

Which Factors Influence the Price of Health Insurance?

You work as a Data Scientist for the New York Times. They are doing a story titled, "Which Factors Influence the Price of Health Insurance?". Many factors that affect how much you pay for health insurance are not within your control. Nonetheless, it's good to have an understanding of what they are. Hence, you have collected data about individuals for the new story. Your data contains basic factors about the

Overall, you want to find how well the factors you collected can predict the individual insurance costs billed by the health insurance company.

The first step of any machine learning project is to understand your data. The dataset contains the variables/features below. You should look at each feature to understand what type of feature it is, i.e., is it a categorical feature, binary, numeric, etc?

This is what you will predict:

charges: Individual medical costs billed by health insurance What type of item are you predicting? (e.g., categorical, numeric, binary, ordinal, etc.)

```
In [3]: 1 # NOTES
        2
        3 # What is Machine Learning?  $f(x) = y$ 
        4
        5 # Features and Scikit-Learn
        6
        7 ## Methods of construction features
        8
        9 ### DictVectorizer (Tabular Data)
       10
       11 ### CountVectorizer (Text Data)
       12
```

```

In [4]: 1 import csv
        2 from sklearn.feature_extraction import DictVectorizer
        3 from sklearn.model_selection import train_test_split
        4
        5 dataset = [] # Features go here
        6 y = [] # What you want to predict goes here
        7 with open('insurance.csv') as iFile:
        8     iCSV = csv.reader(iFile, delimiter=',')
        9     header = next(iCSV)
        10    for row in iCSV:
        11        item = {}
        12        item['age'] = float(row[0])
        13        item['sex'] = row[1]
        14        item['bmi'] = float(row[2])
        15        item['children'] = float(row[3])
        16        item['smoker'] = row[4]
        17        item['region'] = (row[5])
        18        dataset.append(item)
        19        y.append(float(row[6]))
        20
        21 # Here the lists are split into a 80% training portion and a 20% validation
        22 train_dataset, val_dataset, train_y, val_y = train_test_split(dataset, y,
        23
        24 # DictVectorizer will take a list of dictionaries and convert it into a matrix
        25 vec = DictVectorizer()
        26
        27 train_matrix = vec.fit_transform(train_dataset) # .fit_transform() should
        28 val_matrix = vec.transform(val_dataset) # .transform() should only be applied

```

```

In [5]: 1 # Here is what one dictionary looks like in the dataset
        2 dataset[0]

```

```

Out[5]: {'age': 19.0,
         'sex': 'female',
         'bmi': 27.9,
         'children': 0.0,
         'smoker': 'yes',
         'region': 'southwest'}

```

```

In [6]: 1 print(vec.feature_names_) # This list contains the feature names for each
        2
        3 ['age', 'bmi', 'children', 'region=northeast', 'region=northwest', 'region=southeast',
        4 'region=southwest', 'sex=female', 'sex=male', 'smoker=no', 'smoker=yes']

```

```

In [7]: 1 print(train_matrix.toarray()[0:2,:]) # This prints the first two rows of the training matrix
        2
        3 [[46.    19.95  2.    0.    1.    0.    0.    1.    0.    1.    0. ]
        4 [47.    24.32  0.    1.    0.    0.    0.    1.    0.    1.    0. ]]

```

```
In [9]: 1 import numpy as np
2 # Modify the weights below to achieve the LOWEST MSE as possible.
3 weights = np.array([0., # 'age'
4                     0., # 'bmi'
5                     0., # 'children'
6                     0., # 'region=northeast'
7                     0., # 'region=northwest'
8                     0., # 'region=southeast'
9                     0., # 'region=southwest'
10                    0., # 'sex=female'
11                    0., # 'sex=male',
12                    0., # 'smoker=no'
13                    0.]) # 'smoker=yes'
```

```
In [10]: 1 from sklearn.metrics import mean_squared_error
2 predictions = val_matrix.dot(weights)
3
4 print("MSE:", mean_squared_error(val_y, predictions))
```

MSE: 323425978.93488324

```
In [11]: 1 # Finding the weights manually is HARD.
2 # Here we will use scikit-learn to find the weights for us. We are trai
3 # Run the cells below to see how well it performs and what weights it f
4 from sklearn.linear_model import LinearRegression
5
6 clf = LinearRegression(fit_intercept=False)
7
8 clf.fit(train_matrix, train_y)
9
10 lr_predictions = clf.predict(val_matrix)
11 print("MSE:", mean_squared_error(val_y, lr_predictions))
```

MSE: 33596625.74155255

```
In [12]: 1 print("LR Weights:\n", "\n".join([f"{x:.5f}" for x in clf.coef_]))
```

```
LR Weights:
 257.01791
 337.34778
 425.47128
 342.99371
-27.60490
-315.85397
-467.22874
-224.54077
-243.15313
-12059.18785
11591.49395
```

*****Note that a lower MSE means the model is better (i.e., your prediction is closer to the real amount), How

You suspect that the Linear Regression model is overfitting (e.g., weights are too large). Moreover,

you think other models may provide better predictive ability. Choosing a model is only half the battle. In scikit-learn trying a different model is as simple as importing and using a different classifier object. But, to get the best performance from any given model, you need to tune all of the "knobs" available within the model. You can find all the "knobs" a model contains by looking at the scikit learn documentation under the section "Parameters"

https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.Ridge.html (https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.Ridge.html)

You are going to experiment with the "alpha" parameter. The alpha parameter will control overfitting in the Ridge Regression model. Try the values 0, .00001, .0001, .001, .01, .1, 1, 10, 100. Which value gives the best results. You should try this manually, by setting the alpha value to the different values and seeing the results. What do you find? Specifically, how does the alpha value impact the parameters? Also, what alpha gives the best performance?

```
In [13]: 1 # Try the different alphas in this cell
          2 from sklearn.linear_model import Ridge
          3
          4 clf = Ridge(alpha=0, fit_intercept=False)
          5
          6 clf.fit(train_matrix, train_y)
          7
          8 rr_predictions = clf.predict(val_matrix)
          9 print("MSE:", mean_squared_error(val_y, rr_predictions))
         10 print("\nRidge Weights:\n", "\n".join([f"{x:.5f}" for x in clf.coef_]))
```

MSE: 35516055.0927728

Ridge Weights:

```
262.40381
331.86322
325.48535
173.70912
-474.11913
469.62882
-707.97476
-912.67017
373.91421
-11921.84742
11383.09146
```

Trying every hyper parameter one-by-one is time consuming. Instead, let us do it using GridSearchCV. Note that GridSearchCV will use cross-validation applied to the training dataset. So, the exact results may be different that what was found in the previous cell. Why is that?

```
In [14]: 1 from sklearn.model_selection import GridSearchCV
2
3 params = {"alpha": [0, .00001, .0001, .001, .01, .1, 1, 10, 100, 1000]}
4
5 rr = Ridge(fit_intercept=False)
6
7 clf = GridSearchCV(rr, param_grid=params, cv=10, scoring='neg_mean_squa
8
9 clf.fit(train_matrix, train_y)
10
11 rr_predictions = clf.predict(val_matrix)
12 print("MSE:", mean_squared_error(val_y, rr_predictions))
13 print("Best Alpha:", clf.best_params_['alpha'])
14 print("Cross-Validation Score:", -clf.best_score_)
15 print("\nRidge Weights:\n", "\n".join([f"{x:.5f}" for x in clf.best_est
```

MSE: 35516054.80993609

Best Alpha: 1e-05

Cross-Validation Score: 38399017.57406919

Ridge Weights:

262.40380
331.86321
325.48535
173.70911
-474.11912
469.62880
-707.97474
-912.67014
373.91420
-11921.84709
11383.09115

Final Notes

```
1
2 So, where do you go from here? You can try different models such as:
3
4 Lasso: https://scikit-learn.org/stable/modules/generated/sklearn.linear\_model.Lasso.html
5 ElasticNet: https://scikit-learn.org/stable/modules/generated/sklearn.linear\_model.ElasticNet.htm
6 Random Forest Regression: https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestRegres
7 SVM Regression: https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVR.html
```

8 Each method has its own hyperparameters you must test to find what works best. Another option is to explore the use of "Feature Engineering". How will the model perform if you remove one or more features (columns)? What if you transform columns in some non-trivial way, e.g., square the values in a column (e.g., $\text{age} = \text{age}^2$) or combine values via interaction terms (e.g., $\text{age} * \text{gender} = \text{Male}$). Overall, the combinations are endless. If you had access to the data at a specific company, then you could also try to collect more specific data, e.g., income, family history, etc. In the end, this is a creative endeavor as much as it is a technical one.

In []:

1