

# Predicting Credit Card Default Using Machine Learning

Jesse Mailhot | Mentor: Shmuel Naaman | Springboard Data Science Career Track

---

1/21/2019



# Contents

---

- ❖ Introduction
- ❖ Dataset
- ❖ Data Wrangling
- ❖ Data Analysis
- ❖ Data Modeling
- ❖ Analysis Results



# Introduction

---

In this project, I aim to solve the following problems:

- ❖ Which features are the best predictor of default?
- ❖ What is the relationship between default and other features?
  - ❖ Marital status, education level, gender, etc...



# Dataset

---

## Credit Card Default

- ❖ Demographic data: sex, marital status, education level, age
- ❖ Payment data: credit limit, bill amount, payment amount, payment delay
- ❖ April - September, 2005
- ❖ 30,000 Taiwanese credit card customers.

Source: UCI Machine Learning Repository



# Data Wrangling

---

- ❖ Change SEX to 0 & 1
- ❖ Change unlabeled EDUCATION values to 'other'
- ❖ Change unlabeled MARRIAGE values to 'other'
- ❖ Binarize PAY\_0 - PAY\_6



# Data Analysis

---

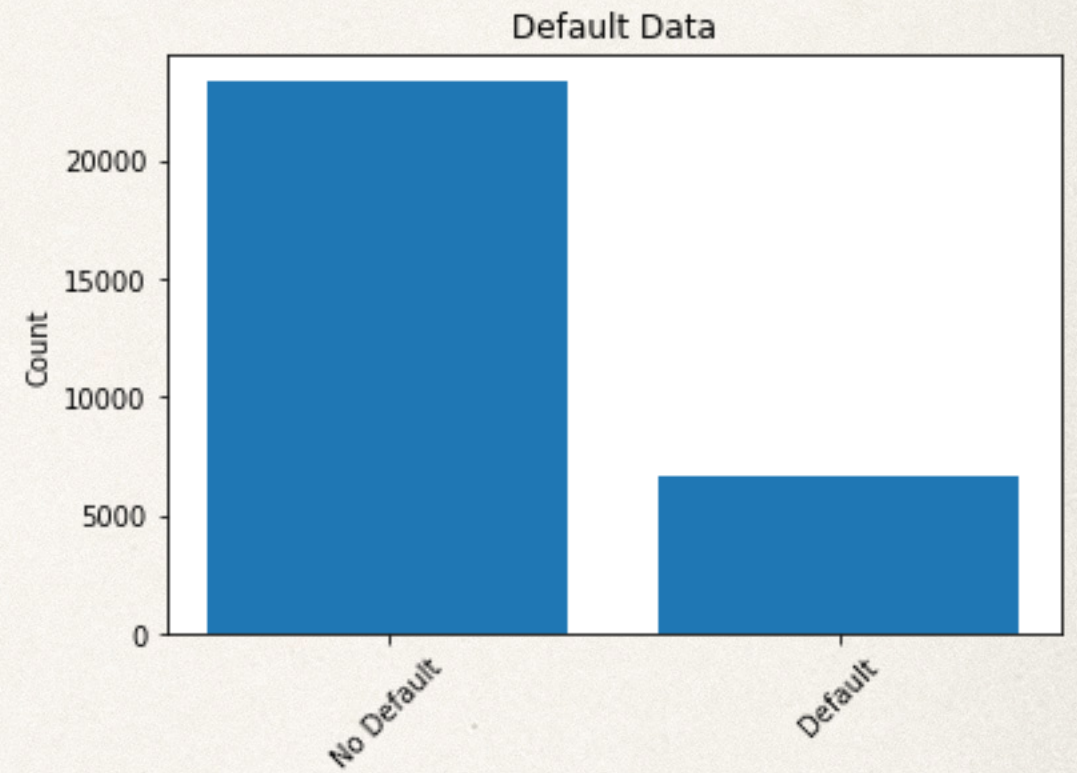
- ❖ Investigate each variable in relation to default
- ❖ Look for trends and outliers



# Default

---

- ❖ Unbalanced dataset
- ❖ 22.12% default percentage

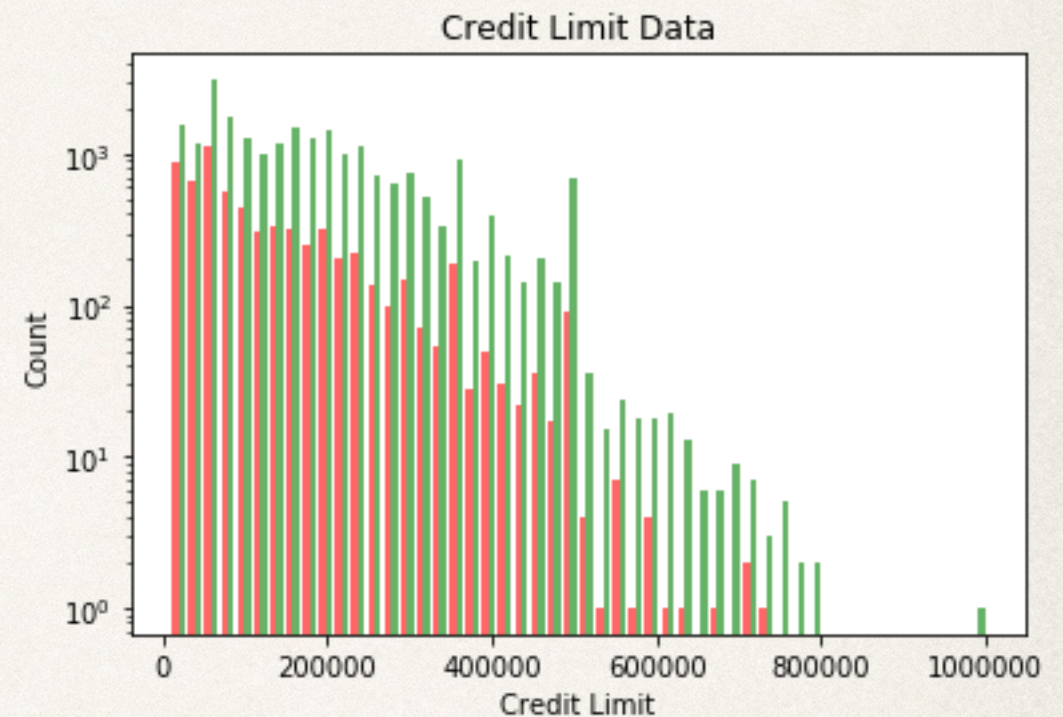




# Credit Limit

---

- ❖ Mean credit limit for defaulting clients is NT\$47990 lower than non-defaulting clients

