

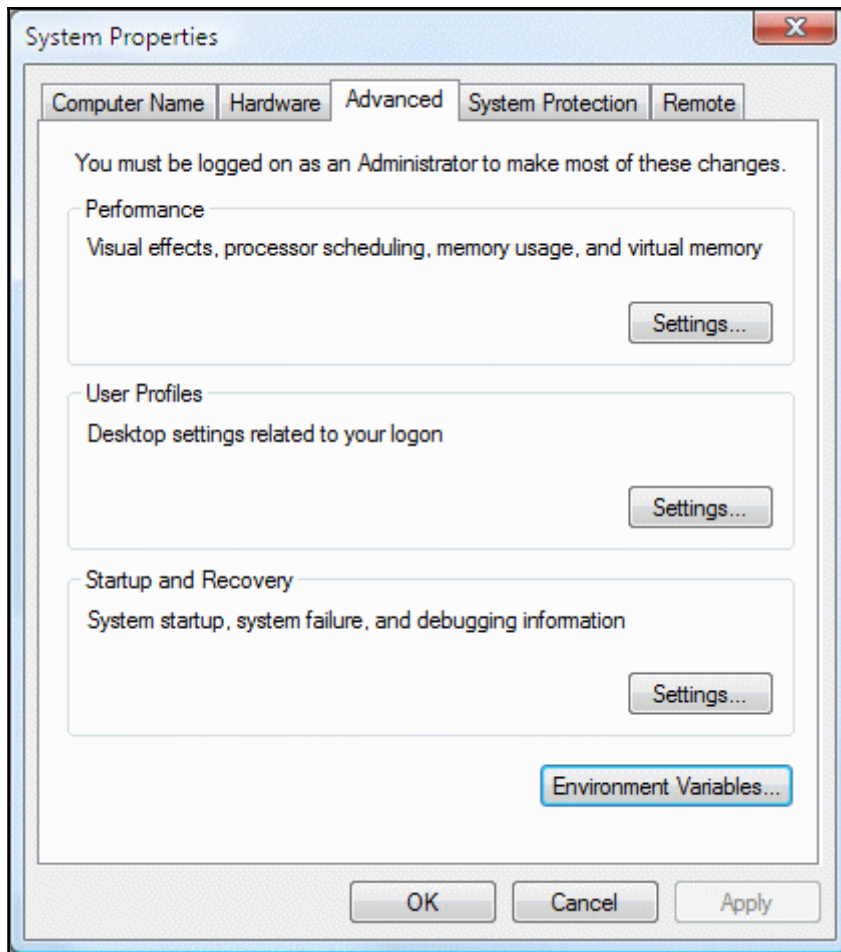
# 1

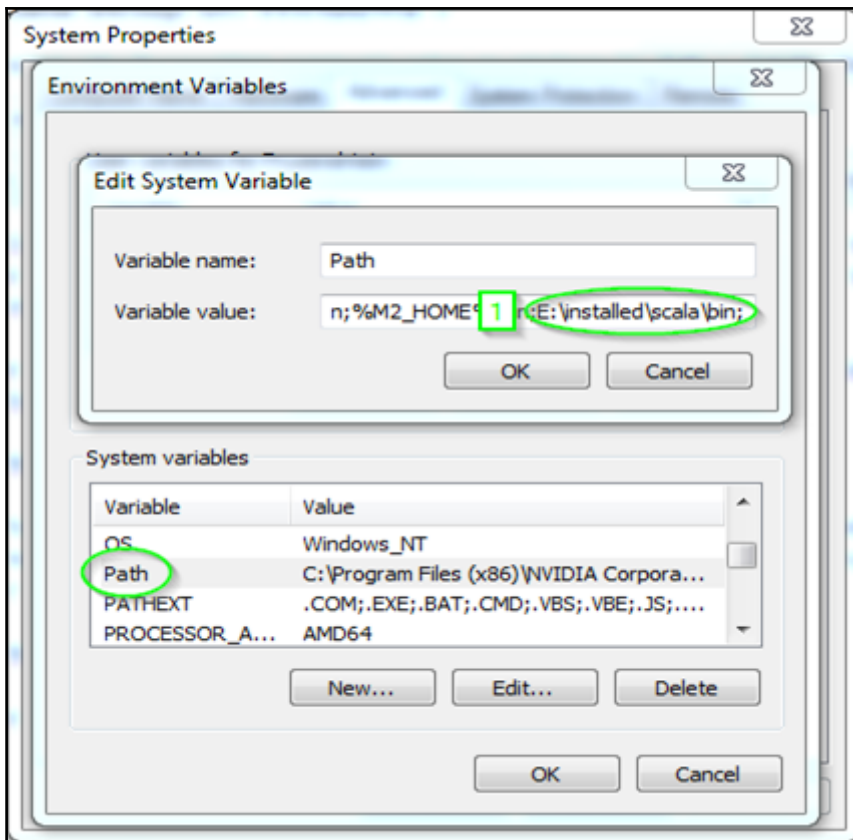
## Graphic Bundle

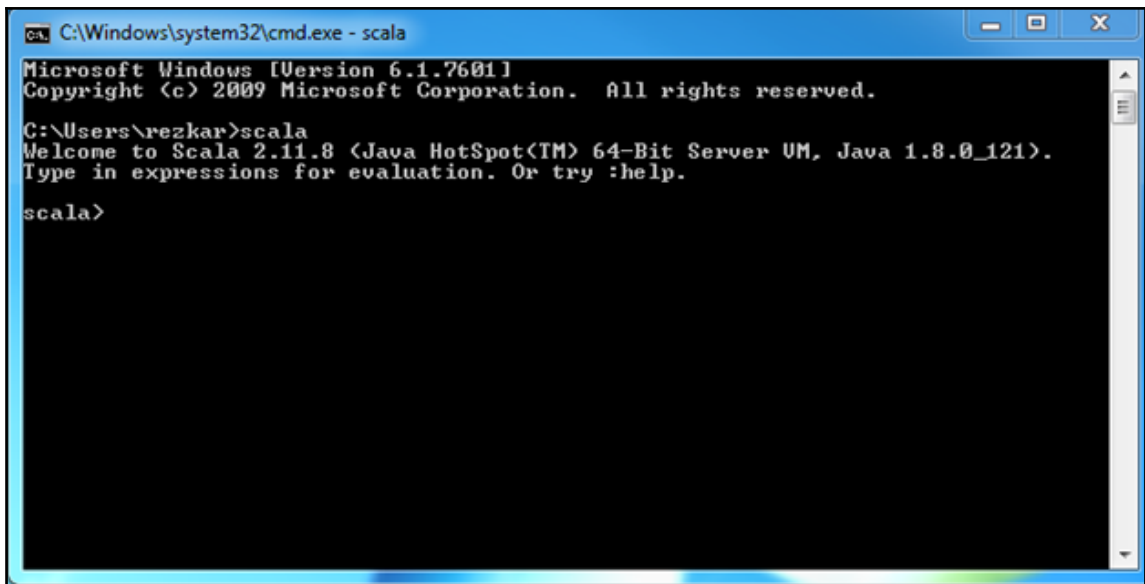
### Chapter 1: Introduction to Scala

InfoWorld Scorecard	Documentation and help system (15.0%)	Ease of use (30.0%)	Plug-in ecosystem (25.0%)	Java features (30.0%)	Overall Score (100%)
Eclipse 3.6	8.0	6.0	10.0	8.0	7.9 ★★★★☆
JetBrains IntelliJ IDEA 9.0.3	7.0	9.0	8.0	9.0	8.5 ★★★★☆
NetBeans 6.9	8.0	8.0	8.0	8.0	8.0 ★★★★☆
Oracle JDeveloper Studio 11g (11.1.1.3.0)	9.0	8.0	5.0	8.0	7.4 ★★★★☆

<b>Archive</b>	<b>System</b>	<b>Size</b>
<a href="#">scala-2.11.8.tgz</a>	Mac OS X, Unix, Cygwin	27.35M
<a href="#">scala-2.11.8.msi</a>	Windows (msi installer)	109.35M
<a href="#">scala-2.11.8.zip</a>	Windows	27.40M
<a href="#">scala-2.11.8.deb</a>	Debian	76.02M
<a href="#">scala-2.11.8.rpm</a>	RPM package	108.16M
<a href="#">scala-docs-2.11.8.txz</a>	API docs	46.00M
<a href="#">scala-docs-2.11.8.zip</a>	API docs	84.21M
<a href="#">scala-sources-2.11.8.tar.gz</a>	Sources	



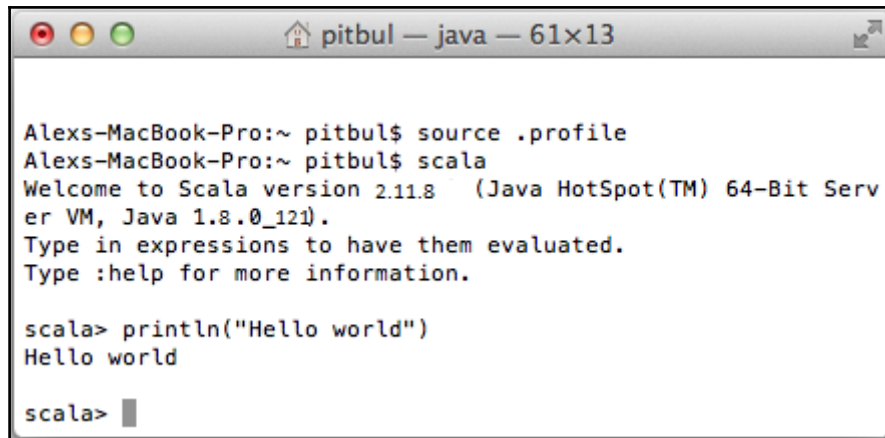




```
C:\Windows\system32\cmd.exe - scala
Microsoft Windows [Version 6.1.7601]
Copyright (c) 2009 Microsoft Corporation. All rights reserved.

C:\Users\rezkar>scala
Welcome to Scala 2.11.8 (Java HotSpot(TM) 64-Bit Server VM, Java 1.8.0_121).
Type in expressions for evaluation. Or try :help.

scala>
```



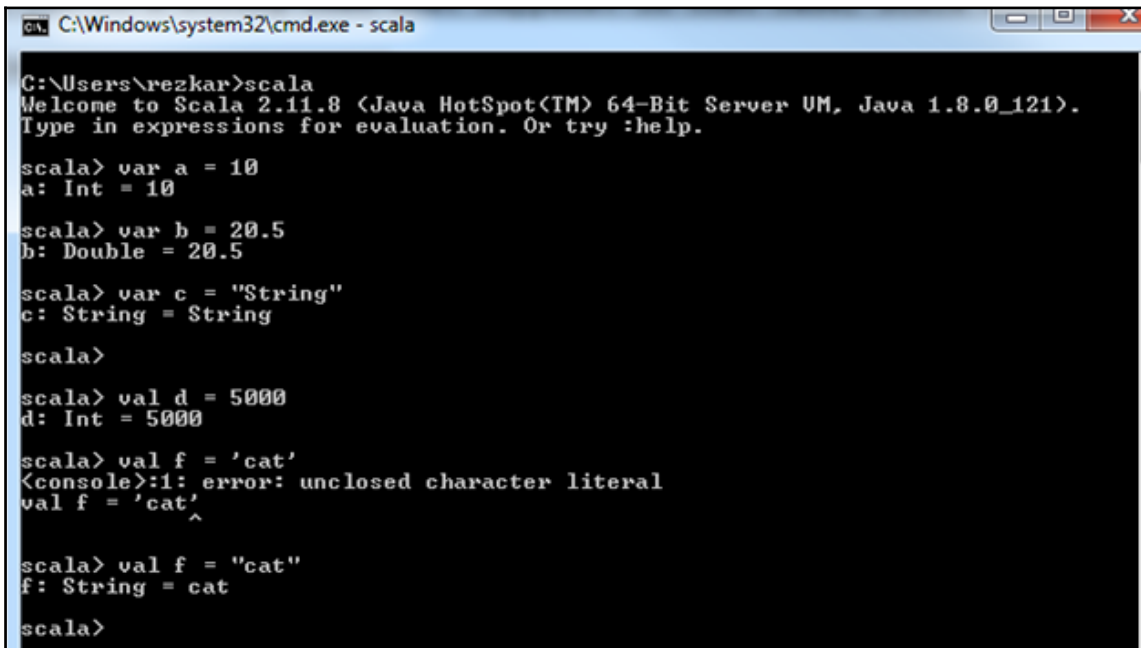
```
pitbul — java — 61x13

Alexs-MacBook-Pro:~ pitbul$ source .profile
Alexs-MacBook-Pro:~ pitbul$ scala
Welcome to Scala version 2.11.8 (Java HotSpot(TM) 64-Bit Server VM, Java 1.8.0_121).
Type in expressions to have them evaluated.
Type :help for more information.

scala> println("Hello world")
Hello world

scala> █
```

```
asif@ubuntu:~$ cd /usr/local/share/
asif@ubuntu:/usr/local/share$ ls
ca-certificates  fonts  man  scala-2.11.8  sgml  texmf  xml
asif@ubuntu:/usr/local/share$ cd ~
asif@ubuntu:~$ echo "export SCALA_HOME=/usr/local/share/scala-2.11.8" >> ~/.bashrc
asif@ubuntu:~$ echo "export PATH=$PATH:$SCALA_HOME/bin" >> ~/.bashrc
asif@ubuntu:~$ source ~/.bashrc
asif@ubuntu:~$ scala
Welcome to Scala version 2.9.2 (Java HotSpot(TM) 64-Bit Server VM, Java 1.8.0_121).
Type in expressions to have them evaluated.
Type :help for more information.
```



```
C:\Windows\system32\cmd.exe - scala
C:\Users\rezkar>scala
Welcome to Scala 2.11.8 (Java HotSpot(TM) 64-Bit Server VM, Java 1.8.0_121).
Type in expressions for evaluation. Or try :help.

scala> var a = 10
a: Int = 10

scala> var b = 20.5
b: Double = 20.5

scala> var c = "String"
c: String = String

scala>

scala> val d = 5000
d: Int = 5000

scala> val f = 'cat'
<console>:1: error: unclosed character literal
val f = 'cat'
        ^

scala> val f = "cat"
f: String = cat

scala>
```

```
scala> val i:Int = "hello"
<console>:11: error: type mismatch;
 found   : String("hello")
 required: Int
    val i:Int = "hello"
           ^

scala> val x = "hello"
x: String = hello

scala> x.re
reduce
reduceLeft
reduceLeftOption
reduceOption
reduceRight
reduceRightOption
regionMatches
replace
replaceAll
replaceAllLiterally
replaceFirst
repr
reverse
reverseIterator
reverseMap

scala> val x = new AnyRef{def helloWord = "Hello, world!"}
x: AnyRef{def helloWord: String} = $anon$1@58065f0c

scala> x.helloWord
def helloWord: String

scala> x.helloWord
warning: there was one feature warning; re-run with -feature for details
res0: String = Hello, world!

scala> _
```

```
scala> trait Logging < override def toString = "Logging " >
defined trait Logging

scala> class A extends Logging < override def toString = "A->" + super.toString
>
defined class A

scala> trait B extends Logging < override def toString = "B->" + super.toString
>
defined trait B

scala> trait C extends Logging < override def toString = "C->" + super.toString
>
defined trait C

scala> class D extends A with B with C < override def toString = "D->" + super.t
oString >
defined class D

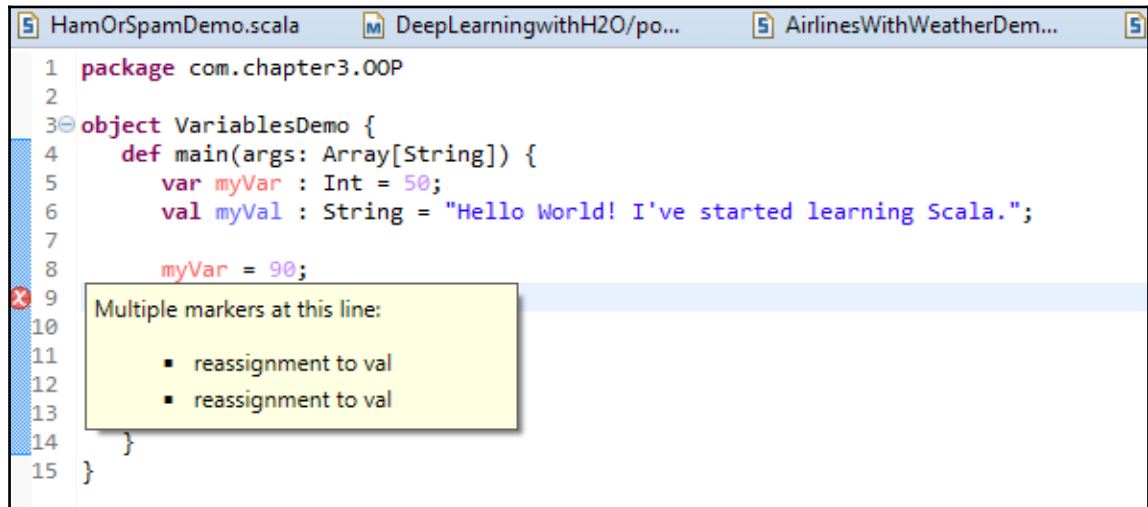
scala> new D()
res4: D = D->C->B->A->Logging

scala> _
```



```
>HelloWorld.scala x  
object HelloWorld{  
    def main(args:Array[String]){  
        println("Hello, world!")  
    }  
}
```

## Chapter 2: Object-Oriented Scala



The screenshot shows an IDE window with three tabs: "HamOrSpamDemo.scala", "DeepLearningwithH2O/po...", and "AirlinesWithWeatherDem...". The code in the first tab is as follows:

```
1 package com.chapter3.OOP
2
3 object VariablesDemo {
4     def main(args: Array[String]) {
5         var myVar : Int = 50;
6         val myVal : String = "Hello World! I've started learning Scala.";
7
8         myVar = 90;
9
10
11     }
12 }
13
14 }
15 }
```

A tooltip is displayed over line 9, indicating a compiler error: "Multiple markers at this line:"

- reassignment to val
- reassignment to val

```
C:\Windows\system32\cmd.exe - scala
C:\Users\rezkar>scala
Welcome to Scala 2.11.8 (Java HotSpot(TM) 64-Bit Server VM, Java 1.8.0_121).
Type in expressions for evaluation. Or try :help.

scala> case class Character(name: String, isHacker: Boolean)
defined class Character

scala> val nail = Character("Nail", true)
nail: Character = Character(Nail,true)

scala> val joyce = nail.copy(name = "Joyce")
joyce: Character = Character(Joyce,true)

scala> println(nail == joyce)
false

scala> println(nail.equals(joyce))
false

scala> println(nail.equals(nail))
true

scala> println(nail.hashCode())
-112671915

scala> println(nail.toString())
Character(Nail,true)

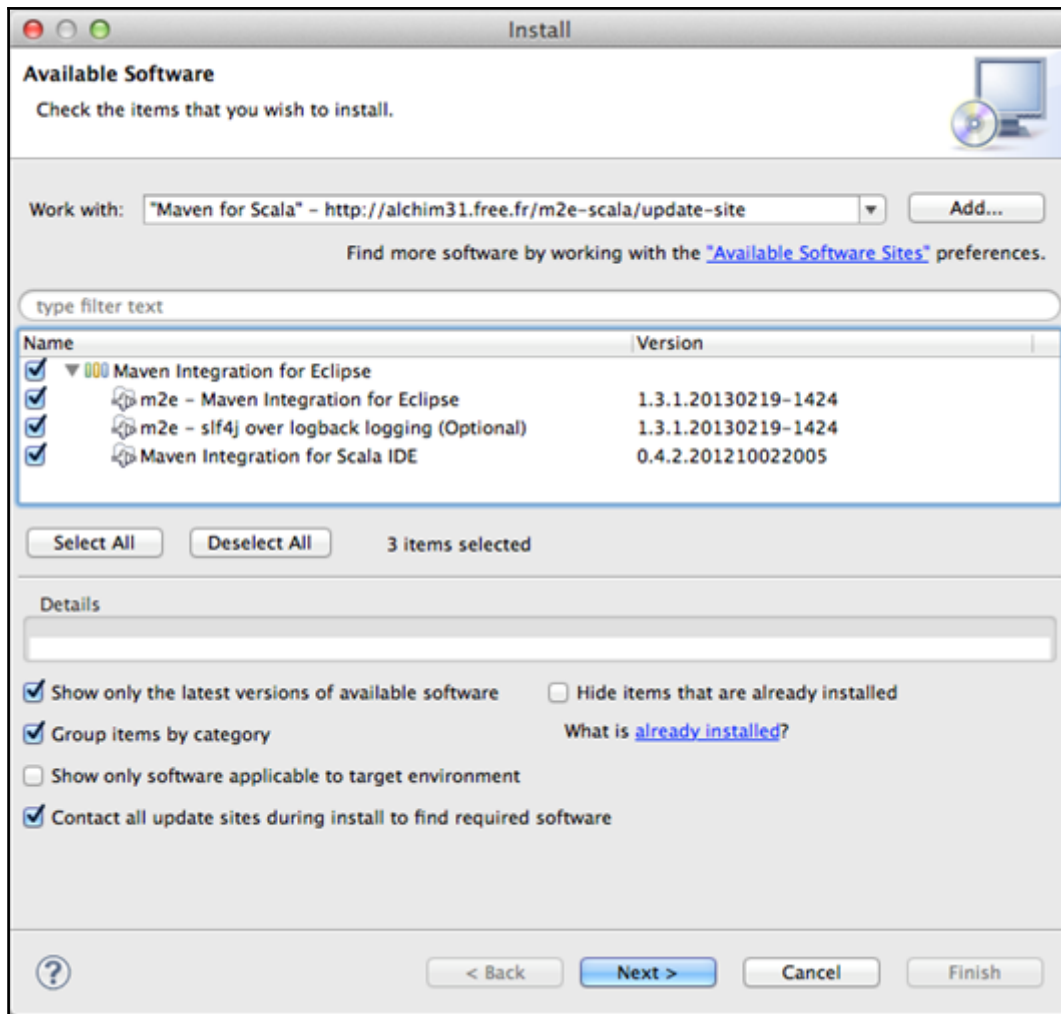
scala>    joyce match {
|         case Character(x, true) => s"$x is a hacker"
|         case Character(x, false) => s"$x is not a hacker"
|     }
res5: String = Joyce is a hacker

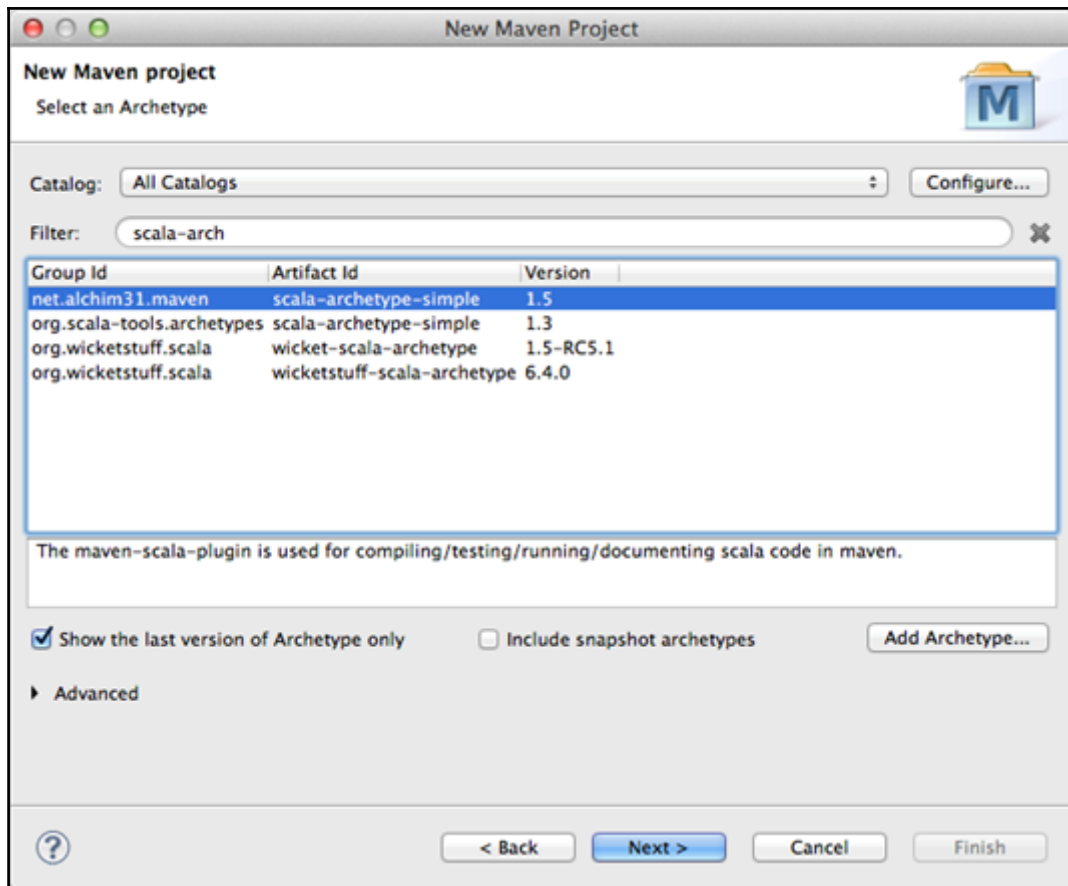
scala> _
```

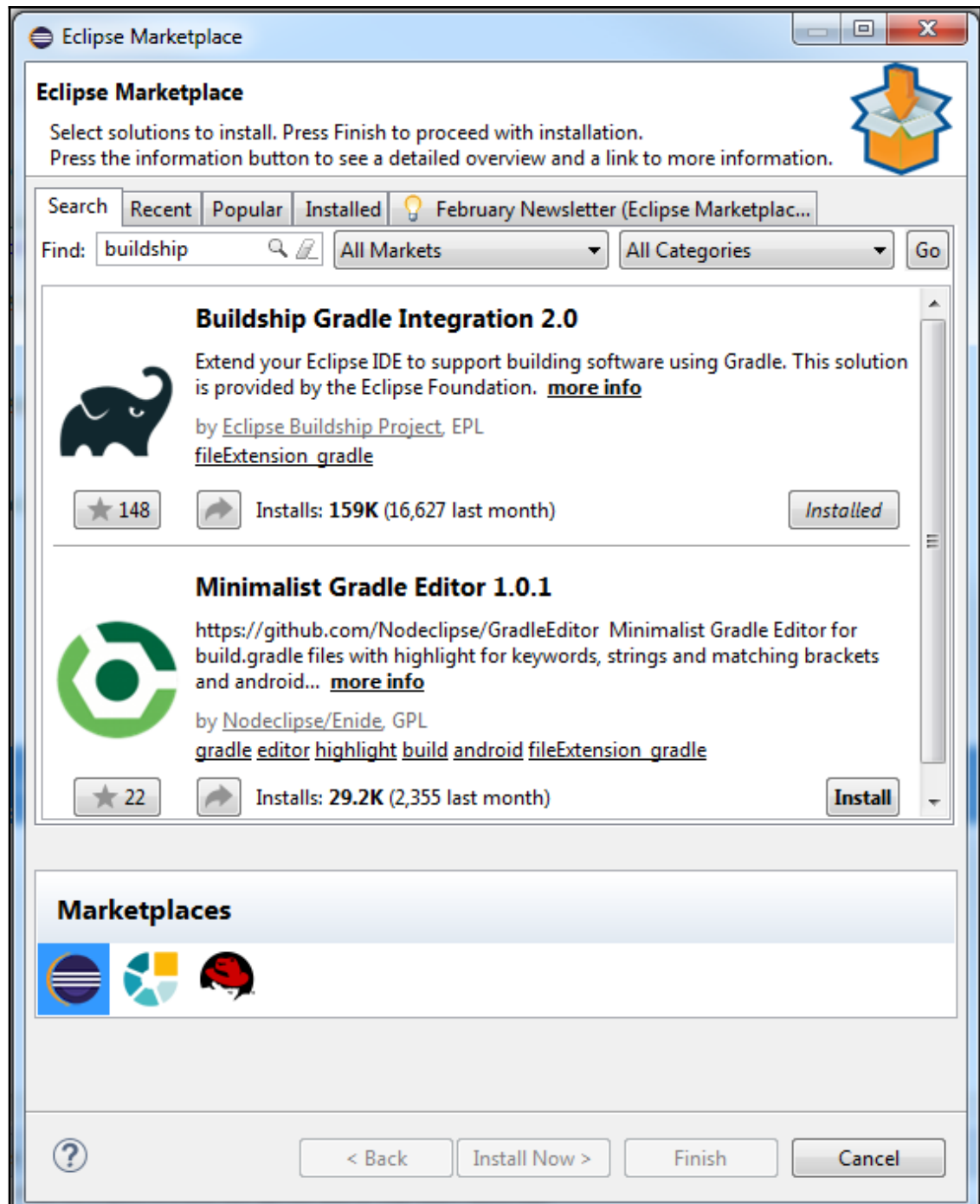
```
scala> object test { def printSomething() = {println("Inside an object")} }
defined object test

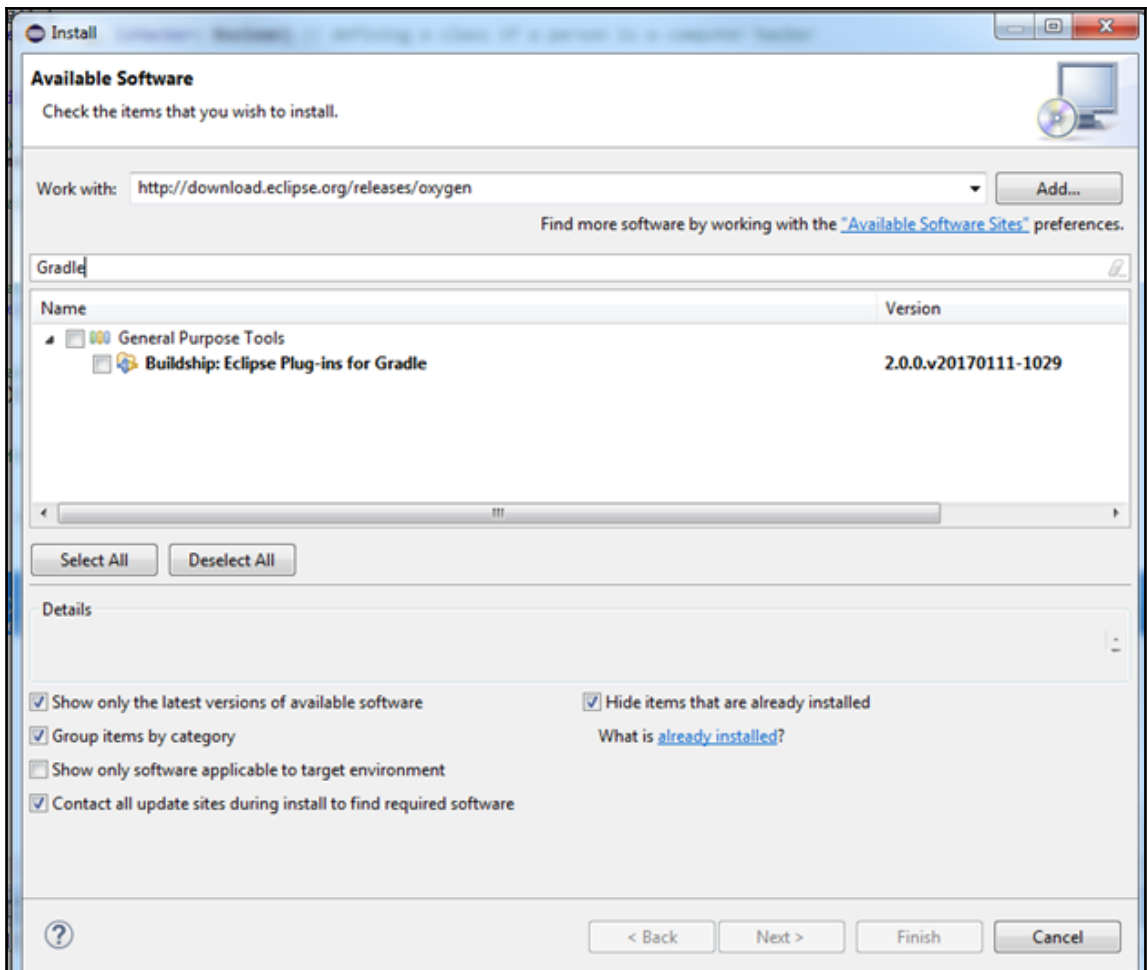
scala> test.printSomething
Inside an object

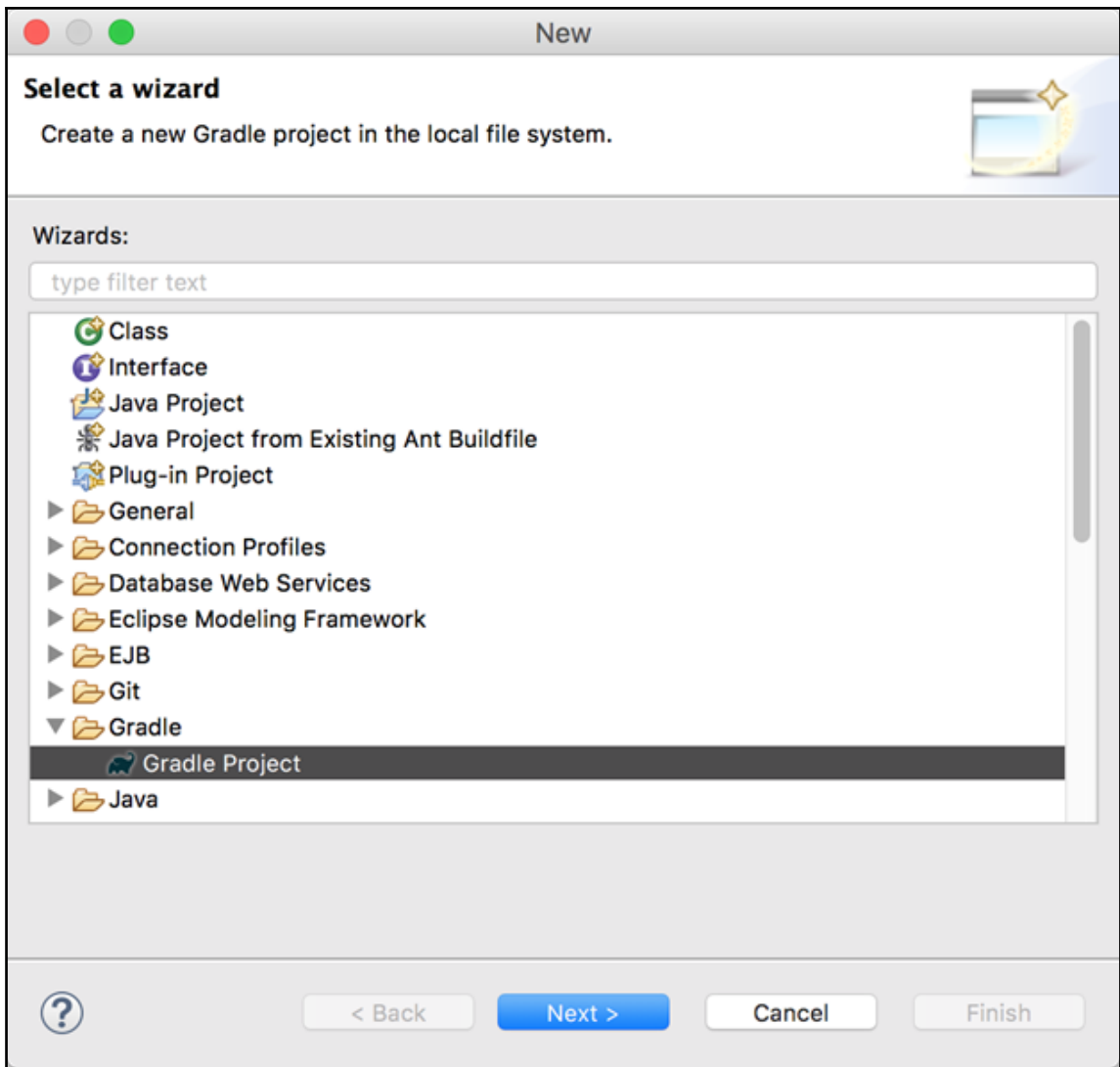
scala> val x = new test()
<console>:11: error: not found: type test
    val x = new test()
                ^
```



