

## Medical Device Proposal (Final Report)

### Problem Statement

TG Medical (USA) Inc. has been a worldwide leader in manufacturing and specialization in disposable gloves. Using the Medicare Durable Medical Equipment, Devices & Supplies dataset (via Data.CMS.gov - <https://data.cms.gov/resources/medicare-durable-medical-equipment-devices-supplies-by-geography-and-service-data-dictionary>), I will be proposing to management a new medical device to add on to our product list and predicting whether it will generate profit for the company.

### Data Wrangling

The raw dataset contains 40267 rows and 18 columns, there were 2 columns with missing values and half of the columns contain strings which will need to be converted into dummy values later.

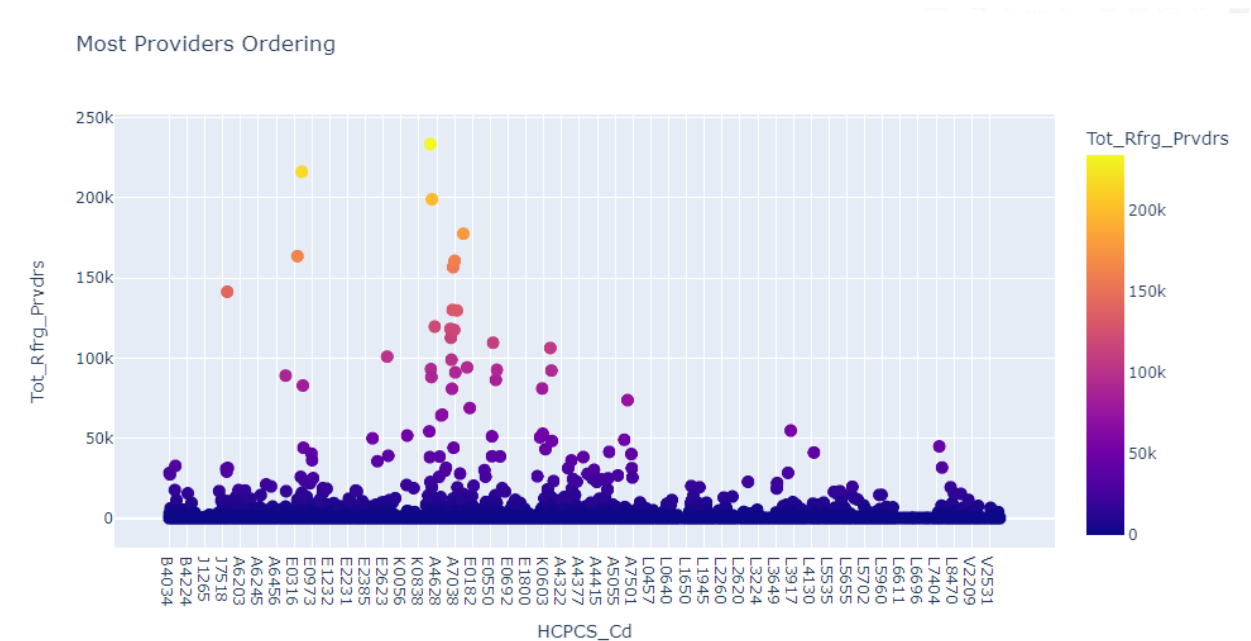
The various columns are:

- Rfrg\_Prldr\_Geo\_Lvl – Geography Level ('State' or 'National')
- Rfrg\_Prldr\_Geo\_Cd – Referring Provider Geography Code
- Rfrg\_Prldr\_Geo\_Desc – Referring Provider Geography Description
- BETOS\_Lvl – BETOS Level (3 classifications – Durable Medical Equipment, Prosthetic and Orthotic Devices, and Drugs and Nutritional Products)
- BETOS\_Cd – BETOS Code (BETOS code that is assigned to the HCPCS code)
- BETOS\_Desc – BETOS Description
- HCPCS\_Cd – HCPCS Code (for the specific product/service)
- HCPCS\_Desc – HCPCS Description
- Suplr\_Rentl\_Ind – Supplier Rental Indicator (if the product/service is a rental)
- Tot\_Rfrg\_Prldr – Number of Referring Providers
- Tot\_Suplrs – Number of Suppliers
- Tot\_Suplr\_Benes – Number of Supplier Beneficiaries
- Tot\_Suplr\_Clms – Number of Supplier Claims
- Tot\_Suplr\_Srvcs – Number of Supplier Services
- Avg\_Suplr\_Sbmtd\_Chrg – Average Supplier Submitted Charges
- Avg\_Suplr\_Mdcr\_Alowd\_Amt – Average Supplier Medicare Allowed Amount
- Avg\_Suplr\_Mdcr\_Pymt\_Amt – Average Supplier Medicare Payment Amount
- Avg\_Suplr\_Mdcr\_Stdzd\_Amt – Average Supplier Medicare Standard Payment Amount

