

Costa Rica Big Data School: Large Scale Data Analysis with Spark

Weijia Xu

Research Scientist, Group Manager

Data Mining & Statistics

Texas Advanced Computing Center

University of Texas at Austin

Dec. 2018

Starting Zeppelin on Wrangler

1. Connect to Wrangler

- Open terminal window or application and run following:
- **ssh USERNAME@wrangler.tacc.utexas.edu**
- If succeed, you should see

```
login1.wrangler(1)$
```

- You are now on the login node on wrangler

2. Then start Zeppelin session on compute node

- **sbatch --reservation=hadoop+TRAINING-OPEN+2552 /data/apps/zeppelin_user/job.zeppelin**

```
login1.wrangler(2)$ sbatch --reservation=hadoop+TRAINING-OPEN+2552 /data/apps/zeppelin_user/job.zeppelin
```

- If succeed, you should see

```
--> Submitting job request to current queue ...
--> Checking available allocation (TRAINING-OPEN)...OK
Submitted batch job 98613
```

Access your Zeppelin session

3. Wait a few minutes for job to run
4. then access zeppelin via web
 - cat zeppelin.out

```
[login1.wrangler(6)$ cat zeppelin.out
```

- At the end of output you may see something like

```
Your applicatin is now running!
Application UI is at http://wrangler.tacc.utexas.edu:58091
Zeppelin username and password: use your TACC credential
[login1: ~]$
```

- Open URL in web browser
- If ...
 - Nothing happened yet, just wait couple of more minutes.
 - You can login but cannot open/create any notebook, try using in **incognito (or private)** mode or a different browser if possible.

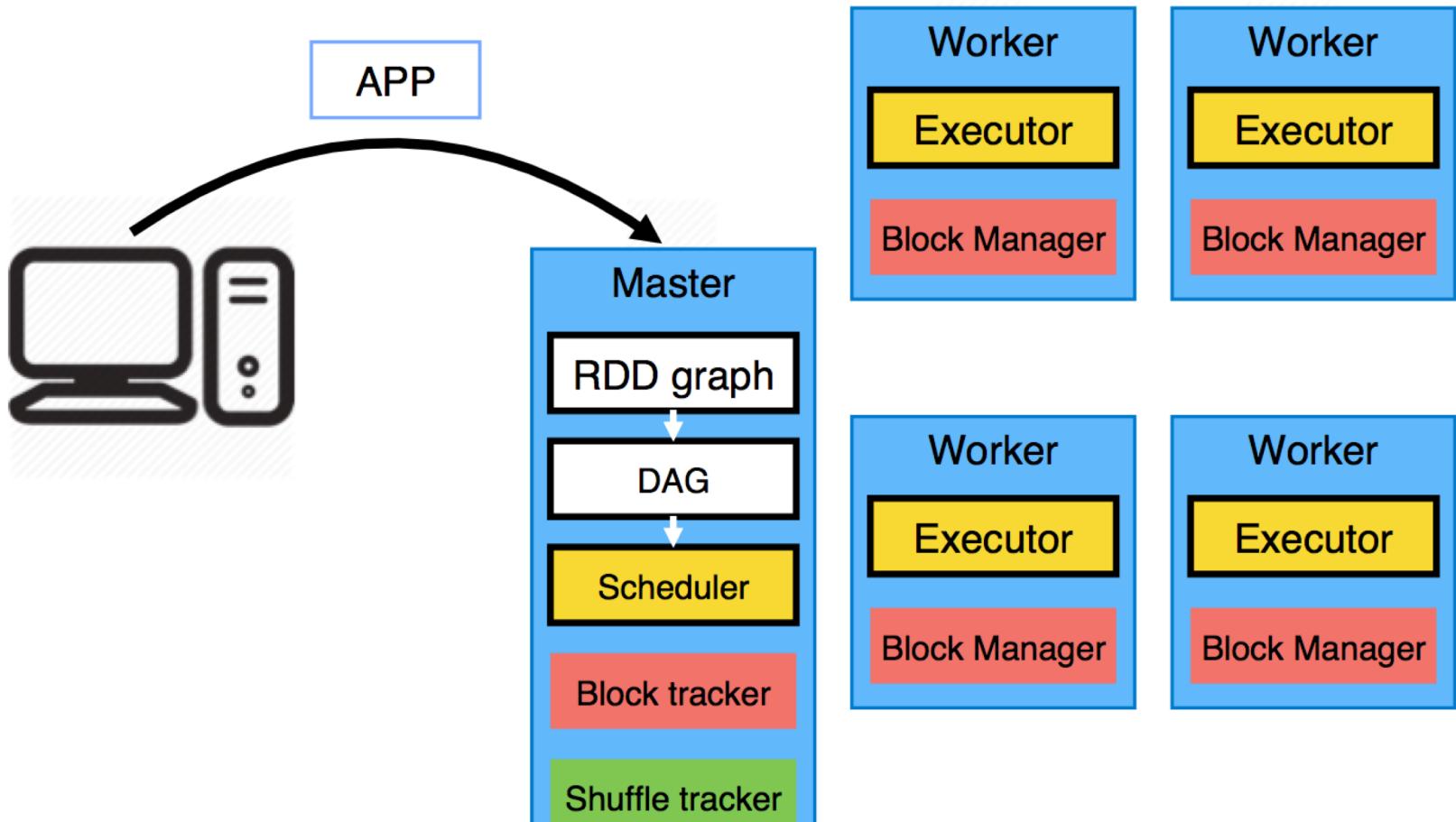
Day 2 Review

- Why Spark is better than Hadoop for analysis
 - RDD
- Scala Basics
 - Immutable variables
 - List and tuple
 - Functions
- A close look of Spark RDD
 - An “upgrade” from key/value pair representation of data
 - Immutable and can be recreated
 - Transformation vs. Action

Data Model in Spark

- RDD
 - Resilient Distributed Dataset
 - Implemented from the beginning of Spark framework
 - Support a number of transformation functions, e.g.
 - Map, `rdd.map(x => x*2)`
 - reduce, `rdd.reduce(_ + _)`
 - filter, `rdd.filter(_ % 3 == 0)`
 - ...
 - Think it as a set of objects stored in memory across nodes, e.g.
`val rdd = sc.parallelize(0 until 10000)`

How does Spark run computation job



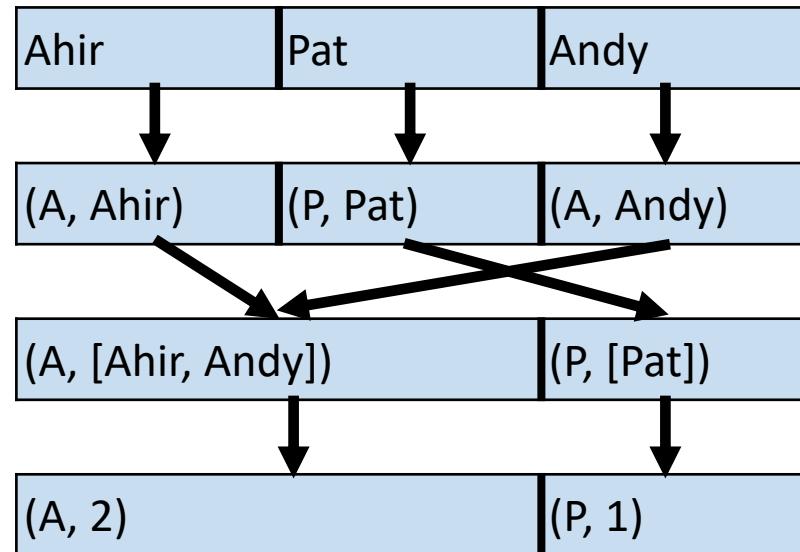
Sample Execution of a Spark Job

Read a file from HDFS and count how many words start with each letter.

1. val lines = sc.textFile("hdfs://names")
2. val kvp = lines.map(name => (name(0), name))
3. val groups = kvp.groupByKey()
4. val res = groups.mapValues(names => names.toSet.size)
5. res.collect

Sample Execution of a Spark Job

```
1. val lines = sc.textFile("hdfs://names")
2. val kvp = lines.map(name => (name(0), name))
3. val groups = kvp.groupByKey()
4. val res = groups.mapValues(names => names.toSet.size)
5. res.collect
```



Create RDDs

```
val lines = sc.textFile("hdfs://names")
```

```
val kvp = lines.map(name => (name(0), name))
```

```
val groups = kvp.groupByKey()
```

```
val res = groups.mapvalues(names => names.toSet.size)
```

```
res.collect
```

HadoopRDD

MapPartitionsRDD

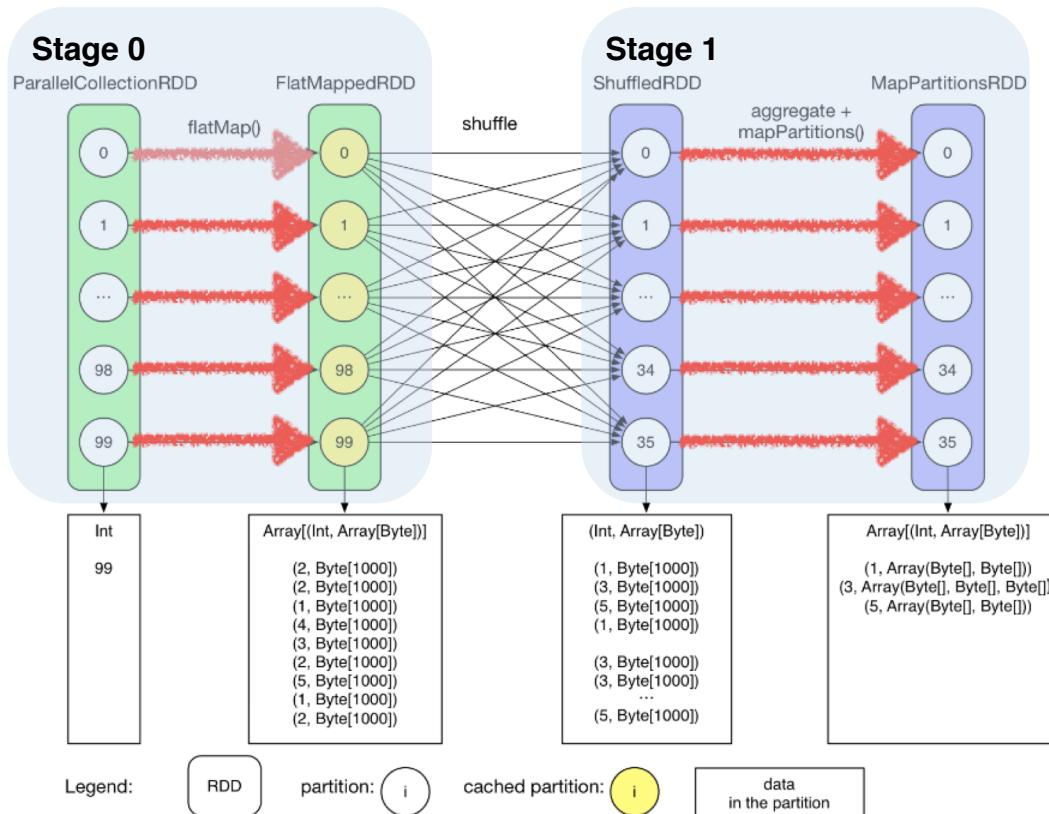
ShuffledRDD

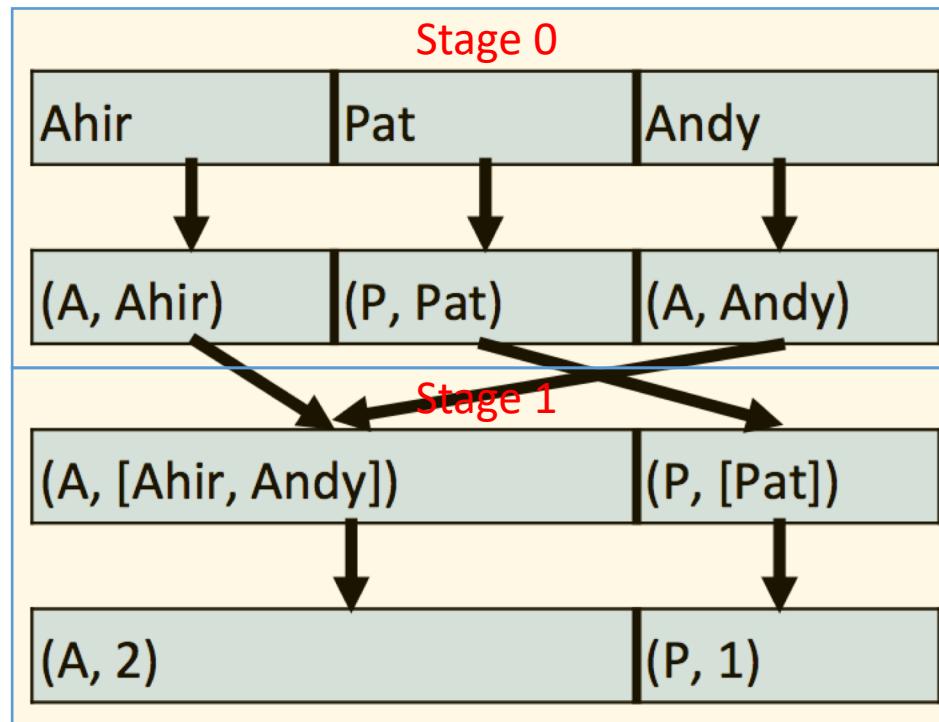
MapPartitionsRDD

collect()

DAG Generation

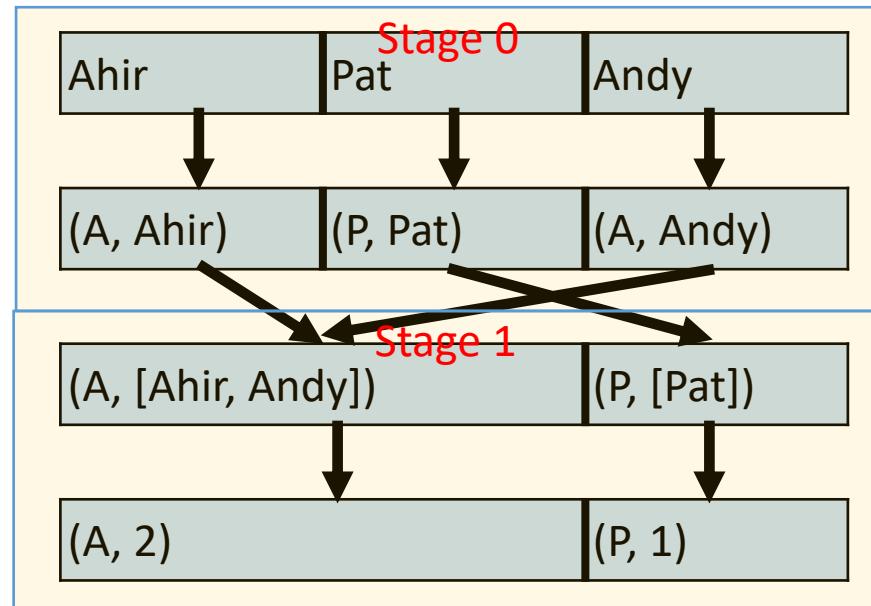
One task per final partition until a shuffle





Sample Execution of a Spark Job

```
1. val lines = sc.textFile("hdfs://names")
2. val kvp = lines.map(name => (name(0), name))
3. val groups = kvp.groupByKey()
4. val res = groups.mapValues(names => names.toSet.size)
5. res.collect
```



Data Model in Spark

- Dataframe
 - Since Spark 1.3
 - An abstract API built on top of RDD.
 - Have schema to describe the data.
 - Can use off-heap storage for large data.
 - Think it as a table stored across data nodes
- Dataset
 - Since Spark 1.6
 - An API to combine advantages of both RDD and DataFrame
- As of Spark 2.X, Dataset and DataFrame APIs are merged.
DataFrame is a Dataset[Row]

Working with File in Spark

- One of the advantage of the DataFrame is easily read structured data file. e.g. reading a csv file
 - ```
val df = spark.read.format("csv")
 .options("header", true)
 .load("/tmp/data/mtcars.csv")
```
- The results can easily be shown as a table with show

```
val df=spark.read.format("csv").option("header", true).load("/tmp/data/mtcars.csv")
df.show()
```

|   | model             | mpg  | cyl | displ | hp  | drat | wt    | qsec  | vs | am | gear | carb |
|---|-------------------|------|-----|-------|-----|------|-------|-------|----|----|------|------|
| 1 | Mazda RX4         | 21   | 6   | 160   | 110 | 3.9  | 2.62  | 16.46 | 0  | 1  | 4    | 4    |
| 1 | Mazda RX4 Wag     | 21   | 6   | 160   | 110 | 3.9  | 2.875 | 17.02 | 0  | 1  | 4    | 4    |
| 1 | Datsun 710        | 22.8 | 4   | 108   | 93  | 3.85 | 2.32  | 18.61 | 1  | 1  | 4    | 1    |
| 1 | Hornet 4 Drive    | 21.4 | 6   | 258   | 110 | 3.08 | 3.215 | 19.44 | 1  | 0  | 3    | 1    |
| 1 | Hornet Sportabout | 18.7 | 8   | 360   | 175 | 3.15 | 3.44  | 17.02 | 0  | 0  | 3    | 2    |
| 1 | Valiant           | 18.1 | 6   | 225   | 105 | 2.76 | 3.46  | 20.22 | 1  | 0  | 3    | 1    |
| 1 | Duster 360        | 14.3 | 8   | 360   | 245 | 3.21 | 3.57  | 15.84 | 0  | 0  | 3    | 4    |

# Working with File in Spark

- Several common format can be easily read and write to/from Dataframe
  - text
  - JSON
  - parquet
  - ORC
  - JDBC
  - ...
- Parquet and ORC
  - Columnar store file.
  - Data are compressed and stored as binary, could be 60~80% smaller than text file format.
  - Store the data schema as well.

# Read and Write Files with DataFrame

- Write as a JSON file

```
df.write.json("cars.json")
```

- Write as a Parquet file

```
df.write.parquet("cars.parquet")
```

- Write as a delimited file

```
df.write.option("delimiter","\t").csv("cars.tab")
```

- Load from JSON file

```
val df_json = spark.read.json("cars.json")
```

- Load from Parquet file

```
val df_parquet = spark.read.parquet("cars.parquet")
```

# Notes about I/O with Spark

- Default file system
  - When use with hadoop cluster mode, the default file system is inside hdfs. e.g.  
“/tmp/data” refers path within the hdfs
  - When use in local mode, the default is the local file system.
  - To explicitly specify file system uses prefix file:/// or hdfs:/// with the path
- e.g.
  - read from Linux file system

```
val df=spark.read.format("csv")
 .option("header", true)
 .load("file:///work/00791/xwj/DMS/R-training/RDataMining/mtcars.csv")
```

- read from hdfs file system

```
val df=spark.read.format("csv")
 .option("header", true)
 .load("/tmp/data/mtcars.csv")
```

# Notes about I/O with Spark

- Path to File and Path to Folder
  - Both file name and directory path can be used as variable for reading.
  - Will automatically read all the files within the directory when just path to directory is given
  - The output path is always treated as directory

```
| hadoop fs -ls
drwxr-xr-x - xwj hadoop 0 2017-05-03 23:46 cars.json
drwxr-xr-x - xwj hadoop 0 2017-05-03 23:46 cars.parquet
drwxr-xr-x - xwj hadoop 0 2017-05-04 00:02 cars.tab
```

- Save Modes
  - Default mode will report error when save to existing data files
  - Can change to different mode with *mode()* function with parameters  
*error, append, overwrite, ignore*

# RDD vs. DataFrame

- From DataFrame to RDD with `rdd` function, e.g.

```
val car_rdd = df.rdd
car_rdd: org.apache.spark.rdd.RDD[org.apache.spark.sql.Row] = MapPartitio
res87: Array[org.apache.spark.sql.Row] = Array([Mazda RX4,21,6,160,110,3.
258,110,3.08,3.215,19.44,1,0,3,1], [Hornet Sportabout,18.7,8,360,175,3.15
3.69,3.19,20,1,0,4,2], [Merc 230,22.8,4,140.8,95,3.92,3.15,22.9,1,0,4,2],
0,3,3], [Merc 450SL,17.3,8,275.8,180,3.07,3.73,17.6,0,0,3,3], [Merc 450SL
0,3,3])
```

- From RDD to DataFrame with `toDF`

```
val rdd = sc.parallelize(0 until 10000)
val rdd_df = rdd.toDF
rdd_df.show

rdd: org.apache.spark.rdd.RDD[Int] = ParallelCollectionRDD[142] at p
rdd_df: org.apache.spark.sql.DataFrame = [value: int]
+---+
| value|
+---+
| 0|
| 1|
| 2|
| 3|
| 4|
| 5|
| 6|
| 7|
| 8|
| 9|
| 10|
| 11|
| 12|
| 13|
| 14|
| 15|
| 16|
| 17|
| 18|
| 19|
| 20|
| 21|
| 22|
| 23|
| 24|
| 25|
| 26|
| 27|
| 28|
| 29|
| 30|
| 31|
| 32|
| 33|
| 34|
| 35|
| 36|
| 37|
| 38|
| 39|
| 40|
| 41|
| 42|
| 43|
| 44|
| 45|
| 46|
| 47|
| 48|
| 49|
| 50|
| 51|
| 52|
| 53|
| 54|
| 55|
| 56|
| 57|
| 58|
| 59|
| 60|
| 61|
| 62|
| 63|
| 64|
| 65|
| 66|
| 67|
| 68|
| 69|
| 70|
| 71|
| 72|
| 73|
| 74|
| 75|
| 76|
| 77|
| 78|
| 79|
| 80|
| 81|
| 82|
| 83|
| 84|
| 85|
| 86|
| 87|
| 88|
| 89|
| 90|
| 91|
| 92|
| 93|
| 94|
| 95|
| 96|
| 97|
| 98|
| 99|
| 100|
| 101|
| 102|
| 103|
| 104|
| 105|
| 106|
| 107|
| 108|
| 109|
| 110|
| 111|
| 112|
| 113|
| 114|
| 115|
| 116|
| 117|
| 118|
| 119|
| 120|
| 121|
| 122|
| 123|
| 124|
| 125|
| 126|
| 127|
| 128|
| 129|
| 130|
| 131|
| 132|
| 133|
| 134|
| 135|
| 136|
| 137|
| 138|
| 139|
| 140|
| 141|
| 142|
| 143|
| 144|
| 145|
| 146|
| 147|
| 148|
| 149|
| 150|
| 151|
| 152|
| 153|
| 154|
| 155|
| 156|
| 157|
| 158|
| 159|
| 160|
| 161|
| 162|
| 163|
| 164|
| 165|
| 166|
| 167|
| 168|
| 169|
| 170|
| 171|
| 172|
| 173|
| 174|
| 175|
| 176|
| 177|
| 178|
| 179|
| 180|
| 181|
| 182|
| 183|
| 184|
| 185|
| 186|
| 187|
| 188|
| 189|
| 190|
| 191|
| 192|
| 193|
| 194|
| 195|
| 196|
| 197|
| 198|
| 199|
| 200|
| 201|
| 202|
| 203|
| 204|
| 205|
| 206|
| 207|
| 208|
| 209|
| 210|
| 211|
| 212|
| 213|
| 214|
| 215|
| 216|
| 217|
| 218|
| 219|
| 220|
| 221|
| 222|
| 223|
| 224|
| 225|
| 226|
| 227|
| 228|
| 229|
| 230|
| 231|
| 232|
| 233|
| 234|
| 235|
| 236|
| 237|
| 238|
| 239|
| 240|
| 241|
| 242|
| 243|
| 244|
| 245|
| 246|
| 247|
| 248|
| 249|
| 250|
| 251|
| 252|
| 253|
| 254|
| 255|
| 256|
| 257|
| 258|
| 259|
| 260|
| 261|
| 262|
| 263|
| 264|
| 265|
| 266|
| 267|
| 268|
| 269|
| 270|
| 271|
| 272|
| 273|
| 274|
| 275|
| 276|
| 277|
| 278|
| 279|
| 280|
| 281|
| 282|
| 283|
| 284|
| 285|
| 286|
| 287|
| 288|
| 289|
| 290|
| 291|
| 292|
| 293|
| 294|
| 295|
| 296|
| 297|
| 298|
| 299|
| 300|
| 301|
| 302|
| 303|
| 304|
| 305|
| 306|
| 307|
| 308|
| 309|
| 310|
| 311|
| 312|
| 313|
| 314|
| 315|
| 316|
| 317|
| 318|
| 319|
| 320|
| 321|
| 322|
| 323|
| 324|
| 325|
| 326|
| 327|
| 328|
| 329|
| 330|
| 331|
| 332|
| 333|
| 334|
| 335|
| 336|
| 337|
| 338|
| 339|
| 340|
| 341|
| 342|
| 343|
| 344|
| 345|
| 346|
| 347|
| 348|
| 349|
| 350|
| 351|
| 352|
| 353|
| 354|
| 355|
| 356|
| 357|
| 358|
| 359|
| 360|
| 361|
| 362|
| 363|
| 364|
| 365|
| 366|
| 367|
| 368|
| 369|
| 370|
| 371|
| 372|
| 373|
| 374|
| 375|
| 376|
| 377|
| 378|
| 379|
| 380|
| 381|
| 382|
| 383|
| 384|
| 385|
| 386|
| 387|
| 388|
| 389|
| 390|
| 391|
| 392|
| 393|
| 394|
| 395|
| 396|
| 397|
| 398|
| 399|
| 400|
| 401|
| 402|
| 403|
| 404|
| 405|
| 406|
| 407|
| 408|
| 409|
| 410|
| 411|
| 412|
| 413|
| 414|
| 415|
| 416|
| 417|
| 418|
| 419|
| 420|
| 421|
| 422|
| 423|
| 424|
| 425|
| 426|
| 427|
| 428|
| 429|
| 430|
| 431|
| 432|
| 433|
| 434|
| 435|
| 436|
| 437|
| 438|
| 439|
| 440|
| 441|
| 442|
| 443|
| 444|
| 445|
| 446|
| 447|
| 448|
| 449|
| 450|
| 451|
| 452|
| 453|
| 454|
| 455|
| 456|
| 457|
| 458|
| 459|
| 460|
| 461|
| 462|
| 463|
| 464|
| 465|
| 466|
| 467|
| 468|
| 469|
| 470|
| 471|
| 472|
| 473|
| 474|
| 475|
| 476|
| 477|
| 478|
| 479|
| 480|
| 481|
| 482|
| 483|
| 484|
| 485|
| 486|
| 487|
| 488|
| 489|
| 490|
| 491|
| 492|
| 493|
| 494|
| 495|
| 496|
| 497|
| 498|
| 499|
| 500|
| 501|
| 502|
| 503|
| 504|
| 505|
| 506|
| 507|
| 508|
| 509|
| 510|
| 511|
| 512|
| 513|
| 514|
| 515|
| 516|
| 517|
| 518|
| 519|
| 520|
| 521|
| 522|
| 523|
| 524|
| 525|
| 526|
| 527|
| 528|
| 529|
| 530|
| 531|
| 532|
| 533|
| 534|
| 535|
| 536|
| 537|
| 538|
| 539|
| 540|
| 541|
| 542|
| 543|
| 544|
| 545|
| 546|
| 547|
| 548|
| 549|
| 550|
| 551|
| 552|
| 553|
| 554|
| 555|
| 556|
| 557|
| 558|
| 559|
| 560|
| 561|
| 562|
| 563|
| 564|
| 565|
| 566|
| 567|
| 568|
| 569|
| 570|
| 571|
| 572|
| 573|
| 574|
| 575|
| 576|
| 577|
| 578|
| 579|
| 580|
| 581|
| 582|
| 583|
| 584|
| 585|
| 586|
| 587|
| 588|
| 589|
| 590|
| 591|
| 592|
| 593|
| 594|
| 595|
| 596|
| 597|
| 598|
| 599|
| 600|
| 601|
| 602|
| 603|
| 604|
| 605|
| 606|
| 607|
| 608|
| 609|
| 610|
| 611|
| 612|
| 613|
| 614|
| 615|
| 616|
| 617|
| 618|
| 619|
| 620|
| 621|
| 622|
| 623|
| 624|
| 625|
| 626|
| 627|
| 628|
| 629|
| 630|
| 631|
| 632|
| 633|
| 634|
| 635|
| 636|
| 637|
| 638|
| 639|
| 640|
| 641|
| 642|
| 643|
| 644|
| 645|
| 646|
| 647|
| 648|
| 649|
| 650|
| 651|
| 652|
| 653|
| 654|
| 655|
| 656|
| 657|
| 658|
| 659|
| 660|
| 661|
| 662|
| 663|
| 664|
| 665|
| 666|
| 667|
| 668|
| 669|
| 670|
| 671|
| 672|
| 673|
| 674|
| 675|
| 676|
| 677|
| 678|
| 679|
| 680|
| 681|
| 682|
| 683|
| 684|
| 685|
| 686|
| 687|
| 688|
| 689|
| 690|
| 691|
| 692|
| 693|
| 694|
| 695|
| 696|
| 697|
| 698|
| 699|
| 700|
| 701|
| 702|
| 703|
| 704|
| 705|
| 706|
| 707|
| 708|
| 709|
| 710|
| 711|
| 712|
| 713|
| 714|
| 715|
| 716|
| 717|
| 718|
| 719|
| 720|
| 721|
| 722|
| 723|
| 724|
| 725|
| 726|
| 727|
| 728|
| 729|
| 730|
| 731|
| 732|
| 733|
| 734|
| 735|
| 736|
| 737|
| 738|
| 739|
| 740|
| 741|
| 742|
| 743|
| 744|
| 745|
| 746|
| 747|
| 748|
| 749|
| 750|
| 751|
| 752|
| 753|
| 754|
| 755|
| 756|
| 757|
| 758|
| 759|
| 760|
| 761|
| 762|
| 763|
| 764|
| 765|
| 766|
| 767|
| 768|
| 769|
| 770|
| 771|
| 772|
| 773|
| 774|
| 775|
| 776|
| 777|
| 778|
| 779|
| 780|
| 781|
| 782|
| 783|
| 784|
| 785|
| 786|
| 787|
| 788|
| 789|
| 790|
| 791|
| 792|
| 793|
| 794|
| 795|
| 796|
| 797|
| 798|
| 799|
| 800|
| 801|
| 802|
| 803|
| 804|
| 805|
| 806|
| 807|
| 808|
| 809|
| 8010|
| 8011|
| 8012|
| 8013|
| 8014|
| 8015|
| 8016|
| 8017|
| 8018|
| 8019|
| 8020|
| 8021|
| 8022|
| 8023|
| 8024|
| 8025|
| 8026|
| 8027|
| 8028|
| 8029|
| 8030|
| 8031|
| 8032|
| 8033|
| 8034|
| 8035|
| 8036|
| 8037|
| 8038|
| 8039|
| 8040|
| 8041|
| 8042|
| 8043|
| 8044|
| 8045|
| 8046|
| 8047|
| 8048|
| 8049|
| 8050|
| 8051|
| 8052|
| 8053|
| 8054|
| 8055|
| 8056|
| 8057|
| 8058|
| 8059|
| 8060|
| 8061|
| 8062|
| 8063|
| 8064|
| 8065|
| 8066|
| 8067|
| 8068|
| 8069|
| 8070|
| 8071|
| 8072|
| 8073|
| 8074|
| 8075|
| 8076|
| 8077|
| 8078|
| 8079|
| 8080|
| 8081|
| 8082|
| 8083|
| 8084|
| 8085|
| 8086|
| 8087|
| 8088|
| 8089|
| 8090|
| 8091|
| 8092|
| 8093|
| 8094|
| 8095|
| 8096|
| 8097|
| 8098|
| 8099|
| 80100|
| 80101|
| 80102|
| 80103|
| 80104|
| 80105|
| 80106|
| 80107|
| 80108|
| 80109|
| 80110|
| 80111|
| 80112|
| 80113|
| 80114|
| 80115|
| 80116|
| 80117|
| 80118|
| 80119|
| 80120|
| 80121|
| 80122|
| 80123|
| 80124|
| 80125|
| 80126|
| 80127|
| 80128|
| 80129|
| 80130|
| 80131|
| 80132|
| 80133|
| 80134|
| 80135|
| 80136|
| 80137|
| 80138|
| 80139|
| 80140|
| 80141|
| 80142|
| 80143|
| 80144|
| 80145|
| 80146|
| 80147|
| 80148|
| 80149|
| 80150|
| 80151|
| 80152|
| 80153|
| 80154|
| 80155|
| 80156|
| 80157|
| 80158|
| 80159|
| 80160|
| 80161|
| 80162|
| 80163|
| 80164|
| 80165|
| 80166|
| 80167|
| 80168|
| 80169|
| 80170|
| 80171|
| 80172|
| 80173|
| 80174|
| 80175|
| 80176|
| 80177|
| 80178|
| 80179|
| 80180|
| 80181|
| 80182|
| 80183|
| 80184|
| 80185|
| 80186|
| 80187|
| 80188|
| 80189|
| 80190|
| 80191|
| 80192|
| 80193|
| 80194|
| 80195|
| 80196|
| 80197|
| 80198|
| 80199|
| 80200|
| 80201|
| 80202|
| 80203|
| 80204|
| 80205|
| 80206|
| 80207|
| 80208|
| 80209|
| 80210|
| 80211|
| 80212|
| 80213|
| 80214|
| 80215|
| 80216|
| 80217|
| 80218|
| 80219|
| 80220|
| 80221|
| 80222|
| 80223|
| 80224|
| 80225|
| 80226|
| 80227|
| 80228|
| 80229|
| 80230|
| 80231|
| 80232|
| 80233|
| 80234|
| 80235|
| 80236|
| 80237|
| 80238|
| 80239|
| 80240|
| 80241|
| 80242|
| 80243|
| 80244|
| 80245|
| 80246|
| 80247|
| 80248|
| 80249|
| 80250|
| 80251|
| 80252|
| 80253|
| 80254|
| 80255|
| 80256|
| 80257|
| 80258|
| 80259|
| 80260|
| 80261|
| 80262|
| 80263|
| 80264|
| 80265|
| 80266|
| 80267|
| 80268|
| 80269|
| 80270|
| 80271|
| 80272|
| 80273|
| 80274|
| 80275|
| 80276|
| 80277|
| 80278|
| 80279|
| 80280|
| 80281|
| 80282|
| 80283|
| 80284|
| 80285|
| 80286|
| 80287|
| 80288|
| 80289|
| 80290|
| 80291|
| 80292|
| 80293|
| 80294|
| 80295|
| 80296|
| 80297|
| 80298|
| 80299|
| 80300|
| 80301|
| 80302|
| 80303|
| 80304|
| 80305|
| 80306|
| 80307|
| 80308|
| 80309|
| 80310|
| 80311|
| 80312|
| 80313|
| 80314|
| 80315|
| 80316|
| 80317|
| 80318|
| 80319|
| 80320|
| 80321|
| 80322|
| 80323|
| 80324|
| 80325|
| 80326|
| 80327|
| 80328|
| 80329|
| 80330|
| 80331|
| 80332|
| 80333|
| 80334|
| 80335|
| 80336|
| 80337|
| 80338|
| 80339|
| 80340|
| 80341|
| 80342|
| 80343|
| 80344|
| 80345|
| 80346|
| 80347|
| 80348|
| 80349|
| 80350|
| 80351|
| 80352|
| 80353|
| 80354|
| 80355|
| 80356|
| 80357|
| 80358|
| 80359|
| 80360|
| 80361|
| 80362|
| 80363|
| 80364|
| 80365|
| 80366|
| 80367|
| 80368|
| 80369|
| 80370|
| 80371|
| 80372|
| 80373|
| 80374|
| 80375|
| 80376|
| 80377|
| 80378|
| 80379|
| 80380|
| 80381|
| 80382|
| 80383|
| 80384|
| 80385|
| 80386|
| 80387|
| 80388|
| 80389|
| 80390|
| 80391|
| 80392|
| 80393|
| 80394|
| 80395|
| 80396|
| 80397|
| 80398|
| 80399|
| 80400|
| 80401|
| 80402|
| 80403|
| 80404|
| 80405|
| 80406|
| 80407|
| 80408|
| 80409|
| 80410|
| 80411|
| 80412|
| 80413|
| 80414|
| 80415|
| 80416|
| 80417|
| 80418|
| 80419|
| 80420|
| 80421|
| 80422|
| 80423|
| 80424|
| 80425|
| 80426|
| 80427|
| 80428|
| 80429|
| 80430|
| 80431|
| 80432|
| 80433|
| 80434|
| 80435|
| 80436|
| 80437|
| 80438|
| 80439|
| 80440|
| 80441|
| 80442|
| 80443|
| 80444|
| 80445|
| 80446|
| 80447|
| 80448|
| 80449|
| 80450|
| 80451|
| 80452|
| 80453|
| 80454|
| 80455|
| 80456|
| 80457|
| 80458|
| 80459|
| 80460|
| 80461|
| 80462|
| 80463|
| 80464|
| 80465|
| 80466|
| 80467|
| 80468|
| 80469|
| 80470|
| 80471|
| 80472|
| 80473|
| 80474|
| 80475|
| 80476|
| 80477|
| 80478|
| 80479|
| 80480|
| 80481|
| 80482|
| 80483|
| 80484|
| 80485|
| 80486|
| 80487|
| 80488|
| 80489|
| 80490|
| 80491|
| 80492|
| 80493|
| 80494|
| 80495|
| 80496|
| 80497|
| 80498|
| 80499|
| 80500|
| 80501|
| 80502|
| 80503|
| 80504|
| 80505|
| 80506|
| 80507|
| 80508|
| 80509|
| 80510|
| 80511|
| 80512|
| 80513|
| 80514|
| 80515|
| 80516|
| 80517|
| 80518|
| 80519|
| 80520|
| 80521|
| 80522|
| 80523|
| 80524|
| 80525|
| 80526|
| 80527|
| 80528|
| 80529|
| 80530|
| 80531|
| 80532|
| 80533|
| 80534|
| 80535|
| 80536|
| 80537|
| 80538|
| 80539|
| 80540|
| 80541|
| 80542|
| 80543|
| 80544|
| 80545|
| 80546|
| 80547|
| 80548|
| 80549|
| 80550|
| 80551|
| 80552|
| 80553|
| 80554|
| 80555|
| 80556|
| 80557|
| 80558|
| 80559|
| 80560|
| 80561|
| 80562|
| 80563|
| 80564|
| 80565|
| 80566|
| 80567|
| 80568|
| 80569|
| 80570|
| 80571|
| 80572|
| 80573|
| 80574|
| 80575|
| 80576|
| 80577|
| 80578|
| 80579|
| 80580|
| 80581|
| 80582|
| 80583|
| 80584|
| 80585|
| 80586|
| 80587|
| 80588|
| 80589|
| 80590|
| 80591|
| 80592|
| 80593|
| 80594|
| 80595|
| 80596|
| 80597|
| 80598|
| 80599|
| 80600|
| 80601|
| 80602|
| 80603|
| 80604|
| 80605|
| 80606|
| 80607|
| 80608|
| 80609|
| 80610|
| 80611|
| 80612|
| 80613|
| 80614|
| 80615|
| 80616|
| 80617|
| 80618|
| 80619|
| 80620|
| 80621|
| 80622|
| 80623|
| 80624|
| 80625|
| 80626|
| 80627|
| 80628|
| 80629|
| 80630|
| 80631|
| 80632|
| 80633|
| 80634|
| 80635|
| 80636|
| 80637|
| 80638|
| 80639|
| 80640|
| 80641|
| 80642|
| 80643|
| 80644|
| 80645|
| 80646|
| 80647|
| 80648|
| 80649|
| 80650|
| 80651|
| 80652|
| 80653|
| 80654|
| 80655|
| 80656|
| 80657|
| 80658|
| 80659|
| 80660|
| 80661|
| 80662|
| 80663|
| 80664|
| 80665|
| 80666|
| 80667|
| 80668|
| 80669|
| 80670|
| 80671|
| 80672|
| 80673|
| 80674|
| 80675|
| 80676|
| 80677|
| 80678|
| 80679|
| 80680|
| 80681|
| 80682|
| 80683|
| 80684|
| 80685|
| 80686|
| 80687|
| 80688|
| 80689|
| 80690|
| 80691|
| 80692|
| 80693|
| 80694|
| 80695|
| 80696|
| 80697|
| 80698|
| 8
```

# Convert RDD[Object] to DataFrame

```
case class Person(id: Int, name: String)
val person = sc.parallelize(Seq(Person(1, "Mike"),
 Person(2, "Smith"),
 Person(3, "Brooke")))
val person_df = person.toDF
person_df.show

defined class Person
person: org.apache.spark.rdd.RDD[Person] = ParallelCollectionRDD[160] c
person_df: org.apache.spark.sql.DataFrame = [id: int, name: string]
+---+----+
| id| name|
+---+----+
1	Mikel
2	Smith
3	Brookel
+---+----+
```

# Common Functions with DataFrame

- `df.show(int n)`
  - Display a number of rows from the DataFrame
  - The default value will only show 20 rows.
- `df.printSchema`
  - Display schema of the DataFrame in the stdout
- `df.describe(cols)`
  - Return a DataFrame with summary statistics of the selected columns.
  - Default will run on all columns.

```
| df.describe("mpg").show
```

```
+-----+-----+
|summary| mpg|
+-----+-----+
| count| 321
| mean | 20.090624999999996
| stddev| 6.026948052089103
| min | 10.4
| max | 33.9
+-----+-----+
```

```
| df.printSchema
```

```
root
|-- model: string (nullable = true)
|-- mpg: string (nullable = true)
|-- cyl: string (nullable = true)
|-- disp: string (nullable = true)
|-- hp: string (nullable = true)
|-- drat: string (nullable = true)
|-- wt: string (nullable = true)
|-- qsec: string (nullable = true)
|-- vs: string (nullable = true)
|-- am: string (nullable = true)
|-- gear: string (nullable = true)
|-- carb: string (nullable = true)
```

# Dataset and DataFrame

- DataFrame and Dataset shares the same API
- DataFrame is a special types of Dataset[row]
- Most functions working with RDD will also work with Dataset, some may not work directly with DataFrame due to type conversion

```
val rdd = sc.parallelize(0 until 100)
val rdd_df = rdd.toDF
val rdd_ds = rdd.toDS

rdd.filter(_ < 10).collect
rdd_df.filter("value < 10").show
rdd_ds.filter("value < 10").show
rdd_ds.filter(_ < 10).show

rdd.map(_ * 2).collect
rdd_ds.map(_ * 2).show
rdd_df.map(_ * 2).show //this line will fail
rdd_df.select('value * 2).show
```

# RDD vs Dataframe

Data Frame supports use of “column name”.  
filter example

```
val car_rdd = df.rdd
car_rdd.filter(_.(1).asInstanceOf[String].toDouble > 20).collect
```

```
car_rdd: org.apache.spark.rdd.RDD[org.apache.spark.sql.Row] = MapPartitionsRDD[58] at rdd at <console>:27
res100: Array[org.apache.spark.sql.Row] = Array([Mazda RX4,21,6,160,110
,3.9,2.62,16.46,0,1,4,4], [Mazda RX4 Wag,21,6,160,110,3.9,2.875,17.02,0
,1,4,4], [Datsun 710,22.8,4,108,93,3.85,2.32,18.61,1,1,4,1], [Hornet 4
Drive,21.4,6,258,110,3.08,3.215,19.44,1,0,3,1], [Merc 240D,24.4,4,146.7
,62,3.69,3.19,20,1,0,4,2], [Merc 230,22.8,4,140.8,95,3.92,3.15,22.9,1,0
,4,2], [Fiat 128,32.4,4,78.7,66,4.08,2.2,19.47,1,1,4,1], [Honda Civic,3
0.4,4,75.7,52,4.93,1.615,18.52,1,1,4,2], [Toyota Corolla,33.9,4,71.1,65
,4.22,1.835,19.9,1,1,4,1], [Toyota Corona,21.5,4,120.1,97,3.7,2.465,20.
01,1,0,3,1], [Fiat X1-9,27.3,4,79,66,4.08,1.935,18.9,1,1,4,1], [Porsche
914-2,26,4,120.3,91,4.43,2.14,16.7,0,1,5,2], [Lotus Europa,30.4,4,95.1,
113,3.77,1.513,16.9,1,1,5,2], [Volvo 142E,21.4,4,121,109,4.11,2.78,18.6
,1,1,4,2])
```

```
df.filter("mpg>20").show
```

| model          | mpg  | cyl | displ | hp  | drat |
|----------------|------|-----|-------|-----|------|
| Mazda RX4      | 21   | 6   | 160   | 110 | 3.9  |
| Mazda RX4 Wag  | 21   | 6   | 160   | 110 | 3.9  |
| Datsun 710     | 22.8 | 4   | 108   | 93  | 3.85 |
| Hornet 4 Drive | 21.4 | 6   | 258   | 110 | 3.08 |
| Merc 240D      | 24.4 | 4   | 146.7 | 62  | 3.69 |
| Merc 230       | 22.8 | 4   | 140.8 | 95  | 3.92 |
| Fiat 128       | 32.4 | 4   | 78.7  | 66  | 4.08 |
| Honda Civic    | 30.4 | 4   | 71.1  | 65  | 3.15 |
| Toyota Corolla | 33.9 | 4   | 71.1  | 65  | 3.19 |
| Toyota Corona  | 21.5 | 4   | 120.1 | 97  | 3.7  |
| Fiat X1-9      | 27.3 | 4   | 79    | 66  | 4.08 |
| Porsche 914-2  | 26   | 4   | 120.3 | 91  | 4.43 |
| Lotus Europa   | 30.4 | 4   | 95.1  | 113 | 3.77 |
| Volvo 142E     | 21.4 | 4   | 121   | 109 | 4.11 |

# Spark SQL

- A Structured Query Language Syntax.
- It follows HiveQL syntax. Similar to standard database SQL but not identical.
- Can be used from command line, spark-shell, pyspark, sparkR shell or over JDBC drive
- Can be used in the form of function call or as SQL statement

# Basic Spark SQL Support Functions

## Select

```
df.select("model", "mpg", "wt").show
```

```
+-----+-----+-----+
| model | mpg | wt |
+-----+-----+-----+
Mazda RX4	21	2.62
Mazda RX4 Wag	21	2.875
Datsun 710	22.8	2.32
Hornet 4 Drive	21.4	3.215
Hornet Sportabout	18.7	3.44
Valiant	18.1	3.46
Duster 360	14.3	3.57
Merc 240D	24.4	3.19
Merc 230	22.8	3.15
Merc 280	19.2	3.44
Merc 280C	17.8	3.44
Merc 450SE	16.4	4.07
Merc 450SL	17.3	3.73
Merc 450SLC	15.2	3.78
Cadillac Fleetwood	10.4	5.25
```

```
df.select($"model", $"mpg" * 1.6).show
```

```
+-----+-----+
| model | (mpg * 1.6) |
+-----+-----+
Mazda RX4	33.61
Mazda RX4 Wag	33.61
Datsun 710	36.48000000000004
Hornet 4 Drive	34.24
Hornet Sportabout	29.92
Valiant	28.96000000000004
Duster 360	22.88000000000003
Merc 240D	39.04
Merc 230	36.48000000000004
Merc 280	30.72
Merc 280C	28.48000000000004
Merc 450SE	26.24
Merc 450SL	27.68000000000003
Merc 450SLC	24.32
```

# Basic Spark SQL Support Function

```
df.select("model", "mpg", "wt", "cyl")
 .filter("cyl > 4").show
```

```
+-----+-----+-----+
| model| mpg | wt | cyl |
+-----+-----+-----+
Mazda RX4	21	2.62	6
Mazda RX4 Wag	21	2.875	6
Hornet 4 Drive	21.4	3.215	6
Hornet Sportabout	18.7	3.44	8
Valiant	18.1	3.46	6
Duster 360	14.3	3.57	8
Merc 280	19.2	3.44	6
Merc 280C	17.8	3.44	6
Merc 450SE	16.4	4.07	8
Merc 450SL	17.3	3.73	8
Merc 450SLC	15.2	3.78	8
Cadillac Fleetwood	10.4	5.25	8
Lincoln Continental	10.4	5.42	8
Chrysler Imperial	14.7	5.345	8
```

```
df.groupBy("cyl").count().show
```

```
+---+---+
| cyl | count |
+---+---+
8	14
6	7
4	11
+---+---+
```

# Basic SQL Syntax

|                     |           |
|---------------------|-----------|
| SELECT              | Select()  |
| col1, col2, ....    |           |
| FROM                | ?         |
| table1, table2, ... |           |
| [WHERE]             | filter()  |
| condition1 AND OR   |           |
| condtion2 ....      |           |
| [GROUP BY]          | groupBy() |
| col1,...            |           |
| [ORDER BY]          | sort()    |
| col1,...            |           |

# Running SQL Statement

- Register the DataFrame as a table
- Using `spark.sql(SQL_STATEMENT)`

```
df.createOrReplaceTempView("cars")
spark.sql("SELECT * FROM cars WHERE cyl = 6")
 .show

+-----+-----+-----+-----+-----+-----+-----+-----+
| model | mpg | cyl | displ | hp | drat | wt | qsec | vs | am | gear | carb |
+-----+-----+-----+-----+-----+-----+-----+-----+
Mazda RX4	21	6	160	110	3.9	2.62	16.46	0	1	4	4
Mazda RX4 Wag	21	6	160	110	3.9	2.875	17.02	0	1	4	4
Hornet 4 Drive	21.4	6	258	110	13.08	3.215	19.44	1	0	3	1
Valiant	18.1	6	225	105	2.76	3.46	20.22	1	0	3	1
Merc 280	19.2	6	167.6	123	13.92	3.44	18.3	1	0	4	4
Merc 280C	17.8	6	167.6	123	13.92	3.44	18.9	1	0	4	4
Ferrari Dino	19.7	6	145	175	13.62	2.77	15.5	0	1	5	6
+-----+-----+-----+-----+-----+-----+-----+-----+
```

# Running SQL Statement

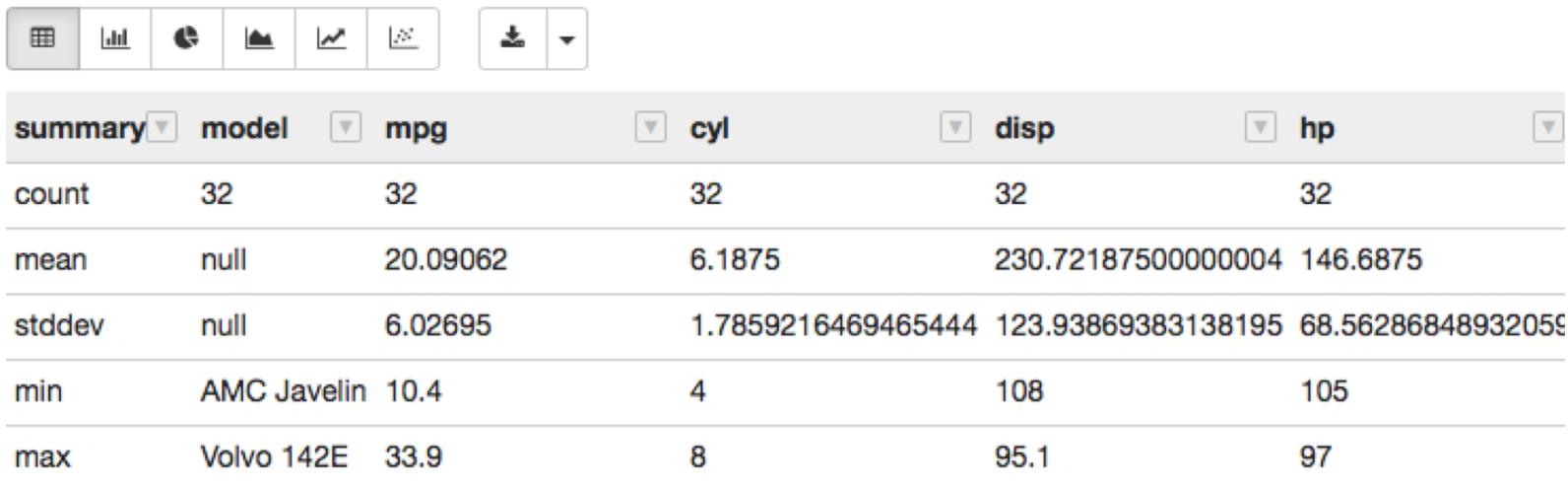
```
spark.sql("SELECT cyl, avg(mpg) as avg_mpg, count(1) as count "+
 "from cars "+
 "group by cyl")
.show
```

```
+---+-----+---+
|cyl| avg_mpg|count|
+---+-----+---+
8	15.10000000000003	14
6	19.74285714285714	7
4	26.663636363636364	11
+---+-----+---+
```

# Using Zeppelin

Zeppelin has some nice feature built-in to work with data frame. Such as display a DataFrame better with [z.show](#)

```
z.show(df.describe())
```

A screenshot of the Zeppelin interface. At the top, there is a toolbar with various icons for data manipulation and visualization. Below the toolbar is a table representing the output of the `df.describe()` command. The table has a header row with columns for `summary`, `model`, `mpg`, `cyl`, `disp`, and `hp`. The data rows show statistical summary values for each column: count (32), mean (20.09062), stddev (6.02695), min (AMC Javelin 10.4), and max (Volvo 142E 33.9).

| summary | model       | mpg      | cyl                | disp               | hp                |
|---------|-------------|----------|--------------------|--------------------|-------------------|
| count   | 32          | 32       | 32                 | 32                 | 32                |
| mean    | null        | 20.09062 | 6.1875             | 230.72187500000004 | 146.6875          |
| stddev  | null        | 6.02695  | 1.7859216469465444 | 123.93869383138195 | 68.56286848932059 |
| min     | AMC Javelin | 10.4     | 4                  | 108                | 105               |
| max     | Volvo 142E  | 33.9     | 8                  | 95.1               | 97                |

# Using Zeppelin

In addition to show a DataFrame in the table format, Zeppelin also provide a few charting options.



# Spark's Machine Learning Library (MLlib)

- A library implementing a common set of machine learning algorithms
  - Scalable
  - Efficient
  - Easy to use API
  - Not including all machine learning algorithm
    - and probably, never will.

# MLlib API

- **RDD-based API**
  - spark.mllib
  - Original version from Spark 1.X
  - Based on operations on RDD, so is further development
  - In “maintenance mode” since Spark 2.0, and still usable
  - And probably will go away in Spark 3.X ?
- **DataFrame-based API**
  - spark.ml
  - Now, the primary API for Mllib
  - More user friendly, based on operations on Dataframe
  - More uniform across different methods,
  - ML pipeline

# ML Pipeline

- Probably steal from scikit-learn ☺
- A representation of entire workflow in practical machine learning and analysis.
- Combine multiple processing steps together.
  - Data processing,
  - Algorithms
  - Evaluations
- An unified interface
  - Transformer, Estimator, Parameter

# Transformer

- A Transformer has a function called
  - `transform()`
- Convert an input dataframe to a new dataframe
  - Feature selection
    - map values in a column to a new set of values
    - Appending new column
    - Assemble multiple columns into a feature vector
  - Inference
    - Learning model take an input dataframe to make prediction

# Estimator

- Abstract concept of an algorithm
  - e.g. a learning algorithm
  - But not limited to learning algorithm
  - May also be used for feature extraction e.g. word2vec
- Implements function called
  - `fit()`
- Take a Dataframe as input and output a *model*
  - model is a \_\_\_\_\_ ?

# ML Pipeline Example

```
val cars = spark.read.format("csv").option("header", true).load("/tmp/data/mtcars.csv")
 .selectExpr("model", "mpg + 0.0 as mpg", "disp + 0.0 as disp",
 "hp + 0.0 as hp", "drat + 0.0 as drat", "wt + 0.0 as wt",
 "cyl + 0.0 as label")

val training= cars.sample(false, 0.8)
val test= cars.except(training)

val assembler = new VectorAssembler()
 .setInputCols(Array("mpg", "disp", "hp", "drat", "wt"))
 .setOutputCol("features")

val lr = new LogisticRegression()
 .setMaxIter(10)
 .setRegParam(0.2)
 .setElasticNetParam(0.0)

val pipeline = new Pipeline().setStages(Array(assembler, lr))
val lrModel = pipeline.fit(training)

val result = lrModel.transform(test).select('model, 'label, 'prediction)
result.show
```

Prepare datasets

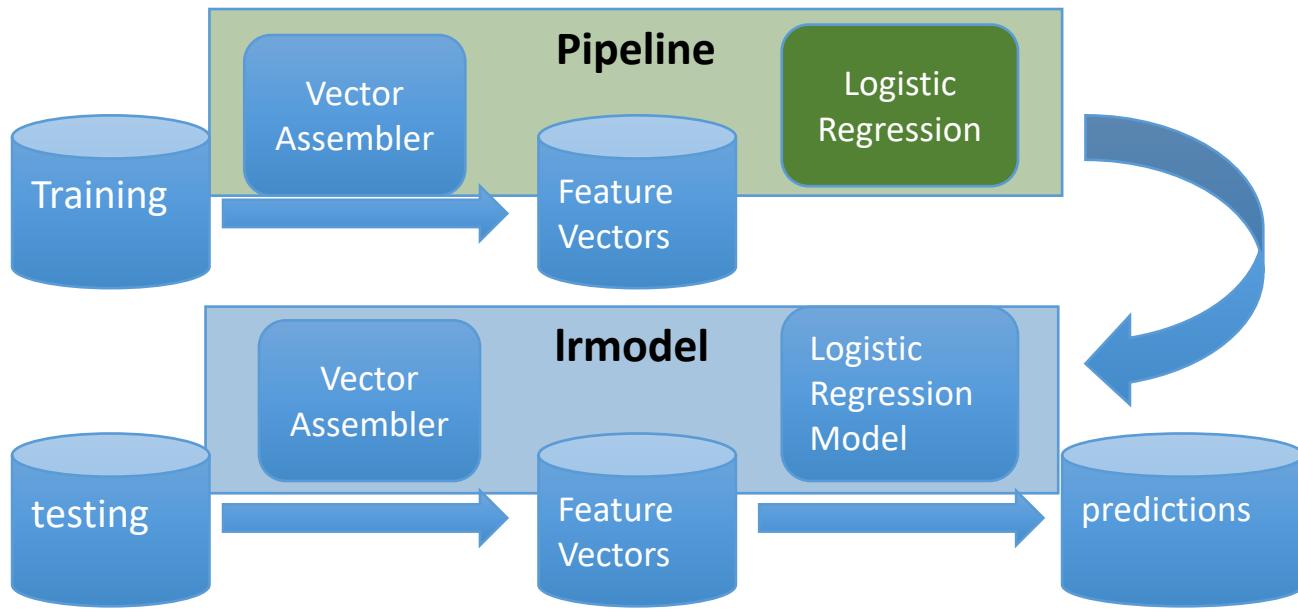
Assemble feature vectors

Estimator

Define analysis

Run analysis

See results



# Features of MLlib

- Data transformation, feature extraction
  - > 30 methods implemented
- Classification and Regression methods
  - ~ 20 methods implemented
- Clustering
  - 6 methods
- Other methods
  - Collaborative filtering
  - Frequent Pattern Mining
  - Dimension reduction
  - ....