

01

02

03

04

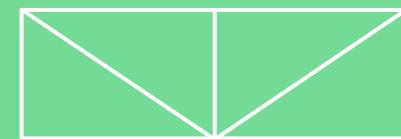
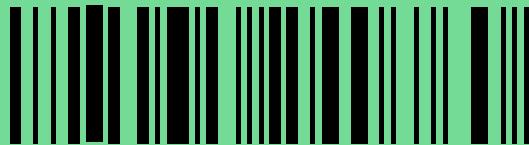
05

06

# GG3

## Good Games For Everyone

Jerome Sim, Joey Kang, Tay Yi Li, Adriel Chen



# TABLE OF CONTENTS

01

## BACKGROUND

Problem Statement,  
Big Picture, Goals

## PROCESSING

Web Scraping, Data  
Processing, EDA

02

03

## MODELLING

Model Analysis  
& Evaluation

## CONCLUSION

Recommendation,  
Going Forward

04

01

02

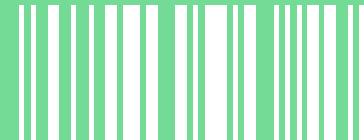
03

04

05

06

01



# BACKGROUND

About GGE and our data science goals

# ABOUT US

Good Games for Everyone is a PC building company, specialising in constructing high-performance gaming terminals for our customers. Adriel, Joey, and Yi Li make up their Data Science Division.



