

# Preprocessing the Data in Apache Spark

# Steps in preprocessing

- Deploy data to the cluster
- Creating building RDD
- Verifying the data thru. sampling
- Cleaning data: For example
  - Converting the datatype
  - Filing missing values
- Other steps: integration, reduction, transformation, discretization

# The Data Set

A	B	C	D	E	F	G	H	I	J	K	L	
id_1	id_2	cmp_fnar	cmp_fnar	cmp_lnar	cmp_lnar	cmp_sex	cmp_bd	cmp_bm	cmp_by	cmp_plz	is_match	
37291	53113	0.833333	?		1 ?		1	1	1	1	0	TRUE
39086	47614		1 ?		1 ?		1	1	1	1	1	TRUE
70031	70237		1 ?		1 ?		1	1	1	1	1	TRUE
84795	97439		1 ?		1 ?		1	1	1	1	1	TRUE
36950	42116		1 ?		1	1	1	1	1	1	1	TRUE
42413	48491		1 ?		1 ?		1	1	1	1	1	TRUE
25965	64753		1 ?		1 ?		1	1	1	1	1	TRUE
49451	90407		1 ?		1 ?		1	1	1	1	0	TRUE
39932	40902		1 ?		1 ?		1	1	1	1	1	TRUE
46626	47940		1 ?		1 ?		1	1	1	1	1	TRUE
48948	98379		1 ?		1 ?		1	1	1	1	1	TRUE
4767	4826		1 ?		1 ?		1	1	1	1	1	TRUE
45463	69659		1 ?		1 ?		1	1	1	1	1	TRUE
11367	13169		1 ?		1 ?		1	1	1	1	1	TRUE
10782	89636		1 ?		1 ?		1	0	1	1	1	TRUE
26206	39147		1 ?		1 ?		1	1	1	1	1	TRUE
16662	27083		1	1	1 ?		1	1	1	1	1	TRUE
18823	30204		1	1	1 ?		1	1	1	1	1	TRUE
38545	85418		1 ?		1 ?		1	1	1	1	1	TRUE
28963	39172		1 ?		1 ?		1	1	1	1	1	TRUE
38918	44578		1 ?		1 ?		1	1	1	1	1	TRUE
23499	40452		1 ?		1 ?		1	1	1	1	1	TRUE
46469	94155		1 ?		1 ?		1	1	1	1	1	TRUE
25816	29809		1 ?		1 ?		1	1	1	1	1	TRUE
28777	39482		1 ?		1 ?		1	1	1	1	1	TRUE
18914	18916		1 ?		1 ?		1	1	1	1	1	TRUE
15105	16971		1 ?		1 ?		1	1	1	1	1	TRUE
75313	76006		1 ?		1 ?		1	1	1	1	1	TRUE
29624	90371		1	1	1 ?		1	1	1	1	1	TRUE
30127	30732		1 ?		1 ?		1	1	1	1	1	TRUE

# Deploy the data to the cluster

- Distributed computing requires the file distributed across the cluster
- Transfer the local data files to hdfs
  - `$ hdfs dfs -mkdir linkage`
  - `$ hdfs dfs -put block_*.csv linkage`

# Creating

- Create a RDD (Resilient Distributed Dataset) from text file
  - `val rawblocks = sc.textFile("hdfs:///user/yxie2/linkage2")`
- Create RDD from external databases
  - `val hiveContext = new org.apache.spark.sql.hive.HiveContext(sc)`
  - `val test_enc_orc = hiveContext.sql("select * from test_enc_orc")`
- Spark is a lazy execution

# Sampling data

- Sample:
  - `$val head = rawblocks.first()`
  - `$val top10 = rawblocks.take(10)`
- View data:
  - Printing on to client console
    - `$head.foreach(println)`

# The Data Set

A	B	C	D	E	F	G	H	I	J	K	L	
id_1	id_2	cmp_fnar	cmp_fnar	cmp_lnar	cmp_lnar	cmp_sex	cmp_bd	cmp_bm	cmp_by	cmp_plz	is_match	
37291	53113	0.833333	?		1 ?		1	1	1	1	0	TRUE
39086	47614		1 ?		1 ?		1	1	1	1	1	TRUE
70031	70237		1 ?		1 ?		1	1	1	1	1	TRUE
84795	97439		1 ?		1 ?		1	1	1	1	1	TRUE
36950	42116		1 ?		1	1		1	1	1	1	TRUE
42413	48491		1 ?		1 ?		1	1	1	1	1	TRUE
25965	64753		1 ?		1 ?		1	1	1	1	1	TRUE
49451	90407		1 ?		1 ?		1	1	1	1	0	TRUE
39932	40902		1 ?		1 ?		1	1	1	1	1	TRUE
46626	47940		1 ?		1 ?		1	1	1	1	1	TRUE
48948	98379		1 ?		1 ?		1	1	1	1	1	TRUE
4767	4826		1 ?		1 ?		1	1	1	1	1	TRUE
45463	69659		1 ?		1 ?		1	1	1	1	1	TRUE
11367	13169		1 ?		1 ?		1	1	1	1	1	TRUE
10782	89636		1 ?		1 ?		1	0	1	1	1	TRUE
26206	39147		1 ?		1 ?		1	1	1	1	1	TRUE
16662	27083		1	1	1 ?		1	1	1	1	1	TRUE
18823	30204		1	1	1 ?		1	1	1	1	1	TRUE
38545	85418		1 ?		1 ?		1	1	1	1	1	TRUE
28963	39172		1 ?		1 ?		1	1	1	1	1	TRUE
38918	44578		1 ?		1 ?		1	1	1	1	1	TRUE
23499	40452		1 ?		1 ?		1	1	1	1	1	TRUE
46469	94155		1 ?		1 ?		1	1	1	1	1	TRUE
25816	29809		1 ?		1 ?		1	1	1	1	1	TRUE
28777	39482		1 ?		1 ?		1	1	1	1	1	TRUE
18914	18916		1 ?		1 ?		1	1	1	1	1	TRUE
15105	16971		1 ?		1 ?		1	1	1	1	1	TRUE
75313	76006		1 ?		1 ?		1	1	1	1	1	TRUE
29624	90371		1	1	1 ?		1	1	1	1	1	TRUE
30127	30732		1 ?		1 ?		1	1	1	1	1	TRUE

- Pre-Process the data – Part II: Structuring Data with Tuples and Case Classes
  - The records in the head array are all strings of comma-separated fields
  - To make it a bit easier to analyze this data, we will need to parse these strings into a structured format that converts the different fields into the correct data type



```
def parse(line: String) ={
  val pieces = line.split(',')
  val id1 = pieces(0).toInt
  val id2 = pieces(1).toInt
  val scores = pieces.slice(2,11).map(_.toDouble)
  val matched = pieces(11).toBoolean
  (id1, id2, scores, matched)
}
val lines= sc.textFile("hdfs://...")
val flines = lines.map(line => parse(line) )
```

- Needs to redefine the toDouble method

```
def toDouble(s: String) = {  
    if ("?".equals(s)) Double.NaN else s.toDouble  
}
```

```
def parse(line: String) = {  
    val pieces = line.split(',')  
    val id1 = pieces(0).toInt  
    val id2 = pieces(1).toInt  
    val scores = pieces.slice(2,11).map(toDouble)  
    val matched = pieces(11).toBoolean  
    (id1, id2, scores, matched)  
}
```

```
val flines = lines.map(line => parse(line) )  
val filtered_data = flines.filter(tup => tup(4) > 0)
```

# Transformation

Transformation	Meaning
<code>map(func)</code>	Return a new distributed dataset formed by passing each element of the source through a function <i>func</i> .
<code>filter(func)</code>	Return a new dataset formed by selecting those elements of the source on which <i>func</i> returns true.
<code>flatMap(func)</code>	Similar to map, but each input item can be mapped to 0 or more output items (so <i>func</i> should return a Seq rather than a single item).
<code>mapPartitions(func)</code>	Similar to map, but runs separately on each partition (block) of the RDD, so <i>func</i> must be of type <code>Iterator&lt;T&gt; =&gt; Iterator&lt;U&gt;</code> when running on an RDD of type T.
<code>mapPartitionsWithIndex(func)</code>	Similar to <code>mapPartitions</code> , but also provides <i>func</i> with an integer value representing the index of the partition, so <i>func</i> must be of type <code>(Int, Iterator&lt;T&gt;) =&gt; Iterator&lt;U&gt;</code> when running on an RDD of type T.
<code>sample(withReplacement, fraction, seed)</code>	Sample a fraction <i>fraction</i> of the data, with or without replacement, using a given random number generator <i>seed</i> .
<code>union(otherDataset)</code>	Return a new dataset that contains the union of the elements in the source dataset and the argument.
<code>intersection(otherDataset)</code>	Return a new RDD that contains the intersection of elements in the source dataset and the argument.
<code>distinct([numTasks])</code>	Return a new dataset that contains the distinct elements of the source dataset.
<code>groupByKey([numTasks])</code>	When called on a dataset of (K, V) pairs, returns a dataset of (K, Iterable<V>) pairs. <b>Note:</b> If you are grouping in order to perform an aggregation (such as a sum or average) over each key, using <code>reduceByKey</code> or <code>aggregateByKey</code> will yield much better performance. <b>Note:</b> By default, the level of parallelism in the output depends on the number of partitions of the parent RDD. You can pass an optional <code>numTasks</code> argument to set a different number of tasks.
<code>reduceByKey(func, [numTasks])</code>	When called on a dataset of (K, V) pairs, returns a dataset of (K, V) pairs where the values for each key are aggregated using the given reduce function <i>func</i> , which must be of type <code>(V, V) =&gt; V</code> . Like in <code>groupByKey</code> , the number of reduce tasks is configurable through an optional second argument.
<code>aggregateByKey(zeroValue)(seqOp, combOp, [numTasks])</code>	When called on a dataset of (K, V) pairs, returns a dataset of (K, U) pairs where the values for each key are aggregated using the given combine functions and a neutral "zero" value. Allows an aggregated value type that is different than the input value type, while avoiding unnecessary allocations. Like in <code>groupByKey</code> , the number of reduce tasks is configurable through an optional second argument.
<code>sortByKey([ascending], [numTasks])</code>	When called on a dataset of (K, V) pairs where K implements Ordered, returns a dataset of (K, V) pairs sorted by keys in ascending or descending order, as specified in the boolean <i>ascending</i> argument.
<code>join(otherDataset, [numTasks])</code>	When called on datasets of type (K, V) and (K, W), returns a dataset of (K, (V, W)) pairs with all pairs of elements for each key. Outer joins are supported through <code>leftOuterJoin</code> , <code>rightOuterJoin</code> , and <code>fullOuterJoin</code> .
<code>cogroup(otherDataset, [numTasks])</code>	When called on datasets of type (K, V) and (K, W), returns a dataset of (K, (Iterable<V>, Iterable<W>)) tuples. This operation is also called <code>groupWith</code> .
<code>cartesian(otherDataset)</code>	When called on datasets of types T and U, returns a dataset of (T, U) pairs (all pairs of elements).

# Action

Action	Meaning
<b>reduce</b> ( <i>func</i> )	Aggregate the elements of the dataset using a function <i>func</i> (which takes two arguments and returns one). The function should be commutative and associative so that it can be computed correctly in parallel.
<b>collect</b> ()	Return all the elements of the dataset as an array at the driver program. This is usually useful after a filter or other operation that returns a sufficiently small subset of the data.
<b>count</b> ()	Return the number of elements in the dataset.
<b>first</b> ()	Return the first element of the dataset (similar to <code>take(1)</code> ).
<b>take</b> ( <i>n</i> )	Return an array with the first <i>n</i> elements of the dataset.
<b>takeSample</b> ( <i>withReplacement</i> , <i>num</i> , [ <i>seed</i> ])	Return an array with a random sample of <i>num</i> elements of the dataset, with or without replacement, optionally pre-specifying a random number generator seed.
<b>takeOrdered</b> ( <i>n</i> , [ <i>ordering</i> ])	Return the first <i>n</i> elements of the RDD using either their natural order or a custom comparator.
<b>saveAsTextFile</b> ( <i>path</i> )	Write the elements of the dataset as a text file (or set of text files) in a given directory in the local filesystem, HDFS or any other Hadoop-supported file system. Spark will call <code>toString</code> on each element to convert it to a line of text in the file.
<b>saveAsSequenceFile</b> ( <i>path</i> ) (Java and Scala)	Write the elements of the dataset as a Hadoop SequenceFile in a given path in the local filesystem, HDFS or any other Hadoop-supported file system. This is available on RDDs of key-value pairs that implement Hadoop's <code>Writable</code> interface. In Scala, it is also available on types that are implicitly convertible to <code>Writable</code> (Spark includes conversions for basic types like <code>Int</code> , <code>Double</code> , <code>String</code> , etc).
<b>saveAsObjectFile</b> ( <i>path</i> ) (Java and Scala)	Write the elements of the dataset in a simple format using Java serialization, which can then be loaded using <code>SparkContext.objectFile()</code> .
<b>countByKey</b> ()	Only available on RDDs of type (K, V). Returns a hashmap of (K, Int) pairs with the count of each key.
<b>foreach</b> ( <i>func</i> )	Run a function <i>func</i> on each element of the dataset. This is usually done for side effects such as updating an <a href="#">Accumulator</a> or interacting with external storage systems. <b>Note:</b> modifying variables other than Accumulators outside of the <code>foreach()</code> may result in undefined behavior. See <a href="#">Understanding closures</a> for more details.