**Insurance Data Analysis Outline**  Introduction Data Overview Exploratory Data Analysis Feature Engineering Hypothesis Testing Results and Conclusion Introduction Insurance companies have a tough task at determining premiums for their customers. While the health care law in the United States as an example does have some rules for the companies to follow to determine premiums, it is really up to the companies on what factors they want to hold more weight to. So, what are the most important factors? how much statistical importance do they hold? and between LASSO, Ridge and Linear regression which model has the highest accuracy? To answer this problematic there is nothing better than the most interpretable model, linear regression, but before going into modeling, we will try in this report to uncover some insights about our data using EDA and Hypothesis Testing. **Data Overview** The insurance dataset was took from here and this is a brief description: "Machine Learning with R by Brett Lantz is a book that provides an introduction to machine learning using R. As far as I can tell, Packt Publishing does not make its datasets available online unless you buy the book and create a user account which can be a problem if you are checking the book out from the library or borrowing the book from a friend. All of these datasets are in the public domain but simply needed some cleaning up and recoding to match the format in the book." Columns age: age of primary beneficiary sex: insurance contractor 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 region: the beneficiary's residential area in the US, northeast, southeast, southwest, northwest. charges: Individual medical costs billed by health insurance Acknowledgements The dataset is available on GitHub here **Exploratory Data Analysis** import pandas as pd import seaborn as sns import numpy as np import matplotlib.pyplot as plt sns.set\_style('whitegrid') !pip install wget import wget file = wget.download('https://raw.githubusercontent.com/stedy/Machine-Learning-with-R-datasets/master/insurance Requirement already satisfied: wget in c:\users\optimist\anaconda3\lib\site-packages (3.2) 100% [......] 54288 / 54288 data = pd.read csv(file) In [4]: data.head() Out[4]: bmi children smoker region charges sex 19 female 27.900 yes southwest 16884.92400 0 18 male 33.770 no southeast 1725.55230 male 33.000 2 28 no southeast 4449.46200 3 no northwest 21984.47061 33 male 22.705 male 28.880 32 3866.85520 no northwest data.describe() bmi children charges age **count** 1338.000000 1338.000000 1338.000000 1338.000000 39.207025 30.663397 1.094918 13270.422265 mean 14.049960 6.098187 1.205493 12110.011237 std 18.000000 15.960000 0.000000 1121.873900 min 25% 27.000000 26.296250 0.000000 4740.287150 50% 39.000000 30.400000 1.000000 9382.033000 **75**% 51.000000 34.693750 2.000000 16639.912515 64.000000 53.130000 5.000000 63770.428010 max data.info() <class 'pandas.core.frame.DataFrame'> RangeIndex: 1338 entries, 0 to 1337 Data columns (total 7 columns): # Column Non-Null Count Dtype 1338 non-null int64 age object 1338 non-null 1338 non-null float64 bmi children 1338 non-null int64 4 smoker 1338 non-null object 5 region 1338 non-null object 6 charges 1338 non-null float64 dtypes: float64(2), int64(2), object(3) memory usage: 73.3+ KB ax = sns.boxplot(data = data['charges'], orient = 'Vertical') ax.set ylabel('Charges') ax.set\_title('Charges Boxplot') Out[7]: Text(0.5, 1.0, 'Charges Boxplot') Charges Boxplot 60000 50000 40000 30000 20000 10000 At first glance, it seems that the charges column has many outliers, before dealing with them we'll check the charges by smoking status ax = sns.boxplot(x='smoker', y='charges', data=data) ax.set\_title('Charges Boxplot by smoking status') Text(0.5, 1.0, 'Charges Boxplot by smoking status') Charges Boxplot by smoking status 60000 50000 40000 30000 20000 10000 0 smoker Now that huge amount of outliers seems be justified, since smokers are subjects to greater risk, which should be reflected in the amount of charges