# Insurance charge estimation based on regression analysis

Dake Li[1] and Likun Lin[2]

[1]Department of Mathematics, Fudan University
[2]Department of Mathematics, City University of Hong Kong

December 1, 2023

**Abstract**

In this article, we selected a set of data from Kaggle website about medical insurance charge, provided with BMI, age, whether smoking or not. By analysing this dataset, we aimed to predict the corresponding insurance charges for new BMI-age-smoking data. We mainly set up a linear regression model. However, in this process, we found that there was a significant difference between the data respect to smokers and non-smokers in the sample. Therefore, we classified the smokers and non-smokers into two categories and established prediction models separately.

## 1 Motivation and Background

The role of medical insurance in society is multifaceted. In terms of individuals, we want to build up a suitable expectation of the expenses on ourselves' medical insurance. Therefore, given personal information, a prediction model of the insurance charge is very useful. As taught in this course, linear regression model is a very powerful tool for people to estimate the relation between two variables. Since every function has a best linear approximation with its first-order derivatives in a neighbourhood, linear regression is especially useful to model some relations locally. In this article, we mainly built up a linear regression model to predict one's medical insurance charge in the future, based on his or her BMI, age, and whether smoking or not.

## 2 Objectives

The main objective of this article is to set up a prediction model of medical insurance charges, based on one's BMI, age, and whether smoking.

# 3  Data Collections and Data Preprocessing

## 3.1  Data Collection

We fetch the dataset from the website Kaggle (Data source: https://www.kaggle.com/datasets/mirichoi0218/insurance). This dataset contains a bunch of samples of the medical insurance charges, together with the customers' personal information. The original data has too many features, such as the customer's sex, region, have children yet or not and so on, but in order to have clearer insights, we want to take only a few features into consideration.



Figure 1: Feature Importance (given by xgboost)

By XGBoost analysis, we conclude that the most improtant two features are age and BMI. In addition, from our life experience, whether smoking or not is also non-negligible. Therefore, we still label our selected data with "smoker" or "non-smoker".



Figure 2: Part of Original Data



Figure 3: Part of Selected Data

## 3.2 Data Preprocessing

Firstly, it is observed that the data label of 'smoker' and 'non-smoker' led to a significant gap in the insurance charge, just as figure 1 and figure 2 show.
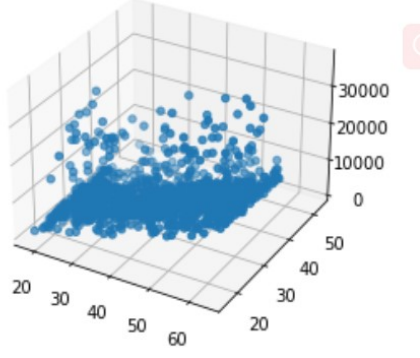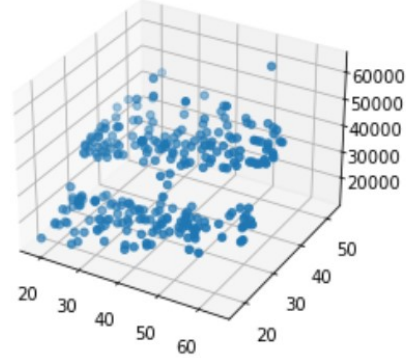


Figure 4: Nonsmoker's condition

Figure 5: Smoker's condition

According to this observation, we intend to develop different prediction models for each of the condition.

For each category(smoker or non-smoker), we randomly separate the raw data into two parts, that is, 70% of the raw data to be the training set, and the remaining 30% to be the test set.

# 4 Clustering and Classification

## 4.1 Smokers' condition

Firstly, we begin with the category of smokers.

By observation of the BMI-charge(figure 3) and the age-charge(figure 4) figure, it is reasonable to believe that the category of smokers can be also divided into two parts, where each part possesses different properties.
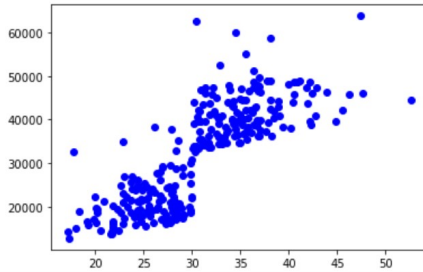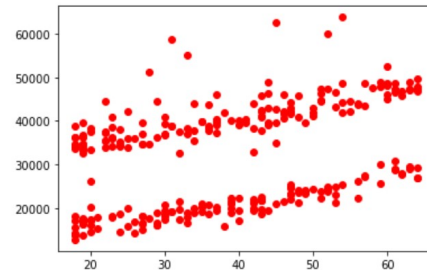


Figure 6: BMI-charges figure

Figure 7: Age-charge figure

The following pair-plot figure provides more evidence that the dataset of smoker category has some intrinsic partition.
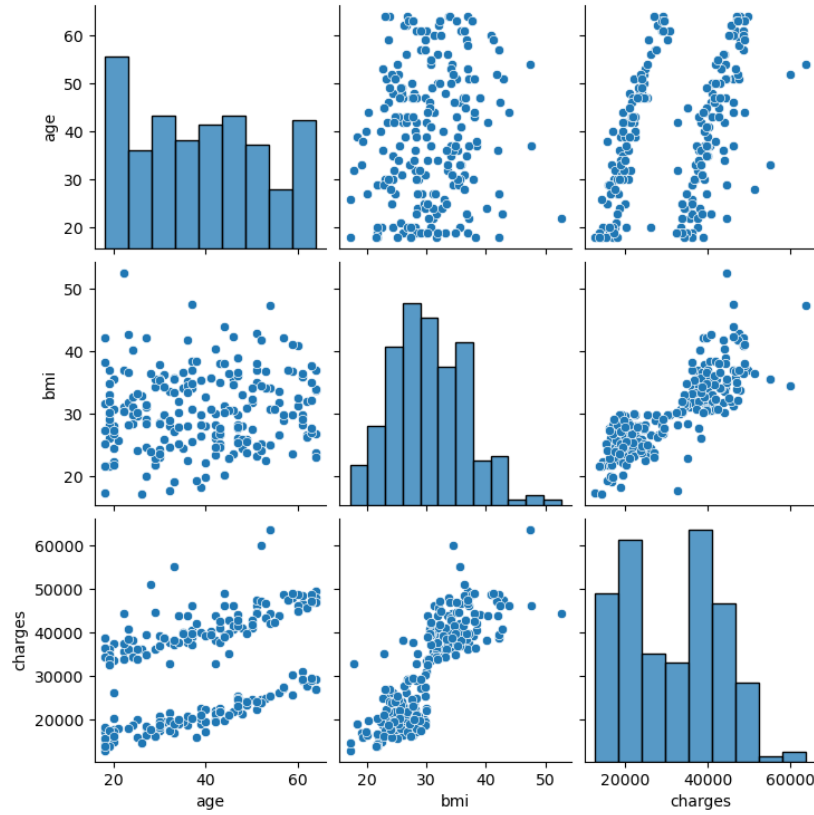
Figure 8: Pair-plot: age, BMI, charges

We conclude that the category of smoker can be further divided into two parts. The next idea is to apply linear regression on each part one by one.

The specific steps are:

1. Do clustering on the training set(70% of the whole smoker data).

2. Label the clusters by "class 1" and "class 2".

3. Using the obtained label, train a SVM model to classify.

4. Do linear regression on the each cluster.

5. For test set, we use the trained SVM to assign the label and apply the linear regression.

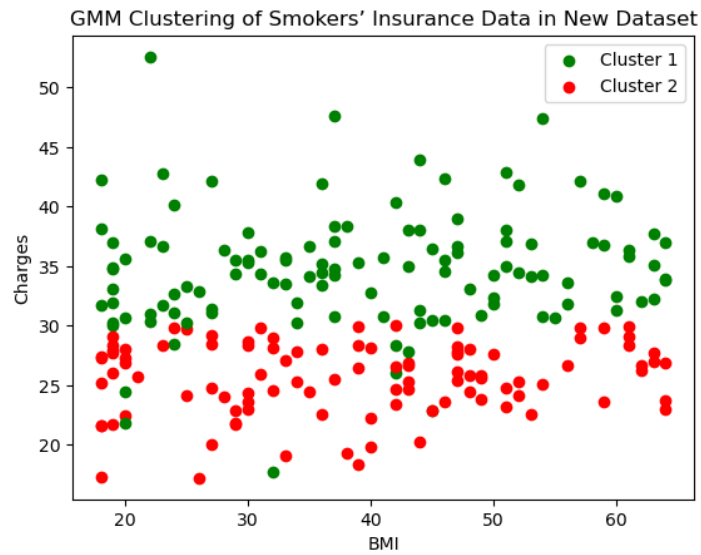The following figures show the result of out clustering:

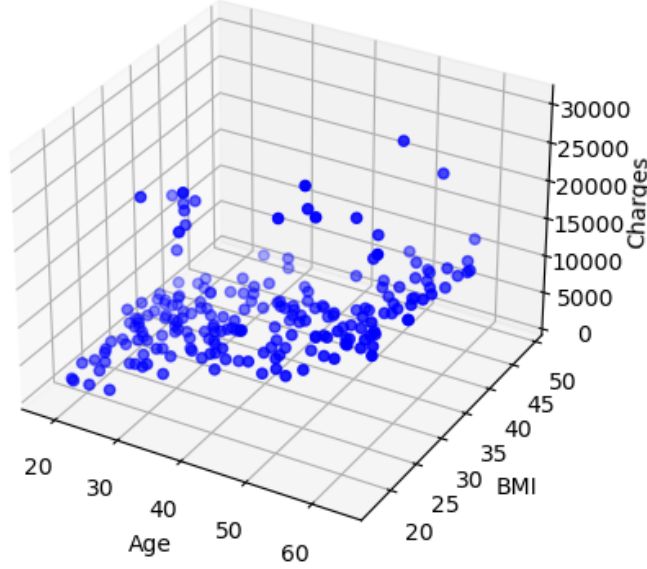Figure 9: GMM Clustering of Smokers Category



Figure 10: 3D plot for clustering result of smokers

According to the above analysis, the smoking population data can be divided into two clusters: high charge and low charge.

## 4.2 Non-smokers' condition

By observation of the scatter plot, we can still find two clusters in the non-smoker data.



3D Scatter Plot of Original Data

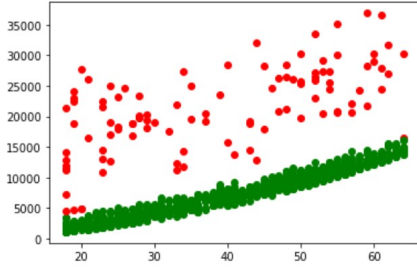We apply the same steps shown in the smokers' condition,



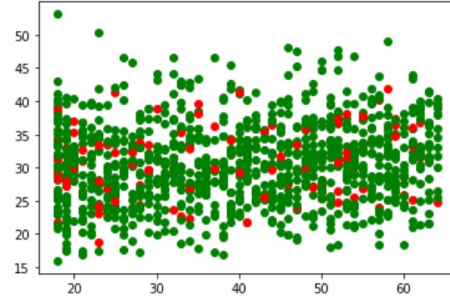Figure 11: GMM result under nonsmoker condition

Figure 12: GMM cluster result on age bmi of non-smoker condition

# 5 Ordinary Linear Regression model

We use the ordinary least squares method to do the linear regression. To be more specific, suppose $(x_1^i, x_2^i, y^i) \in \mathbb{R}^2 (i = 1, 2, \ldots, N)$ are our data points, where $x_1^i$ denotes age, $x_2^i$ denotes BMI, and $y^i$ denotes charge. Our goal is to find a vector $\mathbf{w} = (b_0, b_1, b_2) \in \mathbb{R}^3$ such that $\mathbf{w}$ minimizes the loss function:

$$L(\mathbf{w}) = \sum_{i=1}^{N} \left[ y^i - (b_0 + b_1 x_1^i + b_2 x_2^i) \right]^2 = \|\mathbf{y} - X\mathbf{w}\|_2^2$$

6

where

$$\mathbf{y} = \begin{pmatrix} y^1 \\ y^2 \\ \vdots \\ y^N \end{pmatrix} \in \mathbb{R}^N, \quad X = \begin{pmatrix} 1 & x_!^1 & x_2^1 \\ 1 & x_!^2 & x_2^2 \\ \vdots & \vdots & \vdots \\ 1 & x_!^N & x_2^N \end{pmatrix} \in \mathbb{R}^{N \times 3}.$$

To minimize the loss function, we seek for stationary point. Let $\nabla L(\mathbf{w}) = 0$ we get

$$\mathbf{w}^* = (X^T X)^{-1} X^T \mathbf{y}.$$

The result of linear regression for smokers category:
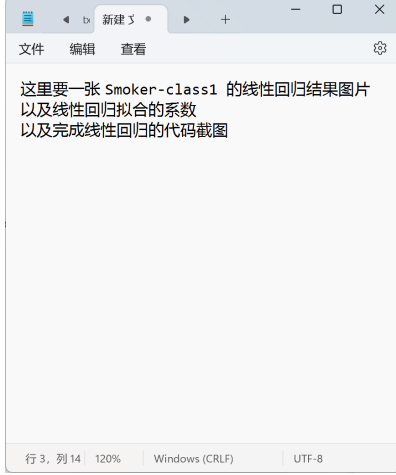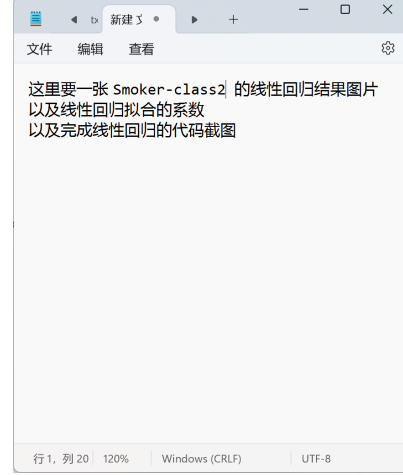


Figure 13: Smoker:class 1



Figure 14: Smoker:class2

The result of linear regression for non-smokers category:
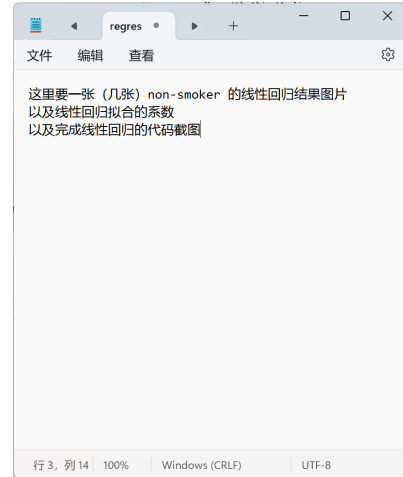


Figure 15: Non-smoker



Figure 16: Non-smoker

# 6 Conclusion

This paper focused on the relationship between 'charge' and its independent variables 'smoker', 'BMI', and 'age'. Our discussion is based on the fact that different

'smoker' labels would lead to a significant difference in the value of charges. Further experiments proved that 'smokers' can be parted into different classes for their charges with respect to 'age' and 'BMI', while 'non-smokers' does not produce such a separating effect so only one class is adopted in the result. Generally speaking, the charge of the above three classes all possess a strong positive linear relationship with age. The result shows that our model may produce more errors when predicting the non-smoker category. It is suggested that relevant data, e.g. family status, is lacking in the present model, and thus further research can be conducted with aims of observing the relation between family and charge.

# 7    Appendix

Data source:
https://www.kaggle.com/datasets/mirichoi0218/insurance
.csv version:
https://drive.google.com/file/d/1_VQk7_KoigrBS3NEPzBX5rdbI6CE4w43/view?usp=sharing
Python Prepossessing code:
See attachment