

SCREENSHOTS

```
File Edit View Run Kernel Help
Lesson_3_Project.py: Storing Test Results.py: Health Insurance Cost.py: Python 3

171: import pandas as pd
172: import numpy as np
173: import scipy as sp
174: import sklearn as sk
175: import matplotlib.pyplot as plt
176: from sklearn.model_selection import cross_val_score, KFold
177: from sklearn import model_selection
178: from sklearn import linear_model
179: from sklearn.metrics import mean_squared_error, mean_absolute_error

180: insurance=pd.read_csv('insurance.csv')
181: insurance.info()

182: (Class 'pandas.core.frame.DataFrame')
183: RangeIndex: 1338 entries, 0 to 1337
184: Data columns (total 8 columns):
185: age          1338 non-null int64
186: sex          1338 non-null int64
187: bmi          1338 non-null float64
188: children     1338 non-null int64
189: smoker       1338 non-null int64
190: region       1338 non-null int64
191: charges      1338 non-null float64
192: insuranceid   1338 non-null int64
193: dtypes: float64(1), int64(7)
194: memory usage: 81.7 KB

195: def map_smoking(column):
196:     mapped=[]
197:     for row in column:
198:         if row=="yes":
199:             mapped.append(1)
200:         else:
201:             mapped.append(0)
202:     return mapped
203: insurance["smoker_new"]=map_smoking(insurance["smoker"])

204: column_name=[col for col in insurance.select_dtypes(include="object")]
```


```
File Edit View Run Kernel Help
Lesson_3_Project.py: Storing Test Results.py: Health Insurance Cost.py: Python 3

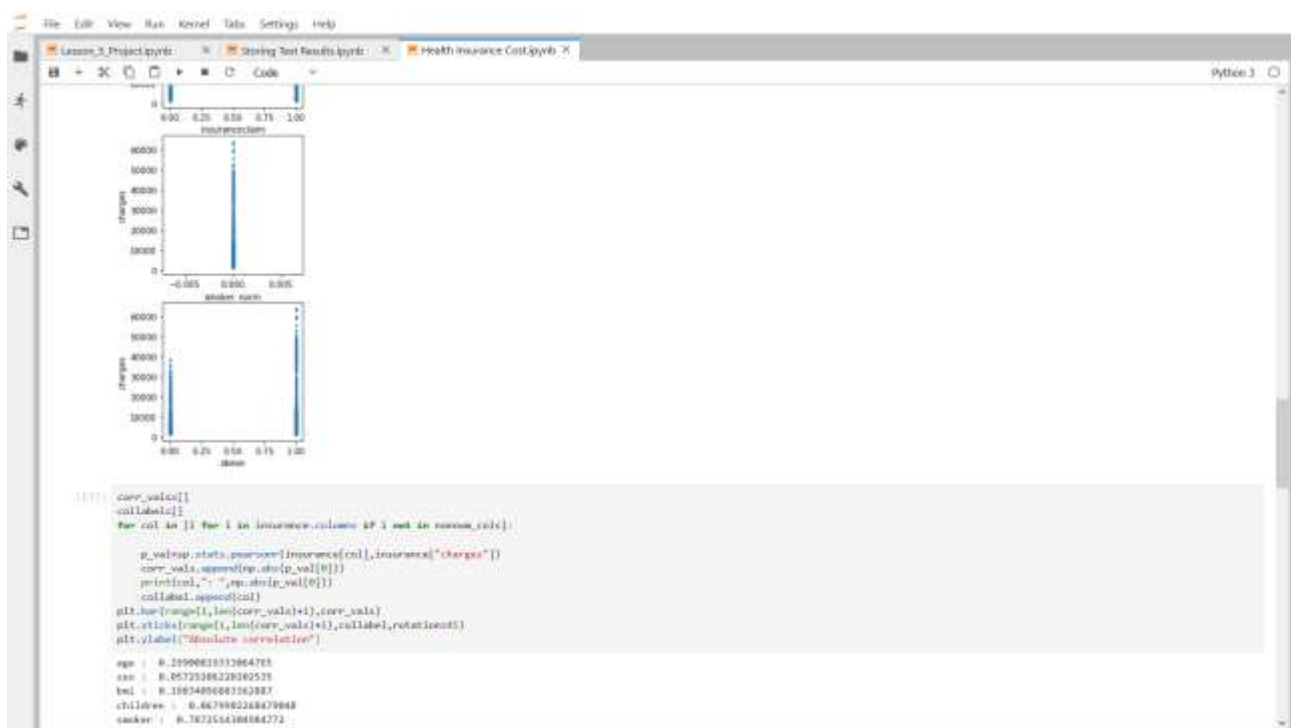
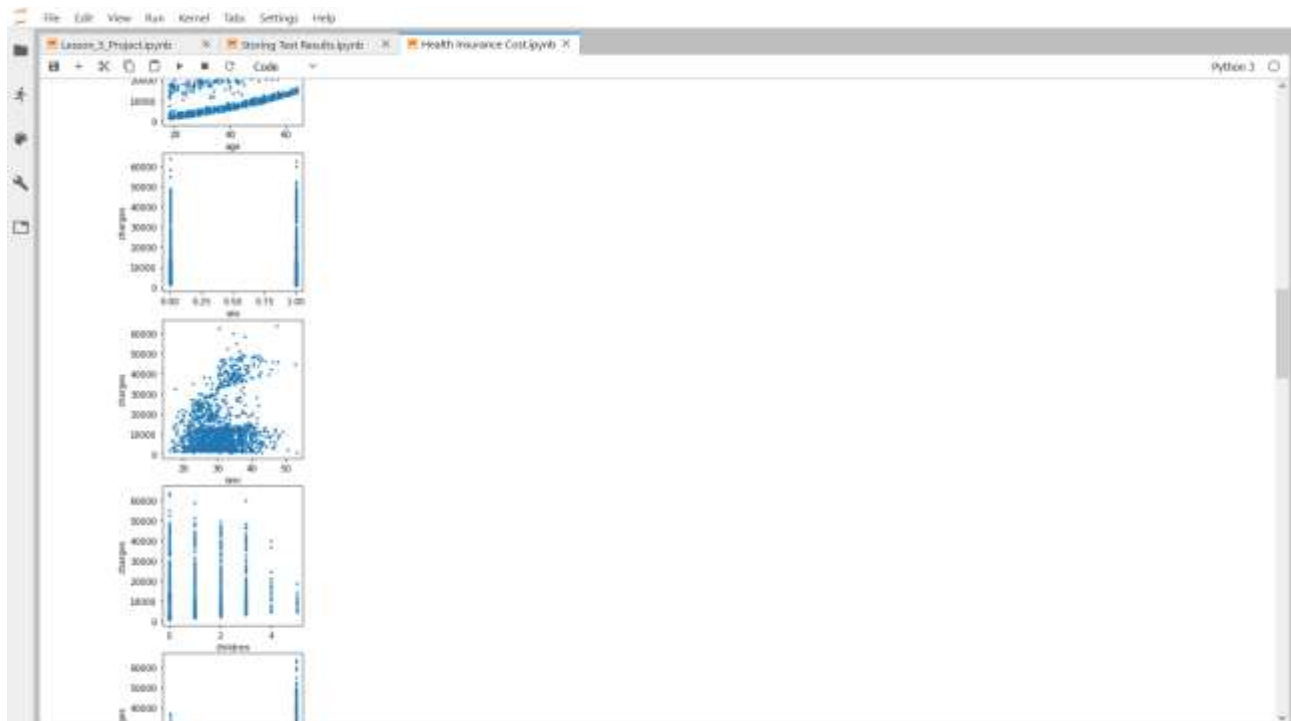
205: def map_obese(column):
206:     mapped=[]
207:     for row in column:
208:         if row>30:
209:             mapped.append(1)
210:         else:
211:             mapped.append(0)
212:     return mapped
213: insurance["obese"]=map_obese(insurance["bmi"])

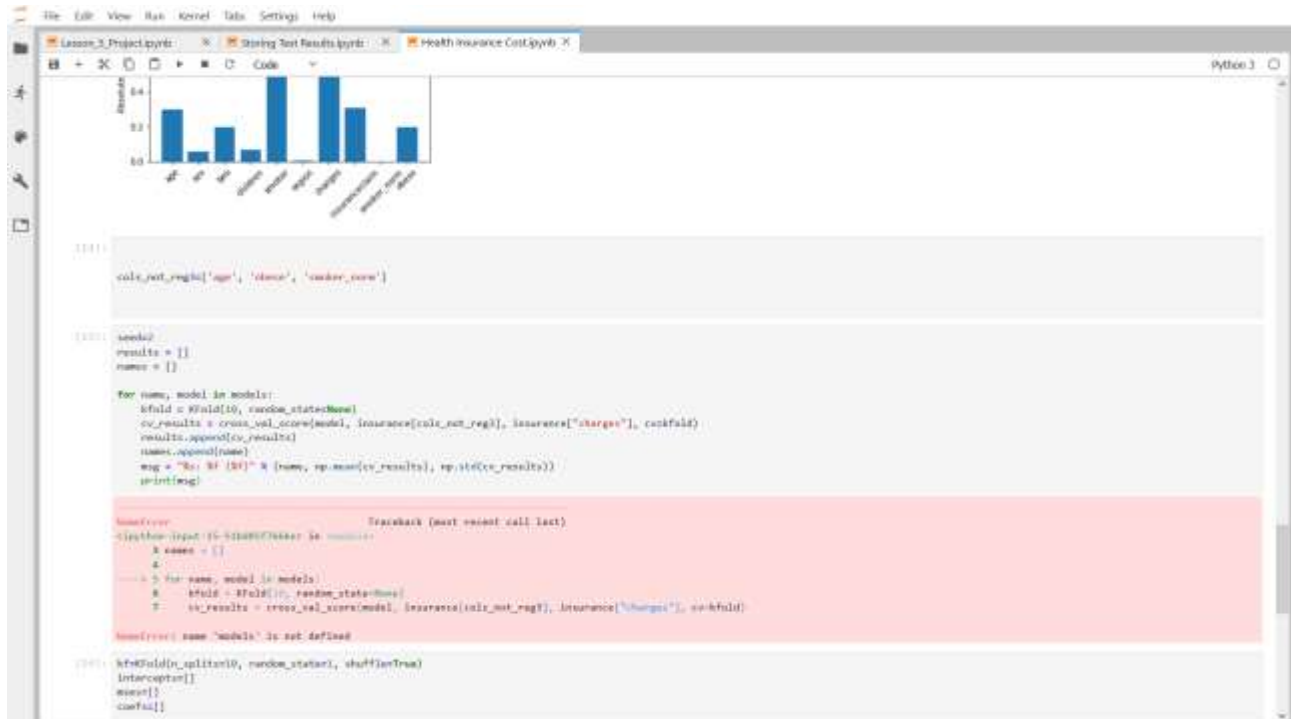
214: insurance.head(5)

215:    age  sex    bmi  children  smoker  region    charges  insuranceid  smoker_new  obese
0    19    0   27.900        0         1     0  16894.52400          1          0         0
1    18    1   31.770        1         0     2   1725.55200          2          0         1
2    38    1   33.000        3         0     2   4446.46200          3          0         1
3    31    1   22.705        0         0     1  21804.47081          4          0         0
4    32    1   26.880        0         0     1   3866.85020          5          0         0

216: column=list(insurance.columns)
217: fig=plt.subplots(column[1],figsize=(5,10))
218: ax[0].set_ylabel("charges")
219: p_value=[]
220: for ind,col in enumerate(1 for i in insurance.columns if i not in ["smoker","region","charges","sex_new"]):
221:     ax[ind].scatter(insurance[col],insurance.charges,xr0)
222:     ax[ind].set_ylabel(col)
223:     ax[ind].set_xlabel("charges")
224: plt.show()

225: 
```





```

import matplotlib.pyplot as plt
%matplotlib inline
from sklearn.model_selection import cross_val_score, KFold
from sklearn import model_selection
from sklearn import linear_model
from sklearn.metrics import mean_squared_error, mean_absolute_error

```

Loading the data

```

insurance=pd.read_csv("insurance2.csv")
insurance.info()
def map_smoking(column):
    mapped=[]

    for row in column:

        if row=="yes":
            mapped.append(1)
        else:
            mapped.append(0)

    return mapped
insurance["smoker_norm"]=map_smoking(insurance["smoker"])
nonnum_cols=[col for col in insurance.select_dtypes(include=["object"])]
def map_obese(column):
    mapped=[]
    for row in column:
        if row>30:
            mapped.append(1)
        else:
            mapped.append(0)
    return mapped
insurance["obese"]=map_obese(insurance["bmi"])
insurance.head(5)

```

We now explore the relation between the features given and the insurance costs

```

colnum=len(insurance.columns)-3
fig,ax=plt.subplots(colnum,1,figsize=(3,25))
ax[0].set_ylabel("charges")
p_vals={}
for ind,col in enumerate([i for i in insurance.columns if i not in
["smoker","region","charges","sex_norm"]]):

    ax[ind].scatter(insurance[col],insurance.charges,s=5)
    ax[ind].set_xlabel(col)
    ax[ind].set_ylabel("charges")
plt.show()

corr_vals=[]
collabel=[]
for col in [i for i in insurance.columns if i not in nonnum_cols]:

```

```

    p_val=sp.stats.pearsonr(insurance[col],insurance["charges"])
    corr_vals.append(np.abs(p_val[0]))
    print(col,": ",np.abs(p_val[0]))
    collabel.append(col)
plt.bar(range(1,len(corr_vals)+1),corr_vals)
plt.xticks(range(1,len(corr_vals)+1),collabel,rotation=45)
plt.ylabel("Absolute correlation")

cols_not_reg3=['age', 'obese', 'smoker_norm']
seed=2
results = []
names = []

for name, model in models:
    kfold = KFold(10, random_state=None)
    cv_results = cross_val_score(model, insurance[cols_not_reg3],
insurance["charges"], cv=kfold)
    results.append(cv_results)
    names.append(name)
    msg = "%s: %f (%f)" % (name, np.mean(cv_results), np.std(cv_results))
    print(msg)

kf=KFold(n_splits=10, random_state=1, shuffle=True)
intercepts=[]
mses=[]
coefs=[]

for train_index, test_index in kf.split(insurance[cols_not_reg3]):

    lr=linear_model.LinearRegression()

    lr.fit(insurance[cols_not_reg3].iloc[train_index],insurance["charges"].iloc[train_index])
    lr_predictions=lr.predict(insurance[cols_not_reg3].iloc[test_index])

    lr_mse=mean_squared_error(insurance["charges"].iloc[test_index],lr_predictions)

    intercepts.append(lr.intercept_)

    coefs.append(lr.coef_)
    mses.append(lr_mse)

rmses=[x**.5 for x in mses]
avg_rmse=np.mean(rmses)
avg_intercept=np.mean(intercepts)
age_coefs=[]
obesity_coefs=[]
smoking_coefs=[]
for vals in coefs:
    #print vals[0]
    age_coefs.append(vals[0])
    obesity_coefs.append(vals[1])
    smoking_coefs.append(vals[2])
age_coef=np.mean(age_coefs)

```

```
obesity_coef=np.mean(obesity_coefs)
smoking_coef=np.mean(smoking_coefs)
print("a: ",age_coef," b: ",obesity_coef," c: ",smoking_coef," intercept:
",avg_intercept)
```