



Coursera IBM-- Supervised Machine Learning: Regression Final Project

By: Liam Webster

6/13/2022

Objective:

The main objective of this analysis was to explore the predictability of individual medical expenditure via a select few parameters. Specifically factors such as age, bmi, smoking history, and number of kids was taken into account. With further investigation this model could then be used both in predictability or interpretation applications. Insurance companies would be interested in the predictability applications so they could accurately quote new customers. Consumers would be interested in the interpretation applications so they could address their attributes/risks and lower their expected medical expenditure.

Data:

The dataset chosen comes from ACME Insurance Inc. The dataset has 1338 rows/observations and 7 columns/features-- BMI, number of children, smoking history, region, and individual expenditure.

	age	sex	bmi	children	smoker	region	charges
0	19	1	27.900	0	1	southwest	16884.92400
1	18	0	33.770	1	0	southeast	1725.55230
2	28	0	33.000	3	0	southeast	4449.46200
3	33	0	22.705	0	0	northwest	21984.47061
4	32	0	28.880	0	0	northwest	3866.85520

...

HARSH SINGH. (2022, March). Medical Insurance Payout, 1.0. Retrieved June 12th 2022 from <https://www.kaggle.com/datasets/harshsingh2209/medical-insurance-payout>.

Approach:

First a general inspection of the data set was performed-- via plots and manual inspection. The couple of rows containing NaN data were dropped. The 'sex' and 'smoker' column were manually encoded-- 1 for female/0 for male and 1 for smoker/0 for non-smoker, respectively. A copy of the data frame was made and with said copy the 'region' column was one-hot encoded; with the other copy the 'region' column was dropped. Thus resulting into two clean data frames one with the original columns except 'region' is dropped, referred to as data frame 1. And one with the original columns except 'region' is one-hot encoded referred to as data frame 2.

Models:

First training and test splits were created with a 70 : 30 ratio respectively. These were declared globally and were used in each of the regression models. The first model trained was a simple linear regression model. Data frame 1 resulted in a 0.767 R2 score while data frame 2 resulted in a 0.771 R2 score. As this difference is quite negligible in comparison to the complexity introduced by one-hot encoding the 'region' column further testing on data frame 2 was halted. The next model trained was a linear regression model with standard scaling. This model was less accurate and produced an R2 score of 0.694. This is likely due to already low variance in the observations. The next model trained was a linear regression model with polynomial features. A degree 20 polynomial feature resulted in a negative R2 score, but a degree 2 polynomial feature resulted in an R2 score of 0.845. The next model trained was a Ridge and Lasso regularization model which after tuning both resulted in an R2 score of 0.845. The last model trained was an ElasticNet model that resulted in an R2 score of 0.769.

Key Findings:

The best model-- minimizing variance and bias --was the linear regression model with degree 2 polynomial featuring. It resulted in the highest R2 score while also relaying good interpretability. As can be seen in the graphs in Appendix A; the polynomial transformation helped to normalize the right skewed data.

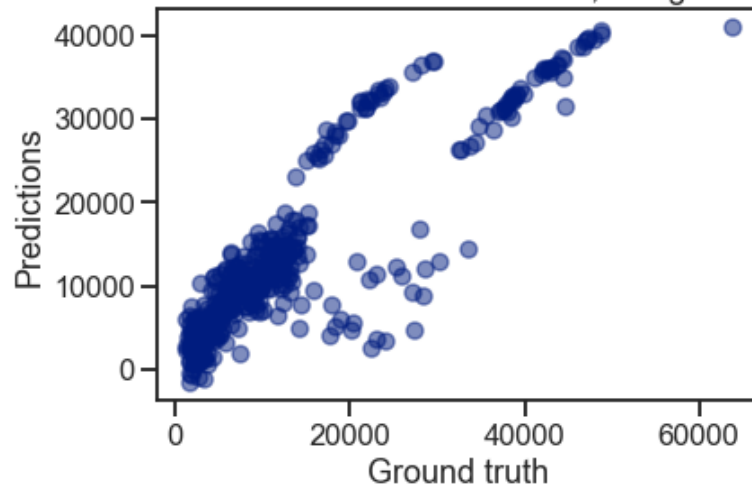
The degree 2 polynomial featuring accounted for a diminishing return trend that the linear model couldn't account for. There was a flattening out in medical expenditure even as the risk factors continued to increase. Each of the 6 factors contributed in similar magnitude to the model. Thus Lasso or Ridge was not able to zero any of the coefficients without underfitting the model. Thus it can be concluded that the base linear regression model was not complex enough and underfit the data and a degree 20 polynomial feature introduced too much complexity and overfit the data. A degree 2 polynomial feature introduced the optimum amount of complexity, positioning the model in the minimization point of bias and variance error.

Further Research:

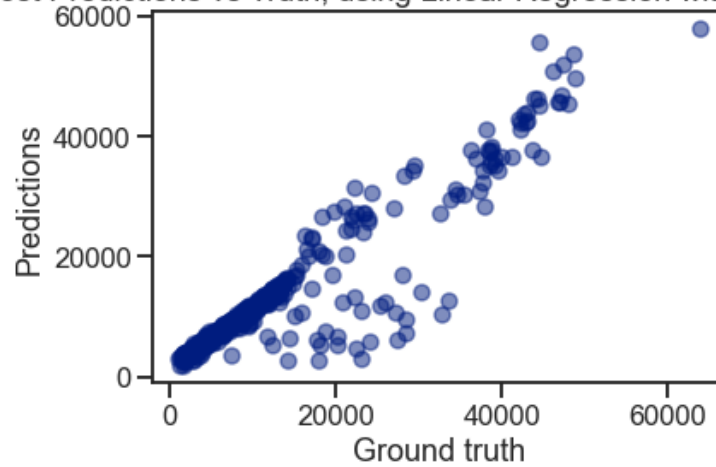
To further this research the next steps to take would be to further validate the findings. Possibly using k-fold cross validation on our degree 2 polynomial linear regression model. Trying different scaling methods such as log scaling could present new findings. Another step to take would be to reintroduce the one-hot encoded data frame into our degree 2 polynomial linear regression model.

Appendix A:

Medical Insurance Cost Predictions vs Truth, using Linear Regression



Medical Insurance Cost Predictions vs Truth, using Linear Regression with Degree 2 Polynomial Features



Appendix B:

