

rec03

October 19, 2022

1 CS 1656 – Introduction to Data Science

1.1 Instructor: Xiaowei Jia / Teaching Assistant: Mert Inan

1.1.1 Additional credits: E. Karageorgos, Xiaoting Li, Tahereh Arabghalizi, Evangelos Karageorgos, Zuha Agha, Anatoli Shein, Phuong Pham

1.2 Recitation 4: Regression and Decision Trees

In this recitation, we will learn how to do regression and classification with decision trees using scikit-learn python package.

```
[1]: import pandas as pd
import numpy as np
from sklearn import linear_model, tree, metrics
import matplotlib.pyplot as plt
%matplotlib inline
```

1.3 Linear Regression

LinearRegression fits a linear model with coefficients $w = (w_1, \dots, w_p)$ to minimize the residual sum of squares between the observed responses in the dataset, and the responses predicted by the linear approximation.

LinearRegression will take in its fit method arrays X, y and will store the coefficients w of the linear model in its `coef_` member.

We will now go through an example of linear regression on bike sharing dataset.

```
[2]: df = pd.read_csv('http://data.cs1656.org/bike_share.csv')
df.head()
```

```
[2]:   instant  season  hr  holiday  weekday  workingday  weathersit  temp  \
0         1       1   0        0         6           0           1  0.24
1         2       1   1        0         6           0           1  0.22
2         3       1   2        0         6           0           1  0.22
3         4       1   3        0         6           0           1  0.24
4         5       1   4        0         6           0           1  0.24

      temp_feels  hum  windspeed  cnt
0         0.2879  0.81         0.0   16
```

1	0.2727	0.80	0.0	40
2	0.2727	0.80	0.0	32
3	0.2879	0.75	0.0	13
4	0.2879	0.75	0.0	1

The attributes of the dataset are as follows:

- instant: record index
- season : season (1:spring, 2:summer, 3:fall, 4:winter)
- hr : hour (0 to 23)
- holiday : weather day is holiday or not (extracted from <http://dchr.dc.gov/page/holiday-schedule>)
- weekday : day of the week (0 to 6)
- workingday : if day is neither weekend nor holiday is 1, otherwise is 0.
- + weathersit :
 - 1: Clear
 - 2: Misty, Cloudy
 - 3: Light Snow, Light Rain
 - 4: Heavy Rain, Ice Pallets
- temp : Normalized temperature in Celsius. The values are divided by 41 (max)
- temp_feels: Normalized feeling temperature in Celsius. The values are divided by 50 (max)
- hum: Normalized humidity. The values are divided by 100 (max)
- windspeed: Normalized wind speed. The values are divided by 67 (max)
- cnt: count of total rental bikes including both casual and registered

Our target variable, y , is *cnt*. We will use a single attribute as input feature for this example and will select *temp* as our input feature X . You will be using all attributes in one of your tasks.

1.3.1 Subsample

As our dataset consists of more than 17000 rows, we will randomly subsample our dataset to select 1000 rows.

```
[3]: df_subsample = df.sample(1000)
df_subsample.head()
```

```
[3]:
```

	instant	season	hr	holiday	weekday	workingday	weathersit	temp	\
2853	2854	2	19	0	2	1	1	0.68	
1642	1643	1	23	0	0	0	1	0.26	
14565	14566	3	2	0	2	1	1	0.68	
15383	15384	4	4	1	1	0	1	0.32	
449	450	1	18	0	4	1	2	0.26	

	temp_feels	hum	windspeed	cnt
2853	0.6364	0.57	0.3582	365
1642	0.2576	0.65	0.2239	28
14565	0.6364	0.82	0.1343	5
15383	0.3333	0.87	0.1343	7
449	0.2576	0.56	0.1940	185

1.3.2 Train & Test Split

We will split the subsample into 90% training set and 10% test set by slicing the first 900 rows for training and using the rest for testing. As fit takes numpy arrays as input we will use *values* function to convert our Dataframe column into numpy array and use double brackets in order to make the arrays two-dimensional.

```
[4]: train = df_subsample.iloc[0:900]
      train_x = train[['temp']].values
      train_y = train[['cnt']].values

      test=df_subsample.iloc[900:]
      test_x = test[['temp']].values
      test_y = test[['cnt']].values
      print (type(train_x), type(train_y), type(test_x), type(test_y))

<class 'numpy.ndarray'> <class 'numpy.ndarray'> <class 'numpy.ndarray'> <class
'numpy.ndarray'>
```

