

Use Cases in Big Data Software and Analytics

Vol. 1, Fall 2017

Bloomington, Indiana

Monday 4th December, 2017, 12:30

Editor:
Gregor von Laszewski
Department of Intelligent Systems
Engineering
Indiana University
laszewski@gmail.com

Contents

1 Preface	7
1.0.1 Disclaimer	7
1.0.2 Citation	7
1.1 List of Papers	8
2 Biology	11
3 Business	11
2 hid106	Status: 0%
A Music Recommendation System	
Shiqi Shen, Qiaoyi Liu	11
3 hid302	Status: 100%
Recipe Ingredients Analysis	
Sushant Athaley	11
4 hid310	Status: 100%
Gerrymandering Detection Using Data Analysis	
Kevin Duffy	29
4 Edge Computing	29
5 Education	29
6 Energy	29
7 Environment	29
8 Government	29
9 Health	29
5 hid327	Status: 99%
How Big Data will Help Improve People's Health Worldwide	
Paul Marks	29
6 hid335	Status: 95%
Using Machine Learning Classification of Opioid Addiction for Big Data Health Analytics	
Sean Shiverick	45
10 Lifestyle	74

11 Machine Learning	74
7 hid343	Status: 100 %
Income Prediction Using Machine Learning Techniques	
Borga Edionse Usifo	74
12 Media	103
13 Physics	103
14 Security	103
8 hid316	Status: 100%
Edge Analytics for Weather Monitoring and Forecasting	
Robert Gasiewicz	103
15 Sports	103
16 Technology	103
9 hid306	Status: 100%; 12/3/2017
Predicting Housing Prices - Kaggle Competition	
Murali Cheruvu, Anand Sriramulu	103
10 hid315	Status: 80% - finishing up .txt
TBI - A Journey Beyond Contact Sports	
Garner, Jeffry	118
11 hid337	Status: Dec 04 17 100%
IoT and Big Data Analytics for Equipment Predictive Health Management (PHM)	
Ashok Reddy Singam, Anil Ravi	118
17 Text	133
18 Theory	133
19 Transportation	133
20 TBD	133
12 hid101	Status: Dec 04 17 0%
TBD	
Huiyi Chen	133
13 hid102	Status: Dec 04 17 0%
None	
Gregor von Laszewski	133
14 hid104	Status: 0%
TBD	
Jones, Gabriel and Millard, Mathew	133
15 hid105	Status: Dec 04 17 0%
None	
Gregor von Laszewski	133

16	hid107		Status: Dec 04 17 0%
	None		
	Gregor von Laszewski	133
17	hid109		Status: 0%
	A Music Recommendation System		
	Shiqi Shen, Qiaoyi Liu	133
18	hid201		Status: not started
	edge computing application using mqtt and raspberry pi		
	Arnav, Arnav	133
19	hid202		Status: 0%
	NA		
	Himani Bhatt	145
20	hid203		Status: not started
	TBD		
	Chandwani, Nisha	145
21	hid204		Status: 0%
	Benchmarking a BigData Docker deployment		
	Chaturvedi, Dhawal	145
22	hid205		Status: 0%
	Benchmarking a BigData Docker deployment		
	Gregor von Laszewski	145
23	hid208		Status: Dec 04 17 0%
	None		
	Gregor von Laszewski	145
24	hid209		Status: Dec 04 17 0%
	None		
	Gregor von Laszewski	145
25	hid210		Status: Dec 04 17 0%
	None		
	Gregor von Laszewski	145
26	hid211		Status: unkown
	Continuous motion tracking using Deep Neural Networks and Recurrent Neural Networks		
	Khamkar, Ajinkya	145
27	hid212		Status: Dec 04 17 0%
	None		
	Gregor von Laszewski	145
28	hid214		Status: 0%
	Benchmarking a BigData Docker deployment		
	Junjie Lu	145
29	hid215		Status: Dec 04 17 0%
	None		
	Gregor von Laszewski	145
30	hid216		Status: Dec 04 17 0%
	None		
	Gregor von Laszewski	145

31	hid218		Status: unkown	
	None			
	Geng Niu			145
32	hid219	Fake Review Detection on Yelp Dataset	Status: Dec 04 17 0%	
	Gregor von Laszewski			145
33	hid224	Big Data Analytics in Detection of DDoS (Distributed Denial-of-Service) attacks	Status: 0%	
	Rawat, Neha			164
34	hid225		Status: Dec 04 17 0%	
	None			
	Gregor von Laszewski			164
35	hid230	big data with natural language processing	Status: unkown	
	Yuanming Huang			164
36	hid231	Big Data Analytics on Food Products around the world	Status: not started	
	Vegi, Karthik			164
37	hid232	Benchmarking a BigData Docker deployment	Status: 0%	
	Gregor von Laszewski			186
38	hid234	Big Data Applications in the Travel Industry and its Potential in Improving Travel Accessibility	Status: 10%	
	Weixuan Wang			186
39	hid235	Big Data analytics in predict house price	Status: unkown	
	Yujie Wu			192
40	hid237	Analyzing everyday challenges of people with visual impairments	Status: 0%	
	Tousif Ahmed			192
41	hid304		Status: Dec 04 17 0%	
	None			
	Gregor von Laszewski			192
42	hid308	Benchmarking a BigData Docker deployment	Status: 0%	
	Pravin Deshmukh			192
43	hid312	To be decided	Status: not yet started	
	Neil Eliason			192
44	hid313	The Impact of Clinical Trial Results on Pharmaceutical Stock Performance	Status: 90%	
	Tiffany Fabianac			192
45	hid319	Facial Recognition and Object Detection using Raspberry Pi Robot Car	Status: 30%	
	Mani Kumar Kagita			202

46 hid320		Status: 100% Dec 03 2017
	Real Estate Big Data Analysis	
	Elena Kirzhner	202
47 hid323		Status: Dec 04 17 0%
	None	
	Gregor von Laszewski	202
48 hid324		Status: Dec 04 17 0%
	None	
	Gregor von Laszewski	202
49 hid326		Status: unkown
	None	
	Mohan Mahendrakar	202
50 hid331		Status: Dec 04 17 0%
	Big Data Applications in Predicting Hospital Readmissions	
	Tyler Peterson	202
51 hid332		Status: 100%
	Big Data Analytics to Reduce Health Care in the United States	
	Judy Phillips	215
52 hid333		Status: 100%
	IoT and Big Data Analytics for Equipment Predictive Health Management type: latex	
	Anil Ravi, Ashok Reddy Singam	234
53 hid334		Status: 0%
	Sentiment based Macroeconomic Investing using GDELT	
	Peter Russell	234
54 hid336		Status: Dec 04 17 0%
	None	
	Gregor von Laszewski	234
55 hid339		Status: Dec 2 2017 100%
	Diagnosis of Coronary Artery Disease Using Big Data Analysis	
	Hady Sylla	234
56 hid340		Status: Dec 04 17 0%
	None	
	Gregor von Laszewski	234
57 hid341		Status: 0%
	Not submitted	
	Tibenkana, Jacob	234
58 hid342		Status: 0%
	TBD	
	Nsikan Udojen	234
59 hid346		Status: unkown
	NOt submitted	
	Zachary Meier	234
60 hid348		Status: 80%
	Not submitted	
	Budhaditya Roy	234

Chapter 1

Preface

1.0.1 Disclaimer

The papers provided are contributed by students of the i523 class thought at Indiana University in Fall of 2017. The students were educated in plagiarizm and we hope that all papers meet the high standrads provided by the policies set at Indiana University in regards to plagiarizm. In case you notice any issues, please contact Gregor von Laszewski (laszewski@gmail.com) so we cn address the issue with the student.

1.0.2 Citation

The proceedings is at this time available as a draft. To cite this proceedings you can use the following citation entry:

```
@Book{las17-i523,
  editor = {Gregor von Laszewski},
  title = {Use Cases in Big Data Software and Analytics},
  publisher = {Indiana University},
  year = {2017},
  volume = {1},
  series = {i523},
  address = {Bloomington, IN},
  edition = {1},
  month = dec,
  url={https://github.com/laszewski/laszewski.github.io/raw/master/papers/vonLaszewski-i
} }
```

Contributors to the volume can cite their contribution as follows. They just need to *FILLIN* the missing information

```
@InBook{las17-,
  author = {FILLIN},
```

```

editor =      {Gregor von Laszewski},
title =       {Use Cases in Big Data Software and Analytics},
chapter =     {FILLIN},
publisher =   {Indiana University},
year =        {2017},
volume =      {1},
series =      {i523},
address =     {Bloomington, IN},
edition =     {1},
month =       dec,
url={https://github.com/laszewski/laszewski.github.io/raw/master/papers/vonLaszewski-i
pages =       {FILLIN},
}

```

1.1 List of Papers

HID	Author	Title
101	Huiyi Chen	TBD
0	Gregor von Laszewski	None
104, 216	Jones, Gabriel and Millard, Mathew	TBD
0	Gregor von Laszewski	None
109, 106	Shiqi Shen, Qiaoyi Liu	A Music Recommendation System
0	Gregor von Laszewski	None
109, 106	Shiqi Shen, Qiaoyi Liu	A Music Recommendation System
hid111	error: yaml	A Music Recommendation System
hid201	error: yaml	edge computing application using mqtt and raspberry pi
202	Himani Bhatt	NA
0	Chandwani, Nisha	TBD
0	Chaturvedi, Dhawal	Benchmarking a BigData Docker deployment
205	Gregor von Laszewski	Benchmarking a BigData Docker deployment
0	Gregor von Laszewski	None
0	Gregor von Laszewski	None
0	Gregor von Laszewski	None
211	Khamkar, Ajinkya	Continuous motion tracking using Deep Neural Networks and Recurrent Neural Networks
0	Gregor von Laszewski	None
hid213	error: yaml	None
214	Junjie Lu	Benchmarking a BigData Docker deployment
0	Gregor von Laszewski	None
0	Gregor von Laszewski	None
218	Geng Niu	None
219	Gregor von Laszewski	Fake Review Detection on Yelp Dataset
224	Rawat, Neha	Big Data Analytics in Detection of DDoS (Distributed Denial-of-Service) attacks

0	Gregor von Laszewski	None
hid228	error: yaml	None
hid229	error: yaml	None
230	Yuanming Huang	big data with natural language processing
0	Vegi, Karthik	Big Data Analytics on Food Products around the world
hid232	error: yaml	Benchmarking a BigData Docker deployment
233	Wang, Jiaan and Chaturvedi, and Dhawal 204	Big Data in Safe Driver Prediction
hid234	error: yaml	Big Data Applications in the Travel Industry and its Potential in Improving Travel Accessibility
235	Yujie Wu	Big Data analytics in predict house price
hid236	error: yaml	Big Data analytics in predict house price
237	Tousif Ahmed	Analyzing everyday challenges of people with visual impairments
301	Gagan Arora	Importance of Big data in predicting stock price
302	Sushant Athaley	Recipe Ingredients Analysis
0	Gregor von Laszewski	None
hid305	error: yaml	None
306, 338	Murali Cheruvu, Anand Sriramulu	Predicting Housing Prices - Kaggle Competition
308	Pravin Deshmukh	Benchmarking a BigData Docker deployment
hid309	error: yaml	Benchmarking a BigData Docker deployment
310	Kevin Duffy	Gerrymandering Detection Using Data Analysis
hid311	error: yaml	Gerrymandering Detection Using Data Analysis
hid312	error: yaml	To be decided
313	Tiffany Fabianac	The Impact of Clinical Trial Results on Pharmaceutical Stock Performance
hid314	error: yaml	The Impact of Clinical Trial Results on Pharmaceutical Stock Performance
315	Garner, Jeffry	TBI - A Journey Beyond Contact Sports
316	Robert Gasiewicz	Edge Analytics for Weather Monitoring and Forecasting
hid318	error: yaml	Edge Analytics for Weather Monitoring and Forecasting
319	Mani Kumar Kagita	Facial Recognition and Object Detection using Raspberry Pi Robot Car
320	Elena Kirzhner	Real Estate Big Data Analysis
hid321	error: yaml	Real Estate Big Data Analysis
0	Gregor von Laszewski	None
0	Gregor von Laszewski	None
325	J. Robert Langlois	The importance of data sharing and the replication of the sciences, but what about data archiving?
218	Mohan Mahendrakar	None
327	Paul Marks	How Big Data will Help Improve People's Health Worldwide
hid328	error: yaml	How Big Data will Help Improve People's Health Worldwide
hid329	error: yaml	How Big Data will Help Improve People's Health Worldwide
330	Janaki Mudvari Khatiwada	Big Data Analytics in Monitoring Outdoor Air Quality
hid331	error: yaml	Big Data Applications in Predicting Hospital Readmissions
332	Judy Phillips	Big Data Analytics to Reduce Health Care in the United States

333, 337	Anil Ravi, Ashok Reddy Singam	IoT and Big Data Analytics for Equipment Predictive Health Management type: latex
334	Peter Russell	Sentiment based Macroeconomic Investing using GDELT
335	Sean Shiverick	Using Machine Learning Classification of Opioid Addiction for Big Data Health Analytics
0	Gregor von Laszewski	None
hid337	error: yaml	IoT and Big Data Analytics for Equipment Predictive Health Management (PHM)
hid338	error: yaml	IoT and Big Data Analytics for Equipment Predictive Health Management (PHM)
332	Hady Sylla	Diagnosis of Coronary Artery Disease Using Big Data Analysis
0	Gregor von Laszewski	None
341	Tibenkana, Jacob	Not submitted
342	Nsikan Udoyen	TBD
343	Borga Edionse Usifo	Income Prediction Using Machine Learning Techniques
345	Ross Wood	Agricultural Data Science: Then, Now, and Beyond
346	Zachary Meier	NOt submitted
347	Jeramy Townsley	Mapping Police Killing of Citizens in the United States
348	Budhaditya Roy	Not submitted

Recipe Ingredient Analysis

Sushant Athaley
Indiana University
sathaley@iu.edu

ABSTRACT

Food is the unavoidable part of day to day of human life. Ingredients play a major role or are the basic requirement in preparation of any kind of food. We can find the humongous list of ingredients getting used across globally along with other details which constitute to big data. We explore ingredients getting used in various recipes across the globe to understand most used ingredient, key ingredients of various cuisine and the relationship between the ingredients to find out closely related ingredients which can always provide great dish if used together.

KEYWORDS

i523, hid302, big data, ingredient, recipe, analysis, python, gephi

1 INTRODUCTION

Ingredients are vital for human existence as well as for food or restaurant industry. We use it every day for cooking and food industry uses it to produce consumable for their customers. Ingredient inspires chefs to come up with new culinary artistry. So what do we know about this essential element of the life and what data tell us? Ingredients come in different size, color, shape, flavor, nutrition, taste, texture, grows in specific weather conditions and this provides a great opportunity for various analysis which can be useful for the human being as well as business industries. So main focus of this study is on the ingredients used in various recipes across the cuisines. This study evaluates recipe ingredient dataset from Kaggle [7] to analyze most used ingredients, key cuisine ingredients and ingredient relationship.

This study is organized as follows, section *Ingredient* defines ingredient and its various characteristics. section *Ingredieint Analytics* describes various analytics which can be performed on the ingredient with some examples. Section *Project* describes the aim of this study. Section *technologies* provides information on the tools and technologies used for this project. Section *Methodology* covers overall process carried out in this project. Section *Dataset* describes data structure used along with loading process and data findings. Section *Analysis and Findings* describes various analysis carried out on the data and the visual representation of the analysis. Section *Shortcomings* captures shortcomings of the project. Section *Future Work* talks about what else can be done with this dataset which is not covered in the current scope of the project. Section *Conclusion* concludes the study.

2 INGREDIENT

Food is defined as “Edible or potable substance (usually of animal or plant origin), consisting of nourishing and nutritive components such as carbohydrates, fats, proteins, essential mineral and vitamins, which (when ingested and assimilated through digestion) sustains life, generates energy, and provides growth, maintenance,

and health of the body” [2]. Thus food is the basic necessity for human for the sustainability. Food can be eaten raw, cooked or processed. As human race evolved over the period of time, the way we eat food is also evolved. Food cooking is just not the basic necessity but its an art and science in today’s era. Food preparation consists of various cooking techniques, tools, and ingredients to make it palatable or edible by humans. The ingredient is by far the most important part of any food or recipe preparation. The recipe consists of the list of ingredients and the set of instruction to cook particular food dish [5]. An ingredient is defined as “Any of the foods or substances that are combined to make a particular dish” [9]. Ingredients impart various flavors, aroma, texture, and color to the cooking dish. Ingredients are mostly derived from vegetables, fruits, nuts, grains, living organisms, herbs, flowers, and spices. It comes in both solid and liquid forms. Another characteristic of ingredients is the nutritional value they provide which is essential for the human body.

3 INGREDIENT ANALYTICS

Ingredients characteristics and the combination of other related data provides various opportunities to analyze ingredient in different ways. Analysis of the flavors present in ingredient can provide us with the categorization of the different ingredient by the flavor profile which can be helpful in deciding substitute ingredient if a certain ingredient is not present or pairing ingredient from different flavor categories to construct the dish as per the taste required [1]. This analysis also helps to understand which ingredients cannot be used together. A similar analysis is carried out to correlate ingredient across recipes to come up with top 50 combinations of ingredients which can be used together [6]. Flavourspace application provides functionality to search recipe based on the ingredients, suggests alternate ingredient if not present, adjust the recipe as per the taste which is a good example of big data analytics in food industry [10]. Foodpairing application takes another approach to form the connection between unfamiliar ingredients and provides information on how to use such ingredient to make a dish, this is very helpful in terms of sustainability as we can use ingredient which is ample available but not in use due to the absence of information on using such ingredients [8].

Another study conducted on most used ingredient provides insight that sugar, oil, pepper, and salt are most commonly occurring ingredient, among spices clove, in vegetable onion, garlic , and tomatoes, butter in milk product, eggs followed by chicken in the animal product are the most used ingredient in the categories [3]. This information can help in better planning and sourcing of such ingredients which are in high demand.

Ingredient nutrition analysis can help find out nutrition of the food prepared by those ingredients. This would be helpful in menu planning where nutrition information is the key factor such as school, hospitals or any other dietary program [4].

Recipe cost is calculated by including the cost of the ingredient used in that recipe. Ingredient cost as per the quantity used in recipe provides base information to calculate the price of any recipe. This ingredient cost analysis provides an avenue to reduce the cost of the recipe by using substitute ingredient of lesser cost. This can also help in household budget to keep in check as well as make restaurant industry profitable.

Ingredient used in recipe can provide insight into type of weather received by that cuisine as ingredient can grow in certain weather condition. This can help chef locally source the ingredient and maintain local agriculture sustainability.

4 PROJECT

This project study is conducted to analyze ingredients getting used in various recipes across the cuisines to find out

- Most used ingredients across cuisines or globally
- Key ingredients used by cuisines
- Ingredient relationship or connection to understand the related ingredients

4.1 Technologies

Technologies and tools used in this projects are

- Python version 3.6 is used for data load and processing
- Gephi 0.9.2 for visualization
- Spyder 3.0 as a Python IDE

4.2 Methodology

The first step was to source the data. We were interested in the dataset which provides recipe information along with the ingredient used in the recipe. Since we wanted to analyze distribution across cuisines, data should also contain cuisine tagging. This dataset can be generated by pulling recipe data from various online applications or pick from publicly available datasets. We finalized publicly available dataset at Kaggle application satisfying need for this project.

Figure 1 shows methodology used for this project to analyze ingredient data.

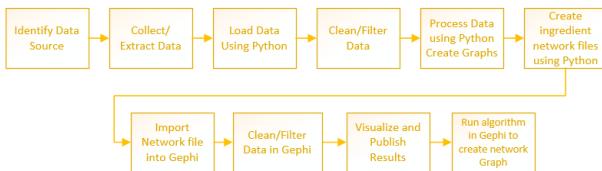


Figure 1: Flowchart of the Methodology to Analyze Ingredients

DataSet is loaded through Python script and further processed to clean the data. This cleaned data then processed to analyze ingredient distribution across cuisine and per cuisine. Gephi software is used to analyze the relationship and to find out the ingredient modularity. The Python script is used to create the network files required by the Gephi tool. Gephi requires Nodes and relationship in terms of Edges between the nodes for the analysis. The Python

script is used to create Node and Edges file in excel format so that it can be imported into Gephi. Distinct ingredients used in recipe becomes the nodes. Edges or relationship between ingredients is derived by relating ingredients appearing in the same recipe. All ingredient in the same recipe is considered related to each other. Network files created by Python are imported in Gephi to produce the graph for the visualization. Gephi tools data laboratory is used to clean up the data and filters are applied to provide usable network visualization.

4.3 Dataset

The dataset for this study is sourced from Kaggle application [7]. This dataset is publicly available and featured in “What’s Cooking?” competition. This dataset is in JSON format and of 12MB size. This dataset contains recipe id, cuisine and list of ingredients as described in Figure 2. This dataset contains total 39774 recipes

```
{
  "id": 24717,
  "cuisine": "indian",
  "ingredients": [
    "tumeric",
    "vegetable stock",
    "tomatoes",
    "garam masala",
    "naan",
    "red lentils",
    "red chili peppers",
    "onions",
    "spinach",
    "sweet potatoes"
  ]
},
```

Figure 2: Ingredient Data Structure

across various cuisines. We used two different methods to load this data. Cuisine and ingredient analysis is done by loading data into *pandas dataframe* and to analyze ingredient relationship data has been loaded into *json* object. Figure 3 shows the code for data loading used in this project.

```
#read the ingredient data using pandas
dfTrain = pd.read_json('./data/train.json')
```

```
#load data using json
dataFilePath='./data/train.json'
with open(dataFilePath) as data_file:
    data = json.load(data_file)
```

Figure 3: Data Loading

Ingredient extraction from the data structure and processing was challenging as ingredients are listed comma separated for each recipe. Also, ingredient list can vary by recipe and there is no

proper structure. Another issue with the ingredient list is ingredient appears in various forms but it's the same ingredient which gives duplicate data. For example, salt appears as salt, kosher salt, Morton Salt, sea salt, table salt, Himalayan salt, fine sea salt, low sodium salt, fine salt. This is the same ingredient but come across in recipe as a different ingredient and getting counted as a separate ingredient in the analysis. Some ingredients are listed along with measures like (10 oz.) frozen chopped spinach, (10 oz.) frozen chopped spinach, thawed and squeezed dry, (14.5 oz.) diced tomatoes and getting counted as a separate ingredient. Some ingredients are listed along with the brand name like KRAFT Reduced Fat Shredded Mozzarella Cheese, Johnsonville Smoked Sausage, Johnsonville Mild Italian Sausage Links etc and also constitutes to the ingredient list. This variation makes difficult to get the proper ingredient list for the analysis. Extensive work is needed to clean and correct the noisy data so that proper analysis can be carried out. This correction process is not carried out as part of this project.

Certain ingredients like salt or water etc should be avoided from the analysis as those are not the ingredient we are looking for the analysis. We tried to clean such elements during ingredient relationship analysis but we had little success as those ingredients are present in the dataset in various forms.

4.4 Analysis and Findings

4.4.1 Recipe Distribution By Cuisine. We first analyze entire dataset to understand the total number of recipes and their distribution across various cuisines. We use Pythons Panda library to get the total recipe count as 39774 and plot the distribution. Figure 4 shows number of recipes per cuisine. Dataset is heavily dominated by Italian cuisine followed by Mexican cuisine and with very few recipes from Russian and Brazilian cuisines. This also highlights another shortcoming of the dataset that it doesn't have equal representation of all cuisines which might give us biased analysis.

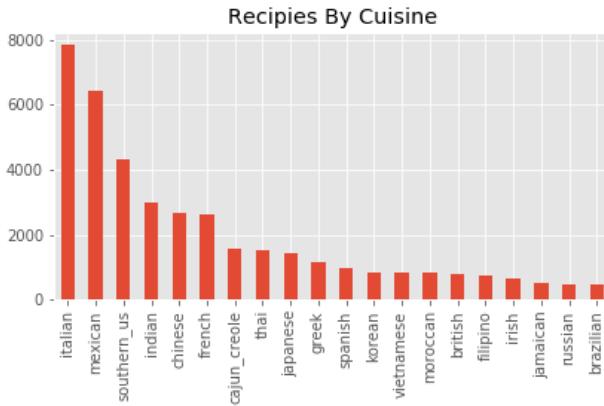


Figure 4: Recipe Distribution By Cuisine

Table1 describes recipe count for every cuisine.

4.4.2 Most Used Ingredients All Cuisines. The second analysis is carried out to understand top 20 ingredients getting used across cuisine or globally. Ingredient *Salt* is obvious topper followed by

Table 1: Recipe Count By Cuisine

Cuisine	Recipe Count
brazilian	467
british	804
cajun creole	1546
chinese	2673
filipino	755
french	2646
greek	1175
indian	3003
irish	667
italian	7838
jamaican	526
japanese	1423
korean	830
mexican	6438
moroccan	821
russian	489
southern us	4320
spanish	989
thai	1539
vietnamese	825

Oil and Onions. This also proves our craving for saltiness and fat. Figure 5 shows top 20 ingredient across cuisines.

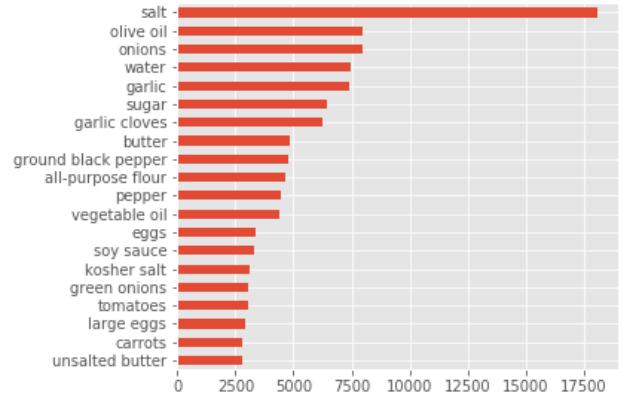


Figure 5: Top 20 Ingredients

4.4.3 Ingredients Distribution By Cuisines. The third analysis is carried out to understand key ingredient for each cuisine. These key ingredients define those cuisines and provide unique test characterized by that cuisine. We limited ingredient list to top 10 to get the key ingredients for each cuisine. Figure 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25 shows top 10 key ingredient used in the corresponding cuisines.

4.4.4 Ingredients Relationship. Forth analysis is carried out to understand the relationship between the ingredient to find out ingredient clusters. This analysis helps us understand the ingredient combinations which can be used together to provide great dish

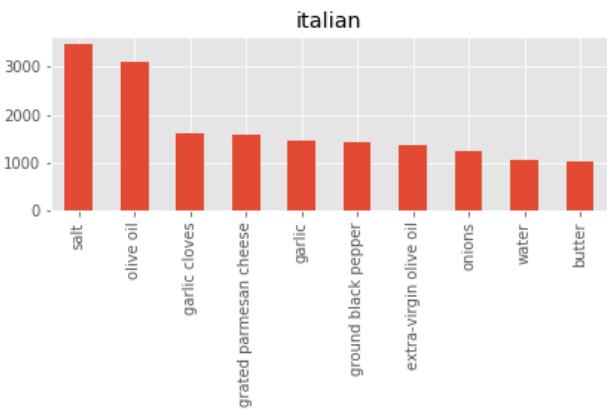


Figure 6: Top 10 Ingredients

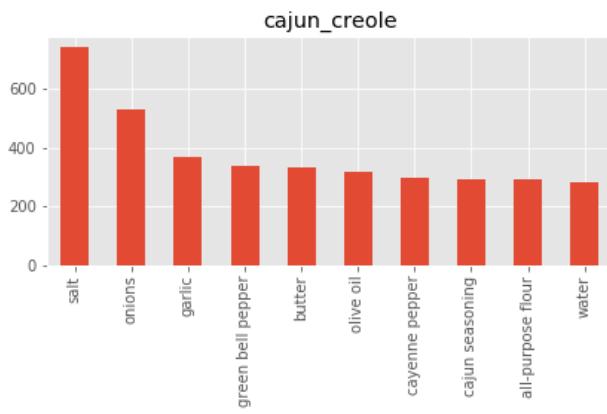


Figure 9: Top 10 Ingredients

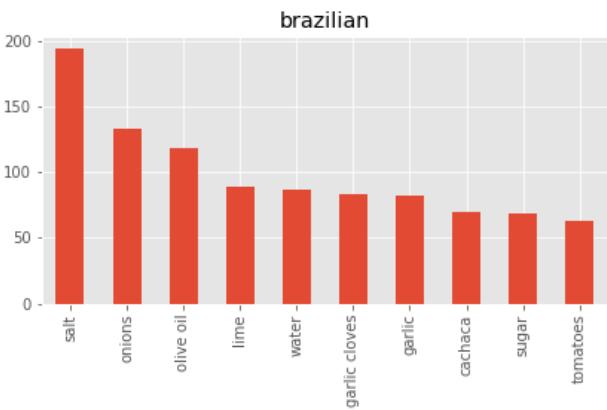


Figure 7: Top 10 Ingredients

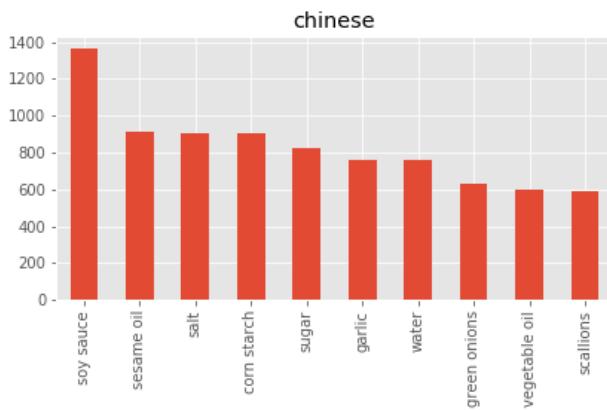


Figure 10: Top 10 Ingredients

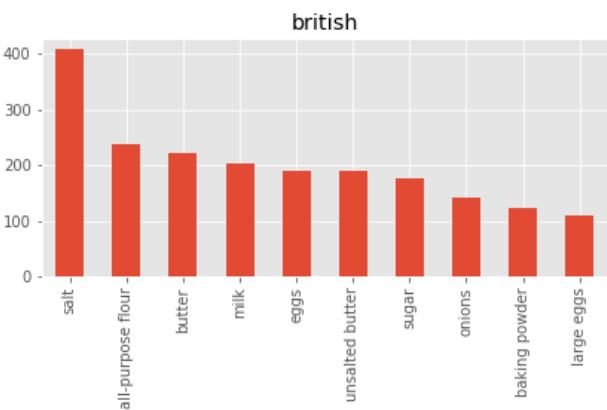


Figure 8: Top 10 Ingredients

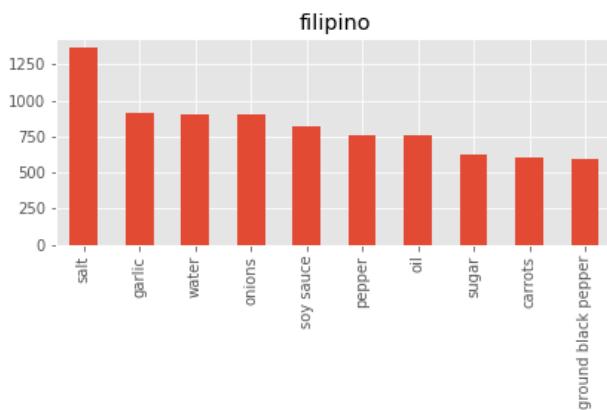


Figure 11: Top 10 Ingredients

every time. This model can be used to predict ingredients for certain recipe based on the cluster. We used Gephi tool to analyze and produce the graph for this analysis. Gephi accepts network

structure in terms of Node and Edge relationship. We created this network using python by relating all ingredients present in the recipe with each other. Ingredients become the node and source

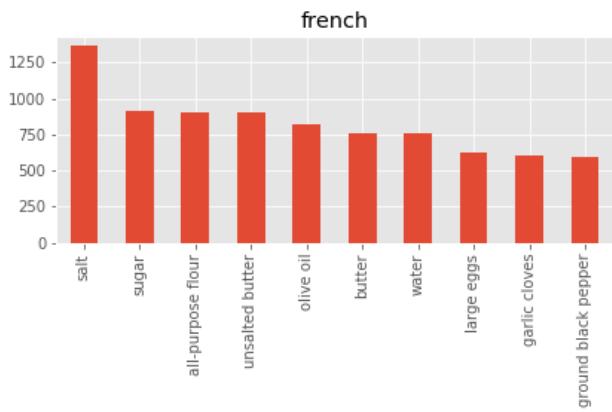


Figure 12: Top 10 Ingredients

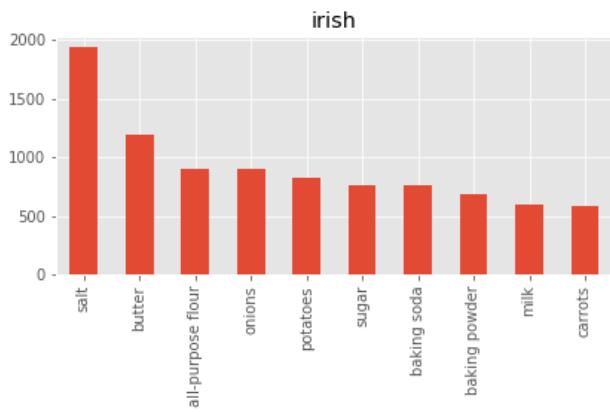


Figure 15: Top 10 Ingredients

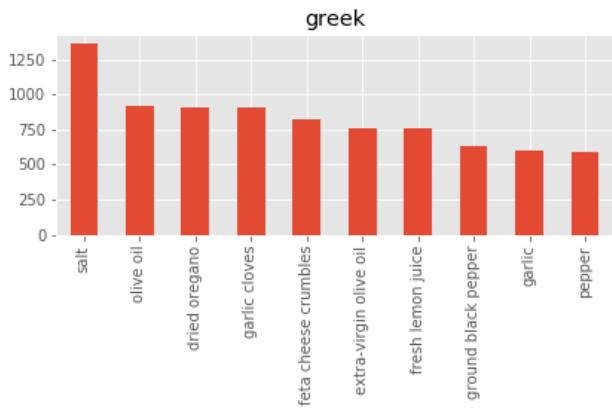


Figure 13: Top 10 Ingredients

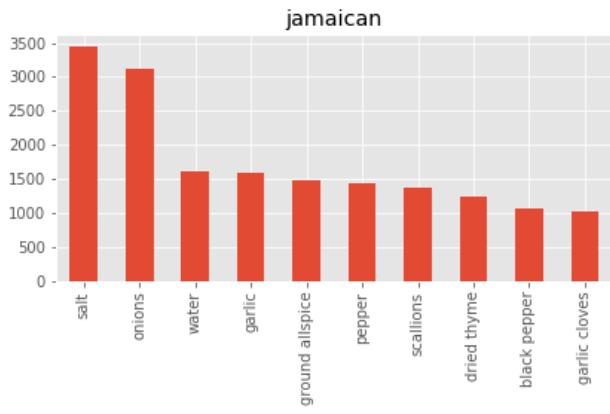


Figure 16: Top 10 Ingredients

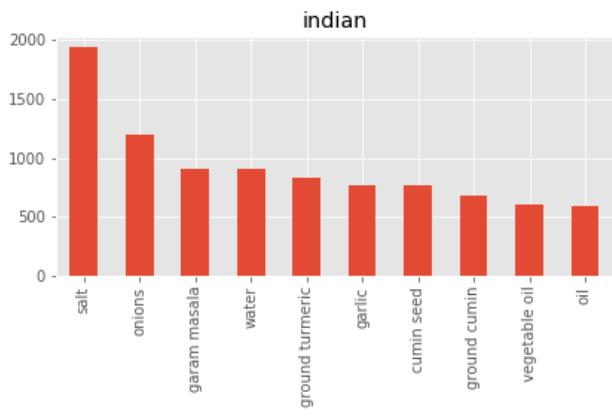


Figure 14: Top 10 Ingredients

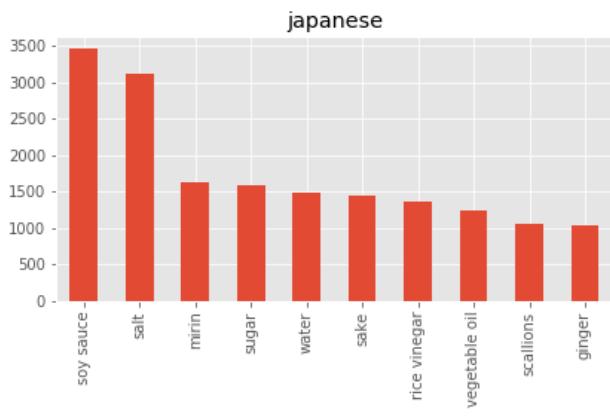


Figure 17: Top 10 Ingredients

and target nodes become the edges. These network files generated in excel spreadsheet are converted to CSV format and imported into the Gephi tool. Import created 5405 Nodes and 290828 edges

for processing and analysis. Force Atlas 2 layout present in Gephi has been applied to the network which brings nodes with higher weights and shared connections closer to each other. We also used

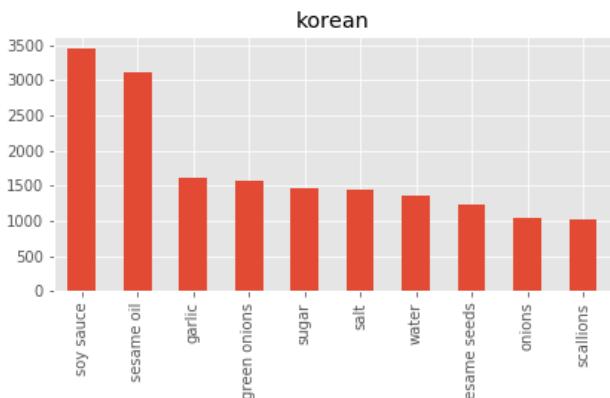


Figure 18: Top 10 Ingredients

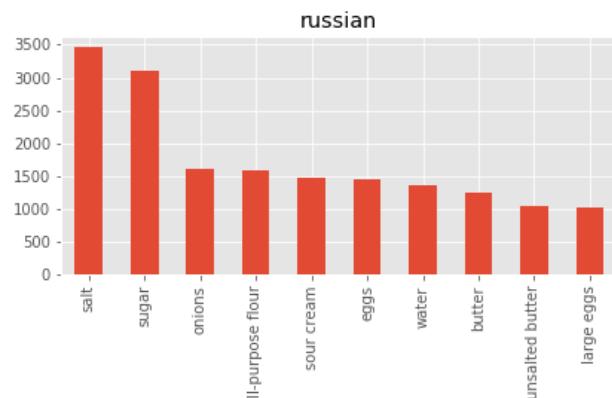


Figure 21: Top 10 Ingredients

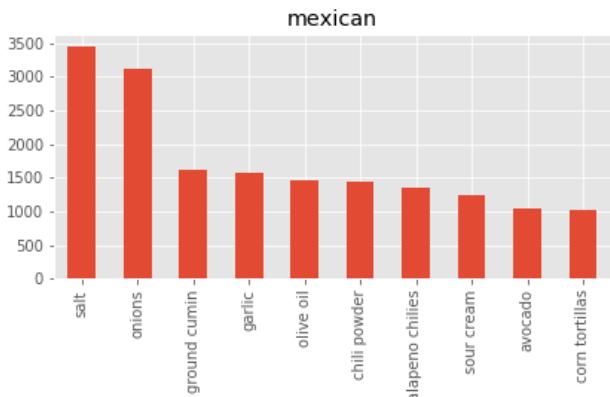


Figure 19: Top 10 Ingredients

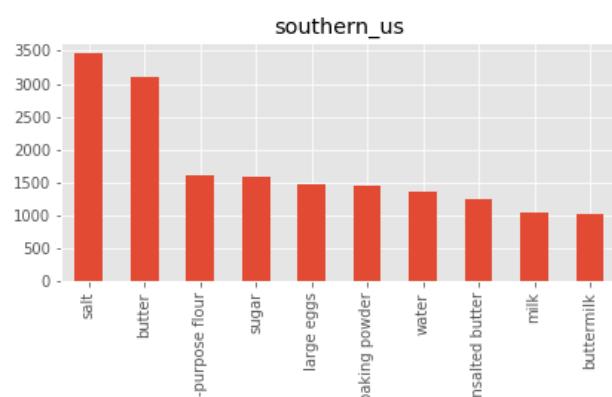


Figure 22: Top 10 Ingredients

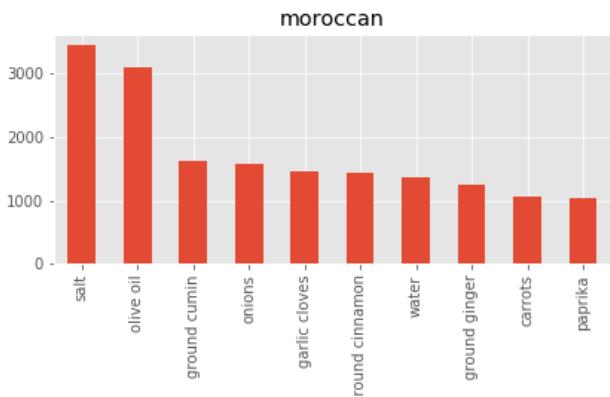


Figure 20: Top 10 Ingredients

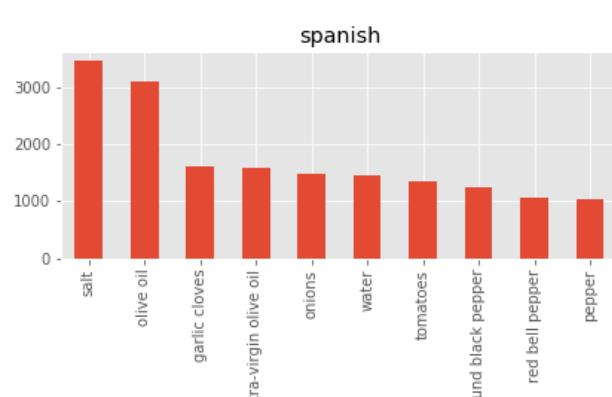


Figure 23: Top 10 Ingredients

Gephi Data Laboratory to clean up duplicate or unwanted nodes. Filtering based on Degree Range and Edge Weight has been applied to data to reduce node and edges to get the graph which can be used

for analysis and avoid crashing Gephi due to large data. Modularity statistic uncovered 5 ingredient clusters which can be identified

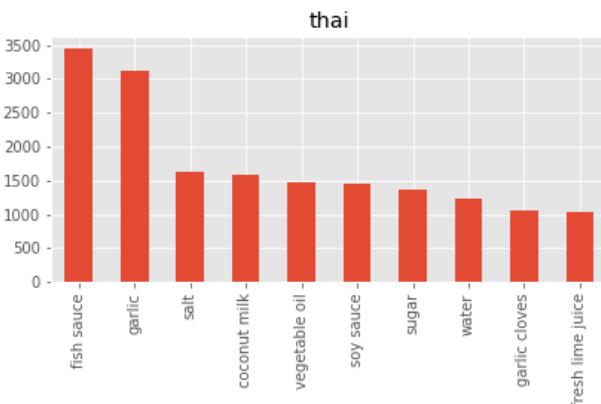


Figure 24: Top 10 Ingredients

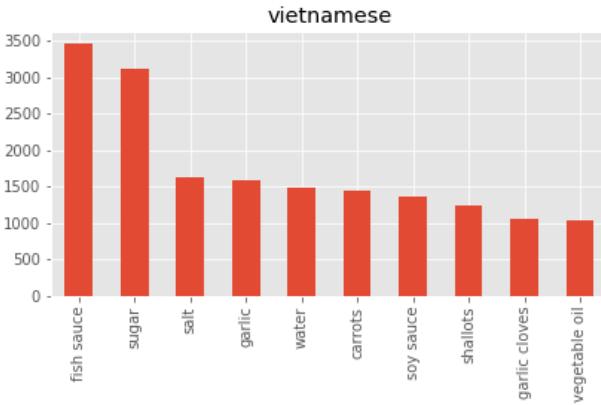


Figure 25: Top 10 Ingredients

by different color in the graph. This cluster can approximately relate to the cuisines present in our dataset and confirms our earlier analysis of ingredient by cuisine.

- Orange - Mexican
- Brown - Indian
- Blue - Chinese
- Green - Italian
- Gray - Southern US

Figure 26 shows ingredient cluster of more than 1000 nodes.
 Figure 27 shows ingredient cluster of around 100 nodes.

4.5 Shortcomings

Improper documentation of ingredient names in the dataset reduces the correctness of this analysis. In absence of proper ingredient name and duplication of ingredient name prevents getting exact ingredient weight into the analysis. A dataset with uniform ingredient name can help this analysis to achieve its best. If we don't find proper ingredient name then this analysis needs to include extensive data cleaning process which can be considered an improvement to this project.

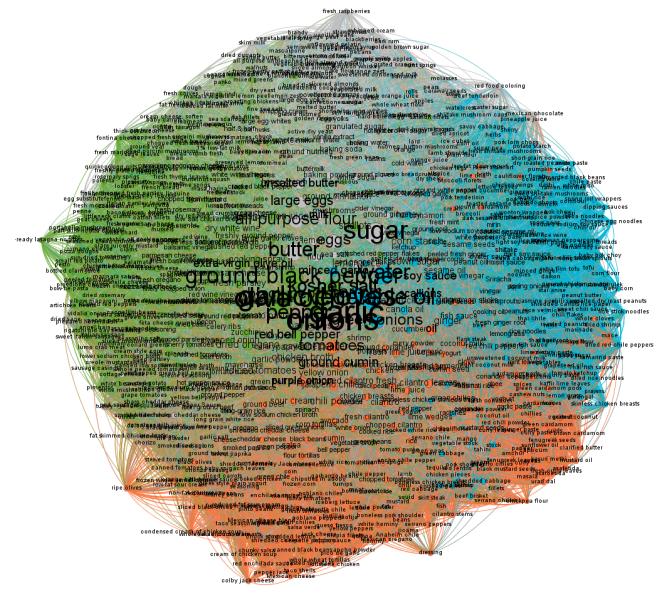


Figure 26: Ingredient Cluster

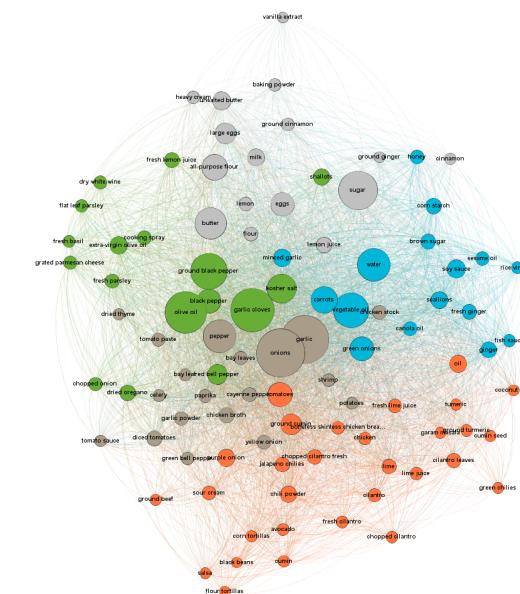


Figure 27: ingredient Cluster 100 Nodes

Network file creation algorithm can be enhanced further by considering the number of recipes for the ingredient to provide additional weight to the relationship which can provide the stronger bond between the ingredients.

4.6 Future Work

This dataset can be analyzed to find out ingredient overlap between various cuisine and can provide insight into the influence of one cuisine on another. Usually, geographically neighboring cuisines are influenced by each other as they share common ingredients.

5 CONCLUSION

This project shows most used ingredient, ingredient distribution by cuisine and predictive ingredient relationship model as per the goal of the project. We also show various opportunities present with ingredient data analysis and role of big data analytics. We prove human craving for salty and fatty food as salt and oil are most used ingredient across cuisines as per the analysis. We understand now based on our analysis key ingredient of any cuisine. Ingredient cluster shows why those ingredients are the base of certain cuisine and recipe of those ingredients always turn out delicious. We also crave for the good data so that we can provide more accurate analysis of the ingredients. Ingredient analysis has potential not only to help restaurant and food industry but it can help with our social responsibility of sustainability and understanding different cuisines and culture. As food industries interest grows in big data analytics, we will continue to see more evaluations of the ingredients.

ACKNOWLEDGMENTS

The author would like to thank Dr. Gregor von Laszewski for his support and suggestions in this project. The author would also like to thank Kaggle application for hosting ingredient dataset which is used in this project and various online resource which helped understand Python and Gephi.

REFERENCES

- [1] Bagrow James P. Barabsi Albert-Lszl Ahn Yong-Yeol, Ahnert Sebastian E. 2011. Flavor network and the principles of food pairing. *Scientific Reports* 1, 196 (2011). <https://www.nature.com/articles/srep00196#supplementary-information>
- [2] businessdictionary. 2017. Food. web. (2017). <http://www.businessdictionary.com/definition/food.html>
- [3] Usashi Chatterjee, Vinit Kumar, and Devika P. Madalli. 2016. Formalizing Food Ingredients for Data Analysis and Knowledge Organization. 10 (07 2016), 289–309.
- [4] S. M. Church. 2015. The importance of food composition data in recipe analysis. web. (2015). <http://onlinelibrary.wiley.com/doi/10.1111/nbu.12125/abstract>
- [5] collinsdictionary. 2017. Recipe. web. (2017). <https://www.collinsdictionary.com/us/dictionary/english/recipe>
- [6] inkhorn82. 2014. A Delicious Analysis. web. (2014). <https://www.r-bloggers.com/a-delicious-analysis-aka-topic-modelling-using-recipes/>
- [7] kaggle. 2015. What's Cooking? web. (2015). <https://www.kaggle.com/c/whats-cooking/data>
- [8] Bernard Lahousse. 2016. Using Big Data to Transform Unfamiliar Ingredients Into Tasty Recipes. web. (2016). <https://foodtechconnect.com/2016/04/20/big-food-data-recipes-from-unfamiliar-ingredients/>
- [9] oxforddictionaries. 2017. Ingredient. web. (2017). <https://en.oxforddictionaries.com/definition/ingredient>
- [10] Matthew Robinson. 2015. Big Data Analytics and Food Come Together At Flavourspace. web. (2015). <http://www.theculinaryexchange.com/food-innovation/big-data-analytics-and-food-come-together-at-flavourspace/>

A ISSUES

DONE:

Example of done item: Once you fix an item, change TODO to DONE

A.1 Assignment Submission Issues

DONE:

Do not make changes to your paper during grading, when your repository should be frozen.

A.2 Uncaught Bibliography Errors

DONE:

Missing bibliography file generated by JabRef

DONE:

Bibtex labels cannot have any spaces, _ or & in it

DONE:

Citations in text showing as [?]: this means either your report.bib is not up-to-date or there is a spelling error in the label of the item you want to cite, either in report.bib or in report.tex

A.3 Formatting

DONE:

Incorrect number of keywords or HID and i523 not included in the keywords

DONE:

Other formatting issues

A.4 Writing Errors

DONE:

Errors in title, e.g. capitalization

DONE:

Spelling errors

DONE:

Are you using *a* and *the* properly?

DONE:

Do not use phrases such as *shown in the Figure below*. Instead, use *as shown in Figure 3*, when referring to the 3rd figure

DONE:

Do not use the word *I* instead use *we* even if you are the sole author

DONE:

Do not use the phrase *In this paper/report we show* instead use *We show*. It is not important if this is a paper or a report and does not need to be mentioned

DONE:

If you want to say *and* do not use & but use the word *and*

DONE:

Use a space after . , :

DONE:

When using a section command, the section title is not written in all-caps as format does this for you

\section{Introduction} and NOT \section{INTRODUCTION}

A.5 Citation Issues and Plagiarism

DONE:

It is your responsibility to make sure no plagiarism occurs.
The instructions and resources were given in the class

DONE:

Claims made without citations provided

DONE:

Need to paraphrase long quotations (whole sentences or longer)

DONE:

Need to quote directly cited material

A.6 Character Errors

DONE:

Erroneous use of quotation marks, i.e. use “quotes”, instead of ” ”

DONE:

To emphasize a word, use *emphasize* and not “quote”

DONE:

When using the characters & # % - put a backslash before them so that they show up correctly

DONE:

Pasting and copying from the Web often results in non-ASCII characters to be used in your text, please remove them and replace accordingly. This is the case for quotes, dashes and all the other special characters.

DONE:

If you see a ?gure and not a figure in text you copied from a text that has the fi combined as a single character

A.7 Structural Issues

DONE:

Acknowledgement section missing

DONE:

Incorrect README file

DONE:

In case of a class and if you do a multi-author paper, you need to add an appendix describing who did what in the paper

DONE:

The paper has less than 2 pages of text, i.e. excluding images, tables and figures

DONE:

The paper has more than 6 pages of text, i.e. excluding images, tables and figures

DONE:

Do not artificially inflate your paper if you are below the page limit

A.8 Details about the Figures and Tables

DONE:

Capitalization errors in referring to captions, e.g. Figure 1, Table 2

DONE:

Do use *label* and *ref* to automatically create figure numbers

DONE:

Wrong placement of figure caption. They should be on the bottom of the figure

DONE:

Wrong placement of table caption. They should be on the top of the table

DONE:

Images submitted incorrectly. They should be in native format, e.g. .graffle, .pptx, .png, jpg

DONE:

Do not submit eps images. Instead, convert them to PDF

DONE:

The image files must be in a single directory named "images"

DONE:

In case there is a powerpoint in the submission, the image must be exported as PDF

DONE:

Make the figures large enough so we can read the details. If needed make the figure over two columns

DONE:

Do not worry about the figure placement if they are at a different location than you think. Figures are allowed to float. For this class, you should place all figures at the end of the report.

DONE:

In case you copied a figure from another paper you need to ask for copyright permission. In case of a class paper, you must include a reference to the original in the caption

DONE:

Remove any figure that is not referred to explicitly in the text (As shown in Figure ..)

DONE:

Do not use *textwidth* as a parameter for *includegraphics*

DONE:

Figures should be reasonably sized and often you just need to add *columnwidth*

e.g.

```
/includegraphics[width=\columnwidth]{images/myimage.pdf}  
re
```

bibtext report

=====

This is BibTeX, Version 0.99d (TeX Live 2016)

The top-level auxiliary file: report.aux

The style file: ACM-Reference-Format.bst

Database file #1: report.bib

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barab while executing---line 6169 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barab while executing---line 6169 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barab while executing---line 6169 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barab while executing---line 6169 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barab while executing---line 6169 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barab while executing---line 6169 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barab while executing---line 6169 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barab while executing---line 6261 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barab while executing---line 6261 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barab while executing---line 6261 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barab while executing---line 6457 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barab while executing---line 6457 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barab while executing---line 6457 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barab while executing---line 6457 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barab while executing---line 6457 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barab while executing---line 6457 of file ACM-Reference-Format.bst

Warning--page numbers missing in both pages and numpages fields in Ahn2011

Warning--no journal in Chatterjee2016

(There were 16 error messages)

make[2]: *** [bibtex] Error 2

latex report

```
=====
[2017-12-04 12.23.22] pdflatex report.tex
This is pdfTeX, Version 3.14159265-2.6-1.40.17 (TeX Live 2016) (preloaded format=pdflatex)
Missing character: ""
Typesetting of "report.tex" completed in 2.4s.
./README.yml
 32:12   error    trailing spaces  (trailing-spaces)
 42:81   error    line too long (81 > 80 characters)  (line-length)
 43:81   error    line too long (83 > 80 characters)  (line-length)
 46:81   error    line too long (82 > 80 characters)  (line-length)
 47:81   error    line too long (84 > 80 characters)  (line-length)
 48:81   error    line too long (102 > 80 characters) (line-length)
 49:81   error    line too long (96 > 80 characters)  (line-length)
 54:1    error    trailing spaces  (trailing-spaces)
```

Compliance Report

```
name: Sushant Athaley
hid: 302
paper1: Nov 3 2017 100%
paper2: 100%
project: 100%
```

```
yamlcheck
```

```
wordcount
```

```
(null)
wc 302 project (null) 2935 report.tex
wc 302 project (null) 3564 report.pdf
wc 302 project (null) 361 report.bib
```

```
find "
```

```
-----  
80: "id": 24717,  
81: "cuisine": "indian",  
82: "ingredients": [  
83: "tumeric",  
84: "vegetable stock",  
85: "tomatoes",  
86: "garam masala",  
87: "naan",  
88: "red lentils",  
89: "red chili peppers",  
90: "onions",  
91: "spinach",  
92: "sweet potatoes"  
106: dataFilePath="./data/train.json"
```

passed: False

find footnote

```
-----
```

passed: True

find input{format/i523}

```
-----
```

passed: False

find input{format/final}

```
-----
```

passed: False

floats

66: Figure \ref{f:methodology} shows methodology used for this project
to analyze ingredient data.

67: \begin{figure}[!ht]

68: \centering\includegraphics[width=\columnwidth]{images/methodology.
PNG}

69: \caption{Flowchart of the Methodology to Analyze Ingredients
}\label{f:methodology}

76: The dataset for this study is sourced from Kaggle application
\cite{www-kaggle}. This dataset is publicly available and featured
in ‘‘What’s Cooking?’’ competition. This dataset is in JSON format
and of 12MB size. This dataset contains recipe id, cuisine and
list of ingredients as described in Figure \ref{c:data-structure}.

77: \begin{figure}[htb]

96: \caption{Ingredient Data Structure}\label{c:data-structure}

98: This dataset contains total 39774 recipes across various cuisines.
We used two different methods to load this data. Cuisine and
ingredient analysis is done by loading data into \emph{pandas
dataframe} and to analyze ingredient relationship data has been
loaded into \emph{json} object. Figure \ref{c:data-loading} shows
the code for data loading used in this project.

99: \begin{figure}[htb]

110: \caption{Data Loading}\label{c:data-loading}

120: We first analyze entire dataset to understand the total number of
recipes and their distribution across various cuisines. We use
Pythons Panda library to get the total recipe count as 39774 and
plot the distribution. Figure
\ref{f:Number_of_recipes_by_cuisine} shows number of recipes per
cuisine. Dataset is heavily dominated by Italian cuisine followed
by Mexican cuisine and with very fewer recipes from Russian and
Brazilian cuisines. This also highlights another shortcoming of
the dataset that it doesn’t have equal representation of all
cuisines which might give us biased analysis.

121: \begin{figure}[!ht]

122: \centering\includegraphics[width=\columnwidth]{images/Number_of_r
ecipes_by_cuisine.png}

123: \caption{Recipe Distribution By Cuisine
}\label{f:Number_of_recipes_by_cuisine}

126: Table\ref{t:recipecount} describes recipe count for every
cuisine.

127: \begin{table}[htb]

130: \label{t:recipecount}

159: The second analysis is carried out to understand top 20

ingredients getting used across cuisine or globally. Ingredient \emph{Salt} is obvious topper followed by \emph{Oil} and \emph{Onions}. This also proves our craving for saltiness and fat. Figure \ref{f:Ingredient_Distribution} shows top 20 ingredient across cuisines.

160: \begin{figure} [!ht]
161: \centering\includegraphics[width=\columnwidth]{images/Ingredient_Distribution.png}
162: \caption{Top 20 Ingredients }\label{f:Ingredient_Distribution}

166: The third analysis is carried out to understand key ingredient for each cuisine. These key ingredients define those cuisines and provide unique test characterized by that cuisine. We limited ingredient list to top 10 to get the key ingredients for each cuisine. Figure \ref{f:italian_10_most_used_ingredients}, \ref{f:brazilian_10_most_used_ingredients}, \ref{f:british_10_most_used_ingredients}, \ref{f:cajun_creole_10_most_used_ingredients}, \ref{f:chinese_10_most_used_ingredients}, \ref{f:filipino_10_most_used_ingredients}, \ref{f:french_10_most_used_ingredients}, \ref{f:greek_10_most_used_ingredients}, \ref{f:indian_10_most_used_ingredients}, \ref{f:irish_10_most_used_ingredients}, \ref{f:jamaican_10_most_used_ingredients}, \ref{f:japanese_10_most_used_ingredients}, \ref{f:korean_10_most_used_ingredients}, \ref{f:mexican_10_most_used_ingredients}, \ref{f:moroccan_10_most_used_ingredients}, \ref{f:russian_10_most_used_ingredients}, \ref{f:southern_us_10_most_used_ingredients}, \ref{f:spanish_10_most_used_ingredients}, \ref{f:thai_10_most_used_ingredients}, \ref{f:vietnamese_10_most_used_ingredients} shows top 10 key ingredient used in the corresponding cuisines.

167: \begin{figure} [!ht]
168: \centering\includegraphics[width=\columnwidth]{images/italian_10_most_used_ingredients.png}
169: \caption{Top 10 Ingredients }\label{f:italian_10_most_used_ingredients}

172: \begin{figure} [!ht]
173: \centering\includegraphics[width=\columnwidth]{images/brazilian_10_most_used_ingredients.png}
174: \caption{Top 10 Ingredients }\label{f:brazilian_10_most_used_ingredients}

177: \begin{figure} [!ht]
178: \centering\includegraphics[width=\columnwidth]{images/british_10_

```

    most_used_ingredients.png}
179: \caption{Top 10 Ingredients
} \label{f:british_10_most_used_ingredients}
182: \begin{figure}[!ht]
183: \centering\includegraphics[width=\columnwidth]{images/cajun_creole_10_most_used_ingredients.png}
184: \caption{Top 10 Ingredients
} \label{f:cajun_creole_10_most_used_ingredients}
187: \begin{figure}[!ht]
188: \centering\includegraphics[width=\columnwidth]{images/chinese_10_most_used_ingredients.png}
189: \caption{Top 10 Ingredients
} \label{f:chinese_10_most_used_ingredients}
192: \begin{figure}[!ht]
193: \centering\includegraphics[width=\columnwidth]{images/filipino_10_most_used_ingredients.png}
194: \caption{Top 10 Ingredients
} \label{f:filipino_10_most_used_ingredients}
197: \begin{figure}[!ht]
198: \centering\includegraphics[width=\columnwidth]{images/french_10_most_used_ingredients.png}
199: \caption{Top 10 Ingredients
} \label{f:french_10_most_used_ingredients}
202: \begin{figure}[!ht]
203: \centering\includegraphics[width=\columnwidth]{images/greek_10_most_used_ingredients.png}
204: \caption{Top 10 Ingredients
} \label{f:greek_10_most_used_ingredients}
207: \begin{figure}[!ht]
208: \centering\includegraphics[width=\columnwidth]{images/indian_10_most_used_ingredients.png}
209: \caption{Top 10 Ingredients
} \label{f:indian_10_most_used_ingredients}
212: \begin{figure}[!ht]
213: \centering\includegraphics[width=\columnwidth]{images/irish_10_most_used_ingredients.png}
214: \caption{Top 10 Ingredients
} \label{f:irish_10_most_used_ingredients}
217: \begin{figure}[!ht]
218: \centering\includegraphics[width=\columnwidth]{images/jamaican_10_most_used_ingredients.png}
219: \caption{Top 10 Ingredients
} \label{f:jamaican_10_most_used_ingredients}
222: \begin{figure}[!ht]
223: \centering\includegraphics[width=\columnwidth]{images/japanese_10_most_used_ingredients.png}

```

```

224: \caption{Top 10 Ingredients
    }\label{f:japanese_10_most_used_ingredients}
227: \begin{figure}[!ht]
228: \centering\includegraphics[width=\columnwidth]{images/korean_10_m
ost_used_ingredients.png}
229: \caption{Top 10 Ingredients
    }\label{f:korean_10_most_used_ingredients}
232: \begin{figure}[!ht]
233: \centering\includegraphics[width=\columnwidth]{images/mexican_10_
most_used_ingredients.png}
234: \caption{Top 10 Ingredients
    }\label{f:mexican_10_most_used_ingredients}
237: \begin{figure}[!ht]
238: \centering\includegraphics[width=\columnwidth]{images/moroccan_10
_most_used_ingredients.png}
239: \caption{Top 10 Ingredients
    }\label{f:moroccan_10_most_used_ingredients}
242: \begin{figure}[!ht]
243: \centering\includegraphics[width=\columnwidth]{images/russian_10_
most_used_ingredients.png}
244: \caption{Top 10 Ingredients
    }\label{f:russian_10_most_used_ingredients}
247: \begin{figure}[!ht]
248: \centering\includegraphics[width=\columnwidth]{images/southern_us
_10_most_used_ingredients.png}
249: \caption{Top 10 Ingredients
    }\label{f:southern_us_10_most_used_ingredients}
252: \begin{figure}[!ht]
253: \centering\includegraphics[width=\columnwidth]{images/spanish_10_
most_used_ingredients.png}
254: \caption{Top 10 Ingredients
    }\label{f:spanish_10_most_used_ingredients}
257: \begin{figure}[!ht]
258: \centering\includegraphics[width=\columnwidth]{images/thai_10_mos
t_used_ingredients.png}
259: \caption{Top 10 Ingredients
    }\label{f:thai_10_most_used_ingredients}
262: \begin{figure}[!ht]
263: \centering\includegraphics[width=\columnwidth]{images/vietnamese_
10_most_used_ingredients.png}
264: \caption{Top 10 Ingredients
    }\label{f:vietnamese_10_most_used_ingredients}
278: Figure \ref{f:ingredient_modularity} shows ingredient cluster of
more than 1000 nodes.
279: \begin{figure}[!ht]
280: \centering\includegraphics[width=\columnwidth]{images/ingredient_

```

```
    modularity.png}
281: \caption{Ingredient Cluster }\label{f:ingredient_modularity}
284: Figure \ref{f:ingredient_modularity100} shows ingredient cluster
     of around 100 nodes.
285: \begin{figure}[!ht]
286: \centering\includegraphics[width=\columnwidth]{images/ingredient_
     modularity100.png}
287: \caption{ingredient Cluster 100 Nodes
     }\label{f:ingredient_modularity100}
```

figures 27

tables 1

includegraphics 25

labels 28

refs 9

floats 28

True : ref check passed: (refs >= figures + tables)

True : label check passed: (refs >= figures + tables)

True : include graphics passed: (figures >= includegraphics)

False : check if all figures are refered to: (refs >= labels)

Label/ref check

passed: True

When using figures use columnwidth

[width=1.0\columnwidth]

do not change the number to a smaller fraction

find textwidth

passed: True

below_check

bibtex

label errors

bibtex errors

This is BibTeX, Version 0.99d (TeX Live 2016)

The top-level auxiliary file: report.aux

The style file: ACM-Reference-Format.bst

Database file #1: report.bib

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barabas Istvan", while executing---line 6169 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barabas Istvan", while executing---line 6169 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barabas Istvan", while executing---line 6169 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barabas Istvan", while executing---line 6169 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barabas Istvan", while executing---line 6169 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barabas Istvan", while executing---line 6169 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barabas Istvan", while executing---line 6261 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barabas Istvan", while executing---line 6261 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barabas Istvan", while executing---line 6261 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barabas Istvan", while executing---line 6457 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barabas Istvan", while executing---line 6457 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barabas Istvan", while executing---line 6457 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barabas Istvan", while executing---line 6457 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barabas Istvan", while executing---line 6457 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barabas Istvan", while executing---line 6457 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Ahn Yong-Yeol, Ahnert Sebastian E., Bagrow James P., Barabas Istvan", while executing---line 6457 of file ACM-Reference-Format.bst

Warning--page numbers missing in both pages and numpages fields in Ahn2011

Warning--no journal in Chatterjee2016

(There were 16 error messages)

bibtex_empty_fields

entries in general should not be empty in bibtex

find ""

passed: True

ascii

=====
The following tests are optional
=====

Tip: newlines can often be replaced just by an empty line

find newline

passed: True

cites should have a space before \cite{} but not before the {

find cite {

passed: True

How Big Data will Help Improve People's Health Worldwide

Paul Marks

Indiana University

Online Student

Shepherdsville, Kentucky 40165

pcmarks@iu.edu

ABSTRACT

Aside from people changing their habits, big data analytics may hold the best possibility for the improvement of worldwide health. It will enable the ability to correctly diagnose patients more quickly, even when the patients may not be able to be physically seen by a provider. It will be used to create treatment plans specific to not only an illness, but to the patient's overall health condition and history, demographics, environment, and access to resources. While it may not solve the problem of everyone not having access to the best of care, it can help to make sure everyone can get the best care possible for them. This paper explores the ways in which big data is evolving in the field of healthcare to make these possibilities become realities and looks at some of the social concerns which could hold it back.

KEYWORDS

i523, hid327, healthcare, patient treatment, genomics, diagnosis

1 INTRODUCTION

There have been many advances in big data analytics over the last several years. More and more data is able to be processed in a shorter amount of time. There are also many new sources of data. Data is not what big data is about though. It is about taking data and turning it into information that can be useful. The application of big data can vary, but very few may be more important than the ability to use data for the betterment of people's health across the globe. This is one way in which data science can make a substantial contribution to humanity.

Making this a reality is not, nor will be, a simple task. Health data itself requires the proper handling of the information as it is very sensitive. On one hand people have a right to privacy. On the other, if data is kept isolated, not combined with records from other people, then this limits the ability to gather insight and find breakthroughs. The key is to ensure privacy, but keep the integrity and relationships of the data in order to preserve privacy while gaining insight. The insight gained has endless possibilities.

One issue facing the medical profession today is a lack of trained professionals. The number of patients per healthcare worker around the world can vary from more than six per 1,000 people to less than one half per 1,000[43]. It is easy to see how this one fact greatly impacts the expected lifespan of people. But what if a patient could be examined, diagnosed, and have access to a treatment plan without a human doctor needed? It may sound futuristic, but the technology is being implemented today thanks in part to data analytics.

The impact of big data on healthcare doesn't stop there. The cost of treating 5 percent of the most chronic conditions can consume up

to 50 percent of the money spent on healthcare[42]. One reason for this is prevention, diagnosis, and treatment plans are not optimized. There is not one way to help patients avoid chronic conditions. It is based on many inputs depending on the person, their environment, and other factors. These same aspects impact the effectiveness of treatment plans as well. One size does not fit all. Through analytics many factors are being analyzed along with the results of prior plans to determine which methods would be the most effective. Avoiding a chronic condition not only saves money, but extends a patient's life and improves the quality of their life.

The ability to take many factors into account for a patient goes well beyond chronic conditions. Genomic technology is progressing which is allowing for a person's individual genome to be one of the inputs. Each person on earth has their own specific genome with billions of combinations, some of which directly impact their health and susceptibility to illnesses. Through big data analytics, this type of analysis may one day be commonplace like taking blood pressure and other vital statistics into account.

The discovery of new drugs and how they can be used to treat people is being sped up by the power of big data techniques. Drug research requires an immense amount of information to be correlated and processed. Big data is helping to speed this up and even helps speed up clinical trials by matching the right set of circumstances to provide viable results.

Progress does not always come without drawbacks, and big data analytics in healthcare is no exception.

2 HANDLING THE DATA

2.1 Security

Any use of healthcare data must take into account the ability to protect the data. Therefore a brief understanding of the task must be addressed. Healthcare information usually has two forms of protected information: Personally Identifiable Information (PII) and Protected Health Information (PHI). In order to be able to keep data with this type of information you must follow very strict rules on safeguarding it. The best known regulations are based on the Health Insurance Portability and Accountability Act (HIPAA) of 1996. Among the governmental standards to comply with HIPAA are the Security Control Assessment[19] and Defense Information Systems Agency's Security Technical Implementation Guides[3]. These types of requirements can be costly and require constant changes to remain secure.

Even with the ability to secure the data properly, any company wishing to obtain data must have an approved reason to get the information or the approval of the patients involved. Obtaining approval from each patient in a big data application is not practical.

Data is needed from too many people to obtain approval for each of them. A common way to handle this is through de-identification.

De-identification is the ability to alter the data in such a way that you cannot link health information to a person or identify individuals in the data. However, in order for the data to be useful for analysis it cannot be changed randomly so the links between certain data elements from record to record is lost. For instance, a diagnosis of a specific cancer in a patient must still be able to be linked to treatment data, x-rays, blood tests, etc. from that patient. In other words, de-identification has to be done in such a way that the data integrity remains in place, but the individual's identity is protected. This can become complicated because data elements such as age, sex, and geographical location are important.

Fortunately there are software solutions to assist in the de-identification of medical information. The software is broken into two categories: structured data and free-form text. De-identification of structured data is generally easier. The data has a known set of fields of which the ones which can identify a person and their health are known. These fields are added to the software and algorithms are run against them. The resultant data is useful for analysis, but the identity of any individual is safe. This is because the algorithm changes data in such a manner that it protects the person and the data integrity. Examples of tools in this arena include PARAT from Privacy Analytics, Inc., mu-Argus from the Netherlands national statistical agency, Cornell Anonymization Toolkit (CAT), an anonymization toolkit from the University of Texas at Dallas, sdcMirco from r-project.org[26]. Commercial tools like Privacy Analytics Eclipse claim to de-identify 10 million records per day from a variety of sources[50].

Unstructured data is more complex. The data which needs to be de-identified can be located anywhere within the dataset. This includes the text or metadata attached to images such as x-rays. Vital clinical, diagnosis, treatment, and other medical information is also included throughout unstructured data. Not being able to identify all PHI and PII can cause privacy concerns. Not linking all the correct data together reduces data integrity which reduces the usefulness of the data being studied.

Being able to properly de-identify and link unstructured data is being studied and refined. There are challenges for solutions to the problem. Informatics for Integrating Biology and the Bedside[25] has held challenges to help further solutions for this problem. The most recent was held in 2014. Track 1 of this challenge noted that "Removing protected health information (PHI) is a critical step in making medical records accessible to more people, yet it is a very difficult and nuanced"[25]. The ability to properly de-identify the data is rooted in the ability for the software to perform natural language processing. The focus of the challenge was all eighteen HIPAA defined PHI types[36]. While not as mainstream as de-identifying structured data, the ability to de-identify unstructured data will continue to progress and be solved through commercially available products over time.

2.2 Data Sharing

There are many sources of healthcare data. This is a major hurdle as the data is in different systems which are governed by different entities and used for purposes[33]. Data is stored in claims systems,

clinical settings, pharmacies, and others. It is stored in different formats. These sources may not contain similar key data that allows it to be easily brought together. Individual patients usually have a single provider who is their primary insurer. This data is usually in standard formats. However the same patients may have many providers of care using different systems. While most providers leverage electronic health records, these systems can contain many free-form text fields, images, and other types of fields. These data sets contain a wealth of information, but they are missing data which could be vital such as social, environmental, and community data. Other sources of data which could be useful are habits which people store on themselves such as food and activity tracking they may enter into any number of online applications[12].

While more data is being collected, there are still barriers to sharing it. There are the security and privacy concerns discussed earlier, but also the costs and who pays for them which must be addressed. There are tools and strategies being worked on in the industry to make sharing data across disparate systems possible. So far a widely adopted solution has not emerged[11]. Until such time that it does, data analytics in healthcare will be hampered.

3 BIG DATA IN A CLINICAL SETTING

Being a doctor can be like being a human big data machine at times. They take in many variables, process it against the history of information they have, and come to some sort of conclusion. In many cases there are multiple diagnosis that can be made. In fact sometimes there a lot of diagnosis that can be made. Unfortunately while much of the work is very scientific it does not mean that coming to a conclusion is a precise science.

Different doctors have different backgrounds. They have seen different patients, seen different diseases, studied at different locations, and read different literature. In short, their diagnosis is based off of their experiences. Unfortunately experiences are a form of bias. It is not that someone is doing this on purpose for the betterment or detriment of someone, but it is how our brains are wired. Physicians are not immune to this and it can affect the ability to treat all patients and conditions equally or appropriately[13]. When set up correctly and fine-tuned over time, data analytics can minimize biases.

3.1 Big Data as a Physician Assistant

What if each doctor had the collective knowledge of others? That could make for better and more accurate diagnoses around the globe. A doctor in the United States would have the knowledge of thousands of years of alternative medicine which may only be taught in schools in the far east. Not only is it possible, but big data is making it happen today through technologies such as IBM's Watson.

3.2 IBM's Watson Health

One of the challenges facing doctors today is the ability to keep up with changes in healthcare. Even doctors who specialize in a field cannot keep up with the amount of information that is being published. One estimate is that 8,000 medical journal articles are published each day[57]. This makes medicine a good fit for big

data. Watson Health, IBM's name for their cognitive supercomputer focused on healthcare, is able to ingest millions of pages of information in seconds. This information becomes part of the core information Watson has at its disposal as it assists clinicians by offering recommendations for them to consider. In this way Watson is not the final decision maker, but helps doctors be better at what they do[30].

While Watson is delegated to a physician's assistant currently, it may not always be so. In order to test how accurate it is, IBM tried it on 1,000 patients. In this test Watson and the attending physician agreed 99 percent of the time. In fact, in 30 percent of the cases Watson offered pathways which the physician had not considered. Armed with information like this IBM believes that computer cognitive thinking will be mainstream in the next ten years[57]. Because of advances in other technology areas have been progressing so quickly, it is hard to disagree with them. For instance, computers are now able to instantaneously make decisions that seemed unimaginable just a few years ago which as lead to the realization of autonomous driving vehicles. The question may not be the technology, but if people will accept a diagnosis from a computer program such as Watson.

Watson was also tested to see how examining a patient's entire genome would be more beneficial than simply running a panel which focuses on a limited number of areas most commonly known to be related to the cancer a patient may be experiencing. While the cost of and speed of sequencing a person's genome has been reduced, there is still a lot of work to using this data for a specific diagnosis and treatment plan. Both Watson and team from the New York Genome Center analyzed a patient's genome. Each of them was able to identify gene mutations which would have pointed to a clinical trial or drug which may have been a better match than the treatment the patient received. The difference being that it took the team of physicians approximately 160 hours to come to their conclusion. Watson provided its results in 10 minutes[56]. While not perfect, Watson adds another tool doctors can leverage which would allow them to better diagnose and treat patients.

How does Watson do it? It is actually very similar to how a human doctor works. The patient's symptoms and other information is made available to Watson. From there it deduces the relevant elements and leverages any background information it may have such as patient and family history, labs, x-rays, and other test results. It then accesses other sources of information it has accumulated over time: treatment guidelines, relevant articles and studies, and potentially information from other patients similar to this patient. Watson develops hypotheses and runs them through a process to test its hypotheses and provide a confidence score for each. Watson then provides its recommended treatment options with its confidence rating to the physician[20].

One advantage of Watson, or any such system, is that every time it is used that patient is getting all of its collective knowledge. Today when a patient see a physician they are diagnosed by that physician and maybe one or two other people generally from the same office. However as Watson gets *trained* by specialists in such fields as Oncology, every doctor who uses Watson's assistance becomes as or more knowledgeable than the collective group. This means that each doctor is providing top of the field care even if they are being seen nowhere near a facility that is considered as the best

world[37]. A patient in a country not seen as having world-class healthcare can get diagnosed as if they were at the Sloan-Kettering Cancer Center. It also means that a patient who may be seeing a specialist in one area may be diagnosed with an ailment outside of their field. This can save time in receiving the appropriate diagnosis and subsequent treatment which gives patients the best chance for recovery.

There are obstacles to making Watson available worldwide and that is the ability to understand different languages. Watson knows English, Brazilian Portuguese, Japanese, and Spanish and is learning others. As an example, IBM, the Cleveland Clinic, and Mubadala are teaming up and are building a hospital in the Middle East. The Cleveland Clinic is already a user of Watson Health and is expected to leverage that in the new facility as many chronic conditions in the United States are present in the Middle East as well. To prepare for this, IBM is teaching Watson Arabic[60]. As Watson learns more languages it will be able to be leveraged in areas around the world which that language is spoken allowing for those populations to advance their healthcare knowledge.

Another advantage that Watson has over human physicians is that it never forgets. Even doctors who try to keep up with changes in healthcare, they will never be able to remember information as precisely as Watson. And Watson is also consistent. A single doctor may be mostly consistent, but different doctors will provide different diagnoses given the same input. Watson will not unless it is programmed differently or new knowledge is ingested which can create a more accurate diagnosis. It also does not have bad days, get tired, and is available 24x7x365. Watson's incremental costs, the cost of using it for one or one million patients, is low. IBM has spent billions on it and is continuing to invest, but those costs will be spread out as usage goes up thus making Watson cheaper over time[20].

3.3 Implementing Big Data Diagnostic Systems

Leveraging such technologies can be implemented in various ways. The easiest way is to look at them as another tool in a physicians' tool chest. Once fully implemented the inclusion of big data assisted technologies will be seamless. Clinical information is being collected digitally on an increasing basis. As vital signs, x-rays, diagnostic images, lab results, and even discussions with the patients are collected digitally they will become part of the patient's electronic health record and the overall collective knowledge base. Watson or other software could provide insight to the physician. It may be present a collection of diagnoses scored in likelihood based on the evidence collected so far[29]. It could provide recommendations for next steps or information which could lead to a more complete recommendation.

The idea behind such a system, Watson or any similar tool, is to make physicians better through more accurate diagnosis. It allows for the use of big data without removing the human aspect of medicine. This will help to begin to include the big data and computer health diagnosis to patients who would otherwise not be open to it. For many people their relationship with their doctor is personal. They discuss items with their doctor they do not discuss with anyone else. They may not trust a computer with their health[29]. A non-caring, non-breathing inanimate object cannot

be trusted with something so human. In this implementation a doctor would still be there providing the personal interaction with the patient and thus providing them with the best care including the collective knowledge of the system.

3.4 Replacing Doctors for Routine Visits

Having a doctor meet with a patient initially may not always be required. The ability for big data to leverage healthcare data could lead to helping alleviate the shortage of doctors and nurses in the United States and around the world. Worldwide there is an estimated shortage of skilled health professionals of 17.4 million of which 2.6 million are doctors. The problem does not get much better over time as the estimate for 2030 is over 14 million[45]. It takes a lot of time and money for a student to achieve the level of knowledge to fill these positions. Unless the students are already in the pipeline then there is not a good response to the problem. People cannot switch careers and be a doctor or a nurse in twelve months or some short time-frame.

Adding new big data doctors is simple. It is mostly a hardware problem. Buy the right equipment, install the right software, train the staff, and Dr. Data can see patients. Leveraging automated machines to take vital signs will free up time for staff[15] similar to how checking out via automated tellers at the grocery store has reduced the number of cashiers and baggers needed. A physical office offering virtual doctor's visits could be staffed with people trained on the technology more than medical professionals. They would be there to help make sure that people are using the machines correctly and to wipe down equipment after a patient has used it. A nurse would be there in case certain patients are unable to use the equipment and their information must be taken manually. They could also be there to take blood samples which would be processed by automated machines and included in the patient's profile.

Automated diagnosis systems are in use today in a limited basis. In the United Kingdom the National Health Service has approved the use of Your.MD (an AI powered mobile app) for diagnosis. When people are comfortable using a technology like this it limits the number of more basic cases a doctor has to see and allows them to concentrate on more difficult tasks. Another tool, Ada, learns a user's history, provides an assessment, and adds an option to contact an actual doctor if needed. Babylon Health takes it one step further by adding follow-ups with users to see how they are doing and can even set up a video consultation with a live general practitioner if needed[15].

4 LIMITING EPIDEMICS

Incorporating big data analytics into the healthcare environment has the ability to limit the spread of disease by taking current circumstances outside of the immediate patient into account. In a linked system data from other local, regional, national, and global patients can be leveraged. Are there other patients presenting similar circumstances? Did the other patients provide more details or mention something slightly different? Taking this into account may help to diagnose a specific person and to identify an outbreak of something. Is a disease spreading? Did patients come from a similar location such as a building? By being able to correlate this information immediately there is the potential to stop an outbreak from

spreading thus saving an untold number of patients from pain and suffering and saving healthcare dollars by not having to treat more patients. Epidemics have an economic impact at many levels including "the micro (individual and household), meso (establishment, village or city) and macro (national and international)"[46].

