

Data Transformation Class

```
1  ## Author: Goh Boon Xiang
2  from pyspark.ml.feature import StringIndexer, OneHotEncoder
3  from pyspark.sql import DataFrame, SparkSession
4  from pyspark.sql.functions import udf, col
5  from pyspark.sql.types import ArrayType, StringType
6  from pyspark.ml.feature import Tokenizer, HashingTF, IDF
7  import pyspark.sql.functions as F
8  from nltk.stem import WordNetLemmatizer
9
10 class DataTransformations:
11
12     @staticmethod
13     def tokenize(df: DataFrame, input_col: str, output_col: str) -> DataFrame:
14         tokenizer = Tokenizer(inputCol=input_col, outputCol=output_col)
15         return tokenizer.transform(df)
16
17     @staticmethod
18     def lemmatize_tokens(df: DataFrame, tokens_col: str) -> DataFrame:
19         lemmatizer = WordNetLemmatizer()
20
21         def lemmatize(tokens):
22             return [lemmatizer.lemmatize(token) for token in tokens]
23
24         lemmatize_udf = udf(lemmatize, ArrayType(StringType()))
25         return df.withColumn(tokens_col, lemmatize_udf(col(tokens_col)))
```

```

27     @staticmethod
28     def calculate_tfidf(df: DataFrame, tokens_col: str, raw_features_col: str =
"rawFeatures", features_col: str = "features", num_features: int = 10000) -> DataFrame:
29         hashingTF = HashingTF(inputCol=tokens_col, outputCol=raw_features_col,
numFeatures=num_features)
30         featurized_df = hashingTF.transform(df)
31
32         idf = IDF(inputCol=raw_features_col, outputCol=features_col)
33         idf_model = idf.fit(featurized_df)
34         return idf_model.transform(featurized_df)
35
36     @staticmethod
37     def oversample(df: DataFrame, label_col: str, majority_label: int, minority_label:
int) -> DataFrame:
38         major_df = df.filter(col(label_col) == majority_label)
39         minor_df = df.filter(col(label_col) == minority_label)
40
41         major_count = major_df.count()
42         minor_count = minor_df.count()
43
44         if minor_count > 0:
45             ratio = int(major_count / minor_count)
46             if ratio > 1:
47                 oversampled_df = minor_df.withColumn("dummy",
F.explode(F.array([F.lit(x) for x in range(ratio)]))).drop('dummy')
48                 return major_df.union(oversampled_df)
49             else:
50                 return df
51         else:
52             return df
53
54     @staticmethod
55     def one_hot_encode(df: DataFrame, categorical_col: str) -> DataFrame:
56         """Perform one-hot encoding on the specified categorical column."""
57         # Index the categorical column
58         indexer = StringIndexer(inputCol=categorical_col, outputCol=f"
{categorical_col}_index")
59         indexed_df = indexer.fit(df).transform(df)
60
61         # One-Hot Encode the indexed column
62         encoder = OneHotEncoder(inputCols=[f"{categorical_col}_index"], outputCols=[f"
{categorical_col}_encoded"])
63         encoded_df = encoder.fit(indexed_df).transform(indexed_df)
64
65         return encoded_df

```