

# Machine Learning 2: Prediction of the Individual Medical Cost

Conor Fallon; Dennis Fast; Leonhard Liu; Tassilo Henninger

1 17 2023

## Table of contents

<b>1</b>	<b>Baseline - Linear Regression</b>	<b>8</b>
1.1	Mathematical Overview . . . . .	8
1.2	Hyperparameter Optimization . . . . .	9
1.3	Performance . . . . .	10
<b>2</b>	<b>Algorithm 1 - Regression Tree</b>	<b>11</b>
2.1	Mathematical Overview . . . . .	11
2.2	Hyperparameter Optimization . . . . .	12
2.3	Performance . . . . .	14
<b>3</b>	<b>Algorithm 2 - Neural Networks</b>	<b>16</b>
3.1	Mathematical Overview . . . . .	16
3.2	Initial Model . . . . .	18
3.3	Hyperparameter Optimization . . . . .	20
3.4	Performance . . . . .	21

## Project Goal and Outline

The task at hand is a regression problem. Broadly speaking, we wish to create a series of models, each of which will be able to make predictions on a withheld test set. The dataset in question is a Medical Cost dataset. The variable we wish to predict is the expected insurance premium for a given individual based on the following input variables:

- age
- sex: male or female
- Body mass index, is a measure of one's weight relative to their height
- children: Number of children/dependants also covered under this person's health insurance
- smoker: Whether the person smokes or not
- region: Four locations: northeast, southeast, southwest, and northwest. All refer to the United States
- charges: The premium billed by the health insurance company. This is the predictor variable

Of course, some of these variables seem obvious in how they affect a person's insurance premium; a 98 year-old obese smoker will have a higher insurance premium than a 24 year-old non-smoker with a healthy BMI. However, through some exploratory data analysis as well as by comparing our methods, a more discerning picture of how each variable impacts the value we wish to predict will be discussed.

There will be three models compared to each other: a baseline model, which consists of a simple linear regression; a regression tree; and a neural network. The best model is ultimately decided upon by using the Root Mean Squared error, although other metrics, where suitable, will be used to influence our decision making. For both the Neural Network and the Regression Tree models, we wish to find a 'best' Neural Network and a 'best' Regression Tree with which to work. For this purpose, these models are initially trained only on the train split (as outlined further on), and the 'best' model for each is decided upon based on their performance on a validation split. As such, the selection of the best model in our two algorithm classes is based on the validation error.

It is important to emphasise that this leaves the test data untouched at the stage when these models are created, leaving the test data purely for the final model comparison between the best Neural Net and the best Regression Tree. Then when ultimately comparing the two chosen models, these are retrained on the combined train and validation splits, and finally tested on the test set, which has been withheld for this purpose.

This report will conclude with a discussion of our findings, an examination of what worked and what did not work, and will also mention avenues for future work.

## Dataset and Preprocessing

The dataset we used is included in a Book called “Machine Learning with R” by Brett Lantz and is also available on Kaggle. It can be downloaded via the following links:

- <https://github.com/stedy/Machine-Learning-with-R-datasets>
- <https://www.kaggle.com/datasets/mirichoi0218/insurance>

It contains 1338 rows, 7 columns and there are no missing values. The 6 predictor variables are like in the project outline introduced: age(numeric, integer), sex(female/male), BMI(numeric, float), children(numeric, integer), smoker(boolean) and region(categorical). The outcome variable is “charges” (numeric, float).

The only preprocessing we had to do was encoding the categorical features as factor variables. Next we split our dataset into a 60% train set, a 20% validation and a 20% test set, as it is asked in the project description.

## EDA

One of the objectives of exploratory data analysis is to get a feel for the data you are dealing with by describing the key features of the data and summarizing the results.

The first thing we can do is to take a look at the descriptive statistics via the summary function and the histograms.

- we have individuals ranging from the age of 18 to 64 with a mean and median around 39 ages.
- the sex is almost perfectly balanced as 49.5% of the individuals are female and 50.5% are male.
- the BMI is ranging from 15.96 to 53.13 with a mean and median around 30. The BMI is calculated by the following formula:

$$BMI = \frac{mass_{kg}}{height_m^2}$$

A common use of the BMI is to assess how far an individual’s body weight departs from what is normal for a person’s height. The WHO regards an adult BMI of less than 18.5 as underweight, a BMI of 25 or more as overweight and 30 or more is considered obese. From the summary and the histogram, we can clearly see, that our population of individuals is in mean considered obese and therefore an unhealthy population.

Most of the individuals have no children, the mean is 1 child and the maximum is 5. Regarding the factor feature smoker, the dataset is unbalanced. Only 20.5% of our individuals are smokers, respectively 79.5% are non smokers. The region is well balanced again between the four categories “southwest”, “southeast”, “northwest” and “northeast”.

