

amazon Reviews

Curtis, Joe, Nate, and Ryan

Overview

- ▶ Project Objective
 - ▶ Outline
- ▶ Methods and Tools
 - ▶ AWS
 - ▶ R
 - ▶ Python
- ▶ Build Model
- ▶ Descriptive Analysis

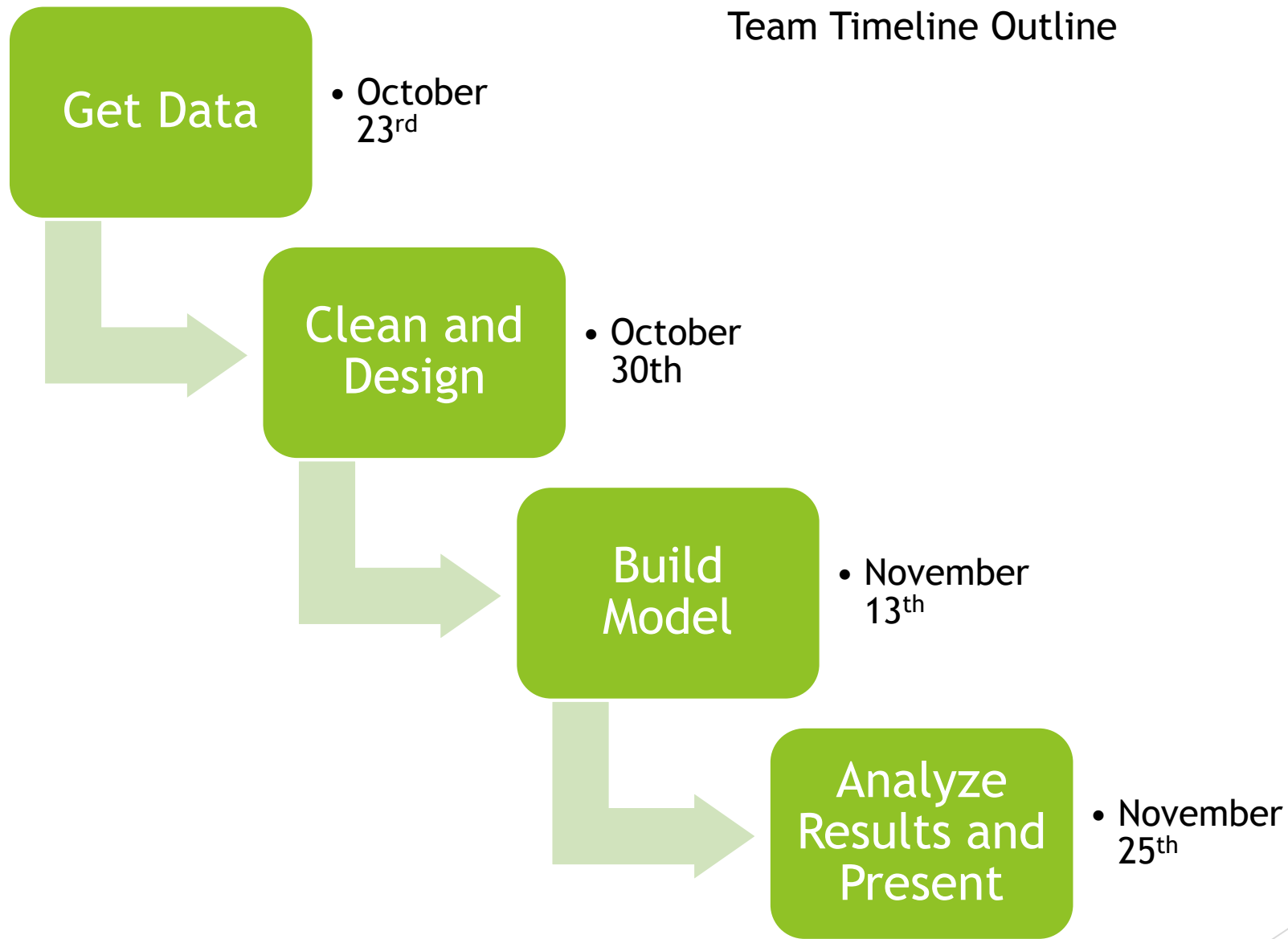
Project Goal

- Summary :

This project will entail web scraping, text mining, and predictive models with the objective of predicting "Review" (y/n) and/or the rating (number of stars). This approach seeks to help sellers target the reviewers most likely to review their product with a high rating, which will also be seen as helpful to other shoppers.

- Key Idea: Want to find text that leads to a good review.

Team Timeline Outline



Team Member	Responsibilities
Ryan	<ul style="list-style-type: none">• Project Manager tasks• Utilize R to pull in Data• AWS• Support Role
Curtis & Nate	<ul style="list-style-type: none">• Utilize Python to pull data• Tasks that require python• Support Role
Joe	<ul style="list-style-type: none">• Model Developer• Interpret Results• Support Role

Configuring AWS

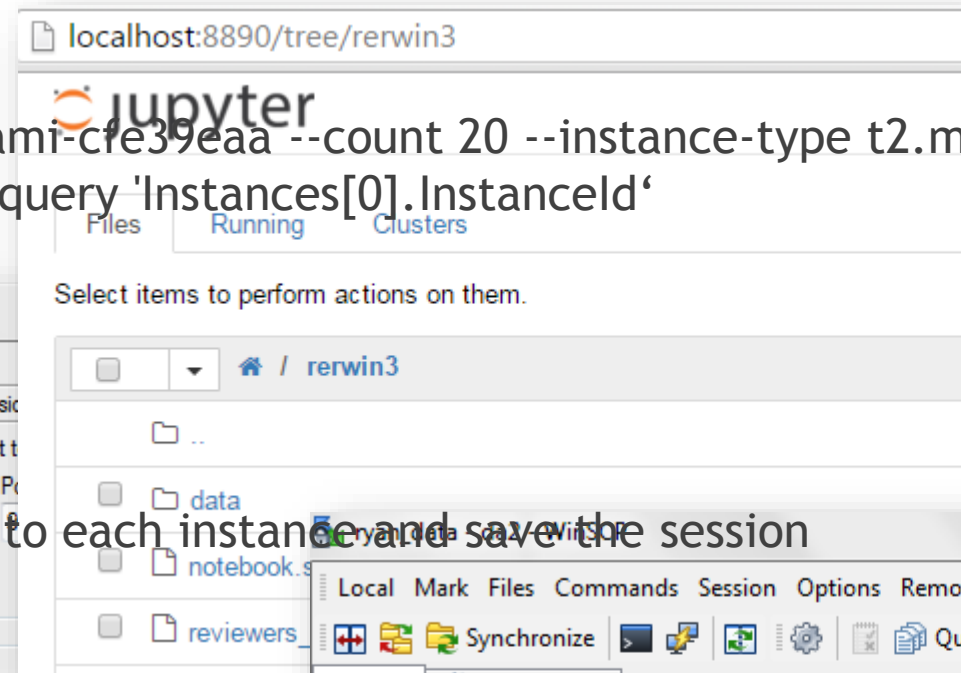
```
da2.eecs.utk.edu - PuTTY
Using username "rerwin3".
Authenticating with public key "rsa-key-20150824"
Welcome to Ubuntu 14.04.3 LTS (GNU/Linux 3.10.0-229.
* Documentation: https://help.ubuntu.com/
Last login: Sun Nov 22 19:03:35 2015 from c-50-142-2
3bd28df80bf
```

PuTTY Configuration

cd docker

./runnew.sh


Then use PuTTY interface to login to each instance and save the session

A screenshot of a Windows File Explorer window. It shows two side-by-side views of directories. The left pane shows the local path 'C:\...\Ryan\Dropbox\RACHEL_RYAN\2_Data\new_data' and contains a list of CSV files. The right pane shows the remote path '/home/rerwin3/AmazonProductReviews/ryan_data' and contains a similar list of CSV files. The files are named 1.csv, 501.csv, 1001.csv, 1501.csv, and 2001.csv. The 'Size' and 'Changed' columns are visible in the right pane.

Name	Size	Type
..		Parent directory
1.csv	11,671 KB	Microsoft Excel spreadsheet
501.csv	9,004 KB	Microsoft Excel spreadsheet
1001.csv	8,748 KB	Microsoft Excel spreadsheet
1501.csv	9,139 KB	Microsoft Excel spreadsheet
2001.csv	9,336 KB	Microsoft Excel spreadsheet

Name	Size	Changed
..		11/22/2015
1.csv	11,671 KB	11/21/2015
501.csv	9,004 KB	11/21/2015
1001.csv	8,748 KB	11/21/2015
1501.csv	9,139 KB	11/21/2015
2001.csv	9,336 KB	11/22/2015



ec2-54-152-30-129.compute-1.amazonaws.com:8787

File Edit Code View Plots Session Build Debug Tools Help

data_last_step.R x call_last_step.R x Untitled1 * x





Source on Save Run Source

```
5
6 # load the list of top reviewers -----
7 reviewers <- read.csv("reviewers_list.csv", stringsAsFactors = F)
8
9
10 # what page to scrape to for reviews
11 reviewers <- reviewers %>%
12   mutate(
13     page_num = ceiling(Reviews/10),
14     page_num = ifelse(page_num <= 10, page_num, 10)
15   )
16
17
18 # grab the url for each person I'm assigned
19 reviewers <- reviewers %>%
20   filter(page_num > 0)
21
22
23 # create the review links
24 review_links <- .create_links(reviewers$url_name,
25                               reviewers$page_num)
26
27
28 # call and store the results -----
29 setwd("~/data")
30
31 # where are my functions
```

Console ~/data/  

times[] <- end_loop

Environment History

  Import Dataset  

Global Environment





Data


review_links	487280 obs.
reviewers	9986 obs. d

Values


end_loop	
file_name	"2001.csv"


Files Plots Packages Help Viewers


 New Folder  Upload  Delete 


 Home


Name


 .Rprofile


 call_last_step.R


 data


 data_last_step.R

 notebook.sh

 products.csv

 R

 reviewers_list.csv

 start.sh

Python Extraction

- ▶ What to Extract
 - ▶ Reviewer ID
 - ▶ Item Price
 - ▶ URL
 - ▶ Item Description
 - ▶ Rating
 - ▶ Summary of Review
 - ▶ Item Category

Troubleshooting

- ▶ Separate Fields
- ▶ Comments
- ▶ VPN
- ▶ Run time

≈ 300,000 Reviews

► Code

1. Class

2. Driver

3. Function

Snip it- Driver

```
#Driver
with open('/home/jhughe39/AmazonProductReviews/python_implementation/reviews4.csv', 'w', newline='') as csvfile:

    # header for csv file
    fn = ['reviewerName', 'reviewerUrl', 'itemNo', 'itemTitle', 'itemUrl', 'itemBrand', 'itemPrice',
          'itemCategory', 'dateReviewed', 'Rating', 'Summary', 'Description']

    # open csv file for writing
    reviewWriter = csv.DictWriter(csvfile, delimiter=',', fieldnames=fn)
    reviewWriter.writeheader()

    f = open("/home/jhughe39/AmazonProductReviews/python_implementation/ReviewURLs.txt", "r")

    for url in f:
        review = get_review(url)
        time.sleep(3)
        if review.itemNo != 'none' and review.reviewerUrl != 'none':
            reviewWriter.writerow(review.__dict__)
```



Scrape Data

Find Most
Common Words
(Cluster Analysis)

