

Table of contents

| | | |
|----------|---|----------|
| 1 | Data types | 1 |
| 1.1 | Local vectors | 2 |
| 1.2 | Local matrices | 2 |
| 2 | Main concepts | 3 |
| 2.1 | Transformer | 4 |
| 2.2 | Estimator | 4 |
| 2.3 | Pipeline | 4 |
| 2.4 | Parameters | 5 |
| 3 | Data preprocessing | 5 |
| 3.1 | Extracting, transformings, and selecting features | 5 |
| 4 | Feature transformations | 6 |
| 4.1 | VectorAssembler | 6 |
| 4.2 | Data Normalization | 7 |
| 4.2.1 | StandardScaler | 8 |
| 4.3 | Categorical columns | 10 |
| 4.3.1 | StringIndexer | 11 |
| 4.3.2 | IndexToString | 12 |
| 4.3.3 | SQLTransformer | 15 |

Spark MLlib is the Spark component providing the machine learning/data mining algorithms

- Pre-processing techniques
- Classification (supervised learning)
- Clustering (unsupervised learning)
- Itemset mining

MLlib APIs are divided into two packages:

- `pyspark.mllib`: It contains the original APIs built on top of RDDs. This version of the APIs is in maintenance mode and will be probably deprecated in the next releases of Spark.
- `pyspark.ml`: It provides higher-level API built on top of DataFrames (i.e, `Dataset<Row>`) for constructing ML pipelines. It is recommended because the DataFrame-based API is more versatile and flexible, also providing the pipeline concept. This is the one explained in this course.

1 Data types

Spark MLlib is based on a set of basic local and distributed data types:

- Local vector
- Local matrix
- Distributed matrix
- ...

DataFrames for ML-based applications contain objects based on these basic data types.

1.1 Local vectors

Local `pyspark.ml.linalg.Vector` objects in MLlib are used to store vectors (in dense and sparse representations) of double values. The MLlib algorithms work on vectors of doubles, used to represent the input records/data (one vector for each input record). Non double attributes/values must be mapped to double values before applying MLlib algorithms.

i Example

Consider the vector of doubles `[1.0, 0.0, 3.0]`. It can be represented

- In dense format as `[1.0, 0.0, 3.0]`
- In sparse format as `(3, [0, 2], [1.0, 3.0])`, where
 - 3 is the size of the vector
 - The array `[0, 2]` contains the indexes of the non-zero cells
 - The array `[1.0, 3.0]` contains the values of the non-zero cells

The following code shows how dense and sparse vectors can be created in Spark

```
1 from pyspark.ml.linalg import Vectors
2
3 # Create a dense vector [1.0, 0.0, 3.0]
4 dv = Vectors.dense([1.0, 0.0, 3.0])
5
6 # Create a sparse vector [1.0, 0.0, 3.0] by specifying
7 # its indices and values corresponding to non-zero entries
8 # by means of a dictionary
9 sv = Vectors.sparse(3, { 0:1.0, 2:3.0 })
```

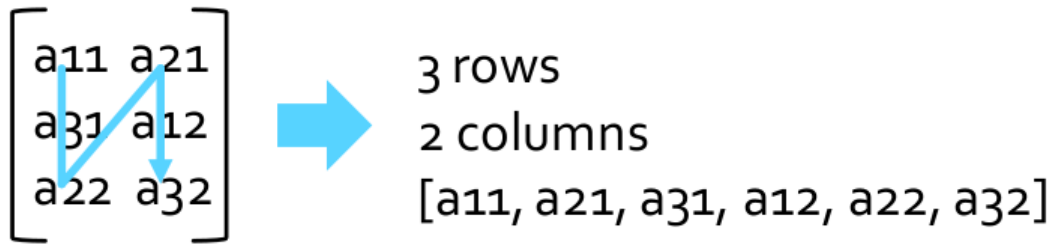
In the sparse vector

| | |
|------------------|--|
| 3 | Size of the vector |
| 2:3.0 | Index and value of a non-empty cell |
| { 0:1.0, 2:3.0 } | Dictionary of <i>index : value</i> pairs |

1.2 Local matrices

Local `pyspark.ml.linalg.Matrix` objects in MLlib are used to store matrices (in dense and sparse representations) of double values. The column-major order is used to store the content of the matrix in a linear way.

Figure 1: Local matrices

**i** Example

The following code shows how dense and sparse matrices can be created in Spark.

```
1 from pyspark.ml.linalg import Matrices
2
3 # Create a dense matrix with two rows and three columns
4 # 3.0 0.0 0.0
5 # 1.0 1.5 2.0
6 dm = Matrices.dense(2,3,[3.0, 1.0, 0.0, 1.5, 0.0, 2.0])
7
8 # Create a sparse version of the same matrix
9 sm = Matrices.sparse(2,3, [0, 2, 3, 4], [0, 1, 1, 1] , [3, 1, 1.5, 2])
```

In the dense matrix vector

| | |
|--------------------------------|------------------------------|
| 2 | Number of rows |
| 3 | Number of columns |
| [3.0, 1.0, 0.0, 1.5, 0.0, 2.0] | Values in column/major order |

In the sparse matrix vector

| | |
|----------------|---|
| 2 | Number of rows |
| 3 | Number of columns |
| [0, 2, 3, 4] | One element per column that encodes the offset in the array of non-zero values where the values of the given column start. The last element is the number of non-zero values. |
| [0, 1, 1, 1] | Row index of each non-zero value |
| [3, 1, 1.5, 2] | Array of non-zero values of the represented matrix |

2 Main concepts

Spark MLlib uses DataFrames as input data: the input of the MLlib algorithms are structured data (i.e., tables), and all input data must be represented by means of tables before applying the MLlib

algorithms; also document collections must be transformed in a tabular format before applying the MLlib algorithms.

The DataFrames used and created by the MLlib algorithms are characterized by several columns, and each column is associated with a different role/meaning

- **label:** the target of a classification/regression analysis;
- **features:** the vector containing the values of the attributes/features of the input record/data points;
- **text:** the original text of a document before being transformed in a tabular format;
- **prediction:** the predicted value of a classification/regression analysis.

2.1 Transformer

A Transformer is an ML algorithm/procedure that transforms one DataFrame into another DataFrame by means of the method `transform(inputDataFrame)`.

i Example 1

A feature transformer might take a DataFrame, read a column (e.g., text), map it into a new column (e.g., feature vectors), and output a new DataFrame with the mapped column appended.

i Example 2

A classification model is a Transformer that can be applied on a DataFrame with features and transforms it into a DataFrame with also the prediction column.

2.2 Estimator

An Estimator is an ML algorithm/procedure that is fit on an input (training) DataFrame to produce a Transformer: each Estimator implements a `fit()` method, which accepts a DataFrame and produces a Model of type **Transformer**.

An Estimator abstracts the concept of a learning algorithm or any algorithm that fits/trains on an input dataset and returns a model

i Example

The Logistic Regression classification algorithm is an Estimator: calling `fit(input training DataFrame)` on it a Logistic Regression Model is built, which is a Model/a Transformer.

2.3 Pipeline

A Pipeline chains multiple Transformers and Estimators together to specify a Machine learning/Data Mining workflow. In a pipeline, the output of a transformer/estimator is the input of the next one.

i Example

A simple text document processing workflow aiming at building a classification model includes several steps

1. Split each document into a set of words;
2. Convert each set of words into a numerical feature vector;
3. Learn a prediction model using the feature vectors and the associated class labels.

2.4 Parameters

Transformers and Estimators share common APIs for specifying the values of their parameters.

In the new APIs of Spark MLlib the use of the pipeline approach is preferred/recommended. This approach is based on the following steps

1. The set of Transformers and Estimators that are needed are instantiated;
2. A pipeline object is created and the sequence of transformers and estimators associated with the pipeline are specified;
3. The pipeline is executed and a model is trained;
4. (optional) The model is applied on new data.

3 Data preprocessing

Input data must be preprocessed before applying machine learning and data mining algorithms

- To organize data in a format consistent with the one expected by the applied algorithms;
- To define good (predictive) features;
- To remove bias (e.g., normalization);
- To remove noise and missing values.

3.1 Extracting, transformings, and selecting features

MLlib provides a set of transformers than can be used to extract, transform and select features from DataFrames

- Feature Extractors (e.g., TF-IDF, Word2Vec)
- Feature Transformers (e.g., Tokenizer, StopWordsRemover, StringIndexer, IndexToString, OneHotEncoderEstimator, Normalizer)
- Feature Selectors (e.g., VectorSlicer)

See the up-to-date list [here](#).

4 Feature transformations

Several algorithms are provided by MLlib to transform features. They are used to create new columns/features by combining or transforming other features. It is possible to perform feature transformations and feature creations by using the standard methods for DataFrames and RDDs.

4.1 VectorAssembler

`VectorAssembler` (`pyspark.ml.feature.VectorAssembler`) is a transformer that combines a given list of columns into a single vector column. It is useful for combining features into a single feature vector before applying ML algorithms.

Given `VectorAssembler(inputCols, outputCol)`

- **inputCols**: the list of original columns to include in the new column of type `Vector`. The following input column types are accepted
 - all numeric types, boolean type, and vector type
 - Boolean values are mapped to 1 (True) and 0 (False)
- **outputCol**: the name of the new output column

When the transform method of `VectorAssembler` is invoked on a `DataFrame` the returned `DataFrame` has a new column (`outputCol`): for each record, the value of the new column is the concatenation of the values of the input columns. It has also all the columns of the input `DataFrame`.

i Example

Consider an input `DataFrame` with three columns: create a new `DataFrame` with a new column containing the concatenation of “colB” and “colC” in a new vector column. Set the name of the new column to “features”.

Original DataFrame

| colA | colB | colC |
|------|------|-------|
| 1 | 4.5 | True |
| 2 | 0.6 | True |
| 3 | 1.5 | False |
| 4 | 12.1 | True |
| 5 | 0.0 | True |

Transformed DataFrame

| colA | colB | colC | features |
|------|------|-------|-------------|
| 1 | 4.5 | True | [4.5, 1.0] |
| 2 | 0.6 | True | [0.6, 1.0] |
| 3 | 1.5 | False | [1.5, 0.0] |
| 4 | 12.1 | True | [12.1, 1.0] |
| 5 | 0.0 | True | [0.0, 1.0] |

Notice that columns of DataFrames can also be vectors.

```
1 from pyspark.mllib.linalg import Vectors
2 from pyspark.ml.feature import VectorAssembler
3
4 # input and output folders
5 inputPath = "data/exampleDataAssembler.csv"
6
7 # Create a DataFrame from the input data
8 inputDF = spark.read.load(
9     inputPath,
10    format="csv",
11    header=True,
12    inferSchema=True
13 )
14
15 # Create a VectorAssembler that combines columns colB and colC
16 # The new vector column is called features
17 myVectorAssembler = VectorAssembler(
18     inputCols = ['colB', 'colC'],
19     outputCol = 'features'
20 )
21
22 # Apply myVectorAssembler on the input DataFrame
23 transformedDF = myVectorAssembler.transform(inputDF)
```

4.2 Data Normalization

MLlib provides a set of normalization algorithms (called scalers)

- StandardScaler
- MinMaxScaler
- Normalizer
- MaxAbsScaler

4.2.1 StandardScaler

`StandardScaler` (`pyspark.ml.feature.StandardScaler`) is an Estimator that returns a Transformer (`pyspark.ml.feature.StandardScalerModel`). `StandardScalerModel` transforms a vector column of an input `DataFrame` normalizing each feature of the input vector column to have unit standard deviation and/or zero mean.

Given `StandardScaler(inputCol, outputCol)`

- `inputCol`: the name of the input vector column (of doubles) to normalize
- `outputCol`: the name of the new output normalized vector column

Invoke the `fit` method of `StandardScaler` on the input `DataFrame` to infer a `StandardScalerModel`. The returned model is a Transformer.

Invoke the `transform` method of `StandardScalerModel` on the input `DataFrame` to create a new `DataFrame` that has a new column (`outputCol`): for each record, the value of the new column is the normalized version of the input vector column. It has also all the columns of the input `DataFrame`.

i Example

Consider an input `DataFrame` with four columns: create a new `DataFrame` with a new column containing the normalized version of the vector column features. Set the name of the new column to “scaledFeatures”.

Original DataFrame

| colA | colB | colC | features |
|------|------|-------|-------------|
| 1 | 4.5 | True | [4.5, 1.0] |
| 2 | 0.6 | True | [0.6, 1.0] |
| 3 | 1.5 | False | [1.5, 0.0] |
| 4 | 12.1 | True | [12.1, 1.0] |
| 5 | 0.0 | True | [0.0, 1.0] |

Transformed DataFrame

| colA | colB | colC | features | scaledFeatures |
|------|------|-------|-------------|----------------|
| 1 | 4.5 | True | [4.5, 1.0] | [0.903, 2.236] |
| 2 | 0.6 | True | [0.6, 1.0] | [0.120, 2.236] |
| 3 | 1.5 | False | [1.5, 0.0] | [0.301, 0.0] |
| 4 | 12.1 | True | [12.1, 1.0] | [2.428, 2.236] |
| 5 | 0.0 | True | [0.0, 1.0] | [0.0, 2.236] |

```
1 from pyspark.mllib.linalg import Vectors
2 from pyspark.ml.feature import VectorAssembler
3 from pyspark.ml.feature import StandardScaler
4 # input and output folders
5 inputPath = "data/exampleDataAssembler.csv"
6
7 # Create a DataFrame from the input data
8 inputDF = spark.read.load(
9     inputPath,
10    format="csv",
11    header=True,
12    inferSchema=True
13 )
14
15 # Create a VectorAssembler that combines columns colB and colC
16 # The new vector column is called features
17 myVectorAssembler = VectorAssembler(
18     inputCols = ['colB', 'colC'],
19     outputCol = 'features'
20 )
21
22 # Apply myVectorAssembler on the input DataFrame
23 transformedDF = myVectorAssembler.transform(inputDF)
24
25 # Create a Standard Scaler to scale the content of features
26 myScaler = StandardScaler(
27     inputCol="features",
28     outputCol="scaledFeatures"
29 )
30
31 # Compute summary statistics by fitting the StandardScaler
32 # Before normalizing the content of the data we need to compute mean and
33 # standard deviation of the analyzed data
34 scalerModel = myScaler.fit(transformedDF)
35
36 # Apply myScaler on the input column features
37 scaledDF = scalerModel.transform(transformedDF)
```

4.3 Categorical columns

Frequently the input data are characterized by categorical attributes (i.e., string columns), and the class label of the classification problem is a categorical attribute. The Spark MLlib classification and regression algorithms work only with numerical values, so categorical columns must be mapped to double values.

4.3.1 StringIndexer

A `StringIndexer` (`pyspark.ml.feature.StringIndexer`) is an Estimator that returns a Transformer of type `pyspark.ml.feature.StringIndexerModel`. `StringIndexerModel` encodes a string column of “labels” to a column of “label indices”: each distinct value of the input string column is mapped to an integer value in $[0, \text{number of distinct values})$.

`StringIndexer(inputCol, outputCol)`

- `inputCol`: the name of the input string column to map to a set of integers
- `outputCol`: the name of the new output column

Invoke the `fit` method of `StringIndexer` on the input `DataFrame` to infer a `StringIndexerModel`. The returned model is a Transformer.

Invoke the `transform` method of `StringIndexerModel` on the input `DataFrame` to create a new `DataFrame` that has a new column (`outputCol`): for each record, the value of the new column is the integer (casted to a double) associated with the value of the input string column. It has also all the columns of the input `DataFrame`.

i Example

Consider an input `DataFrame` with two columns: create a new `DataFrame` with a new column containing the integer version of the string column category. Set the name of the new column to “categoryIndex”.

Original DataFrame

| id | category |
|----|----------|
| 1 | a |
| 2 | b |
| 3 | c |
| 4 | c |
| 5 | a |

Transformed DataFrame

| id | category | categoryIndex |
|----|----------|---------------|
| 1 | a | 0.0 |
| 2 | b | 2.0 |
| 3 | c | 1.0 |
| 4 | c | 1.0 |
| 5 | a | 0.0 |

```

1  from pyspark.mllib.linalg import Vectors
2  from pyspark.ml.feature import StringIndexer
3
4  # input DataFrame
5  df = spark.createDataFrame(
6      [(1,"a"),(2,"b"),(3,"c"),(4,"c"),(5,"a")],
7      ["id","category"]
8  )
9
10 # Create a StringIndexer to map the content of category
11 # to a set of "integers"
12 indexer = StringIndexer(
13     inputCol="category",
14     outputCol="categoryIndex"
15 )
16
17 # Analyze the input data to define the mapping string -> integer
18 indexerModel = indexer.fit(df)
19
20 # Apply indexerModel on the input column category
21 indexedDF = indexerModel.transform(df)

```

4.3.2 IndexToString

`IndexToString` (`pyspark.ml.feature.IndexToString`), which is symmetrical to `StringIndexer`, is a Transformer that maps a column of “label indices” back to a column containing the original “labels” as strings. Classification models return the integer version of the predicted label values. To obtain human readable results, remap those values to the original ones.

Given `IndexToString(inputCol, outputCol, labels)`

- `inputCol`: the name of the input numerical column to map to the original a set of string “labels”;
- `outputCol`: the name of the new output column;
- `labels`: the list of original “labels”/strings; the mapping with integer values is given by the positions of the strings inside labels.

Invoke the transform method of `IndexToString` on the input `DataFrame` to create a new `DataFrame` that has a new column (`outputCol`): for each record, the value of the new column is the original string associated with the value of the input numerical column. It has also all the columns of the input `DataFrame`.

i Example

Consider an input `DataFrame` with two columns: create a new `DataFrame` with a new column containing the integer version of the string column category and then map it back to the string version in a new column.

Original DataFrame

| id | category |
|----|----------|
| 1 | a |
| 2 | b |
| 3 | c |
| 4 | c |
| 5 | a |

Transformed DataFrame

| id | category | categoryIndex | originalCategory |
|----|----------|---------------|------------------|
| 1 | a | 0.0 | a |
| 2 | b | 2.0 | b |
| 3 | c | 1.0 | c |
| 4 | c | 1.0 | c |
| 5 | a | 0.0 | a |

```

1 from pyspark.mllib.linalg import Vectors
2 from pyspark.ml.feature import StringIndexer
3 from pyspark.ml.feature import IndexToString
4
5 # input DataFrame
6 df = spark.createDataFrame(
7     [(1,"a"),(2,"b"),(3,"c"),(4,"c"),(5,"a")],
8     ["id","category"]
9 )
10
11 # Create a StringIndexer to map the content of category
12 # to a set of integers
13 indexer = StringIndexer(
14     inputCol="category",
15     outputCol="categoryIndex"
16 )
17
18 # Analyze the input data to define the mapping
19 # string -> integer
20 indexerModel = indexer.fit(df)
21
22 # Apply indexerModel on the input column category
23 indexedDF = indexerModel.transform(df)
24
25 # Create an IndexToString to map the content of numerical
26 # attribute categoryIndex to the original string value
27 converter = IndexToString(
28     inputCol="categoryIndex",
29     outputCol="originalCategory",
30     labels=indexerModel.labels
31 )
32
33 # Apply converter on the input column categoryIndex
34 reconvertedDF = converter.transform(indexedDF)

```

4.3.3 SQLTransformer

`SQLTransformer` (`pyspark.ml.feature.SQLTransformer`) is a transformer that implements the transformations which are defined by SQL queries. Currently, the syntax of the supported (simplified) SQL queries is

```
1 SELECT attributes, function(attributes) FROM __THIS__
```

Where `__THIS__` represents the `DataFrame` on which the `SQLTransformer` is invoked.

`SQLTransformer` executes an SQL query on the input `DataFrame` and returns a new `DataFrame` associated with the result of the query.

Given `SQLTransformer(statement)`

- **statement:** the SQL query to execute.

When the `transform` method of `SQLTransformer` is invoked on a `DataFrame` the returned `DataFrame` is the result of the executed statement query.

i Example

Consider an input `DataFrame` with two columns: “text” and “id”: create a new `DataFrame` with a new column, called “numWords”, containing the number of words occurring in column “text”.

Original DataFrame

| id | text |
|----|----------------------------------|
| 1 | This is Spark |
| 2 | Spark |
| 3 | Another sample sentence of words |
| 4 | Paolo Rossi |
| 5 | Giovanni |

Transformed DataFrame

| id | text | numWords |
|----|----------------------------------|----------|
| 1 | This is Spark | 3 |
| 2 | Spark | 1 |
| 3 | Another sample sentence of words | 5 |
| 4 | Paolo Rossi | 2 |
| 5 | Giovanni | 1 |


```
1 from pyspark.sql.types import *
2 from pyspark.ml.feature import SQLTransformer
3
4 #Local Input data
5 inputList = [
6     (1,"ThisisSpark"),
7     (2,"Spark"),
8     (3,"Anothersamplesentenceofwords"),
9     (4,"PaoloRossi"),
10    (5,"Giovanni")
11 ]
12
13 # Create the initial DataFrame
14 dfInput = spark.createDataFrame(
15     inputList,
16     ["id","text"]
17 )
18
19 # Define a UDF function that counts the number of words
20 # in an input string
21 spark.udf.register(
22     "countWords",
23     lambda text: len(text.split(" ")),
24     IntegerType()
25 )
26
27 # Define an SQLTransformer to create the columns we are
28 # interested in
29 sqlTrans = SQLTransformer(
30     statement="""
31     SELECT *, countWords(text) AS numLines
32     FROM __THIS__
33     """
34 )
35
36 # Create the new DataFrame by invoking the transform method of
37 # the defined SQLTransformer
38 newDF = sqlTrans.transform(dfInput)
```