

Assignment1

```

1 //Importing library
2 import com.johnsnowlabs.nlp.base._
3 import com.johnsnowlabs.nlp.annotator._
4 import com.johnsnowlabs.nlp.pretrained.PretrainedPipeline
5 import com.johnsnowlabs.nlp.SparkNLP
6 import org.apache.spark.sql.functions._
7 import org.apache.spark.sql.Row

```

Cmd 2

```

1 //Reading the file and filtering
2 val txtfile = sc.textFile("/FileStore/tables/largestttt56613_0-1.txt")
3 val dataDF =
  txtfile.filter(x=>x.length>0).zipWithIndex.toDF("text","_id").select("_id",
  "text") //Filtering
4 dataDF.take(10)

```

txtfile: org.apache.spark.rdd.RDD[String] = /FileStore/tables/largestttt56613_0-1.txt MapPartitionsRDD[143] at textFile at command-2187602793455931:1

Command took 0.34 seconds -- by kxn190007@utdallas.edu at 9/23/2020, 5:57:17 PM on UTD1



Cmd 3

```

1 //Using Library
2 val pipeline = PretrainedPipeline("recognize_entities_dl", "en")
3 val predictions = pipeline.transform(dataDF)
4 pipeline.transform(dataDF).select("entities.result")
5 val output =
  pipeline.transform(dataDF).select(explode($"entities.result"))
6 output.show()

```

► (1) Spark Jobs

-  predictions: org.apache.spark.sql.DataFrame = [_id: long, text: string ... 6 more fields]
-  output: org.apache.spark.sql.DataFrame = [col: string]

```

+-----+
|               col|
+-----+
|Project Gutenberg...|
|           English|
|       Andrew Lang|
|       United States|
|Project Gutenberg...|
|             eBook|
|       United States|

```