Div Head,  
DS  
**Adriel Chen**



Team Lead,  
DS  
**Joey Kang**

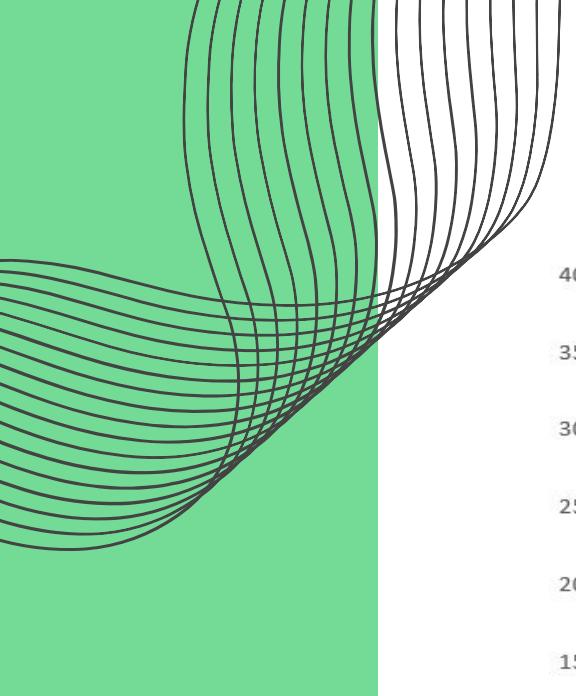


Team Lead,  
DS  
**Tay Yi Li**

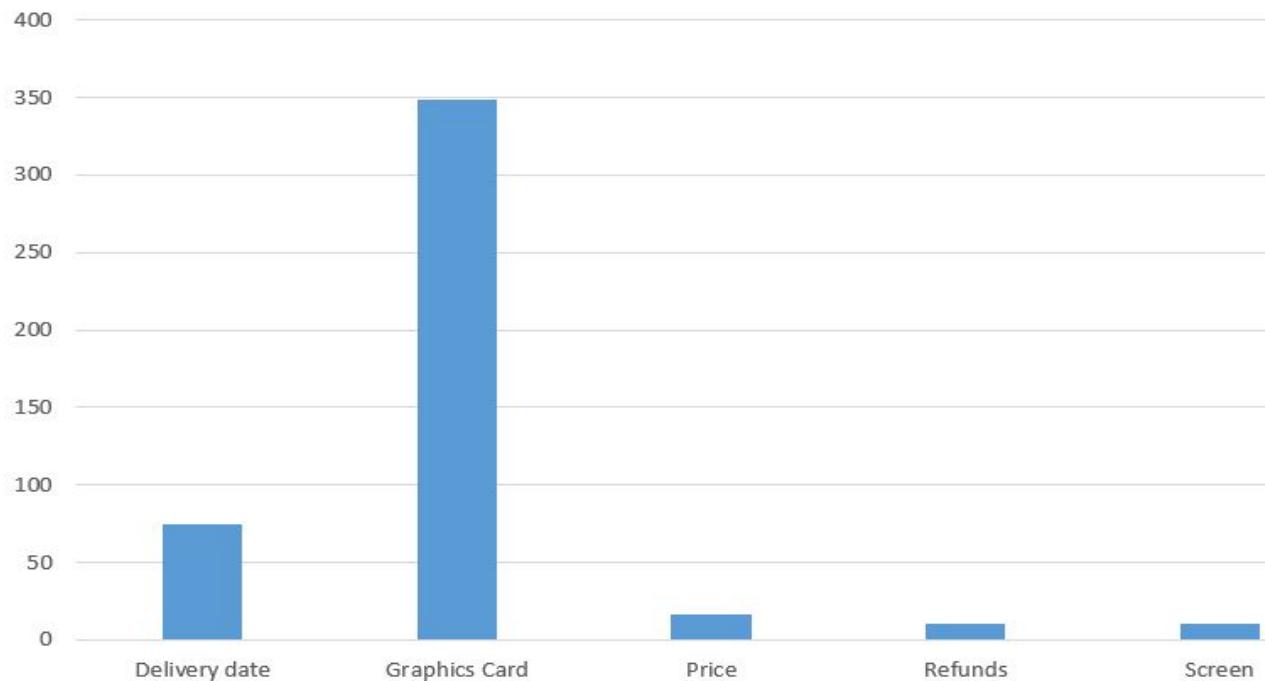


Chief  
Hygienist  
**Jerome Sim**



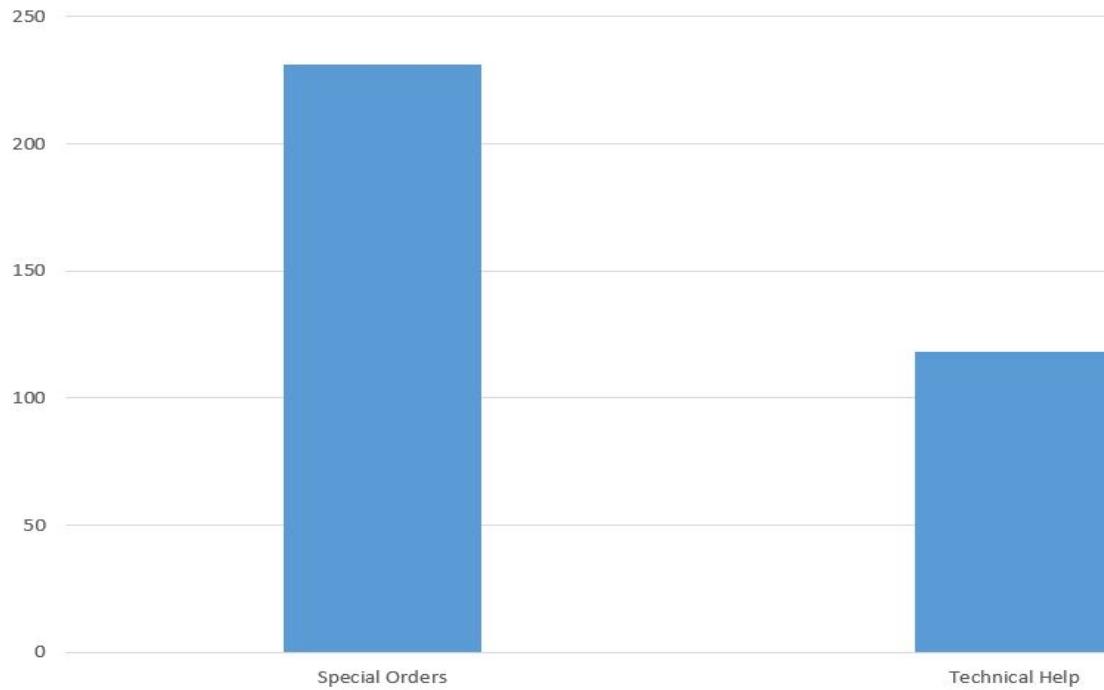


## Count of Customer Service Requests





### Proportion of Question Types (Graphics Card)

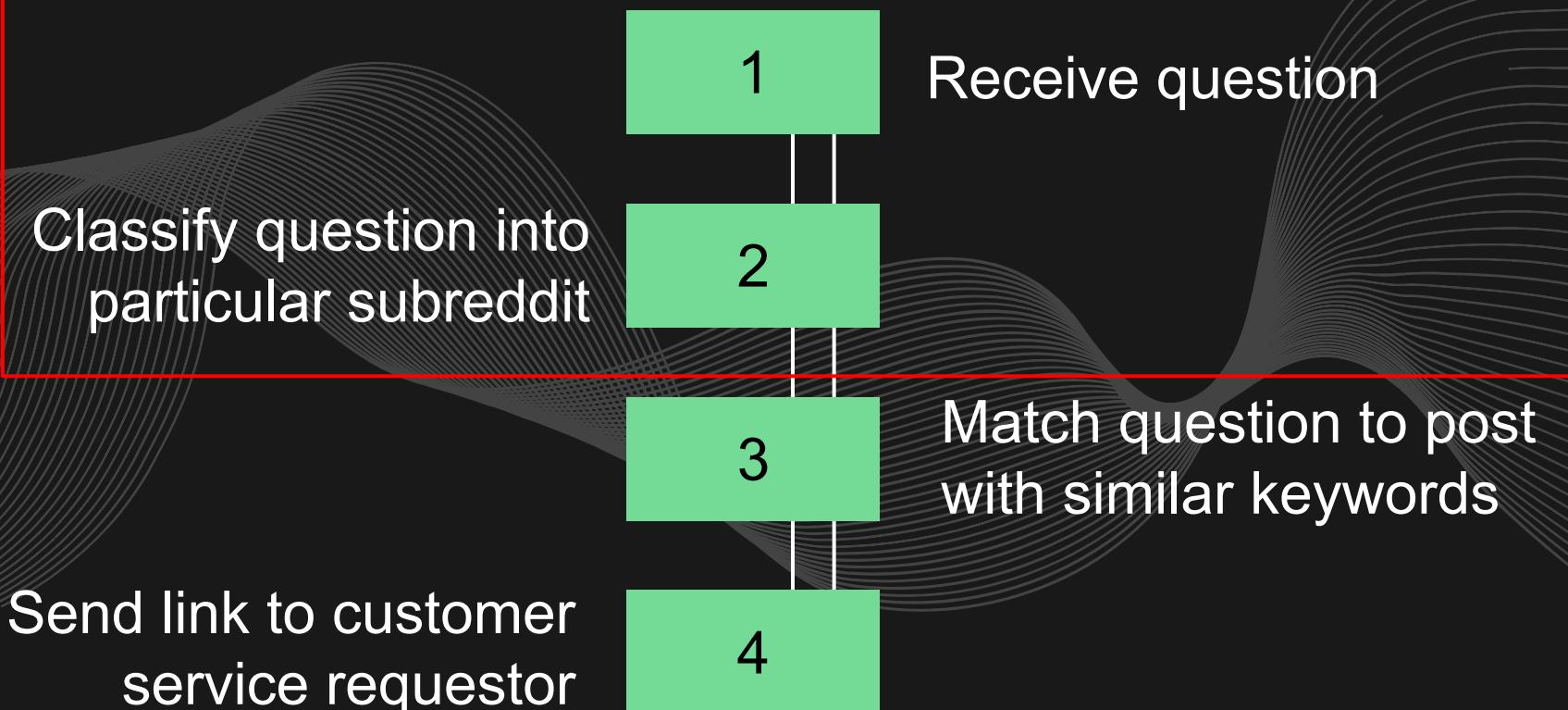


# Problem Statement

Management tasked the Data Science and Software Engineering teams to find a way to field after-hours requests for technical help regarding graphics cards.



# The Big Picture



# Goals

## Aim

To build a classifier that identifies keywords for the bot to determine relevant subreddits accurately.

## KPIs:

- Sensitivity
- Specificity
- Accuracy



01

02

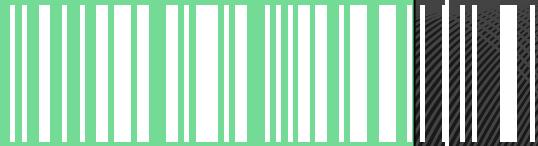
03

04

05

06

02



# DATA PROCESSING

**Web-scraping, Data  
Cleaning and Exploratory  
Data Analysis**

# WEB-SCRAPING

PushShift API to scrape from Reddit's r/Nvidia and r/AMD.

Data:

- 10000 rows each
- Variables: Title, Post's text, Author, Subreddit, Time Stamp, Number of Comments, Number of Crossposts, Post's Score, Total Awards Received by Post.



