

```
In [1]: warnings > ignore warnings
import pandas as pd # data manipulation
numpy = np # linear algebra
matplotlib.pyplot as plt # plots
import seaborn sns # plots

In [2]: warnings.filterwarnings('ignore')

In [3]: df=pd.read_csv('data/bike-sharing-demand.csv')
df=df.drop(['Unnamed: 0'],axis='columns')

In [4]: df.head()
```

	age	sex	bmi	children	smoker	region	charges
0	18	male	33.77	1	no	southeast	1725.5523
1	28	male	31.00	3	no	southeast	4449.4620
2	25	male	26.22	0	no	northeast	2721.3208
3	13	male	34.42	0	no	southeast	2083.4020
4	19	male	24.60	1	no	southeast	1357.2370

```
In [5]: dataset = df.drop(['charges'],axis='columns') df.charges.values
X = dataset'
```

## Split data

```
In [6]: sklearn.model_selection.RandomState(1) train_test_split
X_train X_test Y_train Y_test=train_test_split(X
                                                Y
                                                test_size=0.3
                                                random_state=1)
```

## Data preprocessing

- **MinMaxScaler** consists of adjusting the data on a scale from 0 to 1. With the aim that the variables are comparable to each other.
- **OneHotEncoding** it is used for qualitative variables. For example the geographical location or the color of a car. It consists of creating several fictitious variables according to the number of categorical variables. Where a 1 is indicated if said observation belongs to the class. It's a similar thing with truth tables if you're familiar with programming.
- **PolynomialFeatures** it consists of raising to a certain power. To break the linearity of the data. Since sometimes a linear regression is not enough to solve the problem.

```
In [7]: sklearn.preprocessing.MinMaxScaler OneHotEncoder PolynomialFeatures

In [9]: poly=PolynomialFeatures(degree=3) # degree polynomial
rescale=MinMaxScaler()
poly.transform(x)

x=pd.get_dummies(x.drop_first(['sex','smoker','region'])) # dummy variables
x[['age','bmi','children']] = rescale.fit_transform(x[['age','bmi','children']]) # rescale
x=poly.fit_transform(x) # polynomial transform

In [10]: X_train_poly=poly.transform(X_train)
X_test_poly=poly.transform(X_test)
```

## Creation of the polynomial model

```
In [11]: sklearn.linear_model.LinearRegression

In [12]: poly.fit(X_train_poly,Y_train) # we give the training data

Out[12]: LinearRegression

In [13]: lin=LinearRegression

In [14]: lin.fit(X_train_poly,Y_train) # we give the training data

Out[14]: LinearRegression

In [15]: lin.score(X_train_poly,Y_train) # evaluate with training data

Out[15]: 0.9999999999999999

In [16]: lin.score(X_test_poly,Y_test) # evaluate with testing data

Out[16]: 0.9999999999999999
```

## Data transform

```
In [17]: sklearn.compose.ColumnTransformer
sklearn.pipeline.Pipeline

In [18]: tf_columns=make_column_transformer((MinMaxScaler(),OneHotEncoder(drop='if_binary')),# rescale
                                          (OneHotEncoder(drop='if_binary'),'sex','smoker')) # OneHotEncoding
```

## Mean square error

Measures the average error between the original and predicted values.

```
In [19]: sklearn.metrics.mean_squared_error
```

## Gradient Boosting

It is an algorithm that uses such complex algorithms. That is perfected according to the learning rate assigned by the user. In such a way that each tree will become better than the previous one. And the predicted value is going to average the predictions of the weakest algorithms.

```
In [20]: sklearn.ensemble.GradientBoostingRegressor

In [21]: evaluate_max_depth_lr

estimator_list=[]
mse_train_list=[]
mse_test_list=[]

estimators=np.arange(100,1000,steps=100)
n_estimators=len(estimators)

model=GradientBoostingRegressor(max_depth=max_depth,n_estimators=estimator,learning_rate=lr,random_state=1)
model=Pipeline(['transformer',tf_columns,('model',model)])
model.fit(X_train,Y_train)
pred_train=model.predict(X_train)
pred_test=model.predict(X_test)

mse_train=mean_squared_error(Y_train,pred_train)
mse_test=mean_squared_error(Y_test,pred_test)

estimator_list.append(estimator)
mse_train_list.append(mse_train)
mse_test_list.append(mse_test)

return estimator_list,mse_train_list,mse_test_list

In [22]: dataframe_evaluate_trees_max_depth_lr

n_trees,mse_test,mse_train=evaluate_max_depth_max_depth_lr(lr)

df_evaluate=pd.DataFrame({'n_trees':n_trees,'mse_test':mse_test,'mse_train':mse_train})

df_evaluate
```

first\_evaluate=dataframe\_evaluate\_trees\_max\_depth\_lr=0.9999999999999999

second\_evaluate=dataframe\_evaluate\_trees\_max\_depth\_lr=0.9999999999999999

## Ideal number of estimators


```
In [23]: sns.set_style(style='whitegrid')
fig,ax=plt.subplots(figsize=(10,10))

ax.set_title('Max depth 2')
ax.plot(first_evaluate['n_trees'],first_evaluate['mse_train'],label='Train MSE')
ax.plot(first_evaluate['n_trees'],first_evaluate['mse_test'],label='Test MSE')
ax.set_xlabel('n_estimators')
ax.set_ylabel('mse')

ax_1=plt.title('Max depth 3')
ax_1.plot(second_evaluate['n_trees'],second_evaluate['mse_train'],label='Train MSE')
ax_1.plot(second_evaluate['n_trees'],second_evaluate['mse_test'],label='Test MSE')
ax_1.set_xlabel('n_estimators')
ax_1.set_ylabel('mse')

ax.legend()
ax_1.legend()

plt.show()
```



With a maximum depth of 2 there is less overfitting. The ideal number of estimators for this model ranges from 600 to 650.

```
In [24]: first_evaluate.query('n_estimators<600 and n_estimators>400')

Out[24]:
```

	n_trees	mse_test	mse_train
255	610	2.572382e+06	2.035120e+06
256	612	2.569720e+06	2.030072e+06
257	614	2.567132e+06	2.027725e+06
258	616	2.565719e+06	2.026049e+06
259	618	2.565256e+06	2.017800e+06
260	620	2.558398e+06	2.032292e+06
261	622	2.555005e+06	2.003796e+06
262	624	2.550332e+06	2.004480e+06
263	626	2.547820e+06	2.000770e+06
264	628	2.545979e+06	1.995625e+06
265	630	2.539644e+06	1.989548e+06
266	632	2.535433e+06	1.984598e+06
267	634	2.534456e+06	1.982152e+06
268	636	2.531220e+06	1.979840e+06
269	638	2.528137e+06	1.975120e+06
270	640	2.525122e+06	1.969584e+06
271	642	2.522637e+06	1.965623e+06
272	644	2.519813e+06	1.963009e+06
273	646	2.518178e+06	1.959314e+06
274	648	2.516109e+06	1.957123e+06
275	650	2.509953e+06	1.949535e+06

Starting from the estimator number 630 there is no longer an improvement for the test data.

```
In [27]: gbr=GradientBoostingRegressor(max_depth=
learning_rate=0.1,
n_estimators=650,
random_state=1)

In [42]: gbr_pipeline=Pipeline(['transformer',tf_columns
                               ('gbr',gbr)])

In [43]: gbr_pipeline.fit(X_train,Y_train)
```

```
Out[43]: GradientBoostingRegressor(max_depth=61, max_depth2=
n_estimators=650, random_state=1)

In [44]: gbr_pipeline.score(X_train,Y_train) # evaluate with training data

Out[44]: 0.9999999999999999

In [45]: gbr_pipeline.score(X_test,Y_test) # evaluate with test data

Out[45]: 0.9999999999999999
```

## XGBOOST

It is similar to adaboost. With the difference that xgboost can be used with GPU. Which training will be faster.

```
In [52]: xgboost.XGBRegressor
```