1.3.3 Fit

To fit our linear regression model, apply the following function. Note that the fit function takes numpy array of the format [num_samples,num_features].

```
[5]: # Create linear regression object
      regr = linear_model.LinearRegression()

      # Train the model using the training sets
      regr.fit(train_x, train_y)
```

```
[5]: LinearRegression()
```

Now that we have fit our linear regression model onto the training data, our estimated model coefficients are stored in *coeff* attribute of our model.

```
[6]: # The coefficients
      print('Coefficients: \n', regr.coef_)
```

```
Coefficients:
[[367.19376965]]
```

1.3.4 Predict

We will now use our trained linear regression model to make predictions on our test set. Our model will take temperature attribute, *temp*, of our test data and will make predictions on the count of people who are bike sharing, given by *cnt*.

```
[7]: predict_y = regr.predict(test_x)
      # Printing predicted and actual values side by side for comparison
      np.column_stack((predict_y,test_y))
```

```

[7]: array([[231.90962664, 1.      ],
           [ 70.34436799, 3.      ],
           [129.09537113, 228.     ],
           [ 55.6566172 , 88.      ],
           [253.94125282, 21.      ],
           [136.43924653, 70.      ],
           [202.53412506, 121.     ],
           [121.75149574, 6.       ],
           [246.59737742, 167.     ],
           [202.53412506, 251.     ],
           [ 85.03211878, 15.      ],
           [143.78312192, 142.     ],
           [209.87800046, 724.     ],
           [202.53412506, 83.      ],
           [107.06374496, 10.      ],
           [253.94125282, 47.      ],
           [187.84637428, 471.     ],
           [298.00450517, 236.     ],
           [ 85.03211878, 36.      ],
           [253.94125282, 63.      ],
           [173.15862349, 46.      ],
           [ 63.0004926 , 219.     ],
           [129.09537113, 4.       ],
           [114.40762035, 147.     ],
           [173.15862349, 296.     ],
           [ 85.03211878, 34.      ],
           [261.28512821, 350.     ],
           [129.09537113, 119.     ],
           [151.12699731, 207.     ],
           [202.53412506, 70.      ],
           [283.31675439, 218.     ],
           [275.97287899, 482.     ],
           [224.56575124, 353.     ],
           [180.50249889, 126.     ],
           [ 77.68824338, 75.      ],
           [195.19024967, 3.       ],
           [217.22187585, 190.     ],
           [261.28512821, 617.     ],
           [114.40762035, 81.      ],
           [298.00450517, 273.     ],
           [209.87800046, 770.     ],
           [239.25350203, 133.     ],
           [129.09537113, 238.     ],
           [231.90962664, 148.     ],
           [253.94125282, 141.     ],
           [253.94125282, 296.     ],
           [231.90962664, 142.     ],

```

[121.75149574,	55.],
[202.53412506,	10.],
[114.40762035,	274.],
[275.97287899,	183.],
[165.8147481 ,	12.],
[253.94125282,	436.],
[342.06775753,	210.],
[261.28512821,	340.],
[187.84637428,	334.],
[231.90962664,	562.],
[283.31675439,	371.],
[173.15862349,	301.],
[224.56575124,	563.],
[246.59737742,	5.],
[114.40762035,	139.],
[92.37599417,	15.],
[107.06374496,	25.],
[99.71986956,	520.],
[107.06374496,	55.],
[107.06374496,	3.],
[209.87800046,	274.],
[246.59737742,	19.],
[165.8147481 ,	12.],
[224.56575124,	86.],
[114.40762035,	8.],
[202.53412506,	56.],
[261.28512821,	112.],
[224.56575124,	513.],
[158.47087271,	644.],
[85.03211878,	39.],
[173.15862349,	113.],
[224.56575124,	727.],
[151.12699731,	313.],
[63.0004926 ,	9.],
[173.15862349,	150.],
[165.8147481 ,	33.],
[253.94125282,	77.],
[136.43924653,	138.],
[239.25350203,	235.],
[136.43924653,	44.],
[129.09537113,	380.],
[151.12699731,	21.],
[246.59737742,	686.],
[231.90962664,	92.],
[136.43924653,	30.],
[268.6290036 ,	192.],
[312.69225596,	317.],

```
[151.12699731, 208.      ],
[ 99.71986956, 181.      ],
[195.19024967,  57.      ],
[275.97287899, 624.      ],
[129.09537113, 402.      ],
[ 99.71986956,  34.      ]])
```

