

DS 6003 Assignment 1

Jiangxue Han (jh6rg)

Motivation

- The dataset contains 1000 camera models with 18 features including basic information and price.
- The model is designed to find the relationship between camera's price and camera's property.
- Price is the target feature. Both classification and regression techniques have been implemented. For the classification, a new variable is created indicating whether the price of corresponding camera is above the average price.
- download link: <https://www.kaggle.com/crawford/1000-cameras-dataset/>

Code snippet and explanation

- Create a new feature “isExpensive” indicating the whether the price is above average

```
In [55]: 1 from pyspark.sql.functions import *
2 df = df.withColumn('isExpensive', when(df.Price > 458, 1).otherwise(0)).drop(df.Price)
3 df = df.drop(df.Model).dropna()
4 df

Out[55]: DataFrame[date: bigint, Max_Resolution: double, Low_Resolution: double, Effective_pixels: double, Zoom_wide: double,
Zoom_tele: double, Normal_focus_range: double, Macro_focus_range: double, Storage_included: double, Weight: double, D
imensions: double, isExpensive: int]
```

- Correlation analysis: It seems the features don't have a strong relationship with the target feature.

```
In [20]: 1 print("Pearson's r(Effective pixels,Max resolution) = {}".format(df.corr("Effective_pixels", "Max_Resolution")))
2 print("Pearson's r(Normal focus range,Min resolution) = {}".format(df.corr("Normal_focus_range", "Low_Resolution")))
3 print("Pearson's r(Dimensions,Weight) = {}".format(df.corr("Dimensions", "Weight")))

Pearson's r(Effective pixels,Max resolution) = 0.9538458570187706
Pearson's r(Normal focus range,Min resolution) = -0.12543611622804537
Pearson's r(Dimensions,Weight) = 0.6778848089757021
```

```
In [63]: 1 print("Pearson's r(Date,Price) = {}".format(df.corr("date", "Price")))
2 print("Pearson's r(Max Resolution,Price) = {}".format(df.corr("Max_Resolution", "Price")))
3 print("Pearson's r(Low Resolution,Price) = {}".format(df.corr("Low_Resolution", "Price")))
4 print("Pearson's r(Zoom wide,Price) = {}".format(df.corr("Zoom_wide", "Price")))
5 print("Pearson's r(Normal focus range,Price) = {}".format(df.corr("Normal_focus_range", "Price")))
6 print("Pearson's r(Macro focus range,Price) = {}".format(df.corr("Macro_focus_range", "Price")))
7 print("Pearson's r(Effective pixels,Price) = {}".format(df.corr("Effective_pixels", "Price")))
8 print("Pearson's r(Dimension,Price) = {}".format(df.corr("Dimensions", "Price")))

Pearson's r(Date,Price) = -0.022522217044592785
Pearson's r(Max Resolution,Price) = 0.18420092652935718
Pearson's r(Low Resolution,Price) = 0.15420406063216538
Pearson's r(Zoom wide,Price) = -0.4591033999208599
Pearson's r(Normal focus range,Price) = -0.27385393615399717
Pearson's r(Macro focus range,Price) = -0.12757671236241522
Pearson's r(Effective pixels,Price) = 0.190284007565066
Pearson's r(Dimension,Price) = 0.2642561651698723
```

- Build Linear Regression model

```
7]: 1 df = df.drop(df.Low_Resolution).drop(df.Effective_pixels)

38]: 1 stages = []
2 df = df.withColumnRenamed("Price", "label")
3 assemblerInputs = df.columns[:-1]
4 assembler = VectorAssembler(inputCols=assemblerInputs, outputCol="features")
5 stages += [assembler]

39]: 1 pipeline = Pipeline(stages = stages)
2 pipelineModel = pipeline.fit(df)
3 df3 = pipelineModel.transform(df)
4 selectedCols = ['label', 'features']
5 df3 = df3.select(selectedCols)

10]: 1 # create train/test sets
2 seed = 42
3 (testDF, trainingDF) = df3.randomSplit((0.20, 0.80), seed=seed)
4 print ('training set N = {}, test set N = {}'.format(trainingDF.count(), testDF.count()))

training set N = 843, test set N = 193
```

```
11]: 1 lr = LinearRegression()
2 lrModel = lr.fit(trainingDF)
```

```
12]: 1 predictionsAndLabelsDF = lrModel.transform(testDF)
```

