

```
In [33]: df1['Changed_Credit_Limit'] = df1['Changed_Credit_Limit'].replace('_',np.nan)
df1['Changed_Credit_Limit'] = df1.groupby('Customer_ID')['Changed_Credit_Limit'].fillna(method='ffill')
df1['Changed_Credit_Limit'] = df1.groupby('Customer_ID')['Changed_Credit_Limit'].fillna(method='bfill')
```

```
In [34]: df1['Num_Credit_Inquiries'] = df1['Num_Credit_Inquiries'].replace('.', '')
df1['Num_Credit_Inquiries'] = df1['Num_Credit_Inquiries'].astype('float')
```

```
In [35]: df1['Num_Credit_Inquiries'] = df1['Num_Credit_Inquiries'].apply(lambda x: np.nan if x>15 else x)
df1['Num_Credit_Inquiries'] = df1.groupby('Customer_ID')['Num_Credit_Inquiries'].fillna(method='ffill')
df1['Num_Credit_Inquiries'] = df1.groupby('Customer_ID')['Num_Credit_Inquiries'].fillna(method='bfill')
df1['Num_Credit_Inquiries'] = df1['Num_Credit_Inquiries'].astype('int')
```

```
In [36]: df1['Credit_Mix'] = df1['Credit_Mix'].replace('_',np.nan)

df1['Credit_Mix'] = df1.groupby('Customer_ID')['Credit_Mix'].fillna(method='ffill')
df1['Credit_Mix'] = df1.groupby('Customer_ID')['Credit_Mix'].fillna(method='bfill')
```

```
In [37]: df1['Outstanding_Debt'] = df1['Outstanding_Debt'].str.replace('_', '')
df1['Outstanding_Debt'] = df1['Outstanding_Debt'].astype('float')
```

```
In [38]: def convert_to_months(age_str):
    if pd.isna(age_str) or age_str == 'NA':
        return np.nan
    else:
        years, months = age_str.split(' and ')
        years = int(years.split()[0])
        months = int(months.split()[0])
        return years * 12 + months
    # Apply conversion
df1['Credit_months'] = df1['Credit_History_Age'].apply(convert_to_months)
```

```
In [39]: df1['Credit_months'] = df1.groupby('Customer_ID')['Credit_months'].fillna(method='ffill')
df1['Credit_months'] = df1.groupby('Customer_ID')['Credit_months'].fillna(method='bfill')
```

```
In [40]: df1['Credit_months'] = df1['Credit_months'].astype('int')
```

```
In [41]: df1['Credit_months'].isna().sum()
```

```
Out[41]: 0
```

```
In [42]: def convert_to_years(age_str):
        if pd.isna(age_str) or age_str == 'NA':
            return np.nan
        else:
            years, months = age_str.split(' and ')
            years = int(years.split()[0])

            return years
        # Apply conversion
df1['Credit_years'] = df1['Credit_History_Age'].apply(convert_to_years)
```

```
In [43]: df1['Credit_years'] = df1.groupby('Customer_ID')['Credit_years'].fillna(method='ffill')
df1['Credit_years'] = df1.groupby('Customer_ID')['Credit_years'].fillna(method='bfill')
```

```
In [44]: df1['Amount_invested_monthly'] = df1['Amount_invested_monthly'].str.replace('_', '')
df1['Amount_invested_monthly'] = df1.groupby('Customer_ID')['Amount_invested_monthly'].fillna(method='ffill')
df1['Amount_invested_monthly'] = df1.groupby('Customer_ID')['Amount_invested_monthly'].fillna(method='bfill')
df1['Amount_invested_monthly'] = df1['Amount_invested_monthly'].astype('float')
```

```
In [45]: df1['Payment_Behaviour'] = df1['Payment_Behaviour'].replace('!@9#%8', np.nan)
df1['Payment_Behaviour'] = df1.groupby('Customer_ID')['Payment_Behaviour'].fillna(method='ffill')
df1['Payment_Behaviour'] = df1.groupby('Customer_ID')['Payment_Behaviour'].fillna(method='bfill')
```

```
In [46]: df1['Payment_Behaviour'] = df1.groupby('Customer_ID')['Payment_Behaviour'].fillna(method='ffill')
df1['Payment_Behaviour'] = df1.groupby('Customer_ID')['Payment_Behaviour'].fillna(method='bfill')
```

```
In [47]: df1['Monthly_Balance'].dtypes
```

```
Out[47]: dtype('O')
```

```
In [48]: df1['Monthly_Balance'] = df1.groupby('Customer_ID')['Monthly_Balance'].fillna(method='ffill')
df1['Monthly_Balance'] = df1.groupby('Customer_ID')['Monthly_Balance'].fillna(method='bfill')
df1['Monthly_Balance'] = df1['Monthly_Balance'].replace('__-33333333333333333333333333__ ', np.nan)
df1['Monthly_Balance'] = df1['Monthly_Balance'].astype('float')
```

In [49]: df1.info()

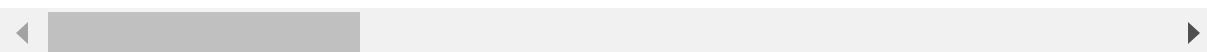
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 100000 entries, 0 to 99999
Data columns (total 29 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   ID                                    100000 non-null  object
1   Customer_ID                          100000 non-null  object
2   Month                                100000 non-null  object
3   Name                                  98630 non-null   object
4   Age                                   100000 non-null  int64
5   SSN                                   100000 non-null  object
6   Occupation                           100000 non-null  object
7   Annual_Income                        100000 non-null  float64
8   Monthly_Inhand_Salary                100000 non-null  float64
9   Num_Bank_Accounts                    100000 non-null  int64
10  Num_Credit_Card                       100000 non-null  int64
11  Interest_Rate                         100000 non-null  float64
12  Num_of_Loan                           100000 non-null  int64
13  Type_of_Loan                          100000 non-null  object
14  Delay_from_due_date                   100000 non-null  int64
15  Num_of_Delayed_Payment                 100000 non-null  int64
16  Changed_Credit_Limit                   100000 non-null  object
17  Num_Credit_Inquiries                   100000 non-null  int64
18  Credit_Mix                             100000 non-null  object
19  Outstanding_Debt                       100000 non-null  float64
20  Credit_Utilization_Ratio               100000 non-null  float64
21  Credit_History_Age                     90970 non-null   object
22  Payment_of_Min_Amount                  100000 non-null  object
23  Total_EMI_per_month                    100000 non-null  float64
24  Amount_invested_monthly                100000 non-null  float64
25  Payment_Behaviour                      100000 non-null  object
26  Monthly_Balance                        99991 non-null   float64
27  Credit_months                          100000 non-null  int64
28  Credit_years                           100000 non-null  float64
dtypes: float64(9), int64(8), object(12)
memory usage: 22.1+ MB
```

In [50]: `df1.head()`

Out[50]:

	ID	Customer_ID	Month	Name	Age	SSN	Occupation	Annual_Income	Monthly_Inh
0	0x1628a	CUS_0x1000	January	Alistair Barrf	17	913-74-1218	Lawyer	30625.940000	27
1	0x1628b	CUS_0x1000	February	Alistair Barrf	17	913-74-1218	Lawyer	30625.940000	27
2	0x1628c	CUS_0x1000	March	Alistair Barrf	17	913-74-1218	Lawyer	30625.940000	27
3	0x1628d	CUS_0x1000	April	Alistair Barrf	17	913-74-1218	Lawyer	30625.940000	27
4	0x1628e	CUS_0x1000	May	Alistair Barrf	17	913-74-1218	Lawyer	30625.940000	27

5 rows × 29 columns



In [51]: `df1['Changed_Credit_Limit'] = df1['Changed_Credit_Limit'].astype('float')`

## Downloading the Cleaned data

In [52]: `# df1.to_csv('cleaned_credit_info.csv')`

## Importing the cleaned data

In [53]: `df1 = pd.read_csv('/content/cleaned_credit_info.csv')`

In [54]: `df1.describe()`

Out[54]:

	Unnamed: 0	Age	Annual_Income	Monthly_Inhand_Salary	Num_Bank_Accou
<b>count</b>	100000.000000	100000.000000	100000.000000	100000.000000	100000.000000
<b>mean</b>	49999.500000	33.311180	176415.701298	4198.771619	5.3689
<b>std</b>	28867.657797	10.764783	1429618.051414	3187.494354	2.5927
<b>min</b>	0.000000	14.000000	7005.930000	303.645417	0.0000
<b>25%</b>	24999.750000	24.000000	19457.500000	1626.761667	3.0000
<b>50%</b>	49999.500000	33.000000	37578.610000	3096.378333	5.0000
<b>75%</b>	74999.250000	42.000000	72790.920000	5961.745000	7.0000
<b>max</b>	99999.000000	56.000000	24198062.000000	15204.633330	10.0000

In [55]: `df1.describe(include='object').T`

Out[55]:

	count	unique	top	freq
<b>ID</b>	100000	100000	0x1628a	1
<b>Customer_ID</b>	100000	12500	CUS_0x1000	8
<b>Month</b>	100000	8	January	12500
<b>Name</b>	98630	10139	Langep	48
<b>SSN</b>	100000	12500	913-74-1218	8
<b>Occupation</b>	100000	15	Lawyer	7096
<b>Type_of_Loan</b>	100000	6260	Not Specified	12816
<b>Credit_Mix</b>	100000	3	Standard	45848
<b>Credit_History_Age</b>	90970	404	15 Years and 11 Months	446
<b>Payment_of_Min_Amount</b>	100000	3	Yes	52326
<b>Payment_Behaviour</b>	100000	6	Low_spent_Small_value_payments	27588

- There is a total of 12500 unique customers.
- The data is from January to August.
- Low\_spent\_Small\_value\_payments has high frequency 27588
- Standard has highest frequency from credit mix categories.