billed to them by the insurance firm. In [9]: ax = sns.boxplot(x='smoker', y='charges', data=data) ax.set title('Charges Boxplot by smoking status') Out[9]: Text(0.5, 1.0, 'Charges Boxplot by smoking status') Charges Boxplot by smoking status 60000 50000 40000 30000 20000 10000 0 smoker Let's explore our data in more detail. ax =sns.countplot(x = 'children', hue='sex', data=data) for rect in ax.patches: ax.text (rect.get\_x() + rect.get\_width() / 2,rect.get\_height() + 0.75,rect.get\_height(),horizontalalignment ax.set\_title('Number of children by gender') Out[10]: Text(0.5, 1.0, 'Number of children by gender') Number of children by gender 289 285 300 female 250 200 150 100 77 50 ax =sns.countplot(x = 'sex', hue='smoker', data=data) for rect in ax.patches: ax.text (rect.get x() + rect.get width() / 2, rect.get height() + 0.75, rect.get height(), horizontalalignment ax.set title('Smoking status by gender') Out[11]: Text(0.5, 1.0, 'Smoking status by gender') Smoking status by gender smoker 517 500 no 400 300 200 159 100 0 female male ax =sns.countplot(x = 'region', data=data) for rect in ax.patches: ax.text (rect.get\_x() + rect.get\_width() / 2,rect.get\_height() + 0.75,rect.get\_height(),horizontalalignment ax.set\_title('Insurance by region') Out[12]: Text(0.5, 1.0, 'Insurance by region') Insurance by region 350 325 325 324 300 250 200 150 100 50 region ax = sns.histplot(data['age'], kde = False) for rect in ax.patches: ax.text (rect.get x() + rect.get width() / 2, rect.get height() + 0.75, rect.get height(), horizontalalignment ax.set title('Insurance destribution by age') Out[13]: Text(0.5, 1.0, 'Insurance destribution by age') Insurance destribution by age 200 194 175 150 116 115 108 125 106 101 100 91 75 50 25 0 40 In [14]: plt.subplots(figsize = (10,7))ax = sns.histplot(data['charges'], kde = False) for rect in ax.patches: ax.text (rect.get\_x() + rect.get\_width() / 2,rect.get\_height() + 0.75,rect.get\_height(),horizontalalignment ax.set\_title('Charges destribution') Out[14]: Text(0.5, 1.0, 'Charges destribution') Charges destribution 200 175 150 140139 125 100 75 50 25 10000 50000 20000 30000 40000 60000 from scipy.stats import shapiro stat, p = shapiro(data['charges']) print('stat=%.3f, p=%.3f' % (stat, p)) **if** p > 0.05: print('Probably Gaussian') print('Probably not Gaussian') stat=0.815, p=0.000Probably not Gaussian plt.subplots(figsize = (10,7)) ax = sns.histplot(np.log(data['charges']), kde = False) for rect in ax.patches: ax.text (rect.get x() + rect.get width() / 2,rect.get height()+ 0.75,rect.get height(),horizontalalignment ax.set title('Charges destribution') Out[16]: Text(0.5, 1.0, 'Charges destribution') 175 150 140 125 118 95 75 75 50 40 25 20 0 7.0 7.5 8.0 8.5 9.0 10.0 10.5 11.0 charges from scipy.stats import shapiro stat, p = shapiro(np.log(data['charges'])) print('stat=%.3f, p=%.3f' % (stat, p)) **if** p > 0.05: print('Probably Gaussian') else: print('Probably not Gaussian') stat=0.983, p=0.000Probably not Gaussian **Feature Engineering** In this section we will: Label some features and prepare them for modeling. • Discretize the BMI variable and uncover some insights about it. Transform the BMI distribution to Gaussian. In [18]: data.head() bmi children smoker region charges 19 female 27.900 yes southwest 16884.92400 male 33.770 1725.55230 southeast 2 28 male 33.000 southeast 4449.46200 22.705 21984.47061 33 male northwest male 28.880 northwest 3866.85520 In [19]: data['smoker'].replace(to replace=['no','yes'], value=[0,1],inplace=True) Out[19]: bmi children smoker region charges age 19 female 27.900 1 southwest 16884.92400 male 33.770 1725.55230 southeast 2 28 male 33.000 southeast 4449.46200 33 male 22.705 0 northwest 21984.47061 32 male 28.880 0 northwest 3866.85520 According to the CDC (Centers of Desease Control and Prevention), adults BMI is related to their weight status as follows: BMI **Weight Status** Below 18.5 Underweight 18.5 - 24.9**Healthy Weight** 25.0 - 29.9Overweight 30.0 and Above Obesity from IPython.display import HTML