

# Predicting Credit Card Fraud with R



# Background

- **Instructor:** John Garcia
  - I teach advanced data analytics at the University of North Texas
- **Project:** We will use R to fit three classification model to a highly imbalanced dataset:



1

**Decision Tree:** Uses a tree-like model of decisions to arrive at a classification prediction.

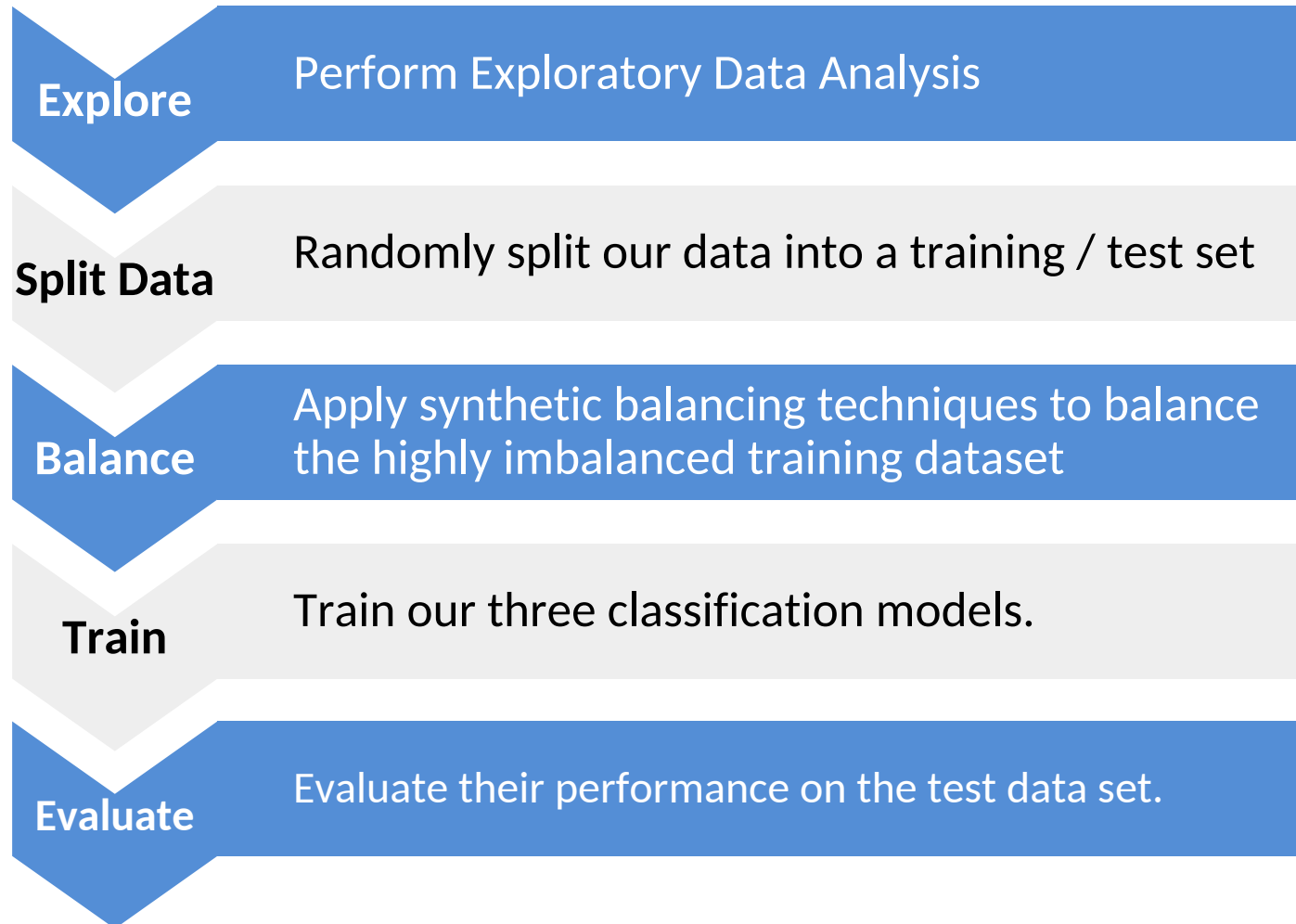
2

**Naïve Bayes classifier:** Uses Bayes' theorem to use probability to arrive at a classification prediction.

3

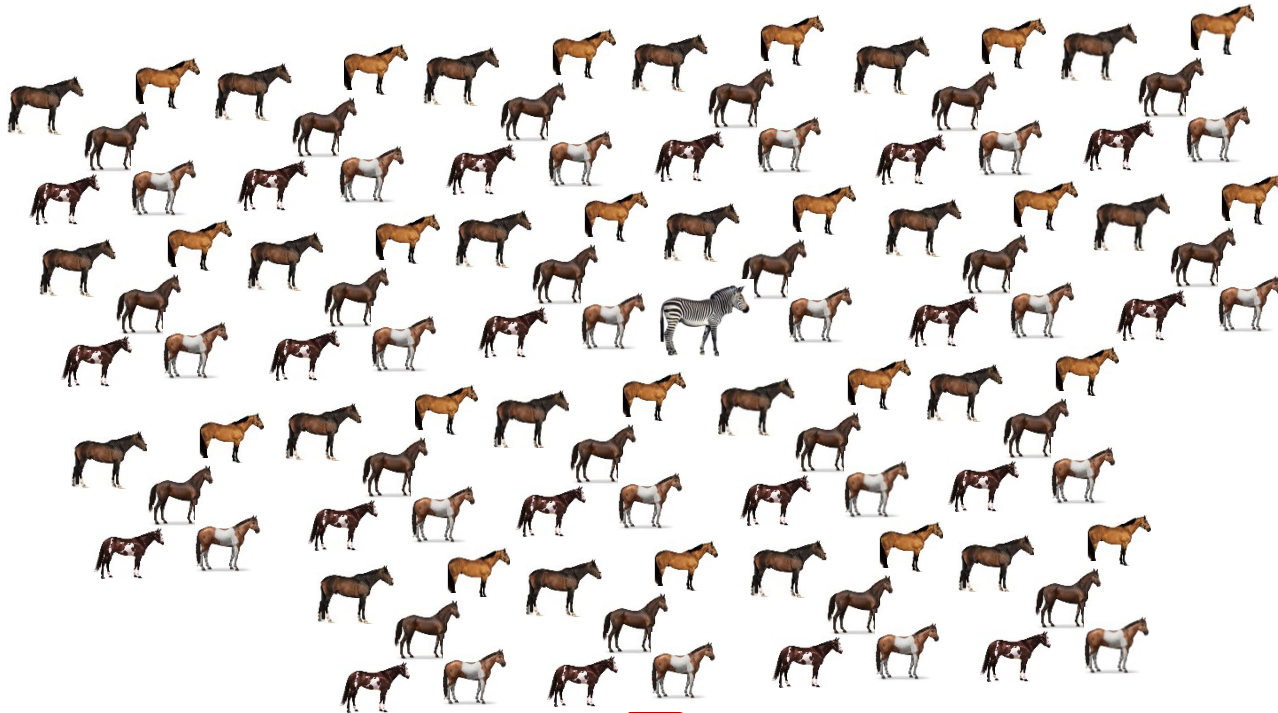
**Linear Discriminant Analysis:** Finds a linear combination of features that is used to separate the classes.

# We will use R to.....



# Why Balance our Dataset?

Dataset: 99 Horses.....1 Zebra



Prediction Algorithm  
All Horses

100 Horses

99% Accuracy



## Predicting Credit Card Fraud with R



### 1 Introduction

07:49



### 2 Exploratory Data Analysis

08:48



### 3 Create Training / Test Datasets

06:29



### 4 Generate Synthetic Samples

08:36

### 5 Train Classifiers on Original Imbalanced Dataset

09:17

### 6 Train Classifiers on SMOTE Balanced Dataset

06:07

### 7 Train Classifiers on ADASYN Balanced Dataset

04:45

### 8 Train Classifiers on DB SMOTE Balanced Dataset

05:56

### 9 Compare All The Trained Models

05:59



## Upload File:



Upload files to your cloud workspace

You can also use any upload service like [ufile.io](https://ufile.io)

## Download File:

Click on a file to download it

Show All

← Administrator

📁 Desktop



📁 Predicting Credit Card Fraud

📄 Predicting\_CC\_Fraud\_WorkingFile.Rmd



## Predicting Credit Card Fraud with R



05:20:48



Predicting\_CC\_Fraud\_WorkingFile.Rmd | creditcardFraud

```
29
30 *** Task 1.1: Import the dataset from Dropbox
31
32 Next, using the "read.csv" function, we will import the credit card fraud dataset and set
33 the class to a factor. This dataset is a subset of the dataset from sourced from
34 https://www.kaggle.com/mlg-ulb/creditcardfraud, which includes anonymized credit card
35 transactions.
36
37
38 #A. Load the dataset
39 creditcardFraud <- read.csv("Predicting Credit Card Fraud/creditcardFraud.csv")
40
41 #B. Change class to factor the as.factor function encodes the vector as a factor or
42 category
43 creditcardFraud$class<-as.factor(creditcardFraud$class)
44
45 *** Task 2: Explore The Data
46
47 * Now that we have downloaded the data we can start the training of the models, but it is
48 important that we first understand and explore our data as it helps us identify potential
49 data quality issues and it provides us the needed context to develop an appropriate
50 model.
51
52 * In this project, we will briefly explore the data and perform a high-level exploratory
53 data analysis (EDA) of the dataset
54
55
56 #A. Structure of the dataset
57
58
59 #B. Missing data?
60
61
62 #C. Check the imbalance in the dataset
63
64
65 #D. Compile histograms for each variable
66 par(mfrow = c(3,5)) #Change setting to view 3x5 charts
67 i <- 1
68 for (i in 1:30)
69 {hist((creditcardFraud[,i]), main = paste("Distribution of ",
70 colnames(creditcardFraud[i]), xlab = colnames(creditcardFraud[i]), col = "light blue")
71 }
72
73 #E. Compute the correlations among the variables
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
```

Environment

- Global Environ
- Data
  - creditcard

Files | Plots

39:3 | Chunk 3 | R Markdown | 00:02 / 08:48 | CC | 1x

- 100% +