# DATA CLEANING

- Kept useful columns
- Removed duplicates, deleted posts, and posts containing photos/videos
- Combined title and Post's text columns
- Concatenated Nvidia and AMD dataframes



# PRE-PROCESSING TEXT

Removing stopwords  
from nltk package

1

Removing Links,  
Punctuations

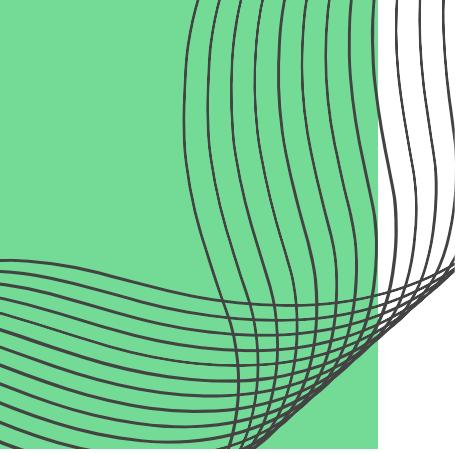
2

Removing additional  
stopwords identified  
through EDA

3

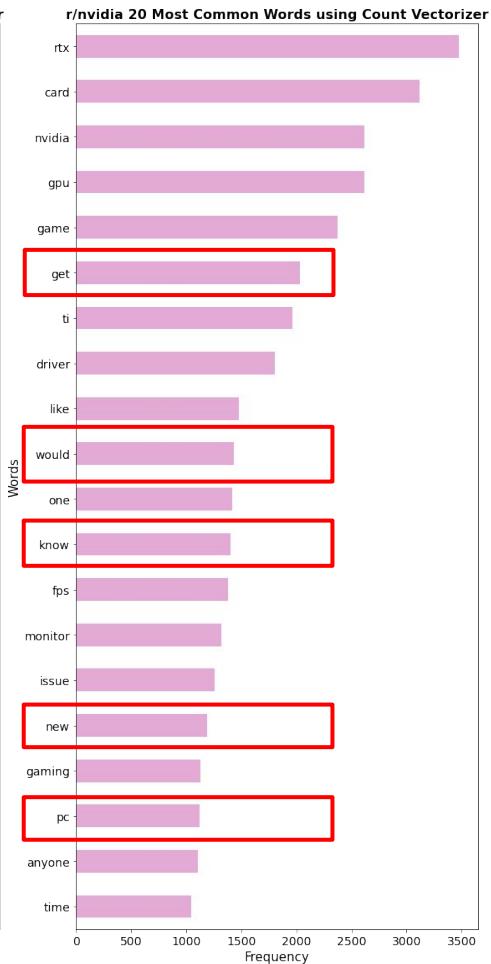
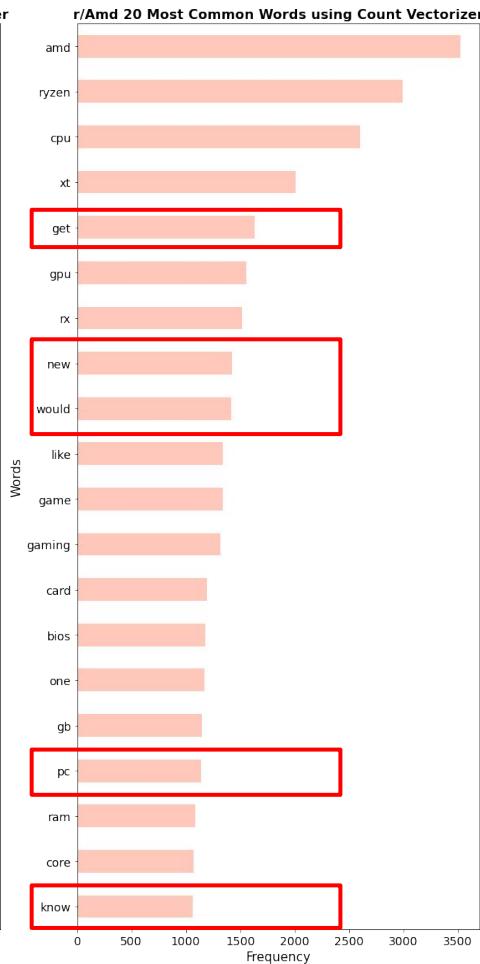
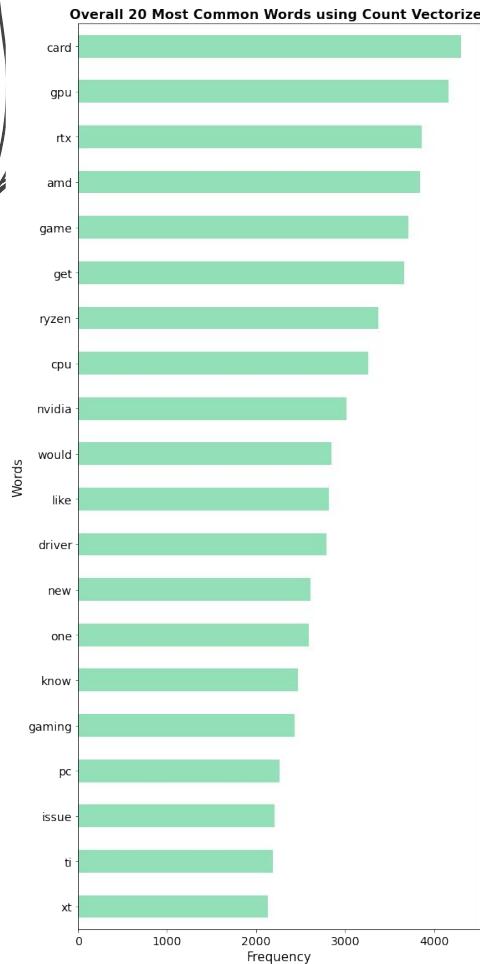
Lemmatization

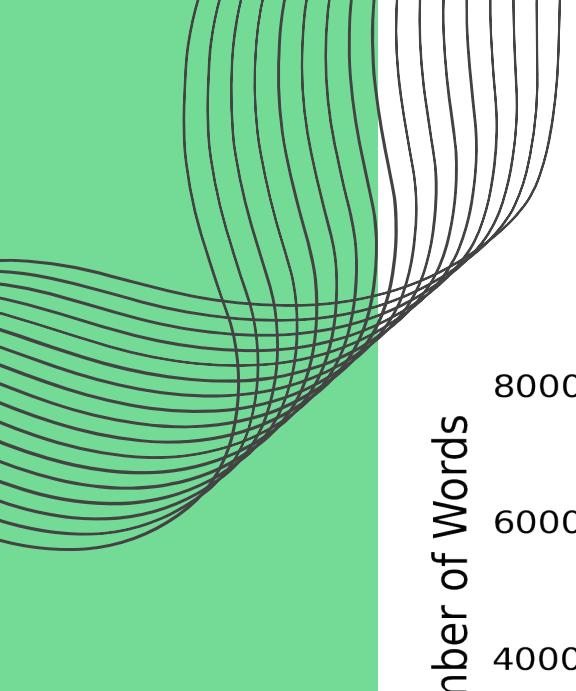
4



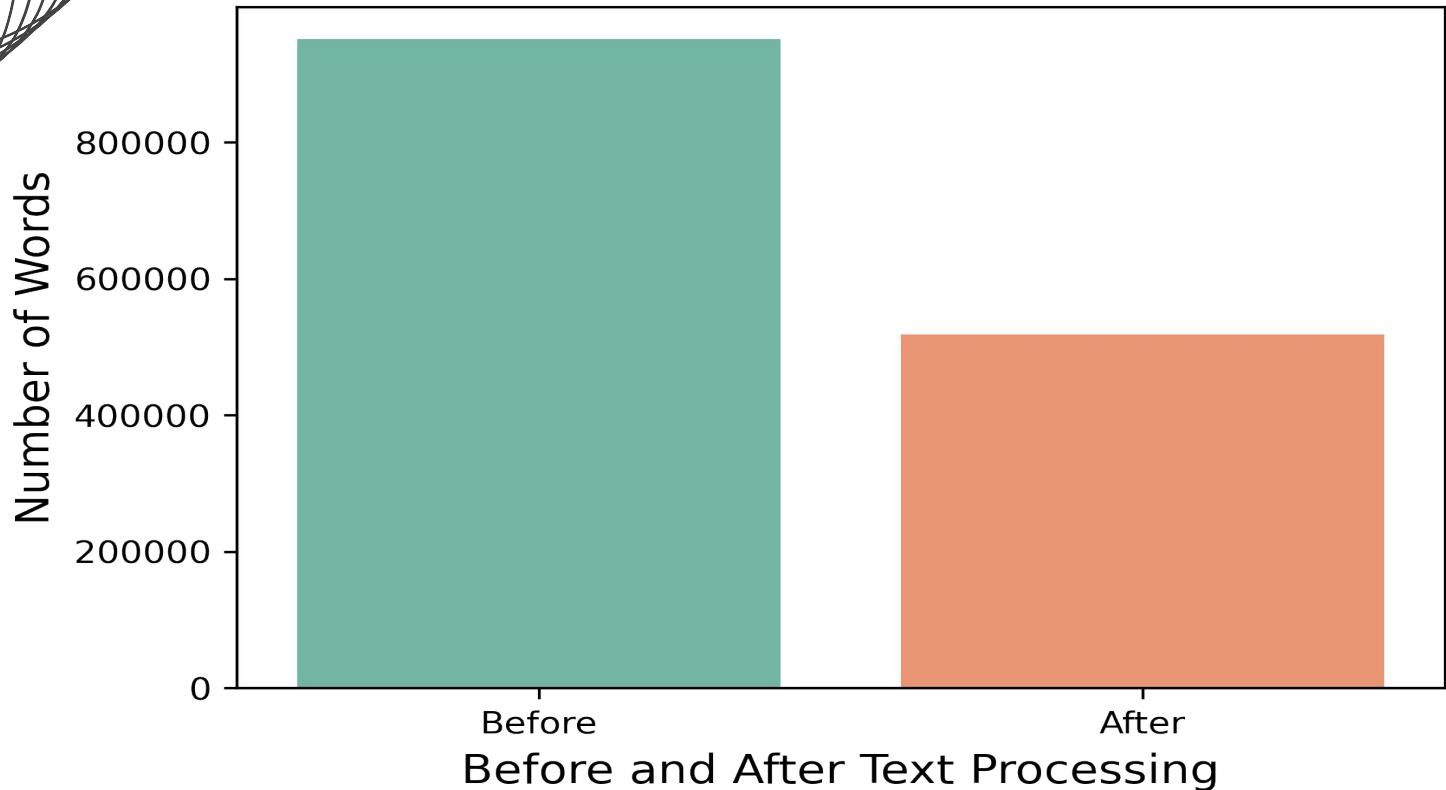
Removed generic words that appear in both subreddits:

- get
- new
- would
- pc
- know

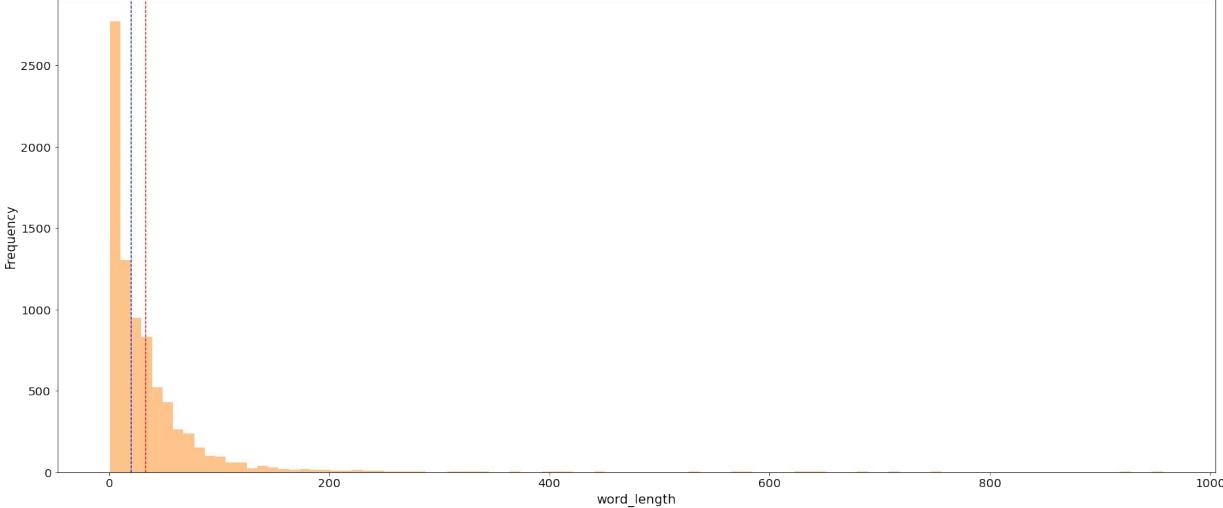
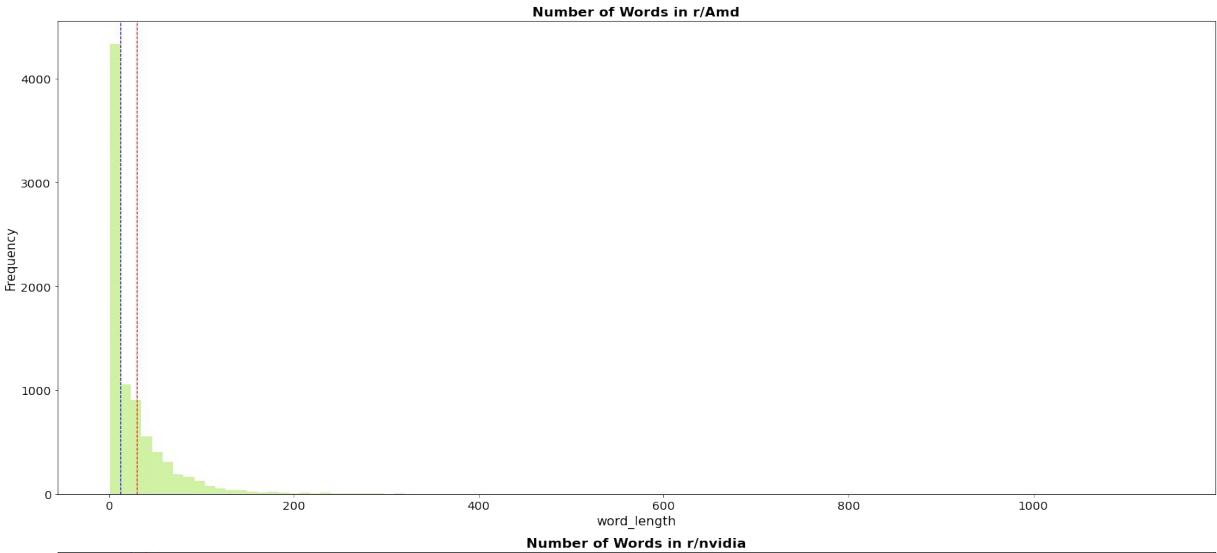




