



# FUNDAMENTALS OF BIG DATA

3 V's →   
 ↗ Volume   
 ↘ Velocity   
 ↘ Variety

Big Data processing systems

Hadoop/MapReduce { - OPEN SOURCE  
 - BATCH PROCESSING  
 Apache Spark { - OPEN SOURCE  
 - BATCH AND REAL-TIME PROCESSING

Features of Spark { - Distributed cluster computing framework  
 - Efficient in memory computations  
 - Fast data processing

Modos de deploy { Local Mode  
 Cluster Mode

Spark está escrito en scala.

SparkContext (SC) { - SparkContext es un entry point a Spark  
 - Un entry point es una manera de conectarse a un cluster de spark.

→ Version → SC-VERSION  
 → Python Version → SC-PYTHON-VERSION  
 → Master (es la URL del cluster o un string local si corre en local) → SC-MASTER  
 Para cargar datos { SC.parallelize([1,2,3,...,1])  
 SC.textFile('file.txt') } → Se le puede añadir un org. "numPartitions: 1000" → Define el número mínimo de particiones

## Spark RDD (Resilient distributed dataset)

Para ver el N° de particiones de un RDD → `rdd.getNumPartitions()`

Transformaciones (nuevos RDD)

- `map()` → crea un nuevo RDD en el muestreo  
 - `filter()` → `rdd.filter(rdd.filter(lambda x: x > 2))`  
 - `flatMap()` → retorna multiples elementos por cada elemento del RDD original  
 - `union()` → une 2 RDDs → `combinados = rdd.b.union(rdd.c)`

`Rdd = sc.parallelize([1...1])`  
`rdd.map = rdd.map(lambda x: x * 2)`

`rdd = sc.parallelize(['hola', 'que', 'tal', 'estas?'])`  
`rdd.flatMap = rdd.flatMap(lambda x: x.split(' '))` → `['hola', 'que', 'tal', 'estas?']`

Acciones

(se escriben al RDD)

- `collect()` → retorna todos los elementos del dataset como un array  
 - `take(N)` → primeros N elementos como un array.  
 - `first()` → Retorna solo el 1° elemento  
 - `count()` → Retorna el N° de elementos que tiene el RDD

`rdd.collect()`

`rdd.take(2)`

`rdd.first()`

`rdd.count()`

Pair RDD

'llave': 'valor'

`rdd.pair = sc.parallelize([(A, 1), (B, 2)])`

Se pueden crear con tuplas como dato de entrada →

si tienen misma llave, se suman el valor.

Tienen la mismas transformaciones y acciones

- `reduceByKey()` → combina valores con la misma key → `rdd.reduce = rdd.reduceByKey(lambda x, y: x + y)`  
 - `sortByKey()` → ordenan por llave  
 - `groupByKey()` → agrupa valores con la misma llave → `rdd.group = rdd.groupByKey()`  
 - `join()` une 2 RDD por el valor de la llave → `rdd.join = rdd._join(rdd2)`

# FUNDAMENTALS OF BIG DATA

## Advanced RDD actions

• `reduce(fcn)` → por ejemplo por sumar todos los valores → `rdd.reduce(lambda x,y: x+y)`

• `saveAsTextFile()` → guarda el RDD en cada partición como un archivo distinto → `rdd.saveAsTextFile('nombre')`

↳ para guardarlo en el único archivo → `rdd.coalesce(1).saveAsTextFile('nombre')`

• `countByKey()` → cuenta los elementos por cada llave solo devuelve como Dict.

Para leer RDD { `collectAsMap()` → retorna un diccionario con los valores {key:val} Deben usarse si hay pocos datos ya que todos se cargan en la memoria.

## Pyspark DataFrames I

SparkSession es el entry point para interactuar con estos DataFrames → Hace lo que SparkContext para los RDD.

Para crear DataFrames {   
 - Desde un RDD → `df = spark.createDataFrame(rdd, schema=(nombre, columnas))`   
 - Desde archivos {   
    csv → `df = spark.read.csv('archivo.csv')`   
    json → `df = " " .json('archivo.json')`   
    txt → `df = " " .txt('archivo.txt')`   
 } Reservar `headerLine` si hay una columna   
 `options`   
 `inferSchema=True` → asigna el tipo de dato a cada columna.

Operaciones {   
 - Transformaciones {   
    - `select(col)` → extrae las columnas del df.   
    - `show()` → muestra las primeras filas.   
    - `filter(condition)` → filtra los valores.   
    - `groupBy(col)` → agrupa por una columna.   
    - `count()`   
    - `orderBy(col)` → Ordena por una columna.   
    - `dropDuplicates()` → Retorna un nuevo df sin duplicados.   
    - `withColumnRenamed('Antiguo nombre', 'nuevo')`   
    - `printSchema()` → Muestra los tipos de datos que contiene.   
    - `columns` → Retorna las columnas del df.   
    - `describe()` → Retorna las estadísticas de las columnas numéricas.   
    - `withColumn('nuevo', columna con viejo)` → Crea una nueva columna a partir de otras.   
    - `drop()` → Para eliminar una columna.   
    - `distinct()` → Valores únicos.   
    - `startswith()` → Compara si empieza por un carácter.   
 }   
 }   
 - `write()` → Guarda el DataFrame en un archivo o base de datos.

Usando consultas SQL {   
 - Primero debemos crear una tabla temporal → `df.createOrReplaceTempView('tabla')`   
 - En seguida un comando SQL → `df.sql('SELECT col1, col2 FROM tabla')`   
 } El resultado siempre es un DataFrame.

# FUNDAMENTALS OF BIG DATA

## PySpark Databases II

Data visualization with df.

- Spark Databases as functions can be plotted, scatter... they 9 queries.
  - `PySpark.dfs.explore`
    - Re an histogram → `hist(df, bins=20, colors=...)`
    - Otras funciones → `distplot(), pandas.histogram()`
  - `WackySpark`
    - Convertir a Wackydf → `h.df = df.toWacky()`
    - Crear un histograma → `h.df.hist()`

Convertir a pandas → `toPandas()`

WackyPySpark.dfs.explore  
WackySpark

## PySpark MLlib

→ SOLO SOPORTA RDDS

Collaborative filtering → Produce recommendations

Clustering → agrupar los datos con características similares

Classification → identificar a que categoría pertenece.

from pyspark.mllib.recommendation import ALS

from pyspark.mllib.clustering import KMeans

from pyspark.mllib.classification import LogisticRegressionWithLBFGS

- `user-user` ⇒ Busca usuarios similares a el usuario objetivo
- `item-item` ⇒ " items " a los items del usuario obj.
- `Rating` ⇒ es una clase que guarda el RDD de rating (user, product, rating)
- `Compartes`
  - `randomSplit()` ⇒ para dividir en train y test
  - `ALS (Alternating Least Squares)`
    - `ratings = sc.parallelize([14, 12, 13])`
    - `model = ALS.train(ratings, rank=50, iterations=50)`
  - `predictAll()` ⇒ genera la salida

from pyspark.mllib.recommendation import Rating  
r: Rating(user=..., product=..., rating=...)

train, test = rdd.randomSplit([0.6, 0.4])

ratings = sc.parallelize([14, 12, 13])  
model = ALS.train(ratings, rank=50, iterations=50)

unrated = sc.parallelize([(1, 1, 1), (1, 1, 1)])  
pred = model.predictAll(unrated)

- `MLlib` contiene 2 tipos de datos específicos
- `Vector`
  - `Dense Vector` ⇒ guarda los enteros en un array de floating point numbers
  - `Sparse Vector` ⇒ solo valores ≠ 0 y sus índices
- `Classification`
  - `HashingTF` ⇒ mapea los x a valores en el vector de x
  - `LogisticRegression` ⇒
    - `data = [LabelPoint(...), ...]`
    - `rdd = sc.parallelize(data)`
    - `train = LogisticRegressionWithLBFGS.train(rdd)`
    - `train.predict(...)`

Label Point ⇒ simple de clasificar  
frase = "Hacker hack que"  
words = frase.split(" ")  
tf = HashingTF(words)  
tf.transform(words)

frase = "Hacker hack que"  
words = frase.split(" ")  
tf = HashingTF(words)  
tf.transform(words)

data = [LabelPoint(...), ...]  
rdd = sc.parallelize(data)  
train = LogisticRegressionWithLBFGS.train(rdd)  
train.predict(...)

# FUNDAMENTALS OF BIG DATA

## PySpark MLlib