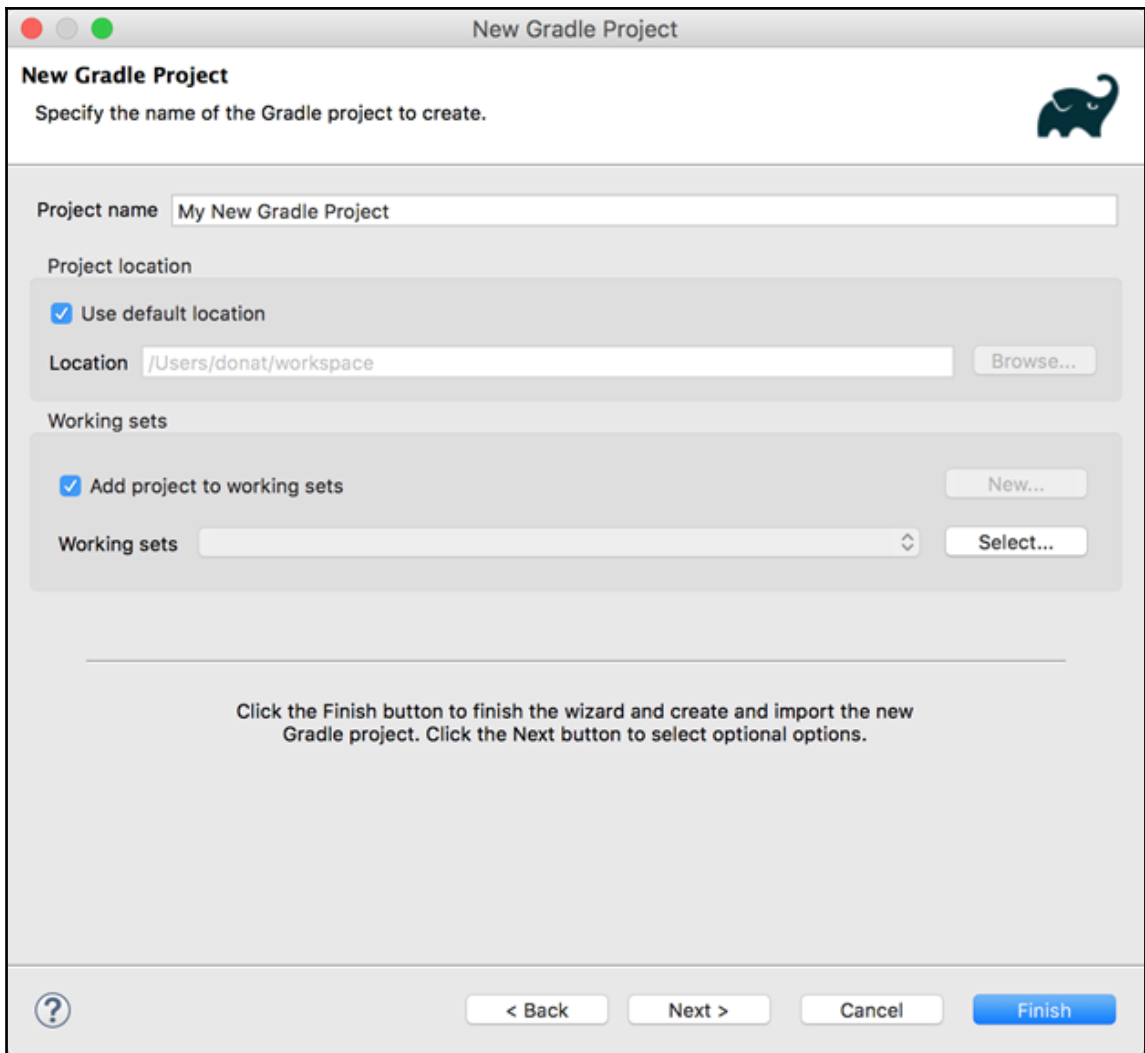


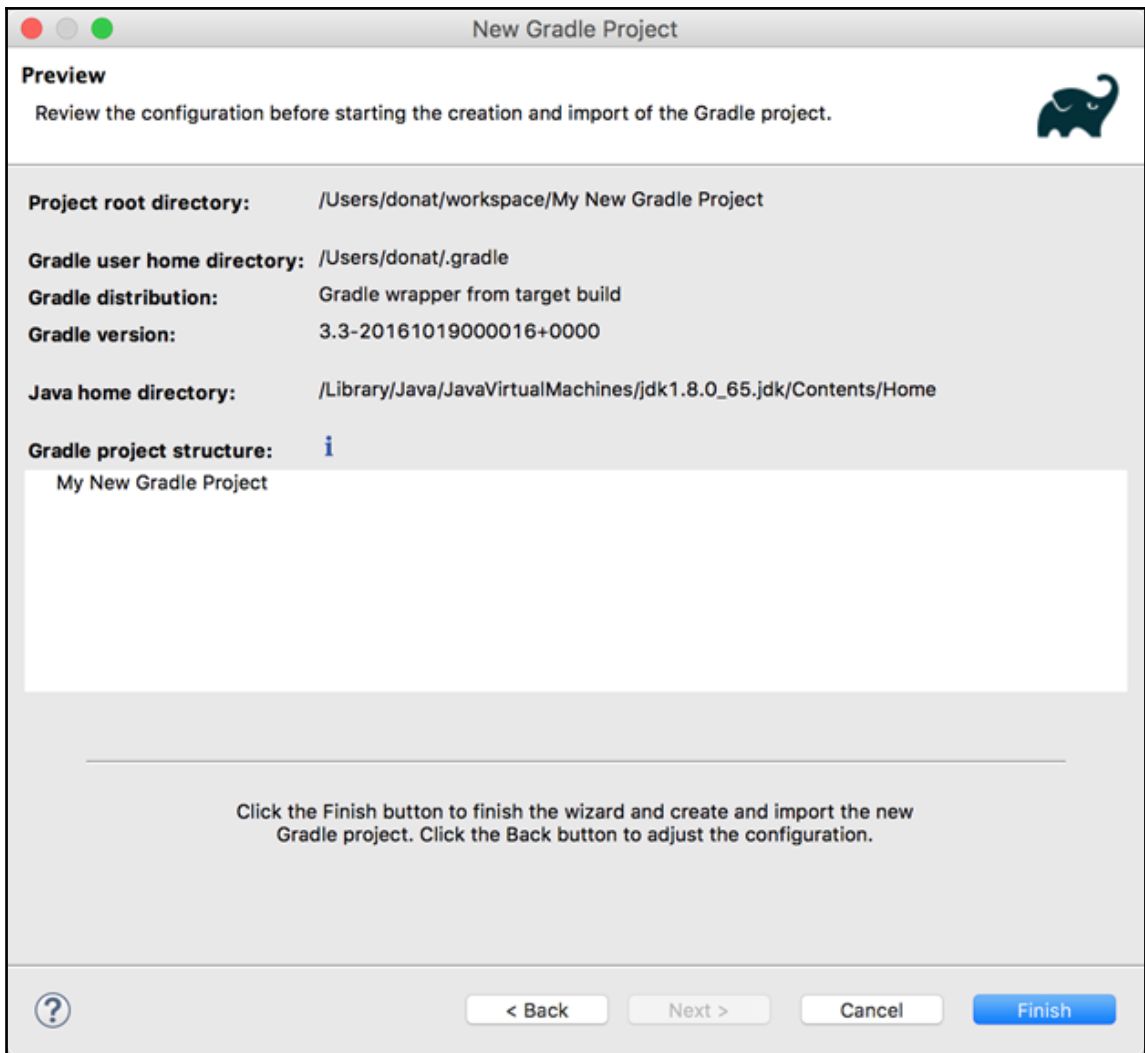


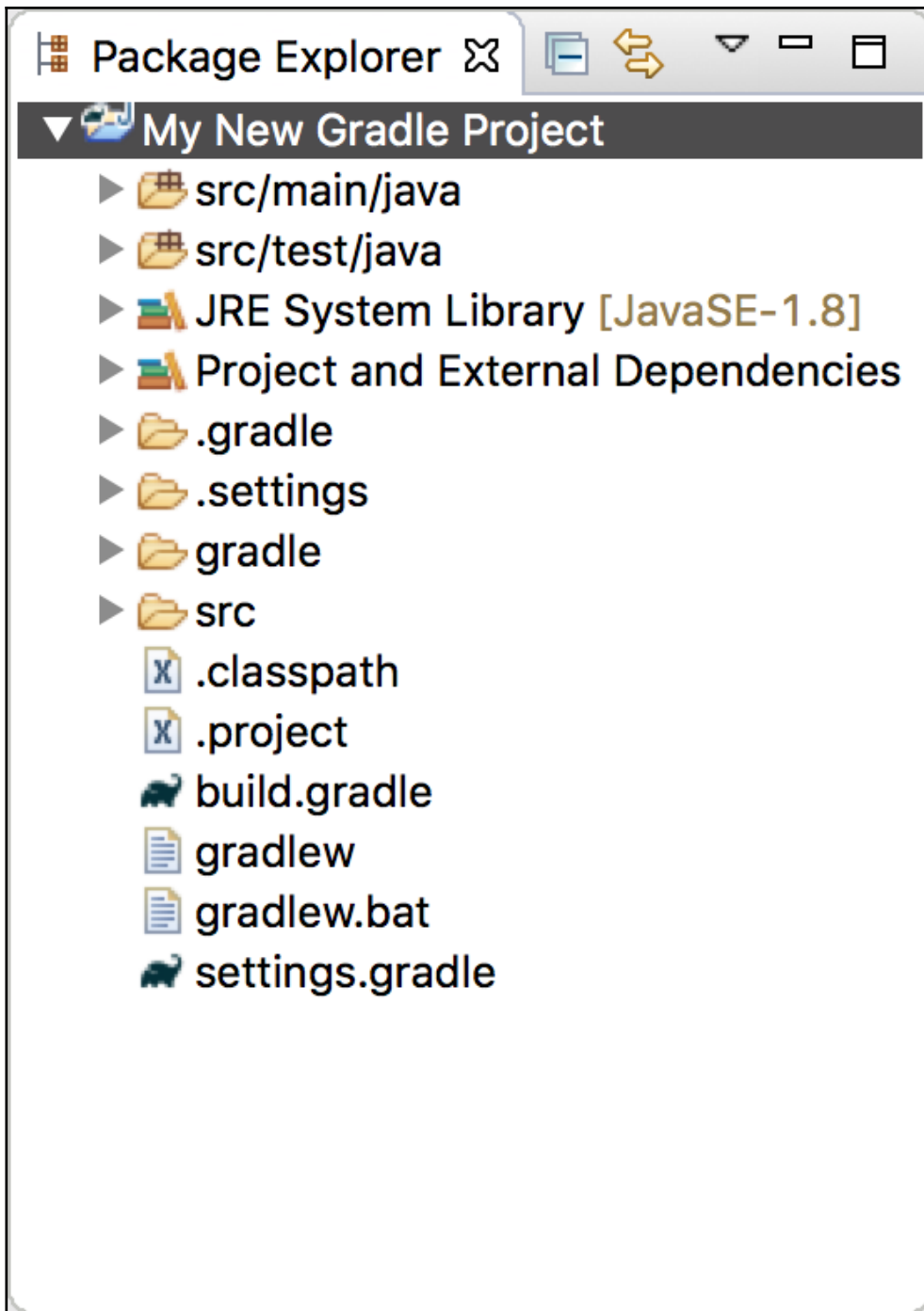




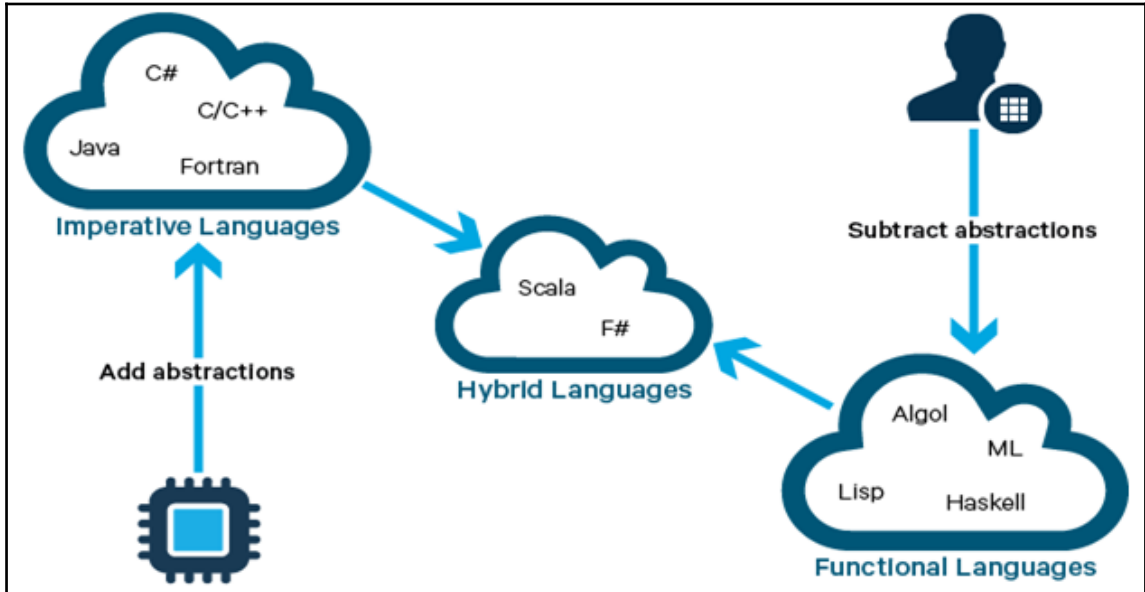








## Chapter 3: Functional programming concepts



```
scala> def quarterMaker(value: Int): Double = value.toDouble/4
quarterMaker: (value: Int)Double

scala> def addTwo(value: Int): Int = value + 2
addTwo: (value: Int)Int

scala> def applyFuncOnRange(begin: Int, end: Int, func: Int => AnyVal): Unit = {
  |   for (i <- begin to end)
  |     println(func(i))
  | }
applyFuncOnRange: (begin: Int, end: Int, func: Int => AnyVal)Unit

scala>
scala>
scala>
scala>
scala>
scala>
```

```
scala> applyFuncOnRange(1,10,quarterMaker)
0.25
0.5
0.75
1.0
1.25
1.5
1.75
2.0
2.25
2.5

scala>

scala>

scala>

scala>

scala>

scala> █
```

```
scala> applyFuncOnRange(1,10,addTwo)
3
4
5
6
7
8
9
10
11
12

scala>

scala>

scala>

scala>

scala>

scala> █
```

```
scala> def bankFee(amount: Double) = amount * 0.05
bankFee: (amount: Double)Double

scala> def TransferMoney(money: Double, bankFee: Double => Double): Double = {
  |   money + bankFee(money)
  | }
TransferMoney: (money: Double, bankFee: Double => Double)Double

scala> TransferMoney(100, bankFee)
res2: Double = 105.0

scala>

scala>

scala>

scala>

scala>

scala>

scala> █
```

```
scala> def TransferMoney(money: Double, bankFee: Double => Double): Double = {
  |   money + bankFee(money)
  | }
TransferMoney: (money: Double, bankFee: Double => Double)Double

scala> TransferMoney(100, (amount: Double) => amount * 0.05)
res12: Double = 105.0

scala> TransferMoney(100, amount => amount * 0.05)
res13: Double = 105.0

scala>

scala>

scala>

scala>

scala>

scala>

scala> █
```

```
scala> def TransferMoney(money: Double) = {
  |   if (money > 1000)
  |   (money: Double) => "Dear customer we are going to add the following amount as Fee
: "+money * 0.05
  |   else
  |   (money: Double) => "Dear customer we are going to add the following amount as Fee
: "+money * 0.1
  | }
TransferMoney: (money: Double)Double => String

scala> val returnedFunction = TransferMoney(1500)
returnedFunction: Double => String = <function1>

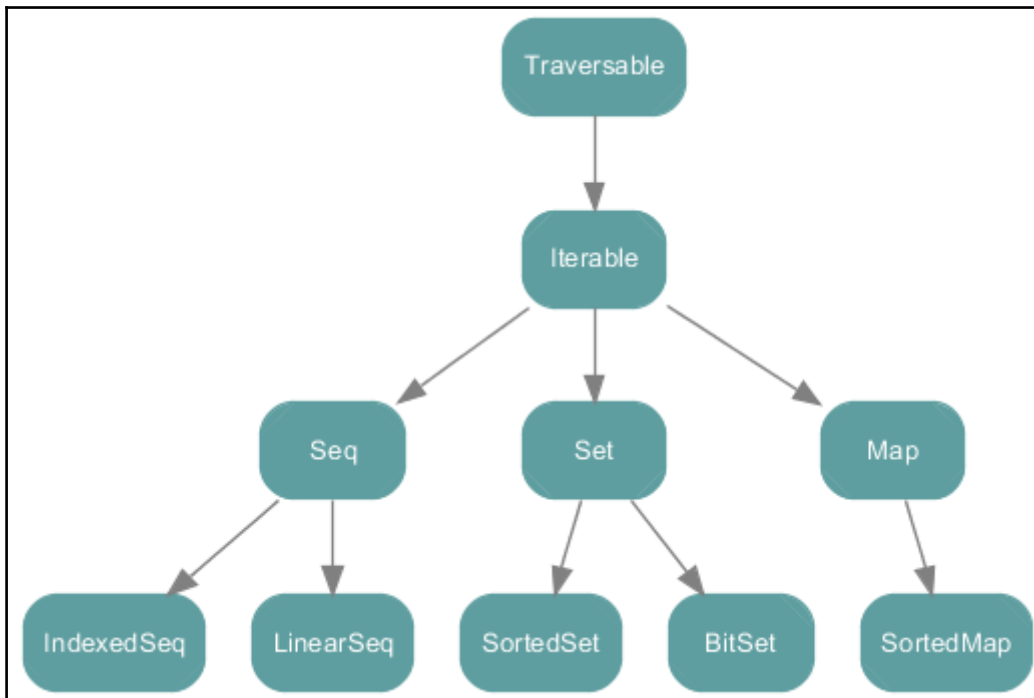
scala>
  | returnedFunction(1500)
res17: String = Dear customer we are going to add the following amount as Fee: 75.0

scala>

scala>

scala>

scala> █
```



```
scala> val evenList = List(2,4,6,8,10)
evenList: List[Int] = List(2, 4, 6, 8, 10)

scala> evenList.map(x => x * 2 )
res18: List[Int] = List(4, 8, 12, 16, 20)

scala>

scala>

scala>

scala>

scala>

scala>

scala>

scala>

scala> █
```

```
scala> def func(x : Int) = if(x > 4) Some(x) else None
func: (x: Int)Option[Int]

scala> evenList.map(x => func(x))
res19: List[Option[Int]] = List(None, None, Some(6), Some(8), Some(10))

scala>

scala>

scala>

scala>

scala>

scala>

scala>

scala>

scala> █
```



```
scala> def around(x : Int) = List(x-1, x, x+1)
around: (x: Int)List[Int]

scala> evenList.map(x => around(x))
res23: List[List[Int]] = List(List(1, 2, 3), List(3, 4, 5), List(5, 6, 7), List(7, 8, 9), List(9, 10, 11))

scala> evenList.flatMap(x => around(x))
res24: List[Int] = List(1, 2, 3, 3, 4, 5, 5, 6, 7, 7, 8, 9, 9, 10, 11)

scala>
scala>
scala>
scala>
scala>
scala>
scala>
scala>
scala>
```

```
scala> val range = List.range(1,10)
range: List[Int] = List(1, 2, 3, 4, 5, 6, 7, 8, 9)

scala> val odds = range.filter(_ % 2 != 0)
odds: List[Int] = List(1, 3, 5, 7, 9)

scala>
scala>
scala>
scala>
scala>
scala>
scala>
scala>
scala>
scala>
```

```
scala> for (x <- 10 until (0, -2))
  |   yield x
res25: scala.collection.immutable.IndexedSeq[Int] = Vector(10, 8, 6, 4, 2)
scala>
scala>
scala>
scala>
scala>
scala>
scala>
scala>
scala>
scala>
```

```
scala> for (x <- 1 to 10 if x % 2 == 0)
  |   yield x
res26: scala.collection.immutable.IndexedSeq[Int] = Vector(2, 4, 6, 8, 10)
scala>
scala>
scala>
scala>
scala>
scala>
scala>
scala>
scala>
scala>
```

```
scala> for (x <- 1 to 10; y <- 1 until x)
  |   yield (x, y)
res29: scala.collection.immutable.IndexedSeq[(Int, Int)] = Vector((2,1), (3,1), (3,2), (4,1)
, (4,2), (4,3), (5,1), (5,2), (5,3), (5,4), (6,1), (6,2), (6,3), (6,4), (6,5), (7,1), (7,2),
(7,3), (7,4), (7,5), (7,6), (8,1), (8,2), (8,3), (8,4), (8,5), (8,6), (8,7), (9,1), (9,2),
(9,3), (9,4), (9,5), (9,6), (9,7), (9,8), (10,1), (10,2), (10,3), (10,4), (10,5), (10,6), (1
0,7), (10,8), (10,9))

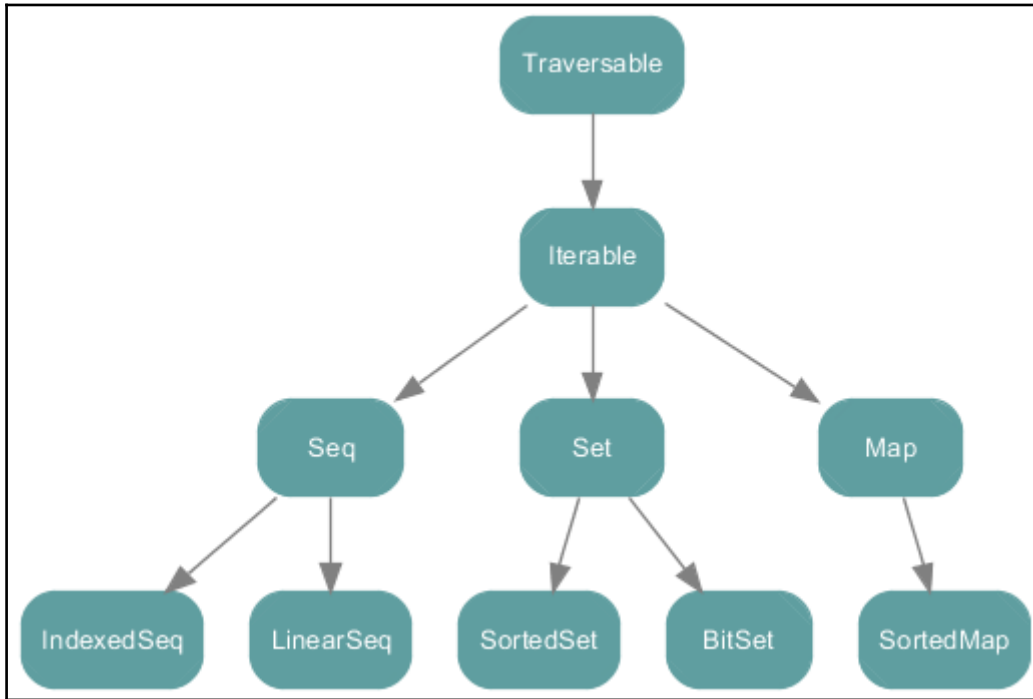
scala>

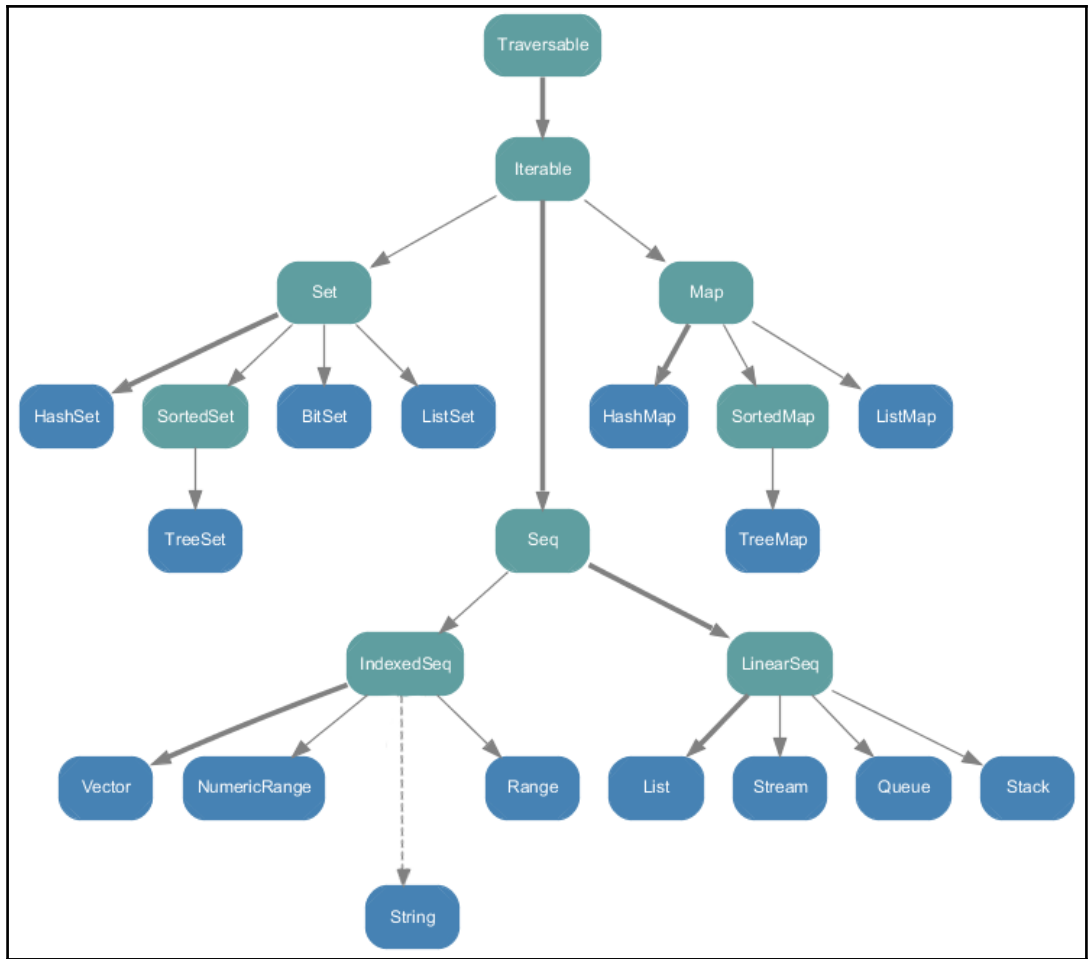
scala> (1 to 10).flatMap(
  |   i => (1 until i).map(
  |     j => (i, j)
  |   )
  | )
res30: scala.collection.immutable.IndexedSeq[(Int, Int)] = Vector((2,1), (3,1), (3,2), (4,1)
, (4,2), (4,3), (5,1), (5,2), (5,3), (5,4), (6,1), (6,2), (6,3), (6,4), (6,5), (7,1), (7,2),
(7,3), (7,4), (7,5), (7,6), (8,1), (8,2), (8,3), (8,4), (8,5), (8,6), (8,7), (9,1), (9,2),
(9,3), (9,4), (9,5), (9,6), (9,7), (9,8), (10,1), (10,2), (10,3), (10,4), (10,5), (10,6), (1
0,7), (10,8), (10,9))

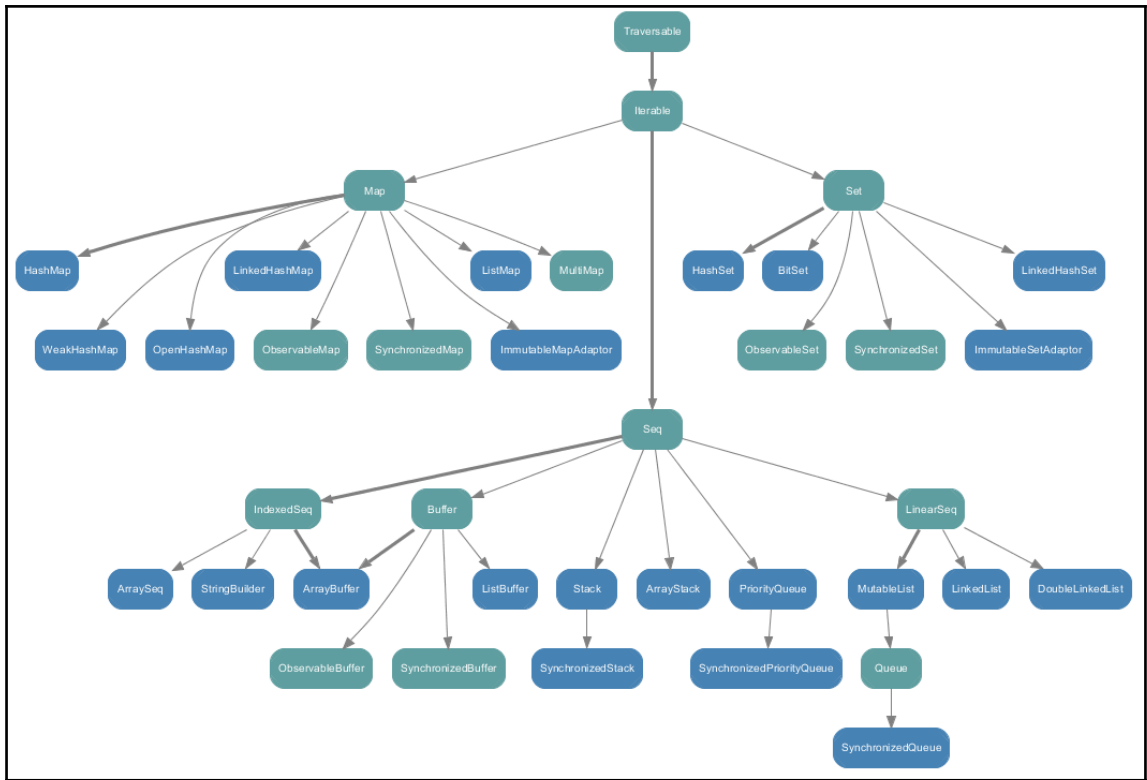
scala>

scala> █
```

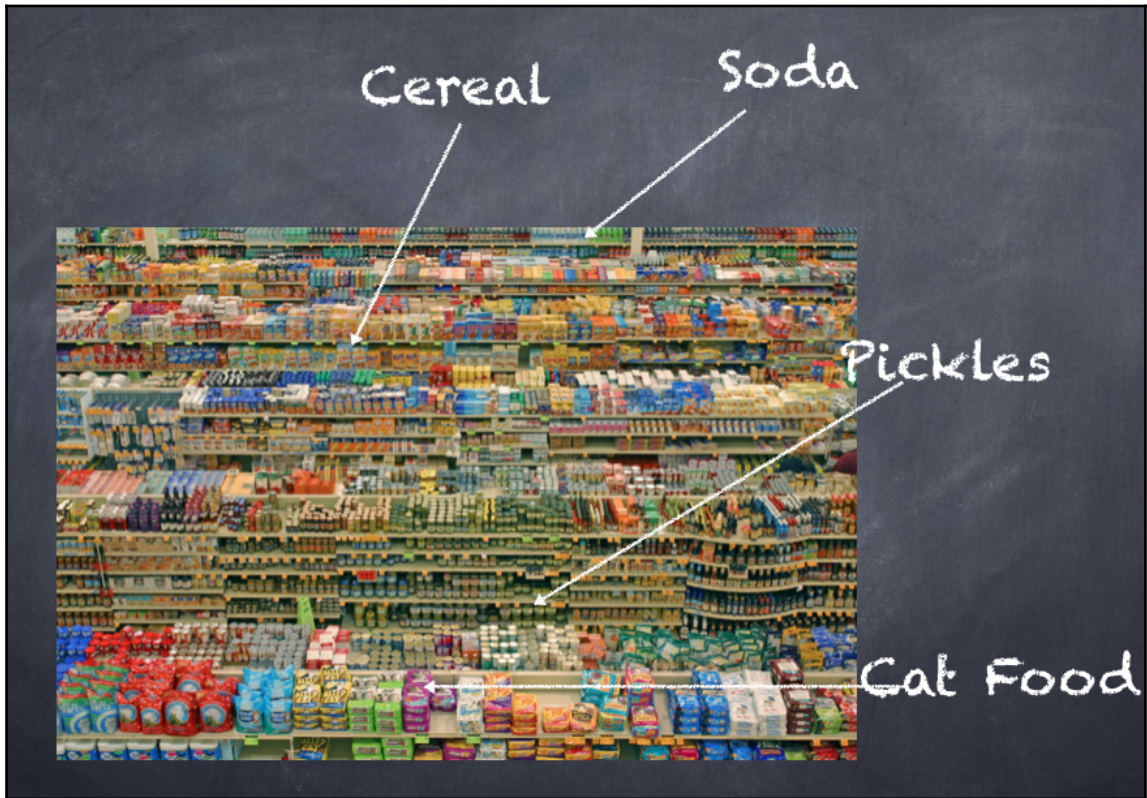
## Chapter 4: Collections APIs







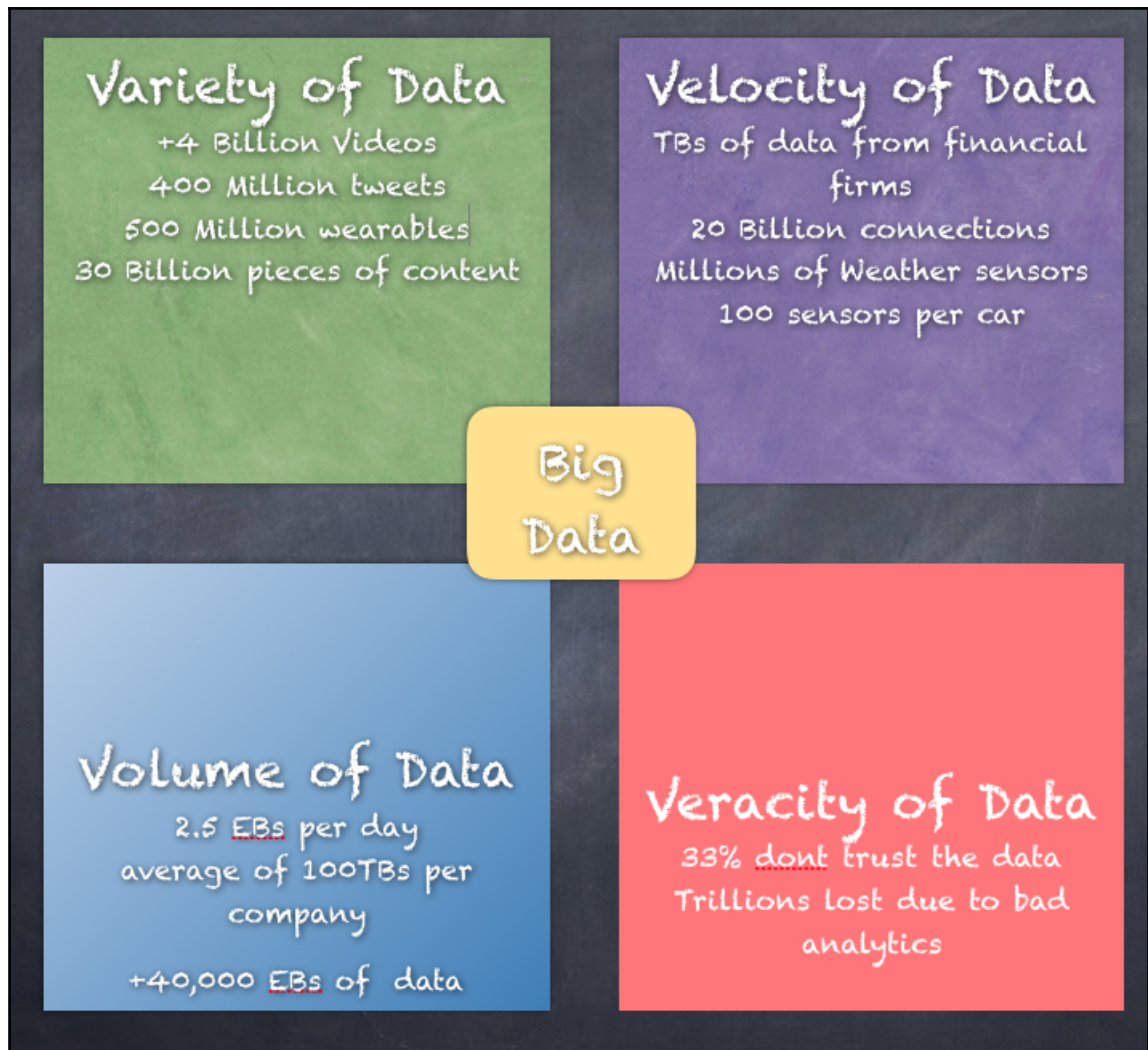
## Chapter 5: Tackle Big Data – Spark Comes to the Party

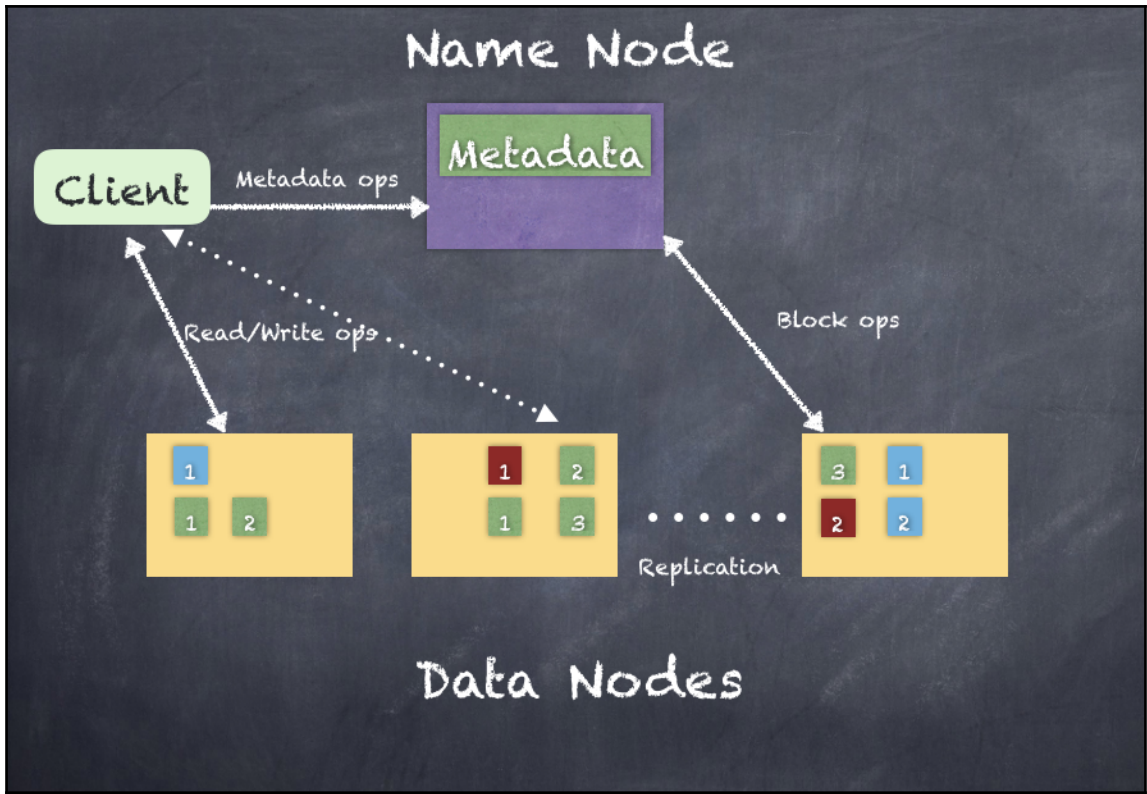


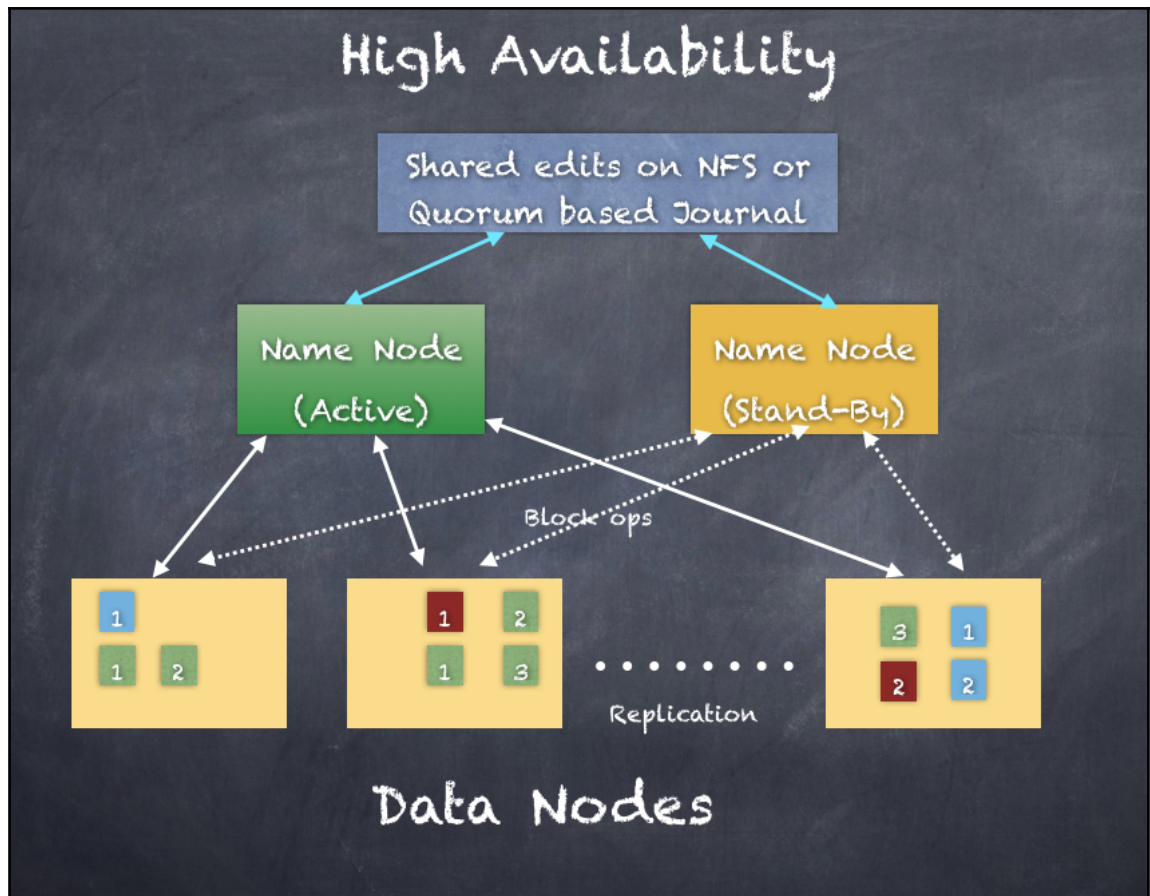


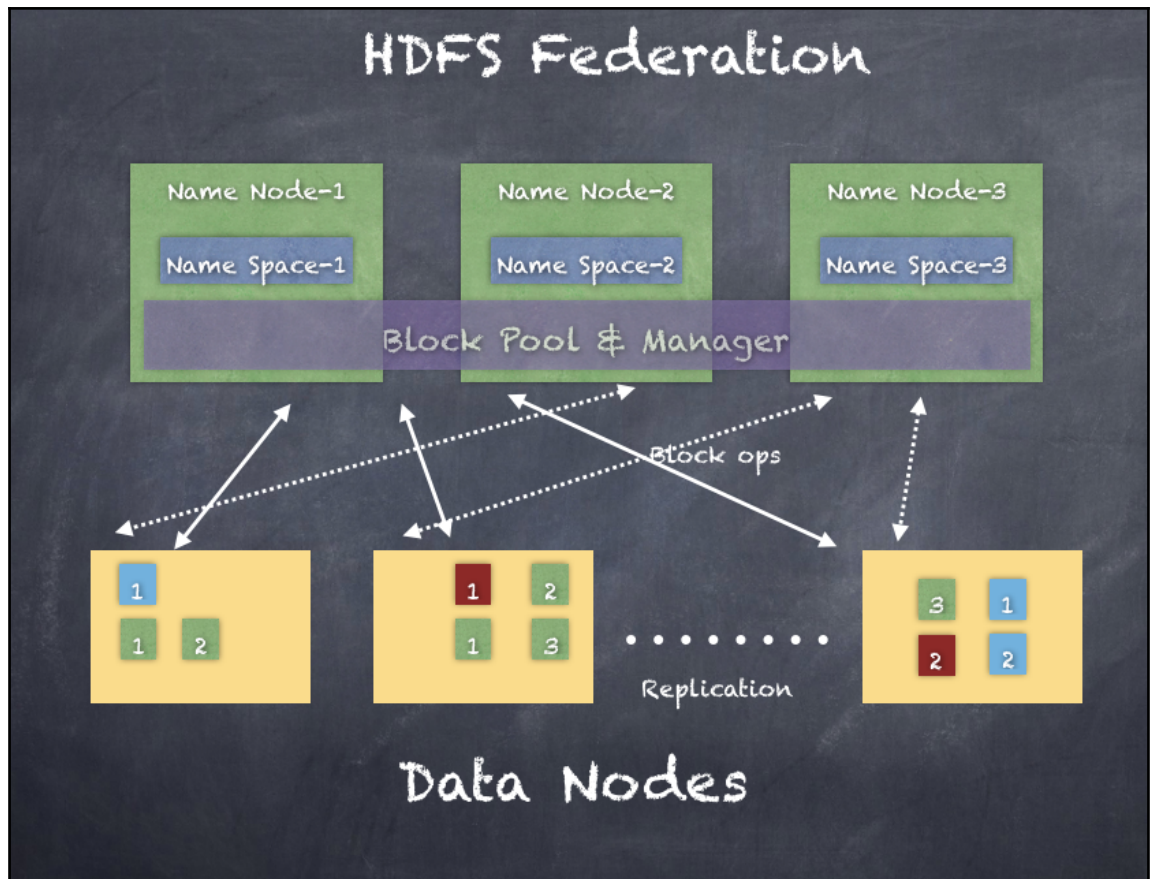
Receipt No	Date	ItemId	Quantity	Cost\$	Sale\$
123	03/01	24	3	9	12
123	03/01	25	2	6	7
124	03/02	12	4	7	9
125	03/02	25	1	3	4