HTML('<html><head><style>h3 {text-align: center;} <div style="text-align:center;"> </style></head><body> <h3>BN **BMI For Adults Widget Body Mass Index (BMI) Calculator for Adults** Calculator What is BMI? **Calculate Your BMI** English | Metric **Height:** 0 | feet | 0 | inch(es) Weight: 0 pounds (8 ounces = .5 pounds)Calculate **Grab This Widget** data.loc[(data['bmi'] < 18.5 ), 'bmi group'] = 'Underweight'</pre> data.loc[(data['bmi'] >= 18.5)&(data['bmi'] < 25), 'bmi group' ] = 'Healthy Weight' data.loc[(data['bmi'] >= 25)&(data['bmi'] < 30), 'bmi group'] = 'Overweight'</pre> data.loc[(data['bmi'] >= 30), 'bmi group'] = 'Obesity' data.head() bmi children smoker sex region charges bmi\_group 16884.92400 female 27.900 southwest Overweight Obesity 18 male 33.770 southeast 1725.55230 2 southeast 28 male 33.000 3 4449.46200 Obesity 3 33 male 22.705 northwest 21984.47061 Healthy Weight Overweight male 28.880 0 3866.85520 32 northwest ax = sns.boxplot(x='bmi\_group', y='charges', data=data) ax.set title('Charges Boxplot by smoking status') Out[22]: Text(0.5, 1.0, 'Charges Boxplot by smoking status') Charges Boxplot by smoking status 60000 50000 40000 30000 20000 10000 0 Underweight Overweight Obesity Healthy Weight bmi\_group mean = data.groupby('bmi group')['charges'].mean() mean = mean.to frame() mean = mean.sort\_values(by='charges') mean charges bmi\_group Underweight 8852.200585 **Healthy Weight** 10409.337709 **Overweight** 10987.509891 **Obesity** 15552.335469 In [24]: plt.plot(mean, 'o', linestyle='-') Out[24]: [<matplotlib.lines.Line2D at 0x19b61947af0>] 15000 14000 13000 12000 11000 10000 9000 Healthy Weight Overweight Underweight Obesity We can see that the mean of the charges is correlated positively with the BMI categories ax = sns.histplot(data['bmi'], kde = True) ax.set title('BMI distribution') Out[25]: Text(0.5, 1.0, 'BMI distribution') BMI distribution 140 120 100 60 40 20 0 Visually we can see that the dansity function curve is chaped like a bell, is it enough to assume that the distribution is normal? **Let's perform Shapiro-Wilk Normality test Shapiro-Wilk Normality Test** Let's test whether the BMI has a Gaussian distribution. Here are the hypothesis. H0: the sample has a Gaussian distribution. H1: the sample does not have a Gaussian distribution. from scipy.stats import shapiro stat, p = shapiro(data['bmi']) print('stat=%.3f, p=%.3f' % (stat, p)) **if** p > 0.05: print('Probably Gaussian') else: print('Probably not Gaussian') stat=0.994, p=0.000Probably not Gaussian Even though we got a bell shaped density function plot, statistically we can't assume that it has a normal distribution. We will try to transform it by using the square root function (sqrt). bmi = np.array(data['bmi']) sqrt bmi = np.sqrt(data['bmi']) ax = sns.histplot(sqrt bmi, kde = True) ax.set title('Transformed BMI distribution') Out[27]: Text(0.5, 1.0, 'Transformed BMI distribution') Transformed BMI distribution 140 120 100 60 40 **Shapiro-Wilk Normality Test** Let's test whether the transformed BMI has a Gaussian distribution. Here are the hypothesis. H0: the sample has a Gaussian distribution. H1: the sample does not have a Gaussian distribution. from scipy.stats import shapiro stat, p = shapiro(sqrt bmi) print('stat=%.3f, p=%.3f' % (stat, p)) **if** p > 0.05: print('Probably Gaussian') else: print('Probably not Gaussian') stat=0.999, p=0.345Probably Gaussian Now we can assume that the variable is normally distributed. # let's calculate the probability of having a healthy BMI from math import sqrt from scipy.stats import norm sqrt\_mean = sqrt\_bmi.mean() sqrt\_sd = sqrt\_bmi.std() upper limit = norm.cdf((sqrt(24.9) - sqrt mean)/(sqrt sd)) lower\_limit = norm.cdf((sqrt(18.5) - sqrt\_mean)/(sqrt\_sd)) prob = upper\_limit - lower\_limit print('The probability that someone chosen randomly has a healthy BMI is {}%'.format(round(prob\*100, 2))) The probability that someone chosen randomly has a healthy BMI is 15.88% # According to the CDC a person reach obesity when the BMI is 30 or above prob1 = 