```
In [56]: ['Age', 'Annual_Income', 'Monthly_Inhand_Salary',  
         'Num_Bank_Accounts', 'Num_Credit_Card', 'Interest_Rate', 'Num_of_Loan',  
         'Delay_from_due_date', 'Num_of_Delayed_Payment', 'Changed_Credit_Limit',  
         'Num_Credit_Inquiries', 'Outstanding_Debt', 'Credit_Utilization_Ratio',  
         'Total_EMI_per_month', 'Amount_invested_monthly', 'Monthly_Balance',  
         'Credit_months', 'Credit_years']
```

```
Out[56]: ['Age',  
         'Annual_Income',  
         'Monthly_Inhand_Salary',  
         'Num_Bank_Accounts',  
         'Num_Credit_Card',  
         'Interest_Rate',  
         'Num_of_Loan',  
         'Delay_from_due_date',  
         'Num_of_Delayed_Payment',  
         'Changed_Credit_Limit',  
         'Num_Credit_Inquiries',  
         'Outstanding_Debt',  
         'Credit_Utilization_Ratio',  
         'Total_EMI_per_month',  
         'Amount_invested_monthly',  
         'Monthly_Balance',  
         'Credit_months',  
         'Credit_years']
```

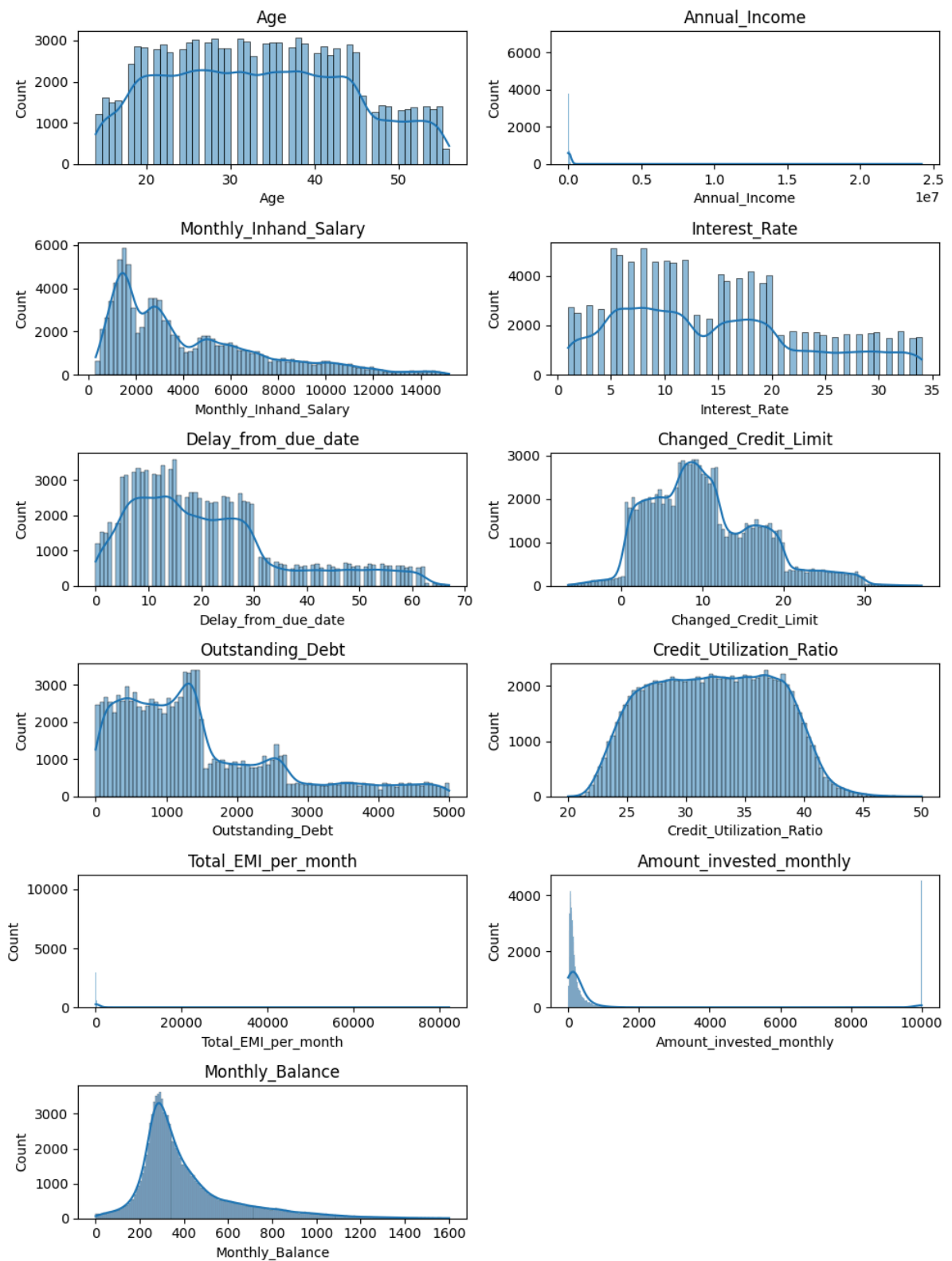
```
In [57]: df1['Customer_ID'].nunique()
```

```
Out[57]: 12500
```

```
In [58]: df1[['Customer_ID', 'Month']][df1['Month']=='August'].count()
```

```
Out[58]: Customer_ID    12500  
         Month          12500  
         dtype: int64
```