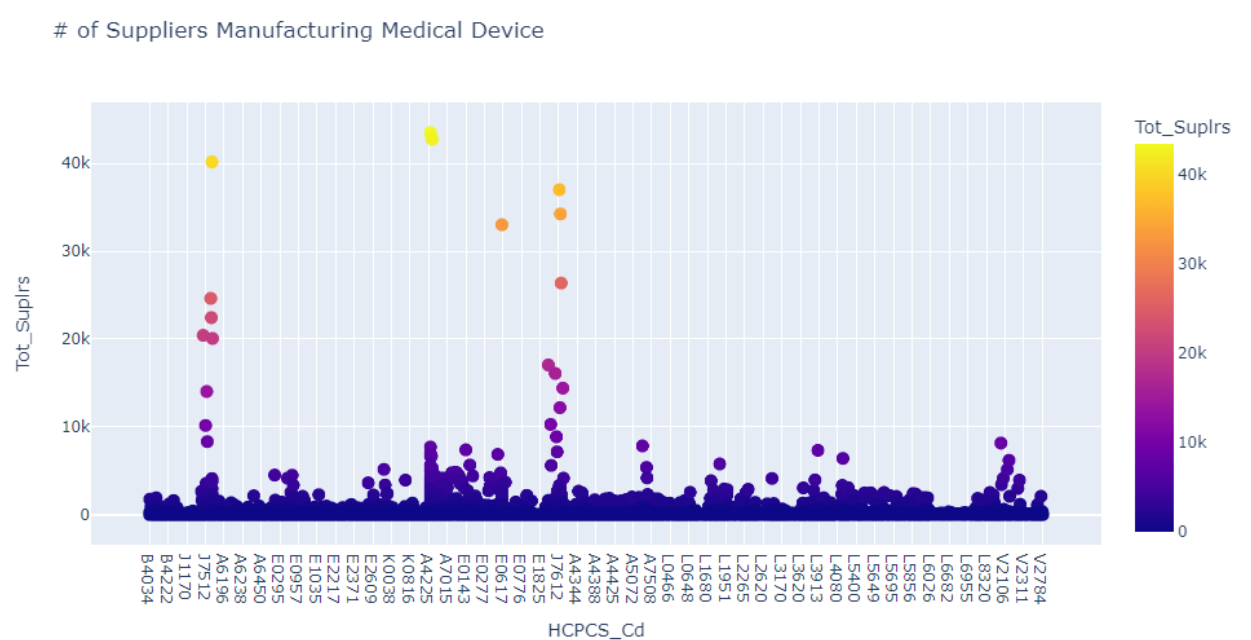
Using the describe() method, we were able to conclude that there are an average of 447 referring providers for an item, an average of 89 suppliers for an item, and an average charge of \$644.40 to a provider. Rfrg\_Prldr\_Geo\_Cd contains categorical values, the mode was used to fill in the missing values. Tot\_Suplr\_Benes contains numerical values, the mean was used to fill in the missing values. There were no duplicated rows and the transformed dataset was saved to a new file 'Wrangled.csv'.

Explanatory Data Analysis

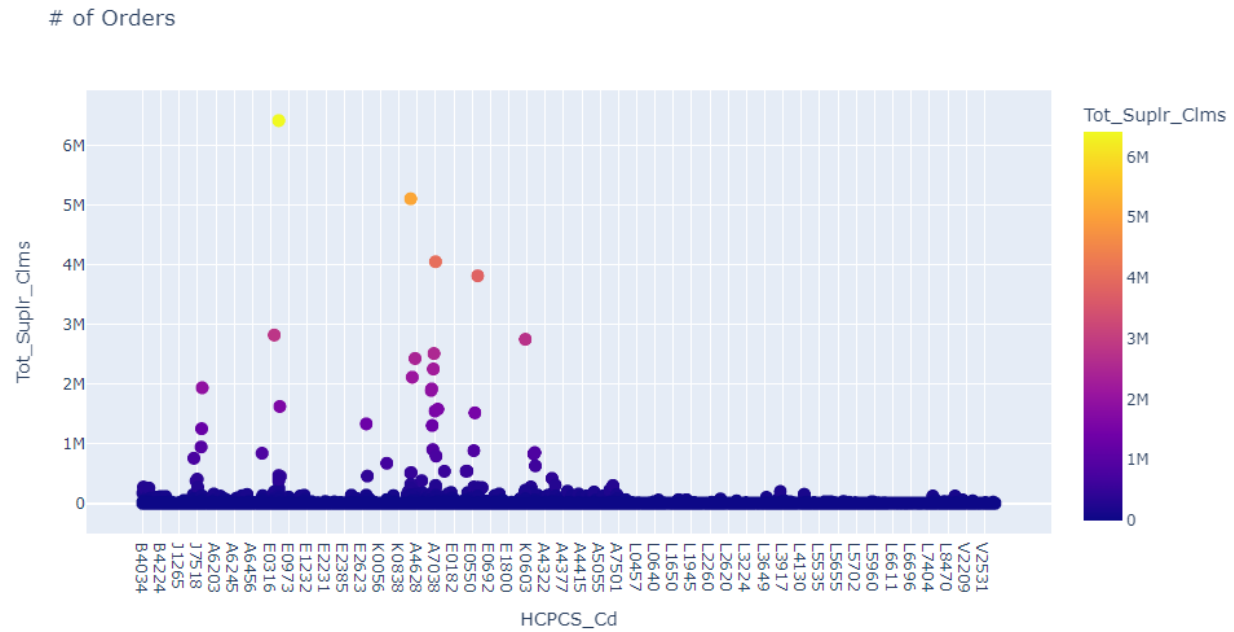
Created several scatter plots to explore the relationship between different columns:



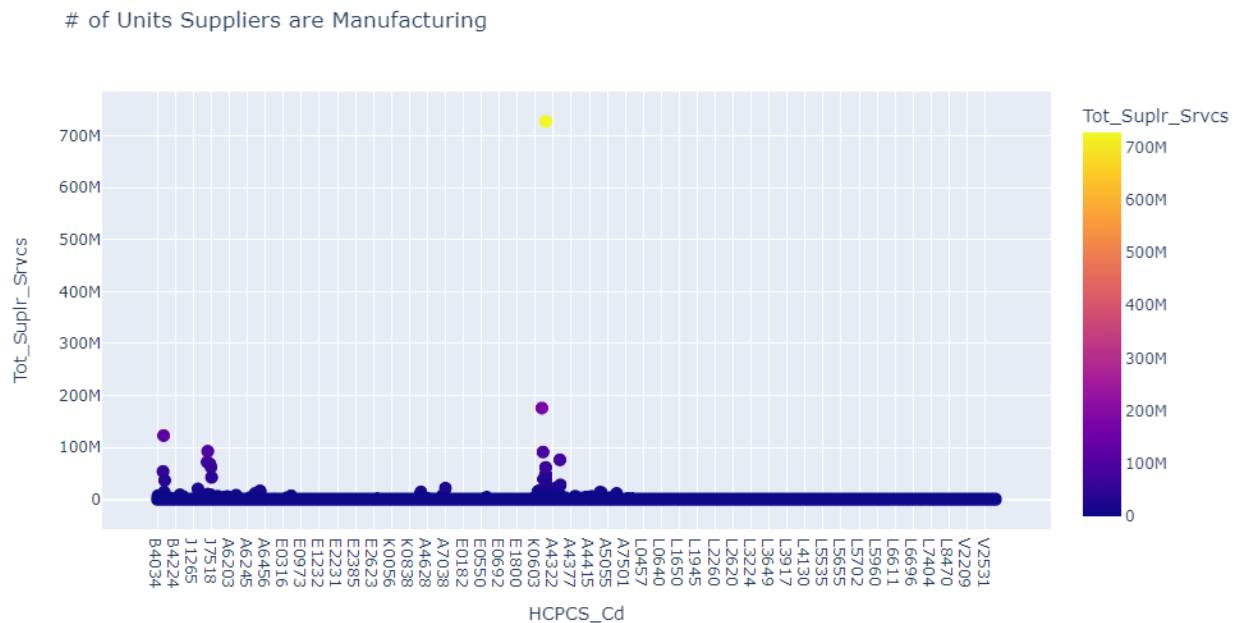
The blood glucose test/reagent strips has the most referring providers (233762 providers) meaning the most providers are ordering this item.



The blood glucose test/reagent strips also have the most suppliers manufacturing it.

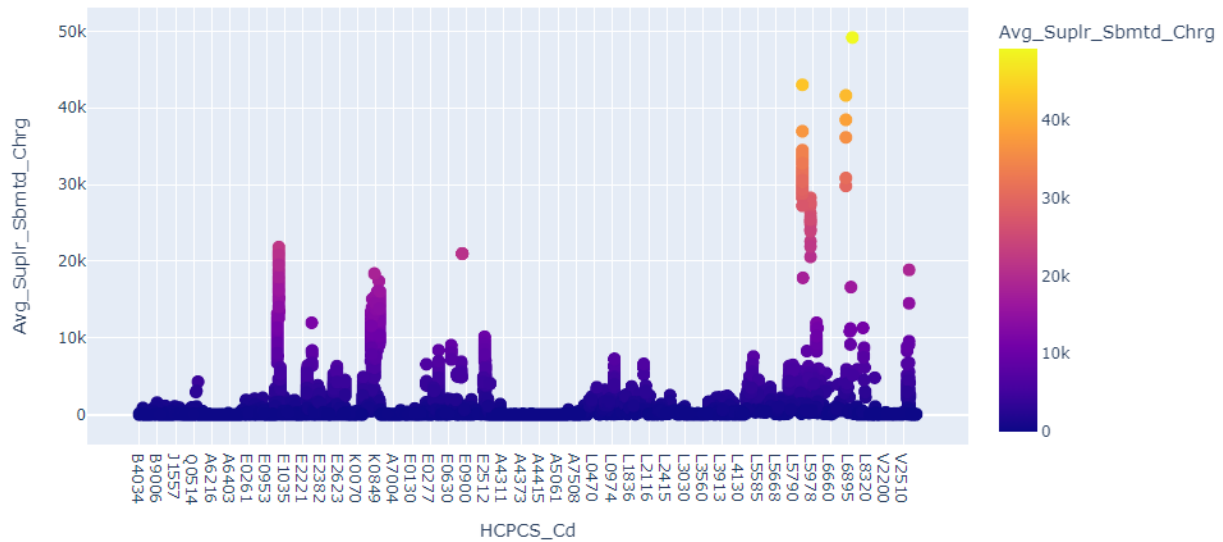


The blood glucose test/reagent strip is the 2<sup>nd</sup> most ordered item, the most ordered item is the oxygen concentrator.



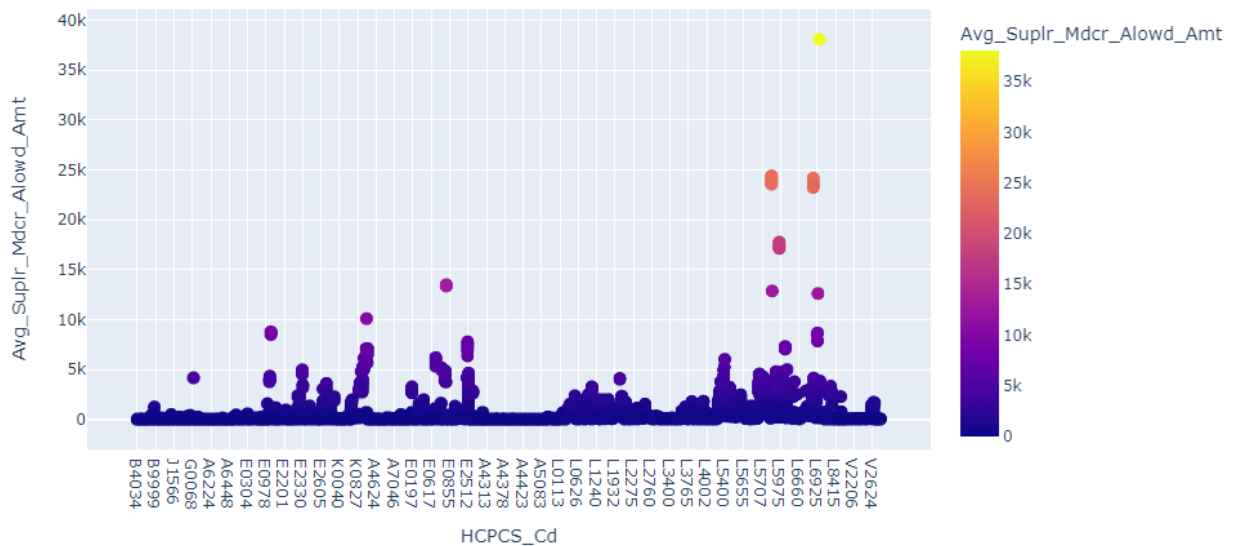
Only under 200M units of the blood glucose test/reagent strips were manufactured, the most manufactured item is the revefenacin inhalation solution.

### Purchasing Price



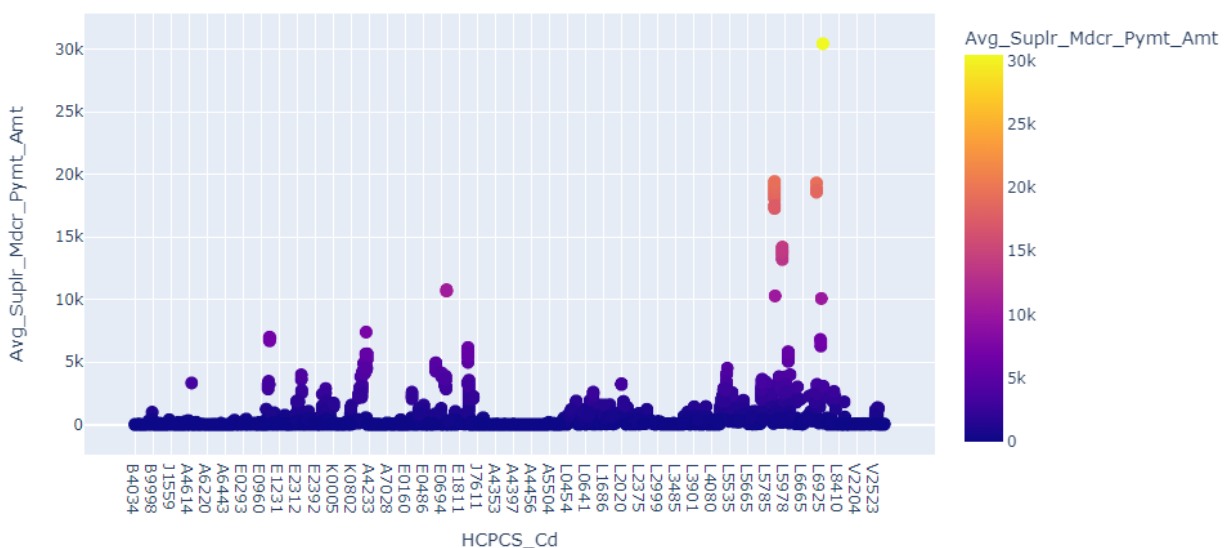
The blood glucose test/reagent strips were being charged on the lower end compared to the rest of the medical devices. The most expensive medical device is the electronic elbow. We also found that the average purchasing price for blood glucose test/reagent strip is \$61.44,

### Avg. Supplier Medicare Allowed Amount



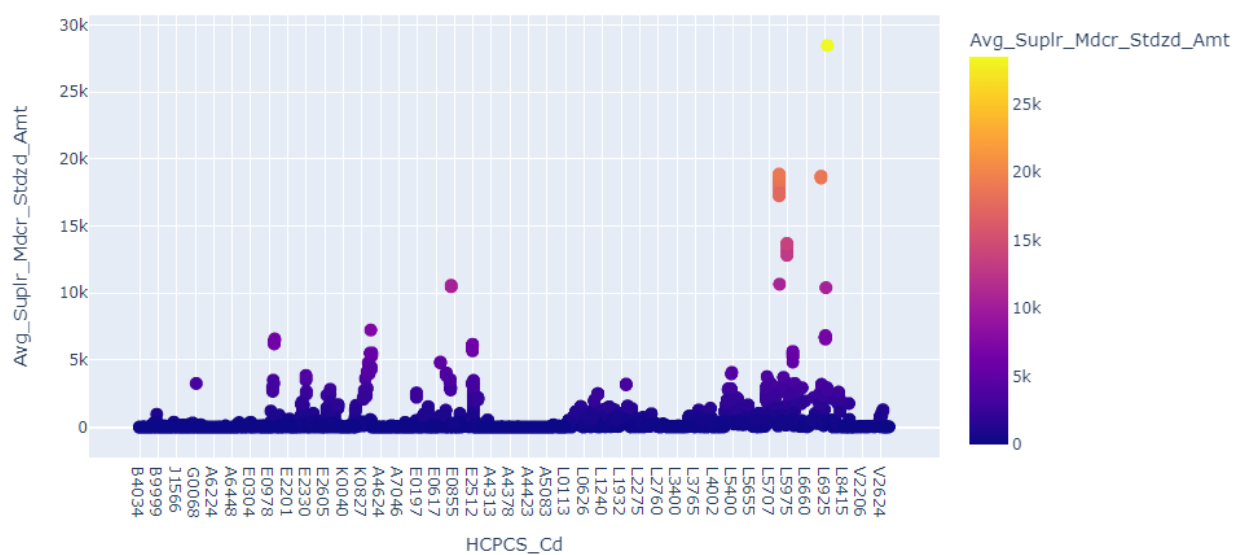
The blood glucose test/reagent strips were being covered by Medicare on the lower end compared to the rest of the medical devices, this makes sense as the blood glucose test/reagent strips do not cost as much as other items. The most covered by Medicare is the electronic elbow which is also the most expensive item.

Avg. Supplier Medicare Payment Amount

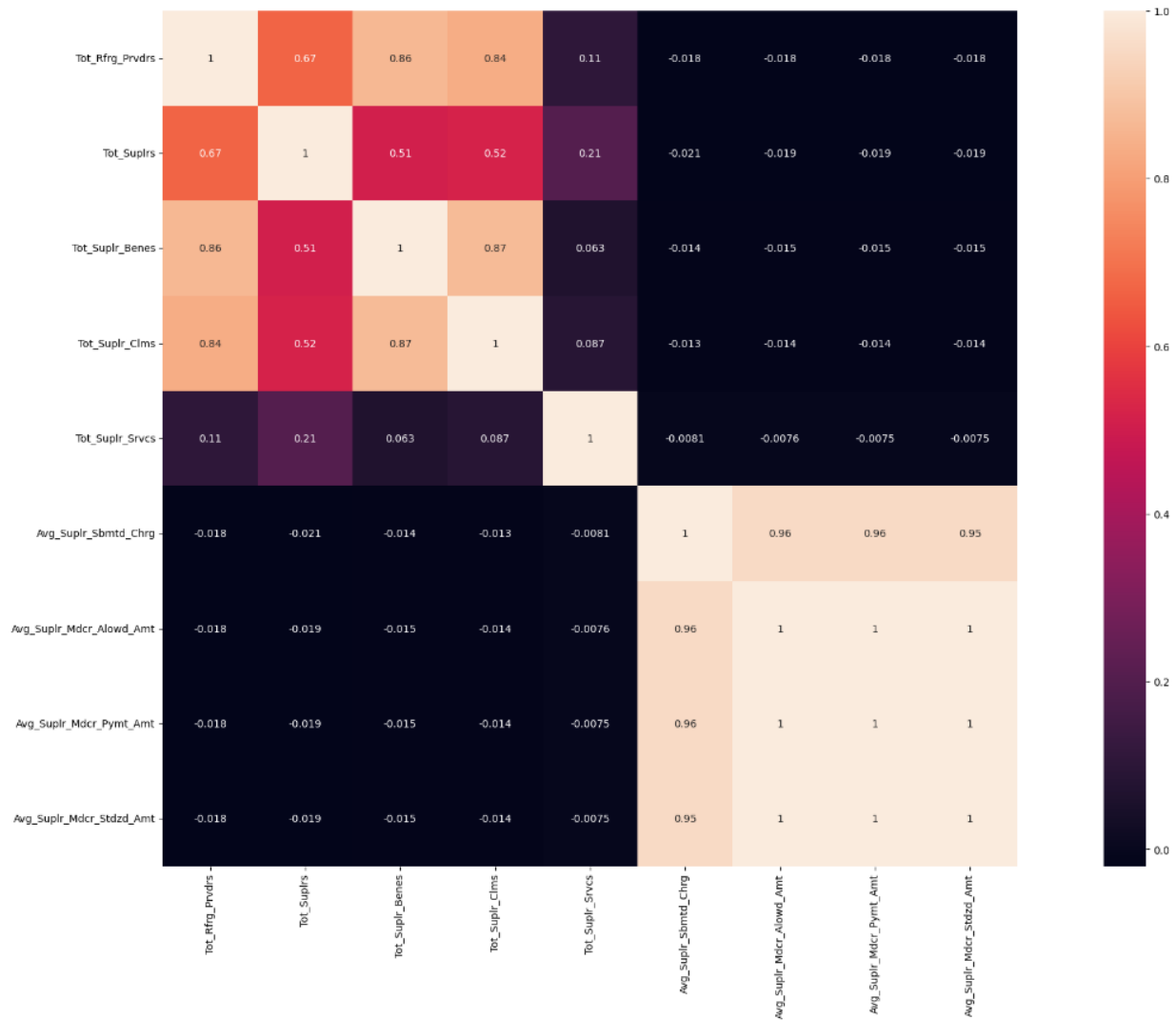


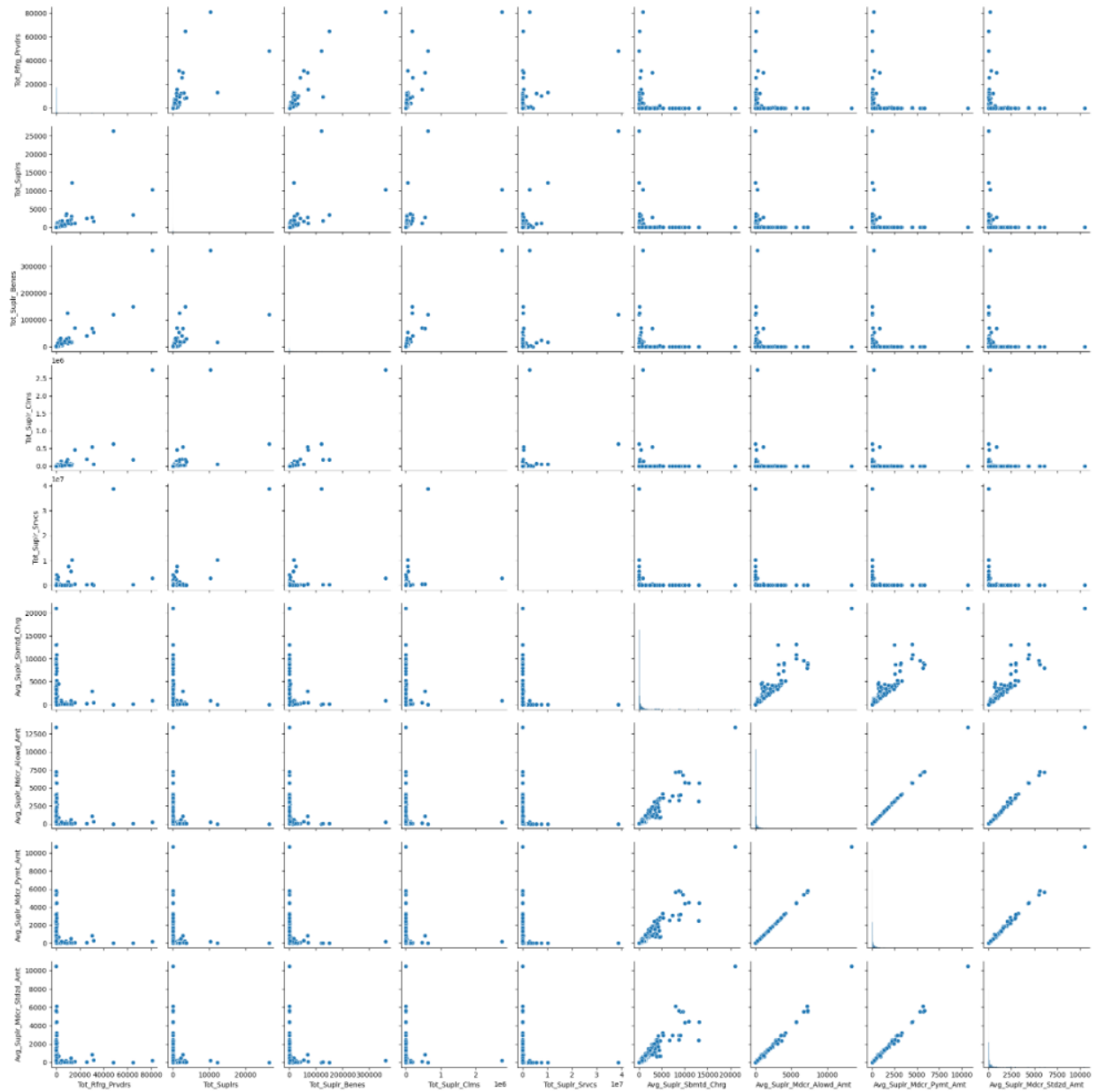
The results for the average supplier Medicare payment amount are the same as the previous distribution.

Avg. Supplier Medicare Standard Payment Amount



The results for the average supplier Medicare payment amount are the same as the previous distribution.





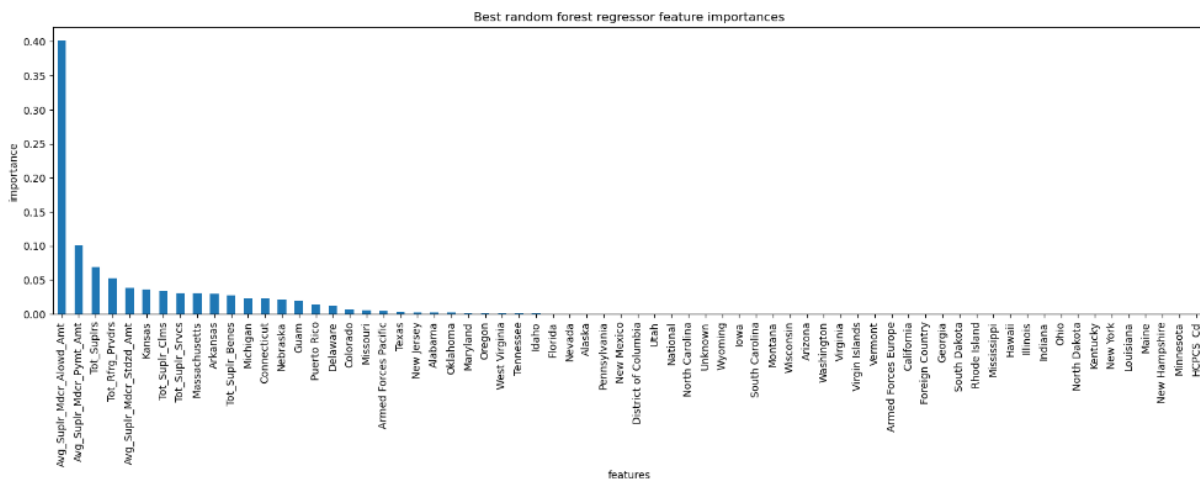
There were a few strong correlations between the Medicare payment columns; however, we did not see any other strong correlations between the other features. The features will need to be further analyzed in the next section.

## Pre-processing and Training Data Development

### Best Estimators:

```
Avg_Suplr_Mdcr_Pymt_Amt      160.274626
Puerto Rico                 -1.384814
Avg_Suplr_Mdcr_Alowd_Amt     -33.559312
Avg_Suplr_Mdcr_Stdzd_Amt     -146.573691
```

The average supplier Medicare payment amount is the biggest positive feature, this makes sense as submitted claims are paid for by Medicare first and may sometimes contribute to the majority of the payment. Average supplier Medicare allowed amount and average supplier Medicare standardized amount is negatively associated with purchasing price, this is suggesting that the higher the Medicare allowed/standardized payment amount then the lower the purchasing price. We can interpret this as some products may be getting paid out by Medicare more than others.



The dominant top four features in common with the linear model are:

- Avg\_Suplr\_Mdcr\_Alowd\_Amt
- Avg\_Suplr\_Mdcr\_Pymt\_Amt
- Tot\_Suplrs
- Tot\_Rfgr\_Prviders

For the Linear Regression model:

- The mean cross-validated Mean Absolute Error (MAE) is approximately 3.42, with a standard deviation of approximately 1.06.
- The MAE calculated directly on the test set using the best estimator from grid search is approximately 3.48.

For the Random Forest mode:

- The mean cross-validated MAE is approximately 4.08, with a standard deviation of approximately 0.51.
- The MAE calculated directly on the test set using the best estimator from grid search is approximately 3.31.



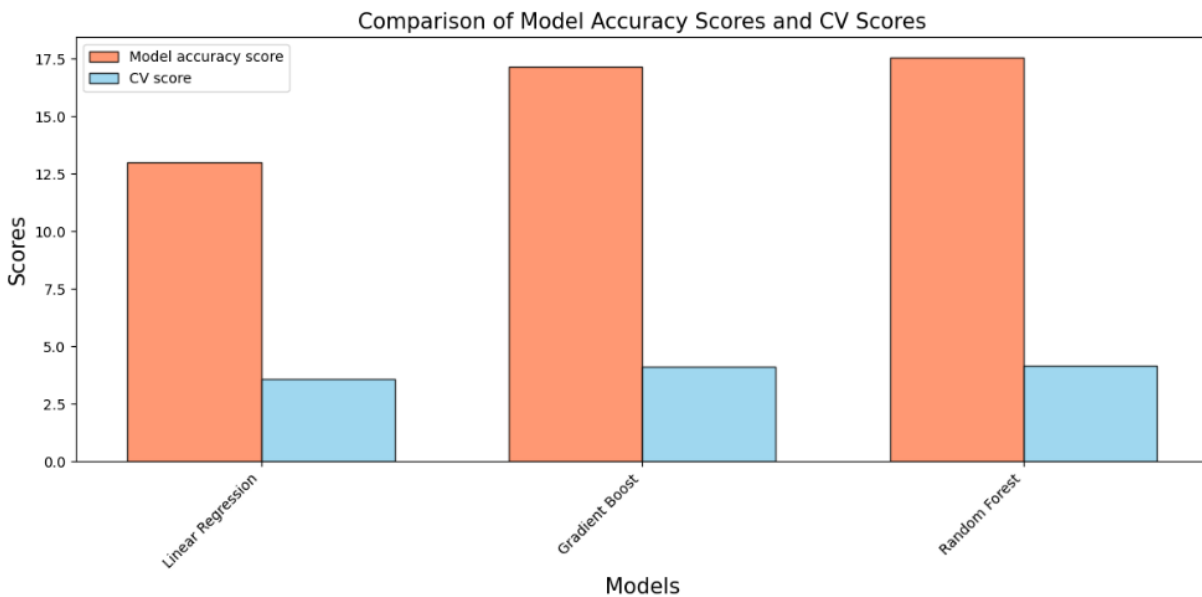
Conclusion: Comparing the mae on the test set, the random forest model appears to perform slightly better than the linear regression model.



There initial rapid improvement in model scores but levels off which indicates that adding more data does not significantly improve model performance beyond a certain point.

## Modeling

I conducted several models – linear regression, gradient boosting regression, and random forest regression.



	Algorithm	Model accuracy score
0	Linear Regression	12.991052
1	Gradient Boost	17.130138
2	Random Forest	17.553441

	Algorithm	CV_Score
0	Linear Regression	3.604310
1	Gradient Boost	4.138857
2	Random Forest	4.189683

Linear Regression appears to perform the best among the three algorithms as it has the lowest score. However, when considering cross-validated scores, Gradient Boost and Random Forest seem to perform better than Linear Regression as they have higher cv\_scores but higher cv\_scores also indicate slightly worse generalization performance compared to Linear Regression.

```
Train MSE: 3.866538869605473
Validation MSE: 5.111203427959549
Test MSE: 3.7461880052871277
Train R2 Score: 0.8463407132970207
Validation R2 Score: 0.8899563587190521
Test R2 Score: 0.7593648624604183
```

The model appears to perform well on the training set as indicated by the relatively low training MSE and high training R2 score. This suggests that the model explains approximately 81% of the variance in the training data. The similar R2 scores for the validation and test sets indicate that the model's performance generalizes well to unseen data.

## Conclusion

- Linear Regression model is the best model for predicting whether the blood glucose test/reagent strip will generate profit for the company, it performs slightly better than Gradient Boosting and Random Forest models.
- A lot of features are not important and have no impact on the target variable.

## Future Work

- Further analysis, such as examining additional evaluation metrics, tuning hyperparameters, or exploring different algorithms to help make a better decision on the best-performing algorithm.
- Integrating data from the company for more accurate predictions.