

Holmusk Data Assignment

Michelle Ng

Problem Statement

- How is costing affected by length stay in the hospital?
- If not, what are the major determinants of cost?

Datasets

- Clinical Data
 - 22 continuous data columns (medical history (7), preop medication (6), symptoms(5), lab_results(3), weight, height)
 - 4 categorical data columns (id, dates (2) , medical_history_3)
- Transaction Data
 - 1 continuous data column (bill id)
 - 2 categorical data columns (patient id, date of admission)
- Demographic Data
 - 5 categorical data columns (patient id, gender, race, resident_status, date of birth)
- Cost Data
 - 2 continuous data columns (bill id, amount)

Missing data and changes in data

- 502 rows of missing data
 - Decided to drop the records, as there is no meaningful way to impute the data
- Mapped errors in recording
 - F, m to female, male
 - India to Indian
 - Etc.
- Date to datetime variables
- Billing data - sum of cost per patient per admission

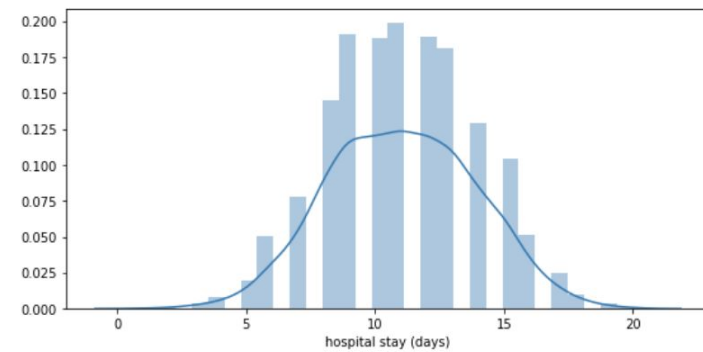
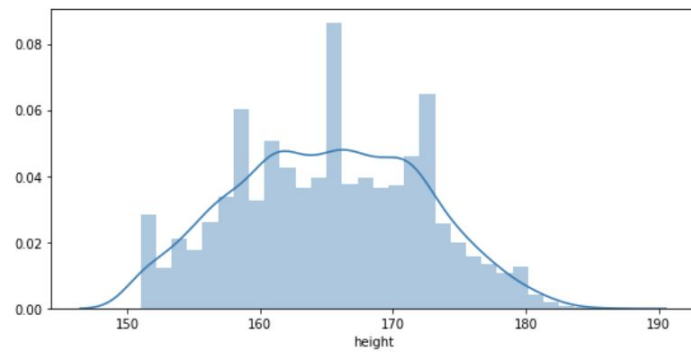
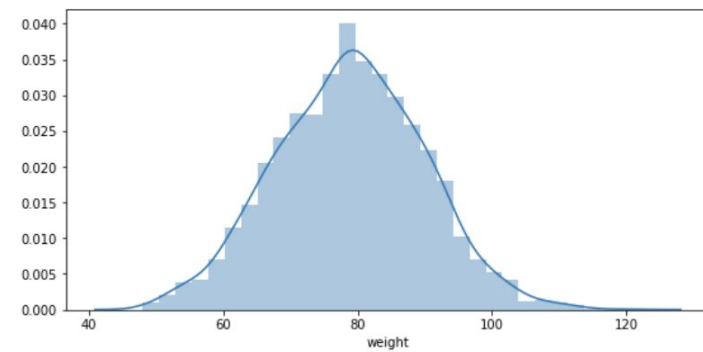
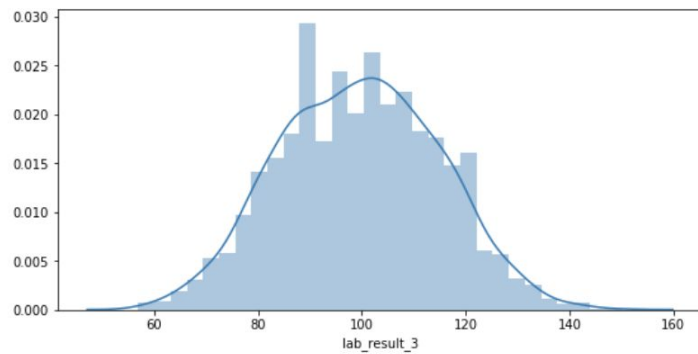
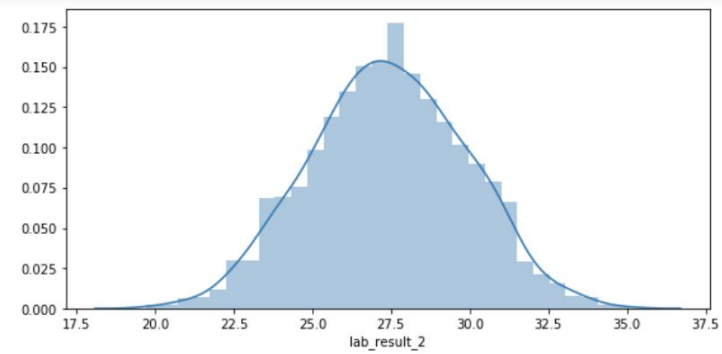
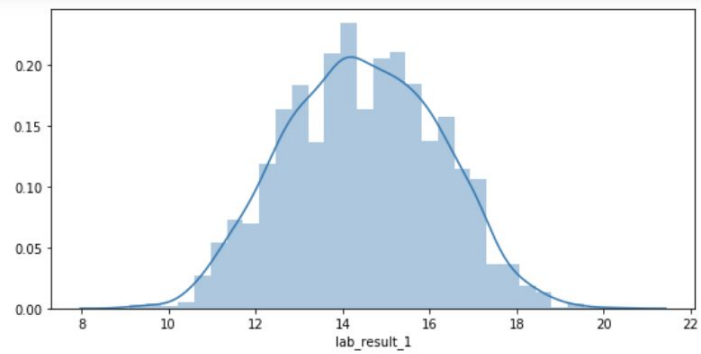
Feature Engineering

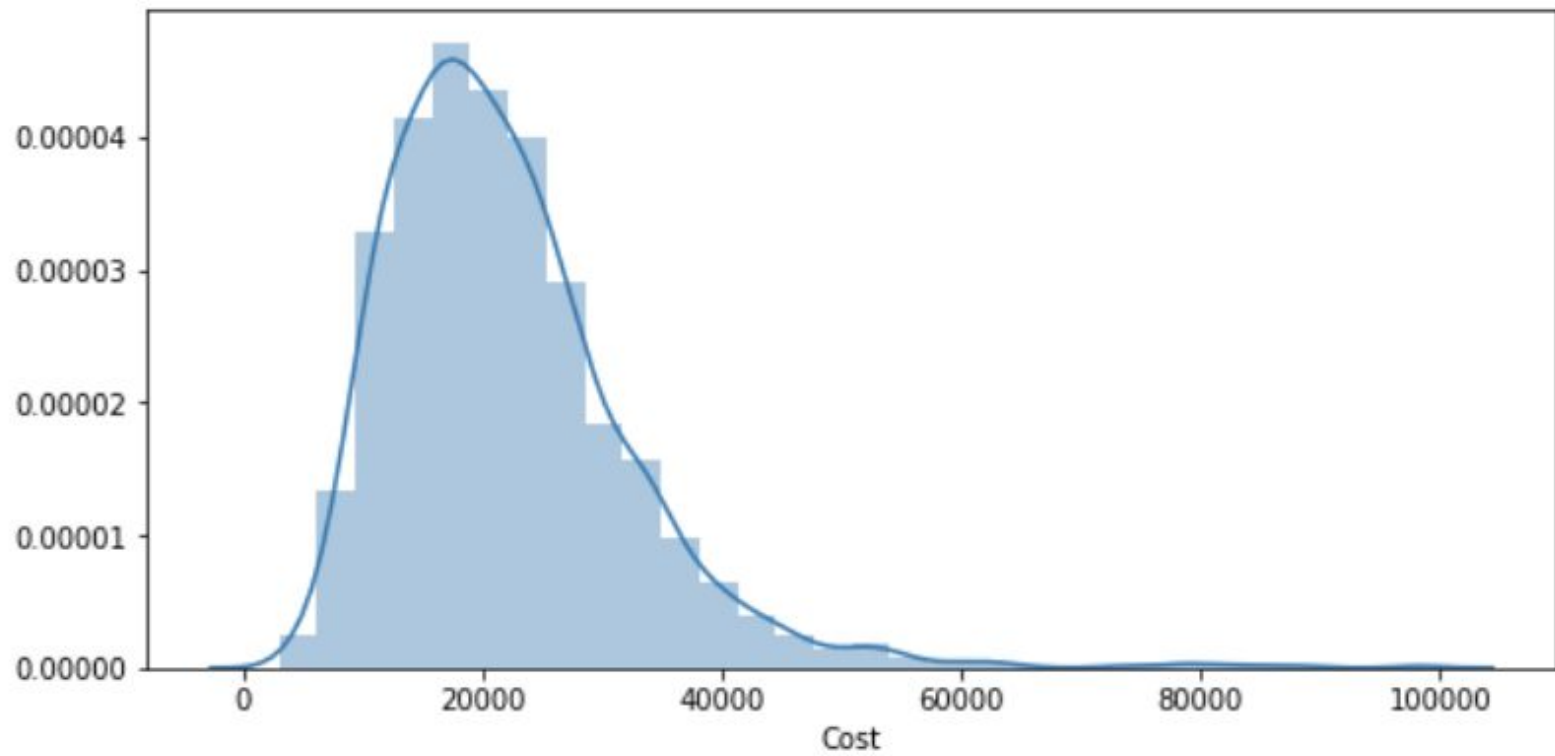
- Calculated hospital stay in days
 - (discharge - admission)
- Calculated age of patient upon admission
 - (admission - date of birth)

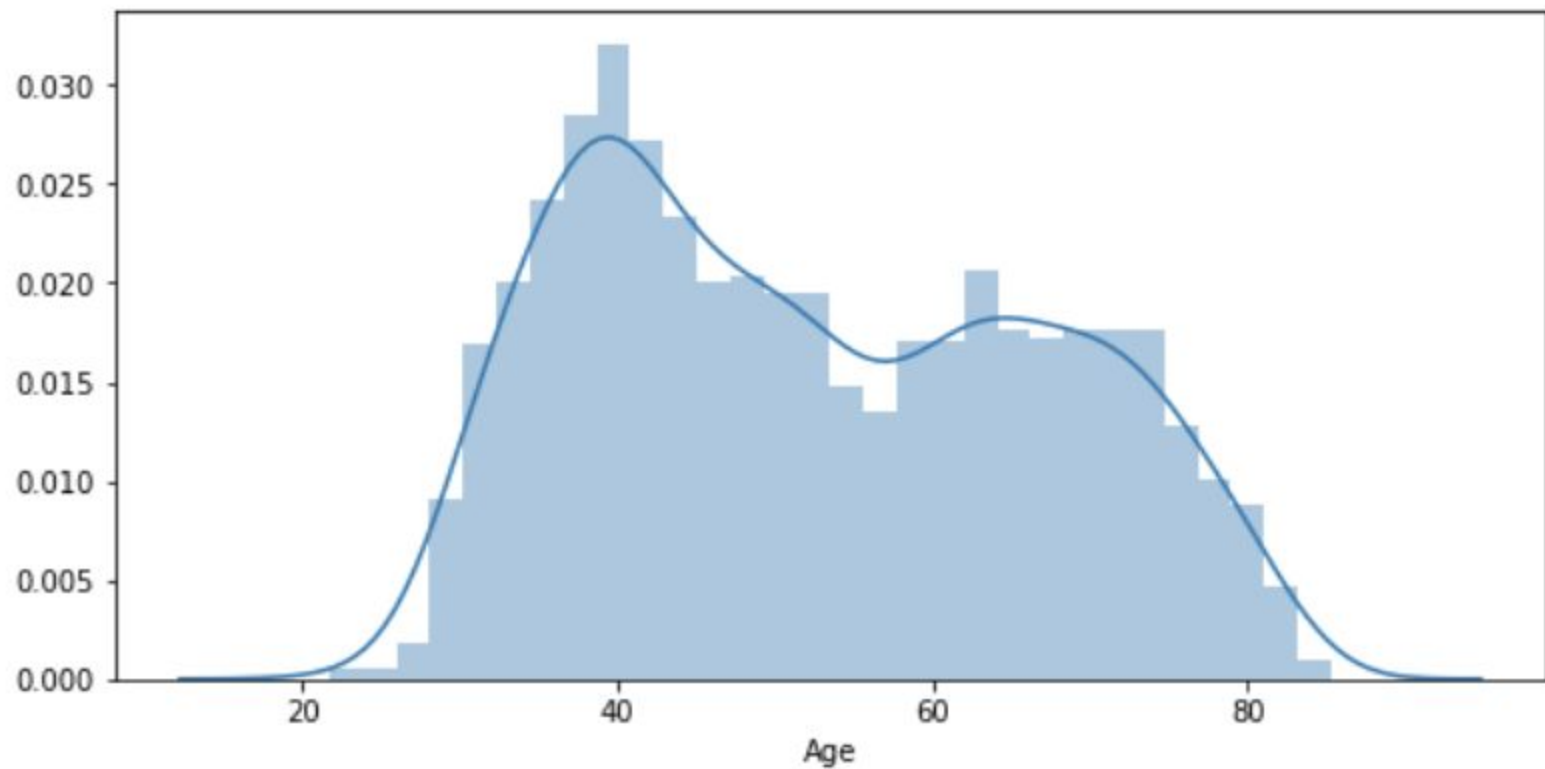
Assumptions

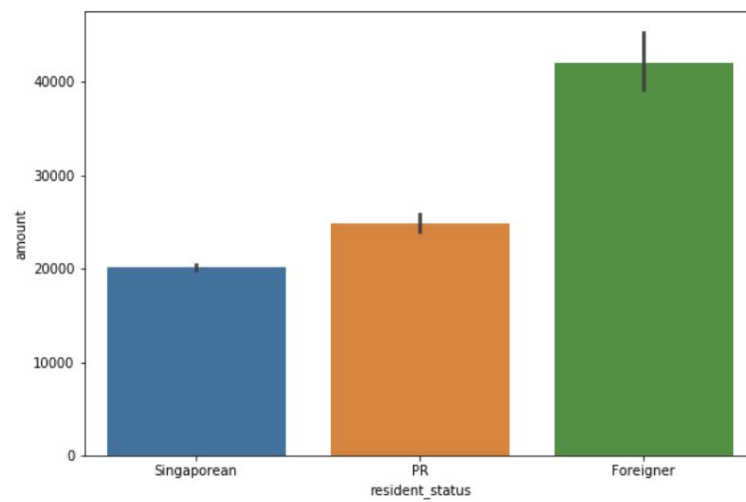
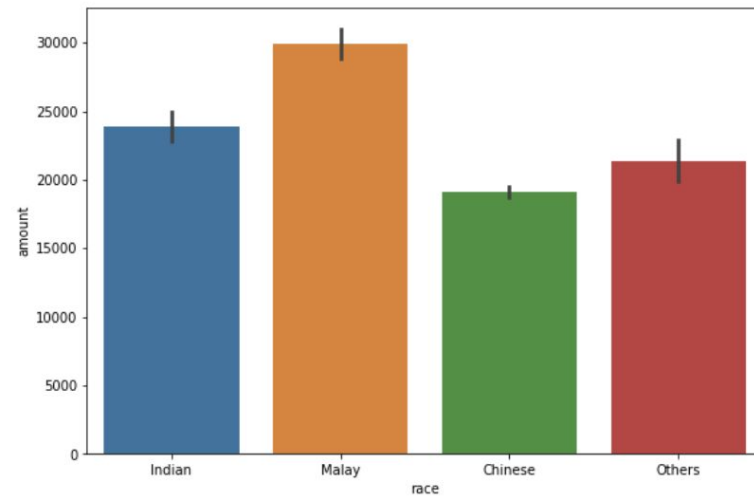
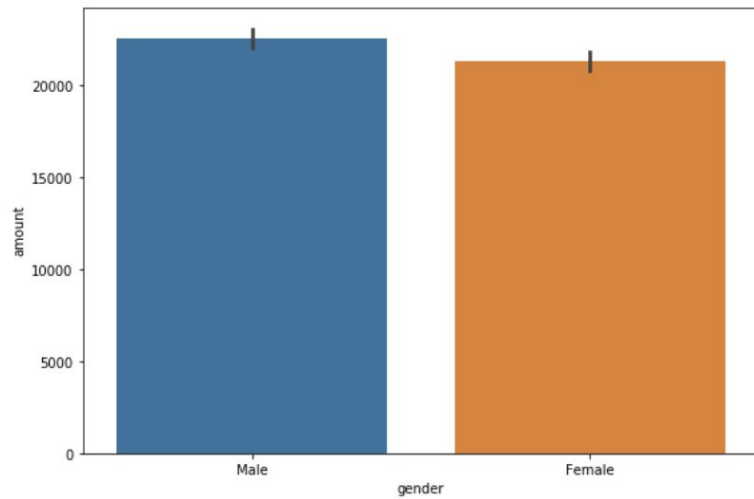
- Central Limit Theorem
 - Large dataset ($n = 2898$)
- Medical history, preop medications, symptoms
 - Dummy coded (0 no, 1 yes)

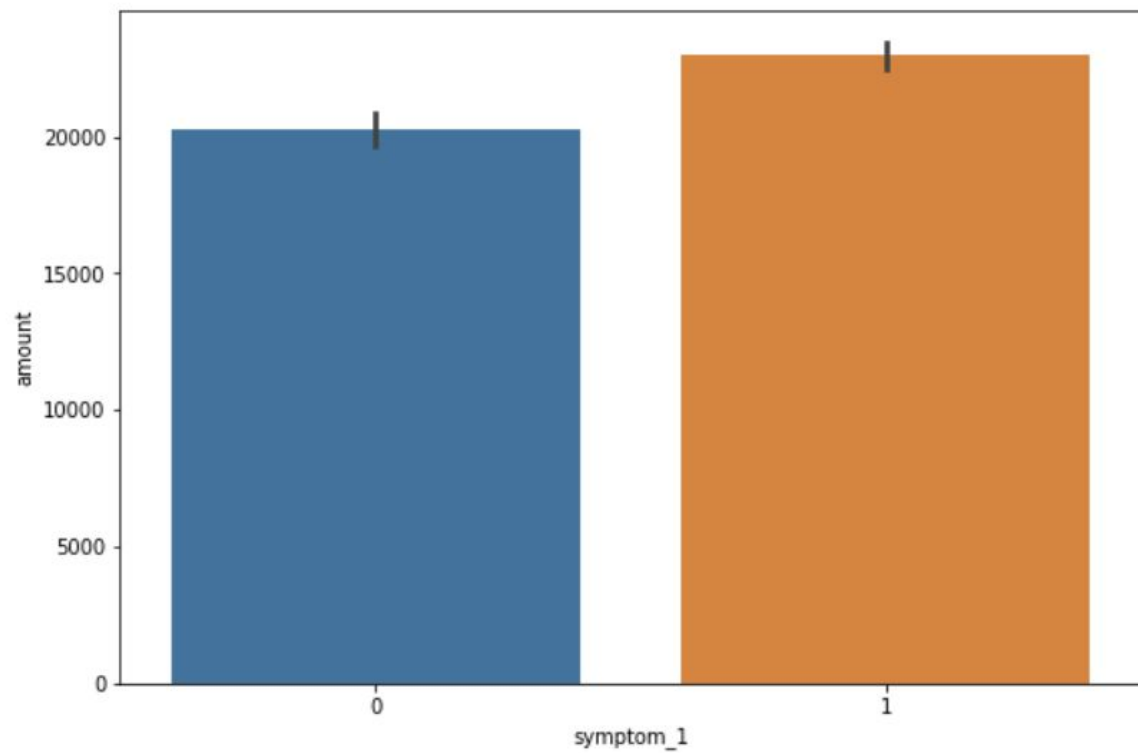
Visualisations

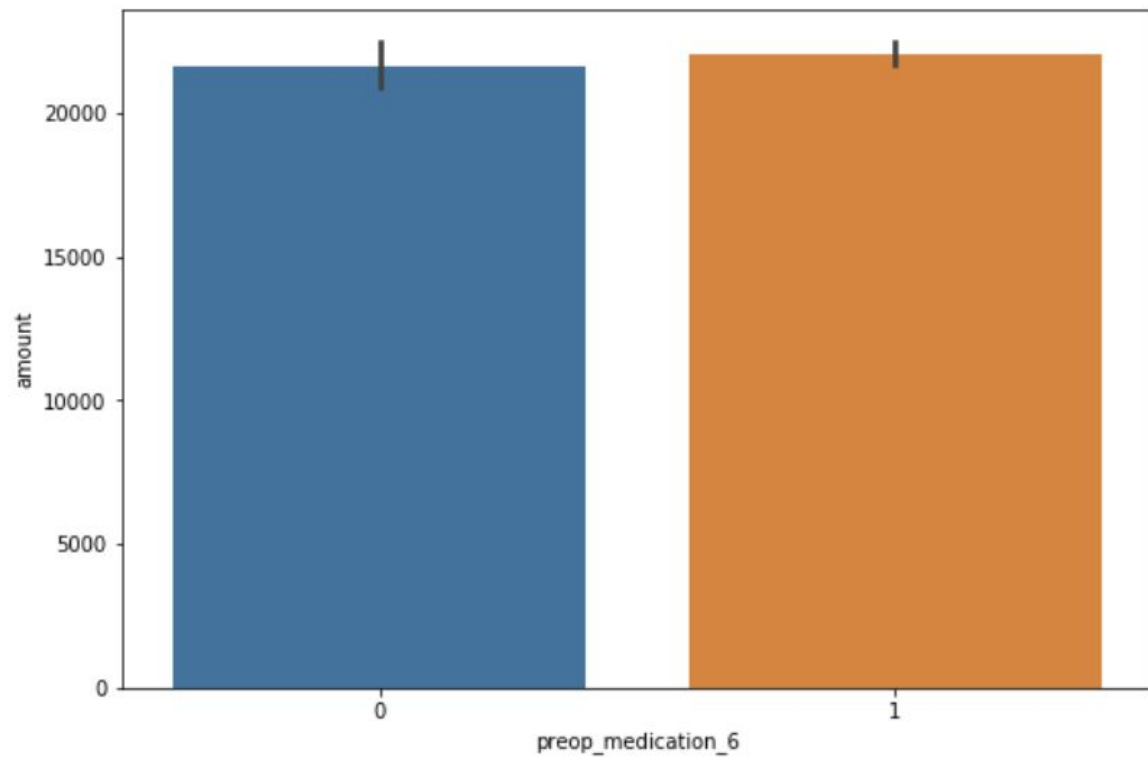






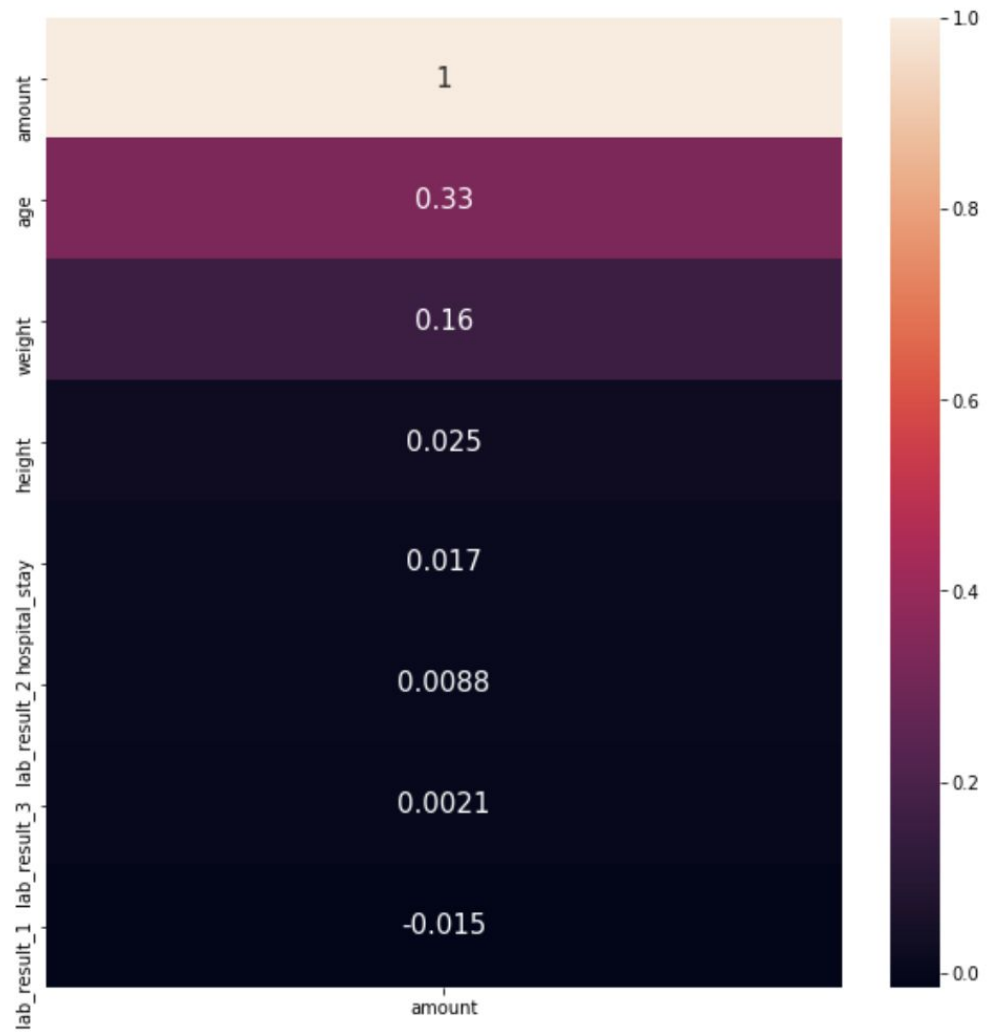






Statistical Tests

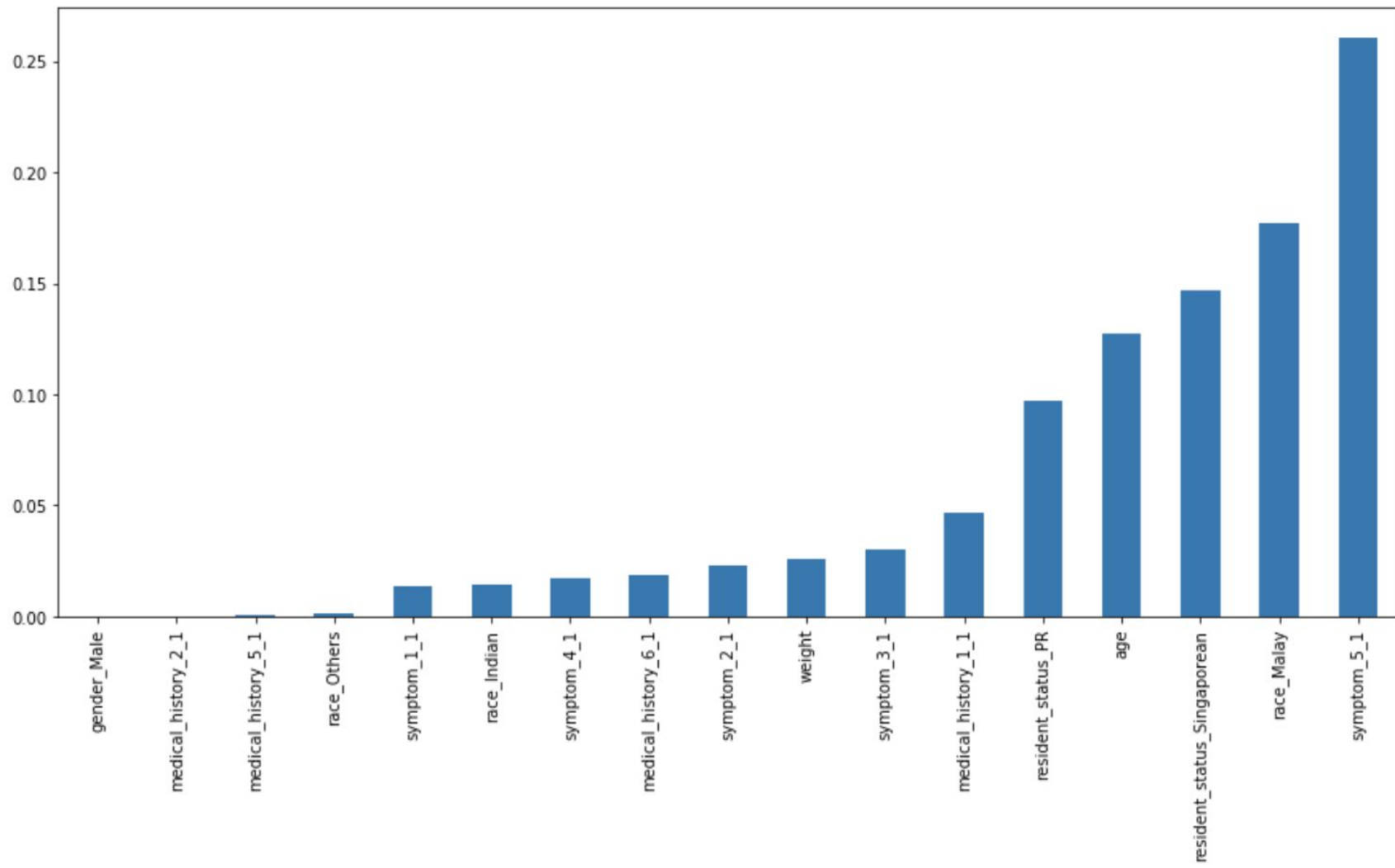
Correlation test



T-test/ANOVA

- Mean cost of treatment compared between categories of the same variable
- Significant difference in cost means:
 - Medical history 1, 2, 5, 6
 - Symptom 1, 2, 3, 4, 5
 - Gender
 - Race
 - Resident status

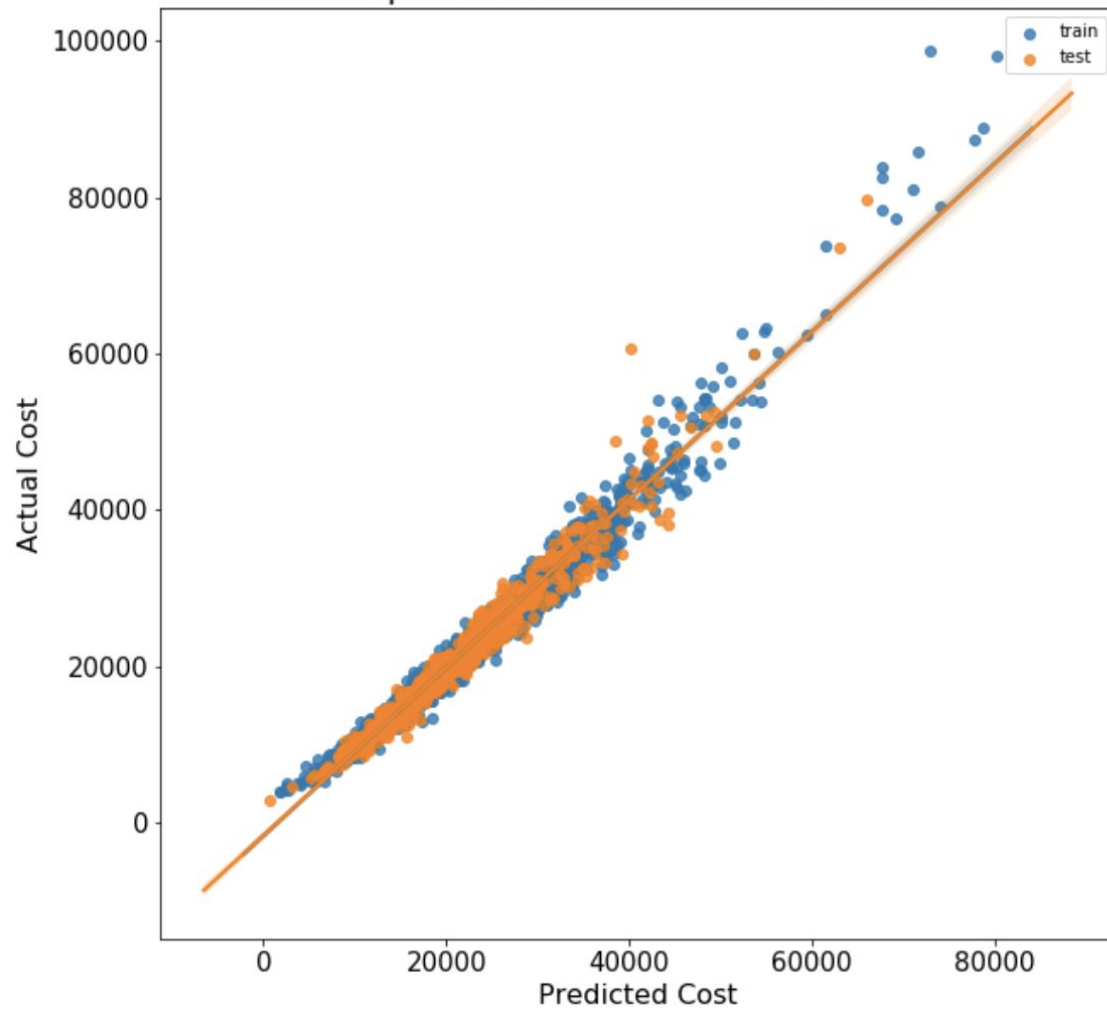
Models	Train Score	Test Score	RMSE
Linear Regression	0.921	0.927	2557.7
Lasso	0.921	0.927	2553.7
Ridge	0.921	0.927	2557.6
ElasticNet	0.751	0.799	4240.7
Random Forest	0.981	0.893	3015.6
Gradient Boost	0.970	0.958	1947.4
AdaBoost	0.532	0.388	7624.6



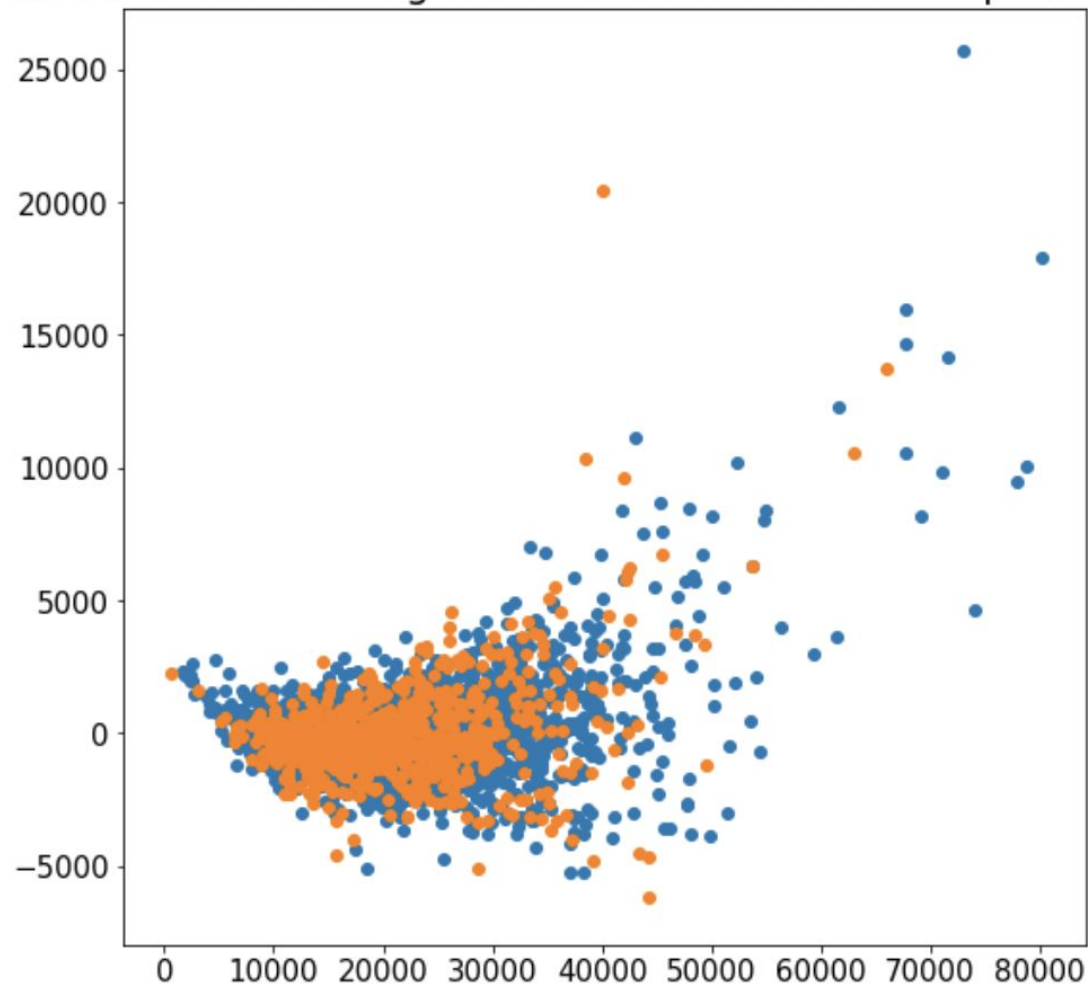
	ranking
symptom_5_1	0.260842
race_Malay	0.176679
resident_status_Singaporean	0.146738
age	0.127064
resident_status_PR	0.097446

Problems with Model

Scatterplot of Predicted Cost vs Actual Cost



Gradient Boosted Regression chart of residuals vs predicted y



Conclusions

- Stay in hospital is not highly correlated to cost of treatment ($r = 0.017$)
- Gradient Boosting model was best able to account for the training and test dataset.
- Symptom 5 is the most important feature, followed by race and age.

Limitations and Improvements

Limitations

- Model is not accurate for extreme ends of cost
- Model is not homoscedastic
- Each treatment is considered as independent from each other

Improvements to be made

- Exploring the effect of previous treatments on the cost
- Exploring other parameters for GridSearch
- Exploring what is the main cause of error at the extreme ends