```
In [59]: numerals_feature = ['Age', 'Annual_Income', 'Monthly_Inhand_Salary',  
                             'Interest_Rate',  
                             'Delay_from_due_date', 'Changed_Credit_Limit',  
                             'Outstanding_Debt', 'Credit_Utilization_Ratio',  
                             'Total_EMI_per_month', 'Amount_invested_monthly', 'Monthly_Balance',  
                             ]  
plt.figure(figsize=(10,15))  
  
for i in range(1, len(numerals_feature) + 1):  
    plt.subplot(7, 2, i)  
    sns.histplot(df1[numerals_feature[i-1]], kde=True)  
    plt.title(numerals_feature[i-1])  
  
plt.tight_layout()  
plt.show()
```



### Insights:

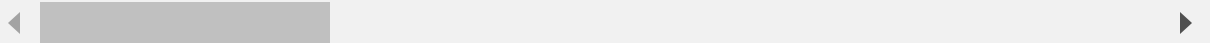
- The distribution of age is relatively uniform, with a slight peak around the mid-20s to early 30s.
- There is a noticeable drop in the number of individuals above the age of 50.
- fewer individuals earning very high monthly inhand salaries.
- The interest rate distribution is fairly uniform, with small peaks around 10-15% and another around 20-25%.
- The delay from the due date has a broad distribution with many individuals experiencing delays of up to 30 days.
- Most changes in credit limits are moderate, with few instances of very high adjustments.
- The distribution shows a high frequency of individuals with outstanding debt in the range of 0 to 2000 units.
- The majority of individuals maintain a credit utilization ratio within the 20% to 40% range, indicating good credit management practices.
- Monthly EMI payments are generally low, suggesting that most individuals do not have heavy debt burdens from EMIs.
- Most individuals invest small amounts monthly, indicating cautious or limited investment behavior.

```
In [60]: num = df1[['Age', 'Annual_Income', 'Monthly_Inhand_Salary',  
                  'Num_Bank_Accounts', 'Num_Credit_Card', 'Interest_Rate', 'Num_of_Loan',  
                  'Delay_from_due_date', 'Num_of_Delayed_Payment', 'Changed_Credit_Limit',  
                  'Num_Credit_Inquiries', 'Outstanding_Debt', 'Credit_Utilization_Ratio',  
                  'Total_EMI_per_month', 'Amount_invested_monthly', 'Monthly_Balance',  
                  'Credit_months', 'Credit_years']]  
  
data = num.corr()
```

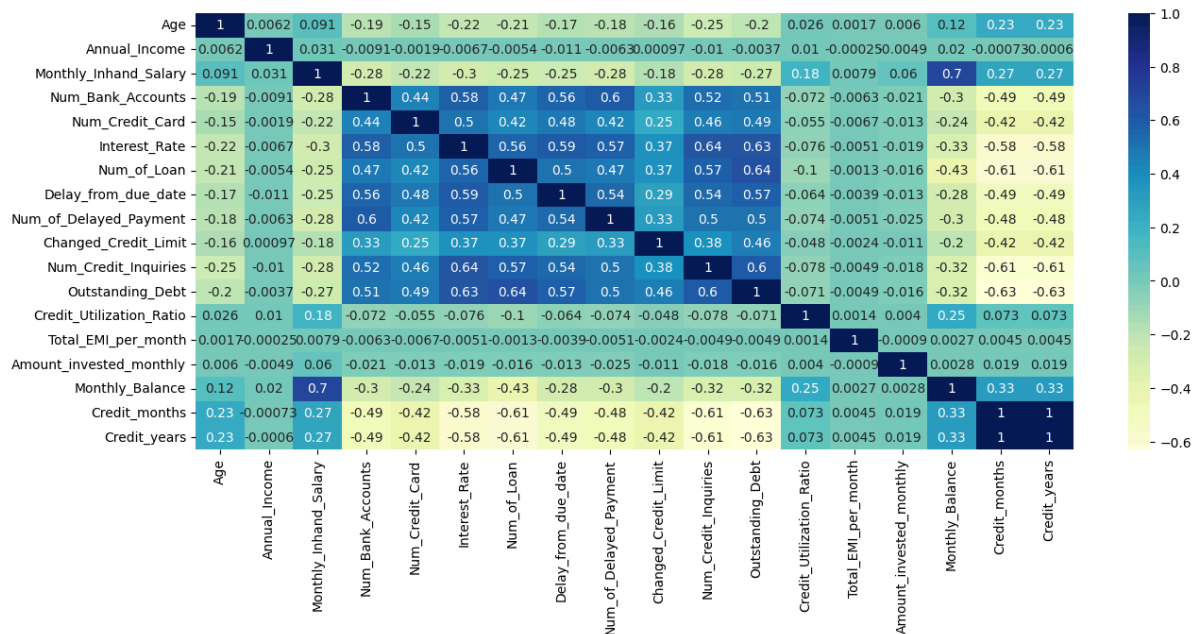
In [61]: data

Out[61]:

	Age	Annual_Income	Monthly_Inhand_Salary	Num_Bank_Account
<b>Age</b>	1.000000	0.006235	0.090623	-0.19033
<b>Annual_Income</b>	0.006235	1.000000	0.030508	-0.00908
<b>Monthly_Inhand_Salary</b>	0.090623	0.030508	1.000000	-0.28324
<b>Num_Bank_Accounts</b>	-0.190335	-0.009087	-0.283243	1.00000
<b>Num_Credit_Card</b>	-0.148567	-0.001941	-0.216958	0.44268
<b>Interest_Rate</b>	-0.217531	-0.006702	-0.301906	0.58429
<b>Num_of_Loan</b>	-0.213355	-0.005440	-0.254155	0.47215
<b>Delay_from_due_date</b>	-0.173892	-0.010670	-0.249300	0.55974
<b>Num_of_Delayed_Payment</b>	-0.183664	-0.006277	-0.284396	0.60035
<b>Changed_Credit_Limit</b>	-0.156400	0.000974	-0.175135	0.33095
<b>Num_Credit_Inquiries</b>	-0.252376	-0.010166	-0.281860	0.52133
<b>Outstanding_Debt</b>	-0.202361	-0.003706	-0.269078	0.50704
<b>Credit_Utilization_Ratio</b>	0.025523	0.010316	0.176081	-0.07176
<b>Total_EMI_per_month</b>	0.001698	-0.000248	0.007949	-0.00634
<b>Amount_invested_monthly</b>	0.005982	-0.004870	0.060218	-0.02060
<b>Monthly_Balance</b>	0.117422	0.019593	0.702747	-0.29980
<b>Credit_months</b>	0.234635	-0.000728	0.271516	-0.48537
<b>Credit_years</b>	0.234569	-0.000602	0.271353	-0.48513



