

Benjamin Laufer

Ecornell Data and Visualization Final Project

Part 1

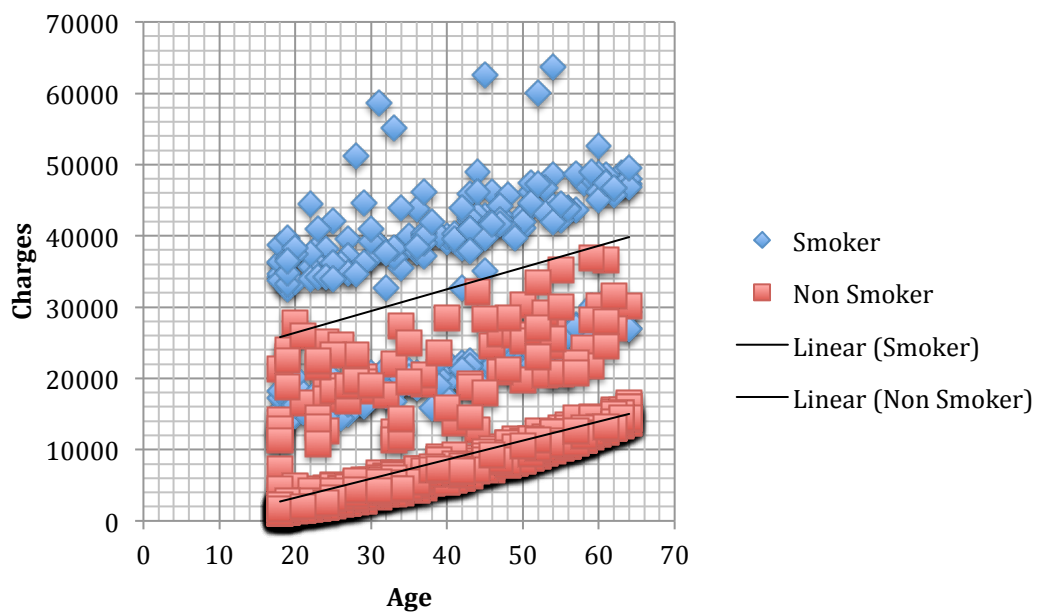
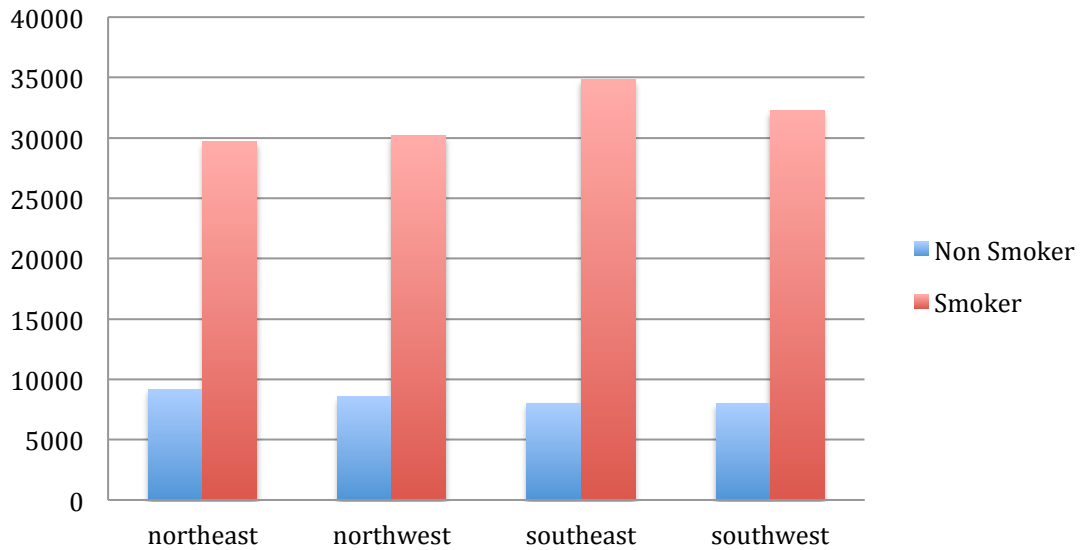
I am examining the affect of smoking on medical charges to help insurance companies better assess what to charge policy-holders. My hypothesis is that smokers spend significantly more on medical costs than non-smokers. The data will be acquired through insurance company looking at new policy-holders and what they were charged in first year. The KPIs I will focus on are the average and standard deviation of what smokers spend on medical bills compared to non-smokers. Factor that will be considered are age, body mass index(BMI) and dependents as I believe all of those will have an effect on medical costs. Considering these factors will ensure that my data is not skewed because of them. The insurance companies will observe new policy-holders and collect the charges over a course the first year. Then I will examine the data to determine if smokers and non-smokers have higher medical costs and if so how great the impact of smoking on charges. To control the sample and make sure it is representative of the population I will select just over 300 people from each region. In order to mitigate bias it is important to consider the factors I mentioned above as that people may lie about whether they are smokers or not

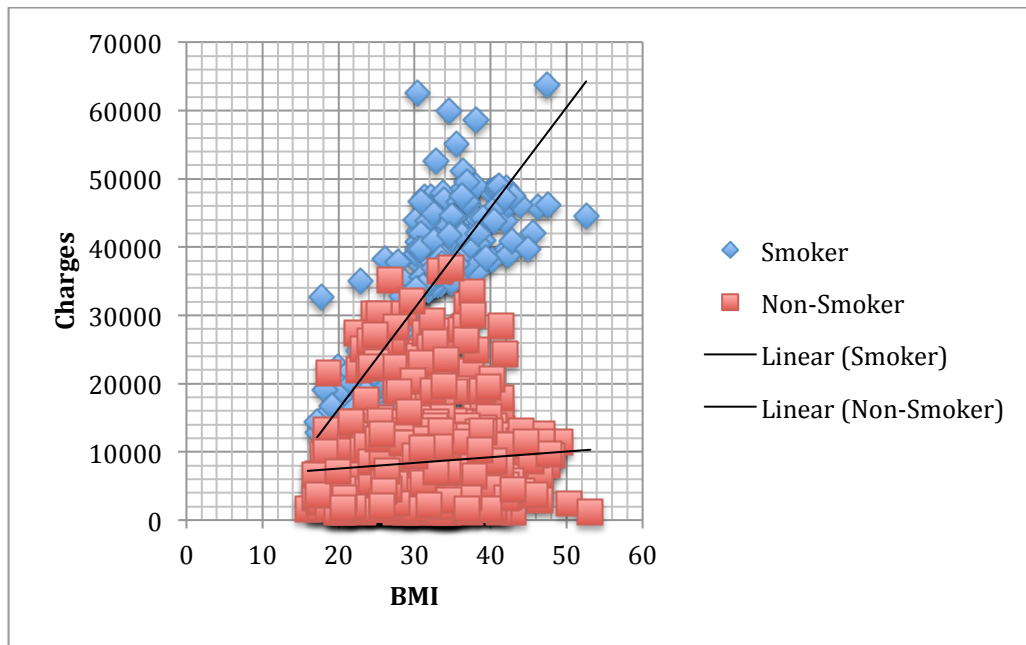
Parts 2 and 3

The summary statistics I will mainly focus on are average and standard deviation of medical charges between smokers and non-smokers. I believe it be suffice and resilient to outliers because of y large sample size

I will use bar charts comparing average charges of smokers to non-smokers, scatter plot graphs showing data on charges to smokers based on weight and age.

Average Charges By Region





	Average	SD	Median	Max	Min
Smoker	32050.23183	11541.54718	34456.34845	63770.42801	12829.4551
Non-Smoker	8434.268298	5993.781819	7345.4053	36910.60803	1121.8739

The chart above shows that on average smokers are charged over \$23,000 more than non-smokers which is about 3.8 times more than what non-smokers pay. The difference was even greater between the medians at over \$27,000. The standard deviation in charges was significantly greater for smokers than non-smokers and their min and max were both significantly higher.

The graphs above strongly display that smokers pay way more than non-smokers from insurance causes. Interestingly, the correlation is about the same for age of smokers and non-smokers. But there is a very strong increase in correlation in medical charges between smokers and BMI. My data does not take into account previous medical conditions or if people lied about being a smoker but the data is nonetheless conclusive that smokers pay way more than non-smokers in medical costs.

*My data for this project was pulled from Kaggle and the original file was entitled Insurance.csv