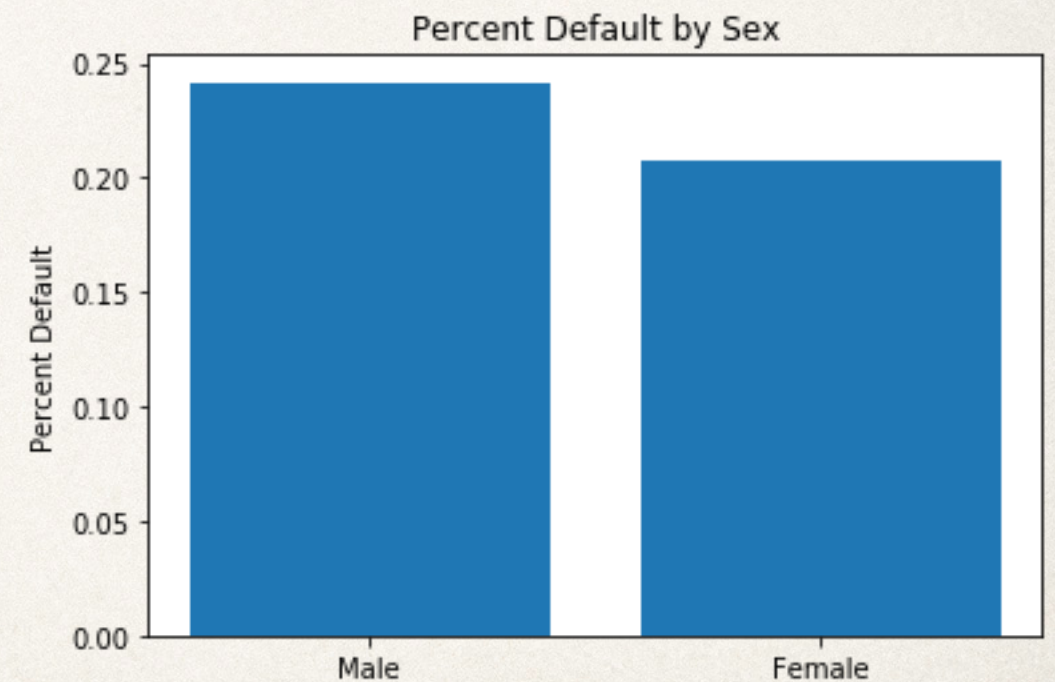
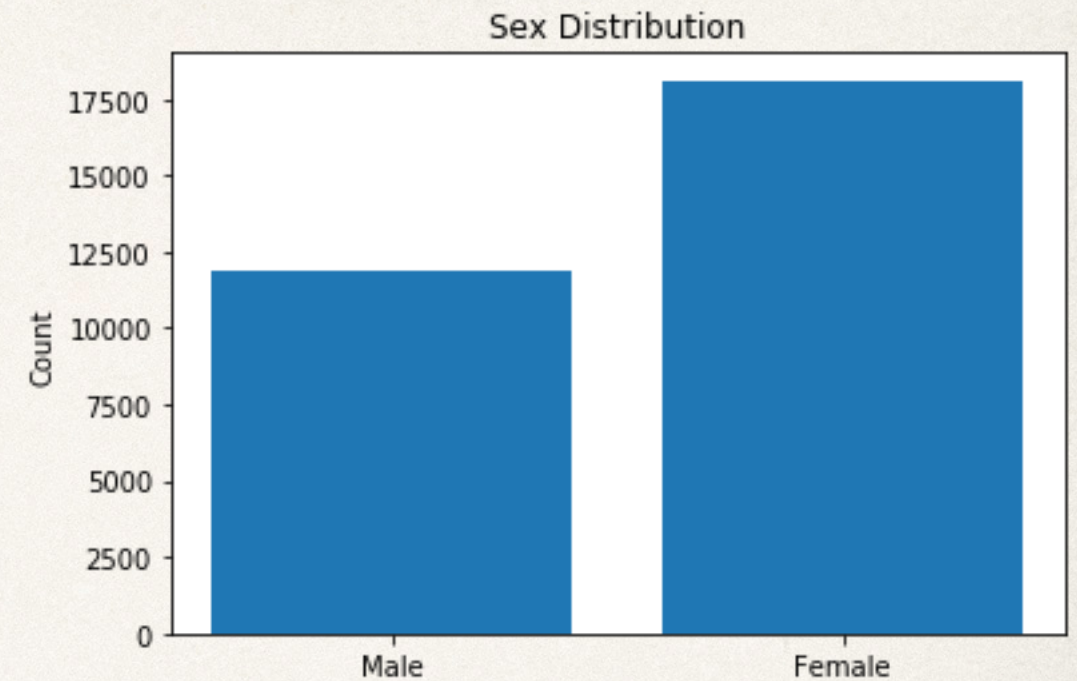




# Sex

---

- ❖ More female clients than male clients
- ❖ Male clients have higher chance of default

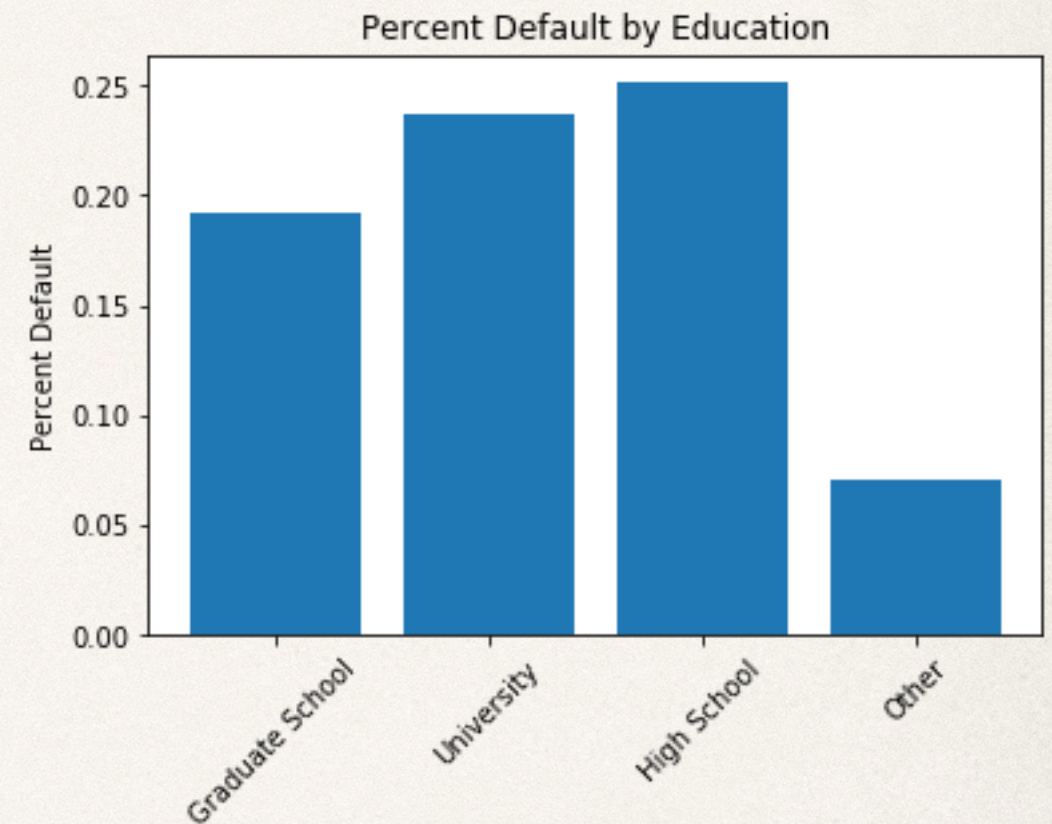




# Education

---

- ❖ Negative correlation between education level and default
- ❖ Low sample size - Other

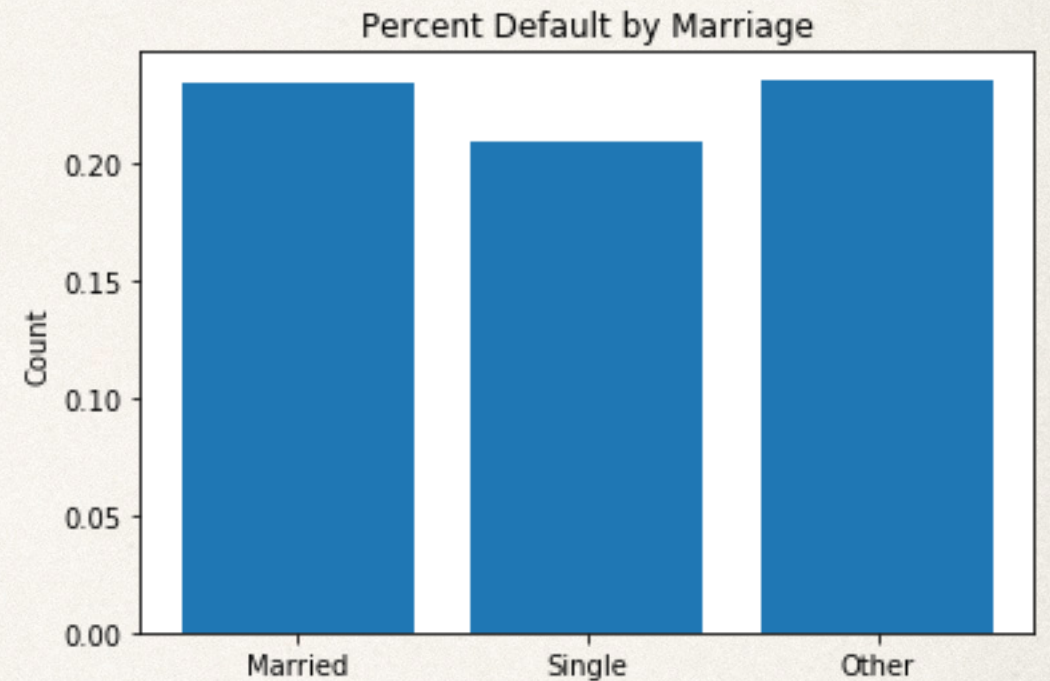




# Marriage

---

- ❖ Single has lower percent default than married
- ❖ Low sample size - Other
  - ❖ Other = divorced?

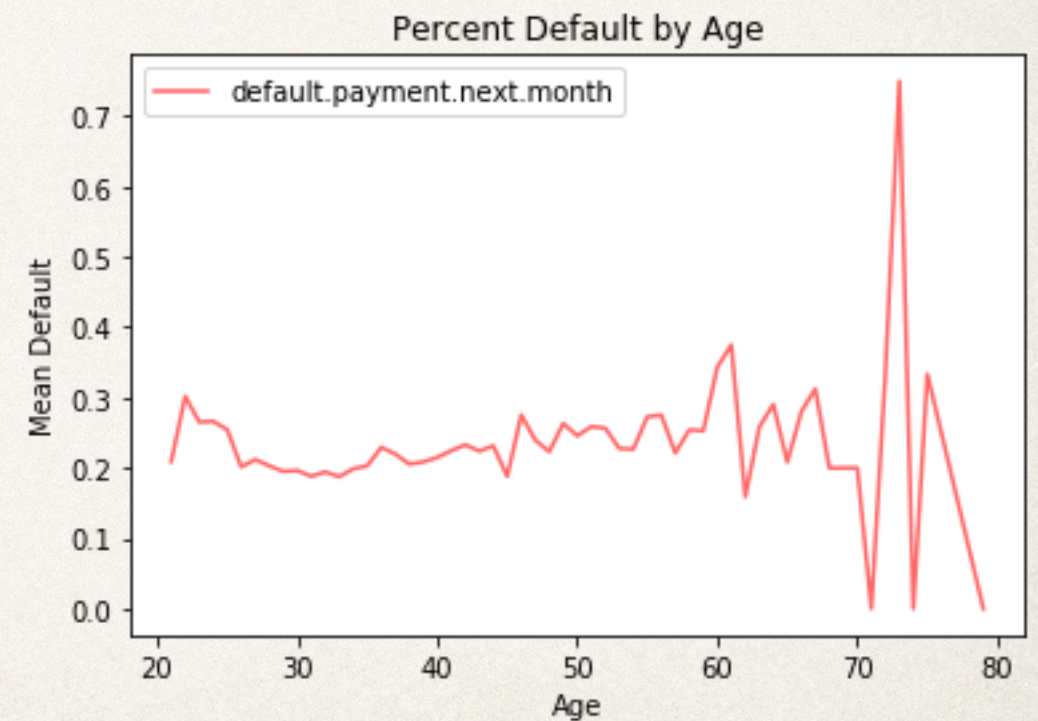
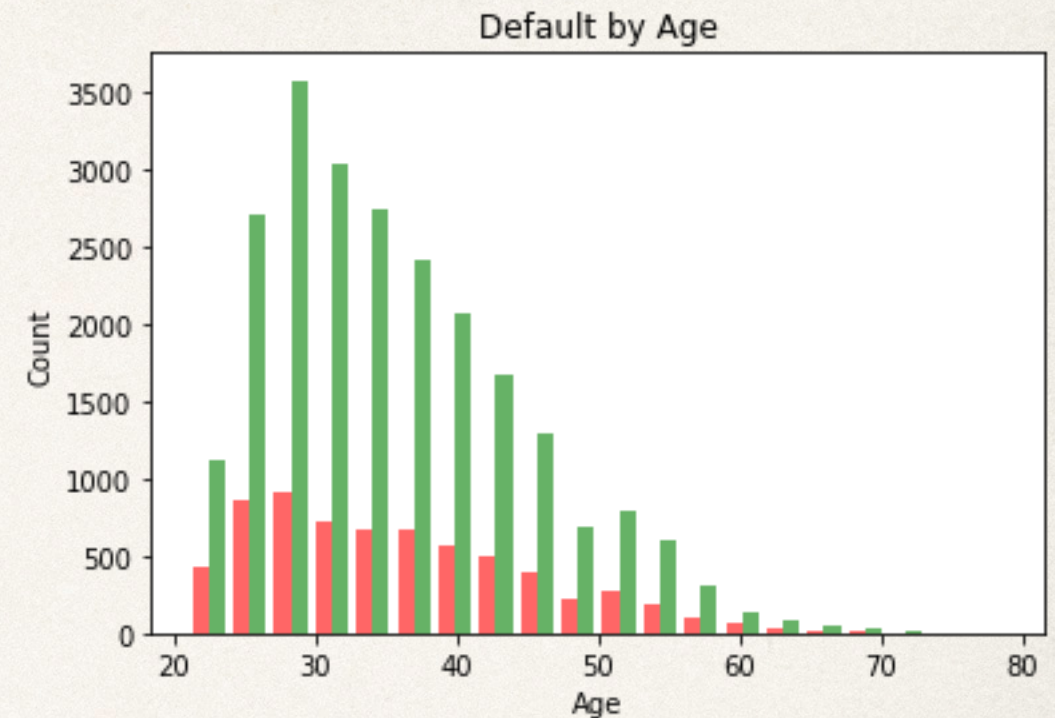




# Age

---

- ❖ Default highest in ages 20-25
- ❖ Decreases after 25
- ❖ Increases after 35
- ❖ Low sample size for ages 50+



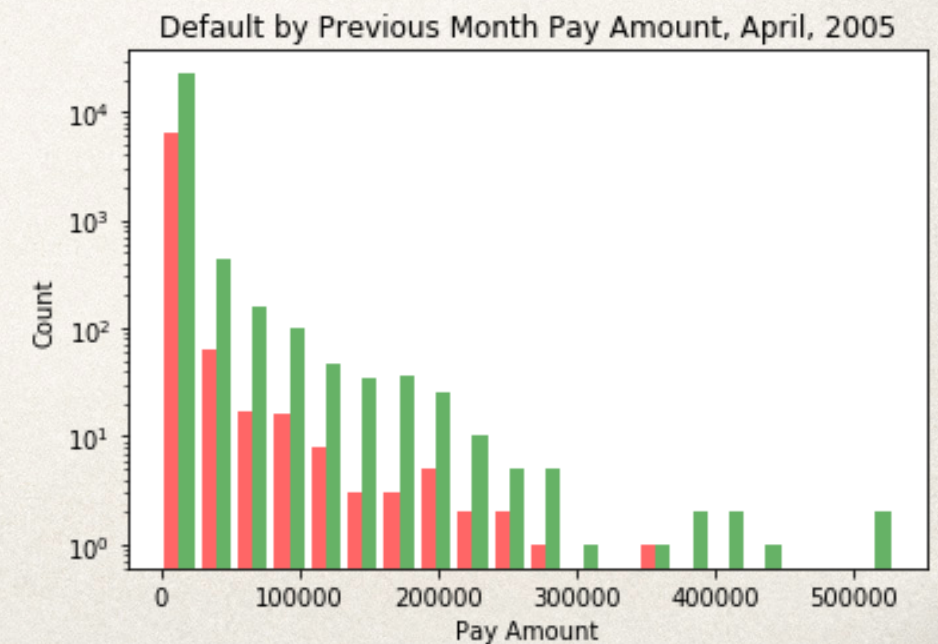
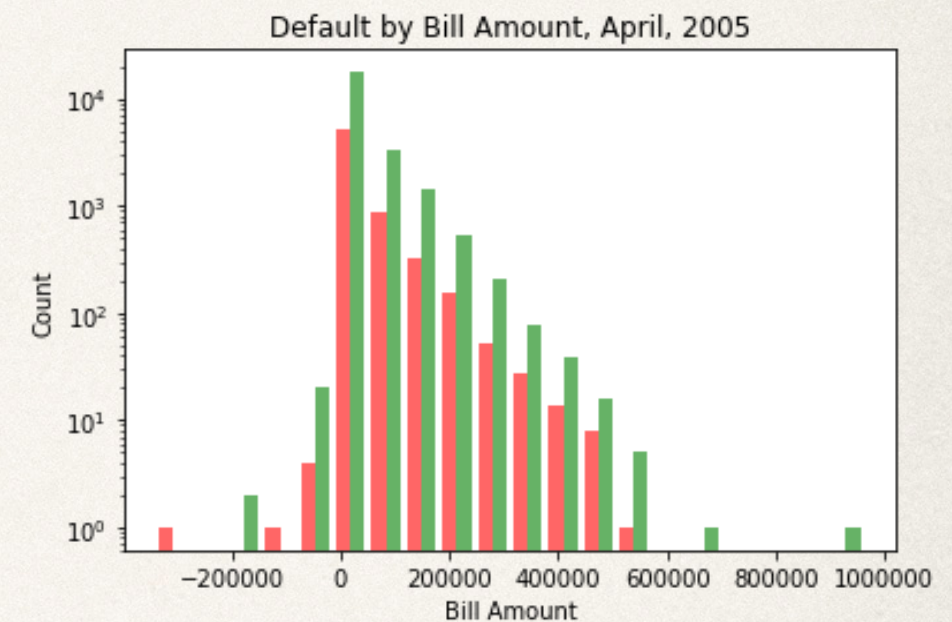
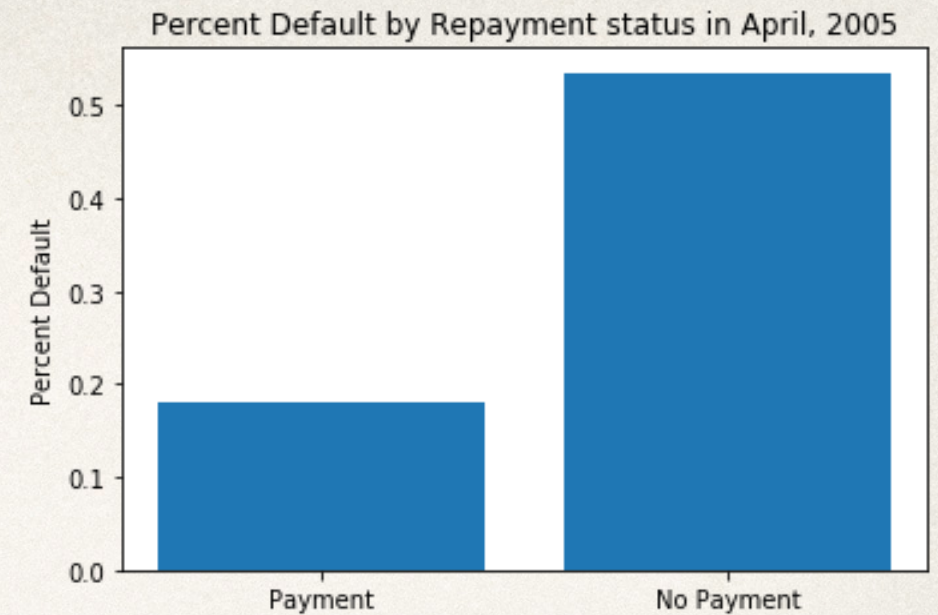


# Pay Data

## April, 2005

---

- ❖ Clients who made payments still defaulted (?)
- ❖ Clients with bill amount 0 or less defaulted (?)
- ❖ Negative correlation between pay amount and default



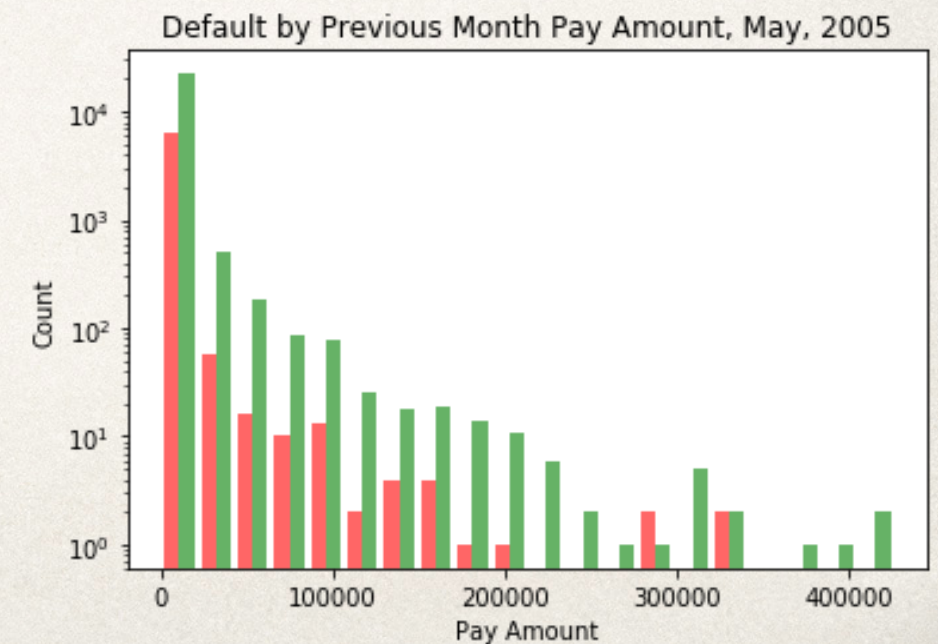
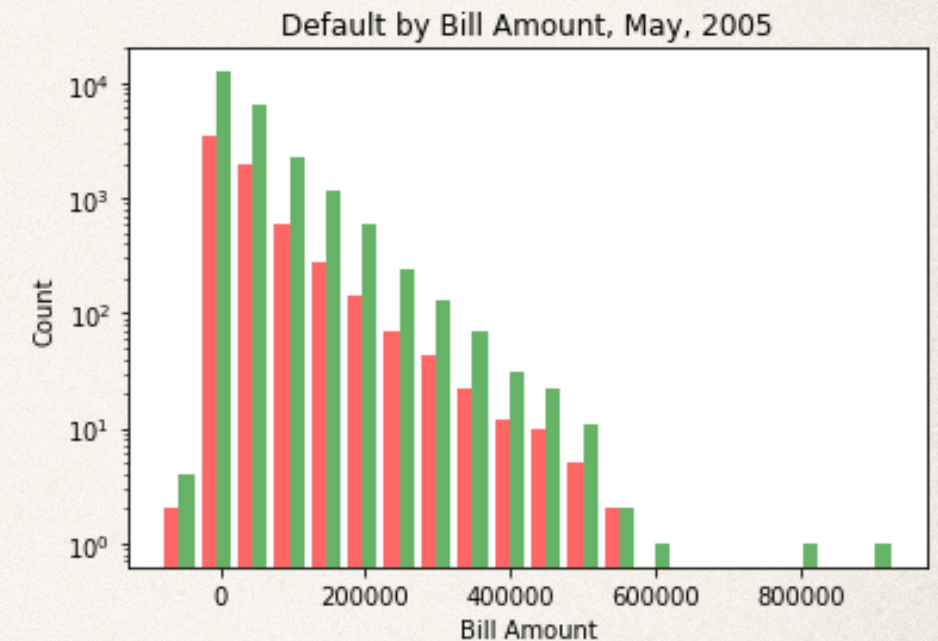
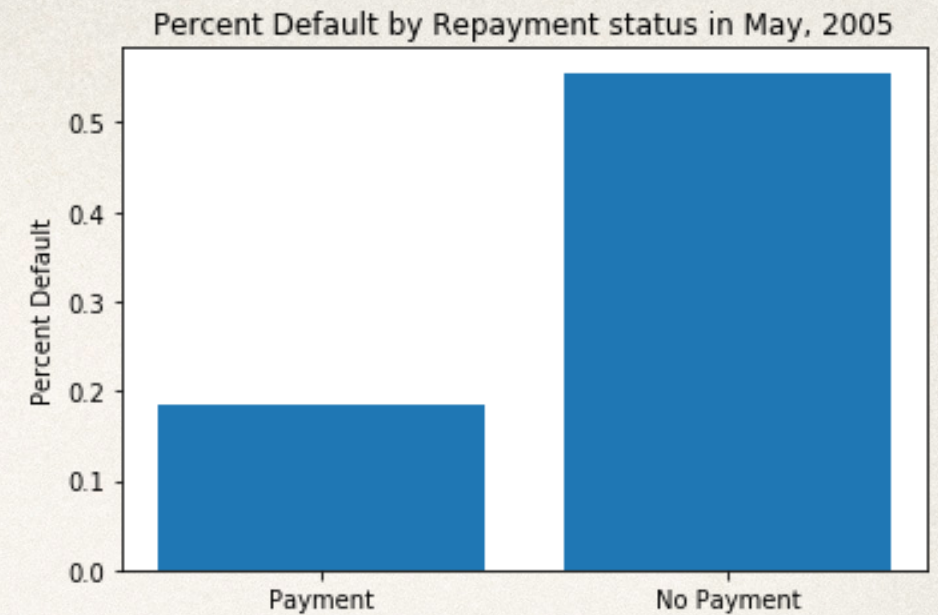