```
In [62]: plt.figure(figsize=(15,6))
sns.heatmap(data, cmap="YlGnBu", annot=True)
plt.show()
```



## Positive Correlations

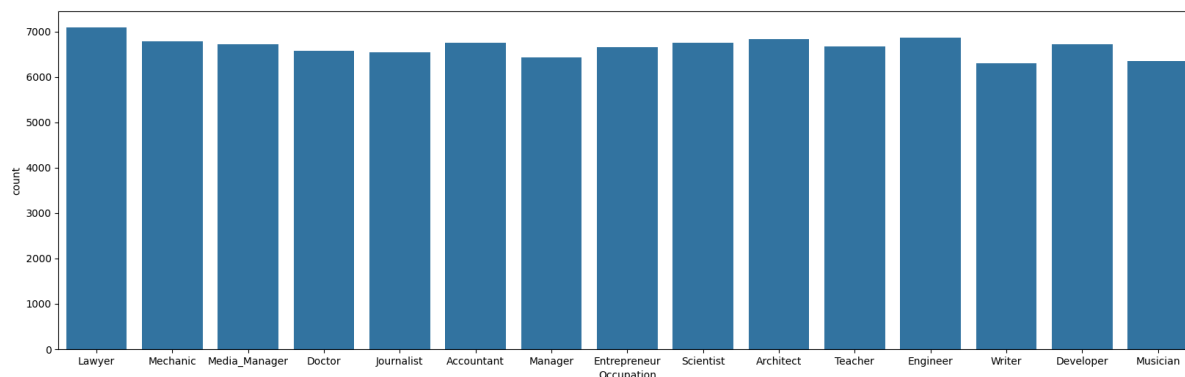
- Higher monthly in-hand salary is strongly associated with a higher monthly balance.
- More bank accounts and credit cards tend to be associated with more loans.
- Higher interest rates are associated with higher outstanding debt.
- A higher credit utilization ratio is somewhat associated with higher outstanding debt.

## Negative Correlations

- Delays in payment are associated with lower monthly balances.
- A higher credit utilization ratio has a small positive relationship with monthly balance, which might indicate that those with higher utilization ratios maintain a slightly higher balance
- More delayed payments are linked to more loans and higher outstanding debt.

```
In [63]: plt.figure(figsize=(20,6))  
sns.countplot(x='Occupation',data=df1)
```

```
Out[63]: <Axes: xlabel='Occupation', ylabel='count'>
```

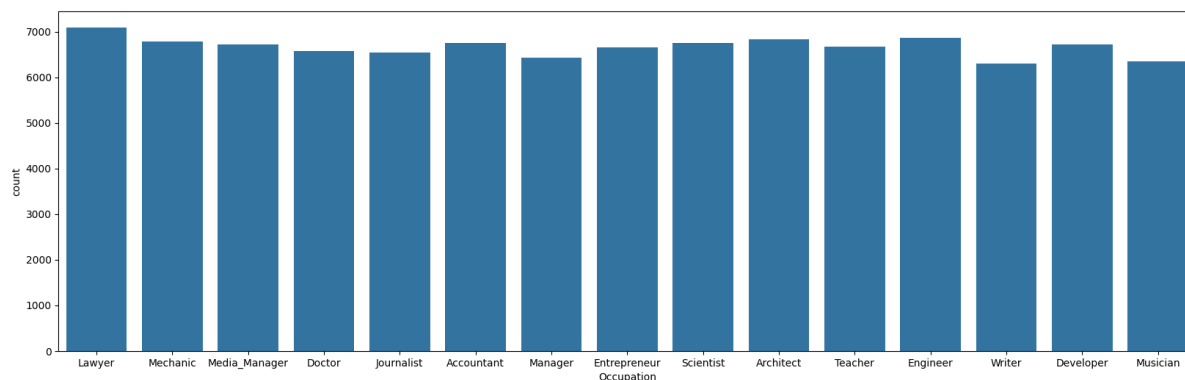


## Insights

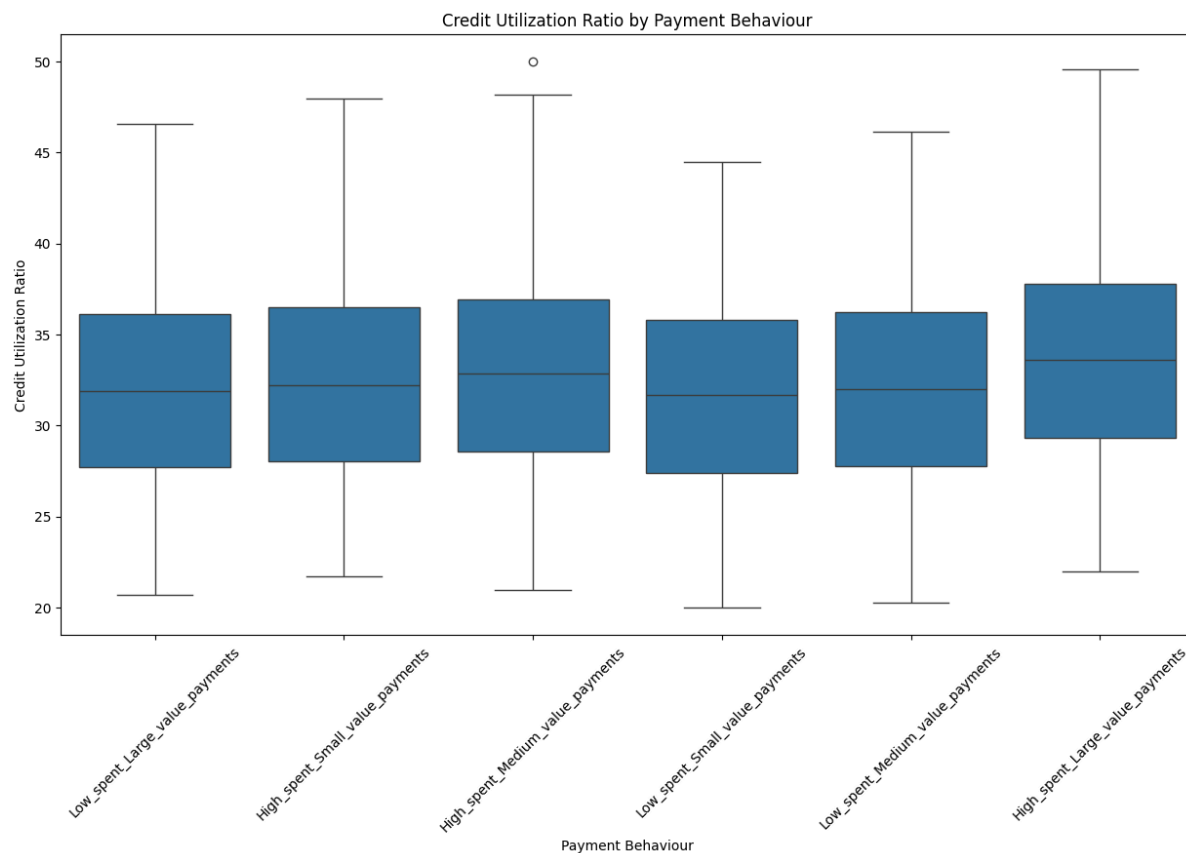
- The distribution is almost uniform across all the categories of occupation

```
In [64]: plt.figure(figsize=(20,6))  
sns.countplot(x='Occupation',data=df1)
```

