

Mushroom Classification



DS 6003 Spark Homework
Runhao Zhao
rz6dg

Motivation



As it is recorded, there are more than 10000 described mushroom species exist worldwide. Many types of mushrooms can be identified as poisonous mushrooms by their shapes, sizes or colors. I am interested in applying various machine learning methods to predict if a mushroom is safe to eat based on multiple observed features. Learning which characteristics spell death and which are most palatable is also important and helpful in identifying edible/poisonous mushrooms. Ultimately, I hope this could be used to help mushroom lovers avoid poisonous mushrooms when spending time outdoors.

Code snippet

```
1 #create a logistic regression model
2 from pyspark.ml.classification import LogisticRegression
3 lr = LogisticRegression(featuresCol = 'features', labelCol = 'label', maxIter=10)
4 lrModel = lr.fit(train)
5 trainingSummary = lrModel.summary
6 roc = trainingSummary.roc.toPandas()
7 plt.plot(roc['FPR'],roc['TPR'])
8 plt.ylabel('False Positive Rate')
9 plt.xlabel('True Positive Rate')
10 plt.title('ROC Curve')
11 plt.show()
```

The codes above are building a logistic regression model on the data set to predict the edibility of a mushroom. The features used in the model include all 28 features in the data set. An ROC curve is created to evaluate the model performance

```
2 for col_name in df1.columns:
3     indexer = StringIndexer(inputCol=col_name, outputCol="new_"+col_name)
4     model = indexer.fit(df1)
5     df1 = model.transform(df1)
```

The codes above are looping through each column and applying label encoding to transform non-numerical labels to numerical labels.

Visualizations

