

```

1
2 import numpy as np # linear algebra
3 import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
4

```

We have a file called hour.csv which contain hourly rental data of 2 different years on hourly basis for all 24 hours. It contain 17379 entries none of which is nan, means no value from this dataset is missing. we have object dteday which shows year-month-date.

General information about different features can be found by **info()** method on pandas data frame. So we have first converted our csv file to pandas dataframe using **read\_csv()** method.

```

1 total_data=pd.read_csv("hour.csv")
2 total_data.head()

```

	instant	dteday	season	yr	mnth	hr	holiday	weekday	workingday	weathersit	temp
0	1	2011-01-01	1	0	1	0	0	6	0	1	0.24
1	2	2011-01-01	1	0	1	1	0	6	0	1	0.22
2	3	2011-01-01	1	0	1	2	0	6	0	1	0.22
3	4	2011-01-01	1	0	1	3	0	6	0	1	0.24

Since date is not given in total\_data by help of lambda function and split method we are converting string object into list object, which will have year month and date as its value at 0,1 and 2 index and then choose value at 2nd index which is date. and then we insert this date value in our total\_data using **total\_data['feature\_name']=value**

And then as mentioned in question we take data upto first 19 days into train dataset and remaining in the test dataset.

```

1 date=total_data.dteday.apply(lambda x: x[0:].split('-')[2])
2 date=date.astype('int64')
3 total_data['date']=date
4 total_data.head()

```

```

5 train_d=total_data[(total_data['date']<=19)]
6 test_d=total_data[(total_data['date']>19)]

```

Here head() method prints first 5 rows of our dataframe to peek how data looks like.

```
1 train_d.head()
```

	instant	dteday	season	yr	mnth	hr	holiday	weekday	workingday	weathersit	temp
0	1	2011-01-01	1	0	1	0	0	6	0	1	0.24
1	2	2011-01-01	1	0	1	1	0	6	0	1	0.22
2	3	2011-01-01	1	0	1	2	0	6	0	1	0.22
3	4	2011-01-01	1	0	1	3	0	6	0	1	0.24

```
1 test_d.head()
```

	instant	dteday	season	yr	mnth	hr	holiday	weekday	workingday	weathersit	te
431	432	2011-01-20	1	0	1	0	0	4	1	1	0.
432	433	2011-01-20	1	0	1	1	0	4	1	1	0.
433	434	2011-01-20	1	0	1	2	0	4	1	1	0.
434	435	2011-01-20	1	0	1	3	0	4	1	1	0.

▼ conclusions from data:-

This info() method tells us that there are 10886 rows in our **train\_d** dataframe which contain first 19 date data from total\_data. It has 18 columns out of which 13 are of type int, 4 are of type float and 1 object

```
1 train_d.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 10886 entries, 0 to 17092
Data columns (total 18 columns):
#   Column          Non-Null Count  Dtype
---  -
0   instant         10886 non-null  int64
1   dteday          10886 non-null  object
2   season          10886 non-null  int64
3   yr              10886 non-null  int64
4   mnth            10886 non-null  int64
5   hr              10886 non-null  int64
6   holiday         10886 non-null  int64
7   weekday         10886 non-null  int64
8   workingday      10886 non-null  int64
9   weathersit       10886 non-null  int64
10  temp            10886 non-null  float64
11  atemp           10886 non-null  float64
12  hum             10886 non-null  float64
13  windspeed       10886 non-null  float64
14  casual          10886 non-null  int64
15  registered      10886 non-null  int64
16  cnt             10886 non-null  int64
17  date            10886 non-null  int64
dtypes: float64(4), int64(13), object(1)
memory usage: 1.6+ MB
```

**describe()** method gives us information about quantitative aspects of int and float data like mean,variance,quantiles and so on.

```
1 train_d.describe()
```

instant season yr mnth hr holiday

Our **test\_d data** frame contain 6493 values which are left after taking first 19 dates of month in train\_d dataframe.

```
1 test_d.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 6493 entries, 431 to 17378
Data columns (total 18 columns):
#   Column                Non-Null Count  Dtype  
---  -
0   instant               6493 non-null  int64  
1   dteday                6493 non-null  object  
2   season                6493 non-null  int64  
3   yr                    6493 non-null  int64  
4   mnth                  6493 non-null  int64  
5   hr                    6493 non-null  int64  
6   holiday                6493 non-null  int64  
7   weekday               6493 non-null  int64  
8   workingday            6493 non-null  int64  
9   weathersit             6493 non-null  int64  
10  temp                  6493 non-null  float64 
11  atemp                 6493 non-null  float64 
12  hum                   6493 non-null  float64 
13  windspeed             6493 non-null  float64 
14  casual                6493 non-null  int64  
15  registered             6493 non-null  int64  
16  cnt                   6493 non-null  int64  
17  date                  6493 non-null  int64  
dtypes: float64(4), int64(13), object(1)
memory usage: 963.8+ KB
```

```
1 test_d.describe()

instant season yr mnth hr holiday
count 6493.000000 6493.000000 6493.000000 6493.000000 6493.000000 6493.000000 6493.000000
mean 8945.809795 2.493300 0.503619 6.565070 11.555367 0.029108 3.000000
std 4991.272309 1.091258 0.500025 3.429462 6.912526 0.168123 2.000000
min 432.000000 1.000000 0.000000 1.000000 0.000000 0.000000 0.000000
25% 4770.000000 2.000000 0.000000 4.000000 6.000000 0.000000 1.000000
50% 9122.000000 3.000000 1.000000 7.000000 12.000000 0.000000 3.000000
75% 13477.000000 3.000000 1.000000 10.000000 18.000000 0.000000 5.000000
max 17379.000000 4.000000 1.000000 12.000000 23.000000 1.000000 6.000000
```

to get insight from data we will plot count values for year,season,month,date and time values matplotlib.

```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
```

now using groupby method on our dataframe we will find mean values of count feature for different different groups of year,season,month and date values.

For printing mean values of count for year 0 and 1 we have used groupby() method, this method groups data as per given input feature and then we obtain value of mean of count using **train\_d.groupby('yr')['cnt'].mean()**

Or we can find mean by explicitly writing **(train\_d['season']==1)** which give us all rows with season 1 and then finding mean of count in those rows.

Mean per year:-

- 0 144.223349
- 1 238.560944

Mean per Season:-

- season1: 116.34326135517499
- season2: 215.25137211855105
- season3: 234.417124039517
- season4: 198.98829553767374

```
1 years=train_d.groupby('yr')['cnt'].mean()
2 print(years)
3 print('\n')
4
5 season1=train_d[(train_d['season']==1)]['cnt'].mean()
6 print("season1: ",season1)
7 season2=train_d[(train_d['season']==2)]['cnt'].mean()
8 print("season2: ",season2)
9 season3=train_d[(train_d['season']==3)]['cnt'].mean()
10 print("season3: ",season3)
11 season4=train_d[(train_d['season']==4)]['cnt'].mean()
12 print("season4: ",season4)
```

yr

```
0    144.223349
1    238.560944
Name: cnt, dtype: float64
```

```
season1: 116.34326135517499
season2: 215.25137211855105
season3: 234.417124039517
season4: 198.98829553767374
```

Similarly we can

```
1 months_mean=train_d.groupby('mnth')['cnt'].mean()
2 print(months_mean)
```

```
mnth
1     90.366516
2    110.003330
3    148.169811
4    184.160616
5    219.459430
6    242.031798
7    235.325658
8    234.118421
9    233.805281
10   227.699232
11   193.677278
12   175.614035
Name: cnt, dtype: float64
```

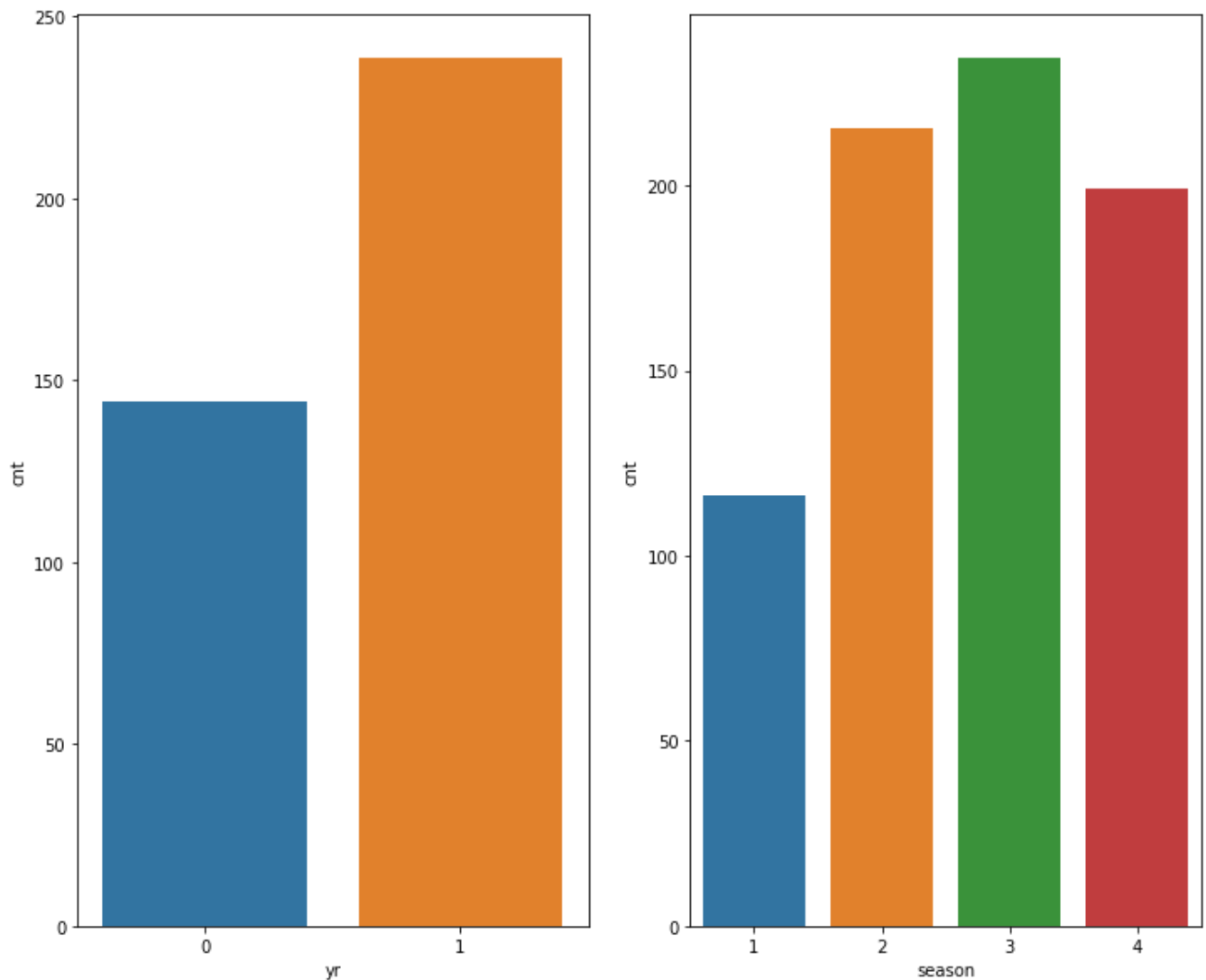
```
1 dates_mean=train_d.groupby('date')['cnt'].mean()
2 print(dates_mean)
3 #mean count for each day
```

```
date
1    180.333913
2    183.910995
3    194.696335
4    195.705575
5    189.765217
6    189.860140
7    183.773519
8    179.041812
9    187.897391
10   195.183566
11   195.679577
12   190.675393
13   194.160279
14   195.829268
15   201.527875
16   191.353659
17   205.660870
18   192.605684
```

```
19 192.311847
Name: cnt, dtype: float64
```

And as asked we can also plot these mean values of count for different groups using matplotlib. subplot provide us convenience to plot different plots simultaneously  $2 \times 2 = 4$  plots.

```
1 %matplotlib inline
2 fig = plt.figure(figsize=[12,10])
3 x1 = fig.add_subplot(1,2,1)
4 x1 = sns.barplot(x='yr',y='cnt',data=train_d.groupby('yr')['cnt'].mean().reset_index())
5
6 x2 = fig.add_subplot(1,2,2)
7 x2 = sns.barplot(x='season',y='cnt',data=train_d.groupby('season')['cnt'].mean().reset_index())
```

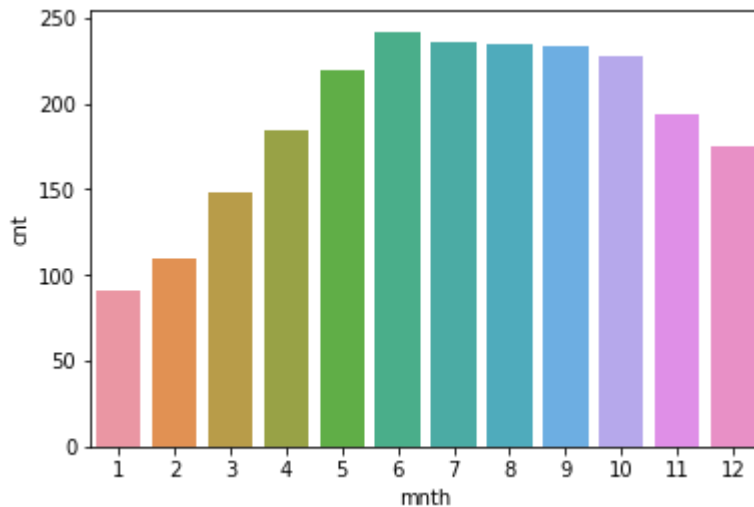


```
1 x3 = fig.add_subplot(1,1,1)
```

```

2 x3 = sns.barplot(x='mnth',y='cnt',data=train_d.groupby('mnth')['cnt'].mean()).
3

```

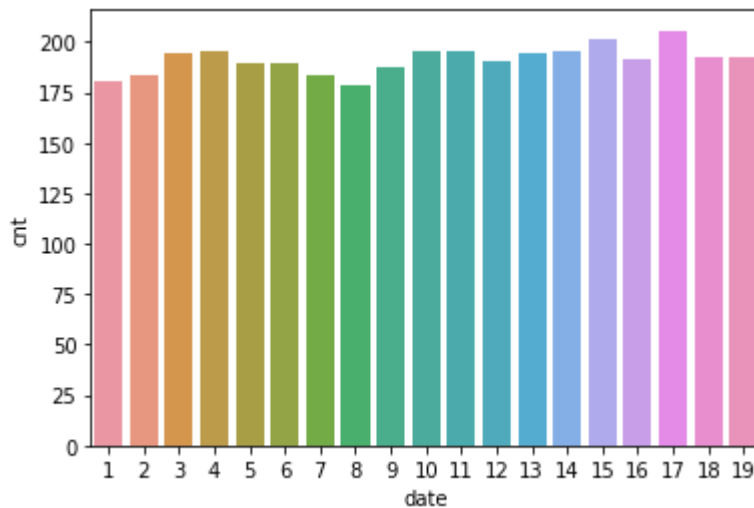


```

1 x4 = fig.add_subplot(1,1,1)
2 x4 = sns.barplot(x='date',y='cnt',data=train_d.groupby('date')['cnt'].mean()).

```

/usr/local/lib/python3.6/dist-packages/ipykernel\_launcher.py:1: MatplotlibDeprecationWarning:   
 """Entry point for launching an IPython kernel.

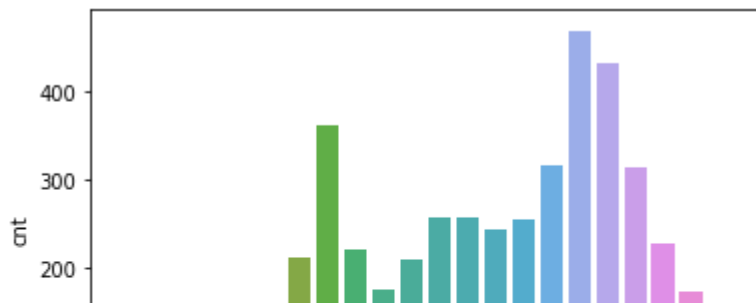


```

1 x5 = fig.add_subplot(1,1,1)
2 x5 = sns.barplot(x='hr',y='cnt',data=train_d.groupby('hr')['cnt'].mean()).rese
3
4

```





Now its time to find poission regression root mean square error values for both and train and test data and we will do this with the help of nd arrays which provide various matrix functionality.



and we will take count as Y for both of our train and test data.

```
1 Y_train=train_d['cnt']
2 Y_test=test_d['cnt']
```

```
1 Y_train = np.array(Y_train).reshape(-1,1)
2 Y_test = np.array(Y_test).reshape(-1,1)
3 Y_train
```

```
array([[ 16],
       [ 40],
       [ 32],
       ...,
       [168],
       [129],
       [ 88]])
```

We will do some manipulation on dataframe so its better to have a copy of our train and test data and use these for linear regression. we are aslo applying one hot encoding on holiday, weathersit, and season.

```
1 posrgsn_train=train_d.copy()
2 posrgsn_test=test_d.copy()
3
4 posrgsn_train=pd.get_dummies(posrgsn_train,columns=['holiday','weathersit','s
5 posrgsn_test=pd.get_dummies(posrgsn_test,columns=['holiday','weathersit','s
```

For training L1 and L2 norm models we take simple training model data and randomly assign 80% of data to new testing data and remaining 20% data to validation data.

```
1 np.random.seed(10)
2 mask=np.random.rand(len(posrgsn_train)) < .8
```

```

3 Rposrgsn_train=posrgsn_train[mask].copy()
4 Rposrgsn_validation=posrgsn_train[~mask].copy()

```

Now from this New (for Regularization-R) training and testing data assing RYvalidation and RYtrain

```

1 RY_train=Rposrgsn_train['cnt']
2 RY_validation=Rposrgsn_validation['cnt']
3
4 RY_train = np.array(RY_train).reshape(-1,1)
5 RY_validation = np.array(RY_validation).reshape(-1,1)

```

Remove all other features such as registered,casual and dteday. Also remove features like working day, atemp which are highly correlated with holiday and temperature.

```

1
2 posrgsn_train.drop(['casual','registered','dteday','cnt','instant','workingda
3 posrgsn_test.drop(['casual','registered','dteday','cnt','instant','workingday
4 Rposrgsn_train.drop(['casual','registered','dteday','cnt','instant','workingc
5 Rposrgsn_validation.drop(['casual','registered','dteday','cnt','instant','wor

```

```

1 posrgsn_test.info()

```

```

<class 'pandas.core.frame.DataFrame'>
Int64Index: 6493 entries, 431 to 17378
Data columns (total 18 columns):
#   Column                Non-Null Count  Dtype
---  -
0   yr                    6493 non-null   int64
1   mnth                 6493 non-null   int64
2   hr                   6493 non-null   int64
3   weekday              6493 non-null   int64
4   temp                 6493 non-null   float64
5   hum                  6493 non-null   float64
6   windspeed            6493 non-null   float64
7   date                 6493 non-null   int64
8   holiday_0            6493 non-null   uint8
9   holiday_1            6493 non-null   uint8
10  weathersit_1          6493 non-null   uint8
11  weathersit_2          6493 non-null   uint8
12  weathersit_3          6493 non-null   uint8
13  weathersit_4          6493 non-null   uint8
14  season_1              6493 non-null   uint8
15  season_2              6493 non-null   uint8
16  season_3              6493 non-null   uint8
17  season_4              6493 non-null   uint8

```

```
dtypes: float64(3), int64(5), uint8(10)
```

```
1 posrgsn_train.head()
```

	yr	mnth	hr	weekday	temp	hum	windspeed	date	holiday_0	holiday_1	weathersit_
0	0	1	0	6	0.24	0.81	0.0	1	1	0	
1	0	1	1	6	0.22	0.80	0.0	1	1	0	
2	0	1	2	6	0.22	0.80	0.0	1	1	0	
3	0	1	3	6	0.24	0.75	0.0	1	1	0	
4	0	1	4	6	0.24	0.75	0.0	1	1	0	

We will normalize our data to so that at the end we can make some sense of our rmse values and to be consistent with our features and in some ml methods it is very critical for better optimization.

```
1 X_train=np.array(posrgsn_train.values)
2
3 X_test=np.array(posrgsn_test.values)
4
```

```
1 print(X_train)
```

```
[[ 0.  1.  0. ...  0.  0.  0.]
 [ 0.  1.  1. ...  0.  0.  0.]
 [ 0.  1.  2. ...  0.  0.  0.]
 ...
 [ 1. 12. 21. ...  0.  0.  1.]
 [ 1. 12. 22. ...  0.  0.  1.]
 [ 1. 12. 23. ...  0.  0.  1.]]
```

```
1 def prediction(w, X):
2     y_hat = np.exp(np.matmul(X, w))
3     return y_hat
4
5 def gradient(X, y, y_hat):
6     gradient = np.divide(np.matmul(np.transpose(X), np.subtract(y_hat, y)), 1
7     #gradient = np.matmul(np.transpose(X), np.subtract(y_hat, y))
8     return gradient
9
10 def gradient_l1(X, y, y_hat, w, reg_const):
11     gradient_values = gradient(X, y, y_hat)
12     m = len(w)
13     for j in range(len(w)):
14         if j == 0:
15             gradient_values[j] = gradient_values[j] + reg_const
```

```

14         if w[j,0] > 0:
15             gradient_values[j,0] += (reg_const / m)
16         else:
17             gradient_values[j,0] -= (reg_const / m)
18
19     return gradient_values
20
21 def gradient_l2(X, y, y_hat, w, reg_const):
22     gradient_values_l2 = gradient(X, y, y_hat) + reg_const * w
23     return gradient_values_l2
24
25
26 def gradient_descent(X, y, alpha = 0.01, iterations = 50000, reg_const =0, c
27     w = np.zeros((len(X[1,:]), 1))
28     for i in range(iterations):
29         y_hat = prediction(w, X)
30         gradient_value = gradient(X, y, y_hat)
31         w = np.subtract(w, alpha * gradient_value)
32     return w
33
34 def gradient_descent1(X, y, alpha = 0.01, iterations = 50000, reg_const =0, c
35     w = np.zeros((len(X[1,:]), 1))
36     for i in range(iterations):
37         y_hat = prediction(w, X)
38         gradient_value = gradient_l1(X, y, y_hat,w,reg_const)
39         w = np.subtract(w, alpha * gradient_value)
40     return w
41
42 def gradient_descent2(X, y, alpha = 0.01, iterations = 50000, reg_const =0, c
43     w = np.zeros((len(X[1,:]), 1))
44     for i in range(iterations):
45         y_hat = prediction(w, X)
46         gradient_value = gradient_l2(X, y, y_hat,w,reg_const)
47         w = np.subtract(w, alpha * gradient_value)
48     return w
49
50
51

```

training X and training Y, now calculate wcap using gradient descent.

```

1 print(X_train)
2 print(Y_train)
3 wcap= gradient_descent(X_train, Y_train, iterations = 10000)

[[ 0.  1.  0. ...  0.  0.  0.]

```

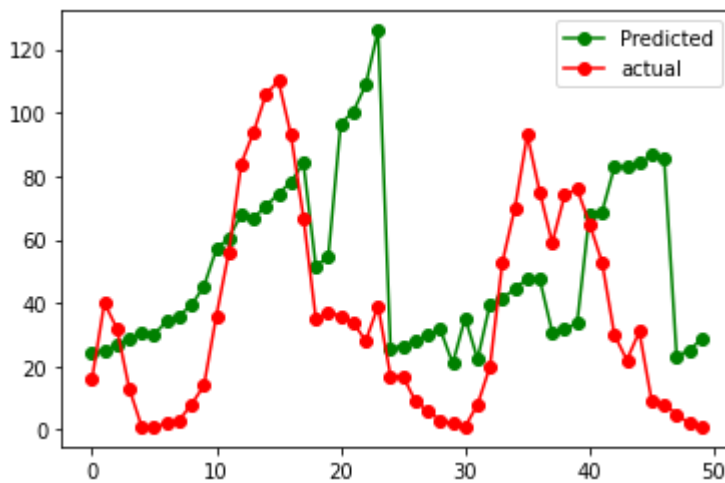
```
[ 0.  1.  1. ...  0.  0.  0.]
[ 0.  1.  2. ...  0.  0.  0.]
...
[ 1. 12. 21. ...  0.  0.  1.]
[ 1. 12. 22. ...  0.  0.  1.]
[ 1. 12. 23. ...  0.  0.  1.]]
[[ 16]
 [ 40]
 [ 32]
 ...
 [168]
 [129]
 [ 88]]
```

```
1 Y_train_hat= np.exp(np.matmul(X_train, wcap))
```

For poisson regression, training rmse and plot

```
1 rmse_train = np.sqrt(np.mean((Y_train - Y_train_hat)**2))
2 print(rmse_train)
3 fig = plt.figure()
4 predicted = plt.plot(range(50), Y_train_hat[:50], 'go-', label='Predicted')
5 actual = plt.plot(range(50), Y_train[:50], 'ro-', label="actual")
6 plt.legend()
7 plt.show()
8
```

155.90952773356904



For poisson regression test rmse and plot

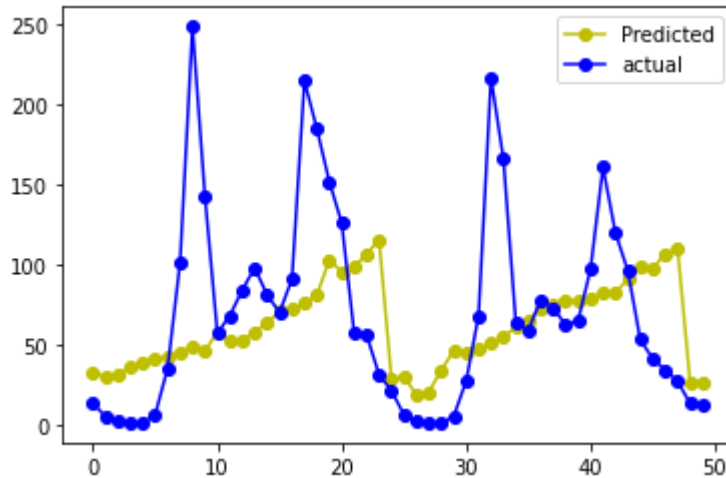
```
1 Y_test_hat= np.exp(np.matmul(X_test, wcap))
2 rmse_test = np.sqrt(np.mean((Y_train - Y_train_hat)**2))
3 print(rmse_test)
4 fig = plt.figure()
```

```

5 predicted = plt.plot(range(50), Y_test_hat[:50], 'yo-', label='Predicted')
6 actual = plt.plot(range(50), Y_test[:50], 'bo-', label="actual")
7 plt.legend()
8 plt.show()
9

```

155.90952773356904



copying train data for regularization and to convert it into 2 different sets ,new training set and validation set.  $10886 \times 0.8 = 8708$  new training set and remaining 20% validation set which is 2177.

```
1 Rposrgsn_validation.info()
```

```

<class 'pandas.core.frame.DataFrame'>
Int64Index: 2117 entries, 11 to 17083
Data columns (total 18 columns):
 #   Column                Non-Null Count  Dtype  
---  -
 0   yr                    2117 non-null  int64  
 1   mnth                  2117 non-null  int64  
 2   hr                    2117 non-null  int64  
 3   weekday               2117 non-null  int64  
 4   temp                  2117 non-null  float64 
 5   hum                   2117 non-null  float64 
 6   windspeed             2117 non-null  float64 
 7   date                  2117 non-null  int64  
 8   holiday_0             2117 non-null  uint8  
 9   holiday_1             2117 non-null  uint8  
10  weathersit_1           2117 non-null  uint8  
11  weathersit_2           2117 non-null  uint8  
12  weathersit_3           2117 non-null  uint8  
13  weathersit_4           2117 non-null  uint8  
14  season_1              2117 non-null  uint8  
15  season_2              2117 non-null  uint8  
16  season_3              2117 non-null  uint8  
17  season_4              2117 non-null  uint8  
dtypes: float64(3), int64(5), uint8(10)
memory usage: 169.5 KB

```

Now let us first obtained train and validation Y values from our our newly created test and validation set data.

```
1
2 Rposrgsn_validation.head()
3
```

	yr	mnth	hr	weekday	temp	hum	windspeed	date	holiday_0	holiday_1	weathersit
<b>11</b>	0	1	11	6	0.36	0.81	0.2836	1	1	0	
<b>14</b>	0	1	14	6	0.46	0.72	0.2836	1	1	0	
<b>18</b>	0	1	18	6	0.42	0.88	0.2537	1	1	0	
<b>30</b>	0	1	7	0	0.40	0.76	0.1940	2	1	0	
<b>32</b>	0	1	9	0	0.38	0.76	0.2239	2	1	0	

```
1 RX_train=np.array(Rposrgsn_train.values)
2 RX_validation=np.array(Rposrgsn_validation.values)
```

```
1 print(RX_train.shape)
2 print(RY_train.shape)
3 print(RX_validation.shape)
4 print(RY_validation.shape)
```

```
(8769, 18)
(8769, 1)
(2117, 18)
(2117, 1)
```

```
1
2 lambda_values=[0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1,3,10]
3
```

I will use different values of lambda ranging from 0.001 to 10 and select the best value of it. so lambda values will be 0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1, 5, 10

checking different lambda values and appending these to array

```
1 Rrmse_validation_array=[]
2 for lambda_v in lambda_values:
```

```

3     Rw=gradient_descent1(RX_train,RY_train,iterations = 5000, reg_const =lamt
4     RY_validation_hat = np.exp(np.matmul(RX_validation, Rw))
5     Rrmse_validation = np.sqrt(np.mean((RY_validation - RY_validation_hat)**2
6     Rrmse_validation_array.append(Rrmse_validation)
7
8
9
10

```

finding optimized lambda and then

```

1  min_index=Rrmse_validation_array.index(min(Rrmse_validation_array))
2  op_lambda=lambda_values[min_index]
3  Rwcap=gradient_descent1(RX_train,RY_train,reg_const=op_lambda,iterations = 30
4
5

```

```

1  Y_test_hat = np.exp(np.matmul(X_test, Rwcap))
2  Rrmse_test = np.sqrt(np.mean((Y_test - Y_test_hat)**2))

```

```

1  print(Rrmse_test)

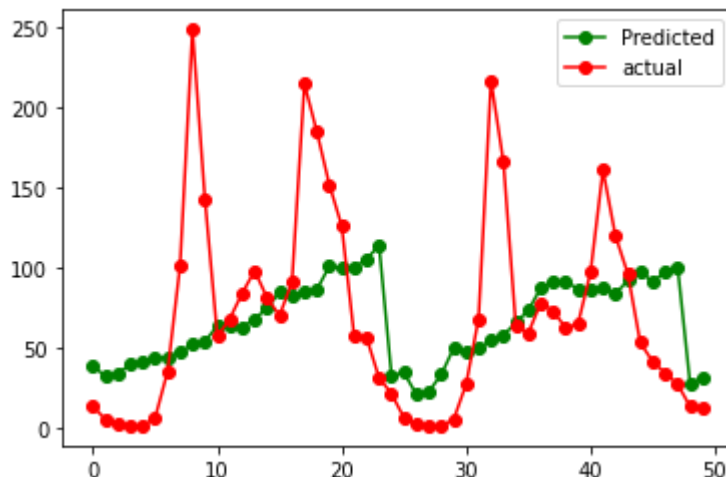
```

235.09398549717537

```

1  fig = plt.figure()
2  predicted = plt.plot(range(50), Y_test_hat[:50], 'go-',label='Predicted')
3  actual = plt.plot(range(50), Y_test[:50], 'ro-',label="actual")
4  plt.legend()
5  plt.show()

```



```

1  R2rmse_validation_array=[]

```



```

2 for lambda_v in lambda_values:
3     Rw=gradient_descent2(RX_train,RY_train,iterations = 5000, reg_const =lamt
4     RY_validation_hat = np.exp(np.matmul(RX_validation, Rw))
5     R2rmse_validation = np.sqrt(np.mean((RY_validation - RY_validation_hat)**
6     R2rmse_validation_array.append(R2rmse_validation)
7
8

```

```

1 min_index=R2rmse_validation_array.index(min(R2rmse_validation_array))
2 op_lambda=lambda_values[min_index]
3 Rwcap=gradient_descent1(RX_train,RY_train,reg_const=op_lambda,iterations = 40

```

```

1 Y_test_hat = np.exp(np.matmul(X_test, Rwcap))
2 R2rmse_test = np.sqrt(np.mean((Y_test - Y_test_hat)**2))

```

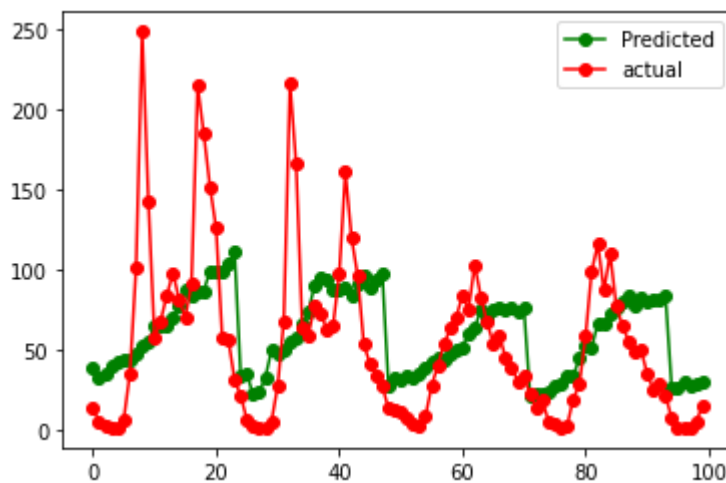
```
1 print(R2rmse_test)
```

216.98590436876498

```

1 fig = plt.figure()
2 predicted = plt.plot(range(100), Y_test_hat[:100], 'go-',label='Predicted')
3 actual = plt.plot(range(100), Y_test[:100], 'ro-',label="actual")
4 plt.legend()
5 plt.show()

```



1