```
Out[64]: <Axes: xlabel='Occupation', ylabel='count'>
```



```
In [65]: plt.figure(figsize=(15, 8))
sns.boxplot(x='Payment_Behaviour', y='Credit_Utilization_Ratio', data=df1)
plt.title('Credit Utilization Ratio by Payment Behaviour')
plt.xlabel('Payment Behaviour')
plt.ylabel('Credit Utilization Ratio')
plt.xticks(rotation=45)
plt.show()
```

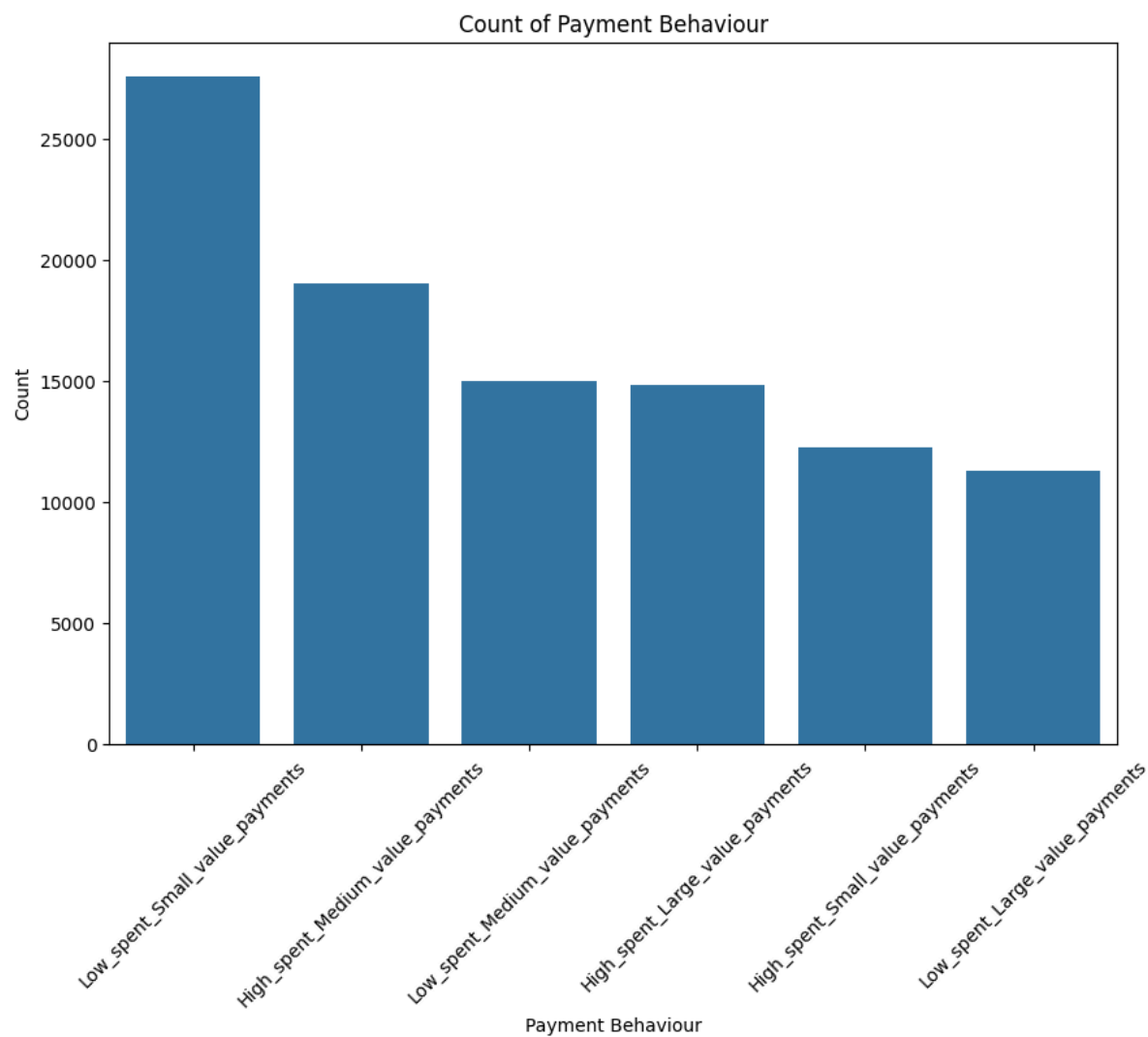


## Insights

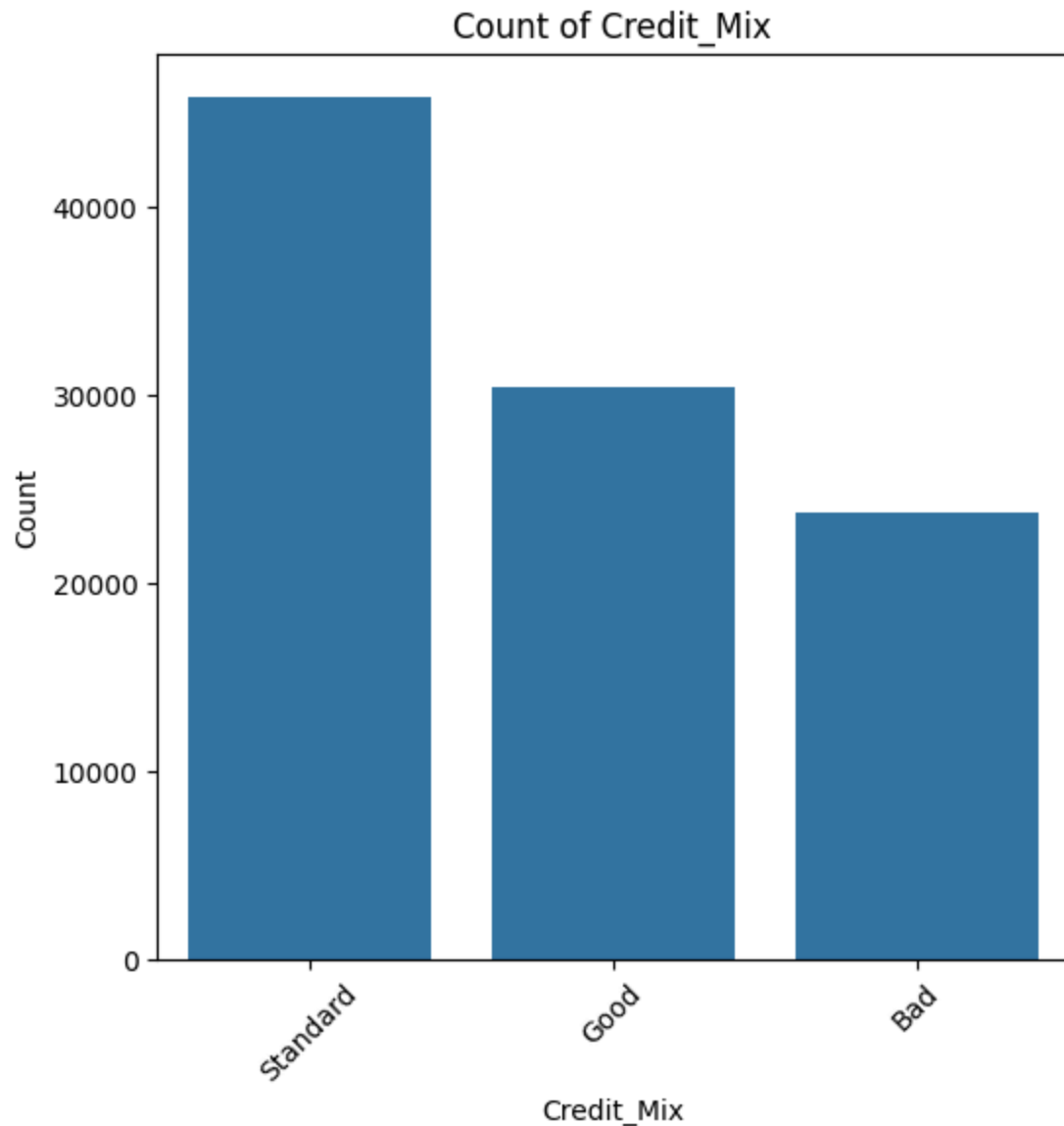
- All payment behavior categories exhibit similar median Credit Utilization Ratios, which appear to be around 30-35.



```
In [66]: plt.figure(figsize=(10, 7))
sns.countplot(x='Payment_Behaviour', data=df1, order = df1['Payment_Behaviour'].value_counts().index)
plt.title('Count of Payment Behaviour')
plt.xlabel('Payment Behaviour')
plt.ylabel('Count')
plt.xticks(rotation=45)
plt.show()
```



```
In [67]: plt.figure(figsize=(6, 6))
sns.countplot(x='Credit_Mix', data=df1, order = df1['Credit_Mix'].value_counts().index)
plt.title('Count of Credit_Mix')
plt.xlabel('Credit_Mix')
plt.ylabel('Count')
plt.xticks(rotation=45)
plt.show()
```



- Low spent smaller values payments and standard credit mix have max frequencies

```
In [68]: df1.groupby(['Credit_Mix'])['Num_Credit_Card'].value_counts().reset_index()
```

Out[68]:

	Credit_Mix	Num_Credit_Card	count
0	Bad	7	4057
1	Bad	8	4017
2	Bad	5	3997
3	Bad	6	3985
4	Bad	10	3864
5	Bad	9	3819
6	Bad	4	29
7	Good	5	6144
8	Good	4	6071
9	Good	3	5668
10	Good	7	4146
11	Good	6	3962
12	Good	1	2185
13	Good	2	2151
14	Good	8	43
15	Good	0	14
16	Standard	6	8982
17	Standard	7	8826
18	Standard	5	8763
19	Standard	4	8258
20	Standard	3	7897
21	Standard	10	1134
22	Standard	8	1012
23	Standard	9	934
24	Standard	2	42

```
In [69]: df1['Credit_Utilization_Category'] = pd.cut(df1['Credit_Utilization_Ratio'], bins=[0, 10, 30, 50, 70, 100], labels=['0-10%', '10-30%', '30-50%', '50-70%', '70-100%'])
credit_mix_vs_credit_utilization = pd.crosstab(df1['Credit_Mix'], df1['Credit_Utilization_Category'], normalize=True)
credit_mix_vs_credit_utilization*100
```

```
Out[69]:
```

	Credit_Utilization_Category	10-30%	30-50%
	Credit_Mix		
	Bad	9.459000	14.309000
	Good	9.958000	20.426000
	Standard	16.945000	28.903000

- **Credit Utilization (10-30%):**

- Bad credit mix: 9.459
- Good credit mix: 9.958
- Standard credit mix: 16.945

- **Credit Utilization (30-50%):**

- Bad credit mix: 14.309
- Good credit mix: 20.426
- Standard credit mix: 28.903

```
In [70]: df1['Payment_Behaviour'].unique()
```

```
Out[70]: array(['Low_spent_Large_value_payments',
                'High_spent_Small_value_payments',
                'High_spent_Medium_value_payments',
                'Low_spent_Small_value_payments',
                'Low_spent_Medium_value_payments',
                'High_spent_Large_value_payments'], dtype=object)
```

## Calculation Of Credit Score:

-

- 1. Payment History (35%):
  - Payment\_of\_Min\_Amount
  - Num\_of\_Delayed\_Payment
  - Delay\_from\_due\_date
- 1. Accounts Owed (30%):
  - Outstanding\_Debt
  - Credit\_Utilization\_Ratio
  - Total\_EMI\_per\_month
- 1. Length of Credit History (15%)
  - Credit\_History\_Age\_Months
- 1. Credit Mix (10%)
  - Credit\_Mix
  - Num\_Bank\_Accounts
  - Num\_Credit\_Card
  - Num\_of\_Loan
- 1. New Credit (10%) -Num\_Credit\_Inquiries

## Assigning Scaling points to categories

```
In [71]: # Create a new column for Credit History Age in months
df1['Credit_History_Age_Months'] = (df1['Credit_years'] * 12) + df1['Credit_months']

# Map Payment_Behaviour to numerical scores
payment_behaviour_map = {
    'Low_spent_Large_value_payments': 1,
    'High_spent_Small_value_payments': 2,
    'High_spent_Medium_value_payments': 3,
    'Low_spent_Small_value_payments': 4,
    'Low_spent_Medium_value_payments': 5,
    'High_spent_Large_value_payments': 6
}
df1['Payment_Behaviour_Score'] = df1['Payment_Behaviour'].map(payment_behaviour_map)

# Define the scoring function for Credit Utilization Ratio
def credit_utilization_score(ratio):
    if ratio <= 0.10:
        return 1.0 # Highest points
    elif ratio <= 0.30:
        return 0.5 # Medium points
    else:
        return 0.0 # Lowest points
```

## Aggregating Data on customer Level

```
In [72]: # Aggregating the data for each customer
df_aggregated = df1.groupby('Customer_ID').agg({
    'Payment_of_Min_Amount': 'last',
    'Num_of_Delayed_Payment': 'sum',
    'Delay_from_due_date': 'sum',
    'Outstanding_Debt': 'mean',
    'Credit_Utilization_Ratio': 'mean',
    'Total_EMI_per_month': 'mean',
    'Credit_History_Age_Months': 'last',
    'Num_Credit_Inquiries': 'sum',
    'Num_Bank_Accounts': 'mean',
    'Num_Credit_Card': 'mean',
    'Num_of_Loan': 'mean',
    'Credit_Mix': 'last',
    'Payment_Behaviour_Score': 'last'
}).reset_index()

# Apply the scoring function
df_aggregated['Credit_Utilization_Ratio_Score'] = df_aggregated['Credit_Utiliz
ation_Ratio'].apply(credit_utilization_score)
```

```
In [73]: df_aggregated.head(5)
```

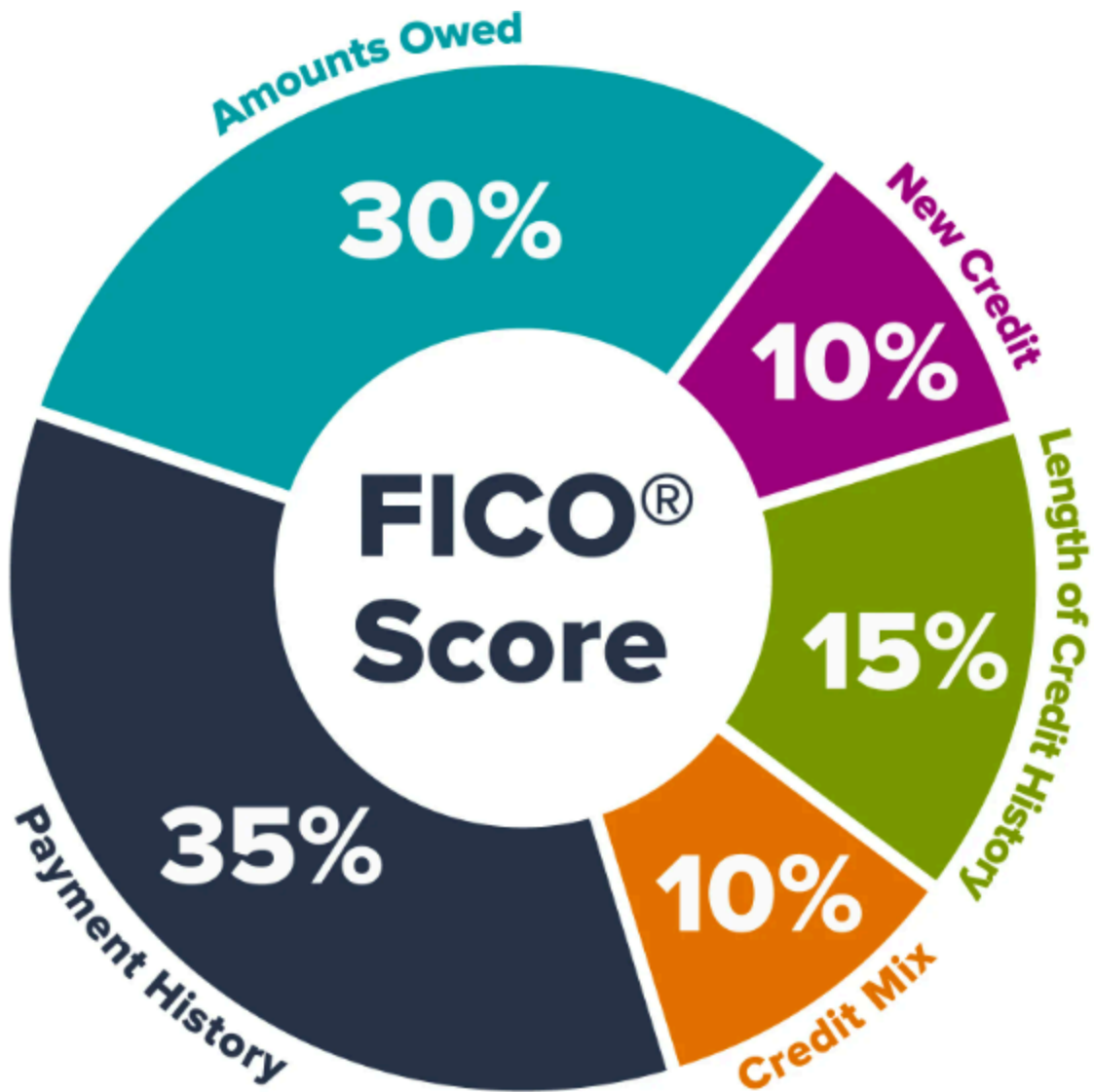
Out[73]:

	Customer_ID	Payment_of_Min_Amount	Num_of_Delayed_Payment	Delay_from_due_date	Outst
0	CUS_0x1000	Yes	200	498	
1	CUS_0x1009	Yes	141	58	
2	CUS_0x100b	No	59	108	
3	CUS_0x1011	Yes	114	218	
4	CUS_0x1013	No	68	100	

## Creating new features by scaling the features

```
In [74]: # Define the normalization functions for other factors
def min_max_scaling(column):
    return (column - column.min()) / (column.max() - column.min())

# Apply min-max scaling to the relevant columns
df_aggregated['Num_of_Delayed_Payment_scaled'] = min_max_scaling(df_aggregated[
    'Num_of_Delayed_Payment'])
df_aggregated['Delay_from_due_date_scaled'] = min_max_scaling(df_aggregated['D
    elay_from_due_date'])
df_aggregated['Outstanding_Debt_scaled'] = min_max_scaling(df_aggregated['Ous
    tanding_Debt'])
df_aggregated['Total_EMI_per_month_scaled'] = min_max_scaling(df_aggregated['T
    otal_EMI_per_month'])
df_aggregated['Credit_History_Age_Months_scaled'] = min_max_scaling(df_aggrega
    ted['Credit_History_Age_Months'])
df_aggregated['Num_Credit_Inquiries_scaled'] = min_max_scaling(df_aggregated
    ['Num_Credit_Inquiries'])
df_aggregated['Num_Bank_Accounts_scaled'] = min_max_scaling(df_aggregated['Num
    _Bank_Accounts'])
df_aggregated['Num_Credit_Card_scaled'] = min_max_scaling(df_aggregated['Num_C
    redit_Card'])
df_aggregated['Num_of_Loan_scaled'] = min_max_scaling(df_aggregated['Num_of_Lo
    an'])
df_aggregated['Payment_Behaviour_Score_scaled'] = min_max_scaling(df_aggregate
    d['Payment_Behaviour_Score'])
```



Applying factor-based normalization to credit score data.



```

In [75]: # Define functions to normalize each factor
def normalize_payment_history(row):
    return (row['Payment_of_Min_Amount'] == 'Yes') * 0.7 + (1 - row['Num_of_De
layed_Payment_scaled']) * 0.2 + (1 - row['Delay_from_due_date_scaled']) * 0.1

def normalize_amounts_owed(row):
    return (1 - row['Outstanding_Debt_scaled']) * 0.5 + (row['Credit_Utilizati
on_Ratio_Score']) * 0.3 + (1 - row['Total_EMI_per_month_scaled']) * 0.2

def normalize_length_of_credit_history(row):
    return row['Credit_History_Age_Months_scaled']

def normalize_new_credit(row):
    return (1 - row['Num_Credit_Inquiries_scaled'])

def normalize_credit_mix(row):
    credit_mix_score = 0
    if row['Credit_Mix'] == 'Good':
        credit_mix_score = 1
    elif row['Credit_Mix'] == 'Standard':
        credit_mix_score = 0.5
    elif row['Credit_Mix'] == 'Bad':
        credit_mix_score = 0
    return (row['Num_Bank_Accounts_scaled'] * 0.25 +
            row['Num_Credit_Card_scaled'] * 0.25 +
            row['Num_of_Loan_scaled'] * 0.25 +
            credit_mix_score * 0.25)

def normalize_payment_behaviour(row):
    return row['Payment_Behaviour_Score_scaled'] * 0.1

# Normalize each factor for all rows
df_aggregated['Payment_History_Score'] = df_aggregated.apply(normalize_payment
_history, axis=1)
df_aggregated['Amounts_Owed_Score'] = df_aggregated.apply(normalize_amounts_ow
ed, axis=1)
df_aggregated['Length_of_Credit_History_Score'] = df_aggregated.apply(normaliz
e_length_of_credit_history, axis=1)
df_aggregated['New_Credit_Score'] = df_aggregated.apply(normalize_new_credit,
axis=1)
df_aggregated['Credit_Mix_Score'] = df_aggregated.apply(normalize_credit_mix,
axis=1)
df_aggregated['Payment_Behaviour_Score'] = df_aggregated.apply(normalize_payme
nt_behaviour, axis=1)

# Calculate the final credit score
df_aggregated['Credit_Score'] = (df_aggregated['Payment_History_Score'] * 0.35
+
                                df_aggregated['Amounts_Owed_Score'] * 0.30 +
                                df_aggregated['Length_of_Credit_History_Scor
e'] * 0.15 +
                                df_aggregated['New_Credit_Score'] * 0.10 +
                                df_aggregated['Credit_Mix_Score'] * 0.10 +
                                df_aggregated['Payment_Behaviour_Score'] * 0.
05) # Adjust weight if needed

```

```
df_aggregated['Credit_Score'] = 300 + (df_aggregated['Credit_Score'] * 550)

df_aggregated.to_csv('credit_scores_aggregated.csv', index=False)
```

## factor wise score and credit score.

```
In [76]: df_aggregated[['Customer_ID', 'Payment_History_Score', 'Amounts_Owed_Score', 'Length_of_Credit_History_Score', 'New_Credit_Score', 'Credit_Mix_Score', 'Credit_Score']]
```

Out[76]:

	Customer_ID	Payment_History_Score	Amounts_Owed_Score	Length_of_Credit_History_Score
0	CUS_0x1000	0.712008	0.543378	0.3042
1	CUS_0x1009	0.854888	0.679023	0.9292
2	CUS_0x100b	0.222732	0.596958	0.4570
3	CUS_0x1011	0.848860	0.651865	0.4570
4	CUS_0x1013	0.215782	0.534822	0.5176
...	...	...	...	...
12495	CUS_0xff3	0.215230	0.576840	0.5088
12496	CUS_0xff4	0.882958	0.623470	0.5467
12497	CUS_0xff6	0.268995	0.665334	0.7310
12498	CUS_0xffc	0.805038	0.568134	0.3851
12499	CUS_0xffd	0.876433	0.528541	0.5467

12500 rows × 7 columns

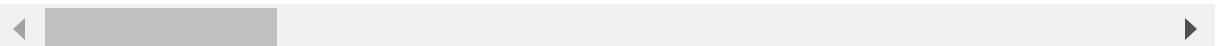