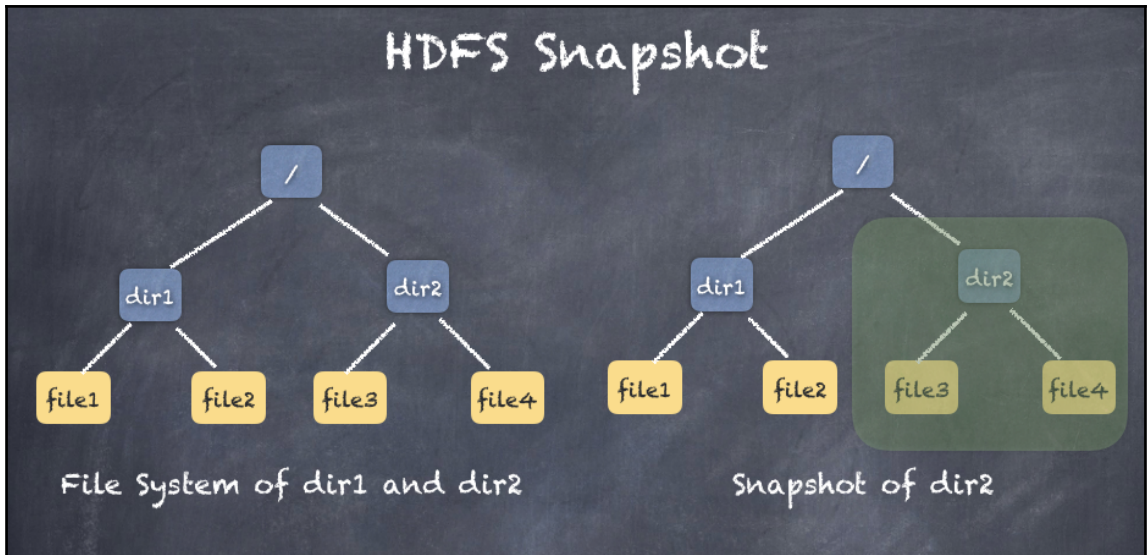


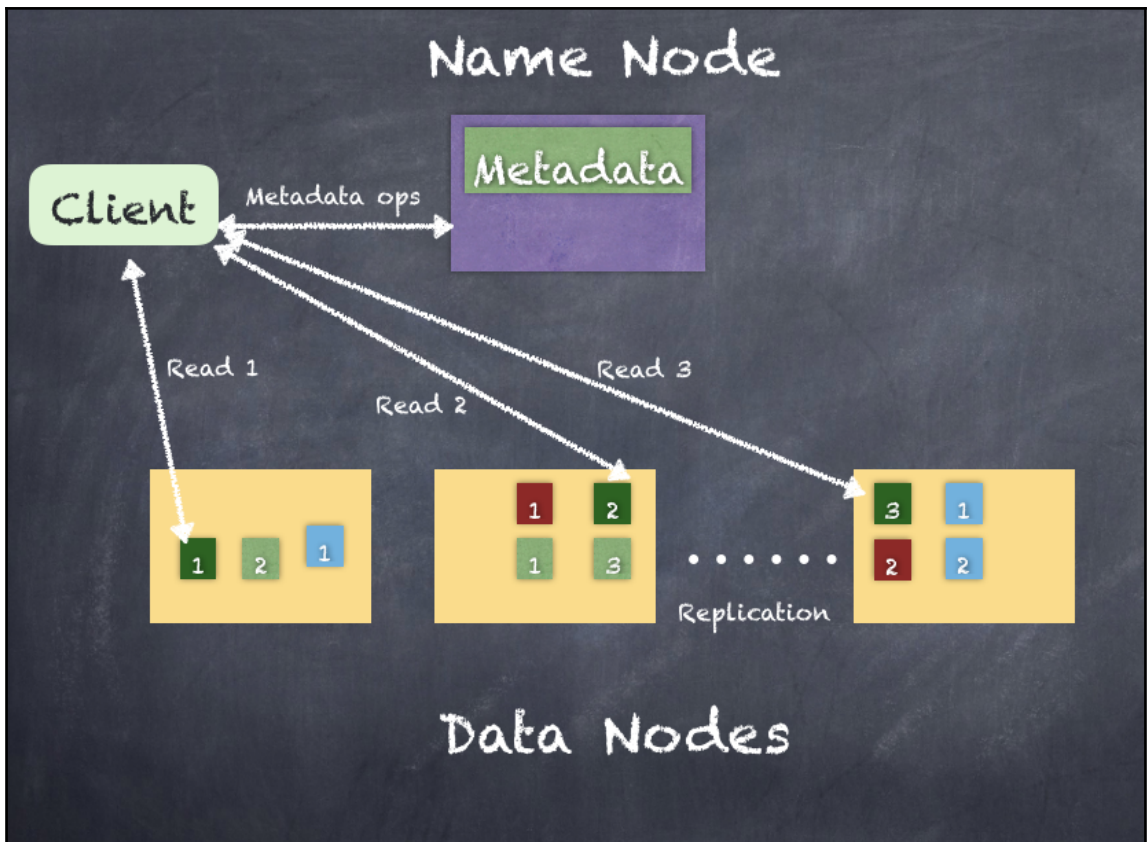


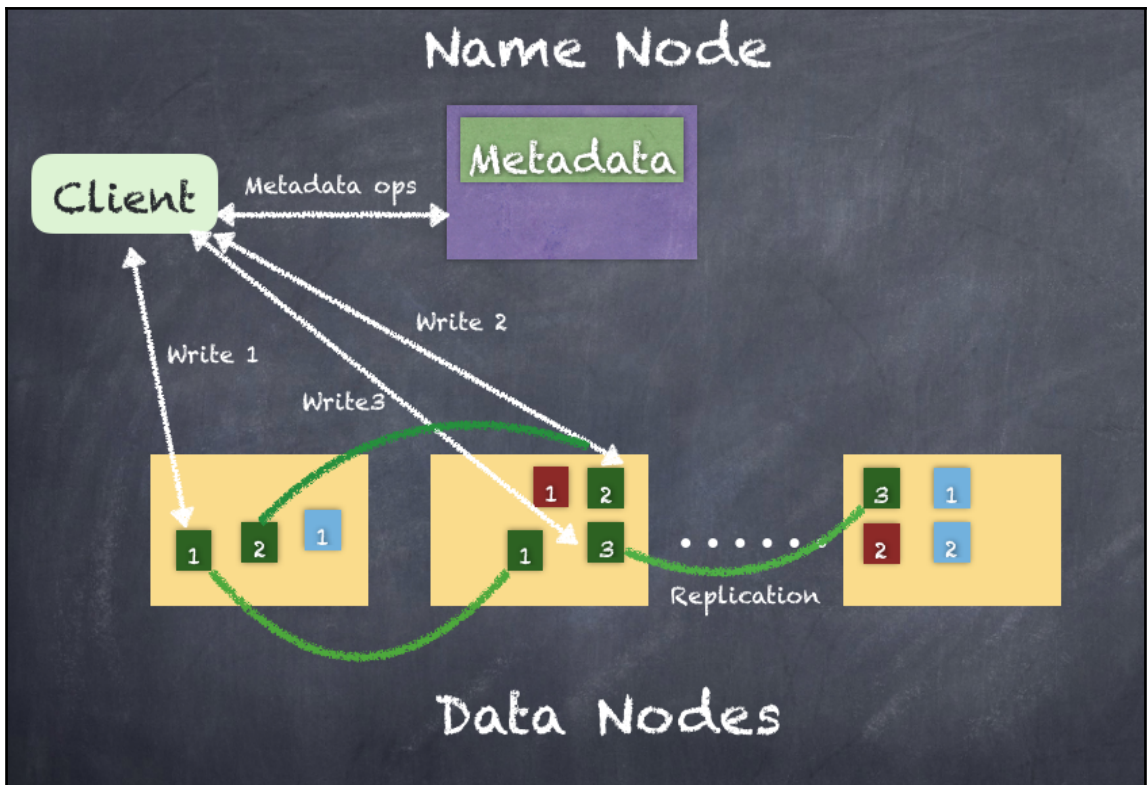








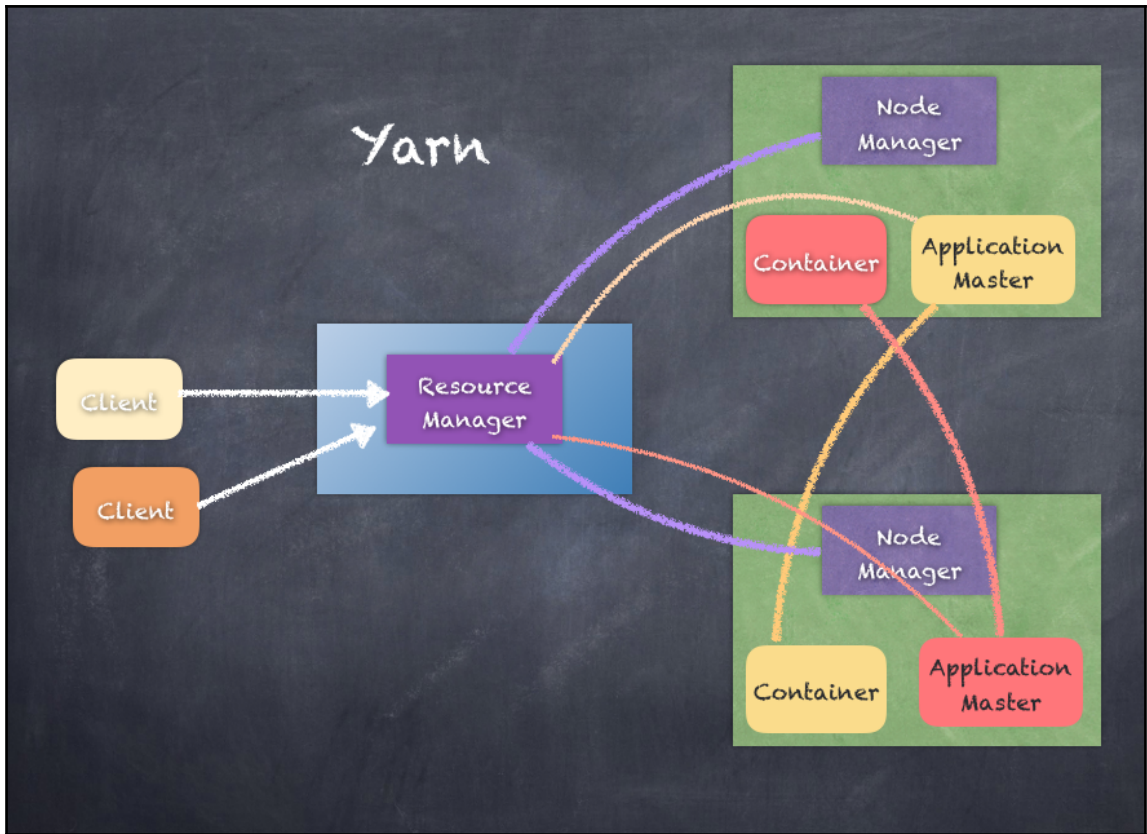


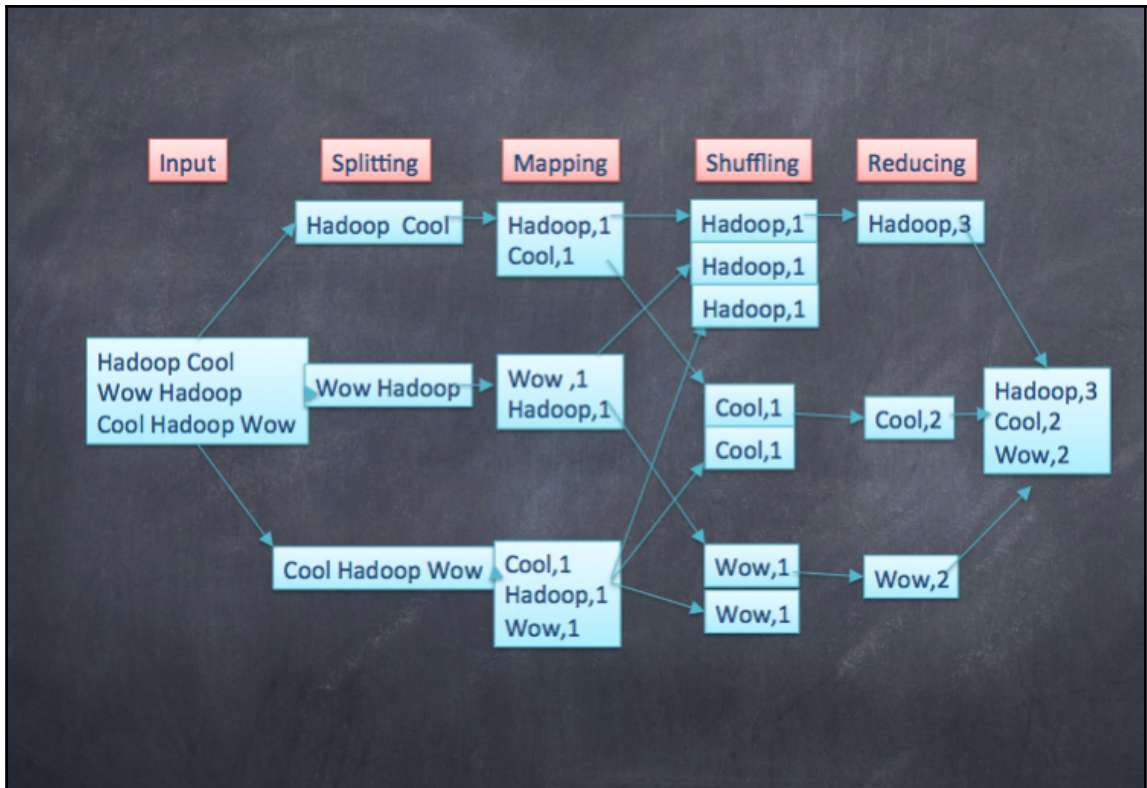


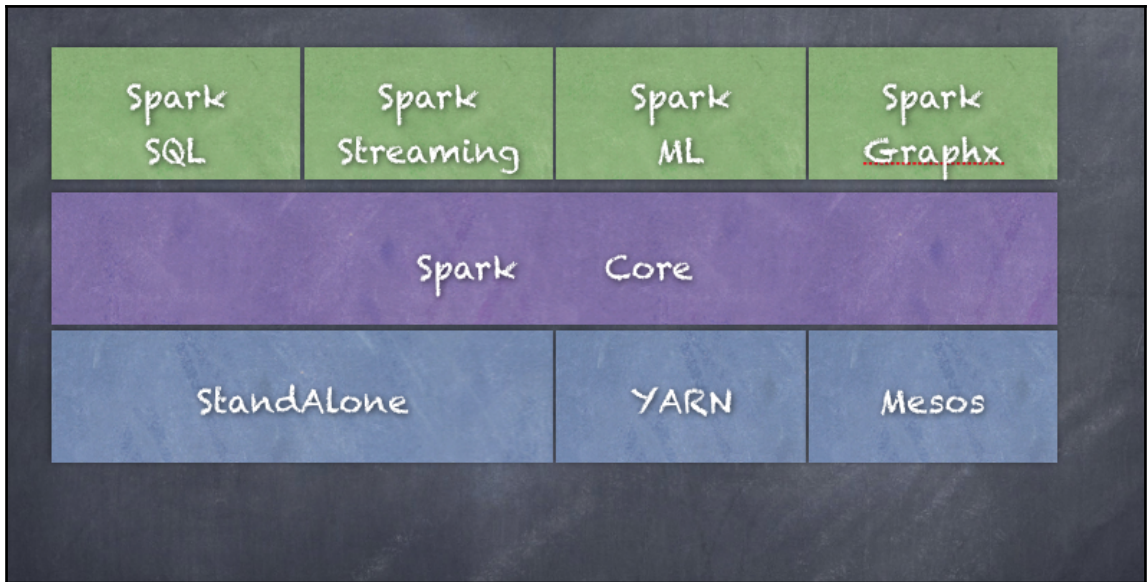
# Sorting the fruits by type

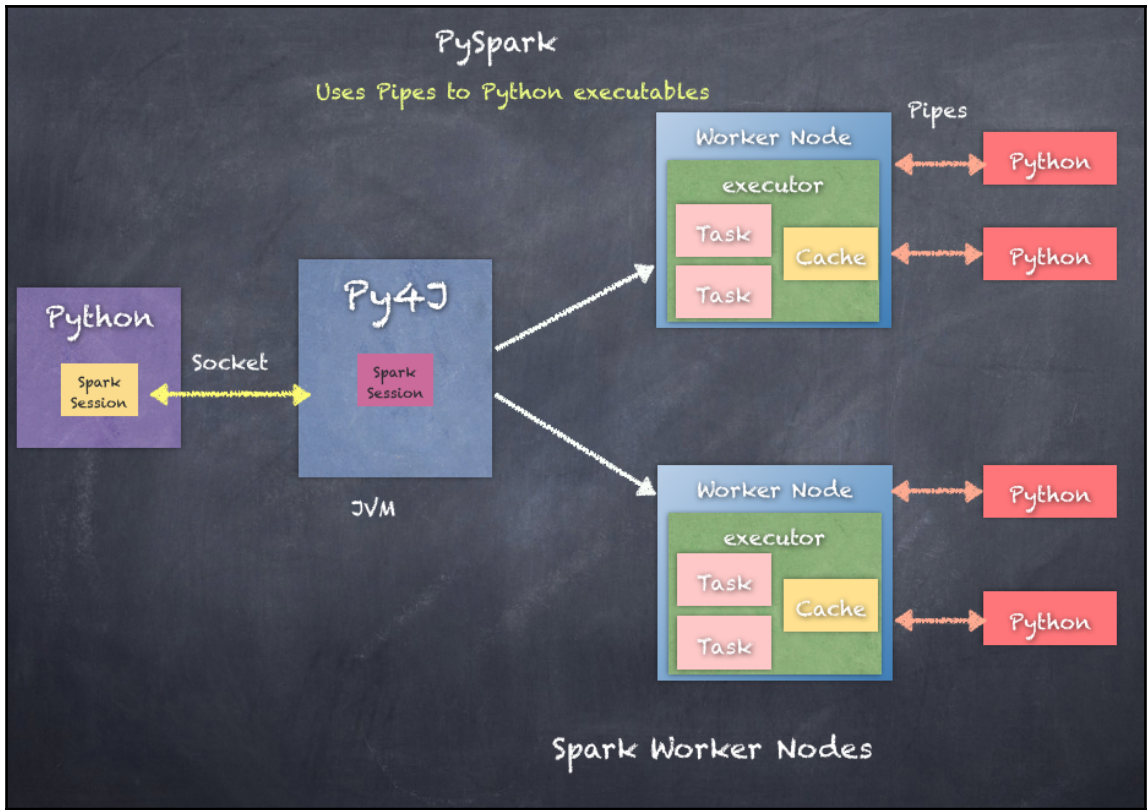


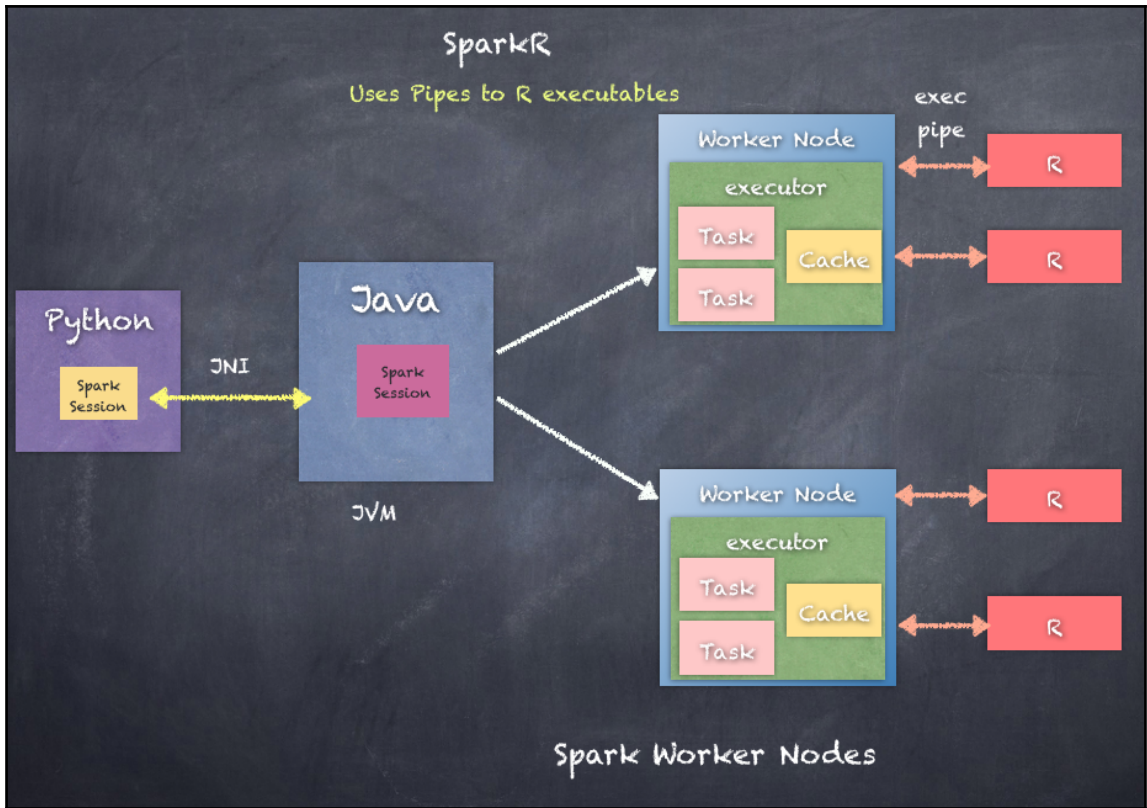




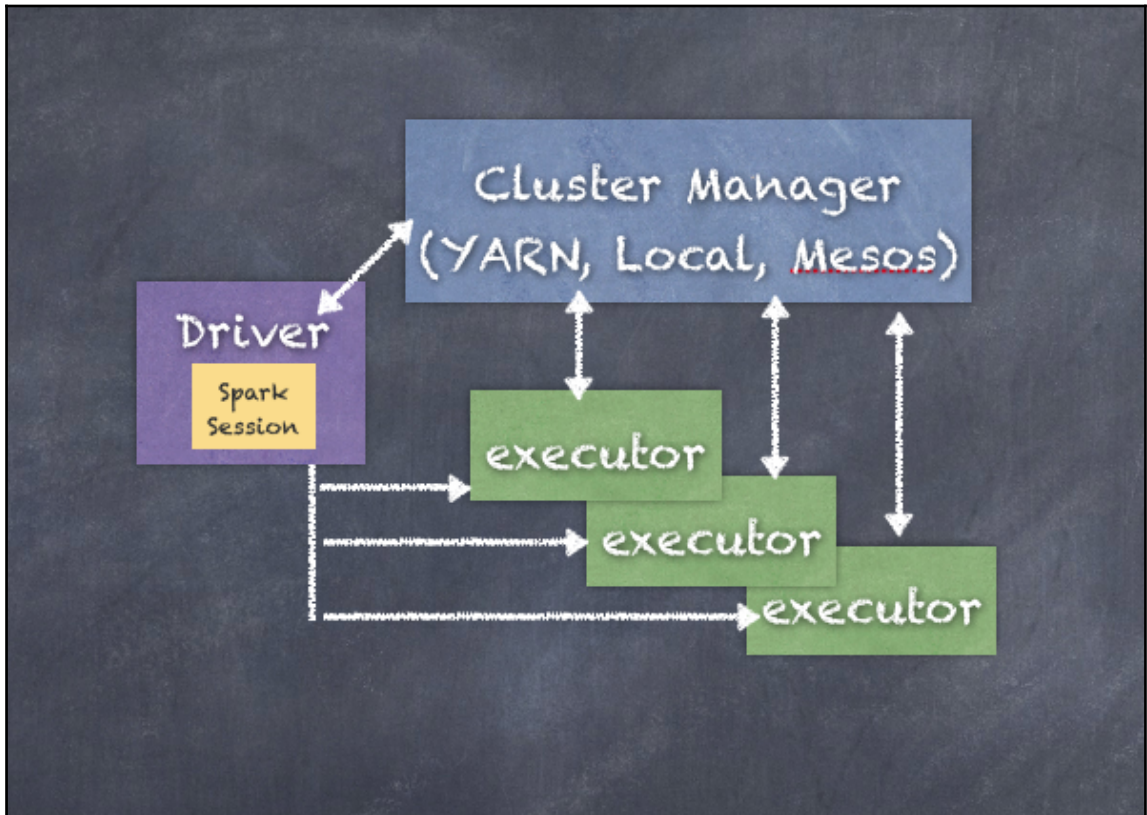


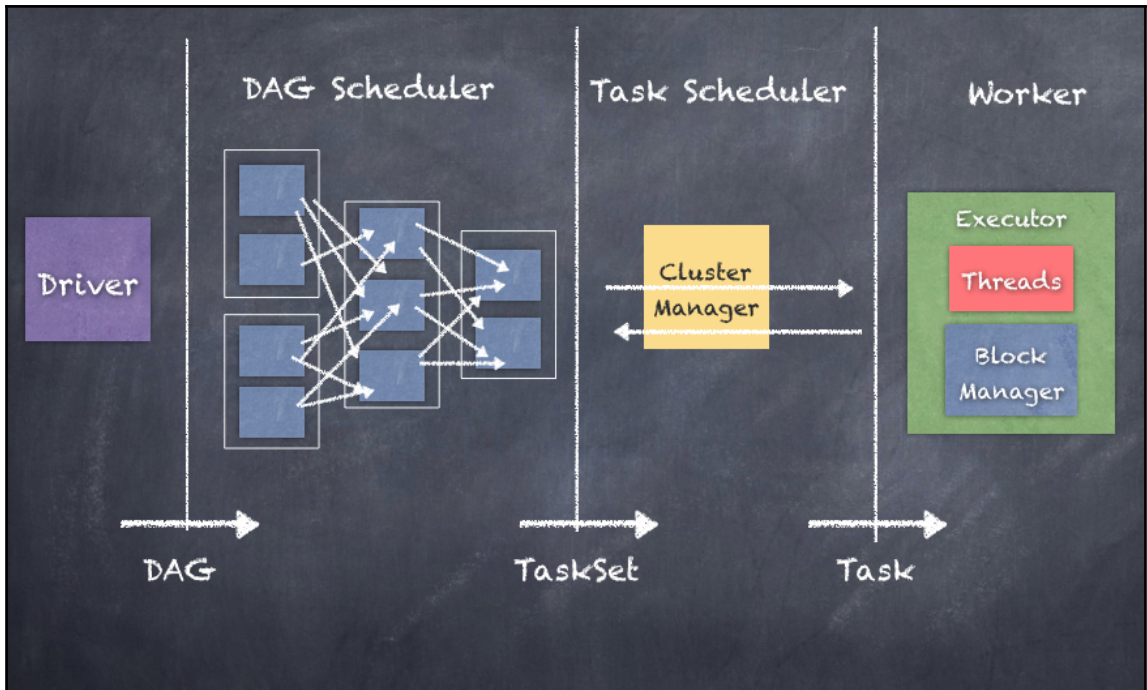


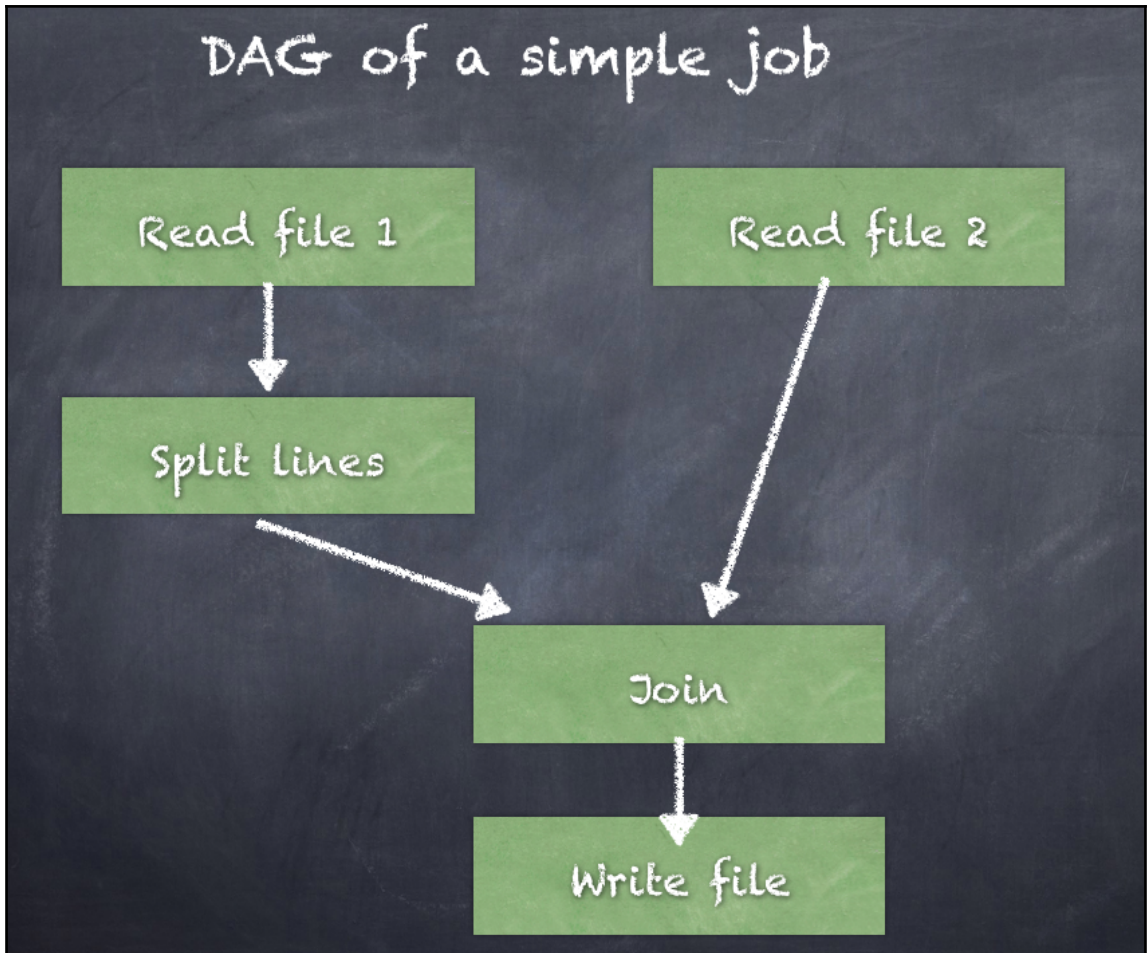




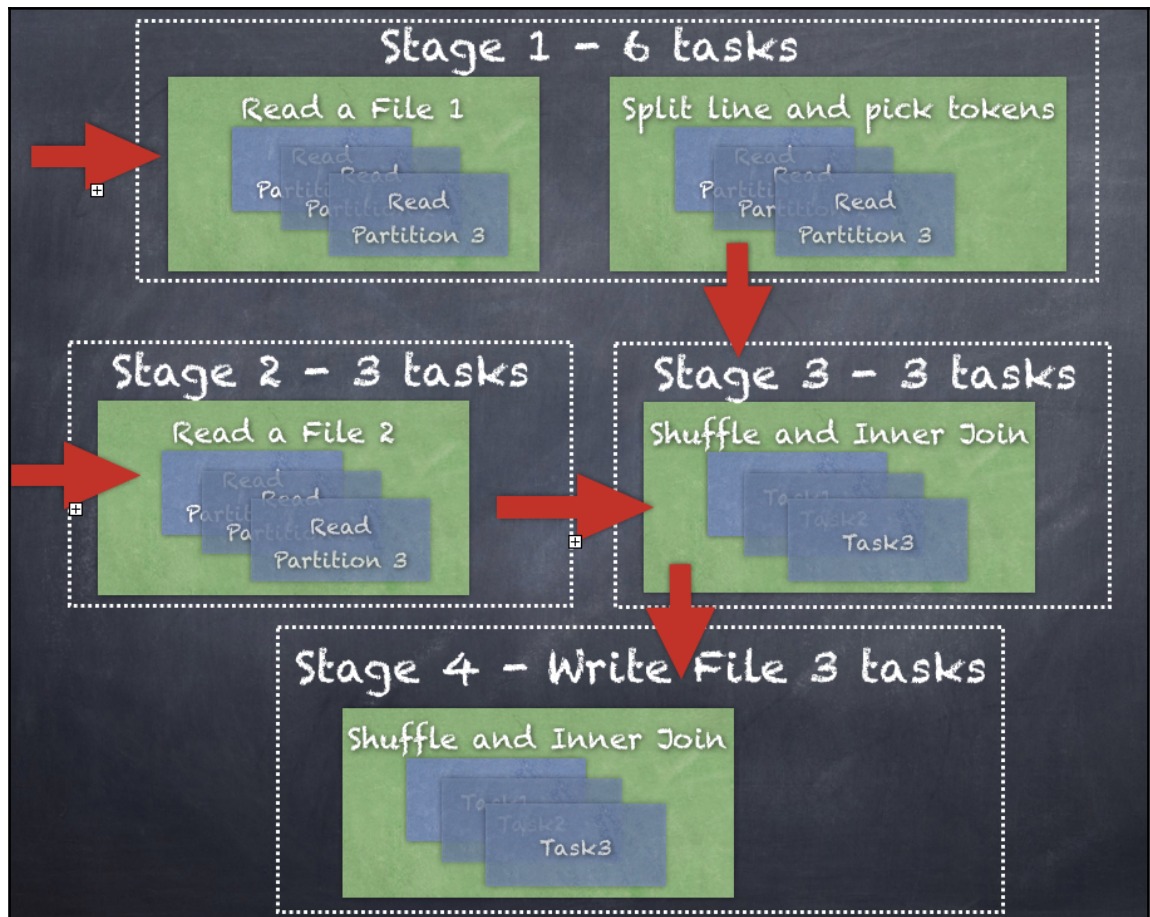
## Chapter 6: Start Working with Spark REPL and RDDs

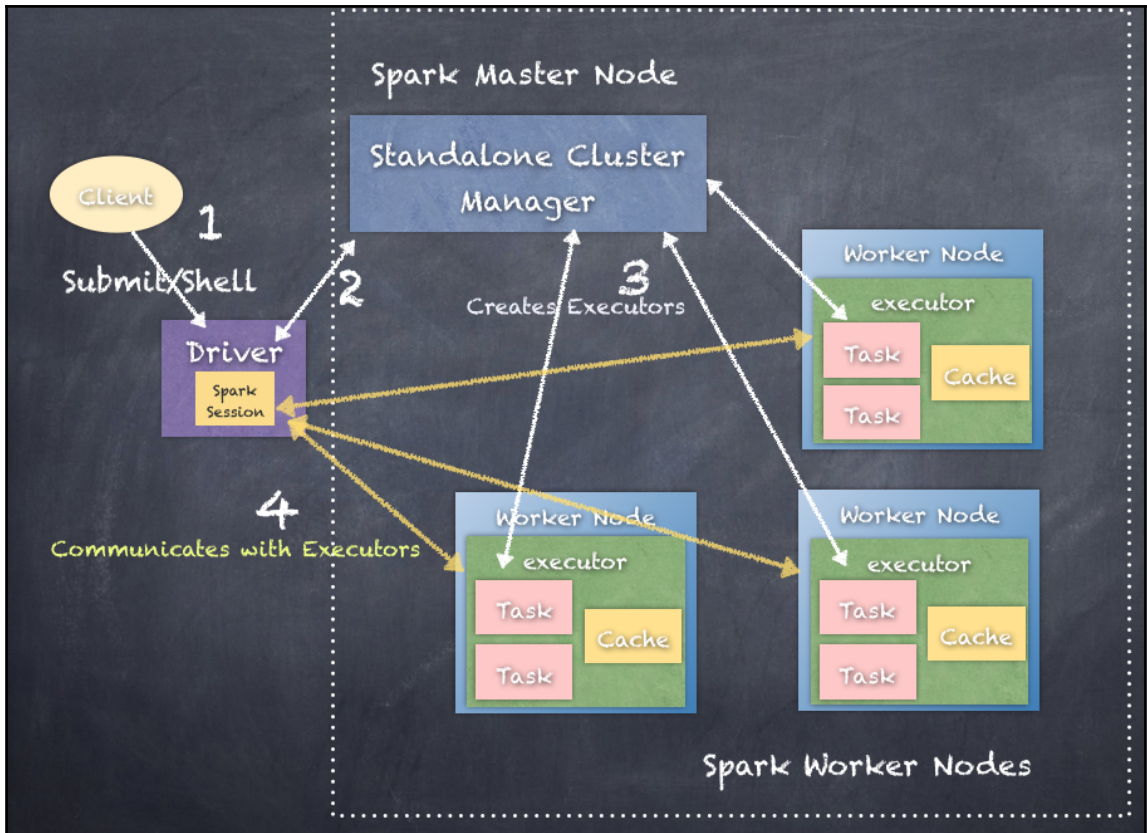














← → ↻ 🏠 spark.apache.org/downloads.html

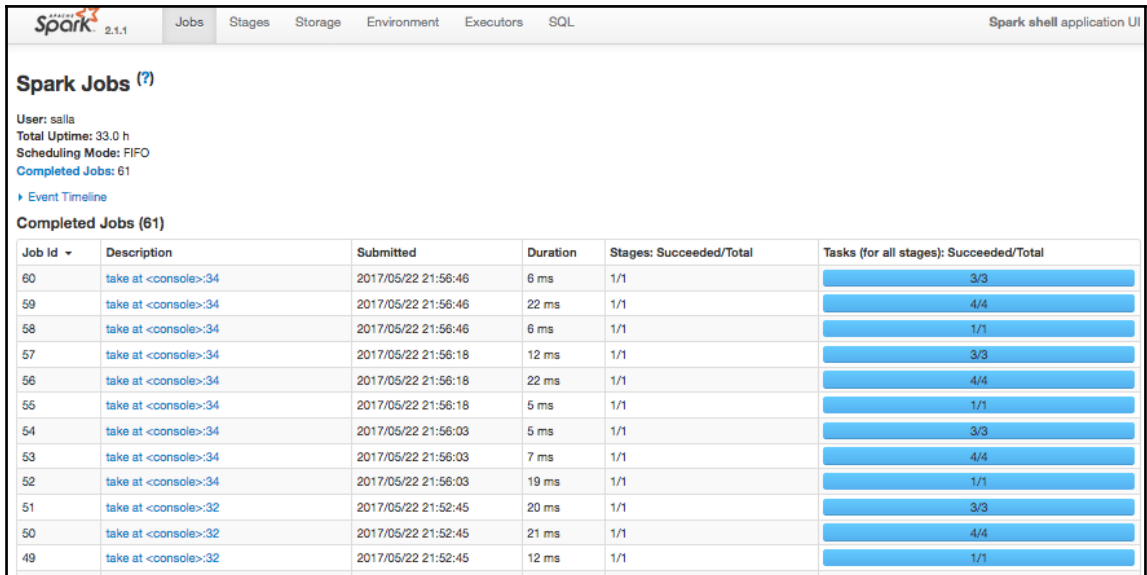
**APACHE Spark™** Lightning-fast cluster computing

Download Libraries Documentation Examples Community Developers

## Download Apache Spark™

1. Choose a Spark release: 2.1.1 (May 02 2017)
2. Choose a package type: Pre-built for Apache Hadoop 2.7 and later
3. Choose a download type: Direct Download
4. Download Spark: [spark-2.1.1-bin-hadoop2.7.tgz](#)
5. Verify this release using the [2.1.1 signatures and checksums](#) and [project release KEYS](#).

*Note: Starting version 2.0, Spark is built with Scala 2.11 by default. Scala 2.10 users should download the Spark source package and build with Scala 2.10 support.*



spark 2.1.1 Jobs Stages Storage Environment Executors SQL Spark shell application UI

## Spark Jobs (?)

User: salla  
 Total Uptime: 33.0 h  
 Scheduling Mode: FIFO  
 Completed Jobs: 61  
[Event Timeline](#)

### Completed Jobs (61)

Job Id	Description	Submitted	Duration	Stages: Succeeded/Total	Tasks (for all stages): Succeeded/Total
60	take at <console>:34	2017/05/22 21:56:46	6 ms	1/1	3/3
59	take at <console>:34	2017/05/22 21:56:46	22 ms	1/1	4/4
58	take at <console>:34	2017/05/22 21:56:46	6 ms	1/1	1/1
57	take at <console>:34	2017/05/22 21:56:18	12 ms	1/1	3/3
56	take at <console>:34	2017/05/22 21:56:18	22 ms	1/1	4/4
55	take at <console>:34	2017/05/22 21:56:18	5 ms	1/1	1/1
54	take at <console>:34	2017/05/22 21:56:03	5 ms	1/1	3/3
53	take at <console>:34	2017/05/22 21:56:03	7 ms	1/1	4/4
52	take at <console>:34	2017/05/22 21:56:03	19 ms	1/1	1/1
51	take at <console>:32	2017/05/22 21:52:45	20 ms	1/1	3/3
50	take at <console>:32	2017/05/22 21:52:45	21 ms	1/1	4/4
49	take at <console>:32	2017/05/22 21:52:45	12 ms	1/1	1/1

2.1.1

Jobs
Stages
Storage
Environment
Executors
SQL

Spark shell application UI

### Executors

**Summary**

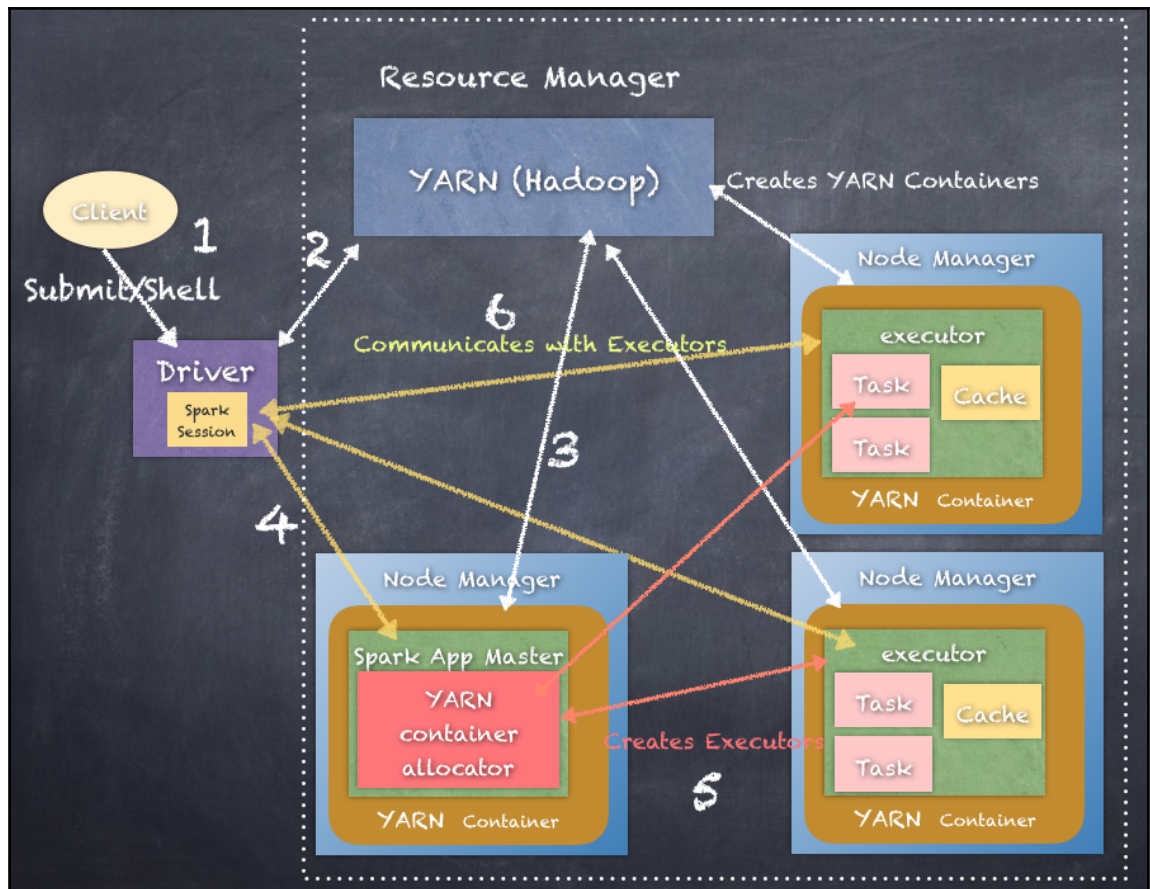
	RDD Blocks	Storage Memory	Disk Used	Cores	Active Tasks	Failed Tasks	Complete Tasks	Total Tasks	Task Time (GC Time)	Input	Shuffle Read	Shuffle Write
Active(1)	4	93.6 KB / 384.1 MB	0.0 B	8	0	0	141	141	5 s (0 ms)	46.3 KB	0.0 B	0.0 B
Dead(0)	0	0.0 B / 0.0 B	0.0 B	0	0	0	0	0	0 ms (0 ms)	0.0 B	0.0 B	0.0 B
Total(1)	4	93.6 KB / 384.1 MB	0.0 B	8	0	0	141	141	5 s (0 ms)	46.3 KB	0.0 B	0.0 B

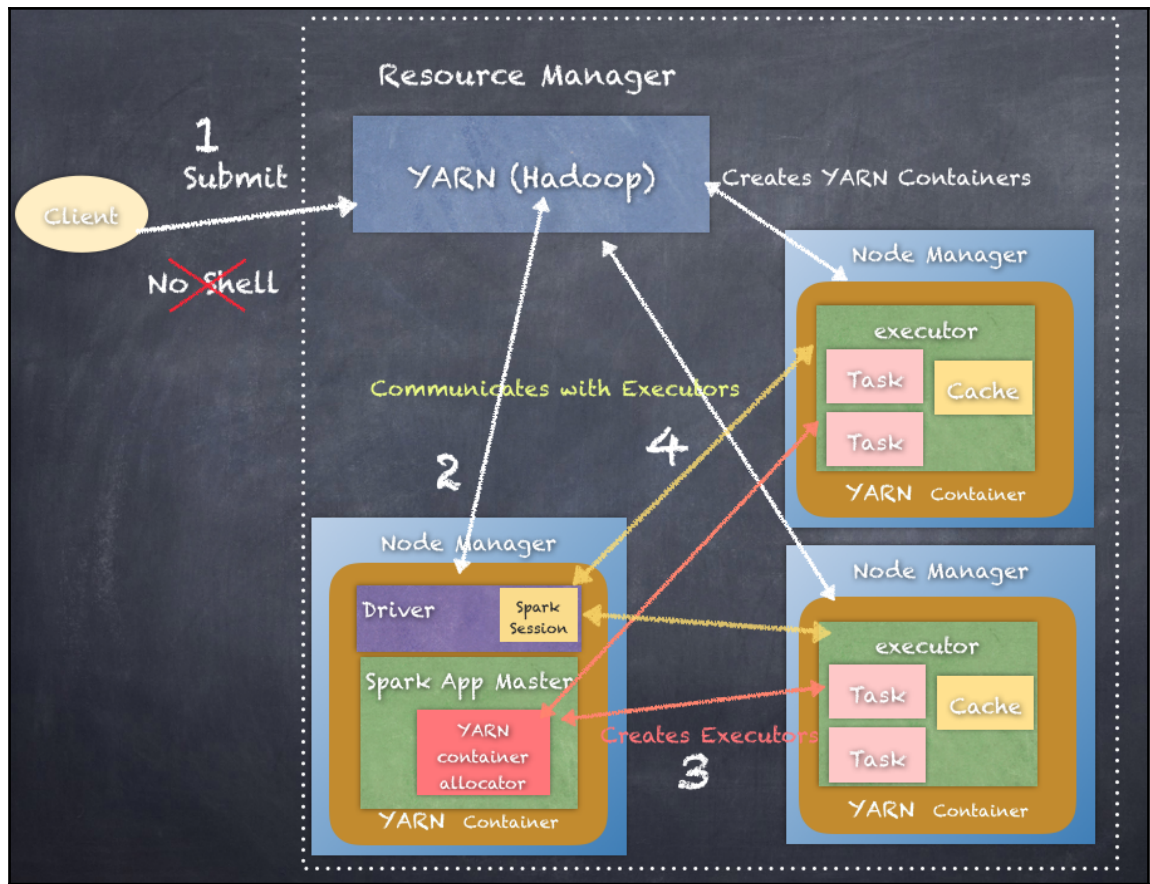
**Executors**

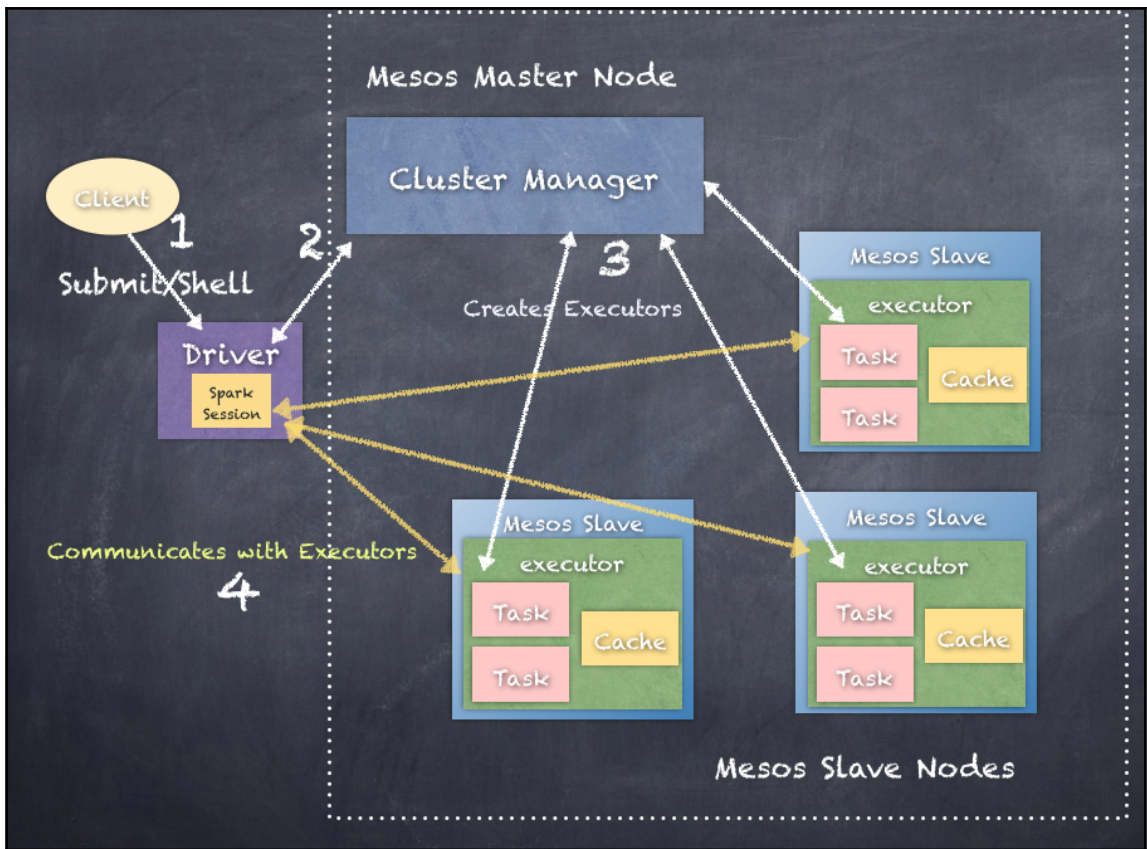
Show  entries Search:

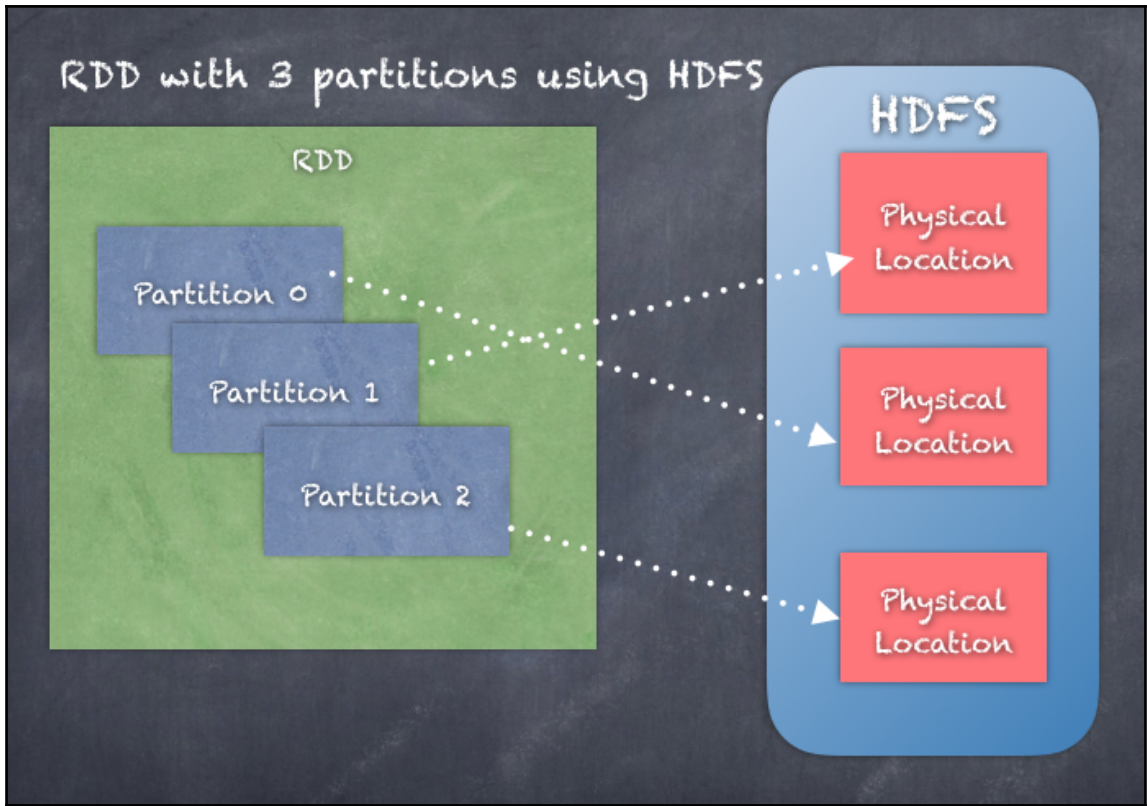
Executor ID	Address	Status	RDD Blocks	Storage Memory	Disk Used	Cores	Active Tasks	Failed Tasks	Complete Tasks	Total Tasks	Task Time (GC Time)	Input	Shuffle Read	Shuffle Write	Thread Dump
driver	192.168.1.253:60407	Active	4	93.6 KB / 384.1 MB	0.0 B	8	0	0	141	141	5 s (0 ms)	46.3 KB	0.0 B	0.0 B	<a href="#">Thread Dump</a>

Showing 1 to 1 of 1 entries [Previous](#) [Next](#)

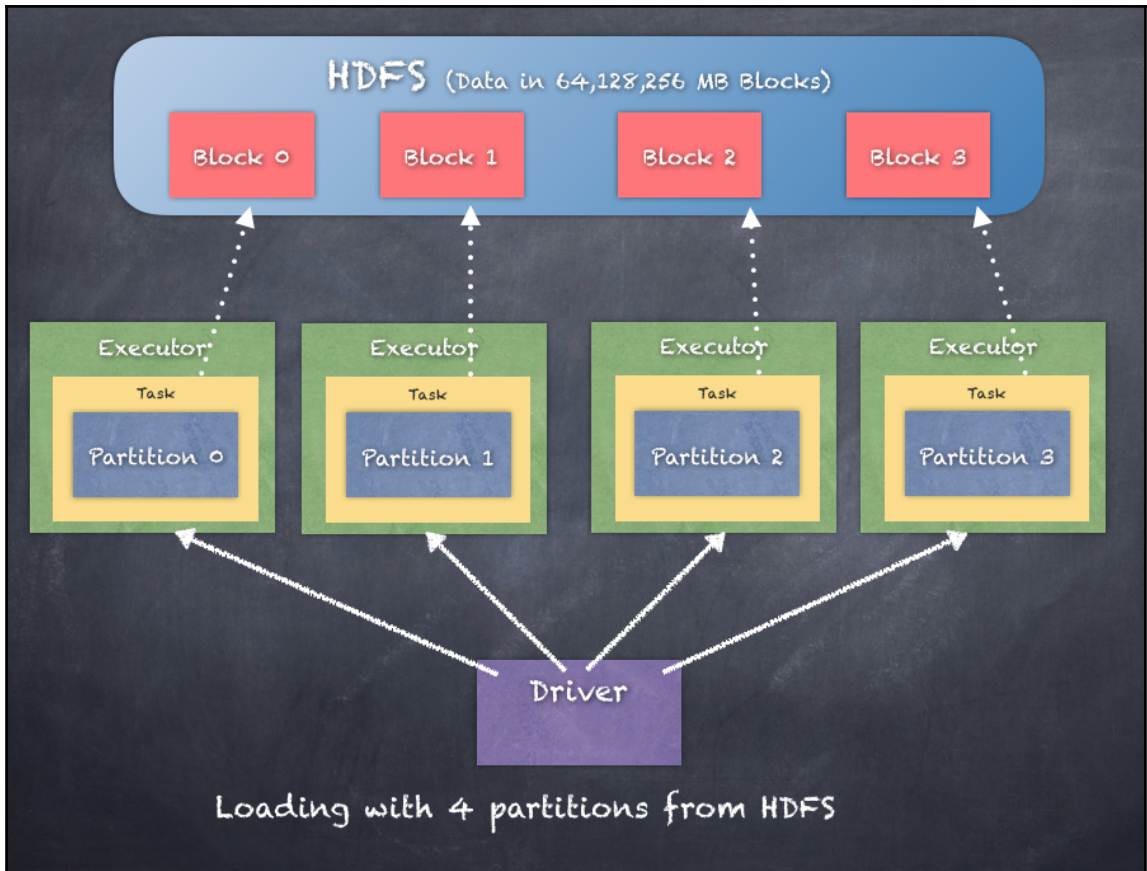


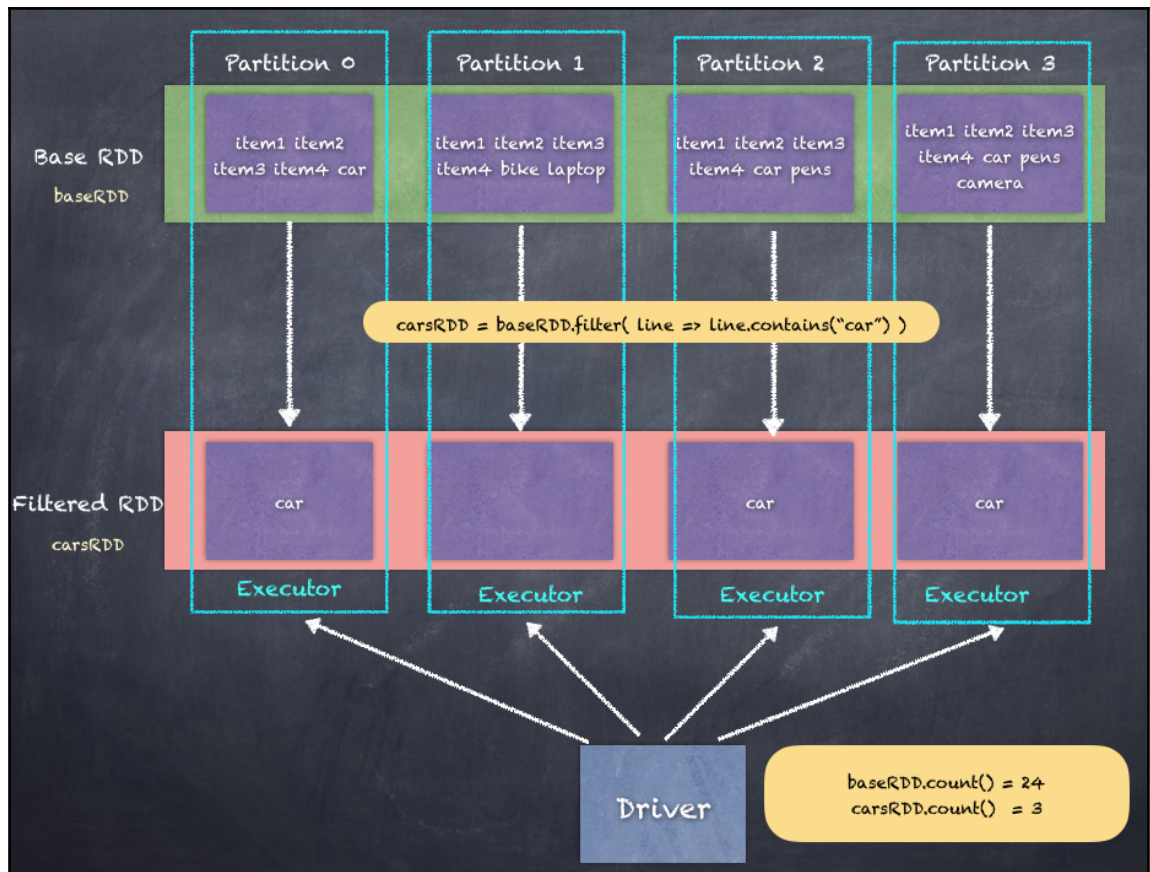


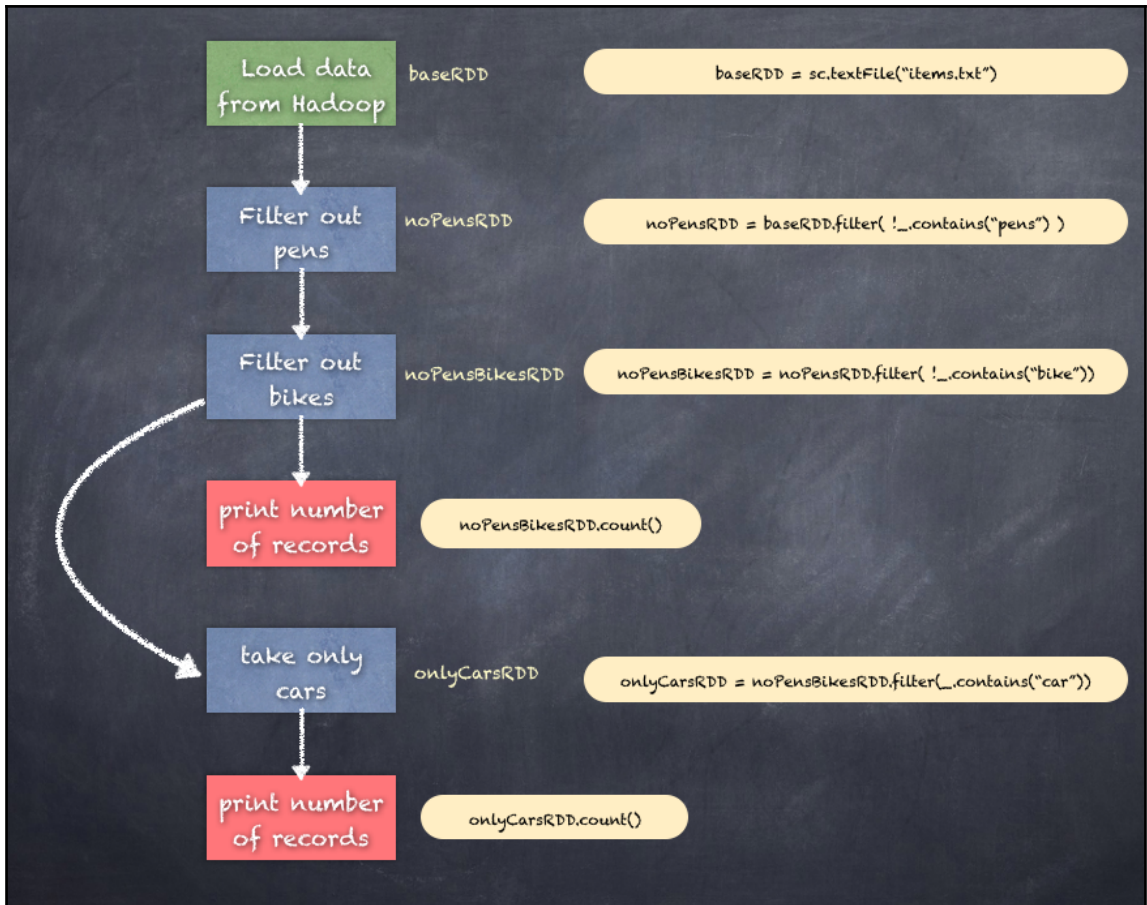


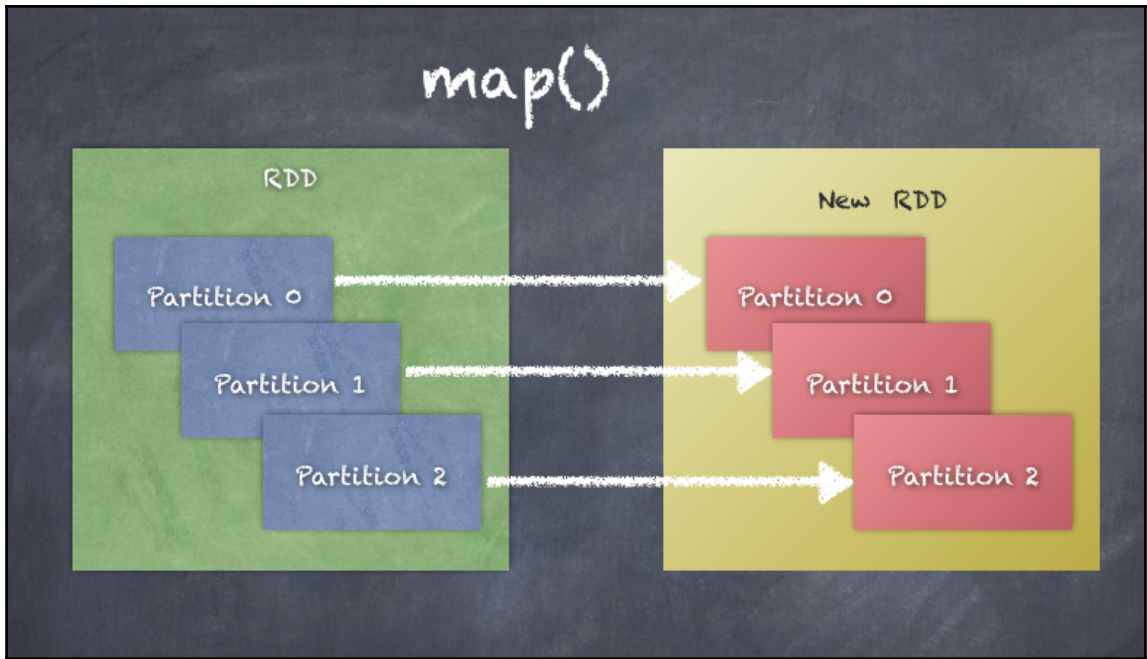


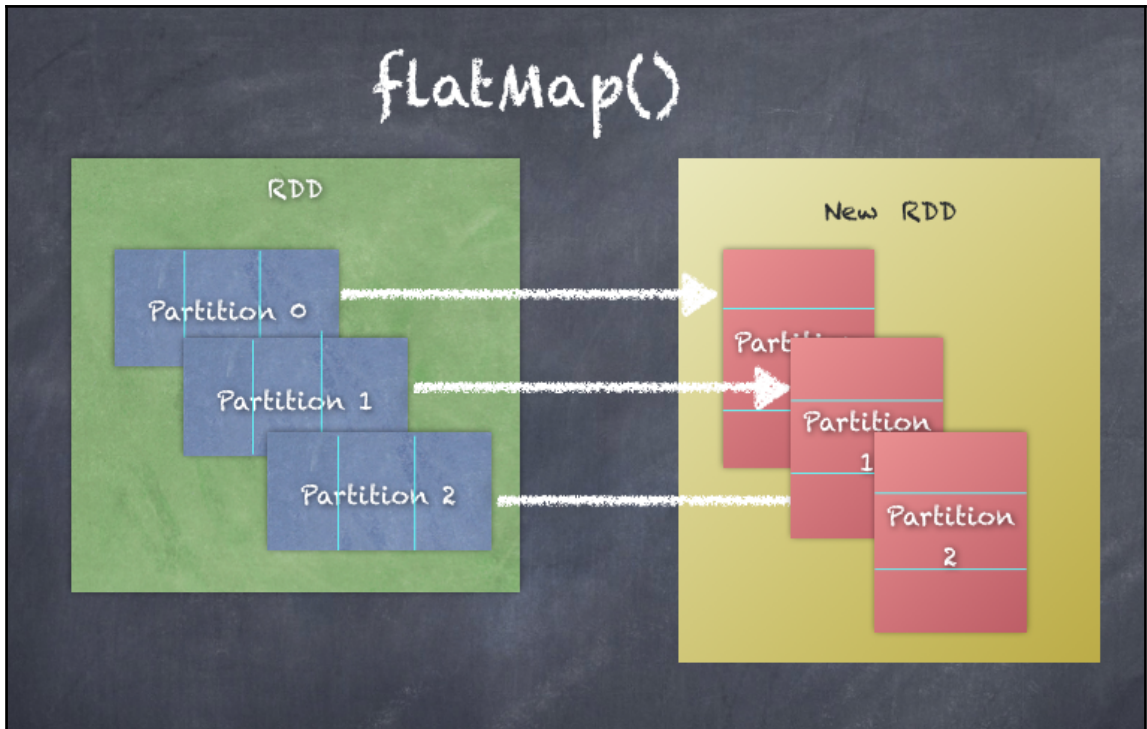


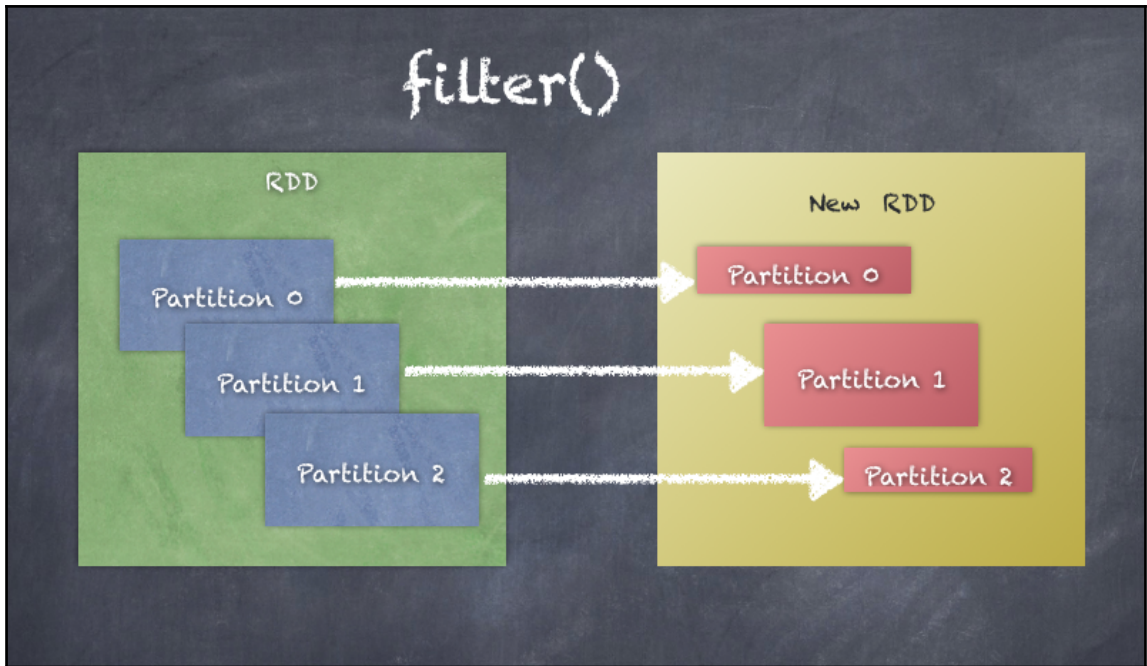


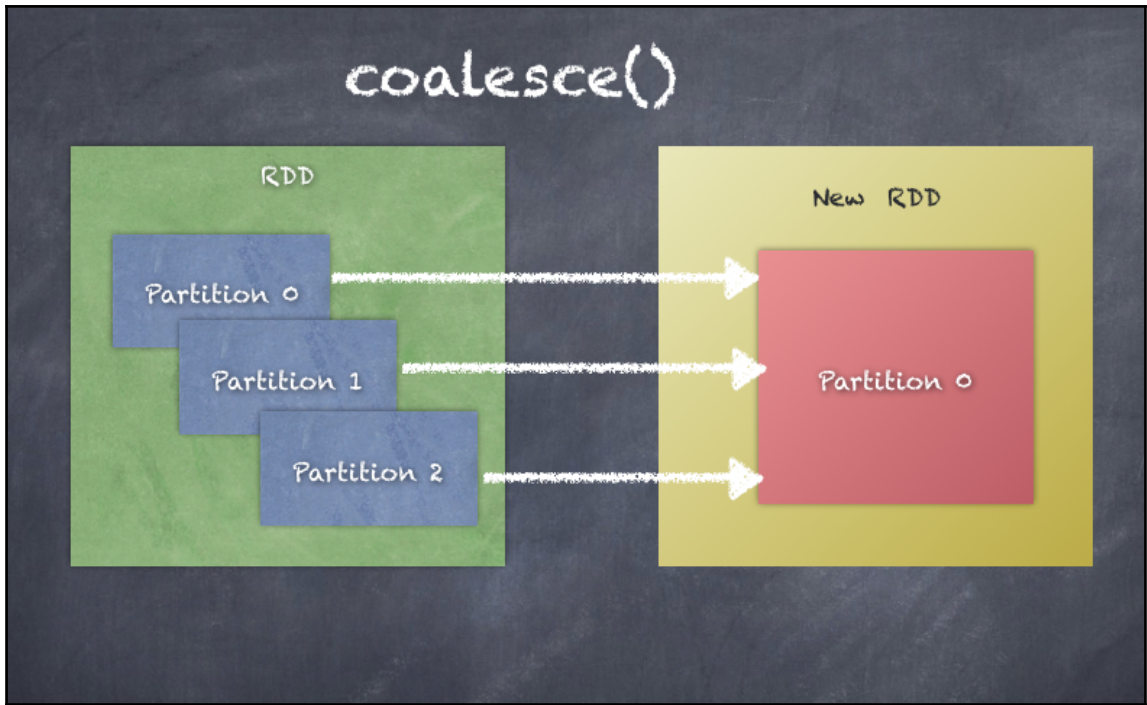


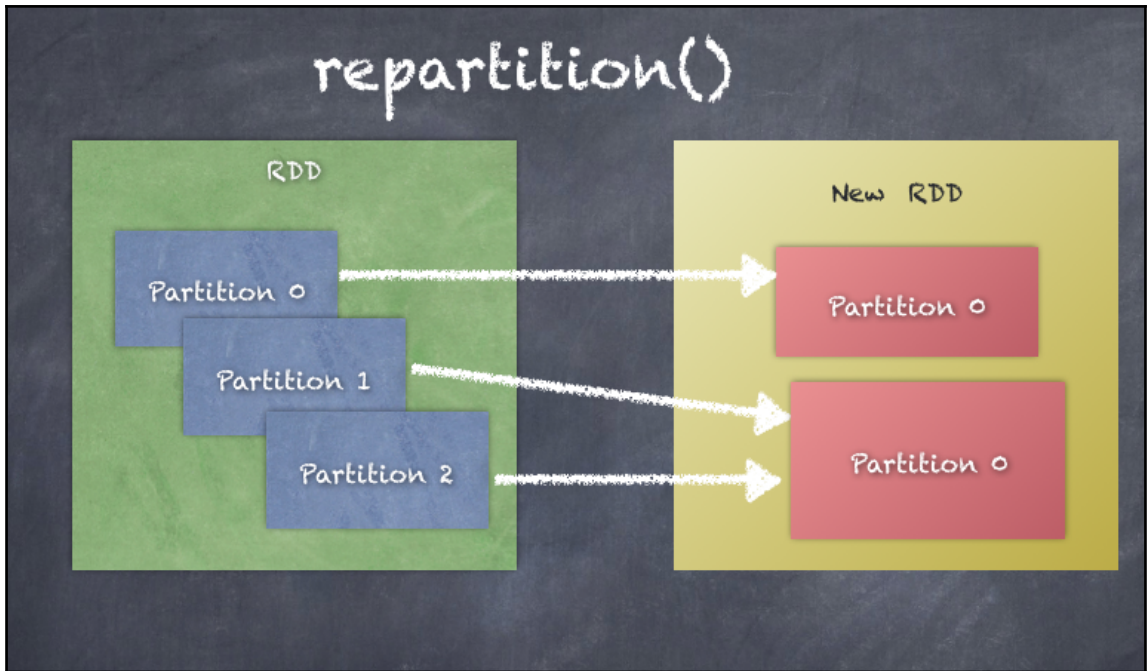




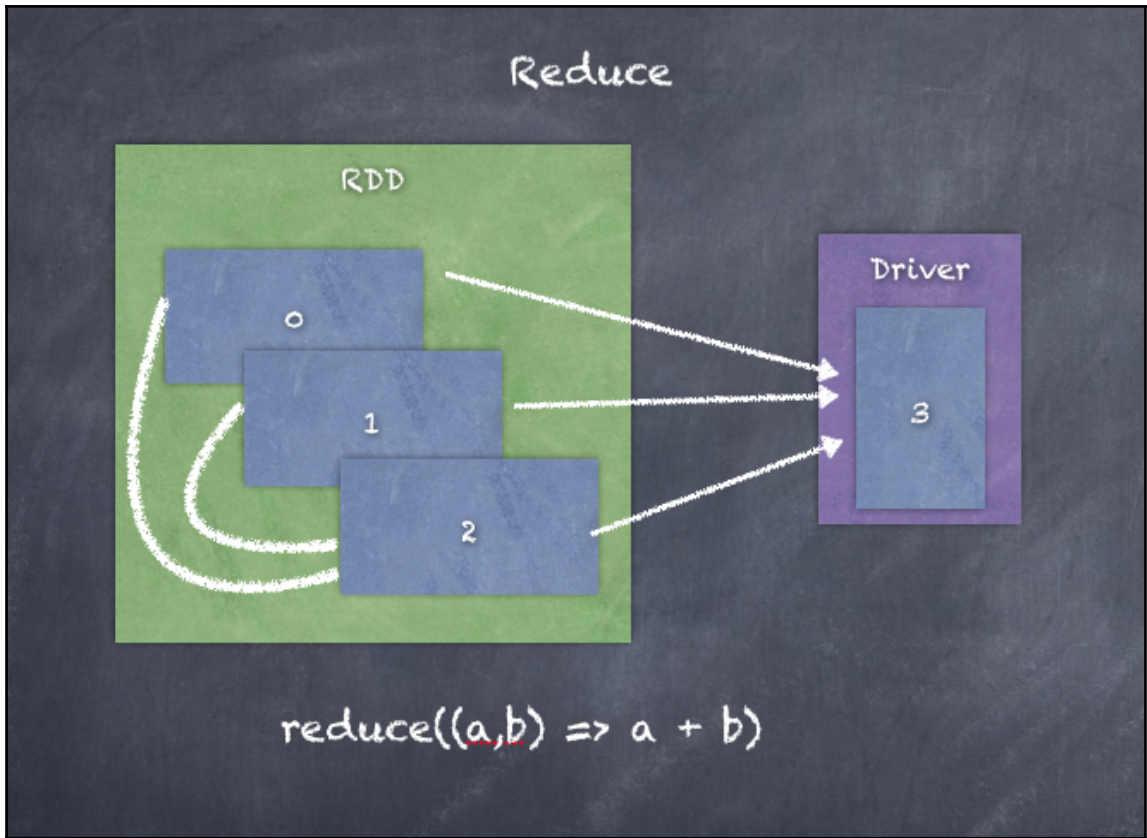


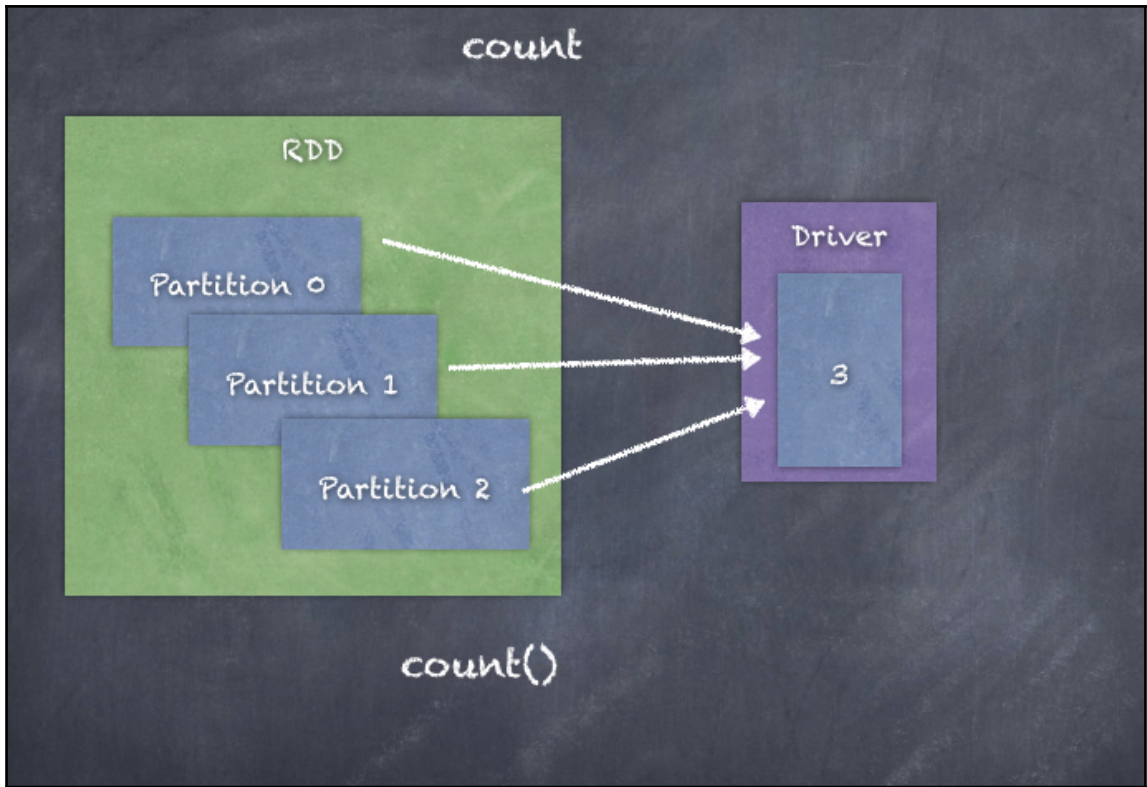


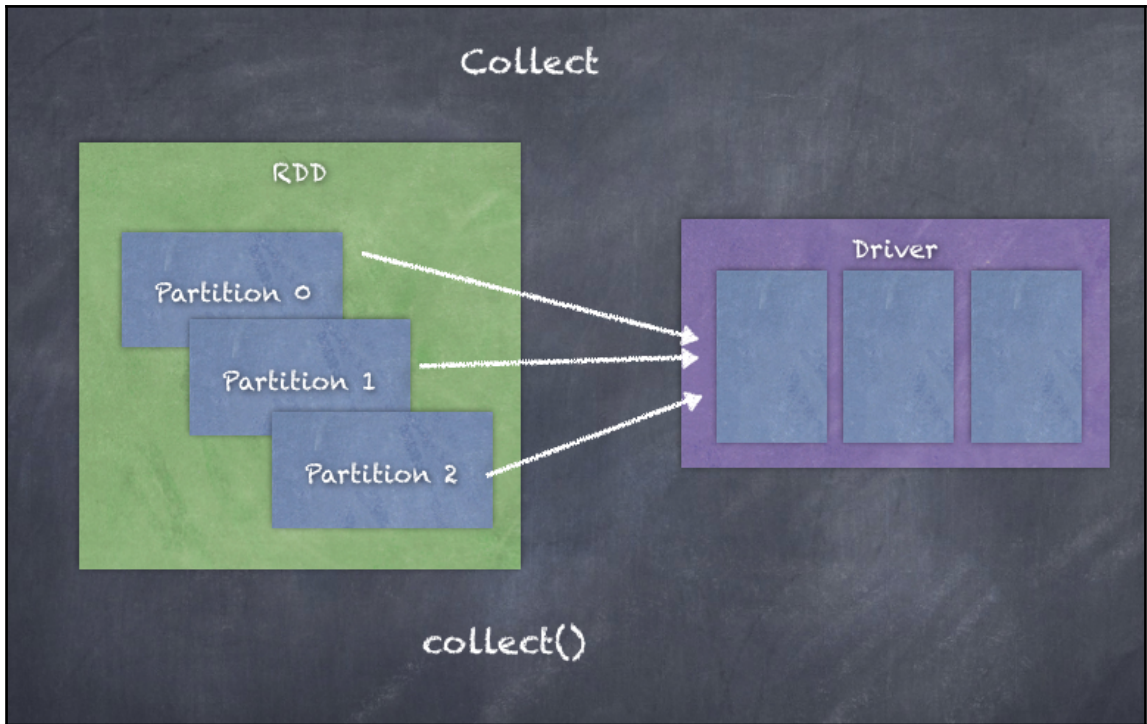












# Caching

The screenshot shows the Spark UI 'Stages' tab. A table lists two completed stages. Stage 0, 'count at <console>:27', has a duration of 0.6 s. Stage 1, 'count at <console>:27', has a duration of 59 ms. A red arrow points from the 'Before caching' text to the 0.6 s duration of Stage 0. A green arrow points from the 'After caching' text to the 59 ms duration of Stage 1.

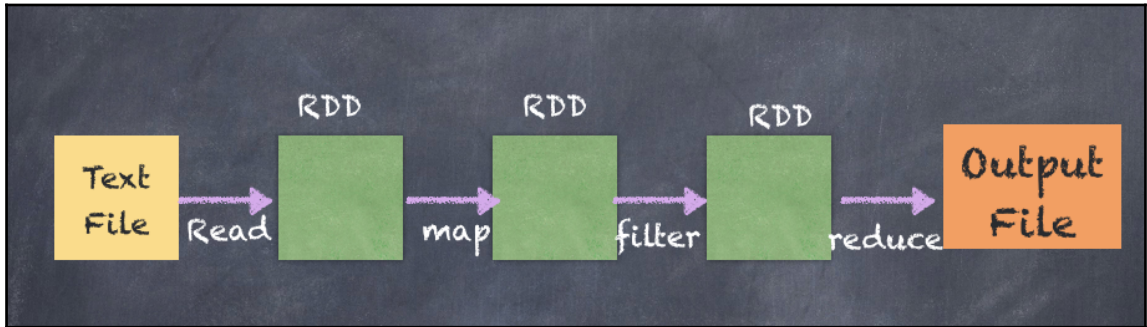
Stage Id	Description	Submitted	Duration	Tasks: Succeeded/Total	Input
1	count at <console>:27	2017/06/30 09:55:57	59 ms	8/8	
0	count at <console>:27	2017/06/30 09:55:47	0.6 s	8/8	

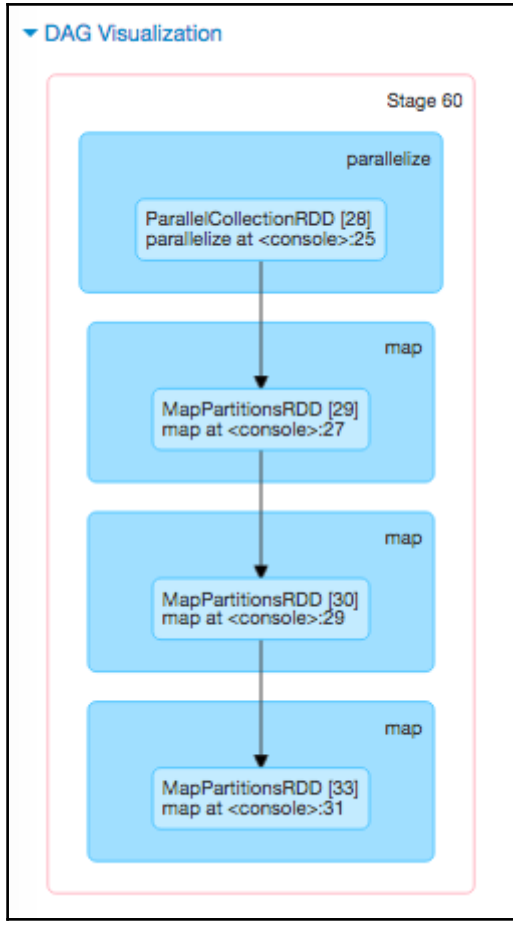
## Storage

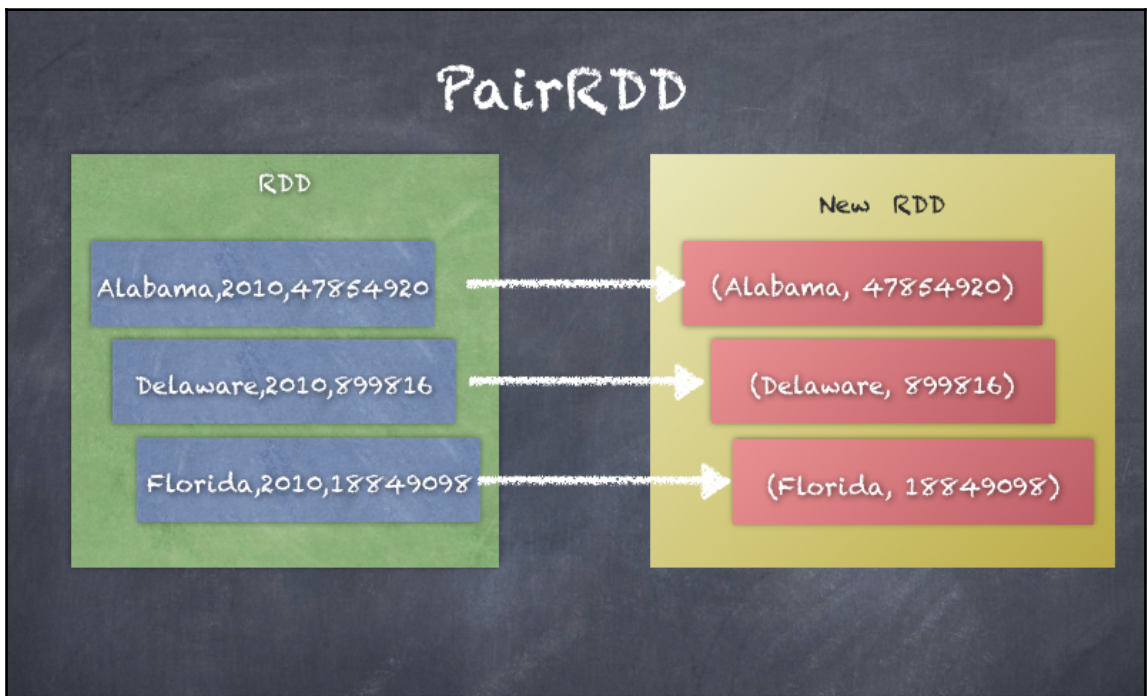
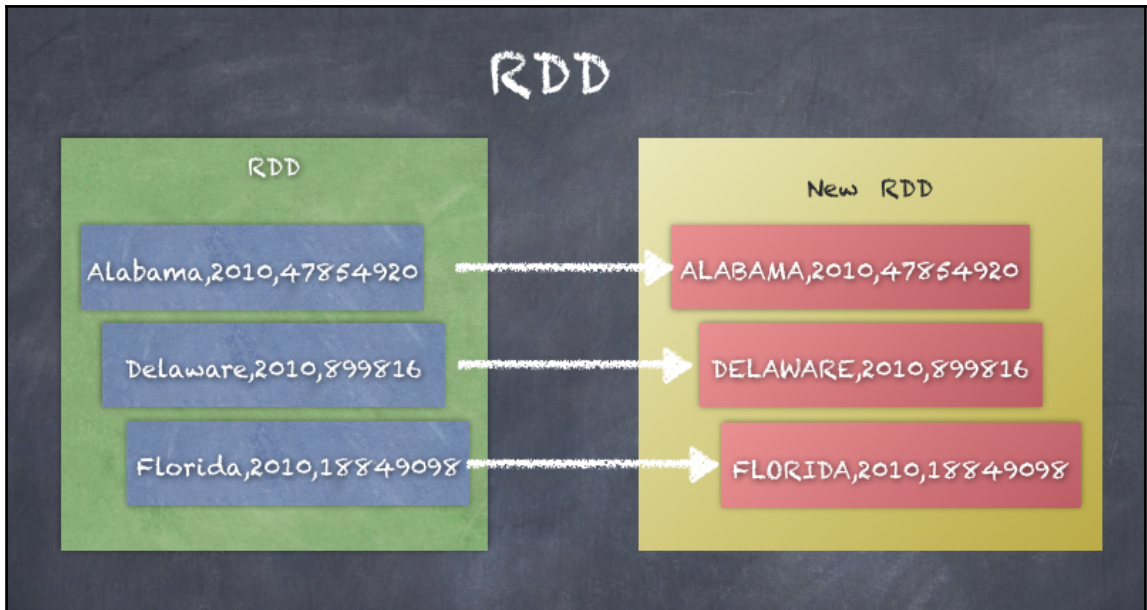
Storage Tab shows the Cache

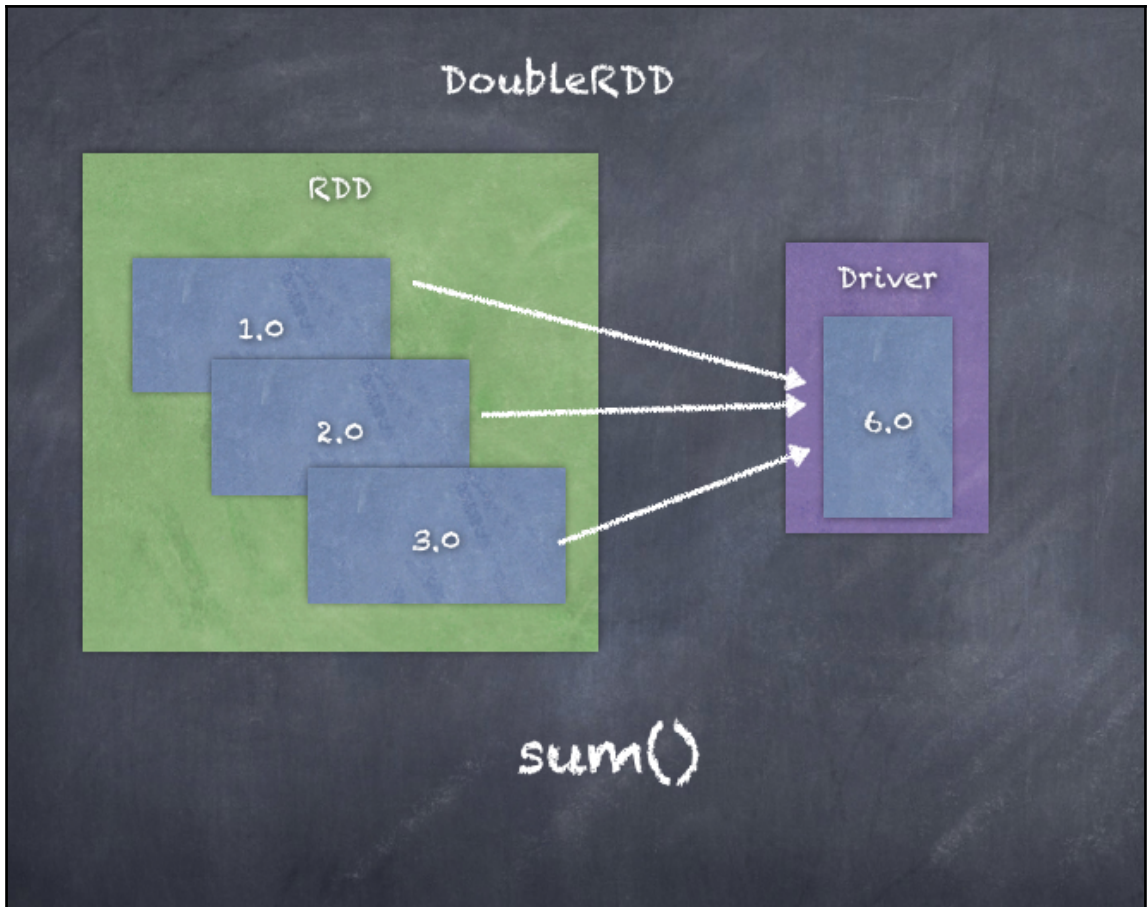
RDD Name	Storage Level	Cached Partitions	Fraction Cached	Size in Memory	Size on Disk
ParallelCollectionRDD	Memory Deserialized 1x Replicated	8	100%	176.0 B	0.0 B

## Chapter 7: Special RDD Operations

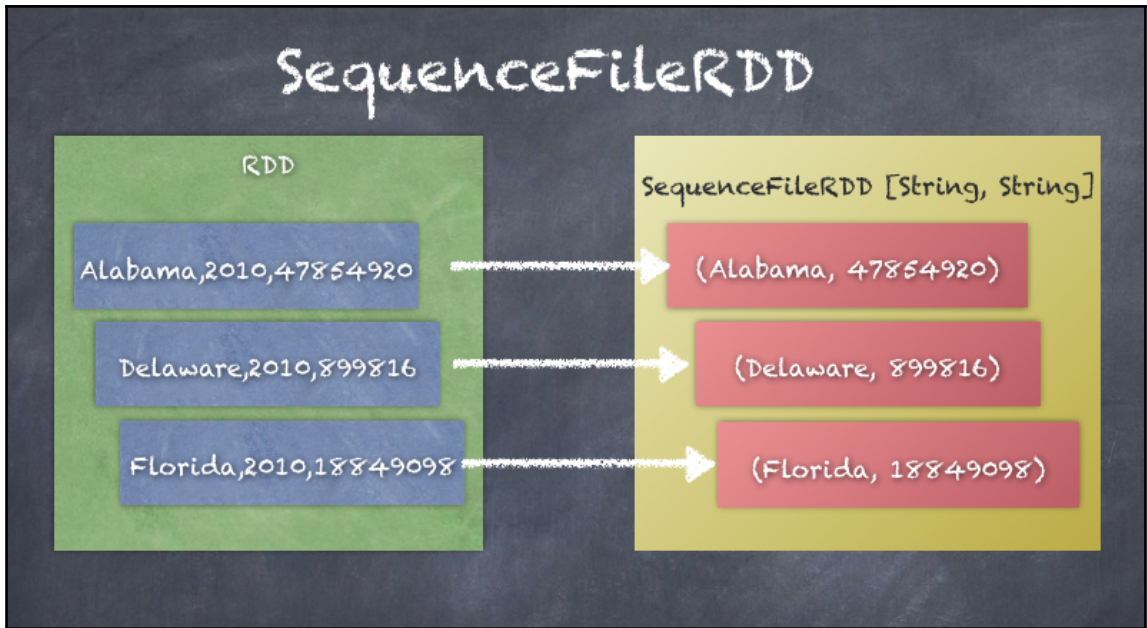


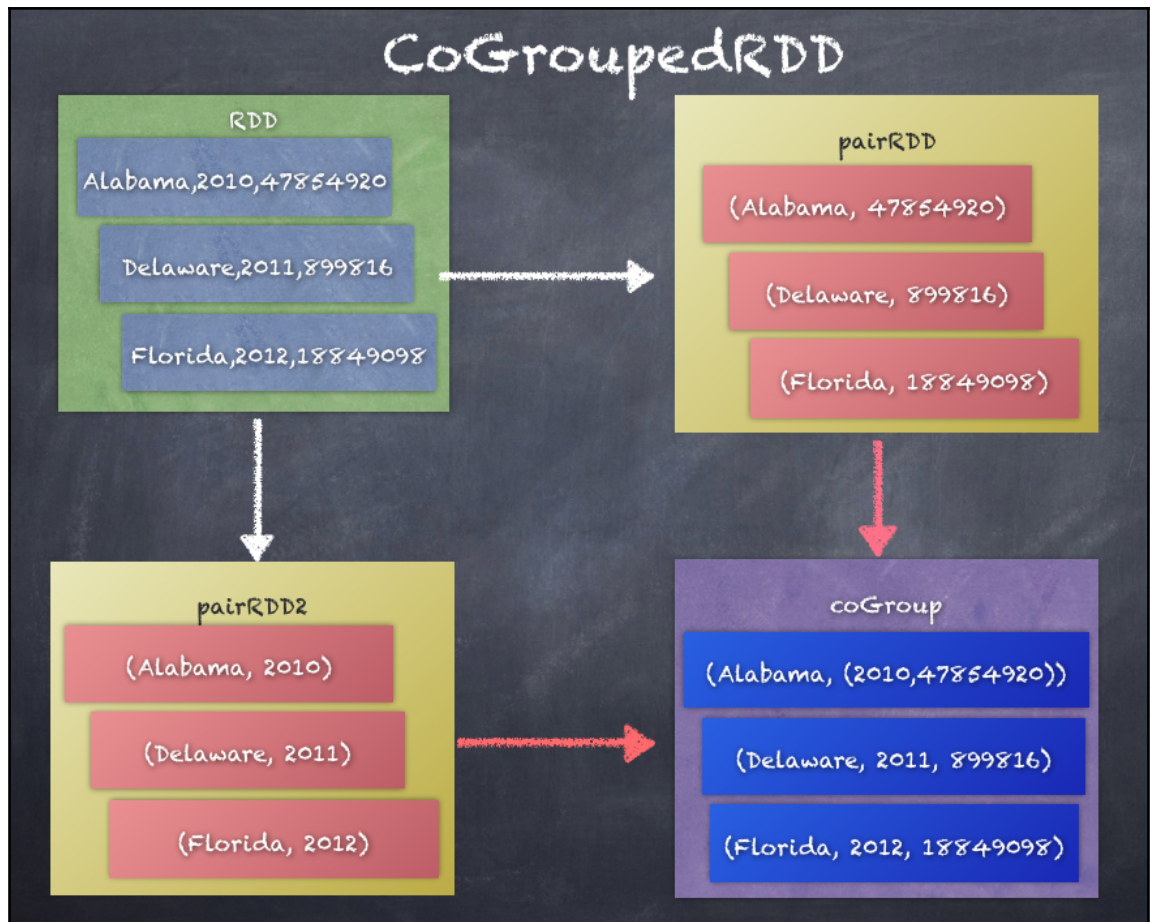


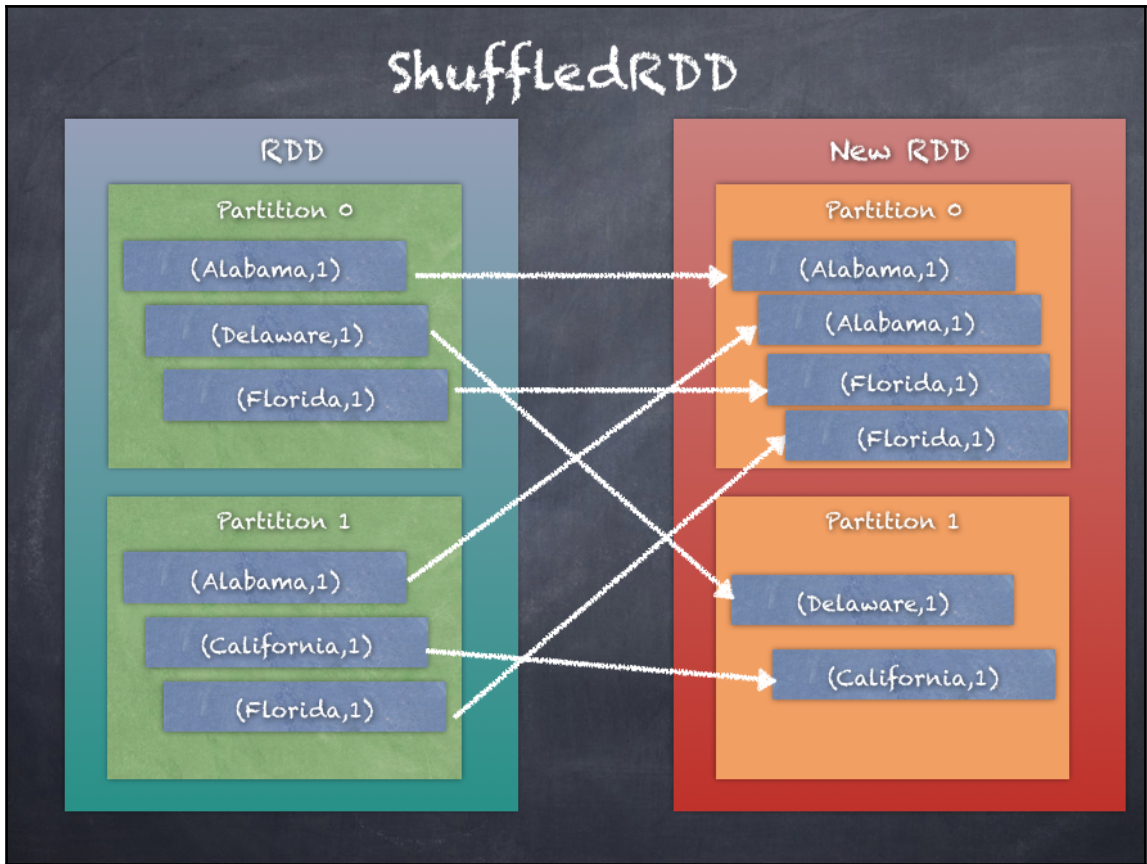


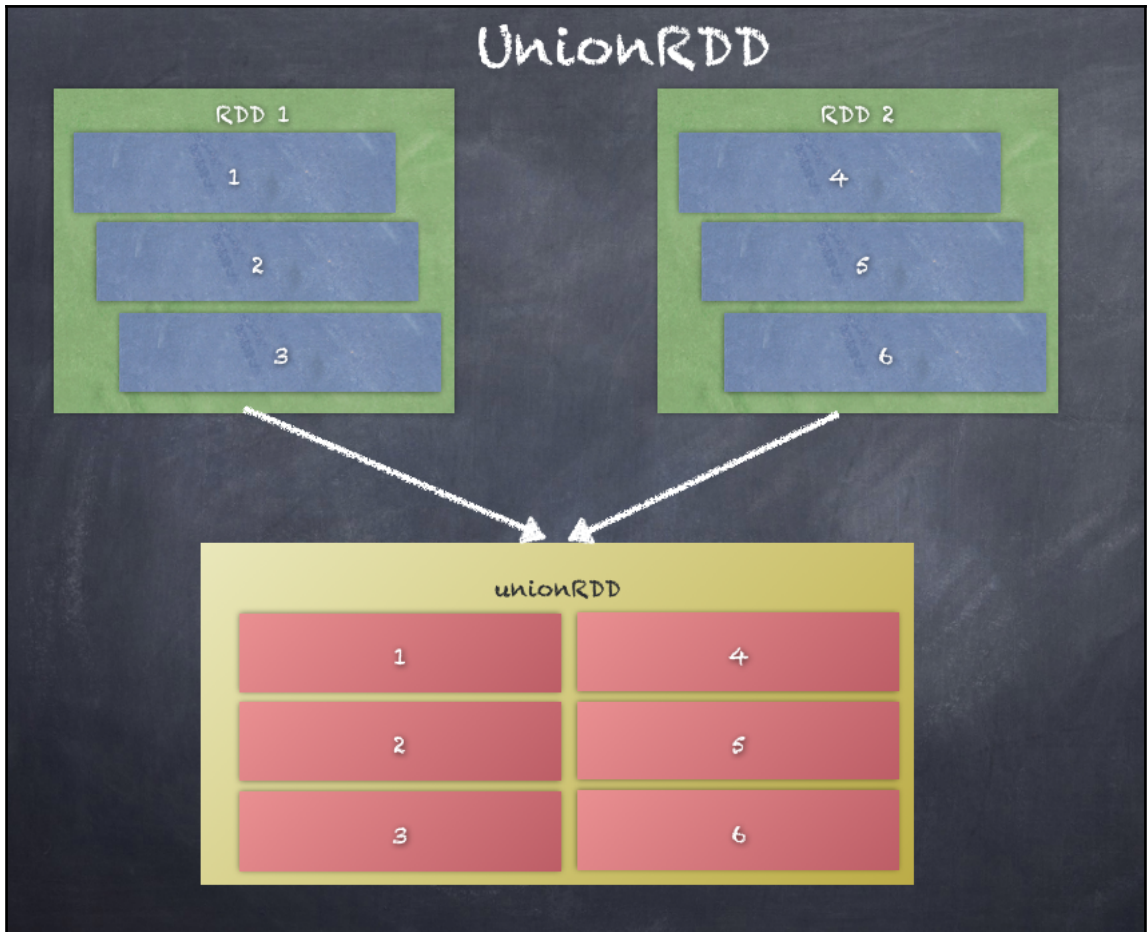


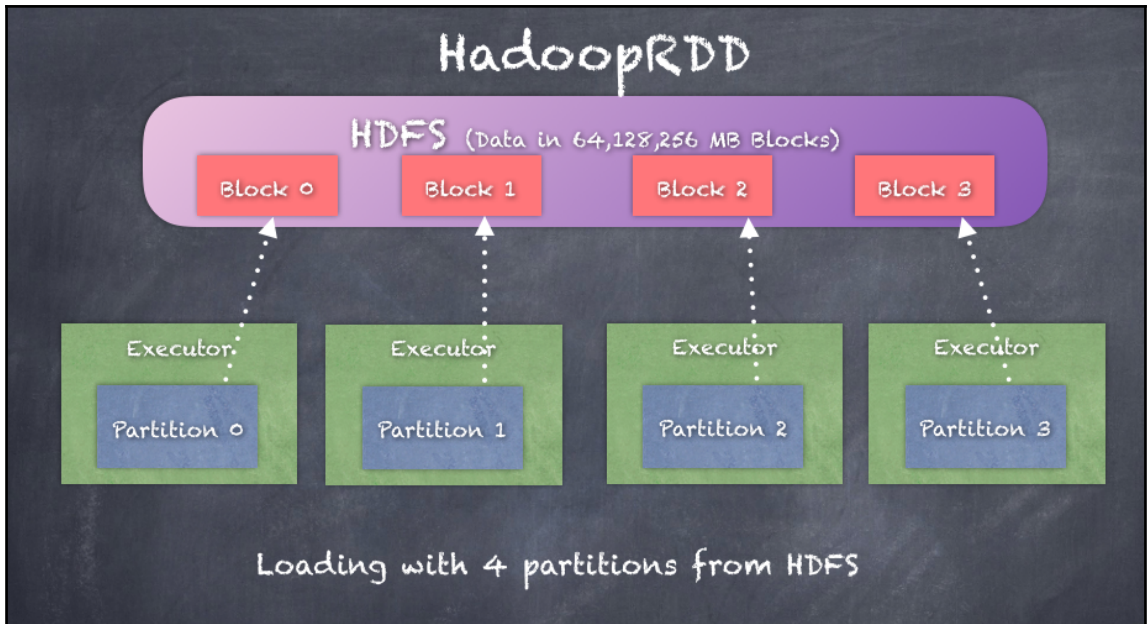


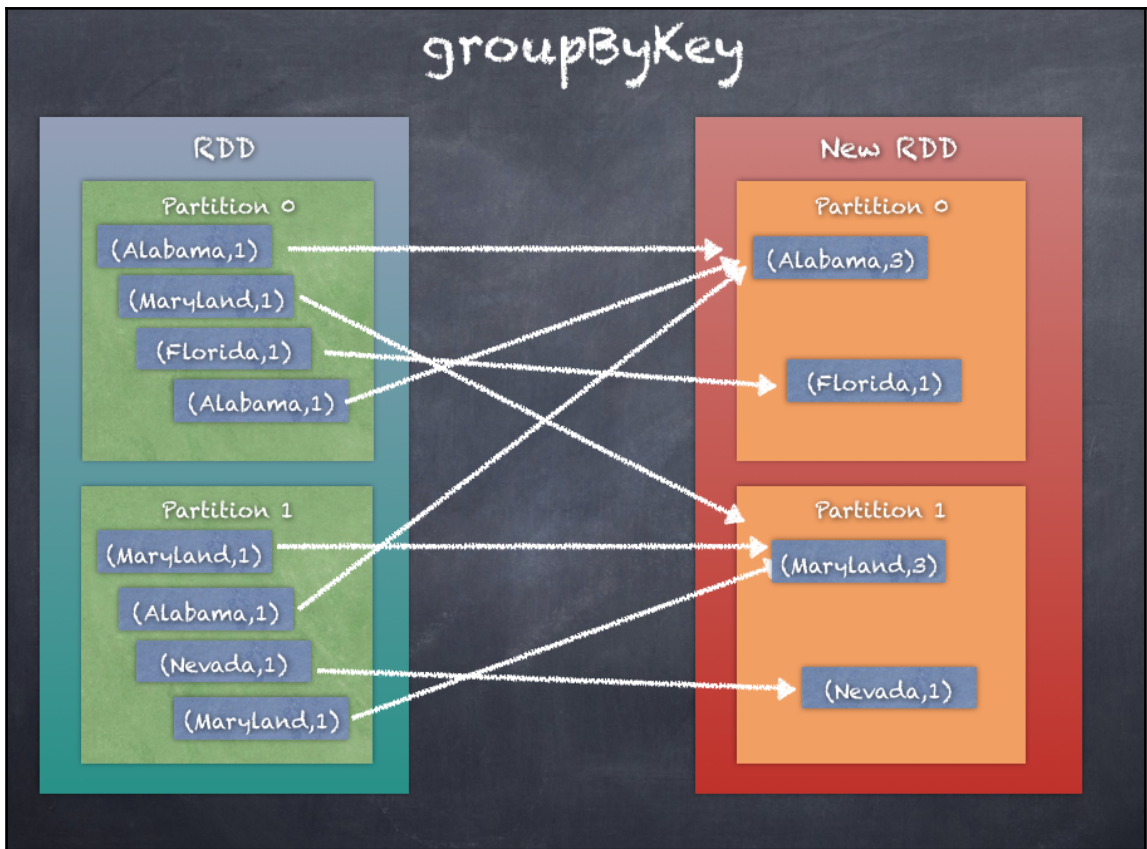


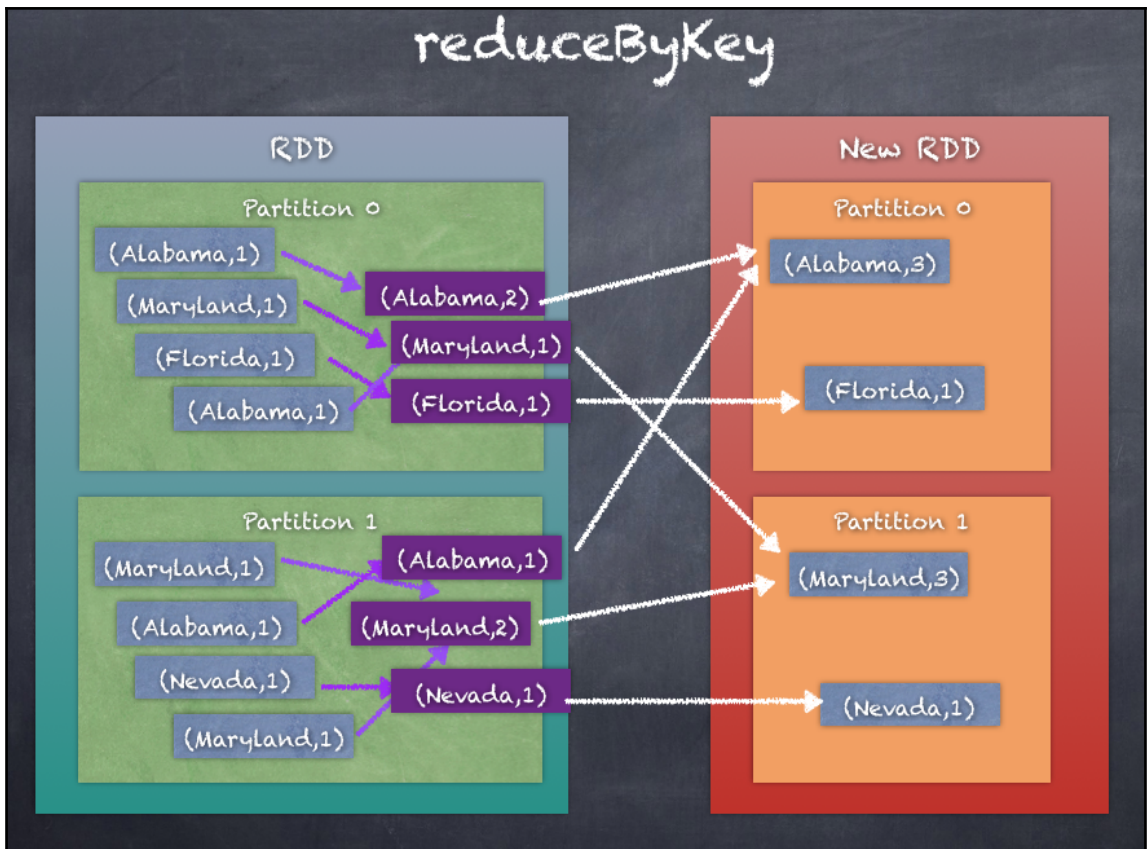


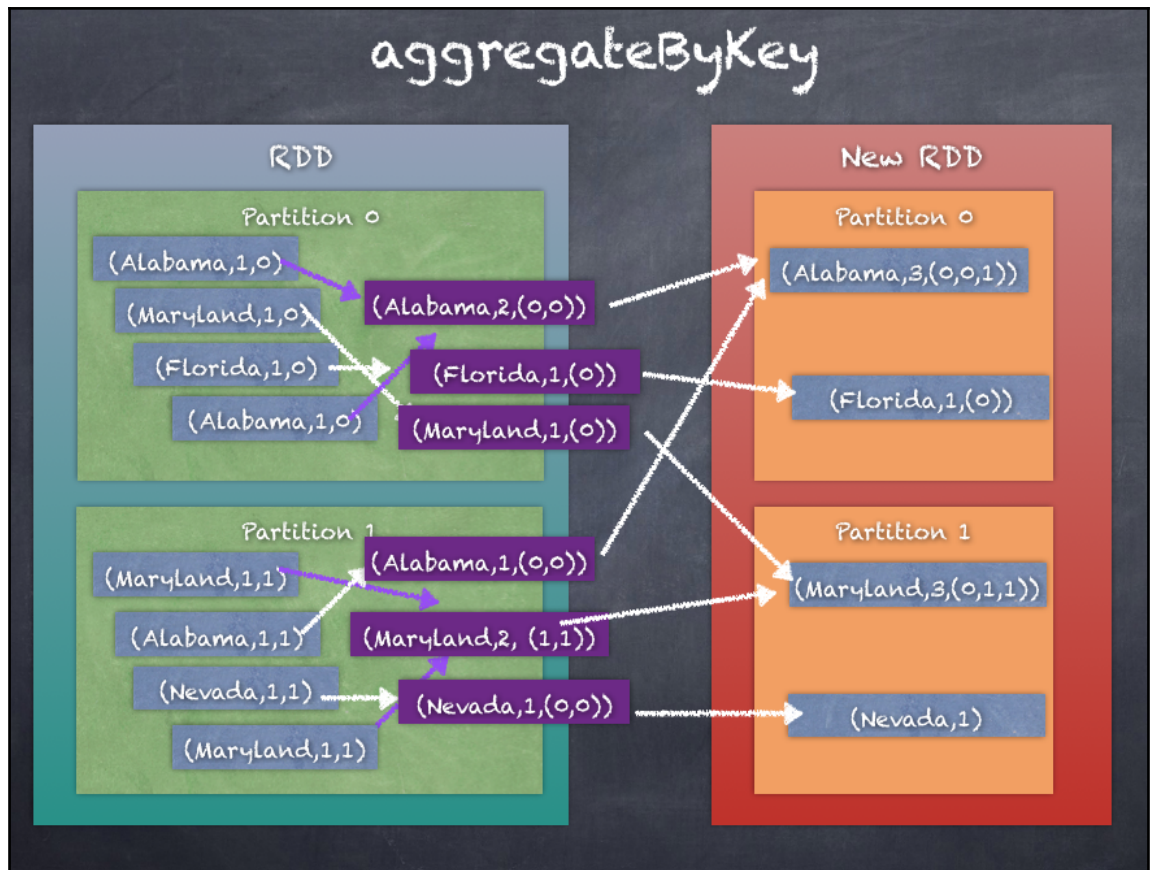




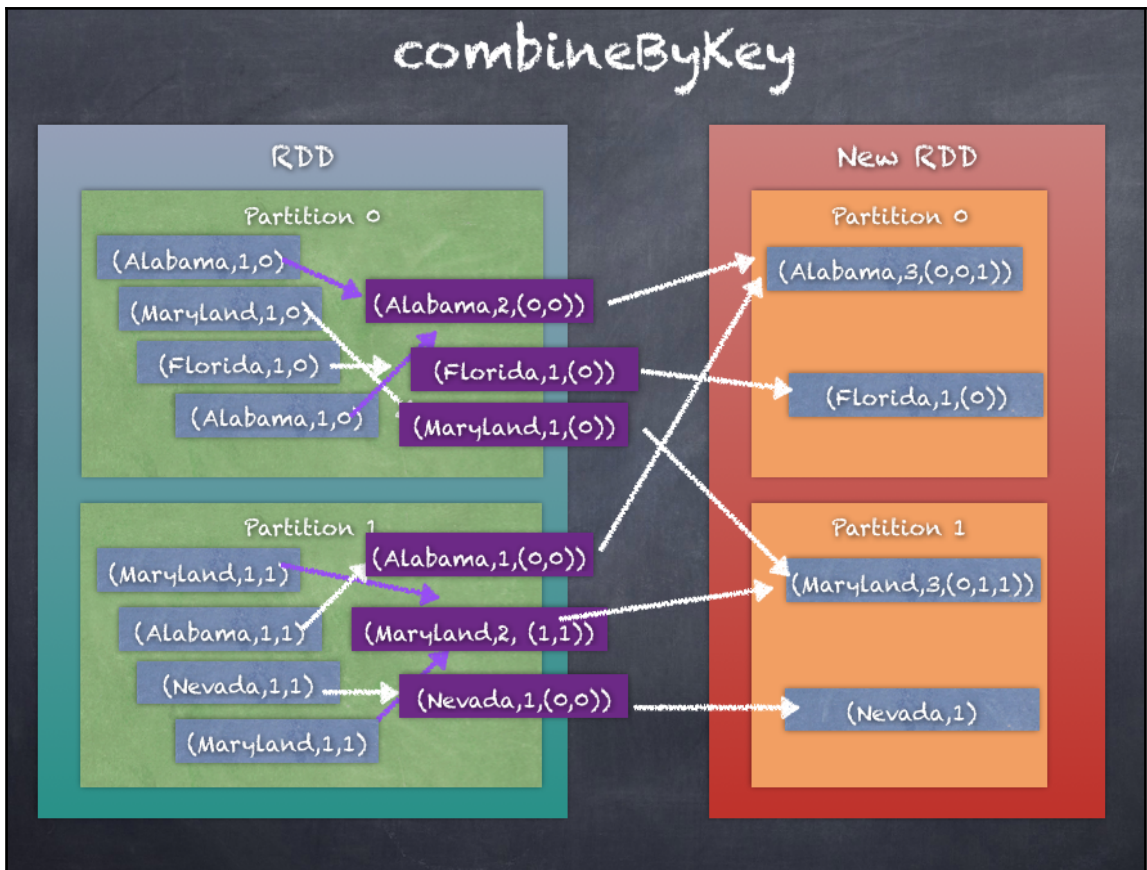


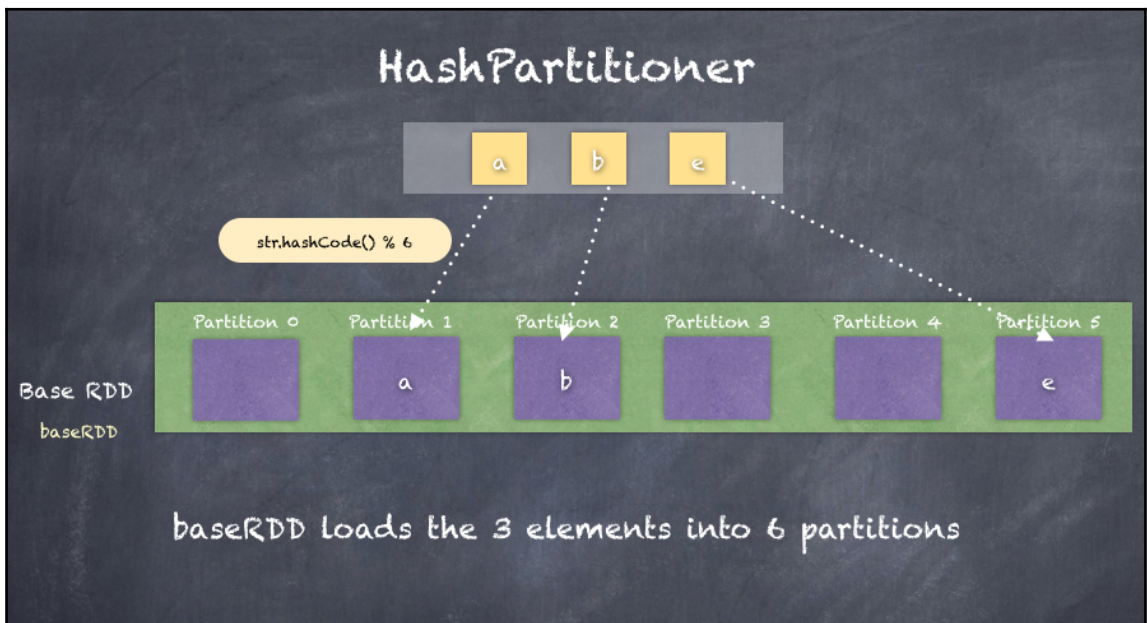
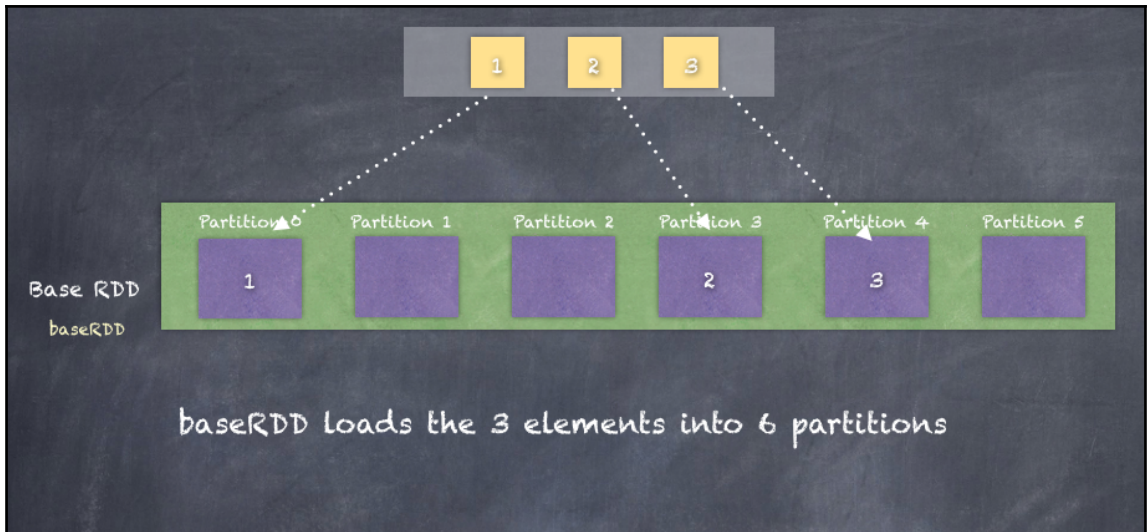


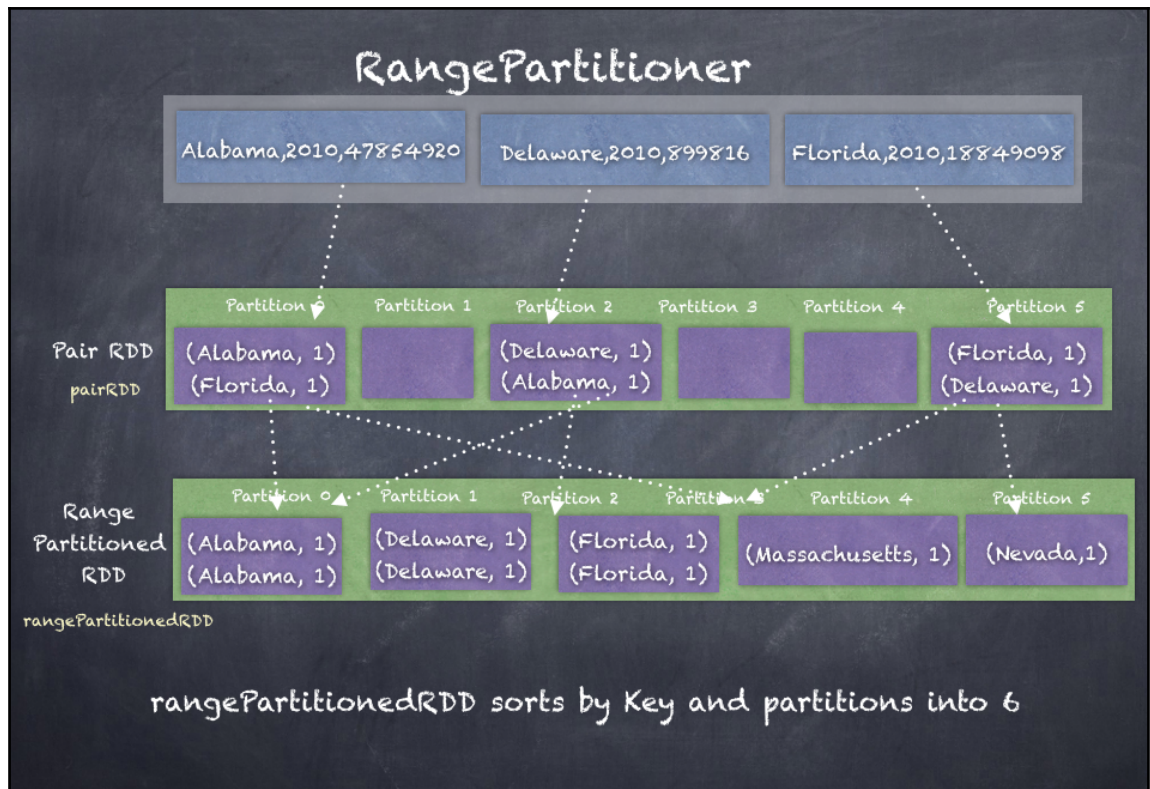


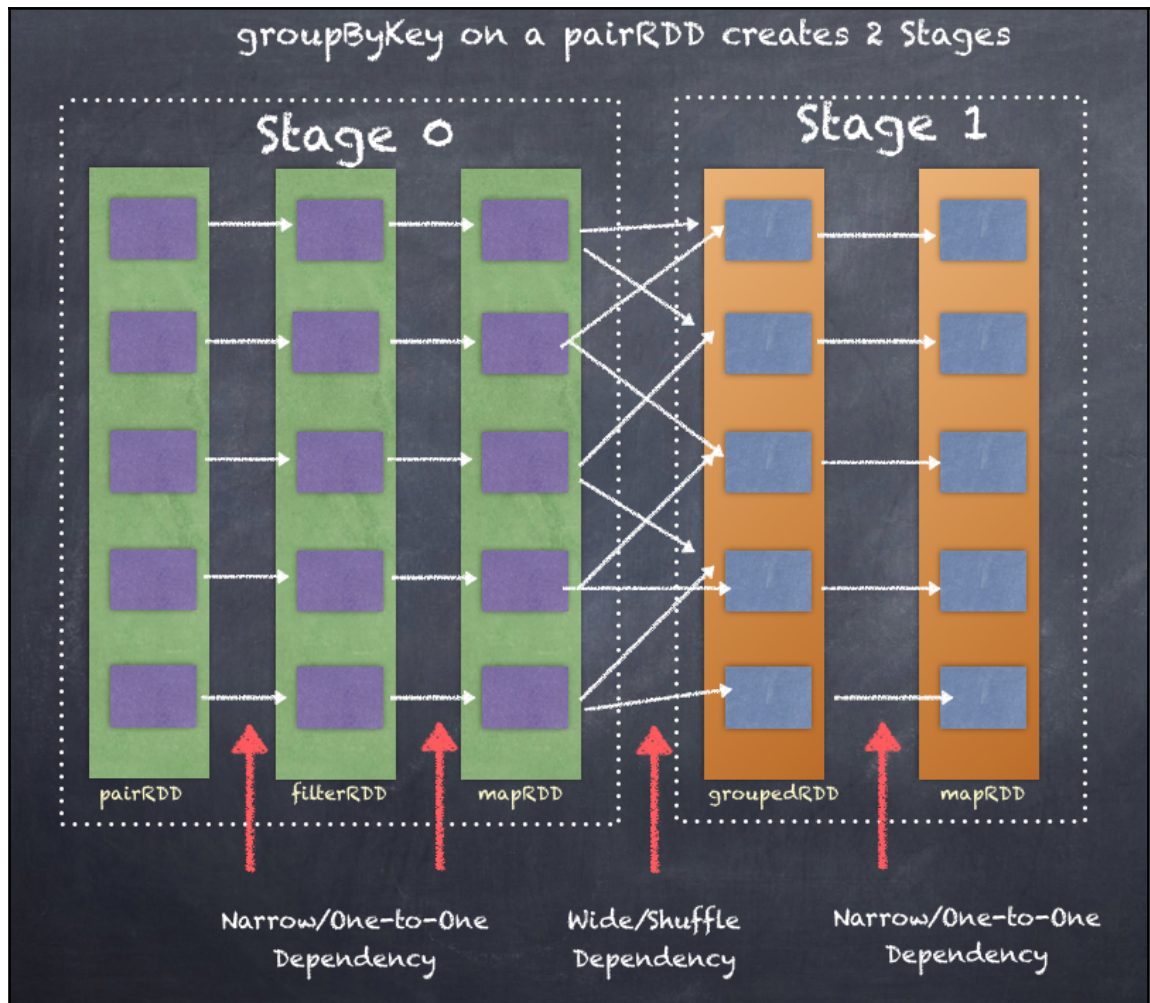


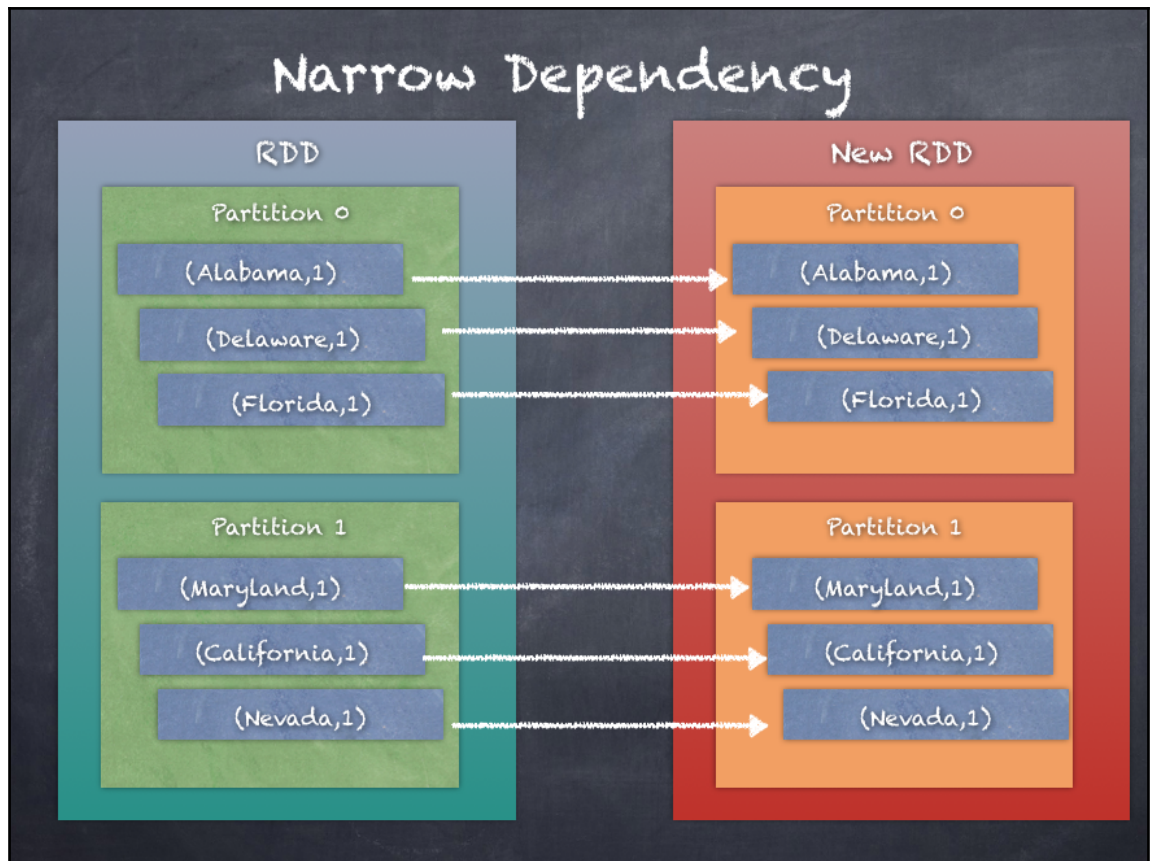


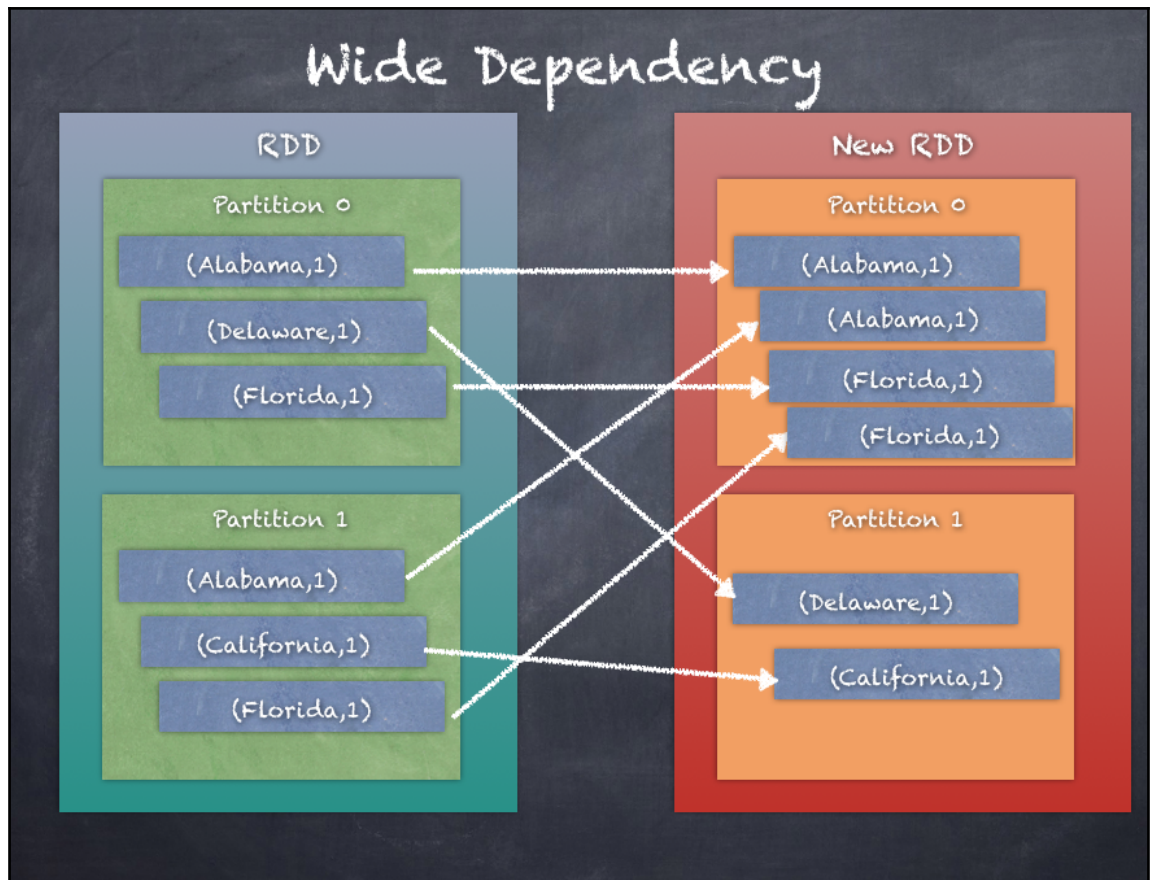


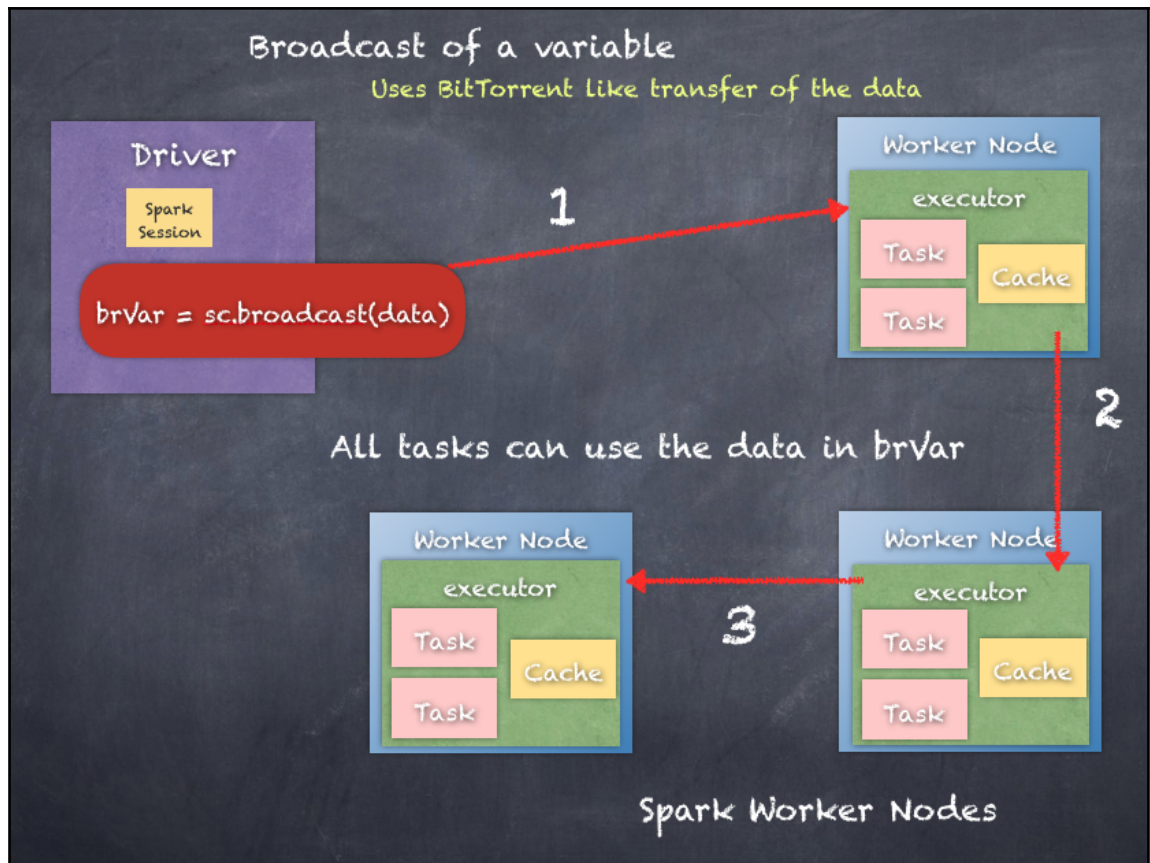








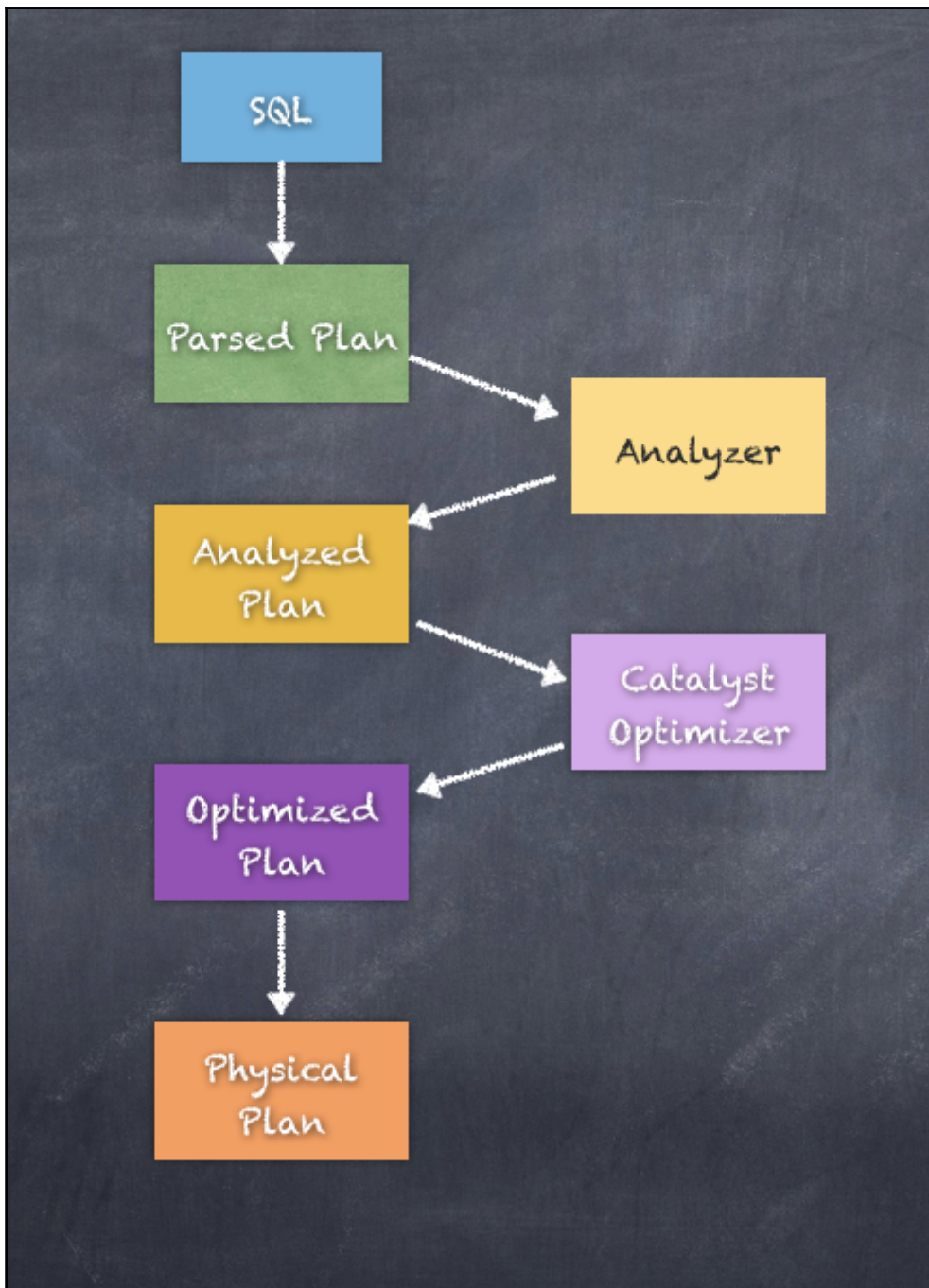




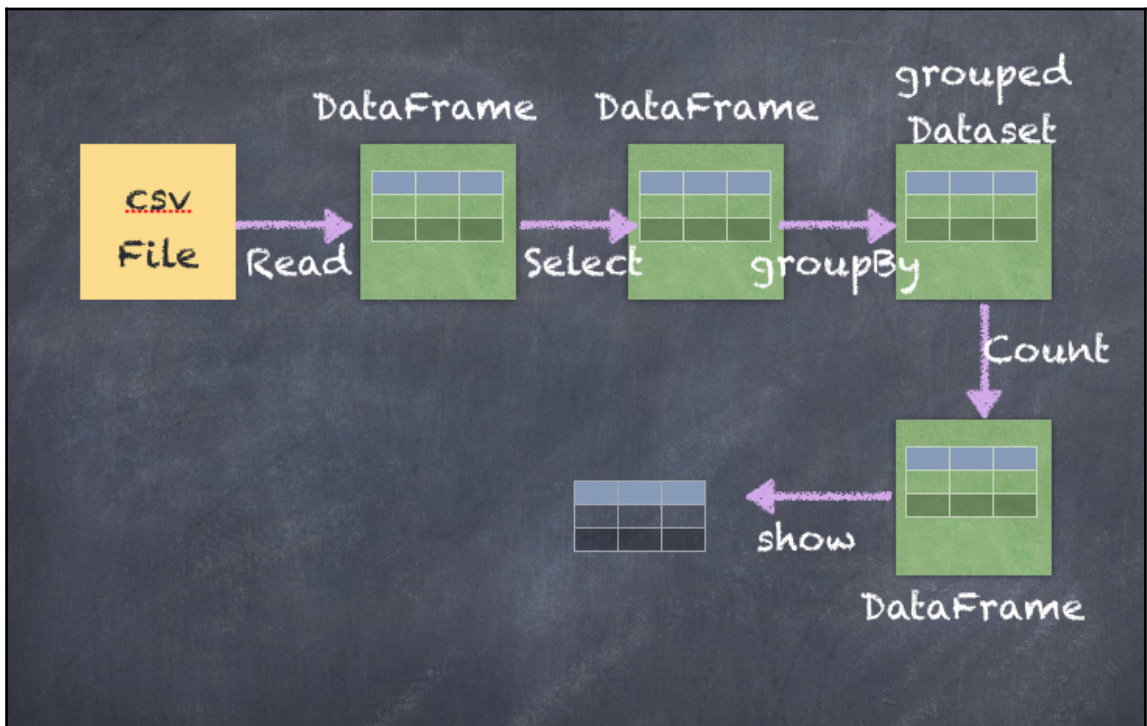
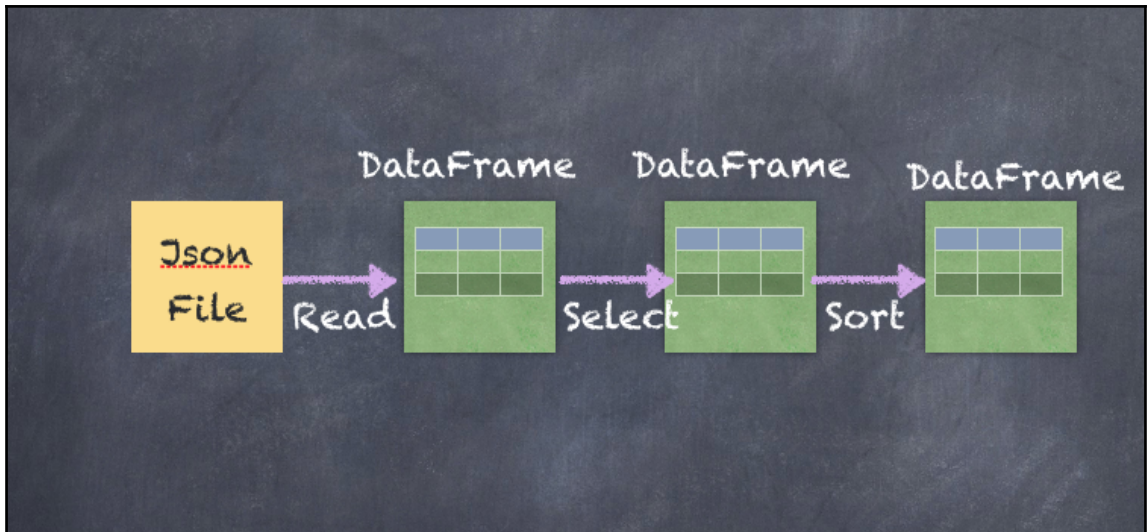
# **Chapter 8: Introduce a Little Structure**

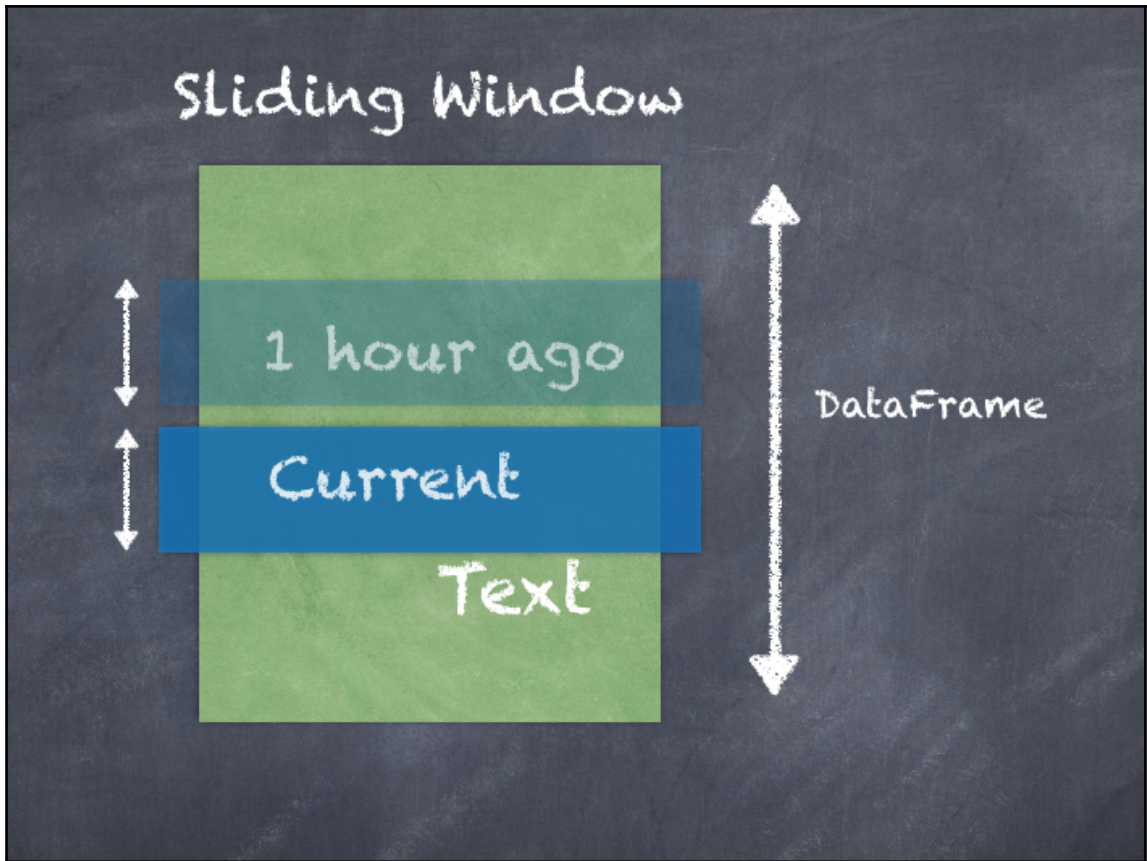
## **SparkSQL**

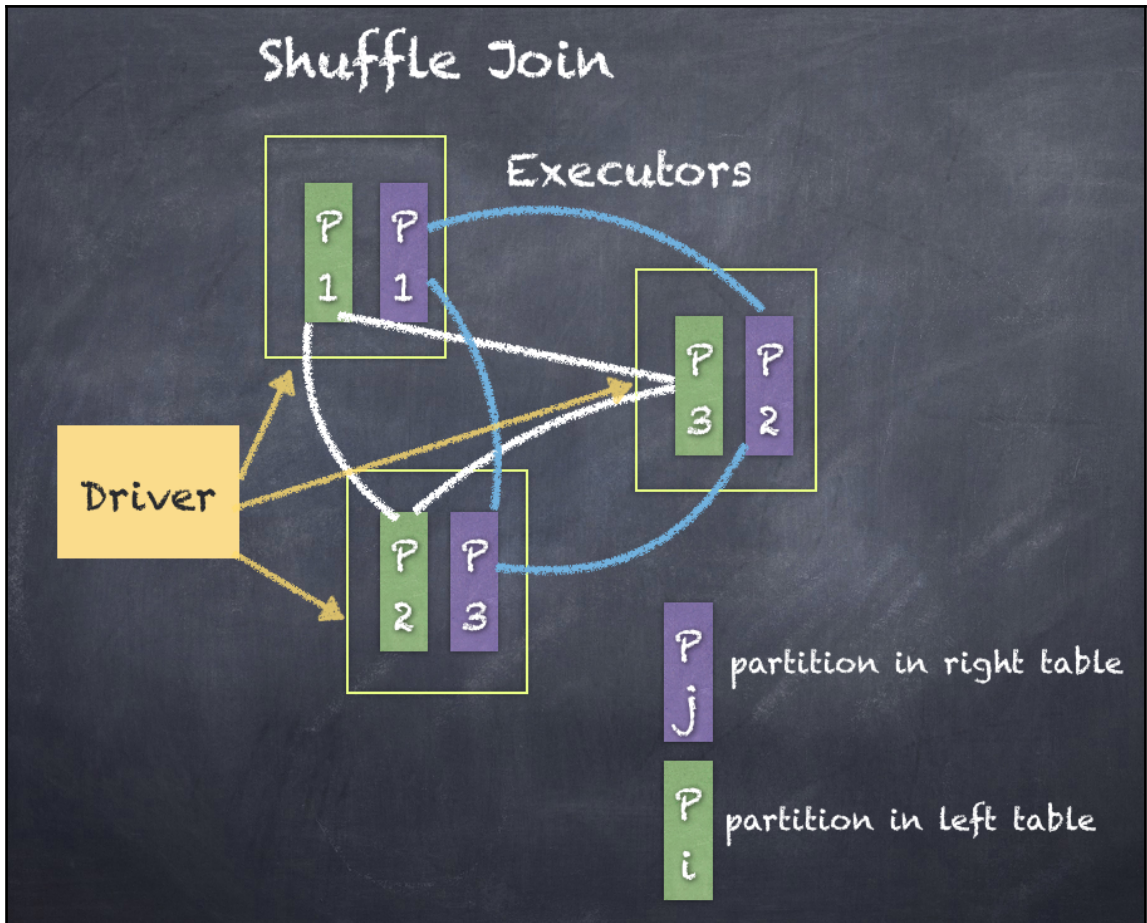


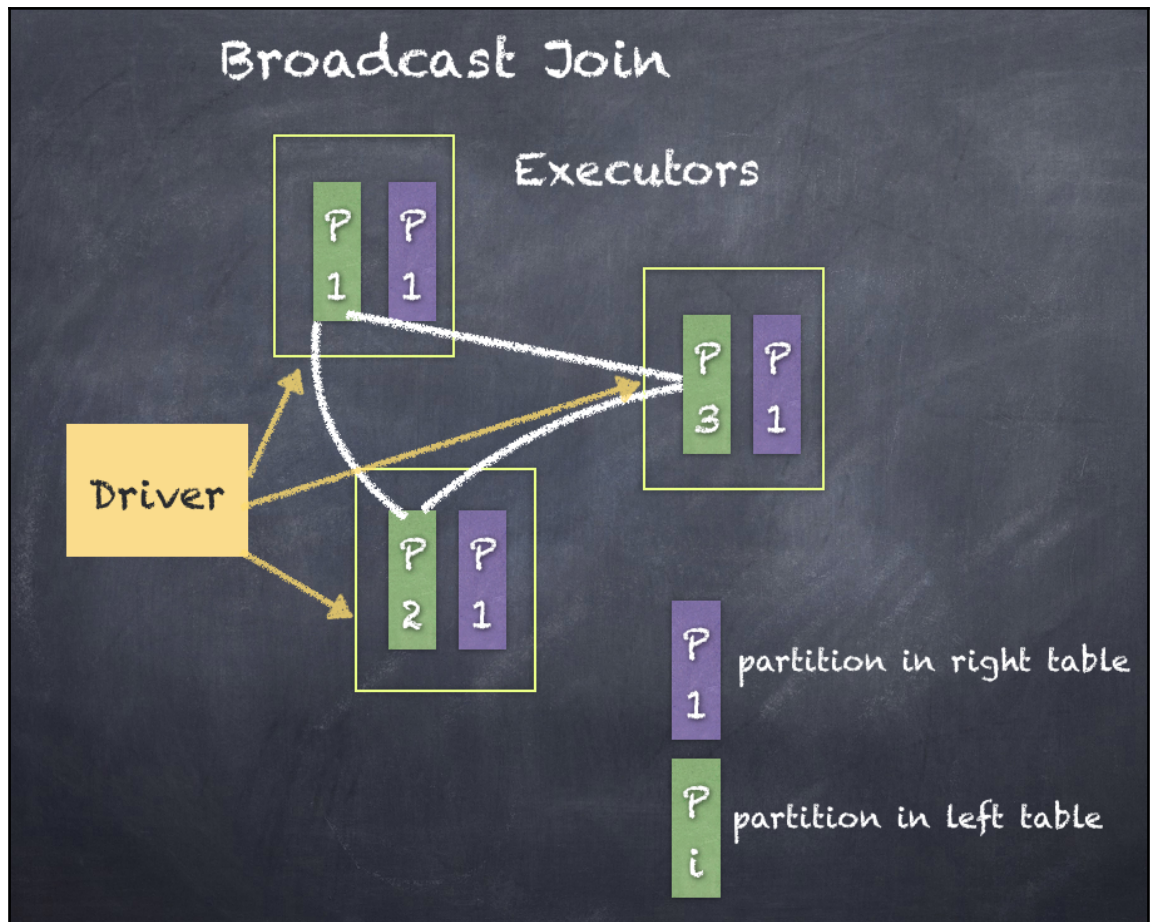


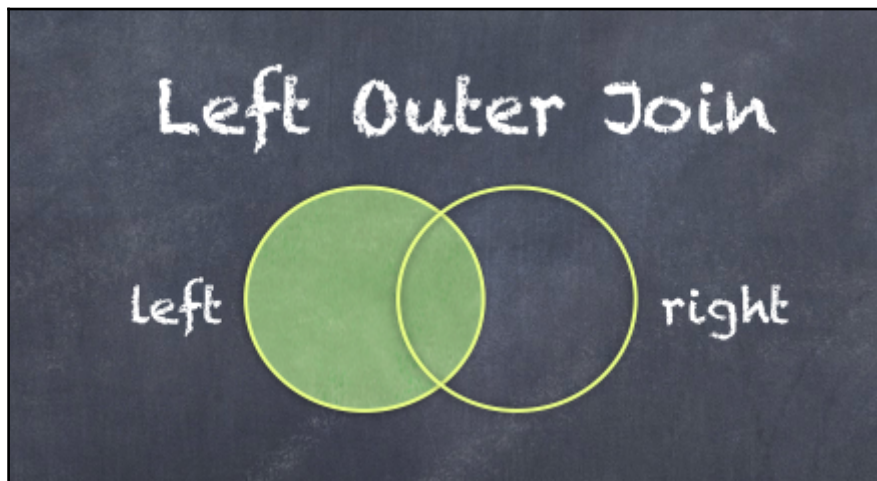
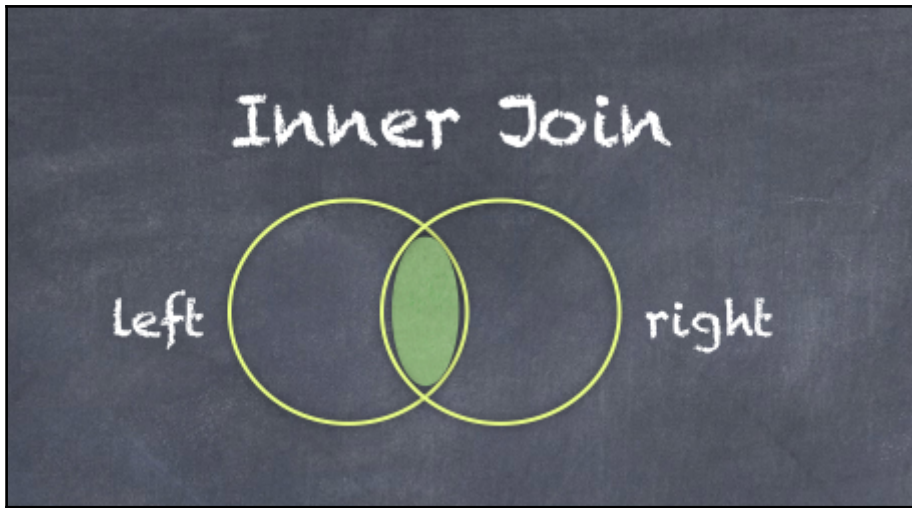


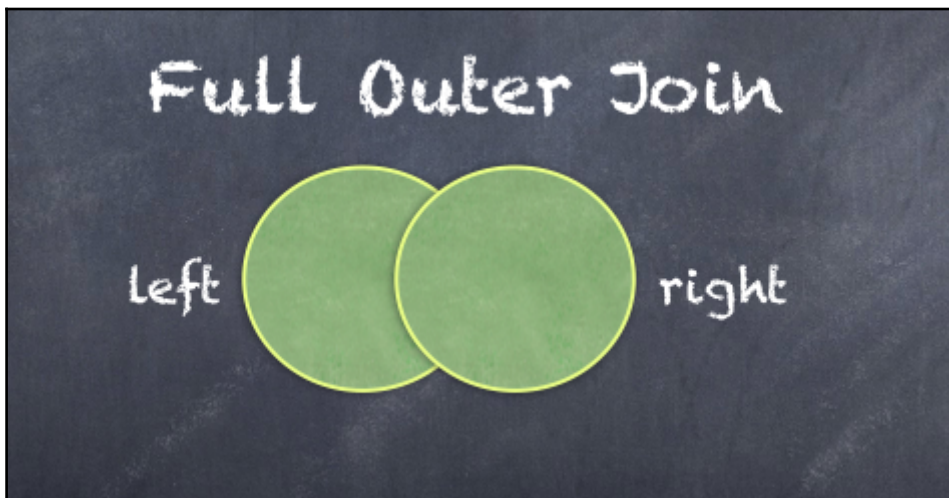
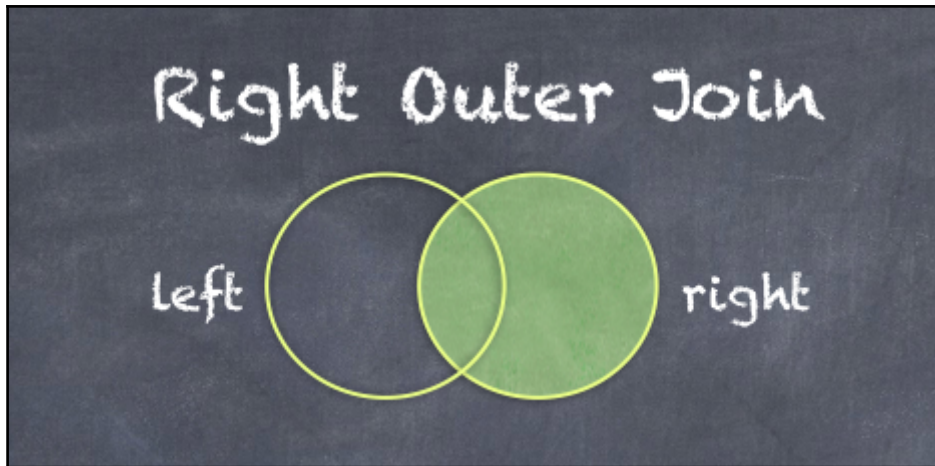






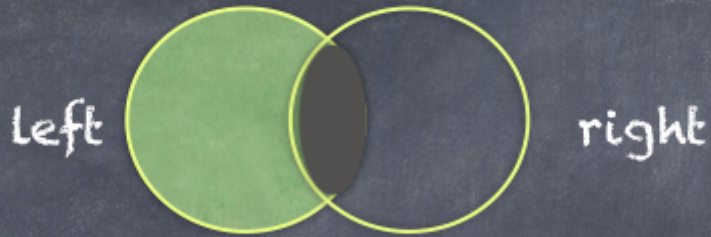




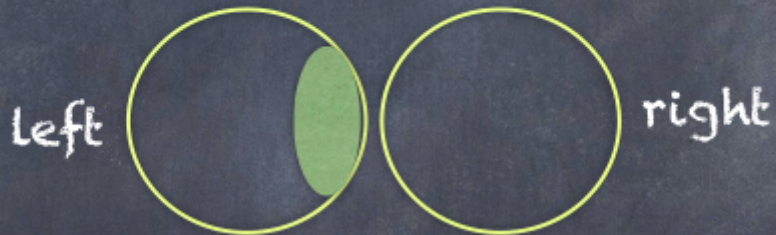


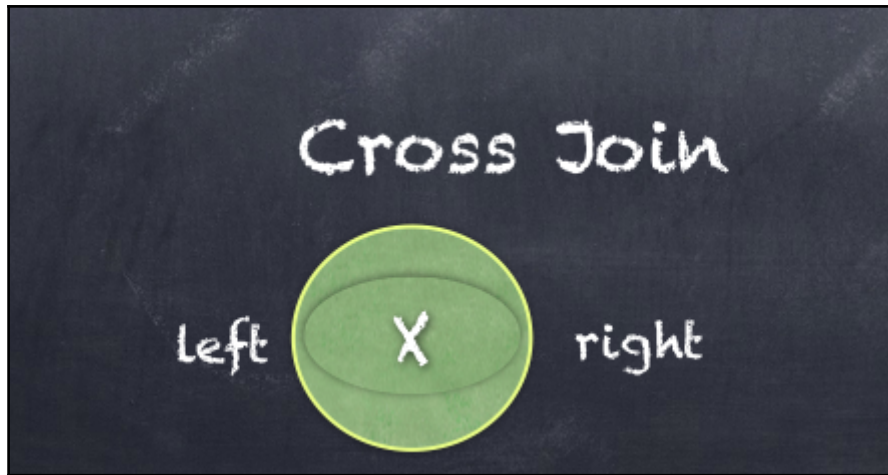


## Left Anti Join

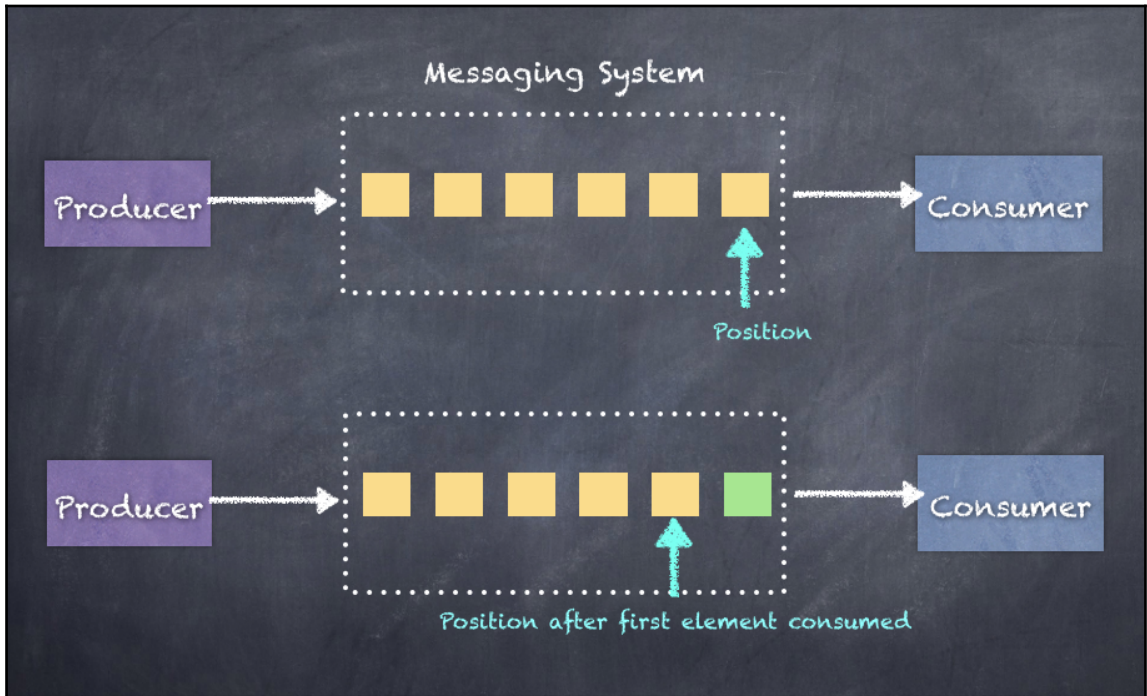


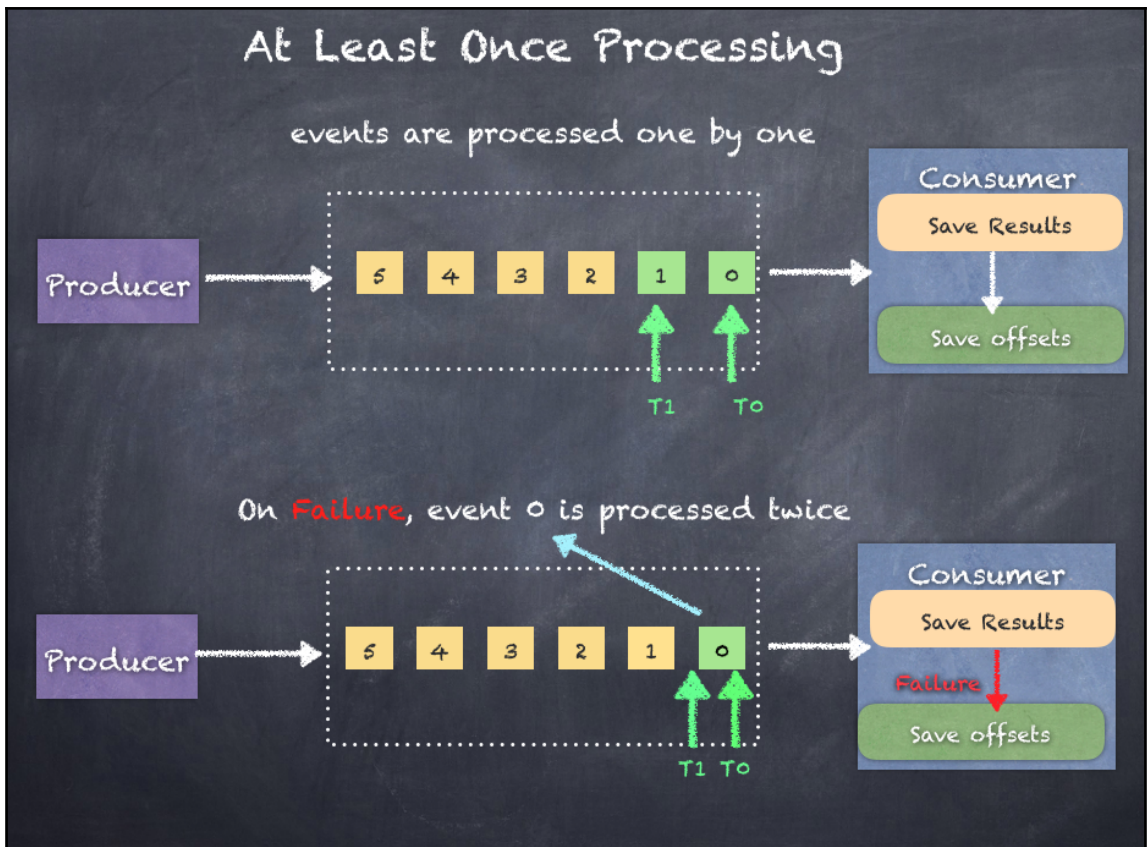
## Left Semi Join

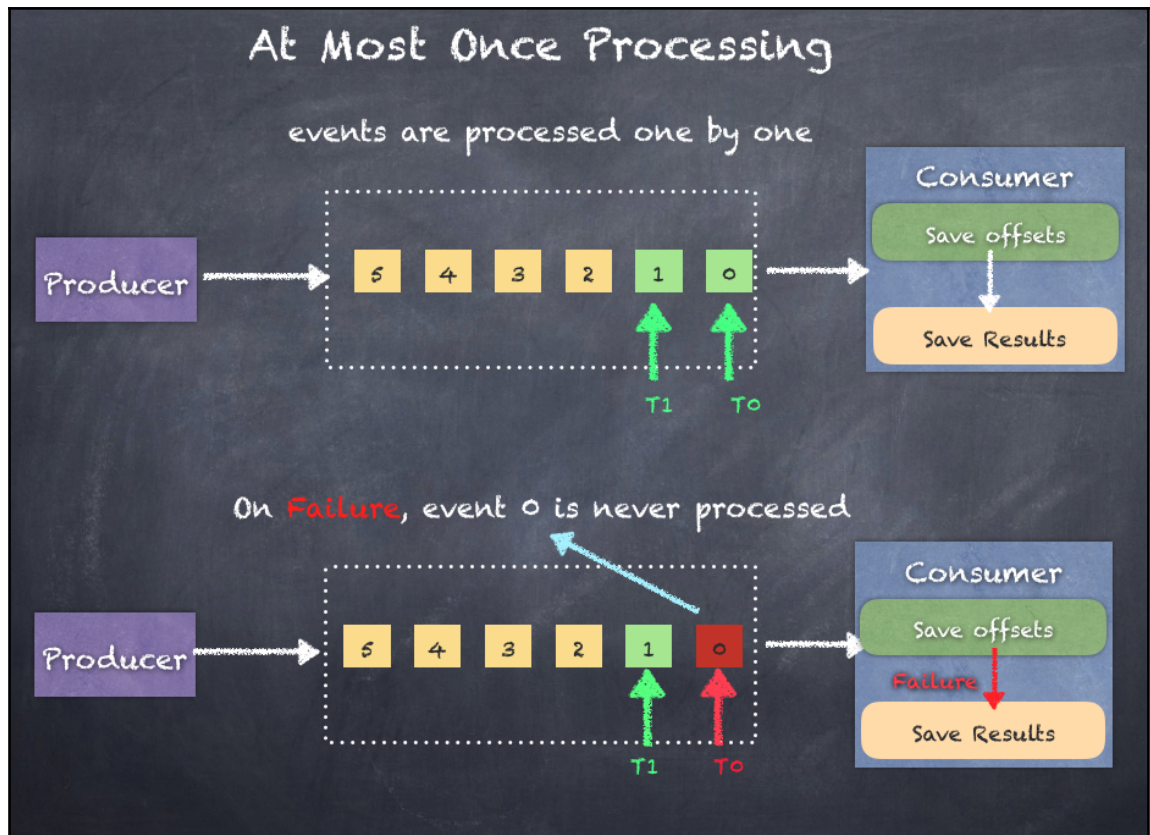


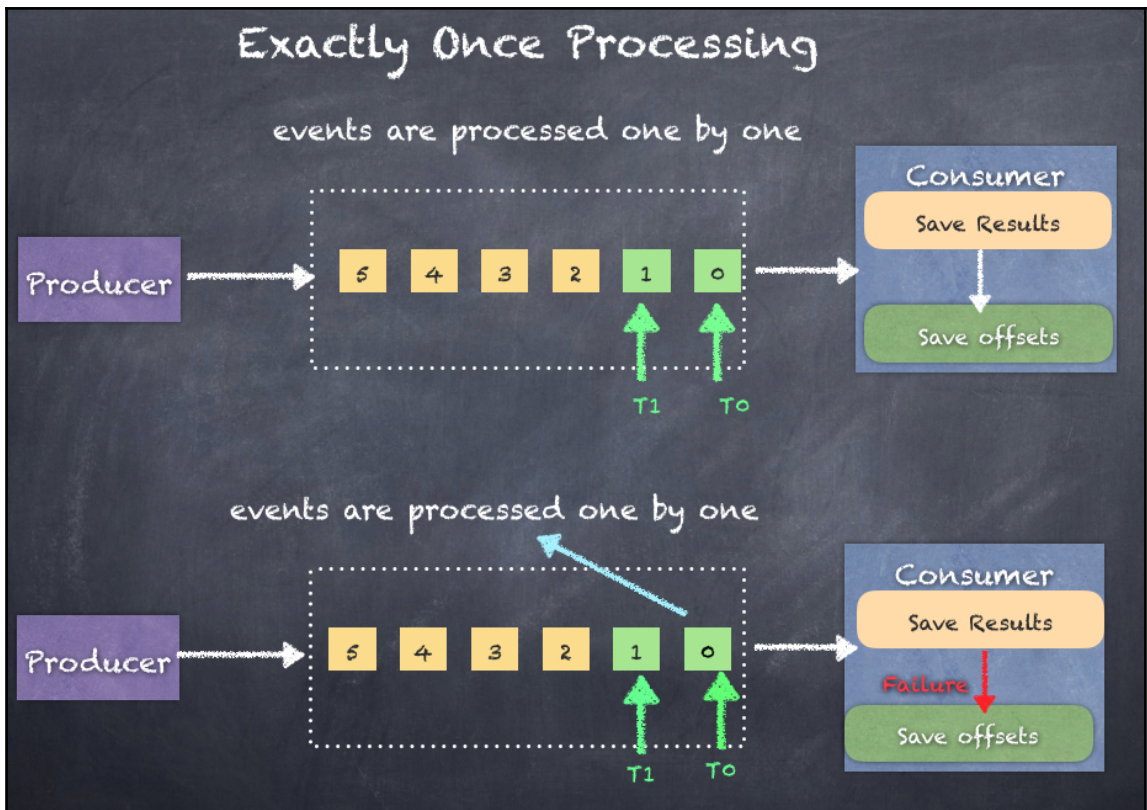


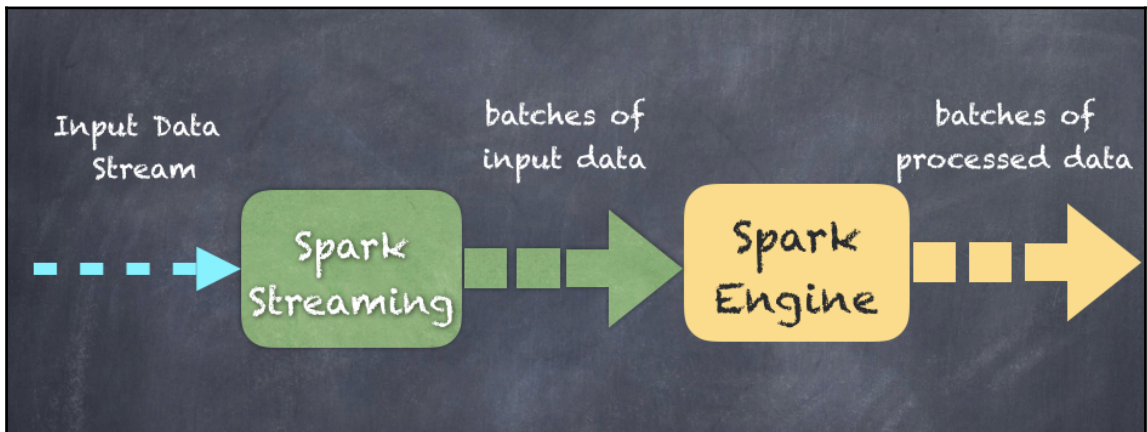
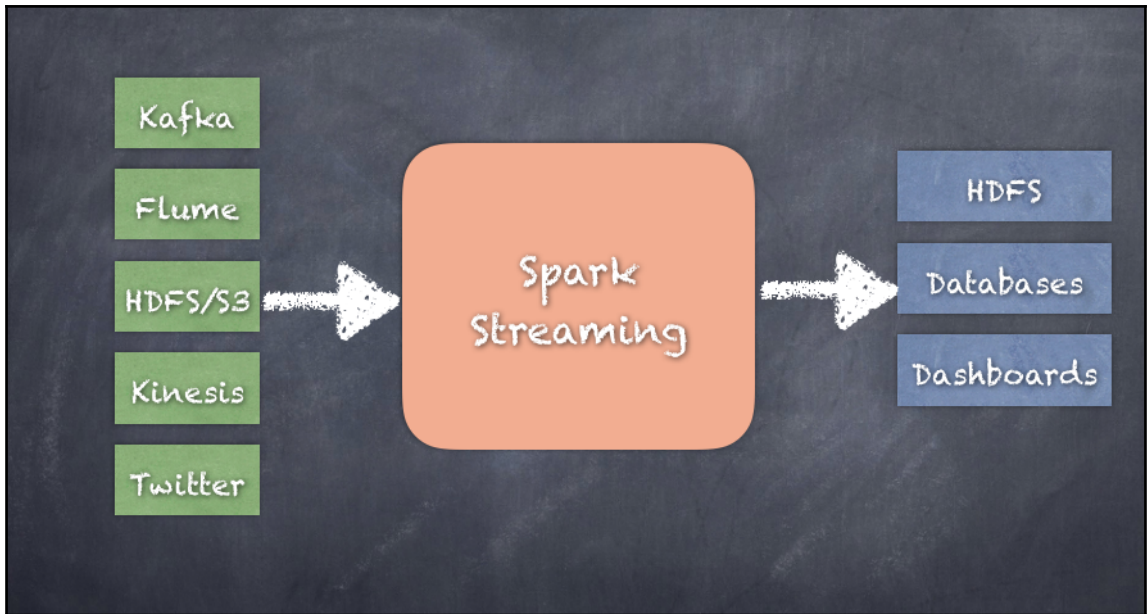
# Chapter 9: Stream Me Up Scotty - Spark Streaming

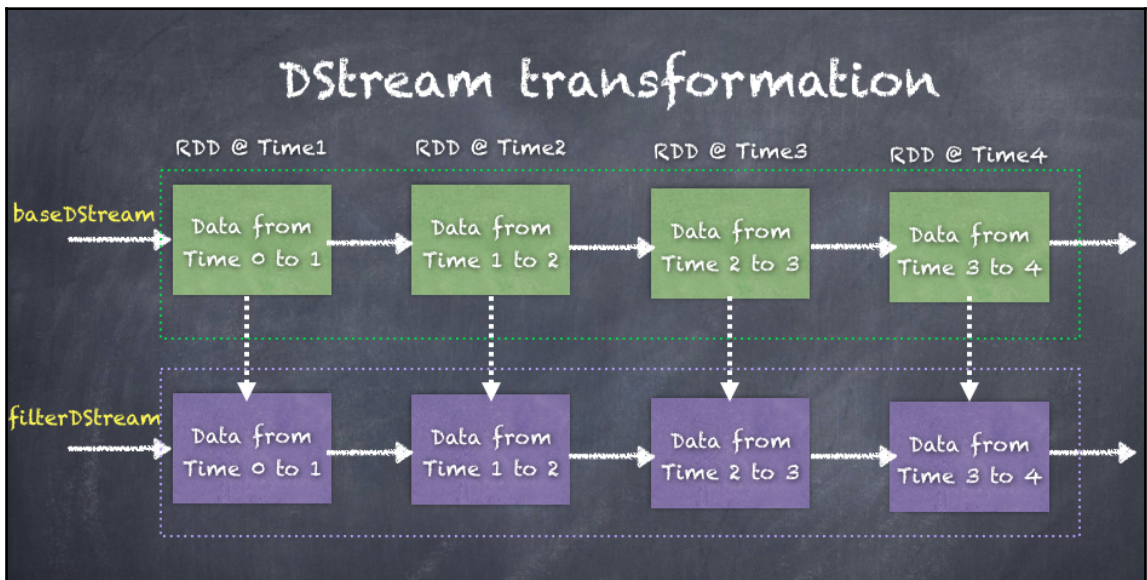
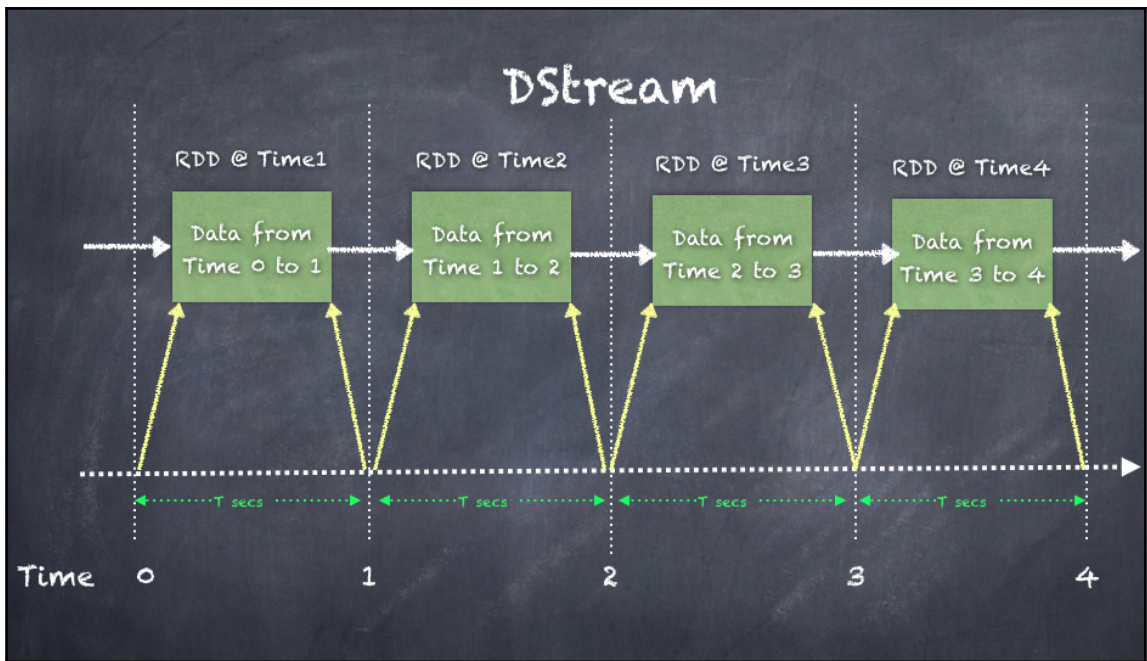




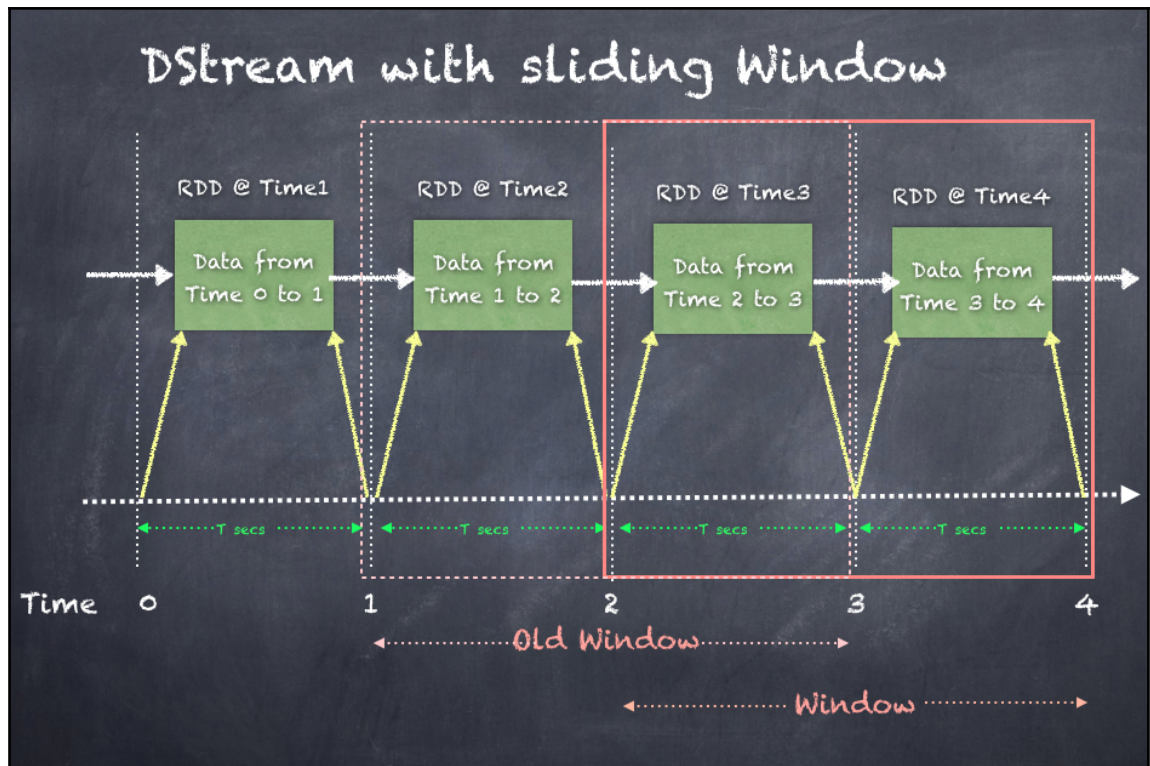


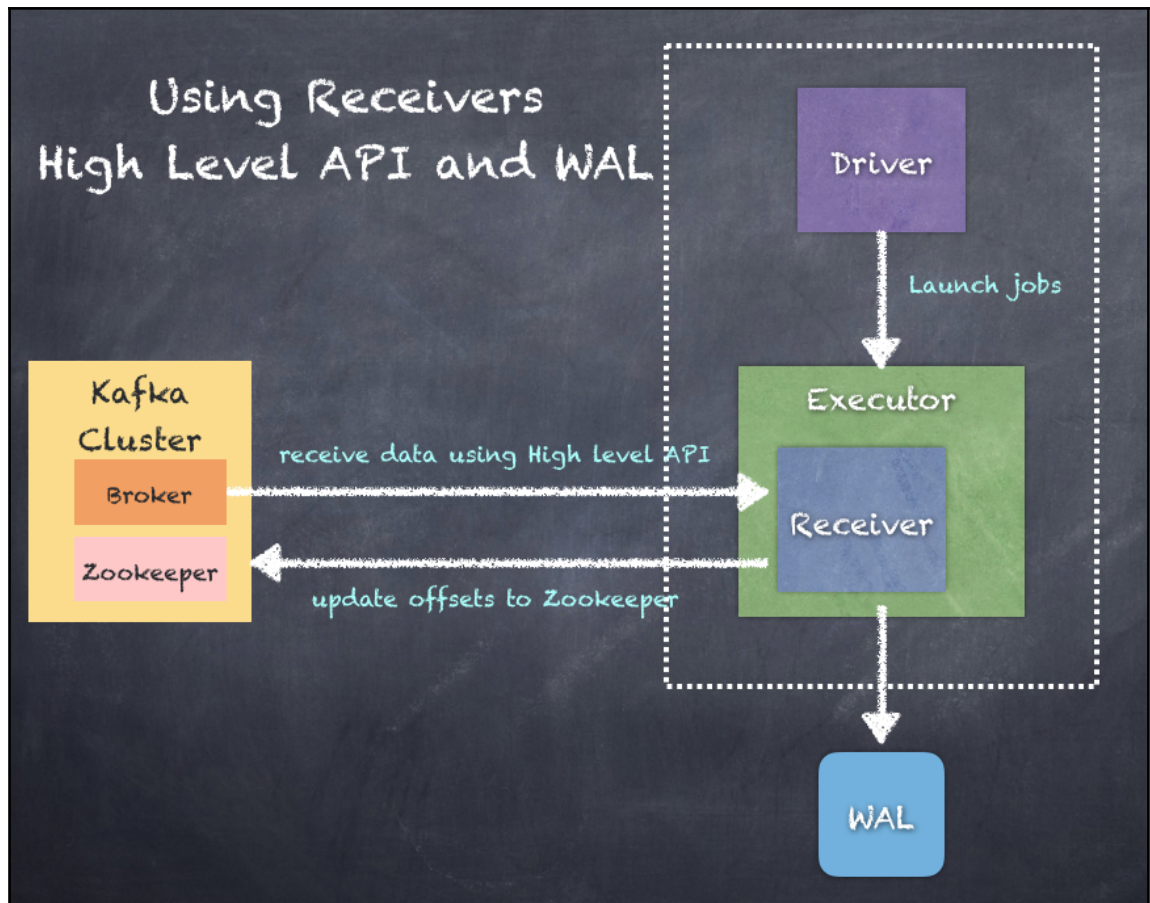


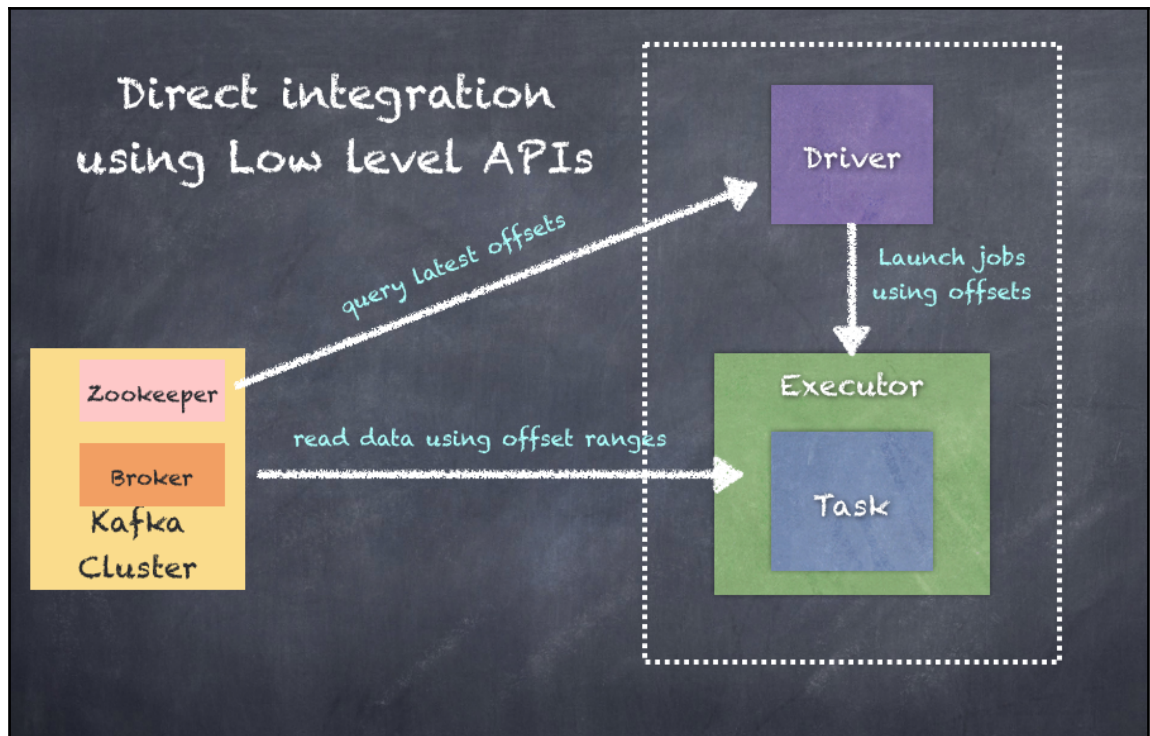




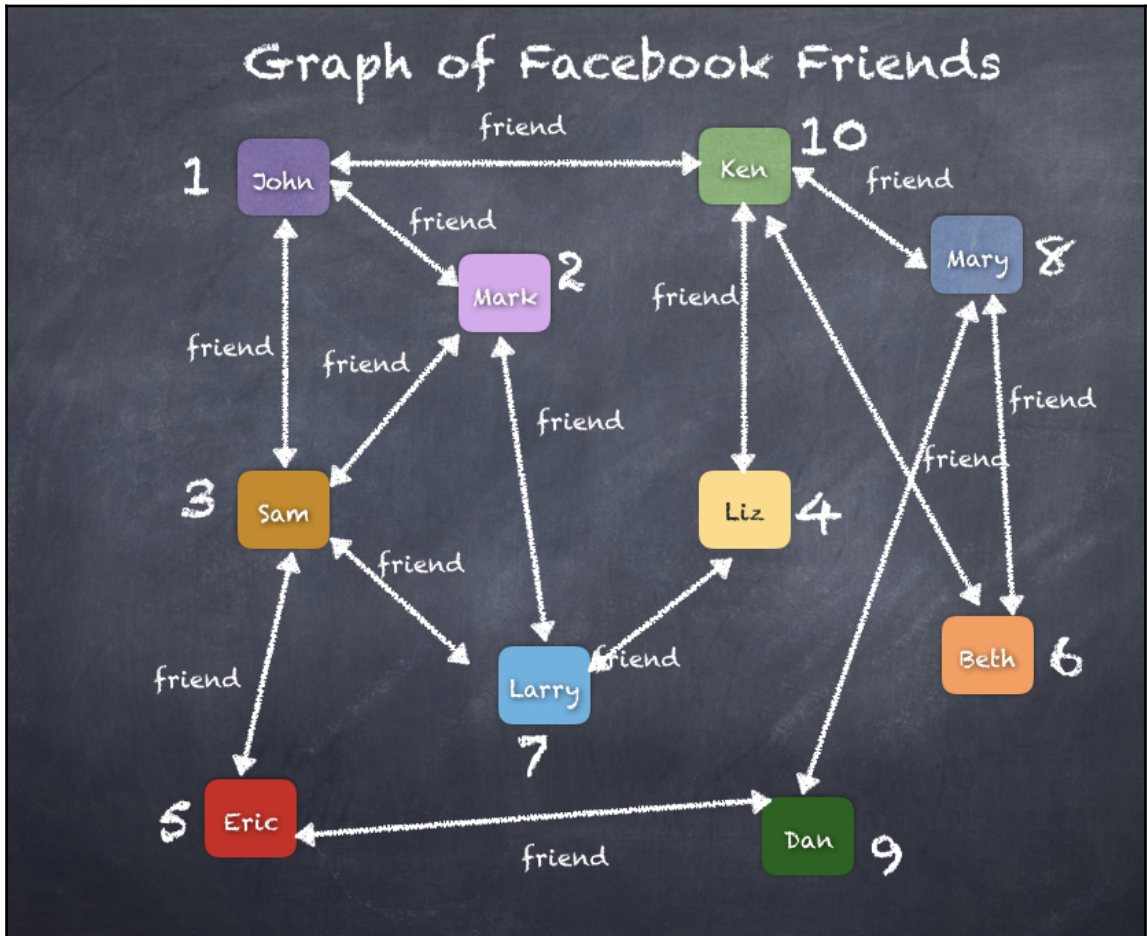


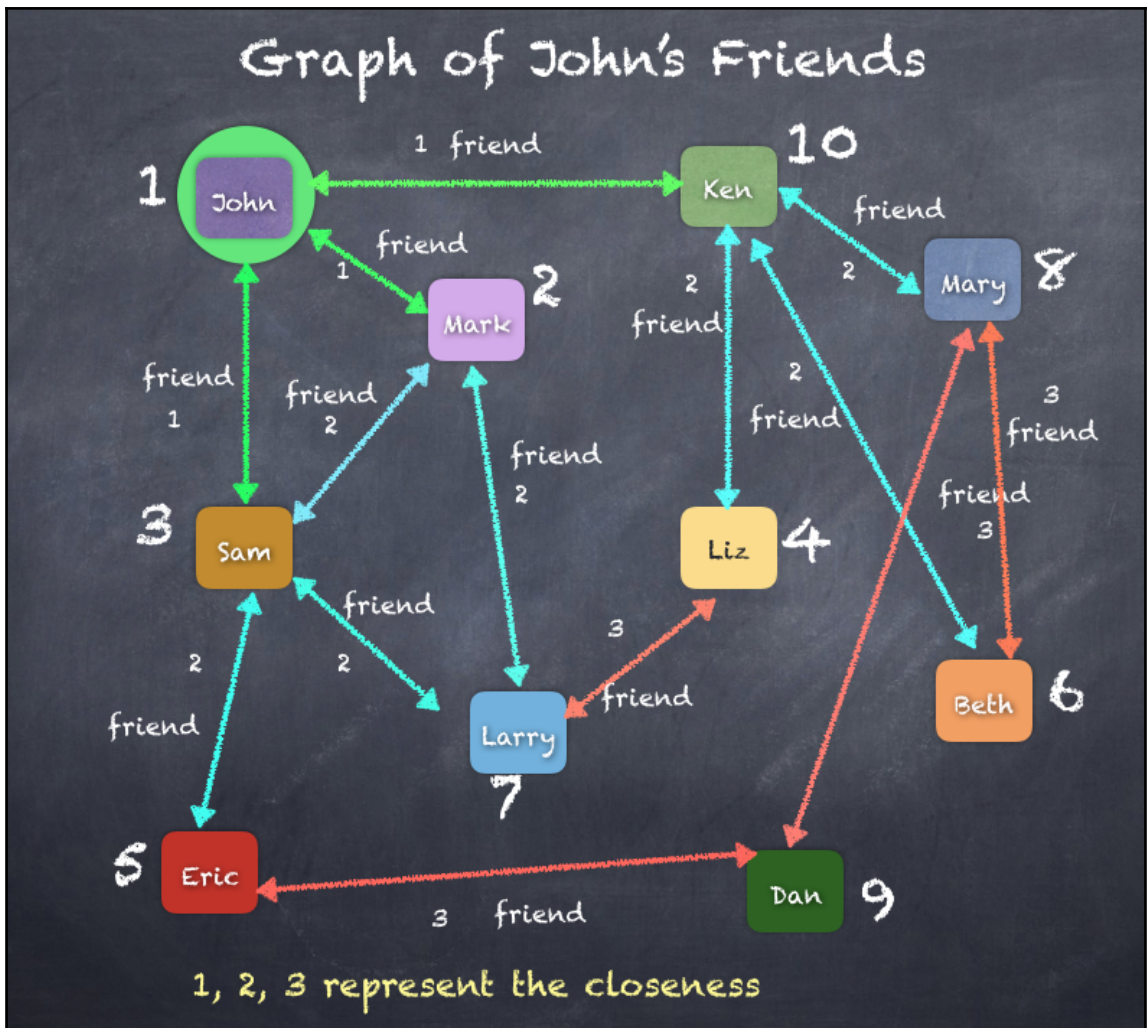






# Chapter 10: Everything is Connected - GraphX





# Graph's Vertices and Edges

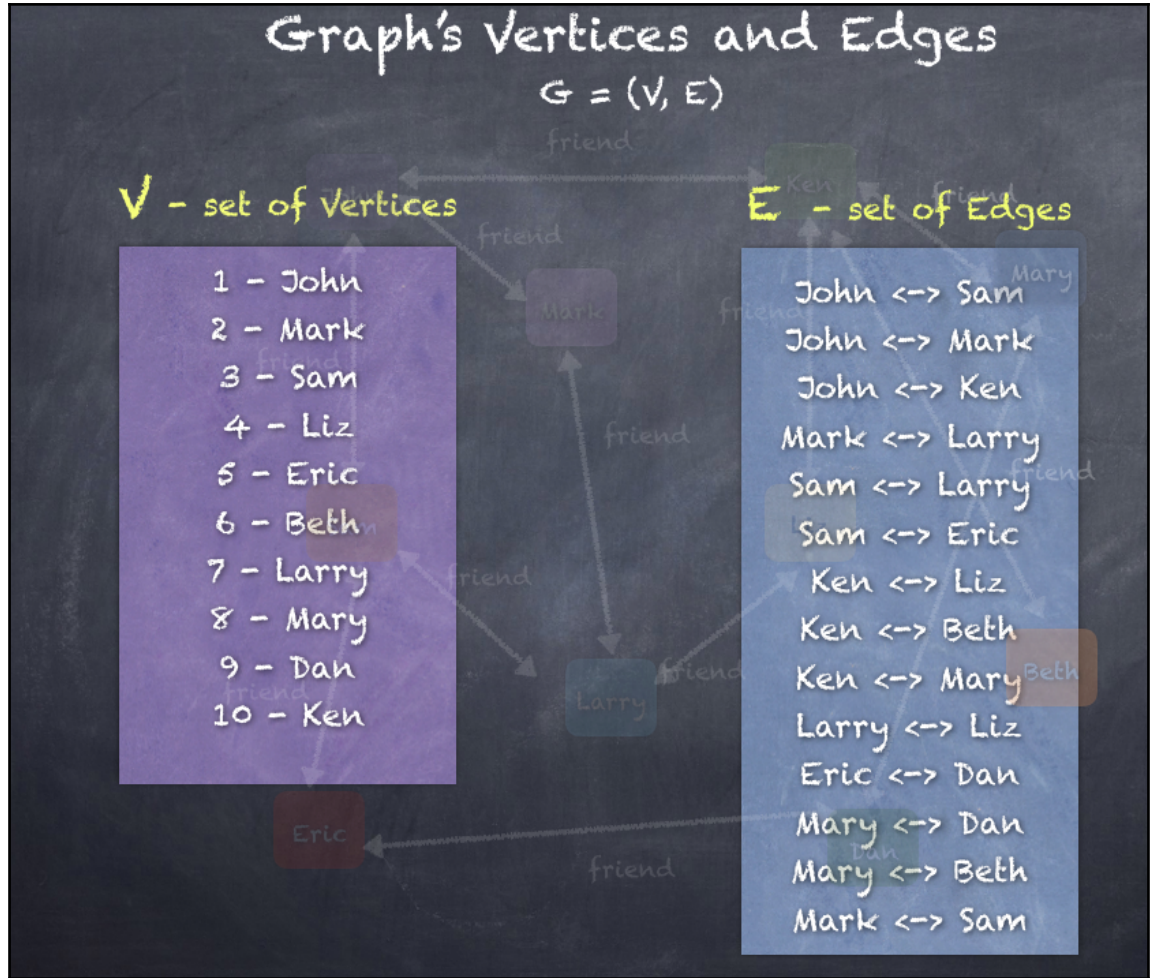
$$G = (V, E)$$

**V** - set of Vertices

- 1 - John
- 2 - Mark
- 3 - Sam
- 4 - Liz
- 5 - Eric
- 6 - Beth
- 7 - Larry
- 8 - Mary
- 9 - Dan
- 10 - Ken

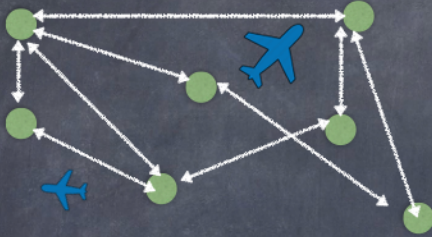
**E** - set of Edges

- John ↔ Sam
- John ↔ Mark
- John ↔ Ken
- Mark ↔ Larry
- Sam ↔ Larry
- Sam ↔ Eric
- Ken ↔ Liz
- Ken ↔ Beth
- Ken ↔ Mary
- Larry ↔ Liz
- Eric ↔ Dan
- Mary ↔ Dan
- Mary ↔ Beth
- Mark ↔ Sam

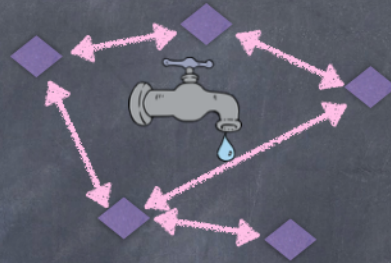


# Graphs in real-life

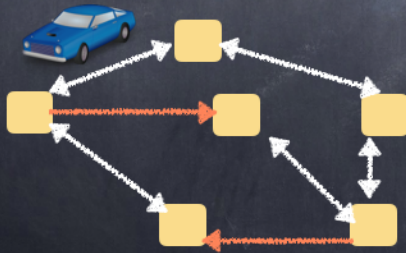
## Flight Routes



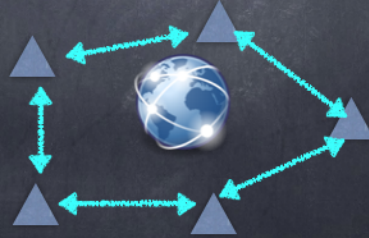
## Water Pipelines

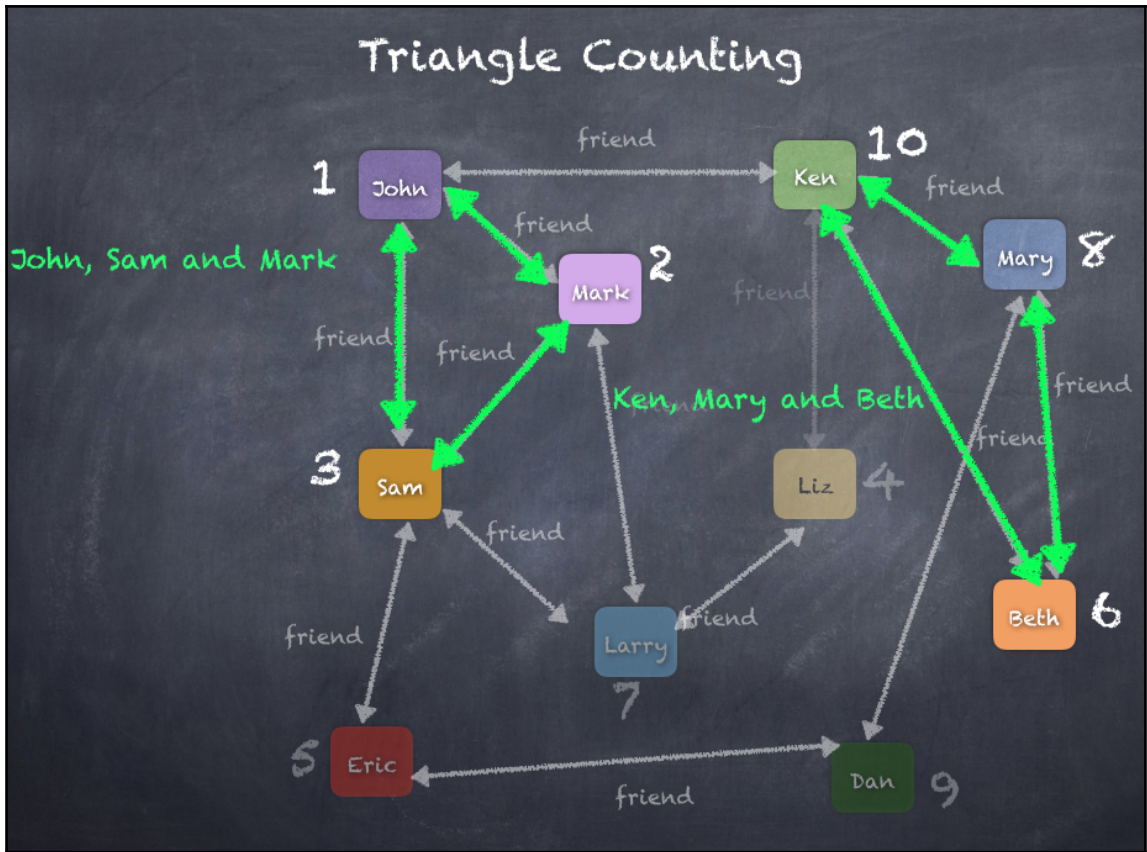


## GPS

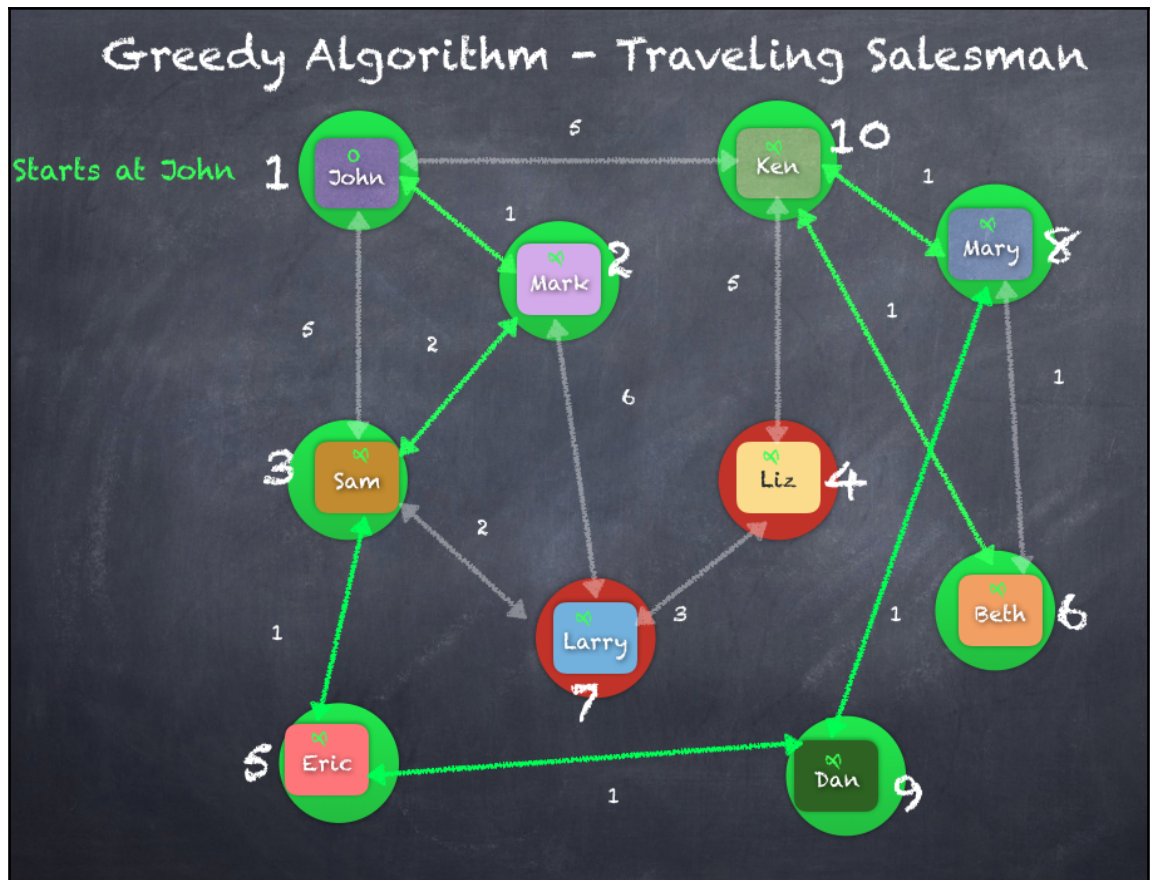


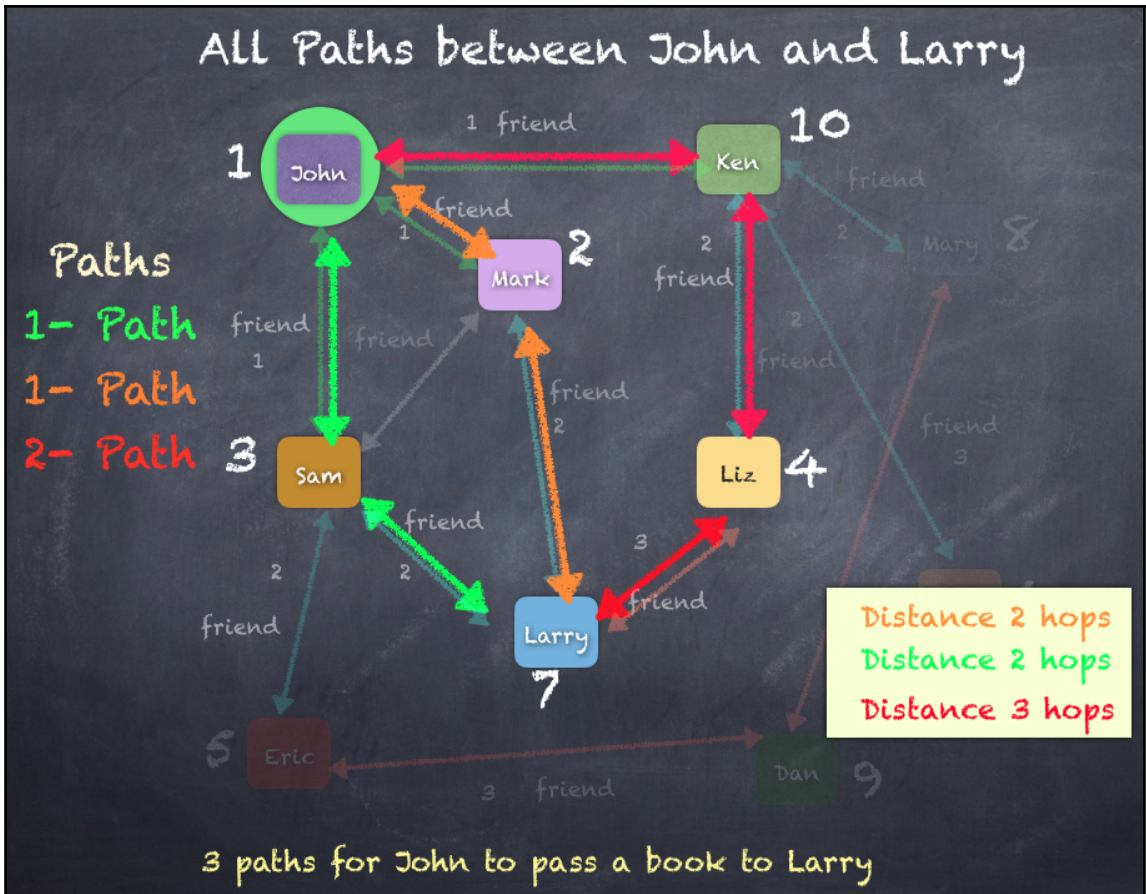
## Internet Routers

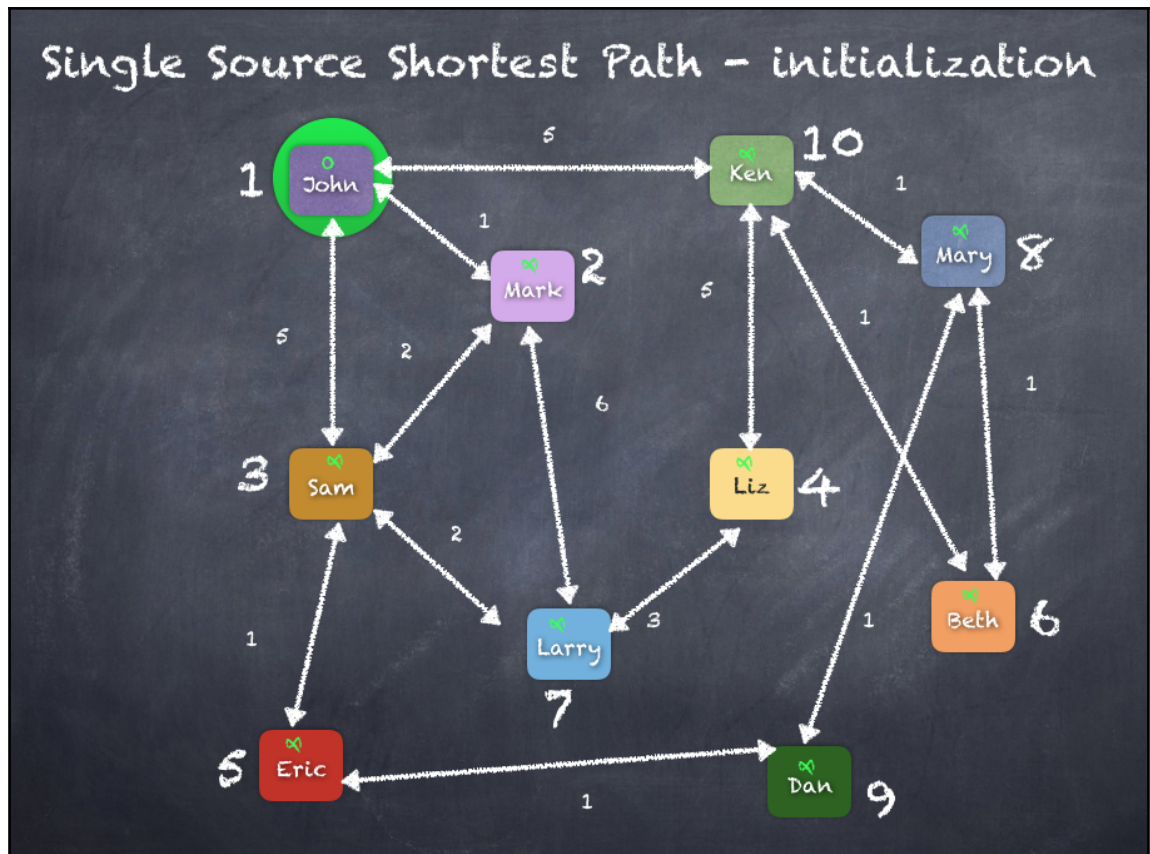


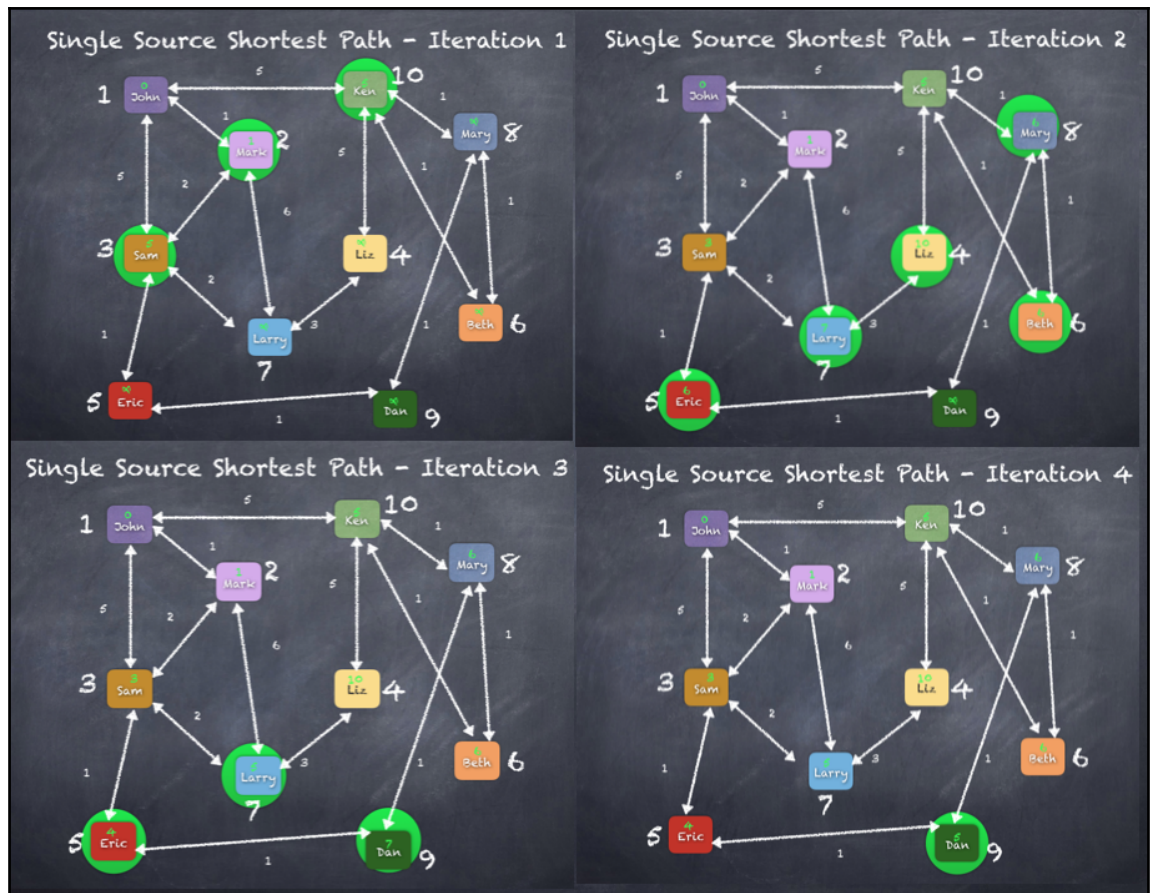


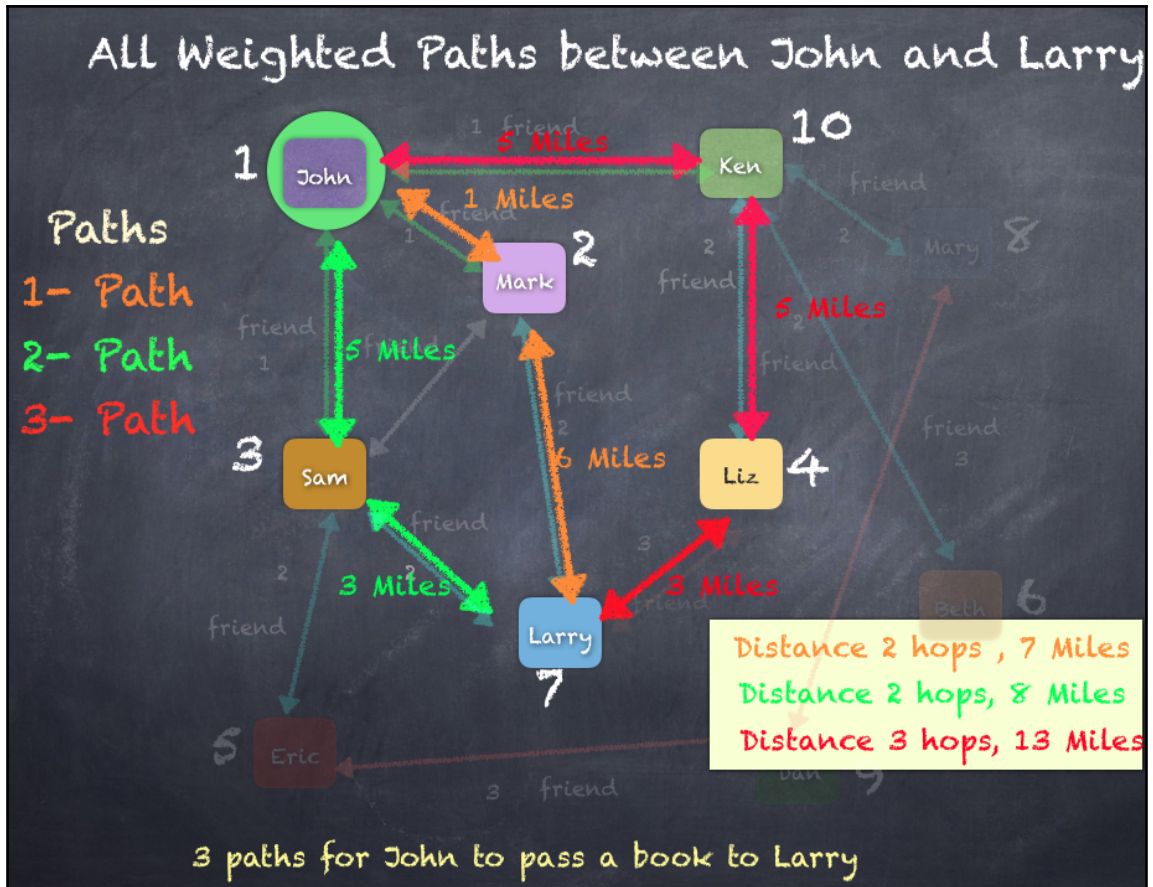


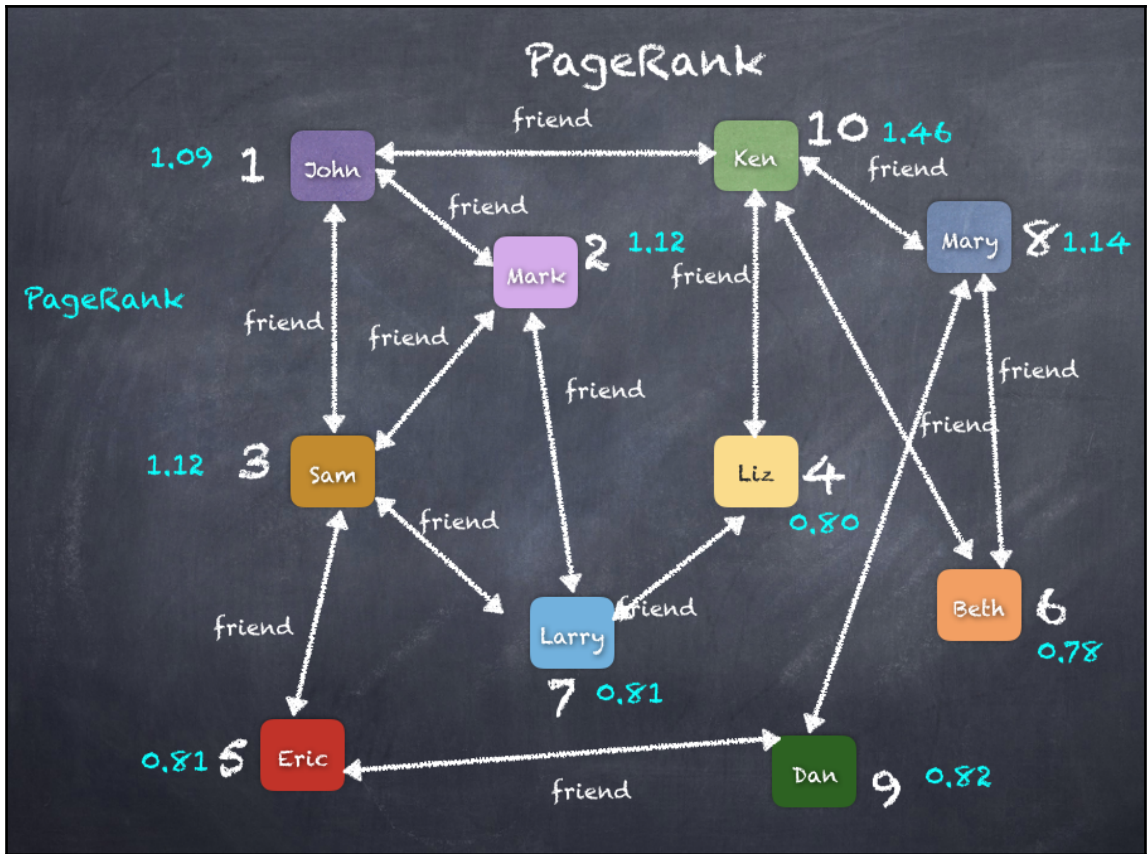




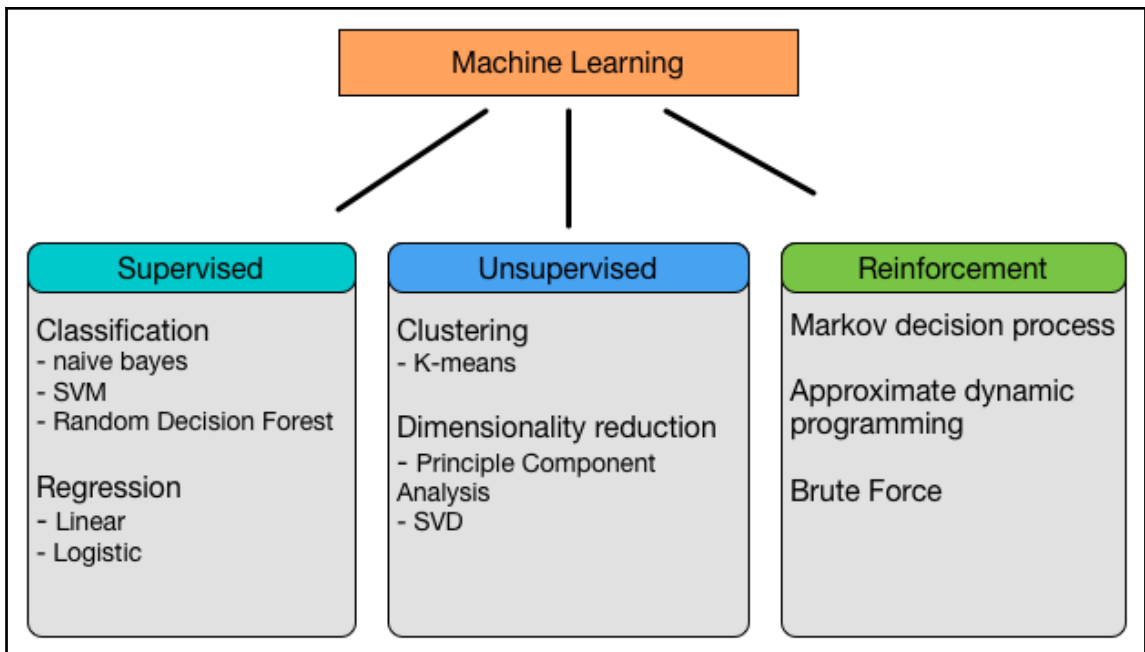
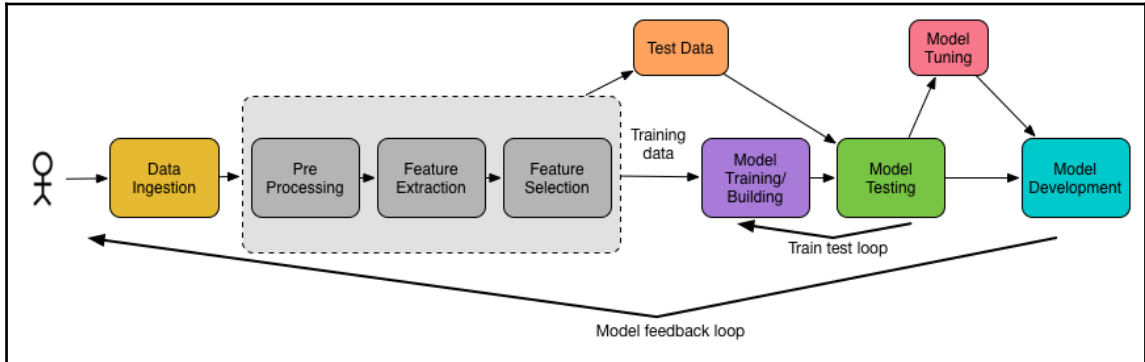


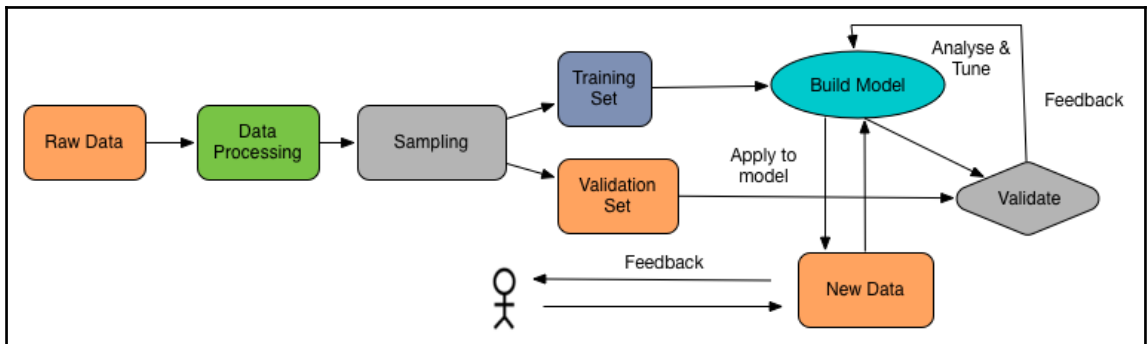
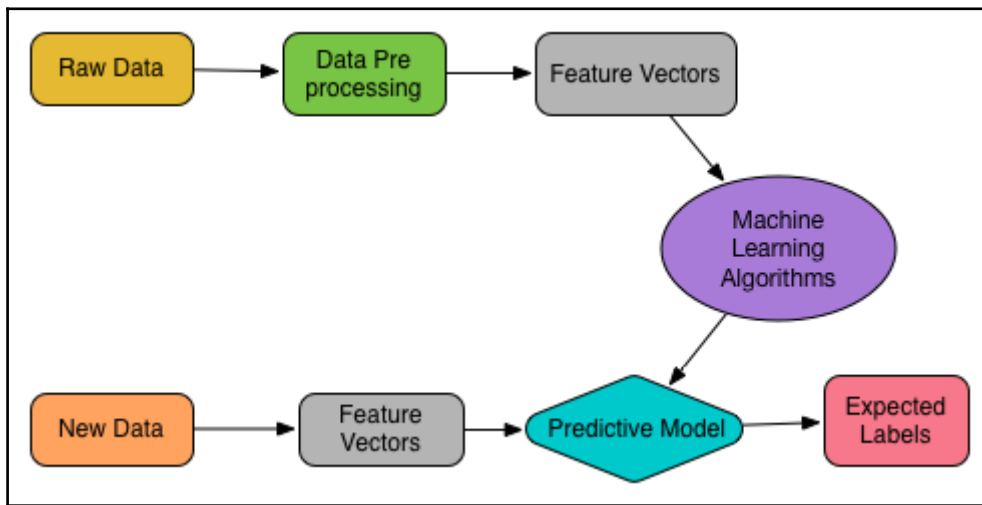
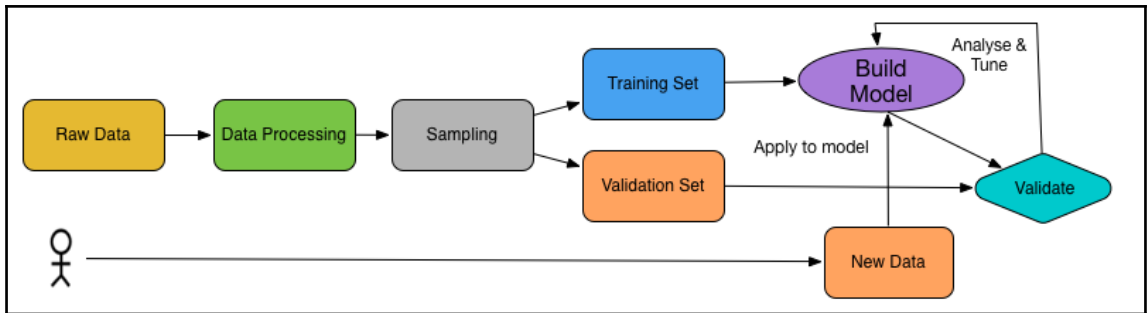




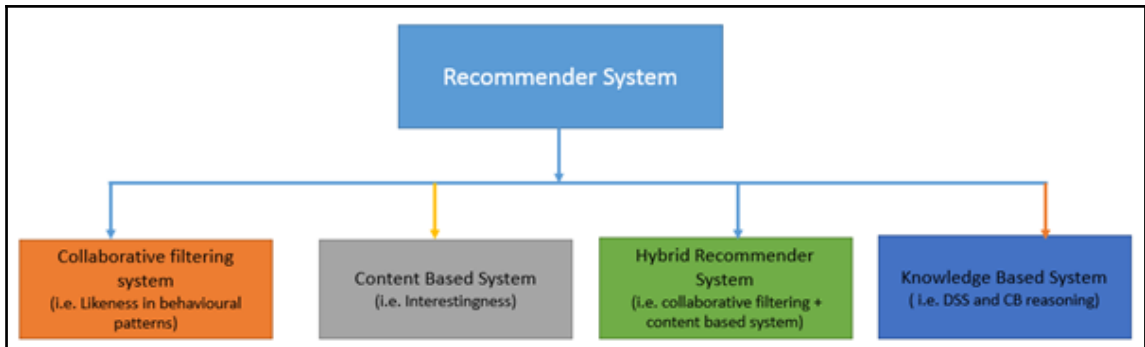


# Chapter 11: Learning Machine Learning - Spark MLlib and Spark ML









id	name
0	[Jason, David]
1	[David, Martin]
2	[Martin, Jason]
3	[Jason, Dael]
4	[Dael, Martin]
5	[Moahmed, Jason]
6	[David, David]
7	[Jason, Martin]

id	name	features
0	[Jason, David]	(3, [0, 1], [1.0, 1.0])
1	[David, Martin]	(3, [1, 2], [1.0, 1.0])
2	[Martin, Jason]	(3, [0, 2], [1.0, 1.0])
3	[Jason, Dael]	(3, [0], [1.0])
4	[Dael, Martin]	(3, [2], [1.0])
5	[Moahmed, Jason]	(3, [0], [1.0])
6	[David, David]	(3, [1], [2.0])
7	[Jason, Martin]	(3, [0, 2], [1.0, 1.0])

sentence	words	tokens
Tokenization, is the process of enchanting words, from the raw text	[tokenization, is, the, process, of, enchanting, words, from, the, raw, text]	9
If you want, to have more advance tokenization, RegexTokenizer, is a good option	[, if, you, want, to, have, more, advance, tokenization, regextokenizer, is, a, good, option]	11
Here, will provide a sample example on how to tokenize sentences	[, here, will, provide, a, sample, example, on, how, to, tokenize, sentences]	11
This way, you can find all matching occurrences	[this, way, you, can, find, all, matching, occurrences]	7

sentence	words	tokens
Tokenization, is the process of enchanting words, from the raw text	[tokenization, is, the, process, of, enchanting, words, from, the, raw, text]	11
If you want, to have more advance tokenization, RegexTokenizer, is a good option	[if, you, want, to, have, more, advance, tokenization, regextokenizer, is, a, good, option]	13
Here, will provide a sample example on how to tokenize sentences	[here, will, provide, a, sample, example, on, how, to, tokenize, sentences]	11
This way, you can find all matching occurrences	[this, way, you, can, find, all, matching, occurrences]	8