5 INSURANCE

The option of having fully automated doctors' visits could alter the insurance market as well. Health insurance is about numbers. Actuaries spend time estimating the health of the consumers they cover and many other factors to determine what premium rates to set[39]. Insurers make a profit by taking in more money than the costs to administrate the plans and the cost of paying for claims combined. To reduce the costs of claims they set predetermined prices for services rendered by hospitals, physicians, and sometimes pharmacy companies. The lower they can drive the cost of the claims they cover the less they charge or the more money they make. Charging less can result in making more money as well as more people may choose to purchase coverage from that insurer.

By creating options for autonomous doctor's visits or telemedicine an insurance company could save money. The more methods can be deployed which can reduce overall healthcare costs, the less people will pay. There are multiple ways in which this can be included to reduce health insurance premiums, a high cost item for most people in the United States and other countries. Insurers can work with healthcare providers who leverage this technology to create a reimbursement policy that is less for services such as telemedicine[34]. They could also offer plans to potential customers which require basic treatments to take place with autonomous or telemedicine options before they go to a doctor's office. This would offer an economic advantage to people which in turn can not only lower costs, but help to increase the adoption of new technologies.

Such a system is not for everyone or every condition. The idea is not to replace all doctor's visits, but to allow those who are comfortable to take advantage of lower cost coverage. It will encourage younger people to keep insurance if it is made more affordable. Currently the highest rate of not having insurance in the United States is when someone can no longer be covered as part of their parents plan, starting around the age of 25[6].

6 PORTABILITY

More importantly than lowering the cost of healthcare or making seeing a doctor more convenient is the ability to make exceptional healthcare available almost anywhere. Big data using an automated doctor can have an impact on under-served areas the like of which no one has ever seen. Today there are people who do not have access to healthcare of any kind. When they get sick they may not have a place to turn. In developed countries the number of patients per doctor is generally in the low hundreds. In poor, *third world countries* the number of patients per doctor is in the thousands or tens of thousands[27]. There are people who try to help, such as Doctors Without Borders, by making visits to these areas to provide some support but it does not reach a level anywhere near what people in some countries have available to them. If each doctor could multiply their impact with technology then the under-served would be helped more. As technology advances so people could

be seen by experts without one being physically present then even more people could be seen.

7 PATIENT DATA COLLECTION

7.1 Actual Data vs. Circumstantial Insight

The more valid data which can be collected on patients the better big data will be able to help improve treatment for people around the world. The more accurate the data, the more accurate the analysis and results will be. Fortunately technology is helping in this area as well. Many people around the world have access to devices which monitor different aspects of our daily lives. Hundreds of millions of people around the world have purchased wearable devices, many of which can be used to monitor activity and inactivity[2]. By the end of next year it is expected that over one-third of people in the world will own a smart phone which can also track this type of activity[1]. While they are not seen as a medical device, they can help to track activity which is useful for diagnosis and treatment. They are another input into the data about a patient which can be used to more accurately gather information. Today doctors rely on a patient to answer questions about their level of activity. With such a devices they can get a more accurate picture.

These devices are useful for more than just activity levels. They also provide insight into areas of people's lives they are not really able to answer accurately such as how they sleep. Many people may sleep they sleep well or not so well, but in fact they are basing this more on how they feel than how much rest and how good of rest they get. Activity trackers are able to track sleep patterns as well. They actively monitor your inactivity. When used correctly a wearer pushes a button to indicate they are going to sleep and when they get up in the morning. The monitor is then able to track how long it takes for someone to get into a motionless/restful state. It continues to track them throughout the night recording if they move around, get up, etc. Getting good sleep is a key element of maintaining overall health[54].

More advanced features of activity trackers include the ability to monitor vital signs like heart rates. They can be extremely important to a diagnosis providing input similar to a mini stress test. This is especially true if a person exercises, such as during jogging. The device can monitor how far a person is moving and their associated heart-rate. By gathering this information, the data can be fed into patient's profile when they visit a doctor (virtually or physically) instead of having to wait for a patient to get a test done and receive that feedback. Shortening the time to collect data and accurately analyze the patient can be the difference between life and death.

One aspect of activity trackers which must be noted is their accuracy and consistency. This is something big data can help with as well. Steps from person to person are not of consistent stride, tracker accuracy changes from device to device, heart rate monitors vary, and sleep are not be tracked similarly across all products and types of activities[55]. Big data can help normalize this input so that it can become a reliable input. Analysis has been done on different monitors to see how accurate they are. In order to bring them into health analysis more tests can be performed to get an accurate picture of how the devices correlate to the actual distances walked and level of sleep.

Activity trackers are only the beginning. *Wearable technology* is an expanding field which is enhancing the collection of passive data. Sensors are being built into clothing which track more accurately and include more types of data[21]. This includes information like breathing rate and muscle activity. They not only collect more types of data, but can wirelessly transmit the data via Bluetooth[32]. This means they can create a more accurate picture based on electronic data which can be used as an input. The more this type of technology becomes commonplace, the more data which can be fed into a patient's health record and the collection of health information.

7.2 Follow-Up Visits

All of these devices also have the ability to not only be used in diagnosis, but in the monitoring of treatment plans. Is the patient exercising as they say they are? Is a medicine or other corrective action helping them to lower their heart rate or get more restful sleep? It can also help to notify the patient or doctor when they are exceeding a prescribed level of respiration or heart rate. This can trigger an alert for a patient if they are at risk or even that they may need to seek treatment. These levels will not only be set based on standards, but patient specific information[5]. They can also take into account the environment the person is in. Are they in a hot location or one with high allergy levels which could negatively impact them? This is what separates the treatment plans of today with those of tomorrow. Use the technology to more accurately collect data on the patient, use it to create a diagnosis, monitor the patient using the technology, feed that data back into the patient's health record, and adjust as needed based on factual information.

Beyond the use of commercially available monitoring systems, there are devices which collect data similar to the information collected by a physician. Simple systems such as a blood pressure monitors are common. Many other pieces of equipment can be prescribed by a physician for home monitoring. These systems not only collect information, but are able to digitally transmit the data so that it can be automatically analyzed with other sources of information. A patient will get feedback without having to visit a doctor[5]. This helps to close another gap in healthcare which affect many people: not following up with their doctor. Missing these visits can negatively impact the patient. By easing the ability to be monitored, automating the data collection, and instantly analyzing that data will lead to better overall prognosis.

Big data will also help to change people's habits. By using the data collected a picture of potential outcomes can be made for a patient to contemplate. Instead of generalities, patients will receive advice based on their medical history, other patients like them, treatment plans, and other inputs based on the variables specific to the patient's circumstances. It can show a patient how they impact their recovery based on what they are doing or not doing. For instance if they miss taking their medicines on time, do not lose weight, continue to smoke, or whatever other variables they are in control of and how it affects their specific recovery or health status. Showing them in advance may give them the motivation they need to follow the plan more closely. Throughout their treatment the model can be updated based on the patient's actual adherence to the plan. This provides another feedback loop for the patient to

course correct their habits if they have not been following it as outlined[5].

Not only will big data help to diagnosis patients more accurately, but it will also allow for the customization of treatment plans at levels not available today. Instead of relying on more general treatment plans, patients will have their plans customized by their specific set of circumstances. Demographic information about the patient will be used to compare to historical plans and outcomes of patients most closely related to their characteristics. This includes not only the patients themselves, but the environments they live in. Pollution, weather, access to ongoing care, income (the patient may have to work whereas a long period of rest would be better) and other circumstances will be variables which may not be controllable by the patient, but can be used to help treat them. The plan will not necessarily be the best treatment course, not everyone has the access to the best care or the ability to abide by it, but will instead be the best plan for them and their circumstances. Each patient will be able to maximize their chances of recovering or otherwise leading the most normal life possible.

8 ACCESS TO HEALTHCARE

It is estimated that over 400 million people do not have access to basic healthcare around the world and others are forced into extreme poverty because of what they pay for healthcare[47]. Through tools referred to as telemedicine, these numbers can be lowered. Telemedicine itself is the ability for people to get evaluated, diagnosed, and treated while the physician is not located where they are. When combined with a mobile diagnostic unit a patient can get similar care to someone who is seen at a clinic[52]. As advances in automated solutions such as IBM Watson evolve, there could be a day when these remote services are performed in very remote areas where communication with a physician would be technically challenging.

9 COST SAVINGS

Another reason why big data will be helping with healthcare more and more in the future is the most basic of reasons: Economics. Regardless of the country or political system, there is always an economic element which must be addressed. No country, no system has an endless supply of any services or funds. Because of that ideas which make the most economic sense have a better chance to be adopted. The economics of automating healthcare with big data analytics will reach a tipping point as time progresses.

Simply put, healthcare is getting more and more expensive every year and computing resources become cheaper every year. Worldwide the per capita expense of healthcare has risen from \$661 to \$1,059 (numbers in United States Dollars or USD) in the last 10 years[22]. That is a 60.21% increase in one decade. The average per capita may seem low to some but that is due to it being worldwide number. Many countries spend almost nothing on healthcare per capita while others spend thousands. For instance, in 2004 Vietnam spent \$30 USD per capita and \$142 USD in 2014. This is a 373% increase, but in total dollars it is still a fraction of \$6,369 (2004) and \$9,403 spent in the United States[22].

In contrast to this the cost of computing power has decreased year over year. Computer power is not as straightforward to analyze, but cost trends are easily seen. One way is to compare the cost using a baseline year and showing other years as a percentage of the cost of the baseline. Using December of 1997 as a baseline (100) of cost for computers, the cost of computers and peripherals in January 2004 had dropped to 16.2. In other words, to get the same amount of computer power in 2004 you only had to spend 16.2 cents for every dollar spent in December of 1997. By January of 2014 it had dropped to 4.9. Comparing the 2004 and 2014 numbers, the same ones used above for healthcare spending, the cost of computing had been reduced by 69.75%[41].

A specific component when it comes to big data is the cost of storage. The decline in the cost of storage over time is staggering. In the early 1980's the cost of one gigabyte (GB) of storage was in the hundreds of thousands of dollars. Using early 2004 as our baseline the cost for one GB of storage had dropped to just under \$2.00. By 2014 the cost had declined further to between three and four cents per GB[31]. The speed at which the data can now be retrieved as compared to 2004 is like comparing the speed of light to the speed of sound. Today's storage units are that much faster.

Using this data one can see that as we are able to leverage big data solutions to provide better healthcare we can also begin to slow the incline of healthcare costs and then lower the cost of healthcare over time. Adding a new virtual doctor will not take years of schooling which can cost hundreds of thousands of dollars in some countries. It will be the cost of some piece of common technology and a licensing fee for the software. As with most everything technology based, increasing the volume decreases the cost. So as more and more virtual doctors are brought online the cost of each will decrease.

10 CHRONIC CONDITIONS

Chronic conditions are ones that "are preventable, and frequently manageable through early detection, improved diet, exercise, and treatment therapy"[59]. They are also very expensive to manage and treat. Worldwide in 2010 the total cost of heart disease alone was \$863 billion dollars (USD) and is expected to be \$1.44 trillion by 2030. Between 2011 and 2031 the cost of the top five chronic diseases (cancer, diabetes, mental illness, heart disease, and respiratory disease) will cost \$47 trillion (USD) globally[28].

It is not only the economic impact of chronic diseases that make them a target for big data analysis. Chronic diseases reduce people's quality of life. This cannot be factored into simple terms such as money. Chronic diseases are the cause of 60 percent of deaths worldwide[44]. In a 2002 study it was estimated that 84 percent of deaths were due to chronic diseases in Europe and Central Asia[14]. Chronic disease is so prevalent and impactful to people's lives that it has been labeled as "the most expensive, fastest growing, and most intricate problem facing healthcare providers in every nation on earth[9]." With data like this it is easy to see why advances in chronic diseases is important. The question becomes how do fight them.

10.1 Prevention

The best way to fight chronic disease is to never have one in the first place. The best way to reduce the number of people who get a chronic condition is early intervention. Big data analytics can be used to help with population health management when it comes to chronic diseases. That is by identifying those who are at a high risk of getting one of these costly, harmful conditions[9]. The ability to leverage big data in prevention is a two part process. First risk factors which are modifiable must be identified and then interventions need to be created which will have an impact on changing the factors[7].

Modifiable is the key word in the first aspect of using big data. A key to fighting many chronic conditions is for people to stop behaviors such as smoking, to eat healthier, and to exercise more. However, if it was as easy as letting people know this then there would be a lot less chronic disease already. Big data can take many factors into account and help to create a more precise message for a people with specific risk elements. For instance instead of telling a patient to eat more nutritious foods, by leveraging elements of their specific health factors a doctor can recommend more precise information such as asking them to include a particular dietary nutrient[7]. Big data can also help with the timing of the message. In a survey patients wanted more information from analytics that would have warned them before they developed a chronic condition[10]. When someone is presented with more personalized information (they are on a path and about to reach a point of no return) vs. general (a healthy lifestyle may prevent you having issues years down the road) they are more compelled to heed that information and act upon it.

Newer technologies outside of a clinical setting are helping to add to the data available to analyze and care for patients. Combining data from a patient's activity monitor, fitness tracking website, or food logs into their plan helps to create a feedback cycle for the healthcare provider. Many applications track food by scanning the USB code from the package. Making it simple helps to get people to do things. The easier it is, the more likely they are to do it. Taking this data and combining it with clinical data such as blood labs and vital statistics can show a patient how they are directly impacting their health in a positive or negative manner. It changes the conversation from more of a public service announcement general message to one unique to them.

A special sub-section of patients are very high-cost patients. In the United States there are roughly five percent of patients who account for almost 50 percent of healthcare spending[8]. Identifying these patients and creating intervention plans that work can have an enormous impact on their lives and the cost of healthcare overall. Patients with seemingly similar risk factors may have very different prognoses. Obvious factors such as age, weight, sex, and vital statistics may be the same. In order for big data to help identify the five percent more data is needed. Including mental health data, genetic information, socioeconomic, marital status, living conditions, and even cultural factors into the analysis will allow for better predictions and better ways to intervene which will lead to better outcomes[8].

10.2 Management

Even with the best of preventive measures there will still be too many people with chronic conditions for years and decades to come. Approximately 25 percent of people with chronic conditions have restrictions in what tasks they can perform for themselves, at work, or at school[24]. Because of this big data must also be leveraged to help manage those with chronic conditions. Managing it is not only based on cost, but helping them to live a better quality of life with less trips to the doctors and less admissions to a hospital. Data analytics can help to customize treatment plans to the circumstances of each patient. It can see patterns in patient's data and help to determine better follow up schedules. This could mean the difference between a visit with their doctor or a costly hospitalization[4].

Part of the solution for using big data to help tackle chronic conditions is leveraging new sources of information from technologies such as wearables. As mentioned earlier they allow for real-time data to be collected, combined with other sources of information including that of other patients, and provide better treatment plans for patients. Historically the medical profession had to rely on subjective input from patients when they came in for a visit. How often were they active, did they log information like their heart rate and blood pressure when they should have. With some wearables all this information and more is gathered in real-time and can trigger an alert to a care management professional[24]. This means that changes can be made when they are needed and the patient can get immediate attention, not days or weeks later.

Another issue with chronic care for providers is that patients may have multiple conditions. They may be overweight, have diabetes, and hypertension. This leads a patient to having multiple doctors each working on a specific condition, but no real coordination across the diseases. A treatment for one condition may have a negative impact on the patient because of treatment or drugs prescribed for another condition. And this situation is not unique as there are many patients suffering from the same conditions simultaneously. Big data analytics can bridge this gap. By combining data from multiple sources, patients, and treatments physicians can create a customized treatment plan for a patient to combat all three illnesses in the best manner without adverse interactions[58].

The result of this is that big data can help people see that treatments are tailored to them and are making a difference. Data analytics allows for patient-centric care, not disease-centric care. Patient managers would work with patients providing details on their plan, their results, and will be able to show patients how the care plan affects their quality of life. It can help to create a healthcare environment "where patients are not only engaged in time but see improved health results at affordable costs"[53].

11 GENOMICS (PERSONALIZED HEALTHCARE)

The field of Genomics is investigating how healthcare can be more personal. How diagnosis and treatment plans will be based on a specific person instead of how the factors or ailment is normally seen and treated in the general population. This is essential work because in the United States up to 47 percent of the cost of healthcare is spent on interventions that do not provide any value. While the actual percentage may vary in other countries, this is a worldwide

problem[30]. Any easy way to understand the difference is over the counter medicine. Generally speaking the instructions on a bottle are broken down into children and adults. Following the directions adults will take the same amount of medicine regardless of their age, weight, or overall health.

Genomics aims to make medicine very specific to an individual by breaking down each person's genome. This is only possible through big data as a single person's genome produces a lot of data because it has up to 25,000 genes with three million base pairs. One human genome can produce up to 100 gigabytes of data[18]. And the information from one individual is not what is required for personalized health. It requires genomes from many individuals. The more data available, the better the analysis can be on similarities between people and how they may react to certain treatments. This multiplies 100 gigabytes by thousands, then millions, then hundreds of millions.

Through advances in technology such analysis is possible. In 2003 the first human genome was sequenced. It was only after 13 years and approximately \$3 billion dollars. By 2015 the same work can be done in a few hours at a cost of just over \$1,000[38]. This means that more and more people can have their genomes sequenced and used for analysis and personalized diagnosis and treatment of diseases. As more and more genomes are collected and analyzed treatment can be based on their personal genome and their family traits through family based analysis. This analysis lets doctors see how people may have inherited a propensity to be susceptible to certain diseases based on mutations in their genomes. In addition, through population based analysis environmental and cultural factors can be included. It is estimated that by 2025 over 100 million genomes could be sequenced[23]. Analyzing the details of the building blocks of so many individuals will be a big data challenge which can have an enormous impact on healthcare.

12 DRUG DISCOVERY

Discovering new drugs which can help us live a better life is something like finding a needle in a haystack. Large libraries of molecules have to be examined "against millions of data points spanning chemical, biological, and clinical databases"[16]. This is done looking for relationships between diseases and drugs to see if a particular drug could be used to treat the disease. While the process is not new, this work is the basis of many new drug discoveries, the ability of current big data techniques speeds up the process allowing for drugs to be discovered more quickly[16].

One of the reasons for it being so complicated goes back to the discussion of the human genome: each person is a unique individual. If you have seen a commercial or advertisement for a prescription drug there is always a list, sometimes a very long list, of possible side effects. These are adverse impacts which can range from minor annoyances to death. Part of the challenge of drug discovery is attempting to identify and quantify the impact a drug may have on people. To speed this process healthcare big data has developed solutions such as array-based technologies which are purpose built to combinatorial problems. This lets researchers find patterns in the data more quickly, speeding up the overall process[17].

Once a drug is thought to have a potential positive use it must go through a testing phase before it is approved for use. This can

be a long process which has successes and failures. Big data is being used for "the improvement of clinical trial designs (e.g., endpoints, inclusion/exclusion criteria, etc.)"[35]. This not only allows for potentially a quicker time to market, and thus the ability to help people sooner, but a cost savings without paying for trials which do not produce viable results.

13 INCENTIVES FOR ADOPTION

In the end many of the advances will only be possible if people accept them. So how can this number be influenced? The most logical way to do so is to make the adoption of these advances financially beneficial. People are more willing to take a chance when they can see a hard benefit. Insurance premiums can help to drive this and provide an immediate benefit. Plans could be offered in which a person's primary care is provided by a big data doctor. People would have to consent to having their information stored electronically and compared against the data sets. Visits to physical doctors including for second opinions would be limited. They could even have different reimbursement models similar to preventive tests. Most insurance today covers preventive services at 100 percent and are not subject to a deductible. Electronic visits could be treated similarly. They could be covered at 100 percent, or some number higher than regular doctors visits, and may or may not be subject to a deductible. Leveraging these types of incentives will help to promote the use of advanced analytics in the healthcare field. As usage grows so will the basis of data available to analyze and the ability to create better analysis models.

Another incentive for leveraging big data analytics by physicians is being led by the governments and private insurance. Instead of paying for services as they are performed, alternate payment models are being explored. For instance, in the United States the Centers for Medicaid and Medicare services is creating Alternate Payment Models to stimulate high-quality, cost-efficient care[40]. Physicians are able to earn more income and profit by achieving better outcomes. They will be willing to invest in computer analysis which will help them to diagnose and treat patients better. The financial incentive will drive change in providers' habits which will benefit the healthcare big data analytics and patients.

14 DRAWBACKS

Leveraging big data innovations does not come without hurdles. One of the first is that people are generally slow or not open to change. The more personal the need for change, the less open they are. Organizations (hospitals, physician groups) are no different. Part of being an individual is making choices based on what information you can gather and leveraging your ability to make a determination. This is part of what makes each person unique. It is also how we learn. The more we become dependent on machines, the less we store in our own brains and we stop "building the networks in our brains to solve a whole host of problems.[51]" As those in the healthcare field rely more on technology to diagnosis and treat patients, the less human innovation may be leveraged which can have a detrimental effect over time.

A major complication in big data analytics in any setting is the quality of data. The term emphasis garbage in, garbage out has probably been applied to computer systems since the beginning.

There are techniques used to combat this, but when it comes to people's health it is a bit more important. A portion of the healthcare data used as a base for analysis comes from existing diagnosis and treatment performed by humans. In looking at second opinions for patients, it was estimated that "10% to 62% of second opinions yield a major change in the diagnosis, treatment, or prognosis"[49]. Extrapolating this number to the base of information in big data for analysis means that a significant portion of the data would be different if a patient simply went to a different doctor.

Aside from the data itself, there is the potential for the algorithms behind big data analysis to be biased or having discrimination built into them. There has been a lot of talk about a lack of diversity in the technology world, especially with companies in Silicon Valley. This lack of diversity could become manifested into the analytics behind healthcare analytics. Different cultures and different races have some unique healthcare challenges. With a lack of diversification in key jobs the developers of healthcare systems could under-serve large portions of the world's population due to a lack of understanding of how certain diseases affect their everyday lives. The United States Federal Trade Commission has asked companies in general to look at how representative their big data is and whether their models have built in biases[48]. The fact that healthcare around the world varies based economic factors makes it easy to understand how the data itself can be discriminatory. More wealthy people will be proportionally more represented than the poor thus skewing the data toward conditions afflicting the wealthy.

While big data will help to diagnose patients and create treatment plans, it does not come without its drawbacks. One of the biggest may be innovation. Part of being human is the ability to think of what has not already been done before. As algorithms and data analysis based on the historical variables begin to become more commonplace, there will be a reduction in the human factor of the medical profession. When faced with what can seem like a dire situation, the human mind can think of new options not previously discovered. Trying something which may not seem to have an impact on the surface, but something completely unrelated to any prior decision made can lead to new alternatives. What will a computer do with a patient when it does not see any hope? A human physician may opt to take a risk. It is a well-informed risk with the patient knowing that there are no guarantees. It is easy to assume when an automated course of action without a substantial chance of a positive outcome is encountered that a physician would be able to intervene. This is true for a while, but as more and more of medicine is turned over to computer diagnosis and treatments the pool of capable physicians will shrink. With less people involved the less chance there is that the truly gifted individuals who make strides in the field will even decide to enter the field in the first place. In other words, these individuals may decide on a different career path and their discoveries would be left undiscovered.

15 CONCLUSIONS

Big data is an expanding science in many fields. The ability to digitize, collect, store, and analyze data has never been more than it is today. The type of information that can be used in data analysis is expanding every day as well. Images, videos, and sound are all part of the inputs into big data. Computers are now able to leverage

natural language processing to make inputs that much easier to collect. As this field continues to grow, the ability to leverage it in improving healthcare around the world will grow as well.

We are on the edge of a shift in healthcare for the betterment of humankind. Advances will not be limited to one nation or one class of people. While healthcare may not be universal in its application, not every person will be able to access the same level of care, there will be benefits which can eventually help all people. A mobile unit which can be taken to almost any part of the planet will be able to have the knowledge better than most doctors practicing today. Doctors will have access to new drugs, diagnostic information, and treatment plans than they ever had before. They will be able to leverage new advances in medicine without having to read as many publications as they can. They will have a tool that reads and learns for them and provides that insight on case by case basis.

Through the use of data analysis of sources of data which did not exist a decade or so ago, we will be able to identify when a disease is starting to spread and react, thus limiting its impact. Because of technology people will be spared from suffering and they will never even know it. By understanding the human genome people who may be more susceptible certain diseases can be treated before they take hold. Babies will have their genome sequenced while they are still in their mother's womb. This one aspect of the power of big data, the ability to process and understand a human genome, may be the single largest breakthrough in healthcare. It can provide insight into how each person individually reacts to the world around them and what science can do to make that interaction better. What science can do to help each person avoid potential chronic conditions which are not only financially costly, but that severely reduce their quality of life or end their life. Through advances in big data we will not only live longer, but live better.

ACKNOWLEDGMENTS

The author would like to thank Dr. Gregor von Laszewski for his support throughout this process. By offering an environment in which students were able to explore areas in big data which interested them, we were all able to further our knowledge individually and collectively. This project is similar to big data itself. It brought together various thoughts which could be considered data points into the collection of the class. With access open to all, and potentially future classes, the collection of projects becomes a big data collection unto itself.

REFERENCES

- [1] [n. d.]. Smartphones industry: Statistics and Facts. Online. ([n. d.]). <https://www.statista.com/topics/840/smartphones/>
- [2] [n. d.]. Statistics and Facts on Wearable Technology. Online. ([n. d.]). <https://www.statista.com/topics/1556/wearable-technology/>
- [3] Defense Information Systems Agency. [n. d.]. Security Technical Implementation Guides (STIGs). Online. ([n. d.]). <https://iase.disa.mil/stigs/Pages/index.aspx>
- [4] Rick Altinger. 2017. Five Big Data Solutions to Manage Chronic Diseases. Online. (08 2017). <https://medcitynews.com/2017/08/five-big-data-solutions-manage-chronic-diseases/?rf=1>
- [5] Geoff Appelboom, Elvis Camacho, et al. 2014. Smart Wearable Body Sensors for Patient Self-Assessment and Monitoring. Online. (2014). <https://archpublichealth.biomedcentral.com/track/pdf/10.1186/2049-3258-72-28?site=http://archpublichealth.biomedcentral.com>
- [6] Jessica Barnett and Edward Berchick. 2017. Health Insurance Coverage in the United States: 2016. Online. (09 2017). <https://www.census.gov/content/dam/Census/library/publications/2017/demo/p60-260.pdf>

- [7] Meredith Barrett, Olivier Humblet, et al. 2013. Big Data and Disease Prevention: From Quantified Self to Quantified Communities. *Big Data* 1, 3 (09 2013), 168–175. <https://doi.org/10.1089/big.2013.0027>
- [8] David Bates, Suchi Saria, et al. 2014. Big Data in Health Care: Using Analytics to Identify and Manage High-Risk and High-Cost Patients. *Health Affairs* 33, 7 (2014), 1123–1131. <https://doi.org/10.1377/hlthaff.2014.0041>
- [9] Jennifer Bresnick. 2015. How Healthcare Big Data Analytics Is Tackling Chronic Disease. Online. (06 2015). <https://healthitanalytics.com/news/how-healthcare-big-data-analytics-is-tackling-chronic-disease>
- [10] Jennifer Bresnick. 2016. How Big Data, EHRs, IoT Combine for Chronic Disease Management. Online. (02 2016). <https://healthitanalytics.com/news/how-big-data-ehrs-iot-combine-for-chronic-disease-management>
- [11] Jennifer Bresnick. 2017. Top 10 Challenges of Big Data Analytics in Healthcare. Online. (06 2017). <https://healthitanalytics.com/news/top-10-challenges-of-big-data-analytics-in-healthcare>
- [12] Jennifer Bresnick. 2017. Which Healthcare Data is Important for Population Health Management? Online. (06 2017). <https://healthitanalytics.com/news/which-healthcare-data-is-important-for-population-health-management>
- [13] Elizabeth Chapman, Anna Kaatz, and Molly Carnes. 2013. Physicians and Implicit Bias: How Doctors May Unwittingly Perpetuate Health Care Disparities. *Journal of General Internal Medicine* 28, 11 (11 2013), 1504–1510. <https://doi.org/10.1007/s11606-013-2441-1>
- [14] D'Vera Cohn. 2007. The Growing Global Chronic Disease Epidemic. Online. (05 2007). <http://www.prb.org/Publications/Articles/2007/GrowingGlobalChronicDiseaseEpidemic.aspx>
- [15] Ben Dickson. 2017. How Artificial Intelligence is Revolutionizing Healthcare. Online. (2017). <https://thenextweb.com/artificial-intelligence/2017/04/13/artificial-intelligence-revolutionizing-healthcare/>
- [16] Brian Eastwood. 2016. Bringing Big Data to Drug Discovery. Online. (09 2016). <http://mitsloan.mit.edu/newsroom/articles/bringing-big-data-to-drug-discovery/>
- [17] Suzanne Elvidge. [n. d.]. Digging for Big Data Gold: Data Mining as a Route to Drug Development Success. Online. ([n. d.]). <https://www.clinicalleader.com/doc/digging-for-big-data-gold-data-mining-as-a-route-to-drug-development-success-0001>
- [18] Bonnie Feldman. 2013. Genomics and the Role of Big Data in Personalizing the Healthcare Experience. Online. (08 2013). <https://www.oreilly.com/ideas/genomics-and-the-role-of-big-data-in-personalizing-the-healthcare-experience>
- [19] Centers for Medicare and Medicaid Services. [n. d.]. CMS Information Security and Privacy Overview. Online. ([n. d.]). <https://www.cms.gov/Research-Statistics-Data-and-Systems/CMS-Information-Technology/InformationSecurity/index.html?redirect=/InformationSecurity/>
- [20] Lauren Friedman. 2014. IBM's Watson Supercomputer May Soon be the Best Doctor in the World. Online. (04 2014). <http://www.businessinsider.com/ibms-watson-may-soon-be-the-best-doctor-in-the-world-2014-4>
- [21] Malaria Gokey. 2016. Why smart clothes, not watches, are the future of wearables. Online. (01 2016). <https://www.digitaltrends.com/wearables/smart-clothing-is-the-future-of-wearables/>
- [22] World Bank Group. [n. d.]. Health Expenditure per Capita (current US\$). Online. ([n. d.]). https://data.worldbank.org/indicator/SH.XPD.PCAP?end=2014&name_desc=true&start=2004&view=chart
- [23] Karen He, Dongliang Ge, and Max He. 2017. Big Data Analytics for Genomic Medicine. *International Journal of Molecular Sciences* 18, 2 (02 2017). <https://doi.org/10.3390/ijms18020412>
- [24] Scalable Health. 2017. Managing Chronic Conditions using Big Data. Online. (03 2017). https://www.scalablehealth.com/Resources/WP/SS_Chronic_Illness_ThoughtPaper.pdf
- [25] Partners Healthcare. 2014. 2014 i2b2/UTHealth Shared-Tasks and Workshop on Challenges in Natural Language Processing for Clinical Data. Online. (2014). <https://www.i2b2.org/NLP/HeartDisease/>
- [26] CHEO Research Institute. [n. d.]. What De-Identification Software Tools are There? Online. ([n. d.]). <http://www.ehealthinformation.ca/faq/de-identification-software-tools/>
- [27] Frank Jacobs. [n. d.]. The Patients Per Doctor Map of the World. Online. ([n. d.]). <http://bigthink.com/strange-maps/185-the-patients-per-doctor-map-of-the-world>
- [28] Kate Kelland. 2011. Chronic Disease to Cost \$47 Trillion by 2030: WEF. Online. (09 2011). <https://www.reuters.com/article/us-disease-chronic-costs-chronic-disease-to-cost-47-trillion-by-2030-wef-idUSTRE78H2IY20110918>
- [29] Bijan Khosravi. 2016. Will You Trust AI To Be Your New Doctor? Online. (03 2016). <https://www.forbes.com/sites/bijankhosravi/2016/03/24/will-you-trust-ai-to-be-your-new-doctor-a-five-year-outcome/#3629545b3724>
- [30] MS Kohn, J Sun, et al. 2014. IBM's Health Analytics and Clinical Decision Support. *Yearbook of Medical Informatics* 9, 1 (2014), 154–162. <https://doi.org/10.15265/IY-2014-0002>
- [31] Matthew Komorowski. 2014. A History of Storage Cost. Online. (03 2014). <http://www.mkomo.com/cost-per-gigabyte-update>
- [32] Max Langridge and Luke Edwards. 2017. Best Smart Clothes: Wearables to Improve Your Life. Online. (10 2017). <http://www.pocket-lint.com/news/131980-best-smart-clothes-wearables-to-improve-your-life>
- [33] Mona Lebied. 2017. 9 Examples of Big Data Analytics in Healthcare That Can Save People. Online. (05 2017). <https://www.datapine.com/blog/big-data-examples-in-healthcare/>
- [34] KJ Lee. 2017. Here's How to Reduce Healthcare Costs. Online. (05 2017). <http://medicaleconomics.modernmedicine.com/medical-economics/news/heres-how-reduce-healthcare-costs?page=0,1>
- [35] Lada Leyens, Matthias Reumann, et al. 2017. Use of Big Data for Drug Development and for Public and Personal Health and Care. *Genetic Epidemiology* 41, 1 (01 2017), 51–60. <https://doi.org/10.1002/gepi.22012>
- [36] Zengjian Liu, Yangxin Chen, et al. 2015. Automatic De-Identification of Electronic Medical Records using Token-Level and Character-Level Conditional Random Fields. *Journal of Biomedical Informatics* 58 (12 2015), S47–S52. <https://doi.org/10.1016/j.jbi.2015.06.009>
- [37] Laura Lorenzetti. 2016. Here's How IBM Watson Health is Transforming the Health Care Industry. Online. (04 2016). <http://fortune.com/ibm-watson-health-business-strategy/>
- [38] Sid Nair. 2015. How Advanced Genomics, Big Data will Enable Precision Medicine. Online. (09 2015). <https://healthitanalytics.com/news/how-advanced-genomics-big-data-will-enable-precision-medicine>
- [39] American Academy of Actuaries. 2016. Drivers of 2017 Health Insurance Premium Changes. Online. (05 2016). <https://www.actuary.org/content/drivers-2017-health-insurance-premium-changes-0>
- [40] Department of Health and Human Services. [n. d.]. APMs Overview. Online. ([n. d.]). <https://ppq.cms.gov/apms/overview>
- [41] United States Department of Labor. [n. d.]. Long-Term Price Trends for Computers, TVs, and Related Items. Online. ([n. d.]). <https://www.bls.gov/opub/ted/2015/long-term-price-trends-for-computers-tvs-and-related-items.htm>
- [42] Optum. [n. d.]. Data Rich, Insight Poor. Online. ([n. d.]). https://cdn-aem.optum.com/content/dam/optum3/optum/en/images/infographics/Game_changer.Track.Two.04_Data_Rich_Insight_Poor_Infog/Images_2016.pdf
- [43] World Health Organization. [n. d.]. Density of Physicians (Total Number per 1000 Population): Latest Available Year. Online. ([n. d.]). http://www.who.int/gho/health_workforce/physicians_density/en/
- [44] World Health Organization. [n. d.]. Chronic Diseases and Health Promotion. Online. ([n. d.]). <http://www.who.int/chp/en/>
- [45] World Health Organization. [n. d.]. Global Health Observatory (GHO) data. Online. ([n. d.]). http://www.who.int/gho/health_workforce/en/
- [46] World Health Organization. 2005. Evaluating the Costs and Benefits of National Surveillance and Response Systems. Online. (2005). http://www.who.int/csr/resources/publications/surveillance/WHO_CDS_EPR_LYO_2005_25.pdf
- [47] World Health Organization and World Bank. 2015. New Report Shows that 400 Million do not have Access to Essential Health Services. Online. (06 2015). <http://www.who.int/mediacentre/news/releases/2015/uhc-report/en/>
- [48] Out-Law.com. 2016. Use of Big Data Can Lead to 'harmful exclusion, discrimination' fit? FTC. Online. (01 2016). https://www.theregister.co.uk/2016/01/08/use_of_big_data_can_lead_to_harmful_exclusion_or_discrimination_us_regulator/
- [49] Velma Payne, Hardeep Singh, et al. 2014. Patient-Initiated Second Opinions: Systematic Review of Characteristics and Impact on Diagnosis, Treatment, and Satisfaction. *Mayo Clinic Proceedings* 89, 5 (05 2014), 687–696. <https://doi.org/10.1016/j.mayocp.2014.02.015>
- [50] Inc. Privacy Analytics. [n. d.]. Privacy Analytics Eclipse. Online. ([n. d.]). <https://privacy-analytics.com/software/privacy-analytics-eclipse/>
- [51] John Robison. 2009. Is Technology Making us Dumber? Online. (11 2009). <https://www.psychologytoday.com/blog/my-life-as-aspergers/200911/is-technology-making-us-dumber>
- [52] Sameer Sawarkar. 2013. Remote Healthcare Solution. Online. (2013). http://www.who.int/health/resources/compendium_ehealth2013.7.pdf
- [53] Abhinav Shashank. 2016. Chronic Care Management Marries Big Data. Online. (12 2016). <http://blog.innovaccer.com/chronic-care-management-marries-big-data/>
- [54] Alyssa Sparacino. 2013. 11 Surprising Health Benefits of Sleep. Online. (07 2013). <http://www.health.com/health/gallery/0,,20459221,00.html#go-ahead-snooze--1>
- [55] Caitlin Stackpool, John Porcari, et al. 2015. ACE-sponsored Research: Are Activity Trackers Accurate? Online. (01 2015). <https://www.acefitness.org/education-and-resources/professional/prosource/january-2015/5216/ace-sponsored-research-are-activity-trackers-accurate>
- [56] Eliza Strickland. 2017. IBM Watson Makes a Treatment Plan for Brain-Cancer Patient in 10 Minutes; Doctors Take 160 Hours. Online. (08 2017). <https://spectrum.ieee.org/the-human-os/biomedical/diagnostics/ibm-watson-makes-treatment-plan-for-brain-cancer-patient-in-10-minutes-doctors-take-160-hours>
- [57] Tom Sullivan. 2017. Cognitive Computing will Democratize Medicine, IBM Watson Officials Say. Online. (04 2017). <http://www.healthcareitnews.com/news/cognitive-computing-will-democratize-medicine-ibm-watson-officials-say>

- [58] Ann Tinker. 2017. How to Improve Patient Outcomes for Chronic Diseases and Comorbidities. Online. (2017). <https://www.healthcatalyst.com/how-to-improve-chronic-diseases-comorbidities>
- [59] Partnership to Fight Chronic Disease. [n. d.]. The Growing Crisis of Chronic Disease in the United States. Online. ([n. d.]). <https://www.fightchronicdisease.org/sites/default/files/docs/GrowingCrisisofChronicDiseaseintheUSfactsheet.81009.pdf>
- [60] Jonathan Vanian. 2015. IBM's Watson Supercomputer is Learning Arabic in Move to Middle East. Online. (07 2015). <http://fortune.com/2015/07/14/ibm-watson-home-middle-east/>

```
bibtext report
```

```
This is BibTeX, Version 0.99d (TeX Live 2016)
The top-level auxiliary file: report.aux
The style file: ACM-Reference-Format.bst
Database file #1: report.bib
Warning--no key, author in StatistaWearable
Warning--no author, editor, organization, or key in StatistaWearable
Warning--to sort, need author or key in StatistaWearable
Warning--no key, author in StatistaPhones
Warning--no author, editor, organization, or key in StatistaPhones
Warning--to sort, need author or key in StatistaPhones
Warning--no key, author in StatistaPhones
Warning--no key, author in StatistaPhones
Warning--no key, author in StatistaWearable
Warning--no key, author in StatistaPhones
Warning--empty author in StatistaPhones
Warning--empty year in StatistaPhones
Warning--no key, author in StatistaWearable
Warning--no author, editor, organization, or key in StatistaWearable
Warning--empty author in StatistaWearable
Warning--empty year in StatistaWearable
Warning--empty year in DISA
Warning--empty year in ClinicalLeader
Warning--empty year in CMS
Warning--empty year in WoldBankPerCapita
Warning--page numbers missing in both pages and numpages fields in PMC5343946
Warning--empty year in eHealthInfo
Warning--empty year in BigThink
Warning--empty year in CMSAPM
Warning--empty year in CompPrices
Warning--empty year in Optum
Warning--empty year in WHODensity
Warning--empty year in WHOChronicDisease
Warning--empty year in WHOGHO
Warning--empty year in PrivacyAnalytics
Warning--empty year in FightChronicDisease
(There were 32 warnings)
```

```
bibtext _ label error
```

```
bibtext space label error
```

```
=====
```

```
bibtext comma label error
```

```
=====
```

```
latex report
```

```
=====
```

```
[2017-12-04 12.23.45] pdflatex report.tex
```

```
This is pdfTeX, Version 3.14159265-2.6-1.40.17 (TeX Live 2016) (preloaded format=pdflatex)
```

```
Missing character: ""
```

```
Missing character: ""
```

```
Missing character: ""
```

```
Typesetting of "report.tex" completed in 1.0s.
```

```
=====
```

```
Compliance Report
```

```
=====
```

```
name: Marks, Paul
```

```
hid: 327
```

```
paper1: 100% 10/25/2017
```

```
paper2: 100% 11/06/17
```

```
project: 99%
```

```
yamlcheck
```

```
-----
```

```
wordcount
```

```
-----
```

```
11
```

```
wc 327 project 11 9864 report.tex
```

```
wc 327 project 11 10594 report.pdf
```

```
wc 327 project 11 2551 report.bib
```

```
find "
```

```
-----
```

```
passed: True
```

```
find footnote
```

```
passed: True
```

```
find input{format/i523}
```

```
4: \input{format/i523}
```

```
passed: True
```

```
find input{format/final}
```

```
passed: False
```

```
floats
```

```
figures 0
```

```
tables 0
```

```
includegraphics 0
```

```
labels 0
```

```
refs 0
```

```
floats 0
```

```
True : ref check passed: (refs >= figures + tables)
```

```
True : label check passed: (refs >= figures + tables)
```

```
True : include graphics passed: (figures >= includegraphics)
```

```
True : check if all figures are refered to: (refs >= labels)
```

```
Label/ref check
```

```
passed: True
```

```
When using figures use columnwidth
```

```
[width=1.0\columnwidth]
```

```
do not change the number to a smaller fraction
```

```
find textwidth
```

```
passed: True
```

below_check

bibtex

label errors

bibtex errors

```
This is BibTeX, Version 0.99d (TeX Live 2016)
The top-level auxiliary file: report.aux
The style file: ACM-Reference-Format.bst
Database file #1: report.bib
Warning--no key, author in StatistaWearable
Warning--no author, editor, organization, or key in StatistaWearable
Warning--to sort, need author or key in StatistaWearable
Warning--no key, author in StatistaPhones
Warning--no author, editor, organization, or key in StatistaPhones
Warning--to sort, need author or key in StatistaPhones
Warning--no key, author in StatistaPhones
Warning--no key, author in StatistaPhones
Warning--no key, author in StatistaWearable
Warning--no key, author in StatistaPhones
Warning--no author, editor, organization, or key in StatistaPhones
Warning--empty author in StatistaPhones
Warning--empty year in StatistaPhones
Warning--no key, author in StatistaWearable
Warning--no author, editor, organization, or key in StatistaWearable
Warning--empty author in StatistaWearable
Warning--empty year in StatistaWearable
Warning--empty year in DISA
Warning--empty year in ClinicalLeader
Warning--empty year in CMS
Warning--empty year in WoldBankPerCapita
Warning--page numbers missing in both pages and numpages fields in PMC5343946
Warning--empty year in eHealthInfo
Warning--empty year in BigThink
Warning--empty year in CMSAPM
Warning--empty year in CompPrices
Warning--empty year in Optum
Warning--empty year in WHODensity
Warning--empty year in WHOChronicDisease
```

```
Warning--empty year in WHOGHO
Warning--empty year in PrivacyAnalytics
Warning--empty year in FightChronicDisease
(There were 32 warnings)
```

```
bibtex_empty_fields
```

```
entries in general should not be empty in bibtex
```

```
find ""
```

```
passed: True
```

```
ascii
```

```
non ascii found 8217
```

```
=====
The following tests are optional
=====
```

```
Tip: newlines can often be replaced just by an empty line
```

```
find newline
```

```
passed: True
```

```
cites should have a space before \cite{} but not before the {
```

```
find cite {
```

```
passed: True
```

Using Machine Learning Classification of Opioid Addiction for Big Data Health Analytics

Sean M. Shiverick

Indiana University Bloomington

smshiver@indiana.edu

ABSTRACT

Classification of opioid addiction can identify important features relevant for predicting drug abuse and overdose death. Machine learning procedures were applied to data from a large National Survey of Drug Use and Health (NSDUH-2015) to classify individuals for illicit opioid use according to demographic characteristics and mental health attributes (e.g., depression). Classification models of opioid addiction can be extended for big data health analytics to include high-dimensional datasets, data collected over previous years, or expanded to the larger population of patients taking prescription opioid medication. The results seek to raise awareness of risk factors related to opioid addiction among patients and medication prescribers, and help decrease the risk of opioid overdose death.

KEYWORDS

Health Analytics, Machine Learning Classifiers, Opioid Addiction, Big Data, i523, hid335

1 INTRODUCTION

Big Data offers tremendous potential to fuel innovation and transform society. Can this momentum be harnessed to address a serious health crisis such as the opioid overdose epidemic? [7] Health informatics is generating huge amounts of data at a rapid pace, from electronic medical records (EMRs), clinical research data, to population-level public health data [5]. This project considers health analytics from two levels, the research questions being addressed and the data used to answer them. The question of interest in this project is whether opioid dependency and addiction can be predicted from demographic attributes and psychological characteristics. Survey research provides data on a wide range of issues that people may be reluctant to disclose, including mental health disorders, personal medical health concerns, prescription medications, and illicit drug use. Responses to surveys may be biased to some degree, but measures of confidentiality and anonymity help to assure more accurate disclosures. The goal of this project is to use machine learning procedures to classify individuals susceptible to opioid abuse and dependence. Understanding the features that contribute to opioid addiction can identify underlying risk factors and increase awareness of potential opioid abuse for patients and health care providers. The results could be extended to big data from previous years of the opioid crisis and to the larger population of patients taking prescription opioid mediation. Different machine learning classification methods are discussed.

1.1 Opioid Overdose Epidemic

The abuse of prescription opioid medication in the U.S. has become a major health crisis of epidemic proportions [23]. Over 2 million Americans were dependent or abused prescription opioids such as oxycodone or hydrocodone in 2014[3]. Overdose deaths from prescription opioids have quadrupled since 1999, resulting in more than 180,000 deaths between 1999 to 2015 [11]. Drug overdose deaths increased significantly for males and females, between 25-44 years, ages 55 and older, for Non-Hispanic Whites and Blacks, in the Northeast, Midwest, and Southern regions of the U.S. [7]. Mobile health applications can monitor patient medication consumption and provide an early warning system for potential abuse, detecting sudden changes in medications, higher dosages, or rapid escalation of a prescribed dosage [22]. Reliable information about medication dosages can be difficult to obtain based on self-reports. Individuals dependent or addicted to prescription opioids may obtain synthetic opioids such as fentanyl or illicit drugs such as heroin. Because the dosage levels and potency of illicit opioids are largely unknown, there is greater risk of drug overdose death. The sharp increase in overdose deaths due to synthetic opioids (other than methadone) has coincided with the increased availability of illicitly manufactured fentanyl, which is indistinguishable from prescription fentanyl. The findings indicate the opioid overdose epidemic is getting worse, and requires urgent action to prevent opioid dependence, abuse and overdose death. The target group for this project is individuals who reported misusing or abusing prescribed opioid medication who also used heroin, shown in Figure 1.

1.2 Machine Learning Approaches

Machine learning is a set of procedures and automated processes for extracting knowledge from data. The two main branches of machine learning are supervised learning and unsupervised learning. Supervised learning problems involve prediction about a specific target variable or outcome of interest. If a given dataset has no target outcome, unsupervised learning methods can be used to discover underlying structure in unlabeled data. The goal of this project is to classify opioid addiction and focuses on supervised learning. Supervised learning is used to predict a certain outcome from a given input, when examples of input/output pairs are available [10]. A machine learning model is constructed from the training set of input-output pairs, to predict new test data not previously seen by the model. The two major approaches to supervised learning problems are regression and classification. When the target variable to be predicted is continuous, or there is continuity between the outcome (e.g., home values, or income), a regression model is used to test the set of features that predict the target variable. If the target is a class label, set of categorical or binary outcomes (e.g., 'spam' or 'ham', 'benign' or 'malignant'), then classification is used

to predict which class or category label that new instances will be assigned to.

1.3 Classification Algorithms

Comparing the performance of different learning algorithms can be helpful for selecting the best model for a given problem [14]. One of the simplest classification algorithms is K-Nearest-Neighbors (KNN) which takes a set of data points and classifies a new data point based on the distance (e.g., Euclidean, by default) to its nearest neighbors. The main parameter for KNN is the number of neighbors, and k of 3 or 5 neighbors works well. The advantage of the KNN classifier is that it provides a solution that is easy to understand. A limitation of KNN is that it does not perform well with a large number of features (100 or more) or sparse datasets. Several different classification algorithms are considered below.

1.3.1 Logistic Regression Classifier. Logistic regression is a commonly used linear model for classification problems. The decision boundary for the logistic regression classifier is a linear function of the input; a binary classifier separates two classes using along a line, plane, or hyperplane. Linear classification models differ in terms of (1) how they measure how well a particular combination of coefficients and intercept fit the training data, and (2) the type of regularization used [10]. The main parameter for linear classification models is the regularization parameter ‘C’. High values of C correspond to less regularization and the model will fit the training set as best as possible, stressing the importance of each individual data point to be classified correctly. By contrast, with low values of C, the model puts more emphasis on finding coefficient vectors (w) that are close to zero, trying to adjust to the ‘majority’ of data points [10]. In addition, the penalty parameter influences the coefficient values of the linear model. The L2 penalty (Ridge) uses all available features, but pushes the coefficient values toward zero. The L1 penalty (Lasso) sets the coefficient values for most features to zero, and uses only a subset for improved interpretability. This paper uses a logistic regression classifier to predict Heroin use from demographic attributes, mental health, prescription opioids, medication use, misuse, and illicit drug use.

1.3.2 Tree Based Models. Decision tree models are widely used for classification and regression. Tree models “learn” a hierarchy of if-else questions that are represented in the form of a decision tree. Building decision trees proceeds from a root node as the starting point and continues through a series of decisions or choices. Each node in the tree either represents either a question or a terminal node (i.e., leaf) that contains the outcome. Applied to a binary classification task, the decision tree algorithm “learns” the sequence of if-else questions that arrives at the outcome most quickly. For data with continuous features, the decisions are expressed in the form of, “Is feature x larger than value y?” [10] In constructing the tree, the algorithm searches through all possible decisions or tests, and find a solution that is most informative about the target outcome. A decision tree classifier is used for binary or categorical targets, and decision tree regression is used for continuous target outcomes. The recursive branching process of tree based models yields a binary tree of decisions, with each node representing a test that considers a single feature. This process of recursive partitioning

is repeated until each leaf in the decision tree contains only a single target. Prediction for a new data point proceeds by checking which region of the partition the point falls in, and predicting the majority in that feature space. The main advantage of tree based models is that they require little adjustment and are easy to interpret. A drawback is that they can lead to very complex models that are highly overfit to the training data. A common strategy to prevent overfitting is *pre-pruning*, which stops tree construction early by limiting the maximum depth of the tree, or the maximum number of leaves. One can also set the minimum number of points in a node required for splitting. Another approach is to build the tree and then remove or collapse nodes with little information, which is called *post-pruning*. Decision trees work well with features measured on very different scales, or with data that has a mix of binary and continuous features.

1.3.3 Random Forests Classifier. A random forest is a collection of decision trees that are slightly different from the others, which each overfits the data in different ways. The idea behind random forests is that overfitting can be reduced by building many trees and averaging their results. This approach retains the predictive power of trees while reducing overfitting. Randomness is introduced into the tree building process in two ways: (a) selecting a bootstrap sample of the data, and (b) selecting features in each node branch [10, 14]. In building the random forest, we first decide how many trees to build (e.g., 10 or 100), and the algorithm makes different random choices so that each tree is distinct. The bootstrapping method repeatedly draws random samples of size n from the dataset (with replacement). The decision trees are built on these random samples that are the same size as the original data, with some points missing and some data points repeated. The algorithm also selects a random subset of p features, repeated separately each node in the tree, so that each decision at the node branch is made using a different subset of features. These two processes help ensure that all of the decision trees in the random forest are different. The important parameters for the random forests algorithm are the number of sampled data points and the maximum number of features; the algorithm could look at all of the features in the dataset or a limited number. A high value for “maximum features” will produce trees in the random forest that are very similar and will fit the data easily based on the most distinctive features, whereas a low value will produce trees that are very different from each other, and reduces overfitting. Random forests is of the most widely used ML algorithms that works well without very much parameter tuning or scaling of data. A limitation of this approach is that Random forests do not perform well with very high-dimensional, data that is sparse data, such as text data.

1.4 Project Goals

The general idea of the project is that prescription opioid dependency and addiction will in many cases lead to the use of illicit opioids such as heroin or fentanyl. According to this reasoning, it was hypothesized that individuals who report using heroin may also be susceptible to misusing or abusing prescription opioid medications. The goal of the study was to identify the set of features important for predicting opioid addiction. The data used in the project is from the National Survey on Drug Use and Health from 2015 (NSUH-

2015) [1], which is the most recent year available. The NSDUH-2015 is a comprehensive survey that covers all aspects of substance use, misuse, dependency, and abuse, including questions related to both prescription medications (opioids, tranquilizers, sedatives) and illicit drugs (e.g., heroin, cocaine, methamphetamine), drug dependency, addiction, and treatment, demographic measures of education and employment, physical health, depression, and mental health treatment. Several classification models were constructed to classify heroin use in the sample by demographics attributes and mental health characteristics (e.g., adult depression). This method addresses the following issues related to opioid dependency and addiction: (i) Identify factors related to illicit opioid use, (ii) Identify factors related to prescription opioid misuse and abuse, and (iii) Examine the relationship between prescription opioid misuse, abuse and heroin use.

2 METHOD

2.1 Data

Data from the 2015 NSUH was downloaded from the Substance Abuse and Mental Health Data Archive (SAMHDA) [1] URL using the ‘get-data.py’ function written to unzip the data files, extract the data as a Pandas data frame, and write the file to CSV file [4]. The dataset consists of 57,146 observations with 2,666 features representing individual-level responses from a survey of the U.S. population. According to the NSDUH codebook, sampling was weighted across states by population size for a representative distribution selected from 6,000 area segments. The sample design used five state sample size groups drawing more heavily from the eight states with the largest population (e.g., CA, FL, IL, MI, NY, OH, PA, TX) which together account for 48 percent of total U.S. population aged 12 or older. All identifying information was collapsed (e.g., age categories) and state identifiers were removed from the public use file to ensure confidentiality. The NSDUH public-use files do not include geographic location, or demographic variables related to ethnicity or immigration status. The weighted survey screening response rate was 81.94 percent and the weighted interview response rate was 71.2 percent.

2.2 Data Cleaning and Preparation

2.2.1 Data Cleaning. All steps of this analysis was completed in a python interactive notebook [16] based following examples from *Python for Data Analysis* [9]. After saving the NSDUH-2015 as a data frame object, the dataset was subset by columns to include demographic characteristics (e.g., age category, sex, marital status, education, employment status, and category of metropolitan area), measures of physical health (e.g., overall health, STDs, Hepatitis, HIV, Cancer, hospitalization), mental health (e.g., Adult Depression, Emotional Distress, Suicidal Thoughts, Plans), Suicide Attempts, Pain Reliever Medication Use, Misuse, and Abuse (over past year, past month), Prescription Opioid Medications Taken in Past year (e.g., Hydrocodone, Oxycodone, Tramadol, Morphine, Fentanyl, Oxymorphone, Demerol, Hydromorphone), Heroin Use, Abuse (over past year, past month), Tranquilizer Use, Sedative Use, Cocaine Use, Amphetamine and Methamphetamine Use, Hallucinogen Use, Drug Treatment (e.g., Inpatient, Outpatient, Hospital, Mental

Health Clinic, ER, Drug Treatment Status), and Mental Health Treatment History. A codebook was created to provide a complete list of variables included with summaries of response categories [19]. The following steps were taken to detect and remove inconsistencies in the data [13]:

- (1) Remove missing values (i.e., ‘NaN’)
- (2) Recode blanks, non-responses, or legitimate skips (e.g., ‘99’, ‘991’, ‘993’) to zero
- (3) Recode dichotomous responses (e.g., “Yes=1”/“No=2”) so that “No=0”
- (4) Recode categorical variables to be consistent with amount or degree (e.g., “1=low”, “2=med”, “3=high”)
- (5) Rename selected variables for better description (e.g., Adult Major Depressive Episode Lifetime changed from ‘AMDELT’ to ‘DEPMELT’)

2.2.2 Aggregate Variables. Because the majority of features were represented as dichotomous “Yes/No” variables, related features were summed to create aggregated variables. For example, overall health, STD, Hepatitis, HIV, Cancer, and hospitalization were aggregated to create a single health measure. The health measure was recoded so that higher scores indicated better health. Questions related to depression, emotional distress, and suicidal thoughts were summed to create a single variable for mental health (‘MENTHLTH’) with scores ranging from 0 to 9. Responses to pain reliever medication use, misuse, abuse, or dependency, were aggregated to create a single variable of pain reliever misuse or abuse (‘PRLMISAB’). All prescription painkiller medications used in the past year were summed. Similarly, all related responses were summed to create single variables for Tranquilizers, Sedatives, Cocaine, Amphetamines, Hallucinogens, Drug Treatment, and Mental Health Treatment. The target outcome of interest for classification, lifetime heroin use (i.e., “Have you ever used heroin before, at any time?”) is a dichotomous variables. The demographic characteristics and aggregated variables were subset and saved to a new data frame consisting of 2 features and 57,146 observations, which was exported to CSV file.

3 RESULTS

3.1 Exploratory Data Analysis

Of the total sample of N=57,146 respondents, 26,736 were male and 30,410 female; 6,343 individuals reported misusing pain medication at some point; however, only 956 respondents had used heroin (570 males, 386 females). Table 1 shows the raw counts of individual substance use by age group (with the sample size for each age group), listing the ten most commonly used opioid pain medications, self-reported misuse of prescription opioid pain relievers (i.e., PRL Misuse Ever), use of prescription Tranquillizers, Sedatives, and Methadone. In addition, self-reported use of illicit drugs such as heroin, cocaine, amphetamines, methamphetamine, Hallucinogens, including LSD and Ecstasy (MDMA). This summary table shows that substance use seems to be highest for individuals between the ages of 18 to 25 and from 35 to 49 years. Of the prescription relievers, Hydrocodone use was almost double the rate of Oxycodone use for each age group, and was significantly higher than any other prescription opioid medication. Use of prescription Fentanyl and

Demerol, two powerful opioids, and synthetic morphines such as Oxymorphone and Hydromorphone, was very low. The rate of prescription Tranquilizer use was several orders of magnitude higher than Sedative use or Methadone use. Compared to other illicit drugs such as Cocaine, Amphetamines, Hallucinogens, heroin use was not very common in this sample. The highest rates of heroin use were seen between the ages of 18 to 49, and was lowest for respondents in the youngest age group 12 to 17, and individuals over 50.

[Table 1 about here.]

Table 2 shows the frequency of individuals reporting that they had experienced mental health issues such as depression, suicidal thoughts, whether they had received mental health treatment, received treatment from a private therapist, or believed that they needed drug treatment, but had not sought treatment, across each age category. Frequency of depression was not included for respondents between 12 to 17 years, and the measure was of adult depression.

[Table 2 about here.]

Figure 1 shows the proportion of individuals who reported misusing prescription opioid pain relievers and who reported using heroin. The left column of the Figure 1 shows the majority of respondents (89 percent) stated they had never misused prescription opioid pain medication or used heroin, although 10 percent reported misusing opioid pain medication at some point. The right panel of Figure 1 shows that, of those individuals who reported using heroin, the proportion who also reported misusing opioid pain medication was almost twice as large as the proportion of those who only used heroin. This is consistent with the hypothesis that misuse of prescription opioids is linked with heroin use for some individuals.

[Figure 1 about here.]

Figure 2 shows the aggregated measure of Opioid Pain Reliever misuse and abuse plotted against the aggregated measure of Heroin use (which includes misuse, abuse, lifetime use, past year use, 30 day use), with weighted regression lines grouped by size of City/Metropolitan region (from none to large). The largest proportion of the sample who report prescription opioid misuse, abuse, and heroin use is represented by observations from large metropolitan areas (red circles) with large population size. However, a small number of observations from rural or small metropolitan regions (blue and green circles) showed very high rates of prescription opioid misuse and abuse. Regression lines (i.e., line of best fit) shown are weighted by the City/Metro region attribute, with a steeper slope shown for smaller metropolitan regions than large metropolitan regions. The difference in slope may be due to the influence of the small number of outliers who had high degrees of prescription opioid misuse, and heroin use. The plot also shows a clear divide on the y-axis, which separates the sample according to high and low or no prescription opioid misuse, although the continuum of heroin use from no, low, to high is distributed fairly evenly along the x-axis.

[Figure 2 about here.]

Figure 3 shows the pairplots of demographic features including mental health (higher scores equal to more depression), Prescription Opioid Pain Reliever (PRL) Medication (aggregated), Heroin

Use (aggregated measure), and Size of City/Metropolitan region. The top row shows that the majority of the sample reported no mental health concerns, whereas a small proportion of the sample reported depression, emotional distress, or suicidal thoughts. Only few people self-described as high in depression reported low Prescription Opioid PRL misuse and abuse. The plot also reveals that prescription opioid misuse and heroin use were distributed approximately evenly for individuals reporting either low, moderate, or high levels of depression, which suggests that depression was not a factor in predicting opioid misuse. The second row shows a small number of individuals from rural areas or small cities who reported very high levels of prescription opioid misuse, although the majority of respondents misusing or abusing prescription opioid were from large metropolitan areas. As described above, the majority of respondents (about 90 percent of the sample) reported they had never misused prescription opioids. In the second row and third and fourth columns, a natural break is seen between individuals who reported high levels of prescription opioid misuse and abuse and those who reported very low or no opioid misuse. A very small proportion of the entire sample reported both misusing and abusing prescription opioids and using heroin, but this is a group of interest. The last column of the second row shows the individuals reporting high levels of opioid misuse and abuse were distributed evenly across city/metropolitan areas of different sizes, with only slightly higher numbers for small cities or rural areas. As stated above, only few participants reported using heroin, and of these, the majority were from large metropolitan areas. Finally, the sample seems to have slightly higher proportions from small and large metropolitan areas, which is likely due to weighted sampling, which drew more from heavily populated regions.

[Figure 3 about here.]

3.2 Classifier Models of Heroin Use

This analysis classified individuals according to whether they had ever used heroin (i.e., "Heroin Use Ever"). All classifier models were constructed using SciKit Learn [10] using an interactive python jupyter notebook [17]. The features of interest were demographic characteristics, health, mental health (adultdepression), prescription opioid misuse and abuse ('PRLMISEVR', 'PRLMISAB', 'PRLANY'), prescription tranquilizers and sedatives ('TRQLZRS', 'SEDATVS'), illicit drugs ('COCAINE', 'AMPHETMN'), drug treatment ('TRTMENT'), and mental health treatment ('MHTRTMT'). The target variable was Heroin Use ('HEROINEVR'). Next, the dataset was split into the training set and test sets using the 'train-test -split' function in 'sklearn'. Model accuracy for the training set and test set are reported, with different parameter values, and features importance.

3.2.1 Logistic Regression Classifier. Logistic Regression Classification is based on a linear equation that calculates the relative weight of each feature for a categorical target or binary outcome ("yes/no") [14]. The logistic regression classifier was fit to the training data in Scikit-Learn, and the model was validated on the test data. By default, the model applies L2 penalty (Ridge). The training set accuracy was 0.983 and the test set accuracy was 0.984. The parameter 'C' determines the strength of regularization, with higher values of C providing greater regularization. The L1 penalty

(Lasso) limits the values of most coefficients to zero, creating a more interpretable model that uses only a few features. Figure 4 plots the coefficients of logistic regression classifier for heroin use with the L1 Penalty (Lasso) under different values of parameter C. The default setting, C=1.0, provides good performance for train and test sets, but the model is very likely underfitting the test data. Using a higher value of C fits a more ‘flexible’ model and generally gives improved accuracy for both training and tests sets. Using a value of C=100 yielded training set accuracy of 0.98 and test set accuracy of 0.98. Figure 4 shows that the features coefficient values did not change much according to the values of parameter C, and the accuracy values were approximately the same for all values of C. Examination of the coefficients from the logistic regression classifier revealed the three features which were most closely associated with Heroin use were: Prescription Opioid Pain Reliever (PRL) Misuse ever (as predicted), Cocaine Use, and Amphetamine use, respectively.

[Figure 4 about here.]

3.2.2 Decision Tree Classifier. The following analysis used the *Decision Tree Classifier* package in Scikit-Learn, which only does pre-pruning. First, the decision model was build using the default setting of a fully developed tree until all leaves are pure. The ‘random state’ features is fixed to break ties internally. Accuracy on the training set was 0.99 and test set accuracy was 0.974. Without restricting their depth, decision trees can become complex; unpruned trees are prone to overfitting and do not generalize well to new data. Limiting the depth of tree decreases overfitting, which results in lower training set accuracy, but improved performance on the test set. Next, pre-pruning was applied, with a maximum depth of 4, which means the algorithm split on four consecutive questions. Training set accuracy of the pruned tree was 0.985 and test set accuracy was 0.984. Even with a depth of 4, the tree can become a bit complex. Figure 5 shows a partial view of the decision tree classifier of heroin use (the entire tree was too wide to include as a legible Figure), and the full tree image is available in the notebook ‘BDA-Analytics-Classifier-Heroin.ipynb’ [17]. The decision tree shows the top features that the algorithm split on to classify heroin use. One way to interpret a decision tree it by following the sample numbers represented at the test split for each node. The classifier algorithm selected Cocaine Use (aggregated score) as the root node of the decision tree. The branch to the left side of the tree represents samples with a score equal to or less than 1.5 (n=40956), whereas the branch to the right represents samples with a Cocaine Use score greater than 1.5 (n=1903). The second split on the right occurs for Any Prescription Opioid Pain Reliever Use (‘PRLANY’), with n=1443 having a score less than or equal to 3.5, and n=460 respondents with a PRL score greater than 3.5. In other words, of those respondents who reported relatively high Cocaine use, a small portion also reported relatively high Prescription Opioid PRL use. Instead of looking at the whole tree, features importance is a common summary function that rates how important each feature is for the classification decisions made in the algorithm. Each feature is assigned an importance value between 0 and 1; with a value of 1 indicating the feature perfectly predicts the target and a value of 0 meaning that the feature was not used at all. Feature importance values also always sum to 1. A feature may have a

low feature importance value because another feature encodes the same information. The top two important features for classifying Heroin Use were Cocaine Use and Any Prescription Opioid PRL Use, with smaller importance given to Opioid PRL Misuse Ever and Prescription Opioid PRL Misuse and Abuse.

[Figure 5 about here.]

3.2.3 Random Forests Classifier. Random forests is an ensemble approach that builds many trees and averages their results to reduce overfitting. The model was build using the Random Forest Classifier package in Scikit-Learn. The parameters of interest for building random forests are: (a) the number of trees (‘n-estimators’), (b) the number of data points for bootstrap sampling (‘n-samples’), and (c) the maximum number of features considered at each node (‘max-features’). The max-features parameter determines how random each tree is, with smaller values of max-features resulting in trees in the random forest that are very different from each other. This analysis applied a random forest consisting of 100 trees to classify Heroin Use, and the random state was set to zero. The training set accuracy was 0.999 and the test set accuracy was 0.984. Often the default settings for random forests work well, but we can apply pre-pruning as with a single tree, or adjust the maximum number of features. Feature importance for random forests is computed by aggregating the feature importance over trees in the random forest, and random forests gives non-zero importance to more features than a single tree. Typically random forests provide a more reliable measure of feature importance than the feature importance for a single tree. Figure 6 shows the feature importance of the random forests classifier for heroin use with 100 trees. Similar to the single tree, the random forest selected Cocaine Use as the most informative feature in the model, followed by Any PRL Use, which is an aggregated measure of prescription opioid medication use. Following after that, several features were tied for third place of importance, namely Education Level, Overall Health, Age Category, and Pain Reliever Misuse and Abuse. Random forests provides much of the same benefit as decision trees, while compensating for some of their shortcomings of overfitting. Single trees are still useful for visually representing the decision process.

[Figure 6 about here.]

3.2.4 Gradient Boosting Classifier Tree. Gradient boosting machines is another ensemble method that combines multiple decision trees for regression or classification by building trees in a serial fashion, where each tree tries to correct for mistakes of the previous one [10]. Gradient boosted regression trees use strong pre-pruning, with shallow trees of a depth of one to five. Each tree only provides a good estimate of part of the data, but combining many shallow trees (i.e., “weak learners”), the use many simple models iteratively improves performance. In addition to pre-pruning and the number of trees, an important parameter for gradient boosting is the learning rate, which determines how strongly each tree tries to correct for mistakes of previous trees. A high learning rate produces stronger corrections, allowing for more complex models. Adding more trees to the ensemble also increases model complexity. Gradient boosting and random forests perform well on similar tasks and data; it is common to first try random forests and then include gradient boosting to attain improvements in accuracy of the learning

model. This analysis used the Gradient Boosting Classifier from Scikit-Learn to classify Heroin Use, with the default setting of 100 trees of maximum depth of 3, and a learning rate of 0.1. The model was built on the training set and evaluated on the test set, with both training set and test set accuracy equal to 0.984. To reduce overfitting, pre-pruning could be implemented by reducing the maximum depth, or by reducing the learning rate. Figure 7 shows that the feature importance for the gradient boosting classifier tree looks similar to the feature importance for random forests, but the gradient boosting has decreased the importance of many features to zero. Again Cocaine is selected as the most informative features, followed by Any Opioid PRL Use. In addition to Prescription Opioid PRL Misuse and Abuse, the gradient boosting classifier selected Amphetamine Use as an informative feature of Heroin Use.

[Figure 7 about here.]

3.3 Classifier Models of Prescription Opioid Pain Reliever (PRL) Misuse

This section reports results from the same set of classification analyses described above using Prescription Opioid Pain Reliever Misuse ('PRLMISEVR') the target variable. Attributes related to Heroin Use were now included as features (e.g., 'HEROINEVR', 'HEROINUSE', 'HEROINFOY'). The classifier models were built using SciKit Learn in a python notebook [18]. The dataset was split into the training set and test sets using the 'train-test-split' function in sklearn and the target variables was designated. Model accuracy for the training set and test set are reported, for different parameter values, with feature importance.

3.3.1 Logistic Regression Classifier. The logistic regression classifier was fit to the training data using the L1 penalty (Lasso), using different values of the regularization parameter C, and the model was validated on the test data. Higher value of parameter C typically gives improved accuracy for both training and tests sets; however, in this case, the training set accuracy was 0.901 and test set accuracy was 0.903, and these values were consistent for all values of parameter C. Figure 8 plots the coefficients of logistic regression classifier for Prescription Opioid PRL Misuse under different values of C. As shown in Figure 8, the features with the highest coefficient values were Treatment (for substance use), Heroin Use (as predicted), as well as Cocaine and Amphetamine use. This result indicates that Prescription Opioid Misuse is positively related to Drug Treatment, meaning that respondents who reported higher levels of opioids misuse were also in treatment, but that people who were misusing opioid medications were also more likely to have used illicit drugs such as heroin, cocaine, and amphetamine.

[Figure 8 about here.]

3.3.2 Decision Tree Classifier. The Decision Tree Classifier package in Scikit-Learn was used to build the tree model, pre-pruning was applied with a maximum depth of 4, which means the algorithm split on four consecutive questions. The training set accuracy of the pruned tree was 0.902 and test set accuracy was 0.902. Figure 9 shows a partial view of the decision tree classifier of prescription opioid misuse (the full tree is included in the 'BDA-Analytics-Classifier-PRL.ipynb' notebook) [18]. As Figure 9 shows, the decision tree classifier selected Cocaine Use as the root note, that

branched by the test score equal to or less than 0.5 (any Cocaine Use). At the second node, on the branch to the right n=5015 samples were further divided according to heroin use, with n=1913 having a score greater than 0.5 (any Heroin Use). At the third node on the right branch, samples were selected according to Tranquilizer medication use, with n=1419 scoring positively. On the left branch, the second node selected was Drug Treatment, with n=2844 respondents scoring positively that they had received Drug Treatment. Feature importance of the decision tree classifier selected Cocaine Use as the most informative feature for Prescription Opioid PRL Misuse. Following afterwards, Tranquilizer Use, Drug Treatment, and Heroin Use were tied for second place.

[Figure 9 about here.]

3.3.3 Random Forests Classifier. The Random Forest Classifier package in Scikit-Learn was used to classify Prescription Opioid PRL Misuse as the target variable, with 100 trees. The model accuracy for the training set was 0.955 and the test set accuracy was 0.896, which suggests that the model overfit the data. Figure 10 shows the feature importance of the random forests classifier for Prescription Opioid PRL Misuse. As Figure 10 shows, several features were identified as important for classifying Prescription Opioid PRL Misuse. The random forest selected Overall Health as the most informative feature in the model, followed by Cocaine Use, Education Level, Age Category, and Size of City Metropolitan region. Because of the additional features included as important, gradient boosting was performed to clarify the feature importance.

[Figure 10 about here.]

3.3.4 Boosted Gradient Classifier. The Gradient Boosting Classifier from Scikit-Learn was used to classify Prescription Opioid PRL Misuse, using the default setting of 100 trees, of maximum depth of 3, and a learning rate of 0.1. The model accuracy for the training set was 0.894 and accuracy for the test set was 0.893. Gradient boosting typically improves test set accuracy by using many simple models iteratively. In this case, model accuracy for gradient boosting was no better than random forests, and this is because the default parameter settings were used; further parameter tuning is needed to improve model performance. Feature importance was a primary interest for identifying features related to 'prescription opioid abuse. Figure 11 shows the feature importance for the gradient boosting classifier tree. As Figure 11 shows, several features were important for classifying prescription opioid misuse, and contrary to the random forests, gradient boosting selected Tranquilizer use as the most informative feature. Following closely in importance were Heroin Use and Age Category. Tied for fourth place were Cocaine Use and Treatment, with Mental Health (depression) coming in fourth in terms of feature importance. This result illustrates that several features are important for understanding Prescription Opioid Misuse, and the relations among features may be complex.

[Figure 11 about here.]

4 DISCUSSION

4.1 Summary of Findings

4.2 Extension to Big Data

The methods used in this project could be extended to better approximate big data for health analytics in the following ways: (1) Include a larger selection of features from the 2600 attributes in the NSDUH-2015 dataset; (2) Include survey data from previous years (e.g., 2005-2015); and (3) Extend the sample to the population of patients who have been prescribed opioid pain medication. There were many additional features that could have been included in the subset of features included in the project dataset. However, data cleaning and preparation can be a time consuming process, especially for datasets with a large number of features [13]. Additional data from the NSDUH was downloaded from previous years (2012 to 2014), and a preliminary examination of the data revealed inconsistencies in questions and prescription opioid medications that would need to be resolved in order to combine data from multiple years. In addition to data cleaning, there are several steps involved in the consolidation of data from multiple sources into a single dataset, which include extraction, integration, and aggregation of features. Unfortunately, time constraints for the project deadline did not allow for the inclusion of data from previous years into this analysis. A future study could integrate data from different years into the analysis or include data from multiple sources.

4.3 Limitations

To be of any use, diverse and often messy raw data has to be sifted through and effectively organized for further analysis, and

there are legitimate questions about the reliability of self report data from survey research for predicting actual behavior.

The question of Value evaluates the quality of the data as it pertains to intended outcomes, such as limiting the spread of contagion and disease prevention.

An important challenge for making sense of big data is developing analytic tools adequate to handle large volumes of data in real time.

4.4 Drug Abuse, Dependency, and Addiction

Drug addiction has many similar characteristics to other chronic medical illnesses, but there are unique challenges to the treatment of addiction [8, 20]. In drug rehabilitation treatment programs, patients undergo intense detoxification that reduces their drug tolerance, but are then released back into the environments associated with their drug use, putting them at high risk for relapse and potential drug overdose [6]. According to a classical conditioning model of addiction, situational cues or events can elicit a motivational state underlying relapse to drug use. Addictive behavior can be also be reinstated after extinction of dependency by exposure to drug-related cues or stressors in the environment [15].

4.5 Dynamics of Epidemic Spreading

If the prescription opioid crisis is a genuine epidemic, then we can conceive of it in terms of the dynamics of epidemic spreading which have been developed based on models of contagious disease. Epidemic spreading is a dynamic process based on networks of direct

person-to-person contact and indirect exposure via transportation pathways [2], that facilitate the distribution of opioid medications or illicit drugs. Instead of thinking about persons as infected or uninfected by biological contagion, in the opioid drug model, we must consider individuals as dependent, addicted or susceptible to dependence and addiction. Epidemics are quantified in terms of the proportion of the population infected, those yet to be infected, and the rate of transmission. Furthermore, the structure of the contact network can influence epidemic spreading [12]. For example, in the case of simple contagion, weak ties among acquaintances or infrequent associations provide shortcuts between distant nodes that reduce distance within the network [?] which can facilitate the spread of contagion, or in this case drug use. Furthermore, opioid contact networks may have “small world” properties where a small number of nodes or people have a very number of connections that can rapidly transmit contagion throughout the network [?]. It may be possible to apply network analysis to identify underlying structure of the contact network of opioid use and addiction, to identify pathways and points of contact between nodes or person in the spreading use, misuse, and abuse of prescription opioid medications. Future research could apply social network modeling to the opioid crisis in order to identify how drug dependency and addiction are subserved by patterns of social interaction.

5 CONCLUSION

Several machine learning methods were used to classify heroin use and prescription opioid misuse and abuse. The results of this analysis are somewhat inconclusive, given that the direction of these effects is unknown. On the one hand there is evidence that individual who reported having used heroin were also more likely to report misusing or abusing prescription opioids. On the other hand, the proportion of individuals who misused or abused prescription opioids, and also reported using heroin, was twice as large as the proportion who reported only using heroin. A general conclusion is that the results provide partial support for the hypothesis that taking prescription opioids leads to a higher likelihood of illicit opioid use. However, the results did not provide sufficient evidence to rule out the alternative hypothesis that people who have used heroin may have a propensity for opioid use therefore be more likely to become dependent on prescription opioid medications. Given that the number of individuals who reported using heroin in this sample was low, additional data may help to provide evidence to resolve this question. A limitation of survey data is there may be bias in self-reports of illicit drug use, as it is a proscribed and illegal behavior, and therefore the data may underestimate the actual rate of heroin use in the general population. Including additional data from previous years may provide a more tests of these hypotheses.

Machine of opioid abuse can contribute to efforts to address prescription opioid addiction, overdoses, in the following ways:

- (1) Identify factors related to opioid dependency
- (2) Inform consumers of opioid medication as to risk factors
- (3) Increase knowledge of opioid abuse for more informed prescriptions.

ACKNOWLEDGMENTS

The author would like to thank Dr. Gregor von Laszewski, the Teaching Assistants, Juliette Zurick, Miao Jiang, Hungri Lee, Grace Li, Saber Sheybani Moghadam, and others who helped to improve this project and report.

REFERENCES

- [1] Substance Abuse, Center for Behavioral Health Statistics Mental Health Services Administration, and Quality. 2016. *National Survey on Drug Use and Health (NSDUH) 2015*. Online data archive. United States Department of Health and Human Services., Ann Arbor, MI. <https://doi.org/10.3886/ICPSR50011.v1>
- [2] Vittoria Colizza, Alain Barrat, Marc Barthélémy, and Alessandro Vespignani. 2006. The role of the airline transportation network in the prediction and predictability of global epidemics. *Proceedings of the National Academy of Sciences of the United States of America* 103, 7 (2006), 2015–2020. <https://doi.org/10.1073/pnas.0510525103> arXiv:<http://www.pnas.org/content/103/7/2015.full.pdf>
- [3] Centers for Disease Control and Prevention. 2017. Prescription Opioid Overdose Data. online. (Oct. 2017). <https://www.cdc.gov/drugoverdose/data/overdose.html>
- [4] hd1 and yoavram. 2016. Python: Download Returned Zip file from URL. Online. (Feb. 2016). <https://stackoverflow.com/questions/9419162/python-download-returned-zip-file-from-url> Stackoverflow.com.
- [5] M. Herland, T. M. Khoshgoftaar, and R. Wald. 2014. A review of data mining using big data in health informatics. *Journal Of Big Data* 1, 2 (2014). <https://doi.org/10.1186/2196-1115-1-2>
- [6] K. Johnson, A. Isham, D.V. Shah, and D.H. Gustafson. 2011. Potential Roles for New Communication Technologies in Treatment of Addiction. *Current psychiatry reports*. (2011). <https://doi.org/10.1007/s11920-011-0218-y>
- [7] Rose A. Judd, Noah Aleshire, Jon E. Zibbell, and R. Matthew Gladden. 2016. *Increases in Drug and Opioid Overdose Deaths, United States, 2000-2014*. techreport 64(50). Centers for Disease Control and Prevention, Atlanta, GA. <https://www.cdc.gov/mmwr/preview/mmwrhtml/mm6450a3.htm> Morbidity and Mortality Weekly Report (MMWR).
- [8] Lisa A. Marsch. 2012. Leveraging technology to enhance addiction treatment and recovery. *Journal of Addictive Diseases* 31, 3 (2012), 313–318. <https://doi.org/10.1080/10550887.2012.694606>
- [9] Wes McKinney. 2017. *Python for Data Analysis*. O'Reilly Media Inc., Sebastopol, CA. <https://github.com/wesm/pydata-book>
- [10] Andreas C. Müller and Sarah Guido. 2017. *Introduction to Machine Learning*. O'Reilly, Sebastopol, CA. https://github.com/amueller/introduction_to_ml_with_python/
- [11] National Institute on Drug Abuse (NIDA). 2017. *Overdose Death Rates*. Summary. National Institutes of Health (NIH), Washington D.C. <https://www.drugabuse.gov/related-topics/trends-statistics/overdose-death-rates>
- [12] Romualdo Pastor-Satorras and Alessandro Vespignani. 2001. Epidemic Spreading in Scale-Free Networks. *Phys. Rev. Lett.* 86 (Apr 2001), 3200–3203. Issue 14. <https://doi.org/10.1103/PhysRevLett.86.3200>
- [13] E. Rahm and H. H. Hai Do. 2000. *Data cleaning: Problems and current approaches*. techreport 23(4). Bulletin of the Technical Committee on Data Engineering, 1730 Massachusetts Avenue, Washington D.C. <https://s3.amazonaws.com/academia.edu.documents/41858217/A00DEC-CD.pdf?AWSAccessKeyId=AKIAIWOWYYGZ2Y53UL3A&Expires=1511155930&Signature=VWRM7u4KwTp6ZxX5jB%2Bh6wMCbpg%3D&response-content-disposition=inline%3B%20filename%3DAutomaticaly-extracting.structure.from.pdf#page=5>
- [14] Sebastian Raschka and Vahid Mirjalili. 2017. *Ptyhon Machine Learning, Second Edition*. Packt, Birmingham, UK. <https://github.com/rasbt/python-machine-learning-book-2nd-edition>
- [15] Yavin Shaham, Uri Shalev, Lin Lu, Harriet de Wit, and Jane Stewart. 2003. The reinstatement model of drug relapse: history, methodology and major findings. *Psychopharmacology* 168, 1 (01 Jul 2003), 3–20. <https://doi.org/10.1007/s00213-002-1224-8>
- [16] S.M. Shiverick. 2017. BDA Project Data. github. (Nov. 2017). <https://github.com/bigdata-i523/hid335/blob/master/project/BDA-Project-Data.ipynb>
- [17] S.M. Shiverick. 2017. Classification Models of Heroin Use. github. (Nov. 2017). <https://github.com/bigdata-i523/hid335/blob/master/project/BDA-Analytics-Classifier-Heroin.ipynb> Interactive Python Jupyter Notebook.
- [18] S.M. Shiverick. 2017. Classification Models of Prescription Opioid Pain Relievers Misuse. github. (Nov. 2017). <https://github.com/bigdata-i523/hid335/blob/master/project/BDA-Analytics-Classifier-PRL.ipynb> Interactive Python Jupyter Notebook.
- [19] S.M. Shiverick. 2017. Project Codebook for Data Variables from NSDUH-2015. github. (Nov. 2017). <https://github.com/bigdata-i523/hid335/blob/master/project/project-data-codebook.txt>
- [20] J. Swendsen. 2016. Contributions of mobile technologies to addiction research. *Dialogues Clinical Neuroscience* 18, 2 (June 2016), 213–221. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4969708/>
- [21] Jake VanderPlas. 2017. *Python Data Science Handbook*. O'Reilly Media Inc., Sebastopol, CA. <https://jakevdp.github.io/PythonDataScienceHandbook/>
- [22] Upkar Varshney. 2013. Smart medication management system and multiple interventions for medication adherence. *Decision Support Systems* 55, 5 (May 2013), 538–551. <https://doi.org/10.1016/j.dss.2012.10.011>
- [23] Nora D. Volkow, Thomas R. Frieden, Pamela S. Hyde, and Stephen S. Cha. 2014. Medication-Assisted Therapies: Tackling the Opioid-Overdose Epidemic. *New England Journal of Medicine* 370, 22 (2014), 2063–2066. <https://doi.org/10.1056/NEJMmp1402780> PMID: 24758595.

A CODE REFERENCES

All code, notebooks, files, and folders for this project can be found in the i523/hid335/project githup repository: url:<https://github.com/bigdata-i523/hid335/tree/master/project>.

A.1 Download and Extract Data file

The ‘get-data.py’ function was written to download the data, unzip the data files, extract the data, and write the NSDUH-2015 dataset to CSV file [4].

A.2 Data Cleaning and Preparation

Data cleaning and preparation steps was conducted using an interactive python Jupyter Notebook [16] based on examples in Python for Data Analysis [9] and the Python Data Science Handbook [21].

A.3 Exploratory Data Analysis

Exploratory Data Analysis and Visualization was conducted using an interactive python notebook: CITE URL

based on examples from Python for Data Analysis [9], and the Python Data Science Handbook [21].

A.4 Machine Learning Classifier Algorithms

Machine learning classification models included logistic regression classifier, decision Tree classifier, random forests classifier, and gradient boosting classifier were constructed using SciKit Learn [10, 14] using two separate interactive python jupyter notebook, one for classifying Heroin Use as the target variable [17], and another notebook for classifying Prescription Opioid Misuse as the target [18].

LIST OF FIGURES

1	Proportion of Individuals Who Reported Ever Misusing Prescription Opioid Pain Relievers and Proportion Who Reported Using Heroin	10
2	Plot of Opioid Pain Medication Misuse and Abuse and Heroin Use with Regression Slopes Weighted by Metropolitan Area Size	11
3	Pairplots of Mental Health, Prescription Opioid Misuse and Abuse, Heroin Use, and Size of City Metropolitan Area	12
4	Coefficients of Logistic Regression Classifier of Heroin Use (With L1 Penalty and Values of Regularization Parameter C)	13
5	Decision Tree Classification of Heroin Use (Partial View)	14
6	Feature Importance for Random Forests Classifier for Heroin Use	15
7	Feature Importance for Gradient Boosting Classifier for Heroin Use	16
8	Logistic Regression Classification of Prescription Opioid (PRL) Misuse with L2 Penalty	17
9	Decision Tree for Prescription Opioid (PRL) Misuse	18
10	Feature Importance for Random Forest Classifier of Prescription Opioid (PRL) Misuse	19
11	Feature Importance for Gradient Boosted Classifier Tree of Prescription Opioid (PRL) Misuse	20

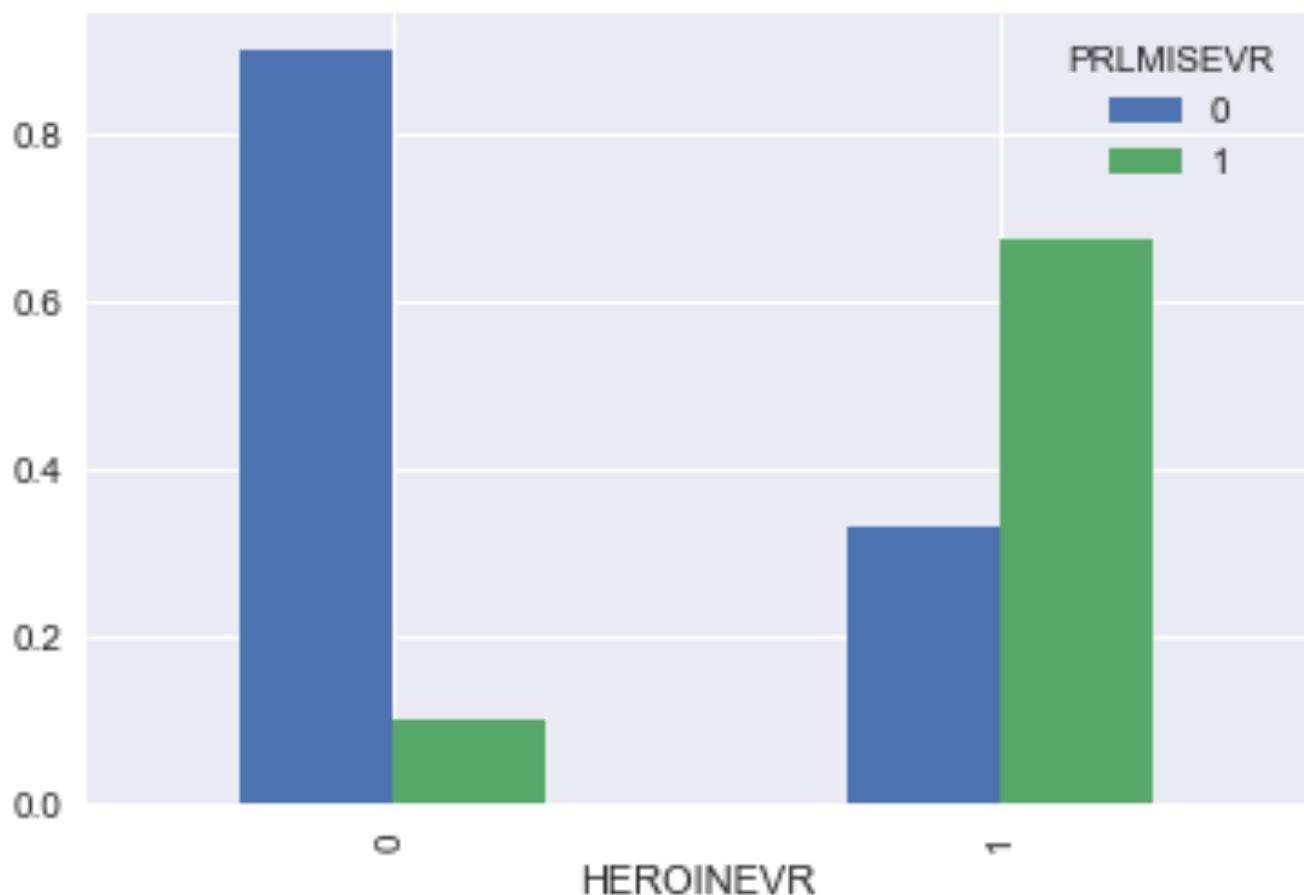


Figure 1: Proportion of Individuals Who Reported Ever Misusing Prescription Opioid Pain Relievers and Proportion Who Reported Using Heroin

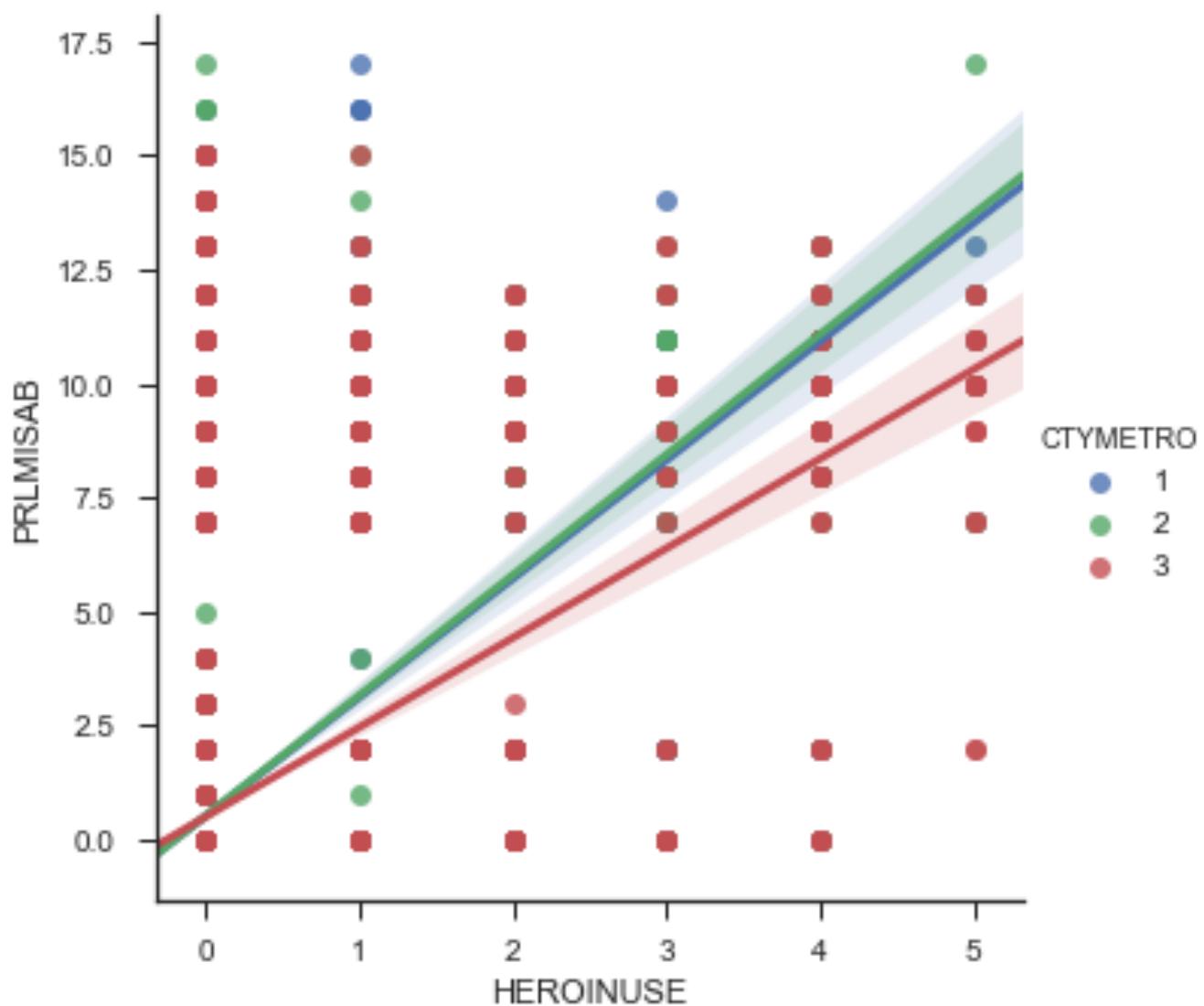


Figure 2: Plot of Opioid Pain Medication Misuse and Abuse and Heroin Use with Regression Slopes Weighted by Metropolitan Area Size

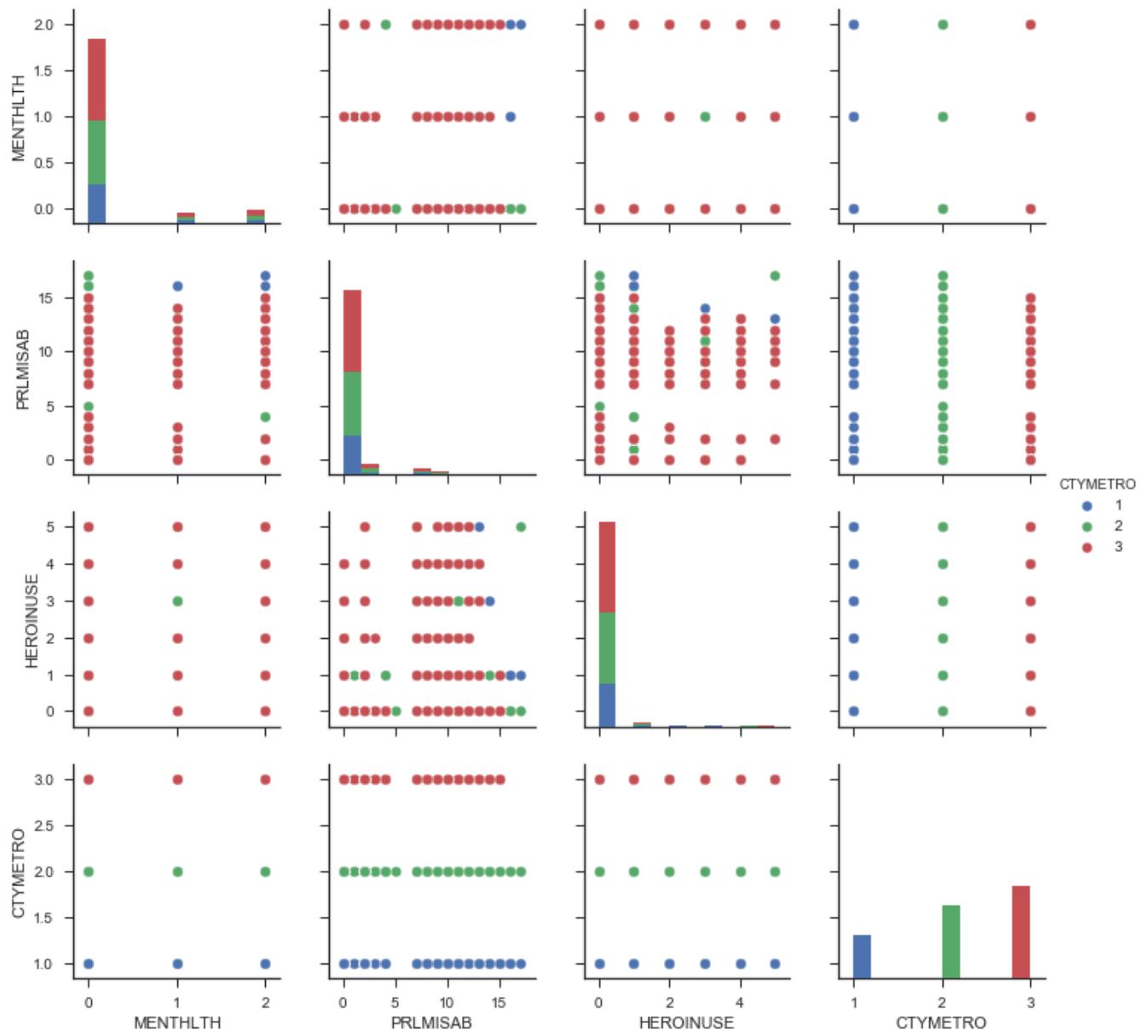


Figure 3: Pairplots of Mental Health, Prescription Opioid Misuse and Abuse, Heroin Use, and Size of City Metropolitan Area

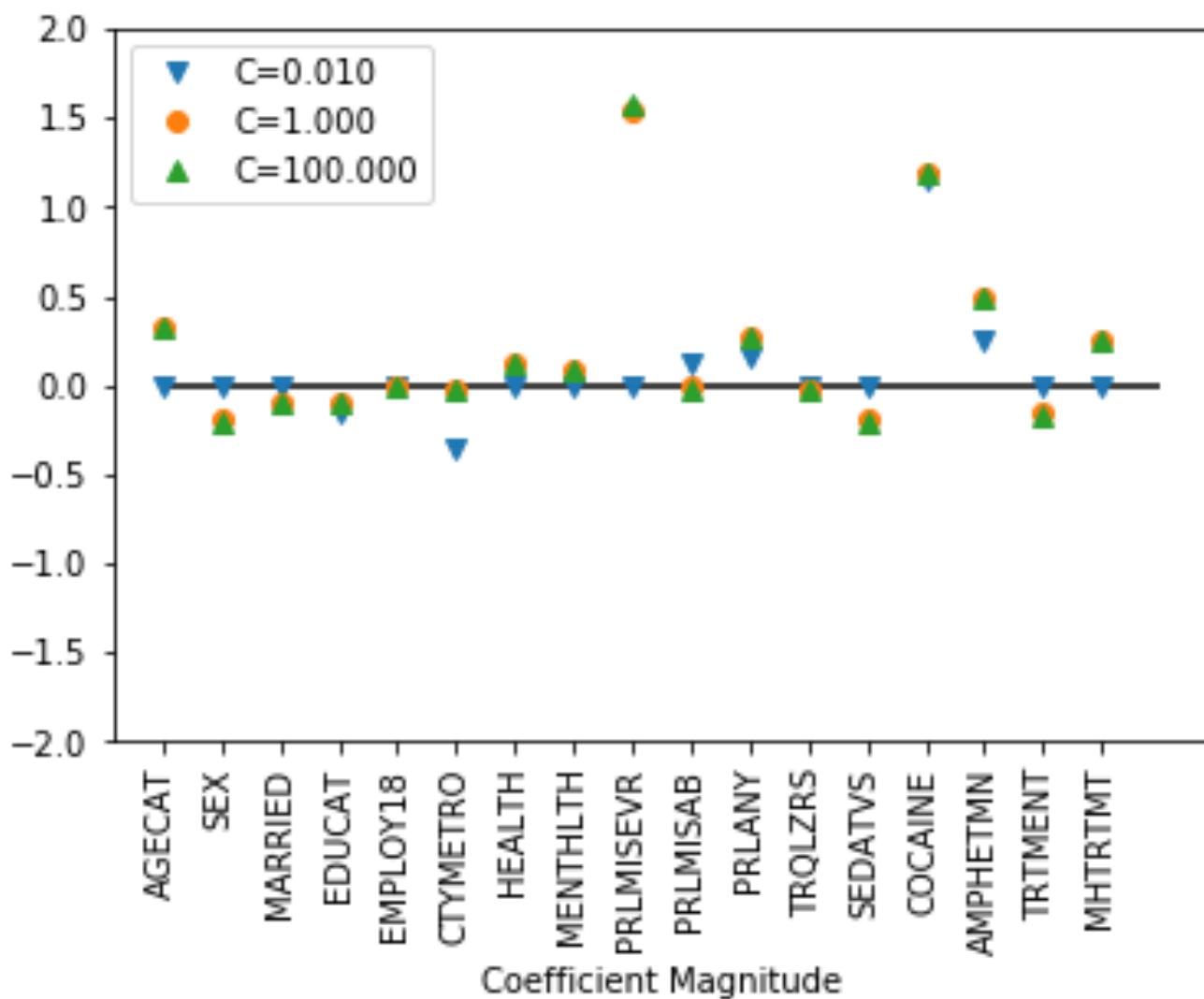


Figure 4: Coefficients of Logistic Regression Classifier of Heroin Use (With L1 Penalty and Values of Regularization Parameter C)

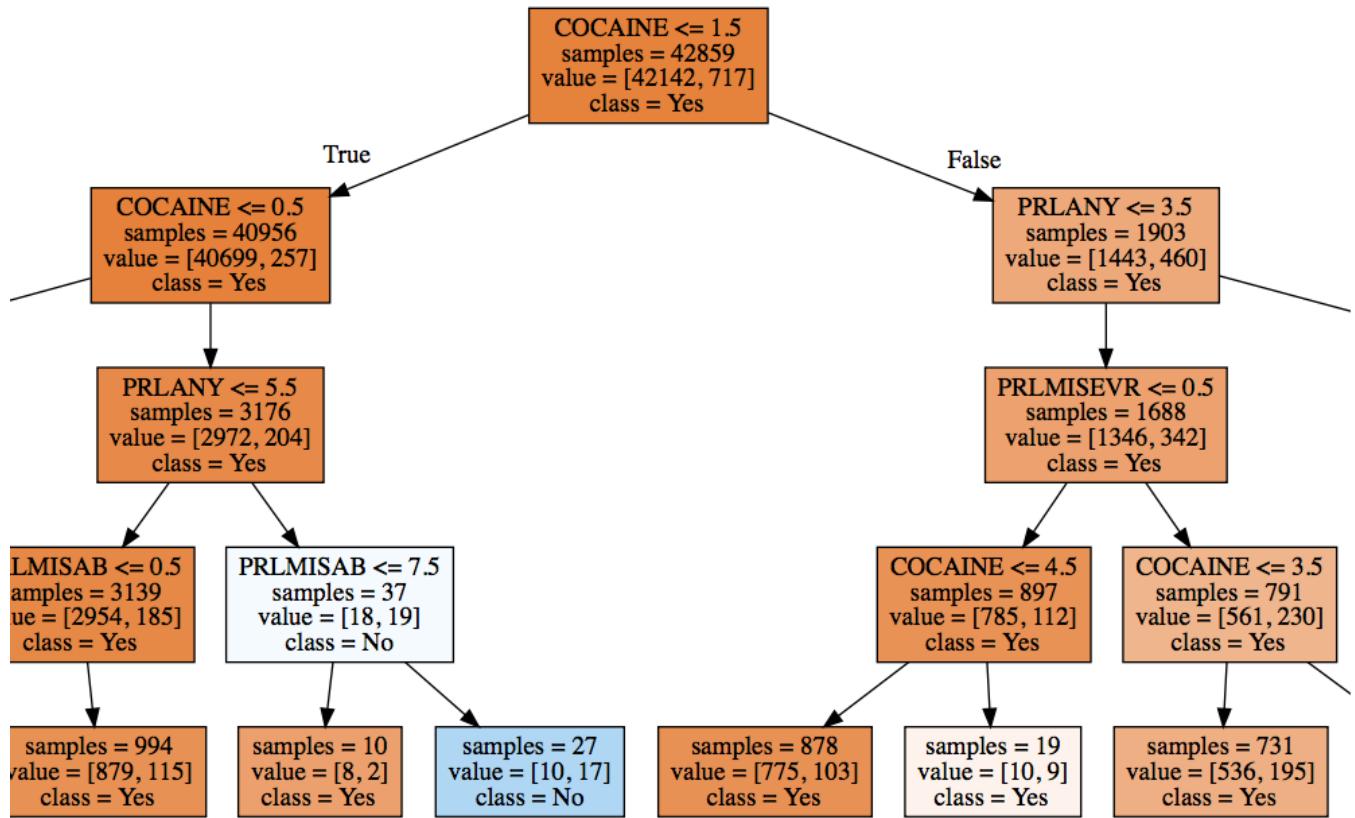


Figure 5: Decision Tree Classification of Heroin Use (Partial View)

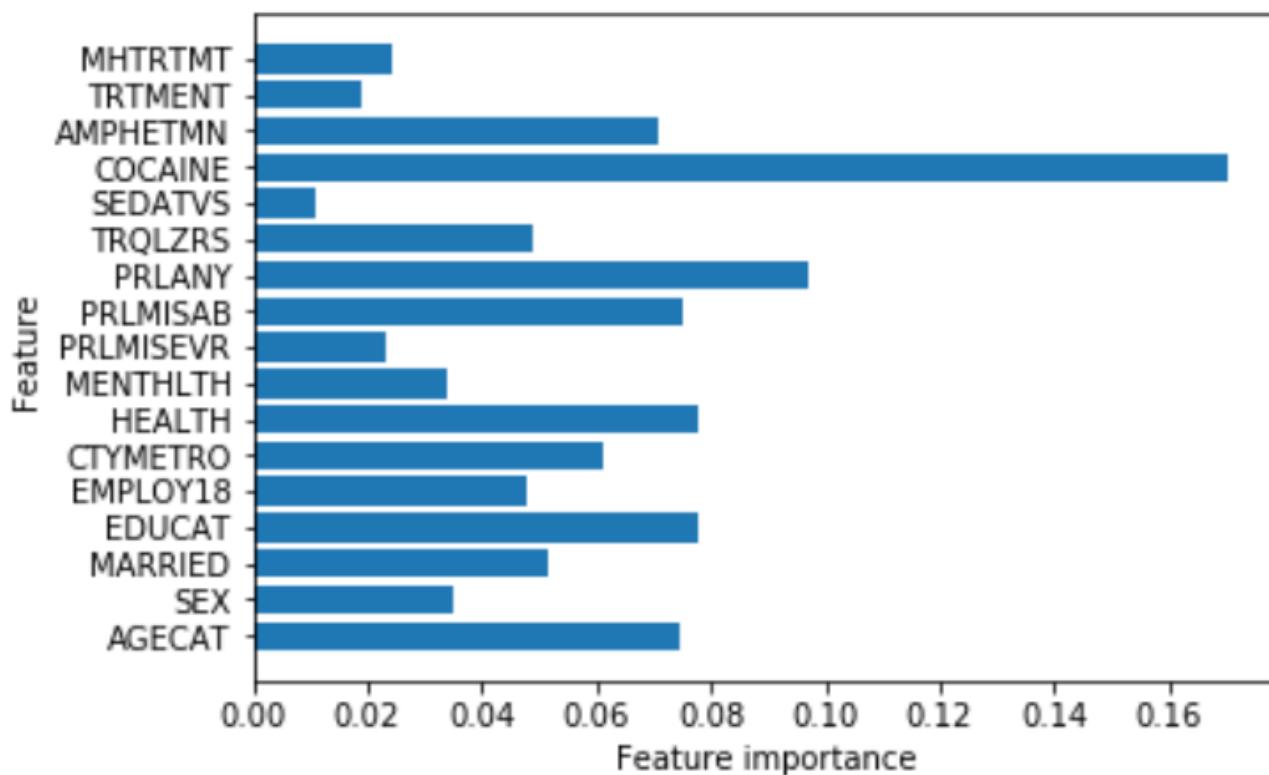


Figure 6: Feature Importance for Random Forests Classifier for Heroin Use

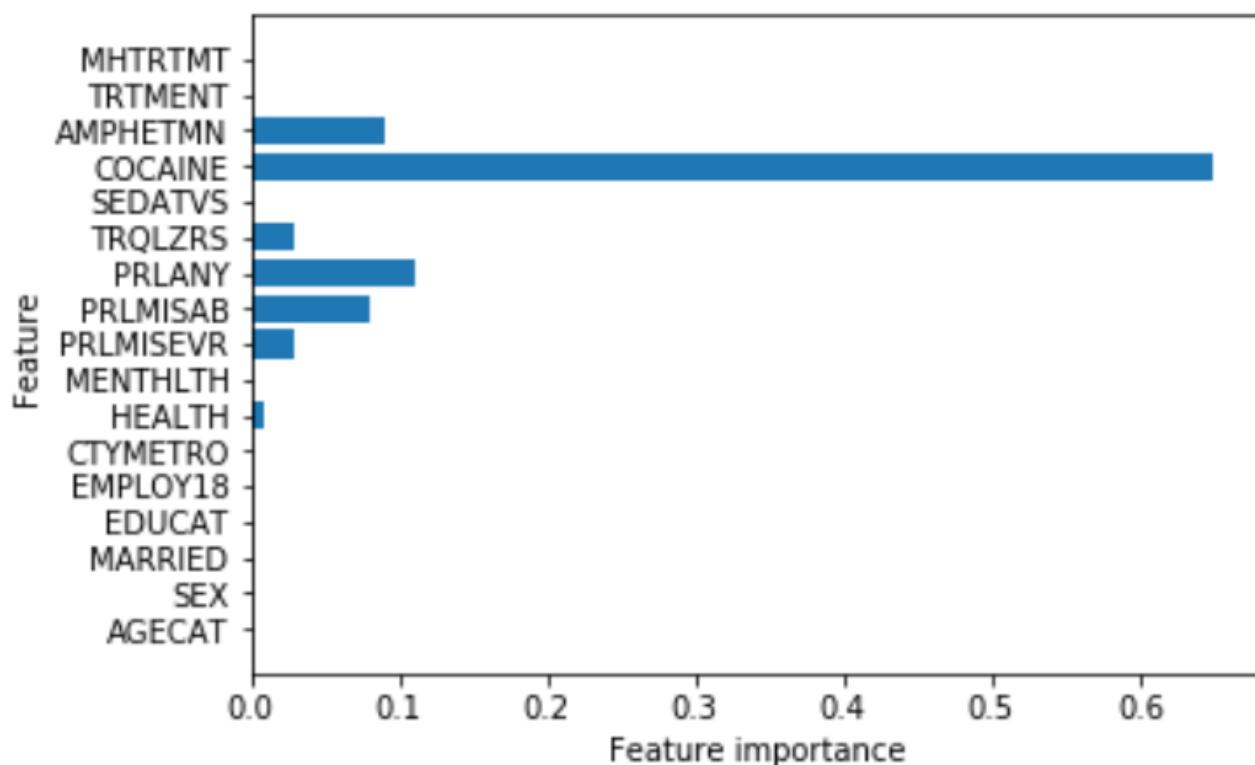


Figure 7: Feature Importance for Gradient Boosting Classifier for Heroin Use

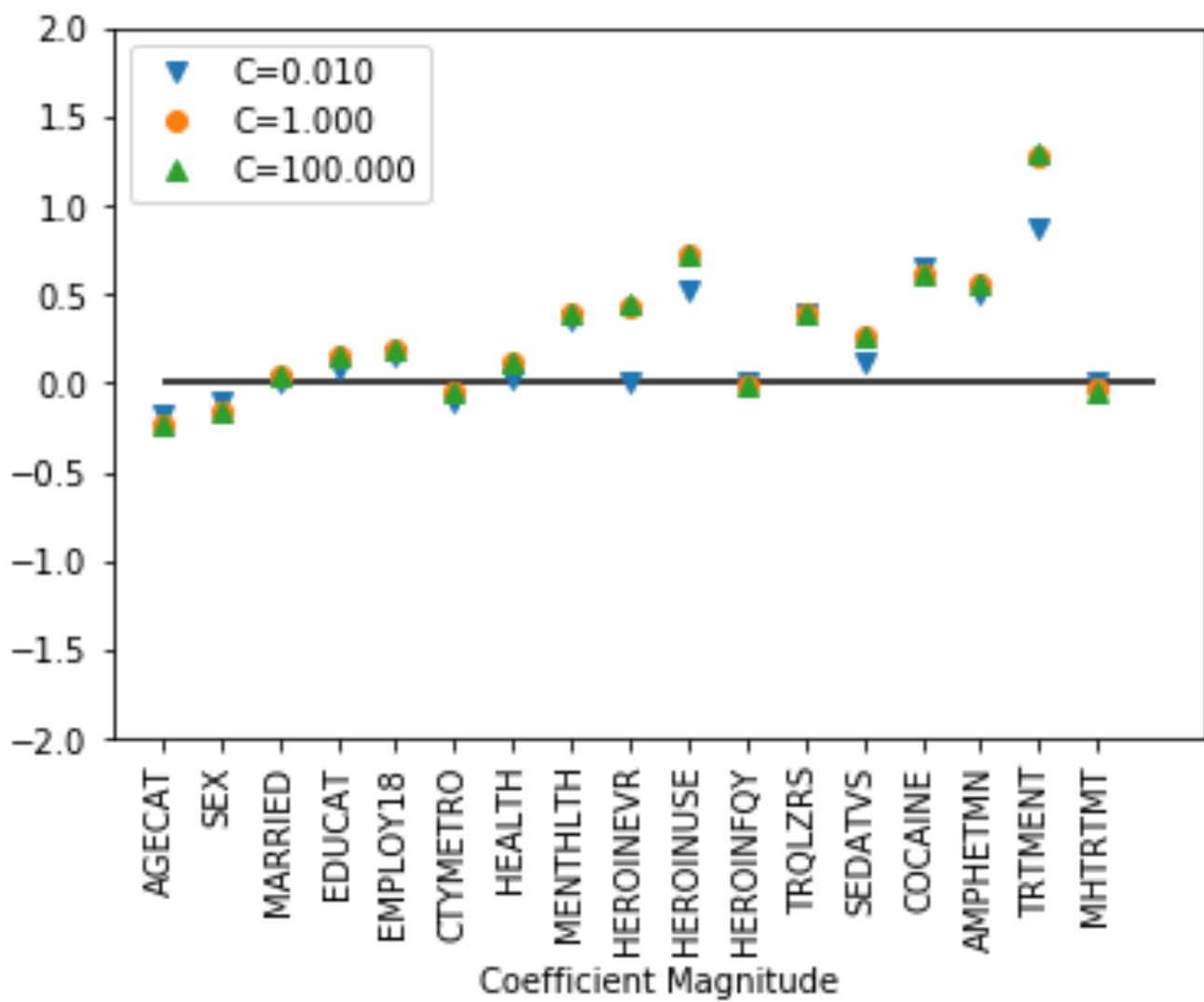


Figure 8: Logistic Regression Classification of Prescription Opioid (PRL) Misuse with L2 Penalty

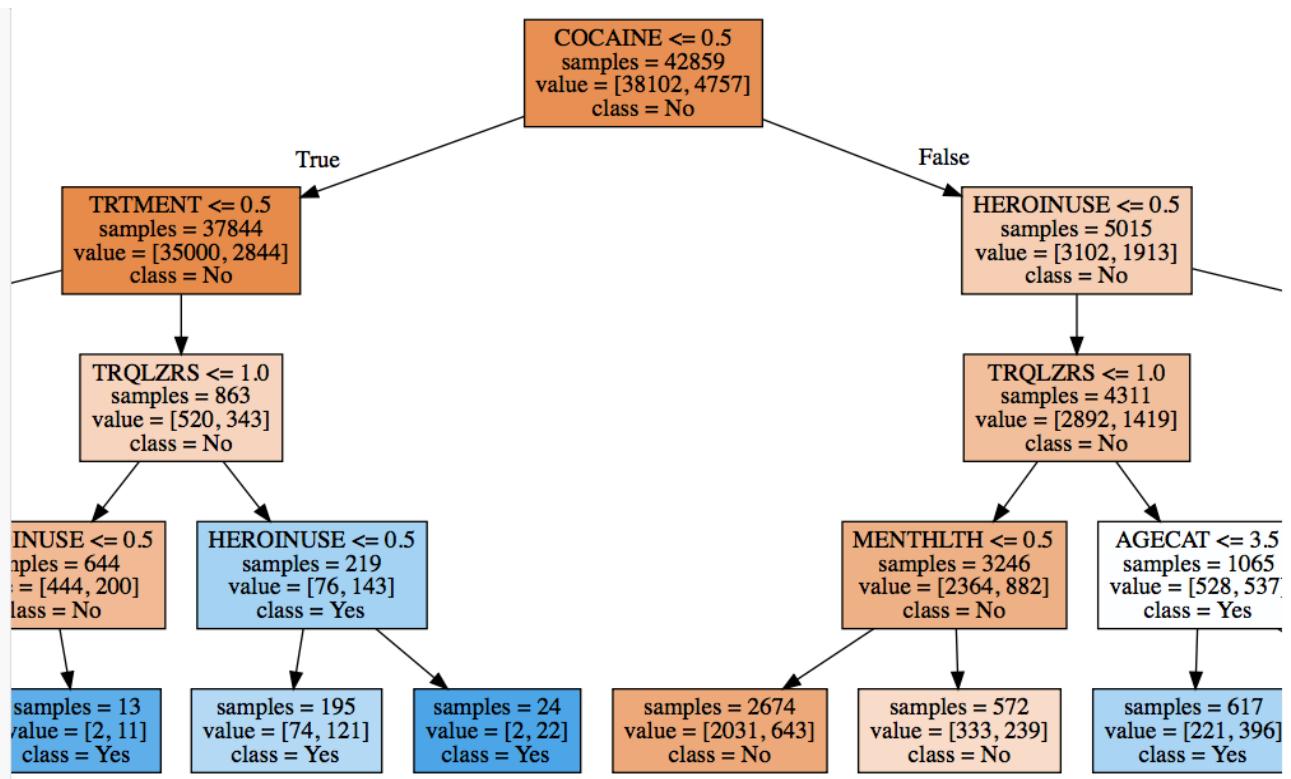


Figure 9: Decision Tree for Prescription Opioid (PRL) Misuse

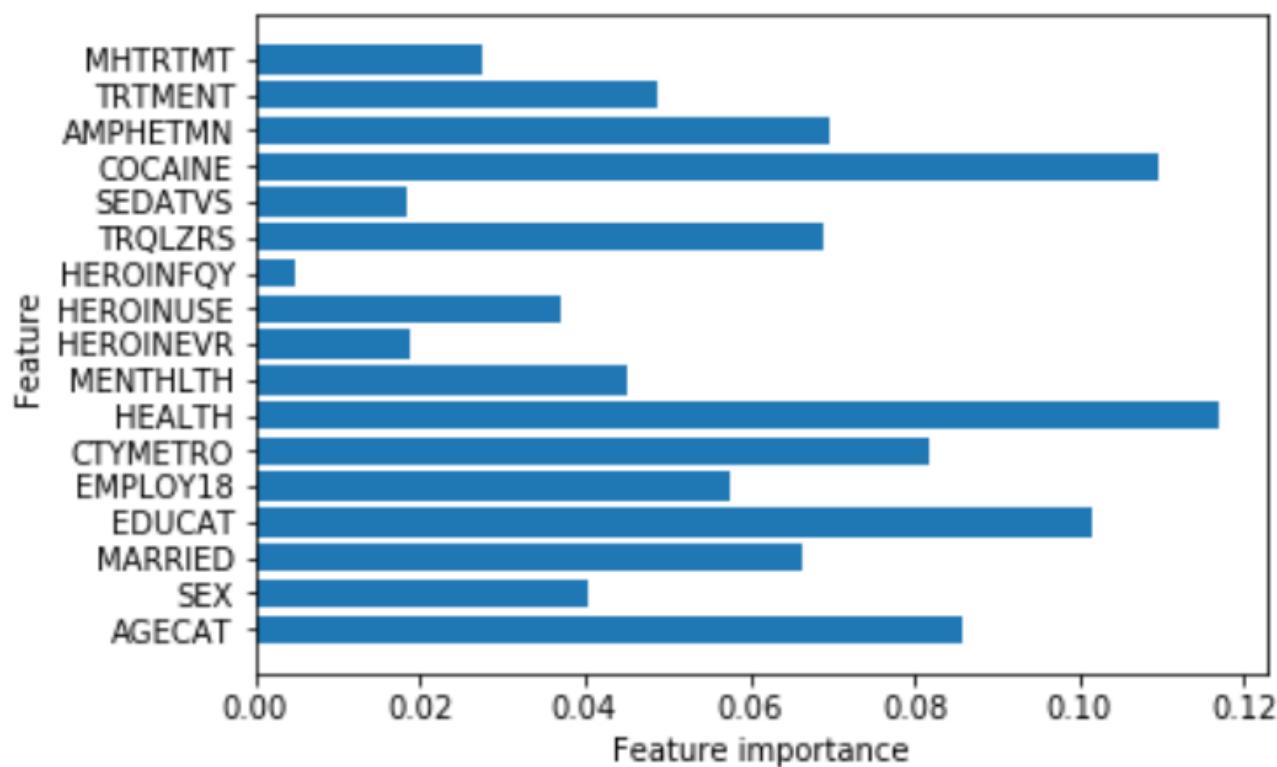


Figure 10: Feature Importance for Random Forest Classifier of Prescription Opioid (PRL) Misuse

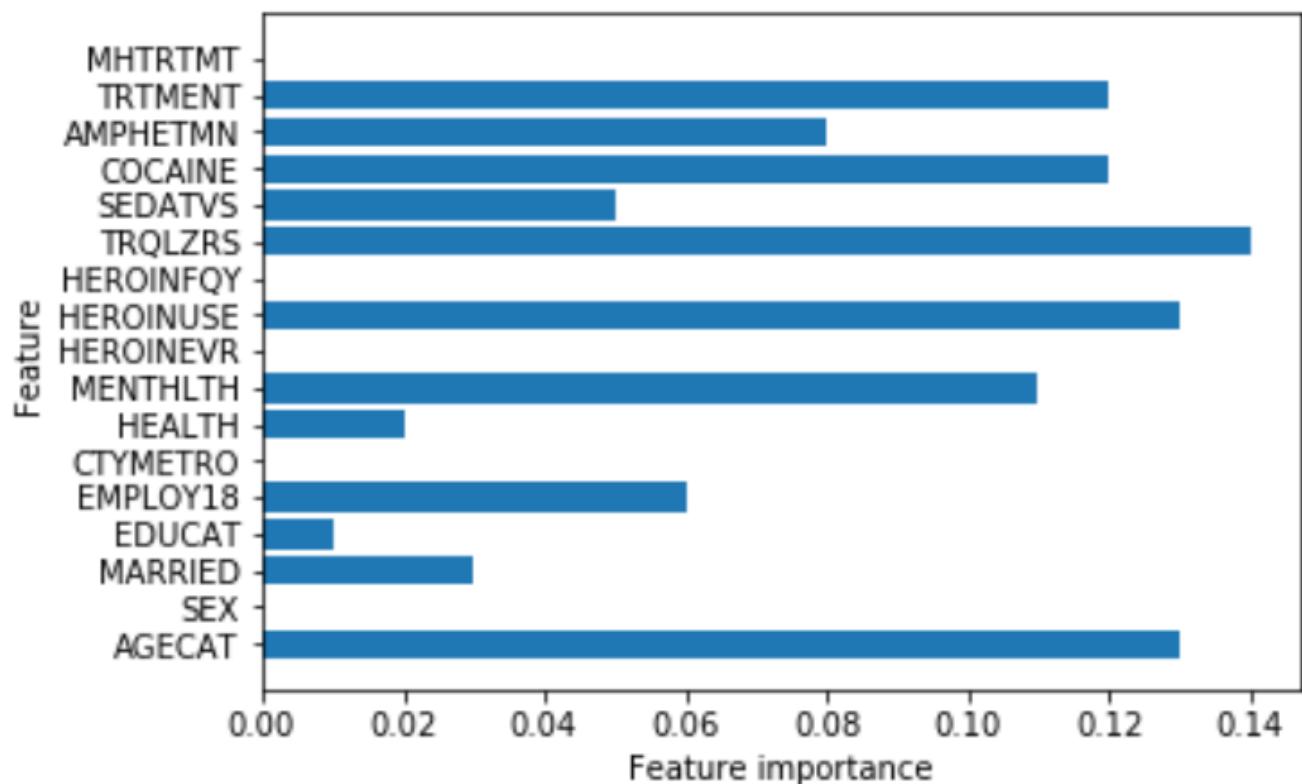


Figure 11: Feature Importance for Gradient Boosted Classifier Tree of Prescription Opioid (PRL) Misuse

LIST OF TABLES

1	Substance Use by Age Group Counts - NSDUH 2015 [1]	22
2	Frequency Table of Mental Health Issues and Treatment NSDUH 2015 [1]	22

Table 1: Substance Use by Age Group Counts - NSDUH 2015 [1]

Age Group	12-17	18-25	26-34	35-49	50+
Sample Size	13585	14553	9084	11169	8755
Oxycodone	545	1632	1132	1345	1044
Hydrocodone	831	2936	2233	2781	2103
Tramadol	241	753	654	829	734
Morphine	251	431	236	313	286
Fentanyl	28	97	81	96	86
Demerol	26	74	49	64	71
Buprenorphine	43	197	167	124	51
Oxymorphone	46	88	57	47	41
Hydromorphone	24	94	107	118	81
PRL Misuse Ever*	798	2127	1475	1343	600
Tranquilizers	405	1469	1064	1405	1153
Sedatives	204	242	157	256	226
Methadone Ever	32	83	96	71	46
Heroin Use Ever*	22	261	259	250	164
Cocaine Use Ever	109	1645	1626	1954	1406
Amphetamines Ever	932	1836	627	383	164
Methamphetamine	42	481	700	898	492
Hallucinogens	450	2660	2020	2127	1197
LSD Use Ever	190	1114	874	1442	907
Ecstasy (MDMA)	199	1867	1403	947	149

Table 2: Frequency Table of Mental Health Issues and Treatment NSDUH 2015 [1]

Age Group	12-17	18-25	26-34	35-49	50+
In Hospital Overnight	730	1149	821	890	1173
Adult Depression	0	2413	1395	1766	967
Suicidal Thoughts	13585	14553	9084	11189	8755
Mental Health Treatment					
Private Therapist	0	592	434	554	311
Treatment Gap*	469	931	321	239	90

```
bibtext report
```

```
=====
This is BibTeX, Version 0.99d (TeX Live 2016)
The top-level auxiliary file: report.aux
The style file: ACM-Reference-Format.bst
Database file #1: report.bib
Warning--I didn't find a database entry for "granovetter73"
Warning--I didn't find a database entry for "watts98"
Warning--page numbers missing in both pages and numpages fields in herland14
Warning--no number and no volume in johnson11
Warning--page numbers missing in both pages and numpages fields in johnson11
(There were 5 warnings)
```

```
bibtext _ label error
```

```
=====
bibtext space label error
```

```
=====
bibtext comma label error
```

```
=====
latex report
```

```
[2017-12-04 12.24.07] pdflatex report.tex
This is pdfTeX, Version 3.14159265-2.6-1.40.17 (TeX Live 2016) (preloaded format=pdflatex)
Label `tab:freq' multiply defined.
p.7   L870  : [granovetter73] undefined
p.7   L874  : [watts98] undefined
Missing character: ""
Missing character: ""
There were undefined citations.
bookmark level for unknown defaults to 0.
The anchor of a bookmark and its parent's must not be the same. Added a new anchor.
The anchor of a bookmark and its parent's must not be the same. Added a new anchor.
There were multiply-defined labels.
Typesetting of "report.tex" completed in 1.1s.
./README.yml
13:5      error      wrong indentation: expected 8 but found 4 (indentation)
15:5      error      wrong indentation: expected 8 but found 4 (indentation)
```

```
39:5      error    wrong indentation: expected 8 but found 4  (indentation)
42:5      error    wrong indentation: expected 8 but found 4  (indentation)
47:80     error    trailing spaces  (trailing-spaces)
48:77     error    trailing spaces  (trailing-spaces)
49:79     error    trailing spaces  (trailing-spaces)
50:35     error    trailing spaces  (trailing-spaces)
54:5      error    duplication of key "chapter" in mapping  (key-duplicates)
57:4      error    wrong indentation: expected 4 but found 3  (indentation)
57:17     error    trailing spaces  (trailing-spaces)
59:4      error    wrong indentation: expected 7 but found 3  (indentation)
61:4      error    wrong indentation: expected 7 but found 3  (indentation)
63:81     error    line too long (94 > 80 characters)  (line-length)
65:81     error    line too long (83 > 80 characters)  (line-length)
65:83     error    trailing spaces  (trailing-spaces)
66:81     error    line too long (86 > 80 characters)  (line-length)
66:86     error    trailing spaces  (trailing-spaces)
67:81     error    line too long (85 > 80 characters)  (line-length)
67:85     error    trailing spaces  (trailing-spaces)
68:79     error    trailing spaces  (trailing-spaces)
69:81     error    line too long (87 > 80 characters)  (line-length)
69:87     error    trailing spaces  (trailing-spaces)
70:81     error    line too long (83 > 80 characters)  (line-length)
70:83     error    trailing spaces  (trailing-spaces)
71:81     error    line too long (81 > 80 characters)  (line-length)
71:81     error    trailing spaces  (trailing-spaces)
72:80     error    trailing spaces  (trailing-spaces)
73:81     error    line too long (82 > 80 characters)  (line-length)
73:82     error    trailing spaces  (trailing-spaces)
74:76     error    trailing spaces  (trailing-spaces)
75:51     error    trailing spaces  (trailing-spaces)
84:33     error    trailing spaces  (trailing-spaces)
```

Compliance Report

name: Sean Shiverick
hid: 335
paper1: 10/25/17 100%
paper2: 100%
project: 95%

yamlcheck

```
wordcount
```

```
(null)
wc 335 project (null) 7694 report.tex
wc 335 project (null) 8042 report.pdf
wc 335 project (null) 997 report.bib
```

```
find "
```

```
passed: True
```

```
find footnote
```

```
passed: True
```

```
find input{format/i523}
```

```
4: \input{format/i523}
```

```
passed: True
```

```
find input{format/final}
```

```
passed: False
```

```
floats
```

```
345: \begin{table}
348: \label{tab:freq}
392: \begin{table}
395: \label{tab:freq}
425: \begin{figure} [!ht]
426: \centering\includegraphics[width=\columnwidth]{images/Figure1.pdf}
    }
429: \label{f:Figure1}
453: \begin{figure} [!ht]
454: \centering\includegraphics[width=\columnwidth]{images/Figure2.pdf}
    }
```

```

457: \label{f:Figure2}
492: \begin{figure} [!ht]
493: \centering\includegraphics[width=\columnwidth]{images/Figure3.pdf}
}
496: \label{f:Figure3}
542: \begin{figure} [!ht]
543: \centering\includegraphics[width=\columnwidth]{images/Figure4.pdf}
}
546: \label{f:Figure4}
591: \begin{figure} [!ht]
592: \centering\includegraphics[width=\columnwidth]{images/Figure5.pdf}
}
594: \label{f:Figure5}
629: \begin{figure} [!ht]
630: \centering\includegraphics[width=\columnwidth]{images/Figure6.pdf}
}
632: \label{f:Figure6}
666: \begin{figure} [!ht]
667: \centering\includegraphics[width=\columnwidth]{images/Figure7.pdf}
}
669: \label{f:Figure7}
706: \begin{figure} [!ht]
707: \centering\includegraphics[width=\columnwidth]{images/Figure8.pdf}
}
711: \label{f:Figure8}
737: \begin{figure} [!ht]
738: \centering\includegraphics[width=\columnwidth]{images/Figure9.pdf}
}
740: \label{f:Figure9}
757: \begin{figure} [!ht]
758: \centering\includegraphics[width=\columnwidth]{images/Figure10.pdf}
}
761: \label{f:Figure10}
786: \begin{figure} [!ht]
787: \centering\includegraphics[width=\columnwidth]{images/Figure11.pdf}
}
790: \label{f:Figure11}

```

```

figures 11
tables 2
includegraphics 11
labels 13
refs 0
floats 13

```

```
True : ref check passed: (refs >= figures + tables)
```

```
True : label check passed: (refs >= figures + tables)
True : include graphics passed: (figures >= includographics)
False : check if all figures are referred to: (refs >= labels)
```

Label/ref check

89: abusing prescribed opioid medication who also used heroin, shown
in Figure 1.

324: however, only 956 respondents had used heroin (570 males, 386
females). Table 1

384: Table 2 shows the frequency of individuals reporting that they
had experienced

413: Figure 1 shows the proportion of individuals who reported
misusing prescription

415: Figure 1 shows the majority of respondents (89 percent) stated
they had never

418: Figure 1 shows that, of those individuals who reported using
heroin, the

433: Figure 2 shows the aggregated measure of Opioid Pain Reliever
misuse and abuse

461: Figure 3 shows the pairplots of demographic features including
mental health

527: few features. Figure 4 plots the coefficients of logistic
regression classifier

533: accuracy of 0.98 and test set accuracy of 0.98. Figure 4 shows
that the

564: Figure 5 shows a partial view of the decision tree classifier of
heroin use

617: feature importance for a single tree. Figure 6 shows the feature
importance

657: or by reducing the learning rate. Figure 7 shows that the feature
importance

695: Figure 8 plots the coefficients of logistic regression classifier
for

697: Figure 8, the features with the highest coefficient values were
Treatment

720: of the pruned tree was 0.902 and test set accuracy was 0.902.
Figure 9 shows

723: notebook) \cite{classifyPRL}. As Figure 9 shows, the decision
tree classifier

748: 0.896, which suggests that the model overfit the data. Figure 10
shows the

748: 0.896, which suggests that the model overfit the data. Figure 10
shows the

750: PRL Misuse. As Figure 10 shows, several features were identified
as important

750: PRL Misuse. As Figure 10 shows, several features were identified

```
    as important
775: prescription opioid abuse. Figure 11 shows the feature importance
for the
775: prescription opioid abuse. Figure 11 shows the feature importance
for the
776: gradient boosting classifier tree. As Figure 11 shows, several
features were
776: gradient boosting classifier tree. As Figure 11 shows, several
features were
passed: False -> labels or refs used wrong
```

When using figures use columnwidth
[width=1.0\columnwidth]
do not change the number to a smaller fraction

```
find textwidth
```

```
passed: True
```

```
below_check
```

WARNING: algorithm and below may be used improperly

```
126: classification algorithms are considered below.
```

```
bibtex
```

```
label errors
```

```
bibtex errors
```

```
This is BibTeX, Version 0.99d (TeX Live 2016)
The top-level auxiliary file: report.aux
The style file: ACM-Reference-Format.bst
Database file #1: report.bib
Warning--I didn't find a database entry for "granovetter73"
Warning--I didn't find a database entry for "watts98"
Warning--page numbers missing in both pages and numpages fields in herland14
Warning--no number and no volume in johnson11
```

```
Warning--page numbers missing in both pages and numpages fields in johnson11
(There were 5 warnings)
```

```
bibtex_empty_fields
```

```
entries in general should not be empty in bibtex
```

```
find ""
```

```
passed: True
```

```
ascii
```

```
=====
The following tests are optional
=====
```

```
Tip: newlines can often be replaced just by an empty line
```

```
find newline
```

```
passed: True
```

```
cites should have a space before \cite{} but not before the {
```

```
find cite {
```

```
passed: True
```

Income Prediction based on Machine Learning Techniques

Borga Edionse Usifo

Indiana University

Bloomington, Indiana 47408

busifo@iu.edu

ABSTRACT

This project takes a closer look to some of the most used supervised learning algorithms in machine learning. We start with the description of the each of the algorithms then we move it to analytics and findings by using that particular algorithm in our data-set. We also provide advantages and disadvantages of each supervised machine learning algorithm for future reference. We mainly focus on our prediction of the income level of individuals by looking at their age, gender, education, location, and other features given by our data-set. We will try each algorithm and try to pick the best features from our data-set to have an optimal prediction.

KEYWORDS

i523, HID343, Machine Learning, Income Prediction, Logistic Regression, Ensemble methods

1 INTRODUCTION

In this project, we try to showcase the performance of the machine learning algorithms on data which we gather from UCI machine learning repository [22]. This data used by Kohavi R. and Becker B. for their research in improving the in Naive Bayes Classifier's accuracy [21].

Data consists of 15 variables, and we try to predict the income of the individuals. To do this prediction task, we first started with data preparation because the data we receive from UCI machine learning repository [22] not fully prepared for any machine learning algorithm. Our first task was the clean the data while applying some statistical techniques to get insights from the dataset. We also used data transformation methods like One-Hot-Encoding[45] to apply logarithmic functions for improving the machine learning algorithms performance before training the data.

Machine Learning algorithms that we discuss in this paper are Gaussian Naive Bayes [46], K Nearest Neighbors [29], Ensemble Methods (Boosting) [8], Support Vector Machines [6], Logistic Regression [34], and Decision Trees [49]. We try to show their weakness, advantages, and their time consumption while training each of them in machine learning algorithms section.

After providing a brief introduction of each of the supervised machine learning algorithms, we will discuss our findings for of each of the algorithms by comparing their accuracy score, F-1 score, recall, and lastly time comparison.

2 IMPORTANCE OF BIG DATA ANALYTICS FOR PREDICTIVE CLASSIFICATION

Importance of big data analytics is getting higher every day since the algorithms become more powerful to predict, classify and cluster any given data set. Importance of our case is any company can be used to predict individuals income to refer them goods in their

income range or governments can provide additional support for the areas that have lower income range. There can be many possible things that can do with this kind of classification predictions.

3 DATA PREPARATION

We first used the pandas [28] to help to load the data in data frame format. This gave us a unique advantage, and faster processing of comma separated values for putting into data frame [48]. Our data consist of 15 variables. Some of these variables are continuous, and some of them are categorical variables, and our target variable was "income" attribute. After putting the data into data frames, we first got a statistical snapshot of continuous variables (age, education, capital gain, capital loss, hours worked) by using the pandas [27] functions as shown in Table 1.

[Table 1 about here.]

3.1 Data Cleaning

After getting a snapshot from income data frame, we recognized that there is a column which has no meaning. The first task was to remove this entire column from our dataset we used pandas drop function for doing this task. After removing this column, we had more concise dataframe to analyze.

Moreover, removing the column we have encountered some missing values which labeled as "question marks" in data frame. In order to remove this values we first changed all the "question mark" values to "NaN" values by using pandas "replace" function [26]. After replacing all the question marks with "NaN" values, we used pandas missing value dropping function to remove all the "NaN" values from our dataset.

Furthermore, we start investigating the types of the variables, and in our case, we found two types of variable one of them labeled as "int64" which stands for integer values, other one labeled as object type of variable. From our previous example especially in "scikit-learn" it is better to use float object rather than "int64" for training the machine learning algorithms. Because their numerical output most of the time is "float64" object. We transferred all the "int64" objects to "float64" objects. This was the last step of the cleaning process.

Our last process is changing the string values to numerical values on our target data which consist of string values ("\$ 50K") for machine learning algorithms to understand this target data we need to transfer it to numerical values. Since we have only two categories, we will assign 1 and 0 as numerical values.

Description	Assigned Value
Individuals who makes more than \$ 50K	1
Individuals who makes at or less than \$ 50K	0

Our shape of the data will also receive impact from changing to numerical. Our number of futures will go from 14 to 103. This is

because we implemented one-hot-encode to our dataset. It is called one hot encoded because we transform the categorical variables into a more acceptable shape for the machine learning algorithms to perform well [45]. In other words “we implement binarization of the category to include as a future to train model [45]”. As we can see in Figure 1 and Figure 2.

[Figure 1 about here.]

[Figure 2 about here.]

4 DATA EXPLORATION

After cleaning the data, we started our data exploration to learn little bit more from our data and make necessary changes if needed before putting into our machine learning algorithms. The first step in this process is getting the total count of the individuals as well as the count of the individuals who are making more than \$50K and less than \$50K which can be seen in below Table 4.

Description	Count
Total Number of Individuals	30162
Individuals who makes more than \$50K	7508
Individuals who makes at or less than \$50K	22654

Moreover, we also look at the statistical values of each of the continuous variable we have. Those values given in Table 2. As we can see we have individuals who’re age ranging from 17 to 90 years old with a mean of 38.58. If we look at the capital gains and capital losses, we have a standard deviation of 7385 and 402 respectively this is also another indication of skew in these variables.

[Table 2 about here.]

We used scatter matrix plot and applied the correlation function to see if we have any reliable correlation between any of the variables. As we can see from the correlation matrix Figure 3 and correlation numbers Figure 4 we do not have the high correlation between any variables. Correlation values range between -1 to 1. The correlation value of 1 is an indication of perfect positive correlation and correlation number -1 indicates a negative correlation between variables [15]. Because of lower correlation values, it will be tough to determine the classification by just looking at the correlations; this indicates we have sophisticated algorithms to determine the relationship between variables to classify individuals incomes.

[Figure 3 about here.]

[Figure 4 about here.]

Furthermore, we also explore the capital gains, capital losses, and hours per week variables which we used a histogram to plot the data into distribution form so we can see how all these attributes distributed. The reason we do the histogram is we want to see any skewness in our data. As shown in the histogram graphs in Figure 5 and Figure 6 in capital gains and capital loss we have highly skewed data which can cause issues later on in our algorithms. We apply a logarithmic function to do highly skewed data to less skewed [24]. Using logarithmic functions adds more value to data from the interpretable standpoint and “it helps to meet the assumptions of inferential statistics [24]”.

[Figure 5 about here.]

[Figure 6 about here.]

Moreover, applying logarithmic function had an impact on distribution. We can see the changes on skew data in Figure 7 after applying logarithmic function.

[Figure 7 about here.]

5 MACHINE LEARNING ALGORITHMS TO CONSIDER

We have multiple algorithms to consider when we are doing the supervised learning. Each algorithm has its benefits and drawbacks. We will consider several supervised machine learning algorithms for our predictions. The application we will use to implement these algorithms will be Python Scikit-Learn library. We will briefly explain each parameter included in these algorithms in Scikit-Learn.

First we’ll look at the Scikit-Learn in Python framework we will go through the advantages in Scikit-Learn how we can implement any machine learning in just couple of simple line of codes in Scikit-Learn.

5.1 Why Scikit-Learn?

Scikit-learn developed by David Cournapeau in 2007. The development came from while he was working on summer code project for Google. After recognized and published by INRIA in 2010 project start the get more attention among worldwide. There are more than 30 active contributors and has secured several sponsorships from big technology companies[17]. “It also has a goal of providing common algorithms to Python users through consistent interface[2]”. Scikit-Learn consists of several elements to make analytical predictions. These elements are shown below[23]:

Supervised Learning Algorithms: One of the most fundamental reason that Scikit-Learn’s popularity comes from highly available supervised learning algorithms. These algorithms vary from regression models to decision trees and many more[23].

Cross Validation: Scikit-Learn includes various techniques to check the accuracy or any statistical measure between training and unseen testing set[23].

Unsupervised Learning Algorithms: Scikit-Learn had also various algorithms to support many unsupervised algorithms some of these include clustering, factor analysis, and neural network analysis[23].

Various example data-sets: Scikit-Learn comes with different data sets included in its package so users can start learning Scikit-Learn without the need of any data-sets[23].

Feature extraction: It has rich feature for extracting images or text from data-sets[23].

Algorithms that we will investigate shown below; we will go more deep analysis on each of these algorithms.

- Gaussian Naive Bayes
- Logistic Regression
- K-Nearest Neighbors (KNN)
- Stochastic Gradient Descent Classifier
- Support Vector Machines
- Decision Trees

5.2 Gaussian Naive Bayes

Naive Bayes bring many beneficial features; it is widely popular among machine learning applications[41]. The popularity of Naive Bayes comes from being able to handle large projects and data-sets faster than most algorithms[41]. It also can handle complex data-sets with categorical and non-categorical inputs [41]. Naive Bayes based on probabilistic classifier of Bayesian theory. It is also a favorite way of doing text categorization [46].

Term naive comes from it is the method of use probability among categories which assumes of independence among given class of attributes as shown in Figure 8. In other words, if we try to classify individuals from their email communications it will not take the order of words into account. Whereas in the English language we can tell the difference between sentence makes sense or not if we randomly re-order our words in the sentences. So it does not understand the text, it only looks at word frequencies as a way to do the classification. This is why it is called “Naive”.

[Figure 8 about here.]

As we state above Naive Bayes derives from Bayesian Theory where the dimensionality of inputs is relatively high. Bayesian Theorem is stated below [16].

$$P(C | X) = \frac{P(X | C) \times P(C)}{P(X)} \quad (1)$$

Naive Bayes Classifier works as follows [16]:

Advantages of Naive Bayes [16]:

- Faster classification time for training data-set.
- Because of independent classification it improves classification performance.
- Performance is relatively good.

Disadvantages of Naive Bayes[16]:

- Often it requires a large number of data-sets to give adequate results.
- On some occasions which are relative to data-sets, it can give less accuracy.

5.3 Logistic Regression

Logistic Regression widely used for predicting “probability of failure in a given system, product, and process [34]“. Logistic Regression also used in natural language analysis, it is an extension of conditional random fields [34]. It works as a classifier which learns the features from the input given and classifies them by multiplying the input value with the weight value [14].

$$P(C | X) = \sum_{i=1}^N W_i \times f_i \quad (2)$$

Main reason that Logistic Regression differs from Linear Regression is output variable for Logistic Regression is binary whereas output variable in Linear Regression is discrete(continuous) [12].

Advantages of Logistic Regression:

- It does not have any assumptions over distribution of classes [18].
- It is fast to train [18].
- Logistic Regression has fast classifying method of unknown data [18].

- We can easily extend to other regression for multiple classes like multinomial regression [18].

Disadvantages of Logistic Regression:

- One of the disadvantages of linear regression is it is not providing flexibility in some instances. What we mean by the “ lack of flexibility is the linear dependency, and linear decision boundary in the instance space is not valid [42]“. This disadvantage can be improved changing from Logistic Regression to Choquistic Regression[42].
- Logistic regression can provide poor results when there are more complex relationships in data [9].
- Logistic models also have over-fitting problems which come from a result of sampling bias [31].
- Because of Logistic Regression’s predictions comes from the independent variable if the researcher includes wrong independent variables then model’s prediction will have no value [31].
- Because it is predictions based on 1 and 0 model will have poor performance when predicting continuous variables [31].

5.4 K-Nearest Neighbors (KNN)

K Nearest neighbor has been primarily studied, and this popularity comes from it has been applied to many applications some of these applications are “spatial databases, pattern recognition, geographic information, image retrieval, computer game, and many other applications [29]“. Due to an increase of mobile devices and people tends to use of applications like navigation K-nearest neighbor found itself another widely used area of location-based services due to an ability to found a target location [29].

Intuition behind the K Nearest Neighbor can be described as follows: “ for a set P of n objects and a querying point q, return the k objects in P that are closest to q [29].“

Advantages of K Nearest Neighbors:

- K Nearest Neighbor is a basic and simple approach to implement [35].
- K Nearest Neighbor can perform well and efficiently with the large amount of data [43].
- K nearest Neighbor also does effectively well with noisy data sets (“if the inverse square of weighted distance used as the distance [43]“). In other words, it is flexible to feature and distance choices [35].

Disadvantages of K Nearest Neighbors:

- K Nearest Neighbor typically require large dataset to perform well [35].
- Time complexity could be high due to computing distance of each query to all training data points [43]. This time might be improved with some indexing (K-D Tree) [43].
- Determining the value of K can be time-consuming [43].
- It can be unclear to know which type of distance to use, as well as which variability to use to get the optimal results [43].
- Switching the different K values can result in the predicted class labels [30].

Many of these disadvantages are improving with the help of parallel distributed computing. Recent improvements in MapReduce framework allows users to run KNN algorithms in the cluster which had a significant effect on reducing the computation time [19].

Another area of improvements on KNN, is to implement different mapping functions such as kernel KNN, kernel difference weighted KNN, adaptive quasi-conformal kernel nearest neighbor, angular similarity, local linear discriminant analysis, and Dempster-Shafer [10].

5.5 Decision Trees

Decision Tree is another widely used algorithm model for classification and regression. Decision Trees uses a recursive split model where each recursive split is identified by each data point; this is an example of non-parametric hierarchical model [13].

Representation of decision trees is as follows; we sort the instances from root to leaf nodes, this sorting gives insights about the classification of the instance, every outcome descending from the root node corresponds to possible values for that variable [33]. We can classify an instance by starting from the root node and checking the attributes labeled on that node and moving down from that node based on attribute given attribute values [33] as shown in Figure 9.

[Figure 9 about here.]

Advantages of Decision Trees:

- Decision Tree applications are easy to interpret and understand [32]. This ease comes from their schematic representation [32]. Interpretation between alternatives can be expressed with single numerical number which is the expected value (EV) [32].
- Decision Trees can handle noisy or incomplete data-sets [32]. In other words it requires little effort of data preparation because of it is flexibility [7].
- It can handle both nominal and numerical variables [32].
- It can be modified easily whenever the new information is available [32].
-

Disadvantages of Decision Trees:

- Because of it is a use of divide and conquer method they can demonstrate good performance if there are few attributes exists when the attributes level goes into large number decision tree become more complex which will result in poor performance [32].
- Decision Trees are also susceptible to training set which can give a result of over-fitting [32]. In other words, it can believe the training set completely which will give an abysmal performance on testing set.
- ID3 and C4.5 decision tree algorithms require discrete values as input data.

5.6 Stochastic Gradient Descent Classifier (SGD)

Stochastic Gradient Descent recently got became more popular because of it is large-scale learning ability in machine learning

problems [11]. It is a useful and straightforward way approach of linear classifiers under convex problems which is Support Vector Machines or Conditional Random Fields [3]. The originality of SGD derives from “Stochastic Approximation” which is a work from Robinson and Manroe [5].

Advantages of Stochastic Gradient Descent:

- One of the advantage of stochastic gradient descent is, it is easy to implement [38].
- Stochastic Gradient Descent is also efficient because of each step only relies on a single derivative which makes the computational cost $1 / n$ than normal gradient descent [37].

Disadvantages of Stochastic Gradient Descent:

- Stochastic Gradient Descent can be required to have many iterations, and it also requires some hyper-parameters [38].
- Feature scaling is a practice which used in the standardization of range of independent variables [47]. SGD also used this feature scaling technique and it can be sensitive to feature scaling [38].
- Another drawback of Stochastic Gradient Descent is while using GPU they are hard to parallelize or distributing them using computer clusters [25].

5.7 Support Vector Machines

Support Vector Machines is fallen under the classification methods in machine learning [6]. It is also a robust classification method that has been widely found itself an area ranging from pattern recognition to text analysis [6].

Fitting a boundary between data points is the principle of the support vector machines. This boundary divides the data points between classes, and each similar data point puts under the same class classification [6]. After training the support vector machines with training data-set, we only need to check whether the test data lies under the boundaries for testing set. Another thing to consider is after it creates the boundaries of the data remaining training data becomes obsolete because we only need the core set of points which supports the boundaries to classify the new data set. This core data points called “support vectors”. It is called vector because of each data point contains a row of observed data values for attributes [6].

[Figure 10 about here.]

Traditionally boundaries are called “hyperplanes” and it is used to describe boundaries in more than three dimensions because they are hard or sometimes impossible to visualize.[7]. Figure 10. Optimality of hyperplane expressed as a linear function which requires maximum distance between the identified classes. It only considers a small number of training example to build this hyperplane. SVM hyperplanes based on “separation of positive (+1) and negative (-1) with the largest margin [39]“.

One of the main characteristic of the machine learning is to generalization. In other words, we want to give a general idea that tends to fit any of our testing datasets optimally. Support vector machines are a perfect regarding generalizations because once the training data fitted by the support vector machines other than support vector data inside the training data becomes redundant

which means that even with the small changes inside the data will not have a significant effect on general boundaries [6].

Advantages of Support Vector Machines:

- Generalizes the data well with the help of boundaries. Which reduces the overfitting [6].
- Classification accuracy in basic support vector machine will yield a 95 percent accuracy with a default settings [6].
- SVM can deliver a unique solution, because of optimality solution is convex. This will give an advantage over Neural Networks which has multiple solutions in local minima [1].

Disadvantages of Support Vector Machines:

- One common disadvantage of SVM, is the lack of transparency because of its non-parametric techniques [1].
- Another biggest disadvantage of SVM is it requires high algorithmic complexity and high level of memory for the large-scale implementations [39].
- According to Burges, biggest limitation of the SVM is in the choice of kernel [4].

5.8 Ensemble Methods

Ensemble methods goes into classification algorithm category, they are learning algorithms which uses weighted vote for its prediction methods, in other words, it is learning rules over a small subset of data then we combine these rules which we learn from the small subset of data to make predictions and/or classification on the testing data [8]. The originality of the Ensemble method comes from Bayesian averaging, but with the recent algorithms include “Bagging, error-correcting, and boosting [8]”.

Bagging refers to simply the looking at data-sets and dividing the data-set to its small subsets then learning the rules of that particular small subset. Next step is combining each learned rule from subsets to apply to more significant data set. Combining method mostly done with averaging the learned rules. Bagging also does better on testing set than standard Linear Regression analysis and linear regression does better on training set especially in third order polynomial [8].

Stacking

Boosting is another method used in Ensemble Methods. The difference from bagging is in boosting we need to pick subsets or examples that we are not good at in other words hardest examples. Then we combine these learned rules with the weighted mean instead of mean used in bagging method.

Boosting is little different than bagging.

Advantages of Ensemble Methods:

- Prediction of the ensemble methods is better than most of the algorithms because of the combining methods intuition makes the model less noisy [36].
- They are more stable than other algorithms. [36]

Disadvantages of Ensemble Methods:

- Over-fitting may cause some disadvantages for ensemble learning but bagging operation will reduce this overfitting [36].

6 FITTING DATA INTO MACHINE LEARNING ALGORITHMS

In this section, we will show the techniques we used on the execution of the prepared data into machine learning algorithms. Before fitting the data into the machine learning algorithms, we split the data into two sets. These sets are the training set and the testing set. We do splitting because of gaining an access of the future data will most likely be hard before future occurs, and because of this fact, it is a good idea to test our model with a dataset which our model has not seen it [40].

We used scikit-learn for splitting data into train and test we saved 20% of data for testing purposes as shown in Table 3 .

[Table 3 about here.]

Furthermore, after splitting the data we put all of our training data into each of the machine learning algorithm to get their prediction results. We also provided code at the beginning and the end of each algorithm to calculate their running time.

Before we move further we need to discuss critical characteristics of a machine learning algorithm. These are;

- Confusion Matrix
- Accuracy
- Recall
- F-1 Score
- Precision

6.0.1 Confusion Matrix: Confusion matrix develops from 4 key elements. These elements are true positive, true negative, false negative, and false positive. As shown in Figure 11 about the constructing a confusion matrix. If we want to build a confusion matrix by targeting individuals who are making more than \$50K our true positive, true negative, false positive, and false negative explained below.

[Figure 11 about here.]

True Positive (TP): We can explain true positive as if the individuals make more than \$50K and our model correctly classifies them as individuals who makes more than \$50K, then this individual is in higher income range, in this case, we call it a true positive [20].

True Negative (TN): Intuition of true negative is if an individual makes less than \$50K and our model correctly classifies them as individuals who makes less than \$50K, then this individual is in lower income range. We call this true negative [20].

False Negative (FN): When an individual makes less than \$50K and our model incorrectly classifies them in higher income range by making a mistake causes a false negative to happen [20].

False Positive (FP): When an individual is making more than \$50K and our model classifies them in lower income range by mistake. This is called false positive [20].

6.0.2 Accuracy: Accuracy answers the question of how good is the model is. In our case this question will be out of all the individuals, how many did the models classify the individuals correctly. The mathematical expression of the accuracy is the ratio between the number of correctly classified points and the number of total points. We can think that if we have high accuracy, our model is

excellent, but this is only where we have identical false positive and false negative values in our dataset [20].

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + FN + TN} \quad (3)$$

6.0.3 Precision. Precision answers the questions of out of all the points predicted to be positive how many of them were actually positive? If we translate this question into our case, we will have out of all the individuals that we are classified as lower income how many were actually have lower income. Higher precision indicates that we have low false positive rate [20]. Mathematical expression of precision is;

$$\text{Precision} = \frac{TP}{TP + FP} \quad (4)$$

6.0.4 Recall (Sensitivity). Recall answers the question of “out of the points that are labeled positive how many of them were correctly predicted is positive ? ”. If we translate this to our case, we will have “out of the points that are labeled higher income how many of them correctly predicted is in higher income range ? ”. Mathematical expression of the recall is;

$$\text{Precision} = \frac{TP}{TP + FN} \quad (5)$$

6.0.5 F-1 Score. The F-1 score is the idea of giving a decision by looking at only one score which will include precision, and recall scores. We cannot just take the average of precision and recall because if either of them is very low. We need a number to be low, even if the other one is not. This will leads us to look at the harmonic mean, and it works as follow. Let's say we have two numbers X and Y. X is smaller than Y, and we have the arithmetic mean, and it always lies between X and Y. It is a mathematical fact that the harmonic mean is always less than the arithmetic mean which is closer to the smaller number than to the higher number. Mathematical expression of F-1 score is;

$$F1Score = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (6)$$

6.1 Results

Now we can look at the results from each of the machine learning algorithm. Results also showed in Table 4 with the visualization of Figure 13. We can also see the running time of the each of the algorithm in Figure 12. Support Vector Machines is the winner for the highest running time for training the algorithm.

[Figure 12 about here.]

[Table 4 about here.]

6.1.1 Naive Bayes. As shown in the Figure 13 we have a comparison of several supervised machine learning algorithms on our dataset. We can see that from the accuracy standpoint Naive Bayes algorithms have the lowest score which means that it did not do a good job for labeling true positives regards to all data but it did a good job in precision standpoint while doing a bad classification from recall standpoint. Two key element for us in this situation is accuracy and f1 score(which consist of precision and recall).

6.1.2 Support Vector Machine. Support Vector Machine is the second best algorithm in our case. This algorithm did very well job on classification it has the second highest accuracy and f1 score.

6.1.3 AdaBoost. As we stated before ensemble algorithms learn from the small portion of the data and combine these learning to do the predictive task. As shown in Figure 13 adaboosting has the highest accuracy score among all the other algorithms. This algorithm should be our first choice to do predictive modeling. We believe that there is still an improvements on accuracy

6.1.4 K-Nearest Neighbors. K-Nearest Neighbor algorithm in our project we set the k value to 5. K Nearest Neighbor algorithm also did a good job by placing itself third in accuracy score.

6.1.5 Decision Tree. Decision Tree is gave a good accuracy but fall behind on f1 score as shown in Figure 13.

[Figure 13 about here.]

7 CONCLUSION

We presented the importance of analytical approach with machine learning algorithms and how they can be used to predict or classify the individuals with many different attributes like age, education, income, etc. We also presented weaknesses and strengths of these algorithms along with their precision, accuracy, recall, and F-1 scores by presenting with the visualizations. We also demonstrated the running time for each algorithm while using big data sets. The source code of this project can found Github website which presented in reference section [44].

ACKNOWLEDGMENTS

The author would like to thank Dr. Gregor von Laszewski for his support and suggestions to write this paper.

REFERENCES

- [1] L. Auria and A. R. Moro. 2008. Support Vector Machines (SVM) as a Technique for Solvency Analysis. Online. http://www.diw.de/english/publications/discussion_papers/27539.html
- [2] L. Ben. 2015. Six Reasons why I recommend scikit-learn. Online. (Oct. 2015). <https://www.oreilly.com/ideas/six-reasons-why-i-recommend-scikit-learn>
- [3] L. Bottou. 2010. Stochastic Gradient Descent. Online. (2010). <http://leon.bottou.org/projects/sgd>
- [4] C. J. C. Burges. 1998. A Tutorial on Support Vector Machines for Pattern Recognition. *Data Mining and Knowledge Discovery* 2, 2 (01 Jun 1998), 121–167. <https://doi.org/10.1023/A:1009715923555>
- [5] N. Deanna, S. Nathan and W. Rachel. 2016. Stochastic gradient descent, weighted sampling, and the randomized Kaczmarz algorithm. *Mathematical Programming* 155, 1 (01 Jan 2016), 549–573. <https://doi.org/10.1007/s10107-015-0864-7>
- [6] B. Deshpande. 2013. When do support vector machines trump other classification methods. Online. (Jan. 2013). <http://www.simafore.com/blog/bid/112816/When-do-support-vector-machines-trump-other-classification-methods>
- [7] B. Deshpande. 2011. 4 key advantages of using decision trees for predictive analytics. Online. (July 2011). <http://www.simafore.com/blog/bid/62333/4-key-advantages-of-using-decision-trees-for-predictive-analytics>
- [8] G. T. Dietterich. n.d.. Ensemble Methods in Machine Learning. (n.d.). <http://web.engr.oregonstate.edu/~tgd/publications/mcs-ensembles.pdf>
- [9] EliteDataScience. 2016. Modern Machine Learning Algorithms: Strengths and Weaknesses. Online. (May 2016). <https://elitedatascience.com/machine-learning-algorithms>
- [10] O. F. Ertugrul and M. E. Tagluk. 2017. A novel version of k nearest neighbor: Dependent nearest neighbor. *Applied Soft Computing* 55, Supplement C (2017), 480 – 490. <https://doi.org/10.1016/j.asoc.2017.02.020>
- [11] M. Fan. n.d.. How and Why to Use Stochastic Gradient Descent? (n.d.). <http://anson.ucdavis.edu/~minjay/SGD.pdf>

- [12] J. Fang. 2013. Why Logistic Regression Analyses Are More Reliable Than Multiple Regression Analyses. *Journal of Business and Economics* 4, 7 (July 2013), 620–633. <http://www.academicstar.us/UploadFile/Picture/2014-6/201461494819669.pdf>
- [13] M. A. Hassan, A. Khalil, S. Kaseb, and M. A. Kasseem. 2017. Potential of four different machine-learning algorithms in modeling daily global solar radiation. *Renewable Energy* 111, Supplement C (2017), 52 – 62. <https://doi.org/10.1016/j.renene.2017.03.083>
- [14] S. T. Indra, L. Wikarsa, and R. Turang. 2016. Using logistic regression method to classify tweets into the selected topics. *2016 International Conference on Advanced Computer Science and Information Systems (ICACSI), Advanced Computer Science and Information Systems (ICACSI), 2016 International Conference on* 1, 385–389 (2016), 385. <http://proxyiub.uits.iu.edu/login?url=https://search-ebscohost.com.proxyiub.uits.iu.edu/login.aspx?direct=true&db=edsee&AN=edsee.7872727&site=eds-live&scope=site>
- [15] Investopedia. n.d. Correlation Coefficient. Online. (n.d.). <https://www.investopedia.com/terms/c/correlationcoefficient.asp>
- [16] D. S. Jadhav and H. P. Channe. 2014. Comparative Study of K-NN, Naive Bayes and Decision Tree Classification Techniques. *International Journal of Science and Research (IJSR)* 5, 1 (Jan. 2014), 1842–1845. <https://www.ijsr.net/archive/v5i1/NOV153131.pdf>
- [17] B. Jason. 2014. A gentle introduction to Scikit-Learn: Python Machine Learning Library. Online. (April 2014). <https://machinelearningmastery.com/a-gentle-introduction-to-scikit-learn-a-python-machine-learning-library/>
- [18] H. Jeff. 2012. Introduction to Machine Learning. Online. (Jan. 2012). http://courses.washington.edu/css490/2012/Winter/lecture_slides/05b_logistic_regression.pdf
- [19] J. Jiaqi and Y. Chung. 2017. Research on K nearest neighbor join for big data. In *2017 IEEE International Conference on Information and Automation (ICIA)*. IEEE, Department of Computer Engineering Wonkwang University Iksan 54538, Korean, 1077–1081. <https://doi.org/10.1109/ICInfa.2017.8079062>
- [20] R. Joshi. 2016. Accuracy, Precision, Recall, and F1 Score: Interpretation of Performance Measures. Online. (Sept. 2016). <http://blog.exsilio.com/all/accuracy-precision-recall-f1-score-interpretation-of-performance-measures>
- [21] R. Kohavi. 1996. Improving the Accuracy of Naive-Bayes Classifiers: A Decision-tree Hybrid. In *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (KDD'96)*. AAAI Press, Silicon Graphics, Inc, 202–207. <http://dl.acm.org/citation.cfm?id=3001460.3001502>
- [22] R. Kohavi and B. Becker. n.d. Predicting whether income exceeds \$50K/yr based on census data. Online. (n.d.). <https://archive.ics.uci.edu/ml/datasets/Census+Income>
- [23] J. Kunal. 2015. Scikit-Learn in python - The most important Machine Learnig Tool I learnt last year. Online. (Jan. 2015). <https://www.analyticsvidhya.com/blog/2015/01/scikit-learn-python-machine-learning-tool/>
- [24] M. D. Lane. n.d. Log Transformations. Online. (n.d.). <http://onlinestatbook.com/2/transformations/log.html>
- [25] V. Q. Le, J. Ngiam, A. Coates, A. Lahiri, B. Prochnow, and A. Y. Ng. 2011. On optimization methods for deep learning. In *International Conference of Machine Learning*. Stanford University, International Conference of Machine Learning, Stanford University, NA. <https://cs.stanford.edu/~acoates/papers/LeNgiCoaLahProNg11.pdf>
- [26] Pandas Library. n.d.. Dataframe replace. Online. (n.d.). <https://pandas.pydata.org/pandas-docs/stable/generated/pandas.DataFrame.replace.html>
- [27] Pandas Library. n.d.. Pandas Dateframe describe. Online. (n.d.). <https://pandas.pydata.org/pandas-docs/stable/generated/pandas.DataFrame.describe.html>
- [28] Pandas Py Data Library. n.d.. Pandas for Python. Online. (n.d.). <https://pandas.pydata.org/>
- [29] L. J. Moon. 2017. Fast k-Nearest Neighbor Searching in Static Objects. *Wireless Personal Communications* 93, 1 (01 Mar 2017), 147–160. <https://doi.org/10.1007/s11277-016-3524-1>
- [30] G. Nick. 2014. KNN. Online. (April 2014). <http://www.nickgillian.com/wiki/pmwiki.php/GRT/KNN>
- [31] R. Nick. NA. The Disadvantages of Logistic Regression. Online. (NA). <http://classroom.synonym.com/disadvantages-logistic-regression-8574447.html>
- [32] C. Petri. 2010. Decison Trees. Online. (2010). <http://www.cs.ubbcluj.ro/~gabis/DocDiplome/DT/DecisionTrees.pdf>
- [33] U. Princeton. NA. Decision Tree Learning. Online. (NA). <http://www.cs.princeton.edu/courses/archive/spr07/cos424/papers/mitchell-decrees.pdf>
- [34] S. A. Raj, L. J. Fernando, and S. Raj. 2017. Predictive Analytics On Political Data. Congress. *World Congress on Computing and Communication Technologies* 10, 1109 (2017), 93–96.
- [35] M. Ray. 2012. Nearest Neighbours: Pros and Cons. Online. (April 2012). <http://www2.cs.man.ac.uk/~raym8/comp37212/main/node264.html>
- [36] S. Ray. 2015. 5 Easy Questions on Ensemble Modeling Everyone Should Know. Online. (Jan. 2015). <https://www.analyticsvidhya.com/blog/2015/09/questions-ensemble-modeling/>
- [37] J. Rie and Z. Tong. 2013. Accelerating Stochastic Gradient Descent using Predictive Variance Reduction. In *Advances in Neural Information Processing Systems 26*, C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger (Eds.). Curran Associates, Inc., Rutgers University, New Jersey, USA, 315–323. <http://papers.nips.cc/paper/4937-accelerating-stochastic-gradient-descent-using-predictive-variance-reduction.pdf>
- [38] Scikitlearn. n.d.. Stochastic Gradient Descent. Online. (n.d.).
- [39] K. N. Shrivastava, P. Saurabh, and B. Verma. 2011. An Efficient Approach Parallel Support Vector Machine for Classification of Diabetes Dataset. *International Journal of Computer Applications in Technology* 36, 6 (Dec. 2011), 19–24. <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.259.3757&rep=rep1&type=pdf>
- [40] D. Steinberg. 2014. Why Data Scientist Split Data into Train and Test. Online. (March 2014). <https://info.salford-systems.com/blog/bid/337783/Why-Data-Scientists-Split-Data-into-Train-and-Test>
- [41] K. B. Tapan. 2015. Naive Bayes vs Logistic Regression: Theory, Implementation and Experimental Validation. *Inteligencia Artificial, Vol 18, Iss 56, Pp 14-30 (2015) 1, 56 (2015), 14.* <http://proxyiub.uits.iu.edu/login.aspx?direct=true&db=edsdobj&AN=edsdobj.0e372b34c5d48bcf72cd437eade1fd1&site=eds-live&scope=site>
- [42] A. F. Tehrani, W. Cheng, and E. Hullermeier. 2011. Choquistic Regression: Generalizing Logistic Regression Using the Choquet Integral. Online. (July 2011). <https://www-old.cs.uni-paderborn.de/fileadmin/Informatik/eim-i-is/PDFs/Talk.EUSFLAT.11.pdf>
- [43] K. Teknomo. 2017. K-Nearest Neighbor Tutorial. Online. (2017). <http://people.revoledu.com/kardi/tutorial/KNN/Strength%20and%20Weakness.htm>
- [44] E. B. Usifo. 2017. Income Prediction. Github. (Dec. 2017). <https://github.com/bigdata-i523/hid343/tree/master/project>
- [45] R. Vasudev. n.d.. What is One Hot Encoding? do you have to use it ? Online. (Aug. n.d.). <https://hackernoon.com/what-is-one-hot-encoding-why-and-when-do-you-have-to-use-it-e3c6186d008f>
- [46] Wikipedia. 2017. Naive Bayes. Online. (Nov. 2017). https://en.wikipedia.org/wiki/Naive_Bayes_classifier
- [47] Wikipedia. NA. Feature Scaling. Online. (NA). https://en.wikipedia.org/wiki/Feature_scaling
- [48] Wikipedia. n.d.. Comma Separated Values. Online. (n.d.). https://en.wikipedia.org/wiki/Comma-separated_values
- [49] Wikipedia. n.d.. Decision Trees. Online. (n.d.). https://en.wikipedia.org/wiki/Decision_tree
- [50] H. Zhang. 2004. *The Optimality of Naive Bayes*. resreport. University of New Brunswick. <http://www.cs.unb.ca/~hzhang/publications/FLAIRS04ZhangH.pdf>

LIST OF FIGURES

1	Example of One Hot Encoding Before [45].	9
2	Example of One Hot Encoding After [45].	10
3	Correlation Matrix [44].	11
4	Scatter Matrix Plot [44].	12
5	Histogram of Capital Gain [44].	13
6	Histogram of Capital Loss [44].	13
7	After Logarithmic Function Applied Histogram of Capital Gain [44].	14
8	Example of Naive Bayes [50].	15
9	Example of Decision Tree Construction[33].	15
10	Example of Shows the Hyperplanes [6].	16
11	Example of Confusion Matrix Construction [20].	16
12	Supervised Learning Algorithm Running Time Results [44].	17
13	Supervised Learning Algorithm Results [44].	18

CompanyName	Categoricalvalue	Price
VW	1	20000
Acura	2	10011
Honda	3	50000
Honda	3	10000

Figure 1: Example of One Hot Encoding Before [45].

Vw	Acura	Honda	Price
1	0	0	20000
0	1	0	10011
0	0	1	50000
0	0	1	10000

Figure 2: Example of One Hot Encoding After [45].

	age	education_num	capital-gain	capital-loss	\
age	1.000000	0.043526	0.080154	0.060165	
education_num	0.043526	1.000000	0.124416	0.079646	
capital-gain	0.080154	0.124416	1.000000	-0.032229	
capital-loss	0.060165	0.079646	-0.032229	1.000000	
hours-per-week	0.101599	0.152522	0.080432	0.052417	
hours-per-week					1.000000
age		0.101599			
education_num		0.152522			
capital-gain		0.080432			
capital-loss		0.052417			
hours-per-week		1.000000			

Figure 3: Correlation Matrix [44].

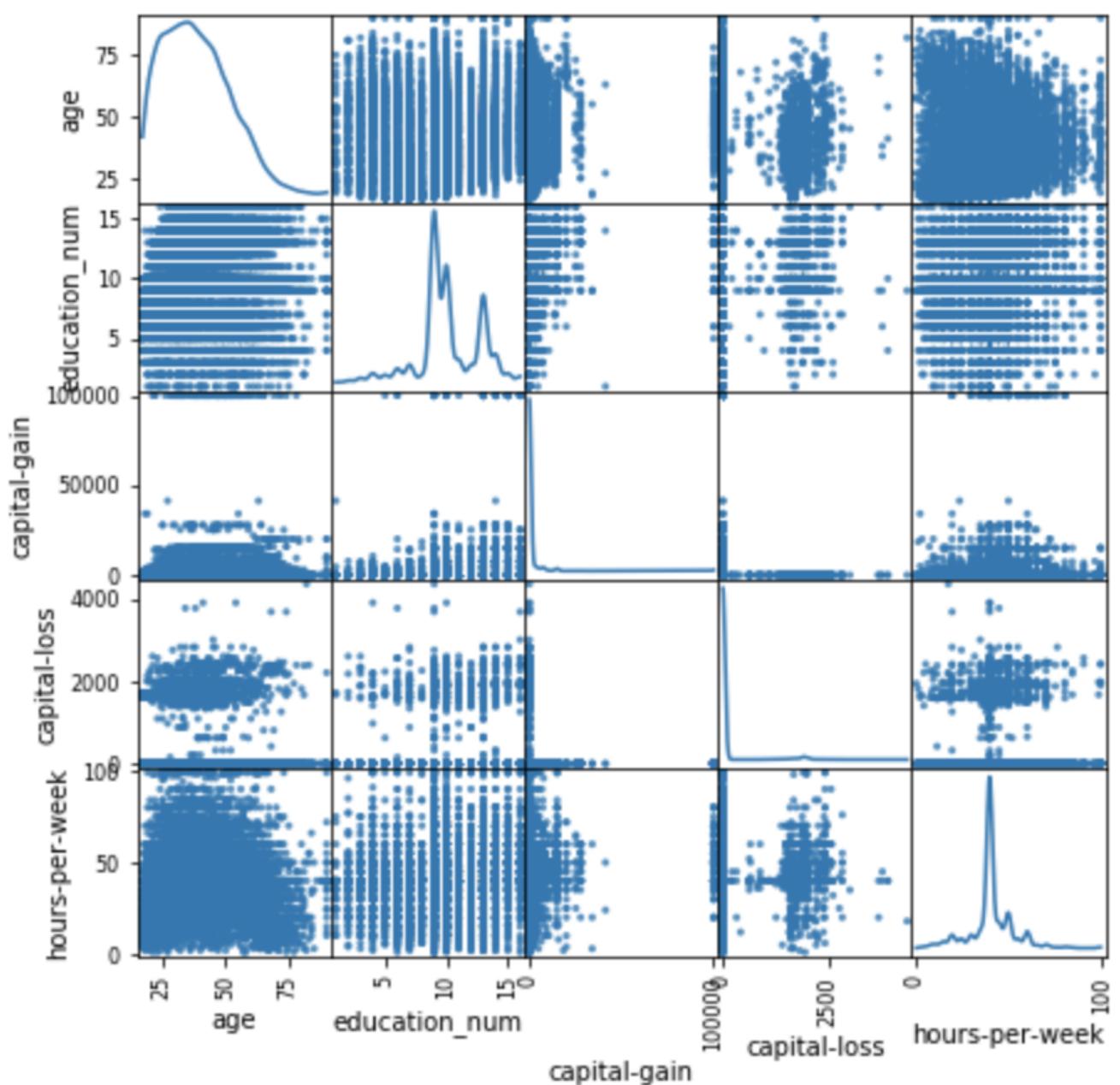


Figure 4: Scatter Matrix Plot [44].

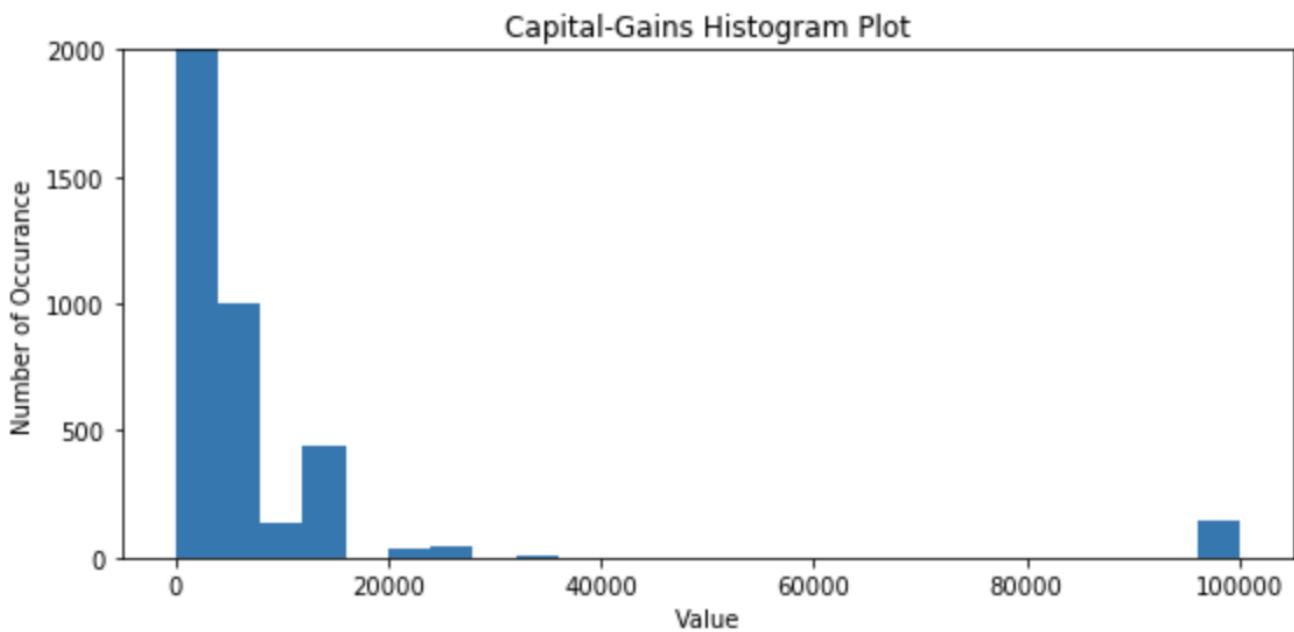


Figure 5: Histogram of Capital Gain [44].

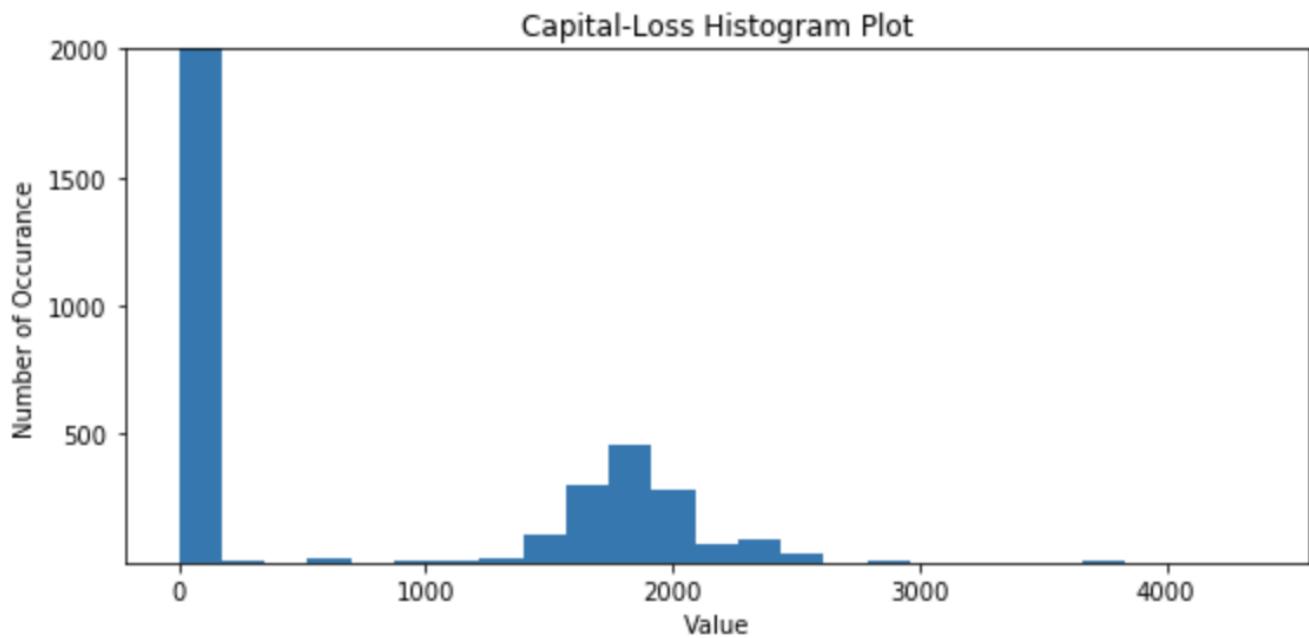


Figure 6: Histogram of Capital Loss [44].

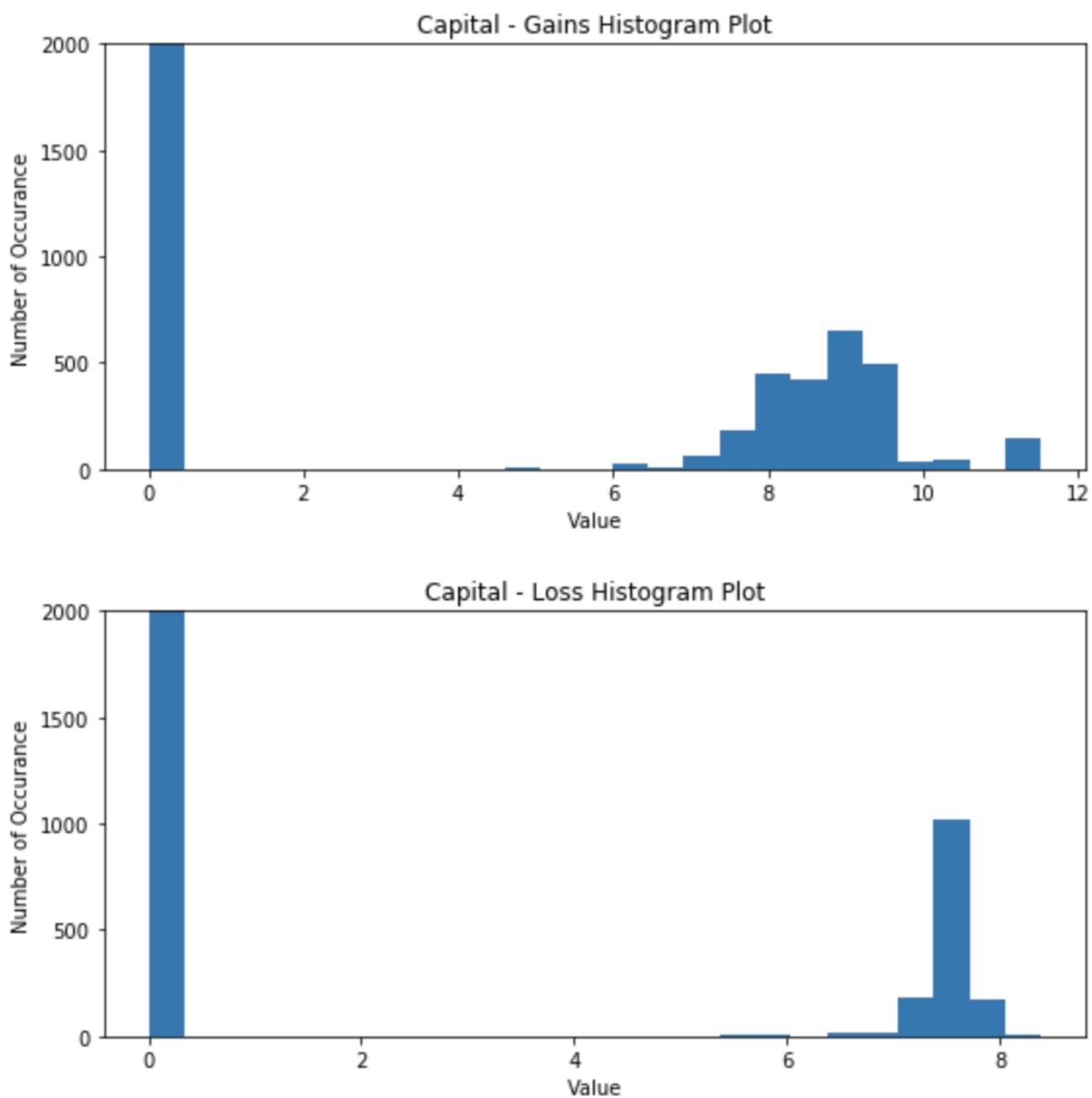


Figure 7: After Logarithmic Function Applied Histogram of Capital Gain [44].

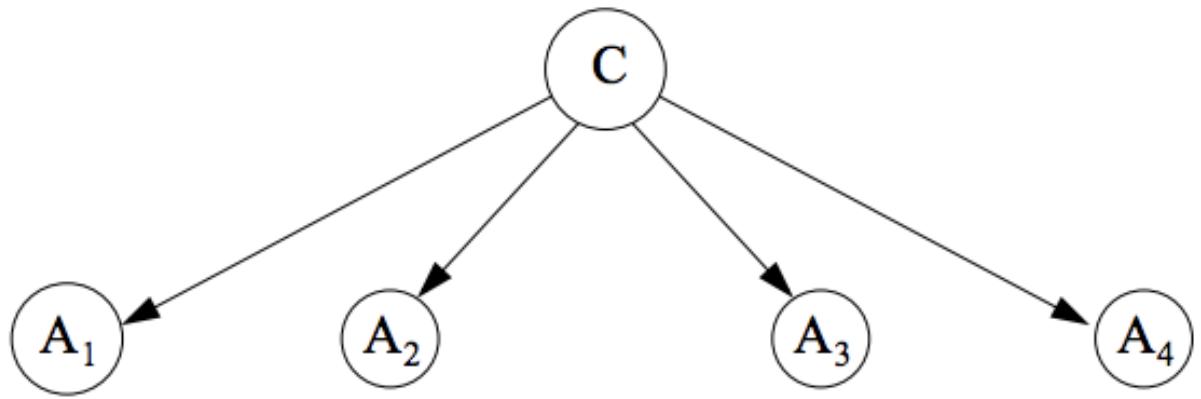


Figure 8: Example of Naive Bayes [50].

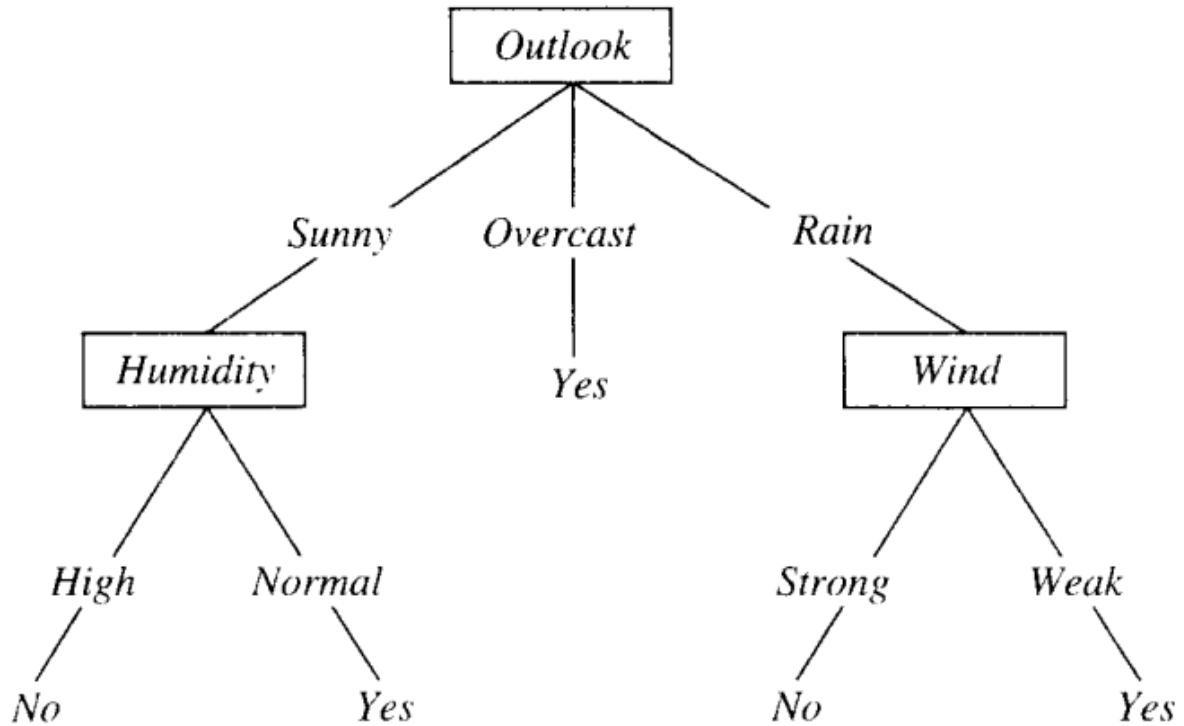


Figure 9: Example of Decision Tree Construction[33].

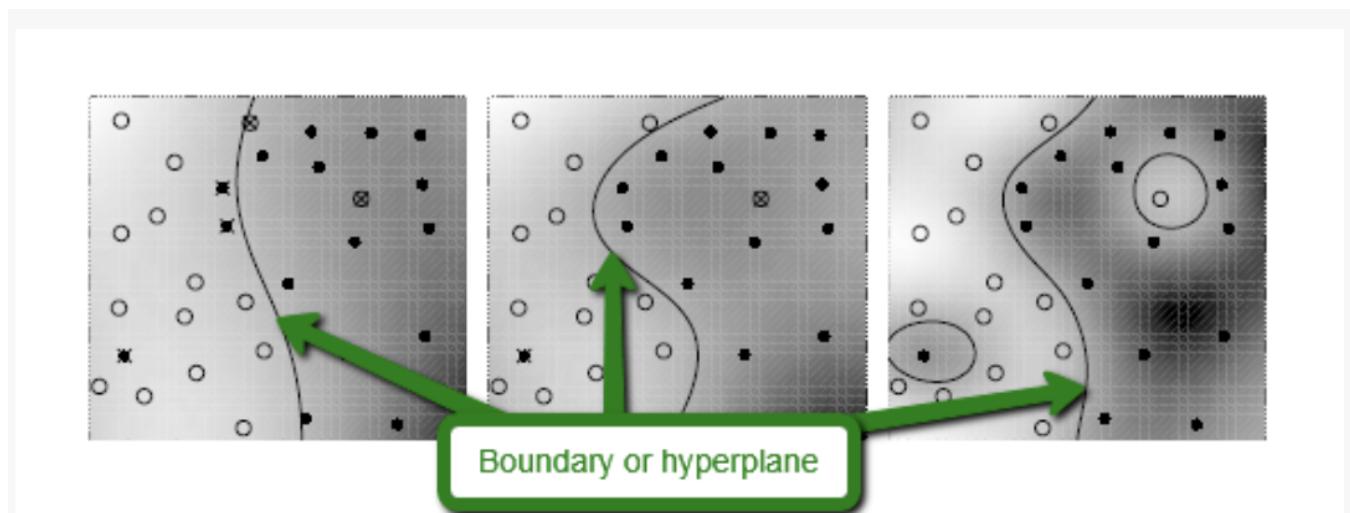


Figure 10: Example of Shows the Hyperplanes [6].

		Predicted class	
		Class = Yes	Class = No
Actual Class	Class = Yes	True Positive	False Negative
	Class = No	False Positive	True Negative

Figure 11: Example of Confusion Matrix Construction [20].

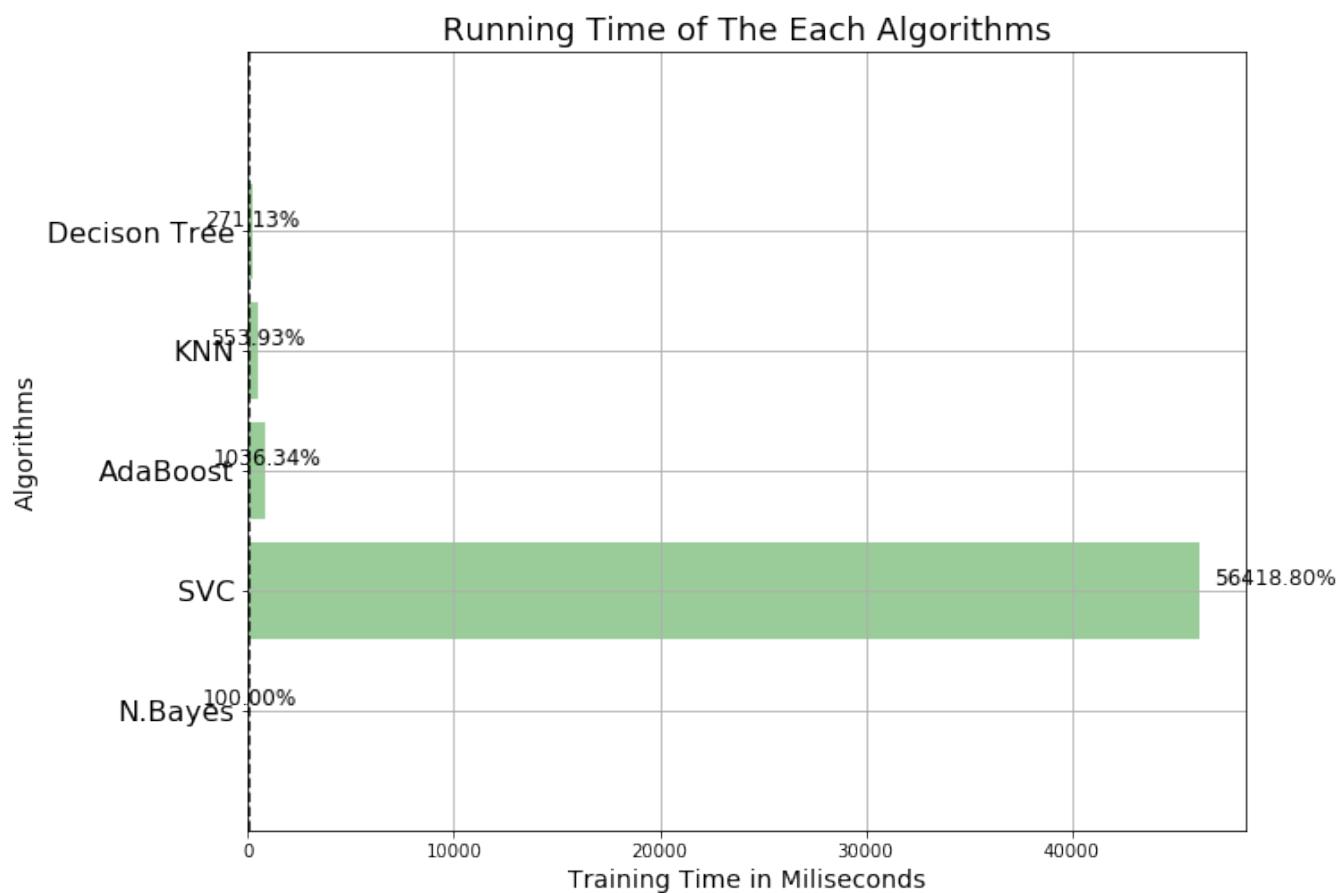


Figure 12: Supervised Learning Algorithm Running Time Results [44].

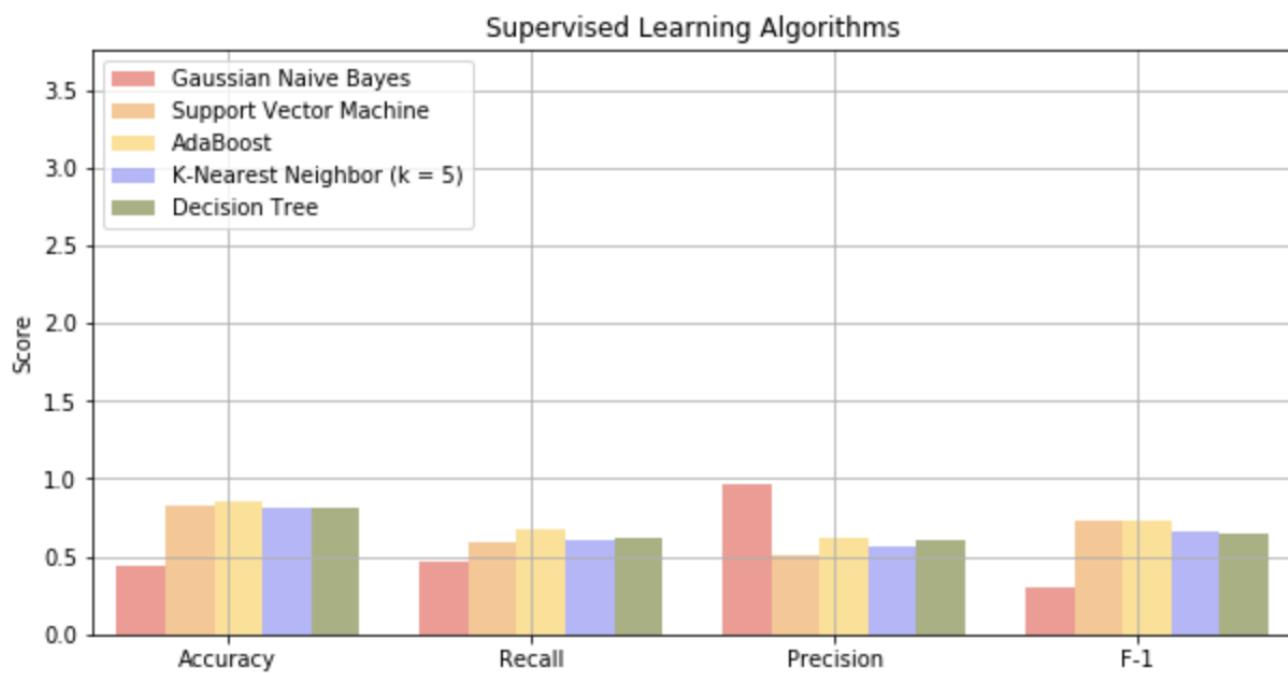


Figure 13: Supervised Learning Algorithm Results [44].

LIST OF TABLES

1	Statistical Summary of The Continuous Variables	20
2	Statistical Summary of Continuous Variables [44].	20
3	Train-Test-Split [44].	20
4	Results of the Algorithms [44].	20

	age	education	cap gain	cap loss	hours
count	32561	32561	32561	32561	32561
mean	38.581	10.08	1077.64	87.303	40.437
std.	13.640	2.572	7385.292	402.960	12.347
min.	17.0	1.0	0	0	1.0
25%	28.0	9.0	0	0	40.0
50%	37.0	10.0	0	0	40.0
75%	48.0	12.0	0	0	45.0
max	90.0	16.0	0	4356.0	99.0

Table 1: Statistical Summary of The Continuous Variables

	Age	Gain	Loss	Hours
Number of Instances	32,561	32,561	32,561	32,561
Mean	38.58	1077.64	87.303	40.437
Standard Deviation	13.640	7385.292	402.960	12.347
Minimum Value	17	0	0	1
25th percentile	28	0	0	40
50th percentile	37	0	0	40
75th percentile	48	0	0	45
Maximum Values	90	99999	4356	99

Table 2: Statistical Summary of Continuous Variables [44].

Splitting the Data	Sample Size
Training	24129
Testing	6033

Table 3: Train-Test-Split [44].

Name	Accuracy	Recall	Precision	F1 Score
Naive Bayes	0.4442	0.4642	0.9680	0.3053
SVC	0.8301	0.5969	0.5056	0.7284
AdaBoost	0.8499	0.6724	0.6189	0.7361
KNN	0.8184	0.6090	0.5682	0.6561
Decision Tree	0.8161	0.6231	0.6109	0.6459

Table 4: Results of the Algorithms [44].

```
bibtext report
```

```
=====
This is BibTeX, Version 0.99d (TeX Live 2016)
The top-level auxiliary file: report.aux
The style file: ACM-Reference-Format.bst
Database file #1: report.bib
```

```
=====
bibtext _ label error
```

```
=====
report.bib:518:@incollection{NIPS2013_4937,
```

```
=====
bibtext space label error
```

```
=====
report.bib:172:@INPROCEEDINGS{knn-chung,
```

```
=====
bibtext comma label error
```

```
=====
latex report
```

```
=====
[2017-12-04 12.24.22] pdflatex report.tex
This is pdfTeX, Version 3.14159265-2.6-1.40.17 (TeX Live 2016) (preloaded format=pdflatex)
bookmark level for unknown defaults to 0.
The anchor of a bookmark and its parent's must not be the same. Added a new anchor.
The anchor of a bookmark and its parent's must not be the same. Added a new anchor.
Typesetting of "report.tex" completed in 2.1s.
```

```
./README.yml
```

```
8:81      error    line too long (188 > 80 characters) (line-length)
8:188     error    trailing spaces (trailing-spaces)
21:81     error    line too long (534 > 80 characters) (line-length)
34:81     error    line too long (666 > 80 characters) (line-length)
34:666    error    trailing spaces (trailing-spaces)
35:12     error    trailing spaces (trailing-spaces)
37:30     error    trailing spaces (trailing-spaces)
42:5      error    duplication of key "type" in mapping (key-duplicates)
```

```
=====
Compliance Report
```

```
name: Usifo, Borga
hid: 343
paper1: 100 %
paper2: 100 %
project: 100 %
```

```
yamlcheck
```

```
wordcount
```

```
20
wc 343 project 20 5600 report.tex
wc 343 project 20 6227 report.pdf
wc 343 project 20 3252 report.bib
```

```
find "
```

```
passed: True
```

```
find footnote
```

```
passed: True
```

```
find input{format/i523}
```

```
4: \input{format/i523}
```

```
passed: True
```

```
find input{format/final}
```

```
passed: False
```

```
floats
```

```
47: We first used the pandas \cite{www-pandas} to help to load the
```

data in data frame format. This gave us a unique advantage, and faster processing of comma separated values for putting into data frame \cite{www-commasep}. Our data consist of 15 variables. Some of these variables are continuous, and some of them are categorical variables, and our target variable was “income” attribute. After putting the data into data frames, we first got a statistical snapshot of continuous variables (age, education, capital gain, capital loss, hours worked) by using the pandas \cite{www-pandas.describe} functions as shown in Table \ref{stats-table}.

```

49: \begin{table}[!ht]
64: \label{stats-table}
88: \par Our shape of the data will also receive impact from changing to numerical. Our number of futures will go from 14 to 103. This is because we implemented one-hot-encode to our dataset. It is called one hot encoded because we transform the categorical variables into a more acceptable shape for the machine learning algorithms to perform well \cite{www-hackernoon}. In other words ‘‘we implement binarization of the category to include as a future to train model \cite{www-hackernoon}’’. As we can see in Figure \ref{fig:one-hot-before} and Figure \ref{fig:one-hot-after}.
90: \begin{figure}[!ht]
92: \includegraphics[width=\columnwidth]{images/one-hot-before.png}
93: \caption{Example of One Hot Encoding Before \cite{www-hackernoon}.}\label{fig:one-hot-before}
96: \begin{figure}[!ht]
98: \includegraphics[width=\columnwidth]{images/one-hot-after.png}
99: \caption{Example of One Hot Encoding After \cite{www-hackernoon}.}\label{fig:one-hot-after}
104: After cleaning the data, we started our data exploration to learn little bit more from our data and make necessary changes if needed before putting into our machine learning algorithms. The first step in this process is getting the total count of the individuals as well as the count of the individuals who are making more than \$50K and less than \$50K which can be seen in below Table \ref{my-label-2}.
115: \label{my-label-2}
118: \par Moreover, we also look at the statistical values of each of the continuous variable we have. Those values given in Table \ref{my-label}. As we can see we have individuals who’re age ranging from 17 to 90 years old with a mean of 38.58. If we look at the capital gains and capital losses, we have a standard deviation of 7385 and 402 respectively this is also another indication of skew in these variables.
120: \begin{table}[!ht]
135: \label{my-label}

```

```

138: \par We used scatter matrix plot and applied the correlation
   function to see if we have any reliable correlation between any
   of the variables. As we can see from the correlation matrix
   Figure \ref{fig:scatter-matrix} and correlation numbers Figure
   \ref{fig:scatter} we do not have the high correlation between any
   variables. Correlation values range between -1 to 1. The
   correlation value of 1 is an indication of perfect positive
   correlation and correlation number -1 indicates a negative
   correlation between variables \cite{www-investopedia}. Because of
   lower correlation values, it will be tough to determine the
   classification by just looking at the correlations; this
   indicates we have sophisticated algorithms to determine the
   relationship between variables to classify individuals incomes.
140: \begin{figure}[!ht]
142: \includegraphics[width=\columnwidth]{images/scatter.png}
143: \caption{Correlation Matrix \cite{Borga2017}.}\label{fig:scatter-
   matrix}
146: \begin{figure}[!ht]
148: \includegraphics[width=\columnwidth]{images/scatter-matrix.png}
149: \caption{Scatter Matrix Plot
   \cite{Borga2017}.}\label{fig:scatter}
152: \par Furthermore, we also explore the capital gains, capital
   losses, and hours per week variables which we used a histogram to
   plot the data into distribution form so we can see how all these
   attributes distributed. The reason we do the histogram is we want
   to see any skewness in our data. As shown in the histogram graphs
   in Figure \ref{fig:Hist-capital} and Figure \ref{fig:loss-
   capital} in capital gains and capital loss we have highly skewed
   data which can cause issues later on in our algorithms. We apply
   a logarithmic function to do highly skewed data to less skewed
   \cite{www-onlinestat}. Using logarithmic functions adds more
   value to data from the interpretable standpoint and ‘‘it helps to
   meet the assumptions of inferential statistics \cite{www-
   onlinestat}’’.
154: \begin{figure}[!ht]
156: \includegraphics[width=\columnwidth]{images/capital-gain.png}
157: \caption{Histogram of Capital Gain
   \cite{Borga2017}.}\label{fig:Hist-capital}
160: \begin{figure}[!ht]
162: \includegraphics[width=\columnwidth]{images/capital-loss.png}
163: \caption{Histogram of Capital Loss
   \cite{Borga2017}.}\label{fig:loss-capital}
166: \par Moreover, applying logarithmic function had an impact on
   distribution. We can see the changes on skew data in Figure
   \ref{fig:Hist-capital-log} after applying logarithmic function.
168: \begin{figure}[!ht]

```

```

170: \includegraphics[width=\columnwidth]{images/logarithmic-
    applied.png}
171: \caption{After Logarithmic Function Applied Histogram of Capital
    Gain \cite{Borga2017}.}\label{fig:Hist-capital-log}
208: \par Term naive comes from it is the method of use probability
    among categories which assumes of independence among given class
    of attributes as shown in Figure \ref{fig:Naive Bayes}. In other
    words, if we try to classify individuals from their email
    communications it will not take the order of words into account.
    Whereas in the English language we can tell the difference
    between sentence makes sense or not if we randomly re-order our
    words in the sentences. So it does not understand the text, it
    only looks at word frequencies as a way to do the classification.
    This is why it is called ‘‘Naive’’.
210: \begin{figure}[!ht]
213: \includegraphics[width=\columnwidth]{Naive-bayes}
214: \caption{Example of Naive Bayes \cite{Zhang}.}\label{fig:Naive
    Bayes}
299: \par Representation of decision trees is as follows; we sort the
    instances from root to leaf nodes, this sorting gives insights
    about the classification of the instance, every outcome
    descending from the root node corresponds to possible values for
    that variable \cite{www-cs.princeton}. We can classify an
    instance by starting from the root node and checking the
    attributes labeled on that node and moving down from that node
    based on attribute given attribute values \cite{www-cs.princeton}
    as shown in Figure \ref{fig:Decision Tree}.
301: \begin{figure}[!ht]
303: \includegraphics[width=\columnwidth]{images/decison_tree.png}
304: \caption{Example of Decision Tree Construction\cite{www-
    cs.princeton}.}\label{fig:Decision Tree}
352: \begin{figure}[!ht]
354: \includegraphics[width=\columnwidth]{images/hyperplane-
    boundary.png}
355: \caption{Example of Shows the Hyperplanes \cite{www-simafore-
    svm}.}\label{fig:Hyperplane}
358: \par Traditionally boundaries are called ‘‘hyperplanes’’ and it
    is used to describe boundaries in more than three dimensions
    because they are hard or sometimes impossible to
    visualize.\cite{www-simafore}. Figure \ref{fig:Hyperplane}.
    Optimality of hyperplane expressed as a linear function which
    requires maximum distance between the identified classes. It only
    considers a small number of training example to build this
    hyperplane. SVM hyperplanes based on ‘‘ separation of positive
    (+1) and negative (-1) with the largest margin \cite{verma-
    ssv}’’.

```

```

403: \par We used scikit-learn for splitting data into train and test  

    we saved 20\% of data for testing purposes as shown in Table  

    \ref{split} .  

405: \begin{table}![ht]  

414: \label{split}  

430: Confusion matrix develops from 4 key elements. These elements are  

    true positive, true negative, false negative, and false positive.  

    As shown in Figure \ref{fig:confusion-matrix} about the  

    constructing a confusion matrix. If we want to build a confusion  

    matrix by targeting individuals who are making more than \$50K  

    our true positive, true negative, false positive, and false  

    negative explained below.  

432: \begin{figure}![ht]  

434: \includegraphics[width=\columnwidth]{images/confusion-matrix.png}  

435: \caption{Example of Confusion Matrix Construction \cite{www-  

    exsilio}.}\label{fig:confusion-matrix}  

477: Now we can look at the results from each of the machine learning  

    algorithm. Results also showed in Table \ref{result-table} with  

    the visualization of Figure \ref{fig:result-algo}. We can also  

    see the running time of the each of the algorithm in Figure  

    \ref{fig:result-time}. Support Vector Machines is the winner for  

    the highest running time for training the algorithm.  

479: \begin{figure}![ht]  

481: \includegraphics[width=\columnwidth]{images/running-time.png}  

482: \caption{Supervised Learning Algorithm Running Time Results  

    \cite{Borga2017}.}\label{fig:result-time}  

485: \begin{table}![ht]  

497: \label{result-table}  

501: As shown in the Figure \ref{fig:result-algo} we have a comparison  

    of several supervised machine learning algorithms on our dataset.  

    We can see that from the accuracy standpoint Naive Bayes  

    algorithms have the lowest score which means that it did not do a  

    good job for labeling true positives regards to all data but it  

    did a good job in precision standpoint while doing a bad  

    classification from recall standpoint. Two key element for us in  

    this situation is accuracy and f1 score(which consist of  

    precision and recall).  

505: As we stated before ensemble algorithms learn from the small  

    portion of the data and combine these learning to do the  

    predictive task. As shown in Figure \ref{fig:result-algo}  

    adaboosting has the highest accuracy score among all the other  

    algorithms. This algorithm should be our first choice to do  

    predictive modeling. We believe that there is still an  

    improvements on accuracy  

510: Decision Tree is gave a good accuracy but fall behind on f1 score  

    as shown in Figure \ref{fig:result-algo}.

```

```
517: \begin{figure}[!ht]
519: \includegraphics[width=\columnwidth]{images/result-score.png}
520: \caption{Supervised Learning Algorithm Results
    \cite{Borga2017}.}\label{fig:result-algo}
```

figures 13
tables 4
includegraphics 13
labels 18
refs 16
floats 17

```
True : ref check passed: (refs >= figures + tables)
False : label check passed: (refs >= figures + tables)
True : include graphics passed: (figures >= includegraphics)
False : check if all figures are referred to: (refs >= labels)
```

Label/ref check
passed: True

When using figures use columnwidth
[width=1.0\columnwidth]
do not change the number to a smaller fraction

find textwidth

passed: True

below_check

WARNING: algorithm and below may be used improperly

```
104: After cleaning the data, we started our data exploration to learn
    little bit more from our data and make necessary changes if
    needed before putting into our machine learning algorithms. The
    first step in this process is getting the total count of the
    individuals as well as the count of the individuals who are
    making more than \$50K and less than \$50K which can be seen in
    below Table \ref{my-label-2}.
```

WARNING: code and below may be used improperly

184: Scikit-learn developed by David Cournapeau in 2007. The development came from while he was working on summer code project for Google. After recognized and published by INRIA in 2010 project start the get more attention among worldwide. There are more than 30 active contributors and has secured several sponsorships from big technology companies\cite{www-machinelearningmystery}. ‘‘It also has a goal of providing common algorithms to Python users through consistent interface\cite{www-oreilly}’’. Scikit-Learn consists of several elements to make analytical predictions. These elements are shown below\cite{www-analyticvidhya}:

WARNING: algorithm and below may be used improperly

184: Scikit-learn developed by David Cournapeau in 2007. The development came from while he was working on summer code project for Google. After recognized and published by INRIA in 2010 project start the get more attention among worldwide. There are more than 30 active contributors and has secured several sponsorships from big technology companies\cite{www-machinelearningmystery}. ‘‘It also has a goal of providing common algorithms to Python users through consistent interface\cite{www-oreilly}’’. Scikit-Learn consists of several elements to make analytical predictions. These elements are shown below\cite{www-analyticvidhya}:

WARNING: algorithm and below may be used improperly

194: \par Algorithms that we will investigate shown below; we will go more deep analysis on each of these algorithms.

bibtex

label errors

518: NIPS2013_4937: do not use underscore in labels:

bibtex errors

This is BibTeX, Version 0.99d (TeX Live 2016)
The top-level auxiliary file: report.aux
The style file: ACM-Reference-Format.bst
Database file #1: report.bib

```
bibtex_empty_fields
```

```
entries in general should not be empty in bibtex
```

```
find ""
```

```
passed: True
```

```
ascii
```

```
=====
The following tests are optional
=====
```

```
Tip: newlines can often be replaced just by an empty line
```

```
find newline
```

```
passed: True
```

```
cites should have a space before \cite{} but not before the {
```

```
find cite {
```

```
passed: True
```

Predicting Housing Prices - Kaggle Competition

Murali Cheruvu, Anand Sriramulu

Indiana University

3209 E 10th St

Bloomington, Indiana 47408

mcheruvu@iu.edu, asriram@iu.edu

ABSTRACT

Apply exploratory data analysis and implement various advanced supervised machine learning algorithms to predict neighborhood housing sale prices found in the sample test dataset. Compare the predicted models and results from these advanced supervised algorithms. Apply ensembled model to achieve better predictions, hence get good score in kaggle competition.

KEYWORDS

i523, hid306, Supervised Learning Algorithms, Exploratory Data Analysis, Kaggle

1 INTRODUCTION

Part of the kaggle competition, two sample data sets are given with 80 attributes (variables) describing various aspects of the residential homes in Ames and Iowa cities. Training dataset contains sale price of the homes, and using this training data set, how accurately we can predict Sale Prices of the homes in the test dataset using preprocessing and thorough data analysis. Many developers used advanced learning algorithms - XGBoost, Lasso and Neural Network, to predict the sale prices in the kaggle competition and achieved better kaggle scores. Kaggle score is a measure to indicate accuracy and the quality of the algorithm. We have applied various exploratory analysis techniques and engineer the features before applying a few advanced supervised learning algorithms.

2 EXPLORATORY DATA ANALYSIS

There are 1460 rows in the training data set and 1459 rows in the test dataset. Out of the 80 variables, 23 are nominal, 23 are ordinal, 14 are discrete, and 20 are continuous. We have combined training and testing datasets for easier analysis. We excluded Id attribute as it does not add value in the modeling. We also removed Sale Price, the target variable, from the training dataset. All attribute details are given in the appendix section as a reference.

2.1 Handling Missing Values

First part of the analysis, we have checked for the missing values in the dataset. we have also identified that there are 6 variables having most of the missing data. All the missing values for the numeric variables are analyzed further to decide whether to delete the instances of the data with missing values or fill them using meaningful data such as median of the corresponding variable.

```
# python code - check for null values
train = pd.read_csv("../data/train.csv")
test = pd.read_csv("../data/test.csv")
```

```
#combine the data sets
alldata = train.append(test)
na = alldata.isnull().sum()
.sort_values(ascending=False)
```

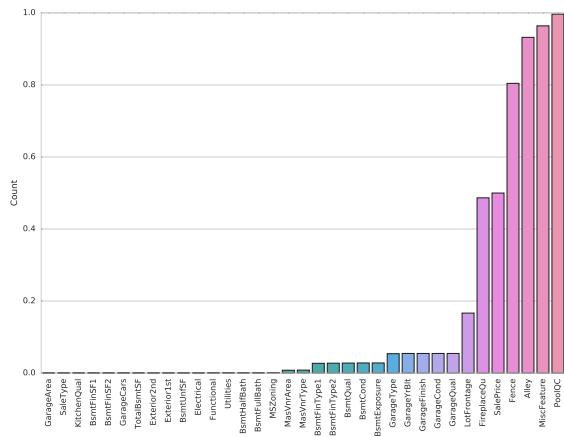


Figure 1: Variables with Missing Values

2.2 Analyze Numerical Variables

There are 37 numeric variables after excluding the *Id* variable. We have analyzed the numeric variables for data patterns such as skewed data and range of the possible values. *over all quality, ground living area, garage cars and garage area* graphs are included as samples.

```
# python code - analyze numeric variables
numerical_features = [f for f in train.columns
if train.dtypes[f] != 'object']

nd = pd.melt(train, value_vars = numerical_features)
plt.figure(figsize = (5,3))
plot = sns.FacetGrid(nd, col='variable', col_wrap=4,
sharex=False, sharey = False)
plot = plot.map(sns.distplot, 'value')
```

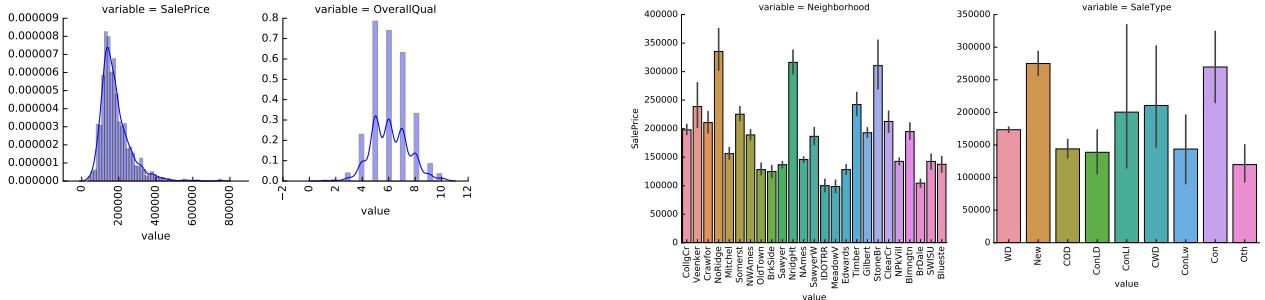


Figure 3: Sample Category Variables

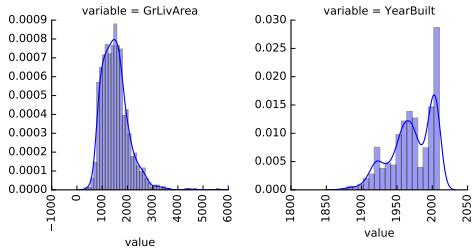


Figure 2: Sample Numerical Variables

2.3 Analyze Categorical Variables

There are 43 categorical variables in the dataset. We have analyzed all categorical variables and found the ways to fill the missing values and also the way we need to convert them into numerical variables. Later in the feature engineering section, we will go through the details. *neighborhood* and *sale type* graphs are included as samples.

```
# python code - analyze numeric variables
cat_features = [f for f in train.columns
if train.dtypes[f] == 'object']
print(cat_features)

def barplot(x,y,**kwargs):
sns.barplot(x=x,y=y)
x = plt.xticks(rotation=90)

plt.figure(figsize = (5,3))

p = pd.melt(train, id_vars='SalePrice',
value_vars=cat_features)

g = sns.FacetGrid (p, col='variable', col_wrap=4,
sharex=False, sharey=False, size=5)

g = g.map(barplot, 'value','SalePrice')
```

2.4 Analyze Correlations

Numpy package offers correlations functionality to analyze the variables that highly positively or negatively correlated with *sale price* and with the other variables. We can visualize a few pair-wise correlation graphs with sale price for detailed analysis. From the correlation-plot we can list the top 10 features those are strongly correlated with the target variable - *sale price*.

```
# python code
corr = alldata[numerical_features].corr()
mask = np.zeros_like(corr)
mask[np.triu_indices_from(mask)] = True
plt.figure(figsize = (15,8))
sns_plot = sns.heatmap(corr, cmap="YlGnBu",
lineweights=.5, mask=mask, vmax=.3)
```

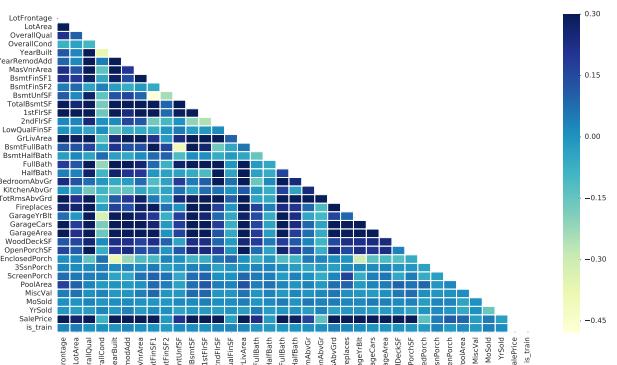


Figure 4: Correlations

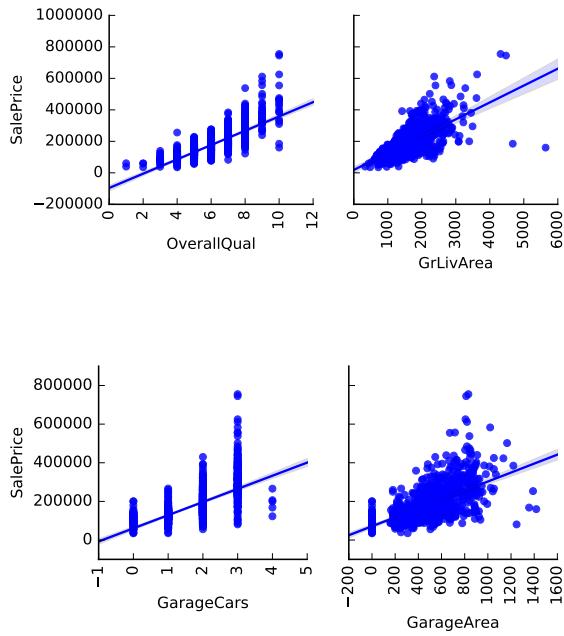


Figure 5: Pair-wise Correlations with Sale Price

- (1) OverallQual: Overall material and finish quality
- (2) GrLivArea: Above ground living area square feet
- (3) GarageCars: Size of garage in car capacity
- (4) GarageArea: Size of garage in square feet
- (5) TotalBsmtSF: Total square feet of basement area
- (6) 1stFlrSF: First Floor square feet
- (7) FullBath: Full bathrooms above grade
- (8) TotRmsAbvGrd: Total rooms above ground
- (9) YearBuilt: Original construction date
- (10) GarageYrBlt: Garage built year

2.5 Skewed Data Analysis

From the numerical analysis, we have identified that there are a few numerical variables need further analysis to identify the skewed data. We did not find any key variables those have skewed more than 75%. However, we wanted to replace the *sale price* with corresponding logarithmic value for the predictive models and later convert it back to the exponential value before submitting to the kaggle competition.

2.6 Outlier Analysis

Continuing with exploratory analysis, We have analyzed the outliers using *Cook's distance* and then further analyzed two key variables - *ground live area* and *garage area* that are in high correlation with *sale price*. We have removed the outlier rows related to these two variables as they are only 8 rows impacting the training dataset.

```
# python code - outlier analysis
import statsmodels.api as sm
from statsmodels.formula.api import ols
```

```
model = ols(formula = "SalePrice ~
GrLivArea + GarageArea", data=train)
fitted = model.fit()
plot = sm.graphics.influence_plot(fitted,
criterion="cooks")
```

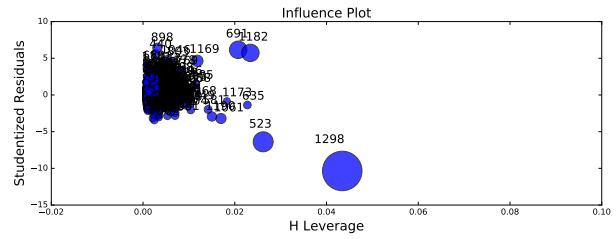


Figure 6: Outliers - Cook's distance

```
# python code - remove outlier rows
# fix all extreme outliers based on outlier analysis
# 8 rows will be deleted
train = train[train.GrLivArea <= 4000]
train = train[train.GarageArea <= 1200]
```

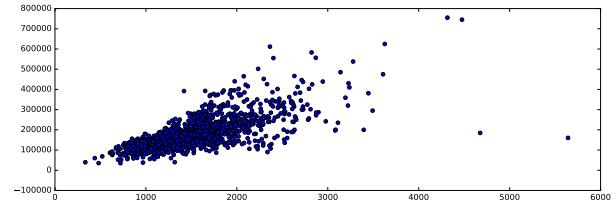


Figure 7: Outliers - Garage Live Area

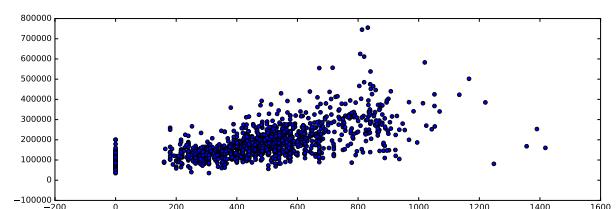


Figure 8: Outliers - Garage Area

2.7 Feature Engineering

Feature engineering is a technique to analyze all the variables those influence target variable for better predictions. Part of feature engineering, we may need to create new features to make the data to be more expressive. One of the key intents, in analyzing categorical variables, is to convert them into numerical factors as most of the machine learning algorithms expect all the variables to be numeric for them to work more effectively. Some of the categorical variables are ordinal. we can use T-shirt sizes: small, medium and large as

an example to explain an ordinal variable. When we convert this category variable into numeric encoding, we need to retain the fact that there is an implicit order within the values. Supposing we give ordinal encoding as - small = 1, medium = 2 and large = 3; we will satisfy the implicit order or weightage and that helps in modeling the system by elevating the importance of this implicit ordering in the values of the ordinal variable.

There are a few other encoding techniques, such as one-hot, binary, polynomial and helmert to factorize categorical variables. We will use ordinal and one-hot encoding techniques for this data set. One-hot encoding converts the category variable into many binary vectors, one new numeric variable for each value in the category. Assume that we have a category variable called signal-light with three possible values: green, yellow and red. We will need to convert these values into numeric - green = 1, yellow = 2 and red = 3. When we apply one-hot encoding on this variable, basically we are creating three new categorical variables - signal-light-green, signal-light-yellow and signal-light-red along with the original variable - signal-light, each is pretty much a binary vector having 1s for all the corresponding values; otherwise 0s. With hot-encoding, we are basically increasing the dimensions in the model. After extensive feature engineering applied on the housing dataset, we have added 228 new features (variables). Following are the python methods to factorize categorical variables and one-hot encoding techniques.

```
# python code - factorize and one-hot
def get_one_hot(df, col_name, fill_val):
if fill_val is not None:
df[col_name].fillna(fill_val, inplace=True)

dummies = pd.get_dummies(df[col_name], prefix="_" + col_name)
df = df.join(dummies)
df = df.drop([col_name], axis=1)
return df
#end def

from sklearn.preprocessing import LabelEncoder

def factorize(df, column, fill_na=None):
le = LabelEncoder()
if fill_na is not None:
df[column].fillna(fill_na, inplace=True)
le.fit(df[column].unique())
df[column] = le.transform(df[column])
return df
#end def
```

3 ALGORITHMS AND METHODOLOGY

Linear regression predicts the target variable using best possible straight line fit to the set predictor variables. The possible best fit is usually the one that minimizes the root mean squared error (RMSE) between the actual and predicted data points. However, with complex problem space such as the housing prices dataset, we have lots of variables relating to the target variable in a non-linear fashion. Trivial supervised learning algorithms will not be

effective to provide accurate *sale price* predictions. To overcome this challenge, we have applied various advanced supervised learning algorithms, such as Support Vector Machine (SVM), Random Forest, Lasso, Ridge, XGBoost and Neural Network, to predict the test data housing prices.

3.1 Support Vector Machine (SVM) Algorithm

Support Vector Machine (SVM) algorithms can be used to solve classification and regression problems. SVM regression relies on kernel functions for modeling the data. SVM creates larger margins between categories of data so that they are linearly separable. SVM handles non-linearly separable data, mainly for regression problems, using kernel functions, such as polynomial, radial basis function (RBF) and sigmoid, to project the data onto a hyperplane. Following python code snippet is the implementation for *sale price* predictions of the housing test dataset.

```
# python code - SVM algorithm
from sklearn.svm import SVR
import matplotlib.pyplot as plt
from sklearn.metrics import mean_squared_error

train = pd.read_csv("train.csv")
test = pd.read_csv("test.csv")
target_vector = train["SalePrice"]

_svm_algo = SVR(kernel = 'rbf', C=1e3, gamma=1e-8)
_svm_algo.fit(train, target_vector)

y_train = target_vector
y_train_pred = _svm_algo.predict(train)

#root mean squared error (RMSE)
rmse_train = np.sqrt(rmse(y_train,y_train_pred))

y_test_pred = _svm_algo.predict(test)
```

3.2 Random Forest Algorithm

Random Forest is an advanced machine learning algorithm for predictive analytics. Random Forest ensembles multiple decision trees to create an additive learning model from the sequence of base models created by each decision tree that worked on a sub-sample dataset. Random Forest models are suitable to handle tabular datasets with hundreds of numeric and categorical features. Along with missing values, non-linear relations between features and the target, will be handled well by random forest algorithms. With proper tuning of hyper-parameters of the random forest algorithm, it can perform well with decent accuracy in the predictions without overfitting the model. Unlike similar regression models, it does not offer feature coefficient information but it provides *feature ranking* functionality very nicely. Following are the random forest algorithm details for the *sale price* predictions implemented in python using *sklearn* package.

```
# python code - random forest algorithm
```

```

from sklearn.ensemble import RandomForestRegressor
from sklearn.metrics import accuracy_score
from sklearn.metrics import confusion_matrix
from sklearn.metrics import mean_squared_error

train = pd.read_csv("train.csv")
test = pd.read_csv("test.csv")
target_vector = train["SalePrice"]

_algo = RandomForestRegressor(n_estimators=100,
oob_score=True, random_state=123456)

model = _algo.fit(train, target_vector)

feat_imp = pd.Series(_algo.feature_importances_,
train.columns).sort_values(ascending=False)

feat_imp[:10].plot(kind='bar',
title='Feature Ranmkngt')
y_train = target_vector
y_train_pred = _algo.predict(train)

#root mean squared error (RMSE)
rmse_train = np.sqrt(rmse(y_train,y_train_pred))
y_test_pred = _algo.predict(test)

```

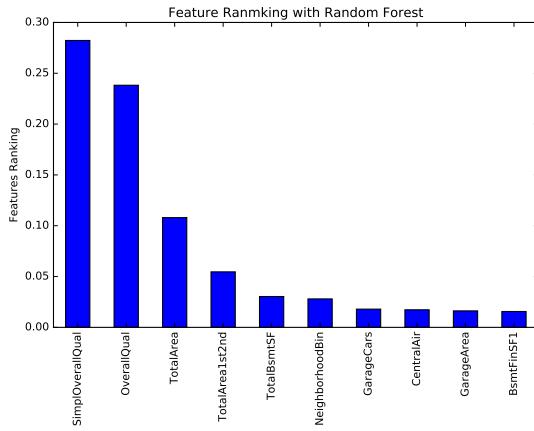


Figure 9: Random Forest Feature Ranking

3.3 Lasso Algorithm

Lasso is a regression model that uses shrinkage to bring data points towards the center, similar to the mean value of all the data points. Lasso stands for Least Absolute Shrinkage and Selection Operator. It is a regularized linear model with penalty term *lambda* to minimize the error. Parameter penalization controls overfitting the input data by shrinking variable coefficients to 0. Essentially this makes the variables no effect in the model, hence reduces the dimensions. Following is the lasso algorithm implementation for *sale price* predictions in python.

```
# python code - lasso algorithm
```

```

from sklearn.linear_model import Lasso
from sklearn.metrics import mean_squared_error

train = pd.read_csv("train.csv")
test = pd.read_csv("test.csv")
target_vector = train["SalePrice"]

#found this best alpha value through cross-validation
_best_alpha = 0.0001

_lasso_algo = Lasso(alpha = _best_alpha,
max_iter = 50000)

model = _lasso_algo.fit(train, target_vector)

y_train = target_vector
y_train_pred = _algo.predict(train)

#root mean squared error (RMSE)
rmse_train = np.sqrt(rmse(y_train,y_train_pred))

y_test_pred = _lasso_algo.predict(test)

```

3.4 Ridge Algorithm

Ridge algorithm is very similar to lasso algorithm with the same goal. While lasso performs *L1 regularization*, ridge applies *L2 regularization* techniques in modeling the predictions. L1 regularization adds penalty to the variables equivalent to *absolute value of the magnitude* of the coefficients, whereas L2 adds the penalty equivalent to *square of the magnitude* of the variable coefficients. Following is the python implementation of the ridge algorithm for the *sale price* predictions.

```

# python code - ridge algorithm
from sklearn.linear_model import Ridge
from sklearn.metrics import mean_squared_error

train = pd.read_csv("train.csv")
test = pd.read_csv("test.csv")
target_vector = train["SalePrice"]

#found this best alpha value through cross-validation
_best_alpha = 0.00099

_ridge_algo = Ridge(alpha = _best_alpha,
normalize = True)

_ridge_algo.fit(train, target_vector)

df = {'features': train.columns.values,
'Coefficients': _ridge_algo.coef_[0]}
coefficients = pd.DataFrame(df)
.coefficients.sort_values(by='Coefficients',
ascending=False)

plt.figure()

```

```

coefficients.iloc[0:10].plot(x=['features'],
kind='bar', title='Top 10 Positive Features')
plt.ylabel('Feature Coefs')
plt.figure()
coefficients.iloc[-10: ].plot(x=['features'],
kind='bar', title='Top 10 Negative Features')
plt.ylabel('Feature Coefs')

y_train = target_vector
y_train_pred = _ridge_algo.predict(train)

#root mean squared error (RMSE)
rmse_train = np.sqrt(rmse(y_train,y_train_pred))

y_test_pred = _ridge_algo.predict(test)

df_predict = pd.DataFrame({'Id': test["Id"],
'SalePrice': np.exp(y_test_pred) - 1.0})

#df_predict = pd.DataFrame({'Id': id_vector,
'SalePrice': sale_price_vector})

df_predict.to_csv('../data/kaggle_python_ridge.csv',
header=True, index=False)

```

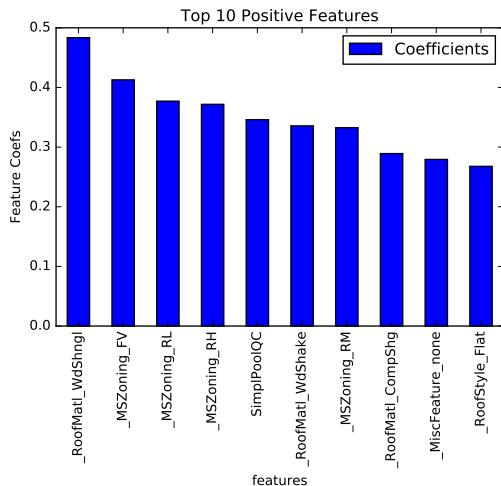


Figure 10: Ridge Algorithm - Top 10 Positive Features

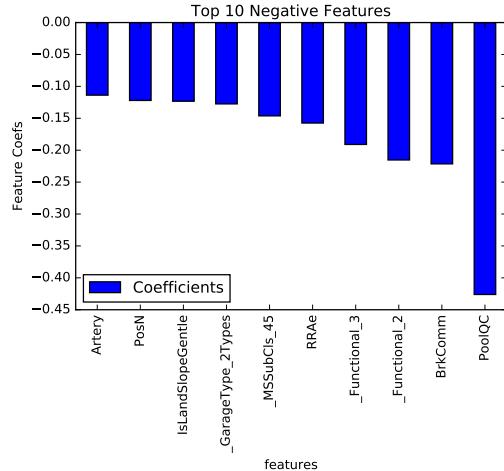


Figure 11: Ridge Algorithm - Top 10 Negative Features

3.5 XGB Boosting Algorithm

XGBoost (eXtreme Gradient Boosting) is one of the Gradient Boosted Machine algorithms. It ensembles (combines) optimized model by taking trained models from all the preceding iterations. XGBoost regularizes the variables (parameters) to reduce the overfit and can work well with variables having missing values. It is empowered with built-in cross validation to reduce the boosting iterations; hence offers better performance along with parallel processing on distributed systems such as Hadoop. By tuning the XGBoost hyper parameters, we can achieve well optimized model that can make more accurate predictions. XGBoost uses *Fscore* to measure the importance of variables. Following table explains the hyper-parameters of XGBoost algorithm and also given the python code implementing XGBoost algorithm for *sale price* predictions.

Table 1: XGBoost Hyper-parameters

Hyper-parameter	Description
Maximum Iterations:	Number of trees in the final model. More the trees, more accuracy.
Maximum Depth:	Depth of each individual tree to control overfitting.
Step Size:	Shrinkage, works similar to learning rate; smaller value takes more iterations.
Column Subsample:	Subset of the columns to use in each iteration

```

# python code - XGBoost algorithm
import xgboost as xgb
from xgboost import XGBClassifier
from xgboost import plot_importance
from sklearn.metrics import mean_squared_error
from sklearn import cross_validation, metrics

```

```

train = pd.read_csv("train.csv")
test = pd.read_csv("test.csv")
target_vector = train["SalePrice"]

_algo = np.random.seed(1234)

_xgb_algo = xgb.XGBRegressor(
    colsample_bytree=0.8,
    colsample_bylevel = 0.8,
    gamma=0.01,
    learning_rate=0.05,
    max_depth=5,
    min_child_weight=1.5,
    n_estimators=6000,
    reg_alpha=0.5,
    reg_lambda=0.5,
    subsample=0.7,
    seed=42,
    silent=1)

_xgb_algo.fit(train, target_vector)

feat_imp = pd.Series(_xgb_algo.booster()
    .get_fscore()
    .sort_values(ascending=False)[0:10]
    .plot(kind='bar',
    title='Top 10 Feature Importances')

y_train = target_vector
y_train_pred = _xgb_algo.predict(train)

#root mean squared error (RMSE)
rmse_train = np.sqrt(rmse(y_train,y_train_pred))

y_test_pred = _xgb_algo.predict(test)

```

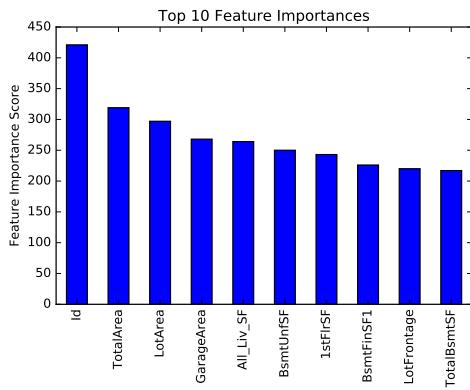


Figure 12: XGBoost Feature Importance

3.6 Neural Network Algorithm

Neural Network is, a *directed graph*, organized by layers and layers are created by number of interconnected neurons (or nodes). Every

neuron in a layer is connected with all the neurons from previous layer; there will be no interaction of neurons within a layer. The performance of a Neural Network is measured using *cost or error function* and the dependent input *weight* variables. *Forward-propagation* and *back-propagation* are two techniques, neural network uses repeatedly until all the input variables are adjusted or calibrated to predict accurate output. During, forward-propagation, information moves in forward direction and passes through all the layers by applying certain weights to the input parameters. *Back-propagation* method minimizes the error in the *weights* by applying an algorithm called *gradient descent* at each iteration step. We have used *TensorFlow* python library to predict the *sale price* of housing dataset using simple feed-forward neural network. TensorFlow uses *tensors*, special multi-dimensional arrays to store the datasets for easier linear algebra and vector calculus operations.

3.7 Model Ensembling

We can create a robust predictive model with better accuracy by merging two or more machine learning algorithms. This technique is called *model ensembling*. Ensembled algorithms may be similar in functionality or may entirely be different from each other. Individual algorithms may not perform great but by ensembling them, the overall system can offer much better performance and accuracy. Variations in the predicting logic in each of these individual algorithms will bring unbiasedness into the unified model. *Bagging*, *boosting* and *stacking* are popular ensembling techniques. Many of the advanced machine learning algorithms use ensembled approaches to achieve accurate classifications or predictions. Random Forest uses bagging, XGBoost uses boosting and Neural Network applies stacking ensembling techniques. For the kaggle submission, we have created an ensembled model by averaging *Sale Price* of the top 3 performing ensembled algorithms - XGBoost, Lasso and Neural Network. Following are the kaggle public scores that we have achieved by submitting all the six algorithms along with the ensembled model. As predicted, ensembled model has scored better compared to the individual algorithms. By applying advanced machine learning algorithms, we have placed our scores within top 15% of the competition. Following table displays each algorithm and the *root-mean-square error* (RMSE) along with the *kaggle score*.

Table 2: Kaggle Submissions

Algorithm	RMSE	Kaggle Score
SVM	0.2069	0.23967
Random Forest	0.0519	0.14607
Ridge	0.0988	0.13687
XGBoost	0.0432	0.13018
Lasso	0.1015	0.12860
Neural Network	0.20	0.12510
Ensemble		0.12011

4 DEVELOPMENT ENVIRONMENT

- Operating Environment: Ubuntu 16.4 through Oracle Virtual Box 5.2

- Programming Language: Python 2.7
- Development Tools/Environment: Jupyter Notebook, Anaconda (data science platform)
- Python Libraries: numpy, pandas, sklearn, matplotlib, seaborn, xgboost and tensorflow
- Repository: git@github.com:bigdata-i523/hid306.git
- Project Folders:
 - Code: all Jupyter notebook files
 - Images: all output images as PDF files
 - Data: all the input and output datasets in CSV format
- Python Jupyter notebook files:
 - 1.1_exploratory_analysis_numerical.ipynb
 - 1.2_exploratory_analysis_categorical.ipynb
 - 1.3_outlier_and_skewed_data_analysis.ipynb
 - 1.4_feature_engineering.ipynb
 - 2.1_algorithm_svm.ipynb
 - 2.2_algorithm_random_forest.ipynb
 - 2.3_algorithm_ridge.ipynb
 - 2.4_algorithm_lasso.ipynb
 - 2.5_algorithm_neural_network_tf.ipynb
 - 2.6_algorithm_xgboost.ipynb
 - 3_ensemble_kaggle_submission.ipynb
- Input Datasets:
 - train.csv
 - test.csv
- Kaggle Submissions:
 - kaggle_python_svm.csv
 - kaggle_python_random_forest.csv
 - kaggle_python_ridge.csv
 - kaggle_python_xgboost.csv
 - kaggle_python_lasso.csv
 - kaggle_python_neural_network.csv
 - kaggle_python_ensemble.csv

5 CONCLUSION

Generally, ensemble models performs better compared to individual algorithms. However, there are a few factors that influence accuracy and performance of the algorithms, such as handcrafted feature engineering, proper cost function with regularized input to address non-linearities in the training datasets and tuning hyper-parameters of the algorithms. While Deep Learning Neural Networks are good for image processing, K-Nearest Neighbor algorithms can handle unsupervised datasets with less complexity. Domain knowledge and algorithm selection play vital role in getting accurate predictions. XGBoost, Random Forest, Lasso and Neural Networks are advanced machine learning algorithms dominating in the data science competitions for classification and regression related tasks. With ensembling and iterative learning techniques, they can scale well and offer better predictions for huge datasets having large number of features.

A KAGGLE HOUSING PRICE DATASET VARIABLES

- Id: Row id
- SalePrice: Sale price of the house in dollars. This is the target variable to predict.

- MSSubClass: The building class
- MSZoning: The general zoning classification
- LotFrontage: Linear feet of street connected to property
- LotArea: Lot size in square feet
- Street: Type of road access
- Alley: Type of alley access
- LotShape: General shape of property
- LandContour: Flatness of the property
- Utilities: Type of utilities available
- LotConfig: Lot configuration
- LandSlope: Slope of property
- Neighborhood: Physical locations within Ames city limits
- Condition1: Proximity to main road or railroad
- Condition2: Proximity to main road or railroad (if a second is present)
- BldgType: Type of dwelling
- HouseStyle: Style of dwelling
- OverallQual: Overall material and finish quality
- OverallCond: Overall condition rating
- YearBuilt: Original construction date
- YearRemodAdd: Remodel date
- RoofStyle: Type of roof
- RoofMatl: Roof material
- Exterior1st: Exterior covering on house
- Exterior2nd: Exterior covering on house (if more than one material)
- MasVnrType: Masonry veneer type
- MasVnrArea: Masonry veneer area in square feet
- ExterQual: Exterior material quality
- ExterCond: Present condition of the material on the exterior
- Foundation: Type of foundation
- BsmtQual: Height of the basement
- BsmtCond: General condition of the basement
- BsmtExposure: Walkout or garden level basement walls
- BsmtFinType1: Quality of basement finished area
- BsmtFinSF1: Type 1 finished square feet
- BsmtFinType2: Quality of second finished area (if present)
- BsmtFinSF2: Type 2 finished square feet
- BsmtUnfSF: Unfinished square feet of basement area
- TotalBsmtSF: Total square feet of basement area
- Heating: Type of heating
- HeatingQC: Heating quality and condition
- CentralAir: Central air conditioning Electrical: Electrical system
- 1stFlrSF: First Floor square feet
- 2ndFlrSF: Second floor square feet
- LowQualFinSF: Low quality finished square feet (all floors)
- GrLivArea: Above grade (ground) living area square feet
- BsmtFullBath: Basement full bathrooms
- BsmtHalfBath: Basement half bathrooms
- FullBath: Full bathrooms above grade
- HalfBath: Half baths above grade
- Bedroom: Number of bedrooms above basement level
- Kitchen: Number of kitchens
- KitchenQual: Kitchen quality

- TotRmsAbvGrd: Total rooms above grade (does not include bathrooms)
- Functional: Home functionality rating
- Fireplaces: Number of fireplaces
- FireplaceQu: Fireplace quality
- GarageType: Garage location
- GarageYrBlt: Year garage was built
- GarageFinish: Interior finish of the garage
- GarageCars: Size of garage in car capacity
- GarageArea: Size of garage in square feet
- GarageQual: Garage quality
- GarageCond: Garage condition
- PavedDrive: Paved driveway
- WoodDeckSF: Wood deck area in square feet
- OpenPorchSF: Open porch area in square feet
- EnclosedPorch: Enclosed porch area in square feet
- 3SsnPorch: Three season porch area in square feet
- ScreenPorch: Screen porch area in square feet
- PoolArea: Pool area in square feet
- PoolQC: Pool quality
- Fence: Fence quality
- MiscFeature: Miscellaneous feature not covered in other categories
- MiscVal: Dollar Value of miscellaneous feature
- MoSold: Month Sold
- YrSold: Year Sold
- SaleType: Type of sale
- SaleCondition: Condition of sale

ACKNOWLEDGMENTS

The authors would like to thank Dr. Gregor von Laszewski and the Teaching Assistants for their support and great suggestions. Authors would also want to thank Kaggle Website and the developers for their valuable information, ideas and contributions.

REFERENCES

- [1] AiO. 2017. House Prices: Advanced Regression Techniques. (Feb. 2017). <https://www.kaggle.com/notapple/detailed-exploratory-data-analysis-using-r>
- [2] Tanner Carbonati. 2017. Detailed Data Analysis & Ensemble Modeling. (Aug. 2017). <https://www.kaggle.com/tannercarbonati/detailed-data-analysis-ensemble-modeling/notebook>
- [3] Yeshwant Chillakuru, Michael Arango, Jack Crum, and Paul Brewster. 2017. Using Neighborhood Level Data to Predict the Residential Sale Price of Properties in Ames, Iowa. (May 2017). <https://rpubs.com/jackcrum/281471>
- [4] Aarshay Jain. 2016. Complete Guide to Parameter Tuning in XGBoost. (March 2016). <https://www.analyticsvidhya.com/blog/2016/03/complete-guide-parameter-tuning-xgboost-with-codes-python/>
- [5] Selva Prabhakaran. [n. d.]. Outlier Treatment. ([n. d.]). <http://r-statistics.co/Outlier-Treatment-With-R.html>
- [6] Siddharth Raina. 2017. Regularized Regression - Housing Pricing. (Jan. 2017). <https://www.kaggle.com/sidraina89/regularized-regression-housing-pricing>
- [7] Kevin Wong. 2016. Predicting Ames House Prices. (Dec. 2016). <http://kevinfo.com/post/predicting-ames-house-prices/>
- [8] Ricky Yue and Jurgen De Jager. 2016. Advanced Regression Modeling on House Prices. (Sept. 2016). <https://nycdatascience.com/blog/student-works/advanced-regression-modeling-house-prices/>

```
bibtext report
```

```
=====
This is BibTeX, Version 0.99d (TeX Live 2016)
The top-level auxiliary file: report.aux
The style file: ACM-Reference-Format.bst
Database file #1: report.bib
Warning--empty year in 1
(There was 1 warning)
```

```
bibtext _ label error
```

```
=====
report.bib:59:@Misc{Aarshay_Jain2016,
```

```
bibtext space label error
```

```
=====
bibtext comma label error
```

```
=====
latex report
```

```
=====
[2017-12-04 12.23.29] pdflatex report.tex
This is pdfTeX, Version 3.14159265-2.6-1.40.17 (TeX Live 2016) (preloaded format=pdflatex)
Typesetting of "report.tex" completed in 1.0s.
./README.yml
41:1     error    trailing spaces  (trailing-spaces)
45:62    error    trailing spaces  (trailing-spaces)
46:67    error    trailing spaces  (trailing-spaces)
47:75    error    trailing spaces  (trailing-spaces)
48:69    error    trailing spaces  (trailing-spaces)
49:64    error    trailing spaces  (trailing-spaces)
53:26    error    trailing spaces  (trailing-spaces)
58:28    error    trailing spaces  (trailing-spaces)
```

```
=====
Compliance Report
=====
```

```
name: Cheruvu, Murali
hid: 306
paper1: 100%; 10/26/2017
paper2: 100%; 11/4/2017
project: 100%; 12/3/2017
```

```
yamlcheck
```

```
wordcount
```

```
9
wc 306 project 9 3869 report.tex
wc 306 project 9 4117 report.pdf
wc 306 project 9 197 report.bib
```

```
find "
```

```
49: train = pd.read_csv("../data/train.csv")
50: test = pd.read_csv("../data/test.csv")
130: sns_plot = sns.heatmap(corr, cmap="YlGnBu",
177: model = ols(formula = "SalePrice ~
178: GrLivArea + GarageArea", data=train)
181: criterion="cooks")
224: dummies = pd.get_dummies(df[col_name], prefix="_" + col_name)
258: train = pd.read_csv("train.csv")
259: test = pd.read_csv("test.csv")
260: target_vector = train["SalePrice"]
288: train = pd.read_csv("train.csv")
289: test = pd.read_csv("test.csv")
290: target_vector = train["SalePrice"]
```

```
326: train = pd.read_csv("train.csv")
327: test = pd.read_csv("test.csv")
328: target_vector = train["SalePrice"]
358: train = pd.read_csv("train.csv")
359: test = pd.read_csv("test.csv")
360: target_vector = train["SalePrice"]
393: df_predict = pd.DataFrame({'Id': test["Id"]},
446: train = pd.read_csv("train.csv")
447: test = pd.read_csv("test.csv")
448: target_vector = train["SalePrice"]

passed: False

find footnote
-----
passed: True

find input{format/i523}
-----
passed: False

find input{format/final}
-----
4: \input{format/final}

passed: True

floats
-----
59: \begin{figure}[H]
61: \includegraphics[width=0.9\columnwidth]{images/missing_values}
62: \caption{Variables with Missing Values} \label{fig:missing_values}
```

```

79: \begin{figure}[H]
81: \includegraphics[width=0.75\columnwidth]{images/num_features_1}
84: \begin{figure}[H]
86: \includegraphics[width=0.75\columnwidth]{images/num_features_2}
87: \caption{Sample Numerical Variables} \label{fig:num_features_2}
114: \begin{figure}[H]
116: \includegraphics[width=1.0\columnwidth]{images/cat_features_1}
117: \caption{Sample Category Variables} \label{fig:cat_features_1}
134: \begin{figure}[H]
136: \includegraphics[width=1.2\columnwidth]{images/correlations}
137: \caption{Correlations} \label{fig:correlations}
140: \begin{figure}[H]
142: \includegraphics[width=0.9\columnwidth]{images/pair_wise_correlations_1}
145: \begin{figure}[H]
147: \includegraphics[width=0.9\columnwidth]{images/pair_wise_correlations_2}
148: \caption{Pair-wise Correlations with Sale Price}
    \label{fig:pair_wise_correlations_2}
184: \begin{figure}[H]
186: \includegraphics[width=.95\columnwidth]{images/outliers}
187: \caption{Outliers - Cook's distance} \label{fig:outliers}
199: \begin{figure}[H]
201: \includegraphics[width=.95\columnwidth]{images/gr_liv_area_outlier}
202: \caption{Outliers - Garage Live Area}
    \label{fig:gr_liv_area_outlier}
205: \begin{figure}[H]
207: \includegraphics[width=.95\columnwidth]{images/garage_area_outlier}
208: \caption{Outliers - Garage Area} \label{fig:garage_area_outlier}
311: \begin{figure}[H]
313: \includegraphics[width=0.85\columnwidth]{images/random_forest_feature_ranking}
314: \caption{Random Forest Feature Ranking}
    \label{fig:random_feature_ranking}
403: \begin{figure}[H]
405: \includegraphics[width=0.8\columnwidth]{images/ridge_feature_ranking_pos}
406: \caption{Ridge Algorithm - Top 10 Positive Features}
    \label{fig:ridge_feature_ranking_pos}
409: \begin{figure}[H]
411: \includegraphics[width=0.8\columnwidth]{images/ridge_feature_ranking_neg}
412: \caption{Ridge Algorithm - Top 10 Negative Features}
    \label{fig:ridge_feature_ranking_neg}

```

```
419: \begin{table}[H]
422: \label{tab:xgb_param}
485: \begin{figure}[H]
487: \includegraphics[width=0.75\columnwidth]{images/xgboost_feature_i
    mportance}
488: \caption{XGBoost Feature Importance} \label{fig:xgb_feature_imp}
498: \begin{table}[H]
500: \label{tab:kaggle}
```

figures 14

tables 2

includegraphics 14

labels 14

refs 0

floats 16

True : ref check passed: (refs >= figures + tables)

True : label check passed: (refs >= figures + tables)

True : include graphics passed: (figures >= includegraphics)

False : check if all figures are referred to: (refs >= labels)

Label/ref check

passed: True

When using figures use columnwidth

[width=1.0\columnwidth]

do not change the number to a smaller fraction

find textwidth

passed: True

below_check

bibtex

label errors

59: Aarshay_Jain2016: do not use underscore in labels:

bibtex errors

```
This is BibTeX, Version 0.99d (TeX Live 2016)
The top-level auxiliary file: report.aux
The style file: ACM-Reference-Format.bst
Database file #1: report.bib
Warning--empty year in 1
(There was 1 warning)
```

bibtex_empty_fields

```
entries in general should not be empty in bibtex
```

find ""

```
passed: True
```

ascii

```
=====
The following tests are optional
=====
```

```
Tip: newlines can often be replaced just by an empty line
```

find newline

```
passed: True
```

```
cites should have a space before \cite{} but not before the {
```

find cite {

```
passed: True
```

IoT and Big Data Analytics for Equipment Predictive Health Management (PHM)

Ashok Reddy Singam
Indiana University
711 N Park Ave
Bloomington, Indiana 47408
asingam@iu.edu

Anil Ravi
Indiana University
711 N Park Ave
Bloomington, Indiana 47408
anilravi@iu.edu

ABSTRACT

The predictive health management (PHM) is an enabling discipline consisting of technologies and methods to assess the reliability of a product in its actual life cycle conditions to determine the advent of failure and mitigate system risk. The PHM system will monitor environmental, operational, and performance related characteristics of the product and gathered data analyzed to assess product health and predict remaining life.

In this application, the industrial rotating equipment such as compressors, vacuum blowers, pumps, and valves etc. are considered to monitor and analyze their operational behavior. The product critical operational parameter data such as vibration, temperature, and load current will be collected from field sensors and analyzed to predict the failure using kNN machine learning classification algorithms. The data will be collected from the field using wireless sensors and stored on the cloud based AWS database server. The product data will be analyzed and made available to all stakeholders to take appropriate preventive actions via web/mobile applications.

KEYWORDS

i523, HID333, HID337, KNN, IoT, Big Data, Analytics

1 INTRODUCTION

The PHM technology can be put within a broader business context by relating it to the Product-Service System (PSS) business model. PSS can be defined as an integrated combination of products and services where the emphasis is put on the ‘sale of use’ rather than the ‘sale of product’. Central to this new business model is a shift from selling a product, and its related spare parts as required, to selling a solution that supports customer needs in the form of a service delivering a fully maintained and useable product [4]. As shown in Figure 1, There are several wireless technologies such as 802.11, cellular, and short distance wireless protocols can be used to collect and send data to the centralized servers. Also, data can be stored in cloud based technologies such as AWS, Microsoft Azure, IBM and Google etc. for processing.

Problem Statement : in the manufacturing operations, automotive and other process industries rotating equipment such as pumps, valves, compressors, and blowers are commonly used equipment for various purposes. These equipment are severely suffered from wear and tear, bearing degradation, shaft misalignment, corrosion, and other mechanical breakdowns. Due to the limitations of wireless enabled sensors based data acquisition it was very difficult to collect this data in the past. Also, due to real-time nature of the data acquisition, it was a huge challenge to store the data

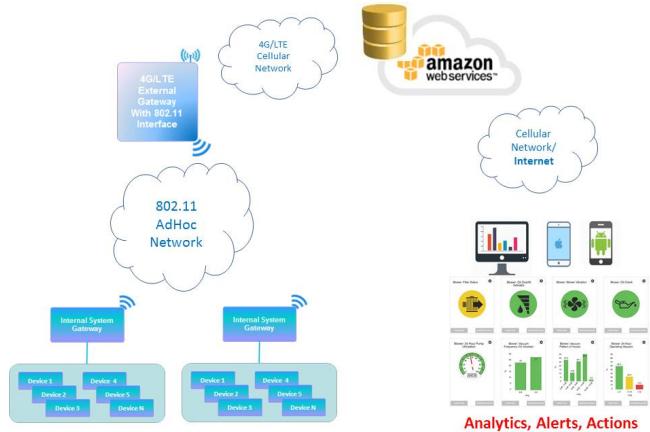


Figure 1: System Architecture

locally and process the information for applying machine learning algorithms. All these technological and infrastructure limitations caused industrial equipment health monitoring had become one of the sector businesses are losing the money due to operation shutdowns and unplanned maintenance etc.

Solution Approach : with the wireless sensors and cloud based server technologies, it has become possible to deploy hundreds of sensors in the manufacturing plant and collect the data and store with minimal costs. Once the data is stored on the servers with high computing power, machine learning algorithms can be used to process the sensor data to predict the equipment failures with reasonable accuracy. This approach has been named as predictive or prognostics health management of the equipment which is widely available in the recent times due to the availability of technological infrastructure.

The PHM generally combines sensing, collecting, storing and analyzing of environmental, operational, and performance related parameters to assess the health of a product and predict remaining useful life. Assessing the health of a product provides information that can be used to meet several critical goals: (1) providing advance warning of failures; (2) minimizing unscheduled maintenance, extending maintenance cycles, and maintaining effectiveness through timely repair actions; (3) reducing the life cycle cost of equipment by decreasing inspection costs, downtime, and inventory; and (4) improving qualification and assisting in the design and logistical support of fielded and future systems [3].

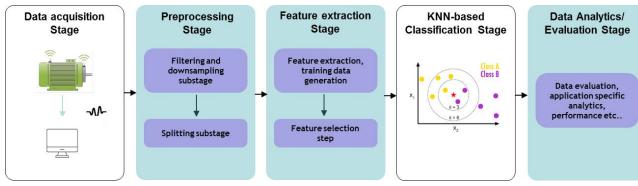


Figure 2: PHM Design Process

The PHM is not a new concept, however, with the advent of sensors, machine learning algorithms, and computing capacity of the servers it has become more prevalent in the recent days. In this application, an attempt has been made to prove the concept of simple PHM implementation and use in real world applications. The application can be re-architected to address more complex products/systems with considerations of scalability, performance, cost and reliability. The limitations of the current application are described in the end of this report.

The parameter monitoring and the analysis of acquired data using prognostic models are fundamental steps for the PHM methods. The sensors are the essential devices used to monitor parameters and obtain long-term accurate information to provide anomaly detection, fault isolation, and rapid failure prediction [3].

Firstly, PHM requires monitoring a large number of product parameters to evaluate the health of a product. Depending on the complexity of the monitored product, it is possible to monitor thousands of parameters in the entire life cycle of the product to provide the information required by PHM. These parameters include operational and environmental loads as well as the performance conditions of the product, for example, temperature, vibration, shock, pressure, acoustic levels, strain, stress, voltage, current, humidity levels, contaminant concentration, usage frequency, usage severity, usage time, power, and heat dissipation. In each case, a variety of monitoring features such as magnitude, variation, peak level, and rate of change may be required in order to obtain characteristics of parameters.

In this application, commonly used equipment in industrial and automobile operations such as air compressors, vacuum blowers, and smart valves are considered for analysis. The critical operational parameters of these products will be collected using applicable sensors from the field and fed to a database at regular intervals.

In general design, the frequency of data collection and storage depends on the number of parameters to be analyzed, cost of the system and operational behavior of the equipment. For this application, since products with rotating parts are considered, the critical parameters that would define the health of the equipment are: input or load current, internal ambient temperature, and vibration of the equipment.

The PHM application design process is shown in Figure 2, which describes various steps of the processes involved. For the implementation of this project, the sensor generated data is simulated using SQL scripts due to development time constraints. However, a detailed step-by-step approach is provided if we need to plug-in the sensor modules in to the application.

Data Acquisition Stage: It is required to have a description of a machine behaving normally that can be used for early detection

of anomalies. This calls for a proper characterization of machine health. As part of this process, various methods are identified to extract health information from vibration measurements and investigate strengths and weaknesses of these methods as health descriptors. This stage will be the core part of PHM application where vibration data were experimentally obtained from a compressor using triaxial accelerometer to collect transverse, longitudinal and vertical axes vibration signals. For the experimental data collection, ACC301A triaxial accelerometer and National Instruments data acquisition system was used. A total of 8 parameters (1) input current (2) input voltage (3) internal ambient temperature (4) external ambient temperature (5) transverse vibration (6) longitudinal vibration (7) vertical axis vibration and (8) acquisition time were captured at 1 second rate, which generated about 65000 records. This data has been analyzed for identifying the feature classification.

Pre Processing stage: during this stage collected data will be filtered and processed for accuracy in order to adapt them to subsequent feature extraction stage. In this application, all the pre-processing has been done manually to validate the accuracy of the data based on the system conditions. Since spectral analysis of vibration signals are not done (one of the limitations for this application, captured in the end), the data generated from the compressor is considered as the primary frequency of the equipment (which is isolated from the rest of the attachments).

Feature extraction and selection stage: during this stage domain specific vibration spectral analysis has been performed but only considered time-domain behavior for various system operational conditions such as increased load, modified input voltage, and modified external ambient temperature etc. Based on the response of the machine vibration to various external conditions were noted down. This data is used to identify the following feature vectors.

- NORMAL OPERATION AT 30 DEG CENTIGRADE
 - OVER CURRENT FAULT OPERATION
 - OVER TEMPERATURE FAULT OPERATION
 - INPUT OVER VOLTAGE FAULT OPERATION
 - ABNORMAL OPERATION AT 30 DEG CENTIGRADE
 - BEARING DEGRADATION OPERATION

kNN classification stage: this stage is core part of the PHM application, which will predict the unknown test data to be classified in to a known label based on the training data set using nearest neighbor algorithm.

Classifier performance evaluation stage: this stage will be used to evaluate the classifier accuracy of prediction. In this application, k-fold cross-validation method has been used to perform the evaluation.

The data is generated and made available in Oracle database on AWS cloud to perform analysis. The application developed in this project will consist of the following components:

- Sensor Data Generator
 - Machine Learning Algorithm
 - Big Data and IoT
 - PHM Dashboard
 - Decision Alerts
 - Application Script

The following sections will describe the architectural and design aspects of the PHM system implementation in detail.

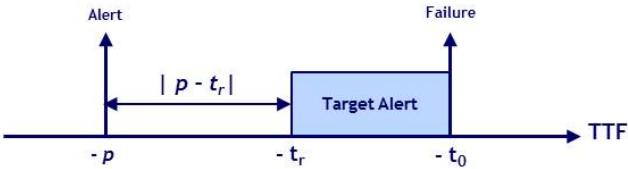


Figure 3: Time relation between alert time and failure time

2 PROGNOSTICS MODEL EVALUATION

[1] The prediction is typically performed only after the *health* of the component or system deteriorates beyond a certain threshold. In this application, faults and failures are identified in the training data set. The faults identified are: Over current fault, over temperature fault and over voltage fault. If over current fault is occurred, the equipment will tend to draw higher current than nominal values which if continued further several times eventually leads to a permanent failure of the equipment. In this application, when motor bearing starts degrading, the first observation will be over current followed by over temperature conditions. Often times, that threshold is tripped because a fault occurs. A fault is a state of a component or system that deviates from the normal state such that the integrity of the component is outside of its required specification. A fault does not necessarily imply that the overall system does not operate anymore; however, the damage that characterizes the fault often grows under the influence of operations to a failure. The latter is the state at which the component or system does not meet its desired function anymore. It is the task of prognostics to estimate the time that it takes from the current time to the failed state, conditional on anticipated future usage. This would give operators access to information that has significant implications on system safety or cost of operations. Where safety is impacted, the ability to predict failure allows operators to take action that preserves the assets either through rescue operation or through remedial action that avert failure altogether. Where minimizing cost of operations is the primary objective, predictive information allows operators to avert secondary damage, or to perform maintenance in the most cost-effective fashion. Often times, there is a mix of objectives that need to be optimized together, sometimes weighted by different preferences.

As emphasized above, predictive models evaluation needs to take domain specificities into account. Such specificities cover two aspects: capability of failure prediction and TTF estimation. From the point of view of TTF, it is desirable that a predictive model can generate alerts in a *targeted* time window prior to a failure. A model that predicts a failure too early leads to non-optimal component use [2] which will impact the reliability or availability of the system.

As shown in the Figure 3, the time to failure prediction will be estimated based on the classified result data set and alert the stakeholders to take relevant actions. The target alert zone will be identified based on the abnormal behavior of the equipment over the period.

3 APPLICATION DESIGN ANALYSIS

The PHM application in this project considered to use rotating equipment temperature, load current and vibration data for analyzing and predicting the future operational behavior. Vibration signals from rotating components are usually analyzed in the frequency domain, because significant peaks in the signal spectrum appear at frequencies that are related to the rotation frequency of the component. In this application, only time domain parameters with peak vibration magnitudes irrespective of the frequency component. The training data set consists of normal, abnormal, and fault conditions vibration patterns describes the system characteristics from which its status can be estimated. The PHM application for industrial equipment machine failure detection problem directly correlates to the pattern classification problem. From the vibration data collected, each accelerometer will output values of X, Y, and Z data then using a KNN we can similarly identify which vibration parameter(s) determines problems in our machines, or *likely to experience failure*. The typical defects or failures that can be detected are: machine imbalance, shaft misalignment, pumps cavitation, structural and rotating looseness, early stage bearing wear, gear teeth problems, and other high-frequency defects.

This application used *Sensor Data Gen* SQL script module to generate the sensor data and store in Oracle database on AWS. This is the critical module as we have not used the real data collection from the field. However, the sensor hardware and necessary environment to generate the data is identified and experimented to work with. A brief description about the hardware is provided in the end of this report.

The PHM application is designed such that the fundamental concepts can be verified to open a discussion on limitations, performance, scalability, ROI and reliability of the system.

The following sections describe the application design components with necessary implementation details:

3.1 Sensor Data Generator

The SQL data generator script is designed to generate training data as well test data for this application with following eleven parameters: Acquisition time, equipment name, part number, serial number, internal ambient temperature, external ambient temperature, input voltage, input current, and vibration data for x, y, and z axes. The following database design architecture followed for Sensor Data Gen module:

- Sensor Data Generator PL SQL Objects
 - Tables
 - * SENSOR TRAIN DATA for storing training data
 - * SENSOR TEST DATA for storing testing data
 - Views
 - * SENSOR TRAIN DATA VIEW
 - * SENSOR TEST DATA VIEW
 - Packages BIG DATA 503 PRJ PKG
 - * Generate Test Set
 - * Generate Train Set
 - * Delete Data Set
 - * Update Test Data Labels

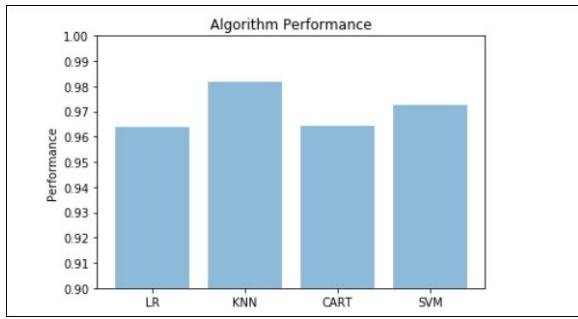


Figure 4: Algorithm performance

3.2 Machine Learning Algorithm

3.2.1 Classifier evaluation. Typical classifier evaluation methods include ROC Curves, Reject Curves, Precision-Recall Curves, and Statistical Tests. The statistical tests consists of following methods to perform evaluation:

- Estimating the error rate of a classifier
- Comparing two classifiers
- Estimating the error rate of a learning algorithm
- Comparing two algorithms

Out of the listed statistical tests, the error rate estimation method is used in this application to evaluate the performance. An experimental data is used to estimate the error rate or accuracy of various classifiers. Then a comparison has been made to choose the classifier to use in the application.

The following list of performance for various classifiers is observed during the accuracy calculation. The same set of training data has been used for all the classifiers, which has resulted the following performance.

- LogisticRregression: 0.963636 (0.044536)
- KNN: 0.981818 (0.036364)
- DecisionTreeClassifier: 0.964394 (0.059656)
- SVM: 0.972727 (0.041660)

Based on the performance shown, kNN has been selected to use for this application.

3.2.2 *k* Nearest Neighbor - *kNN*. [5] In this application, neighbors-based classification is chosen to classify the unknown instance to the known trained labels. Neighbors-based classification does not attempt to construct a general internal model, but simply stores instances of the training data.

Classification is computed from a simple majority vote of the nearest neighbors of each point: a query point is assigned the data class which has the most representatives within the nearest neighbors of the point.

KNN falls in the supervised learning family of algorithms. Informally, this means that we are given a labelled dataset consisting of training observations (x,y) and would like to capture the relationship between x and y . More formally, our goal is to learn a function $h:X \rightarrow Y$ so that given an unseen observations x , $h(x)$ can confidently predict the corresponding output y .

In the classification setting, the K-nearest neighbor algorithm essentially boils down to forming a majority vote between the K

most similar instances to a given unseen observation. Similarity is defined according to a distance metric between two data points. A popular choice is the Euclidean distance given by:

$$d(x, x') = \sqrt{(x_1 - x'_1)^2 + (x_2 - x'_2)^2 + \dots + (x_n - x'_n)^2}$$

But other measures can be more suitable for a given setting and include the Manhattan, Chebyshev and Hamming distance. An alternate way of understanding KNN is by thinking about it as calculating a decision boundary (i.e. boundaries for more than 2 classes) which is then used to classify new points.

In the application design, sci-kit open source python libraries are used for implementing the kNN algorithms. Scikit is built on NumPy, SciPy, and matplotlib. The k-neighbors classification in KNeighborsClassifier is the more commonly used of the two techniques. The optimal choice of the value k is highly data-dependent: in general a larger k suppresses the effects of noise, but makes the classification boundaries less distinct.

The sklearn.neighbors.KNeighborsClassifier class has the following methods, which are used in the application deisgn:

- fit: Fit the model using X as training data and y as target values.
- get params: Fit the model using X as training data and y as target values.
- kneighbors: Finds the K -neighbors of a point.
- kneighbors graph: Computes the (weighted) graph of k -Neighbors for points in X .
- predict: Predict the class labels for the provided data.
- predict_proba: Return probability estimates for the test data X .
- score: Returns the mean accuracy on the given test data and labels.
- set params: Set the parameters of this estimator.

3.2.3 *K-fold cross-validation*. To estimate the test error in the model, a cross-validation approach followed in which a subset of the training set will be holding out from the fitting process. This subset, called the validation set, can be used to select the appropriate level of flexibility of our algorithm. There are different validation approaches that are used in practice, and we will be exploring one of the more popular ones called **k-fold cross validation**. The k-fold cross validation (the k is totally unrelated to K) involves randomly dividing the training set into k groups, or folds, of approximately equal size. The first fold is treated as a validation set, and the method is fit on the remaining $k - 1$ folds. The misclassification rate is then computed on the observations in the held-out fold. This procedure is repeated k times; each time, a different group of observations is treated as a validation set. This process results in k estimates of the test error which are then averaged out.

In this application, an average k-fold cross validation accuracy of 0.99 achieved, which is explained in the appendix section of the report.

3.3 Big Data and IoT

In PHM systems big data is characterized by one or more 3Vs: volume, velocity and variety due to streaming of real-time IoT sensors.

Most of the IoT systems present challenges in combinations of velocity and volume. The important feature of the IoT application is that by observing the behavior of “many things” it will be possible to gain important insights, optimize processes, etc. This requires storing all the events (velocity and volume challenge) to run analytical queries over the stored events and perform analytics (data mining and machine learning) over the data to gain insights. In general PHM applications, data will be collected through field sensors at specific rate which accounts for large amount of data per day in the order of multi-million records. This data will be stored in any NOSQL or RDBMS based database for storage and processing. Since the big data infrastructure is much reliable and available widely from multiple vendors, it would help to build complex PHM systems with large number of feature vectors for classification.

In this application, for the demonstration of the concept, the real vibration data from the compressor equipment has been collected via accelerometer sensors. This vibration data has been analyzed in time domain and established the labels based on the compressor design performance parameters. Later, this data analysis is used to design a SQL script for generating training and test data sets. However, the real-time PHM system will have continuous streaming of data coming from hundreds of devices at faster rates (in the order of milliseconds to tens of seconds). This data needs to be captured by reliable and scalable platforms such as AWS IoT or similar and use the machine learning algorithms to classify the unknown data.

3.4 PHM Dashboard

Once all the test data set has been classified into appropriate labels, the prediction of the failure can be performed based on the trending of the equipment behavior over the period. In order to understand the equipment performance insight, following queries will be used on the classified data:

- Faults Reported by Equipment Part Number
- Faults Reported by Serial Number
- Abnormal Behavior by Equipment Part Number
- Abnormal Behavior by Serial Number over the period range

There can be more application specific information obtained from classified data set to take various decisions. Figure 4 shows various PHM data analytics for this project. Figure 4(a) displays all the serial numbers of equipment 1 with bearing degradation problems. The X axis gives serial numbers while the Y axis gives number of occurrences of bearing degradation for that particular serial number. Similarly Figure 4 (b) and 4 (c) give the details of over temperature and over current faults of various serial numbers.

As part of data visualization, result data file is queried based on the analytics metrics interested. The python matplotlib package has been used to draw the charts as needed for showing the analytics. In real world application, a more sophisticated business intelligence tools such as Tableau, Microsoft BI, and Amazon Quick Sight can be used to show the PHM dashboards. These dashboards are targeted for business users so that they will be able to customize the views, add filters and drill down into specific information as needed.

3.5 Decision Alerts

Once the results data set has been generated by the prediction algorithm, and then based on the analytics queries, PHM system

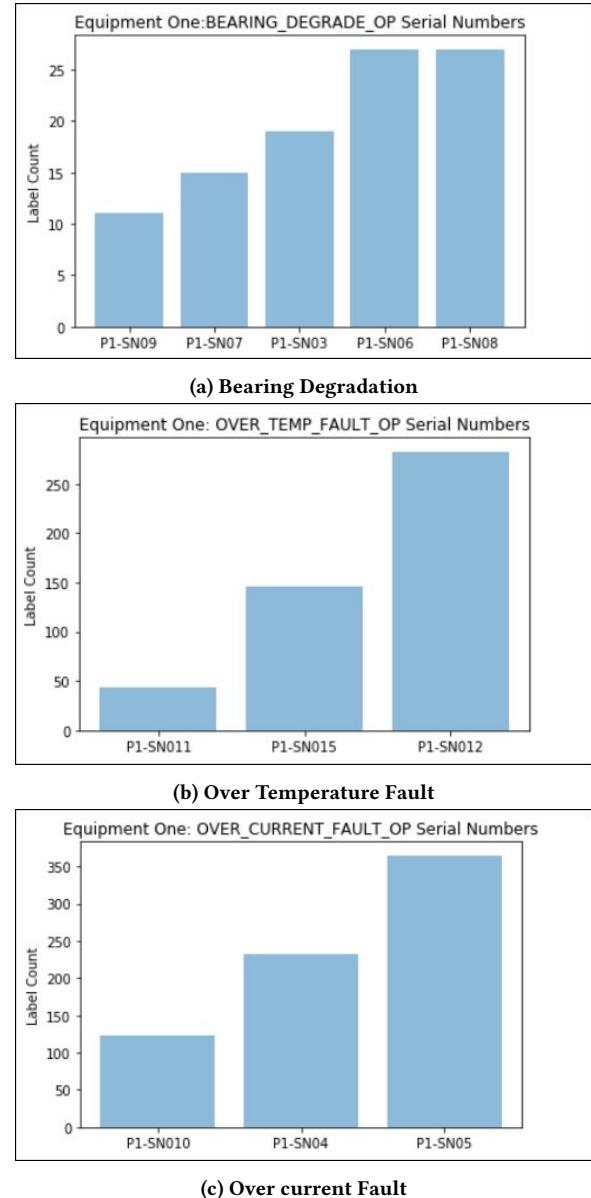


Figure 5: PHM Data Analytics

can send out the alert messages to appropriate stakeholders. The typical messages include the following s minimum:

- SN 10002: Faulted X times on over temperature in last Y days, needs maintenance to clean the filter
- SN 10005: Consistently indicating bearing degradation from last X days, needs lubrication maintenance
- SN 10009: Consistently drawing over current from last X days, needs mechanical load maintenance

In this application all the unseen test data is classified and labeled in the result data file. However, in real world application, along with the dashboards a comprehensive alerting capabilities can be built. The application checks for out of range alert conditions on

selected incoming report parameters, looking for warning or alarm conditions that are higher or lower than expected under normal operating conditions.

3.6 Application Script

All the application specific python modules are captured in a ipynb file which can be run from Jupyter notebook or Github. The location of the file is described in appendix b. This application specific data has been created/generated on AWS cloud database, which will be accessed during the run time by python notebook.

4 APPLICATION LIMITATIONS

The PHM application developed in this project has several limitations. Typically, PHM applications suffer from the prediction accuracy rate which influences the ability to take decisions that will have broader impact on the business operations and financial aspects. However, with the advanced machine learning classifiers and model evaluation methods this can be addressed to achieve reasonable confidence. Following are some of the limitations of this application, which can be addressed and improved in large real-time PHM systems.

Data acquisition hardware : in this application the data is not collected from real-time sensors for voltage, current, temperature and vibration data. There will be inherent accuracy in the raw data generated by SQL script. However, a sample vibration dataset has been collected from the field, which is used as basis to generate the simulated data set.

Feature extraction analysis : the equipment performance parameters of interest need to be down selected from large set of incoming parameter data.

When analyzing vibration data in the time domain only few parameters are available in quantifying the strength of a vibration profile: amplitude, peak-to-peak value, and RMS. The amplitude is valuable for shock events but it does not take into account the time duration and thus the energy in the event. The same is true for peak-to-peak with the added benefit of providing the maximum excursion of the wave, useful when looking at displacement information, specifically clearances. The RMS value is generally the most useful because it is directly related to the energy content of the vibration profile and thus the destructive capability of the vibration.

This requires in-depth domain specific analysis, in this case a detailed mathematical modeling of vibration spectral analysis to precisely select the features and corresponding behavior patterns. Such analytical data should be used for training data feature set. In this application, a primitive approach of time-domain analysis of vibration magnitudes used for determining features. However, in real application these features need to be mathematically analyzed to identify the features that represent the system behavior as close as possible.

Model accuracy and scoring : the kNN algorithm used in this application validated using k-cross fold cross-validation. There are several other model evaluation and scoring methods such as accuracy (or error rate), True Positive Rate (TPR), False Positive Rate (FPR), False Negative Rate (FNR), True Negative Rate (TNR), sensitivity etc. These metrics provide a simple and effective way to measure the performance of a classifier. This application can

be further improved by applying more performance measurement methods to increase the effectiveness of the algorithm design.

Scalability : the application designed in this project is very primitive to understand the basic concepts of PHM and kNN classifier implementation. This application cannot be used for PHM application in business use. To implement real world PHM application, a more comprehensive design needed by considering modularity, service oriented architecture, large number of sensors integration, big data and analytics integration etc.

5 RECOMMENDATIONS

Generality : since each rotating equipment vibrates in a different manner, a monitoring method needs to be retrained for each machine. The training on several repeated measurements on several similar equipment in several operating modes may allow for a more general monitoring method.

Feature extraction and dimensionality : In this application it has been assumed that a proper feature (selection) has been chosen, such that the feature dimensionality is not too high. If the data lies in a subspace, application of an initial dimensionality reduction may be a good idea. It is highly recommended to perform spectral analysis on vibration data and identify various fault frequencies and their sources. This would help to extract the optimized feature vectors for the given application followed by selecting the more relevant ones.

Model evaluation : classifier accuracy and effectiveness will be varied based on the test data set. It is highly recommended to evaluate multiple models with appropriate test data to choose the best classifier for the given application.

Domain specific modeling : It is highly recommended to perform more and more domain-oriented feature vector analysis to meet the needs of predictive model evaluation for PHM applications. Domain-oriented approaches helpful and useful in evaluating classifier for applications. Generic evaluation methods could help developers in investigating overall performance of a model from the statistical viewpoint at the initial stage of model development. Domain-oriented approaches should be further used to evaluate the usefulness and business value [2].

6 CONCLUSION

In this project, the problem statement around industrial rotating equipment maintenance is described and solution principle to address the same using PHM concept is defined, experimented and results are discussed. Since this application is developed to prove the only concept but not the complete solution a section with limitations and recommendations for real world system development is described. Overall, PHM application with kNN classifier algorithm and cross validation accuracy of 0.99 has been implemented, verified and results are analyzed for business decisions.

ACKNOWLEDGMENTS

The authors would like to thank professor Gregor von Laszewski and his team for providing *LaTex* templates and assistance with the *JabRef* tool to organize references.

REFERENCES

- [1] Bhaskar Saha Sankalita Saha Abhinav Saxena, Jose Celaya and Kai Goebel. 2009. Sensor Systems for Prognostics and Health Management. (2009), 16 pages. <https://ntrs.nasa.gov/archive/nasa/casi.ntrs.nasa.gov/20100023445.pdf>
- [2] Jie Liu Kyle R Mulligan Chunsheng Yang, Yanni Zou. 2014. Predictive Model Evaluation for PHM. (2014), 11 pages. [https://www.phmsociety.org/sites/phm_submission/2014/ijphm.14.019.pdf](https://www.phmsociety.org/sites/phmsociety.org/files/phm_submission/2014/ijphm.14.019.pdf)
- [3] Michael H. Azarian Shunfeng Cheng and Michael G. Pecht. 2010. Sensor Systems for Prognostics and Health Management. (2010), 24 pages. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3247731/pdf/sensors-10-05774.pdf>
- [4] Ian Jennions Tonci Grubic and Tim Baines. 2009. The Interaction of PSS and PHM - a mutual benefit case. (2009), 10 pages. [https://www.phmsociety.org/sites/phm_submission/2009/phmc_09_49.pdf](https://www.phmsociety.org/sites/phmsociety.org/files/phm_submission/2009/phmc_09_49.pdf)
- [5] Kevin Zakka. 2016. A Complete Guide to K-Nearest-Neighbors with Applications in Python and R. (2016). <https://kevinzakka.github.io/2016/07/13/k-nearest-neighbor/>

A WORK BREAKDOWN

A.1 HID 333:Anil Ravi

- Identified Project topic.
- Created architecture of the application.
- Ran experimental test to collect vibration data
- Extracted and analyzed feature vectors
- Studied, designed and reviewed kNN algorithm
- Created draft project report
- Reviewed the draft project report.

A.2 HID 337:Ashok Reddy Singam

- Implemented sensor data generation SQL script.
- Implemented kNN algorithm in Python
- Implemented k-fold cross validation design
- Created data analytics charts
- Reviewed the draft project report.

B CODE REFERENCE

All code, notebooks and files for this project can be found in the github repository: <https://github.com/bigdata-i523/hid337/blob/master/project/jupyter>

bibtext report

=====

This is BibTeX, Version 0.99d (TeX Live 2016)
The top-level auxiliary file: report.aux
The style file: ACM-Reference-Format.bst
Database file #1: report.bib
Name 1 in "Tonci Grubic, Ian Jennions, and Tim Baines" has a comma at the end for entry
while executing---line 3085 of file ACM-Reference-Format.bst
Name 1 in "Tonci Grubic, Ian Jennions, and Tim Baines" has a comma at the end for entry
while executing---line 3085 of file ACM-Reference-Format.bst
Name 1 in "Tonci Grubic, Ian Jennions, and Tim Baines" has a comma at the end for entry
while executing---line 3085 of file ACM-Reference-Format.bst
Name 1 in "Abhinav Saxena, Jose Celaya, Bhaskar Saha, Sankalita Saha, and Kai Goebel" has
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Abhinav Saxena, Jose Celaya, Bhaskar Saha, Sankalita Saha,
while executing---line 3085 of file ACM-Reference-Format.bst
Name 1 in "Abhinav Saxena, Jose Celaya, Bhaskar Saha, Sankalita Saha, and Kai Goebel" has
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Abhinav Saxena, Jose Celaya, Bhaskar Saha, Sankalita Saha,
while executing---line 3085 of file ACM-Reference-Format.bst
Name 1 in "Abhinav Saxena, Jose Celaya, Bhaskar Saha, Sankalita Saha, and Kai Goebel" has
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Abhinav Saxena, Jose Celaya, Bhaskar Saha, Sankalita Saha,
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Chunsheng Yang, Yanni Zou, Jie Liu, Kyle R Mulligan" for e
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Chunsheng Yang, Yanni Zou, Jie Liu, Kyle R Mulligan" for e
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Chunsheng Yang, Yanni Zou, Jie Liu, Kyle R Mulligan" for e
while executing---line 3085 of file ACM-Reference-Format.bst
Name 1 in "Abhinav Saxena, Jose Celaya, Bhaskar Saha, Sankalita Saha, and Kai Goebel" has
while executing---line 3131 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Abhinav Saxena, Jose Celaya, Bhaskar Saha, Sankalita Saha,
while executing---line 3131 of file ACM-Reference-Format.bst
Name 1 in "Abhinav Saxena, Jose Celaya, Bhaskar Saha, Sankalita Saha, and Kai Goebel" has
while executing---line 3131 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Abhinav Saxena, Jose Celaya, Bhaskar Saha, Sankalita Saha,
while executing---line 3131 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Chunsheng Yang, Yanni Zou, Jie Liu, Kyle R Mulligan" for e
while executing---line 3131 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Chunsheng Yang, Yanni Zou, Jie Liu, Kyle R Mulligan" for e
while executing---line 3131 of file ACM-Reference-Format.bst
Name 1 in "Tonci Grubic, Ian Jennions, and Tim Baines" has a comma at the end for entry
while executing---line 3131 of file ACM-Reference-Format.bst

```
Name 1 in "Tonci Grubic, Ian Jennions, and Tim Baines" has a comma at the end for entry
while executing---line 3131 of file ACM-Reference-Format.bst
Name 1 in "Abhinav Saxena, Jose Celaya, Bhaskar Saha,Sankalita Saha, and Kai Goebel" has
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Abhinav Saxena, Jose Celaya, Bhaskar Saha,Sankalita Saha,
while executing---line 3229 of file ACM-Reference-Format.bst
Name 1 in "Abhinav Saxena, Jose Celaya, Bhaskar Saha,Sankalita Saha, and Kai Goebel" has
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Abhinav Saxena, Jose Celaya, Bhaskar Saha,Sankalita Saha,
while executing---line 3229 of file ACM-Reference-Format.bst
Name 1 in "Abhinav Saxena, Jose Celaya, Bhaskar Saha,Sankalita Saha, and Kai Goebel" has
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Abhinav Saxena, Jose Celaya, Bhaskar Saha,Sankalita Saha,
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Chunsheng Yang, Yanni Zou, Jie Liu, Kyle R Mulligan" for e
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Chunsheng Yang, Yanni Zou, Jie Liu, Kyle R Mulligan" for e
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Chunsheng Yang, Yanni Zou, Jie Liu, Kyle R Mulligan" for e
while executing---line 3229 of file ACM-Reference-Format.bst
Name 1 in "Tonci Grubic, Ian Jennions, and Tim Baines" has a comma at the end for entry
while executing---line 3229 of file ACM-Reference-Format.bst
Name 1 in "Tonci Grubic, Ian Jennions, and Tim Baines" has a comma at the end for entry
while executing---line 3229 of file ACM-Reference-Format.bst
Name 1 in "Tonci Grubic, Ian Jennions, and Tim Baines" has a comma at the end for entry
while executing---line 3229 of file ACM-Reference-Format.bst
(There were 32 error messages)
make[2]: *** [bibtex] Error 2
```

```
latex report
=====
```

```
[2017-12-04 12.24.14] pdflatex report.tex
This is pdfTeX, Version 3.14159265-2.6-1.40.17 (TeX Live 2016) (preloaded format=pdflatex)
Typesetting of "report.tex" completed in 1.0s.
```

```
=====
Compliance Report
=====
```

```
name: Ashok Reddy Singam
hid: 337
paper1: Nov 01 17 100%
paper2: Nov 06 17 100%
project: Dec 04 17 100%
```

```
yamlcheck
```

```
wordcount
```

```
7
```

```
wc 337 project 7 4977 report.tex
wc 337 project 7 4832 report.pdf
wc 337 project 7 217 report.bib
```

```
find "
```

```
passed: True
```

```
find footnote
```

```
passed: True
```

```
find input{format/i523}
```

```
passed: False
```

```
find input{format/final}
```

```
4: \input{format/final}
```

```
passed: True
```

```
floats
```

```
38: \begin{figure}
39: \includegraphics[width=1.0\columnwidth]{images/system_architecture
    }
40: \caption{System Architecture} \label{fig:Figure1}
64: \begin{figure}
65: \includegraphics[width=1.0\columnwidth]{images/phm_process_1}
66: \caption{PHM Design Process} \label{fig:Figure2}
109: \begin{figure}
```

```

110: \includegraphics[width=1.0\columnwidth]{images/ttf_1}
111: \caption{Time relation between alert time and failure time}
    \label{fig:Figure3}
174: \begin{figure}
175: \includegraphics[width=0.9\columnwidth]{images/algperformance}
176: \caption{Algorithm performance} \label{fig:Figure5}
238: \begin{figure}
241: \includegraphics[width=1.0\columnwidth]{images/DEGRADE}
242: \caption{Bearing Degradation} \label{sfig:sfig4a}
246: \includegraphics[width=1.0\columnwidth]{images/OVERTEMP}
247: \caption{Over Temperature Fault} \label{sfig:sfig4b}
251: \includegraphics[width=1.0\columnwidth]{images/OVERCURR}
252: \caption{Over current Fault} \label{sfig:sfig4c}
255: \caption{PHM Data Analytics} \label{fig:Figure4}

```

figures 5

tables 0

includegraphics 7

labels 8

refs 0

floats 5

True : ref check passed: (refs >= figures + tables)

False : label check passed: (refs >= figures + tables)

False : include graphics passed: (figures >= includegraphics)

False : check if all figures are referred to: (refs >= labels)

Label/ref check

- 34: The PHM technology can be put within a broader business context by relating it to the Product-Service System (PSS) business model. PSS can be defined as an integrated combination of products and services where the emphasis is put on the \lq sale of use \rq rather than the \lq sale of product \rq. Central to this new business model is a shift from selling a product, and its related spare parts as required, to selling a solution that supports customer needs in the form of a service delivering a fully maintained and useable product \cite{Tonci2009}. As shown in Figure 1, There are several wireless technologies such as 802.11, cellular, and short distance wireless protocols can be used to collect and send data to the centralized servers. Also, data can be stored in cloud based technologies such as AWS, Microsoft Azure, IBM and Google etc. for processing.
- 59: The PHM application design process is shown in Figure 2, which describes various steps of the processes involved. For the implementation of this project, the sensor generated data is simulated using SQL scripts due to development time constraints.

However, a detailed step-by-step approach is provided if we need to plug-in the sensor modules in to the application.

- 113: As shown in the Figure 3, the time to failure prediction will be estimated based on the classified result data set and alert the stakeholders to take relevant actions. The target alert zone will be identified based on the abnormal behavior of the equipment over the period.
- 233: There can be more application specific information obtained from classified data set to take various decisions. Figure 4 shows various PHM data analytics for this project. Figure 4(a) displays all the serial numbers of equipment 1 with bearing degradation problems. The X axis gives serial numbers while the Y axis gives number of occurrences of bearing degradation for that particular serial number. Similarly Figure 4 (b) and 4 (c) give the details of over temperature and over current faults of various serial numbers.

passed: False -> labels or refs used wrong

When using figures use columnwidth
[width=1.0\columnwidth]
do not change the number to a smaller fraction

find textwidth

passed: True

below_check

bibtex

label errors

bibtex errors

This is BibTeX, Version 0.99d (TeX Live 2016)

The top-level auxiliary file: report.aux

The style file: ACM-Reference-Format.bst

Database file #1: report.bib

Name 1 in "Tonci Grubic, Ian Jennions, and Tim Baines" has a comma at the end for entry

while executing---line 3085 of file ACM-Reference-Format.bst
Name 1 in "Tonci Grubic, Ian Jennions, and Tim Baines" has a comma at the end for entry
while executing---line 3085 of file ACM-Reference-Format.bst
Name 1 in "Tonci Grubic, Ian Jennions, and Tim Baines" has a comma at the end for entry
while executing---line 3085 of file ACM-Reference-Format.bst
Name 1 in "Abhinav Saxena, Jose Celaya, Bhaskar Saha,Sankalita Saha, and Kai Goebel" has
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Abhinav Saxena, Jose Celaya, Bhaskar Saha,Sankalita Saha,
while executing---line 3085 of file ACM-Reference-Format.bst
Name 1 in "Abhinav Saxena, Jose Celaya, Bhaskar Saha,Sankalita Saha, and Kai Goebel" has
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Abhinav Saxena, Jose Celaya, Bhaskar Saha,Sankalita Saha,
while executing---line 3085 of file ACM-Reference-Format.bst
Name 1 in "Abhinav Saxena, Jose Celaya, Bhaskar Saha,Sankalita Saha, and Kai Goebel" has
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Abhinav Saxena, Jose Celaya, Bhaskar Saha,Sankalita Saha,
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Chunsheng Yang, Yanni Zou, Jie Liu, Kyle R Mulligan" for e
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Chunsheng Yang, Yanni Zou, Jie Liu, Kyle R Mulligan" for e
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Chunsheng Yang, Yanni Zou, Jie Liu, Kyle R Mulligan" for e
while executing---line 3085 of file ACM-Reference-Format.bst
Name 1 in "Abhinav Saxena, Jose Celaya, Bhaskar Saha,Sankalita Saha, and Kai Goebel" has
while executing---line 3131 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Abhinav Saxena, Jose Celaya, Bhaskar Saha,Sankalita Saha,
while executing---line 3131 of file ACM-Reference-Format.bst
Name 1 in "Abhinav Saxena, Jose Celaya, Bhaskar Saha,Sankalita Saha, and Kai Goebel" has
while executing---line 3131 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Abhinav Saxena, Jose Celaya, Bhaskar Saha,Sankalita Saha,
while executing---line 3131 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Chunsheng Yang, Yanni Zou, Jie Liu, Kyle R Mulligan" for e
while executing---line 3131 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Chunsheng Yang, Yanni Zou, Jie Liu, Kyle R Mulligan" for e
while executing---line 3131 of file ACM-Reference-Format.bst
Name 1 in "Tonci Grubic, Ian Jennions, and Tim Baines" has a comma at the end for entry
while executing---line 3131 of file ACM-Reference-Format.bst
Name 1 in "Tonci Grubic, Ian Jennions, and Tim Baines" has a comma at the end for entry
while executing---line 3131 of file ACM-Reference-Format.bst
Name 1 in "Abhinav Saxena, Jose Celaya, Bhaskar Saha,Sankalita Saha, and Kai Goebel" has
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Abhinav Saxena, Jose Celaya, Bhaskar Saha,Sankalita Saha,
while executing---line 3229 of file ACM-Reference-Format.bst
Name 1 in "Abhinav Saxena, Jose Celaya, Bhaskar Saha,Sankalita Saha, and Kai Goebel" has
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Abhinav Saxena, Jose Celaya, Bhaskar Saha,Sankalita Saha,

```
while executing---line 3229 of file ACM-Reference-Format.bst
Name 1 in "Abhinav Saxena, Jose Celaya, Bhaskar Saha,Sankalita Saha, and Kai Goebel" has
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Abhinav Saxena, Jose Celaya, Bhaskar Saha,Sankalita Saha,
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Chunsheng Yang, Yanni Zou, Jie Liu, Kyle R Mulligan" for e
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Chunsheng Yang, Yanni Zou, Jie Liu, Kyle R Mulligan" for e
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Chunsheng Yang, Yanni Zou, Jie Liu, Kyle R Mulligan" for e
while executing---line 3229 of file ACM-Reference-Format.bst
Name 1 in "Tonci Grubic, Ian Jennions, and Tim Baines" has a comma at the end for entry
while executing---line 3229 of file ACM-Reference-Format.bst
Name 1 in "Tonci Grubic, Ian Jennions, and Tim Baines" has a comma at the end for entry
while executing---line 3229 of file ACM-Reference-Format.bst
Name 1 in "Tonci Grubic, Ian Jennions, and Tim Baines" has a comma at the end for entry
while executing---line 3229 of file ACM-Reference-Format.bst
(There were 32 error messages)
```

bibtex_empty_fields

entries in general should not be empty in bibtex

find ""

passed: True

ascii

=====

The following tests are optional

=====

Tip: newlines can often be replaced just by an empty line

find newline

passed: True

cites should have a space before \cite{} but not before the {

```
find cite {
```

```
-----  
passed: True
```

IoT Application Using MQTT and Raspberry Pi Robot Car

Arnav Arnav

Indiana University Bloomington
Bloomington, Indiana 47408, USA
aarnav@iu.edu

ABSTRACT

As the number of connected edge devices increases there is a need for fast communication between these devices, and to analyse the data collected by these devices, which is made possible by the use of a scalable lightweight communication protocol such as MQTT, which is easy to use, data agnostic, and application independent. We look at one such application of the protocol, to control a robot car remotely, over wireless network, navigating with the help of a raspberry pi camera on the car.

KEYWORDS

i523, HID201, Edge Computing, Raspberry Pi, MQTT, Robot Car, IoT

1 INTRODUCTION

As the number of edge devices increases, and sensor networks become more and more common in Internet of Things (IoT) applications, the need arises to allow these resource constrained devices to communicate with each other in a power efficient and secure manner. In many cases these devices may not be able to process traditional HTTP requests efficiently, and as the number of devices increases, sending an HTTP request to each of the devices in order to get data may not be efficient [2][6].

Message Queue Telemetry Transport (MQTT) is a lightweight machine to machine (M2M) messaging protocol, that uses a client/server based publish/subscribe model and is ideal for IoT applications. The protocol has been designed on top of TCP/IP protocol for us in situations where network bandwidth and available memory are limited [25][15]. The Eclipse Paho Project currently provides support for MQTT [4]. MQTT clients are available for various languages like Python, C, and Lua.

We look at one such application here that uses MQTT for communication between a raspberry pi and a desktop. The raspberry pi controls the stepper motors of the robot car according to the message it receives over mqtt, and drives the car accordingly. Another program running on the raspberry pi uses the raspberry pi onboard camera to capture pictures and send them back to the desktop to help in navigation. Thus we create a simple robot car that can be used remotely for monitoring purposes. The robot car can be controlled from anywhere in the world, as long as both the controlling device (desktop) and the raspberry pi can connect to the MQTT broker.

We can use multiple such cars and controlling devices to control the cars independently or from a common device to drive multiple cars together, thus controlling a swarm of cars. As these cars may be using different platforms like raspberry pi or arduino, Using MQTT allows us to write the controller program independent of the subscriber programs running on the different robot cars and

even in different languages. All that is needed to control a car is that the subscriber can understand the messages sent by the controller.

2 RELATED WORK

There have been many edge computing applications that involve robot cars or swarm of cars.

[19] provides an example of a raspberry pi car that uses distance sensor, and face detection on the raspberry pi 2. The car is controlled over wifi and is built using the GoPiGo robot car kit [8]

Zheng Wang used raspberry pi in [24] to build a sophisticated self driving car that can detect stop signs and traffic signals and drive appropriately on a small test track. The car has a camera and a distance sensor that stream data to a TCP server running on a desktop. The system uses Haar Cascades provided in opencv to detect objects like stop signs and traffic signals and a trained neural network which uses the image to predict the direction in which the car should move. The distance is calculated using the image from the raspberry pi camera with the help of a monocular vision method proposed by Chu, Ji, Guo, Li and Wang in 2004 [9].

As the part of Eclipse IoT open challenge [1] built a robot car that is controlled using the Constrained Application Protocol (CoAP) which snaps images and communicates the images over MQTT

OpenHAB provides a vendor neutral platform that allows users to integrate various home automation systems and provides an application interface to control those devices [16]. It allows integration of various devices with MQTT.

The FloodNet project at University of Southampton [7] aims at "providing a pervasive, continuous embedded monitoring presence". The system is intelligent and obtains "environmental self-awareness and resilience to ensure robust transmission of data", ensuring data quality and allowing exploration of environments in new ways. The project uses MQTT for communicating data from the sensors on field to visualization and simulation applications.

As a part of IBM's Extreme Blue projects, Say it Sign it [23] is a sophisticated, innovative speech to sign language translation system. The application uses speech recognition and renders an avatar that signs the corresponding words in British sign Language, using MQTT and microbroker for communication.

3 TECHNOLOGIES AND HARDWARE

The project uses MQTT to communicate between a controller running on a desktop and a raspberry pi that drives the robot car with the help of stepper motors. We describe these technologies in detail.

3.1 MQTT

MQTT works via a publish-subscribe model that contains 3 entities: (1) a publisher, that sends a message, (2) a broker, that maintains

queue of all messages based on topics and (3) multiple subscribers that subscribe to various topics they are interested in [20].

This allows for decoupling of functionality at various levels. The publisher and subscriber do not need to be close to each other and do not need to know each others identity. They need only to know the broker, as the publisher and the subscribers do not have to be running either at the same time nor on the same hardware [12].

MQTT implements a hierarchy of topics that are related to all messages. These topics are recognised by strings separated by a forward-slash (/), where each part represents a different topic level. This is a common model introduced in file systems but also in internet URLs.

A topic looks therefore as follows: *topic-level0/topic-level1/topic-level2*.

All subscribers subscribe to different topics via the broker. Subscribing to *topic-level0* allows the subscriber to receive all messages that are associated with topics that start with *topic-level0*.

This is different from traditional message queues as the message is forwarded to multiple subscribers, and allows for a more flexible approach with the help of topics [12]. The basic steps in an MQTT client application include connecting to the broker, subscribing to some topics, waiting for messages and performing the appropriate action when a certain message is received [25].

MQTT allows the publisher and subscriber to respond to messages with the help of callbacks that are executed on different events, in a non-blocking manner. The paho-mqtt package for python provides callbacks methods like `on-connect()`, `on-message()` and `on-disconnect()`, which are fired when the connection to the broker is complete, a message is received from the broker, and when the client is disconnected from the broker respectively. These methods are used in conjunction with the `loop-start()` and `loop-end()` methods which start and end an asynchronous loop that listens for these events and fires the relevant callbacks, allowing the clients to perform other tasks [5].

MQTT has been designed to be flexible and options are provided to easily change the quality of service (QoS) as required by the application. Three basic levels of QoS are supported by the protocol, Atmost-once (QoS level 0), Atleast-once (QoS level 1) and Atmost-once (QoS level 2) [13][5].

The QoS level of 0 can be used in applications where some dropped messages may not affect the application. Under this QoS level, the broker forwards a message to the subscribers only once and does not wait for any acknowledgement [13] [5].

The QoS level of 1 can be used in situations where the delivery of all messages is important and the subscriber can handle duplicate messages. Here the broker keeps on resending the message to a subscriber after a certain timeout until the first acknowledgement is received. A QoS level of 2 should be used in cases where all messages must be delivered and no duplicate messages should be allowed. In this case the broker sets up a handshake with the subscriber to check for its availability before sending the message [13] [5].

The MQTT specification uses TCP/IP to deliver the messaged to the subscribers, but it does not provide any form of security by default to make it useful for resource constrained IoT devices. “It allows the use of username and password for authentication,

but by default this information is sent as plain text over the network, making it susceptible to man-in-the middle attacks” [18] [14]. Therefore, in sensitive applications some form of additional security measures are recommended which may include network layer security with the use of Virtual Private Networks (VPNs), Transport Layer Security, or application layer security [14].

Transport Layer Security (TLS) and Secure Sockets Layer (SSL) are cryptographic protocols that establish the identity of the server and client with the help of a handshake mechanism which uses trust certificates to establish identities before encrypted communication can take place [3]. If the handshake is not completed for some reason, the connection is not established and no messages are exchanged [14]. “Most MQTT brokers provide an option to use TLS instead of plain TCP and port 8883 has been standardized for secured MQTT connections” [18].

Using TLS/SSL security however comes at an additional cost. If the connections are short-lived then most of the time can be spent in the handshake itself, which may take up few kilobytes of bandwidth. In case the connections are short-lived, temporary session IDs and session tickets can be used to resume a session instead of repeating the handshake process. If the connections are long term, the overhead of the handshake is negligible and TLS/SSL security should be used [18][14].

3.2 Raspberry Pi

The raspberry pi is a credit card sized development board that was developed by Eben Upton with the goal to create a low cost device that can be used for education and prototyping [17]. Since its creation the board has been adapted for various different projects by educators hobbyists and in the industry [22]. The board is developed as open hardware except for the Broadcom chip that controls the main components of the board, and most raspberry pi projects are available openly with detailed documentation.

The board’s Broadcom system on chip consists of an ARM processor and it can be used just like a normal computer by connecting a monitor, a keyboard and a mouse. The raspberry pi can communicate to other devices with the help of wifi and bluetooth and is capable of accessing the internet. All this put together makes the raspberry pi a very useful device [22].

The raspberry pi comes in various models, Model A+, which is one of the smallest form factors, raspberry pi2 Model B, raspberry pi3 Model B and Model B+ that have more gpio pins. The raspberry pi 3 Model B is the newest design and consists of on board wifi and bluetooth, eliminating the need to use usb wifi and bluetooth attachments. It has a 1.2 GHz ARM 8 microprocessor, 1 GB RAM, a dual core Videocore IV GPU, and 40 general purpose input and output (GPIO) pins. The board has an ethernet port and four USB ports and an HDMI port to connect to a monitor [11][10].

The raspberry pi Zero is the development board that has the smallest form factor. Even though the raspberry pi zero includes no ethernet or USB ports, and does not come with GPIO pins soldered on, its small size and cost effectiveness make it extremely useful in applications such as IoT where space is constrained [21].

The raspberry pi uses a micro SD card to boot and various operating systems, that support the ARM architecture can be used. The most common operating systems are Raspbian, a derivative

of the Debian linux, and Pidora, a derivative of Fedora. There are other operating systems centered around using the raspberry pi for various purposes, like openELEC and RaspBMC, which make it easy to use raspberry pi as a multimedia center. For users who want non-linux operating system, RISC OS may be a good choice. The raspberry pi foundation provides new users the opportunity to try out various operating systems with the help of their New Out Of The Box Software (NOOBS), which allows the users to pick which operating system they want to use [17].

3.3 Stepper Motors

Stepper motors are brushless motors that divide the complete rotation into a number of parts known as steps. The motor consists of electromagnetic coils and a rotating core that aligns itself according to the combined magnetic effect of the coils. The stepper motor can move from one step to another and remain in a single step based on which coils are turned on. The torque of the motor can be increased or decreased with the current supplied to the coils, and the speed of rotation can be controlled by setting the time interval between switching the coils on and off [26].

3.4 OpenCV

4 ARCHITECTURE

5 RESULTS

5.1 Setup Instructions

To run the application successfully on both the raspberry pi and the desktop, it must be ensured that all the required libraries are installed. A Makefile has been provided that can do this on both the raspberry pi and the desktop.

* First, the motors should be connected to the raspberry pi correctly. The program uses the raspberry pi GPIO pins, and assumes that for the left motor, the pins IN1, IN2, IN3, IN4 are connected to GPIO pins 7, 11, 13, and 15, and for the right motor, they are connected to GPIO pins 8, 10, 12, 16, as shown in the connection diagram

[Figure 1 about here.]

* On the raspberry pi, dependencies for openCV need to be installed. Since the openCV is not available in pip for the arm processor in raspberry pi, we it must be installed from source. This takes a few hours on the raspberry pi. To complete the setup including installation of a MQTT client and opencv on the raspberry pi, clone the repository from github on the raspberry pi and navigate to the code folder, open the terminal and run the command

make setup_pi

* Next, install opencv and an MQTT client and MQTT broker on the desktop. For this, clone the repository from github, navigate into the code folder and run the command

make setup_server

* Note the IP address of the desktop so that we can connect to the MQTT server running on it. Connect the raspberry pi and the desktop on the same wireless network.

* To run the code on the desktop, run the command

make run_server [IP address of the MQTT broker]

* Finally to run the code on the raspberry pi, run

make run_pi [IP address of the MQTT broker]

* Now the raspberry pi car can be controlled by typing in W, A, S, or D keys on the desktop in the terminal where the program is running.

* The program can be stopped on both the raspberry pi and the desktop by running

make kill

5.2 Observations

It was observed that the communication between the raspberry pi and the desktop controller application is pretty seamless. The robot car responds without any observable delays when the network is strong. When the network is weak, however, some delays may be observed. The delay becomes more evident in the case of the images sent by the raspberry pi back to the desktop when the network is not strong.

Using the stepper motors, it is difficult to set a how much a motor should turn when it receives a message. If the motor is not allowed to turn long enough, then between two messages the motor will be idle and if it is turned longer than the interval between two messages, there can be conflicts if in response to each of the messages the subscriber running on the raspberry pi tries to set a different step on the motor. Therefore, the movement can seem a little jerky at times.

However, this is not a problem with 360 degrees continuous servo motors. Since the continuous servo motors use pulse width modulation, the speed and direction of rotation can be controlled by sending a square wave with different duty cycles depending on the motor. Since, the motor can be stopped and started easily, there are no conflicts even if the motor is allowed to turn longer than the interval between two messages. However, the motor would respond to the two messages one after the other.

Thus the raspberry pi robot car can be successfully controlled over wifi using MQTT for communication

5.3 Improvements

The project can be improved in various ways. Firstly, even though the deployment with makefile is easy, installing opencv on raspberry pi takes around 4 hours. This can be avoided if we use docker for deployment on the raspberry pi. Two separate images would be needed however one for the processor on the desktop and another one for the arm 8 processor on the raspberry pi.

Machine learning can be incorporated, by collecting the images and the corresponding messages that were sent to the raspberry pi and use it to train a neural network, which could then be used to drive the robot car autonomously. This would be complicated however since car needs to be driven for a long time to get enough data for the neural network to perform well regardless of the surroundings.

Many different sensors could be added to help improve the monitoring capability of the car, and get more information about the environment. If many controlling devices and cars are present, the cars may be controlled in groups and other functionality added to behave as a swarm of cars to complete tasks collaboratively.

6 CONCLUSION

MQTT is a fast and reliable data agnostic and platform independent protocol that allows communication between devices. Raspberry pi

is small but powerful development board that allows users to build prototypes easily and can be used in various applications because of the significantly powerful arm 8 microprocessor. OpenCv is an open source library for computer vision that is optimised to perform operations on images efficiently and is commonly used in computer vision applications. All these technologies were used to build a robot car, controlled via MQTT over a wireless network. MQTT allows us to easily scale up the number of such cars if needed.

ACKNOWLEDGMENTS

The authors would like to thank Dr. Gregor von Laszewski for giving the opportunity to work on this project and for providing the necessary hardware to complete the project. The author would also like to thank the associate instructors of the class for their help and for answering questions on piazza which helped everyone.

REFERENCES

- [1] bitreactive. 2015. The Raspberry Pi Eclipse IoT Car. bitreactive website. (March 2015). <http://www.bitreactive.com/remote-controlled-raspberry-pi-car-part-3-2/>
- [2] Paul Caponetti. 2017. Why MQTT is the Protocol of Choice for the IoT. xively.com blog website. (August 2017). <http://blog.xively.com/why-mqtt-is-the-protocol-of-choice-for-the-iot/>
- [3] Ian Craggs. 2013. MQTT security: Who are you? Can you prove it? What can you do? IBM developer works website. (March 2013). https://www.ibm.com/developerworks/community/blogs/c565c720-fe84-4f63-873f-607d87787327/entry/mqtt_security?lang=en
- [4] eclipse. [n. d.]. mqtt broker. eclipse mosquitto website. ([n. d.]). <https://mosquitto.org/>
- [5] eclipse paho. [n. d.]. Python Client - documentation. eclipse paho website. ([n. d.]). <https://www.eclipse.org/paho/clients/python/docs/>
- [6] hivemq. [n. d.]. intrawebsite mqtt. hivemq website. ([n. d.]). <https://www.hivemq.com/blog/mqtt-essentials-part-1-introducing-mqtt>
- [7] University of Southampton IAM group. 2005. FloodNet. IAM group website. (April 2005). <http://www.iam.ecs.soton.ac.uk/projects/297.html>
- [8] Dexter Industries. 2017. GoPiGo Build and Program Your Own Robot. dexter industries website. (2017). <https://www.dexterindustries.com/gopigo3/>
- [9] Chu Jiangwei, Ji Lisheng, Guo Lie, Wang Rongben, et al. 2004. Study on method of detecting preceding vehicle based on monocular camera. In *Intelligent Vehicles Symposium, 2004 IEEE*. IEEE, 750–755.
- [10] jwatson. 2016. Raspberry Pi Models Comparison Chart Poster. element14 community website. (June 2016). <https://www.element14.com/community/docs/DOC-82195/l/raspberry-pi-models-comparison-chart-poster-free-download>
- [11] makershed.com. 2016. Raspberry pi comparison chart. makershed.com website. (2016). <https://www.makershed.com/pages/raspberry-pi-comparison-chart>
- [12] Hive mq. [n. d.]. MQTT Essentials Part 2: Publish & Subscribe. HiveMQ website. ([n. d.]). <https://www.hivemq.com/blog/mqtt-essentials-part2-publish-subscribe>
- [13] Hive MQ. [n. d.]. MQTT Essentials Part 6: Quality of Service 0, 1 & 2. Hivemq website. ([n. d.]). <https://www.hivemq.com/blog/mqtt-essentials-part-6-mqtt-quality-of-service-levels>
- [14] Hive MQ. [n. d.]. MQTT Security Fundamentals: TLS / SSL. hive mq website. ([n. d.]). <https://www.hivemq.com/blog/mqtt-security-fundamentals-tls-ssl>
- [15] Mqtt. [n. d.]. Mqtt official website. mqtt official website. ([n. d.]). <http://mqtt.org/>
- [16] OpenHab. 2017. What is openHAB?. openhab website. (November 2017). <https://www.openhab.org/introduction.html>
- [17] opensource.com. 2015. What a Raspberry Pi. opensource.com website. (March 2015). <https://opensource.com/resources/raspberry-pi>
- [18] Todd Ouska. 2016. Transport-level security tradeoffs using MQTT. iot design website. (February 2016). <http://iotdesign.embedded-computing.com/guest-blogs/transport-level-security-tradeoffs-using-mqtt/>
- [19] pythonprogramming.net. 2014. Robotics with Python Raspberry Pi and GoPiGo Introduction. pythonprogramming.net. (April 2014). <https://pythonprogramming.net/robotics-raspberry-pi-tutorial-gopigo-introduction/>
- [20] random nerds tutorial. [n. d.]. What is MQTT and How It Works. random nerds website. ([n. d.]). <https://randomnerdtutorials.com/what-is-mqtt-and-how-it-works/>
- [21] raspberrypi.org. 2015. Raspberry Pi Zero: the 5 dollar computer. raspberrypi.org. (November 2015). <https://www.raspberrypi.org/blog/raspberry-pi-zero/>
- [22] raspberrypi.org. 2015. What is a Raspberry pi. raspberrypi.org website. (May 2015). <https://www.raspberrypi.org/help/what-%20is-a-raspberry-pi/>
- [23] IBM research. 2007. IBM Research Demonstrates Innovative 'Speech to Sign Language' Translation System. IBM website. (September 2007). <http://www-03.ibm.com/press/us/en/pressrelease/22316.wss>
- [24] Zheng Wang. 2015. Self Driving RC Car. Zheng Wang wordpress website. (August 2015). <https://zhengludwig.wordpress.com/projects/self-driving-rc-car/>
- [25] Wikipedia. 2017. MQTT – Wikipedia, The Free Encyclopedia. (November 2017). <https://en.wikipedia.org/w/index.php?title=MQTT&oldid=808683219> [Online; accessed 6-November-2017].
- [26] Wikipedia. 2017. Stepper motor – Wikipedia, The Free Encyclopedia. (2017). https://en.wikipedia.org/w/index.php?title=Stepper_motor&oldid=811220740 [Online; accessed 4-December-2017].

A ISSUES

DONE:

Example of done item: Once you fix an item, change TODO to DONE

A.1 Assignment Submission Issues

Do not make changes to your paper during grading, when your repository should be frozen.

A.2 Uncaught Bibliography Errors

Missing bibliography file generated by JabRef

Bibtex labels cannot have any spaces, _ or & in it

Citations in text showing as [?]: this means either your report.bib is not up-to-date or there is a spelling error in the label of the item you want to cite, either in report.bib or in report.tex

A.3 Formatting

Incorrect number of keywords or HID and i523 not included in the keywords

Other formatting issues

A.4 Writing Errors

Errors in title, e.g. capitalization

Spelling errors

Are you using *a* and *the* properly?

Do not use phrases such as *shown in the Figure below*. Instead, use *as shown in Figure 3*, when referring to the 3rd figure

Do not use the word *I* instead use *we* even if you are the sole author

Do not use the phrase *In this paper/report we show* instead use *We show*. It is not important if this is a paper or a report and does not need to be mentioned

If you want to say *and* do not use & but use the word *and*

Use a space after . , :

When using a section command, the section title is not written in all-caps as format does this for you

\section{Introduction} and NOT \section{INTRODUCTION}

A.5 Citation Issues and Plagiarism

- It is your responsibility to make sure no plagiarism occurs.
- The instructions and resources were given in the class
- Claims made without citations provided
- Need to paraphrase long quotations (whole sentences or longer)
- Need to quote directly cited material

Do not submit eps images. Instead, convert them to PDF

The image files must be in a single directory named "images"

In case there is a powerpoint in the submission, the image must be exported as PDF

Make the figures large enough so we can read the details. If needed make the figure over two columns

Do not worry about the figure placement if they are at a different location than you think. Figures are allowed to float. For this class, you should place all figures at the end of the report.

In case you copied a figure from another paper you need to ask for copyright permission. In case of a class paper, you must include a reference to the original in the caption

Remove any figure that is not referred to explicitly in the text (As shown in Figure ..)

Do not use textwidth as a parameter for includegraphics

Figures should be reasonably sized and often you just need to add columnwidth

e.g.

/includegraphics[width=\columnwidth]{images/myimage.pdf}
re

A.6 Character Errors

- Erroneous use of quotation marks, i.e. use "quotes", instead of "
- To emphasize a word, use *emphasize* and not "quote"
- When using the characters & # % - put a backslash before them so that they show up correctly
- Pasting and copying from the Web often results in non-ASCII characters to be used in your text, please remove them and replace accordingly. This is the case for quotes, dashes and all the other special characters.
- If you see a ffigure and not a figure in text you copied from a text that has the fi combined as a single character

A.7 Structural Issues

- Acknowledgement section missing
- Incorrect README file
- In case of a class and if you do a multi-author paper, you need to add an appendix describing who did what in the paper
- The paper has less than 2 pages of text, i.e. excluding images, tables and figures
- The paper has more than 6 pages of text, i.e. excluding images, tables and figures
- Do not artificially inflate your paper if you are below the page limit

A.8 Details about the Figures and Tables

- Capitalization errors in referring to captions, e.g. Figure 1, Table 2
- Do use *label* and *ref* to automatically create figure numbers
- Wrong placement of figure caption. They should be on the bottom of the figure
- Wrong placement of table caption. They should be on the top of the table
- Images submitted incorrectly. They should be in native format, e.g. .graffle, .pptx, .png, jpg

LIST OF FIGURES

1 Example caption

7

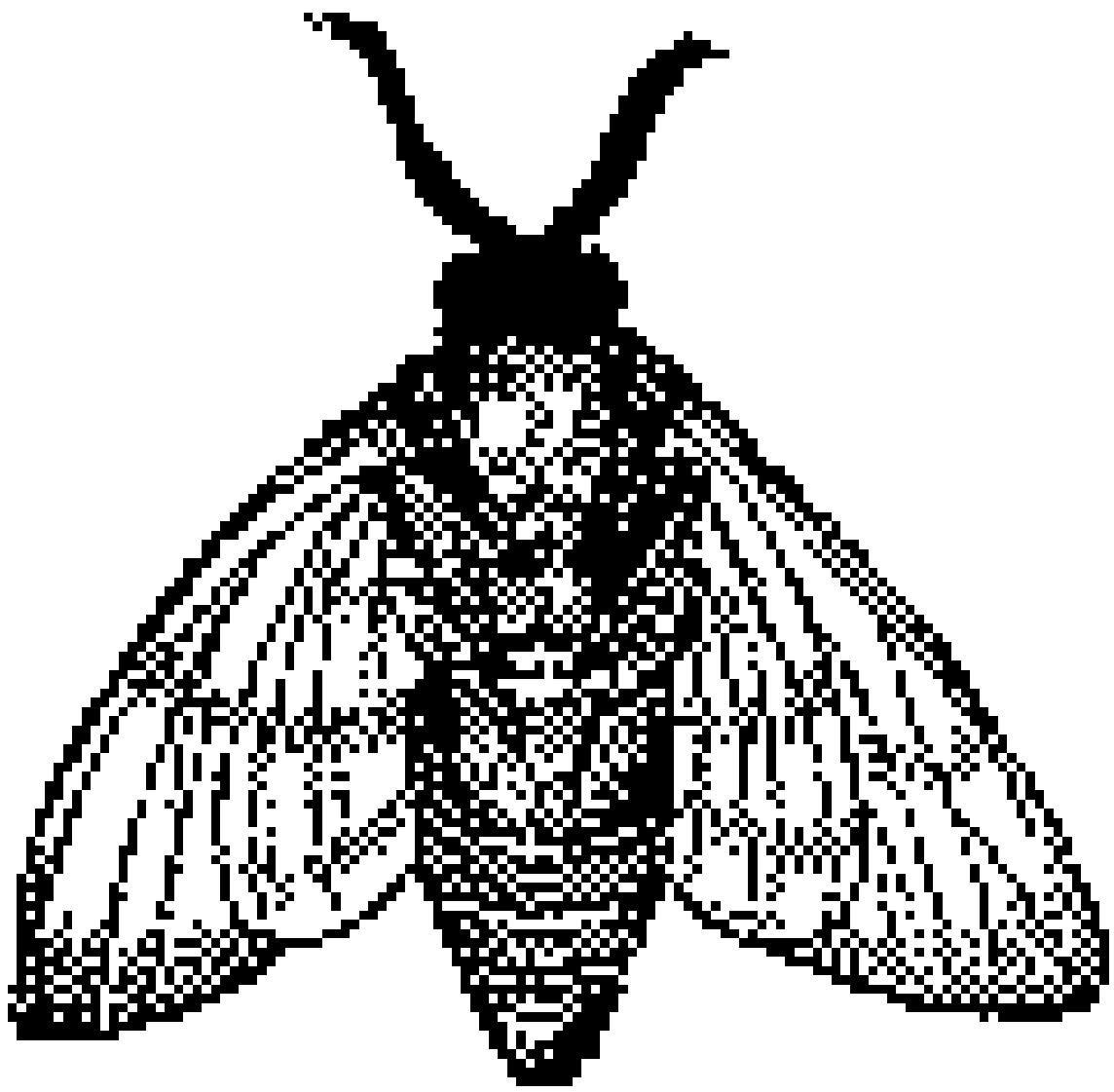


Figure 1: Example caption

```
bibtext report
```

```
=====
This is BibTeX, Version 0.99d (TeX Live 2016)
The top-level auxiliary file: report.aux
The style file: ACM-Reference-Format.bst
Database file #1: report.bib
Warning--empty year in eclipse-mosquitto
Warning--empty year in python-paho-mqtt
Warning--empty year in hivemq-website
Warning--empty publisher in monocular
Warning--empty address in monocular
Warning--empty year in hivemq-details
Warning--empty year in hivemq-qos
Warning--empty year in mqtt-sec-ssl
Warning--empty year in mqtt-official
Warning--empty year in how-mqtt-works
(There were 10 warnings)
```

```
bibtext _ label error
```

```
=====
bibtext space label error
```

```
=====
bibtext comma label error
```

```
=====
latex report
```

```
=====
[2017-12-04 12.22.46] pdflatex report.tex
This is pdfTeX, Version 3.14159265-2.6-1.40.17 (TeX Live 2016) (preloaded format=pdflatex)
Missing character: ""
bookmark level for unknown defaults to 0.
The anchor of a bookmark and its parent's must not be the same. Added a new anchor.
Typesetting of "report.tex" completed in 1.2s.
```

```
./README.yml
9:81      error    line too long (83 > 80 characters) (line-length)
10:81     error    line too long (81 > 80 characters) (line-length)
11:81     error    line too long (82 > 80 characters) (line-length)
12:81     error    line too long (81 > 80 characters) (line-length)
13:81     error    line too long (81 > 80 characters) (line-length)
27:81     error    line too long (81 > 80 characters) (line-length)
27:81     error    trailing spaces (trailing-spaces)
28:81     error    line too long (83 > 80 characters) (line-length)
29:80     error    trailing spaces (trailing-spaces)
30:81     error    line too long (83 > 80 characters) (line-length)
30:83     error    trailing spaces (trailing-spaces)
31:81     error    line too long (83 > 80 characters) (line-length)
31:83     error    trailing spaces (trailing-spaces)
32:81     error    line too long (89 > 80 characters) (line-length)
33:81     error    line too long (89 > 80 characters) (line-length)
34:81     error    line too long (89 > 80 characters) (line-length)
34:89     error    trailing spaces (trailing-spaces)
55:5      error    no new line character at the end of file (new-line-at-end-of-file)
55:1      error    trailing spaces (trailing-spaces)
```

Compliance Report

```
name: Arnav, Arnav
hid: 201
paper1: 20th Oct 2017 100%
paper2: 80%
project: not started
```

```
yamlcheck
```

```
wordcount
```

```
7
wc 201 project 7 3314 report.tex
wc 201 project 7 4123 report.pdf
wc 201 project 7 1507 report.bib
```

```
find "
```

passed: True

find footnote

passed: True

find input{format/i523}

4: \input{format/i523}

passed: True

find input{format/final}

passed: False

floats

222: \%ref{f:connection}
225: \begin{figure}[!ht]
226: \centering\includegraphics[width=\columnwidth]{images/fly.pdf}
227: \caption{Example caption}\label{f:fly}

figures 1
tables 0
includegraphics 1
labels 1
refs 1
floats 1

True : ref check passed: (refs >= figures + tables)
True : label check passed: (refs >= figures + tables)
True : include graphics passed: (figures >= includegraphics)
True : check if all figures are referred to: (refs >= labels)

Label/ref check
passed: True

When using figures use columnwidth
[width=1.0\columnwidth]

do not change the number to a smaller fraction

find textwidth

passed: True

below_check

bibtex

label errors

bibtex errors

```
This is BibTeX, Version 0.99d (TeX Live 2016)
The top-level auxiliary file: report.aux
The style file: ACM-Reference-Format.bst
Database file #1: report.bib
Warning--empty year in eclipse-mosquitto
Warning--empty year in python-paho-mqtt
Warning--empty year in hivemq-website
Warning--empty publisher in monocular
Warning--empty address in monocular
Warning--empty year in hivemq-details
Warning--empty year in hivemq-qos
Warning--empty year in mqtt-sec-ssl
Warning--empty year in mqtt-official
Warning--empty year in how-mqtt-works
(There were 10 warnings)
```

bibtex_empty_fields

entries in general should not be empty in bibtex

find ""

passed: True

ascii

The following tests are optional

Tip: newlines can often be replaced just by an empty line

find newline

passed: True

cites should have a space before \cite{} but not before the {

find cite {

passed: True

Unsupervised Learning For Detecting Fake Online Reviews

Syam Sundar Herle
Indiana University
Bloomington, IN 47408, USA
syampara@iu.edu

ABSTRACT

Nowadays decision making done by organization and individuals are

KEYWORDS

HID 219, Opinionated Spamming

1 INTRODUCTION

1.1 The Need to Modernize Global Record Keeping

Contracts, transactional records, and verification systems are part of the foundational core of the global economy. However, as Iansiti and Lakhani [25] explain, these tools have not modernized to keep up with the needs of the rapidly evolving global economy and are “like a rush-hour gridlock trapping a Formula 1 car.” Records and transactions are still being managed as they were in the 20th century which creates broad consequences for nearly every industry including supply chain and healthcare.

In supply chain, data management methods for records and logistics are usually inconsistent across the different levels of a supply chain [3]. The outdated record management method encourages redundant data to be stored at the same organization as well as across the supply chain which increases IT maintenance costs and decreases trust and transparency [18]. These issues prevent a tertiary party, like the government, to effectively scrutinize records.

Outdated data management processes also negatively impact healthcare. In the USA in 2014 healthcare fraud cost an estimated \$272 billion [12], and in 2016, healthcare data breaches impacted over 27 million patients [9]. Today, medical data management is stifled by antiquated technology that limits patients’ ability to manage and control access to their electronic medical records [13]. In addition, pharmaceutical supply chains are enervated by current record-keeping technologies. Transactional records are rarely shared across pharmaceutical supply chain organizations which consequently increases inventory levels [35]. As a result, total healthcare cost and the opportunity for counterfeit drugs increases [48]. In addition, verification systems are often independent among supply chain retailers and prescribers. The lack coordination opens the door for “doctor shopping” and greater prescription medication abuse [14].

1.2 Rational Exuberance for Blockchain

Blockchain, “an open, distributed ledger that can record transactions between two parties efficiently and in a verifiable and permanent way” [25], has the potential to resolve these and other fundamental problems of the global economy by overcoming many

of the antiquated shortcomings of the traditional means of managing and verifying contracts and transactions. However, like TCP/IP in the 1970s and 1980s, blockchain is an immature technology that faces numerous challenges to mass adoption. In spite of its current limitations, blockchain is already seeing promising applications in various industries extending beyond just finance including healthcare and supply chain. One particularly exciting use case sits at the intersection of healthcare and supply chain for a more secure distribution system for opioid medications that could potentially mitigate the opioid crisis.

2 BLOCKCHAIN OVERVIEW

2.1 The Blockchain Framework

Blockchain is a foundational technology comprised of numerous technological processes and entities. Some of the most significant pieces follow.

2.1.1 Node. Nodes are the individual units connected to the blockchain network. They are computers with adequate software to maintain a blockchain. The blockchain network connects all the nodes and can read and write data to a block [47] [28].

2.1.2 Block. Blocks are the group of records, bundled together by nodes. They follow a specific set of rules and have limited size. Blocks are also linked to the last generated block, thus forming a chain [47].

2.1.3 Smart Contracts. Smart contracts are the codes with time stamps to represent a contract [47]. Iansiti and Lakhani [25] believe that “smart contracts” may be the most transformative blockchain application at the moment,” because they allow for automatic payments whenever contract conditions are met.

2.1.4 Submit Transaction. In case of a new transaction submission to the network, an individual node circulates it to all the other nodes in the network [47]. The main purpose of circulation is approval.

2.1.5 Transaction Approval. When a transaction is submitted and circulated in the network, each node verifies it. Invalid transactions are deleted [47].

2.1.6 Consensus. For multiple systems to work in a distributed network, they must have an agreement. Such a structure is useful in case of fault tolerance when those agreed set of protocols help to restore the data [47].

2.2 Data Types in Blockchain

There are three major types of data stored on a blockchain, namely un-encrypted, encrypted and hashed [18].

2.2.1 Un-encrypted Data. All the organizations have read access to the un-encrypted data. Such data is fully transparent and facilitates immediate dispute resolution [18].

2.2.2 Encrypted Data. The encrypted data can only be read by the organizations with the access to such data. This means an organization should have a decryption key to read to read the encrypted data. Encrypted data provides restricted access but is also stored in every node in the blockchain. In case of a dispute, the decryption key could be used by different organizations to rectify the entry or deletion of any record [18].

2.2.3 Hash Data. Hash data is also a hidden data, where hash keys act like fingerprints to represent changes or entry for any data record. Each organization can easily confirm their hash keys. Breaking the hash key is nearly impossible. Only the hash key is in the blockchain while the record data is stored off-chain by individual organizations. Data could be revealed, in case of a dispute, by the respective organization [18].

2.3 Benefits of Blockchain

Blockchain's framework and data types provide such broad-ranging benefits that blockchain has been proposed as the "cure" to solve many of the world's problems. This exuberance stems from the fundamental benefits that are cornerstones to nearly every industry.

2.3.1 Trust. Blockchains enable parties that do not know each other to trust each other. No single organization is trusted to maintain the records. Instead, all organizations must approve the contents of the record in order to avoid disputes. Therefore, records should have a time stamp and an origin proof. Normally, a third party facilitates this requirement. Blockchains can provide an alternate solution, where organizations jointly manage the records and preventing corruption by a single organization [18].

2.3.2 Access. Blockchains allow for greater control over what information is and is not accessible. The technology enforces identical data to be stored by each organization. When one copy is updated, all the other copies are also updated. This eliminates the need for a third party to facilitate management of records [26]. Alternatively, different levels of read and write access could be provided to different organizations. Although some meta data should be stored in the public ledger.

2.3.3 Redundancy and Security. Blockchain also assists in providing security by disallowing redundancy at the same node. In the areas of logistics and inventory data, blockchain provides a new approach to supply chain management. The core logic of blockchain does not allow duplicate entries to be created in the same place [3]. A unique inventory can have a single entry with multiple updates, but not duplication. This prevents the organizations from creating false information. In the example of a drug inventory, the shipment status for a batch of drugs will be updated for everyone, everywhere. Each entry could be tracked back to its origin [3].

2.3.4 Transparency. Transparency in a business helps to grow trust among organizations. Sharing information can improve relationships among these organizations. Without blockchain, transparency is hard to achieve. Blockchains can help improve the visibility of contracts, legal documents as well as other inter-organization

data [47]. Organizations are not obligated to show all of their data. Some access can be provided to data that could be useful to other organizations and a shared collection of records can also be stored and managed by co-operation from different organizations.

2.3.5 Low Transaction Costs. Through by-passing third-party verification systems such as brokers, lawyers, or banks, blockchain could significantly reduce transaction costs. Not only will this lower costs for existing transactions, it could open up the market for micro-payments [25].

2.4 Challenges to Blockchain Mass Adoption

While it indeed has the potential to help a wide variety of the world's problems, it should not be viewed as a panacea. Blockchain is not mature enough to support mass-market adoption and faces numerous challenges. Rabah [42] states that to be effective, blockchain needs to overcome its shortcomings of lacking standard protocols, unclear regulation, large energy and computing power consumption, privacy, cultural adoption, and high initial capital requirements. Tapscott and Tapscott [51] agree that its current technical infrastructure is not sufficient, its energy consumption and computational requirements are not sustainable, and user-friendly systems have yet to be designed that would allow for mass market adoption.

Society would have to dismantle many technological, governance, organizational, and cultural barriers to create new foundations for a new world economy that relies heavily on blockchain [25]. This will come at the cost of some existing societal norms, core business functions, and people's jobs [25] [42].

2.5 Technology Adoption Lifecycle

Iansiti and Lakhani [25] argue that the process for mass adoption of blockchain may take longer than expected but will follow a fairly predictable technology adoption pattern that parallels the adoption of TCP/IP (transmission control protocol / internet protocol). TCP/IP started as *single-use* and matured to *localized uses, substitutions, and transformations*. It was introduced as a *single-use* in 1972 for e-mail in ARPAnet, a precursor to commercial internet for the US Department of Defense. Met with skepticism, this technology slowly gained traction among some firms in the 1980s and early 1990s for *localized use* and did not become mainstream until the emergence of World Wide Web in the mid-1990s. This then paved the road for infrastructure companies to provide the necessary hardware and software to establish "plumbing" systems for the internet. Once the technical infrastructure was mature enough, companies then developed businesses that *substituted* existing services with online services (such as Amazon books instead of Borders). Finally, a wave of companies created *transformative* applications that fundamentally changed service experiences (such as Napster in the music industry or Skype in telecommunications).

Similarly, blockchain was also launched for a *single use* in 2009 for Bitcoin, a virtual currency. Blockchain has matured to extend beyond cryptocurrencies and is now being applied for various *localized uses* including in healthcare and supply chain. It took over 30 years for TCP/IP to realize its potential, and blockchain will likewise require decades to mature into a revolutionary economic force. However, companies can start planning for this revolution today

and implement blockchains that follow seven design principles [25] [51].

2.6 Seven Design Principles for Blockchain

Tapscott and Tapscott [51] in their book *Blockchain Revolution* propose seven design principles that, when appropriately applied, can help blockchain move down the technology adoption lifecycle and create more honest, cost-effective, and accountable systems.

2.6.1 Networked integrity. Because all organizations on the blockchain must approve updates, “Participants can exchange value directly with the expectation that the other party will act with integrity.” [51].

2.6.2 Distributed Power. Since the blockchain is distributed across a broad network, it cannot be dismantled by authoritarian power, hackers, or other bad actors. There are no single points of failure and the blockchain can still perpetuate even if numerous nodes are compromised [51].

2.6.3 Value as Incentive. Blockchains can align incentives of individual participants with the interests of the entire blockchain. This minimizes organization problems and conflicts of interests [51].

2.6.4 Security. Blockchains can protect against hackers, malware, ransomware, and identity theft by using a variety of security features. Public key infrastructures, private keys, public keys, and verification methods verify participant activities and prevent bad actors from overriding the network [51].

2.6.5 Privacy. Blockchains can and should provide participants with the freedom to expose as little or as much information about themselves as they desire. This allows a participant to act anonymously when desired or to share sensitive information with only appropriate parties when needed [51].

2.6.6 Rights Preserved. To protect against counterfeit items, a blockchain can serve as a public ledger of ownership [51].

2.6.7 Inclusion. Currently, access to certain financial services is limited to those who are deemed “creditworthy”. Blockchains can and should have significantly lower bars of entry that are not managed by banking institutions so that even a poor rural farmer on a remote corner of Earth who isn’t creditworthy, could participate in the blockchain [51].

3 BLOCKCHAIN APPLICATIONS

3.1 Supply Chain

Blockchain, being a public ledger, can be used in different domains with slight variation in its core attributes. While the general implementation says that the data of a single block is public to all the nodes, different sets of access rights could be provided to different classes of users. Such implementation of blockchain could be applied to a supply chain network.

A supply chain requires the involvement of various parties helping each other. This is generally a one-to-one chain network. Often, each organization uses different technologies for record keeping. Record keeping could involve any information ranging from direct

communications to logistics. Trust is an important issue between organizations. Most of the organizations in a supply chain keep individual records, which are not public to other organizations in the supply chain. Organizations share some information like contracts or notarized data. An efficient management of such shared data can be accomplished using a blockchain. The blockchain provides the ability to collect, record, and notarize different types of shared data [18].

Blockchain could also facilitate storing and maintaining logistics data. Such an application could be useful in the field of healthcare, where the government wants to monitor the supply of drugs. An ideal scenario for this would be to mitigate issues like the opioid crisis. Blockchain technology could simplify storage and management of trusted information. It could provide easy access of such critical public sector information to government organizations while providing data security [50]. Blocks comprise of the data records. When these blocks are added to the chain, they become immutable. This means they cannot be deleted or changed by a single organization [50]. A consensus has to be reached by a majority of the organizations for changing any record. Such a feature helps to maintain the security of the records by eliminating data corruption. Each block is verified and managed using some shared protocols. This process can be automated to allow ease of data entry. Two use cases are for counterfeit detection and data analysis.

3.2 Healthcare

Representing over 17% of the United States’ GDP, healthcare costs continue to soar [14]. More effective data management could address many of healthcare’s fundamental issues, and according to a 2011 McKinsey report [33], more effective health data management could save \$300 billion annually. Current innovations focus on placing patients at the center, privacy and access, completeness of information, and cost [14]. Three interesting applications of blockchain for healthcare are in claims adjudication, cyber security and healthcare IoT, and electronic medical records [9].

3.2.1 Claims Adjudication and Fraud Prevention. In 2014, the Economist estimated that the United States wasted \$272 billion dollars on healthcare fraud [12]. Blockchain could not only minimize fraudulent billing; but, by automating claims adjudication and billing processes, obviate the need for administrative and transactional costs through third parties. Gem Health and Capital One are developing a blockchain-based solution for healthcare claims management [9].

3.2.2 Cyber Security and Healthcare IoT. In 2016, there were 450 reported health data breaches, impacting 27 million patients. Hacking and ransomware were responsible for 27% of these breaches. Each additional connected medical device serves as a potential entry point for bad actors. With an estimated 20-30 billion healthcare IoT devices by 2020, blockchain could secure these devices and protect confidential data. Telstra, IBM, and Tierion are three companies that are developing cyber security solutions for connected healthcare devices [9].

3.2.3 Electronic Medical Records. Beleaguered by stifled technology development, limited ownership control by patients, fragmented information systems, and risks of electronic protected

health information hacking, electronic medical records have perhaps the most important use cases for blockchain [57]. Blockchain can provide interoperability of healthcare information, improved security, patient-centric control, and immutable records [9]. Three examples of blockchain-based EMRs include MedRec, Medicalchain, and the Estonian eHealth Foundation. First, by leveraging smart contracts on the Ethereum blockchain, MedRec is a prototype system that provides patients with “one-stop-shop access to their medical history” and shows promise to give ownership of health information back to the patients who can selectively share access through a modern API interface in a secure manner [13]. Second, Medicalchain is a permissioned blockchain distributed on networks of international healthcare providers that allow patients to transfer medical records across national borders [14]. Third, a data security company called Guardtime is using its Keyless Signature Infrastructure system in partnership with the Estonian eHealth Foundation to store Estonian health records on a blockchain.

4 THE OPIOID CRISIS

4.1 Addiction Risk

Since the late 1990s, pharmaceutical companies have downplayed the addictive risk of opioids [38]. However, the addictive nature of prescribed opioid painkillers increases the “potential for unforeseen adverse events for the patient, including overdose, experience of physiological dependence and subsequent withdrawal, addiction, and negative impacts on functioning” [54]. Patients with wholesome medical intentions often fall victim to the pills’ addictive nature. Misuse and eventual abuse of prescribed opioid painkillers are common: 21%-29% of patients prescribed opioids for chronic pain misuse them while 7.8%-11.7% develop an addiction [54]. Moreover, an opioid addiction often serves as a gateway to other illegal drug use. With similar highs, prescription opioid addicts often transition to heroin, an illicit street-made opioid, since it is cheaper and easier to obtain. In fact, 4%-6% of patients using prescribed opioids develop a heroin addiction [38]. Whereas, 75% of heroin users began their opioid addiction with prescription opioids [7].

Despite these risks, opioids are still prescribed at alarming rates. In fact, the United States, with about 5% of the world’s population, consumed 80% of the world’s opioid prescriptions from 2001-2010 [54]. Between 1999 and 2015 the amount of prescribed opioids painkillers such as codeine, fentanyl, oxycodone, Demerol, and Vicodin quadrupled. In the same time period, opioid-related deaths also quadrupled.

4.2 Health Impact

The epidemic has become so severe that in October 2017 President Trump was forced to declare it “a national health emergency” [37]. With no signs of stopping, this epidemic is burgeoning across America killing nearly 91 people a day [45].

In 2015, 33,091 Americans died from an opioid overdose with rural white males at the greatest risk of an opioid overdose. White Americans (27,056) died the most, followed by black Americans (2,741), and Hispanic American (2,507). Generally the middle-aged population was most at risk with the following percent distributions by age group [17]:

- Aged 0-24: 10% of the opioid-related deaths

- 25-34: 26%
- 35-44: 23%
- 45-54: 23%
- 55+: 19%

Males die nearly twice as frequently from an opioid overdose, representing 65% deaths compared with 35% for females [17].

4.3 Financial Impact

The health impacts are the primary reason for concern, but the financial liability associated with the epidemic is also increasing. The estimated financial impact of the crisis grew from \$55.7 billion in 2007 [2] to \$78.5 billion in 2013 [15]. Of the total economic burden, roughly 25% or \$20 billion is conveyed to the public sector [15]. Partitioned between workplace, healthcare, and criminal justice costs, the overall financial burden will continue to rise until a reversal in current trends.

Opioid drug makers are also exposed to significant financial and legal liabilities as lawsuits accusing pharmaceutical companies of deceptive marketing are commonplace. After a U.S. Justice Department probe in 2007, the maker of OxyContin pleaded guilty to federal charges and paid \$634.5 million. In later cases, OxyContin maker Purdue Pharma LP settled two additional cases for a combined \$43.5 million. Since then governments litigating the culpability of opioid drug makers include “South Carolina, Oklahoma, Mississippi, Ohio, Missouri and New Hampshire as well as cities and counties in California, Illinois, Ohio, Oregon, Tennessee and New York” [43]. In a suit filed in April 2017 against the three largest drug retailers in the USA - CVS, Walgreens, and Walmart - lawyers for plaintiffs Cherokee Nation claim that the “Defendants turned a blind eye to the problem of opioid diversion and profited from the sale of prescription opioids to the citizens of the Cherokee Nation in quantities that far exceeded the number of prescriptions that could reasonably have been used for legitimate medical purposes” [21].

4.4 Responses to Mitigate the Crisis

The private sector, government, and academia alike recognize the importance of solving this crisis and are implementing strategies to help mitigate the opioid crisis.

4.4.1 Private Sector. Drug retailers are taking immediate action. In September 2017, CVS pharmacy announced actions to limit patient supply of prescription opioids to seven days, to restrict the strength of opioids dispensed for first time patients and to install 750 more in-store drug disposal kiosks [5] [19].

A longer-term private sector solution is through the use of radio frequency identification (RFID) technology as a method to improve supply chain security [52] [56]. RFID tracking tags are small microchips that are either printed, etched, stamped, or vapor-deposited onto product labels and are intended to replace barcodes. RFID can be read without direct line of sight and at distances up to 30 feet. Research shows that RFID tags have the potential to reduce costs, increase transparency, and identify counterfeit lots. RFID tags have many advantages over current barcode tracking methods. RFID tags can hold up to 32,000 alphanumeric characters compared to just 20 in a barcode. RFID tags have a much higher upfront cost but decrease total supply chain cost due to the timely

process to scan each individual barcode. And unlike RFID tags, barcodes are susceptible to wear and tear and are easily replicated. RFID technology also has its flaws. In addition to the higher upfront cost, each tag costs between 5-10 US cents, significantly higher than bar codes. Moreover, they are vulnerable to electromagnetic interference and poor manufacturing, are larger, and require a much larger IT infrastructure [52] [27]. From a security and transparency perspective, RFID technology is a good option to conform to The Drug Quality and Security Act [1].

4.4.2 Government. Through policy and politics, the federal government is attempting to find solutions to the epidemic. In the same address President Trump declared the opioid epidemic a national health crisis, he proposed “really tough, really big, really great advertising” [11]. Tom Price of the U.S. Department of Health and Human Services outlined a more detailed federal long-term plan including, “improving access to treatment and recovery services, promoting use of overdose-reversing drugs, strengthening our understanding of the epidemic through better public health surveillance, providing support for cutting edge research on pain and addiction, and advancing better practices for pain management”[41]. Additionally, President Trump’s Commission on Combating Drug Addiction and the Opioid Crisis repeatedly mentions “data sharing” as a method to cope and limit the opioid crisis [37].

Multiple studies indicate that states with strong prescription drug monitoring programs (PDMPs) show a significant reduction in the number of opioid-related deaths [39] [40]. Evidence suggests that 72% of physicians were aware of their states’ PDMPs in 2015, but only 52% used their services. Physicians noted difficulties understanding the data formats and retrieval systems as the main barriers to continual use of PDMPs [46]. As a result, low registration rates are common in the 49 states that offer some form PDMPs [20].

Increasing access to Naloxone, an opioid antagonist that rapidly reverses the opioid overdose damage, may be the most important immediate solution to reducing opioid-related deaths [20]. Between 1998 and 2014, 52,283 naloxone kits were distributed among the 30 states with naloxone distribution programs resulting in 26,453 overdose reversals [20]. 27 states have “third-party prescription” laws that allow physicians to prescribe Naloxone to family and friends of individuals with an opioid addiction [20]. To further reduce opioid-related deaths states must reduce malpractice liability for physicians prescribing Naloxone and make Naloxone available without a prescription [20].

In addition, states have started to pass legislation protecting Good Samaritans. As of 2014, 23 states had laws protecting cooperating bystanders, from low-level misdemeanors and drug possession. Without these laws, bystanders are subject to criminal charges and even murder if it is proven they supplied the deadly drugs. Consequently, these laws are necessary to encourage immediate life-saving calls to 911 [4] [20].

Other solutions states should consider is access to medical marijuana, as Pardo [39] found that states with legal medical marijuana dispensaries have lower opioid-related deaths.

4.4.3 Academia. Academic research is helping to propose effective solutions to the opioid crisis. For example, Indiana University announced plans to commit \$50 million and 70 researchers to find solutions that lead to a decline in opioid-related deaths [44]. In

a similar proposal, researchers at the Network for Public Health Law, Boston University, and Northeastern University proposed a four-step solution including “improving clinical decision making and access to evidence-based treatment, investing in comprehensive public health approaches, and re-focusing law enforcement response ”[10].

5 AN OVERVIEW OF PHARMACEUTICAL SUPPLY CHAINS

5.1 Network Nodes

Forward facing supply chain activities occur before a customer purchase. In a pharmaceutical supply chain, forward facing nodes includes manufacturers, warehouses, distributors, and retailers. Reverse facing supply chain activities occur after the sale and include collecting, recycling, redistributing, and disposing of unwanted medications.

5.1.1 Primary Manufacturing. Produces the main active ingredient [49].

5.1.2 Secondary Manufacturing. Often at a different geographic location for tax and labor reasons, secondary manufacturers combine the active ingredients produced by primary manufacturers and adding excipient substances. Secondary manufacturers produce distribution ready SKU medications through one or more of the following processes: granulation, compression, coating, quality control, and packaging [49].

5.1.3 Market Warehouses and Distribution Centers. Due to the cost of setup and cleaning, it is common for primary manufacturers to produce a years’ worth of active ingredients for a particular medication in one batch. This strategy creates a lot of excess finished and work-in-progress inventory that is stored in warehouse and distribution centers [49].

5.1.4 Wholesalers. Roughly 80% of demand flows through wholesalers. The industry is highly competitive and consolidated. The largest five wholesalers accounted for roughly 45% of industry revenue [49] [23].

5.1.5 Pharmacies and Hospitals. The last node on the pharmaceutical supply chain before medications are distributed at a patient level. Major retailers include pharmacies CVS, Walgreens, Walmart, and Rite Aid and hospital systems such as Community Health Systems, Hospital Corporation of America, and Ascension Health [49].

5.1.6 Reverse Supply Chain. The reverse supply chain is often overlooked as a key component of the pharmaceutical supply chain network. Few people take their unwanted medications to proper collection sites. Instead, medications are discarded in the trash and sewage. In fact, in 2003 the world disposed of at least \$760 million worth of prescription medications [24]. By 2014, this number ballooned to an estimated \$5 billion [32]. The roughly 10 million unused and unexpired prescription medications could be recycled and reused, but instead improper disposal leads to dangerous compounds in sewage effluent, surface water, and even drinking water

[24] [32]. Hua, Tang, and Wu [24] suggest a combination of government subsidies, penalties, and marketing to encourage drug makers to collect unwanted and expired medications.

5.2 Weaknesses

The nature of the current pharmaceutical production and supply chain system creates multiple weaknesses.

5.2.1 Lead Time. Lead times, the time it takes between manufacturing and end sale, can take up to 300 days [49]. As a result, high safety stocks are needed to react to future demand.

5.2.2 High Service Levels. The necessity for on-time pharmaceutical products forces retailers to maintain high service levels, the targeted rate of stock-outs. In many cases and especially in hospitals, patient health relies on having the right medication at the right time. A failure to meet this immediate demand could lead to fatal consequences [31] [24].

5.2.3 Imbalance of Information. Another major disadvantage is the lack of collaboration between raw material suppliers, manufacturers, warehouses, wholesalers, and retailers. “The problem is that the different decision-makers do not have access to the same information regarding the state of the entire supply chain network, and in addition they usually operate under different objective functions” [48]. In this decentralized method, manufacturers have a difficult time forecasting demand. In addition, an imbalance of information between supply chain nodes increases cost and stock-outs. However, Nematollahi, Hosseini-Motlagh, and Heydari [35] found that collaborative decision making through information sharing can increase economic benefits for the entire supply chain while also increasing drug fill rate.

5.2.4 Manufacturing Strategy. The mixture of manufacturers ‘push’ strategy and retailers ‘pull’ strategy, results in high safety stocks. At any given point, there is usually 4 to 24 weeks of finished goods that has yet to be delivered to patients [49].

5.2.5 Large Network. Medications pass through several nodes before they are delivered to the market. Safety and security issues face organization conflicts as the capital cost to prevent theft and mismanagement is not equally spread across the supply chain. The number of nodes also increases the likelihood for counterfeits to enter the market. Between each node, medications are shipped and handled between multiple parties and often times across national and state borders [49].

5.2.6 Government Regulation. The government heavily regulates pharmaceutical supply chains to ensure a safe and steady supply of medications. The Drug Quality and Security Act [1] signed by President Barack Obama in 2013 introduced new regulations for the manufacturing and the distribution of pharmaceutical products. The policy mandates the creation of systems to trace lot-level transactions and systems to verify product legitimacy. In addition, any company within the supply chain must obtain federal licensure and authenticate the licensure of their trading partners. These required changes place immense financial pressure on pharmaceutical companies, drug distributors, and prescribers to develop sustainable supply chain solutions. The 2023 deadline gives pharmaceutical

companies time to test and implement the most sustainable and practical solution [16].

5.2.7 Counterfeits. High inventory levels increase supply chain cost, the potential for theft, and the introduction of counterfeits. It is estimated that 10% of the worldwide pharmaceuticals are counterfeit and approaching 25% in developing countries [30]. Pharmaceutical companies lose an estimated \$200 billion annually due to counterfeit drugs [9].

6 BLOCKCHAIN’S POTENTIAL TO MITIGATE THE OPIOID CRISIS

6.1 Moving Opioid Distribution onto the Blockchain

Blockchain can mitigate the opioid crisis through more secure opioid distribution. The 2013 federal passing of The Drug Quality and Security Act [1] provides pharmaceutical supply chain organizations with the necessary regulatory incentives to quickly move onto the blockchain.

The first step to moving opioid distribution onto the blockchain rests in the initial infrastructure investment plan for development and maintenance. The next step is to establish the policies and security clearances of each organization [6]. Once these critical questions are answered, an opioid distribution blockchain would be similar to blockchains in other industries. Each blockchain would start with the genesis node created by the raw material supplier. From there on, each additional downstream node would timestamp an additional hash. When the opioid eventually reaches the patient, the block would contain information on all involved supply chain nodes with timestamps and distribution information including prescribing physician and pharmacist.

The Hyperledger design principles [8] [53], Tapscott and Tapscott’s seven design principles for blockcahn [51] and BlockSci [29] analysis protocols should be included in the design of the blockchain.

6.2 Benefits

6.2.1 Cost Savings. As a proactive cost saving maneuver, drug makers and retailers can move onto the supply chain to prevent future litigation [36]. In addition, blockchain automation saves time and operating costs [53].

6.2.2 Reducing Lead Times. Collaborative record-sharing is the foundation and ultimate strength of blockchain technology. Nematollahi, Hosseini-Motlagh, and Heydari [35] show that collaborative record-sharing among pharmaceutical nodes increases both the social and economic effectiveness of the supply chain. The economic benefits realized through the reduction of the total supply chain inventory levels also decreases lead times.

6.2.3 Collaborative Information Sharing. Blockchain technology has the potential to reduce the opioid epidemic through transparent and decentralized record keeping. In particular, blockchain has the potential to identify prescription drug fraud. Currently without blockchain, opioid addicts can take advantage of the incomplete feedback between physicians and pharmacists by “doctor shopping”,

modifying, and duplicating prescriptions [14]. With pharmaceutical records on the blockchain, this type of activity is easily identifiable.

Blockchain can reduce illegal opioid prescribing and distribution. In the current centralized record keeping system, the U.S. Drug Enforcement Administration (DEA) relies the Controlled Substances Act of 1970, that requires drug companies to report unusually large or otherwise suspicious orders [22]. Drug makers on the other hand claim their responsibility to report is too vague. As a result, identifying “pill mills” is unnecessarily difficult and time consuming. The DEA’s pharmaceutical unit has 600 investigators [22]. With blockchain, record keeping is standardized and accessible to all parties with the correct cryptographic keys.

6.2.4 Post-Sale Opioid Collection. Blockchain technology can also increase the usefulness of post-sale opioid collection. Current medication packaging lacks 2D DataMatrix barcodes making it nearly impossible to identify historical information such as who is returning their medication, who prescribed and sold the medication, and when the medication was prescribed and returned [55]. Blockchain can trace this information leading to better post-sale analysis. In turn, this information can be studied to improve prescribing methodology.

6.2.5 Counterfeit Detection. Blockchain can reduce the high prevalence of illicit counterfeit drugs. Blocks are immutable, that is once a block is created it cannot be deleted or erased [14]. In addition, each batch of product can be traced back to its origin. This means that each batch will have a block of code associated with it. If a batch does not have its presence in the blockchain, then it can be deemed as a counterfeit [50]. Furthermore, blocks with abnormal distribution patterns can be flagged and removed from the supply chain. Creating illicit blocks is easily identifiable as all new blocks must be approved by all parties on the blockchain.

6.2.6 Data Analysis. Academic institutions and researchers should have access to superkeys to analyze trend analysis to provide predictions and usage patterns over various locations and times of the year [34]. Data analysis can provide both a descriptive and predictive overview of the opioid supply chain.

6.3 Outlook

Moving forward, blockchain must overcome multiple adoption risks. Blockchain monopolization in the pharmaceutical supply chain will reduce the effectiveness, safety, and security of the system. Future governmental regulations can prevent mergers and acquisitions in this industry. In addition, future quantum computing power may be strong enough to break cryptographic keys. Investing in security is necessary for future blockchain success. Lost keys will result in irretrievable data; thus, developing a system to overcome this issue is critical. Lastly, blockchain technology is only as good as its users. Encouraging accurate and timely data entry will ultimately define the usability of blockchain in prescription drug distribution [14].

Nonetheless, blockchain can play a role in mitigating the deadly opioid epidemic by providing cost savings, reducing lead times, facilitating information sharing, facilitating post-sale opioid collection, detecting counterfeit, and providing rich data for analysis.

7 CONCLUSION

Although still in its infancy, blockchain has the potential to be just as transformative as TCP/IP. Early and potential applications in healthcare and supply chain suggest that blockchain is indeed moving along the path of technology adoption. Because blockchain is a low-cost solution for supply chain management and provides security and transparency, it can be used for digital data and communication to overall the distribution of controlled substances such as opioids but this model has yet to be tested. But the computation and infrastructure cost for the entire model is low and should be tested to develop a proof of concept system that leverages blockchain to more securely distribute prescription opioids. A prototype model of blockchain can be developed which emulates the current structure of a pharmaceutical supply chain. Such a model can be vital to test out the flaws of blockchain and how to accurately tailor it to the specific use case.

ACKNOWLEDGMENTS

The authors would like to thank Dr. Gregor von Laszewski for his support and suggestions to write this paper.

REFERENCES

- [1] 113th Congress. 2013. H.R.3204 - Drug Quality and Security Act. (Nov. 2013). <https://www.congress.gov/bill/113th-congress/house-bill/3204> Sponsor Rep. Fred Upton.
- [2] Howard G. Birnbaum, Alan G. White, Matt Schiller, Tracy Waldman, Jody M. Cleveland, and Carl L. Roland. 2011. Societal Costs of Prescription Opioid Abuse, Dependence, and Misuse in the United States. *Pain Medicine* 12, 4 (2011), 657–667. <https://doi.org/10.1111/j.1526-4637.2011.01075.x>
- [3] Paul Brody. 2017. How Blockchain Revolutionizes Supply Chain Management. (Aug. 2017). <http://www.digidatamag.com/finance/2017/08/23/how-the-blockchain-revolutionizes-supply-chain-management-05306209>
- [4] Scott Burris, Joanna Norland, and Brian R Edlin. 2001. Legal aspects of providing naloxone to heroin users in the United States. *International Journal of Drug Policy* 12, 3 (2001), 237 – 248. <http://www.sciencedirect.com/science/article/pii/S0953595901000809>
- [5] Shamard Charles. 2017. CVS to Limit Opioid Prescriptions to 7-Day Supply. (Sept. 2017). <https://www.nbcnews.com/storyline/americas/heroin-epidemic/cvs-limit-opioid-prescriptions-7-day-supply-n803486>
- [6] K. Christidis and M. Devetsikiotis. 2016. Blockchains and Smart Contracts for the Internet of Things. *IEEE Access* 4 (06 2016), 2292–2303. <https://doi.org/10.1109/ACCESS.2016.2566339>
- [7] Theodore Cicero, Matthew Ellis, Hilary L Surrott, and Steven Kurtz. 2014. The Changing Face of Heroin Use in the United States A Retrospective Analysis of the Past 50 Years. *JAMA psychiatry* 71 (05 2014), E1–E6.
- [8] Sharon Cocco and Gari Singh. 2017. Top 6 technical advantages of Hyperledger Fabric for blockchain networks. (Aug. 2017). <https://www.ibm.com/developerworks/cloud/library/cl-top-technical-advantages-of-hyperledger-fabric-for-blockchain-networks/index.html>
- [9] Reenita Das. 2017. Does Blockchain Have A Place In Healthcare? Technical Report. Forbes. <https://www.forbes.com/sites/reenitadas/2017/05/08/does-blockchain-have-a-place-in-healthcare/#5ebcaa6dc1c31>
- [10] Corey Davis, Traci Green, and Leo Beletsky. 2017. Action, Not Rhetoric, Needed to Reverse the Opioid Overdose Epidemic. *Journal of Law, Medicine & Ethics* 45 (2017), 20 – 23. <http://proxyub.uits.iu.edu/login?url=https://search.ebscohost.com.proxyub.uits.iu.edu/login.aspx?direct=true&db=aph&AN=122737813&site=ehost-live&scope=site>
- [11] JULIE HIRSCHFELD DAVIS. 2017. Trump Declares Opioid Crisis a Public Health Emergency but Requests No Funds. (oct 2017). https://www.washingtonpost.com/investigations/cherokee-nation-sues-drug-firms-retailers-for-flooding-communities-with-opioids/2017/04/20/03d04a74-2519-11e7-b503-9d616bd5a305_story.html?utm_term=.ee0423b994ba
- [12] The Economist. 2014. *The 272 billion dollar swindle*. Technical Report. The Economist, <https://www.economist.com/news/united-states/21603078-why-thieves-love-americas-health-care-system-272-billion-swindle>.
- [13] Ariel Ekblaw, Asaf Azaria, Thiago Vieira, and Andrew Lippman. 2016. *MedRec: Medical Data Management on the Blockchain*. Technical Report. pubpub.org.

- [14] Mark A. Engelhardt. 2017. Hitching Healthcare to the Chain: An Introduction to Blockchain Technology in the Healthcare Sector. *Technology Innovation Management Review* 7 (10/2017 2017), 22–34. <https://doi.org/10.22215/timreview/1111>
- [15] Curtis Florence, Chao Zhou, Feijun Luo, and likang xu. 2016. The Economic Burden of Prescription Opioid Overdose, Abuse, and Dependence in the United States, 2013. *Medical Care* 54 (01 2016), 901–906.
- [16] U.S. Food and Drug Administration. 2014. Title II of the Drug Quality and Security Act. (Dec. 2014). <https://www.fda.gov/Drugs/DrugSafety/DrugIntegrityandSupplyChainSecurity/DrugSupplyChainSecurityAct/ucm427033.htm>
- [17] Kaiser Family Foundation. 2015. Opioid Overdose Deaths by Race/Ethnicity. (july 2015). <https://www.kff.org/other/state-indicator/opioid-overdose-deaths-by-raceethnicity/?currentTimeframe=0&sortModel=%7B%22colId%22:%22Location%22,%22sort%22:%22asc%22%7D>
- [18] Gideon Greenspan. 2016. Four genuine blockchain use cases. (May 2016). <https://www.multichain.com/blog/2016/05/four-genuine-blockchain-use-cases>
- [19] Claire Hansen. 2017. CVS to Limit Opioid Prescriptions. (Sept. 2017). <https://www.usnews.com/news/national-news/articles/2017-09-22/cvs-to-enforce-new-limits-on-opioid-prescriptions>
- [20] Kathryn Hawk, Federico E Vaca, and Gail D'Onofrio. 2015. Reducing Fatal Opioid Overdose: Prevention, Treatment and Harm Reduction Strategies. *The Yale journal of biology and medicine* 88 (09 2015), 235–245.
- [21] Scott Higham and Lenny Bernstein. 2017. Cherokee Nation sues drug firms, retailers for flooding communities with opioids. (april 2017). https://www.washingtonpost.com/investigations/cherokee-nation-sues-drug-firms-retailers-for-flooding-communities-with-opioids/2017/04/20/03d04a74-2519-11e7-b503-9d616bd5a305_story.html?utm_term=.ee0423b994ba
- [22] Scott Higham and Lenny Bernstein. 2017. THE DRUG INDUSTRY IS TRIUMPH OVER THE DEA. (Oct. 2017). https://www.washingtonpost.com/graphics/2017/investigations/dea-drug-industry-congress/?utm_term=.86b20fd58ff4
- [23] Hoovers. 2017. Drug Wholesalers. (2017). <http://subscriber.hoovers.com/H-industry360/overview.html?industryId=1493>
- [24] Mei-na Hua, Hua-jun Tang, and Zi-lin Wu. 2016. Analysis of a pharmaceutical reverse supply chain based on unwanted medications categories in household. In *Industrial Engineering and Engineering Management (IEEM), 2016 IEEE International Conference on*. IEEE, IEEE, Bali, Indonesia, 1493–1497.
- [25] Marco Iansiti and Karim Lakhani. 2017. *The Truth About Blockchain*. Technical Report. Harvard Business Review.
- [26] Marco Iansiti and Karim R. Lakhani. 2017. The Truth About Blockchain. (feb 2017). <https://hbr.org/2017/01/the-truth-about-blockchain>
- [27] RFID Journal. 2017. How much does an RFID tag cost today? (Aug. 2017). <http://www.rfidjournal.com/faq/show?85>
- [28] Kost De Serves N. Chilton B Kakavand, H. 2017. The Blockchain Revolution: An Analysis Of Regulation And Technology Related To Distributed Ledger Technologies. (april 2017). <http://www fintechconnective com/wpcontent/uploads/2016/11/Luther-Systems-DLA-Piper-Article-onBlockchain-Regulation-and-Technology-SK.pdf>
- [29] Harry A. Kalodner, Steven Goldfeder, Alishah Chator, Malte Moser, and Arvind Narayanan. 2017. BlockSci: Design and applications of a blockchain analysis platform. *CoRR* abs/1709.02489 (oct 2017), 1–14. <http://arxiv.org/abs/1709.02489>
- [30] Theodore Kelesidis, Iosif Kelesidis, Petros I. Rafailidis, and Matthew E. Falagas. 2007. Counterfeit or substandard antimicrobial drugs: a review of the scientific evidence. *Journal of Antimicrobial Chemotherapy* 60, 2 (2007), 214–236. <https://doi.org/10.1093/jac/dkm109> arXiv:/oup/backfile/content_public/journal/jac/60/2/10.1093.jac.dkm109/1/dkm109.pdf
- [31] Peter Kelle, John Woosley, and Helmut Schneider. 2012. Pharmaceutical supply chain specifics and inventory solutions for a hospital case. *Operations Research for Health Care* 1, 2 (2012), 54 – 63. <https://doi.org/10.1016/j.orhc.2012.07.001>
- [32] Jeanne Lenzer. 2014. US could recycle 10 million unused prescription drugs a year. *BMJ* 349 (2014), g7677. <https://doi.org/10.1136/bmj.g7677>
- [33] James Manyika, Michael Chui, Brad Brown, Jacques Bughin, Richard Dobbs, Charles Roxburgh, and Angela Hung Byers. 2011. *Big data: The next frontier for innovation, competition, and productivity*. Technical Report. McKinsey Global Institute, <https://www.mckinsey.com/business-functions/digital-mckinsey/our-insights/big-data-the-next-frontier-for-innovation>.
- [34] Steven McKie. 2015. The Blockchain Meets Big Data and Real-time Analysis. (June 2015). <https://bitcoinmagazine.com/articles/blockchain-meets-big-data-realtime-analysis-1435183048>
- [35] Mohammadreza Nematollahi, Seyyed-Mahdi Hosseini-Motlagh, and Jafar Heydari. 2017. Economic and social collaborative decision-making on visit interval and service level in a two-echelon pharmaceutical supply chain. *Journal of Cleaner Production* 142, Part 4 (2017), 3956 – 3969. <https://doi.org/10.1016/j.jclepro.2016.10.062>
- [36] Yuki Noguchi. 2017. 41 States To Investigate Pharmaceutical Companies Over Opioids. (Sept. 2017). <https://www.npr.org/sections/thetwo-way/2017/09/19/552135830/41-states-to-investigate-pharmaceutical-companies-over-opioids>
- [37] Commission on Combating Drug Addiction and the Opioid Crisis. 2016. Commission Interim Report. (june 2016). <https://www.whitehouse.gov/sites/whitehouse.gov/files/ondcp/commission-interim-report.pdf>
- [38] National Institute on Drug Abuse. 2017. Opioid Crisis. (june 2017). <https://www.drugabuse.gov/drugs-abuse/opioids/opioid-crisis#one>
- [39] Bryce Pardo. 2017. Do more robust prescription drug monitoring programs reduce prescription opioid overdose? *Addiction* 112, 10 (2017), 1773–1783. <https://doi.org/10.1111/add.13741> ADD-16-0812.R1.
- [40] Stephen W. Patrick, Carrie E. Fry, Timothy F. Jones, and Melinda B. Buntin. 2016. Implementation Of Prescription Drug Monitoring Programs Associated With Reductions In Opioid-Related Death Rates. *Health Affairs* 35, 7 (2016), 1324–1332. <https://doi.org/10.1377/hlthaff.2015.1496>
- [41] Tom Price. 2017. Strategy for Fighting Opioid Crisis. (april 2017). <https://www.hhs.gov/about/leadership/secretary/speeches/2017-speeches/secretary-price-announces-hhs-strategy-for-fighting-opioid-crisis/index.html> Tom Price's remarks at the National Rx Drug Abuse and Heroin Summit.
- [42] Kefa Rabah. 2017. Overview of Blockchain as the Engine of the 4th Industrial Revolution. *Mara Research Journal of Business & Management-ISSN: 2519-1381*, 1, 1 (2017), 125–135.
- [43] Nate Raymond. 2017. S. Carolina sues OxyContin maker Purdue over deceptive marketing. (Aug. 2017). <https://www.reuters.com/article/south-carolina-purduepharma/s-carolina-sues-oxycontin-maker-purdue-over-deceptive-marketing-idUSL2N1L10S3>
- [44] Shari Rudavsky. 2017. Indiana has an opioid crisis. See what the state's leading university is doing to help. (Oct. 2017). <https://www.indystar.com/story/news/2017/10/10/state-has-opioid-crisis-see-what-its-leading-university-pledges-50-million-address-opioid-crisis/747151001/>
- [45] David F Scholl L Rudd RA, Seth P. 2016. Increases in Drug and Opioid Involved Overdose Deaths United States. *MMWR Morb Mortal Wkly Rep* 2016, 65:1445fi??1452 (May 2016). <https://dx.doi.org/10.15585/mmwr.mm65501e1>
- [46] Lainie Rutkow, Lydia Turner, Eleanor Lucas, Catherine Hwang, and G. Caleb Alexander. 2015. Most Primary Care Physicians Are Aware Of Prescription Drug Monitoring Programs, But Many Find The Data Difficult To Access. *Health Affairs* 34, 3 (2015), 484–492. <https://doi.org/10.1377/hlthaff.2014.1085>
- [47] Krystsina Sadouskaya. 2017. *Adoption of Blockchain Technology in Supply Chain and Logistics*. Master's thesis. Mikkeli University of Applied Sciences.
- [48] Nihar Sahay and Marianthi Iterapetritou. 2013. Centralized vs. Decentralized Supply Chain Management Optimization. (11 2013). <https://aiche.confex.com/aiche/2013/webprogram/Paper319958.html>
- [49] Nilay Shah. 2004. Pharmaceutical supply chains: key issues and strategies for optimisation. *Computers & Chemical Engineering* 28, 6 (2004), 929 – 941. <https://doi.org/10.1016/j.compchemeng.2003.09.022> FOCAPO 2003 Special issue.
- [50] Axel Doemeyer Steve Cheng, Matthias Daub and Martin Lundqvist. 2017. Using blockchain to improve data management in the public sector. (Feb. 2017). <https://www.mckinsey.com/business-functions/digital-mckinsey/our-insights/using-blockchain-to-improve-data-management-in-the-public-sector>
- [51] Don Tapscott and Alex Tapscott. 2016. *Blockchain Revolution: How the Technology behind Bitcoin is changing Money, Business, and the World*. Penguin Random House LLC, 375 Hudson St, New York, New York 10014.
- [52] Douglas Taylor. 2014. RFID in the Pharmaceutical Industry: Addressing Counterfeits with Technology. *Journal of Medical Systems* 38, 11 (12 Oct 2014), 141. <https://doi.org/10.1007/s10916-014-0141-y>
- [53] TheLinuxFoundationProject. 2017. Revolutionizing the Supply Chain. (2017). <https://www.hyperledger.org/projects/sawtooth/seafood-case-study>
- [54] Kevin Vowles, Mindy Mcenteer, Peter Siyahhan Julnes, Tessa Frohe, John Ney, and David N van der Goes. 2015. Rates of opioid misuse, abuse, and addiction in chronic pain. *Pain* 156, 4 (04 2015), 569–576.
- [55] Dan Walles. 2017. Track and trace is on the way. Is your drug supply chain ready? (June 2017). <https://medcitynews.com/2017/06/track-and-trace-are-you-ready/>
- [56] David C. Wyld. 2008. Genuine medicine?: Why safeguarding the pharmaceutical supply chain from counterfeit drugs with RFID is vital for protecting public health and the health of the pharmaceutical industry. *Competitiveness Review* 18, 3 (2008), 206–216. <https://doi.org/10.1108/10595420810905984>
- [57] Ben Yuan, Wendy Lin, and Colin McDonnell. 2016. *Blockchains and electronic health records*. Technical Report. MIT.

A ISSUES

DONE:

Example of done item: Once you fix an item, change DONE to DONE

A.1 Assignment Submission Issues

DONE:

Do not make changes to your paper during grading, when your repository should be frozen.

A.2 Uncaught Bibliography Errors

DONE:

Missing bibliography file generated by JabRef

DONE:

Bibtex labels cannot have any spaces, _ or & in it

DONE:

Citations in text showing as [?]: this means either your report.bib is not up-to-date or there is a spelling error in the label of the item you want to cite, either in report.bib or in report.tex

A.3 Formatting

DONE:

Incorrect number of keywords or HID and i523 not included in the keywords

DONE:

Other formatting issues

A.4 Writing Errors

DONE:

Errors in title, e.g. capitalization

DONE:

Spelling errors

DONE:

Are you using *a* and *the* properly?

DONE:

Do not use phrases such as *shown in the Figure below*. Instead, use *as shown in Figure 3*, when referring to the 3rd figure

DONE:

Do not use the word *I* instead use *we* even if you are the sole author

DONE:

Do not use the phrase *In this paper/report we show* instead use *We show*. It is not important if this is a paper or a report and does not need to be mentioned

DONE:

If you want to say *and* do not use & but use the word *and*

DONE:

Use a space after . , :

DONE:

When using a section command, the section title is not written in all-caps as format does this for you

\section{Introduction} and NOT \section{INTRODUCTION}

A.5 Citation Issues and Plagiarism

DONE:

It is your responsibility to make sure no plagiarism occurs. The instructions and resources were given in the class

DONE:

Claims made without citations provided

DONE:

Need to paraphrase long quotations (whole sentences or longer)

DONE:

Need to quote directly cited material

A.6 Character Errors

DONE:

Erroneous use of quotation marks, i.e. use “quotes”, instead of ” ”

DONE:

To emphasize a word, use *emphasize* and not “quote”

DONE:

When using the characters & # % - put a backslash before them so that they show up correctly

DONE:

Pasting and copying from the Web often results in non-ASCII characters to be used in your text, please remove them and replace accordingly. This is the case for quotes, dashes and all the other special characters.

DONE:

If you see a ffigure and not a figure in text you copied from a text that has the fi combined as a single character

A.7 Structural Issues

DONE:

Acknowledgement section missing

DONE:

Incorrect README file

DONE:

In case of a class and if you do a multi-author paper, you need to add an appendix describing who did what in the paper

DONE:

The paper has less than 2 pages of text, i.e. excluding images, tables and figures

DONE:

The paper has more than 6 pages of text, i.e. excluding images, tables and figures

DONE:

Do not artificially inflate your paper if you are below the page limit

A.8 Details about the Figures and Tables

DONE:

Capitalization errors in referring to captions, e.g. Figure 1, Table 2

DONE:

Do use *label* and *ref* to automatically create figure numbers

DONE:

Wrong placement of figure caption. They should be on the bottom of the figure

DONE:

Wrong placement of table caption. They should be on the top of the table

DONE:

Images submitted incorrectly. They should be in native format, e.g. .graffle, .pptx, .png, jpg

DONE:

Do not submit eps images. Instead, convert them to PDF

DONE:

The image files must be in a single directory named "images"

DONE:

In case there is a powerpoint in the submission, the image must be exported as PDF

DONE:

Make the figures large enough so we can read the details. If needed make the figure over two columns

DONE:

Do not worry about the figure placement if they are at a different location than you think. Figures are allowed to float. For this class, you should place all figures at the end of the report.

DONE:

In case you copied a figure from another paper you need to ask for copyright permission. In case of a class paper, you must include a reference to the original in the caption

DONE:

Remove any figure that is not referred to explicitly in the text (As shown in Figure ..)

DONE:

Do not use *textwidth* as a parameter for *includegraphics*

DONE:

Figures should be reasonably sized and often you just need to add *columnwidth*

e.g.

```
/includegraphics[width=\columnwidth]{images/myimage.pdf}  
re
```

bibtext report

=====

This is BibTeX, Version 0.99d (TeX Live 2016)
The top-level auxiliary file: report.aux
The style file: ACM-Reference-Format.bst
Database file #1: report.bib
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3085 of file ACM-Reference-Format.bst
Name 1 in "Steve Cheng, Matthias Daub, Axel Domeyer, and Martin Lundqvist" has a comma a
while executing---line 3085 of file ACM-Reference-Format.bst
Name 1 in "Steve Cheng, Matthias Daub, Axel Domeyer, and Martin Lundqvist" has a comma a
while executing---line 3085 of file ACM-Reference-Format.bst
Name 1 in "Steve Cheng, Matthias Daub, Axel Domeyer, and Martin Lundqvist" has a comma a
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Rudd RA, Seth P, David F, Scholl L" for entry opsis10
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Rudd RA, Seth P, David F, Scholl L" for entry opsis10
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Rudd RA, Seth P, David F, Scholl L" for entry opsis10
while executing---line 3085 of file ACM-Reference-Format.bst
Name 2 in "Kathryn Hawk and Federico E Vaca, and Gail D'Onofrio" has a comma at the end
while executing---line 3085 of file ACM-Reference-Format.bst
Name 2 in "Kathryn Hawk and Federico E Vaca, and Gail D'Onofrio" has a comma at the end
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3131 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3131 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3131 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3131 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3131 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3131 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Rudd RA, Seth P, David F, Scholl L" for entry opsis10
while executing---line 3131 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Rudd RA, Seth P, David F, Scholl L" for entry opsis10
while executing---line 3131 of file ACM-Reference-Format.bst
Name 1 in "Steve Cheng, Matthias Daub, Axel Domeyer, and Martin Lundqvist" has a comma a
while executing---line 3131 of file ACM-Reference-Format.bst
Name 1 in "Steve Cheng, Matthias Daub, Axel Domeyer, and Martin Lundqvist" has a comma a
while executing---line 3131 of file ACM-Reference-Format.bst
Name 2 in "Kathryn Hawk and Federico E Vaca, and Gail D'Onofrio" has a comma at the end
while executing---line 3229 of file ACM-Reference-Format.bst
Name 2 in "Kathryn Hawk and Federico E Vaca, and Gail D'Onofrio" has a comma at the end
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Rudd RA, Seth P, David F, Scholl L" for entry opsis10
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Rudd RA, Seth P, David F, Scholl L" for entry opsis10
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Rudd RA, Seth P, David F, Scholl L" for entry opsis10
while executing---line 3229 of file ACM-Reference-Format.bst
Name 1 in "Steve Cheng, Matthias Daub, Axel Domeyer, and Martin Lundqvist" has a comma a
while executing---line 3229 of file ACM-Reference-Format.bst

```
Name 1 in "Steve Cheng, Matthias Daub, Axel Domeyer, and Martin Lundqvist" has a comma a
while executing---line 3229 of file ACM-Reference-Format.bst
Name 1 in "Steve Cheng, Matthias Daub, Axel Domeyer, and Martin Lundqvist" has a comma a
while executing---line 3229 of file ACM-Reference-Format.bst
(There were 44 error messages)
make[2]: *** [bibtex] Error 2
```

```
latex report
=====
```

```
[2017-12-04 12.22.53] pdflatex report.tex
This is pdfTeX, Version 3.14159265-2.6-1.40.17 (TeX Live 2016) (preloaded format=pdflatex)
Missing character: ""
Typesetting of "report.tex" completed in 1.4s.
./README.yml
37:81    error    line too long (85 > 80 characters) (line-length)
37:85    error    trailing spaces (trailing-spaces)
38:81    error    line too long (85 > 80 characters) (line-length)
38:85    error    trailing spaces (trailing-spaces)
39:81    error    line too long (87 > 80 characters) (line-length)
40:81    error    line too long (89 > 80 characters) (line-length)
41:81    error    line too long (89 > 80 characters) (line-length)
41:89    error    trailing spaces (trailing-spaces)
42:53    error    trailing spaces (trailing-spaces)
54:1     error    trailing spaces (trailing-spaces)
```

```
=====
Compliance Report
=====
```

```
name: Syam Sundar Herle Parampali Sreenath
hid: 219
```

```
paper1: 100%
paper2: 100%
```

```
yamlcheck
```

```
wordcount
```

```
10
wc 219 project 10 5929 report.tex
wc 219 project 10 7733 report.pdf
wc 219 project 10 3283 report.bib
```

```
find "
```

```
161: Increasing access to Naloxone, an opioid antagonist that rapidly
      reverses the opioid overdose damage, may be the most important
      immediate solution to reducing opioid-related deaths
      \cite{Hawk01}. Between 1998 and 2014, 52,283 naloxone kits were
      distributed among the 30 states with naloxone distribution
      programs resulting in 26,453 overdose reversals \cite{Hawk01}. 27
      states have "third-party prescription" laws that allow physicians
      to prescribe Naloxone to family and friends of individuals with
      an opioid addiction \cite{Hawk01}. To further reduce opioid-
      related deaths states must reduce malpractice liability for
      physicians prescribing Naloxone and make Naloxone available
      without a prescription \cite{Hawk01}.
```

```
passed: False
```

```
find footnote
```

```
passed: True
```

```
find input{format/i523}
```

```
4: \input{format/i523}
```

```
passed: True
```

```
find input{format/final}
```

passed: False

floats

figures 0
tables 0
includegraphics 0
labels 0
refs 0
floats 0

True : ref check passed: (refs >= figures + tables)
True : label check passed: (refs >= figures + tables)
True : include graphics passed: (figures >= includegraphics)
True : check if all figures are refered to: (refs >= labels)

Label/ref check
passed: True

When using figures use columnwidth
[width=1.0\columnwidth]
do not change the number to a smaller fraction

find textwidth

passed: True

below_check

bibtex

label errors

bibtex errors

This is BibTeX, Version 0.99d (TeX Live 2016)

The top-level auxiliary file: report.aux

The style file: ACM-Reference-Format.bst

Database file #1: report.bib

Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3085 of file ACM-Reference-Format.bst
Name 1 in "Steve Cheng, Matthias Daub, Axel Domeyer, and Martin Lundqvist" has a comma a
while executing---line 3085 of file ACM-Reference-Format.bst
Name 1 in "Steve Cheng, Matthias Daub, Axel Domeyer, and Martin Lundqvist" has a comma a
while executing---line 3085 of file ACM-Reference-Format.bst
Name 1 in "Steve Cheng, Matthias Daub, Axel Domeyer, and Martin Lundqvist" has a comma a
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Rudd RA, Seth P, David F, Scholl L" for entry opsis10
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Rudd RA, Seth P, David F, Scholl L" for entry opsis10
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Rudd RA, Seth P, David F, Scholl L" for entry opsis10
while executing---line 3085 of file ACM-Reference-Format.bst
Name 2 in "Kathryn Hawk and Federico E Vaca, and Gail D'Onofrio" has a comma at the end
while executing---line 3085 of file ACM-Reference-Format.bst
Name 2 in "Kathryn Hawk and Federico E Vaca, and Gail D'Onofrio" has a comma at the end
while executing---line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3131 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3131 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3131 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3131 of file ACM-Reference-Format.bst

Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3131 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3131 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Rudd RA, Seth P, David F, Scholl L" for entry opsis10
while executing---line 3131 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Rudd RA, Seth P, David F, Scholl L" for entry opsis10
while executing---line 3131 of file ACM-Reference-Format.bst
Name 1 in "Steve Cheng, Matthias Daub, Axel Domeyer, and Martin Lundqvist" has a comma a
while executing---line 3131 of file ACM-Reference-Format.bst
Name 1 in "Steve Cheng, Matthias Daub, Axel Domeyer, and Martin Lundqvist" has a comma a
while executing---line 3131 of file ACM-Reference-Format.bst
Name 2 in "Kathryn Hawk and Federico E Vaca, and Gail D'Onofrio" has a comma at the end
while executing---line 3229 of file ACM-Reference-Format.bst
Name 2 in "Kathryn Hawk and Federico E Vaca, and Gail D'Onofrio" has a comma at the end
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Kakavand, H., Kost De Serves, N., Chilton, B" for entry pa
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Rudd RA, Seth P, David F, Scholl L" for entry opsis10
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Rudd RA, Seth P, David F, Scholl L" for entry opsis10
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Rudd RA, Seth P, David F, Scholl L" for entry opsis10
while executing---line 3229 of file ACM-Reference-Format.bst
Name 1 in "Steve Cheng, Matthias Daub, Axel Domeyer, and Martin Lundqvist" has a comma a
while executing---line 3229 of file ACM-Reference-Format.bst
Name 1 in "Steve Cheng, Matthias Daub, Axel Domeyer, and Martin Lundqvist" has a comma a
while executing---line 3229 of file ACM-Reference-Format.bst
Name 1 in "Steve Cheng, Matthias Daub, Axel Domeyer, and Martin Lundqvist" has a comma a
while executing---line 3229 of file ACM-Reference-Format.bst

(There were 44 error messages)

bibtex_empty_fields

entries in general should not be empty in bibtex

find ""

15: note = "",

42: note = "",

76: note = "",

436: pages = "",

437: doi = "",

439: note = "",

passed: False

ascii

=====

The following tests are optional

=====

Tip: newlines can often be replaced just by an empty line

find newline

passed: True

cites should have a space before \cite{} but not before the {

find cite {

passed: True

Big Data Analytics on Food Products Around the World

Karthik Vegi

Indiana University Bloomington
College Mall Apartments
Bloomington, Indiana 47401
kvegi@iu.com

Nisha Chandwani

Indiana University Bloomington
Park Doral Apartments
Bloomington, Indiana 47408
nchandwa@iu.edu

ABSTRACT

Food is one of the basic necessities of human-being. It helps us gain energy to recharge our body to do the daily activities like moving, playing, and thinking. From being a cave man to producing wide variety of foods, we have come a long way. The civilizations shaped the food habits of the world and there is a lot of variance in the food habits across countries. We analyze the *Open Food Facts* database that gathers information on food products from around the world to unearth some food habits of the world and we predict the food grade based on the nutrition facts of the food products.

KEYWORDS

i523, hid231, hid203, big data, food habits, food products, nutrition

1 INTRODUCTION

Open Food Facts is a non-profit initiative started by Stephane Gignard and run by thousands of volunteers around the world. Any person around the world can contribute to the database by simply scanning a product using a mobile app which is made available to IOS and Android. This massive database of food products opens up a lot of opportunities to analyze the food products around the world and understand their food habits. We are particularly interested in the consumption of nutrients that come along with the food items across the world, the composition of different fat content, and the prediction of nutrition grade based on the nutrients.

2 ANALYSIS OF FAT CONTENT IN FOOD

Fat is definitely a nutrient that the body needs and are essential to aid in cell growth, help with energy generation, maintaining body temperature, protect organs, help absorb other essential nutrients that aid in producing energy, improve blood cholesterol level, help reduce inflammation in case of injury, and help in storing energy that can be used for survival when you go without food for few days [1]. But, we do need to keep a track of the consumption because anything that is remotely excess leads to a variety of serious health issues [1].

2.1 Dietary Fats

There are different types of fat if? some are good and some are bad and some needs to be taken within a certain limit [1].

2.1.1 *Saturated Fat.* More intake of saturated fats results in the cholesterol levels in the blood which increases the risk of heart related diseases [1]. The American Heart Association suggests around 5 percent of daily calories from foods containing saturated fat [1]. Meat, cheese and milk are some of the sources of saturated fat [1].

2.1.2 *Trans Fat.* Any type of trans fat whether it is natural or artificial is not good [1]. The reason why food manufacturers use trans-fat is because they are less expensive, can be produced artificially, easy to use with other ingredients, last for a long time and also aid in improving the taste of the food [1]. Trans fats raise the bad fat levels and decreases the good fat levels [1]. The American Heart suggests to completely cut off trans-fat from the diet [1].

2.1.3 *Monounsaturated Fat.* Monounsaturated fats have a good effect on the body when taken within limit [1]. They help reduce the bad cholesterol levels in the blood and thereby decrease the risk of heart diseases [1]. They also help in gaining vitamin E which is a good nutrient that acts as antioxidant [1]. Olive oil, avocados, and sesame oil are some of the sources for monounsaturated fats [1].

2.1.4 *Polyunsaturated Fat.* Polyunsaturated fats have a good effect on the body when taken within limit [1]. They help reduce the bad cholesterol levels in the blood and thereby decrease the risk of heart diseases [1]. They also provide some nutrients that are essential for the body [1]. Soybean oil and sunflower oil are some of the sources for polyunsaturated fats [1].

2.2 Data Cleaning and Transformation

To make the analysis more interesting, the top 20 countries with most value counts for the attributes have been considered. The countries with names combined with other countries were also cleaned in the process. The data set was analyzed was missing values and the attributes with more than 60 percent missing values were removed from the analysis to add consistency. Only the columns that are meaningful in the analysis were retained and the rest were removed from further analysis.

We then display the top 5 countries as a pie-chart and the 5 countries are namely United States, France, Switzerland, Germany, and Spain as shown in Figure 1.

[Figure 1 about here.]

We then impute all the null values with zeroes and we then check the dietary fat content in the foods and check the top countries with fat content using a histogram. The analysis with respect to the fat countries is as follows

2.3 Data Analysis

The top 5 countries with most fat content in the food items are Serbia, United States, Switzerland, Germany, and Sweden as shown in Figure 2.

[Figure 2 about here.]

The top 5 countries with most saturated fat content in the food items are Serbia, United States, Germany, France and Switzerland as shown in Figure 3

[Figure 3 about here.]

The top 5 countries with most trans-fat content in the food items are United States, Brazil, Canada, Australia, Russia, and Serbia as shown in Figure 4

[Figure 4 about here.]

The top 5 countries with most cholesterol content in the food items are United States, Canada, Portugal, Brazil, France, and Italy as shown in Figure 5

[Figure 5 about here.]

3 ANALYSIS OF SUGAR AND SALT CONTENT

4 NUTRITION GRADE LABELLING SYSTEM

France recently took a decision to implement a nutri-score system which will use a color coding mechanism to label the food products that will help consumers know the nutrition grade of the product [2]. The World Health Organization regional office for Europe as a part of its 5 year action plan from 2015-2020 recommends a labelling mechanism for the consumers to know about the quality of the food products at a first glance [2]. This will not only make it easier for the consumers to pick healthier options but it will also regulate food manufacturers to resort to healthier ingredients instead of going for low cost artificial or less healthier ingredients [2].

France after United Kingdom became the second country to implement this system to indicate the main ingredients like fat, salt and sugar content in the food items [2]. France made use of an evidence based system to study different labelling systems to arrive at the best one [2]. By implementing this system, the World Health Organization will keep a check on the growing number of diet related diseases in the Europe region [2]. Europe being the largest consumer of cheese wants to regulate the ingredients that go into the manufacturing process so that people are well informed about their food choices [2].

4.1 Nutrition Grade Prediction as a Big Data Problem

We build a predictive classification model to predict the food nutrition grade based on the ingredients of the food. The goal is to apply various machine learning algorithms to the problem at hand, measure the prediction accuracy to compare and contrast the different algorithms, and arrive at the best algorithm that suits the given data and the problem. This problem can be solved using Big Data and Machine Learning techniques given the size and the complexity of the data.

5 MACHINE LEARNING

Machine Learning is a field in which we train computers in a way that they can learn from the input data [4]. The ideology is that

computers use the training data that is made available to them, learn from it, build a model and use this experience to build knowledge that can be applied on new unseen data [4]. A wonderful example to demonstrate machine learning is the application to detect spam emails where the machine builds knowledge from previously seen emails which are marked as spam, checks new emails to see if they match the historic spam emails and label them as spam or non-spam [4].

5.1 Types of Machine Learning Algorithms

There are primarily two types of machine learning algorithms which are descriptive models and predictive models [4]. A *Descriptive Model* is described as the analysis done and insights gained from slicing and dicing the data in new and interesting ways [4]. One example of descriptive model is pattern discovery that is often used in market basket analysis where transnational purchase details are analyzed [4]. A *Predictive Model* on the other hand involves predicting one value using one or more variables [4]. The learning algorithms tried to build a model that captures the relationship between a response variable and dependent target variables [6].

5.2 Types of Learning

Unsupervised Learning is the process where there is no explicit training data to learn from, so there is simply no mechanism where the machine can learn from previously available data [4]. The same email example can be looked at in a different way where we now want to do anomaly detection in emails [4]. Here the main goal is to detect unusual messages from the bunch of messages and we do not have experience of previous data [4].

Supervised Learning in contrast is the process of gaining knowledge or expertise from the training data which can be applied to future unseen data [4]. Here the model is first trained by using a bulk of training examples and this model is applied against testing data to measure the accuracy [4]. The variable that we need to predict is identified which is called the response variable and the variables that are used to predict the response variables, called the predictor variables are identified [4]. If the existing variables are not sufficiently giving the accuracy that is expected, a method called feature engineering is done where new variables are derived by combining existing variables [4].

6 PREDICTION ANALYSIS

Prediction analysis is the process of working on a large data set using a combination of statistical, data mining and machine learning algorithms to predict the outcome based on past data [4]. There are primarily two types of prediction analysis in machine learning, namely regression and classification [6]. In regression, we try to predict a continuous variable from the predictor variables [6]. A good example of regression is to predict the housing prices from different parameters like year of construction, location, amenities, number of bed rooms etc [6]. Here the response variable is continuous and it is not predefined [6]. Classification on the other hand tries to predict a categorical variable in which we assign each record with a predefined label or a class [6].

Classification is the task of assigning each data record to a predefined class [6]. In machine learning, classification is categorized as a supervised learning technique [6]. This problem has applications in various fields like spam detection, medical applications, astronomy and banking to identify fraudulent transactions from genuine transactions [6]. It is the task as coming up with a model which is essentially a function that maps every data record to a class label [6].

The task at hand is a classification problem since we are trying to predict the food nutrition grade of the products based on the ingredients that go into the product. For this problem we are considering only the data for the country France, since the nutrition grade is available for most food products from the country. Another reason is that France is the first countries in the region to come up with the idea of adding a color coded label to the food products mentioning the nutrition grade. In the subsequent sections we discuss the machine learning techniques used to solve this problem.

6.1 K Nearest Neighbors

6.1.1 Overview. Some of the classification algorithms in machine learning work on the principle of eager learning that involves a two step process where first a model is built from the training data and the model is applied on testing data [6]. In contrast, K nearest neighbors is a lazy learning algorithm where the process of modeling the training data is not done until the test examples are classified [6]. *Rote Classifier* is a good example of lazy learning algorithm memorizes the entire training data to perform classification but has the drawback of not being able to map every test example against the training example [6]. K nearest neighbors algorithm overcomes this drawback by finding all the records that are closest or nearest to the training records [6].

The nearest neighbour puts each attribute list as a data point in the n-dimensional space, given n the number of attributes [6]. Once we have the training examples, we take each test example and compute its distance to the training example classes and assign a class label [6]. Any of the popular distance measures among Euclidean distance, Manhattan distance, Minkowski distance and Mahalanobis distance can be used [6]. The k denotes the k closest points to the test example [6]. Figure 6 shows the structure of the data [6].

[Figure 6 about here.]

6.1.2 Support in Python. KNeighborsClassifier is available in the scikit learn python library.

6.2 Logistic Regression

Logistic regression or logit regression is a special type of regression analysis where the response variable that we need to predict is a categorical variable [6]. Typically logistic regression models the response variable to take two values 1 or 0, pass or fail, win or lose [6]. Logistic regression that takes more than two values for the response variable is called multinomial logistic regression [6]. Here the probability of the response variable to take a categorical value is modelled as a function of the predictor variables [6].

Like a lot of machine learning algorithms, logistic regression works by making a lot of assumptions which should be taken care as a part of the data cleaning and transformation process [4]. It does not assume a linear relationship between the response variables and predictor variables [4]. Since it applies a log transformation on the predicted probabilities, it can handle a variety of relationship between the predictor variables [4]. If the predictor variables are multivariate normal, the algorithm achieves best result although it works even if they are not [4]. Stepwise method must be used in the logistic regression to ensure that we are neither overfitting nor underfitting the data [4]. A very important assumption to be noted in logistic regression is that the each attribute list must be independent, in the sense the data records must not be derived from a before-after setup experiment [4]. It also requires a decently large sample size to work on [4].

6.2.1 Support for Python. LogisticRegression is available in the scikit learn python library.

6.3 Random Forest Classifier

Random forest is a ensemble classification algorithm which is very powerful [6]. Ensemble method is a special process to improve the accuracy of the prediction [6]. The classification algorithms we have seen so far predict the response variable using a single classifier on the test data but ensemble methods use multiple classifiers in tandem and aggregate the predictions to boost the accuracy by a huge margin [6]. Using a combination method, the ensemble method derives a set of base classifiers from the training data and on each iteration takes a vote of all the base classifiers to arrive at a result [6].

Random forest is an ensemble method which works very well for classification problems [6]. It combines the predictions made by multiple classifiers where each classifier independently works on the training data and casts its vote [6]. Unlike methods like AdaBoost which generates values based on independent random vectors using a varied probability distribution, random forest generates values based on fixed probability distribution [6].

6.3.1 Rationale for Random Forest. Consider an example, where we have 25 base classifiers and each base classifier has an error rate of 0.35 [6]. As discussed, the random forest takes the majority vote given by the base classifiers [6]. The model makes a wrong prediction if half or more base classifiers predict wrong, if not the accuracy is improved with an error rate of 0.06 which is far better than using just a single classifier [6].

6.3.2 Support for Python. RandomForestClassifier is available in the scikit learn python library.

7 EXPERIMENTS AND RESULTS

7.1 Algorithm

The problem at hand is to correctly identify the nutrition grade of the food item. The possible labels are, *a* to *e*, with *a* being the best and *e* being the worst grade for a food item. For this task, we have used machine learning techniques that help in predicting the label of each food item. Before getting into the details of each step of the

method, we first present a concise version of the algorithm used for this task:

- (1) Select all the records for country, France. Drop records where nutrition grade is not populated.
- (2) Separate the predictors from the response variable in order to perform data cleaning and data transformation steps.
- (3) Check for missing values in the predictors obtained in the step above. Drop columns with more than 60% missing values.
- (4) Impute the missing values with 0 for remaining columns.
- (5) After imputing the missing values, standardize all the numerical predictors using standard scaler.
- (6) Check for the correlation between different numerical predictors. Drop one predictors from each pair of predictors that show high correlation.
- (7) Combine the pre-processed predictors and the response variable in a single dataframe.
- (8) Divide the data obtained in step above into training and test data using stratified sampling.
- (9) Train different classifiers on the training data and check the performance of each classifier on the test data.

7.2 Data set

For the classification problem, we selected the records for country France.

Number of examples: 123,961

Number of variables: 12

Response variables: *Nutrition Grade*

Predictor variables: *Energy per 100g, Fat per 100g, Saturated Fat per 100g, Carbohydrates per 100g, Sugars per 100g, Fiber per 100g, Proteins per 100g, Salt per 100g, Trans-fat per 100g, Sodium per 100g*

7.3 Python Packages Used

The following Python packages were used to solve the classification problem:

- Pandas: Provides high performance data structures for data analysis and data munging
- Matplotlib: Plotting library that helps embedding plots into applications using GUI
- Seaborn: Visualization package based on matplotlib used for drawing high level statistical graphics
- Scikit-learn: Toolbox with solid implementation of machine learning and other algorithms
- Scipy: Package that supports scientific computing with modules for linear algebra and integration

7.4 Data Cleaning

7.4.1 Step 1: Data Sparsity. Data sparsity refers to the situation where a lot of attributes have missing values which is an advantage in some cases because you only need to store and analyze the data that is available to you and save on computation time and storage [6]. We first check the data value counts for each country. United States, France, Switzerland, Germany and Spain come as the top 5 countries with most data. Since the food nutrition grade was

implemented in France, it has most products with the data available so for this classification problem we use the data filtered on France.

7.4.2 Step 2: Handling Missing Values. Missing values is a common scenario and they can be handled in different ways. You could choose to eliminate the data objects with missing values but at the expense of missing some critical analysis [6]. Estimating the missing values is also a good way to handle them, especially when the data comes from time series etc, where you could possibly interpolate the missing values from the ones that are closer to it [6]. Ignoring the missing values is another technique which can be applied for tasks like clustering where the similarity can be calculated using the attributes other than the missing ones [6].

The data set was first analyzed to check the missing values in all the columns. The threshold limit has been set at 60 percent. All the columns with missing values more than 60 percent were removed from the analysis to make the result more consistent. Once the columns were removed, the data set has to be re-indexed to maintain the order. Only the columns that are important for the prediction task has been retained from the original data set. In this case, all the ingredients which are primarily the predictor variables were included. The missing values in the response variable also need to be taken care of. Removing the records with missing values for the response variable proved to be the best option after trying out various things.

Imputation was used to handle the null values in the predictor variables. Imputation can be done in a variety of ways with by imputing the missing values by calculating the mean and the mode or just replacing them with 0. Since all the predictor variables have numeric values, all the null values have been replaced with 0. To ensure the imputation process has been done correctly, the sum of missing values is calculated.

7.4.3 Step 3: Outlier Treatment. Outliers are data objects with quite distinct characteristics from the other data records [6]. There is a considerable difference between anomalies and outliers, where anomalies refer to data records that have bad data which is noise and need to be ignored, anomalies often contain interesting aspects and can lead to some good analysis [6]. In applications like *Fraud Detection*, anomalies could be of utmost importance [6]. The outliers in the data have been looked at by using box plots and have been handled as a part of the data cleaning process.

7.5 Exploratory Data Analysis

For exploratory data analysis, we used the Seaborn package along with Matplotlib for visualizations. The measure of spread that is the range and variance of the values is a good way to understand the different aspects of the predictor variables. Box-plots are a method of visualization to look at the distribution of values for a numerical attribute [6]. The box plots show the percentiles where the lower and upper ends of the box indicate 25th and 75th percentile, the line inside the box indicates the 50th percentile, the tails indicate the 10th and 90th percentile respectively [6].

7.5.1 Bi-variate box-plots. Bi-variate box-plots go beyond univariate box plots by showing the relationship between the predictor variable and the response variable [6]. We look at the bi-variate box-plots for each of the important predictor variables namely saturated fat, polyunsaturated fat, sugars and salt and the response variable which is the nutrition grade. Figure 7 shows the bi-variate box plots.

[Figure 7 about here.]

By looking at the box plots, we can understand some important aspects of how the response variable is related to the predictor variables. We see that as the saturated fat content increases, the food grade decreases and as the polyunsaturated fat content increases the nutrition grade is better. When the sugar levels increase the health quotient of the food comes down and the energy levels behave in a interesting manner where the energy for the nutrition grade A is higher, the energy slightly increases with the nutrition grade. While increase in energy does not necessarily imply that the nutrition quality is high because there are a lot of instant energy foods that have a lot of additives but they are often rated low when it comes to health.

7.5.2 Correlation. Correlation between data objects is the measure of the linear relationship between the attributes of the object that are continuous variables [6]. Correlation analysis is the process of finding of the correlations between the different predictor variables and helps find collinearity problem [4]. The relationship could be either linear or non-linear given the data [6]. The correlation coefficient can range anywhere between -1 and 1, where 1 indicates a positive correlation and -1 indicates negative correlation [4]. Correlation plot visually shows the correlation coefficient between the variables in a nicely laid out plot. Figure 8 shows the correlation plot.

[Figure 8 about here.]

By looking at the correlation plot, we can see that sugars, fat, energy are positively in correlation with the nutrition grade. They will play an important role in the prediction algorithm. Also, sodium and salt are highly correlated with each other and this may lead to collinearity problem if not handled. Collinearity is the state where the independent variables are highly correlated with each other which can add a lot of noise to the data [5]. Some of the problems because of collinearity are that the regression coefficient may not be estimated correctly and also makes it very difficult to explain the response variables using the predictor variables [5]. So we remove sodium from the predictor variables and proceed to the next step.

7.5.3 Data Transformation. Data transformation refers to the transformation that is applied to the variables [6]. For each data object, we apply a transformation function to all the attributes of the object to ensure that the attributes do not have a lot of variance in the data [6]. This process is also called standardization since we are applying a standard function to make sure all the attributes fall within a given range [6]. There are different methods that can be applied to achieve scaling namely log transformation, absolute value, square root transformation [6].

We use the method called normalization where all the values fall in between the range 0 and 1. To achieve this, we use the prepossessing package from sklearn which provides utility functions and transformer classes to change raw data into a standard representation. A lot of machine learning algorithms work well on standard data. If some of the variables have extreme values, they might dominate the model function and might disturb the estimation parameter.

There was a massive improvement in the prediction accuracy of the algorithms before and after data scaling which proves the importance of data standardization with respect to machine learning algorithms.

7.6 Data Sampling

In a supervised machine learning approach, the model is trained on one sample of the data and later tested on a different sample of the data. Thus, in order to test the performance of the nutrition grade classifier, the data for the country France was divided into two samples, training and testing. There are various ways to achieve this split or sampling of the data. Some of these sampling methods are:

- **Simple Random Sampling:** This is one of the simplest sampling techniques. In this technique, every data point has an equal chance of being selected. In other words, it works similar to a lottery system where every outcome has an equal probability. The biggest advantage of this technique is the ease of implementation and its unbiased nature while generating the sample. However, random sampling might not always result in a sample that can represent the true population. It generally works well when we have huge data to sample from.
- **Stratified Sampling:** This technique is a more sophisticated method of sampling data. Stratified sampling generates a sample such that the proportion of each class in the sample is same as that in the true population. In this technique, the entire population is divided into groups or strata. The next step is to randomly select data points from each strata such that the final sample has the same proportion for each strata as that present in the true population. Thus, the sample generated by this technique is a good representative of the true population. Stratified sampling is a very useful technique when the classes in the data are highly imbalanced.

For our classifier, we chose to divide the data for France into training and test samples using stratified sampling technique. The strata or groups were created based on the response variable, i.e., food grade. This ensured that the training and test data had the same proportion for each food grade.

7.7 Data Modeling

Once the data was divided into training and test data, the next step was to train different classifiers and tune their respective parameters for better accuracy. We implemented three different models for classifying the food grade. Each of these models along with their parameters are:

- K Nearest Neighbors (kNN): For kNN, the grade of a food item in test data is classified by first finding the k most similar food items in the training data. It then takes vote (food grade label) from each of these neighbors and based on the majority vote, the food item in the test data is assigned a food grade. Thus, one of the most important parameter for kNN is k , i.e., the number of neighbors to consider from the training data. We tried different k values and found that $k=3$ gave the best accuracy.
- Logistic Regression: For logistic regression, one of the important parameters is penalty. This parameter specifies the kind of regularization to be applied. This parameter can take two possible values, l_1 regularization and l_2 regularization. Both these values penalize high magnitude of the co-efficients of the predictors in order to prevent the model from over-fitting. For our model, we have used l_2 regularization as it works well even in the presence of highly correlated features.
- Random Forest: For random forest, there are many parameters, such as the number of trees in the forest, the maximum depth of the trees, maximum number of features to consider at each split, the minimum number of samples required in a sub-tree to qualify for further split, the minimum number of samples required to qualify as a leaf node, etc. For our data, we have kept most of the models at their default values except for the number of estimators or trees in the forest. We have set this value to 100 as the classifier produced very high accuracy with 100 trees in the forest.

7.8 Evaluation Metrics and Results

There are various evaluation metrics for assessing the performance of classifiers. Some of these evaluation metrics are [3]:

- Accuracy: This metric gives the proportion of the total number of correctly classified instances
- Precision: This gives the proportion of the true positive instances from the total instances classified as positive
- Recall: This gives the proportion of the positive instances that are correctly classified
- F-Measure: This gives the harmonic mean between precision and the recall values
- Confusion Matrix: This is a useful way of checking the accuracy of the classifier. It clearly shows the number of instances correctly classified for each label. Thus, if we know that the classes in the data are not well-balanced, it's always a good idea to check the confusion matrix along with accuracy. Consider a case where 95% of the instances belong to class A and only 5% of the instances belong to class B. If a classifier is trained on a dataset with such imbalance, there is a high chance that the classifier would return label A for each test instance. The classifier would be still be able to correctly classify 95% of the test instances resulting in 95% accuracy. Thus, this is a case where accuracy can be misleading and thus a quick look at the confusion matrix can help understand the problem with the classifier. For such a case, the confusion matrix will clearly show

that all the instances of the minority class, B, have been misclassified.

For our model, we used accuracy as well as confusion matrix for evaluating the results. The confusion matrix did not show any serious issues for any of the classifiers. The accuracy for each of the three classifiers was:

- (1) Logistic Regression: With l_2 penalty, the accuracy of logistic regression was 78.9%. Figure 9 shows the confusion matrix.

[Figure 9 about here.]

- (2) K Nearest Neighbors: With k as 3, the accuracy of kNN was 95.74%. Figure 10 shows the confusion matrix.

[Figure 10 about here.]

- (3) Random Forest: With number of trees as 100, the accuracy of random forest classifier was 99.68%. Figure 11 shows the confusion matrix.

[Figure 11 about here.]

Thus, we obtained the best results with Random Forest classifier.

8 CONCLUSION

This concludes

ACKNOWLEDGMENTS

This project was undertaken as a part of the course objective for I523: Big Data Applications and Analytics at Indiana University, Bloomington. We would like to thank Prof. Gregor von Laszewski and all the TAs for their help, support and suggestions.

REFERENCES

- [1] American Heart Association. 2017. Dietary Fats. Webpage. (March 2017). <https://healthyforgood.heart.org/eat-smart/articles/dietary-fats>
- [2] World Health Organization Europe. 2017. Labelling systems to guide consumers to healthier options. Webpage. (March 2017). <http://www.euro.who.int/en/countries/france/news/news/2017/03/france-becomes-one-of-the-first-countries-in-region-to-recommend-colour-coded-front-of-pack>
- [3] M Hossain and MN Sulaiman. 2015. A review on evaluation metrics for data classification evaluations. *International Journal of Data Mining & Knowledge Management Process* 5, 2 (2015), 1.
- [4] Shai Shalev-Shwartz and Shai Ben-David. 2014. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press, New York, USA.
- [5] Statistics Solutions. 2017. Multicollinearity. Webpage. (March 2017). <http://www.statisticssolutions.com/multicollinearity/>
- [6] Pang-Ning Tan, Michael Steinbach, and Vipin Kumar. 2006. *Introduction to Data Mining*. Pearson, Boston, USA.
- [7] Karthik Vegi and Nisha Chandwani. 2017. Code base - Analysis on food products around the world. github. (Dec. 2017). <https://github.com/bigdata-i523/hid231/tree/master/project/code>

LIST OF FIGURES

1	Top 5 countries [7]	8
2	Top 5 countries with most fat content [7]	9
3	Top 5 countries with most saturated fat content [7]	10
4	Top 5 countries with most trans-fat content [7]	11
5	Top 5 countries with most cholesterol content [7]	12
6	K nearest neighbors algorithm[6]	13
7	Bi-variate box plots [7]	13
8	Correlation Plot [7]	14
9	Confusion matrix for Logistic Regression [7]	15
10	Confusion matrix for K Nearest Neighbors [7]	16
11	Confusion matrix for Random Forest [7]	17

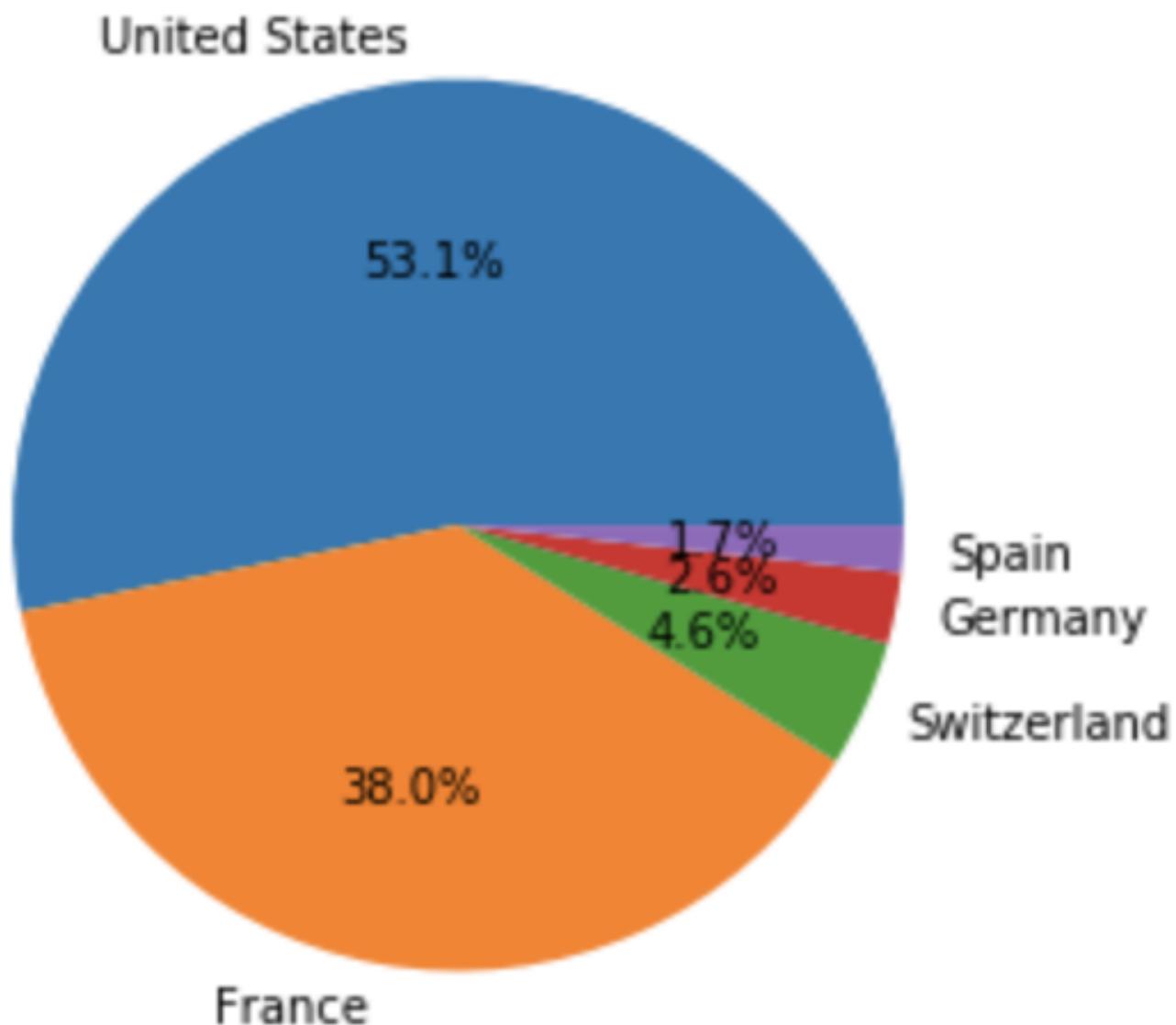


Figure 1: Top 5 countries [7]

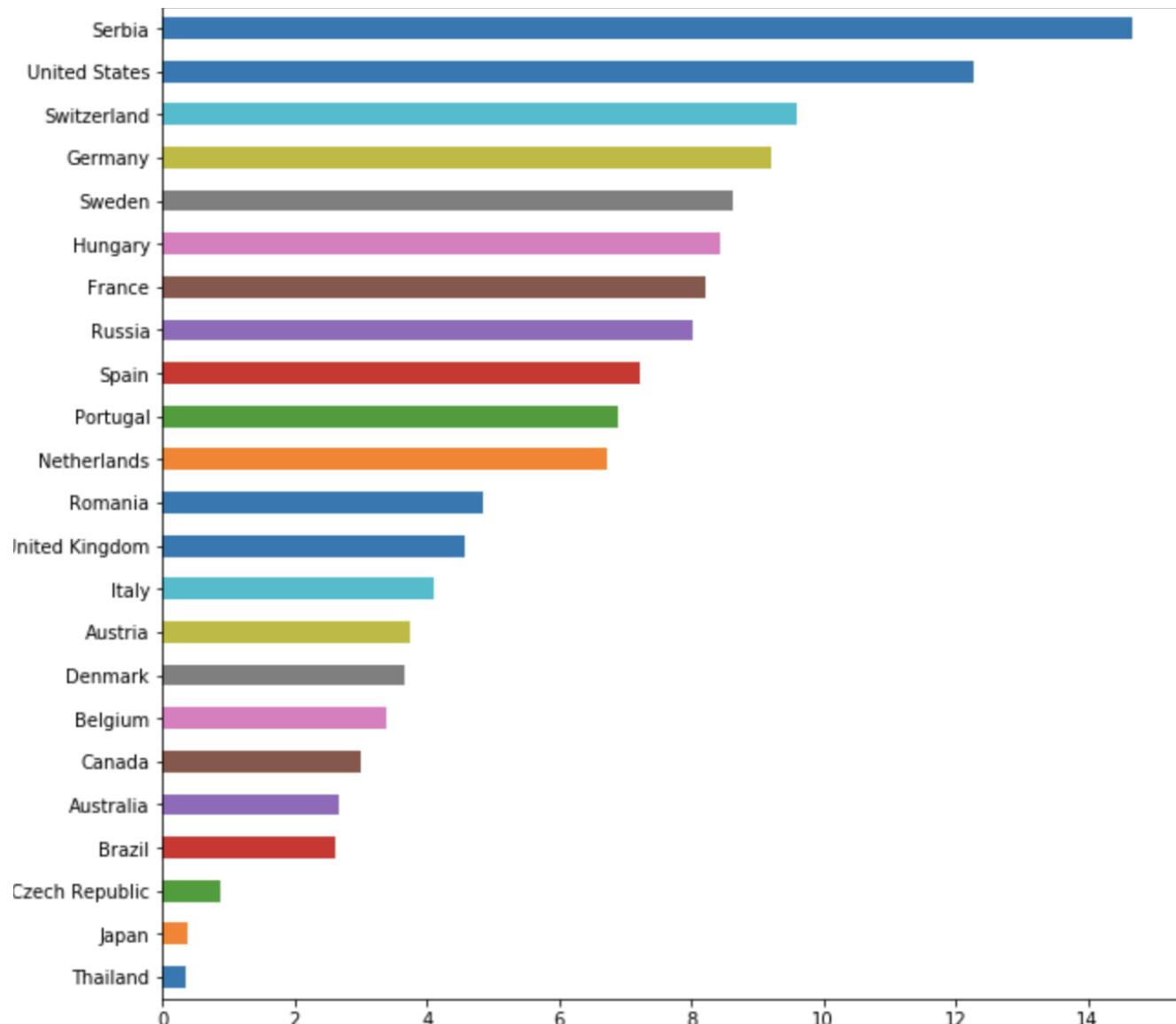


Figure 2: Top 5 countries with most fat content [7]

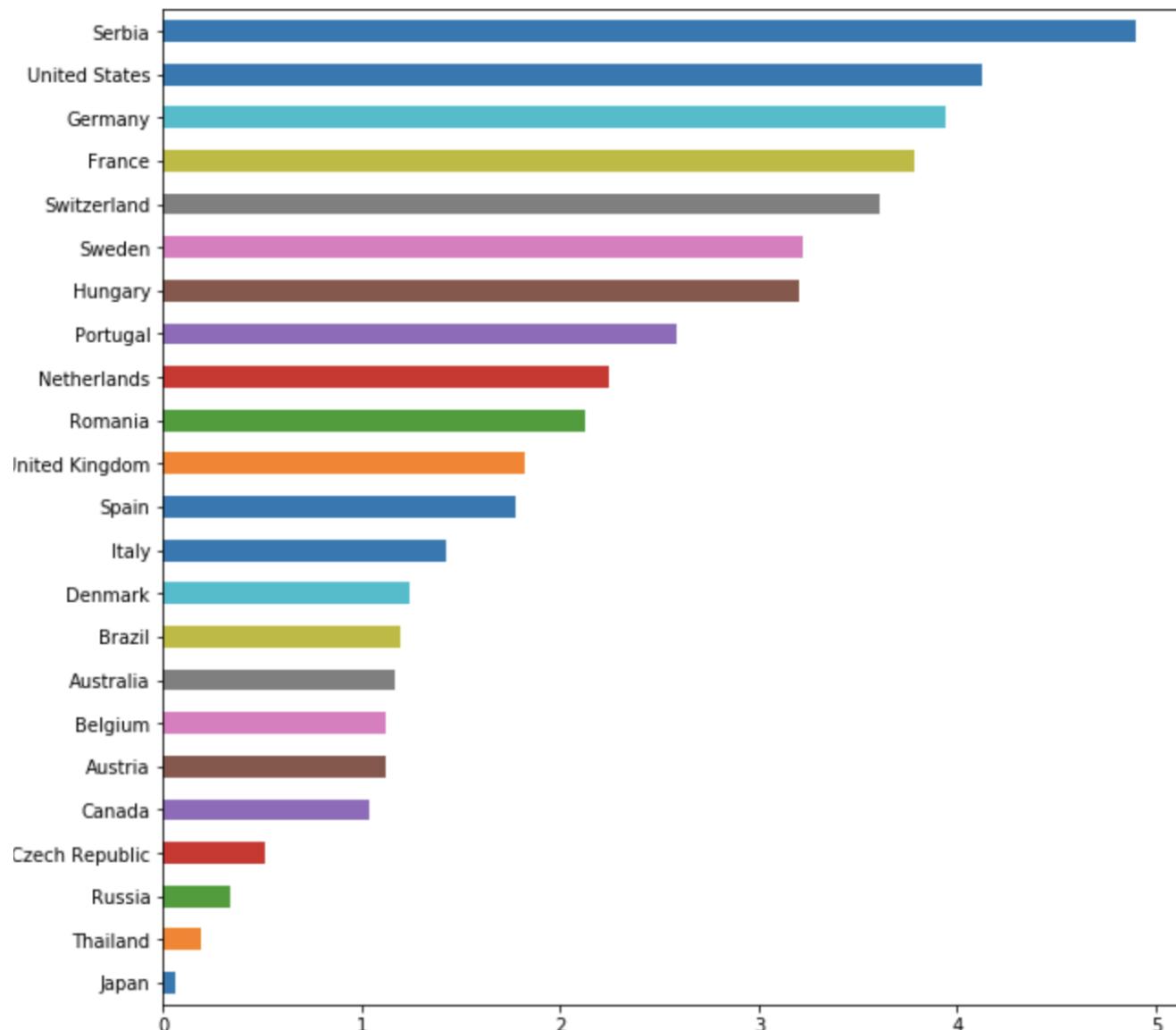


Figure 3: Top 5 countries with most saturated fat content [7]

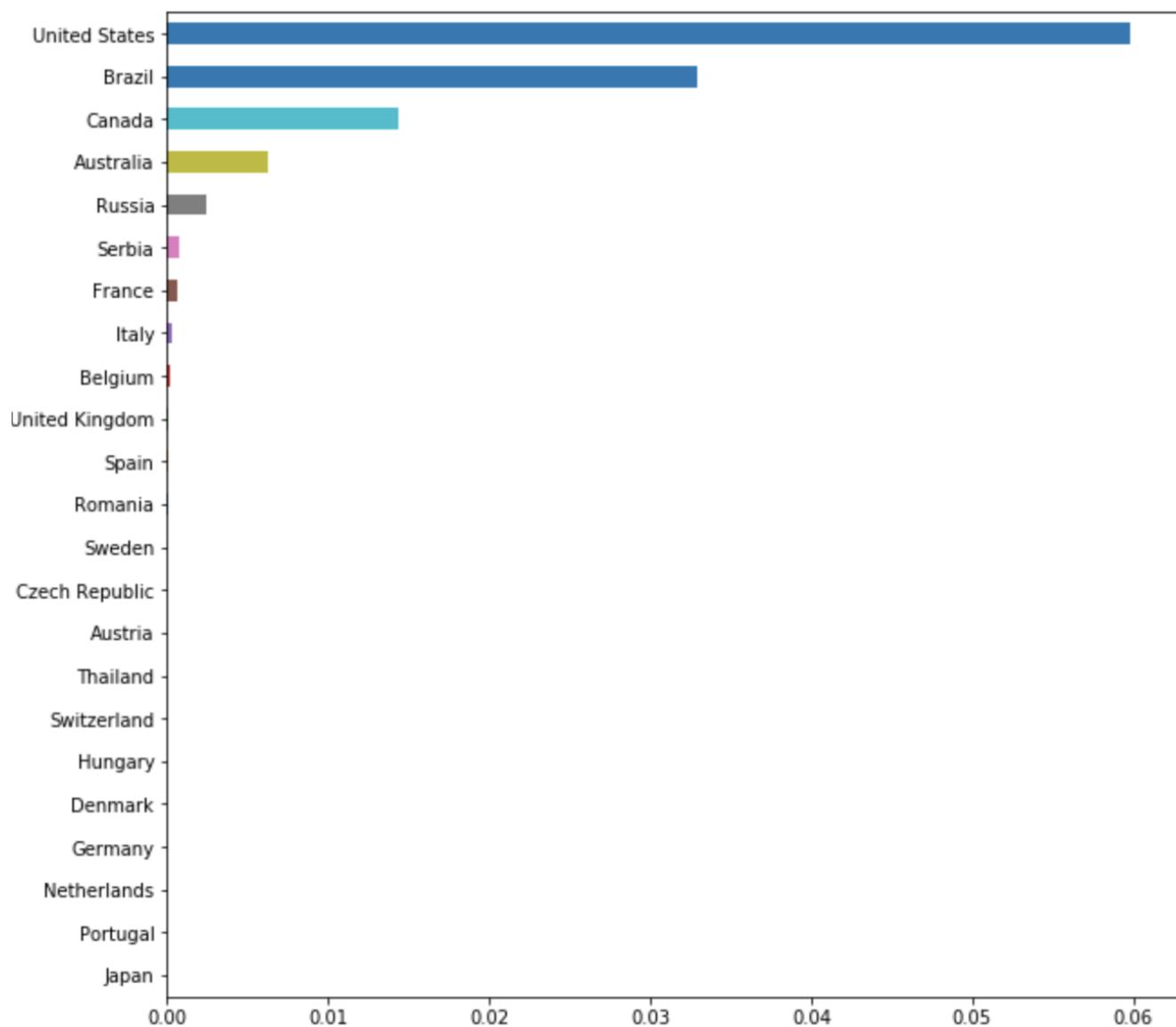


Figure 4: Top 5 countries with most trans-fat content [7]

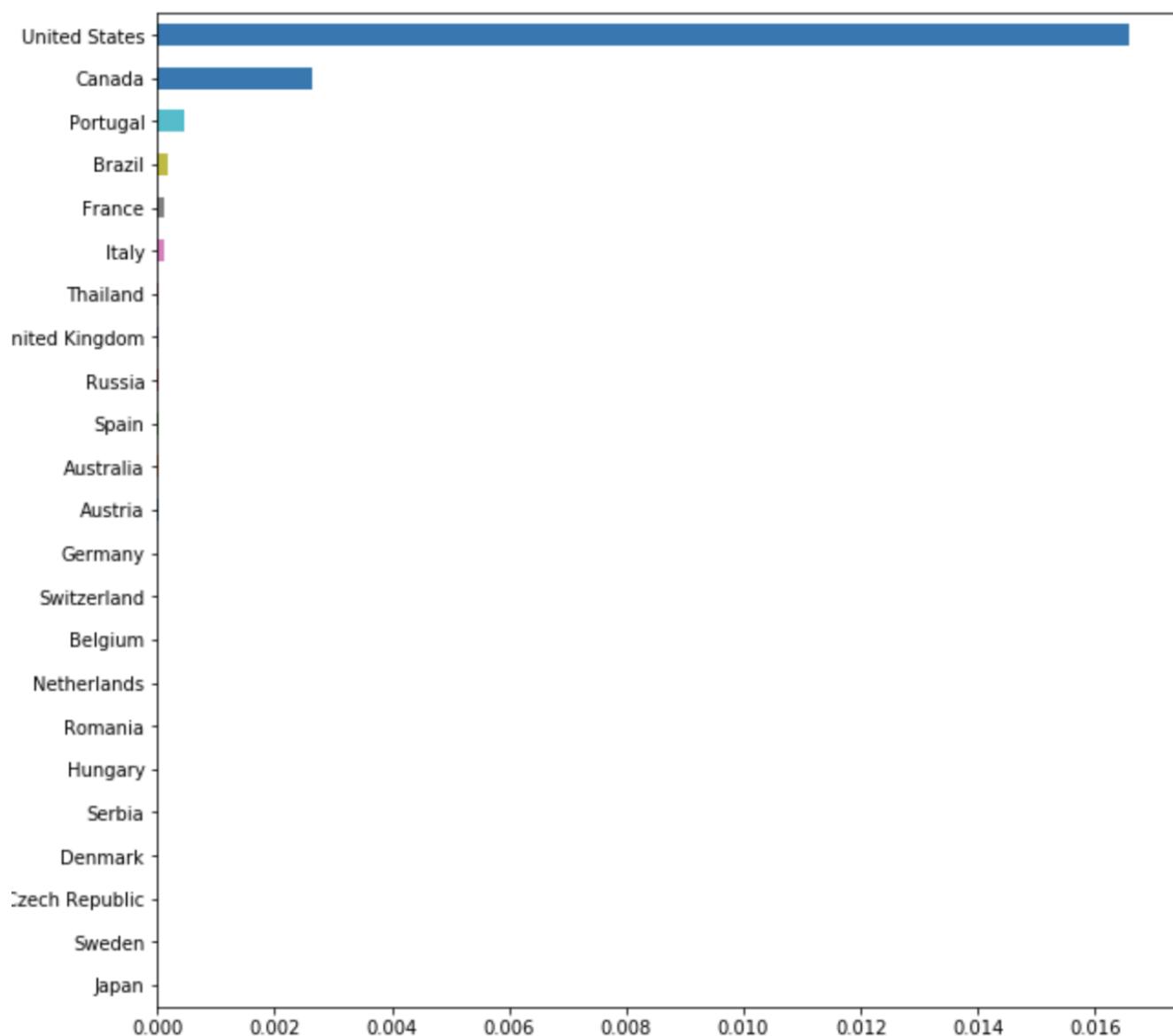


Figure 5: Top 5 countries with most cholesterol content [7]

Algorithm 5.2 The k -nearest neighbor classification algorithm.

- 1: Let k be the number of nearest neighbors and D be the set of training examples.
 - 2: **for** each test example $z = (\mathbf{x}', y')$ **do**
 - 3: Compute $d(\mathbf{x}', \mathbf{x})$, the distance between z and every example, $(\mathbf{x}, y) \in D$.
 - 4: Select $D_z \subseteq D$, the set of k closest training examples to z .
 - 5: $y' = \operatorname{argmax}_v \sum_{(\mathbf{x}_i, y_i) \in D_z} I(v = y_i)$
 - 6: **end for**
-

Figure 6: K nearest neighbors algorithm[6]

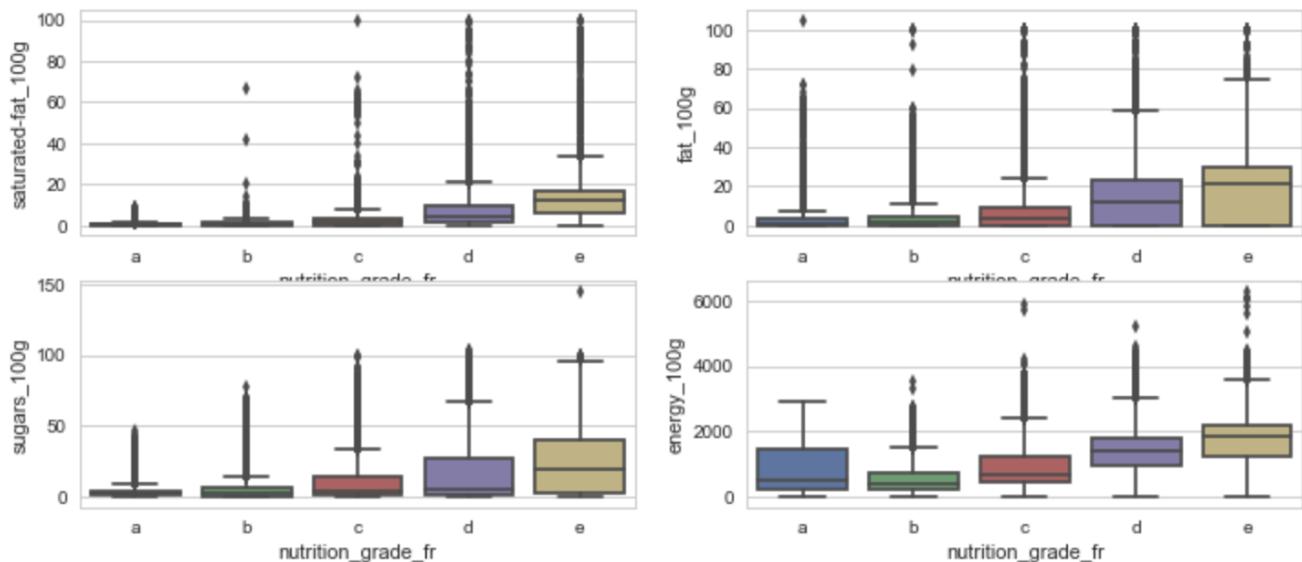


Figure 7: Bi-variate box plots [7]

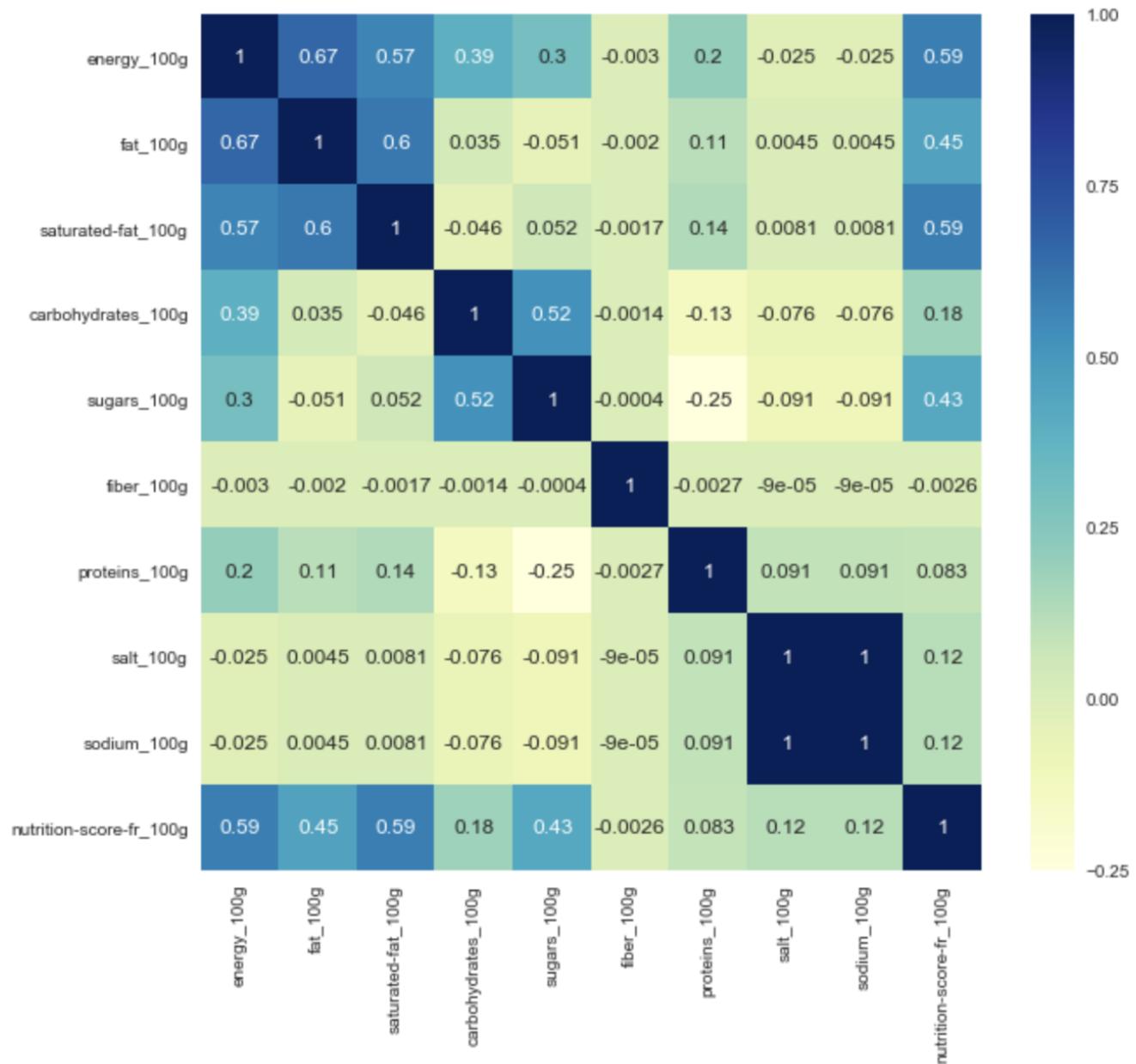


Figure 8: Correlation Plot [7]

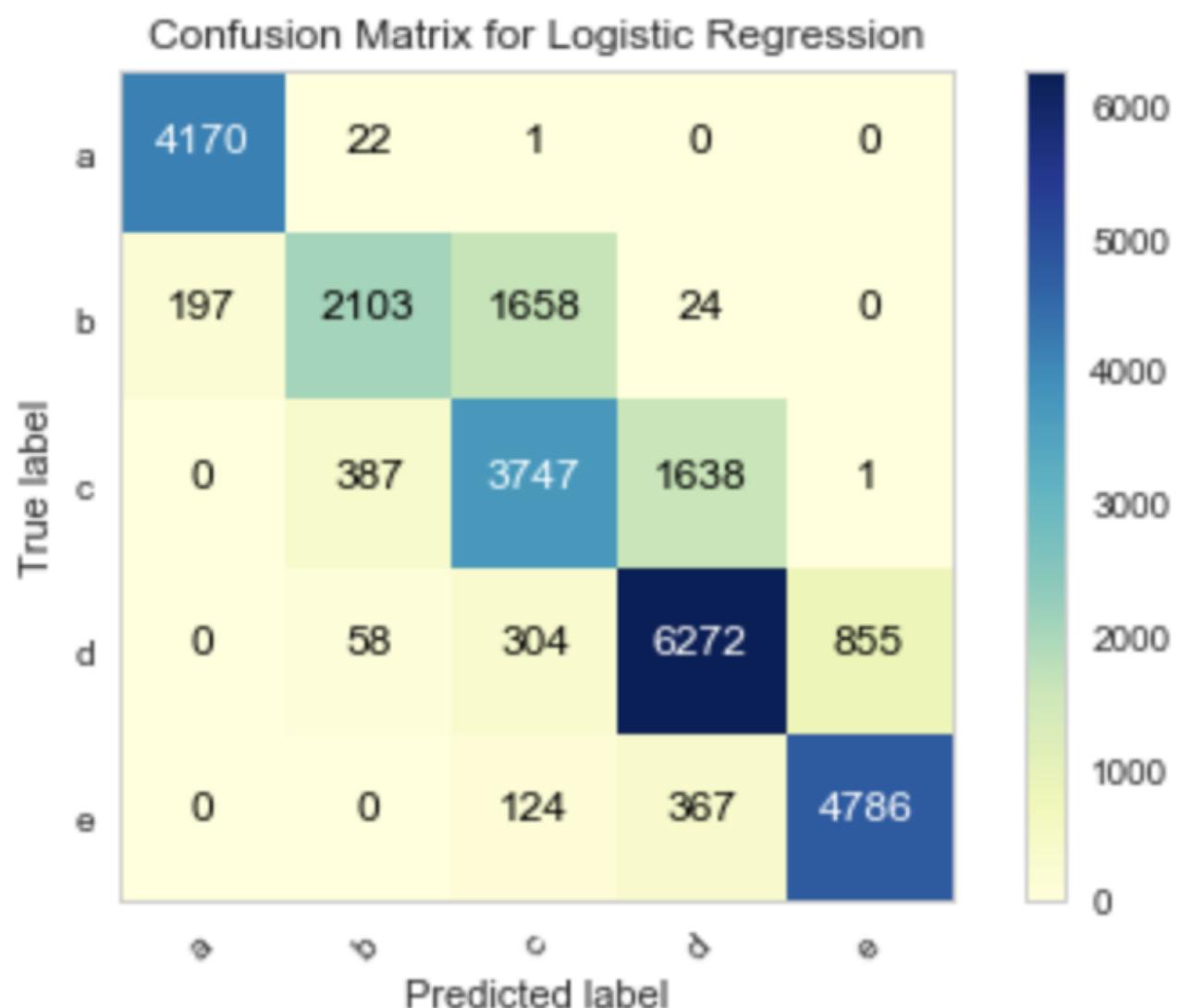


Figure 9: Confusion matrix for Logistic Regression [7]

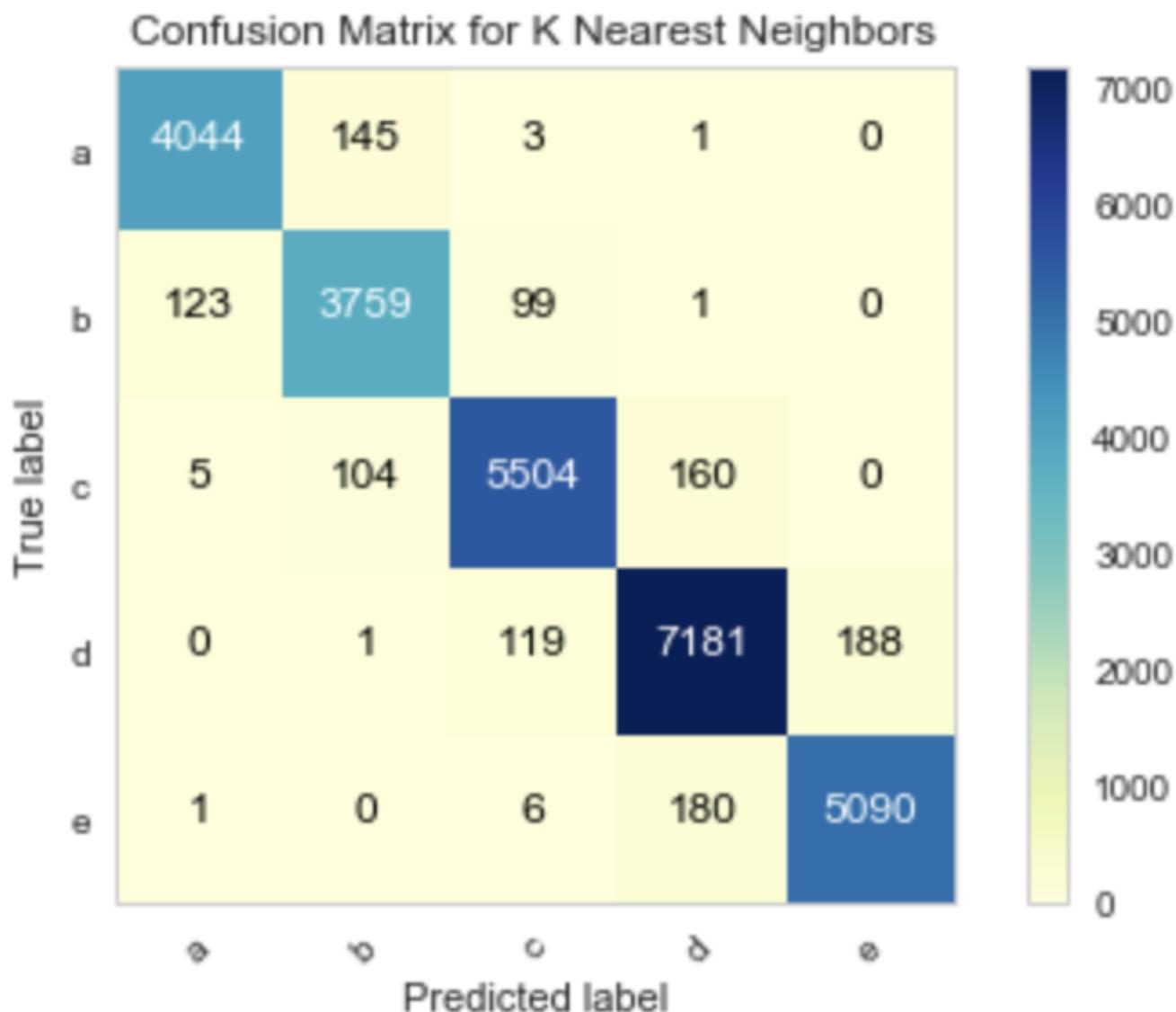


Figure 10: Confusion matrix for K Nearest Neighbors [7]

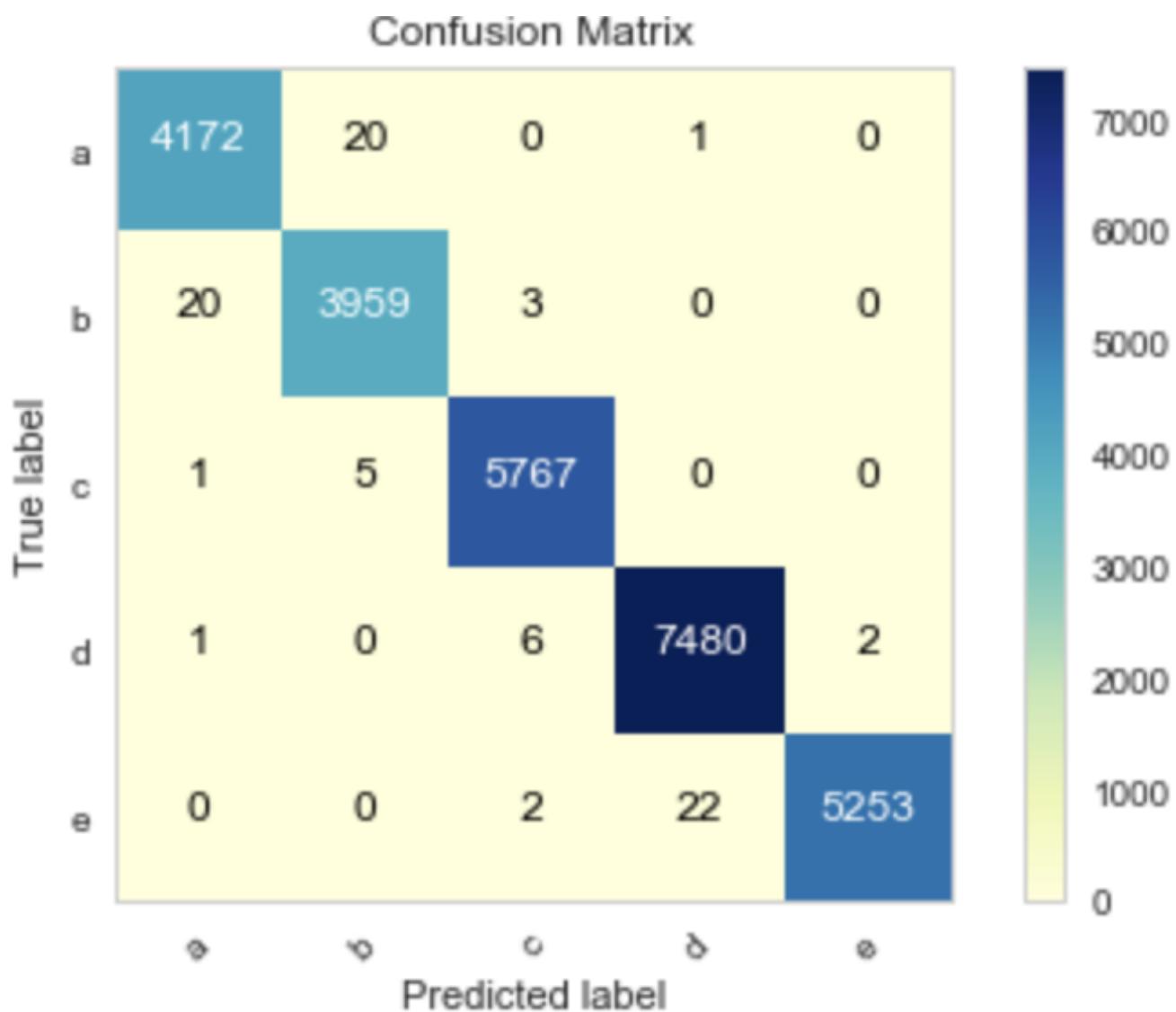


Figure 11: Confusion matrix for Random Forest [7]

```
bibtext report
```

```
=====
```

```
This is BibTeX, Version 0.99d (TeX Live 2016)
The top-level auxiliary file: report.aux
The style file: ACM-Reference-Format.bst
Database file #1: report.bib
```

```
bibtext _ label error
```

```
=====
```

```
bibtext space label error
```

```
=====
```

```
bibtext comma label error
```

```
=====
```

```
latex report
```

```
=====
```

```
[2017-12-04 12.23.02] pdflatex report.tex
This is pdfTeX, Version 3.14159265-2.6-1.40.17 (TeX Live 2016) (preloaded format=pdflatex)
Missing character: ""
bookmark level for unknown defaults to 0.
The anchor of a bookmark and its parent's must not be the same. Added a new anchor.
Typesetting of "report.tex" completed in 2.0s.
```

```
=====
```

```
Compliance Report
```

```
=====
```

```
name: Vegi, Karthik
hid: 231
paper1: Oct 29 17 100%
paper2: 100%
```

```
yamlcheck
```

```
=====
```

```
wordcount
```

```
-----  
17  
wc 231 project 17 5507 report.tex  
wc 231 project 17 5499 report.pdf  
wc 231 project 17 203 report.bib  
  
find "  
-----  
  
passed: True  
  
find footnote  
-----  
  
passed: True  
  
find input{format/i523}  
-----  
  
4: \input{format/i523}  
  
passed: True  
  
find input{format/final}  
-----  
  
passed: False  
  
floats  
-----  
  
65: We then display the top 5 countries as a pie-chart and the 5  
countries are namely United States, France, Switzerland, Germany,  
and Spain as shown in Figure \ref{fig:Fig7}.  
67: \begin{figure}  
68: \includegraphics[width=1.0\textwidth]{images/fig7.png}  
70: \label{fig:Fig7}  
76: The top 5 countries with most fat content in the food items are  
Serbia, United States, Switzerland, Germany, and Sweden as shown  
in Figure \ref{fig:Fig8}.  
78: \begin{figure}  
79: \includegraphics[width=1.0\textwidth]{images/fig8.png}  
81: \label{fig:Fig8}  
84: The top 5 countries with most saturated fat content in the food  
items are Serbia, United States, Germany, France and Switzerland
```

```

as shown in Figure \ref{fig:Fig9}
86: \begin{figure}
87: \includegraphics[width=1.0\textwidth]{images/fig9.png}
89: \label{fig:Fig9}
93: The top 5 countries with most trans-fat content in the food items
    are United States, Brazil, Canada, Australia, Russia, and Serbia
    as shown in Figure \ref{fig:Fig10} \\
95: \begin{figure}
96: \includegraphics[width=1.0\textwidth]{images/fig10.png}
98: \label{fig:Fig10}
101: The top 5 countries with most cholesterol content in the food
    items are United States, Canada, Portugal, Brazil, France, and
    Italy as shown in Figure \ref{fig:Fig11} \\
103: \begin{figure}
104: \includegraphics[width=1.0\textwidth]{images/fig11.png}
106: \label{fig:Fig11}
142: The nearest neighbour puts each attribute list as a data point in
    the n-dimensional space, given n the number of attributes
    \cite{book-tan}. Once we have the training examples, we take each
    test example and compute its distance to the training example
    classes and assign a class label \cite{book-tan}. Any of the
    popular distance measures among Euclidean distance, Manhattan
    distance, Minkowski distance and Mahalanobis distance can be used
    \cite{book-tan}. The k denotes the k closest points to the test
    example \cite{book-tan}. Figure \ref{fig:Fig1} shows the
    structure of the data \cite{book-tan}.
144: \begin{figure}
145: \includegraphics[width=1.0\textwidth]{images/fig1.png}
147: \label{fig:Fig1}
221: \subsubsection{Bi-variate box-plots} Bi-variate box-plots go
    beyond uni-variate box plots by showing the relationship between
    the predictor variable and the response variable \cite{book-tan}.
    We look at the bi-variate box-plots for each of the important
    predictor variables namely saturated fat, polyunsaturated fat,
    sugars and salt and the response variable which is the nutrition
    grade. Figure \ref{fig:Fig2} shows the bi-variate box plots. \\
223: \begin{figure}
224: \includegraphics[width=1.0\textwidth]{images/fig2.png}
226: \label{fig:Fig2}
232: Correlation between data objects is the measure of the linear
    relationship between the attributes of the object that are
    continuous variables \cite{book-tan}. Correlation analysis is the
    process of finding of the correlations between the different
    predictor variables and helps find collinearity problem
    \cite{book-shai}. The relationship could be either linear or non-
    linear given the data \cite{book-tan}. The correlation

```

coefficient can range anywhere between -1 and 1, where 1 indicates a positive correlation and -1 indicates negative correlation \cite{book-shai}. Correlation plot visually shows the correlation coefficient between the variables in a nicely laid out plot. Figure \ref{fig:Fig3} shows the correlation plot. \\
 234: \begin{figure}\\
 235: \includegraphics[width=1.0\textwidth]{images/fig3.png}\\
 237: \label{fig:Fig3}\\
 282: \item Logistic Regression: With \$l_2\$ penalty, the accuracy of logistic regression was 78.9%. Figure \ref{fig:Fig4} shows the confusion matrix. \\
 284: \begin{figure}\\
 285: \includegraphics[width=1.0\textwidth]{images/fig4.png}\\
 287: \label{fig:Fig4}\\
 290: \item K Nearest Neighbors: With k as 3, the accuracy of kNN was 95.74%. Figure \ref{fig:Fig5} shows the confusion matrix. \\
 292: \begin{figure}\\
 293: \includegraphics[width=1.0\textwidth]{images/fig5.png}\\
 295: \label{fig:Fig5}\\
 298: \item Random Forest: With number of trees as 100, the accuracy of random forest classifier was 99.68%. Figure \ref{fig:Fig6} shows the confusion matrix. \\
 300: \begin{figure}\\
 301: \includegraphics[width=1.0\textwidth]{images/fig6.png}\\
 303: \label{fig:Fig6}

```

figures 11
tables 0
includegraphics 11
labels 11
refs 11
floats 11
  
```

```

True : ref check passed: (refs >= figures + tables)
True : label check passed: (refs >= figures + tables)
True : include graphics passed: (figures >= includegraphics)
True : check if all figures are referred to: (refs >= labels)
  
```

```

Label/ref check
passed: True
  
```

When using figures use columnwidth
 [width=1.0\columnwidth]
 do not change the number to a smaller fraction

```
find textwidth
-----
68: \includegraphics[width=1.0\textwidth]{images/fig7.png}
79: \includegraphics[width=1.0\textwidth]{images/fig8.png}
87: \includegraphics[width=1.0\textwidth]{images/fig9.png}
96: \includegraphics[width=1.0\textwidth]{images/fig10.png}
104: \includegraphics[width=1.0\textwidth]{images/fig11.png}
145: \includegraphics[width=1.0\textwidth]{images/fig1.png}
224: \includegraphics[width=1.0\textwidth]{images/fig2.png}
235: \includegraphics[width=1.0\textwidth]{images/fig3.png}
285: \includegraphics[width=1.0\textwidth]{images/fig4.png}
293: \includegraphics[width=1.0\textwidth]{images/fig5.png}
301: \includegraphics[width=1.0\textwidth]{images/fig6.png}

passed: False
```

```
below_check
-----
```

```
bibtex
-----
```

```
label errors
```

```
bibtex errors
```

```
This is BibTeX, Version 0.99d (TeX Live 2016)
The top-level auxiliary file: report.aux
The style file: ACM-Reference-Format.bst
Database file #1: report.bib
```

```
bibtex_empty_fields
```

```
entries in general should not be empty in bibtex
```

```
find ""
```

```
passed: True
```

```
ascii
```

```
non ascii found 8211
```

```
=====
The following tests are optional
=====
```

```
Tip: newlines can often be replaced just by an empty line
```

```
find newline
```

```
passed: True
```

```
cites should have a space before \cite{} but not before the {
```

```
find cite {
```

```
passed: True
```

Big Data Applications in the Travel Industry and its Potential in Improving Travel Accessibility

Weixuan Wang
Indiana University Bloomington
School of Public Health
Bloomington, Indiana 47405
wangweix@indiana.edu

ABSTRACT

This is my abstract

KEYWORDS

i523, HID234

1 INTRODUCTION

Big Data Applications in the Travel Industry and its Potential in Improving Travel Accessibility. Here is the introduction [1].

2 TRANSPORTATION AND BIG DATA

3 TRAVEL AND BIG DATA

4 PROMISE OF BIG DATA IN TRAVEL ACCESSIBILITY

5 CONCLUSION

Put here an conclusion. Conlcusions and abstracts must not have any citations in the section.

ACKNOWLEDGMENTS

The authors would like to thank Dr. Gregor von Laszewski for his support and suggestions to write this paper.

REFERENCES

- [1] Ian Editor (Ed.). 2007. *The title of book one* (1st. ed.). The name of the series one, Vol. 9. University of Chicago Press, Chicago. <https://doi.org/10.1007/3-540-09237-4>

```
bibtext report
```

```
=====
This is BibTeX, Version 0.99d (TeX Live 2016)
The top-level auxiliary file: report.aux
The style file: ACM-Reference-Format.bst
Database file #1: report.bib
Warning--empty chapter and pages in editor00
(There was 1 warning)
```

```
bibtext _ label error
```

```
=====
report.bib:22:@Inbook{las_gergor00,
```

```
bibtext space label error
```

```
=====
bibtext comma label error
```

```
=====
latex report
```

```
[2017-12-04 12.23.13] pdflatex report.tex
```

```
This is pdfTeX, Version 3.14159265-2.6-1.40.17 (TeX Live 2016) (preloaded format=pdflatex)
Typesetting of "report.tex" completed in 0.8s.
```

```
./README.yml
```

```
9:81      error    line too long (86 > 80 characters)  (line-length)
24:81     error    line too long (119 > 80 characters)  (line-length)
25:81     error    line too long (125 > 80 characters)  (line-length)
25:125    error    trailing spaces  (trailing-spaces)
26:81     error    line too long (118 > 80 characters)  (line-length)
26:118    error    trailing spaces  (trailing-spaces)
32:28     error    trailing spaces  (trailing-spaces)
33:81     error    line too long (91 > 80 characters)  (line-length)
38:81     error    line too long (82 > 80 characters)  (line-length)
38:82     error    trailing spaces  (trailing-spaces)
39:81     error    line too long (82 > 80 characters)  (line-length)
39:82     error    trailing spaces  (trailing-spaces)
40:81     error    line too long (87 > 80 characters)  (line-length)
41:81     error    line too long (86 > 80 characters)  (line-length)
41:86     error    trailing spaces  (trailing-spaces)
```

```
42:81    error    line too long (84 > 80 characters)  (line-length)
42:84    error    trailing spaces  (trailing-spaces)
43:81    error    line too long (86 > 80 characters)  (line-length)
43:86    error    trailing spaces  (trailing-spaces)
45:1     error    trailing spaces  (trailing-spaces)
51:12   error    too many spaces after colon  (colons)
51:81    error    line too long (108 > 80 characters)  (line-length)
```

Compliance Report

```
name: Weixuan Wang
hid: 234
paper1: Oct 22 2017 100%
paper2: Nov 9 2017 100%
project: 10%
```

```
yamlcheck
```

```
wordcount
```

```
1
wc 234 Project 1 125 content.tex
wc 234 Project 1 131 report.pdf
wc 234 Project 1 144 report.bib
```

```
find "
```

```
passed: True
```

```
find footnote
```

```
passed: True
```

```
find input{format/i523}
```

```
4: \input{format/i523}
```

```
passed: True
```

```
find input{format/final}
```

```
passed: False
```

```
floats
```

```
figures 0
tables 0
includegraphics 0
labels 0
refs 0
floats 0
```

```
True : ref check passed: (refs >= figures + tables)
```

```
True : label check passed: (refs >= figures + tables)
```

```
True : include graphics passed: (figures >= includegraphics)
```

```
True : check if all figures are refered to: (refs >= labels)
```

```
Label/ref check
```

```
passed: True
```

```
When using figures use columnwidth
```

```
[width=1.0\columnwidth]
```

```
do not change the number to a smaller fraction
```

```
find textwidth
```

```
passed: True
```

```
below_check
```

```
bibtex
```

```
label errors
```

```
22: las_gergor00: do not use underscore in labels:
```

```
bibtex errors
```

```
This is BibTeX, Version 0.99d (TeX Live 2016)
The top-level auxiliary file: report.aux
The style file: ACM-Reference-Format.bst
Database file #1: report.bib
Warning--empty chapter and pages in editor00
(There was 1 warning)
```

```
bibtex_empty_fields
```

```
entries in general should not be empty in bibtex
```

```
find ""
```

```
3: author = "",
```

```
15: chapter = "",
```

```
16: pages = "",
```

```
17: number = "",
```

```
18: type = "",
```

```
19: month = "",
```

```
20: note = "",
```

```
36: chapter = "",
```

```
37: pages = "",
```

```
38: number = "",
```

```
39: type = "",
```

```
40: month = "",
```

```
41: note = "",
```

```
passed: False
```

```
ascii
```

```
=====
The following tests are optional
=====
```

```
Tip: newlines can often be replaced just by an empty line
```

```
find newline
```

```
=====
passed: True
```

```
cites should have a space before \cite{} but not before the {
```

```
find cite {
```

```
=====
passed: True
```

The Impact of Clinical Trial Results on Pharmaceutical Stock Performance

Tiffany Fabianac

Indiana University

Bloomington, Indiana 47408, USA

tifabi@iu.edu

ABSTRACT

While many relate stock market trading to gambling, successful traders have turned stock picking into a science. The likes of Warren Buffet tell us that successful stock buying is all in the research. So what kind of research aids in the prediction of companies within the highly volatile pharmaceutical market? The use of available, open-source APIs and Google Alerts are used to explore if clinical trial results can directly impact stock performance in small, mid, and large cap pharmaceutical companies. Key words and/or phrases in results and related news articles are identified as possible predictors of market effect. As well as a comparison to already established analyst ratings from Barclays, Goldman, Morningstar, or others which have already been shown to impact stock performance.

KEYWORDS

Big Data, HID313, i523, Stock Market, Pharmaceutical

1 INTRODUCTION

A “stock” is a piece of ownership in a company. Offering stocks for sale provides capital to the selling company in exchange for a stake in the company. A stock market is a collection of exchanges where trading of stocks takes place [5]. Evidence of early stock markets date back to the fourteenth century with the offering of state loan stocks throughout Italy. Even prior to the organization of stock markets, price fluctuations for goods such as wheat and barley were tracked by early economists. The first “modern” stock market appeared in Amsterdam in the seventeenth century where the volume of stocks traded and the fluidity in which they were traded reached a new high [2].

The biggest stock markets in the world are currently the New York Stock Exchange (NYSE), the National Association of Securities Dealers Automated Quotations (NASDAQ), and the London Stock Exchange. NYSE did stuff... NASDAQ began as an all-electric equities exchange in 1971 and today provides trading, technology, and information services for financial markets. Today [8].

Throughout the history of markets, prices have been tracked and insightful traders have attempted to predict and capitalize on price fluctuation. The age of computers opened new doors for stock analysis and trend prediction to facilitate capital gains for traders. Financial companies like Goldman Sachs and JPMorgan Chase & Co. have hired mathematicians, statisticians, and trade analysts since the early days of trading in an effort to predict the market in a consistent manner. Once an algorithm is established and used consistently the algorithm itself but be considered as a variable that could effect the prediction outcome [4].

A major complexity in creating algorithms for the stock market is that the market tends to follow the erratic emotions and feelings

of humans. If computers were running the market, making trade decisions based on logic and reason, then the market would be much more stable. The volatility of human emotions about money and stocks creates tremendous volatility in the market. The revolution of social media has provided a means of measuring the mood of possible traders [1].

An Application Programming Interface (API) acts as the middle-man between the requesting service and the performing service. When a user or system submits a request the request is passed to the API which translates it for the processing system then returns the results in a receivable format.

1.1 Pharmaceutical Sector

The pharmaceutical industry

Like the financial sector trying to predict the stock market, the pharmaceutical industry has devoted resources to developing prediction algorithms and machine learning systems. The efforts of drug manufacturers are to create a system that consistently predicts or aids in identifying drug targets.

1.2 Clinical Trials

A clinical trial is a planned experiment involving patients with the intent to elucidate an appropriate or effective treatment option(s) for the population of patients afflicted with the same medical condition. A big concern with clinical trials is that inferences are made for the entire population of patients from a relatively small sample size [9]. One of the first clinical trials recorded was carried out in the eighteenth century to evaluate six treatments on twelve patients with scurvy. Two patients that were given oranges and lemons recovered very quickly. Fisher introduced the concept of randomization in the nineteenth century [3].

Clinical trials have four defined phases. Phase I trials identify how well a drug is tolerated by determining the maximally tolerated dose (MTD) on a very small sample size. Phase I trials have very simple experimental designs as the only intent is to examine toxicity. Phase II explores biological activity or effect on a small patient sample size. The design of a Phase II trial is dependent on the design of the Phase I trial as both share the intent to evaluate adverse events. Phase III trials follow the design of Phase II trials but on a bigger sample size with the intent to solidify a treatment's effectiveness in clinical practice [3].

Clinical trial designs have been very slow to evolve due to

1.3 Established Analyst Ratings

1.4 Data Resources

NASDAQ's website provides historical stock performance data that can be exported as a Comma-Separated Values (CSV) file. The

disadvantage of NASDAQ's free export service is that each stock must be exported separately. The free quote service can be accessed at <http://www.nasdaq.com/quotes/>. NASDAQ provides API services for subscribers starting at \$5,000 per year [7]. Access to NASDAQ's API services can also be granted through corporate sponsorship. NASDAQ's free CSV export services were used to collect initial project data. In example, the stock history for Celsion Corporation during the week of August 21, 2017 is shown.

```
date ,close ,volume ,open ,high ,low
2017/08/25,1.3700,179097.0000,1.3600,1.4100,1.3000
2017/08/24,1.3600,149832.0000,1.3100,1.3600,1.2810
2017/08/23,1.3100,223451.0000,1.2500,1.3300,1.2430
2017/08/22,1.2800,164594.0000,1.3200,1.3200,1.2400
2017/08/21,1.3300,169037.0000,1.3300,1.3700,1.2800
```

Exports such as this one offered by NASDAQ and API interfaces for stock data are provided by numerous companies. The Yahoo! Finance API is explored below and the Google Finance API was used to perform the stock data extraction for the analysis presented. Additional resources such as stock tracking apps and free exports are available. CSV exports such as the one listed above can be downloaded from Google Finance, Yahoo! Finance, and many others. This publication does not provide a complete list of available resources, but attempts to present a few for comparison.

Python.org provides a python module to pull stock data from Yahoo! Finance [10]. The package can be installed through Git by cloning the Git directory where the package is available: <https://github.com/finance.git>. To install the python package without Git the tape archive can be downloaded from <https://pypi.python.org/pypi/yahoo-finance>. While Yahoo! Finance is a great resource, the API does not function consistently. To install the python package without Git, the tape archive can be downloaded from <https://pypi.python.org/pypi/yahoo-finance>. Tape archives allow for compression of multiple files which can be restored to their original format using the tar command in the command line [6]. Apply the tar options: z - filter archive through gzip, x - extract an archive file, and f - filename of archive, use "cd" to change the current working directory, and then install the python module using the package management command "pip":

```
tar -zxf yahoo-finance-1.4.0.tar.gz
cd yahoo-finance
pip install yahoo-finance
```

2 METHODS

2.1 Data Collection

Data collection was initiated with the use to Google Alerts. Google allows for alerts to be configured from <https://www.google.com/alerts>. Gmail users can configure these alerts to be sent through email when news or other types of articles are released to the web. The Google Alerts for this project were: "Phase III Trial", "Phase 3 Trial", and "Meets Primary End Point". When these phrases are detected by google, the link to the webpage and a short description are sent via email to the configured email address. On busy days, an excess of 100 alerts were received for this data. On slow days, only a couple alerts were received. Only very infrequently were no messages received.

To collect data from the received Google Alerts without too much manual clicking, Gmail has an available API which allows users to pull data from a Gmail account. To start using the Gmail API a user must first configure their Authentication credentials through Google's developer console. Once credentials are received the form of a JSON file, the Google Client Library can be installed using pip to install google-api-python-client. The Google Development team has provided a quickstart file which facilitates the first authentication run. Running this quick start guide will open a browser window and prompt the user to log into a Gmail account. The user then accepts the authorization and can run the Gmail API from command line or other compilers.

Headlines of the received alerts, usually the title of the article and the first couple of lines, are referred to as "Snippets" by Google's Gmail API. This project pulled only the Snippets and the date from the Google Alerts. The Snippets do not contain the whole article but may still provide enough evidence of sentiment for further analysis and prediction of the associated stock. Unfortunately, no solution was identified for extracting the appropriate stock symbols from the Snippets so this task had to be performed manually.

The Python code that calls the Gmail API and writes a csv from the data starts, after calling all needed libraries, with defining the scope of the authorization. Google mail can be opened with a Read-only or Modify authentication. Next, the credentials are established by the JSON file received during the API authentication setup. This JSON must be saved in the same directory as the code being run. The code sets the variables for User ID and Label then runs an execution command calling the Messages.List API, which looks like this:

```
GMAIL.users().messages().list(userId='me', labelIds=[INBOX], q='from:googlealerts-noreply@google.com before:2017/11/24').execute()
```

Google has defined the user ID "me" as the global for the authenticated account in use. The label ID "INBOX" designates that the messages will be pulled from the inbox folder, but any other folder could be called here as well as a collection of labels that Google has defined such as "UNREAD". The "q" designates a query. The query will return only messages from the Google Alerts email address which have been received by the twenty-fourth of November 2017. This data was selected so that all returned records would have atleast five market days of stock prices to compare. This execution returns a dictionary which contains message IDs for all the messages that matched the query.

The next step is to "get" the messages with the use of the Messages.Get API. While looping through the dictionary of message ID from the defined query, the script retrieves the Date and Snippet for each. Additional options could return the Sender, Receiving Email, Email body, among others. The syntax is shown here:

```
GMAIL.users().messages().get(userId='me', id=m_id).execute()
```

The user ID is the same as described previously with the ID being the current message ID within the loop. This execute command returns a dictionary which is parsed from "payload" to "headers" to extract the Date. The Snippet is also grabbed from the message dictionary and along with the Date, passed to a final list to be written to a .csv file.

[Figure 1 about here.]

Figure 1 shows the entire code to extract Google Alerts data using the Google provided Gmail API.

The Python package pandas is an incredible resource. The Pandas package has a resource for getting stock market data from free online sources such as Yahoo! mentioned above and Google. To install this package through Git, simply clone the directory, use the “Change Directory” command “cd” to change the current working directory, and installing the python module as follows:

```
$ git clone git://github.com/pydata/pandas-datareader.git  
$ cd pandas-datareader  
$ python setup.py install
```

If the Python setup returns the error: “python: command not found” run the following with the path to the python installation:

```
$ PATH="$PATH:/c/Python27"
```

Pandas-datareader and many other packages can also be installed via pip. In example, many additional packages are needed to run a python script using pandas-datareader. These packages can be configured all at once or one at a time as follows:

```
pip -m install --user numpy scipy matplotlib ipython jupyter pandas sympy  
nose urllib3 chardet idna
```

Unlike the NASDAQ export, using Google as a data source for pandas-datareader requires each attribute to be called separately. This means calling the Close Price, Open Price, High Price, etc individually and joining them through code. Also, unlike NASDAQ’s export but this time in a positive light, multiple tickers can be passed together. This allows for all historical data to be pulled for many stocks with a single code.

[Figure 2 about here.]

Figure 2 shows the code to combine the data produced by the Google Alert mining and available historic stock price data.

2.2 Data Analysis

3 RESULTS

3.1 Comparison to Established Analyst Ratings

4 CONCLUSION

ACKNOWLEDGMENTS

The author would like to thank Dr. Gregor von Laszewski and the teaching assistants of the Fall 2017 i523 course for their support and suggestions to write this paper.

REFERENCES

- [1] Johan Bollen, Huina Mao, and Xiaojun Zeng. 2011. Twitter mood predicts the stock market. *Journal of Computational Science* 2, 1 (2011), 1 – 8. <https://doi.org/10.1016/j.jocs.2010.12.007>
- [2] F. Braudel. 1982. *Civilization and Capitalism, 15th-18th Century, Vol. II: The Wheels of Commerce*. University of California Press. <https://books.google.com/books?id=WPDbSXQsvGIC>
- [3] L.M. Friedman, C. Furberg, and D.L. DeMets. 1998. *Fundamentals of Clinical Trials*. Springer. <https://books.google.com/books?id=yzxT0Zh3X3IC>
- [4] Thomas Hellstrom and Kenneth Holmstrom. 1997. *Predict the stock market*. techreport. Department of Mathematics and Physics, Mälardalen University, Sweden.
- [5] Investopedia. 2017. Stock Market. Website. (09 2017). <https://www.investopedia.com/terms/s/stockmarket.asp?lgl=rira-layout>
- [6] LINFO. 2006. The tar Command. website. (07 2006). <http://www.linfo.org/tar.html>
- [7] NASDAQ. 2017. NASDAQ DataOnDemand Subscription Plans. Website. (2017). <https://www.nasdaqdod.com/Shop/ProductConfig.aspx?product=webservices&service=NASDAQDataOnDemand>
- [8] NASDAQ. 2017. NASDAQ’s Story. website. (2017). <http://business.nasdaq.com/discover/nasdaq-story/index.html>
- [9] S.J. Pocock. 2013. *Clinical Trials: A Practical Approach*. Wiley. <https://books.google.com/books?id=TxbTBQAAQBAJ>
- [10] Python.org. 2016. yahoo-finance 1.4.0. Website. (11 2016). <https://pypi.python.org/pypi/yahoo-finance>

LIST OF FIGURES

- 1 The Google API Python code calls the Gmail APIs Messages.list which lists reduced properties of Gmail messages and Messages.Get which returns the messages themselves. Lists is used to query the messages that are wanted based on the defined criteria: userId=me, labelIds=INBOX], q=from:googlealerts-noreply@ google.com. Get then retrieves the messages identified in using List and returns the messages content for Date and Snippet. 5
- 2 This Python script takes in the Date, Stock Ticker Symbol, and Snippet from the Google API .csv that was produced using both manual mining of the stock symbols and the python script provided for getting the Date and Snippet from Gmail. This code returns a modified .csv which lists an "L" for stocks that did not increase by 10% in five days and a "W" for stocks that increased by atleast 10%. It also prints the stocks that increased by atleast 10% along with the highest price over 5 days, the starting price on the day that the Google Alert was received, and the percent change. 5

your code here

Figure 1: The Google API Python code calls the Gmail APIs `Messages.list` which lists reduced properties of Gmail messages and `Messages.Get` which returns the messages themselves. `Lists` is used to query the messages that are wanted based on the defined criteria: `userId=me, labelIds=INBOX], q=from:googlealerts-noreply@ google.com`. `Get` then retrieves the messages identified in using `List` and returns the messages content for Date and Snippet.

your code here

Figure 2: This Python script takes in the Date, Stock Ticker Symbol, and Snippet from the Google API .csv that was produced using both manual mining of the stock symbols and the python script provided for getting the Date and Snippet from Gmail. This code returns a modified .csv which lists an “L” for stocks that did not increase by 10% in five days and a “W” for stocks that increased by atleast 10%. It also prints the stocks that increased by atleast 10% along with the highest price over 5 days, the starting price on the day that the Google Alert was received, and the percent change.

```
bibtext report
```

```
=====
This is BibTeX, Version 0.99d (TeX Live 2016)
The top-level auxiliary file: report.aux
The style file: ACM-Reference-Format.bst
Database file #1: report.bib
Warning--empty address in Braudel
Warning--empty address in Friedman
Warning--can't use both author and editor fields in Pocock
Warning--empty address in Pocock
(There were 4 warnings)
```

```
bibtext _ label error
```

```
=====
bibtext space label error
```

```
=====
bibtext comma label error
```

```
=====
latex report
```

```
=====
[2017-12-04 12.23.38] pdflatex report.tex
This is pdfTeX, Version 3.14159265-2.6-1.40.17 (TeX Live 2016) (preloaded format=pdflatex)
bookmark level for unknown defaults to 0.
The anchor of a bookmark and its parent's must not be the same. Added a new anchor.
Typesetting of "report.tex" completed in 1.1s.
./README.yml
69:81     error    line too long (83 > 80 characters)  (line-length)
```

```
=====
Compliance Report
```

```
=====
name: Tiffany Fabianac
hid: 313
paper1: Oct 31 2017 100%
```

```
paper2: 100%
project: 90%
```

```
yamlcheck
```

```
wordcount
```

```
(null)
wc 313 project (null) 2437 report.tex
wc 313 project (null) 2632 report.pdf
wc 313 project (null) 410 report.bib
```

```
find "
```

```
108: Data collection was initiated with the use to Google Alerts.
Google allows for alerts to be configures from
https://www.google.com/alerts. Gmail users can configure these
alerts to be sent through email when news or other types of
articles are released to the web. The Google Alerts for this
project were: "Phase III Trial", "Phase 3 Trial", and
"Meets Primary End Point". When these phrases are detected by
google, the link to the webpage and a short description are sent
via email to the configured email address. On busy days, an
excess of 100 alerts were received for this data. On slow days,
only a couple alerts were received. Only very infrequently were
no messages received.
```

```
158: $ PATH="$PATH:/c/Python27"
```

```
passed: False
```

```
find footnote
```

```
passed: True
```

```
find input{format/i523}
```

```
4: \input{format/i523}
```

```
passed: True
```

```
find input{format/final}
-----
passed: False

floats
-----
130: \begin{figure}[htb]
134: \caption{The Google API Python code calls the Gmail APIs  
Messages.list which lists reduced properties of Gmail messages  
and Messages. Get which returns the messages themselves. Lists is  
used to query the messages that are wanted based on the defined  
criteria: userId=me, labelIds=INBOX], q=from:googlealerts-  
noreply@ google.com. Get then retrieves the messages identified  
in using List and returns the messages content for Date and  
Snippet.}\label{c:googleapi}
137: Figure \ref{c:googleapi} shows the entire code to extract Google  
Alerts data using the Google provided Gmail API.
173: \begin{figure}[htb]
177: \caption{This Python script takes in the Date, Stock Ticker  
Symbol, and Snippet from the Google API .csv that was produced  
using both manual mining of the stock symbols and the python  
script provided for getting the Date and Snippet from Gmail. This  
code returns a modified .csv which lists an ‘‘L’’ for stocks that  
did not increase by 10\% in five days and a ‘‘W’’ for stocks that  
increased by atleast 10\%. It also prints the stocks that  
increased by atleast 10\% along with the highest price over 5  
days, the starting price on the day that the Google Alert was  
received, and the percent change.}\label{c:stock}
180: Figure \ref{c:stock} shows the code to combine the data produced  
by the Google Alert mining and available historic stock price  
data.
```

```
figures 2
tables 0
includegraphics 0
labels 2
refs 2
floats 2
```

```
True : ref check passed: (refs >= figures + tables)
True : label check passed: (refs >= figures + tables)
True : include graphics passed: (figures >= includegraphics)
True : check if all figures are refered to: (refs >= labels)
```

```
Label/ref check
passed: True
```

```
When using figures use columnwidth
[width=1.0\columnwidth]
do not change the number to a smaller fraction
```

```
find textwidth
```

```
passed: True
```

```
below_check
```

```
bibtex
```

```
label errors
```

```
bibtex errors
```

```
This is BibTeX, Version 0.99d (TeX Live 2016)
The top-level auxiliary file: report.aux
The style file: ACM-Reference-Format.bst
Database file #1: report.bib
Warning--empty address in Braudel
Warning--empty address in Friedman
Warning--can't use both author and editor fields in Pocock
Warning--empty address in Pocock
(There were 4 warnings)
```

```
bibtex_empty_fields
```

```
entries in general should not be empty in bibtex
```

```
find ""
```

passed: True

ascii

The following tests are optional

Tip: newlines can often be replaced just by an empty line

find newline

passed: True

cites should have a space before \cite{} but not before the {

find cite {

passed: True

Big Data Applications in Predicting Hospital Readmissions

Tyler Peterson

Indiana University - School of Informatics, Computing, and Engineering

711 N. Park Avenue

Bloomington, Indiana 47408

typeter@iu.edu

ABSTRACT

Hospital readmissions occur when a patient is discharged from a hospital and subsequently readmitted to a hospital within a short time frame. Hospitals are held accountable and penalized for readmissions that occur within 30 days of the initial inpatient stay. In 2016, nearly 2,600 hospitals were penalized \$528 million collectively for readmissions. Machine learning is increasingly being used to build models that predict if a patient has a high probability of being readmitted, which allows hospital staff to prioritize resources around high-risk patients and potentially prevent the otherwise likely readmission. Healthcare providers possess every-growing stores of medical data that are essential for building accurate predictive models. While most of this information is private and not widely available for research, there are a few public datasets that researchers can use to build models and gain a better understand what kind of information is significant in the task of identifying high-risk patients. One such dataset includes over 100,000 patient admissions that occurred at 130 US hospitals between 1999 and 2008 and includes many features that can be used to build models. Open-source Python tools such as scikit-learn, pandas and matplotlib have tools necessary for preparing, modeling and visualizing data. These tools can be used to define algorithms that describe the problem of hospital readmissions by creating classifiers that assign samples based on the probability of readmission. Machine learning techniques such as logistic regression, support vector machines and decision trees are ideal for modeling data for classification problems, and these open-source tools include methods for assessing and optimizing the algorithms. The effectiveness of machine learning in classifying patients for risk of readmission is a growing topic of study and implementation of tools for assisting healthcare providers will likely continue to increase.

KEYWORDS

hid331, i523, Big Data, Hospital Readmissions, Machine Learning, Classification, Python

1 INTRODUCTION

Hospital readmissions are problematic for both patients and healthcare providers. Even a single hospital admission for a patient can be an inconvenient, expensive and anxiety-inducing major life event. For a patient to be subsequently readmitted to the hospital, the patient again experiences the negative aspects of being in a hospital, along with a diminished quality of life that accompanies a recurrent disease or medical issue. Healthcare providers are increasingly being held accountable and often penalized for an inability to keep recently discharged patients from being readmitted. It has been estimated that nearly 1 in 5 Medicare patients discharged from a

hospital will be readmitted within 30 days [4] The Hospital Readmission Reduction Program (HRRP), which originated in 2013 as a provision in the Affordable Care Act, serves as an example of an initiative that punishes hospitals for readmissions by administering financial penalties on hospitals with disproportionately high readmission rates among Medicare beneficiaries [1]. The HRRP levies a reduction in Medicare reimbursement, and uses the 'all-cause' definition for readmissions, which means that a subsequent hospital stay that occurs for any reason within 30 days of the initial stay counts against the hospital [1]. The program focuses on patients initially admitted with a heart attack, heart failure, pneumonia, chronic obstructive pulmonary disease, a coronary artery bypass graft procedure or a hip/knee replacement procedure [1]. If a hospital's risk-adjusted readmission rate is higher than the national average, then that hospital will be penalized. Further, the excessiveness of the rate is considered as well, ensuring that providers with the worst readmission rates have proportionately higher penalties [1]. In 2016, the US government penalized 79 percent of US hospitals, which amounts to 2,597 institutions [8]. The penalties for those readmissions, applied to the 2017 fiscal year reimbursements, amounted to \$528 million nationally, \$108 million higher than the previous year [8]. Effectively this means that the care provided to readmitted patients is uncompensated care, which still requires valuable resources such as medical supplies, pharmaceuticals, the occupancy of hospital beds and the attention of medical staff. HRRP has had the intended effect of bringing increased attention to readmissions, and some healthcare providers are leveraging their ever-increasing medical data stores to better understand their patients. Several organizations are using machine learning to identify high-risk patients. Assessing patients for the likelihood of readmission presents a binary classification problem, where a model's goal is to come to one of two conclusions on each case. The model analyzes each patient and the patient's accompanying attributes and concludes either that the patient will be readmitted or will not be readmitted.

1.1 Applying Machine Learning to Hospital Readmissions

There are several studies pertaining to the effectiveness of using machine learning to build predictive models that address this problem. A 2011 study conducted a systematic review of the topic and found 26 studies discussing predictive models related to hospital readmissions. These models were created using administrative claims data, electronic medical record (EMR) data, or a combination of each type of dataset [3]. Administrative claims data is primarily gathered for billing purposes and contains information about procedures, diagnoses, length of hospital stay and location of care [7]. The advantage of this type of data is that it typically describes large

populations and is inexpensive to acquire because it's already gathered for billing [4]. EMRs contain the basic information contained in administrative claim data, and also include lab data, image data and the results of various diagnostic tests, as well as social and behavioral information. Of the 26 studies reviewed by this paper, only 4 reported an area under the curve (AUC) value greater than 0.70, indicating that the other 22 models performed relatively poorly at classifying high-risk patients. Interestingly, 3 of the 4 studies with a moderately high AUC built models with clinical information found in EMRs in addition to administrative claims data, which suggests that the rich information available in EMRs adds discriminative power to the predictive models [4]. One study that demonstrates the power of incorporating EMR data was conducted at Mount Sinai Health System in New York, NY. Mount Sinai developed a model to predict readmissions among patients with heart failure, which is the top cause of readmission among Medicare beneficiaries [9]. To build the model, Mount Sinai leveraged their EMR system to mine 4,205 patient attributes, including 1,763 diagnosis codes, 1,028 medications, 846 laboratory measurements, 564 surgical procedures, and 4 types of vital signs. The study used a cohort of 1,068 patients, 178 of whom were readmitted within 30 days [9]. The model achieved a prediction accuracy rate of 83.19 percent and an AUC value of 0.78. Commenting on this outcome, Mount Sinai said that the model would benefit from the inclusion of several years of data from several different hospital sites [9]. In other words, even more data is needed to further improve the accuracy of the model.

2 ANALYSIS

Though the data used by institutions to build models is not widely available, there are a few public datasets that can be used by machine learning practitioners to better understand how predictive modeling techniques can be applied to the task of predicting readmissions. One such dataset comes from the Cerner Corporation's Health Facts database, which is comprised of comprehensive clinical EMR records voluntarily provided by hospitals across the United States [10]. Researchers extracted a subset of 101,766 encounters from the nearly 74 million records in the Health Facts database for the purpose of studying diabetic inpatient encounters. The admissions span 10 years from 1999 to 2008, and occurred at 130 different hospitals across the United States. The researchers used the following criteria to narrow down the dataset: 1) the encounter is an inpatient encounter 2) it was a diabetic encounter, meaning at least one diabetic diagnosis code was associated with the episode of care 3) The length of stay was between 1 and 14 days 4) the patient had at least one lab test and 5) the patient was administered at least one medication [10]. This dataset is now publicly available on the UCI Machine Learning Repository. The dataset contains 55 attributes, or features, that are potentially related to hospital readmissions, including diagnoses defined by ICD-9-CM codes, in-hospital procedures, hospital characteristics, individual provider information, lab data, pharmacy data, and demographic data, such as age, gender and race. Each patient encounter record also has a label indicating whether or not the patient was readmitted within 30 days. Since the dataset includes these labels, supervised machine learning techniques can be used, as opposed to unsupervised machine learning techniques. Further, since the model will need to

predict whether or not a patient will be readmitted within 30 days, this is a binary classification problem. The model will classify each patient encounter as highly likely of readmission within 30 days or not likely of readmission within 30 days.

3 OVERVIEW OF SUPERVISED MACHINE LEARNING

Several algorithms can be used for supervised classification problems, including logistic regression, support vector machines (SVM) and decision trees. An open-source Python library called scikit-learn has modules for training and evaluating models built using each of these three types of algorithms. Though there are several types of algorithms, there are several fundamentals of machine learning that apply to all predictive modeling techniques.

3.1 Minimization of Error

The goal of a machine learning algorithm is to minimize the error made in the predictions. The general form of this concept can be represented by the formula:

$$Y = f(x) + \epsilon$$

where Y is the actual outcome associated with the sample, X represents the attributes associated with each sample, $f(X)$ is a function that represents the systematic information X provides about Y , and ϵ is the error term describing the differences between the predicted value returned by $f(X)$ and the actual value represented by Y [6]. A perfect prediction means $f(X)$ equals Y and ϵ equals zero. In reality, the error term will rarely be zero, so each prediction yields a certain amount of error. The prediction accuracy for each sample is evaluated by this formula, and sum of the error terms from each evaluation represents the magnitude of error made by the model. The goal is make the sum of errors as low as possible [6]. The error term is minimized through optimization of $f(X)$, which is intended to describe the patterns that exist between the independent variables, represented by X , and the dependent variable, represented by Y . Said differently, the equation describes the relationship between the features and the label. The way that this function describes this relationship is through coefficient weights. Each feature in the dataset is paired with a numerical weight that accentuates or diminishes the impact of a feature on the predicted outcome. The way in which these coefficients can be interpreted differs by which algorithm is being used, but the intuition remains the same: the coefficients are adjusted to highlight the important features in the dataset. Once the coefficients are determined, the model has been fit to the data.

3.2 Training Set vs. Test Set

The coefficient weights of the model are defined by analyzing the samples in a dataset. In a practical sense, the value of a model depends on its ability to accurately predict the outcomes of new samples that were unseen at the time the model was determined [3]. A model that performs well when making predictions with new data is said to generalize well. A machine learning practitioner will want to have confidence in the model's ability to generalize before deploying the model to make predictions in real-time, and will not necessarily have a new dataset of previously unseen observations

to run through the model. To get around this, the original dataset is often split into two parts. The first part of the dataset is referred to as the training set and is used to determine the coefficient weights. The second part of the dataset is referred to as the test set, and this set is run through the model derived from the training set. The accuracy of the predictions on the test set is compared to the accuracy of the predictions on the training set to determine the extent to which the model generalizes [3]. A model that has high training accuracy, but low test accuracy, is said to be overfitting the data. This means that the model, in its efforts to minimize ϵ , has become too complex and focuses too closely on the samples in the training data set. By chasing patterns in the training data caused more so by random chance than by the true characteristics of X , the model no longer generalizes to the unseen samples in the test set [6][3]. An overfit model describes characteristics in the training data that are not in the test data, leading to poor predictions on the test set. A model can also underfit the data, which means the model is failing to capture the relationship between Y and X and will likely perform poorly on both the training and test datasets.

3.3 The Bias/Variance Trade-off

Bias and Variance are two important components related to training models using machine learning. Variance describes the extent to which a model changes due to small adjustments in the training data. Since the training data used to fit a model can vary, it is reasonable to expect that a model will change when different samples are selected into the training dataset, but ideally the model changes only slightly [6]. If a model is quite complex and is overfitting the training data, then slight changes in the training samples can have a large effect on the coefficient weights. Low variance is preferable [6]. Bias refers to the error that occurs when trying to describe a phenomenon using a derived model. For example, if a machine learning technique assumes a linear relationship between the independent and dependent variables, but the relationship is highly non-linear, then the selected technique will result in high bias [6]. The chosen approach to the problem is not well-suited to the task at hand and will make predictions that are far from reality. As the method for training a model becomes more complex and able to fit to the perceived important important in the training data, variance will increase and bias will decrease. The model will become more flexible and therefore more sensitive to variations in the training data, but will reduce bias by reducing the prediction error. The important part of the relationship between these two components is that as a model becomes more complex, the bias decreases more rapidly than the variance increases, so the trade-off of increasing variance while decreasing bias leads to a net gain in improvement of the model [6]. However, there is a point at which the model becomes too complex and the net gain begins to disappear. Increased model complexity leads to significantly higher variance without appreciable improvement in bias [6].

3.4 Model Evaluation

Several statistics can be used for evaluating model accuracy. For classification problems, a basic technique for evaluation is the confusion matrix. Figure 1 shows the general framework of a confusion matrix which shows the counts of each type of prediction and the

accuracy of that prediction. A true positive is an outcome that is predicted to be positive and is positive in reality [2]. A true negative is an outcome that is predicted to be negative and is negative in reality [2]. These are the preferred responses. In the context of hospital readmissions, a true positive is a prediction that a patient in the test dataset, according to the trained model, will be readmitted to the hospital within 30 days, and this occurs in reality. A true negative is a prediction that a patient in the test dataset will not be readmitted, and this occurs in reality. On the other hand, a false positive is an outcome that is predicted to be positive but is negative in reality [2]. A false negative is an outcome that is predicted to be negative but is positive in reality. These are errors in prediction [2]. If a healthcare provider acts on a false positive, that could mean that a patient, who without intervention would not have been readmitted within 30 days, received resources and attention that were not necessary. In the case of a false negative, this means a patient who eventually did get readmitted within 30 days could have benefited from additional attention and resource from a healthcare team. These four components - true positives, true negative, false positives, and false negatives - can be combined to create more nuanced metrics. Two of those metrics are sensitivity and specificity. Sensitivity refers to the true positive detection rate. This is the percentage of positive occurrences that are successfully identified [2]. Specificity is the true negative detection rate. This is the percentage of negative occurrences that are successfully identified [2]. In the context of readmissions, low sensitivity means many patients who eventually get readmitted are not predicted to be high-risk before the readmission occurs. Low specificity means that many patients who would not otherwise be readmitted are predicted to be readmitted. There is a trade-off between sensitivity and specificity, and an improvement in one often causes the other to worsen. Preference toward sensitivity or specificity often depends on the cost of incorrect predictions. A patient who otherwise would not be readmitted who is predicted to be high-risk is the type of case that will incur unnecessary resources. While this requires healthcare providers to invest resources that are not needed, the readmission is nevertheless avoided and there are potentially other benefits achieved by the hospital, such as increased satisfaction of the patient and their family. On the other hand, a patient who eventually gets readmitted but was not identified beforehand will likely be costly to a hospital in a couple ways. The provider must dedicate resources to stabilizing and healing the patient, while also incurring penalties if this type of readmission occurs frequently. If the expense of an unexpected readmission is higher than the expense of deploying unnecessary resources to low-risk patients, then a model that favors higher sensitivity at the expense of lower specificity is preferable. Sensitivity and specificity can be assessed in tandem by the receiver operating characteristic (ROC) curve, which is quite useful for evaluating supervised classification models. The ROC curve plots the true positive rate against the false positive rate (100 minus the true negative rate) for varying decision thresholds. This illustrates the trade-off between sensitivity and specificity and can provide guidance on which decision threshold is appropriate for the task [2]. ROC curves are often leveraged to evaluate the performance of models by calculating the area under the ROC curve, also known as the AUC. The goal is to maximize the

AUC value, and that value points to the optimal balance between sensitivity and specificity [2].

3.5 Data Preparation

The data needs to be cleaned up so that the model can evaluate the features properly. For logistic regression and support vector machines, the data must be numerical. Columns such as 'num_procedures' and 'num_lab_procedures' contain numerical data, and are ready to use as-is. Other columns such as 'A1Cresult' includes values such as ", ", and ". These values must be encoded to work properly. The Python library Pandas has a function called 'get_dummies' will create columns that include two values, 0's and 1's, for each unique value in a column. In the case of the column 'A1Cresult', this process will yield X columns, one for each unique value. For each observation, a 1 will appear in the column corresponding to value of the original feature. Consideration must be given to problems that can arise from transforming categorical variables in this manner.

For decision tree methods, the get dummies, change multi level categorical variables delete columns and why - all 1's, all 0's redundancy between med specialty and diagnosis categories cleanup of diagnosis codes multicollinearity scale 0 to 1

only include one admit per patient, samples need to be independent

4 LOGISTIC REGRESSION

4.1 Logistic Regression - Intuition

Logistic regression models the probability that a sample belongs to a certain class given the feature values of the sample [3]. This probability can be represented as:

$$p(x) = Pr(Y = 1|X)$$

In the context of predicting hospital readmissions, this translates to 'the likelihood that a patient will be readmitted within 30 days of discharge given the patient's characteristics.' To determine the probability, logistic regression utilizes the logistic function, which takes in the coefficient weights and feature responses for each sample and returns a the probability - a number between 0 and 1[3]. In the case of logistic regression involving multiple features, the model takes the form:

$$f(x) = \frac{e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}{1 + e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}$$

The model is fit to the data by adjusting the coefficient weights using a method called maximum likelihood. The intuition of this process is that the estimates for the coefficients are set such that the predicted probability of a certain outcome corresponds as closely as possible to the actual label of that sample. This means that the ideal coefficient weights, when plugged into the logistic function, return a number close to one for the readmitted patients and a number close to zero for the patients not readmitted [3].

4.2 Logistic Regression - Preprocessing

4.3 Logistic Regression - Execute Analysis

fitting, changing parameters

4.4 Logistic Regression - Evaluate Analysis

classification reports, f score predict method predict proba decision function

L2 vs L1

5 DECISION TREES

5.1 Decision Trees - Intuition

Tree-based machine learning techniques are conceptually simple strategies for modeling data. Decision trees segment the observations based on a series of if/else statements, ultimately leading to a class prediction [3]. The sequence of if/else rules can be visualized as an upside-down tree. The first node of the tree is called the root node and contains all observations. The root node is split into two nodes by applying an if/else rule to a specific feature in the observations.[5]. An example if/else rule might be 'patient received hemoglobin A1c test', and each observations is separated into different nodes depending on whether than observations received a hemoglobin A1c test. Each node that splits into additional nodes is referred to as a nonterminal node [5]. A node that does not further split that data is referred to as a terminal node, which predicts a class for all observations that have followed each branch into that particular node [5]. Decision trees use a top-down, greedy approach to segment data called recursive binary splitting [?]. The 'top-down' part of the term means that all observations are grouped together at the start of the tree and then subsequently split into separate branches. The approach is referred to as greedy because the if/else question chosen at each branch is the best split possible at that particular branch. This is a potentially short-sighted strategy, hence the name greedy, because it is possible that a rule that creates the best split at one particular node may not benefit the model in the long term as well as the second or third best split [6]. Recursive binary splitting determines the best split by iterating through each feature and splitting the observations by various if/else rules. The combination of feature and if/else rule that yields the lowest classification error rate, and therefore is most informative way to separate the observations by class, is chosen as the rule for that particular node [6] [3]. The classification error rate can be determined by dividing the number of observations that are in the minority in a region by the total number of observations in that region. The two resulting nonterminal nodes can be split by another recursive binary split, over and over, until each leads to a terminal that scores a classification error rate of zero, but a tree of this complexity will likely be overfitting the training data [6]. There are several strategies for preventing overfitting, such as pruning and limiting the depth of the tree to a certain number of terminal nodes.

5.2 Decision Trees - Preprocessing

5.3 Decision Trees - Execute Analysis

5.4 Decision Trees - Evaluate Analysis

6 SUPPORT VECTOR MACHINES

6.1 Support Vector Machines - Intuition

A support vector machine (SVM) is a classifier that can be used to separate data using linear and non-linear boundaries. There

are a few subtypes of SVMs. A maximal margin classifier seeks to perfectly separate the observations into their respective classes by using a linear boundary. As this is often not possible given the distribution of the observations, a more lenient type of SVM is the support vector classifier (SVC). An SVC is a classifier that resigns to the fact that a few observations will be misclassified but strives to achieve an overall better separating boundary [6]. In two dimensional space, SVCs determine a line the best separates the classes. In three dimensional space, SVCs separate the classes using a plane. In four dimensional space and beyond, SVCs use a 'hyperplane' to draw a boundary between the classifiers. The term 'support vector' comes from the concept that the boundary is drawn in a way that, in addition to minimizing error, it maximizes the distance between the boundary and the nearest observations [6]. These observations are the support vectors. If the observations are not linearly-separable, an SVM can be used to advance upon SVCs by drawing non-linear boundaries [6]. In all subtypes of SVMs, unseen observations are classified based on which side of the hyperplane that the observations appears.

6.2 Support Vector Machines - Preprocessing

6.3 Support Vector Machines - Execute Analysis

6.4 Support Vector Machines - Evaluate Analysis

7 MODEL OPTIMIZATION

we don't know the test error, CV can help

Cross validation is a strategy for estimating the test error by cycling through the training data [6] GridsearchCV

8 HOW TO IMPROVE ANALYSIS

Additional features, socioeconomic status(SES)

additional studies that includes SES, do they improve?

9 PITFALLS

curse of dimensionality

10 INCORPORATING BY THE BEDSIDE

bedside alerts, discharge planning, case management team assignment, home care

11 PREVIOUS ANALYSES

flaws in methodology additional data points

12 FIGURES

In Figure 1 we show a fly. Please note that because we use just columwidth that the size of the figure will change to the column-width of the paper once we change the layout to final. CHnaging the layout to final should not be done by you. All figures will be listed at the end.

[Figure 1 about here.]

When copying the example, please do not check in the images from the examples into your images directory as you will not need

them for your paper. Instead use images that you like to include. If you do not have any images, do not dreate the images folder.

13 CONCLUSION

This is my conclusion

ACKNOWLEDGMENTS

The authors would like to thank Dr. Gregor von Laszewski for his support and suggestions to write this paper.

REFERENCES

- [1] Cristina Boccuti and Gisella Casillas. 2017. Aiming for Fewer Hospital U-turns: The Medicare Hospital Readmissions Reduction Program. Online. (March 2017). <http://files.kff.org/attachment/Issue-Brief-Fewer-Hospital-U-turns-The-Medicare-Hospital-Readmission-Reduction-Program>
- [2] Christopher M Florkowski. 2008. Sensitivity, Specificity, Receiver Operating Characteristic (ROC) Curves and Likelihood Ratios: Communicating the Performance of Diagnostic Tests. *Clinical Biochemistry Review* 29 (August 2008), S83–S87. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2556590/>
- [3] Sarah Guido and Andreas Müller. 2017. *Introduction to Machine Learning with Python* (1st edition ed.). O'Reilly Media.
- [4] Danning He, Simon C Mathews, Anthony N Kalloo, and Susan Huffless. 2013. Mining High-dimensional Administrative Claims Data to Predict Early Hospital Readmissions. *Journal of Informatics in Health and Biomedicine* 21, 2 (March 2013), 272–279. <https://doi.org/doi.org/10.1136/amiainjnl-2013-002151>
- [5] A.J. Izenman. 2008. *Modern Multivariate Statistical Techniques*. Springer Science and Business Media, LLC. https://doi.org/10.1007/978-0-387-78189-1_9
- [6] Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. 2015. *An Introduction to Statistical Learning*. Springer Science and Business Media. <https://doi.org/DOI10.1007/978-1-4614-7138-7>
- [7] Paul LaBrec. 2016. Analyze this! Administrative claims data or EHR data in health services research? Online. (January 2016). <https://www.3mhisinsideangle.com/blog-post/analyze-this-administrative-claims-data-or-ehr-data-in-health-services-research/>
- [8] Jordan Rau. 2016. Medicare's Readmission Penalties Hit New High. Online. (August 2016). <https://khn.org/news/more-than-half-of-hospitals-to-be-penalized-for-excess-readmissions/amp/>
- [9] Khader Shameer, Kipp W Johnson, Alexandre Yahia, Riccardo Miotto, Li Li, Doran Ricks, Jebakumar Jebakaran, Patricia Kovatch, Partha P Sengupta, Annette Gelijns, Alan Moskowitz, Bruce Darrow, David Reich, Andrew Kasarskis, Nicholas P Tattonetti, Sean Pinney, and Joel T Dudley. 2016. Predictive Modeling of Hospital Readmission Rates Using Electronic Medical Record-Wide Machine Learning: A Case-Study Using Mount Sinai Heart Failure Cohort. In *PSB*, Pacific Symposium on Biocomputing (Ed.), Vol. 22. Pacific Symposium on Biocomputing, 276–287. <https://www.ncbi.nlm.nih.gov/pubmed/27896982>
- [10] Beata Strack, Jonathan P DeShazo, Chris Gennings, Juan Olmo, Sebastian Ventura, Krzysztof J Cios, and John N Clore. 2014. Impact of HbA1c Measurement on Hospital Readmission Rates: Analysis of 70,000 Clinical Database Patient Records. *BioMed Research International* 2014 (April 2014), 1–11. <https://doi.org/dx.doi.org/10.1155/2014/781670>

A ISSUES

DONE:

Example of done item: Once you fix an item, change TODO to DONE

A.1 Assignment Submission Issues

Do not make changes to your paper during grading, when your repository should be frozen.

A.2 Uncaught Bibliography Errors

Missing bibliography file generated by JabRef

Bibtex labels cannot have any spaces, _ or & in it

Citations in text showing as [?]: this means either your report.bib is not up-to-date or there is a spelling error in the label of the item you want to cite, either in report.bib or in report.tex

A.3 Formatting

Incorrect number of keywords or HID and i523 not included in the keywords

Other formatting issues

A.4 Writing Errors

Errors in title, e.g. capitalization

Spelling errors

Are you using *a* and *the* properly?

Do not use phrases such as *shown in the Figure below*. Instead, use *as shown in Figure 3*, when referring to the 3rd figure

Do not use the word *I* instead use *we* even if you are the sole author

Do not use the phrase *In this paper/report we show* instead use *We show*. It is not important if this is a paper or a report and does not need to be mentioned

If you want to say *and* do not use & but use the word *and*

Use a space after . , :

When using a section command, the section title is not written in all-caps as format does this for you

\section{Introduction} and NOT \section{INTRODUCTION}

A.5 Citation Issues and Plagiarism

It is your responsibility to make sure no plagiarism occurs.
The instructions and resources were given in the class

Claims made without citations provided

Need to paraphrase long quotations (whole sentences or longer)

Need to quote directly cited material

A.6 Character Errors

Erroneous use of quotation marks, i.e. use “quotes”, instead of ” ”

To emphasize a word, use *emphasize* and not “quote”

When using the characters & # % _ put a backslash before them so that they show up correctly

Pasting and copying from the Web often results in non-ASCII characters to be used in your text, please remove them and replace accordingly. This is the case for quotes, dashes and all the other special characters.

If you see a ffigure and not a figure in text you copied from a text that has the fi combined as a single character

A.7 Structural Issues

Acknowledgement section missing

Incorrect README file

In case of a class and if you do a multi-author paper, you need to add an appendix describing who did what in the paper

The paper has less than 2 pages of text, i.e. excluding images, tables and figures

The paper has more than 6 pages of text, i.e. excluding images, tables and figures

Do not artificially inflate your paper if you are below the page limit

A.8 Details about the Figures and Tables

Capitalization errors in referring to captions, e.g. Figure 1, Table 2

Do use *label* and *ref* to automatically create figure numbers

Wrong placement of figure caption. They should be on the bottom of the figure

Wrong placement of table caption. They should be on the top of the table

Images submitted incorrectly. They should be in native format, e.g. .graffle, .pptx, .png, jpg

Do not submit eps images. Instead, convert them to PDF

The image files must be in a single directory named ”images”

In case there is a powerpoint in the submission, the image must be exported as PDF

Make the figures large enough so we can read the details. If needed make the figure over two columns

Do not worry about the figure placement if they are at a different location than you think. Figures are allowed to float. For this class, you should place all figures at the end of the report.

In case you copied a figure from another paper you need to ask for copyright permission. In case of a class paper, you must include a reference to the original in the caption

Remove any figure that is not referred to explicitly in the text (As shown in Figure ..)

Do not use *textwidth* as a parameter for *includegraphics*

Figures should be reasonably sized and often you just need to add *columnwidth*

e.g.

/includegraphics[width=\columnwidth]{images/myimage.pdf}
re

LIST OF FIGURES

1 Example caption

8

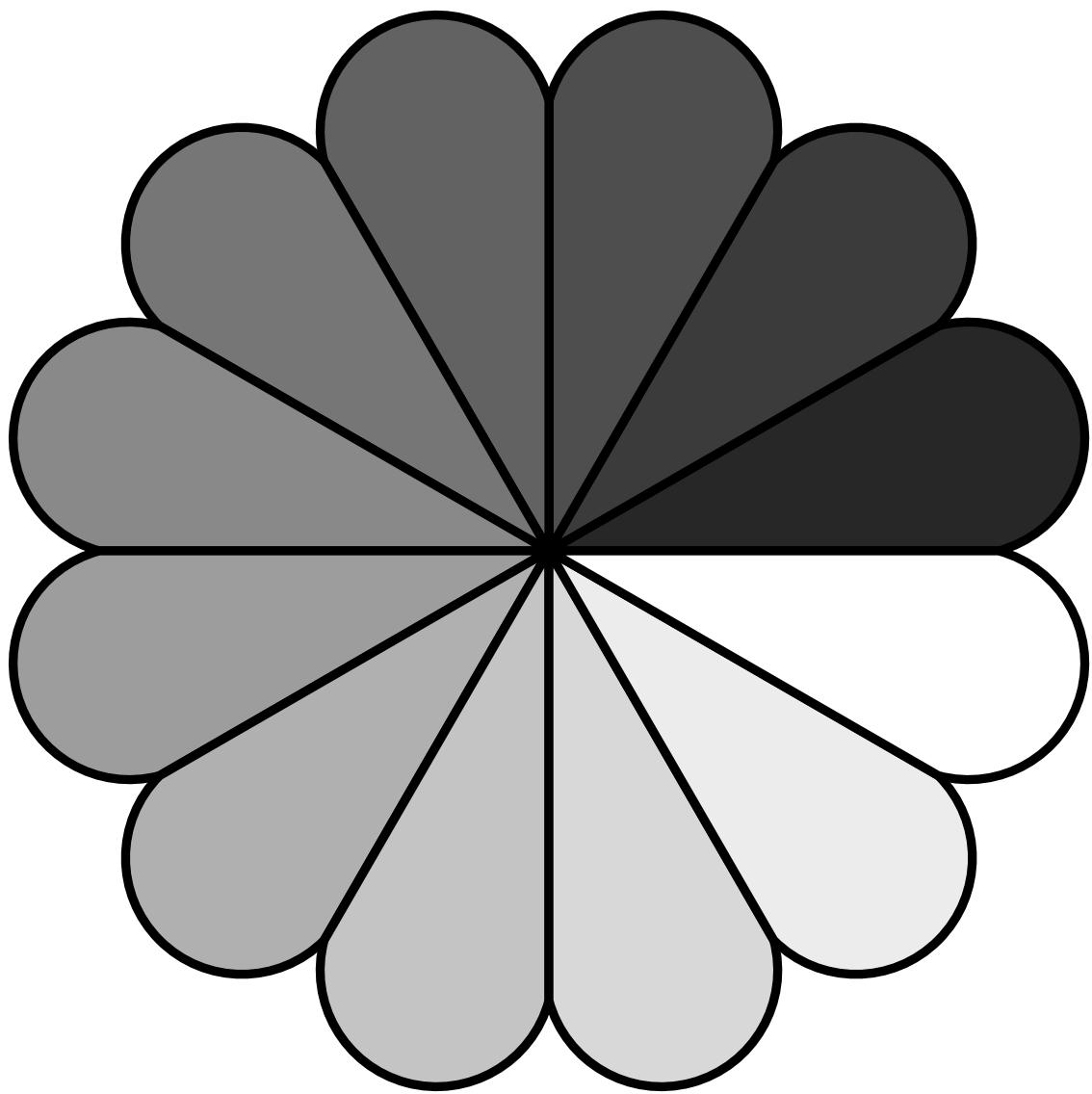


Figure 1: Example caption

```
bibtext report
```

```
=====
This is BibTeX, Version 0.99d (TeX Live 2016)
The top-level auxiliary file: report.aux
The style file: ACM-Reference-Format.bst
Database file #1: report.bib
Warning--I didn't find a database entry for "03"
Warning--can't use both author and editor fields in cite08
Warning--empty address in cite08
Warning--unrecognized DOI value [doi.org/10.1136/amiajnl-2013-002151]
Warning--empty address in cite13
Warning--empty address in cite03
Warning--unrecognized DOI value [DOI 10.1007/978-1-4614-7138-7]
Warning--empty publisher in cite01
Warning--empty address in cite01
Warning--unrecognized DOI value [dx.doi.org/10.1155/2014/781670]
(There were 10 warnings)
```

```
bibtext _ label error
```

```
=====
bibtext space label error
```

```
=====
bibtext comma label error
```

```
=====
latex report
```

```
[2017-12-04 12.23.54] pdflatex report.tex
This is pdfTeX, Version 3.14159265-2.6-1.40.17 (TeX Live 2016) (preloaded format=pdflatex
p.4 L135 : [03] undefined
Missing character: ""
There were undefined citations.
bookmark level for unknown defaults to 0.
```

The anchor of a bookmark and its parent's must not be the same. Added a new anchor.
Typesetting of "report.tex" completed in 1.4s.

Compliance Report

```
name: Tyler Peterson
hid: 331
paper1: Oct 22 17 100%
paper2: Nov 6 17 100%
project: Dec 04 17 0%
```

```
yamlcheck
```

```
wordcount
```

```
8
wc 331 project 8 4476 report.tex
wc 331 project 8 5281 report.pdf
wc 331 project 8 779 report.bib
```

```
find "
```

```
passed: True
```

```
find footnote
```

```
passed: True
```

```
find input{format/i523}
```

```
4: \input{format/i523}
```

```
passed: True
```

```
find input{format/final}
```

```
passed: False
```

```
floats
```

```
190: In Figure \ref{f:fly} we show a fly. Please note that because we
     use
196: \begin{figure}[!ht]
197: \centering\includegraphics[width=\columnwidth]{images/rosette.pdf
     }
198: \caption{Example caption}\label{f:fly}
```

```
figures 1
```

```
tables 0
```

```
includegraphics 1
```

```
labels 1
```

```
refs 1
```

```
floats 1
```

```
True : ref check passed: (refs >= figures + tables)
```

```
True : label check passed: (refs >= figures + tables)
```

```
True : include graphics passed: (figures >= includegraphics)
```

```
True : check if all figures are referred to: (refs >= labels)
```

```
Label/ref check
```

```
77: Several statistics can be used for evaluating model accuracy. For
     classification problems, a basic technique for evaluation is the
     confusion matrix. Figure 1 shows the general framework of a
     confusion matrix which shows the counts of each type of prediction
     and the accuracy of that prediction. A true positive is an outcome
     that is predicted to be positive and is positive in reality
     \cite{cite12}. A true negative is an outcome that is predicted to
     be negative and is negative in reality \cite{cite12}. These are
     the preferred responses. In the context of hospital readmissions,
     a true positive is a prediction that a patient in the test
     dataset, according to the trained model, will be readmitted to the
     hospital within 30 days, and this occurs in reality. A true
     negative is a prediction that a patient in the test dataset will
     not be readmitted, and this occurs in reality.
```

```
passed: False -> labels or refs used wrong
```

```
When using figures use columnwidth
```

```
[width=1.0\columnwidth]
```

```
do not change the number to a smaller fraction
```

```
find textwidth
```

```
passed: True
```

```
below_check
```

```
bibtex
```

```
label errors
```

```
bibtex errors
```

```
This is BibTeX, Version 0.99d (TeX Live 2016)
The top-level auxiliary file: report.aux
The style file: ACM-Reference-Format.bst
Database file #1: report.bib
Warning--I didn't find a database entry for "03"
Warning--can't use both author and editor fields in cite08
Warning--empty address in cite08
Warning--unrecognized DOI value [doi.org/10.1136/amiajnl-2013-002151]
Warning--empty address in cite13
Warning--empty address in cite03
Warning--unrecognized DOI value [DOI 10.1007/978-1-4614-7138-7]
Warning--empty publisher in cite01
Warning--empty address in cite01
Warning--unrecognized DOI value [dx.doi.org/10.1155/2014/781670]
(There were 10 warnings)
```

```
bibtex_empty_fields
```

```
entries in general should not be empty in bibtex
```

```
find ""
```

```
passed: True
```

ascii

=====
The following tests are optional
=====

Tip: newlines can often be replaced just by an empty line

find newline

passed: True
cites should have a space before \cite{} but not before the {

find cite {

passed: True

Big Data Analytics Role in Reducing Healthcare Costs in the United States

Judy Phillips
Indiana University
PO BOX 4822
Bloomington, Indiana 47408
judkphil@iu.edu

ABSTRACT

In the United States more money is spent on health care than in any other industrialized country in the world. Yet, health care access is often problematic and health care quality indicators are lower or mediocre as compared to other countries with similar economic status. Insights offered by Big Data Analytics can find solutions that will significantly lower costs and improve delivery of health care in the United States. These solutions have the potential to save billions of dollars in health care costs and to improve the quality of care for millions of Americans.

KEYWORDS

I523, HID332, health care costs, predictive analytics, electronic health records, big data

1 INTRODUCTION

Health care spending in the United States greatly exceeds the spending of other industrialized countries. Americans spend 3 trillion dollars annually on health care. Health expenditures currently account for 17.6 percent of the Gross National Product (GDP) and are expected to increase at an average rate of 5.8 percent through 2025. Health care spending has exceeded growth of the Gross National Product (GDP) in 42 of the previous 50 years [2]. Health spending threatens the nation's fiscal health [29]. Despite the excessive spending, the United States ranks among the worst on measures of health care quality, health access equity, and quality of life [22]. Policy makers do not know how to respond.

Big data analytics has the potential to help manage and address some of the cost issues while simultaneously improving patient health outcomes. Big Data ability gives us the ability to combine and analyze data from a wide variety of sources in ways that have never before been possible. This new information is providing new and valuable insights into ways to provide more effective and efficient patient care. The associations, patterns, and trends in big data may hold the key to reducing expenditures, improving care, and saving lives [29]. The information is being used to achieve more accurate and timely diagnoses, better match treatment plans to patient needs, and predict and identify at-risk patients and populations [22]. Mobile applications are being used to monitor patient care in real time. Big data can reduce health care waste, improve coordination of care, expose fraud and abuse, and to speed up the research and development pipeline.

The cost savings estimates are substantial. McKinsey and Company estimates that Big data analytics has the potential to reduce health care costs in the United States by 12 to 17 percent. This

equals to a savings of between 348 to 493 billion dollars annually [6].

Some of the tools and methodologies that big data uses to introduce efficiencies into the American health care system include: Outcome based reimbursement methodologies, electronic health records, medical device monitoring, predictive analytics, evidence based medicine, genomic analysis, and claim prepayment fraud analysis. Big data technologies are adding value and improving efficiency in almost every area of health care including clinical decision support, administration, pharmaceutical research and development, and population health management.

2 COMPARISON TO OTHER COUNTRIES

According to the Organization for Economic Cooperation and Development (OECD), the United States spends 2.5 times per person than the average of OECD related industrialized nations. In 2016, the United States spent 9822 dollars per person annually on health care. In comparison, the average amount spent per person among all OECD nations was 4033 dollars. The next highest spender was Switzerland at 7919 dollars per person [28]. The average spending as a percentage of Gross National Product (GDP) among OECD nations was 9 percent. Switzerland was again the next highest spender at 12 percent of their Gross National Product (GDP) being spent on health care. According to a McKinsey and Company analysis, the United States spends 600 billion dollars more annually than the estimated benchmark amount as calculated based upon the country's size and wealth as compared to other OECD related nations [18].

The United States lags in many standard indicators of health quality. According to a Commonwealth Fund study of 11 developed countries in 2013, the United States ranked fifth in quality and worst in infant mortality. The United States also ranked last in the prevention of deaths from treatable conditions such as strokes, diabetes, high blood pressure and treatable cancers. The average life expectancy in the United States is 76.3 years. The average life expectancy among all OECD countries is 77.9 years. The incidence of obstetric trauma is 9.6 per 100,000 births in the United States compared to 5.7 incidents per 100,000 in other countries. The statistics for preventable hospital admissions also compare poorly in comparison to other nations. In the United States the hospital admission rate for asthma and COPD was 262 per 100,000 in comparison to the average of 236 per 100,000. Thirty eight percent of the population in the United States is obese. The average obesity rate in other countries is nineteen percent. The United States has fewer physicians and hospitals. In the United States, there are 2.6 practicing physicians and 2.8 hospital beds per 1000 population.

This compares to an average of 3.4 physicians and 4.7 hospital beds on the average in the other countries [28].

The United States has material problems with health care access. Most other OEDC countries have achieved almost universal insurance coverages. On the average, 98 percent of persons in OEDC countries have health insurance. In the United States only 90 percent have health insurance. In addition, cost sharing requirements often make access additionally prohibitive. In 2016, 22.3 percent of the persons in the United States had skipped a medical consultation due to cost concerns. In comparison, the average percentage of individuals who had skipped medical visits due to cost in OEDC nations was 10.5 percent. In the United States 11.6 percent of the population had skipped taking a prescribed medication due to cost in 2016. This compares to an average of 7.1 percent of the population in other OEDC countries who reported foregoing foregone a prescribed medication due to cost [28].

3 HEALTH COST DRIVERS

Why is health care so much more expensive in the United States than it is anywhere else in the world? Some of the contributing factors include: the basic health care economic payment structure, inefficient and wasteful use of resources, medical errors, lack of transparency within the system, unnecessary administrative costs, and fraud and abuse.

3.1 Health Care Payment Structure

Many of the cost issues can be contributed to the complex, un-coordinated, multi-payer payment structure. Private insurance companies, Medicaid, and Medicare are the primary payers. An individual's eligibility by payer is dependent upon factors such as employment status, income level, age, and whether or not they are disabled. Most citizens obtain private insurance through their employment. Individuals who are 65 years of age or older or disabled are eligible for Medicare. These individuals may also purchase private Medicare Supplement insurance on their own to pay expenses that Medicare does not cover. Low income individuals may be eligible for Medicaid. If an individual is not eligible for any of these programs, he can purchase individual health insurance from a private insurance company on his own. However, individual health insurance is expensive. According to data from E-care, in 2016, the average monthly premium for an individual was 393 dollars per month. The average cost for family coverage was 1021 dollars per month [39]. In addition, individual insurance policies often include fairly high cost sharing features. Even though subsidies are available through the Affordable Care Act to offset some of these costs, many people choose to forego insurance entirely due to the prohibitive expense.

The system is inefficient and flawed because the basic economic concepts such as supply and demand and competition do not work in this sector. This is because none of the players are incentivized to manage or reduce costs [3]. Consumers do not manage medical utilization because it is being paid for by a third party, the insurance company. Insurance coverage thus insulates patients from the true costs of medical care [3]. Providers are not incentivized to provide efficient, cost effective care. Most providers are paid via a traditional fee for service methodology. That is, providers are

paid for each service that they provide. Traditional fee for service provider payment methodologies that reward health caregivers for quantity instead of quality often result in overutilization of unnecessary tests and treatment procedures. The structure is such that it encourages the production of inefficient and low value services [3]. Insurance companies pass the cost of services on to the consumers in the form of higher premiums year after year. The cost inflation cycle goes on and on.

Administrative waste is another result of the complexity of the United States multi payer payment structure. Each payer has their own rules and standards. Benefit and coverage options can vary dramatically among individuals even within the same insurance company. According to the OEDC 2008 estimates, the United States spends 7.3 percent of health care expenses on administrative activities. This is more than any other country. Comparatively, Germany spends 5.6 percent, Canada spends 2.6 percent and France spends 1.9 percent [28]. Administrative activities include transaction related activities such as billing and claims payment, and regulatory compliance such as those required to comply with government and nongovernment accreditation and regulation including licensing requirements.

3.2 Clinical and Operational Waste

McKinsey and Company estimates that clinical waste amounts to 273 dollars annually [29]. According to the Congressional budget office, 30 percent of United States spending is wasteful or not necessary [8]. There are two types of waste: operational, and clinical [3].

Operational waste results from duplication of services or inefficient production processes. An example would be a duplicate medical service because of lost medical records or the same service already being provided by another caregiver [3].

Clinical waste is created by the creation of low value outputs or care that is not optimally managed. One type of clinical waste is the spending on goods and services that provide marginal or no health benefit over less costly alternatives. Some clinical waste is the result of the uncertainty in the science of medicine. An example would be when a patient is misdiagnosed or when the treatment protocol is uncertain [3]. Other types of clinical waste may be symptoms of a flawed fee for service payment structure. These may include such things as over screening, excessive office visits, or the use of branded instead of generic drugs. Another example is when a newer or more modern treatment is marketed and sold even when it does not provide a better outcome as compared to the traditional treatment. An example was a 2 million dollar prostate cancer machine that was being marketed in 2014. It made the price of the procedure significantly more, but it did nothing to improve the health outcome [8]. Other examples of types of treatment that are the result of clinical waste include avoidable emergency room use, unnecessary hospital admissions, and excessive antibiotic use [3].

3.3 Medical Errors

Medical errors cost the United States system between 17 and 29 billion dollars annually [3]. This amount could be as much as 1 trillion dollars a year if lost productivity is taken into account [27].

This compares to an estimate of 750 million in Canada [3]. The Institute of Medicine estimates that preventable medical errors claim between 44000 and 98000 lives in hospitals each year [3].

3.4 Fraud and Abuse

The National Healthcare Anti-Fraud Association estimates losses due to health care fraud at 80 billion dollars annually. Other industry sources estimate fraudulent related losses to be around 200 billion. This accounts for approximately 2 to 3 percent of total health care spending. Research indicates that only 5 percent of these losses are ever recovered [10].

4 BIG DATA

Big data refers to electronic data sets that are so large and complex that they cannot be managed with traditional hardware and software. A report delivered to the United States Congress in August 2012 defines big data as large volumes of high velocity, complex, variable data that require progressive techniques and technologies to capture, storage, distribution, manage, and analyze the information. Big data characteristics include variety, velocity, and veracity, and volume [29]. Health care data is big data because it involves the processing of overwhelmingly large complex data sets, from a wide variety of sources and a very rapid speed [29]. In addition, the data is extremely difficult to sort, organize, and decipher [11]. Recent advances in Big Data technology gives us the ability to capture, share and store healthcare data at an unprecedented pace.

4.1 Volume

The health care industry has always generated large amounts of data. Data is needed for record keeping, compliance and regulatory reporting and patient care. Historically, this data has been stored in hard copy format. Now, more and more data is being created and stored digitally. In 2011, there were estimated to be 150 Exabytes of health related data. The amount of health related big data is growing rapidly. It is expected to soon reach the zettabyte scale and then soon after that, into the yottabytes [29].

4.2 Velocity

Traditionally, health care data has been static: for example, paper files, x-ray films, and prescriptions [29]. Ironically, in many medical situations, the speed of the response can mean the difference between life and death. Increasingly, more and more of the data is being collected in real time and at a rapid pace. For example, medical monitoring devices information collect data continuously, and can support immediate response [29].

4.3 Variety

There is an enormous variety of data being collected. The data is in multimedia including images, video, text, numerical, multimedia, paper, and electronic records. Formats include structured, unstructured, and semi-structured. Sources of data include patients, physicians, hospitals, laboratories, research companies, insurance companies, and government agencies. Data comes from web and social media such as Facebook, twitter, health plan websites and smart phone applications. Machine to machine data comes from patient sensors. Biometric data is available such as fingerprints,

genetics, hand writing information and imagining reports [29]. Physicians generate electronic medical records, physician notes, and medical correspondence. Pharmaceutical companies maintain research and development information in medical databases. The United States government houses databases concerning clinical drug trials. Data is collected by the United States Centers Disease Control and Prevention [6].

4.4 Veracity

The characteristic Veracity addresses whether the information is credible and error free. Veracity is extremely important in health care because life or death decisions on being based upon the information provided. There is a particular concern because interpretations of unstructured data such as physician notes could be incorrect or imprecise. Big data architecture, platforms, methodologies and tools are designed to take into account the uncertainties of big data analytics [29].

4.5 Unstructured Big Data

Unstructured data now makes up about 80 percent of the health care information that is available and is growing exponentially. Sources of unstructured data include: medical devices, physician and nurses notes, and medical correspondence. Being able to access to this information is an invaluable resource for improving patient care and increasing efficiency [22]. Big data technology gives us the ability to capitalize and make use of the valuable clinical information that is unstructured [15].

Traditional databases have well defined structures. The data exists in a table and column format, tables have well defined schemas, and each piece of data is stored within its own well defined space. Big data is not like that at all. Data is extracted from the source systems in its raw format. Massive amounts of this data are stored in a somewhat chaotic fashion in a distributed file system. For example, the Hadoop Distributed File System (HDFS) stores data in directories of files in a hierarchical form. The convention is to store files in 64 Megabyte files in the data nodes using a high degree of compression [15].

Big data is raw data. Big data is not cleansed or transformed in any way. No business rules are applied. The approach is to transform and apply business rules or bind the data semantically as late in the process as possible. In other words, the approach is to bind as close to the application layer as possible [15].

Big unstructured data is less expensive than traditional databases. Most traditional relational databases require propriety software that is associated with expensive licensing and maintenance agreements. Relational databases also need significant specialized resources for design, administration, and maintenance. Because of its unstructured format and open source concept, big unstructured data is much less expensive to own and operate. Big data needs little design work and is easy to maintain. A Hadoop cluster is built using inexpensive commodity hardware and runs on traditional disk drives using a direct attached (DAS) configuration instead of an expensive storage area (SAN). The practice of storage redundancy makes the configuration more tolerable to hardware failures. Hadoop clusters are designed so that they are able to rebuild failed nodes easily [15].

Big unstructured data is more difficult to use. Traditional relational database users are able to access the data using a simple structured query language (SQL) that uses a sophisticated query engine that has been optimized to extract the data. Unstructured data is much more difficult to query. A sophisticated data user, such as a data scientist may be needed to manipulate the data. However tools are being developed to solve this problem. One tool is SparkSQL. This tool leverages conventional SQL for querying and works by converting SQL queries into MapReduce jobs. Another example is Microsoft Polybase which can join data from Hadoop and traditional databases and return a single result set [15].

To summarize, advances in Big Data technology, including data management of unstructured datasets and cloud computing are facilitating the development of platforms for more effectively capturing, storing, and manipulating large data sets sourced from multiple sources [29].

4.6 Big Data Trends for Healthcare

The costs for storing and parallel processing are decreasing [22]. Previously, we had to choose what data to capture and store because storage costs were so high. Now we can capture and store everything [17]. The use of the Internet of Things is growing. Internet connected technology is everywhere and has become a common and accepted part of our culture. For example, wearable fitness devices are continuously generating health information and sending it to the cloud.

Another trend is the establishment of standards and incentives in the industry that encourage the digitization and sharing of health care data. The Health Insurance Portability and Accountability Act (HIPAA) establishes national standards for electronic healthcare transactions for the submission of claims. Claims are the documents that health providers submit to insurance companies to get paid. Such standards encourage the widespread use of Electronic Document exchange. These standards have made it possible to effectively and easily share and exchange medical information between providers and insurance companies [22]. Medicare and Medicaid have set up Electronic Medical Record (EHR) incentive programs to encourage professionals and hospitals to adopt and demonstrate meaningful use of EHRs. The Affordable Health Care Act (ACA) encourages the shift from fee for service to value based payment structures by financing initiatives to test new payment models [33].

5 VALUE BASED REIMBURSEMENT

One of the most important strategies that we can take to reduce health care in the United States is to change the way that we reimburse providers from the traditional fee for service methodology to outcome based reimbursement. McKinsey and Company estimates that this strategy alone could reduce health care spending in the United States by 1 trillion dollars over the next decade [23]. This will also mitigate medical inflation because it will automatically promote preventative care and discourage the use of low value expensive technologies. Other benefits include: improved care coordination and the reduction of redundant care. All of this results in better health outcomes, and enhanced patient satisfaction.

With the fee for service payment structure providers are paid a fee for each and every service that they perform. This tends to

encourage overutilization instead of the efficient use of medical resources. The United States tends to perform more and more expensive diagnostic services and treatment services than any other country in the world. The United States is well known for over testing and over treatment [26]. Hospitals are rewarded for preventable readmissions. Physicians are rewarded as much for a failed medical procedure as they are for a successful one. It is up to each individual physician to determine what tests and treatment services to order. From a clinical perspective, many of these tests are not medically necessary. This is a wasteful use of resources.

The goal of value based reimbursement structures are to align payment incentives with the administration of efficient, high quality medical care. Basing provider reimbursement on performance and patient outcomes encourages providers to work towards optimizing patient health instead of just providing more health care services. Caregivers are also incentivized to be more innovative and to search for ways to improve health care delivery [5].

Many payers, including private health insurance companies, Medicare, and Medicaid are starting to base reimbursement on value based incentives. The Affordable Health Care Act includes provisions to encourage the development and adoption of more effective care delivery models. Some payers are also starting to reward pharmaceutical companies by basing reimbursements on drug effectiveness [18]. Systems that have been adopted to date include: patient centered medical homes, episode based payments, global payments, shared savings programs, value based contracting, and population models, including accountable care organizations.

In the patient centered home model, the primary care physician coordinates the patients care and is rewarded for improving quality and reducing costs for individual patients. Another value based system is a population model that rewards providers for improving the health of the entire population [20]. An example of this type of program is an Accountable Care Organization (ACO). In Accountable Care Organizations, groups of doctors, hospitals, and other providers work together to provide coordinated care for patients. In Medicare supported Accountable Care Organizations, providers share in Medicare savings when they deliver high quality care and manage costs wisely [7].

Big Data Analytics can play an integral role in the development and testing of new payment model methodologies. The development and adoption of such models are still in the infancy stage. Big Data Analytics has the potential to provide information that will result in innovative payment structure and reward insights. Big data can also play a role developing clinical best practices and in identifying reasons for unjustified clinical variability in current practices.

Big Data will help to support the implementation of models that have already been adopted. Value based health care depends upon quality data collection and precise data analytics [20]. First, the data must be collected and analyzed in order to define what defines quality care. Big Data is collected and analyzed in order to establish clinical guidelines that promote a more rational use of specific diagnostic tests and treatment protocols. Second, this information must be made available to health care givers in a format that they can use for day to day clinical decision making. This is often in the form of a cloud based integration platform [20]. Next, data must be collected on an ongoing basis to provide feedback indicating

whether the providers are meeting the defined standards and if not, what can be done to improve performance. In addition, the same data can benefit future patients when data analytics are taken beyond the initial reporting and are used to develop care protocols for entire patient populations [20].

One example is in which big data is being used to track and modify provider behavior is at Memorial Care, a six hospital system in Fountain Valley, California. Memorial Care uses physician performance analytics to analyze performance of hospital doctors and outpatient providers. So far, such tracking has resulted in the reduction 280 dollars per hospital stay for the average adult patient. This equates to a 13.8 million annual dollar savings for the Fountain Valley Hospital system [9].

6 ELECTRONIC HEALTH RECORDS

An Electronic Medical Record (EMR) is a digitized version of a patients medical chart. Whereas, an electronic medical record (EMR) typically includes information from one health provider, an electronic health record (EHR) includes information from multiple providers and documents all of the available information about the patient. The objective is to provide in one place, an electronic record of a patients health. This enables the sharing of information between providers. An electronic health record (EHR) contains medical history, diagnosis, medications, immunizations dates, allergy information, radiology images, and test results [36]. These records are made available to providers in real time. Electronic health record (EHR) systems often include electronic prescription subscribing systems. Also, they can include and be integrated with evidence based tools that help providers make immediate decisions about patients care. For example, an Electronic Health record system can also automatically check for problems such as medication conflicts and notify clinicians with alerts [13].

Electronic Health Records (EHRs) improve patient health care in so many ways. Physicians have better organized, more accessible, and more complete information about the patient. A clinicians ability to make an accurate diagnosis is improved. Easily accessible patient information reduces medical errors and unnecessary tests. There is a reduction in the incidence of duplicate tests. Coordination of care is improved because every caregiver is made aware of simultaneous care that is being provided by other caregivers. It easier to communicate critical clinical information to all applicable providers in a timely fashion. Because information is made available to providers in real time, there is a drastic reduction in the probability of errors caused by such things as allergic reactions or drug interactions, especially in emergency situations. Because electronic subscribing allows physicians to communicate directly with the pharmacies, prescriptions are no longer lost or misread [13]. Preventative care improves because it is easier to track and manage when patients are due for vaccinations and screenings. It becomes possible to track prescriptions to determine if a patient has been following doctors orders [34]. Productivity is increased, overlap care is reduced, and coordination of care is enhanced [5]. In general, electronic health records (EHRs) improve quality of care enhance patient safety, and contribute to better outcomes [13].

Electronic Health records (EHRs) have significantly improved the ability to treat chronically ill patients. In the past, providers

had to limit the decisions to the amount of information that was available to them at the time. The planning of care of a chronically diseased patient that had many symptoms was often mismanaged or delayed. Electronic health records (EHRs) enable the physicians to facilitate personalized treatment for these patients in a way that has never before been possible [5]. Providers have a comprehensive record of historical treatments, diagnostic data, medical history, and meticulous medical information all in one place [?]. The result is more efficient and effective treatment for chronically ill patients. There is a reduction in the number of potential side effects and an increase the patients quality of life all at a much reduced cost. [5].

Electronic health records (EHRs) also save money by reducing administrative costs. They reduce transcription costs and eliminate chart storage and access costs.

Between 2001 and 2014 Electronic Health record (EHR) usage in physician offices rose from 20 percent to 82 percent. According to Health Information Technology for Economic and Clinical Health (HITECH) research, electronic health records are being used in 94 percent of hospitals in the United States [34]. This amount of data that is being collected by large health systems and treatment centers around the country is massive [31].

7 PREDICTIVE ANALYTICS

7.1 Definition

Predictive analytics is the process of learning from historical data in order to make predictions about the future. The objective of predictive health analytics is to provide insights that enable personalized medical care for each individual patient [30]. Traditionally, physicians have always used predictive analytics, as they have always provided health care based upon what they know about the medical history of each individual patient. Predictive Health analytics seeks to supplement that knowledge with software tools that enable physicians to make more informed choices about the patients treatment based upon data from population cohorts [31]. Patients are directed to specific treatment plans based upon their specific conditions as compared to other patients in a similar cohort. This additional knowledge has the potential to provide physician with the information they need to provide a more effective treatment plans [31]. This becomes especially important for patients with complex medical histories who are suffering from multiple conditions [34]. Predictive analytics can also improve the accuracy of diagnosing patient conditions, better match treatments with outcomes, and better predict the specific patients at risk for disease [34].

Predictive analytics takes advantage of disparate data sources including: clinical, claims, research, sensors, social media, and genomic analysis.

Predictive analytics has the potential to materially reduce health care costs and improve patient care. Insights provided can in clinical decision support, prevent hospital readmission preventions, aid in adverse incidence avoidance, and help chronic disease management. In addition, predictive analytics can identify treatments and programs that do not deliver demonstrable benefits or that cost too much [29]. Some predictive models reduce readmissions by identifying environmental of lifestyle factors that increase risk

or trigger adverse events so that treatment plans can be adjusted according. [29].

7.2 Patient Profile Analytics

Patient Profile Analytics is a specific type of predictive analysis in which patient profiles are developed to identify individuals who may be at risk for developing a disease and who could benefit from proactive management, such as lifestyle modifications. For example, patient profile analytics can be used to identify patients who may be at risk for developing diabetes.

7.3 Risk Stratification

One area in which predicting patients at risk can yield the greatest results is in identifying the patients who are at the greatest risk for the most adverse outcomes or costliest diseases [29]. Risk stratification is a methodology that can be used to identify and track the sickest and potentially costliest patients. The tool ranks or stratifies patients by potential risk and flags high risk cases for additional management. A risk stratification predictive tool takes into account risk factors such as missed doctors appointments in addition the symptoms. The tool enables doctors to intervene earlier to avoid hospital admissions and costly treatment [9].

7.4 Predictive Analytic Examples

Hundreds of thousands of dollars are spent on cancer care. Big data can be used to develop individualized, personalized cancer care programs. There is a web based application, which was sponsored by the National Cancer Institute that uses data from the Prostate, Lung, Colorectal, and Ovarian Cancer Screening trial together with patient risk factor and demographic data to help develop patient specific treatment regimens [6].

Congestive heart failure accounts for more medical spending than any other diagnosis. The earlier this condition is diagnosed, the easier it is to treat and to avoid dangerous and expensive complications. However, early manifestation is difficult to recognize and can easily be missed by physicians [22]. Machine learning algorithms have the ability to take into account many more factors than doctors alone. Predictive modeling and machine learning using large sample sizes can identify nuances and patterns that were previously impossible to see. As a result, machine learning models in the form of predictive analytics substantially improved clinicians ability to accurately diagnose persons with congestive heart failure [34].

Optum labs has developed a database with the electronic health records of over 30 million patients. They use the database to develop predictive analytic tools, the objective of which is to help doctors make Big data informed decisions that will improve patients treatment [22].

Parkland Hospital in Dallas, Texas uses predictive modeling to identify high risk patients in the coronary care unit and to predict likely outcomes when the patients are sent home. To date, Parkland has reduced readmissions for Medicare patients with heart failure by 31 percent. This equates to a 500000 dollar annual savings for this one hospital [9].

8 INTERNET CONNECTED MEDICAL DEVICES

Internet connected medical devices are becoming more affordable and are being used more and more commonly. Gartner, the analysis firm, estimates that there will be more than 25 billion connected health devices by the year 2020 [15]. These devices collect data in real time and send information into the cloud. Devices include blood pressure monitors, pulse oximeters, glucose monitors, and electronic scales [15]. Some of these devices are being used as preventive care devices. Other devices are being used by health care providers to aid in the monitoring of patient conditions. Big Data is required because the process involves the capture and analysis of large volumes of fast moving data from in hospital and in home devices in real time.

8.1 Preventative Care

Millions of people are using mobile technology help live healthier lifestyles. Smart phone applications together with wearable devices such as Fitbit, Jawbone, and Samsung Gear Fit are designed to track the wearers exercise and activity levels [12]. Measures that are typically tracked include: the number of steps taken, number of calories burned, and number of stairs climbed. The objective is to encourage the users to take a more active role in their own health and wellbeing by being more physically active. Such devices can provide individuals with the information that they need to make more informed decisions, better manage their health, and to more easily track and adopt healthier behaviors [3]. In the future, it is conceivable that it will be routine to share this information with personal physicians and that it will be incorporated into regular health care management.

An individuals data can be uploaded from the device to the cloud where it is aggregated with information from other users [15]. In an initiative between Apple and IBM, a big data platform is being developed that will allow iPhone and Apple Watch users to share their data with IBMs Watson Health cloud health care analytics service. The information will use the combination of real time activity information in combination with biometric data to discover new medical insights [12].

8.2 Medical Monitoring

Remote monitoring enable medical professional to monitor a patient remotely using various technological devices. The devices can be worn by patients with health conditions at home and in medical facilities to stream data continuously to provide real time remote patient monitoring. The devices can improve care by giving patients the ability to self-manage their conditions. Processing of real time events can be supplemented with machine learning algorithms to help provide physicians with information they need to make lifesaving interventions [22]. Patient care tends to be more proactive as patient vital signs are can be monitored constantly [22]. Medical alerts can be sent to care providers such that they immediately aware of changes in a patients condition and can respond accordingly. Devices are often used for adverse risk prediction. Remote monitoring is typically used to monitor conditions such as heart disease, diabetes mellitus, and asthma. One example of the

use of personal devices in patient care is pediatricians monitoring asthmatics to identify environmental triggers for attacks [6].

Real time systems analysis improves patient care while simultaneously reducing health care costs [5]. The devices are especially advantageous to individuals who reside in remote areas. Other advantages include: a reduced incidence of severe events, improved in patient safety, and high patient satisfaction levels.

9 PUBLIC HEALTH

Data science is being used in cities throughout the United States to predict and impede potential public health issues before they even start. For example, the Chicago Department of Public Health is modeling a program to target lead exposure in children. Information is collected from multiple sources such as, home inspection records, assessor values, health records, and census data. Predictive analytic algorithms then determine which houses have the highest potential risk. This information is then being incorporated into Electronic health records (EHRs) to automatically alert physicians to possible lead exposure risk concerning their pediatric and pregnant patients. Chicago has similar programs in place for food protection and tobacco control [14].

In San Diego, California the public health department routinely gathers big data health related information and publishes it on a user friendly web site. Information is gathered from sources such as marketing companies, mobile apps and demographic data. The data includes everything from vegetable consumption to diabetes occurrences. In one initiative, Live Well, the information was able to reduce the obesity rates at a local elementary school by 5 percent. A project that is currently in progress is the study and analysis of areas that have high rates of Alzheimers [19].

10 TRANSPARENCY

In the United States, health care price information is rarely made available to the health care consumers when they receive the care. Patients usually become aware of the costs when they receive the bill. The price of health procedures can vary radically by provider. Prices can even vary by payer for the same provider. In one study, it was estimated that consumers paid 10 to 17 percent less when they were given access to comparative price data. According a paper that was published by the American Economic Journal Economic Policy, if patients had access to price data and were willing to shop around, they could be pay significantly less for everything from routine screenings to knee surgery [2]. This tended to work best for consumers who had to pay for at least some portion of their own care.

Online pricing is a potential Big Data solution. Health related price web sites provide approximate prices for health services and procedures in fairly transparent formats. Online resources are now being made available by insurers, government agencies, internet companies and medical care providers. National insurers such as Anthem, United Health group, Humana, Aetna, and Cigna offer pricing tools to their customers. Some states, including New Hampshire, Maine, Oregon, and Massachusetts publish health pricing websites. The internet company Healthcarebluebook.com publishes information for all consumers in the United States [35].

The trend towards pay for performance reimbursement agreements will also help the cost transparency issue. This is because these pricing structures encourage health care providers to share information [5].

11 EVIDENCE BASED MEDICINE

Evidence based medicine (EBM) is an approach to medical practice that emphasizes the use of evidence from well designed and well conducted research to optimize decision making [37]. Evidence based medicine is an approach that supplements a clinicians knowledge, which may be limited by knowledge gaps or bias, with the formal and explicit information such as scientific literature or best practice methodology. Evidence based medicine eliminates guesswork for health care providers. Instead of having to rely only on their own personal judgement, providers can base treatment and protocols on credible scientific data [5].

Big Data analytics supports the research and development of evidence based best practice treatment protocols. Structured and unstructured data from a variety of sources is combined and big data algorithms are applied. Sources may include electronic medical records, financial and operational data, clinical data, and genomic data [29]. The aggregating individual data sets into big data sets enable analysis for conditions that typically have small populations. An example is the study of individuals with gluten allergies [18].

12 DRUG COSTS

It is a well known fact that drugs in the United States are priced higher than they are in other countries. There are many complicated contributing factors. One factor is lack of price regulation. Another factor is the economic structure of the health care system. Because the system includes multiple payers, there is no one payer with the power to effectively negotiate with the pharmaceutical companies as there are in other economies. Therefore, drug companies typically set drug prices at whatever the market will bear. Newly developed drugs usually have higher price tags. Big Data analytics cannot fix all of the problems with the drug market, but there are some areas in which it may have an impact: medication therapy management capabilities, drug comparison technology, and pharmaceutical research and development process improvements [4].

12.1 Medication Therapy Management

Big data analytics can play a significant role in improving the Medication Therapy Management process. Adverse drug events cost billions of dollars and result in thousands of patient deaths. Physicians and pharmacist are often overwhelmed to the point of not having the time to implement appropriate drug therapies. Drug therapies are becoming more difficult to manage as more patients are taking multiple medications. Big Data cloud analytics are helping clinicians better co manage drug therapies, and to identify drug interactions, adverse side effects, and additive toxicities in real time. The results include a reduction in the number of patient deaths, emergency room visits, hospital admissions, and hospital readmissions [9].

12.2 Comparison of Competitor Drugs

In the research, there tends to be a lot of information about individual drugs. However, there is not much information about how drugs perform in comparison to their competitors. There needs to be more drug comparative information so that physicians are better informed about the true benefits of prescribing a more costly medication as compared to a less expensive or generic drug [4]. Big data technology can play a role in making such comparisons easier to accomplish.

12.3 Pharmaceutical Research and Development

Big Data can help to streamline the Pharmaceutical Research and development process. As a result, important drugs can be delivered to the market more quickly and the cost of drug development will be reduced.

Big data can enhance the process of identifying appropriate patients to enroll in the clinical trials. First, multiple sources are now available from which to select patients. For example, social media can be incorporated into the selection process and used in addition to physician information. Secondly, the participate selection criteria can include more inclusive factors, such as genetic information. This will enable better targeting of potential trial subjects which will result in more pertinent information, while at the same time shorting trail times and reducing expenses [24].

Trial can be monitored and tracked in real time. Real time trial monitoring can decrease the number of safety and operational issues. The result is the avoidance of potentially costly issues such as adverse events or unnecessary delays [24].

Electronically captured data can improve communication. Information can be shared easily between functions and external parties. All interested individuals can have access to the data at the same time including all departments, external partners, physicians, and contract research organizations (CROs). This will replace the issue of having rigid departmental data silos that hinder interaction [24].

Genomic and proteomic data can be used to speed drug development by providing the capability to better target treatments based upon genetic indicators [17].

13 ADMINISTRATIVE COSTS

According to the Institute of Medicine (IOM), the United States spends 361 billion annually on health care administration. This is more than twice our total spending on heart disease and three times our spending on cancer. Also according to the IOM, fully half of these expenditures are unnecessary [9].

One way that providers can save money is to digitize billing processes such as benefit verification, denial management, and claims submission. A benefit verification that is done electronically costs 49 cents per patient. Comparatively, the same process done manually costs 8 dollars. It is estimated that providers could save 9.4 dollars annually by transitioning to electronic processing [21].

One example in which digitized processes are being used to streamline billing processes effectively is at the Phoenix Childrens Hospital in Arizona. They use a tool that automatically converts the clinical notes in the electronic health record (EHR) system to billable diagnostic codes [21].

14 FRAUD AND ABUSE

Common types of fraud and abuse include: billing for services that are not rendered, billing for more expensive procedures than were actually delivered, and the performance of unnecessary services.

In the past, the process of identifying misrepresented claims was tedious and time consuming. Big Data analytics makes it possible to easily identify and tag such claims. According to an article by RevCycle Intelligence, when there is repeated misrepresentation of some key fact or event, patterns are created in the data that can be detected by comparing the information to legitimate claims [10]. Anthem Health Insurance, one of the nations biggest insurance payers, uses big data and machine learning algorithms to tag suspicious claims as the claims are being processed. Tagged claims are then sent to clinical coding experts for review. The objective is to identify and address fraudulent claims before they are actually paid [10].

The Center for Medicare and Medicaid Services used predictive data analytics to identify and recover 210.7 million [22] in health care fraud in 2015. They did this by assigning risk scores to claims and providers via algorithms. This enabled the identification of abnormal billing patterns in claim submissions [10].

United Healthcare realized a 2200 percent return on their investment in a Hadoop Big Data platform that was used to identify and tag inaccurate claims using a systemic and repeatable methodology [22].

Other uses of Big Data analytics in fighting fraud and abuse include: identifying links between providers to access whether an identified unethical activity is being practiced by related providers, identification of a hospitals overutilization of services in a short time period, recognizing patients who are receiving health care services from different hospitals in different locations at the same time, and detecting prescriptions that are filled for the same patient in multiple locations at the same time. Big Data analytics can also utilize machine learning algorithms combined with historical information to detect trends in anomalies and suspicious data patterns.

15 GENOMICS ANALYTICS

Big data is playing a major role in the field of genomics and precision medicine. These technologies are helping clinicians choose the best treatment plan for individuals based upon their genetic makeup. Combining data from electronic health records (EHRs), clinical trials, and genetic testing gives researchers information to develop more effective treatments for complex diseases such as cancer and diabetes [25], and HIV. Genetic testing that has been made possible by the mapping of the human genome will cut costs and improve survival rates [1].

One area in which genomics can have a dramatic impacts is in pharmaceuticals management. In the United States, 300 million dollars are spent annually on pharmaceuticals. Studies indicate that between 20 to 75 percent of patients are not responsive to prescribed drug therapies. This can often be contributed to incorrect dosing or drug mismatches. However, 50 percent of the time it is because of a molecular mismatch between the patient and the drug. According to Alan Mertz, president of the American Clinical Laboratory Association, an estimated 30 to 110 billion can be saved

by using genetic test to select a drug that is a precise match for the genetics of the patient. By using each patients unique genomic profile, therapy can become more targeted and the instances of inappropriate care will be reduced [1].

For breast cancer patients, genetic testing can identify which 30 percent of women of an overabundance of the HER2 protein. Regular chemotherapy will not help these women, but a drug called Herceptin does. Having this information not only provides doctors with the information they need to prescribe the correct medication, it enables thousands of women avoid needless harsh, expensive chemotherapy treatment. As a result, genetic testing has been shown to reduce the risk of death by 33 percent and the risk of recurrence by 52 percent for breast cancer patients. The resulting savings are estimated to be 24 thousand dollars per patient [1].

Genetic tests can help physicians select the appropriate drug for patients with metastatic colon cancer. According to one estimate, 700 million dollars could be saved annually be obtaining this information before administering treatment [1].

According to a 2006 Brookings/AEI estimate, using genetic tests to determine the appropriate dose of the blood thinner, warfarin, could save the United States 1.1 billion dollars annually. According to a study in June 2010 by the Journal of American College of Cardiology, this test could reduce hospital admissions that are caused by inaccurate dosages by 31 percent [1].

Genomic technology is also good for the United States economy. According to Battelle, a global research organization, human genome sequencing projects generated 796 billion in economic output, 244 billion in personal income and 3.8 million job-years of employment in the United States [1].

The process of gene sequencing continues becomes more efficient and cost effective. It is expected to become a regular part of medical care in the near future [15].

16 TELEMEDICINE

Telemedicine is receiving medical treatment and advice remotely, on a computer over the internet with a physician [12]. Telemedicine has been in the market for 40 years, but the with availability of internet connected technology such as smartphones, wireless devices, and video conferences, it is becoming commonplace. It is primarily used for initial diagnosis, remote patient monitoring, and medical education. However, it is also being used for more complicated care such as telesurgery. Telesurgery is a technique in which doctors perform surgery via robots with the assistance of high speed real time data delivery technology [34].

Telemedicine is especially beneficial to patients who live in rural communities who may have to travel long distances to see a doctor or specialist. Telemedicine also gives doctors who are located in multiple locations the ability to discuss and share information. Telemedicine facilitates medical education by giving caregivers the ability to observe and be trained by subject experts no matter where their location.

Telemedicine has the potential to significantly reduce costs by reducing the number of outpatient and hospital visits [38].

17 USE CASES

Valence Health has built a data lake that they use as their primary data repository using a MapR Converged Data Platform. The system includes 3000 inbound data feeds and contains 45 different types of data including: lab test results, patient vitals, prescriptions, immunizations, pharmacy benefits, claims information from doctors and hospitals. The system reports dramatically better system performance than legacy system technology. For example, previously, it took 22 hours to process 20 million laboratory records. Now the processing time for the same number of records is 20 minutes. In addition, the new system requires less hardware [22].

The National Institute of Health developed a data lake which combines data sets from separate institutions. Now that all of the data is housed in the same location, analysis is more efficient and can be more easily shared [22].

United Healthcare uses Hadoop to maintain a platform with tools that they use to analyze information generated from claims, prescriptions, provider contracts, plan subscriber, and review information [22].

Novartis, a global healthcare company, uses Hadoop and Apache Spark to build a workflow system that aids in the integration, processing, and analysis of Next Generation Sequencing research as it relates to Genomic Analytics [22].

18 CHALLENGES

One of the most compelling challenges is clinicians willingness and ability to change behavior based upon the information provided by the data. Studies have shown that it takes more than a decade of compelling clinical evidence before a new finding becomes common clinical practice. Therefore, we need to do a better job of working with clinicians on finding ways to use the data to provide higher quality care [17].

In health care, the privacy, security, and confidentiality of the patient is paramount [15]. Big data technology has inconsistent security technology. The Health Insurance Portability and Accountability Act (HIPPA) is a federal law that was passed in 1996 that sets a national standards to protect the confidentiality of medical records and personal health information. The HIPAA law is applicable to any component of the information can be used to identify a person. The protections apply to both electronic and non-electronic forms of information [32]. HIPAA regulations make it a federal offense to breach patient security. It is important to work with vendors who understand the importance of security [15]. Liason Technologies is one company that provides solutions to the healthcare and life sciences industry that has experience meeting the HIPAA security requirements [22].

Health care data has inconsistent formatting and definitional issues [17]. There is proliferation of data formats and data representations. There are inconsistent variable definitions. A value may have different meanings for different groups. For example, a cohort definition for an asthmatic patient often differs from one group of clinicians to another [16]. Big data has the challenge of bringing all of this information together.

Another issue is lack of technical experts. The manipulation and extraction of data from often unstructured data sets require special knowledge. There have been some recent changes in tooling that

will make it easier for individualized with less specialized skills to manipulate the data. For example, Big data is starting to use include SQL as a tools for querying and data manipulation. Examples are Microsoft Polybase, Impala, and SQL Hadoop [15].

19 CONCLUSION

Big data analytics has huge potential to save the United States billions of dollars in health care costs while drastically improving health outcomes. Vast amounts of information is being captured, stored and combined in ways that offer insights have never before been possible. Innovative Big data tools are reducing medical waste, decreasing medical errors, fighting fraud, and keeping people healthier. Value based reimbursement solutions have the potential to revolutionize the health delivery system in the United States by motivating providers to find ways to deliver the best possible medical care with the most economical use of resources. The development of most of these tools is only in the preliminary stage. Therefore, we are only beginning to realize some of the potential benefits. Big data really does have the potential to bend the cost curve. Big data in health care is here to stay.

ACKNOWLEDGMENTS

The author would like to thank Dr. Gregor von Laszewski and the teaching assistants in the Data Science department at Indiana University for their support and suggestions to write this paper.

REFERENCES

- [1] American Clinical Library Association. 2011. Genetic Testing Can Help the United States Cut Costs and Improve Care. Web page as article. (2011). <https://www.prnewswire.com/news-releases/genetic-testing-can-help-the-us-cut-costs-and-improve-health-care-126105103.html>
- [2] American Economic Association. 2017. Would Price Transparency Lower Health-care Costs. Web page as article. (Feb. 2017). <https://www.aeaweb.org/research/health-care-price-transparency>
- [3] Effros Rachel M Palar Kartika Keeler Emmett B Bentleu, Tanya. 2018. Waste in the US Health System - A conceptual framework. *The Milbank Quarterly* 86 (Dec. 2018), 629–659. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2690367/>
- [4] Business Insider. 2016. Why the Price of Prescription drugs in the US is Out of Control. Web page as article. (Aug. 2016). <http://www.businessinsider.com/why-the-us-pays-more-for-prescription-drugs-2016-8>
- [5] Christian Ofori Boateng. 2016. Top 3 Ways Big Data Helps Decrease the Cost of Health Care. Web page as Article. (Nov. 2016). <https://go.christiansteven.com/top-3-ways-big-data-helps-decrease-the-cost-of-health-care>
- [6] CIO. 2015. How Big Data can save 400 billion in healthcare costs. Web page as Article. (Oct. 2015). <https://www.cio.com/article/2993986/big-data-how-big-data-can-help-save-400-billion-in-healthcare-costs.html>
- [7] CMS Centers for Medicare and Medicaid Services. 2017. Accountable Care Organizations. Web page. (Nov. 2017). <https://www.cms.gov/Medicare/Medicare-Fee-for-Service-Payment/ACO/>
- [8] Consumer Reports. 2014. Why is Healthcare so Expensive. Web page. (2014). <https://www.consumerreports.org/cro/magazine/2014/11/it-is-time-to-get-mad-about-the-outrageous-cost-of-health-care/index.htm>
- [9] DataFloq. 2016. Five ways Big Data in reducing healthcare costs. Web page as article. (March 2016). <https://datafloq.com/read/5-ways-big-data-reducing-healthcare-costs/89>
- [10] Datameer. 2017. The Role of Big Data in Preventing Healthcare Fraud, Waste, and Abuse. Web page as article. (2017). <https://www.datameer.com/company/datameer-blog/role-big-data-preventing-healthcare-fraud-waste-abuse/>
- [11] Digitalist. 2016. Can Big Data Analytics Save Billions in Healthcare Costs. Web page as Article. (Feb. 2016). <http://www.digitalistmag.com/resource-optimization/2016/02/29/big-data-analytics-save-billions-in-healthcare-costs-04037289>
- [12] Forbes. 2015. How Big Data is changing Healthcare. Web page as Article. (April 2015). <https://www.forbes.com/sites/bernardmarr/2015/04/21/how-big-data-is-changing-healthcare/#427274d12873>
- [13] Forbes. 2016. How an Electronic Health Record can Save Time, Money and Lives. Web page. (Dec. 2016). <https://www.forbes.com/sites/robertpearl/2016/12/01/how-an-electronic-health-record-can-save-time-money-and-lives/2/#4445b8275f57>
- [14] Harvard Business Review. 2014. How Cities are Using Analytics to Improve Public Health. Web page as article. (2014). <https://hbr.org/2014/09/how-cities-are-using-analytics-to-improve-public-health>
- [15] Health Catalyst. 2017. Big Data in Healthcare Made Simple: Where it Stands Today and Where its Going. Web page as Article. (Oct. 2017). <https://www.healthcatalyst.com/big-data-in-healthcare-made-simple>
- [16] Health Catalyst. 2017. Five Reasons Healthcare Data is so Complex. Web page as article. (Nov. 2017). <https://www.healthcatalyst.com/>
- [17] Health Catalyst. 2017. Hadoop in Healthcare A no nonsense Q and A. Web page as article. (Nov. 2017). <https://www.healthcatalyst.com/Hadoop-in-healthcare>
- [18] Kayyali, Basel, Knott, David, Kuiken, Steve Van. 2013. McKinsey on Healthcare. Web page as Article. (2013). <http://healthcare.mckinsey.com/big-data-revolution-us-healthcare>
- [19] KQED Science. 2015. How San Diego is Using Big Data to Improve Public Health. Web page as article. (Aug. 2015). <https://ww2.kqed.org/futureofyou/2015/08/19/how-san-diego-is-using-big-data-to-improve-public-health/>
- [20] Liaison. 2017. Value Based Healthcare - The patient is the Center but Data is the Key. Web page as blog. (2017). <https://www.liaison.com/blog/2017/06/22/value-based-healthcare-patient-center-data-key/>
- [21] Managed Healthcare Executive. 2017. Five ways to reduce healthcare administrative costs. Web page as article. (2017). <http://managedhealthcareexecutive.modernmedicine.com/managed-healthcare-executive/news/five-ways-reduce-healthcare-administrative-costs>
- [22] McDonald, Carol. 2016. How Big Data is Reducing Costs and Improving Outcomes in Healthcare. Web page as Article. (2016). <https://mapr.com/blog/reduce-costs-and-improve-health-care-with-big-data/>
- [23] McKinsey and Company. 2013. The Trillion Dollar Prize. Web page as article. (Feb. 2013). <https://healthcare.mckinsey.com/sites/default/files/the-trillion-dollar-prize.pdf>
- [24] McKinsey and Company. 2017. How Big Data can Revolutionize pharmaceutical R and D. Web page as article. (Nov. 2017). <https://www.mckinsey.com/industries/pharmaceuticals-and-medical-products/our-insights/how-big-data-can-revolutionize-pharmaceutical-r-and-d>
- [25] Pacient. 2017. How Big Data Can Improve Health Care. Web page as article. (Nov. 2017). <https://pacient.care/decks/privacy-technology/health-technology/how-big-data-can-improve-healthcare>
- [26] PBSO News Hour. 2012. Health Costs: How the US Compares with Other Countries. Web page as Article. (Oct. 2012). <https://www.pbs.org/newshour/health/health-costs-how-the-us-compares-with-other-countries>
- [27] Practice Fusion. 2017. EHR Adoption Rates 20 Must see stats. Web page as Article. (March 2017). <https://www.practicefusion.com/blog/ehr-adoption-rates/>
- [28] OECD Publishing. 2017. *Health at a Glance 2017*. OECD, Paris. http://dx.doi.org/10.1787/health_glance-2017-en
- [29] Raghupathi Viju Raghupathi, Wullianallur. 2014. Big Data Analytics in Healthcare Promise and Potential. *Springer Health Information Science and Systems* 2 (Feb. 2014), 2–3. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4341817/>
- [30] Rock Health. 2017. The Future of Personalized Healthcare: Predictive Analytics. Web page. (Nov. 2017). <https://rockhealth.com/reports/predictive-analytics/>
- [31] Search Technologies. 2017. Using Big Data Predictive Analytics to Improve Healthcare. Web page as article. (2017). <https://www.searchtechnologies.com/blog/predictive-analytics-in-healthcare>
- [32] Stephen B Thacker. 2003. HIPAA Privacy Rule and Public Health. *CDC* 52 (April 2003), 1–12. <https://www.cdc.gov/mmwr/preview/mmwrhtml/m2e411a1.htm>
- [33] The Common Wealth Fund. 2017. The Affordable Care Act Payment and Delivery System reforms: A progress report. Web page as article. (Feb. 2017). <http://www.commonwealthfund.org/publications/issue-briefs/2015/may/aca-payment-and-delivery-system-reforms-at-5-years>
- [34] The datapine blog. 2017. Nine examples of Big Data Analytics in Healthcare that Can Save People. Web page. (May 2017). <https://www.datapine.com/blog/big-data-examples-in-healthcare/>
- [35] The Wall Street Journal. 2017. How to Research Medical Prices. Web page as article. (Nov. 2017). <http://guides.wsj.com/health/health-costs/how-to-research-health-care-prices/>
- [36] US Department of Health and Human Resources. 2017. EHR Basics. Web page. (2017). <https://www.healthit.gov/providers-professionals/learn-ehr-basics>
- [37] Wikipedia. 2017. Evidence Based Medicine. Web page. (Nov. 2017). https://en.wikipedia.org/wiki/Evidence-based_medicine
- [38] Wikipedia. 2017. Telemedicine. Web page. (Nov. 2017). <https://en.wikipedia.org/wiki/Telemedicine>
- [39] Zane Benefits. 2017. FAQ - How much does Individual Insurance cost. Web page. (Nov. 2017). <https://www.zanebenefits.com/blog/bid/97380/faq-how-much-does-individual-health-insurance-cost>

bibtext report

=====

This is BibTeX, Version 0.99d (TeX Live 2016)
The top-level auxiliary file: report.aux
The style file: ACM-Reference-Format.bst
Database file #1: report.bib.bib
Warning--string name "april" is undefined
--line 15 of file report.bib.bib
Warning--string name "june" is undefined
--line 33 of file report.bib.bib
Warning--string name "june" is undefined
--line 108 of file report.bib.bib
Warning--string name "sept" is undefined
--line 117 of file report.bib.bib
Warning--string name "sept" is undefined
--line 126 of file report.bib.bib
Warning--string name "sept" is undefined
--line 135 of file report.bib.bib
Warning--string name "sept" is undefined
--line 186 of file report.bib.bib
Warning--string name "july" is undefined
--line 249 of file report.bib.bib
Warning--string name "sept" is undefined
--line 330 of file report.bib.bib
Warning--string name "april" is undefined
--line 348 of file report.bib.bib
Warning--I didn't find a database entry for "www-google-christion"
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing--line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing--line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing--line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing--line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing--line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing--line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing--line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing--line 3085 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing--line 3085 of file ACM-Reference-Format.bst


```
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing---line 3229 of file ACM-Reference-Format.bst
(There were 40 error messages)
make[2]: *** [bibtex] Error 2
```

latex report

```
=====
```

```
[2017-12-04 12.24.00] pdflatex report.tex
This is pdfTeX, Version 3.14159265-2.6-1.40.17 (TeX Live 2016) (preloaded format=pdflatex
p.5 L139 : [www-google-christion] undefined
There were undefined citations.
Typesetting of "report.tex" completed in 1.0s.
```

```
=====
```

Compliance Report

```
=====
```

```
name: Judy Phillips
hid: 332
paper1: Oct 31 2017 100%
paper2: 100%
project: 100%
```

yamlcheck

```
-----
```

wordcount

```
-----
```

```
10  
wc 332 project 10 8955 report.tex  
wc 332 project 10 9322 report.pdf  
wc 332 project 10 1551 report.bib
```

```
find "
```

```
passed: True
```

```
find footnote
```

```
passed: True
```

```
find input{format/i523}
```

```
4: \input{format/i523}
```

```
passed: True
```

```
find input{format/final}
```

```
passed: False
```

```
floats
```

```
figures 0  
tables 0  
includegraphics 0  
labels 0  
refs 0  
floats 0
```

```
True : ref check passed: (refs >= figures + tables)  
True : label check passed: (refs >= figures + tables)  
True : include graphics passed: (figures >= includegraphics)  
True : check if all figures are refered to: (refs >= labels)
```

```
Label/ref check  
passed: True
```

```
When using figures use columnwidth  
[width=1.0\columnwidth]  
do not change the number to a smaller fraction
```

```
find textwidth
```

```
passed: True
```

```
below_check
```

```
bibtex
```

```
label errors
```

```
bibtex errors
```

```
This is BibTeX, Version 0.99d (TeX Live 2016)  
The top-level auxiliary file: report.aux  
The style file: ACM-Reference-Format.bst  
Database file #1: report.bib.bib  
Warning--string name "april" is undefined  
--line 15 of file report.bib.bib  
Warning--string name "june" is undefined  
--line 33 of file report.bib.bib  
Warning--string name "june" is undefined  
--line 108 of file report.bib.bib  
Warning--string name "sept" is undefined  
--line 117 of file report.bib.bib  
Warning--string name "sept" is undefined  
--line 126 of file report.bib.bib  
Warning--string name "sept" is undefined  
--line 135 of file report.bib.bib  
Warning--string name "sept" is undefined  
--line 186 of file report.bib.bib  
Warning--string name "july" is undefined  
--line 249 of file report.bib.bib  
Warning--string name "sept" is undefined  
--line 330 of file report.bib.bib
```



```
while executing---line 3131 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing---line 3131 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing---line 3131 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing---line 3131 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing---line 3229 of file ACM-Reference-Format.bst
Too many commas in name 1 of "Bentleu,Tanya, Effros, Rachel M, Palar, Kartika, Keeler, E
while executing---line 3229 of file ACM-Reference-Format.bst
(There were 40 error messages)
```

bibtex_empty_fields

entries in general should not be empty in bibtex

find ""

passed: True

ascii

=====
The following tests are optional
=====

Tip: newlines can often be replaced just by an empty line

find newline

passed: True

cites should have a space before \cite{} but not before the {

find cite {

138: Electronic Health Records (EHRs) improve patient health care in so many ways. Physicians have better organized, more accessible, and more complete information about the patient. A clinician's ability to make an accurate diagnosis is improved. Easily accessible patient information reduces medical errors and unnecessary tests. There is a reduction in the incidence of duplicate tests. Coordination of care is improved because every caregiver is made aware of simultaneous care that is being provided by other caregivers. It is easier to communicate critical clinical information to all applicable providers in a timely fashion. Because information is made available to providers in real time, there is a drastic reduction in the probability of errors caused by such things as allergic reactions or drug interactions, especially in emergency situations. Because electronic prescribing allows physicians to communicate directly with the pharmacies, prescriptions are no longer lost or misread \cite{www-google-elec}. Preventative care improves because it is easier to track and manage when patients are due for vaccinations and screenings. It becomes possible to track prescriptions to determine if a patient has been following doctors' orders \cite{www-google-datapine}. Productivity is increased, overlap care is reduced, and coordination of care is enhanced \cite{www-}

google-christian}. In general, electronic health records (EHRs) improve quality of care enhance patient safety, and contribute to better outcomes \cite{www-google-elec}.

passed: False

bibtext report

```
=====
This is BibTeX, Version 0.99d (TeX Live 2016)
The top-level auxiliary file: report.aux
The style file: ACM-Reference-Format.bst
Database file #1: report.bib
Warning--entry type for "NPR" isn't style-file defined
--line 69 of file report.bib
Warning--no key, author in NPR
Warning--to sort, need author or key in NPR
Warning--no key, author in NPR
Warning--no key, author in NPR
Warning--no key, author in NPR
Warning--empty author in NPR
(There were 7 warnings)
```

bibtext _ label error

bibtext space label error

bibtext comma label error

Big Data for Edge Computing

Ben Trovato

Institute for Clarity in Documentation
P.O. Box 1212
Dublin, Ohio 43017-6221
trovato@corporation.com

G.K.M. Tobin

Institute for Clarity in Documentation
P.O. Box 1212
Dublin, Ohio 43017-6221
webmaster@marysville-ohio.com

Gregor von Laszewski

Indiana University
Smith Research Center
Bloomington, IN 47408, USA
laszewski@gmail.com

ABSTRACT

This paper provides a sample of a L^AT_EX document which conforms, somewhat loosely, to the formatting guidelines for ACM SIG Proceedings.

KEYWORDS

Big Data, Edge Computing i523

1 INTRODUCTION

Put here an introduction about your topic. We just need one sample reference so the paper compiles in LaTe_X so we put it here [?].

2 FIGURES

In Figure 1 we show a fly. Please note that because we use just columwidth that the size of the figure will change to the column-width of the paper once we change the layout to final. Changing the layout to final should not be done by you. All figures will be listed at the end. Please do not use phrases such as as shown in the figure below.

[Figure 1 about here.]

When copying the example, please do not check in the images from the examples into your images directory as you will not need them for your paper. Instead use images that you like to include. If you do not have any images, do not create the images folder.

7 LONG EXAMPLE

If you like to see a more elaborate example, please look at report-long.tex.

8 CONCLUSION

Put here an conclusion. Conlcusions and abstracts must not have any citations in the section.

ACKNOWLEDGMENTS

The authors would like to thank Dr. Gregor von Laszewski for his support and suggestions to write this paper.

REFERENCES

We include an appendix with common issues that we see when students submit papers. One particular important issue is not to use the underscore in bibtex labels. Sharelatex allows this, but the proceedings script we have does not allow this.

When you submit the paper you need to address each of the items in the issues.tex file and verify that you have done them. Please do this only at the end once you have finished writing the paper. To do this change TODO with DONE. However if you check something on with DONE, but we find you actually have not executed it correctly, you will receive point deductions. Thus it is important to do this correctly and not just 5 minutes before the deadline. It is better to do a late submission than doing the check in haste.

3 TABLES

In case you need to create tables, you can do this with online tools (if you do not mind sharing your data) such as <https://www.tablesgenerator.com/> or other such tools (please google for them). They even allow you to manage tables as CSV.

or generate them by hand while using the provided template in Table?. Not ethat the caption is before the tabular environment.

[Table 1 about here.]

4 QUOTES

Do not use "these quotes" but use these "these quotes".

5 LABELS

Do not use Figure 1 user the ref for the figure while using its label

6 FOOTNOTES

Footnotes must be avoided in papers. All URLs must be included as full references/citations and used with the \cite command ¹.

¹do not use footnotes

LIST OF FIGURES

1 Example caption

3

Figure 1: Example caption

LIST OF TABLES

1 My caption

5

Table 1: My caption

1	2	3
4	5	6
7	8	9

```
bibtext report
```

```
=====
```

```
This is BibTeX, Version 0.99d (TeX Live 2016)
The top-level auxiliary file: report.aux
The style file: ACM-Reference-Format.bst
Database file #1: report.bib
Warning--I didn't find a database entry for "editor00"
(There was 1 warning)
```

```
bibtext _ label error
```

```
=====
```

```
bibtext space label error
```

```
=====
```

```
bibtext comma label error
```

```
=====
```

```
latex report
```

```
[2017-12-04 12.24.39] pdflatex report.tex
This is pdfTeX, Version 3.14159265-2.6-1.40.17 (TeX Live 2016) (preloaded format=pdflatex
p.1 L57 : [editor00] undefined
p.1 L85 : 't:mytable' undefined
Empty 'thebibliography' environment.
There were undefined citations.
bookmark level for unknown defaults to 0.
The anchor of a bookmark and its parent's must not be the same. Added a new anchor.
The anchor of a bookmark and its parent's must not be the same. Added a new anchor.
There were undefined references.
Typesetting of "report.tex" completed in 0.8s.
```

```
=====
```

```
Compliance Report
```

```
=====
```

```
name: Budhaditya Roy
hid: 348
paper1: 100% Oct 25 17
```

```
paper2: 100%
project: 80%
```

```
yamlcheck
```

```
wordcount
```

```
(null)
wc 348 project (null) 579 content.tex
wc 348 project (null) 551 report.pdf
wc 348 project (null) 0 report.bib
```

```
find "
```

```
103: Do not use "these quotes" but use these ``these quotes''.
```

```
passed: False
```

```
find footnote
```

```
113: \footnote{do not use footnotes}.
```

```
passed: False
```

```
find input{format/i523}
```

```
4: \input{format/i523}
```

```
passed: True
```

```
find input{format/final}
```

```
passed: False
```

```
floats
```

```
62: In Figure \ref{f:fly} we show a fly. Please note that because we
use
```

```
69: \begin{figure}[!ht]
70: % \centering\includegraphics[width=\columnwidth]{images/fly.pdf}
71: \caption{Example caption}\label{f:fly}
86: or generate them by hand while using the provided template in
    Table\ref{t:mytable}. Not ethat
89: \begin{table}[htb]
92: \label{t:mytable}
```

```
figures 1
tables 1
includegraphics 1
labels 2
refs 2
floats 2
```

```
True : ref check passed: (refs >= figures + tables)
True : label check passed: (refs >= figures + tables)
True : include graphics passed: (figures >= includegraphics)
True : check if all figures are refered to: (refs >= labels)
```

Label/ref check

```
106: Do not use Figure 1 user the ref for the figure while using its
      label
passed: False -> labels or refs used wrong
```

```
When using figures use columnwidth
[width=1.0\columnwidth]
do not change the number to a smaller fraction
```

find textwidth

passed: True

below_check

WARNING: figure and below may be used improperly

67: figure below.

bibtex

```
label errors
```

```
bibtex errors
```

```
This is BibTeX, Version 0.99d (TeX Live 2016)
The top-level auxiliary file: report.aux
The style file: ACM-Reference-Format.bst
Database file #1: report.bib
Warning--I didn't find a database entry for "editor00"
(There was 1 warning)
```

```
bibtex_empty_fields
```

```
entries in general should not be empty in bibtex
```

```
find ""
```

```
passed: True
```

```
ascii
```

```
=====
The following tests are optional
=====
```

```
Tip: newlines can often be replaced just by an empty line
```

```
find newline
```

```
passed: True
```

```
cites should have a space before \cite{} but not before the {
```

```
find cite {
```

```
passed: True
```