1-norm.cdf((sqrt(30) - sqrt mean)/(sqrt sd)) print('The probability that someone chosen randomly has obesity is {}%'.format(round(prob1\*100, 2))) The probability that someone chosen randomly has obesity is 52.36%**Hypothesis Testing** Since the bmi variable follows the Normal distribution, it's worth to test: whether the bmi differ by gender. whether the bmi differ by region. data.groupby('sex')['bmi'].mean() Out[31]: sex female 30.377749 30.943129 Name: bmi, dtype: float64 T-Test We would like to test whether there is a difference in the BMI by gender. let's state our hypothesis: H0: the means of the samples are equal. H1: the means of the samples are unequal. female bmi = data.loc[data['sex'] == 'female']['bmi'] male bmi = data.loc[data['sex'] == 'male']['bmi'] # Example of the Student's t-test from scipy.stats import ttest ind stat, p = ttest ind(female bmi, male bmi) print('stat=%.3f, p=%.3f' % (stat, p)) **if** p > 0.05: print('Statistically the means are the same') print('Statistically the means are not the same') stat = -1.697, p = 0.090Statistically the means are the same region\_mean = data.groupby('region')['bmi'].mean() region\_mean.to\_frame() bmi region northeast 29.173503 northwest 29.199785 southeast 33.355989 **southwest** 30.596615 **ANOVA Test** To test whether the means are equal by region we will use the one-way ANOVA test. Let's state our hypothesis. H0: the means of the samples are equal. H1: one or more of the means of the samples are unequal. In [34]: from scipy.stats import f\_oneway northwest\_bmi = data.loc[data['region'] == 'northwest']['bmi'] northeast\_bmi = data.loc[data['region'] == 'northeast']['bmi'] southwest\_bmi = data.loc[data['region'] == 'southwest']['bmi'] southeast\_bmi = data.loc[data['region'] == 'southeast']['bmi'] stat, p = f\_oneway(northwest\_bmi, northeast\_bmi, southwest\_bmi, southeast\_bmi) print('stat=%.3f, p=%.3f' % (stat, p)) **if** p > 0.05: print('Statistically the means are equal') print('Statistically there is at least one mean that is different from the others') stat=39.495, p=0.000Statistically there is at least one mean that is different from the others Results We got insightful visualizations about our data's distribution and variables. We saw how important the outliers from the charges variable, and why it's convenient to keep them since they represent the smoking and obese people who obviously are exposed to greater risk. After transforming the BMI variable, we were able to calculate the probabilities of chosing randomly a healthy person and an obese person which evaluated respectively to 15.88% and 52.36%. Finally, we saw how the BMI means by gender are statistically equivalent, while they differ by region. Next we will try to predict charges using a famous Machine Learning technique which is Regression. **Machine Learning** Now let's check our data again data.head() bmi children smoker bmi\_group age sex region charges southwest 16884.92400 19 female 27.900 0 Overweight male 33.770 Obesity 18 southeast 1725.55230 male 33.000 2 28 3 4449.46200 Obesity southeast northwest 21984.47061 male 22.705 Healthy Weight 3 33 0 32 male 28.880 northwest 3866.85520 Overweight Now we ended up with an interesting form, since the BMI variable is repeated in a descrete form within the BMI\_group variable. We avoid using both variables in the same model and compare which one comes with better results. The **first model** (model 1) will use the bmi variable and the **second model** (model 2) will use the bmi\_group variable. For both data sets, we will fit and compare the results of the following models: • Linear regression without scaling the variables Linear regression scaling the variables Ridge regression Lasso regression target = data['charges'] Model 1 data1 = data.drop(['bmi\_group','charges'], axis = 1) data1.head() age sex bmi children smoker region 19 female 27.900 0 1 southwest 18 male 33.770 southeast 28 male 33.000 southeast 33 22.705 northwest male 32 male 28.880 0 northwest data.dtypes

] : _	<pre>bmi_group object dtype: object  # we will apply one hot encoding to categorical variables data1 = pd.get_dummies(data1, drop_first=True, columns=['sex', 'region']) data1.head()</pre>
	age         bmi         children         smoker         sex_male         region_northwest         region_southeast         region_southwest           0         19         27.900         0         1         0         0         1           1         18         33.770         1         0         1         0         1         0           2         28         33.000         3         0         1         0         1         0           3         33         22.705         0         0         1         1         0         0           4         32         28.880         0         0         1         1         0         0           from sklearn.preprocessing import StandardScaler, PolynomialFeatures           from sklearn.model_selection import KFold, cross_val_predict           from sklearn.linear model import LinearRegression, Lasso, Ridge
: [	<pre>from sklearn.metrics import r2_score from sklearn.pipeline import Pipeline  kf = KFold(shuffle=True, random_state=555, n_splits=5)  We will use grid search cv for our parameter tuning and extract meaningful results for model comparison.</pre>
•	<pre>("linear_regression", LinearRegression())])  params = {     'polynomial_featuresdegree': [1, 2, 3] }  grid = GridSearchCV(estimator, params, cv=kf)  grid.fit(datal, target)  GridSearchCV(cv=KFold(n_splits=5, random_state=555, shuffle=True),</pre>
	<pre>best_score = grid.best_score_ alpha_degree = (None, grid.best_params_['polynomial_featuresdegree']) grid.best_score_, grid.best_params_  (0.80498278137892, {'polynomial_featuresdegree': 2})  y_predict = grid.predict(data1) R2_score = r2_score(target, y_predict) R2_score  0.8263982356893842</pre>
	<pre>lin_reg = [best_score, alpha_degree, R2_score]  grid.best_estimatornamed_steps['linear_regression'].coef_  array([-2.42776896e+14, -9.48522278e+00, 1.21766541e+03, 2.14934236e+03, -1.16106515e+04, 1.57063411e+01, -4.35292528e+02, 1.09503035e+03, -1.41419043e+03, 3.42281159e+00, -8.41032041e-01, -1.31590000e+01, 6.36059624e+00, -3.24930299e+00, 2.40093767e+01, 3.3266655e+01, 5.25883999e+01, -1.59068005e+01, -5.79741546e+01, 1.48347083e+03, -4.68431718e+01, -6.06927929e+01, -1.76246304e+02, -2.97008648e+00, 1.98220737e+02, 3.18172890e+02, -1.64280590e+01, 8.58725959e+02, 9.74944317e+02, -1.20970446e+03, -1.16106515e+04, 2.72831080e+03, 4.48508446e+02, 4.85329068e+02, -3.75695394e+03, 1.57063411e+01, -5.25585778e+01, -2.77421424e+02, 2.71545681e+03, -4.35292528e+02, 0.00000000e+00, 0.00000000e+00, 1.09503035e+03, 0.00000000e+00, -1.41419043e+03])</pre>
	<pre>Model 1 Linear Regression scaled variables  from sklearn.model_selection import GridSearchCV  # Same estimator as before estimator = Pipeline([("scaler", StandardScaler()),</pre>
	<pre>grid.fit(data1, target)  GridSearchCV(cv=KFold(n_splits=5, random_state=555, shuffle=True),</pre>
	<pre>y_predict = grid.predict(data1) R2_score = r2_score(target, y_predict) R2_score  0.847724191182748  lin_reg_scaled = [best_score, alpha_degree, R2_score]  grid.best_estimatornamed_steps['linear_regression'].coef_  array([ 8.91874812e+11,</pre>
	-2.08290804e+15, 8.01572493e+02, 5.02655252e+01, -9.12860186e+01, -2.89492861e+01, 1.23563197e+02, 1.04043340e+02, 3.30772521e+02, 2.74900267e+02, -2.35873652e+02, 2.00301188e+01, 3.64025048e+03, 1.86750312e+01, -3.69470407e+01, -3.82637836e+02, -2.05459569e+02, -1.37265727e+02, -1.81197729e+02, -1.36556818e+02, 1.24082500e+02, -1.28228573e+02, -2.07371051e+02, -2.66757659e+15, -1.17041019e+01, -7.48406225e+01, -2.10032229e+02, 1.58595345e+02, 2.22164568e+15, 8.86137498e+01, 1.59852476e+02, 6.70488298e+01, 1.13560489e+16, -2.99947801e+16, 4.09218254e+16, 3.38851846e+16, 3.45233683e+16, 3.86691196e+16])  Model 1 Ridge Regression  N.B: While alpha ranges used in models such as Lasso and Ridge regression may look inconsistent, we will just overlook this those ranges contain the parameters that produce good results.
