SC1015 Mini Project

Medical Insurance in US

Members:

Eugene Ian Teo-Aldridge (U2222253B) Jiang Yifei(U2221092F) Huang Shu yang (U2210979A)

Table of contents

01 02 03

Problem Definition Data Preparation Exploratory Data & Dataset & Cleaning Analysis

04 05

Machine Learning
Techniques Re

Insights & Recommendations

O1 Problem Definition & Dataset

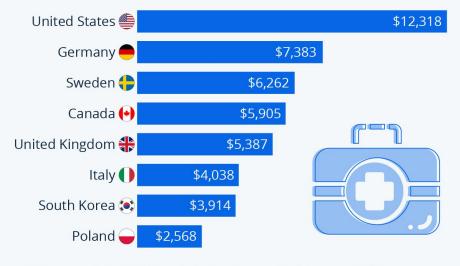
Expensive healthcare in US

Most common cause of personal bankruptcies in US

Medical insurance protects against high medical costs

The U.S. Has the Most Expensive Healthcare in the World

Per-capita health expenditure in selected countries in 2021



Includes government and private/compulsory and voluntary spending Source: OECD











Medical Insurance providers

Many different medical insurance plans of various premium prices

Slight mismatch in premium prices can result in loss of customers



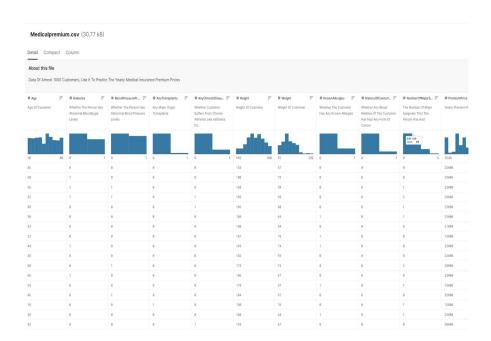
How can we develop an adaptive premium pricing model?

For medical insurance companies to remain competitive

For medical insurance policy holders to choose the best value for money plans

Dataset: Medical Premium in the US by Tejashvi

Many numerical and categorical variables for the different health and demographic factors

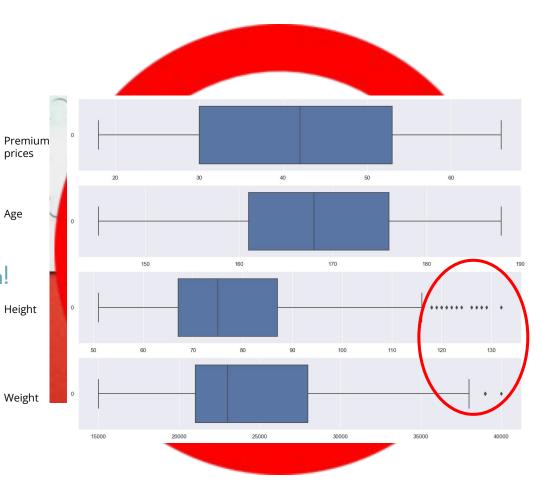


O2 Data preparation and cleaning

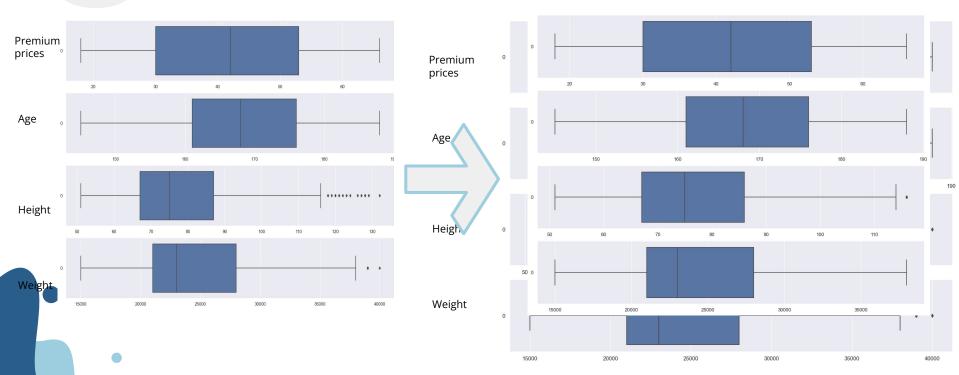
Cleaning the data

Numberial started data outliers

No need to fill in using imputation!