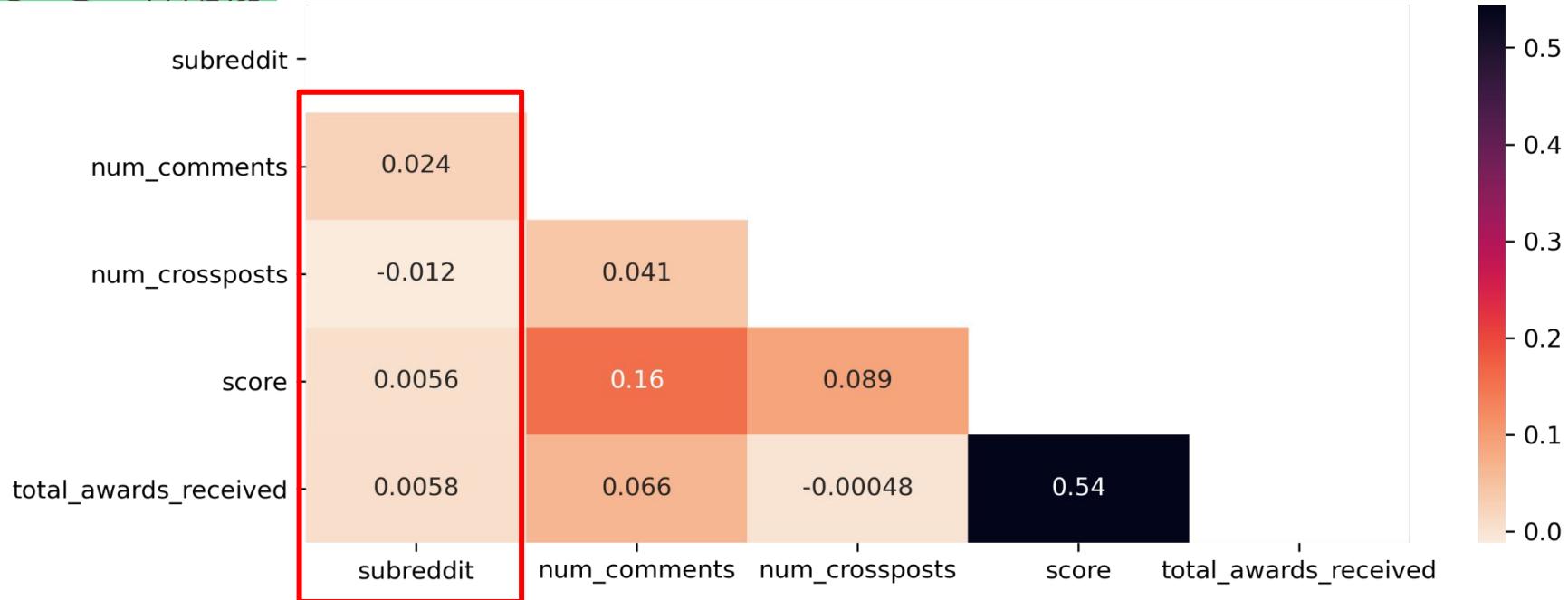
## Number of Words After Text Processing (almost 50% drop)



- Distribution of word length is **positively skewed** in both subreddits
- Average word length for Amd (31) and Nvidia (33) is comparable



# No correlation between any of the numerical columns with target variable (i.e. subreddit)



01

02

03

04

05

06



03

# MODELLING

**Model evaluation and analysis**

# MODELS EVALUATION

## Multinomial Naive Bayes

- 1. Test Accuracy
- 2. AUC
- 3. Sensitivity & Specificity

1

2

3

## Decision Tree

- 1. Test Accuracy
- 2. AUC
- 3. Sensitivity & Specificity

## Logistic Regression

- 1. Test Accuracy
- 2. AUC
- 3. Sensitivity & Specificity

# 1. Decision Tree

TEST ACCURACY



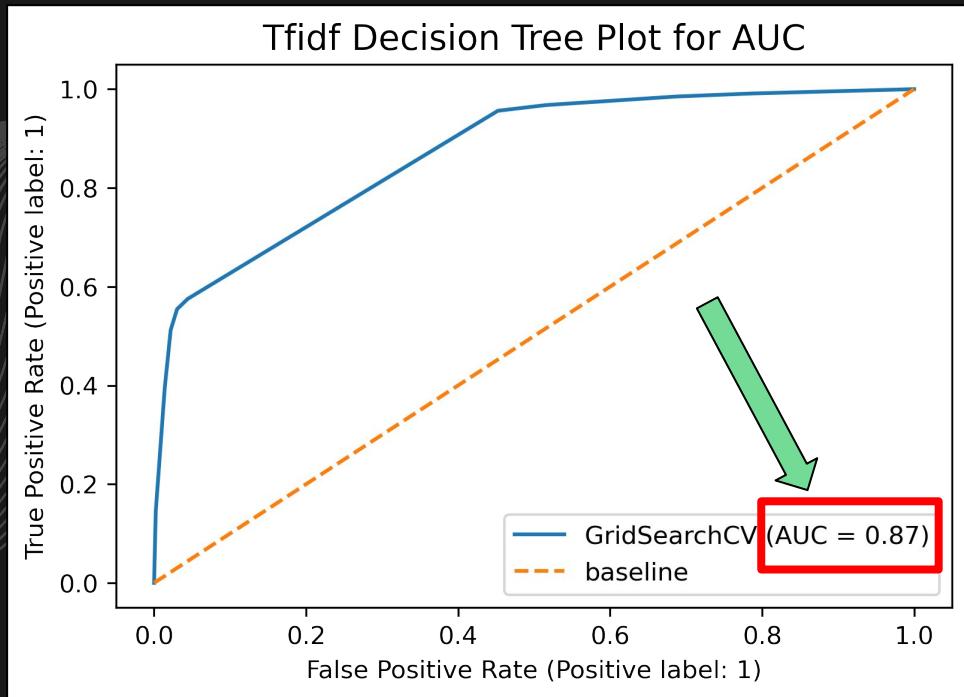
$(\text{True Positives} + \text{True Negatives}) / \text{Total Predictions}$



Overall ability to  
classify posts into  
the correct  
subreddit

# 1. Decision Tree

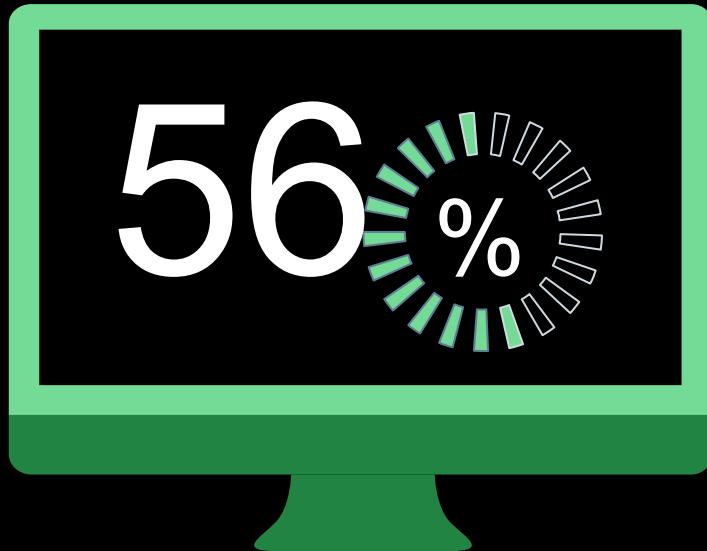
## Area Under Curve



Good Model  
↓  
Larger AUC  
↓  
Better at  
distinguishing  
between Positives  
& Negatives

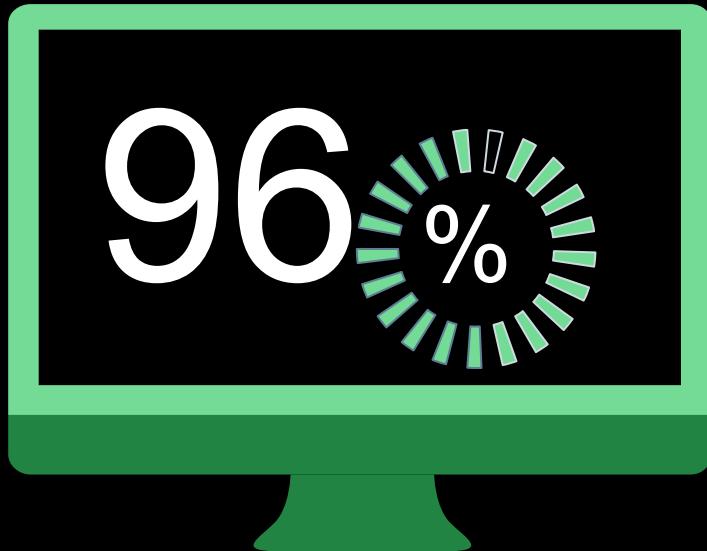
# 1. Decision Tree

SENSITIVITY



**Ratio of True Positives to All  
Actual Positives**  
Ability to correctly classify r/AMD posts

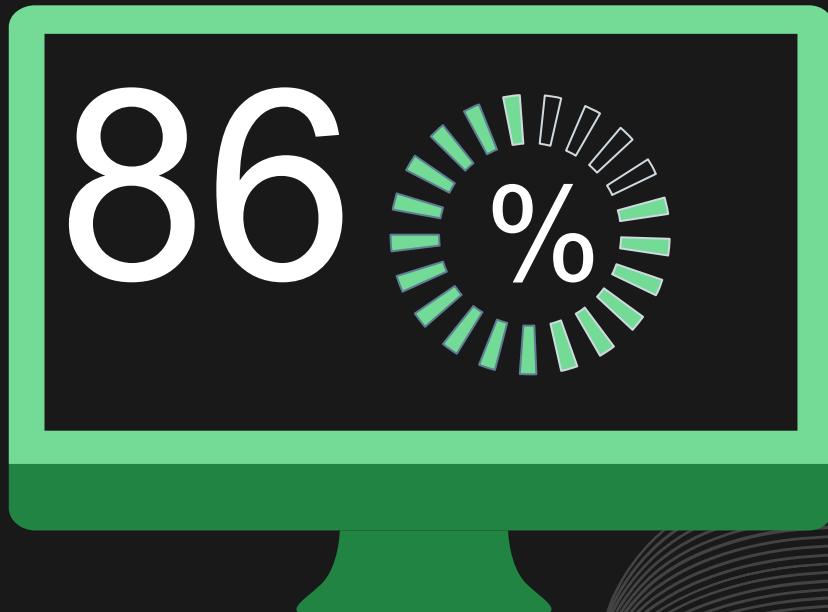
SPECIFICITY



**Ratio of True Negatives to All  
Actual Negatives**  
Ability to correctly classify r/NVIDIA posts

## 2. Multinomial Naive Bayes

TEST ACCURACY



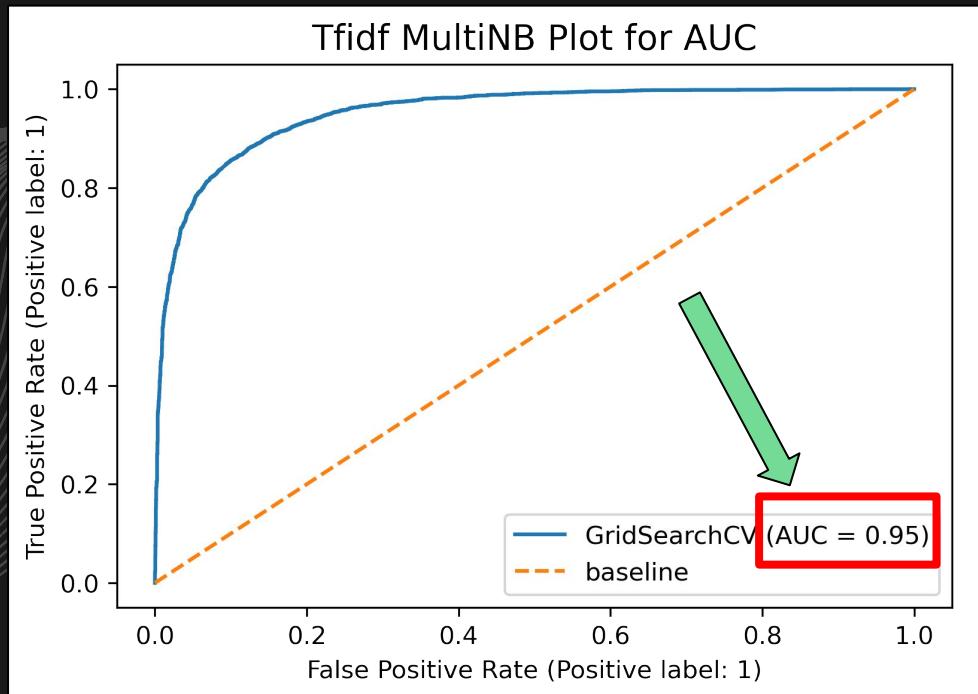
$(\text{True Positives} + \text{True Negatives}) / \text{Total Predictions}$



Overall ability to  
classify posts into  
the correct  
subreddit

# 2. Multinomial Naive Bayes

## Area Under Curve



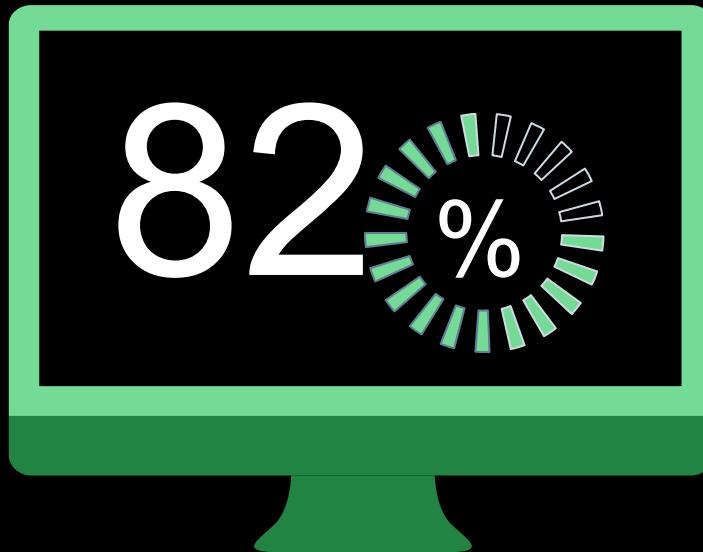
Good Model

Larger AUC

Better at  
distinguishing  
between Positives  
& Negatives

## 2. Multinomial Naive Bayes

SENSITIVITY



**Ratio of True Positives to All  
Actual Positives**  
Ability to correctly classify r/AMD posts

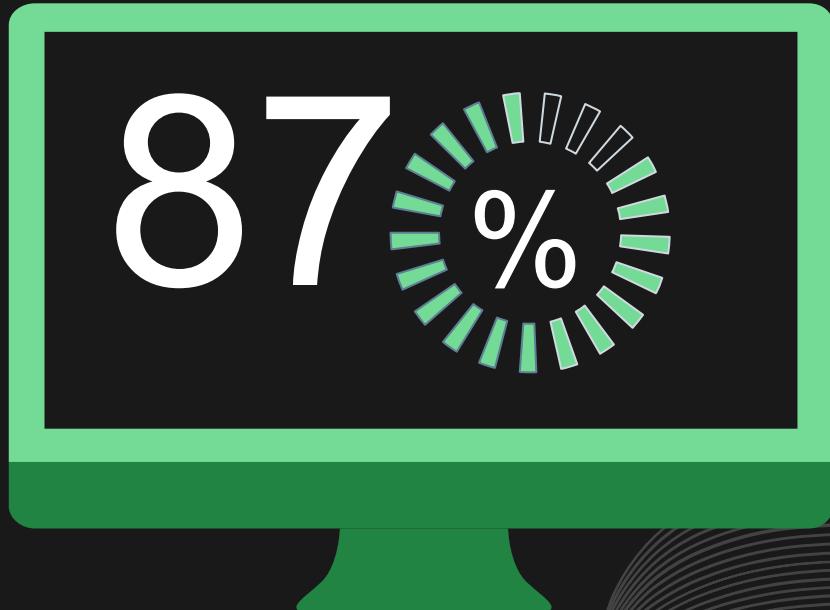
SPECIFICITY



**Ratio of True Negatives to All  
Actual Negatives**  
Ability to correctly classify r/NVIDIA posts

# 3. Logistic Regression

TEST ACCURACY



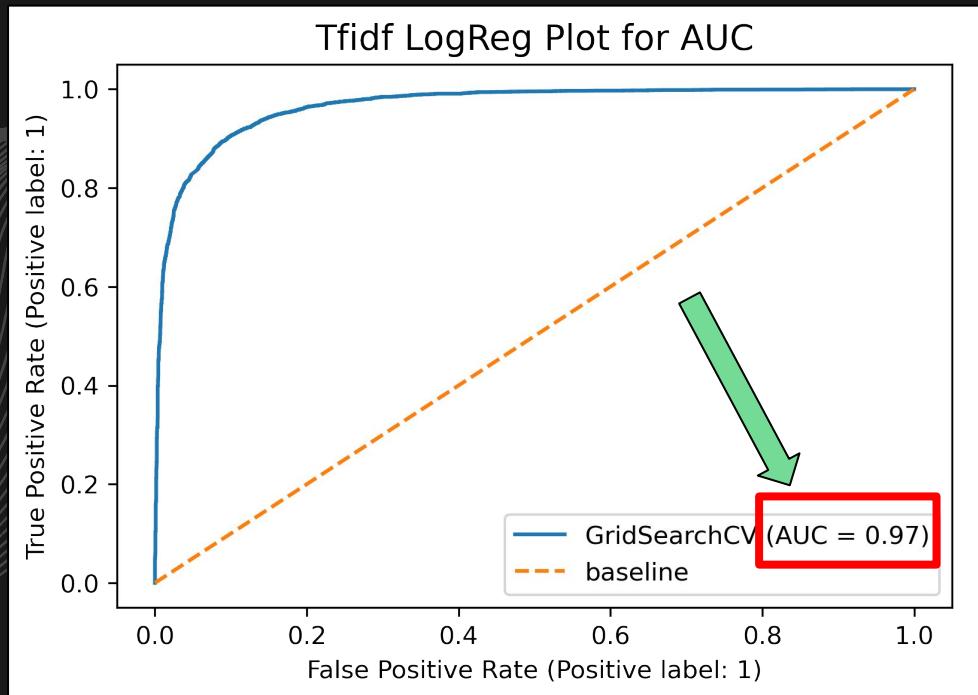
$(\text{True Positives} + \text{True Negatives}) / \text{Total Predictions}$