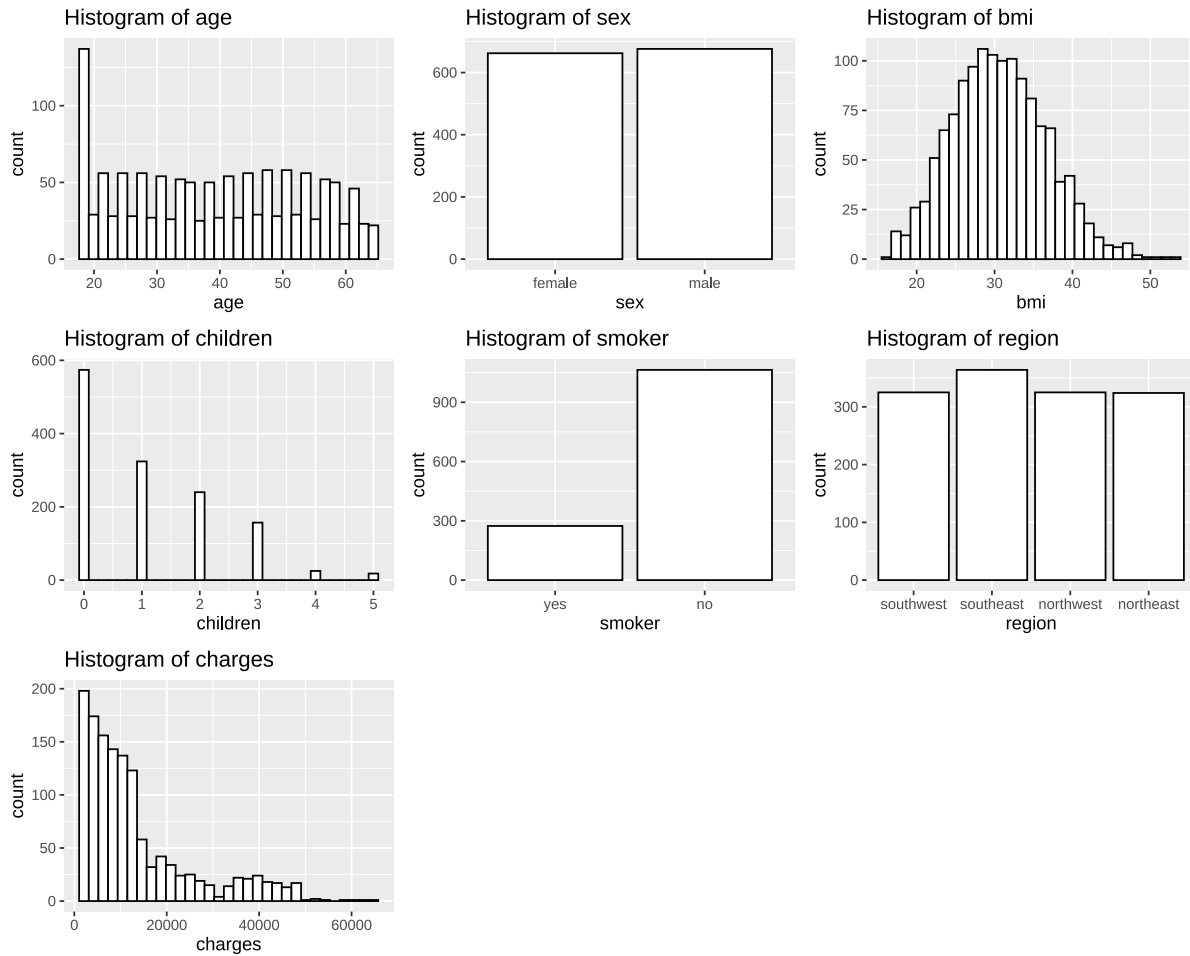
Last we have our predictor variable charges. The minimum of a health insurance charge is 1122\$, the median is 9382\$, the mean 13270\$ and the maximum 63770\$. From the difference of median and mean we can already tell that the distribution is skewed to the right (positively skewed). This means that there are a few large values that are pulling the mean higher, but the majority of the data is concentrated around the median. Looking at the boxplot for the charge we can confirm that. We can inspect a long tail of high values and a shorter tail of low values. We have no symmetric distribution.

```
summary(df.insurance)
```

age	sex	bmi	children	smoker
Min. :18.00	female:662	Min. :15.96	Min. :0.000	yes: 274
1st Qu.:27.00	male :676	1st Qu.:26.30	1st Qu.:0.000	no :1064
Median :39.00		Median :30.40	Median :1.000	
Mean :39.21		Mean :30.66	Mean :1.095	
3rd Qu.:51.00		3rd Qu.:34.69	3rd Qu.:2.000	
Max. :64.00		Max. :53.13	Max. :5.000	

region	charges
southwest:325	Min. : 1122
southeast:364	1st Qu.: 4740
northwest:325	Median : 9382
northeast:324	Mean :13270
	3rd Qu.:16640
	Max. :63770



Next, we want to better understand the relationships between our different variables in the dataset. Therefore we use the correlation matrix to get insights into our data regarding the linearity and their respective strengths. A positive correlation means that as one variable increases, the other variable also increases, and a negative correlation means that as one variable increases, the other variable decreases. A correlation value of 1 or -1 indicates a perfect linear relationship between the variables. A correlation value of 0 indicates that there is no linear relationship between the variables.

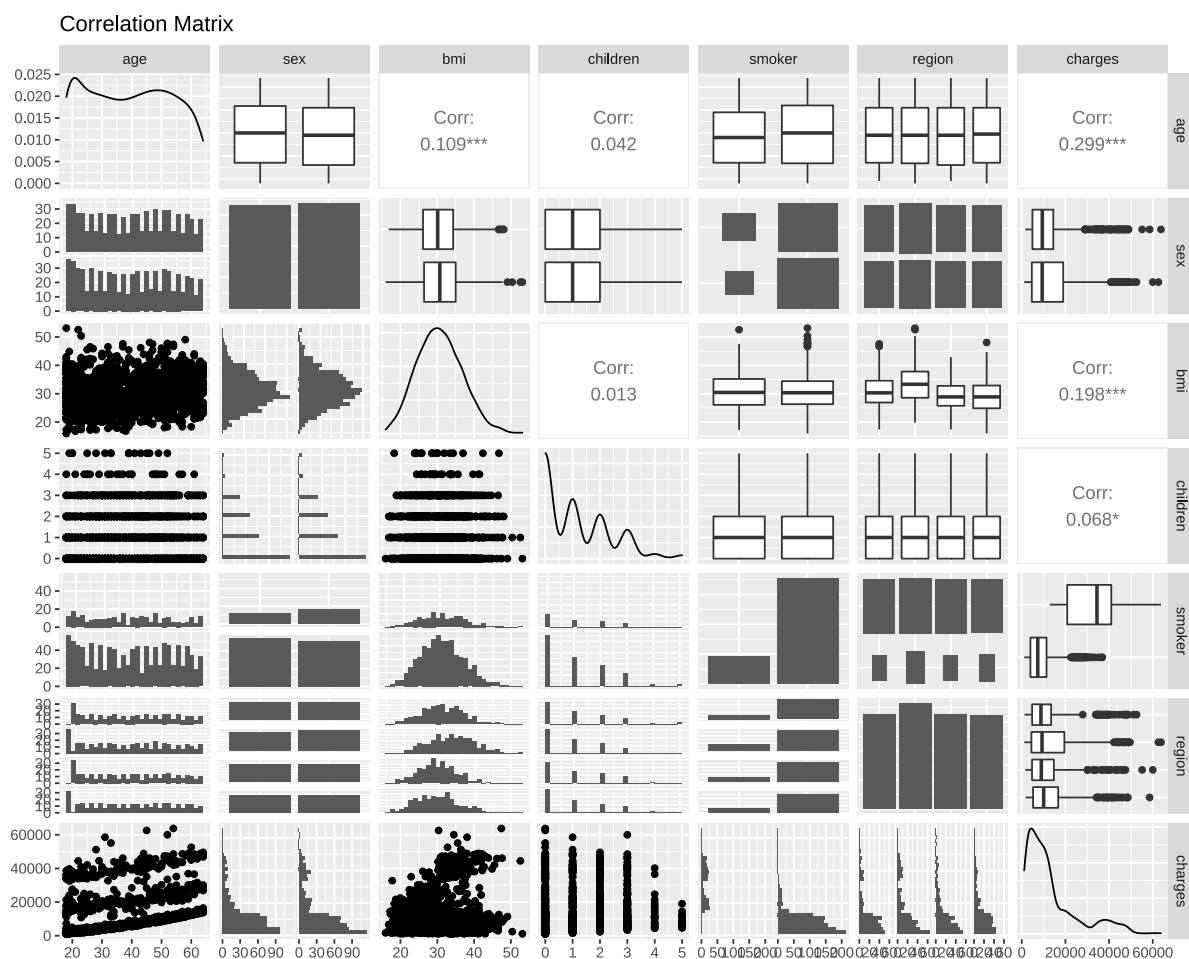
Note however, that correlation does not imply causality, it just indicates that two variables are related. That means the correlation between two variables may be high, but that the correlation is due to a relationship between the two variables and a third variable. This concept is called multicollinearity.

This has to be kept in mind for our baseline model. If we create a linear model and want to interpret the coefficients on how they influence our charge, multicollinearity can be a problem, as interpretability assumes that you can only change the value of one explanatory variable and not the others at the same time. This of course is only true if there are no correlations

between the explanatory variables. If this independence does not hold, we have a problem of multicollinearity.

Multicollinearity can result in the coefficients swinging wildly based on which other independent variables are in the model. Therefore the coefficients become very sensitive to small changes in the model and can not be easily interpreted. One way to detect multicollinearity is to examine the correlation matrix and look for high correlation coefficients between the multiple explanatory variables.

As the correlation metric can only be calculated for numeric values, we get the correlation scores for the features age, bmi, children and charges. Charges has the strongest correlation with Age (+0.299), then bmi (+0.198) and last children with only +0.068. The highest correlation between explanatory variables is +0.109 between bmi and age. As this is quite low, we don't have a problem of multicollinearity between the numeric explanatory variables.



One way to assess the influence of categorical variables on an outcome variable is through a one-way ANOVA, which compares the means of a numerical outcome variable across different

levels of a categorical variable. Thereby we could determine if there is a significant difference in the means of the outcome variable between the different levels of the categorical variable.

However, the ANOVA assumes that the data is normally distributed, independent and that the variances of the groups are equal. We already found that outcome “charges” is positively skewed and not normally distributed. Thereby, the assumption of the ANOVA fails. We can however use the non-parametric alternatives called Kruskal-Wallis test.

The output of Kruskal-Wallis test statistics gives us the p-value, through which we can determine the significance of the test. A small p-value (typically less than 0.05) indicates strong evidence against the null hypothesis, which in this case is that there is no difference in the medians of the outcome variable between the different levels of the categorical variable.

As only the p-value of the categorical variable “smoker” is below the significance level of 5%, we can only reject the null hypothesis for smoker and therefore say, there is a significant difference in the medians of the outcome variable between being a smoker and not. We can confirm the result also by looking at the boxplot, split up by category. This is also visible in small boxplots in the correlation matrix.

Kruskal-Wallis rank sum test

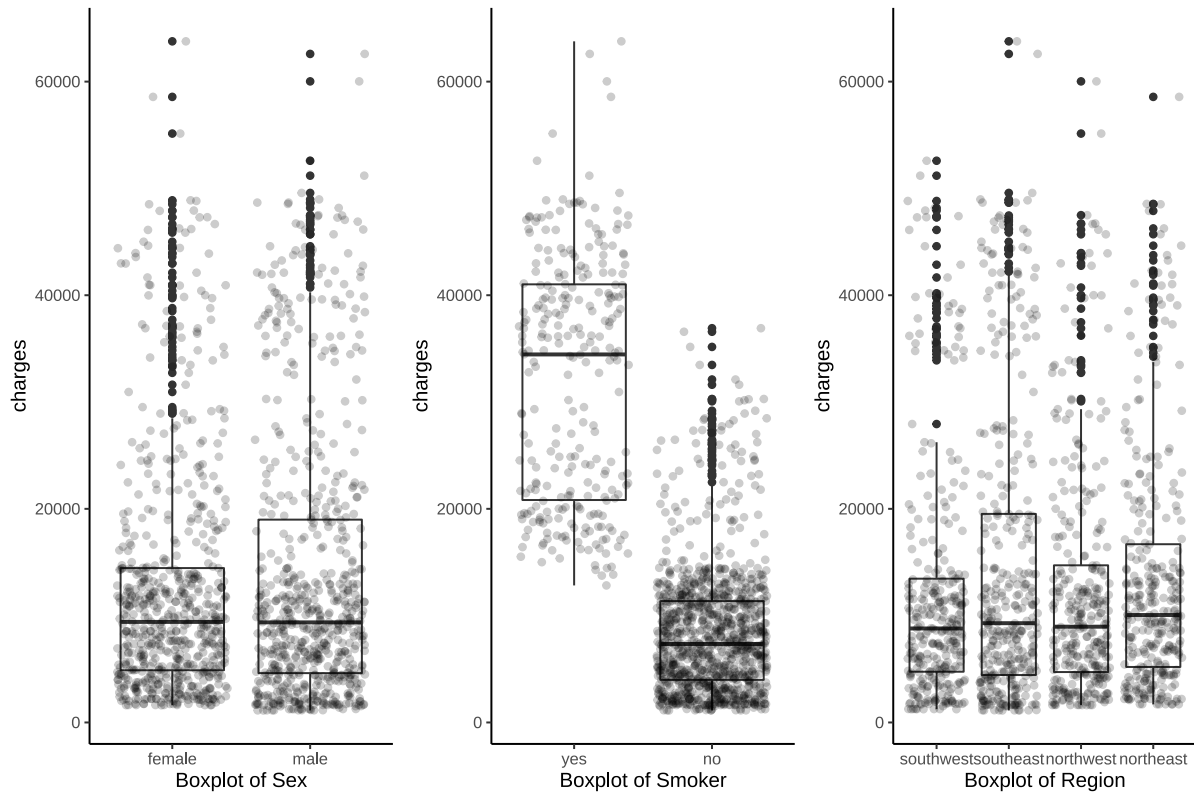
```
data: df.insurance$charge and df.insurance$sex
Kruskal-Wallis chi-squared = 0.1204, df = 1, p-value = 0.7286
```

Kruskal-Wallis rank sum test

```
data: df.insurance$charge and df.insurance$smoker
Kruskal-Wallis chi-squared = 588.52, df = 1, p-value < 2.2e-16
```

Kruskal-Wallis rank sum test

```
data: df.insurance$charge and df.insurance$region
Kruskal-Wallis chi-squared = 4.7342, df = 3, p-value = 0.1923
```



Forecasting

## 1 Baseline - Linear Regression

In order to check the performance of the chosen models, we first created a simple linear regression model as a baseline model, with 'charges' as the outcome variable and all other columns ('age', 'sex', 'BMI', 'children', 'smoker', and 'region') as predictor variables.

### 1.1 Mathematical Overview

The mathematical notation of a linear model represents a linear relationship between the outcome variable 'charges' and the predictor variables 'age', 'sex', 'BMI', 'children', 'smoker', and 'region'. The equation states that the outcome variable 'charges' is a linear combination of the predictor variables, with each predictor variable multiplied by a specific coefficient (beta). By default the categorical features are dummy encoded, to convert the categorical variable into multiple binary features, and then use these binary features as independent variables in the linear regression model. This allows the model to learn separate regression coefficients for each category of the original categorical variable.



$$\text{charges} = \beta_0 + \beta_1 * \text{age} + \beta_2 * \text{sexmale} + \beta_3 * \text{BMI} + \beta_4 * \text{children} + \beta_5 * \text{smokeryes} + \beta_6 * \text{regionnorthwest} + \beta_7 * \text{regionsoutheast} + \beta_8 * \text{regionsouthwest}$$

In this equation, ‘charges’ is represented as charges, the coefficients are represented as  $\beta_0$ ,  $\beta_1$ ,  $\beta_2$ ,  $\beta_3$ ,  $\beta_4$ ,  $\beta_5$ ,  $\beta_6$ ,  $\beta_7$ ,  $\beta_8$  and predictor variables ‘age’, ‘sex’, ‘BMI’, ‘children’, ‘smoker’, and ‘region’ are represented as age, sexmale, BMI, children, smokeryes, regionnorthwest, regionsoutheast and regionsouthwest respectively. Note, that sexfemale is indirectly encoded in sexmale being 0, smokerno is indirectly encoded in smokeryes being 0 and regionnortheast in all other region variables being 0.

The equation can be used to make predictions about the outcome variable ‘charges’ based on the values of the predictor variables. For example, if we know the values of ‘age’, ‘sex’, ‘BMI’, ‘children’, ‘smoker’, and ‘region’ for a specific individual, we can plug those values into the equation and solve for ‘charges’. The result of this calculation would be an estimate of the individual’s ‘charges’.

If our linear model has good predictability, we can also interpret the coefficients on how they influence the outcome. As we found no strong multicollinearity during EDA, we can do that in the following. This is called regression analysis.

## 1.2 Hyperparameter Optimization

Since we use the linear model as a baseline, we don’t conduct the hyperparameter optimization on that model.

```
lm1 <- lm(charges ~ ., train)
summary(lm1)
```

Call:

```
lm(formula = charges ~ ., data = train)
```

Residuals:

Min	1Q	Median	3Q	Max
-11078	-3021	-1085	1448	29854

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-11854.00	1282.48	-9.243	< 2e-16 ***
age	259.98	15.47	16.805	< 2e-16 ***
sexmale	42.63	429.42	0.099	0.92095
bmi	327.40	36.92	8.867	< 2e-16 ***
children	477.80	174.61	2.736	0.00635 **
smokeryes	23594.95	536.18	44.005	< 2e-16 ***

```

regionnorthwest    242.35    597.11    0.406    0.68494
regionsoutheast    -683.86    615.71   -1.111    0.26705
regionsouthwest   -1391.54    612.99   -2.270    0.02347 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 6042 on 793 degrees of freedom
Multiple R-squared:  0.7468,    Adjusted R-squared:  0.7442
F-statistic: 292.3 on 8 and 793 DF,  p-value: < 2.2e-16

```

### 1.3 Performance

Finally, we predict the outcome on the validation set and calculate RMSE and R2 scores.

Model	$RMSE_{valid}$	$R^2_{valid}$
Linear Model (Baseline)	5914.225	0.7953

The R2 score of 0.79 is good enough for interpreting the coefficients. From the model summary we can see that the following explanatory variables are statistically significant:

- Age has a positive effect on the charges. One unit change on age results in an increase of the charges of **259.98\$**
- BMI has a positive effect on the charges. One unit change on BMI results in an increase of the charges of **327.40\$**
- Children has a positive effect on the charges. One unit change on Children results in an increase of the charges of **477.80\$**
- Smoker has a positive effect on the charges. Being a smoker results in an increase of the charges of **23594.95\$**

## 2 Algorithm 1 - Regression Tree

### 2.1 Mathematical Overview

A regression tree is a specific class of tree which will predict a numerical dependent variable based on a number of explanatory variables. They can deal with the situation in which the explanatory variables interact with one another (BMI and age are in all likelihood in some way dependent).

This can be difficult to model through simpler models like linear regression; however, by partitioning the data into smaller regions, the interactions between the explanatory variables, potentially very complex, can become more manageable. In a tree model, this partitioning and sub-partitioning (i.e. recursively partitioning each partition) is done by the tree.

We want a set of predictor variables  $X_1, \dots, X_P$ , into some number  $J$  distinct and mutually exclusive regions (partitions)  $R_1, \dots, R_J$ . In this description, each observation that falls into the region  $R_j$  will have the same value assigned to it. These regions are decided upon as follows: we want to find a selection of regions  $R_1, \dots, R_J$  such that the residual squared error (RSS) is minimised. I.e. minimise:

$$\sum_{j=1}^J \sum_{i \in R_j} (y_i - \hat{y}_{R_j})^2$$

where  $\hat{y}_{R_j}$  is the mean value for the predictor variable in that given region based on the training data.

The splitting is done in a top-down, greedy manner, where the splitting starts at the top of the tree (i.e. all data-points belong to the same region), then this is split in 2, and this process is repeated for each split in the tree. It is greedy in the sense that it does not look ahead steps into the future; it simply predicts the best split for the split in question. The actual splitting is done as follows.

A predictor variable  $X_j$  and a splitting point  $s$  such that the splitting regions  $X|X_j < s$  and  $X|X_j \geq s$  provides the smallest possible RSS. This is then repeated for the other predictor variables until the predictor variable which gave the value of  $s$  which gave the smallest RSS is found and this is chosen as the variable for this split in the tree. We do this until a stopping criteria is found, one which will ideally neither underfit nor overfit the tree.

In our case, we are examining the complexity parameter in order to do this. Of course, if we completely overfit the tree, there will be the smallest RSS; however, we want to penalise overfitting in some way. The approach used here is penalised least squares (PLS).

Let us call our full tree  $T_0$ . We want to find a pruned tree  $T_\alpha \subset T_0$  such that

$$PLS(\alpha) = \sum_{j \in |T_\alpha|}^J \sum_{i \in R_j} (y_i - \hat{y}_{R_j})^2 + \alpha |T_\alpha|$$

with  $\alpha$  greater than or equal to zero, and the other variables as defined already.  $\alpha$  is the complexity parameter. It can be found using  $K$ -fold cross validation.

The rule of thumb that is used by the rpart package is to take the largest value of the complexity parameter that is within 1 standard deviation of the smallest value out of the cross validation setup.

## 2.2 Hyperparameter Optimization

This is how the rpart regression tree works out of the box: the standard hyperparameter settings form a baseline of the regression trees. Rpart uses 10-fold cross validation as its default to find the complexity parameter. The tree looks as follows. To explain how to interpret this output, the first split is made on whether or not one is a smoker; if one is not a smoker, the next split is based on age; otherwise, the next split is made on BMI.

n= 802

```
node), split, n, deviance, yval
      * denotes terminal node
```

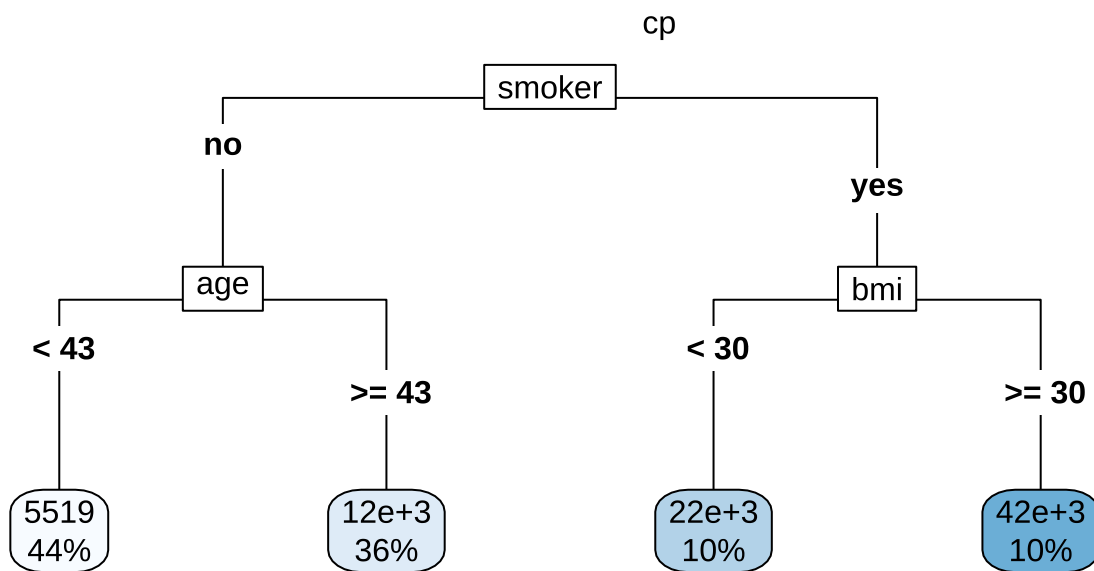
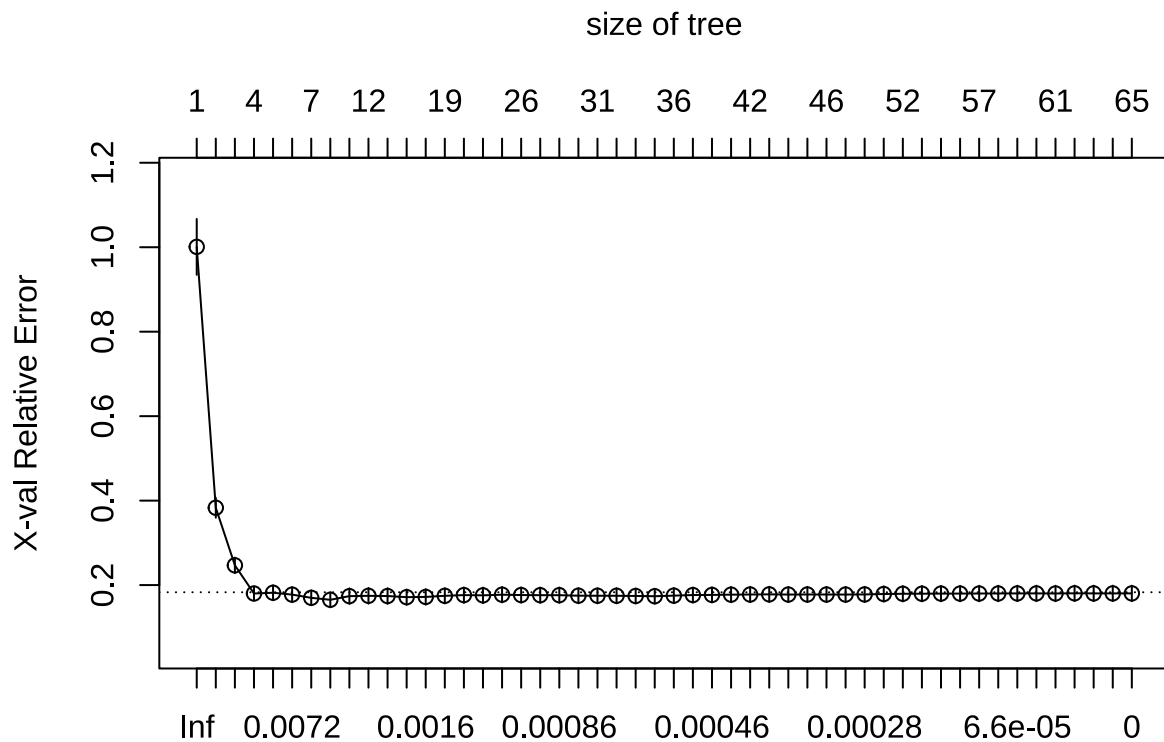
```
1) root 802 114336300000 13335.670
  2) smoker=no 641 23372910000 8626.180
    4) age< 42.5 355 8058769000 5518.603 *
    5) age>=42.5 286 7630553000 12483.490 *
  3) smoker=yes 161 20143560000 32085.890
    6) bmi< 30.01 77 1671806000 21685.210 *
    7) bmi>=30.01 84 2507060000 41619.850 *
```

The next tree to be formed will be a full tree (i.e. totally overfitted) as discussed above, with complexity parameter equaling 0, which will then be pruned back.

We can see in the below plot how the error decreases as the complexity parameter decreases.

Based on this plot a suitable complexity parameter is selected, and this quite interpretable decision tree is shown.

Note, there exist other error measures for the splitting criteria and this will be mentioned in the section on future work. Note that the  $cp=0.019$  was gotten from reading the `cpmatrix` and applying the rule of picking a `cp` within 1sd of the smallest error; `cpmatrix` has been commented out to avoid it being printed in the report. The splitting rule, now shown in an easy-to-read plot, is similar to the one outputted by rpart's default settings.



Let us predict and get some metrics as was done for the baseline model. Do this for the validation set.

This is initially done on the full, overfitted tree, then on the default one, and then on the one with the selected complexity parameter

Model	$RMSE_{test}$	$R2_{test}$
Full, Overfitted Tree	4770.110	0.8644
Default Tree	5076.103	0.8479
Pruned Tree	5076.103	0.8479

Note that the default method and the one found via our rule of thumb find the same value and thus the same RMSE and R2 values.

We will do some quick hyperparameter optimisation to finding the optimal complexity parameter via grid search. We get the following:

- $cp = 0.06$
- $RMSE_{test} = 5076.103$
- $R2_{test} = 0.8479$

It clearly finds a different  $cp$  than before; however, it provides the same RMSE and R2 as the other two trees discussed. This is perhaps unsurprising as they all are creating the same tree (compare first and last plots). Hence this tree seems to offer us the best solution.

## 2.3 Performance

For a single tree you can interpret a fitted model by inspecting the text output of the tree or the tree diagram itself. We could also calculate the Increase in Node Purity and Mean decrease in Accuracy, but that is most often done for ensemble methods, as you e.g can not plot all trees of a random forest. Given that Smoking, BMI and Age are the only variables considered by our ‘best’ model, we can consider that these variables are the most important.

Finally, we predict the outcome of the final tree, found via hyperparameter optimisation, on the validation set and calculate RMSE and R2 scores.

Model	$RMSE_{valid}$	$R2_{valid}$
Regression Tree	5076.103	0.8479

The final tree itself looks as follows:

```
printcp(rpart(charges~.,train))
```

Regression tree:

```
rpart(formula = charges ~ ., data = train)
```

Variables actually used in tree construction:

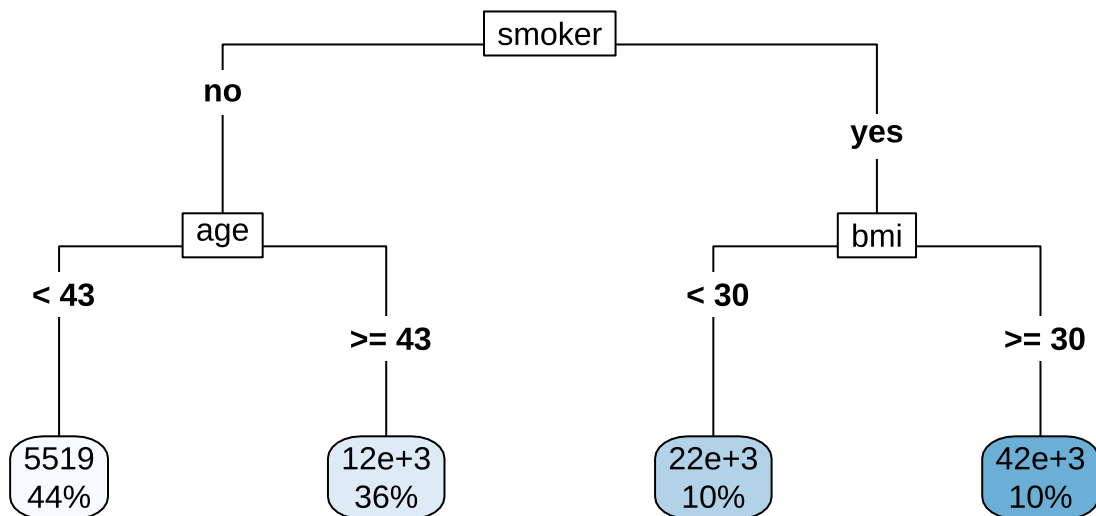
```
[1] age    bmi    smoker
```

Root node error:  $1.1434e+11/802 = 142563982$

n= 802

	CP	nsplit	rel error	xerror	xstd
1	0.619399	0	1.00000	1.00447	0.066367
2	0.139629	1	0.38060	0.38362	0.023519
3	0.067202	2	0.24097	0.24615	0.018445
4	0.010000	3	0.17377	0.17898	0.016712

```
rpart.plot(rpart(charges~.,train),type=5)
```



## 3 Algorithm 2 - Neural Networks

### 3.1 Mathematical Overview

A neural network is a type of machine learning model that is inspired by the structure and function of the human brain. It is composed of layers of interconnected nodes, or “neurons,” that process and transmit information. Neural networks are capable of handling complex and non-linear relationships between input and output variables.

A neural network can be represented mathematically as a series of layers, with each layer performing a set of mathematical operations on the input it receives. The output of one layer serves as the input to the next layer. The first layer, called the input layer, receives the input variables, and the last layer, called the output layer, produces the output variable(s). In between the input and output layers, there are one or more hidden layers that perform computations on the input data.

The mathematical notation for a neural network can be quite complex as it depends on the architecture of the network, the number of layers and the activation functions used.

In general, a neural network can be represented mathematically as:

$$y = f(W_n * f(W_{n-1} * f(\dots f(W_1 * x + b_1) + b_2) \dots + b_{n-1}) + b_n)$$

Where:

- $y$  = output variable
- $x$  = input variable
- $f$  = activation function (such as sigmoid, ReLU, etc.)
- $W$  = weight matrix
- $b$  = bias term

The mathematical notation for the neural network model can be represented as a series of layers, with each layer performing a set of mathematical operations on the input it receives.

The output of one layer serves as the input to the next layer.

The mathematical notation for each hidden layer would be:

$$h_i = f(W_i * h_{i-1} + b_i)$$

Where:

- $h_i$  is the output of the  $i$ -th hidden layer
- $W_i$  is the weight matrix for the  $i$ -th hidden layer
- $b_i$  is the bias term for the  $i$ -th hidden layer



- $f$  is the activation function

$f(x)$  is the output of the function

The mathematical notation for the input layer is:

$$h_0 = W_0 * x + b_0$$

Where:

- $h_0$  is the output of the input layer
- $W_0$  is the weight matrix for the input layer
- $x$  is the input data
- $b_0$  is the bias term for the input layer

The mathematical notation for the output layer is:

$$y = W_{n+1} * h_n + b_{n+1}$$

Where:

- $y$  is the output variable
- $h_n$  is the output of the last hidden layer
- $W_{n+1}$  is the weight matrix for the output layer
- $b_{n+1}$  is the bias term for the output layer

In our model, we use ReLU (Rectified Linear Unit) as activation function, which is defined mathematically as:

$$f(x) = \max(0, x)$$

where:

- $x$  is the input to the function

To train the model properly, an optimizer function should be defined, which adjust the model's parameters so as to minimize the cost function. We use MSE loss, which measures the difference between the model's predictions and the true labels, which mathematical notation is denoted as:

$$MSE = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2$$

where:

- $m$  is the number of examples in the training set
- $y^{(i)}$  is the true label of the  $i$ -th example
- $\hat{y}^{(i)}$  is the predicted label of the  $i$ -th example by the model

In our model, we use the Adam optimizer, which has the following mathematical notation:

$$\theta_{t+1} = \theta_t - \frac{\alpha}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t$$

Where:

- $\theta$  is the set of parameters to optimize
- $\alpha$  is the learning rate
- $\hat{m}_t$  is the biased first moment estimate
- $\hat{v}_t$  is the biased second raw moment estimate
- $\epsilon$  is a small positive constant that is added to the denominator to prevent division by zero

## 3.2 Initial Model

### Python libraries

To train the neural network model, we use Python and the following libraries:

- **pandas** and **numpy** for wrangling data
- **sklearn** for preprocessing data
- **torch** for defining the neural network model
- **pytorch\_lightning** for running and logging the training
- **optuna** for hyperparameter optimization

### Data preprocessing

For data preprocessing we transformed the categorical variables using **LabelEncoder** and scaled the numerical data using **MinMaxScaler** from **sklearn** library:

### Dataset Architecture

To run the training process, we have to define custom data classes first, **trainDataset**, **valDataset**, and **testDataset** that are used to load and preprocess the data for training, validation and testing respectively.

### Model Architecture

The **RegressionModel** class takes four parameters:

- **layers**: an integer representing the number of hidden layers in the network

- num\_neurons: an integer representing the number of neurons in each hidden layer
- learning\_rate: a float representing the learning rate for the optimizer
- batch\_size: an integer representing the batch size to use during training

First, we test the model with predefined parameter set. Therefore, we are using the **RegressionModel** class defined earlier to create an instance of the model

First, to check the performance of the neural network model, we train it for 50 epochs with the following parameters:

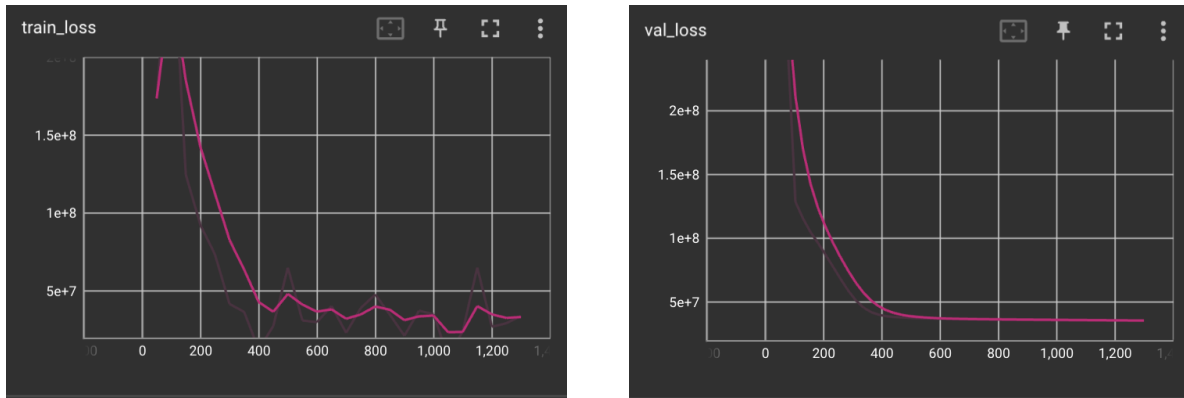
- num\_neurons: 64
- layers: 3
- learning\_rate: 0.001
- batch\_size: 32

Before the training, we can check how many weights will be used in the model:

	Name	Type	Params
0	input	Linear	448
1	hidden	Sequential	12.5 K
2	output	Linear	65
3	relu	ReLU	0

Total: 13.0 K Trainable params

Below you can see the train and validation losses during the training.



After the training of the initial model, we received the following MSE scores:

Model	$MSE_{train}$	$MSE_{valid}$
Neural Network (initial)	34806204	35508452

### 3.3 Hyperparameter Optimization

The HPO process on the neural network model was performed using **Optuna** library in Python.

“Optuna is an automatic hyperparameter optimization software framework, particularly designed for machine learning. It features an imperative, define-by-run style user API... the user of Optuna can dynamically construct the search spaces for the hyperparameters.”<sup>1</sup>

In this project, we are using Torch Lightning and Optuna to optimize a neural network. Torch Lightning handles the “model” which is the neural network, while Optuna handles tuning the hyperparameters. Optuna can be used to search for the optimal hyperparameters of a network and the results are used to fine-tune the network’s performance.

Optuna uses various different sampling algorithms to narrow the search space. Some of the algorithms are for example Grid search, random search, tree-structured Parzen Estimator, CMA-ES based algorithm, and quasi-Monte Carlo sampling algorithm. The default sampler is the Tree-structured Parzen Estimator, which is also being used in this project.

Tree-structured Parzen estimator is a type of bayesian estimator, where we fit a surrogate model to prior observations to estimate the hyperparameter space. The Key idea is instead of  $P(y|x)$ , we model  $P(x|y)$ , where  $x$  is a single hyperparameter and  $y$  is the loss. In TPE we maintain two surrogate models, one for the best cases (minimizing loss) and one for the worst cases (maximizing the loss). With each new observation, the model is fitted to the best and worst-case scenarios, giving us the most probable hyperparameters that result in the best and worst-case scenarios.<sup>2</sup>

Optuna gives the user control of the surrogate model (or surrogate) by setting a hierarchy of parameters and the tree structure. This allows the user to control the model’s fidelity, which is important when searching for optimal hyperparameters. Besides Sampling algorithms, Optuna also provides pruning algorithms for the early stopping of inefficient search trials. Pruning algorithms are useful when the search space is large, as they help to focus on only the efficient search trials while discarding the inefficient ones early on.

Since the dataset chosen for this project is rather on a smaller scale of things, we defined a rather small search space so no pruning was needed.

The steps in this project are as follows: - Train the model with default hyperparameters - Set up hyperparameter space - Initiate the Optuna

<sup>1</sup><https://optuna.readthedocs.io/en>

<sup>2</sup><https://www.jair.org/index.php/jair/article/view/13188>

Once Optuna has finished the trial runs and made a suggestion of the optimal hyperparameters, a final model with the best parameters is trained, validated, and finally tested.

Optuna and Torch Lightning can be combined to produce an efficient neural network model which can provide accurate and reliable results for machine learning tasks.

It is best used when the complexity of the problem and the size of the dataset are large.

We let Optuna optimize 4 parameters for us:

- ‘num\_neurons’: number of neurons in each layer, between 16 and 256
- ‘batch\_size’: amount of data points used at once, between 32 and 512
- ‘layers’: amount of layers, between 1 and 5
- ‘learning\_rate’: learning rate of the Adam optimizer, between 1e-5 and 1e-1

As the result of the HPO, the the following parameters were found to perform the best on the validation set:

- ‘num\_neurons’: 235
- ‘batch\_size’: 314
- ‘layers’: 5
- ‘learning\_rate’: 0.0748

### 3.4 Performance

Finally, we predict the outcome on the validation set and compare MSE values.

Model	$MSE_{valid}$
Neural Network (initial)	35508452
Neural Network (optimized)	18894440

We can see that the optimized model could achieve almost 50% improvement in terms of MSE value, which indicates that the performance of the model has been improved a lot.

## Future Work and Discussion

In this study, we evaluated the performance of three different models: a linear model, a regression tree, and a neural network. Our results showed that the regression tree outperformed the linear model on both of the chosen metrics. The neural network, in turn, outperformed both the linear model and the regression tree.

The linear model has the advantage of being the easiest to interpret of all three models, even though its performance is not as good as the other two models. The difference in performance between the regression tree and the linear model is significant, but the increase in performance of the neural network over the tree model is relatively small. This suggests that the regression tree may be a more suitable model for this task because it offers greater explainability in terms of how it can be visualized, easily explained to non-experts, and understood. The neural network, on the other hand, has some of the “black box” problems that some neural networks are prone to; it is difficult to “see inside” the model to understand how it is making its decisions. Note however, that the regression tree is also not a perfect model, as it only predicts 4 values, as we saw in the tree plot. This is not ideal for predicting a continuous outcome.

For future work, we could attempt to improve the interpretability of the neural network by using model-agnostic tools such as SHAP values. Additionally, we could improve the current regression tree by including other parameters into the hyperparameter optimization, testing different splitting criteria such as gini index vs. entropy. If that doesn’t improve performance, we could also build multiple trees and use the ensemble method Random Forest. However, this would reduce interpretability again.

In conclusion, the choice of model depends on the desired balance between performance and interpretability. While the neural network is the model that showed the best performance, the regression tree may be a more suitable choice for this task because of its greater explainability.