# Pay Data

## May, 2005

---

- ❖ Clients who made payments still defaulted (?)
- ❖ Clients with bill amount 0 or less defaulted (?)
- ❖ Negative correlation between pay amount and default



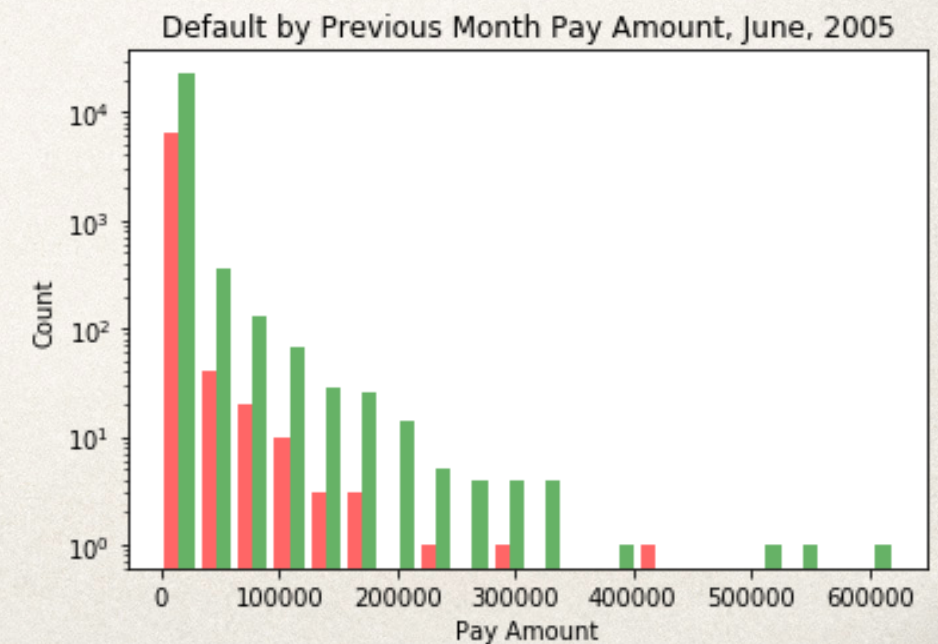
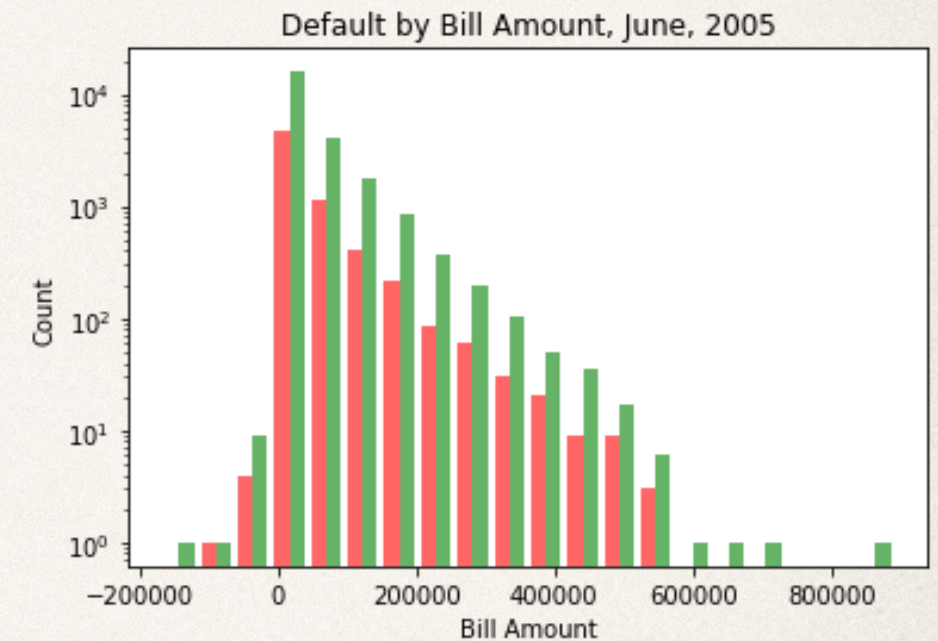
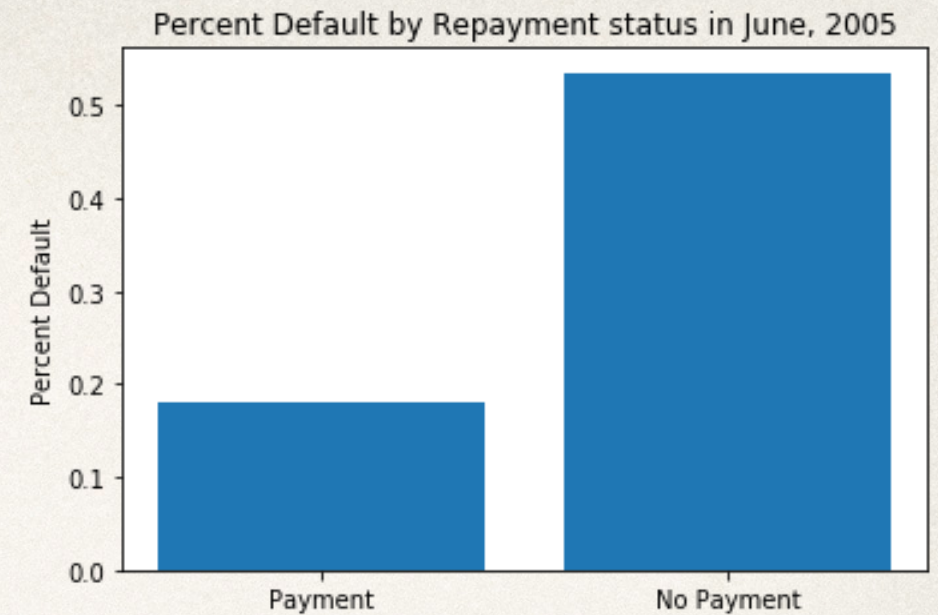


# Pay Data

## June, 2005

---

- ❖ Clients who made payments still defaulted (?)
- ❖ Clients with bill amount 0 or less defaulted (?)
- ❖ Negative correlation between pay amount and default

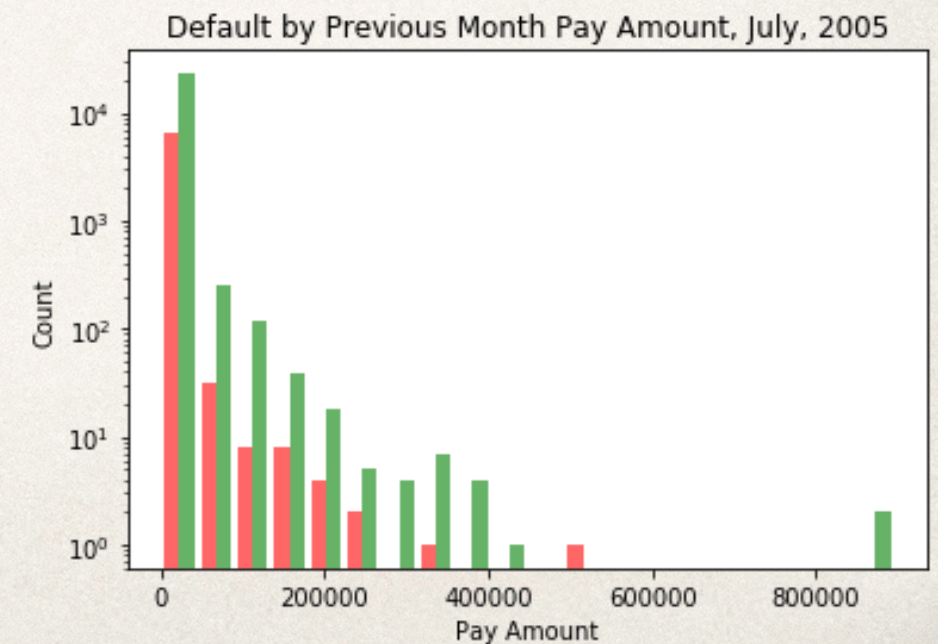
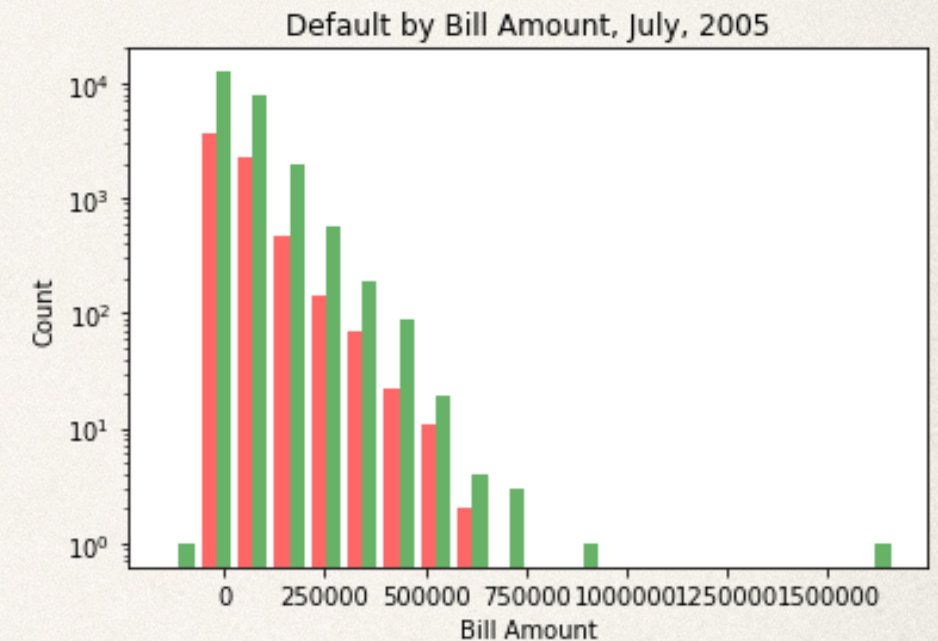
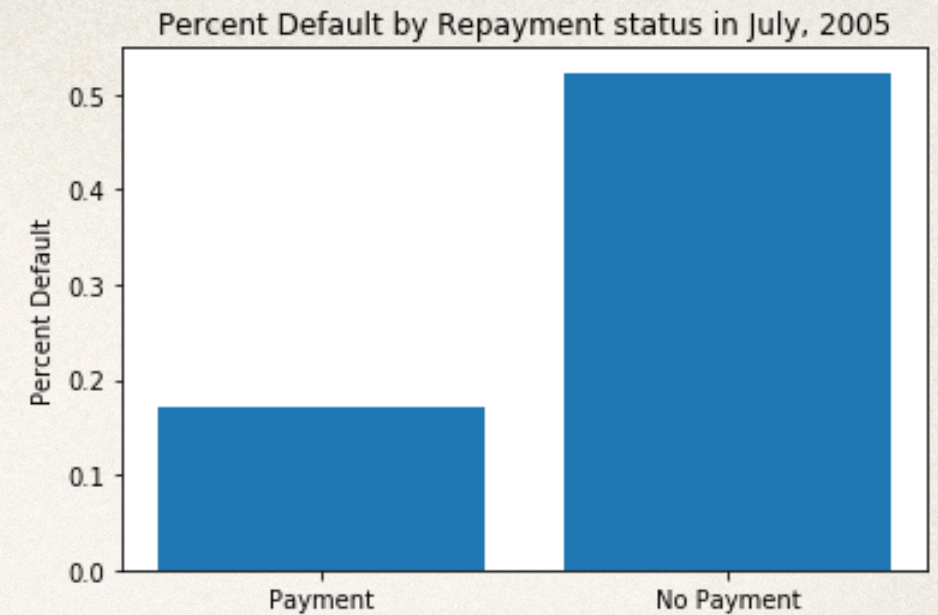




# Pay Data July, 2005

---

- ❖ Clients who made payments still defaulted (?)
- ❖ Clients with bill amount 0 or less defaulted (?)
- ❖ Negative correlation between pay amount and default



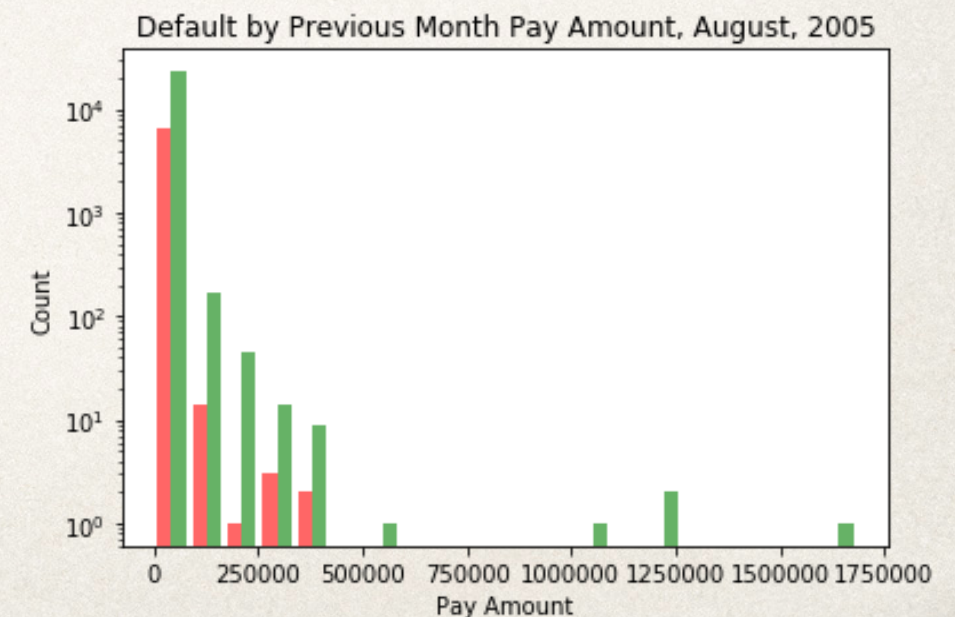
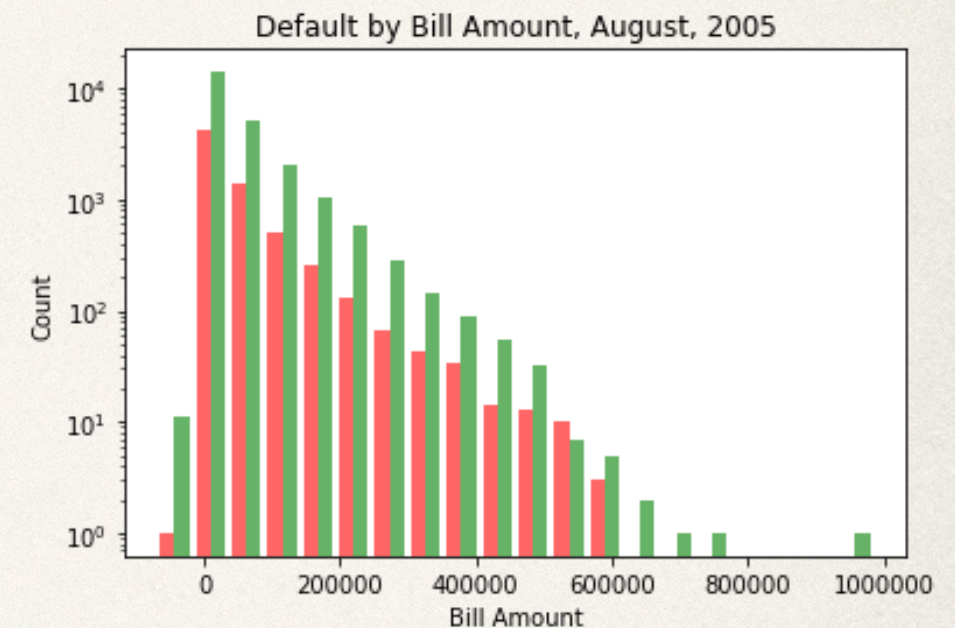
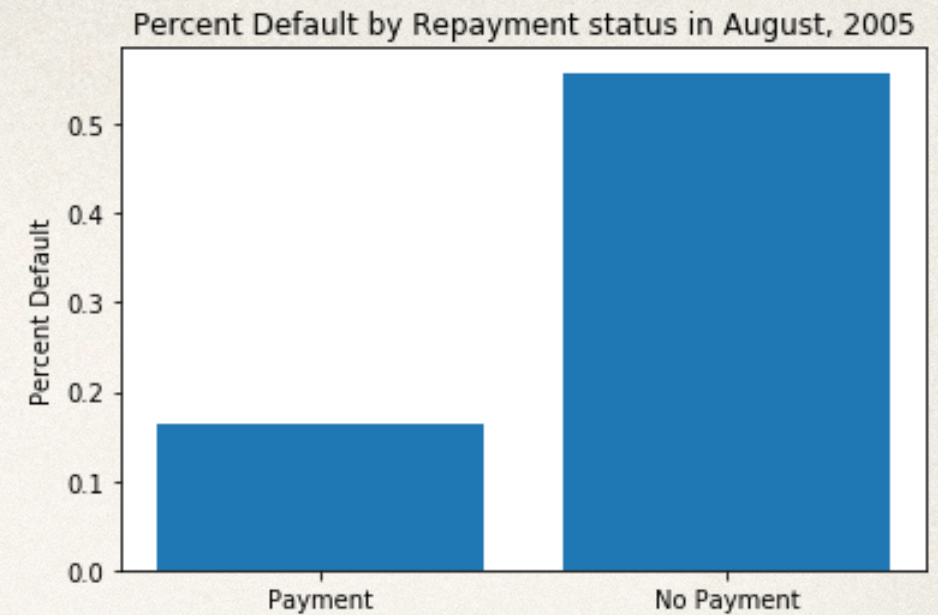


# Pay Data

## August, 2005

---

- ❖ Clients who made payments still defaulted (?)
- ❖ Clients with bill amount 0 or less defaulted (?)
- ❖ Negative correlation between pay amount and default



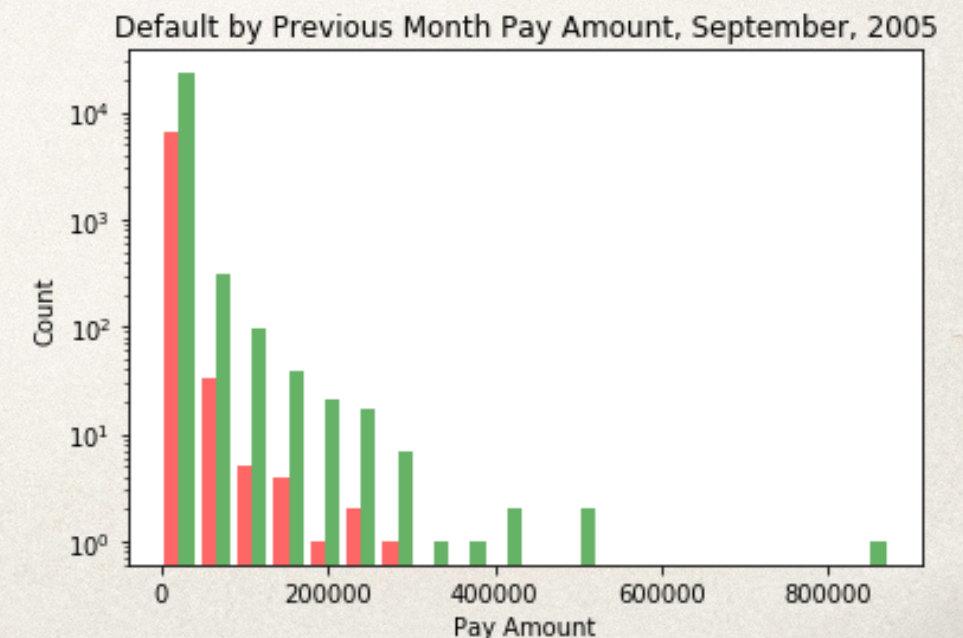
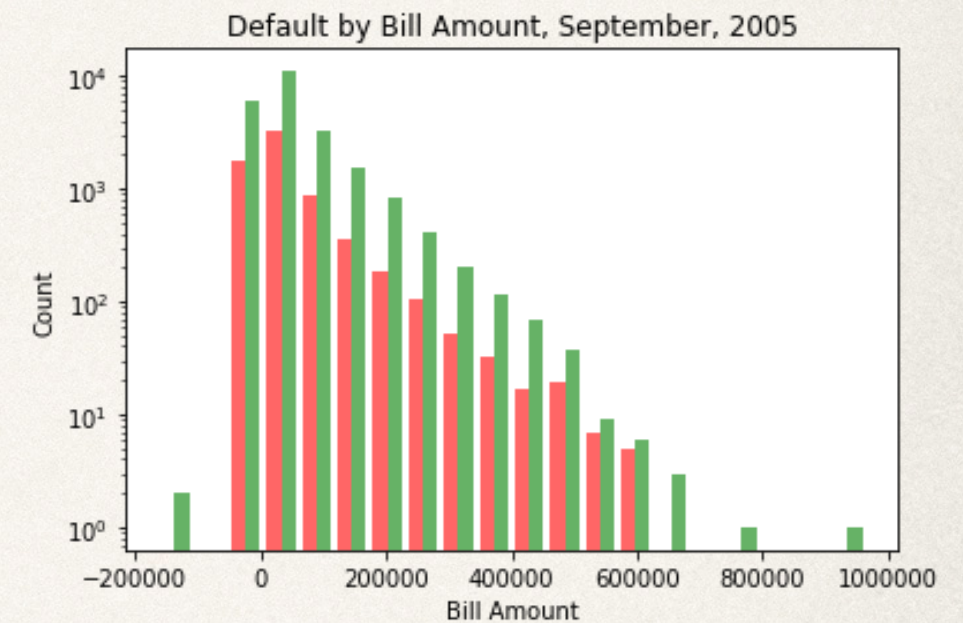
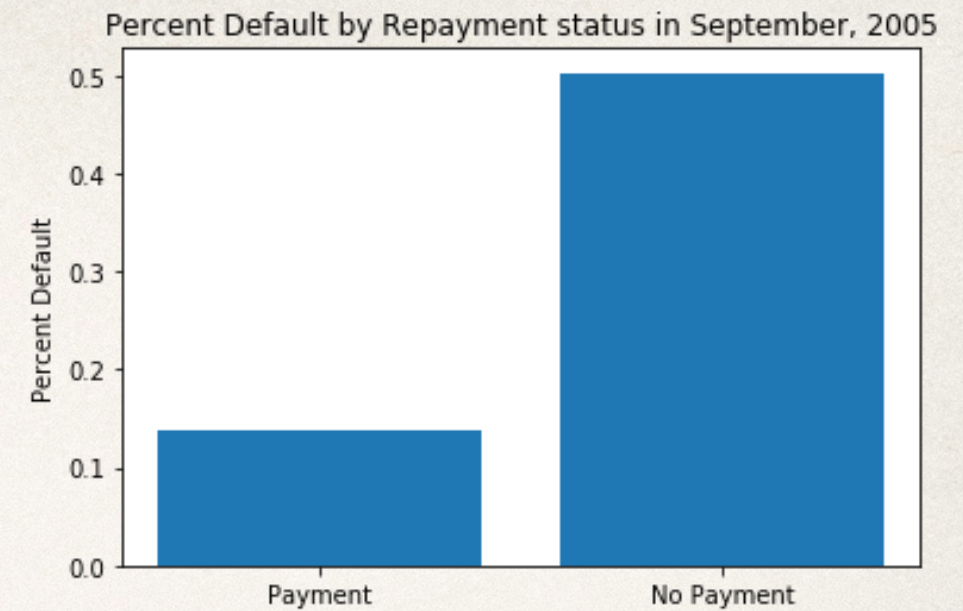


