

[1] Install the required Libraries

▾ [3] Import the required pthon packages

```
1 import pandas as pd
2 import numpy as np
3 from numpy import arange
4 import matplotlib.pyplot as plt
5 %matplotlib inline
6 import seaborn as sns
7 import warnings
8 warnings.filterwarnings('ignore')
9 from IPython.display import display, HTML
10 import panel as pn
11 from colorama import Fore, Style, Back
12 from sklearn.model_selection import train_test_split
13 from sklearn.preprocessing import OneHotEncoder
14 from sklearn.ensemble import IsolationForest
15 from sklearn.neighbors import LocalOutlierFactor
16 from sklearn.preprocessing import StandardScaler
17 import tensorflow as tf
18 from sklearn.ensemble import RandomForestRegressor
19 import statsmodels.api as sm
20 import pickle
21 from tabulate import tabulate
22 from sklearn.ensemble import BaggingRegressor
23 from sklearn.ensemble import GradientBoostingRegressor
24 from sklearn.ensemble import AdaBoostRegressor
25 import xgboost as xgb
26 from sklearn import tree
27 from sklearn import metrics
28 import joblib
29 from typing import List, Union
```

▾ [4] Extra Settings

```
1 pd.set_option('display.max_rows', 50)
2 pd.options.display.float_format = '{:,.2f}'.format
3 pal = ['#C060A1', '#46C2CB', '#ECB365', '#FF0000', '#0039A6', '#414A4C']
4 plt.rcParams['figure.figsize'] = (8, 4)
5 plt.rcParams['font.size'] = 8
6 plt.rcParams['text.color'] = '#0039A6'
```

▾ [5] Backup Functions

```
1 # Define CSS styling
```

```

2 css = """
3 <style>
4     table {
5         width: 50%;
6         text-align: center;
7         'font-size': '12pt';
8     }
9 }
10 th, td {
11     padding: 8px;
12     text-align:center;
13 }
14 th {
15     background-color: #C70039;
16     color: white;
17 }
18 }
19 td {
20     text-align:center!important;
21     background-color: #EEEEEE;
22     color: black;
23 }
24 td::after {
25     content: attr(data-value " :.2f");
26 }
27 </style>
28 """
29 def printout(data, show_index=False, cells_colors=False):
30     """
31     Function to display a pandas DataFrame with custom formatting.
32     Removes trailing zeros and formats float numbers to 2 decimal places.
33     """
34     data = data.applymap(lambda x: "{:.2f}".format(x).rstrip('0').rstrip('.') if isinstance(x, float) else x)
35     styled_df = HTML(css + data.to_html(index=show_index))
36     display(styled_df)
37
38 def encode_categoricals(data):
39     """
40     Function to encode categorical columns in a pandas DataFrame using factorize method.
41     """
42     categorical_columns = data.select_dtypes(include=['object']).columns
43     for column in categorical_columns:
44         data[column], _ = pd.factorize(data[column])
45     return data
46

```

▼ [6] Dataset

Introduction

The dataset used in this section is the insurance.csv dataset, which contains 1338 observations and 7 features. The features are:

- Independent Features :

1. Age: Age of the customer.
2. BMI: Body Mass Index of the individual.
3. Children: Number of children individual has.
4. Smoker: Is that individual a smoker or not.
5. Region: Which region/part individual lives in.

- Dependent Feature

6. Charges: The charges column is our target/independent feature.

The dataset contains 4 numerical features (age, BMI, children, and charges) and 3 categorical features (smoker, region, and charges). The goal is to use this dataset to predict the insurance charges of individuals.

TO-DO List

- Load the dataset
- Get the dimensionality (no. of features (columns) and observations(rows))
- DTypes and names of features
- Descriptive statistics (Max, min, mean, median and std for each continuous feature)
- Declare the target variable (feature)

▼ [6 : 1] Load the Dataset

```

1 # Function to print colored output
2 def print_output(message, data):
3     print(message)
4     printout(data)
5
6 # Load the dataset
7 df = pd.read_csv('/content/gdrive/MyDrive/Colab_Notebooks/Medical Health Insurance Cost Prediction /us_insurance.csv')
8
9 # Print top 5 observations
10 print_output("The Top 5 observations of Dataset: \n", df.head())
11
12 # Print tail 5 observations
13 print_output("\n\n The Tail 5 observations of Dataset: \n", df.tail())

```

The Top 5 observations of Dataset:

age	sex	bmi	children	smoker	region	charges
19	female	27.9	0	yes	southwest	1684.92
18	male	33.77	1	no	southeast	1725.55
28	male	33	3	no	southeast	1449.46
33	male	22.7	0	no	northwest	1984.47
32	male	28.88	0	no	northwest	1866.86

▼ [6 : 2] Get the dimensionality

THE TOP 5 OBSERVATIONS OF DATASET.

```
1 print(Fore.BLUE + Style.BRIGHT +
2       f"The no. of observations :"+
3       Fore.LIGHTRED_EX +
4       f"{len(df)} \n\n" +
5       Fore.BLUE +
6       "The no. of features :"+
7       Fore.LIGHTRED_EX +
8       f"{len(df.columns)}\n\n" +
9       Fore.BLUE +
10      "The Dataset shape is : " +
11      Fore.LIGHTRED_EX +
12      f"{df.shape}"
13      + Style.RESET_ALL)
14
```

The no. of observations :1338

The no. of features :7

The Dataset shape is : (1338, 7)

▼ [6 : 3] DTypes and names of features

```
1 def print_dtype(df):
2     """Prints the name and dtype of each column in the DataFrame"""
3     for col in df.columns:
4         print(f"{Fore.BLUE}{Style.BRIGHT}The feature name and dtype | {Fore.GREEN}{col} --> {df[col].dtype}{Style.RESET_ALL}")
5
6 def get_features(df):
7     """Returns a list of feature names, removing any single quotes"""
8     return [name.replace("'", "") for name in df.columns]
9
10 def get_cat_features(df):
11     """Returns a list of categorical feature names"""
12     return df.select_dtypes(include=['object', 'category']).columns.tolist()
13
14 def get_num_features(df, cat_features):
15     """Returns a list of numerical feature names"""
16     return df.drop(columns=cat_features).columns.tolist()
```

```

17
18 def print_features(df):
19     """Prints the feature names, categorical feature names, and numerical feature names"""
20     features = get_features(df)
21     cat_features = get_cat_features(df)
22     num_features = get_num_features(df, cat_features)
23
24     print(f"{Style.BRIGHT}\n\nThe dataset features : {Fore.MAGENTA}{'', '.join(features)} \n{Fore.BLACK}The categorical features : {Fore.MAGENTA}{'', '.join(cat_features)}{Fore
25
26 # Call the functions
27 print_dtype(df)
28 print_features(df)

```

The feature name and dtype | age --> int64
 The feature name and dtype | sex --> object
 The feature name and dtype | bmi --> float64
 The feature name and dtype | children --> int64
 The feature name and dtype | smoker --> object
 The feature name and dtype | region --> object
 The feature name and dtype | charges --> float64

The dataset features : age, sex, bmi, children, smoker, region, charges
 The categorical features : sex, smoker, region
 The continuous features : age, bmi, children, charges

▼ [6 : 4] Descriptive Statistics

```

1 print("The Continuous Variables (features) \n")
2 printout(df.describe(), show_index=True)
3 print("\nThe Categorical Variables (features) \n")
4 printout(df.describe(include="object"), show_index=True)

```

The Continuous Variables (features)

age	bmi	children	charges
-----	-----	----------	---------

Zoom in on the highest and lowest charges in the dataset 

```

1 def print_max_min_charges(df):
2     # Get the max and min charges in the dataset
3     max_charge = df.loc[df['charges'].idxmax()]
4     min_charge = df.loc[df['charges'].idxmin()]
5
6     # Combine both observations into a dataframe
7     max_min_charges = pd.DataFrame({'The Highest Charges': max_charge, 'The Lowest Charges': min_charge})
8
9     # Print the dataframe with index
10    printout(max_min_charges, show_index=True)
11
12 print_max_min_charges(df)

```

	The Highest Charges	The Lowest Charges
age	54	18
sex	female	male
bmi	47.41	23.21
children	0	0
smoker	yes	no
region	southeast	southeast
charges	63770.43	1121.87

Insights :

You can see the notable differences in charges as well as in other features like age, BMI, and Smoker by looking at the highest and lowest table. It's still early to determine or recognize the importance of predictors. However, further analysis and research are needed to fully understand their impact on insurance charges.

The person with the highest rate of charges is a smoker, 54 years old, and has a 47 BMI (that means he's obese; according to the World Health Organization (WHO), a BMI of less than 18.5 is considered underweight, a BMI between 18.5 and 24.9 is considered normal weight, a BMI between 25 and 29.9 is considered overweight, and a BMI of 30 or higher is considered obese), while the individual with lower insurance charges looks younger, smokes less, and has a healthier lifestyle overall.

▼ [6 : 5] Declare the target variable (feature)

Target Variable : is the variable that the model is trying to predict or estimate. It is also known as the dependent variable or response variable. In regression models like this, the target variable is typically a continuous variable, such as a numeric value or a real number, and the goal is to find a relationship between the target variable and one or more predictor variables, also known as independent variables or features.

The purpose of model is to predict the insurance **charges** (target variable) by using some predictor variables (*age, sex, bmi, children, smoker, region*)

```
1 def format_charges(value):
2     return '{:.2f}$'.format(value)
3
4 def print_charge_stats(label, value, color):
5     print(f"{Style.BRIGHT}\n{label} : {color}{value}{Style.RESET_ALL}")
6
7 charges_stats = df['charges'].agg(['max', 'min', 'mean']).apply(format_charges)
8
9 print_charge_stats("The maximum Charges", charges_stats[0], Fore.GREEN)
10 print_charge_stats("The minimum Charges", charges_stats[1], Fore.MAGENTA)
11 print_charge_stats("The average of Charges", charges_stats[2], Fore.RED)
```

The maximum Charges : **63770.43\$**

The minimum Charges : **1121.87\$**

The average of Charges : **13270.42\$**

▼ [7] Explatory Data Analysis (EDA)

EDA is an essential step in regression modeling that involves understanding the distribution of variables and identifying relationships. Various visualization libraries like seaborn, matplotlib, and plotly are used to create different plots and graphs to gain insights into the dataset. Some significant questions to ask and answer during the EDA process for a health insurance prediction project with the given dataset include:

1. What is the distribution of the target variable (charges)?
2. What is the distribution of each attribute in the dataset, including age, sex, BMI, children, region, and smoker status?
3. Are there any missing or null values in the dataset? If so, how are they distributed across the attributes?
4. What is the correlation between each attribute and the target variable (charges)? Which attributes have the strongest correlation with the target variable?
5. Are there any significant differences in the distribution and correlation of the attributes based on different categories or subgroups (e.g., age, sex, smoker status, region)?
6. Are there any outliers in the dataset? If so, how do they affect the distribution and correlation of the attributes?
7. What insights can be gained from visualizing the data using scatter plots, histograms, box plots, and other visualization techniques?
8. Are there any significant interactions between the attributes that affect the prediction of the target variable?

Answering these questions can help to identify patterns, trends, and relationships in the data, and can provide valuable insights for feature engineering and model selection.

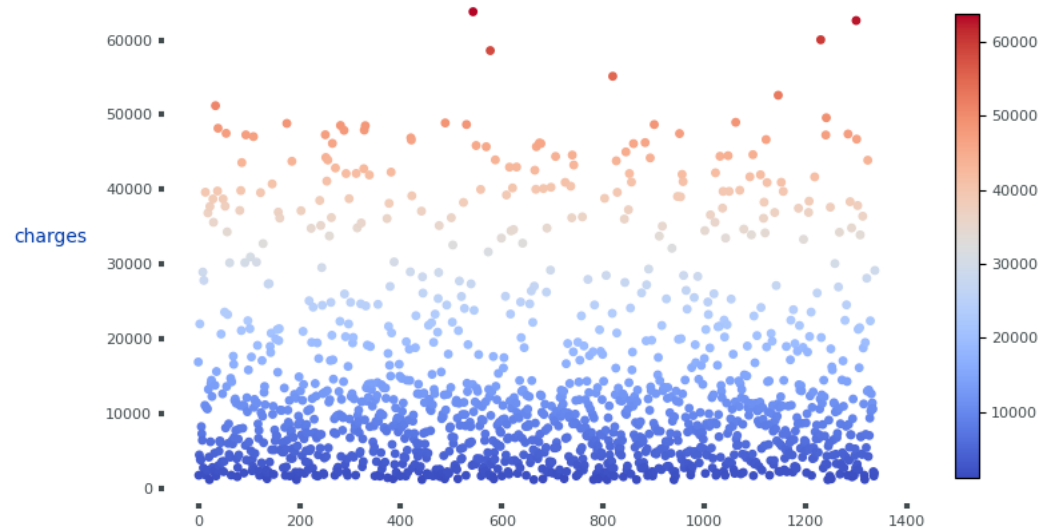
▼ [7 : 1] Distribution of Target Variable (charges)

```
1 plt.figure(figsize=(9, 5))
2 plt.scatter(df.index, df['charges'], c = df['charges'], cmap='coolwarm', s=15)
```

```

3 cbar = plt.colorbar(pad=0.05,shrink=0.9)
4 for t in cbar.ax.get_yticklabels():
5     t.set_fontsize(8)
6     t.set_color("#414A4C")
7 ax= plt.gca()
8 for spine in ax.spines.values():
9     spine.set_linewidth(0)
10 ax.set_ylabel('charges', fontsize=10, rotation=360, labelpad=30, color = "#0039A6")
11 ax.tick_params(axis='x', colors='#414A4C',length=3, width=3, labels=8) # Set the color of the x-axis tick marks and labels
12 ax.tick_params(axis='y', colors='#414A4C', length=3, width=3, labels=8)
13
14 plt.show()

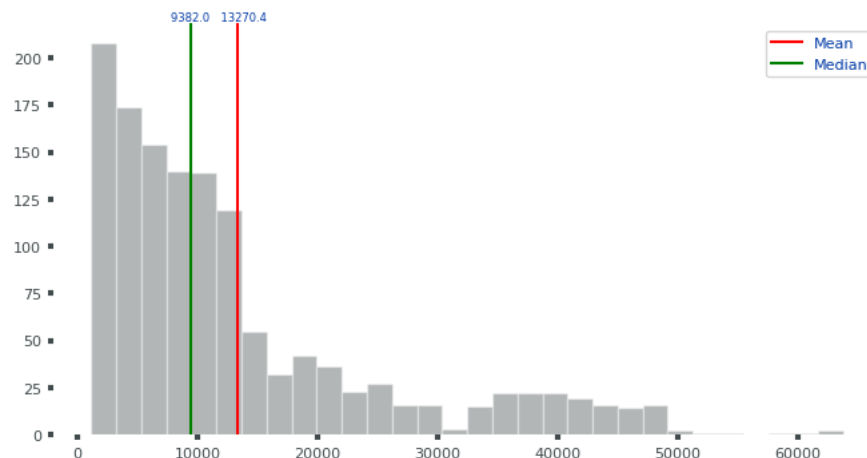
```



```

1 plt.figure(figsize=(8, 4))
2 sns.distplot(df['charges'], kde=False,color='#414A4C', hist_kws={'edgecolor': 'white'})
3 plt.axvline(x=np.mean(df['charges']), color = 'red', label='Mean',)
4 plt.text(np.mean(df['charges']), plt.ylim()[1], f"      {np.mean(df['charges']):.1f}", ha='center', va='bottom', size=6)
5 plt.text(np.mean(df['charges']), plt.ylim()[1], f"{np.median(df['charges']):.1f}      ", ha='right', va='bottom', size=6)
6
7 plt.axvline(x=np.median(df['charges']), color = 'green', label='Median')
8 plt.legend(loc="upper right")
9 ax= plt.gca()
10 for spine in ax.spines.values():
11     spine.set_linewidth(0)
12 ax.set_xlabel('_____ charges _____', fontsize=10, rotation=360, labelpad=30, color = "#0039A6")
13 ax.tick_params(axis='x', colors='#414A4C',length=3, width=3, labels=8) # Set the color of the x-axis tick marks and labels
14 ax.tick_params(axis='y', colors='#414A4C', length=3, width=3, labels=8)
15
16 plt.show()

```

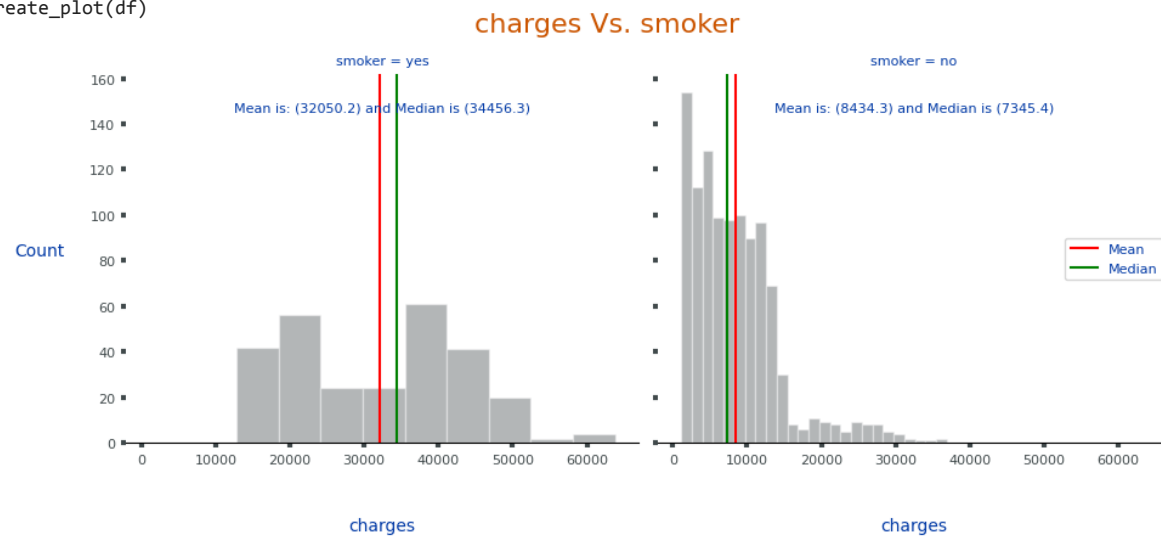
▼ [7 : 2] Target Variable Vs. Independent Variables

```

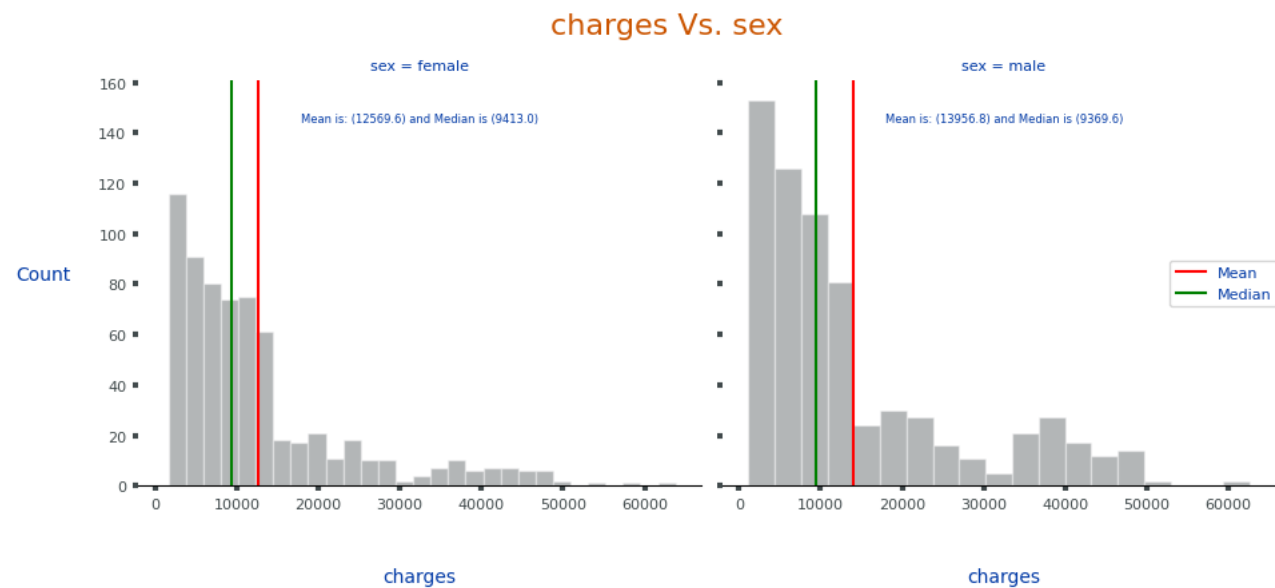
1 plt.rcParams['text.color'] = '#CC5500'
2 def create_plot(data, var_1="smoker", var_2 = "charges", metrics=True, metrics_size= 8):
3
4     g = sns.FacetGrid(data, col=var_1)
5
6     g.add_legend()
7     g.legend.get_title().set_fontsize(30)
8     g.fig.set_size_inches(10, 4.5)
9
10    g.fig.suptitle(f'{var_2} Vs. {var_1}', fontsize=16, color = "#CC5500")
11
12    for ax in g.axes.flat:
13        ax.set_xlabel(var_2, fontsize=10, color='#0039A6', labelpad=30)
14        ax.set_ylabel('Count', fontsize=10, color='#0039A6', rotation = 360, labelpad=30)
15        ax.tick_params(axis='x', colors='#414A4C',length=3, width=3, labels=8)
16        ax.tick_params(axis='y', colors='#414A4C', length=3, width=3, labels=8)
17    g.map(sns.distplot, var_2, color='#414A4C', kde=False,hist_kws={'edgecolor': 'white'})
18
19    if metrics == True:
20        for category, ax in g.axes_dict.items():
21            # Get the count of each level
22            mean = df[df[var_1] == category][var_2].mean()
23            median = df[df[var_1] == category][var_2].median()
24            ax.axvline(x=mean, color = 'red', label='Mean')
25            ax.axvline(x=median, color = 'green', label='Median')
26
27            ax.text(0.5, 0.9, f"Mean is: ({mean :.1f}) and Median is ({median :.1f})",
28                  transform=ax.transAxes, ha='center', color="#0039A6", size = metrics_size)
29        plt.legend(loc="center right")
30    for ax in g.axes.flat:
31        ax.spines['left'].set_visible(False)
32

```

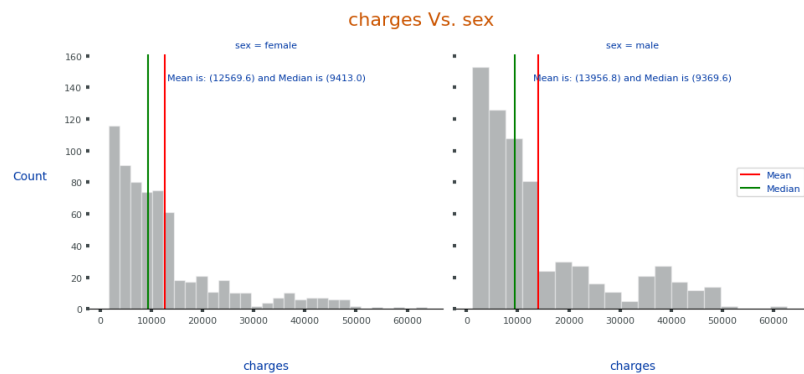
```
33 return plt.show()
34 create_plot(df)
```



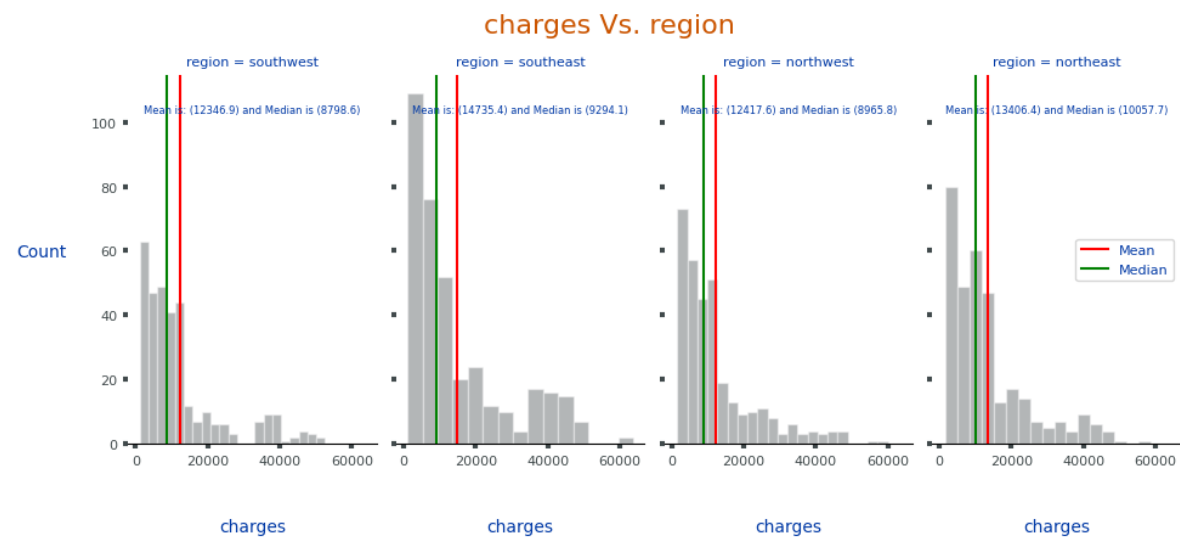
```
1 create_plot(data=df,var_1='sex', metrics = True, metrics_size = 6)
```



```
1 create_plot(data=df,var_1='sex', metrics = True)
```



```
1 create_plot(data=df,var_1='region', metrics = True, metrics_size=6)
```



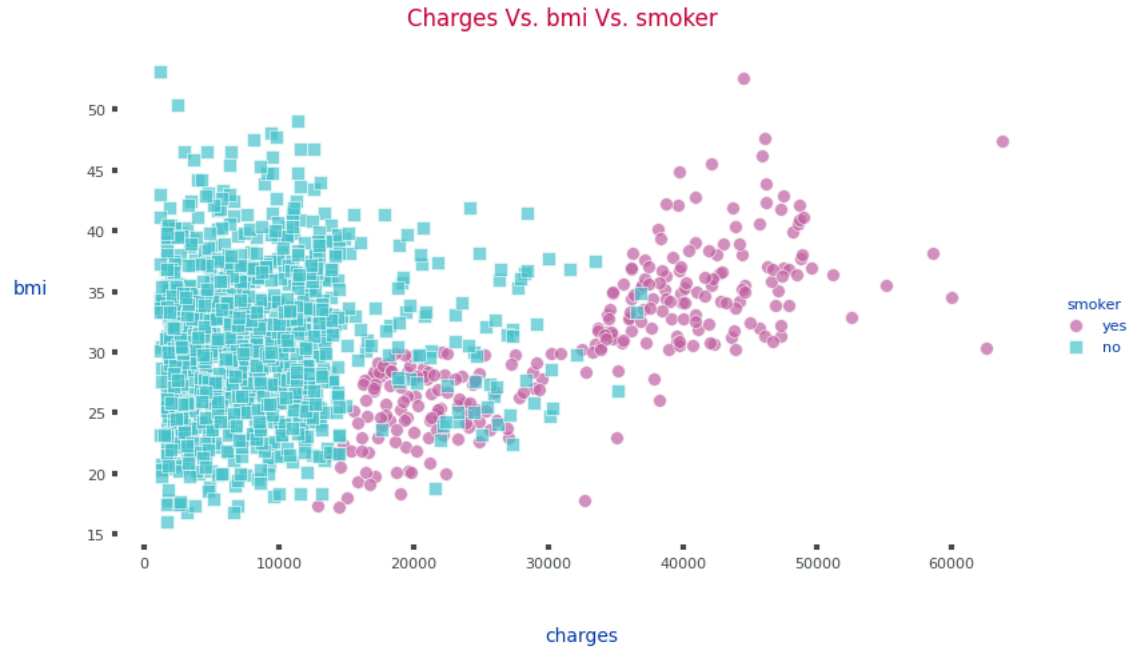
```
1 def bivariate_plot(data, hue, var):
2     # Initialize the FacetGrid object
3     g = sns.FacetGrid(data, hue=hue, palette=pal,
4                       hue_kws={'marker': ['o', 's', 'D', 'v', '^', '*']}, height=5)
5
6     # Set figure size
7     g.fig.set_size_inches(8, 5)
8
```

```

9     # Set title of the figure
10    g.fig.suptitle(f'Charges Vs. {var} Vs. {hue}', fontsize=12, color="#C70039")
11
12    for ax in g.axes.flat:
13        # Set visibility of left and bottom spines
14        ax.spines['left'].set_visible(False)
15        ax.spines['bottom'].set_visible(False)
16
17        # Set the labels for the y-axis and x-axis
18        ax.set_ylabel(var, fontsize=10, rotation=360, labelpad=30, color="#0039A6")
19        ax.set_xlabel('charges', fontsize=10, rotation=360, labelpad=30, color="#0039A6")
20
21        # Set the color of the x-axis and y-axis tick marks and labels
22        ax.tick_params(axis='x', colors='#414A4C', length=3, width=3, labelsiz=8)
23        ax.tick_params(axis='y', colors='#414A4C', length=3, width=3, labelsiz=8)
24
25    # Add a scatterplot to the grid
26    g.map(sns.scatterplot, "charges", var, s=50, alpha=.7)
27
28    # Add a legend to the grid
29    g.add_legend()
30
31    # Show the plot
32    plt.show()

```

```
1 bivariate_plot(df, hue='smoker', var = "bmi")
```



```

1 def barplot(data, var, hue):
2     plt.rcParams['font.size'] = 8
3     g = sns.barplot(data=data, x=var, y='charges', hue=hue, palette=pal, width=0.6)
4     g.set_title(f' charges Vs. {var} (Hue : {hue})', fontsize=12, color="#C70039", y=1.1)
5     fig, ax = plt.gcf(), plt.gca()
6
7     setup_figure_and_axes(fig, ax)
8     setup_labels_and_ticks(ax, var)
9     annotate_bars(ax, data, var, hue)
10    setup_legend(ax)
11
12    plt.show()
13
14
15 def setup_figure_and_axes(fig, ax):
16     sns.move_legend(ax, "upper right", bbox_to_anchor=(1.1, 0.9))
17     fig.set_figwidth(8)
18     fig.set_figheight(4)
19     for spine in ax.spines.values():
20         spine.set_linewidth(0)
21     ax.spines['left'].set_visible(False)
22     ax.spines['bottom'].set_visible(False)
23
24
25 def setup_labels_and_ticks(ax, var):
26     ax.set_ylabel("charges", fontsize=10, rotation=360, labelpad=30, color="#0039A6")
27     ax.set_xlabel(var, fontsize=10, rotation=360, labelpad=30, color="#0039A6")
28     ax.tick_params(axis='x', colors='#414A4C', length=3, width=3, labelsize=8)
29     ax.tick_params(axis='y', colors='#414A4C', length=3, width=3, labelsize=8)
30
31
32 def annotate_bars(ax, data, var, hue):
33     v = h = 0
34     for i, p in enumerate(ax.patches):
35         unique_var_values = data[var].unique().tolist()
36         hue_value_counts = data[data[var] == unique_var_values[v]][hue].value_counts()
37         cat, total = hue_value_counts[h], len(data[data[var] == unique_var_values[v]])
38         precentage = "{:.2f} %".format((100 * cat) / total)
39         height = p.get_height()
40         ax.text(p.get_x() + p.get_width() / 2., height + 0.1 * height, precentage, ha="center", fontsize=6, color='#0039A6')
41         v += 1
42         if v >= data[var].nunique():
43             v, h = 0, h + 1
44
45
46 def setup_legend(ax):
47     plt.setp(ax.get_legend().get_texts(), fontsize='6', fontweight='bold', color="#0039A6")
48     plt.setp(ax.get_legend().get_title(), fontsize='8', color="#C70039")

```

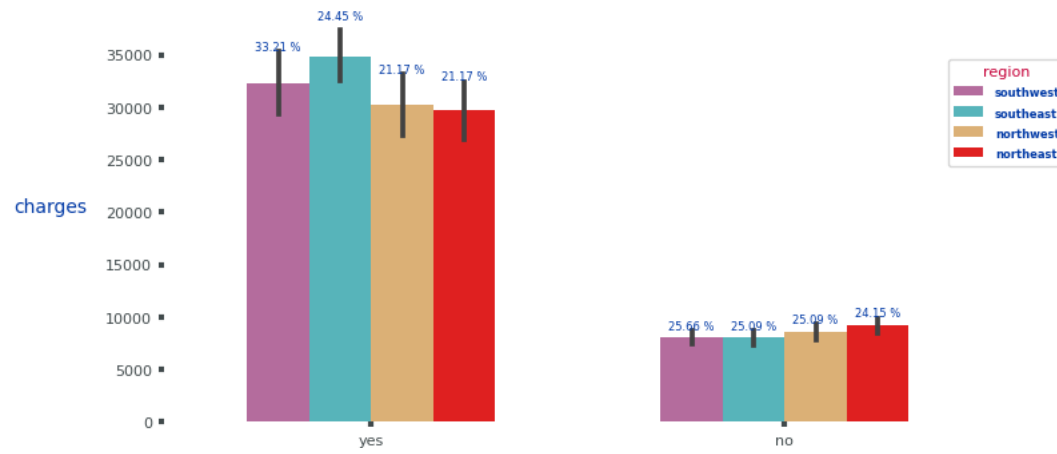


```

1 barplot(df, "smoker", "region")

```

charges Vs. smoker (Hue : region)

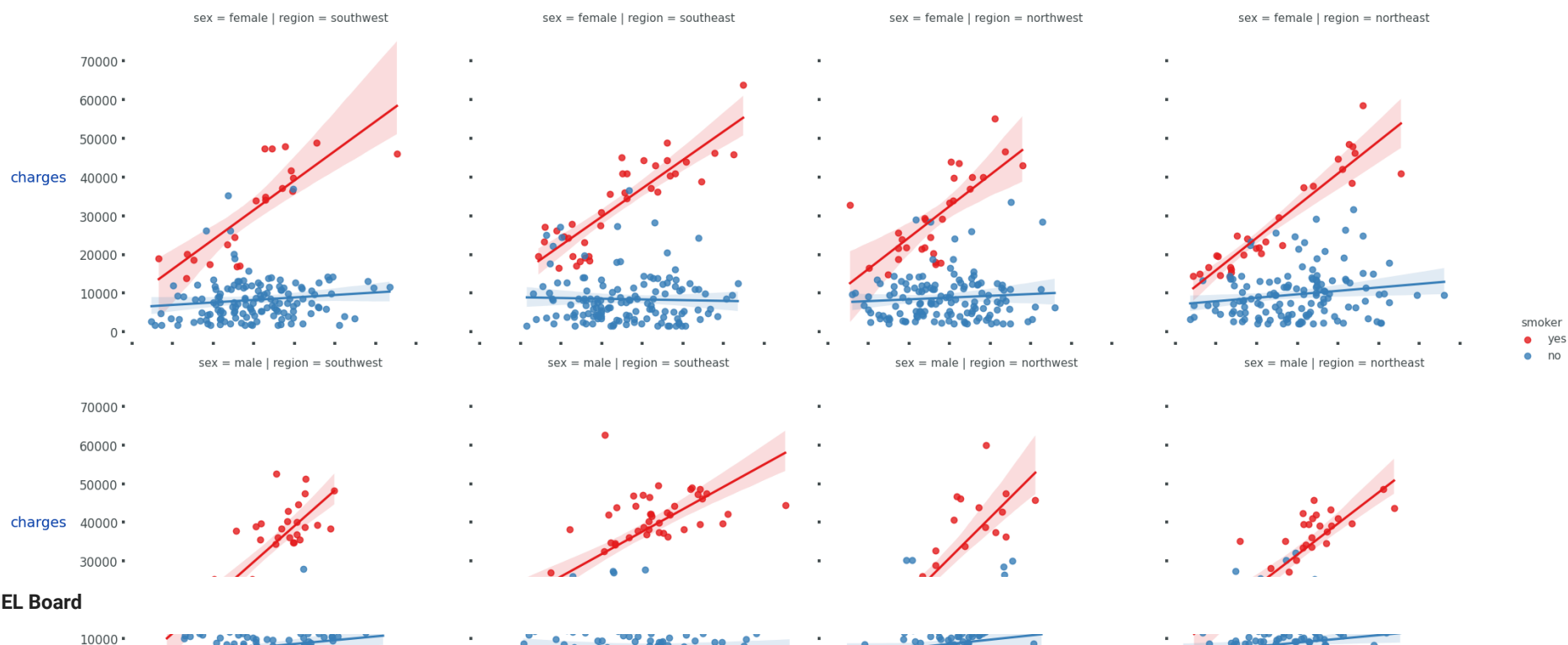


```

1 def plot_regression(df):
2     plot_config = {
3         'font.size': 11,
4         'text.color': "#414A4C"
5     }
6
7     plt.rcParams.update(plot_config)
8
9     g = sns.lmplot(x="bmi", y="charges", row="sex", col="region", hue='smoker', data=df, palette="Set1")
10
11     modify_axes(g)
12
13     plt.show()
14
15
16 def modify_axes(g):
17     g.fig.suptitle('Regression Plot of Charges Vs. BMI considering (Hue : Smoker)', fontsize=20, color = "#C70039", y=1.1)
18
19     for ax in g.axes.flat:
20         ax.spines['left'].set_visible(False)
21         ax.spines['bottom'].set_visible(False)
22
23         ax.set_ylabel("charges", fontsize=14, rotation=360, labelpad=40, color="#0039A6")
24         ax.set_xlabel('BMI', fontsize=14, rotation=360, labelpad=30, color="#0039A6")
25
26         ax.tick_params(axis='x', colors='#414A4C', length=3, width=3, labels=12)
27         ax.tick_params(axis='y', colors='#414A4C', length=3, width=3, labels=12)
28
29
30 plot_regression(df)

```

Regression Plot of Charges Vs. BMI considering (Hue : Smoker)



▼ PANEL Board

```

1 #pn.extension(sizing_mode="stretch_width", design='material', template="fast")
2 #variable_widget_1 = pn.widgets.Select(name="categorical variables", value="smoker", options=df.select_dtypes(include=['object', 'category']).columns.to_list())
3 #variable_widget_2 = pn.widgets.Select(name="continous variables", value="charges", options=df.select_dtypes(exclude=['object', 'category']).columns.to_list())
4 #bound_plot = pn.bind(create_plot,data=df, var_1=variable_widget_1, var_2=variable_widget_2)
5

```

```

1 #app = pn.Column(variable_widget_1, variable_widget_2, bound_plot)
2 #app

```

▼ [7 : 3] Data Preparation and Cleansing

8-Step checklist

- Missing Values :
- Duplicate and Low Variation Data
- Incorrect and Irrelevant Data
- Categorical Data
- Outliers
- Feature Scaling

- Feature Selection
- Validation Split

Missing Values

```

1 def print_missing_values_info(df):
2     """
3     Function to print information about missing values in our dataset.
4
5     Parameters:
6     df (DataFrame): The dataset to check for missing values.
7
8     """
9     # Check missing values in total dataset
10    missing_any = df.isnull().values.any()
11    print(f"{Style.BRIGHT}Any NULL values: {Fore.BLUE}{missing_any}{Style.RESET_ALL}\n")
12
13    # Check missing values per column
14    missing_per_feature = df.isnull().sum().astype(int).to_string(index=True, header=False)
15    print(f"{Style.BRIGHT}NULL values per feature:\n {Fore.GREEN}{missing_per_feature}\n{Style.RESET_ALL}")
16
17 print_missing_values_info(df)

```

Any NULL values: **False**

NULL values per feature:

```

age      0
sex      0
bmi      0
children 0
smoker   0
region   0
charges  0

```

Insights : No missing values in the dataset

Duplicated Data

```

1 # Detect duplicated data
2 dups = df.duplicated()
3
4 # If there are any duplicates, print them out and then remove them
5 if dups.any():
6     printout(df[dups], show_index=True)
7     # Drop the duplicated observations by index
8     df = df.drop(df[dups].index, axis=0).reset_index(drop=True)

```

	age	sex	bmi	children	smoker	region	charges
581	19	male	30.59	0	no	northwest	1639.56

Insights : It can be seen that there is no duplicated data in the dataset except for a single observation. and in our case, the similarity may happen during the data collection process. There are no identical features here, just age, region, smoking habit, sex, and children. However, it is important to note that even though there are no identical features, the presence of observations with duplicated data can still affect the accuracy and reliability of the analysis. Since we're dealing with only one observation, we are fine to proceed with our analysis without any concerns.

Incorrect and Irrelevant Data

Anomaly Detection Algorithms :

Isolation Forest model represents the anomaly score assigned to each data point in your dataset. Anomaly scores indicate how likely a data point is to be an outlier or anomaly. In the Isolation Forest algorithm, lower anomaly scores indicate a higher likelihood of being an anomaly. Anomaly scores usually range between -1 and 1, where values closer to -1 are considered more anomalous.

```

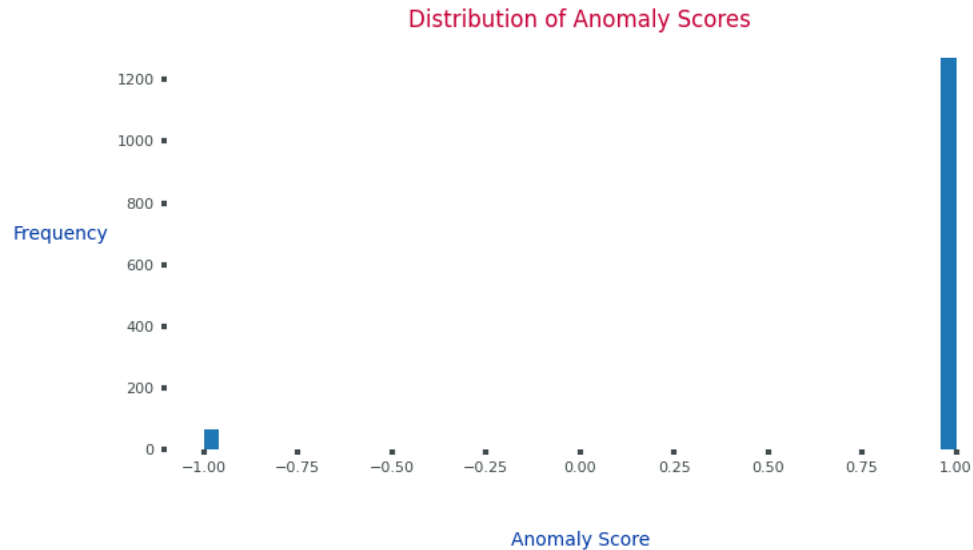
1 def preprocess_data(data):
2     """Preprocesses the data."""
3     data = data.drop('region', axis=1).copy()
4
5     # get dummies of categorical features
6     data = encode_categoricals(data)
7
8     return data
9
10 def detect_anomalies(data):
11     """Detects anomalies in the data."""
12     iso_forest = IsolationForest(contamination=0.05)
13     anomaly_scores = iso_forest.fit_predict(data)
14
15     # Convert anomaly_scores to a Series
16     anomaly_scores_series = pd.Series(anomaly_scores, index=data.index)
17
18     # Filter out anomalies (rows with negative anomaly scores)
19     filtered_df = data[anomaly_scores_series >= 0]
20     filtered_df = filtered_df.reset_index(drop=True)
21
22     return filtered_df, anomaly_scores
23
24 def plot_anomaly_scores(anomaly_scores):
25     """Plots the distribution of anomaly scores."""
26     plt.hist(anomaly_scores, bins=50)
27     ax=plt.gca()
28     for spine in ax.spines.values():
29         spine.set_linewidth(0)
30     ax.set_ylabel('Frequency', fontsize=10, rotation=360, labelpad=30, color = "#0039A6")
31     ax.set_xlabel('Anomaly Score', fontsize=10, rotation=0, labelpad=30, color = "#0039A6")
32     ax.tick_params(axis='x', colors='#414A4C',length=3, width=3, labelsize=8) # Set the color of the x-axis tick marks and labels
33     ax.tick_params(axis='y', colors='#414A4C', length=3, width=3, labelsize=8)
34
35     plt.title('Distribution of Anomaly Scores', fontsize=12, color = "#C70039")
36     plt.show()
37

```

```

38
39 preprocessed_data = preprocess_data(df)
40 filtered_df, anomaly_scores = detect_anomalies(preprocessed_data)
41 plot_anomaly_scores(anomaly_scores)

```



insights : The dataset appears to contain some anomalies that need additional investigation or may even need to be removed, according to the anomaly detection algorithm discussed above. These anomalies could potentially affect the accuracy and reliability of any analysis or predictions made using the dataset. It is crucial to identify the cause of these anomalies and determine whether they are due to errors in data collection, measurement, or other factors before taking any further action. but we have no way to identify the cause, so the anomalous data points will be dropped from the dataset.

Categorical Data

```

1 # Get percentages of Categorical variables
2 for col in get_cat_features(df):
3     printout(pd.DataFrame((df[col].value_counts(normalize=True) * 100).apply(lambda x: "{:.2f}".format(x).rstrip('0').rstrip('.') + " %")), show_index=True)

```

sex	
male	50.49 %
female	49.51 %

```

1 # Count observations per level for all categorical features
2 nums = df.select_dtypes(exclude=['object', 'category']).columns.to_list()
3 for col in get_num_features(df, nums):
4     printout(pd.DataFrame(pd.crosstab(index=df[col], columns = "freq").rename_axis(None, axis=1)), show_index=True)

```

sex		freq
female		662
male		675
smoker		freq
no		1063
yes		274
region		freq
northeast		324
northwest		324
southeast		364
southwest		325

Insights : You can see that there are three categorical features in the dataset, which include sex (male or female), smoker (yes or no), and region (southeast, southwest, northeast, northwest). The male category in sex feature is top with 50.52% (676 males), the positive-smokers category is top with 79.52% (1067 positive-smokers), which looks imbalanced and may require more attention during analysis, and the southeast region is top with 27.2% (regions have roughly equal distributions as follows: southeast (27.20%), southwest (24.29%), northwest (24.29%), and northeast (24.20%)). These categorical features provide valuable insights into the demographics and lifestyle choices of the individuals in the dataset, allowing for further analysis and understanding of potential factors influencing insurance charges.

Feature Selection

USING * Correlation with Random Forest Regressor

```

1 class DataAnalysis:
2     def __init__(self, df, target):
3         self.df = df
4         self.X = df.drop(target, axis=1)
5         self.y = df[target]
6

```

```

7     def calculate_correlation(self):
8         correlation = self.X.corrwith(self.y)
9         return correlation
10
11    def train_model(self, model):
12        model.fit(self.X, self.y)
13        return model
14
15    def get_feature_importances(self, model):
16        feature_importances = model.feature_importances_
17        return feature_importances
18
19    def plot(self, importance):
20        indices = np.argsort(importance)
21        plt.title('\nFeature Importance with Random Forest Regresson method', fontsize=12, color = "#C70039")
22        plt.barh(range(len(indices)), importance[indices], color='#D8D9DA', #4D3C77", "#C70039
23                align='center')
24        ax = plt.gca()
25        for spine in ax.spines.values():
26            spine.set_linewidth(0)
27        ax.spines['left'].set_visible(False)
28        ax.spines['bottom'].set_visible(False)
29        ax.set_ylabel("Features", fontsize=10, rotation=360, labelpad=30, color="#0039A6")
30        ax.set_xlabel("Relative Importance", fontsize=10, rotation=360, labelpad=30, color="#0039A6")
31        ax.tick_params(axis='x', colors='#414A4C', length=3, width=3, labelsize=8)
32        ax.tick_params(axis='y', colors='#414A4C', length=3, width=3, labelsize=8)
33
34        plt.yticks(range(len(indices)), self.X.columns[indices])
35        plt.show()
36
37
38    def display(self, importance, comment, return_importance=False):
39        feature_importance = pd.DataFrame({'Feature': self.X.columns, 'Importance': importance})
40
41        # Sort the DataFrame by importance score
42        feature_importance = feature_importance.sort_values(by='Importance', ascending=False)
43
44        # Print the feature names and importance scores in a pretty format
45        print(f"\n{comment}\n")
46        printout(feature_importance)
47        if return_importance :
48            return feature_importance
49
50
51
52
53 ds = df.copy()
54 ds = encode_categoricals(ds)
55 analysis = DataAnalysis(ds, 'charges')
56
57 correlation = analysis.calculate_correlation()
58 analysis.display(correlation," Correlation scores with train.corrwith(target) method : ")
59
60 model = RandomForestRegressor(random_state=42)
61 trained_model = analysis.train_model(model)

```

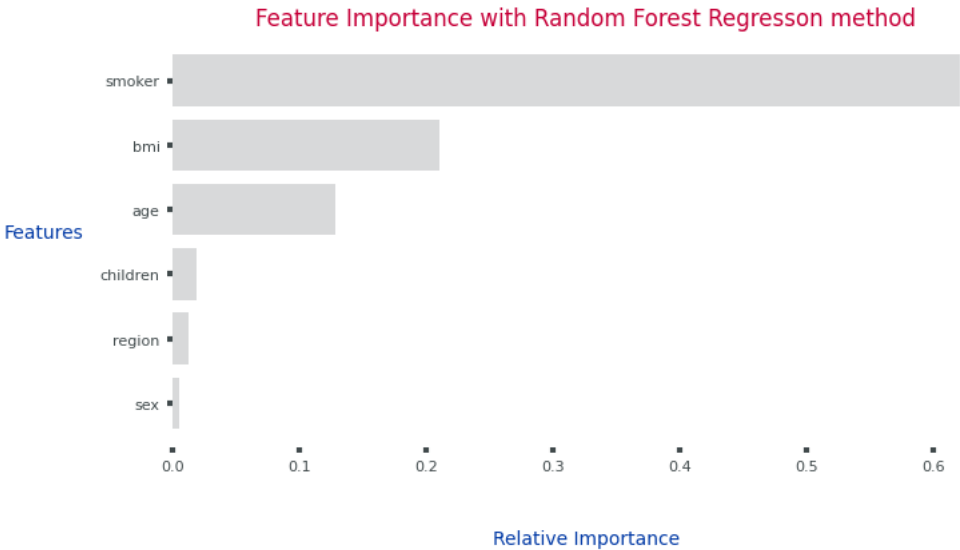
```
62
63 feature_importances = analysis.get_feature_importances(trained_model)
64 feature_importance_df = analysis.display(feature_importances," Correlation scores with Random Forest Regressor method :", return_importance = True)
65 analysis.plot(feature_importances)
```

Correlation scores with train.corrwith(target) method :

Feature	Importance
age	0.3
bmi	0.2
children	0.07
sex	0.06
region	0.01
smoker	-0.79

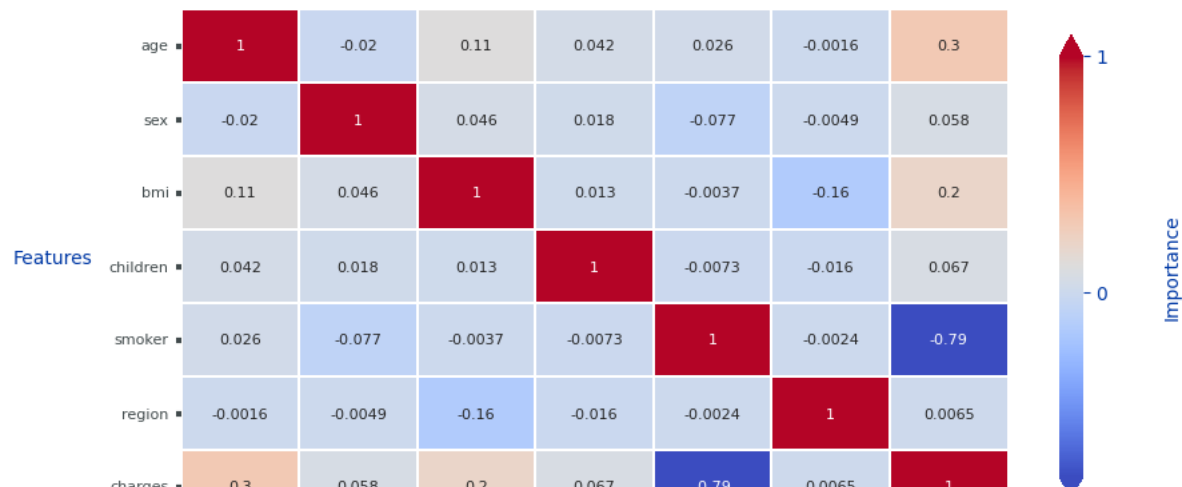
Correlation scores with Random Forest Regressor method :

Feature	Importance
smoker	0.62
bmi	0.21
age	0.13
children	0.02
region	0.01
sex	0.01



USING * Correlation with heatmap

```
1 white',annotfontsize=8, fontsize=10, labelpad=30, axis_color="#0039A6", tick_params={'colors': '#414A4C', 'length': 3, 'width': 3, 'labelsize': 8}):
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30 lor=axis_color)
31 lor=axis_color)
32
33
34
35
36
37
38
39
40 labelpad)
41
42
43
44
45
```



Insights 📊 It looks like the target feature (charges) has a high correlation with the smoker feature in the US Medical Health Insurance Dataset and a low correlation with age and BMI features. This points out that smoking has the potential to be a significant factor in the underlying use case with a scale of 0.79, while age and BMI come next in terms of their impact on the target feature with a scale of 0.3 and 0.2, respectively. Other features like sex, children, and region have no detectable correlation with the target feature. These findings indicate that smoking has a strong influence on the charges. Age and BMI also have some influence, although to a lesser extent. However, factors such as sex, number of children, and region do not seem to have any noticeable impact on the target feature.

Actions 🛠️ We'll drop *region* feature, the lowest correlated variable and this step will help us to build a global app that is not related to the demographic dataset source

```
1 df.drop('region', axis=1, inplace=True)
2 printout(df.head())
```

age	sex	bmi	children	smoker	charges
19	female	27.9	0	yes	16884.92
18	male	33.77	1	no	1725.55
28	male	33	3	no	4449.46
33	male	22.7	0	no	21984.47
32	male	28.88	0	no	3866.86

Outliers

USING : Isolation Forest & Local Outlier Factor

```
1 def format_plot(ax, x_label, y_label):
2     for spine in ax.spines.values():
3         spine.set_linewidth(0)
4     ax.set_ylabel(y_label, fontsize=9, labelpad=20, color="#0039A6")
```

```
5 ax.set_xlabel(x_label, fontsize=10, labelpad=30, color="#0039A6")
6 ax.tick_params(axis='x', colors='#414A4C', length=3, width=3, labelsiz=8)
7 ax.tick_params(axis='y', colors='#414A4C', length=3, width=3, labelsiz=8)
8
9 # Extract the charges column (target variable)
10 ds = df.copy()
11 charges = ds['charges'].values.reshape(-1, 1)
12
13 # Encode categorical features
14 ds = encode_categoricals(ds)
15
16 # Standardize the features
17 scaler = StandardScaler()
18 X_scaled = scaler.fit_transform(ds)
19
20 # Apply Isolation Forest for outlier detection
21 iso_forest = IsolationForest(contamination=0.005)
22 iso_preds = iso_forest.fit_predict(X_scaled)
23
24 # Apply Local Outlier Factor (LOF) for outlier detection
25 lof = LocalOutlierFactor(contamination=0.05)
26 lof_preds = lof.fit_predict(X_scaled)
27
28 # Visualize the results
29 plt.figure(figsize=(10, 4))
30
31 plt.subplot(1, 2, 1)
32 plt.scatter(charges, X_scaled[:, 2], c=iso_preds, cmap='Set1')
33 ax=plt.gca()
34 format_plot(ax, 'Charges', 'Standardized bmi')
35 plt.title('Isolation Forest Outlier Detection', fontsize=11, color='#C70039')
36
37 plt.subplot(1, 2, 2)
38 plt.scatter(charges, X_scaled[:, 2], c=lof_preds, cmap='Set1')
39 ax=plt.gca()
40 format_plot(ax, 'Charges', 'Standardized bmi')
41 plt.title('Local Outlier Factor (LOF) Outlier Detection', fontsize=11, color='#C70039')
42
43 plt.tight_layout()
44 plt.show()
```


Isolation Forest Outlier Detection



Local Outlier Factor (LOF) Outlier Detection



Display outliers in the dataset :

```

1 # Return outlier observations using Isolation Forest
2 iso_outliers = ds[iso_preds == -1]
3 print("Outliers detected by Isolation Forest:\n")
4 printout(iso_outliers, show_index=True)
5
6 # Return outlier observations using LOF
7 lof_outliers = ds[lof_preds == -1]
8 print("\n\nOutliers detected by Local Outlier Factor:\n")
9 printout(lof_outliers, show_index=True)

```

Outliers detected by Isolation Forest:

	age	sex	bmi	children	smoker	charges
438	52	0	46.75	5	1	12592.53
494	21	1	25.7	4	0	17942.11
543	54	0	47.41	0	0	63770.43
549	43	0	46.2	0	0	45863.21
859	37	0	47.6	2	0	46113.51
1046	22	1	52.58	1	0	44501.4
1084	39	0	18.3	5	0	19023.26

Outliers detected by Local Outlier Factor:

	age	sex	bmi	children	smoker	charges
3	33	1	22.7	0	1	21984.47
9	60	0	25.84	0	1	28923.14
32	19	0	28.6	5	1	4687.8
62	64	1	24.7	1	1	30166.62
102	18	0	30.11	0	1	21344.85
115	60	1	28.59	0	1	30260
116	58	1	49.06	0	1	11381.33
140	34	1	22.42	2	1	27375.9
166	20	0	37	5	1	4830.63
172	18	1	15.96	0	1	1694.8
219	24	0	23.21	0	1	25081.77
245	54	1	30.02	0	1	24476.48
286	46	0	48.07	2	1	9432.93
289	52	1	26.4	3	1	25992.82
321	26	0	29.64	4	1	24671.66
340	24	0	27.6	0	1	18955.22
354	18	0	38.28	0	1	14133.04
355	46	1	27.6	0	1	24603.05
379	62	1	31.46	1	1	27000.98
387	50	1	25.36	2	1	30284.64
397	21	1	31.02	0	1	16586.5
430	19	1	33.1	0	1	23082.96
438	52	0	46.75	5	1	12592.53

442	18	1	43.01	0	1	1149.4
443	59	0	36.52	1	1	28287.9
473	47	0	33.34	0	1	20878.78
491	61	0	25.08	0	1	24513.09
520	50	0	27.36	0	1	25656.58
525	18	0	33.88	0	1	11482.63
533	37	1	36.19	0	1	19214.71
539	53	1	31.35	0	1	27346.04
554	25	0	41.33	0	1	17878.9
568	49	0	31.9	5	1	11552.9
573	62	0	36.86	1	1	31620
582	32	0	23.65	1	1	17626.24
598	52	0	37.52	2	1	33471.97
639	33	1	42.4	5	1	6666.24
657	48	0	35.91	1	1	26392.26
659	37	1	46.53	3	1	6435.62
780	18	1	41.14	0	1	1146.8
805	40	0	41.42	1	1	28476.73
846	23	1	50.38	1	1	2438.06
857	25	0	32.23	1	1	18218.16
875	49	0	27.1	1	1	26140.36

Display the descriptive statistics before and after eliminating outliers

```

968      39      0      34.32      5      1      8596.83

1 def print_cleaned_data_stats(cleaned_data, outlier_type):
2     print(f"\n\nCleaned Data after removing {outlier_type} outliers:\n")
3     printout(cleaned_data[nums].describe(), show_index=True)
4
5 # Remove outliers detected by Isolation Forest
6 data_cleaned_iso = ds[iso_preds == 1]
7
8 # Remove outliers detected by Local Outlier Factor
9 data_cleaned_lof = ds[lof_preds == 1]
10
11 print("Original Data before removing outliers:\n")
12 printout(df.describe(), show_index=True)
13
14 # Print cleaned data stats
15 print_cleaned_data_stats(data_cleaned_iso, 'Isolation Forest')
16 print_cleaned_data_stats(data_cleaned_lof, 'Local Outlier Factor')

```

Original Data before removing outliers:

	age	bmi	children	charges
count	1337	1337	1337	1337
mean	39.22	30.66	1.1	13279.12
std	14.04	6.1	1.21	12110.36
min	18	15.96	0	1121.87
25%	27	26.29	0	4746.34
50%	39	30.4	1	9386.16
75%	51	34.7	2	16657.72
max	64	53.13	5	63770.43

Cleaned Data after removing Isolation Forest outliers:

	age	bmi	children	charges
count	1330	1330	1330	1330
mean	39.23	30.61	1.09	13161.19
std	14.05	6.01	1.2	11963.11
min	18	15.96	0	1121.87
25%	27	26.3	0	4724.37
50%	39	30.33	1	9303.3
75%	51	34.59	2	16443.29
max	64	53.13	5	62592.87

Cleaned Data after removing Local Outlier Factor outliers:

	age	bmi	children	charges
count	1270	1270	1270	1270
mean	39.3	30.55	1.08	12934.53
std	13.98	5.93	1.16	12067.22
min	18	16.82	0	1121.87
25%	27	26.22	0	4625.25
50%	39.5	30.3	1	9053.38
75%	51	34.5	2	14585.95

Insights : Taking a closer look at the chart and statistical tables above, it looks like ISO performing better than LOF, but both detectors are not performing as well as they should, as they have a high rate of false positives. This indicates that these detectors are incorrectly identifying a significant number of data points as outliers when they are actually not. It is important to consider alternative outlier detection methods that

can provide more accurate results in our specific case. Additionally, further analysis and evaluation of the data may be required to understand the underlying factors contributing to these false positives and refine our outlier detection approach accordingly.

USING : TensorFlow AutoEncoder

```

1 def standardize_features(data):
2     """Standardize the features of the dataframe."""
3     ds = data.copy()
4     ds = encode_categoricals(ds)
5
6     scaler = StandardScaler()
7     X_scaled = scaler.fit_transform(ds)
8     return X_scaled
9
10 def build_autoencoder(input_dim,X_scaled):
11     """Build and train autoencoder for outlier detection."""
12     autoencoder = tf.keras.Sequential([
13         tf.keras.layers.Input(shape=(input_dim,)),
14         tf.keras.layers.Dense(8, activation='relu'),
15         tf.keras.layers.Dense(input_dim, activation='linear')
16     ])
17     autoencoder.compile(optimizer='adam', loss='mean_squared_error')
18     autoencoder.fit(X_scaled, X_scaled, epochs=50, batch_size=32, verbose=0)
19     return autoencoder
20
21 def calculate_errors(autoencoder, X_scaled):
22     """Calculate reconstruction errors."""
23     reconstructions = autoencoder.predict(X_scaled)
24     reconstruction_errors = np.mean(np.square(X_scaled - reconstructions), axis=1)
25     return reconstruction_errors
26
27 def remove_outliers(data, reconstruction_errors):
28     """Identify and remove outliers."""
29     threshold = np.percentile(reconstruction_errors, 92)
30     outlier_indices = np.where(reconstruction_errors > threshold)[0]
31     data_cleaned_autoencoder = data.drop(outlier_indices)
32     data_cleaned_autoencoder = data_cleaned_autoencoder.reset_index(drop=True)
33     return data_cleaned_autoencoder, threshold
34
35 def display_data(data, data_cleaned_autoencoder, num_features):
36     """Display original and cleaned data."""
37     print("Original Data before removing outliers:\n")
38     printout(data.describe(), show_index=True)
39
40     print("\n\nCleaned Data after removing Autoencoder outliers:\n")
41     printout(data_cleaned_autoencoder[num_features].describe(), show_index=True)
42
43 # Using the helper functions
44 X_scaled = standardize_features(df)
45 input_dim = X_scaled.shape[1]
46 autoencoder = build_autoencoder(input_dim,X_scaled)
47 reconstruction_errors = calculate_errors(autoencoder, X_scaled)

```

```
48 data_cleaned_autoencoder, threshold = remove_outliers(df, reconstruction_errors)
```

```
49 display([data=df, data_cleaned_autoencoder, threshold])
```

Original Data before removing outliers:

	age	bmi	children	charges
count	1337	1337	1337	1337
mean	39.22	30.66	1.1	13279.12
std	14.04	6.1	1.21	12110.36
min	18	15.96	0	1121.87
25%	27	26.29	0	4746.34
50%	39	30.4	1	9386.16
75%	51	34.7	2	16657.72
max	64	53.13	5	63770.43

Cleaned Data after removing Autoencoder outliers:

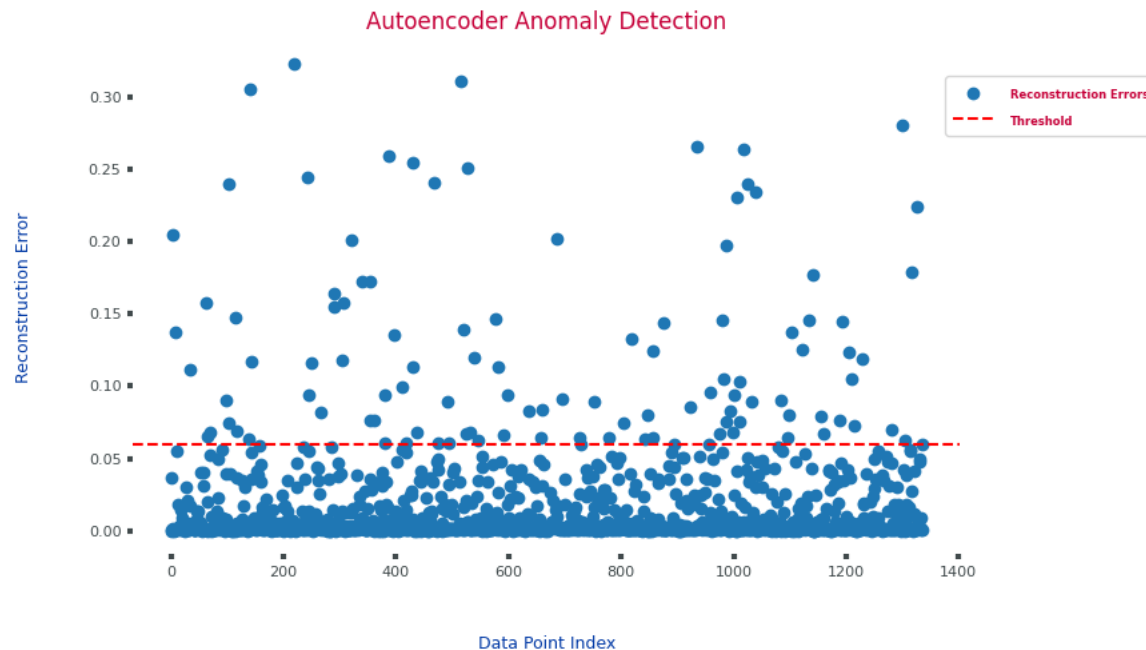
	age	bmi	children	charges
count	1230	1230	1230	1230
mean	39.49	30.81	1.08	12432.59
std	14.01	5.97	1.2	11901.27
min	18	15.96	0	1121.87
25%	27	26.51	0	4464.06
50%	40	30.5	1	8589.57
75%	51	34.8	2	13614.11
max	64	52.58	5	63770.43

```
1 #Plot reconstruction errors
2 fig, ax = plt.subplots(figsize=(8, 5))
3 colors = {'axes': '#0039A6', 'title': '#C70039', 'ticks': '#414A4C'}
4 fontsize = {'labels': 9, 'title': 12, 'legend': 6, 'legend_title': 8}
5 tick_params = {'length': 3, 'width': 3, 'labelsize': 8}
6
7 # Scatter plot
8 ax.scatter(range(len(reconstruction_errors)), reconstruction_errors, cmap='Set1', marker='o', label='Reconstruction Errors')
9
10 # Threshold line
11 ax.axhline(y=threshold, color='red', linestyle='--', label='Threshold')
12
13 # Hide axes spines
14 for spine in ax.spines.values():
15     spine.set_linewidth(0)
16
17 # Set labels
18 ax.set_ylabel('Reconstruction Error', fontsize=fontsize['labels'], labelpad=30, color=colors['axes'])
19 ax.set_xlabel('Data Point Index', fontsize=fontsize['labels'], labelpad=30, color=colors['axes'])
```

```

19 ax.set_xlabel('Data Point Index', fontsize=fontsize['labels'], labelpad=30, color=colors['axes'])
20
21 # Set ticks
22 ax.tick_params(axis='x', colors=colors['ticks'], **tick_params)
23 ax.tick_params(axis='y', colors=colors['ticks'], **tick_params)
24
25 # Set title
26 ax.set_title('Autoencoder Anomaly Detection', fontsize=fontsize['title'], color=colors['title'])
27
28 # Set legend
29 legend = ax.legend(loc='upper right', bbox_to_anchor=(1.25, 0.95))
30 plt.setp(legend.get_texts(), fontsize=fontsize['legend'], fontweight='bold', color=colors['title'])
31 plt.setp(legend.get_title(), fontsize=fontsize['legend_title'], color=colors['title'])
32
33 plt.show()

```



Normalize and Split Data

```

1 # normalize and split
2 def split_data(data, use_log=False):
3     target_var = 'charges'
4     numerical_vars = ['age', 'bmi', 'children']
5     test_size = 0.2
6     random_state_val = 42
7
8     X, y = prepare_data_variables(data, target_var, use_log)
9     X, scaler = preprocess_features(X, numerical_vars)
10

```

```

11     return split_dataset(X, y, test_size, random_state_val, scaler)
12
13 def prepare_data_variables(data, target_var, use_log):
14     X = data.drop(target_var, axis=1)
15     y = np.log10(data[target_var]) if use_log else data[target_var]
16     return X, y
17
18 def preprocess_features(X, numerical_vars):
19     X_num, X_cat = split_numerical_categorical(X, numerical_vars)
20     X_num, scaler = scale_numerical_features(X_num)
21     X_cat = encode_categorical_features(X_cat)
22     return pd.concat([X_num, X_cat], axis=1), scaler
23
24 def split_numerical_categorical(data, numerical_vars):
25     return data[numerical_vars], data.drop(numerical_vars, axis=1)
26
27 def scale_numerical_features(data_num):
28     scaler = StandardScaler()
29     return pd.DataFrame(scaler.fit_transform(data_num), columns=data_num.columns), scaler
30
31 def encode_categorical_features(data_cat):
32     encoder = OneHotEncoder(drop='first', sparse=False)
33     data_cat_encoded = encoder.fit_transform(data_cat)
34     return pd.DataFrame(data_cat_encoded, columns=encoder.get_feature_names_out(data_cat.columns))
35
36 def split_dataset(X, y, test_size, random_state, scaler):
37     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=test_size, random_state=random_state)
38     return scaler, X, y, X_train, X_test, y_train, y_test

1 scaler, X, y, X_train, X_test, y_train, y_test = split_data(data_cleaned_autoencoder)
2 print(X_train.shape, y_train.shape,
3       X_test.shape, y_test.shape)
4
(984, 5) (984,) (246, 5) (246,)

```

▼ [8] Models

```

1 #OLS Regression
2 x2 = sm.add_constant(np.asarray(X_train))
3 #fit linear regression model
4 model = sm.OLS(y_train, x2).fit()
5 model.summary()

```


OLS Regression Results

Dep. Variable: charges **R-squared:** 0.864
Model: OLS **Adj. R-squared:** 0.864
Method: Least Squares **F-statistic:** 1246.
Date: Mon, 04 Sep 2023 **Prob (F-statistic):** 0.00
Time: 16:43:31 **Log-Likelihood:** -9652.2
No. Observations: 984 **AIC:** 1.932e+04
Df Residuals: 978 **BIC:** 1.935e+04
Df Model: 5

Covariance Type: nonrobust

	coef	std err	t	P> t	[0.025	0.975]
const	7496.4314	208.696	35.920	0.000	7086.889	7905.974
x1	3593.0878	140.652	25.546	0.000	3317.073	3869.103
x2	1770.3597	140.415	12.608	0.000	1494.810	2045.910
x3	549.4349	140.664	3.906	0.000	273.396	825.474
x4	-76.3320	282.732	-0.270	0.787	-631.163	478.499

```

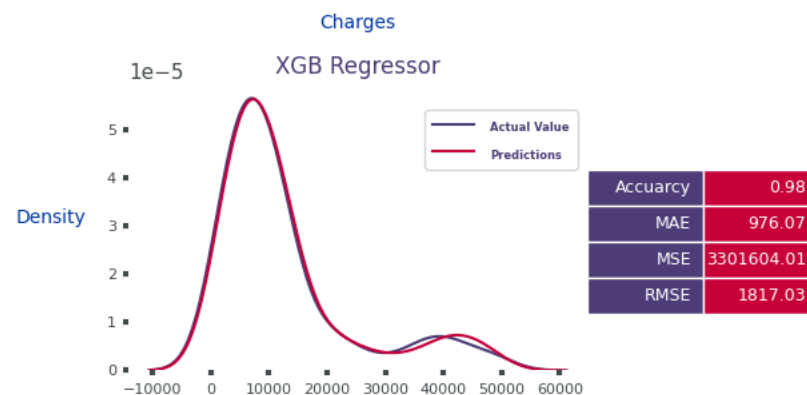
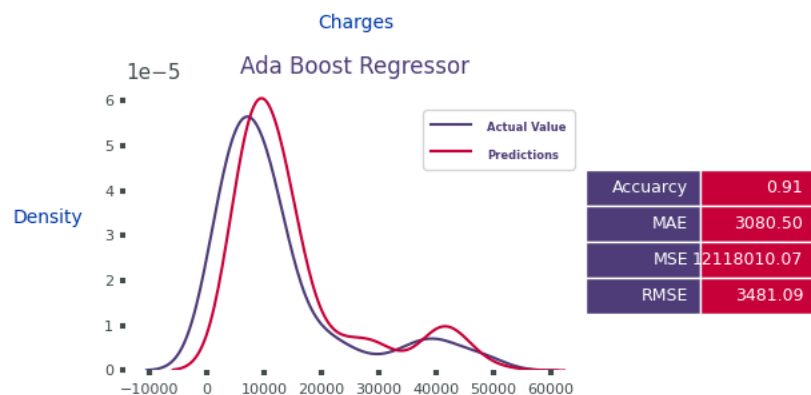
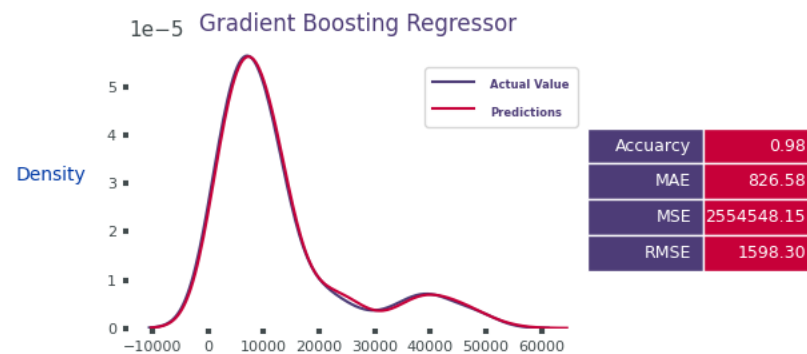
1 #Define a function for model training and prediction
2 def train_and_predict(models, X_train, y_train, X_test, y_test):
3     results, preds = [], []
4     for model_name, model in models:
5         model.fit(X_train, y_train)
6         accuracy = model.score(X_test, y_test)
7         predictions = model.predict(X_test)
8         preds.append([model_name, model, "{:.2f}".format(accuracy), predictions])
9         results.append([model_name, "{:.2f}".format(accuracy)])
10    results.sort(key=lambda x: x[1], reverse=True)
11    return results, preds
12
13 # Define models
14 models = [
15     ("Bagging Regressor", BaggingRegressor(tree.DecisionTreeRegressor(random_state=1))),
16     ("Gradient Boosting Regressor", GradientBoostingRegressor()),
17     ("Ada Boost Regressor", AdaBoostRegressor()),
18     ("XGB Regressor", xgb.XGBRegressor()),
19     ("Random Forest Regressor", RandomForestRegressor(n_estimators=100, criterion='squared_error', min_samples_split=2))
20     # Add more models here
21 ]
22
23 # Train models and make predictions
24 results, preds = train_and_predict(models, X_train, y_train, X_test, y_test)
25
26 # Display results in a table
27 table_headers = ["Model", "Accuracy"]
28 table = tabulate(results, headers=table_headers, tablefmt="pretty")
29 #print(table)
30
31 # Plotting
32 num_plots = len(models)
33 num_rows = int(np.ceil(np.sqrt(num_plots)))
34 num_cols = int(np.ceil(num_plots / num_rows))
35
36 fig, axes = plt.subplots(nrows=num_rows, ncols=num_cols, figsize=(14, 10))
37

```

```

38 for i, (model_name, model, accuracy, predictions) in enumerate(preds):
39     ax = axes[i // num_cols, i % num_cols]
40     sns.distplot(y_test, hist=False, color="#4D3C77", label="Actual Value", ax=ax)
41     sns.distplot(predictions, hist=False, color="#C70039", label="Predictions" , ax=ax)
42     metrics_data = [['Accuracy', accuracy],
43                     ['MAE', "{:.2f}".format(metrics.mean_absolute_error(y_test, predictions))],
44                     ['MSE', "{:.2f}".format(metrics.mean_squared_error(y_test, predictions))],
45                     ['RMSE', "{:.2f}".format(np.sqrt(metrics.mean_squared_error(y_test, predictions)))]
46     table = ax.table(cellText=metrics_data, edges='closed', bbox=[1,0.2,0.5,0.5],
47                     cellColours=[["#4D3C77", "#C70039"]]*4)
48     table.auto_set_font_size(False)
49     table.set_fontsize(9)
50     table.scale(1, 2)
51     for cell in table.get_celld().values():
52         cell.get_text().set_color('white')
53         cell.set_edgecolor('white')
54     ax.set_title(f"{model_name}", fontsize=12, color = "#4D3C77")
55     for spine in ax.spines.values():
56         spine.set_linewidth(0)
57     ax.spines['left'].set_visible(False)
58     ax.spines['bottom'].set_visible(False)
59
60     # Setting the labels for the axis
61     ax.set_ylabel("Density", fontsize=10, rotation=360, labelpad=30, color="#0039A6")
62     ax.set_xlabel("Charges", fontsize=10, rotation=360, labelpad=30, color="#0039A6")
63     tick_params={'colors': '#414A4C', 'length':3, 'width':3, 'labelsize':8}
64     # Setting the tick parameters
65     ax.tick_params(axis='x', **tick_params)
66     ax.tick_params(axis='y', **tick_params)
67     legend = ax.legend(loc='upper right', bbox_to_anchor=(1, 0.95))
68     plt.setp(legend.get_texts(), fontsize=fontsize['legend'], fontweight='bold', color="#4D3C77")
69
70 # Remove empty subplots
71 for i in range(len(models), num_rows * num_cols):
72     fig.delaxes(axes.flatten()[i])
73
74 plt.tight_layout()
75 plt.savefig('models.png')
76 plt.show()

```



```
1 # save the iris classification model as a pickle file
2 #model_pkl_file = "gbr_for_mhip_98.pkl"
3
4 #with open(model_pkl_file, 'wb') as file:
5 #    pickle.dump(model, file)
```

```
1 #joblib.dump(scaler, 'scaler.pkl')
```

Pass new data

```
1 #
2
3 def load_model(model_path: str, scaler_path: str):
4     """
5     Load the saved model and scaler.
6     """
7     model = joblib.load(model_path)
8     scaler = joblib.load(scaler_path)
9
10    return model, scaler
11
12 def create_dataframe(values: List[Union[int, float]], column_names: List[str]) -> pd.DataFrame:
13     """
```

```

12     Create a pandas DataFrame from the provided values and column names.
13     """
14     data = dict(zip(column_names, values))
15     df = pd.DataFrame([data])
16
17     return df
18
19 def scale_and_format_data(df: pd.DataFrame, scaler) -> pd.DataFrame:
20     """
21     Scale and format the data in the DataFrame.
22     """
23     # Scale the numerical columns
24     df.iloc[:, :3] = scaler.transform(df.iloc[:, :3])
25
26     # Apply formatting function to all values in the DataFrame
27     format_function = lambda x: f'{x:.2f}'
28     df = df.applymap(format_function)
29
30     return df
31
32 def predict_charges(model, data: pd.DataFrame) -> float:
33     """
34     Make predictions on the data.
35     """
36     charges = model.predict(df.values)
37
38     return "{:.2f} {}".format(charges[0])
39
40 # Define constants
41 MODEL_PATH = '/content/gdrive/MyDrive/Colab_Notebooks/Medical Health Insurance Cost Prediction /rfr_for_mhip_98.pkl'
42 SCALER_PATH = '/content/gdrive/MyDrive/Colab_Notebooks/Medical Health Insurance Cost Prediction /scaler.pkl'
43 VALUES = [20, 30.5, 2, 0, 0]
44 COLUMN_NAMES = ['age', 'bmi', 'children', 'sex', 'smoker']
45
46 # Load the model and scaler
47 model, scaler = load_model(MODEL_PATH, SCALER_PATH)
48
49 # Create a DataFrame
50 new_df = create_dataframe(VALUES, COLUMN_NAMES)
51
52 # Scale and format the data
53 new_df = scale_and_format_data(new_df, scaler)
54
55 # Make predictions
56 charges = predict_charges(model, new_df)
57
58 print(charges)

```

2164.86 \$

🔍 🔍 🔍 END 🔍 🔍 🔍

Key Insights

1. According to our analysis of the dataset, the most influential factors influencing medical insurance charges are whether the person is a smoker, their BMI (body mass index), and their age.
2. Compared to non-smokers, smokers frequently face significantly higher medical insurance premiums. This important realization emphasizes how lifestyle decisions have a significant impact on healthcare expenditures.
3. An individual's age and the cost of their insurance are strongly associated. Given the higher risk of health problems as people age, it makes sense that older people tend to have higher medical costs.
4. Higher BMI values are associated with higher medical insurance premiums. This implies that keeping a healthy weight can reduce healthcare costs.
5. We evaluated multiple regression models, including Random Forest Regressor, Gradient Boosting Regressor, Bagging Regressor, AdaBoost Regressor, and XGBRegressor. Among these, the Gradient Boosting Regressor and Random Forest Regressor showed better performance in predicting medical insurance charges.
6. Gradient Boosting Regression model showed the lowest Mean Absolute Error, Mean Square Error, and Root Mean Square Error among all models, indicating its capability to better capture the underlying patterns in the data.

Conclusion

This project highlights the process of estimating medical health insurance charges using numerous regression models. Regression models are statistical tools that analyze the relationship between dependent and independent variables. By utilizing these models, this project aims to provide insights into the factors that influence medical health insurance costs and develop accurate predictions based on the available data. Additionally, it explores the effectiveness of different regression techniques such as linear regression, bagging regression, gradient boosting regression, and random forest regression in capturing the complexity of this domain. In order to provide precise forecasts in the field of estimating medical insurance costs, the project emphasizes the significance of data preparation, model selection, and evaluation.