

Isaac Opiyo Wabwire

Prediction of Loan Status Using Monte Carlo Simulation Abstract

Using the data set provided, we built a simple model using Monte Carlo simulation for predicting the fraction of loans that will default after the 3-year duration of the loan. Our model revealed a 95% confidence interval of $6.61\% \pm 0.01\%$ for Monte-Carlo simulation of $N = 1000$ replicated copies of the data set.

Keywords: Loan status, loan origination, loan charge-off, Monte Carlo simulation, predictive analytics

Introduction: Predicting the status of a loan is an important problem in risk assessment. A bank or financial organization has to be able to estimate the risk involved before granting a loan to a customer. Data Science and predictive analytics play an important role in building models that can be used to predict the probability of loan default. In this project, we are provided with a data set `loan_timing.csv` containing 50000 data points. Each data point represents a loan, and two features are provided as follows:

The column with header “days since origination” indicates the number of days that elapsed between origination and the date when the data was collected.

For loans that charged off before the data was collected, the column with header “days from origination to charge-off” indicates the number of days that elapsed between origination and charge-off. For all other loans, this column is blank.

Project Objective: The goal of this project is to use techniques of data science to estimate what fraction of these loans will have charged off by the time all of their 3-year terms are finished.

Exploratory Data Analysis: The data set was important in python and calculations were performed using python. We plot the following figures:

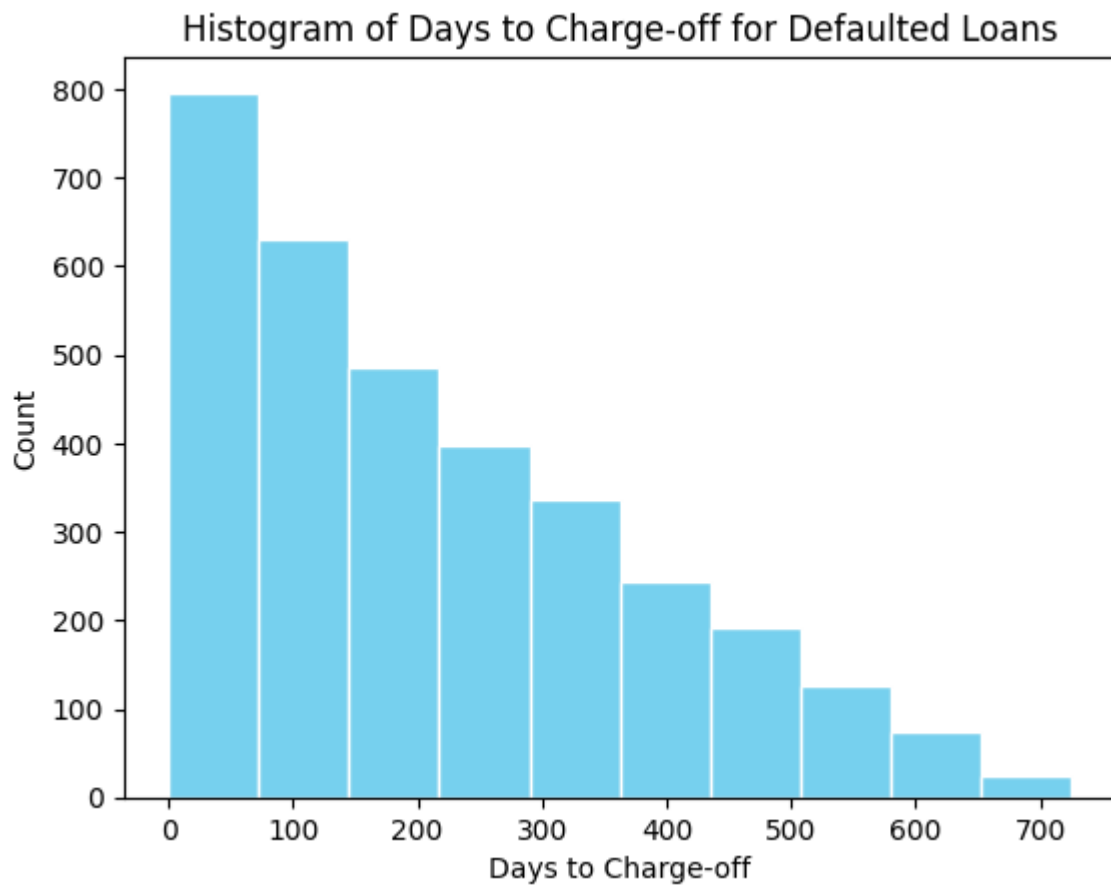


Figure 1: Histogram of days since origination for current loans.

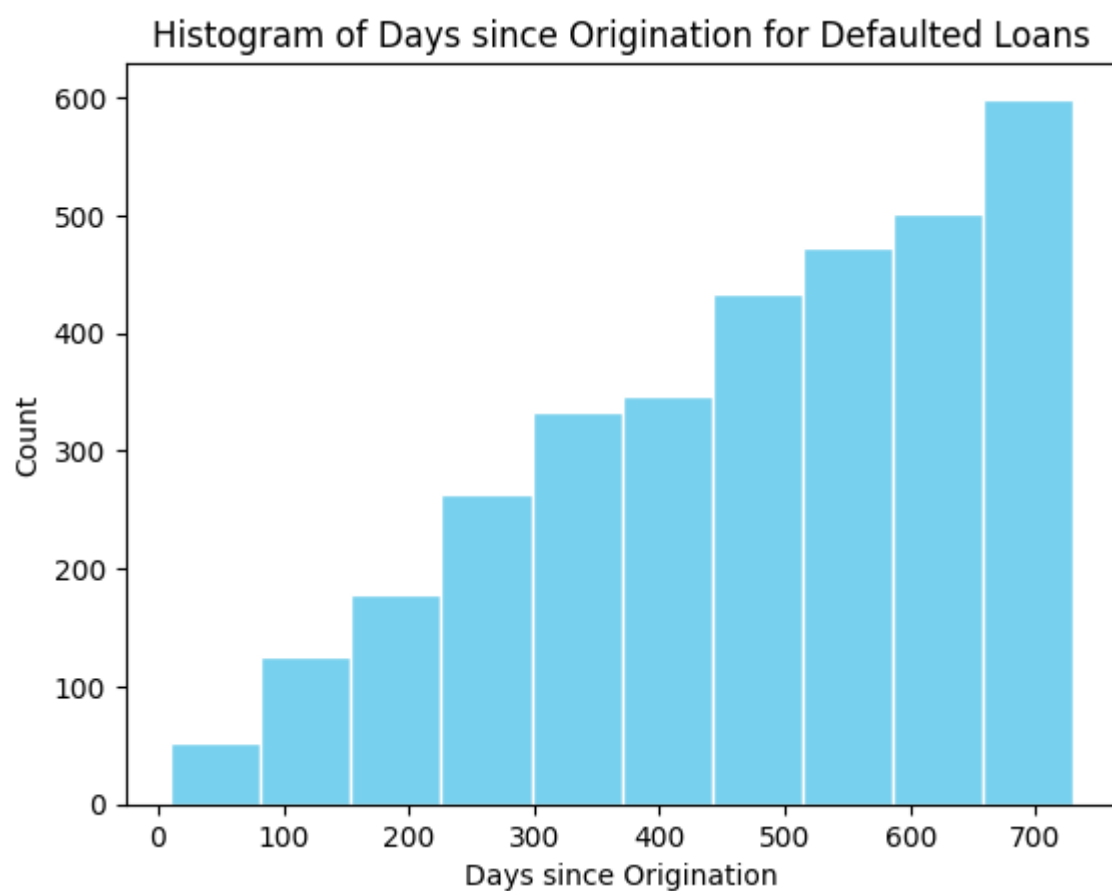


Figure 2: Histogram of days to charge-off for defaulted loans.

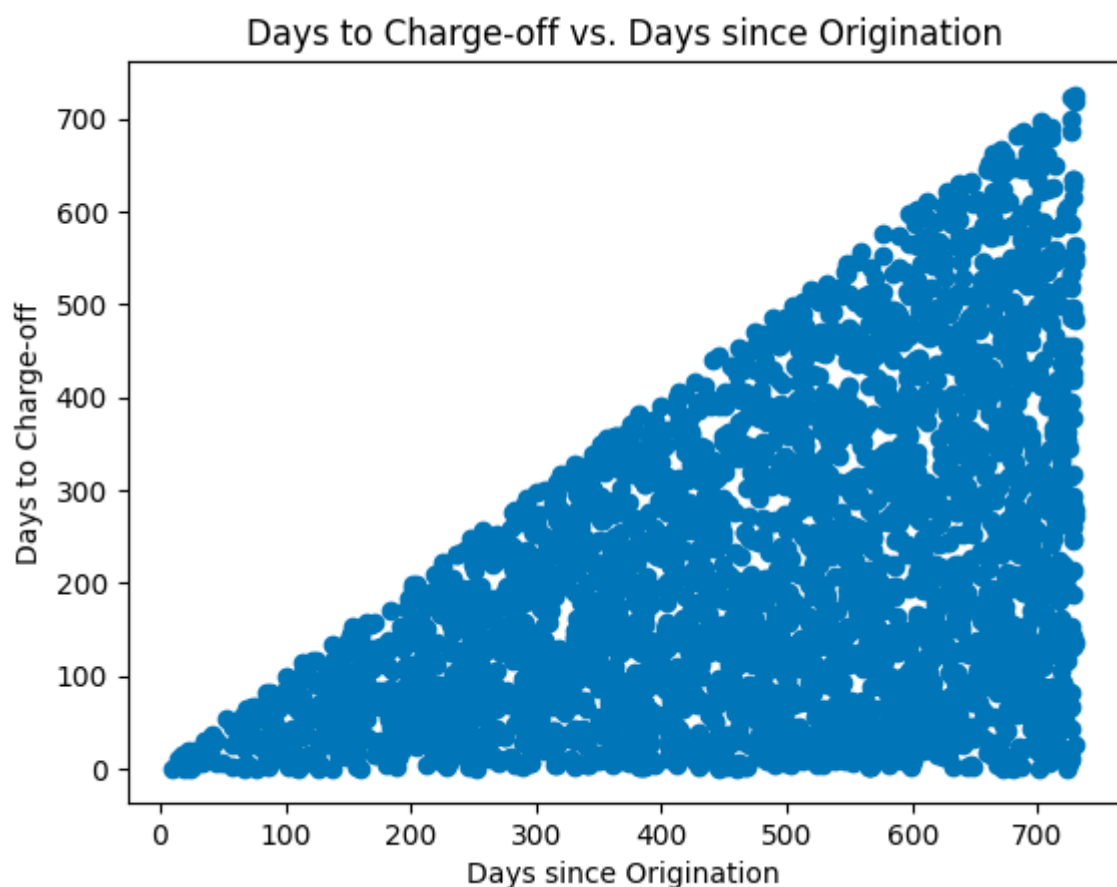


Figure 3: Histogram of days to charge-off for defaulted loans.

Figure 1 shows histogram of active loans, which are uniformly distributed over the days since origination.

From Figure 2, we see that the proportion of loans that are charged off decreases with increasing days from origination to charge-off. This shows that younger loans have a higher probability of defaulting. It also shows that 100% of loans defaulted within 2 years from date of origination.

Figure 3 shows the distribution of defaulted loans as a function of days since origination to the time when data about loan status was collected. The defaulted loans contain a large proportion (71%) of loans that are one year and older. These loans are less likely to default compared to younger loans.