## Selection number of estimators XGBOOST

```
In [53]: evaluate_max_depth_lr

estimator_list=[]
mse_train_list=[]
mse_test_list=[]

estimators=np.arange(100,1000,steps=100)
n_estimators=len(estimators)

model=XGBRegressor(max_depth=max_depth,n_estimators=estimator,learning_rate=lr,random_state=1,verbosity=1)
model=Pipeline(['transformer',tf_columns,('model',model)])
model.fit(X_train,Y_train)
pred_train=model.predict(X_train)
pred_test=model.predict(X_test)

mse_train=mean_squared_error(Y_train,pred_train)
mse_test=mean_squared_error(Y_test,pred_test)

estimator_list.append(estimator)
mse_train_list.append(mse_train)
mse_test_list.append(mse_test)

return estimator_list,mse_train_list,mse_test_list

In [54]: first_evaluate=dataframe_evaluate_trees_max_depth_lr=0.9999999999999999

In [55]: second_evaluate=dataframe_evaluate_trees_max_depth_lr=0.9999999999999999

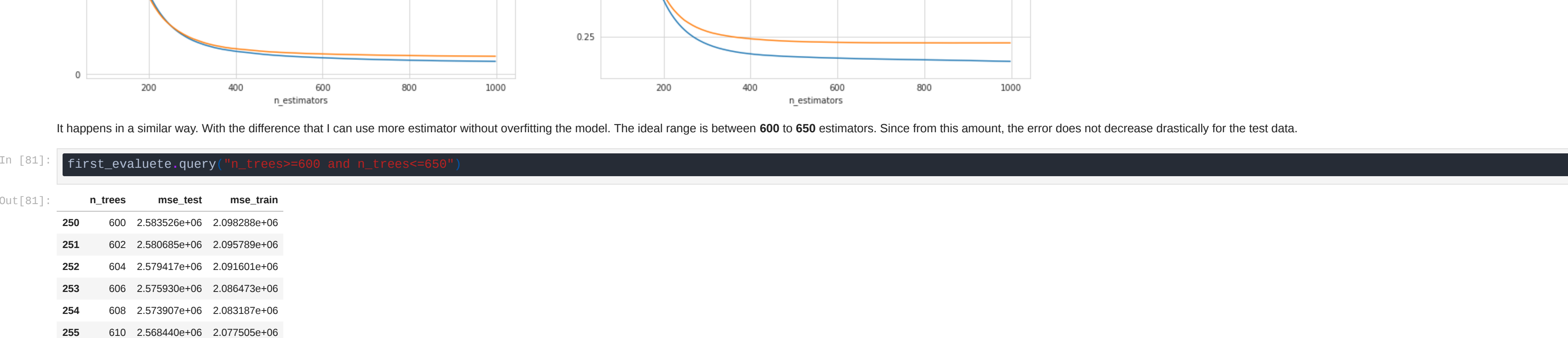
In [56]: sns.set_style(style='whitegrid')
fig,ax=plt.subplots(figsize=(10,10))

ax.set_title('Max depth 2')
ax.plot(first_evaluate['n_trees'],first_evaluate['mse_train'],label='Train MSE')
ax.plot(first_evaluate['n_trees'],first_evaluate['mse_test'],label='Test MSE')
ax.set_xlabel('n_estimators')
ax.set_ylabel('mse')

ax_1=plt.title('Max depth 3')
ax_1.plot(second_evaluate['n_trees'],second_evaluate['mse_train'],label='Train MSE')
ax_1.plot(second_evaluate['n_trees'],second_evaluate['mse_test'],label='Test MSE')
ax_1.set_xlabel('n_estimators')
ax_1.set_ylabel('mse')

ax.legend()
ax_1.legend()

plt.show()
```



It happens in a similar way. With the difference that I can use more estimator without overfitting the model. The ideal range is between 600 to 650 estimators. Since from this amount, the error does not decrease drastically for the test data.

```
In [61]: first_evaluate.query('n_estimators<600 and n_estimators>400')

Out[61]:
```

	n_trees	mse_test	mse_train
260	610	2.558270e+06	2.030200e+06
261	612	2.555005e+06	2.003796e+06
262	614	2.579417e+06	2.019501e+06
263	616	2.571933e+06	2.008473e+06
264	618	2.573907e+06	2.003187e+06
265	620	2.568440e+06	2.017505e+06
266	622	2.560714e+06	2.012123e+06
267	624	2.565027e+06	2.006400e+06
268	626	2.562223e+06	2.001395e+06
269	628	2.558023e+06	2.005995e+06
270	630	2.541228e+06	2.009212e+06
271	632	2.537486e+06	2.002770e+06
272	634	2.540420e+06	2.000200e+06
273	636	2.532498e+06	2.002555e+06
274	638	2.530372e+06	2.003407e+06
275	640	2.528480e+06	2.004488e+06
276	642	2.524949e+06	2.002223e+06
277	644	2.522270e+06	2.001020e+06
278	646	2.518522e+06	2.005413e+06
279	648	2.515875e+06	2.002244e+06
280	650	2.514028e+06	1.998935e+06

We seek 650. Since from this estimator there is no significant improvement.

```
In [40]: xgb_reg=XGBRegressor(max_depth=
learning_rate=0.1,
n_estimators=650,
verbosity=1,
random_state=1)

xgb_pipeline=Pipeline(['transformer',tf_columns,('xgb',xgb_reg)])

In [47]: xgb_pipeline.fit(X_train,Y_train)
```

```
Out[47]: GradientBoostingRegressor(max_depth=61, max_depth2=
n_estimators=650, random_state=1)

In [48]: xgb_pipeline.score(X_train,Y_train)

Out[48]: 0.9999999999999999

In [49]: xgb_pipeline.score(X_test,Y_test)

Out[49]: 0.9999999999999999

In [59]: sklearn.metrics.r2_score
sklearn.model_selection.cross_val_score
```

```
In [73]: Evaluate

__init__(self,model,X_data,y_true):
self.X_data=X_data
self.y_true=y_true

self.model=model
self.predict=self.model.predict(self.X_data)

mse self
mse=mean_squared_error(self.y_true,self.predict)

r2 self
r2=r2_score(self.y_true,self.predict)

cv_score self
cv_score=cross_val_score(self.model,self.X_data,self.y_true,cv=5).mean()
```

## MSE

Measures the average error between the original and predicted values.

```
In [55]: mse_poly_train=Evaluate lin X_train_poly Y_train .mse()
mse_poly_test=Evaluate lin X_test_poly Y_test .mse()

mse_gbr_train=Evaluate gbr_pipeline X_train Y_train .mse()
mse_gbr_test=Evaluate gbr_pipeline X_test Y_test .mse()

mse_xgb_train=Evaluate xgb_pipeline X_train Y_train .mse()
mse_xgb_test=Evaluate xgb_pipeline X_test Y_test .mse()
```

## R<sup>2</sup>

It measures the degree of fit between the original value and the predictions. The closer it is to 1, the closer the original and predicted values will be.

```
In [77]: r2_poly_train=Evaluate lin X_train_poly Y_train .r2()
r2_poly_test=Evaluate lin X_test_poly Y_test .r2()

r2_gbr_train=Evaluate gbr_pipeline X_train Y_train .r2()
r2_gbr_test=Evaluate gbr_pipeline X_test Y_test .r2()

r2_xgb_train=Evaluate xgb_pipeline X_train Y_train .r2()
r2_xgb_test=Evaluate xgb_pipeline X_test Y_test .r2()
```

## Cross Validation

It measures the degree of generalization of the model. It divides the data into several subsets by the amount said by the user. To subsequently obtain the average value of generalization.

```
In [74]: cv_poly=Evaluate lin X_test_poly Y_test .cv_score()
cv_gbr=Evaluate gbr_pipeline X_test Y_test .cv_score()
cv_xgb=Evaluate xgb_pipeline X_test Y_test .cv_score()
```

## Dataframe evaluation

```
In [76]: models_names=['Polynomial Regression','Gradient Boosting','XGBOOST']

mse_train=mse_poly_train,mse_gbr_train,mse_xgb_train
mse_test=mse_poly_test,mse_gbr_test,mse_xgb_test

r2_train=r2_poly_train,r2_gbr_train,r2_xgb_train
r2_test=r2_poly_test,r2_gbr_test,r2_xgb_test

cv=cv_poly,cv_gbr,cv_xgb

evaluate_df=pd.DataFrame({'Model':models_names,
                           'mse_train':mse_train,
                           'mse_test':mse_test,
                           'r2_train':r2_train,
                           'r2_test':r2_test,
                           'cv':cv})
```

```
In [79]: evaluate_df
```

```
Out[79]:
```

	Model	MSE train	MSE test	R2 train	R2 test	CV
0	Polynomial Regression	4.00027e+06	4.40150e+06	0.999999	0.999999	0.999999
1	Gradient Boosting	1.989649e+06	2.539044e+06	0.999549	0.998178	0.979221
2	XGBOOST	1.989653e+06	2.514028e+06	0.998583	0.998187	0.979176

XGBOOST and gradient Boosting have very similar metrics to gradient boosting. But XGBOOST performs better for test data.

```
In [83]: df_test=pd.DataFrame({'true_values':Y_test,
                             'pred_poly':pred_poly_test,
                             'pred_gbr':pred_gbr_test,
                             'pred_xgb':pred_xgb_test,
                             'smoker':X_test['smoker']})

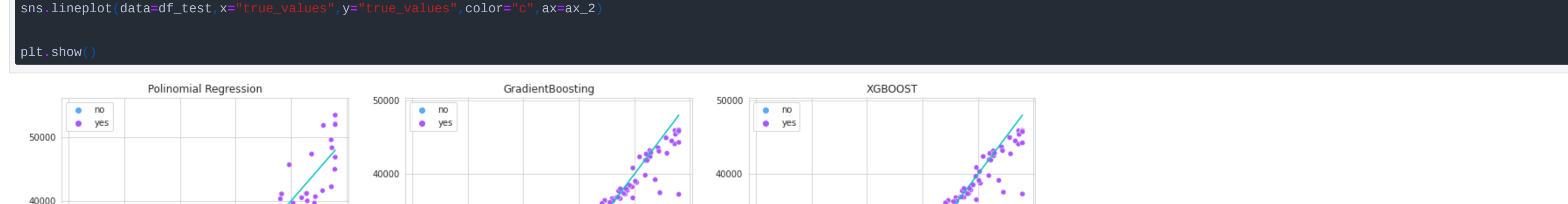
In [90]: fig,ax=plt.subplots(figsize=(10,10))
sns.set_style(style='darkgrid')

ax_0=plt.title('Polynomial Regression')
sns.scatterplot(data=df_test,x='true_values',y='pred_poly',color='red',hue='smoker',ax=ax_0,palette='cool')
sns.lineplot(data=df_test,x='true_values',y='true_values',color='C',ax=ax_0)

ax_1=plt.title('Gradient Boosting')
sns.scatterplot(data=df_test,x='true_values',y='pred_gbr',color='red',hue='smoker',ax=ax_1,palette='cool')
sns.lineplot(data=df_test,x='true_values',y='true_values',color='C',ax=ax_1)

ax_2=plt.title('XGBOOST')
sns.scatterplot(data=df_test,x='true_values',y='pred_xgb',color='red',hue='smoker',ax=ax_2,palette='cool')
sns.lineplot(data=df_test,x='true_values',y='true_values',color='C',ax=ax_2)

plt.show()
```



## Conclusion

- The **Polynomial Regression** is good to predict those users who do not smoke, for smokers it gives not so accurate predictions.
- **Gradient Boosting** It is good for both cases.
- **XGBOOST** Has a better root mean square error than using gradient boosting. In turn, it is faster to train and presents fewer symptoms of overfitting.