	<pre>from sklearn.model_selection import GridSearchCV  # Same estimator as before estimator = Pipeline([("scaler", StandardScaler()),</pre>
	<pre>estimator=Pipeline(steps=[('scaler', StandardScaler()),</pre>
	<pre>grid.best_score_, grid.best_params_  (0.8315159135685688,</pre>
	grid.best_estimatornamed_steps['ridge_regression'].coef_  array([
	<pre>from sklearn.model_selection import GridSearchCV  # Same estimator as before estimator = Pipeline([("scaler", StandardScaler()),</pre>
	GridSearchCV(cv=KFold(n_splits=5, random_state=555, shuffle=True),
	<pre>best_score = grid.best_score_ alpha_degree = (grid.best_params_['lasso_regressionalpha'], grid.best_params_['polynomial_featuresorid.best_score_, grid.best_params_ grid.best_score_, grid.best_params_  (0.8325039849660897, {'lasso_regressionalpha': 73.7410737070725,     'polynomial_featuresdegree': 2})  y_predict = grid.predict(datal) R2_score = r2_score(target, y_predict) R2_score  0.8462302941184598</pre>
	<pre>lasso_reg = [best_score, alpha_degree, R2_score]  grid.best_estimatornamed_steps['lasso_regression'].coef_  array([ 0.00000000e+00,</pre>
	-7.41987842e+01])  lin_reg, lin_reg_scaled, ridge_reg, lasso_reg  ([0.80498278137892, (None, 2), 0.8263982356893842], [0.829965126687363, (None, 2), 0.847724191182748], [0.8315159135685688, (31.62417284363078, 2), 0.8475929725707861], [0.8325039849660897, (73.7410737070725, 2), 0.8462302941184598])  ind = pd.MultiIndex.from_tuples([('Model_1','lin_reg'),('Model_1','lin_reg_scaled'),('Model_1','ridge_reg_olded'), ('Model_1 = pd.DataFrame([lin_reg, lin_reg_scaled, ridge_reg, lasso_reg], index=ind, columns=col)  Model_1.sort_values(by='R2_score', ascending=False)
•	Model_1   lin_reg_scaled   0.829965   (None, 2)   0.847724
	data2 = data.drop(['bmi','charges'], axis = 1) data2.head()    ge
	data2 = pd.get_dummies(data2, drop_first=True, columns=['sex', 'region', 'bmi_group'])  age children smoker sex_male region_northwest region_southeast region_southwest bmi_group_Obesity bmi_group_Overweight  1 18 1 0 1 0 0 1 0 1 0 1  2 28 3 0 1 0 1 1 0 1  3 33 0 0 1 1 0 0 0 0 0  4 32 0 0 1 1 1 0 0 0 0
	<pre>Model 2 Linear regression without variables scaling.  from sklearn.model_selection import GridSearchCV  # Same estimator as before estimator = Pipeline([("polynomial_features", PolynomialFeatures()),</pre>
	<pre>GridSearchCV(cv=KFold(n_splits=5, random_state=555, shuffle=True),</pre>
	R2_score = r2_score(target, y_predict) R2_score  0.869463179501786  lin_reg = [best_score, alpha_degree, R2_score]  grid.best_estimatornamed_steps['linear_regression'].coef_  array([ 6.04903102e-10, -7.54907130e+01, 1.18248353e+03, 5.70261577e+03, -5.78203604e+02, -7.64728279e+02, -6.22084894e+02, -1.27751394e+03, -2.59795927e+02, -9.30999255e+02, -2.32675697e+03, 3.82629595e+00, -2.60248984e+00, -3.12519678e+00, 5.79295742e+00, 9.33377329e+00, 3.77650434e+01, 2.45134252e+01, 1.44011808e+01, 3.04124959e+01,
	5.42459183e+01, -1.25715232e+02, -3.66697107e+02, -6.07464500e+01, 8.28529075e+01, -2.98595600e+02, -3.92325602e+02, 3.43973058e+02, 3.34448761e+02, 6.54277531e+01, 5.70261577e+03, 5.58313536e+01, 9.60416510e+02, 1.41244811e+03, 7.53600840e+02, 2.14151939e+04, 2.75497128e+03, 3.27647654e+03, -5.78203604e+02, 6.55090995e+02, 1.52078133e+02, 4.59355712e+02, -1.30196198e+01, 6.13371513e+02, -1.10584789e+03, -7.64728279e+02, 2.84217094e-13, -9.09494702e-13, 7.53482904e+00, 1.49962931e+02, 3.97784960e+03, -6.22084894e+02, 0.00000000e+00, -1.06905572e+03, -1.48743212e+03, 0.00000000e+00, -1.27751394e+03, -1.66873463e+02, 1.31927023e+03, 2.03486596e+03, -2.59795927e+02, 0.0000000e+00, 0.0000000e+00, -9.30999255e+02, 0.0000000e+00, -2.32675697e+03])  Model 2 Linear Regression scaled variables  from sklearn.model_selection import GridSearchCV
	<pre># Same estimator as before estimator = Pipeline([("scaler", StandardScaler()),</pre>
	<pre>PolynomialFeatures()),     ('linear_regression',</pre>
	<pre>0.8692157125924178  lin_reg_scaled = [best_score, alpha_degree, R2_score]  grid.best_estimatornamed_steps['linear_regression'].coef_  array([-8.75810642e+11, 3.63744781e+03, 1.03190508e+03, 2.15914916e+16,</pre>
	<pre>4.98000000e+02, 1.73750000e+02, 5.11601122e+15, 1.40000000e+02, 2.80000000e+01, 1.01812500e+02, -4.00000000e+00, 1.36000000e+02, -6.05000000e+01, 8.42702824e+15, -5.84734552e+15, 1.98876101e+15, 1.3000000e+01, 4.3000000e+01, 2.08250000e+02, 2.20433851e+15, -2.95215421e+15, -2.33500000e+02, -2.92000000e+02, 3.26645061e+15, 1.47086908e+16, -3.60000000e+01, 2.58000000e+02, 1.17750000e+02, 2.35257926e+14, 5.72869903e+15, -2.66916529e+16, 9.96743827e+15, -1.37807311e+16, -6.69927399e+15])</pre> <pre></pre>
	<pre>params = {     'polynomial_featuresdegree': [1, 2, 3],     'ridge_regressionalpha': np.geomspace(4, 40, 50) }  grid = GridSearchCV(estimator, params, cv=kf)  grid.fit(data2, target)  GridSearchCV(cv=KFold(n_splits=5, random_state=555, shuffle=True),</pre>
	4.82717056, 5.05942087, 5.30284546, 5.557 10.23819169, 10.73078318, 11.24707479, 11.78820681, 12.35537439, 12.94983017, 13.57288709, 14.22592122, 14.91037488, 15.62775975, 16.37966025, 17.16773704, 17.99373068, 18.85946545, 19.76685345, 20.71789872, 21.71470176, 22.75946412, 23.85449327, 25.0022077, 26.20514227, 27.4659538, 28.78742692, 30.17248025, 31.62417284, 33.14571091, 34.74045495, 36.41192712, 38.16381905, 40. ])})  best_score = grid.best_score_ alpha_degree = (grid.best_params_['ridge_regression_alpha'], grid.best_params_['polynomial_features_cgrid.best_score_, grid.best_params_  (0.8522931183263784, {'polynomial_features_degree': 2,     'ridge_regression_alpha': 33.145710914187376})
	<pre>y_predict = grid.predict(data2) R2_score = r2_score(target, y_predict) R2_score  0.8691647800185991  ridge_reg = [best_score, alpha_degree, R2_score]  grid.best_estimatornamed_steps['ridge_regression'].coef_  array([ 0.00000000e+00,  3.55305159e+03,  9.52832090e+02,  3.03869806e+03,</pre>
	-6.38515607e+01, -1.20985580e+01, 3.15551650e+01, 4.47188436e+01, 2.17089492e+02, 1.33797800e+02, 9.43134171e+01, 1.77587528e+02, 7.29364529e+01, -1.66527646e+02, -1.73591115e+02, -3.20365951e+01, 5.22729404e+01, -1.41372450e+02, -1.92103528e+02, 1.93449320e+02, 1.67240480e+02, 2.43282872e+01, 4.44599064e+03, 4.25829442e+01, 1.62901809e+02, 2.73960069e+02, 1.40528771e+02, 4.08980337e+03, 3.41378478e+02, 1.06412670e+02, 5.01748041e+00, 1.25539726e+02, 1.79191822e+01, 8.94024682e+01, -5.65456163e-01, 1.34209160e+02, -8.00317389e+01, -3.59389534e+01, 6.78280633e+01, 1.08075960e+02, -7.36648160e+00, 2.53359111e+01, 1.98897319e+02, -8.95391689e+01, 1.47826486e+02, -2.05279844e+02, -2.66289842e+02, 1.36876593e+01, -1.92849540e+02, -5.07266211e+01, 2.36481332e+02, 7.97069533e+01, -1.89105854e+02, -7.36190834e+02, -1.99667860e+02, -2.84088538e+02, 4.05234144e+01, -3.82032436e+01])  Model 2 Lasso regression
	<pre>from sklearn.model_selection import GridSearchCV  # Same estimator as before estimator = Pipeline([("scaler", StandardScaler()),</pre>
	estimator=Pipeline(steps=[('scaler', StandardScaler()),
	<pre>best_score = grid.best_score_ alpha_degree = (grid.best_params_['lasso_regressionalpha'], grid.best_params_['polynomial_features_orgid.best_score_, grid.best_params_  (0.855033666178441, {'lasso_regressionalpha': 97.64349540069045,     'polynomial_featuresdegree': 2})  y_predict = grid.predict(data2) R2_score = r2_score(target, y_predict) R2_score  0.8663602384916383  lasso_reg = [best_score, alpha_degree, R2_score]</pre>
	grid.best_estimatornamed_steps['lasso_regression'].coef_  array([
	-0. , -0. , 0. , 0. , 0. , 0. , 0. , 111.68362523, -0. , -42.28712766, 0. , -197.58539248, -0. , 157.44565966, 0. , -197.58539248, -0. , 157.44565966, 0. , 0. , -0. , -50.63218696, 0. , 0. , -0. , -18.91177808])  lin_reg, lin_reg_scaled, ridge_reg, lasso_reg  ([0.8517098686270737, (None, 2), 0.869463179501786], [0.8514920934611669, (None, 2), 0.8692157125924178], [0.8522931183263784, (33.145710914187376, 2), 0.8691647800185991], [0.855033666178441, (97.64349540069045, 2), 0.8663602384916383])  ind = pd.MultiIndex.from_tuples([('Model_2','lin_reg'), ('Model_2','lin_reg_scaled'), ('Model_2','ridge_scaled'), ('Mod
•	lin_reg_scaled, ridge_reg, lasso_reg], index=ind, columns=col)   Model_2.sort_values(by='R2_score', ascending=False)    Best_score   (Alpha, Polynomial_degree)   R2_score     Model_2   lin_reg   0.851710   (None, 2)   0.869463     lin_reg_scaled   0.851492   (None, 2)   0.869216     ridge_reg   0.852293   (33.145710914187376, 2)   0.869165     lasso_reg   0.855034   (97.64349540069045, 2)   0.866360    We see that the models using the discretized BMI variable achieved an accuracy close to 87%
	Results   pd.concat ([Model_1], Model_2]   Results   Sort_values (by='R2_score', ascending=False)
	<pre>lin_reg  0.804983</pre>
	### def_importances  The total number of features is: 66  The total number of non null features is: 22     Feature   Coef
	6 region_southwest bmi_group_Overweight 157.445660 7 smoker bmi_group_Overweight 156.346485 8 region_southeast region_southwest 111.683625 9 sex_male bmi_group_Overweight 92.476865 10 smoker region_southeast 85.106855 11 region_northwest bmi_group_Underweight 53.990471 12 children region_northwest 44.757166 13 age region_southeast 42.263804 14 children^2 -12.239369
	15 bmi_group_Underweight^2 -18.911778
	16 region_southeast bmi_group_Overweight -42.287128 17 children region_southwest -43.045228 18 children smoker -45.489261 19 bmi_group_Obesity bmi_group_Overweight -50.632187 20 sex_male -138.626583 21 region_southwest^2 -197.585392  Even thought the LASSO model is ranked fourth, it's still one of the best models due to its interpretability. By this we mean is the most charming because it reduces the magnitude of the coefficients and also nullifies many features in our case 66 - 26 features were nullified.
ii ii ff \\ \\ \\ \\ \\ \\ \\ \\ \\ \\ \\ \\	16 region_southeast bmi_group_Overweight -42.287128  17 children region_southwest -43.045228  18 children smoker -45.489261  19 bmi_group_Obesity bmi_group_Overweight -50.632187  20 sex_male -138.626583  21 region_southwest^2 -197.585392  Even thought the LASSO model is ranked fourth, it's still one of the best models due to its interpretability. By this we mean is the most charming because it reduces the magnitude of the coefficients and also nullifies many features in our case 66 - 20.000 children region_southwest -42.287128  -42.287128  -43.045228  -45.489261  -50.632187  -50.632187  -60.632187