Build Model

- What words lead to a positive review

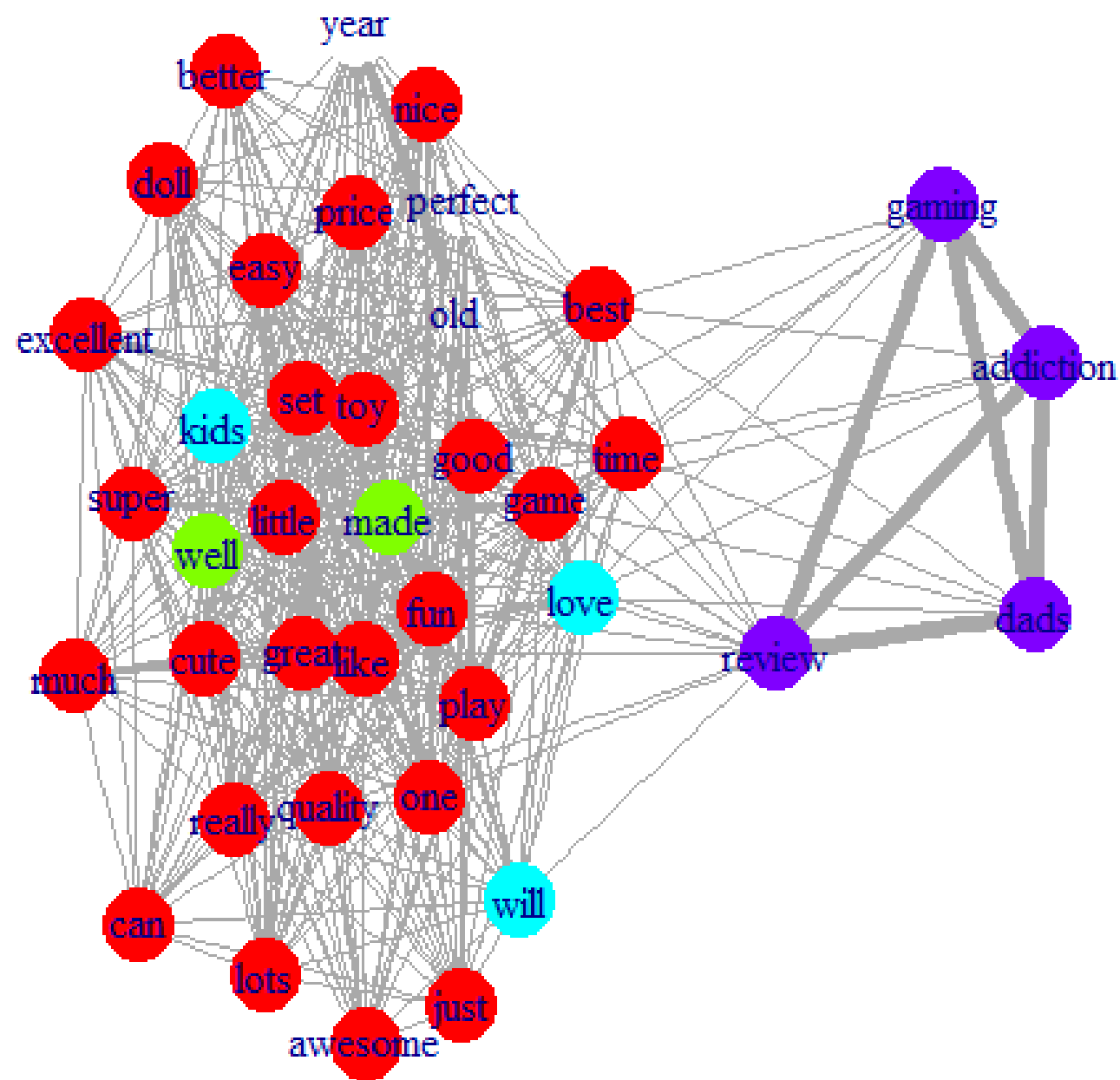
Random Forest Model
Response: 5 Star Rating


Top Item Categories

- ▶ Books
- ▶ Kindle Store
- ▶ Toys and Games
- ▶ Health and Beauty
- ▶ Electronics
- ▶ Amazon Video

Example - Toys and Games

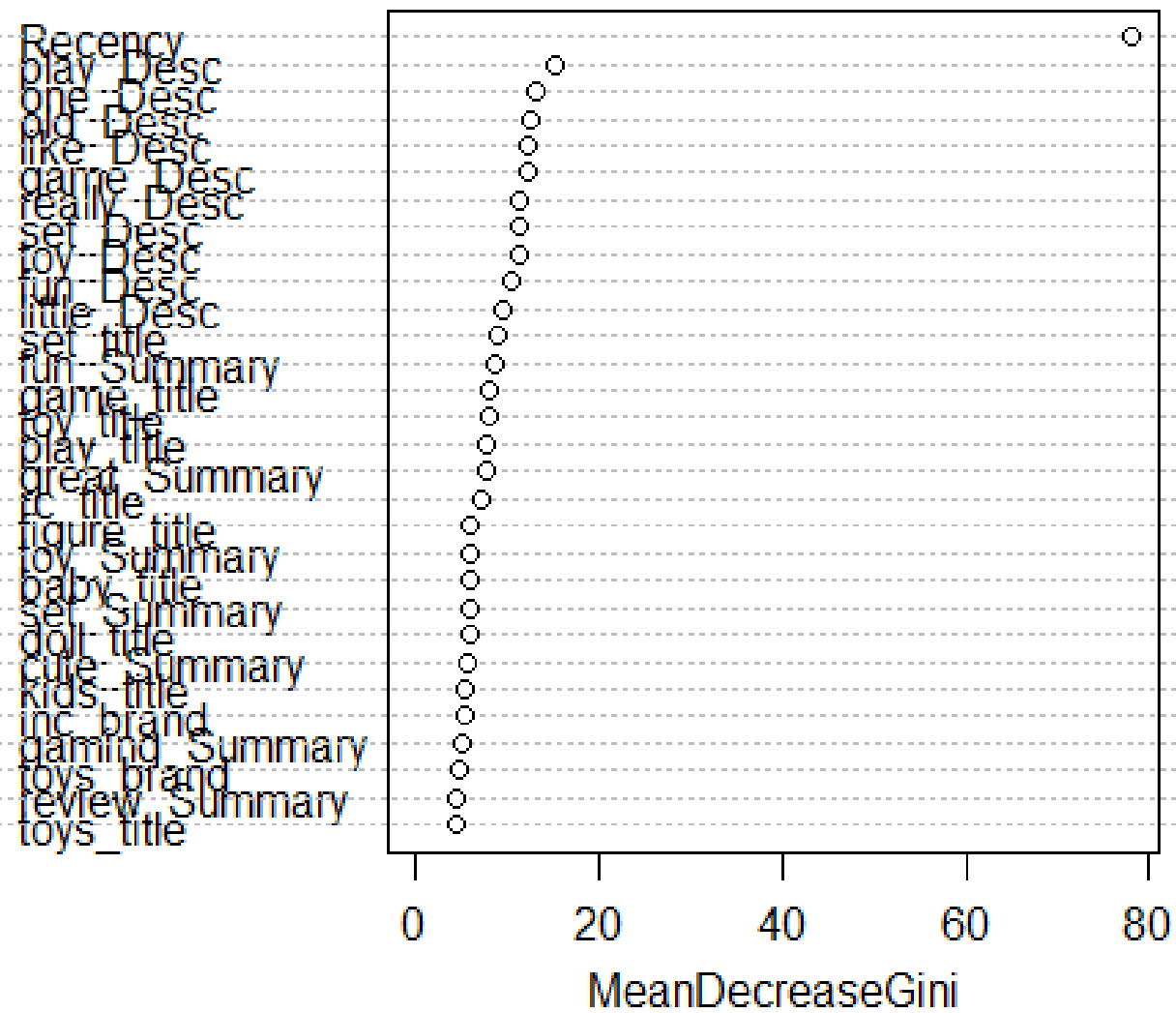
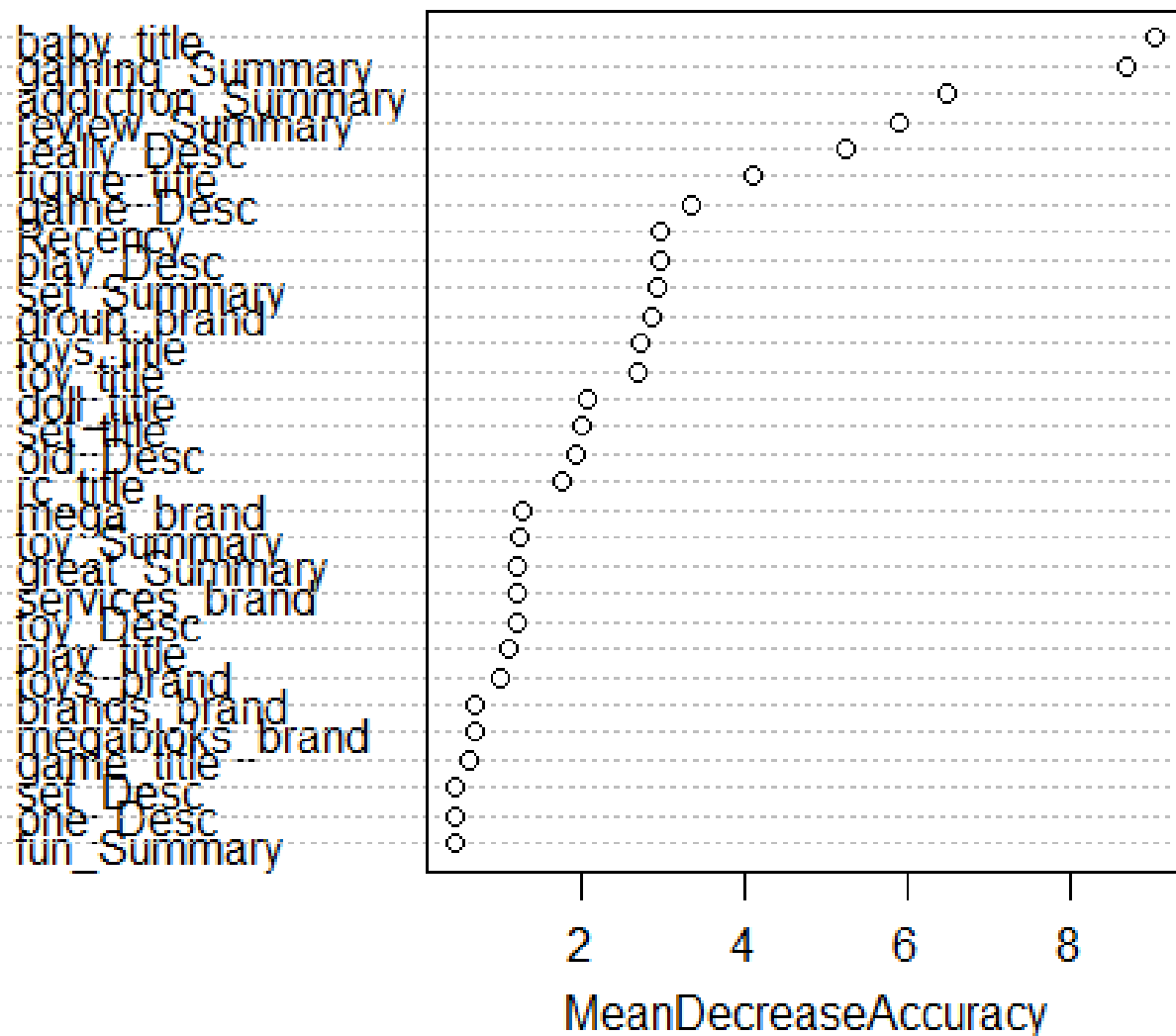
- ▶ Text analysis
 - ▶ Review Summary
 - ▶ Item Description
 - ▶ Brand Names
 - ▶ Product Names
- ▶ Response - 5 Star Rating





awesome
dadscute
best kids play
really doll quality
just old made much
lots t well time price
easy o excellent
s good year will love
nice can little perfect toy
one super great better
like
game review
gamingaddiction

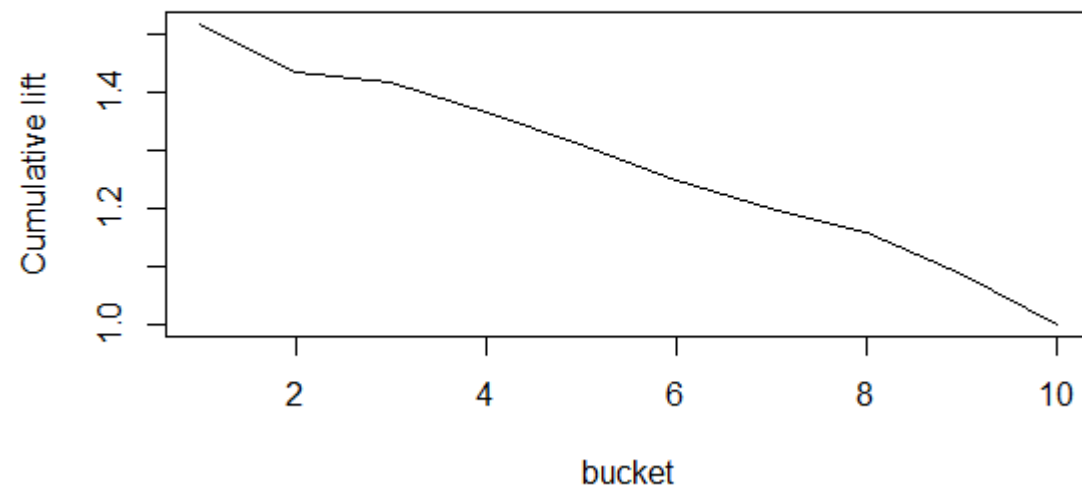
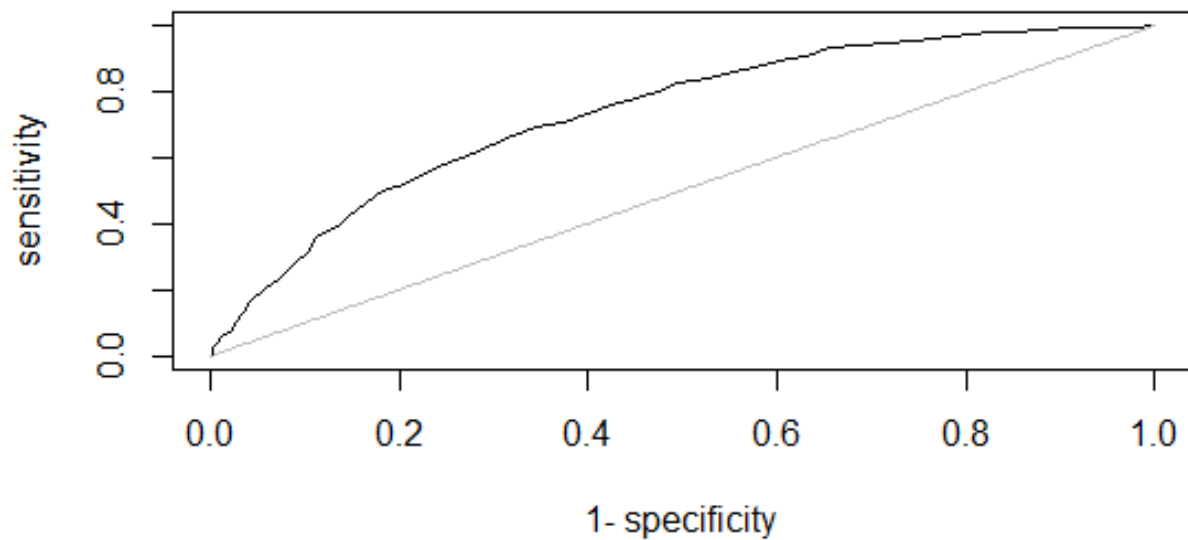
rFmodel2



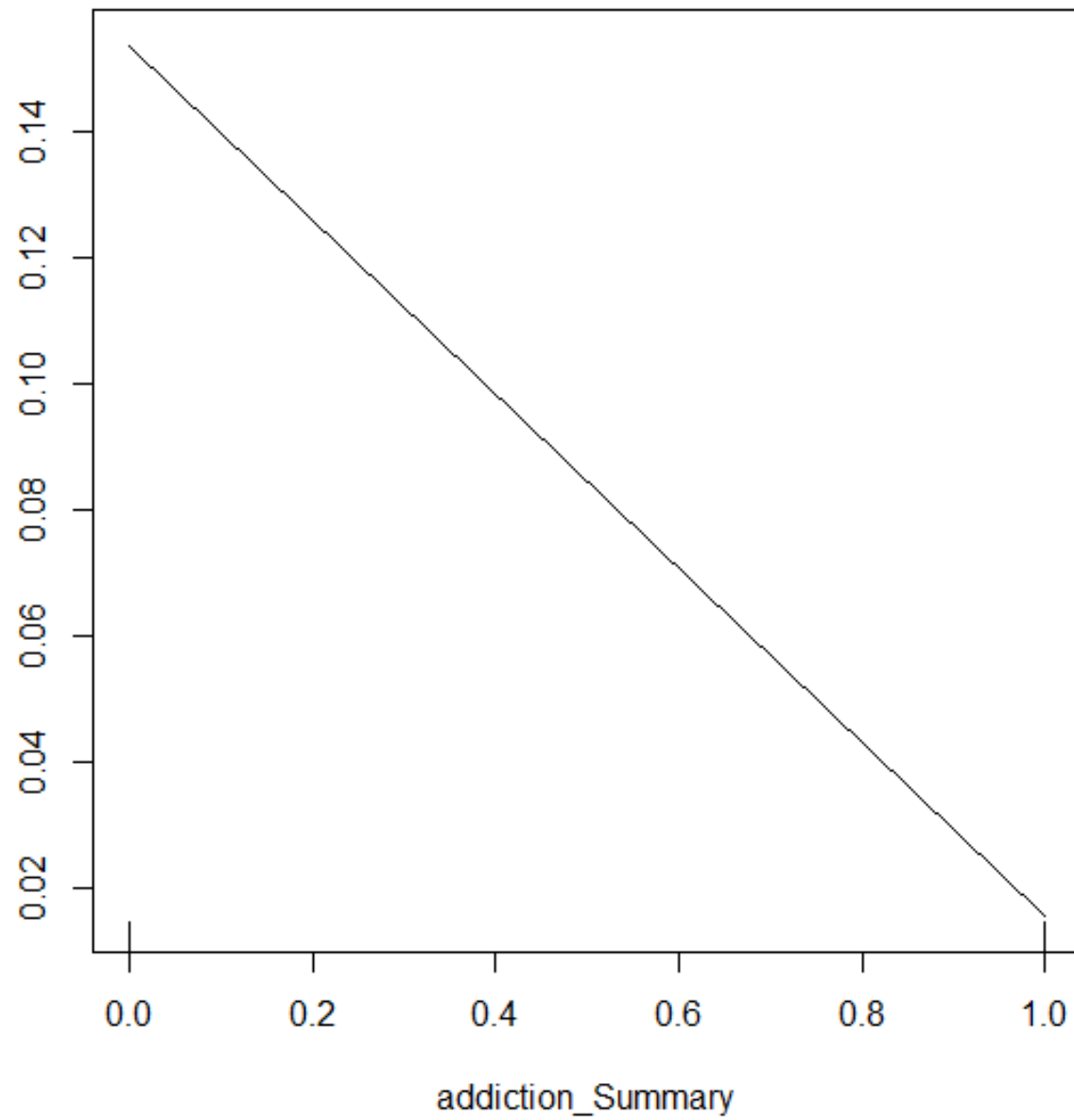
Model Performance

AUC - 0.75 - 0.80

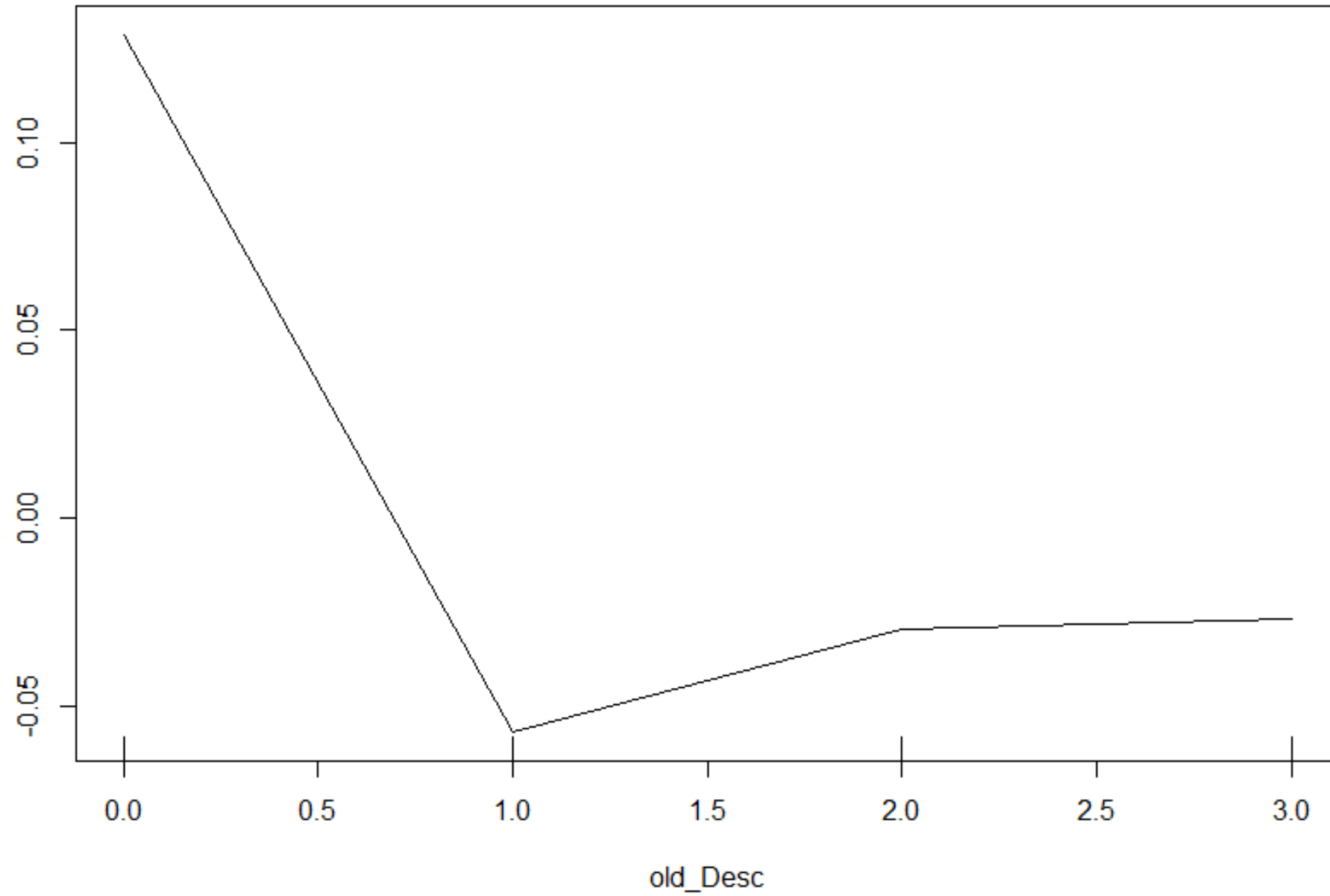
Lift Curve - 1.5



Partial Dependence on addiction_Summary



Partial Dependence on old_Desc



The background features abstract, overlapping green geometric shapes, primarily triangles and polygons, in various shades of green, creating a modern and dynamic visual effect.

End

Any questions?