- Build Logistic Regression model

```
In [152]: 1 df.groupBy().avg('label').collect()

Out[152]: [Row(avg(label)=457.9218146718147)]

In [154]: 1 # Create a new feature "isExpensive" indicating whether Price(label) is above average
2 df = df.withColumn('isExpensive', when(df.label > 458, 1).otherwise(0)).drop(df.label)
3 df

Out[154]: DataFrame[date: bigint, Max_Resolution: double, Zoom_wide: double, Zoom_tele: double, Normal_focus_range: double, Mac
ro_focus_range: double, Storage_included: double, Weight: double, Dimensions: double, isExpensive: int]

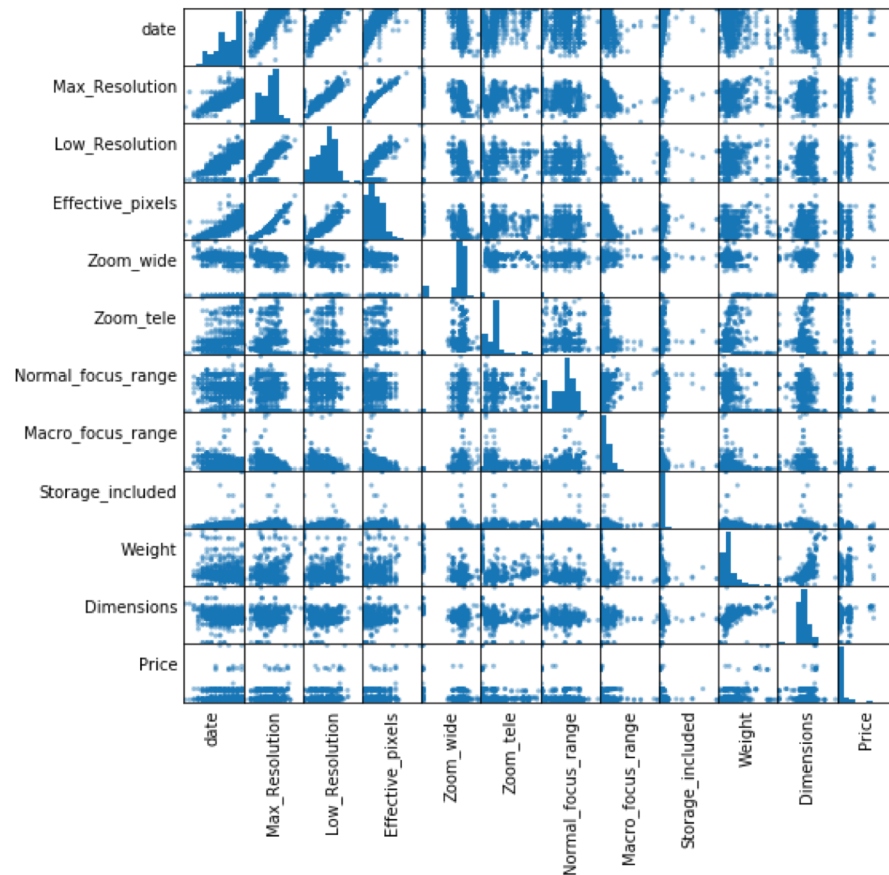
In [161]: 1 stages = []
2 label_stringIdx = StringIndexer(inputCol = 'isExpensive', outputCol = 'label')
3 stages += [label_stringIdx]
4 assemblerInputs = df.columns[:-3]
5 assembler = VectorAssembler(inputCols=assemblerInputs, outputCol="features")
6 stages += [assembler]

In [168]: 1 pipeline = Pipeline(stages = stages)
2 pipelineModel = pipeline.fit(df)
3 df3 = pipelineModel.transform(df)
4 selectedCols = ['label', 'features']
5 df3 = df3.select(selectedCols)

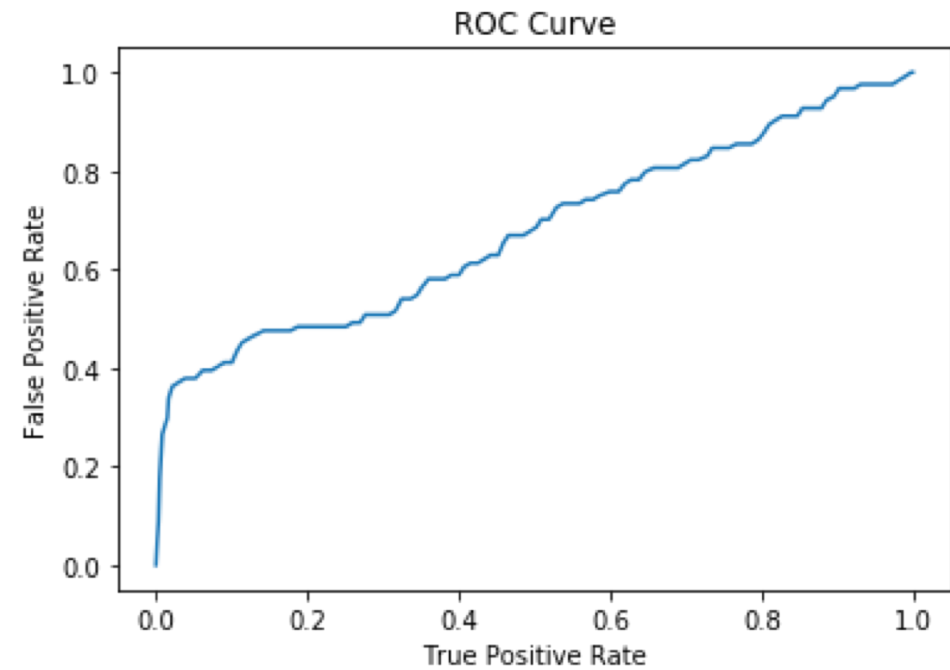
In [169]: 1 trainDF, testDF = df3.randomSplit([0.7, 0.3], seed = 2019)
2 print ('training set N = {}, test set N = {}'.format(train_df.count(), test_df.count()))
```

Visualizations

Correlation map



ROC Curve



AUC: 0.64