1.3.5 Mean Squared Error

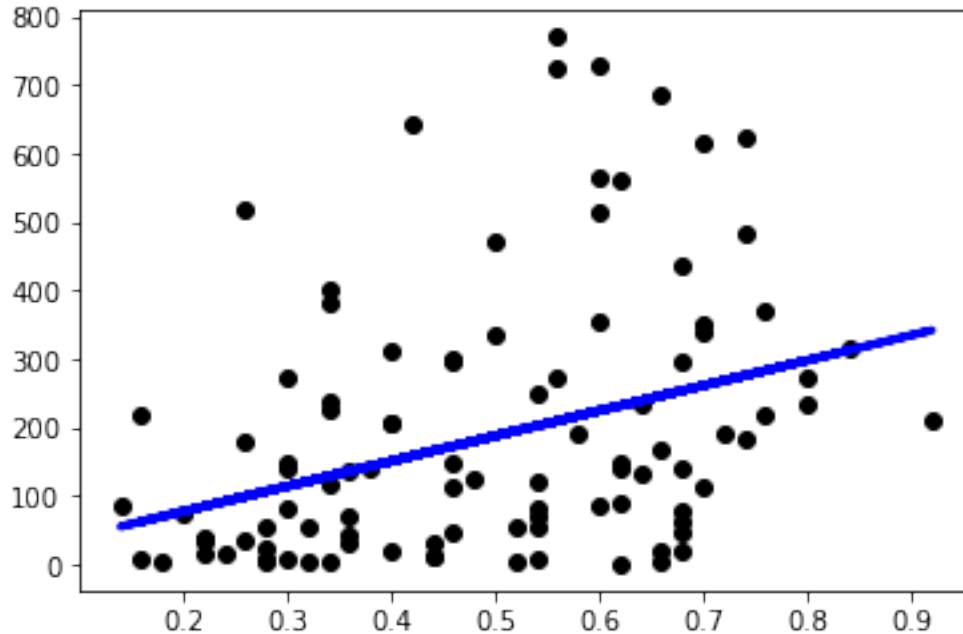
Looks like some of our model's predictions are not good. Now, let's measure the difference between our predicted and actual values by calculating the mean squared error.

```
[8]: meansq_error = np.mean((predict_y - test_y) ** 2)
     print ("Mean squared error: %.2f" % meansq_error)
```

Mean squared error: 33605.67

As expected our mean squared error is high, which means our model is not good. Can we improve it? What if we use more training data? Or more features? ### Plot We can also visualize the difference between our predictions and actual values by plotting.

```
[9]: plt.scatter(test_x, test_y, color='black', linewidth=1)
     plt.plot(test_x, predict_y, color='blue', linewidth=3)
     plt.show()
```



1.4 Decision Trees

Decision Trees (DTs) are a non-parametric supervised learning method used for classification and regression. The goal is to create a model that predicts the value of a target variable by learning simple decision rules inferred from the data features.

We will go through an example of binary classification using decision trees on titanic survival dataset.

```
[10]: dt = pd.read_csv('http://data.cs1656.org/titanic.csv')
      dt.head()
```

```
[10]:
```

	PassengerId	Survived	Pclass	\
0	1	0	3	
1	2	1	1	
2	3	1	3	
3	4	1	1	
4	5	0	3	

	Name	Sex	Age	SibSp	\
0	Braund, Mr. Owen Harris	male	22.0	1	
1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	
2	Heikkinen, Miss. Laina	female	26.0	0	
3	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	
4	Allen, Mr. William Henry	male	35.0	0	

	Parch	Embarked
0	0	S
1	0	C
2	0	S
3	0	S
4	0	S

The attributes of the dataset are as follows: - survival Survival (0 = No; 1 = Yes) - pclass Passenger Class (1 = 1st; 2 = 2nd; 3 = 3rd) - name Name - sex Sex - age Age - sibsp Number of Siblings/Spouses Aboard - parch Number of Parents/Children Aboard - embarked Port of Embarkation (C = Cherbourg; Q = Queenstown; S = Southampton) Our target class variable is *Survived*, whether the passenger survived or not. We will use only a subset of attributes that take discrete values to build our decision tree.

To fit a decision tree model, we will have to convert the categorical values into numerical values. As the only categorical attribute we will use is *Sex*, we will only need to convert that column into numerical values using the following commands.

```
[11]: dt['Sex'] = dt['Sex'].replace(['female', 'male'], [1, 2])
      dt.head()
```

```
[11]:
```

	PassengerId	Survived	Pclass	\
0	1	0	3	
1	2	1	1	
2	3	1	3	

```

3          4          1          1
4          5          0          3

```

	Name	Sex	Age	SibSp	Parch	\
0	Braund, Mr. Owen Harris	2	22.0	1	0	
1	Cumings, Mrs. John Bradley (Florence Briggs Th...	1	38.0	1	0	
2	Heikkinen, Miss. Laina	1	26.0	0	0	
3	Futrelle, Mrs. Jacques Heath (Lily May Peel)	1	35.0	1	0	
4	Allen, Mr. William Henry	2	35.0	0	0	

	Embarked
0	S
1	C
2	S
3	S
4	S

1.4.1 Train & Test Split

We will split our data into train and test set using the first 800 rows for training and the rest for testing.

```

[12]: dt_train_x = dt.iloc[:800][['Pclass', 'Sex', 'SibSp']].values
      dt_train_y = dt.iloc[:800][['Survived']].values

      dt_test_x = dt.iloc[800:][['Pclass', 'Sex', 'SibSp']].values
      dt_test_y = dt.iloc[800:][['Survived']].values

```

1.4.2 Fit

We will now fit our decision tree model onto the training set.

```

[13]: clf = tree.DecisionTreeClassifier()
      clf = clf.fit(dt_train_x, dt_train_y)

```

1.4.3 Predict

```

[14]: dt_predict_y = clf.predict(dt_test_x)
      ## comparing predicted and actual values
      np.column_stack((dt_predict_y, dt_test_y))

```

```

[14]: array([[0, 0],
            [1, 1],
            [0, 1],
            [0, 1],
            [0, 1],
            [0, 0],
            [0, 0],

```


[1, 0],
[0, 0],
[1, 1],
[0, 0],
[0, 0],
[0, 0],
[0, 0],
[0, 0],
[0, 0],
[1, 0],
[0, 0],
[0, 0],
[0, 0],
[1, 1],
[0, 1],
[0, 0],
[1, 1],
[0, 0],
[0, 0],
[0, 0],
[0, 1],
[0, 1],
[1, 1],
[0, 1],
[0, 1],
[0, 0],
[0, 0],
[0, 0],
[1, 1],
[0, 0],
[0, 0],
[0, 1],
[0, 1],
[0, 0],
[0, 0],
[1, 1],
[0, 0],
[0, 0],
[0, 0],
[0, 0],
[0, 0],
[0, 0],
[1, 1],
[0, 0],
[0, 0],
[0, 0],
[1, 1],

```

[1, 0],
[1, 1],
[1, 1],
[0, 1],
[1, 1],
[0, 0],
[0, 0],
[0, 0],
[1, 1],
[0, 0],
[0, 0],
[1, 1],
[1, 1],
[0, 0],
[0, 0],
[0, 1],
[0, 0],
[1, 1],
[0, 0],
[0, 0],
[1, 1],
[1, 1],
[0, 0],
[0, 0],
[0, 0],
[1, 1],
[1, 1],
[0, 0],
[1, 0],
[0, 0],
[0, 0],
[1, 0],
[0, 0],
[1, 1],
[0, 0],
[0, 1],
[0, 0]], dtype=int64)

```

1.4.4 Accuracy

We can measure the accuracy of our prediction by using the following commands.

```

[15]: accuracy = metrics.accuracy_score(dt_test_y,dt_predict_y)
accuracy

```

```

[15]: 0.8021978021978022

```

1.5 Tasks

Task 1

Do linear regression over a sample of 1000 rows of bike share counts, *cnt*, using *weekday*, as input feature. Calculate the mean squared error by using first 900 rows for training and the rest for testing. Return the mean squared error.

[]:

Task 2.1

Repeat Task 1 using all attributes except *instant* (also, scatter plot is not required in this task). Is the mean squared error higher or lower? Is it better to use all attributes?

[]:

Task 2.2

Comparing the results of task 1 and task 2.1, is it better to use all attributes? Why?

Task 3

You will use *bank-data.csv* as input for this task. Use decision trees to do binary classification of *mortgage*{yes,no} using *region*, *sex* and *married* attributes as input features. Use the first 500 rows for training and the rest for testing. Measure the accuracy of your classification. Return the accuracy.

[]: