

CSSM 502 – Course Project Report

1. Research Goal and Data

In this project, I will investigate which predictor has highest impact on medical cost of beneficiaries best for an insurance company. My hypothesis is “smoker” predictor has highest impact.

I decided to develop regression models for ‘insurance_data.csv’ data. I can explain my data as following:

- **Age:** age of primary beneficiary,
- **Sex:** insurance beneficiary gender (female, male),
- **BMI:** body mass index, providing an understanding of body weights that are relatively high or low relative to height, objective index of body weight (kg / m^2) using the ratio of height to weight, ideally 18.5 to 24.9,
- **Children:** number of children covered by health insurance / number of dependents,
- **Smoker:** smoking (yes/no),
- **Medical Cost:** individual medical costs billed by health insurance.

2. Data Preprocessing

My input variables are ‘age’, ‘sex’, ‘bmi’, ‘children’ and ‘smoker’. My target variable is ‘medical cost’. In ‘bmi’ column there are some missing values. I replaced them with most frequent item in ‘bmi’ column. Then I constructed features and target matrices. In features matrix, ‘sex’, ‘smoker’ columns have categorical information. Therefore, I used one-hot encoding on features matrix.

```

1 #read from csv file
2 df = pd.read_csv('insurance_data.csv', delimiter=';')
3 df['bmi'] = df['bmi'].replace(',', '.', regex=True).astype(float)
4 df['medical cost'] = df['medical cost'].replace(',', '.', regex=True).astype(float)
5
6
7 #construct features and target matrices
8 X = df.iloc[:,0:5]
9 Y = df['medical cost']
10
11 #for bmi column there are missing data I replaced missing values with most_frequent in
12 X_processed = SimpleImputer(strategy='most_frequent', missing_values=99)
13 X_processed = X_processed.fit(X[['bmi']])
14 X['bmi'] = X_processed.transform(X[['bmi']])
15
16 # sex and smoker data are categorical, therefore I will apply one-hot encoding
17 # creating instance of one-hot-encoder
18 X_encoded = OneHotEncoder(sparse=True, handle_unknown='ignore')
19 # passing sex and smoker columns
20 X_encoded = pd.DataFrame(X_encoded.fit_transform(X[['sex', 'smoker']]).toarray())
21 # merge encoded columns with other columns
22 X_encoded = X_encoded.join(X['age'])
23 X_encoded = X_encoded.join(X['bmi'])
24 X_encoded = X_encoded.join(X['children'])

```

3. Model Development

I split existing data as train (%80) data and test (default %20) data.

```

1 #split the data as train and test
2 Xtrain, Xtest, Ytrain, Ytest = train_test_split(X_encoded, Y, random_state=0,
3 train_size=0.8, test_size=0.2)

```

First, I used LinearRegression model as regressor. I fitted the model, predicted values for Xtest, calculated accuracy of model with MSE (Mean Squared Error) and MAE (Mean Absolute Error).

```

1 #LinearRegression
2 from sklearn.linear_model import LinearRegression
3 model = LinearRegression(fit_intercept=True) #storing of hyperparameter values
4 model.fit(Xtrain, Ytrain)
5 y_pred_LR = model.predict(Xtest)
6 #calculated MSE to measure model performance
7 MSE_LR = mean_squared_error(Ytest, y_pred_LR)
8 MAE_LR = mean_absolute_error(Ytest, y_pred_LR)
9 MSE_LR, MAE_LR

```

```
(32198271.422296938, 3936.2671465688218)
```

I decided to optimize LinearRegression model. For optimization, I used GridSearchCV module. By help of this module, I tried following hyperparameters:

- **fit_intercept**: [True, False],
- **normalize**: [True, False],
- **copy_X**: [True, False]

GridSearchCV module returned best parameters as following:

- **fit_intercept**: [True],
- **normalize**: [True],
- **copy_X**: [True]

```
1 #For LR optimization, I will use GridSearchCV
2 from sklearn.model_selection import GridSearchCV
3 import numpy as np
4
5 grid_params = {'fit_intercept': [True, False],
6               'normalize' : [True, False],
7               'copy_X' : [True, False]}
8
9 grid = GridSearchCV(LinearRegression(), grid_params, verbose = 1, cv = 7, n_jobs = -1)
10
11 grid.fit(Xtrain,Ytrain)
12 grid.best_params_
```

Fitting 7 folds for each of 8 candidates, totalling 56 fits

```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 8 concurrent workers.
[Parallel(n_jobs=-1)]: Done 56 out of 56 | elapsed: 1.7s finished
```

```
{'copy_X': True, 'fit_intercept': True, 'normalize': True}
```

I run LinearRegression with best hyperparameter values. I fitted the model, predicted values for Xtest, calculated accuracy of model with MSE (Mean Squared Error) and MAE (Mean Absolute Error).

```
1 #run LinearRegression with optimized parameters
2 model_optimized = LinearRegression(fit_intercept=True, normalize=True, copy_X=True) #st
3 model_optimized.fit(Xtrain, Ytrain)
4 y_pred_LR_optimized = model_optimized.predict(Xtest)
5 #calculated MSE to measure model performance
6 MSE_LR_optimized = mean_squared_error(Ytest, y_pred_LR_optimized)
7 MAE_LR_optimized = mean_absolute_error(Ytest, y_pred_LR_optimized)
8 MSE_LR_optimized,MAE_LR_optimized
```

```
(32198271.422296952, 3936.2671465688272)
```

Then I used SVR as regressor. I fitted the model, predicted values for Xtest, calculated accuracy of model with MSE (Mean Squared Error) and MAE (Mean Absolute Error).

```
1 #Support Vector Regressor(SVR)
2 regressor = SVR(kernel = 'poly')
3 regressor.fit(Xtrain, Ytrain)
4 y_pred = regressor.predict(Xtest)
5 #calculated MSE to measure model performance
6 MSE = mean_squared_error(Ytest, y_pred)
7 MAE = mean_absolute_error(Ytest, y_pred)
8 MSE,MAE
```

```
(171644228.42779985, 8150.237918030387)
```

I decided to optimize SVR model. For optimization, I used GridSearchCV module. By help of this module, I tried following hyperparameters:

- **kernel:** ['linear', 'poly', 'rbf', 'sigmoid']
- **gamma:** ['scale', 'auto']

GridSearchCV module returned best parameters as following:

- **kernel:** ['poly']
- **gamma:** ['auto']

```
1 #For SVR optimization, I will use GridSearchCV
2 from sklearn.model_selection import GridSearchCV
3 import numpy as np
4
5 grid_params = {'kernel': ['linear', 'poly', 'rbf', 'sigmoid'],
6               'gamma' : ['scale', 'auto']}
7
8 grid = GridSearchCV(SVR(), grid_params, verbose = 1, cv = 7, n_jobs = -1)
9
10 grid.fit(Xtrain,Ytrain)
11 grid.best_params_
```

Fitting 7 folds for each of 8 candidates, totalling 56 fits

```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 8 concurrent workers.
[Parallel(n_jobs=-1)]: Done 34 tasks | elapsed: 2.0s
[Parallel(n_jobs=-1)]: Done 56 out of 56 | elapsed: 1.7min finished
```

```
{'gamma': 'auto', 'kernel': 'poly'}
```

I run SVR with best hyperparameter values. I fitted the model, predicted values for Xtest, calculated accuracy of model with MSE (Mean Squared Error) and MAE (Mean Absolute Error).

```
1 #run SVR with best parameters
2 regressor_optimized = SVR(kernel = 'poly', gamma = 'auto')
3 regressor_optimized.fit(Xtrain, Ytrain)
4 y_pred_optimized = regressor_optimized.predict(Xtest)
5 MSE_optimized = mean_squared_error(Ytest, y_pred_optimized)
6 MAE_optimized = mean_absolute_error(Ytest, y_pred_optimized)
7 MSE_optimized,MAE_optimized
```

```
(21778183.016043536, 2141.3520245314526)
```

Finally, I used Random Forest Regressor(RFR) as regressor. I fitted the model, predicted values for Xtest, calculated accuracy of model with MSE (Mean Squared Error) and MAE (Mean Absolute Error).

```

1 from sklearn.ensemble import RandomForestRegressor
2 #Random Forest Regressor(RFR)
3 regressor_RFR = RandomForestRegressor(n_estimators = 100, random_state = 0)
4 regressor_RFR.fit(Xtrain, Ytrain)
5 y_pred_RFR = regressor_RFR.predict(Xtest)
6 #calculated MSE to measure model performance
7 MSE_RFR = mean_squared_error(Ytest, y_pred_RFR)
8 MAE_RFR = mean_absolute_error(Ytest, y_pred_RFR)
9 MSE_RFR, MAE_RFR

```

(20354599.438501805, 2588.707771231343)

I decided to optimize RFR model. For optimization, I used GridSearchCV module. By help of this module, I tried following hyperparameters:

- **criterion:** ['mse', 'mae']
- **max_features:** ['auto', 'sqrt', 'log2']
- **bootstrap:** [True, False]

GridSearchCV module returned best parameters as following:

- **criterion:** ['mae']
- **max_features:** ['sqrt']
- **bootstrap:** [True]

```

1 #For RFR optimization, I will use GridSearchCV
2 grid_params = {'criterion' : ['mse', 'mae'],
3               'max_features' : ['auto', 'sqrt', 'log2'],
4               'bootstrap' : [True, False]}
5
6 max_features : {"auto", "sqrt", "log2"}
7
8 grid = GridSearchCV(RandomForestRegressor(), grid_params, verbose = 1, cv = 7, n_jobs
9
10 grid.fit(Xtrain, Ytrain)
11 grid.best_params_

```

Fitting 7 folds for each of 12 candidates, totalling 84 fits

```

[Parallel(n_jobs=-1)]: Using backend LokyBackend with 8 concurrent workers.
[Parallel(n_jobs=-1)]: Done 34 tasks | elapsed: 3.7s
[Parallel(n_jobs=-1)]: Done 84 out of 84 | elapsed: 11.4s finished

```

```
{'bootstrap': True, 'criterion': 'mae', 'max_features': 'sqrt'}
```

I run RFR with best hyperparameter values. I fitted the model, predicted values for Xtest, calculated accuracy of model with MSE (Mean Squared Error) and MAE (Mean Absolute Error).

```

1 #run RFR with best parameters
2 regressor_RFR_optimized = RandomForestRegressor(n_estimators = 100, random_state = 0,
3 regressor_RFR_optimized.fit(Xtrain, Ytrain)
4 y_pred_RFR_optimized = regressor_RFR_optimized.predict(Xtest)
5 #calculated MSE to measure model performance
6 MSE_RFR_optimized = mean_squared_error(Ytest, y_pred_RFR_optimized)
7 MAE_RFR_optimized = mean_absolute_error(Ytest, y_pred_RFR_optimized)
8 MSE_RFR_optimized,MAE_RFR_optimized

```

(16988783.590700284, 2285.703164365672)

4. Model Selection

I summarized performances of regressors in following table. My metrics are MSE (Mean Squared Error) and MAE (Mean Absolute Error). According to these metrics, SVR and RFR are better than Linear Regression. I will choose RFR as a regressor for this data.

Regressor	MSE	MAE
Linear Regression	32.198.271,42	3.936,26
Support Vector Regressor	21.778.183,01	2.141,35
Random Forest Regressor	16.988.783,59	2.285,70

5. Results

Finally, I will look at relation between each single predictor and medical cost parameter. I used Random Forest Regressor model and calculated MSE and MAE values in each scenario. Smoking and medical cost pair has smallest MSE and MAE values. Therefore my hypothesis is true, smoking has highest impact on medical cost.

```

1 #regression between sex and medical cost
2 regressor_RFR_optimized.fit(np.reshape(np.array(Xtrain[0]), (-1, 1)), Ytrain)
3 y_pred1 = regressor_RFR_optimized.predict(np.reshape(np.array(Xtest[0]), (-1, 1)))
4 #calculated MSE to measure model performance
5 MSE_RFR_optimized = mean_squared_error(Ytest, y_pred1)
6 MAE_RFR_optimized = mean_absolute_error(Ytest, y_pred1)
7 MSE_RFR_optimized,MAE_RFR_optimized

```

(17555160.54948175, 8619.733930597014)

```

1 #regression between smoking and medical cost
2 regressor_RFR_optimized.fit(np.reshape(np.array(Xtrain[2]), (-1, 1)), Ytrain)
3 y_pred1 = regressor_RFR_optimized.predict(np.reshape(np.array(Xtest[2]), (-1, 1)))
4 #calculated MSE to measure model performance
5 MSE_RFR_optimized = mean_squared_error(Ytest, y_pred1)
6 MAE_RFR_optimized = mean_absolute_error(Ytest, y_pred1)
7 MSE_RFR_optimized,MAE_RFR_optimized

```

(51329974.80796147, 5342.986025932836)

```

1 #regression between bmi and medical cost
2 regressor_RFR_optimized.fit(np.reshape(np.array(Xtrain['bmi']), (-1, 1)), Ytrain)
3 y_pred1 = regressor_RFR_optimized.predict(np.reshape(np.array(Xtest['bmi']), (-1, 1)))
4 #calculated MSE to measure model performance
5 MSE_RFR_optimized = mean_squared_error(Ytest, y_pred1)
6 MAE_RFR_optimized = mean_absolute_error(Ytest, y_pred1)
7 MSE_RFR_optimized,MAE_RFR_optimized

```

(199270039.43339625, 10178.637050186568)

```

1 #regression between age and medical cost
2 regressor_RFR_optimized.fit(np.reshape(np.array(Xtrain['age']), (-1, 1)), Ytrain)
3 y_pred1 = regressor_RFR_optimized.predict(np.reshape(np.array(Xtest['age']), (-1, 1)))
4 #calculated MSE to measure model performance
5 MSE_RFR_optimized = mean_squared_error(Ytest, y_pred1)
6 MAE_RFR_optimized = mean_absolute_error(Ytest, y_pred1)
7 MSE_RFR_optimized,MAE_RFR_optimized

```

(165838678.944245, 7091.66516902985)

```

1 #regression between children and medical cost
2 regressor_RFR_optimized.fit(np.reshape(np.array(Xtrain['children']), (-1, 1)), Ytrain)
3 y_pred1 = regressor_RFR_optimized.predict(np.reshape(np.array(Xtest['children']), (-1, 1)))
4 #calculated MSE to measure model performance
5 MSE_RFR_optimized = mean_squared_error(Ytest, y_pred1)
6 MAE_RFR_optimized = mean_absolute_error(Ytest, y_pred1)
7 MSE_RFR_optimized,MAE_RFR_optimized

```

(174779559.1861185, 8660.452840671642)