# Pay Data

## September, 2005

---

- ❖ Clients who made payments still defaulted (?)
- ❖ Clients with bill amount 0 or less defaulted
- ❖ Negative correlation between pay amount and default





# Data Modeling

---

- ❖ Engineer 'PIF' feature
- ❖ Split categorical columns into dummy variables
- ❖ Feature selection
- ❖ Model selection
- ❖ Metric selection
- ❖ Model tuning



# Engineering 'PIF'

---

- ❖  $\text{PAY\_AMT}(N) \geq \text{BILL\_AMT}(N - 1)$
- ❖ Drop clients who closed their accounts
  - ❖  $\text{PIF2} = \text{True} \ \& \ \text{default} = \text{True}$



# Split categorical columns

---

- ❖ Marriage, Education
- ❖ Pandas *get\_dummies*



# Feature Selection

---

- ❖ Exclude bill amount, pay amount columns
  - ❖ Described by 'PIF'



# Model Selection

---

- ❖ SciKit Learn
  - ❖ Decision Tree
  - ❖ Support Vector Machine
  - ❖ Random Forest
  - ❖ AdaBoost



# Metric Selection

---

Goal: save the bank money

- ❖ Accuracy:  $(TP + TN) / (FP + FN)$
- ❖ Precision:  $TP / (TP + FP)$
- ❖ Recall:  $TP / (TP + FN)$
- ❖ F1:  $2 * (precision * recall) / (precision + recall)$



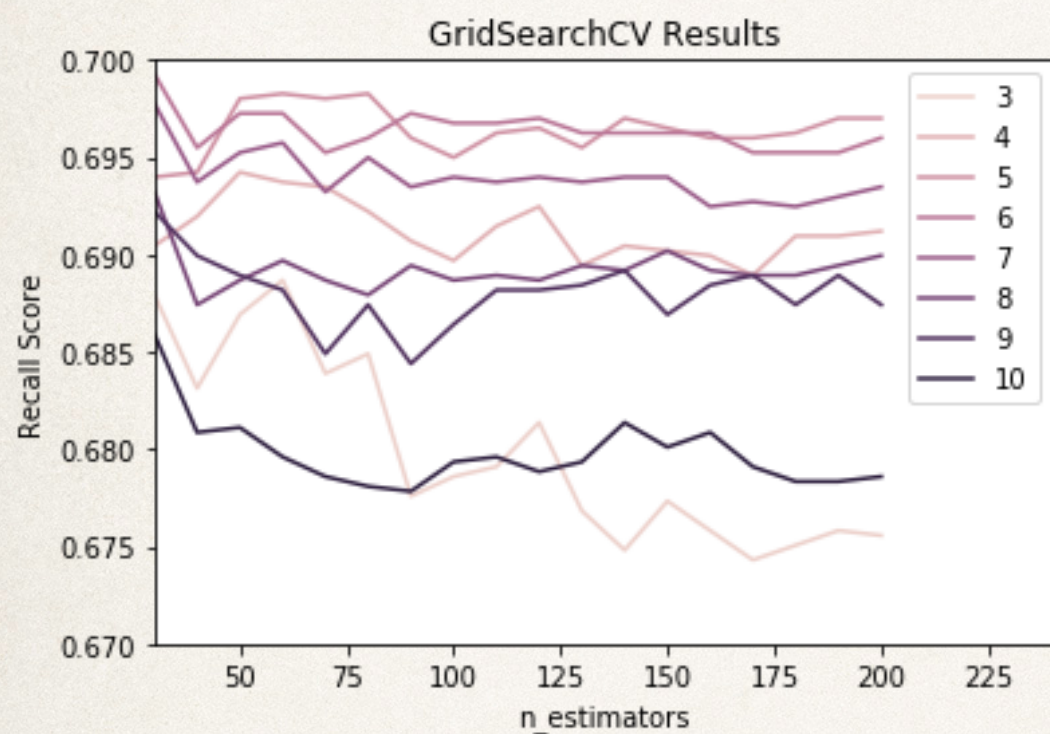
# Model Tuning

---

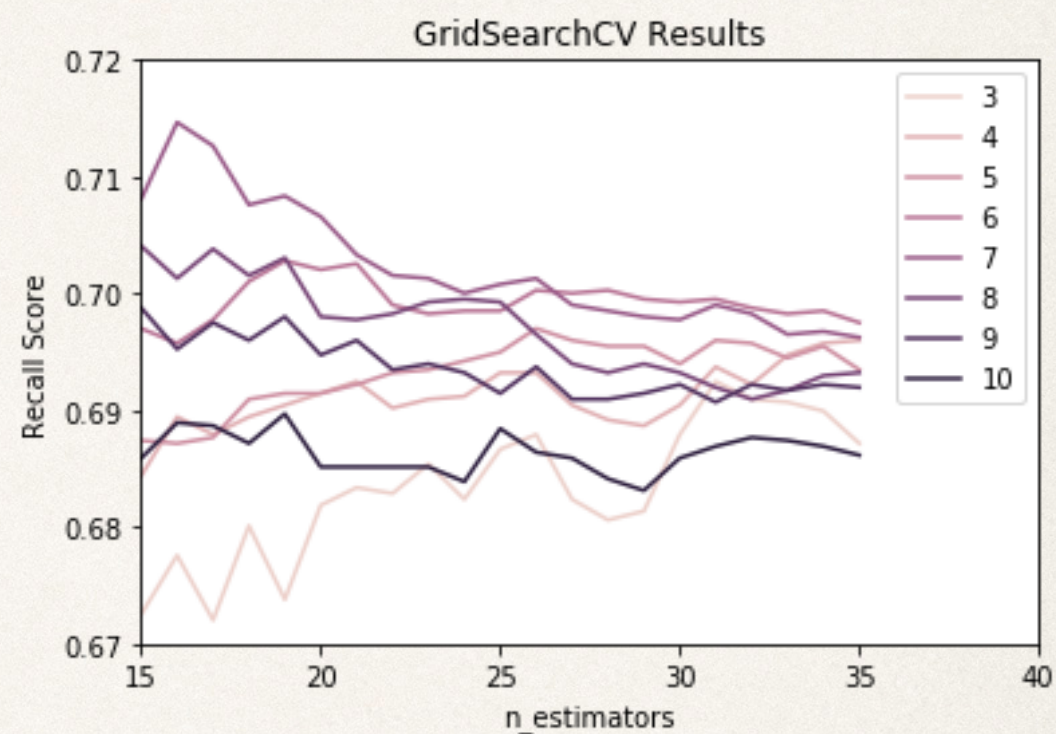
- ❖ SciKit Learn *GridSearchCV*
- ❖ Maximize recall score
- ❖ `max_depth` & `n_estimators`



# Model Tuning



	feature	score
14	PIF_2	0.274671
8	PAY_0	0.183614
9	PAY_2	0.175227
10	PAY_3	0.065750
11	PAY_4	0.062222
15	PIF_3	0.049585



	feature	score
14	PIF_2	0.273711
8	PAY_0	0.225856
9	PAY_2	0.108314
10	PAY_3	0.076283
11	PAY_4	0.067449
0	LIMIT_BAL	0.055806
12	PAY_5	0.045125
15	PIF_3	0.044115



# Analysis Results

---

- ❖ Payment features best indicator of default
- ❖ Demographic features not as clear
  - ❖ Married clients are more likely to default than single
  - ❖ Males more likely to default than females
  - ❖ Default is highest for customers age 20-25