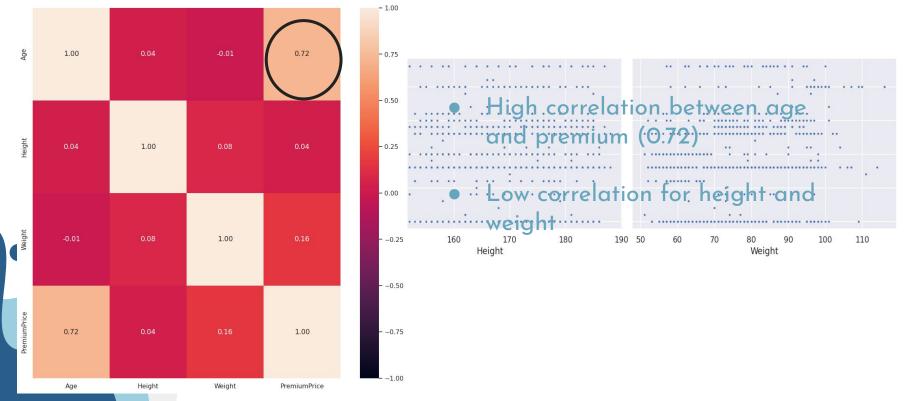
Removing outliers from numeric dataset



21 rows(2% of the dataset removed)

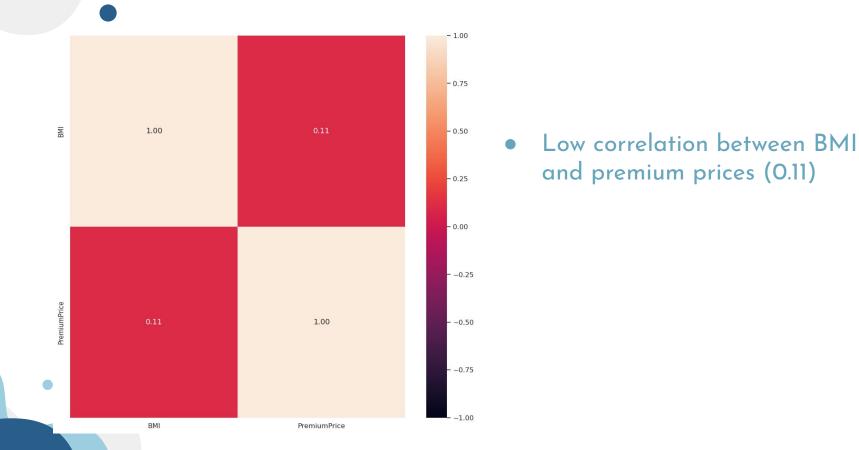
C3 Exploratory Data Analysis

Comparing the numeric variables against premium price



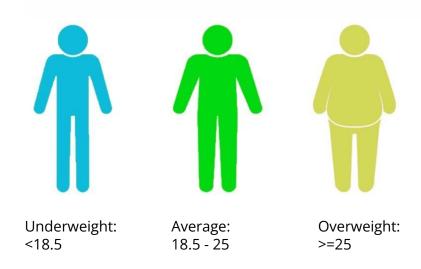
Exploring a new variable for BMI

BMI in relation to Premium prices

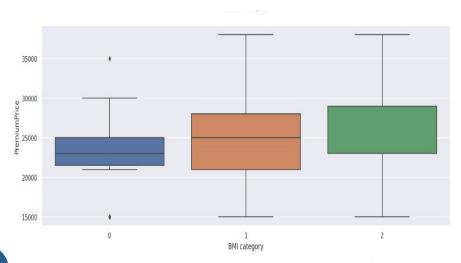


What if there is correlation different categories of BMI?





BMI Categorical findings



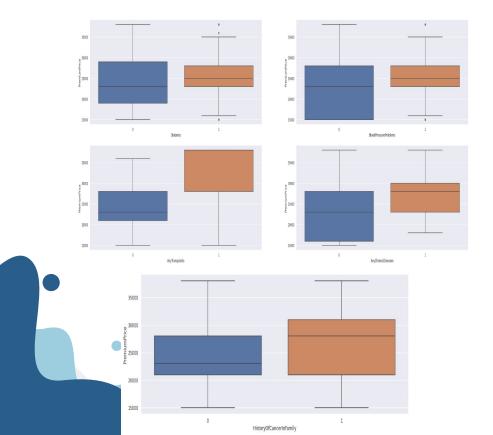
BMI Category premium price medians

Underweight: 23000

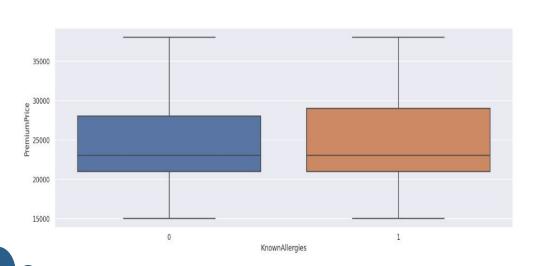
Average: 25000

Overweight: 23000

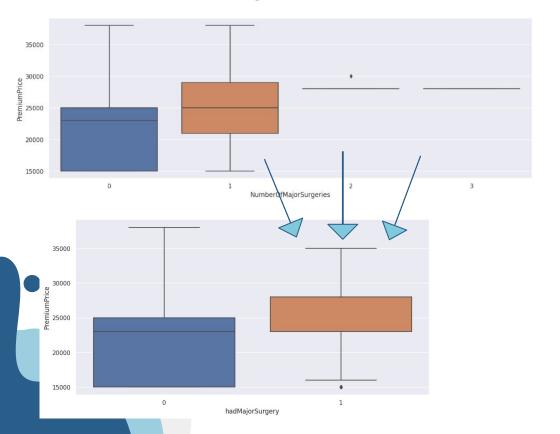
- Still low correlation between the different categories of BMI against premium prices
- Median premium prices lower for both overweight and underweight
- Not much difference in median between each categories
- This goes against our initial hypothesis that being overweight and underweight will increase premium prices



Diabetes,
Blood pressure problems,
Any Transplants,
Any Chronic Diseases and
History of family cancer
has clear relation between the
different categories and
premium prices

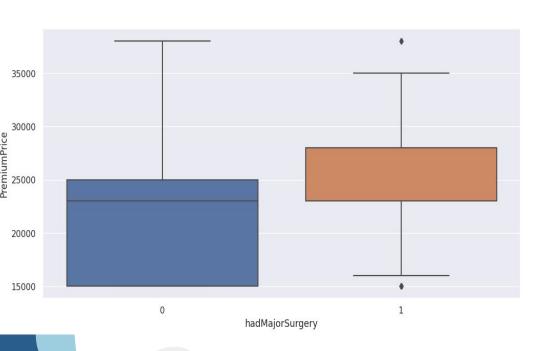


 Known Allergies does not seem to have any obvious relation between its categories and premium prices



 Number of Major Surgeries category for 2 and 3 number of major surgeries have too little data

Merge categories of 1, 2 and 3 together



 Had Major Sugery has a clear relation between its categories and premium prices

O4 Machine learning techniques

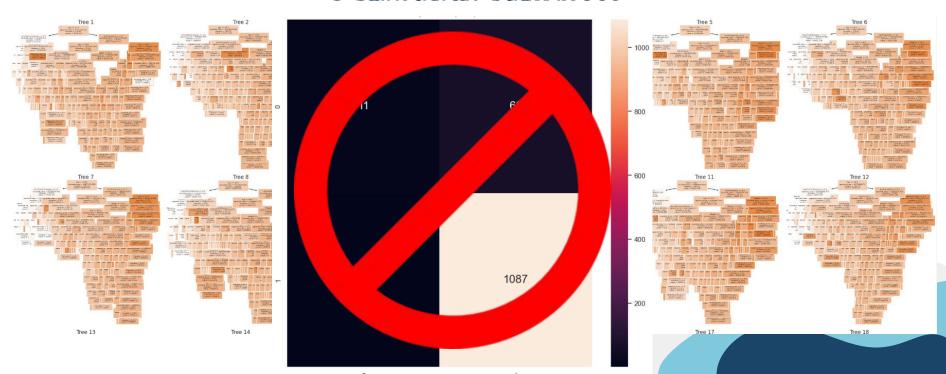
Model 1 - Random forest

Why Random Forest:

- Random forest does not get affected as much by over fitting unlike decision trees
- Random forest is more accurate, which is important when we have so many data we are using to predict premium prices
- Has access to feature importance to measure the importance of various features

Visualisation of result

Ocemf00 and Mahritrees



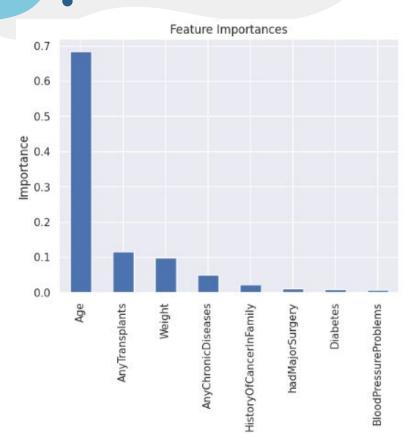
Visualisation and accuracy

Accuracy:

- R² value of 0.9, highly accurate predictions
- Root Mean Squared Error of 2000, fairly low given the average premium price being \$24000
- Low Mean absolute error of 1000, signifying the high accuracy of the prediction



Feature Importance



Importance of features:

- Age (68%)has the highest influence in random forest price prediction
- Of the remaining values, only AnyTransplants (11%), Weight (9%) & AnyChronicDiseases (4%) influence the decision making relatively heavily
- Blood pressure issues, diabetes and having surgery barely influences pricing (<1%)

Model 2 - CatBoost

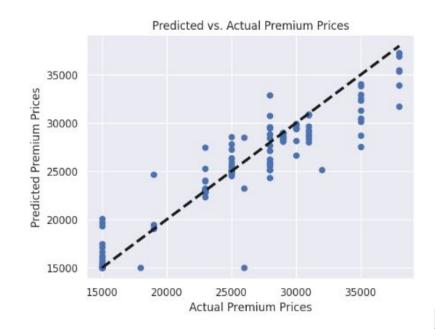
Why CatBoost:

- Specialises in prediction when provided with categorical data, which makes up the majority of our data
- Similarly to random forest, highly accurate and not susceptible to over fitting while also being resistant to noise because of the usage of gradient boosting
- Able to handle numerical data as well, while also able to provide Feature Importance

Visualisation and accuracy

Accuracy:

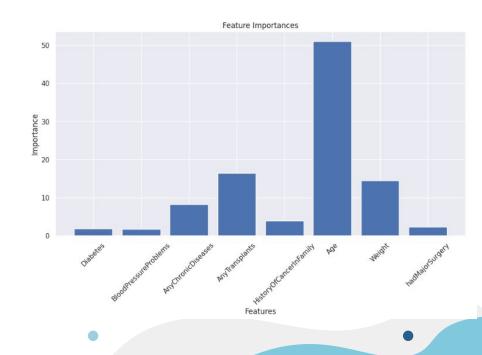
- Very Similar accuracy to Random forest, at 0.88
 R^2 value
- Root Mean Squared Error of 2100 and a Mean absolute error of 1300
- Slightly less accurate than Random forest, which can potentially explained away by the nature of Catboost relying on our categorical data with lower correlation



Feature Importance

Importance of features:

- Age has lower importance this time, at 51%
- Categorical values are more important predictors, with AnyTransplants sitting at 16% and AnyChronicDiseases at 8%
- Overall, less accurate than random forest in our case, but gave more insights into the categorical values importance and correlation



O5 Insights and recommendation

Lessons learnt

New Methods:

- 2 New Machine learning model, the categorical data specialised CatBoost and Random forest
- Feature engineering, making new features (BMI, BMI category) to find new things that may have strong correlations with target variable

Recommendations

What insurance companies should do:

- Understand which factors have most influence in the premium price in the insurance market and set their prices accordingly based on these features
- Remove factors with low importance like diabetes and high blood pressure to make insurance more appealing, expanding the customer base without impacting profitability

Thank you!