Overall ability to  
classify posts into  
the correct  
subreddit

# 3. Logistic Regression

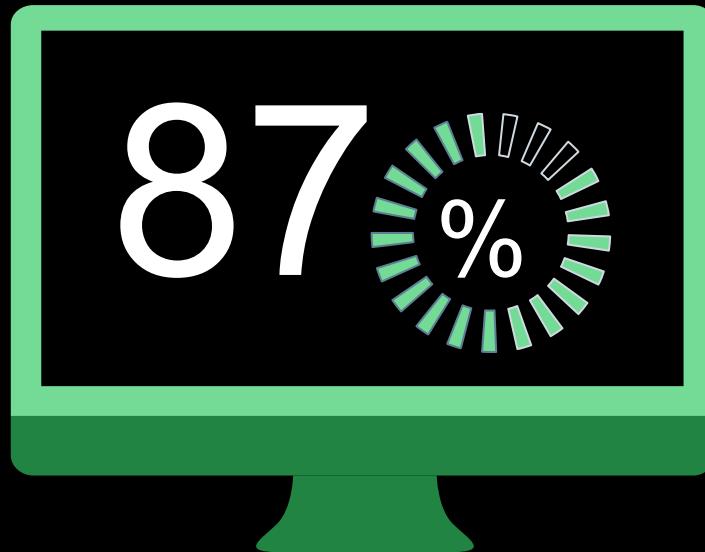
## Area Under Curve



Good Model  
↓  
Larger AUC  
↓  
Better at  
distinguishing  
between Positives  
& Negatives

# 3. Logistic Regression

SENSITIVITY



**Ratio of True Positives to All  
Actual Positives**  
Ability to correctly classify r/AMD posts

SPECIFICITY



**Ratio of True Negatives to All  
Actual Negatives**  
Ability to correctly classify r/NVIDIA posts

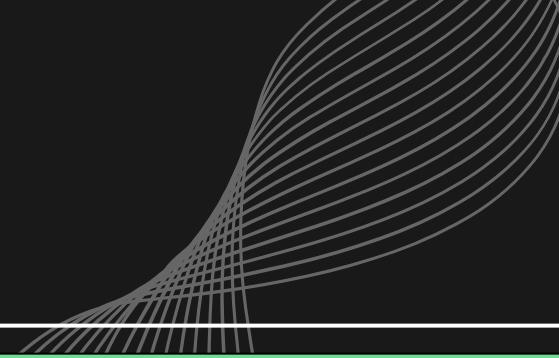
# Models Summary

|                         | TEST ACCURACY | AUC | Sensitivity   Specificity |
|-------------------------|---------------|-----|---------------------------|
| Decision Tree           | 76%           | 87% | 56%   96%                 |
| Multinomial Naive Bayes | 86%           | 95% | 82%   90%                 |
| Logistic Regression     | 87%           | 97% | 87%   88%                 |

Highest!

Highest!

Most balanced and high overall!



**LOGISTIC REGRESSION  
IS THE WAY TO GO**

— CHIEF HYGIENIST



# 04

## CONCLUSION & RECOMMENDATIONS

# 4 STAGE PROCESS



## STAGE 1 & 2 - RECEIVE QN & CLASSIFY

To build a classifier that identifies keywords for the bot to direct customers to subreddits accurately.



## STAGE 3 & 4 – MATCH QN AND SEND LINK

We will be identifying specific posts with solutions and assigning into a dictionary for the software team to direct customers to in the future

# Good enough for Beta Test

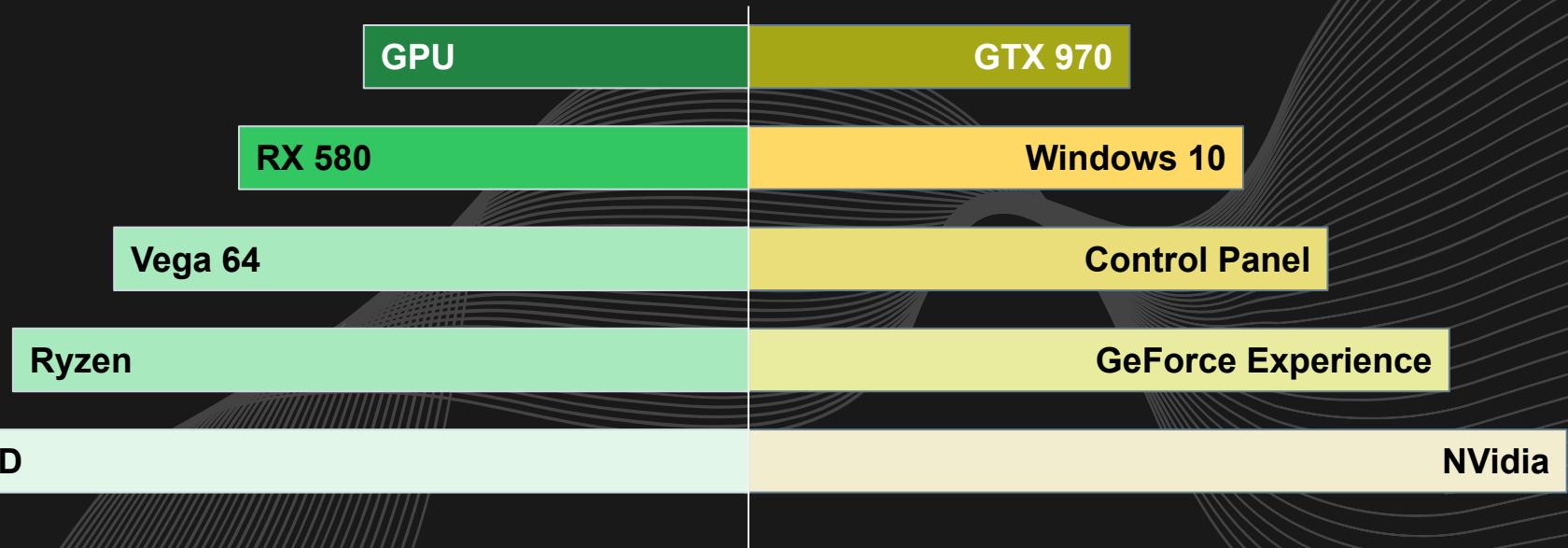
Software engineering can proceed with using our logistic regression model



87  
%

# TOP PHRASES FOR CLASSIFICATION

AMD      NVIDIA



# GOING FORWARD

## RECOMMENDATION

Software team to get Customer ID and Product ID for automatic brand & model imputation



## IMPROVE ACCURACY

Test with more models and tune more hyperparameters

## ONE BOT TO RULE THEM ALL

Complete AI recommendation solution for all customer inquiries



## ADDITIONAL DATA POINTS

Scrape additional solutions from other NVidia and AMD forums



THANK YOU

Any questions?

