

PROJECT #2

Cesar Sandiford

LOAN DEFAULT PREDICTION

This project is about using machine learning to predict loan defaults with given data about customers with a Supervised model.

In times of uncertainty, FinTechs are in big risk when it comes to higher percentage of loan defaults and using machine learning to understand and identify the risk is very important to have in hand so they can make decisions to raise the credit standards in the future and minimize risk.

INSPIRED BY:

★ Personal Saving Rate (PSAVERT)

[DOWNLOAD](#)

Observation:

Nov 2022: **2.4** (+ more)

Updated: Dec 23, 2022

Units:

Percent,
Seasonally Adjusted Annual Rate

Frequency:

Monthly

1Y | 5Y | 10Y | Max

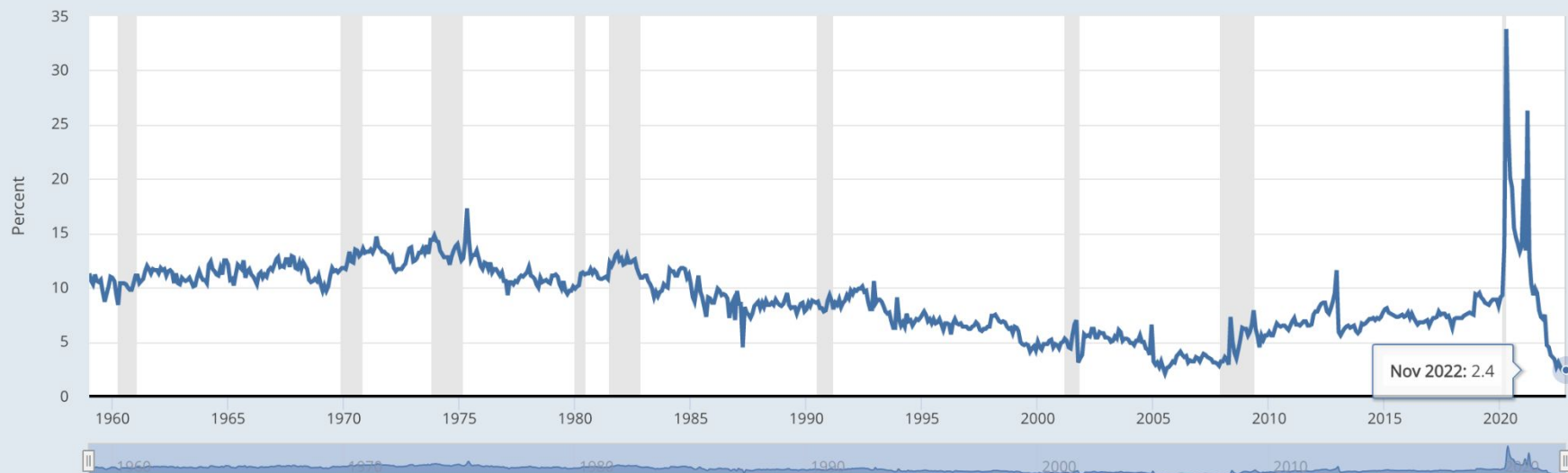
1959-01-01

to

2022-11-01

[EDIT GRAPH](#)

FRED — Personal Saving Rate



Shaded areas indicate U.S. recessions.

Source: U.S. Bureau of Economic Analysis

fred.stlouisfed.org



☆ Consumer Loans: Credit Cards and Other Revolving Plans, All Commercial Banks

(CCLACBW027SBOG)

DOWNLOAD 

Observation:

2023-01-04: **940.0455** (+ more)

Updated: Jan 13, 2023

Units:

Billions of U.S. Dollars,
Seasonally Adjusted

Frequency:

Weekly,
Ending Wednesday

1Y | 5Y | 10Y | Max

2000-06-28

to

2023-01-04

EDIT GRAPH 

FRED — Consumer Loans: Credit Cards and Other Revolving Plans, All Commercial Banks



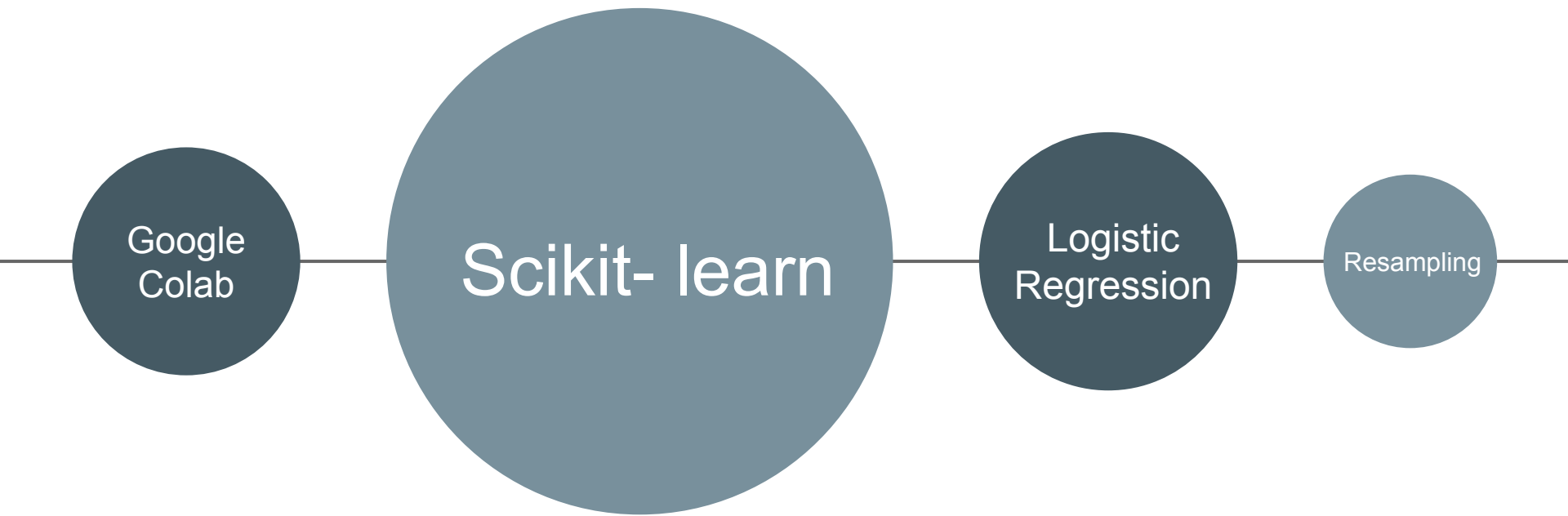
Shaded areas indicate U.S. recessions.

Source: Board of Governors of the Federal Reserve System (US)

fred.stlouisfed.org



SKILLS USED



WALKTHROUGH

IMPORTING LIBRARIES

Numpy

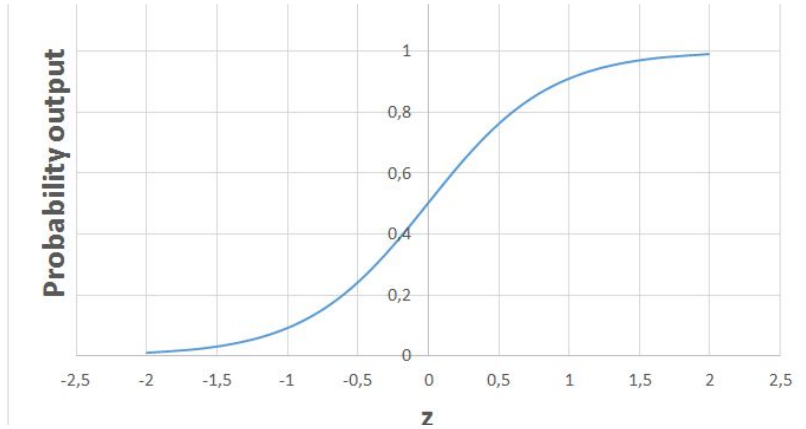
Pandas

Pathlib-Path

Sklearn

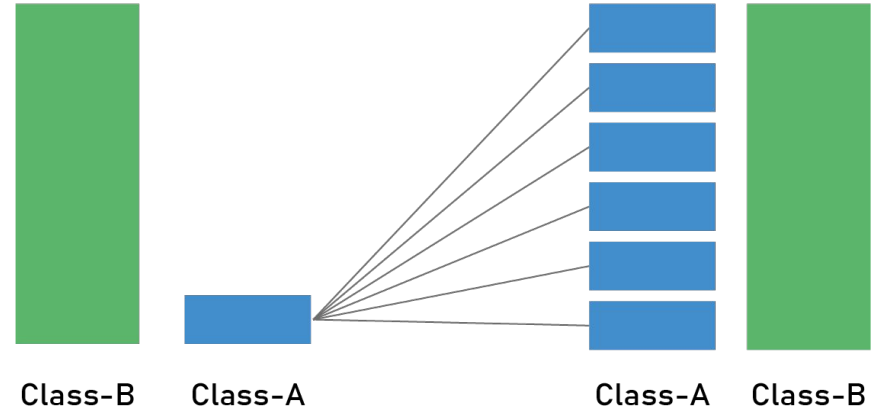
Train-test-split

Classification report



LOGISTIC REGRESSION

Logistic regression can be used to predict default events and model the influence of different variables on a consumer's credit-worthiness. In this paper we use a logistic regression model to predict the creditworthiness of bank customers using predictors related to their personal status and financial history.



RESAMPLING

By resampling our data, we have the chance to create better-performing machine learning models. And, the more resampling techniques that we have, the better our chances are of identifying one that works well for our current data. In this case, the oversampling was the one for this task.

CONCLUSION

By having the Logistic Regression model applied to the dataset, the balanced accuracy score was 58% meaning it was not too accurate. Since the data was imbalanced by the defaults, I gave it a proper chance by oversampling it to increase the accuracy of the mode up to 84%. Huge improvement! F1 score average was 94 vs the resampled average was 90. In conclusion, Oversampling the model brings a better prediction for this situation. Also, predicting defaults can be difficult given that customers could have different lifestyle, emergencies, unexpected situation or money opportunities in life.

CREDITS

KAGGLE:

[HTTPS://WWW.KAGGLE.COM/DATASETS/B934BBD9D19E1C321CF5F121B9B8F9BEA4F1E770EA8D4CFDF445FB34AB80F42C?RESOURCE=DOWNLOAD](https://www.kaggle.com/datasets/b934bbd9d19e1c321cf5f121b9b8f9bea4f1e770ea8d4cfdf445fb34ab80f42c?resource=download)

[HTTPS://WWW.KAGGLE.COM/DATASETS/ITSSURU/LOAN-DATA](https://www.kaggle.com/datasets/itssuru/loan-data)

MODULE 12 CHALLENGE

GOOGLE.COM

STACK OVERFLOW

FEDERAL RESERVE ECONOMIC DATA: [HTTPS://FRED.STLOUISFED.ORG/SERIES/PSAVERT](https://fred.stlouisfed.org/series/PSAVERT) & [HTTPS://FRED.STLOUISFED.ORG/SERIES/CCLACBW027SBOG](https://fred.stlouisfed.org/series/CCLACBW027SBOG)