

```
In [1]: import pandas as pd
df = pd.read_csv("C:/Users/prasa/Desktop/py codes/ds projects/ML/5 One
Hot Encoding/One_Hot_Encoding.csv")
df
```

Out[1]:

	town	area	price
0	monroe township	2600	550000
1	monroe township	3000	565000
2	monroe township	3200	610000
3	monroe township	3600	680000
4	monroe township	4000	725000
5	west windsor	2600	585000
6	west windsor	2800	615000
7	west windsor	3300	650000
8	west windsor	3600	710000
9	robinsville	2600	575000
10	robinsville	2900	600000
11	robinsville	3100	620000
12	robinsville	3600	695000

```
In [3]: dummies=pd.get_dummies(df.town)
dummies
```

Out[3]:

	monroe township	robinsville	west windsor
0	1	0	0
1	1	0	0

	monroe township	robinsville	west windsor
2	1	0	0
3	1	0	0
4	1	0	0
5	0	0	1
6	0	0	1
7	0	0	1
8	0	0	1
9	0	1	0
10	0	1	0
11	0	1	0
12	0	1	0

```
In [5]: merged = pd.concat([df,dummies],axis='columns')
merged
```

Out[5]:

	town	area	price	monroe township	robinsville	west windsor
0	monroe township	2600	550000	1	0	0
1	monroe township	3000	565000	1	0	0
2	monroe township	3200	610000	1	0	0
3	monroe township	3600	680000	1	0	0
4	monroe township	4000	725000	1	0	0
5	west windsor	2600	585000	0	0	1
6	west windsor	2800	615000	0	0	1
7	west windsor	3300	650000	0	0	1
8	west windsor	3600	710000	0	0	1

	town	area	price	monroe township	robinsville	west windsor
9	robinsville	2600	575000	0	1	0
10	robinsville	2900	600000	0	1	0
11	robinsville	3100	620000	0	1	0
12	robinsville	3600	695000	0	1	0

Dummy Variable Trap

When you can derive one variable from other variables, they are known to be multi-colinear. Here if you know values of california and georgia then you can easily infer value of new jersey state, i.e. california=0 and georgia=0. There for these state variables are called to be multi-colinear. In this situation linear regression won't work as expected. Hence you need to drop one column.

NOTE: sklearn library takes care of dummy variable trap hence even if you don't drop one of the state columns it is going to work, however we should make a habit of taking care of dummy variable trap ourselves just in case library that you are using is not handling this for you

```
In [8]: final = merged.drop(['town', 'west windsor'], axis='columns')
        final
```

Out[8]:

	area	price	monroe township	robinsville
0	2600	550000	1	0
1	3000	565000	1	0
2	3200	610000	1	0
3	3600	680000	1	0
4	4000	725000	1	0
5	2600	585000	0	0
6	2800	615000	0	0

	area	price	monroe township	robinsville
7	3300	650000	0	0
8	3600	710000	0	0
9	2600	575000	0	1
10	2900	600000	0	1
11	3100	620000	0	1
12	3600	695000	0	1

```
In [9]: from sklearn.linear_model import LinearRegression
model = LinearRegression()
```

```
In [13]: x = final.drop('price',axis='columns') #Remove price since its dependen
t variable
x
```

Out[13]:

	area	monroe township	robinsville
0	2600	1	0
1	3000	1	0
2	3200	1	0
3	3600	1	0
4	4000	1	0
5	2600	0	0
6	2800	0	0
7	3300	0	0
8	3600	0	0
9	2600	0	1
10	2900	0	1
11	3100	0	1

	area	monroe township	robinsville
12	3600	0	1

```
In [11]: y = final.price
y
```

```
Out[11]: 0    550000
1    565000
2    610000
3    680000
4    725000
5    585000
6    615000
7    650000
8    710000
9    575000
10   600000
11   620000
12   695000
Name: price, dtype: int64
```

```
In [14]: model.fit(x,y)
```

```
Out[14]: LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)
```

```
In [15]: model.predict([[2800,0,1]])
```

```
Out[15]: array([590775.63964739])
```

```
In [16]: model.predict([[3400,0,0]])
```

```
Out[16]: array([681241.66845839])
```

```
In [17]: model.score(x,y) #check accuracy
```

Out[17]: 0.9573929037221873

Using sklearn OneHotEncoder

First step is to use label encoder to convert town names into numbers

```
In [95]: from sklearn.preprocessing import LabelEncoder  
le = LabelEncoder()
```

```
In [96]: dfle = df  
dfle.town = le.fit_transform(dfle.town)  
dfle
```

Out[96]:

	town	area	price
0	0	2600	550000
1	0	3000	565000
2	0	3200	610000
3	0	3600	680000
4	0	4000	725000
5	2	2600	585000
6	2	2800	615000
7	2	3300	650000
8	2	3600	710000
9	1	2600	575000
10	1	2900	600000
11	1	3100	620000
12	1	3600	695000

```
In [97]: X = dfle[['town', 'area']].values
X
```

```
Out[97]: array([[ 0, 2600],
 [ 0, 3000],
 [ 0, 3200],
 [ 0, 3600],
 [ 0, 4000],
 [ 2, 2600],
 [ 2, 2800],
 [ 2, 3300],
 [ 2, 3600],
 [ 1, 2600],
 [ 1, 2900],
 [ 1, 3100],
 [ 1, 3600]], dtype=int64)
```

```
In [98]: y = dfle.price.values
y
```

```
Out[98]: array([550000, 565000, 610000, 680000, 725000, 585000, 615000, 650000,
 710000, 575000, 600000, 620000, 695000], dtype=int64)
```

Now use one hot encoder to create dummy variables for each of the town

```
In [99]: from sklearn.preprocessing import OneHotEncoder
from sklearn.compose import ColumnTransformer
ct = ColumnTransformer([('town', OneHotEncoder(), [0])], remainder = 'passthrough')
```

```
In [100]: X = ct.fit_transform(X)
X
```

```
Out[100]: array([[1.0e+00, 0.0e+00, 0.0e+00, 2.6e+03],
 [1.0e+00, 0.0e+00, 0.0e+00, 3.0e+03],
 [1.0e+00, 0.0e+00, 0.0e+00, 3.2e+03],
 [1.0e+00, 0.0e+00, 0.0e+00, 3.6e+03],
 [1.0e+00, 0.0e+00, 0.0e+00, 4.0e+03],
```

```
[0.0e+00, 0.0e+00, 1.0e+00, 2.6e+03],  
[0.0e+00, 0.0e+00, 1.0e+00, 2.8e+03],  
[0.0e+00, 0.0e+00, 1.0e+00, 3.3e+03],  
[0.0e+00, 0.0e+00, 1.0e+00, 3.6e+03],  
[0.0e+00, 1.0e+00, 0.0e+00, 2.6e+03],  
[0.0e+00, 1.0e+00, 0.0e+00, 2.9e+03],  
[0.0e+00, 1.0e+00, 0.0e+00, 3.1e+03],  
[0.0e+00, 1.0e+00, 0.0e+00, 3.6e+03]])
```

```
In [101]: X = X[:,1:]
```

```
In [102]: X
```

```
Out[102]: array([[0.0e+00, 0.0e+00, 2.6e+03],  
                 [0.0e+00, 0.0e+00, 3.0e+03],  
                 [0.0e+00, 0.0e+00, 3.2e+03],  
                 [0.0e+00, 0.0e+00, 3.6e+03],  
                 [0.0e+00, 0.0e+00, 4.0e+03],  
                 [0.0e+00, 1.0e+00, 2.6e+03],  
                 [0.0e+00, 1.0e+00, 2.8e+03],  
                 [0.0e+00, 1.0e+00, 3.3e+03],  
                 [0.0e+00, 1.0e+00, 3.6e+03],  
                 [1.0e+00, 0.0e+00, 2.6e+03],  
                 [1.0e+00, 0.0e+00, 2.9e+03],  
                 [1.0e+00, 0.0e+00, 3.1e+03],  
                 [1.0e+00, 0.0e+00, 3.6e+03]])
```

```
In [103]: model.fit(X,y)
```

```
Out[103]: LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)
```

```
In [104]: model.predict([[0,1,3400]])
```

```
Out[104]: array([681241.6684584])
```

```
In [105]: model.predict([[1,0,2800]])
```



```
Out[105]: array([590775.63964739])
```

Exercise

At the same level as this notebook on github, there is an Exercise folder that contains carprices.csv. This file has car sell prices for 3 different models. First plot data points on a scatter plot chart to see if linear regression model can be applied. If yes, then build a model that can answer following questions,

- 1) Predict price of a mercedes benz that is 4 yr old with mileage 45000
- 2) Predict price of a BMW X5 that is 7 yr old with mileage 86000
- 3) Tell me the score (accuracy) of your model. (Hint: use `LinearRegression().score()`)