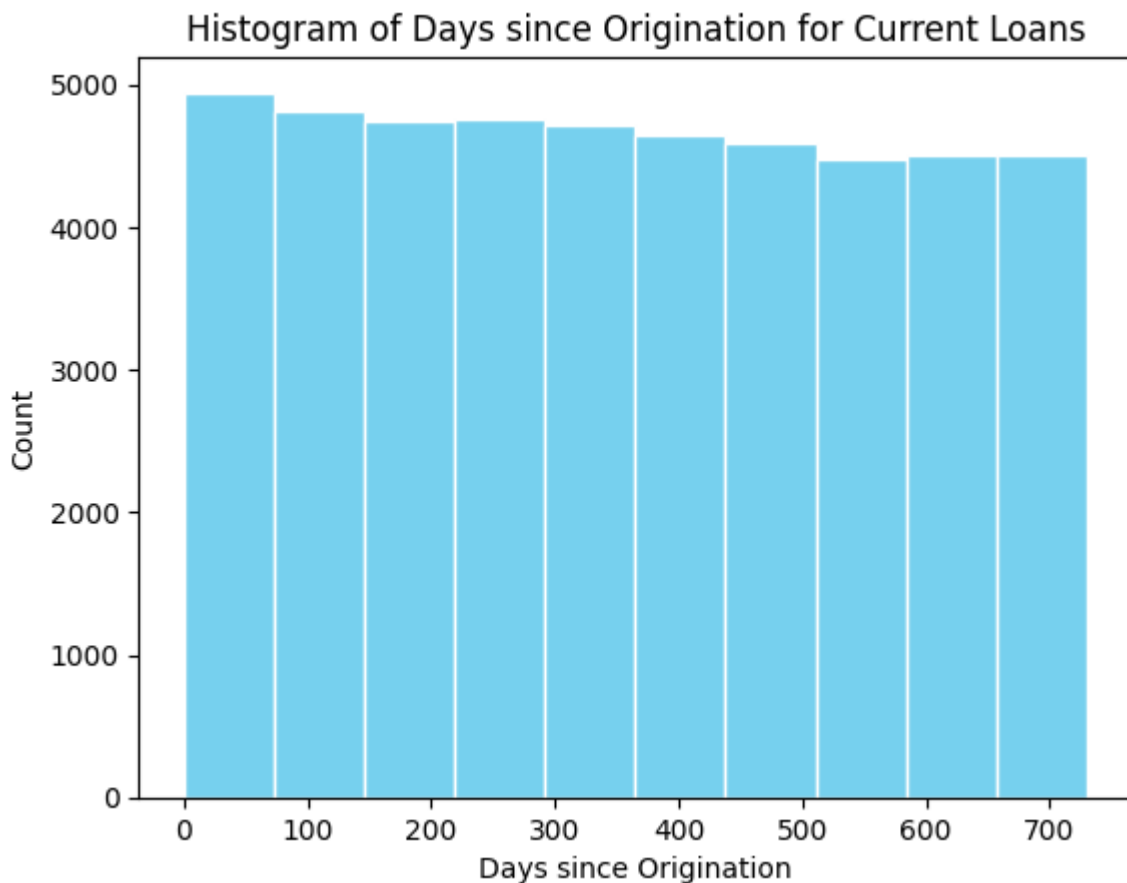


Figure 4: Plot of days to charge-off vs. days since origination for defaulted loans.

Model Selection: Our data set has only 2 features or predictors, and suffers from the problem of prevalence: 93% of the loans have an active status, while 7% has a default status. Use of Linear Regression for predicting fraction of loans that will have charged off after the 3 years loan duration produces a model that is biased towards the active loans.

Figure 3 indicates that relationship between days to charge-off and days since origination for defaulted loans can be simulated using Monte Carlo (MC) simulation. We therefore choose MC simulation as our model for predictive proportion of loans that will default.

Model Calculations: We generate a MC simulation of the defaulted loans, and compare with the original, as shown in figures 4 and 5:

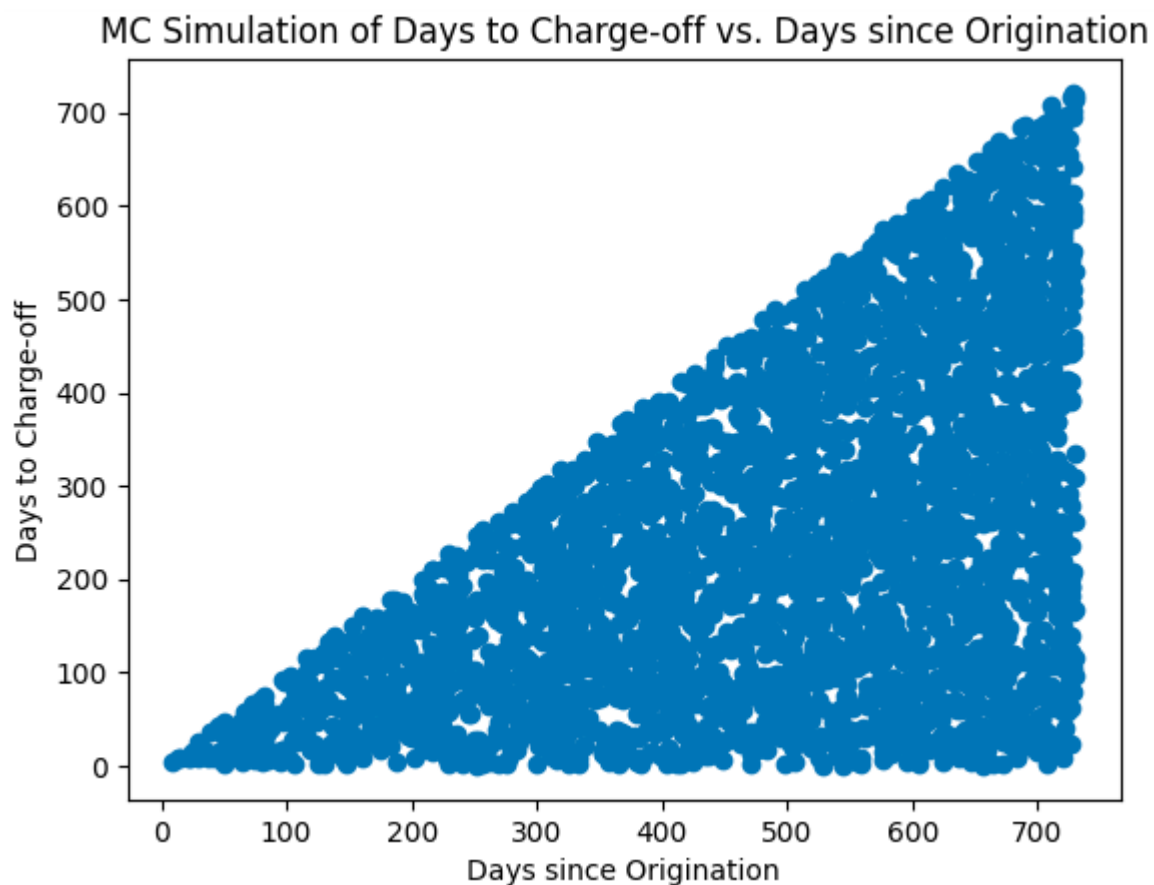


Figure 5: Original and MC simulation of days to charge-off vs. days since origination.

Because there is randomness associated with the charge-off of a loan, we see that MC simulation provides a good approximation for distribution of defaulted loans.

Predictions: Since we have demonstrated that the defaulted loans with charge-off and days since origination in the first 2 years (i.e. 0 to 730 days) can be approximated using a MC simulation, we can predict the fraction of loans that will be charged off by the time all of their 3-year terms are finished using MC simulation.

The total number of charged off loans in our data set is 3,305. This means that there are 46,695 loans that are currently active. Of these active loans, a certain proportion will default over the 3-year period. To estimate the total fraction of defaulted loans, we simulated defaulted loans with charge-off and days since origination covering the entire duration of loan (i.e. 0 to 1095 days), then by appropriate scaling, we computed the fraction of loans that will have charged off after the 3-year term i.e., 1095 days.

By creating 1000 random trials, we obtained the following distribution for the fraction of defaulted loans 3-year term:

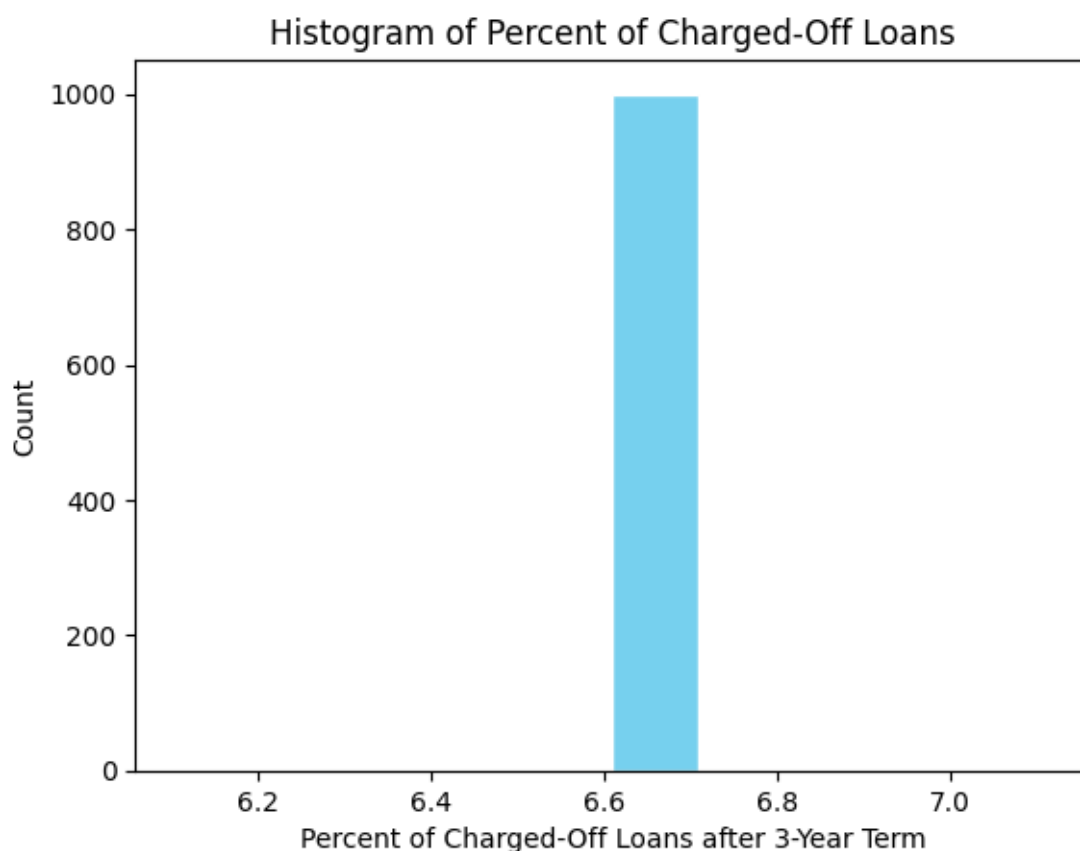


Figure 6: Histogram for fraction of charged off loans after 3-year term using $N = 1000$ samples.

Based on our calculations, the 95% confidence interval for the fraction of loans that will have charged off after the 3-year loan duration is accordingly $6.61\% \pm 0.01\%$.

Conclusions: We have presented a simple model based on the MC simulation for predicting the fraction of loans that will default at the end of the 3-year loan duration period. Different models could be used such as logistic regression, decision tree, etc. I would like to try these different approaches to see if the results are comparable to the MC simulation results.

Appendix: Python Code for Performing Data Analysis

```
# mount user's Google Drive to Google Colab.
from google.colab import drive
drive.mount('/content/drive')

!ls
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

# Specify the file path of the CSV file
file_path = '/content/drive/MyDrive/data/loan_timing.csv'

# Read the CSV file into a pandas DataFrame
df = pd.read_csv(file_path)

df.columns = ["origination", "chargeoff"]

# Partition the dataset into default (charged off) and current loans
default = df.dropna(subset=["chargeoff"])
current = df[df["chargeoff"].isna()]

# Exploratory Data Analysis

# Figure 1: Histogram of days to charge-off for defaulted loans
plt.hist(default["chargeoff"], color="skyblue", edgecolor="white")
plt.xlabel('Days to Charge-off')
plt.ylabel('Count')
plt.title("Histogram of Days to Charge-off for Defaulted Loans")
plt.show()

# Figure 2: Histogram of days since origination for defaulted loans
plt.hist(default["origination"], color="skyblue", edgecolor="white")
plt.xlabel('Days since Origination')
```



```
plt.ylabel('Count')
plt.title("Histogram of Days since Origination for Defaulted Loans")
plt.show()
```

Figure 3: Plot of days to charge-off vs. days since origination for defaulted loans

```
plt.scatter(default["origination"], default["chargeoff"])
plt.xlabel('Days since Origination')
plt.ylabel('Days to Charge-off')
plt.title("Days to Charge-off vs. Days since Origination")
plt.show()
```

Figure 4: Histogram of days since origination for active loans

```
plt.hist(current["origination"], color="skyblue", edgecolor="white")
plt.xlabel('Days since Origination')
plt.ylabel('Count')
plt.title("Histogram of Days since Origination for Current Loans")
plt.show()
```

Monte Carlo Simulation of Defaulted Loans

```
np.random.seed(2)
N = 3 * 365 # Loan duration in days
df_MC = pd.DataFrame({
    'u': np.round(np.random.uniform(0, N, 15500)),
    'v': np.round(np.random.uniform(0, N, 15500))
})
df_MC = df_MC[(df_MC["v"] <= df_MC["u"]) & (df_MC["u"] <= 730) & (df_MC["v"] <= 730)]
```

```
plt.scatter(df_MC["u"], df_MC["v"])
plt.xlabel('Days since Origination')
plt.ylabel('Days to Charge-off')
plt.title("MC Simulation of Days to Charge-off vs. Days since Origination")
plt.show()
```

Predicting the fraction of loans that will charge off by the end of their 3-year terms

```
np.random.seed(2)
```

```
B = 1000
```

```
fractions = []
```

```
for _ in range(B):
```

```
    df2 = pd.DataFrame({
        'u': np.round(np.random.uniform(0, N, 50000)),
        'v': np.round(np.random.uniform(0, N, 50000))
    })
```

```
    df2 = df2[(df2["v"] <= df2["u"]) & (df2["u"] <= 730) & (df2["v"] <= 730)]
```

```
    total = (len(df2) / len(df2[df2["u"] <= 730])) * len(default)
```

```
    fraction = 100.0 * (total / 50000.0)
```

```
    fractions.append(fraction)
```

```
mean_fraction = np.mean(fractions)

print("Mean Fraction of Charged-Off Loans:", mean_fraction)

# Plot the histogram of the percent of charged-off loans
plt.hist(fractions, color="skyblue", edgecolor="white")
plt.xlabel('Percent of Charged-Off Loans after 3-Year Term')
plt.ylabel('Count')
plt.title("Histogram of Percent of Charged-Off Loans")
plt.show()

# Calculate Confidence Interval of Percentage of Defaulted Loans after 3-year term
std_deviation = np.std(fractions)
confidence_interval = (mean_fraction - 2 * std_deviation, mean_fraction + 2 * std_deviation)
print("Confidence Interval:", confidence_interval)
```