```
In [77]: df_aggregated.head()
```

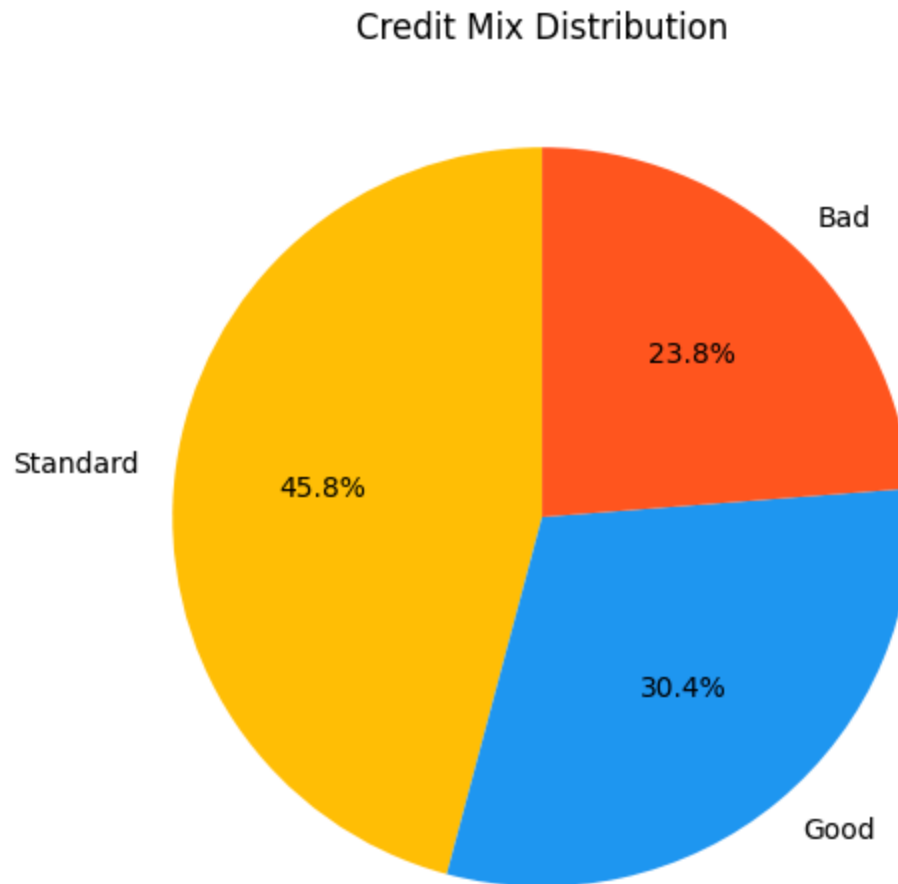
Out[77]:

	Customer_ID	Payment_of_Min_Amount	Num_of_Delayed_Payment	Delay_from_due_date	Outstanding_Amount
0	CUS_0x1000	Yes	200	498	
1	CUS_0x1009	Yes	141	58	
2	CUS_0x100b	No	59	108	
3	CUS_0x1011	Yes	114	218	
4	CUS_0x1013	No	68	100	

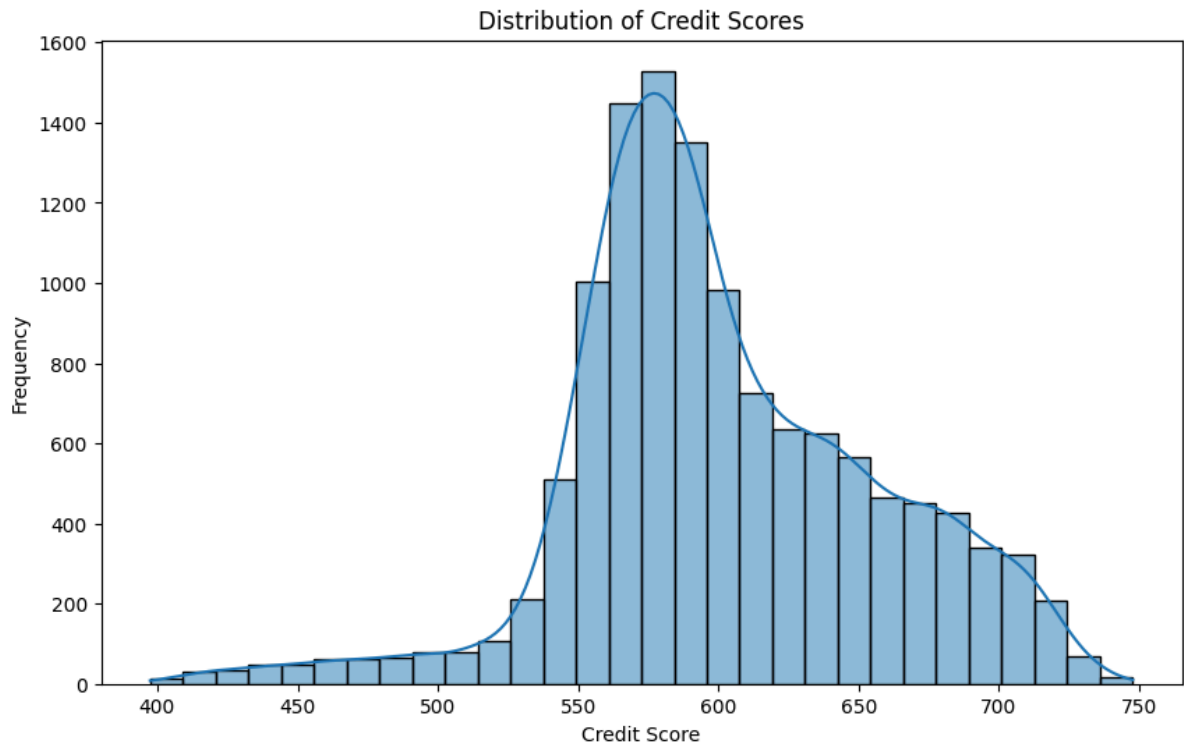
5 rows × 31 columns



```
In [83]: credit_mix_distribution = df_aggregated['Credit_Mix'].value_counts()
plt.figure(figsize=(8, 6))
credit_mix_distribution.plot.pie(autopct='%1.1f%%', startangle=90, colors=['#FFC107', '#2196F3', '#FF5722'])
plt.title('Credit Mix Distribution')
plt.ylabel('')
plt.show()
```



```
In [85]: plt.figure(figsize=(10, 6))
sns.histplot(df_aggregated['Credit_Score'], bins=30, kde=True)
plt.title('Distribution of Credit Scores')
plt.xlabel('Credit Score')
plt.ylabel('Frequency')
plt.show()
```



#### Final Insights:

- The highest frequency of credit scores is around 630-640, indicating that a significant portion of individuals fall within this range.
- The credit scores range from approximately 400 to 750, providing a broad perspective on the creditworthiness of individuals in the dataset.
- There are fewer individuals with scores above 700, indicating that only a small portion of the population has excellent credit.
- The central tendency of the credit scores is around 600-650, which can be considered the average credit score range for the population in this dataset

## Recommendations:

- Provide credit utilization monitoring tools and alerts to help customers maintain their utilization ratios below 30%.
- Target segment (credit score > 700) with premium credit cards, wealth management, and high-interest savings accounts.
- Focus on customers with credit scores between 650-700 by offering balance transfer credit cards, personal loans, and home loan options to meet their financial needs and build loyalty.
- Run targeted marketing campaigns tailored to different customer segments to effectively reach and engage them. -Introduce payment reminders and flexible repayment options to help customers avoid payment delays and improve financial health.

In [ ]: