

Predicting Shipping Prices from Great Britain to India and Bangladesh Using Machine Learning

Zachary J. Fuller

Northwest Missouri State University, Maryville MO 64468, USA
s553830@nwmissouri.edu and zachfuller1@gmail.com

Abstract. In this report, the freight cost, meaning the cost per kg, of shipments from Great Britain to India and Bangladesh are predicted based on a number of input features. Data is cleaned and explored using Pandas, Numpy, and Plotly Express and Plotly Figure Factory. Linear regression and random forest regressor models are created using Sklearn to predict the freight cost of shipments and evaluate the relationships between input variables. The results show that shipping time, gross weight, and India as a receiving country were potentially important features in predicting the freight cost of shipments along these routes. However, only linear regression models performed relatively decently at predicting this target, so the strength of these relationships are somewhat weak.

Keywords: data analytics · data science · machine learning · logistics · shipping · pandas

1 Introduction

This project will cover the domain of logistics. I currently work for a freight auditing company, so I was hoping to potentially expand my knowledge of the industry by working with unfamiliar data. My data comes from Kaggle, as there are a couple of decent data sources on there that are relevant to this topic. My goal was to do some predictive analytics using previously shipped packages as a training set in order to understand the most important input variables to predict the target variable, freight cost of a shipment, which is the cost per kg of that shipment. Predicting the cost to ship something based on multiple factors like weight, destination, origin, and time spent in transit can be accomplished using a regression model. I feel that this could be useful to shippers, as carriers' pricing policies are not very transparent unless you are well versed in logistics, so having a pricing estimator could empower small business owners and online sellers to choose the best shipping option for their needs.

The steps I followed to complete this project were:

1. First, I located my data set.
2. Then, I cleaned the data set using the python packages Pandas and NumPy, removing any fields as necessary and deciding how to handle missing values.
3. Then, I used Sklearn to separate the data into a train and test set.
4. I then evaluated the results. As needed, I retrained and re-tested the data, repeating

steps 2 and 3. This included some more data cleaning to remove outliers and testing different regression models.

5. After training and testing was complete, I used a python plotting package, Plotly, to display the results.

6. Finally, I compiled this report, making sure to examine any possible errors.

The key components of my approach were finding the right data sources, making sure the data was cleaned in a way so as not skew results, and most importantly, the training and testing process. This was completed using multiple linear regression models and a random forest regressor model, so making sure that the data was trained correctly was key for finding the best fit to account for each independent variable.

2 Data Collection

My data set was found on Kaggle [9]. I have also found some supplementary information regarding regression models and their results from a blog called Algosome [2]. The data was found in a few CSV files. I did not have to use any data scraping techniques, as my data set was made available for easy download through the host site. From the "Shipping Optimization Challenge" data set on Kaggle, I used the send time stamp, source country, destination country, freight cost, gross weight, and shipment charges fields. The data set had thousands of columns, so I planned to fill in some null values. Additionally, I considered slicing off the timestamps and only using dates, as it appeared that a lot of this data involves international shipping, and the timestamp could be less important and take up space/time. Ultimately, the approach I settled with was using a different field called shipping time. Finally, I needed to convert the country fields, as they were simply abbreviations in the original data set. I planned to read the CSV files into a Pandas data frame within a Jupyter notebook, merge the necessary tables and fields together, and then perform this cleaning.

3 Data Processing/Cleaning

Curating the data means that we clean and manipulate the data as needed, with the desired result being a more concise and complete data set for further analysis. In the case of my data set, I viewed the data in a Jupyter Notebook using a Pandas data frame. My data was in 2 files with the same structure, so I had to concatenate the 2 files. After examining the number of unique values and any missing values in the data frame, I determined which columns provided useful insight and which ones were bulk that could be dropped. I also had to recreate some of the columns that were stored as hard to understand abbreviations so that they could be more easily understood by myself and any observers of my final project.

I used the python packages Pandas to create and manipulate my data in the form of a data frame. I used different Pandas methods like dropna to drop

rows with missing values, drop to drop unnecessary columns, unique to view the unique values in each column, concat to combine my 2 data sets, isnull and any to find columns with missing values, and map to create a new column based on given values in already existing columns [6]. I also used the python package Plotly, specifically the modules Plotly Express and Figure Factory, as Plotly is my preferred plotting package. The Plotly.express function "histogram" was used to view the distribution of my data [4].

The only column that was missing data was called "shipping time". About 19.7 percent of the values in this column were null. I figured that this would be a very important independent variable in my later model, and I didn't want to feed any poorly imputed data into my model, so I dropped all of the rows that had missing values in this column. The data science phrase "garbage in, garbage out" was a guiding principle in this decision.

After cleaning the data, I was left with 5114 records/rows and 10 columns/attributes.

Below are some important descriptions of the attributes that I kept in my data set:

- shipment id: This was a unique ID given to each shipment in the original data set. I debated whether or not to keep this in or simply refer to individual records with the data frame's index, but I ultimately decided that it may be useful to have this ID for further reference in the future. [9]
- freight cost: The cost/kg of the shipment [9]. This field is the dependent target variable that I was trying to predict.
- gross weight: The weight in kg of the shipment. [9]
- shipment charges: This one was confusing to me at first, but looking back through the documentation on Kaggle [9], I learned that this is the minimum amount that a company would charge for a specific shipment. I did not end up using this in my final analysis, but at time of cleaning, I figured I would keep it in there just in case it was relevant. [9]
- total cost: This was a field I calculated by multiplying the freight cost by the gross weight, as the freight cost is the charge per kg. This was the field that I was originally planning to use as the the dependent variable to predict. However, I soon realized that by predicting this variable that I calculated using other variables in the data set, the original variables used to calculate it would have a nearly perfect co-linear relationship with this variable. This ended up being true when creating some models, where the resulting r2 scores were roughly 0.99. So, I decided to drop this column and just focus on freight cost.
- shipment mode: The way the product is shipped, either by air or ocean. [9]
- shipping company: One of three anonymous shipping companies used. The shipment charges field values are based on this field's value. [9]
- shipping time: Time in days that it took the shipment to reach it's destination. [9]
- sender: The origin country. This was originally stored in a field called source country with a 2-letter abbreviation as a value, so I calculated this field using the Pandas map function to make it easier to understand. However, all of the values here are "Great Britain," so I later ended up dropping this field.

- receiver: The destination country. This was originally stored in a field called destination country with a 2-letter abbreviation as a value, so I calculated this field using the Pandas map function to make it easier to understand.

In the end, the independent variables I decided to use were gross weight, shipping time, shipping company, shipment mode, and receiver. The dependent variable that I was predicting was the freight cost field.

4 Exploratory Data Analysis

Exploratory data analysis (EDA) is the process of investigating data to pull out some key insights, including which are the most important variables and whether or not any outliers exist within the data. EDA allows us to get a better understanding of relationships that may exist in our data before creating our models. It is essential in a data science project because it allows us to check our assumptions before getting to the more complicated model building process, hopefully saving us some time and headaches at that point. If we begin creating and testing our model but don't understand our data, we may consistently get poor results or have to repeat our steps, whereas completing EDA should minimize any redundancies later in the process thanks to a better understanding of the data. Additionally, using EDA, we can spot any errors or outliers in the data and remove those before we get to the next steps. EDA allows us to perfect our data cleaning performed before, as we better understand our data and any flaws that may exist.

Univariate EDA, either graphical or non-graphical (meaning with or without visuals, respectively) is where we examine one variable by itself [8]. I performed univariate graphical analysis by creating histograms of each of my variables to better understand the distribution of each variable.

Multivariate EDA, either graphical or non-graphical, means that we examine multiple variables [8]. I created a few scatter plots to better understand the relationship between some key variables. I also created a correlation matrix to visualize which variables had strong relationships to one another.

For my univariate graphical analysis, I examined each variable that had more than one value on its own histogram, created using Plotly. This allowed me to see whether or not data was skewed in any particular way. I found that most of the data was skewed to the left but with very long tails. I also created box plots with Plotly and found that a lot of my variables had several outliers.

For my multivariate graphical analysis, I created a correlation matrix using Pandas' corr function and Plotly Figure Factory's annotated heat map function, which allowed me to see the relationship between variables. The general result was that there were two very strongly related variables, total cost and gross weight, and the rest were not strongly correlated. However, these two variables being strongly correlated should not be a surprise, as total cost was generated using the gross weight field. Because of this, I decided not to use total cost in my analysis. I also created Plotly scatter plots, which allowed me to get a better

understanding of how three variables related to one another thanks to an x-axis, a y-axis, and the ability to color points on the scatter plot.

Some interesting takeaways were that the gross weight of shipments in my data set were somewhat normally distributed except for having a very long tail, most shipments cost 100,000 dollars or less, air shipments were more common than ocean shipments, and most shipments took less than 6 days to reach their destination. Using multivariate graphical analysis, I found that the shipping time had the strongest relationship with freight cost, but the relationship was rather weak.

```
px.histogram(df, x='freight_cost')
```

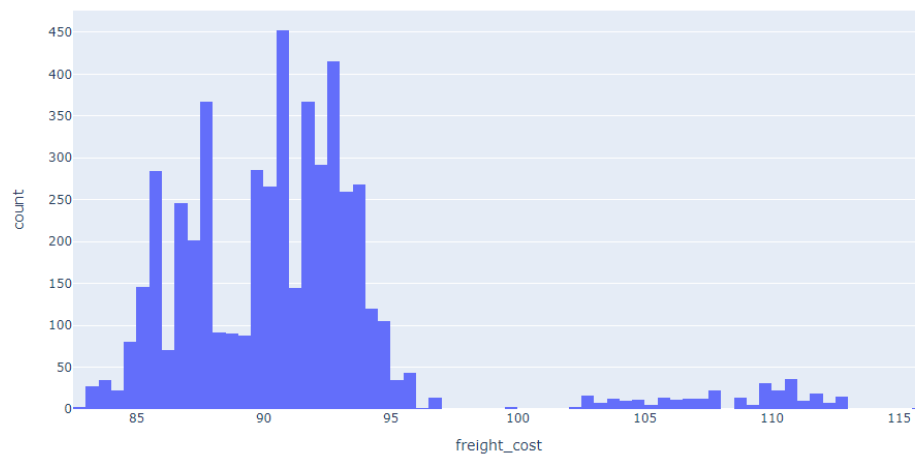


Fig. 1. Freight cost distribution

```
px.histogram(df, x='gross_weight')
```

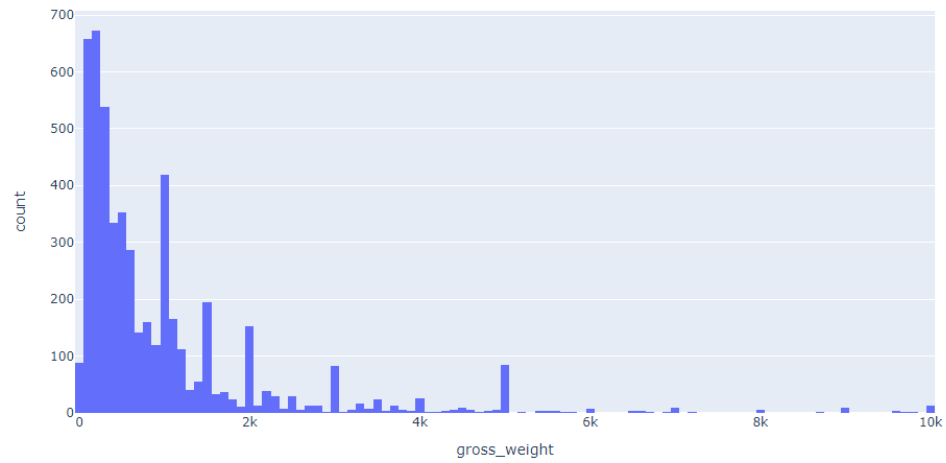


Fig. 2. Gross weight distribution

```
px.histogram(df, x='shipping_time')
```

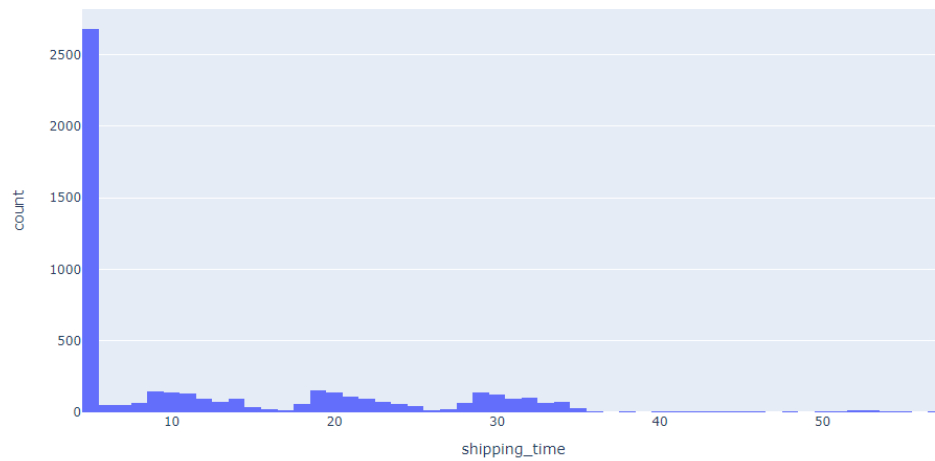


Fig. 3. Shipping time distribution

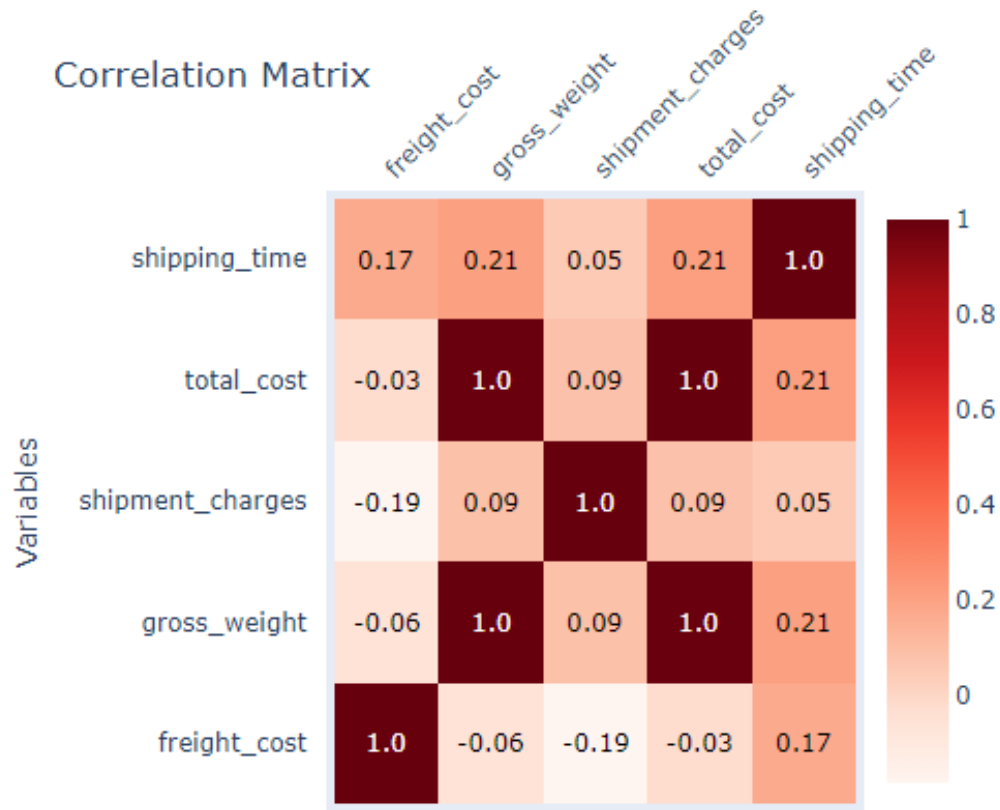


Fig. 4. Correlation matrix

Using box plots (a form of univariate graphical analysis), I found that gross weight and shipping time both had a lot of outliers, which was also apparent in the histograms. Even so, the histograms still looked somewhat normally distributed to me, so I figured that for the time being, it was okay to keep the outliers, as they may be relevant data points. However, when testing my regression model later, I decided to come back and clean the data a bit more by removing these outliers to see how the results changed. It turned out that removing the outliers resulted in poorer performing models, so I decided to keep the outliers in the data set.

```
px.box(df, y='gross_weight', title='gross_weight')
```

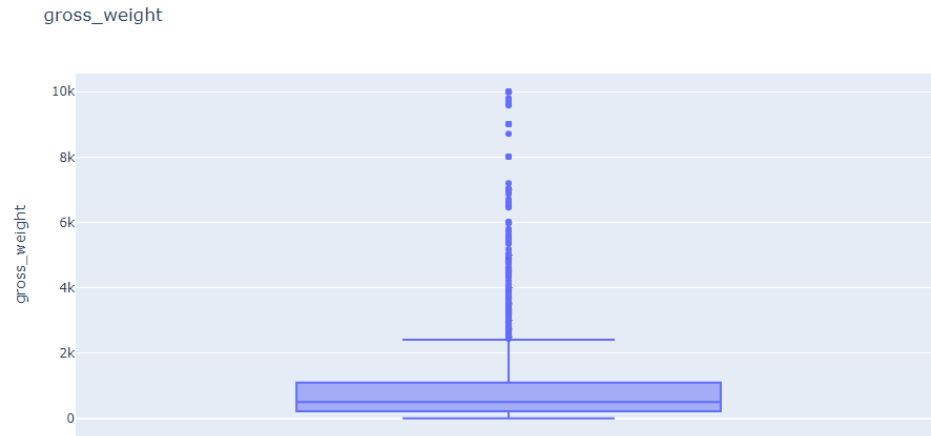


Fig. 5. Gross weight box plot showing a lot of outliers

I created scatter plots (a form of multivariate graphical analysis) and saw that of the two receiver countries, India received the highest cost and weight of shipments, the most expensive and heaviest shipments were shipped through ocean shipping as opposed to air, and anonymous shipping company SC1 handled the heaviest and most expensive shipments.

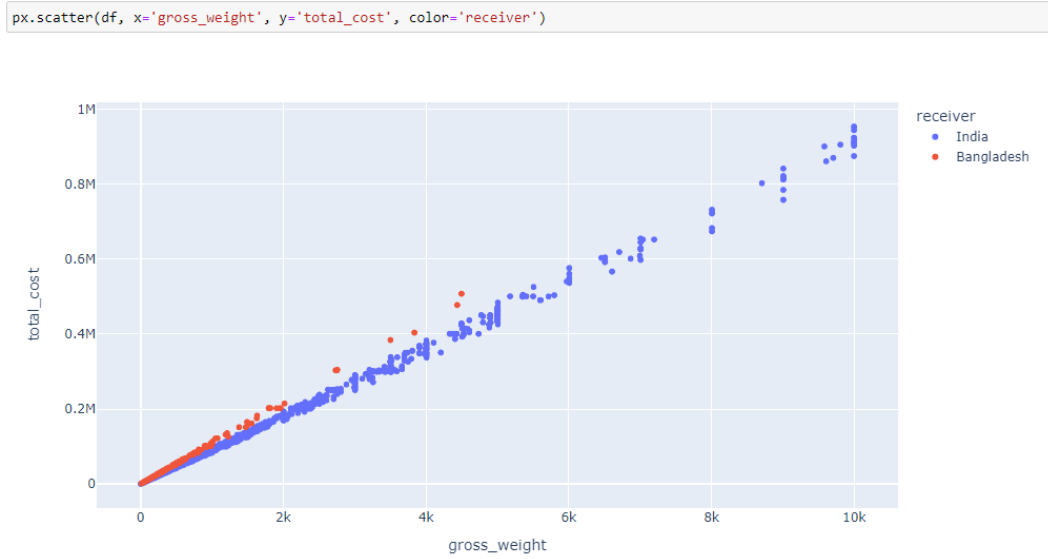


Fig. 6. Comparison of gross weight, total cost, and receiver

Overall, I learned that the most important independent variables would likely be gross weight, shipping time, shipping company, receiver and shipment mode. Freight cost would become the dependent variable so as to avoid the co-linearity of gross weight and total cost becoming a problem when evaluating the models. I would have to clean my data a bit more to drop some unnecessary columns, and I planned to possibly later remove outliers if they became a problem in my model testing.

5 Predictive Analysis

The pipeline to build my predictive model was as follows (some of this is a recap of previous sections):

- Data Collection, where I had to find a suitable data set from the web.
- Data Processing and Cleaning, where some feature engineering occurred and features were selected as relevant, possibly manipulated, turned into new features, or dropped.
- Exploratory Data Analysis occurred and meant that I took a more in depth look at the data using univariate and multivariate graphical analysis, and from here, a bit more data cleaning happened. Relationships between variables were also explored.
- Then, my data was split into training and test data using Sklearn's model selection train test split function [7].
- Next, I fit my model and evaluated the results.

- After this, I cleaned the data a bit more and re-fit the model and re-evaluated the results. This process was repeated a few times until I had a result that I was happy with.

- Finally, I created some visualizations with Plotly to best understand the performance of the models.

I used a linear regression model, which is a simple supervised learning algorithm that tries to find the line of best fit (the regression line) between input variable(s) and a target variable. This can work well with one or two input variables, but I learned that as you add more input variables, the results can be harder to interpret. I also used a random forest regressor model, which is another supervised learning algorithm. This one is based on decision trees. I chose to use this model as well because I was hoping to see if it performed better than my regression model, and a Random Forest Regressor model can help to better understand the relationships between variables [7].

I used Sklearn's model selection train test split function to split my data into train and test sets, with a test size of 0.2, meaning that 20 percent of my data was set aside for testing, leaving 80 percent of it for training. After fitting and evaluating my model a couple of times, I decided to try and add polynomial features instead to see if my results changed. I used Sklearn's preprocessing PolynomialFeatures function to do this, with biases excluded and interaction only set to True. Then I fit used fit transform to fit the polynomials to my input features, and the results of that were used in one of my model's.

The overall implementation and evaluation process of my analysis was as follows:

- Before splitting for training or testing, I created dummy variables to account for categorical variables that I wanted to use as inputs in my models. I used Pandas get dummies function and set drop first to True so as to avoid the "dummy variable trap." According to the blog Algosome, the dummy variable trap occurs when some of the independent variables are co-linear, meaning that they are correlated and can explain one another. [2]

- Then, I split my model into train and test data as explained above, but without polynomials.

- Next, I fit the linear regression model with the training data. To evaluate the data, I examined the mean absolute error, the root mean squared error, the mean squared error, and the r2 score, as well as the coefficients. I noticed that my r2 score was decent at around 0.7, but I was hoping that tweaking the model could result in better performance.

```

Bias is 110.2318807575731
Coefficients [ 1.58805299e-04 -1.69953730e-03 -7.22308374e-01 -9.95959688e-01
-1.71179918e+00 -1.89560588e+01]
MAE is 2.5611475928555136
RMSE is 2.9868129299006
MSE is 8.921051478221408
R^2 0.7075690298573324

```

Fig. 7. Results of first linear regression model run

| | Coefficient |
|-----------------------------|-------------|
| gross_weight | 90.161303 |
| shipping_time | 1.009976 |
| shipping_company_SC2 | 245.030918 |
| shipping_company_SC3 | 52.781520 |
| shipment_mode_Ocean | -297.812438 |
| receiver_India | 0.000000 |

Fig. 8. Coefficients of variables in the first linear regression model run

- I then created some line plots (for numeric variables) and box plots (for the previously string dummy variables) and saw some pretty odd and hard to interpret results.

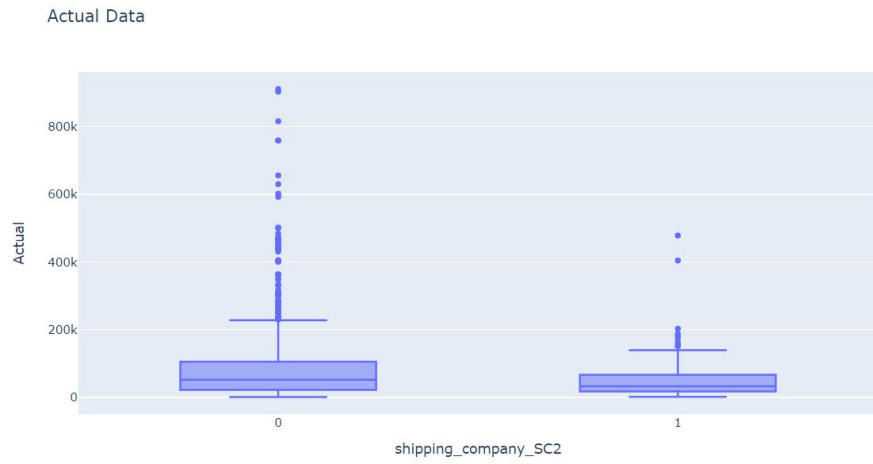


Fig. 9. Results of first linear regression model run - Shipping Company 2 Dummy Variable Box Plot

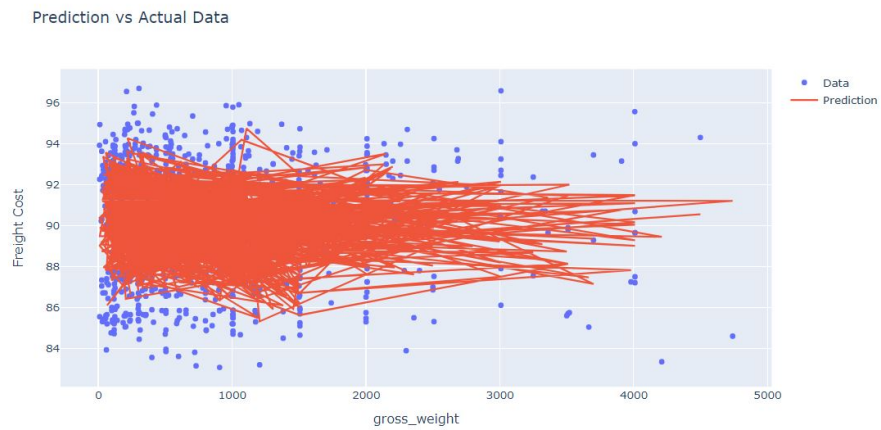


Fig. 10. Results of first linear regression model run - Shipping Time vs Total Cost

- Given the odd results, I performed some more data cleaning by removing outliers from my data using Numpy's `abs` function [3] and Scipy's `stats' zscore` function [5].
- Then, I split my data again and refit the linear regression model. I noticed that my r^2 value dropped drastically to around 0.01, so I realized that the outliers were important and should be kept in the data.

```

Bias is 90.16623395680668
Coefficients [ 5.07291014e-05  2.23837998e-03  4.93094738e-01  1.18532649e-01
-6.11627388e-01  0.00000000e+00]
MAE is 2.481207727749636
RMSE is 2.9080652543565004
MSE is 8.456843523595538
R^2 0.01354182284745098

```

Fig. 11. Results of second linear regression model run

- This is when I created polynomial features and transformed my input variables using them [7].
- Then, I refit my linear regression model again and evaluated the same metrics as before. I saw that the values I evaluated previously had not changed much from my first run of the model, so I knew I needed to take another approach.

```

Bias is 103.61684692447362
Coefficients [-1.53894935e-03  1.31279900e-01  5.68300992e-01 -6.53542150e+00
4.56525084e+00 -1.15430821e+01  2.97556218e-03 -3.44313898e-03
5.68859079e-03  1.36312585e-03 -1.31538106e-03 -3.12747713e-01
9.58547358e-01 -1.55383161e-01  1.50616986e-01 -1.92290628e-12
-1.40186967e+00  1.97017066e+00  0.00000000e+00 -6.53542150e+00
-6.97783122e+00 -1.41576252e-04 -4.58071611e-03 -2.90633968e-03
-6.40014510e-05  0.00000000e+00  5.14752699e-03 -8.59066595e-03
0.00000000e+00  5.68859078e-03  1.58669415e-03  0.00000000e+00
3.59136583e-01 -6.71884297e-01  0.00000000e+00  9.58547357e-01
-1.36046075e-01  0.00000000e+00  0.00000000e+00  0.00000000e+00
0.00000000e+00]
MAE is 2.5645110614566766
RMSE is 3.0059926840971345
MSE is 9.035992016845494
R^2 0.7038012931392317

```

Fig. 12. Results of third linear regression model run with polynomial features added

- Finally, I decided to try a random forest regressor, as I was hoping for some better results than from the linear regression models. I split the data into train and test sets again and fit the random forest regressor model. The results were actually much worse, with an r^2 score of -0.009, but I was able to visualize the importance of each feature. I found that two features, shipping time and gross weight, were the most important in this model, with respective weights of 0.52 and 0.43, which were much higher than the weights of the other variables.

| | Importance |
|-----------------------------|------------|
| shipping_time | 0.529578 |
| gross_weight | 0.438986 |
| shipping_company_SC2 | 0.014471 |
| shipping_company_SC3 | 0.013600 |
| shipment_mode_Ocean | 0.003366 |
| receiver_India | 0.000000 |

Fig. 13. Importance of each input variable in the random forest regressor

- Learning about the importance, I was able to determine that two inputs, shipping time and gross weight, were somewhat important in determining freight cost, but their importance was not very apparent, and just because they held importance in the poorly performing random forest regressor model did not mean that they were important in the better performing linear regression model. In the best performing linear regression model, a dummy variable field called receiver India had the highest absolute value of its coefficient. However, the coefficient of -18.95 was still rather low, so it is hard to determine the importance of this feature. Additionally, the fact that removing outliers from my data set, re-training my model, and fitting the model again led to much poorer results emphasized the fact that the outliers in this data set were actually important in maintaining the predictive power of the model. Even though they may appear to be outliers, their influence is rather important.

- It is important to note that before all of this, I went through the same process using total cost as my target dependent variable. However, I noticed that my r^2 score was around 0.99 after every model's run. This seemed suspicious, leading me to believe there was some over-fitting. Retracing my steps, I realized that because total cost was calculated from gross weight and freight cost, it was bound to have co-linearity with at least one of those features. The relationship between gross weight and total cost was so strong that it skewed the results of every model. Because of this, I had to choose a different target variable, and freight cost seemed to be the most appropriate stand in for total cost. This whole

process was very iterative, and I learned the value of having truly clean data that you understand well.

From my analysis, I learned that only two of my input features were relevant in predicting the target feature. Essentially, shipping time and gross weight were far more important in predicting total cost than any other input feature, but their importance was still minimal. Again, their importance was only highlighted by the random forest regressor model, which did not perform well compared to the linear regression model I first ran. Looking at the coefficients of my first linear regression model, the best performing one, the dummy variable field, receiver India, had the highest absolute value of a coefficient at 18.95, implying greater importance within that model. However, it is difficult to say how much importance can be placed on this variable.

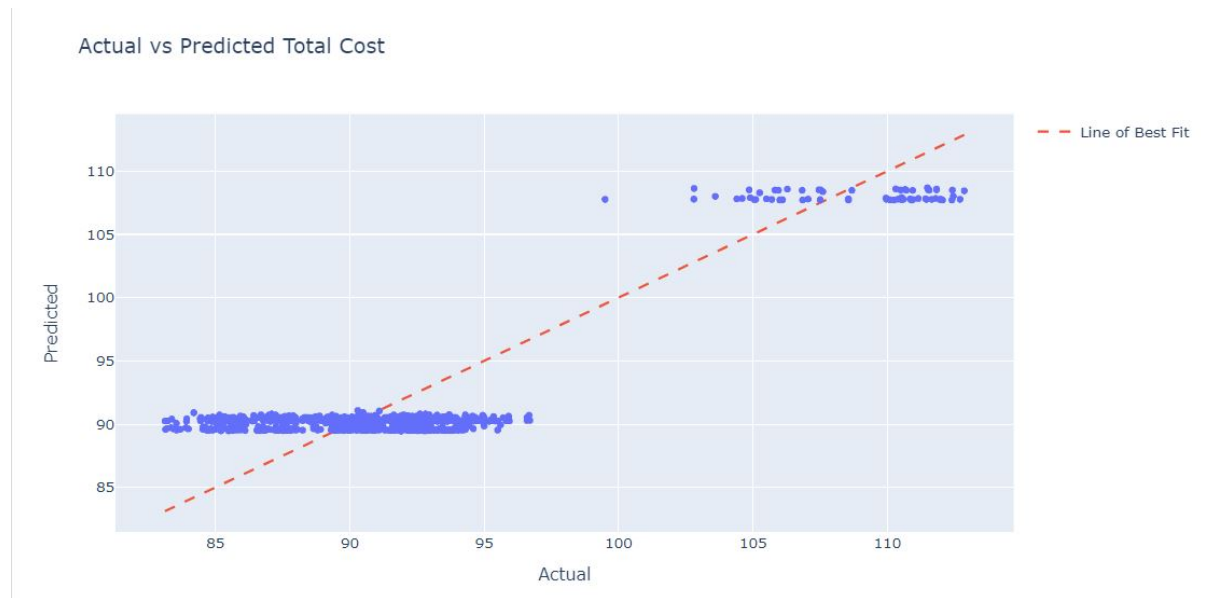


Fig. 14. Actual vs. Predicted Cost of the Freight Cost of Shipments from the original linear regression model

6 Interpretation of Results

The resulting plots from this study included a scatter plot showing the actual values of my target feature (total cost) on the x-axis and the predicted cost of feature on the y-axis, which helped me to best communicate my findings in the section above. Additionally, a table and bar chart showing the importance of each feature helped to visualize the relative importance of each feature in the random forest regressor model. In previous modules, charts that were helpful to

explore the relationship between data variables included scatter plots, box plots, and a heat map to explore correlation.

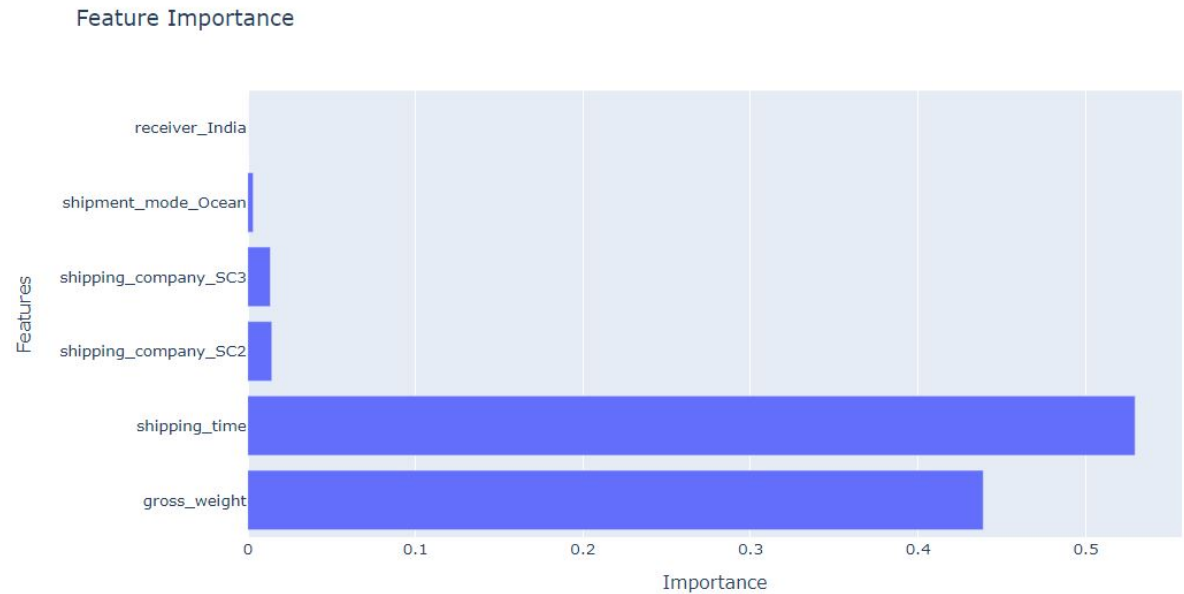


Fig. 15. Feature importance from the random forest regressor model

Using the bar chart of feature importance, I was able to see that two independent variables, shipping time and gross weight, were somewhat important in predicting the freight cost of a shipment, as their bars were much larger than the bars of the other variables. Again, this importance is only reflective of the performance of the random forest regressor model, but we can infer that these features held some weight in all models. Additionally, the scatter plot of actual vs predicted values shows the mediocre predictive power of the linear regression model in this case.

The data lacked significant features, other than shipping time and gross weight, to predict the freight cost of a shipment. However, this does not necessarily mean that there are no other important features that could have influenced the total cost of a shipment. There could have been other data points that simply weren't tracked in this data set that had more importance in predicting the target. Even though the importance of shipping time and gross weight were highlighted by the poorly performing random regressor model, I can not say for certain that these variables are the two most important factors that influence the freight cost of a shipment due to the limited scope of the data set. However, the linear regression model that performed best had a dummy variable field, receiver

India, that had the highest absolute value of a coefficient, with a coefficient of -18.95, showing that it was potentially more important than the other variables in predicting freight cost.

As stated above, shipping time and gross weight had the highest importance scores from the random forest regressor model, showing that they were decently important in one model. Even so, I can not say with certainty that these were important in the linear regression models. The dummy variable, receiver India, had the highest absolute value of coefficients of -18.95 in the best performing linear regression model, showing some importance in deciding the freight cost.

It was interesting to me that all of the coefficients in the best performing linear regression model were so small. No single feature seemed to hold much importance, even though two features, shipping time and gross weight, had much clearer importance in the random forest regressor model. Overall, it seems that a myriad of factors are important in predicting freight cost, and there may not have been enough variables in this data set to best understand the freight cost.

7 Limitations

Although running the linear regression and random forest regressor models displayed the strength of the relationships between shipping time, gross weight, India as a receiving country, and freight cost, the original data set, models, and results are not without their limitations.

For starters, the data set was taken from Kaggle, a site with user submitted data for various purposes. In this case, the data was uploaded for use in a competition [9]. This means that the user who uploaded the data could have picked and chosen the variables to keep, and other important variables could have been left out. Additionally, the scope of this data is rather small, containing information for only two shipment modes (air and ocean), three unnamed shipping companies, one origin country (Great Britain) and two receiver countries that are relatively close to each other when examined at a global scale (India and Bangladesh). Having more specific location information, more companies, more transport modes, information about what was being transported, and more information about the breakdown of charges could have helped these models to have much stronger predictive capabilities. I believe that with more of these data points included, I could have created models that better defined the relationships between some of these variables and total shipment costs.

Two linear regression models performed somewhat well, but given the only decent r^2 scores 0.7 and the lack of clarity as to the actual importance of some input variables on the target variable, it is hard to say what the actual most important factors influence freight cost are.

8 Conclusions

The overall goal of this project was to predict the cost of shipments per kg, referred to in this paper as freight cost, based on a number of input features. Using

tools available in Pandas, Numpy, and Plotly Express, I was able to perform data cleaning and exploratory data analysis on a data set found on Kaggle to identify the most important features for my analysis. Then, Sklearn was used to create and evaluate linear regression models and a random forest regressor model. The results were a few models that only semi-accurately predicated the freight cost of shipments from Great Britain to India and Bangladesh, but the models were not consistent and highlighted only a few features as important indicators: shipping time, gross weight, and India as a receiver, and the importance of these input features are questionable. The results can potentially help determine the cost per kg to ship something along these routes with a moderate degree of accuracy.

□

References

1. Capstone github repo, <https://github.com/HundredDucks/Capstone/tree/master>
2. Dummy variable trap in regression models, <https://www.algosome.com/articles/dummy-variable-trap-regression.html>
3. Numpy user guide, <https://numpy.org/doc/stable/user/index.html>
4. Plotly express in python, <https://plotly.com/python/plotly-express/>
5. Scipy user guide, <https://docs.scipy.org/doc/scipy/tutorial/index.html>
6. User guide, https://pandas.pydata.org/docs/user_guide/index.html
7. User guide, https://scikit-learn.org/stable/user_guide.html
8. What is exploratory data analysis?, <https://www.ibm.com/topics/exploratory-data-analysis>
9. GAUTAM, S.: [1] shipping optimization challenge, <https://www.kaggle.com/datasets/salil007/1-shipping-optimization-challenge?select=train2p.csv>
10. Pollack, N.: 3 ways to improve logistics management using shipping data, <https://transimpact.com/nextsights/3-ways-to-improve-logistics-management-using-shipping-data/>