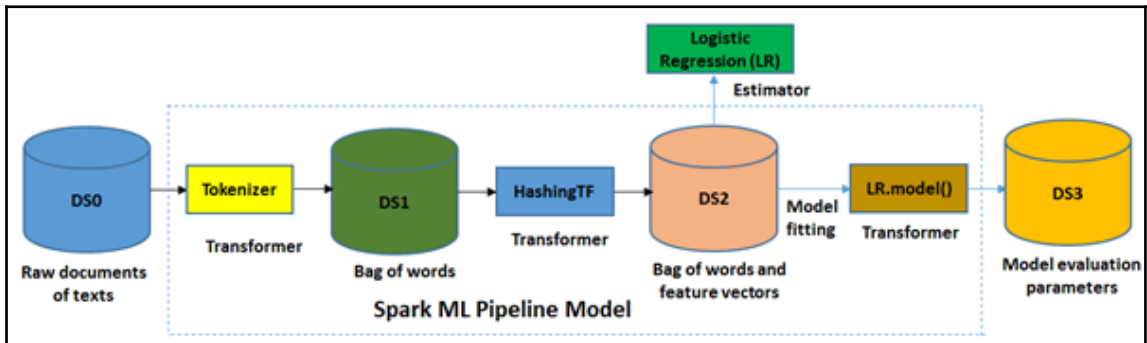
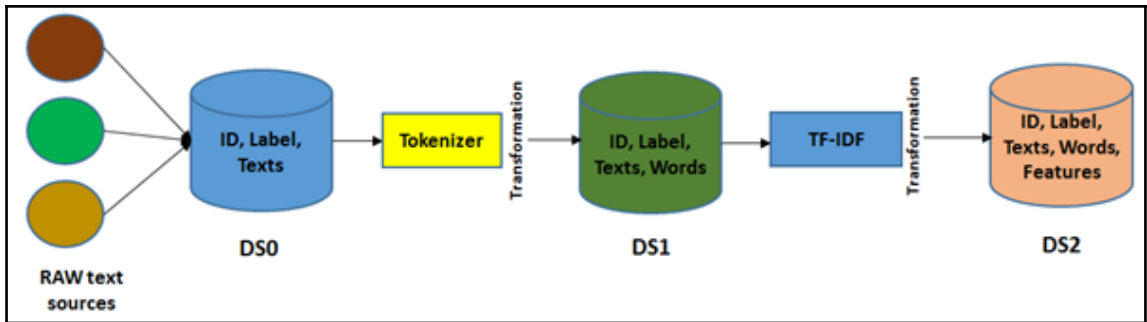
id	filtered
0	[tokenization, process, enchanting, words, raw, text]
1	[want, advance, tokenization, regextokenizer, good, option]
2	[provide, sample, example, tokenize, sentences]
3	[way, find, matching, occurrences]

id	name	address
0	Jason	Germany
1	David	France
2	Martin	Spain
3	Jason	USA
4	Daiel	UK
5	Moahmed	Bangladesh
6	David	Ireland
7	Jason	Netherlands

id	name	address	label
0	Jason	Germany	0.0
1	David	France	1.0
2	Martin	Spain	3.0
3	Jason	USA	0.0
4	Daiel	UK	4.0
5	Moahmed	Bangladesh	2.0
6	David	Ireland	1.0
7	Jason	Netherlands	0.0

id	name	address
0	Jason	Germany
1	David	France
2	Martin	Spain
3	Jason	USA
4	Daiel	UK
5	Moahmed	Bangladesh
6	David	Ireland
7	Jason	Netherlands

id	name	address	categoryIndex	categoryVec
0	Jason	Germany	0.0	(4, [0], [1.0])
1	David	France	1.0	(4, [1], [1.0])
2	Martin	Spain	3.0	(4, [3], [1.0])
3	Jason	USA	0.0	(4, [0], [1.0])
4	Daiel	UK	4.0	(4, [], [])
5	Moahmed	Bangladesh	2.0	(4, [2], [1.0])
6	David	Ireland	1.0	(4, [1], [1.0])
7	Jason	Netherlands	0.0	(4, [0], [1.0])



```

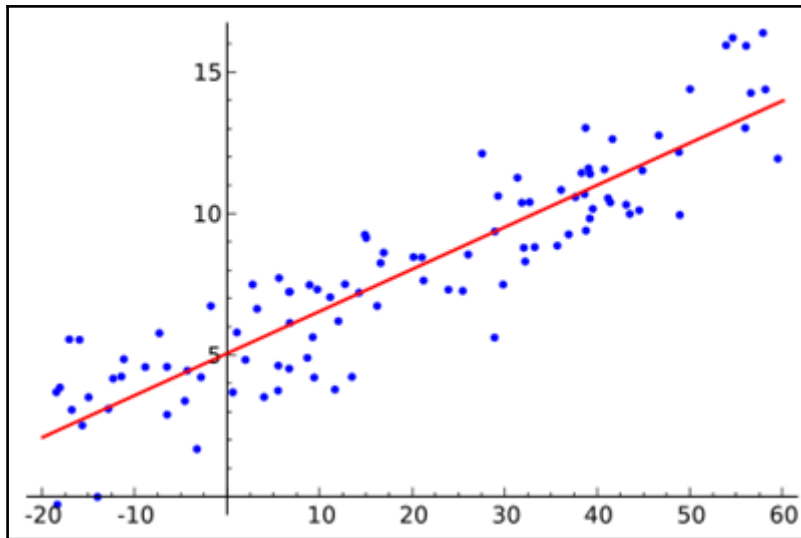
+-----+
| features |
+-----+
|[3.5,2.0,5.0,6.3,5.6,2.4] |
|[4.4,0.1,3.0,9.0,7.0,8.75] |
|[3.2,2.4,0.0,6.0,7.4,3.34] |
+-----+
    
```

```

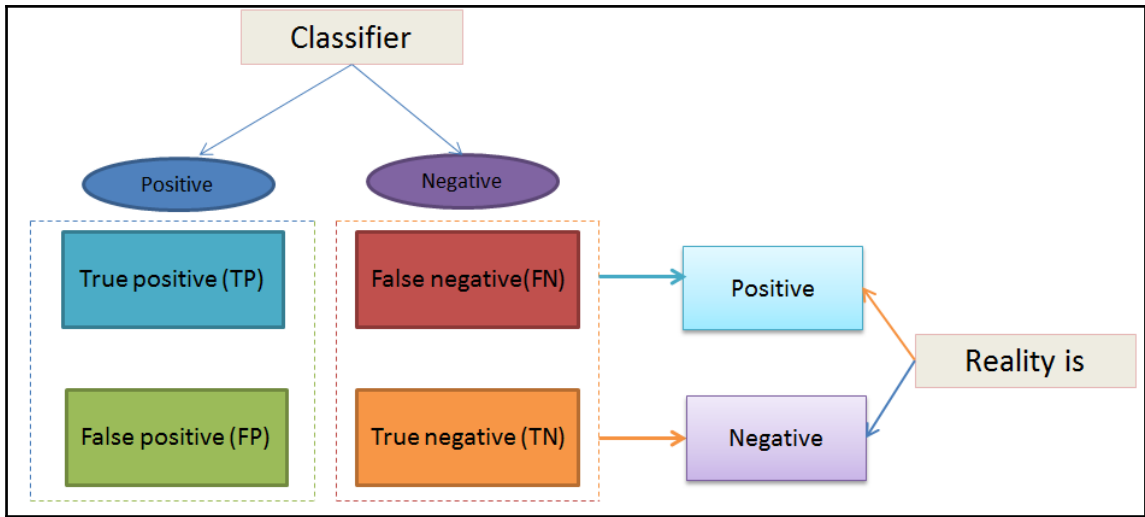
+-----+
| pcaFeatures |
+-----+
| [-5.149253129088702,3.2157431427730385,-6.828533673828153,5.774261462142295] |
| [-12.372614091904445,0.804196667817684,-6.828533673828154,5.774261462142296] |
| [-5.649682494292658,-2.189177804885822,-6.828533673828155,5.7742614621422925] |
+-----+
    
```

label	features
5.0	(780, [152, 153, 154...]
0.0	(780, [127, 128, 129...]
4.0	(780, [160, 161, 162...]
1.0	(780, [158, 159, 160...]
9.0	(780, [208, 209, 210...]
2.0	(780, [155, 156, 157...]
1.0	(780, [124, 125, 126...]
3.0	(780, [151, 152, 153...]
1.0	(780, [152, 153, 154...]
4.0	(780, [134, 135, 161...]
3.0	(780, [123, 124, 125...]
5.0	(780, [216, 217, 218...]
3.0	(780, [143, 144, 145...]
6.0	(780, [72, 73, 74, 99...]
1.0	(780, [151, 152, 153...]
7.0	(780, [211, 212, 213...]
2.0	(780, [151, 152, 153...]
8.0	(780, [159, 160, 161...]
6.0	(780, [100, 101, 102...]
9.0	(780, [209, 210, 211...]
4.0	(780, [129, 130, 131...]
0.0	(780, [129, 130, 131...]
9.0	(780, [183, 184, 185...]
1.0	(780, [158, 159, 160...]
1.0	(780, [99, 100, 101, ...]
2.0	(780, [124, 125, 126...]
4.0	(780, [185, 186, 187...]
3.0	(780, [150, 151, 152...]
2.0	(780, [145, 146, 147...]
7.0	(780, [240, 241, 242...]

only showing top 30 rows



$$RMSD = \sqrt{\frac{\sum_{t=1}^n (\hat{y}_t - y_t)^2}{n}}$$



Metric	Definition
Precision (Positive Predictive Value)	$PPV = \frac{TP}{TP+FP}$
Recall (True Positive Rate)	$TPR = \frac{TP}{P} = \frac{TP}{TP+FN}$
F-measure	$F(\beta) = (1 + \beta^2) \cdot \left( \frac{PPV \cdot TPR}{\beta^2 \cdot PPV + TPR} \right)$
Receiver Operating Characteristic (ROC)	$FPR(T) = \int_T^\infty P_0(T) dT$ $TPR(T) = \int_T^\infty P_1(T) dT$
Area Under ROC Curve	$AUROC = \int_0^1 \frac{TP}{P} d\left(\frac{FP}{N}\right)$
Area Under Precision-Recall Curve	$AUPRC = \int_0^1 \frac{TP}{TP+FP} d\left(\frac{TP}{P}\right)$

Metric	Definition
Confusion Matrix	$C_{ij} = \sum_{k=0}^{N-1} \hat{\delta}(\mathbf{y}_k - \ell_i) \cdot \hat{\delta}(\hat{\mathbf{y}}_k - \ell_j)$ $\begin{pmatrix} \sum_{k=0}^{N-1} \hat{\delta}(\mathbf{y}_k - \ell_1) \cdot \hat{\delta}(\hat{\mathbf{y}}_k - \ell_1) & \dots & \sum_{k=0}^{N-1} \hat{\delta}(\mathbf{y}_k - \ell_1) \cdot \hat{\delta}(\hat{\mathbf{y}}_k - \ell_N) \\ \vdots & \ddots & \vdots \\ \sum_{k=0}^{N-1} \hat{\delta}(\mathbf{y}_k - \ell_N) \cdot \hat{\delta}(\hat{\mathbf{y}}_k - \ell_1) & \dots & \sum_{k=0}^{N-1} \hat{\delta}(\mathbf{y}_k - \ell_N) \cdot \hat{\delta}(\hat{\mathbf{y}}_k - \ell_N) \end{pmatrix}$
Accuracy	$ACC = \frac{TP}{TP+FP} = \frac{1}{N} \sum_{i=0}^{N-1} \hat{\delta}(\hat{\mathbf{y}}_i - \mathbf{y}_i)$
Precision by label	$PPV(\ell) = \frac{TP}{TP+FP} = \frac{\sum_{i=0}^{N-1} \hat{\delta}(\hat{\mathbf{y}}_i - \ell) \cdot \hat{\delta}(\mathbf{y}_i - \ell)}{\sum_{i=0}^{N-1} \hat{\delta}(\hat{\mathbf{y}}_i - \ell)}$
Recall by label	$TPR(\ell) = \frac{TP}{P} = \frac{\sum_{i=0}^{N-1} \hat{\delta}(\hat{\mathbf{y}}_i - \ell) \cdot \hat{\delta}(\mathbf{y}_i - \ell)}{\sum_{i=0}^{N-1} \hat{\delta}(\mathbf{y}_i - \ell)}$
F-measure by label	$F(\beta, \ell) = (1 + \beta^2) \cdot \left( \frac{PPV(\ell) \cdot TPR(\ell)}{\beta^2 \cdot PPV(\ell) + TPR(\ell)} \right)$
Weighted precision	$PPV_w = \frac{1}{N} \sum_{\ell \in L} PPV(\ell) \cdot \sum_{i=0}^{N-1} \hat{\delta}(\mathbf{y}_i - \ell)$
Weighted recall	$TPR_w = \frac{1}{N} \sum_{\ell \in L} TPR(\ell) \cdot \sum_{i=0}^{N-1} \hat{\delta}(\mathbf{y}_i - \ell)$
Weighted F-measure	$F_w(\beta) = \frac{1}{N} \sum_{\ell \in L} F(\beta, \ell) \cdot \sum_{i=0}^{N-1} \hat{\delta}(\mathbf{y}_i - \ell)$

$$\hat{\delta}(x) = \begin{cases} 1 & \text{if } x = 0, \\ 0 & \text{otherwise} \end{cases}$$

$$L(\mathbf{w}; \mathbf{x}, y) := \log \left( 1 + \exp \left( -y \mathbf{w}^T \mathbf{x} \right) \right)$$



$$f(z) = 1 / (1 + e^{-z})$$

cancer_class	thickness	size	shape	madh	epsize	bnuc	bchrom	nNuc	mit
0.0	5.0	1.0	1.0	1.0	2.0	1.0	3.0	1.0	1.0
0.0	5.0	4.0	4.0	5.0	7.0	10.0	3.0	2.0	1.0
0.0	3.0	1.0	1.0	1.0	2.0	2.0	3.0	1.0	1.0
0.0	6.0	8.0	8.0	1.0	3.0	4.0	3.0	7.0	1.0
0.0	4.0	1.0	1.0	3.0	2.0	1.0	3.0	1.0	1.0
1.0	8.0	10.0	10.0	8.0	7.0	10.0	9.0	7.0	1.0
0.0	1.0	1.0	1.0	1.0	2.0	10.0	3.0	1.0	1.0
0.0	2.0	1.0	2.0	1.0	2.0	1.0	3.0	1.0	1.0
0.0	2.0	1.0	1.0	1.0	2.0	1.0	1.0	1.0	5.0
0.0	4.0	2.0	1.0	1.0	2.0	1.0	2.0	1.0	1.0
0.0	1.0	1.0	1.0	1.0	1.0	1.0	3.0	1.0	1.0
0.0	2.0	1.0	1.0	1.0	2.0	1.0	2.0	1.0	1.0
1.0	5.0	3.0	3.0	3.0	2.0	3.0	4.0	4.0	1.0
0.0	1.0	1.0	1.0	1.0	2.0	3.0	3.0	1.0	1.0
1.0	8.0	7.0	5.0	10.0	7.0	9.0	5.0	5.0	4.0
1.0	7.0	4.0	6.0	4.0	6.0	1.0	4.0	3.0	1.0
0.0	4.0	1.0	1.0	1.0	2.0	1.0	2.0	1.0	1.0
0.0	4.0	1.0	1.0	1.0	2.0	1.0	3.0	1.0	1.0
1.0	10.0	7.0	7.0	6.0	4.0	10.0	4.0	1.0	2.0
0.0	6.0	1.0	1.0	1.0	2.0	1.0	3.0	1.0	1.0

only showing top 20 rows

cancer_class	thickness	size	shape	madh	epsize	bnuc	bchrom	nNuc	mit	features
0.0	5.0	1.0	1.0	1.0	2.0	1.0	3.0	1.0	1.0	[5.0,1.0,1.0,1.0,...
0.0	5.0	4.0	4.0	5.0	7.0	10.0	3.0	2.0	1.0	[5.0,4.0,4.0,5.0,...
0.0	3.0	1.0	1.0	1.0	2.0	2.0	3.0	1.0	1.0	[3.0,1.0,1.0,1.0,...
0.0	6.0	8.0	8.0	1.0	3.0	4.0	3.0	7.0	1.0	[6.0,8.0,8.0,1.0,...
0.0	4.0	1.0	1.0	3.0	2.0	1.0	3.0	1.0	1.0	[4.0,1.0,1.0,3.0,...
1.0	8.0	10.0	10.0	8.0	7.0	10.0	9.0	7.0	1.0	[8.0,10.0,10.0,8....
0.0	1.0	1.0	1.0	1.0	2.0	10.0	3.0	1.0	1.0	[1.0,1.0,1.0,1.0,...
0.0	2.0	1.0	2.0	1.0	2.0	1.0	3.0	1.0	1.0	[2.0,1.0,2.0,1.0,...
0.0	2.0	1.0	1.0	1.0	2.0	1.0	1.0	1.0	5.0	[2.0,1.0,1.0,1.0,...
0.0	4.0	2.0	1.0	1.0	2.0	1.0	2.0	1.0	1.0	[4.0,2.0,1.0,1.0,...
0.0	1.0	1.0	1.0	1.0	1.0	1.0	3.0	1.0	1.0	[1.0,1.0,1.0,1.0,...
0.0	2.0	1.0	1.0	1.0	2.0	1.0	2.0	1.0	1.0	[2.0,1.0,1.0,1.0,...
1.0	5.0	3.0	3.0	3.0	2.0	3.0	4.0	4.0	1.0	[5.0,3.0,3.0,3.0,...
0.0	1.0	1.0	1.0	1.0	2.0	3.0	3.0	1.0	1.0	[1.0,1.0,1.0,1.0,...
1.0	8.0	7.0	5.0	10.0	7.0	9.0	5.0	5.0	4.0	[8.0,7.0,5.0,10.0...
1.0	7.0	4.0	6.0	4.0	6.0	1.0	4.0	3.0	1.0	[7.0,4.0,6.0,4.0,...
0.0	4.0	1.0	1.0	1.0	2.0	1.0	2.0	1.0	1.0	[4.0,1.0,1.0,1.0,...
0.0	4.0	1.0	1.0	1.0	2.0	1.0	3.0	1.0	1.0	[4.0,1.0,1.0,1.0,...
1.0	10.0	7.0	7.0	6.0	4.0	10.0	4.0	1.0	2.0	[10.0,7.0,7.0,6.0...
0.0	6.0	1.0	1.0	1.0	2.0	1.0	3.0	1.0	1.0	[6.0,1.0,1.0,1.0,...

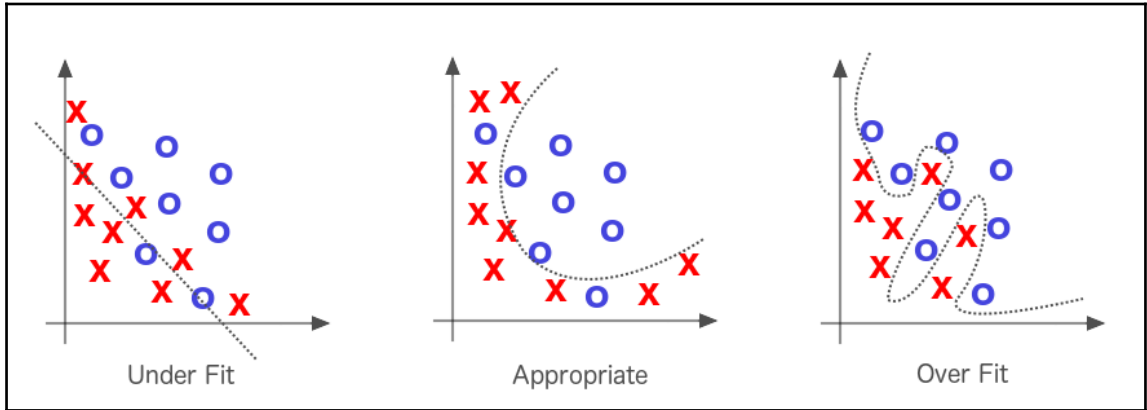
only showing top 20 rows

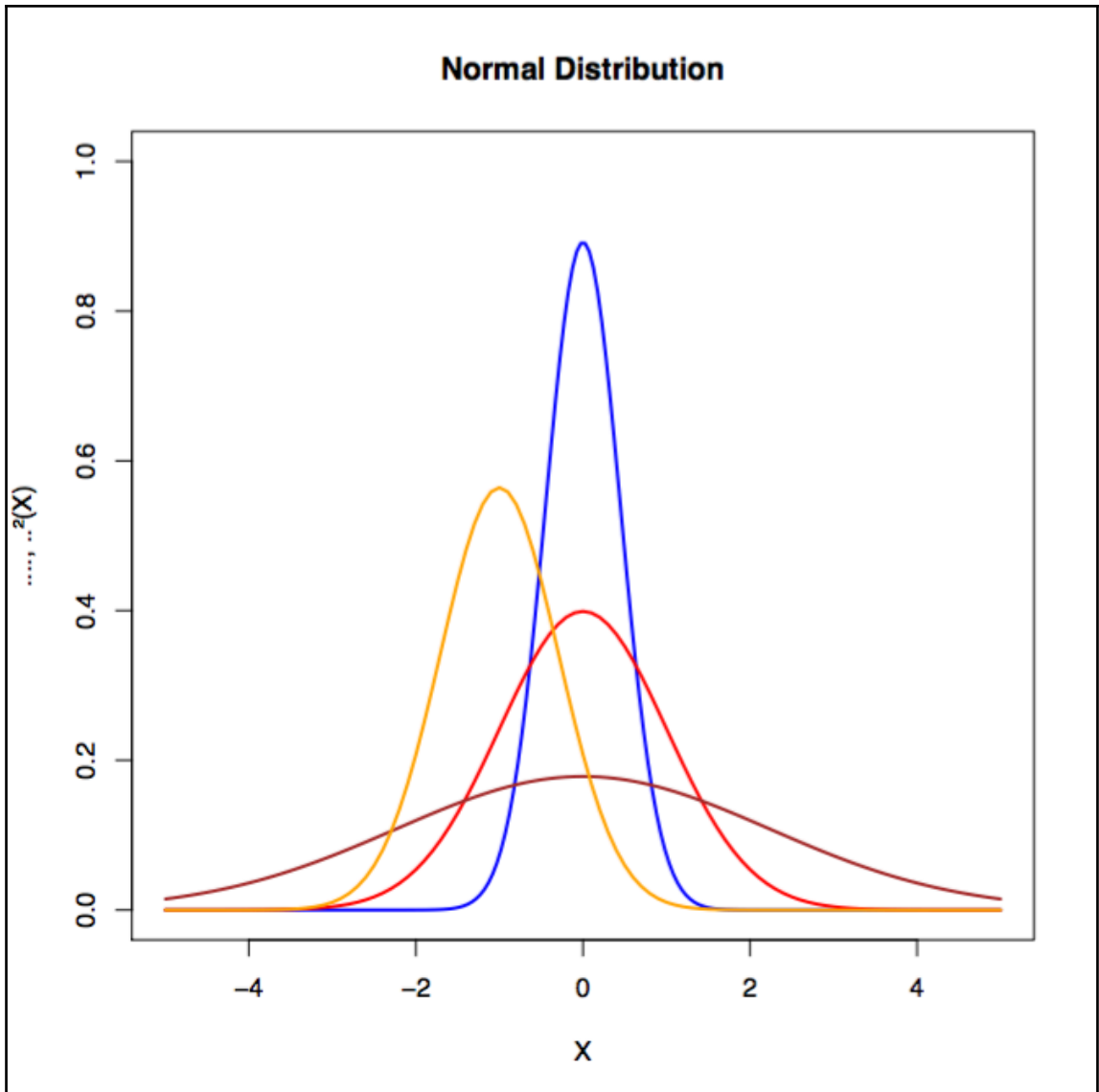
cancer_class	thickness	size	shape	madh	epsize	bnuc	bchrom	nNuc	mit	features	label
0.0	5.0	1.0	1.0	1.0	2.0	1.0	3.0	1.0	1.0	[5.0,1.0,1.0,1.0,...	0.0
0.0	5.0	4.0	4.0	5.0	7.0	10.0	3.0	2.0	1.0	[5.0,4.0,4.0,5.0,...	0.0
0.0	3.0	1.0	1.0	1.0	2.0	2.0	3.0	1.0	1.0	[3.0,1.0,1.0,1.0,...	0.0
0.0	6.0	8.0	8.0	1.0	3.0	4.0	3.0	7.0	1.0	[6.0,8.0,8.0,1.0,...	0.0
0.0	4.0	1.0	1.0	3.0	2.0	1.0	3.0	1.0	1.0	[4.0,1.0,1.0,3.0,...	0.0
1.0	8.0	10.0	10.0	8.0	7.0	10.0	9.0	7.0	1.0	[8.0,10.0,10.0,8....	1.0
0.0	1.0	1.0	1.0	1.0	2.0	10.0	3.0	1.0	1.0	[1.0,1.0,1.0,1.0,...	0.0
0.0	2.0	1.0	2.0	1.0	2.0	1.0	3.0	1.0	1.0	[2.0,1.0,2.0,1.0,...	0.0
0.0	2.0	1.0	1.0	1.0	2.0	1.0	1.0	1.0	5.0	[2.0,1.0,1.0,1.0,...	0.0
0.0	4.0	2.0	1.0	1.0	2.0	1.0	2.0	1.0	1.0	[4.0,2.0,1.0,1.0,...	0.0
0.0	1.0	1.0	1.0	1.0	1.0	1.0	3.0	1.0	1.0	[1.0,1.0,1.0,1.0,...	0.0
0.0	2.0	1.0	1.0	1.0	2.0	1.0	2.0	1.0	1.0	[2.0,1.0,1.0,1.0,...	0.0
1.0	5.0	3.0	3.0	3.0	2.0	3.0	4.0	4.0	1.0	[5.0,3.0,3.0,3.0,...	1.0
0.0	1.0	1.0	1.0	1.0	2.0	3.0	3.0	1.0	1.0	[1.0,1.0,1.0,1.0,...	0.0
1.0	8.0	7.0	5.0	10.0	7.0	9.0	5.0	5.0	4.0	[8.0,7.0,5.0,10.0...	1.0
1.0	7.0	4.0	6.0	4.0	6.0	1.0	4.0	3.0	1.0	[7.0,4.0,6.0,4.0,...	1.0
0.0	4.0	1.0	1.0	1.0	2.0	1.0	2.0	1.0	1.0	[4.0,1.0,1.0,1.0,...	0.0
0.0	4.0	1.0	1.0	1.0	2.0	1.0	3.0	1.0	1.0	[4.0,1.0,1.0,1.0,...	0.0
1.0	10.0	7.0	7.0	6.0	4.0	10.0	4.0	1.0	2.0	[10.0,7.0,7.0,6.0...	1.0
0.0	6.0	1.0	1.0	1.0	2.0	1.0	3.0	1.0	1.0	[6.0,1.0,1.0,1.0,...	0.0

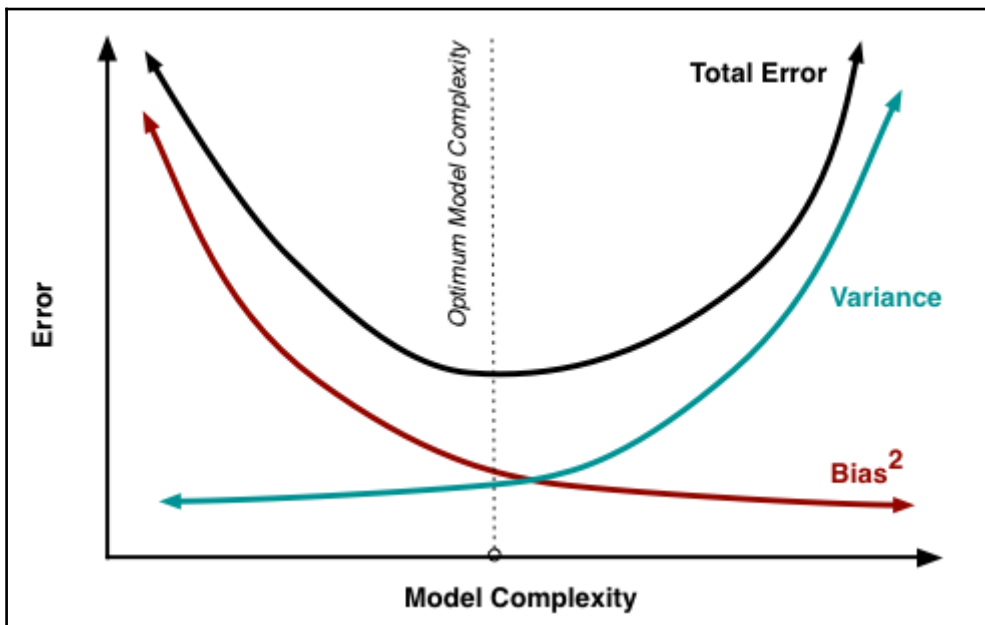
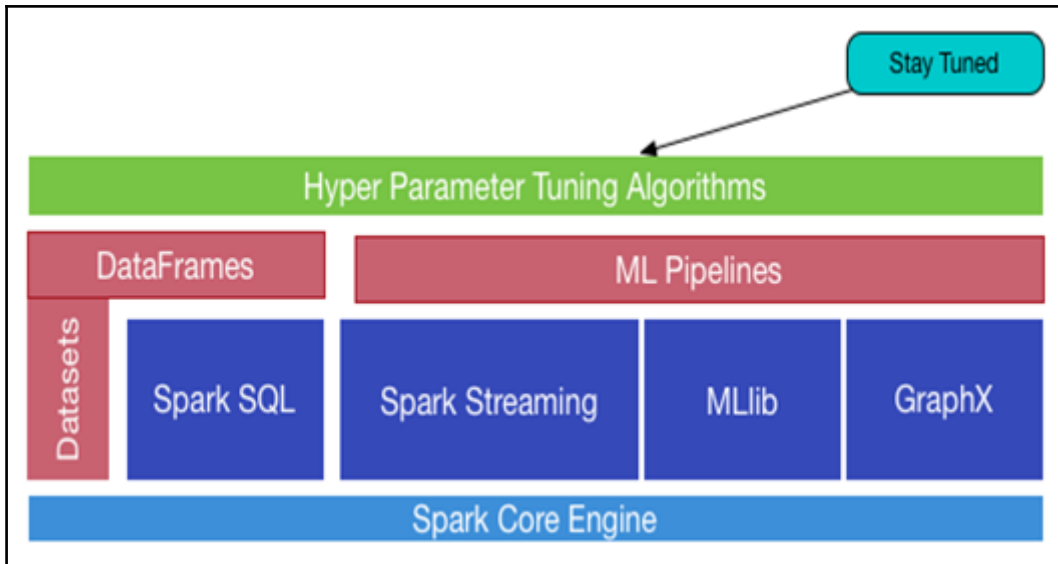
only showing top 20 rows

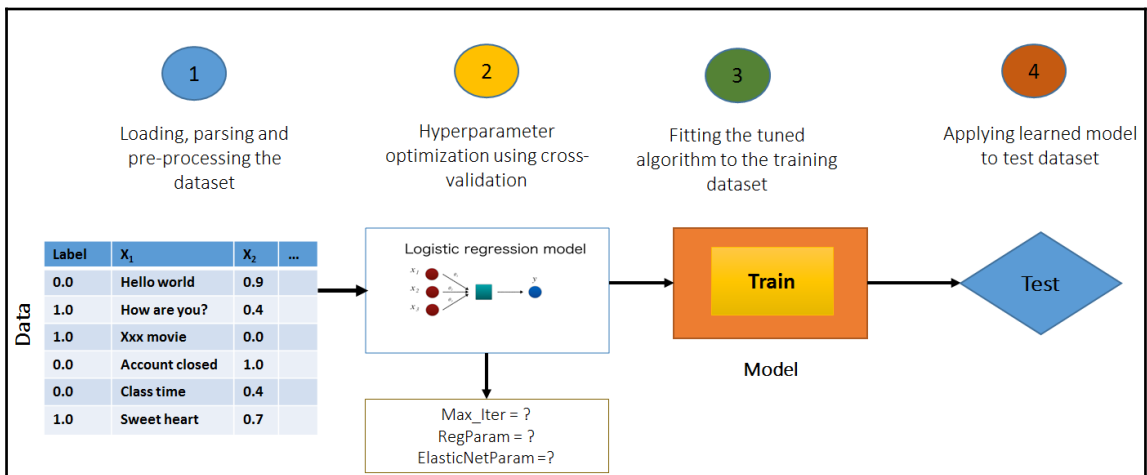
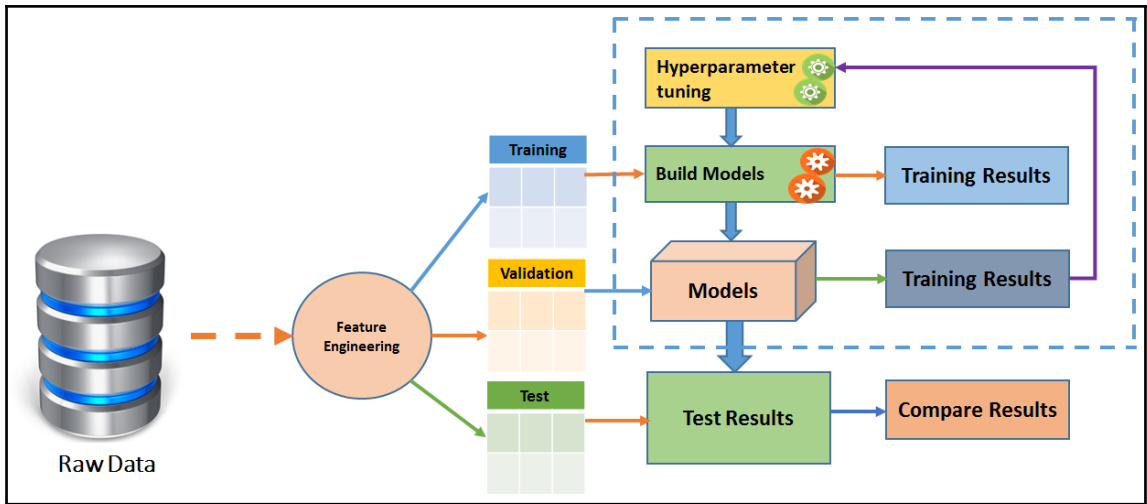


# Chapter 12: Advanced Machine Learning Best Practices









creditability	balance	duration	history	purpose	amount	savings	employment	instPercent	sex	married	guarantors	residenceDuration	assets	age	concCredit	apartment	credits	occupation	dependents	hasPhone	foreign
1.0	0.0	18.0	4.0	2.0	1049.0	0.0	1.0	4.0	1.0	0.0	3.0	1.0	21.0	2.0	0.0	0.0	2.0	0.0	0.0	0.0	
1.0	0.0	9.0	4.0	0.0	2799.0	0.0	2.0	2.0	2.0	0.0	1.0	0.0	36.0	2.0	0.0	1.0	2.0	1.0	0.0	0.0	
1.0	1.0	12.0	2.0	0.0	541.0	1.0	3.0	2.0	1.0	0.0	3.0	0.0	23.0	2.0	0.0	0.0	1.0	0.0	0.0	0.0	
1.0	0.0	12.0	4.0	0.0	2122.0	0.0	2.0	3.0	2.0	0.0	1.0	0.0	39.0	2.0	0.0	1.0	1.0	1.0	0.0	1.0	
1.0	0.0	12.0	4.0	0.0	2171.0	0.0	2.0	4.0	2.0	0.0	3.0	1.0	38.0	0.0	1.0	1.0	1.0	0.0	0.0	1.0	
1.0	0.0	10.0	4.0	0.0	2241.0	0.0	1.0	1.0	2.0	0.0	2.0	0.0	48.0	2.0	0.0	1.0	1.0	1.0	0.0	1.0	
1.0	0.0	8.0	4.0	0.0	3398.0	0.0	3.0	1.0	2.0	0.0	3.0	0.0	39.0	2.0	1.0	1.0	1.0	0.0	0.0	1.0	
1.0	0.0	6.0	4.0	0.0	1361.0	0.0	1.0	2.0	2.0	0.0	3.0	0.0	40.0	2.0	1.0	0.0	1.0	1.0	0.0	1.0	
1.0	3.0	18.0	4.0	3.0	1098.0	0.0	0.0	4.0	1.0	0.0	3.0	2.0	65.0	2.0	1.0	1.0	0.0	0.0	0.0	0.0	
1.0	1.0	24.0	2.0	3.0	3758.0	2.0	0.0	1.0	1.0	0.0	3.0	3.0	23.0	2.0	0.0	0.0	0.0	0.0	0.0	0.0	
1.0	0.0	11.0	4.0	0.0	3905.0	0.0	2.0	2.0	2.0	0.0	1.0	0.0	36.0	2.0	0.0	1.0	2.0	1.0	0.0	0.0	
1.0	0.0	30.0	4.0	1.0	6187.0	1.0	3.0	1.0	3.0	0.0	3.0	2.0	24.0	2.0	0.0	1.0	2.0	0.0	0.0	0.0	
1.0	0.0	6.0	4.0	3.0	1957.0	0.0	3.0	1.0	1.0	0.0	3.0	2.0	31.0	2.0	1.0	0.0	2.0	0.0	0.0	0.0	
1.0	1.0	48.0	3.0	10.0	7582.0	1.0	0.0	2.0	2.0	0.0	3.0	3.0	31.0	2.0	1.0	0.0	3.0	0.0	1.0	0.0	
1.0	0.0	18.0	2.0	3.0	1936.0	4.0	3.0	2.0	3.0	0.0	3.0	2.0	23.0	2.0	0.0	1.0	1.0	0.0	0.0	0.0	
1.0	0.0	6.0	2.0	3.0	2647.0	2.0	2.0	2.0	2.0	0.0	2.0	0.0	44.0	2.0	0.0	0.0	2.0	1.0	0.0	0.0	
1.0	0.0	11.0	4.0	0.0	3939.0	0.0	2.0	1.0	2.0	0.0	1.0	0.0	40.0	2.0	1.0	1.0	1.0	1.0	0.0	0.0	
1.0	1.0	18.0	2.0	3.0	3213.0	2.0	1.0	1.0	3.0	0.0	2.0	0.0	25.0	2.0	0.0	0.0	2.0	0.0	0.0	0.0	
1.0	1.0	36.0	4.0	3.0	2337.0	0.0	4.0	4.0	2.0	0.0	3.0	0.0	36.0	2.0	1.0	0.0	2.0	0.0	0.0	0.0	
1.0	3.0	11.0	4.0	0.0	7228.0	0.0	2.0	1.0	2.0	0.0	3.0	1.0	39.0	2.0	1.0	1.0	1.0	0.0	0.0	0.0	

only showing top 20 rows

creditability	avgbalance	avgamt	avgdur
0.0	0.9033333333333333	3938.1266666666666	24.86
1.0	1.8657142857142857	2985.442857142857	19.207142857142856

summary	balance
count	1000
mean	1.577
stddev	1.257637727110893
min	0.0
max	3.0

creditability	avg(balance)
0.0	0.9033333333333333
1.0	1.8657142857142857

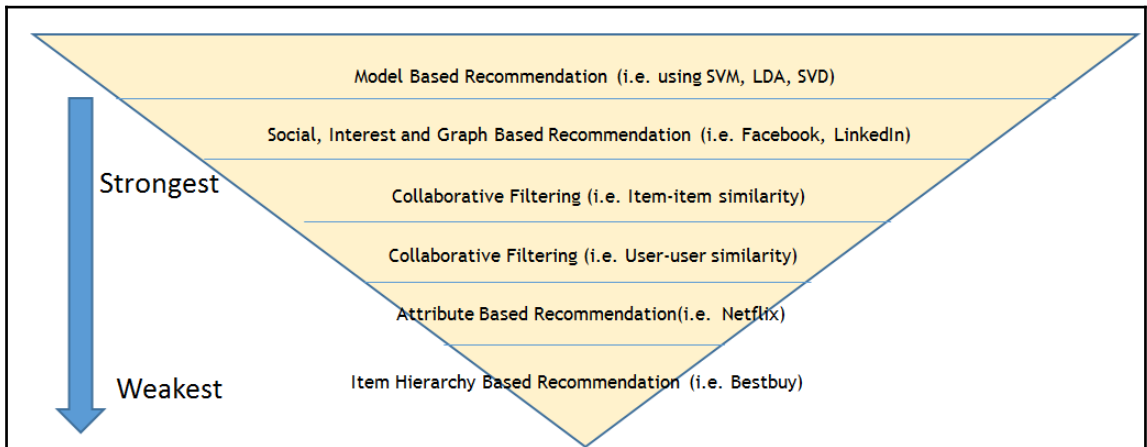


```
+-----+
|           features |
+-----+
| (20, [1, 2, 3, 4, 6, 7, ... |
| (20, [1, 2, 4, 6, 7, 8, ... |
| [1.0, 12.0, 2.0, 9.0... |
| [0.0, 12.0, 4.0, 0.0... |
| [0.0, 12.0, 4.0, 0.0... |
| [0.0, 10.0, 4.0, 0.0... |
| [0.0, 8.0, 4.0, 0.0, ... |
| [0.0, 6.0, 4.0, 0.0, ... |
| [3.0, 18.0, 4.0, 3.0... |
| (20, [0, 1, 2, 3, 4, 5, ... |
| (20, [1, 2, 4, 6, 7, 8, ... |
| [0.0, 30.0, 4.0, 1.0... |
| [0.0, 6.0, 4.0, 3.0, ... |
| [1.0, 48.0, 3.0, 10.... |
| [0.0, 18.0, 2.0, 3.0... |
| [0.0, 6.0, 2.0, 3.0, ... |
| [0.0, 11.0, 4.0, 0.0... |
| [1.0, 18.0, 2.0, 3.0... |
| [1.0, 36.0, 4.0, 3.0... |
| [3.0, 11.0, 4.0, 0.0... |
+-----+
only showing top 20 rows
```

label	features
0.0	(20, [1, 2, 3, 4, 6, 7, ...
0.0	(20, [1, 2, 4, 6, 7, 8, ...
0.0	[1.0, 12.0, 2.0, 9.0...]
0.0	[0.0, 12.0, 4.0, 0.0...]
0.0	[0.0, 12.0, 4.0, 0.0...]
0.0	[0.0, 10.0, 4.0, 0.0...]
0.0	[0.0, 8.0, 4.0, 0.0, ...]
0.0	[0.0, 6.0, 4.0, 0.0, ...]
0.0	[3.0, 18.0, 4.0, 3.0...]
0.0	(20, [0, 1, 2, 3, 4, 5, ...
0.0	(20, [1, 2, 4, 6, 7, 8, ...
0.0	[0.0, 30.0, 4.0, 1.0...]
0.0	[0.0, 6.0, 4.0, 3.0, ...]
0.0	[1.0, 48.0, 3.0, 10....]
0.0	[0.0, 18.0, 2.0, 3.0...]
0.0	[0.0, 6.0, 2.0, 3.0, ...]
0.0	[0.0, 11.0, 4.0, 0.0...]
0.0	[1.0, 18.0, 2.0, 3.0...]
0.0	[1.0, 36.0, 4.0, 3.0...]
0.0	[3.0, 11.0, 4.0, 0.0...]

only showing top 20 rows

label	rawPrediction	probability	prediction
1.0	[21.0,9.0]	[0.7,0.3]	0.0
0.0	[28.9868421052631...]	[0.96622807017543...]	0.0
0.0	[18.0,12.0]	[0.6,0.4]	0.0
0.0	[23.9873417721519...]	[0.79957805907173...]	0.0
0.0	[24.6540084388185...]	[0.82180028129395...]	0.0
0.0	[22.9868421052631...]	[0.76622807017543...]	0.0
0.0	[14.5952380952380...]	[0.48650793650793...]	1.0
0.0	[17.9547224224945...]	[0.59849074741648...]	0.0
0.0	[23.9684210526315...]	[0.79894736842105...]	0.0
0.0	[25.0,5.0]	[0.83333333333333...]	0.0
0.0	[15.5,14.5]	[0.51666666666666...]	0.0
0.0	[22.5,7.5]	[0.75,0.25]	0.0
0.0	[22.9486422749787...]	[0.76495474249929...]	0.0
0.0	[18.0,12.0]	[0.6,0.4]	0.0
0.0	[27.9631948664260...]	[0.93210649554753...]	0.0
0.0	[21.0,9.0]	[0.7,0.3]	0.0
0.0	[24.0,6.0]	[0.8,0.2]	0.0
0.0	[16.0,14.0]	[0.53333333333333...]	0.0
0.0	[23.9921259842519...]	[0.79973753280839...]	0.0
0.0	[14.9890109890109...]	[0.49963369963369...]	1.0



userId	movieId	rating	timestamp
1	16	4.0	1217897793
1	24	1.5	1217895807
1	32	4.0	1217896246
1	47	4.0	1217896556
1	50	4.0	1217896523
1	110	4.0	1217896150
1	150	3.0	1217895940
1	161	4.0	1217897864
1	165	3.0	1217897135
1	204	0.5	1217895786
1	223	4.0	1217897795
1	256	0.5	1217895764
1	260	4.5	1217895864
1	261	1.5	1217895750
1	277	0.5	1217895772
1	296	4.0	1217896125
1	318	4.0	1217895860
1	349	4.5	1217897058
1	356	3.0	1217896231
1	377	2.5	1217896373

only showing top 20 rows

movieId	title	genres
1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy
2	Jumanji (1995)	Adventure Children Fantasy
3	Grumpier Old Men (1995)	Comedy Romance
4	Waiting to Exhale (1995)	Comedy Drama Romance
5	Father of the Bride Part II (1995)	Comedy
6	Heat (1995)	Action Crime Thriller
7	Sabrina (1995)	Comedy Romance
8	Tom and Huck (1995)	Adventure Children
9	Sudden Death (1995)	Action
10	GoldenEye (1995)	Action Adventure Thriller
11	American President, The (1995)	Comedy Drama Romance
12	Dracula: Dead and Loving It (1995)	Comedy Horror
13	Balto (1995)	Adventure Animation Children
14	Nixon (1995)	Drama
15	Cutthroat Island (1995)	Action Adventure Romance
16	Casino (1995)	Crime Drama
17	Sense and Sensibility (1995)	Drama Romance
18	Four Rooms (1995)	Comedy
19	Ace Ventura: When Nature Calls (1995)	Comedy
20	Money Train (1995)	Action Comedy Crime Drama Thriller

only showing top 20 rows

Got 105339 ratings from 668 users on 10325 movies.

title	maxr	minr	cntu
Pulp Fiction (1994)	5.0	0.5	325
Forrest Gump (1994)	5.0	0.5	311
Shawshank Redemption, The (1994)	5.0	0.5	308
Jurassic Park (1993)	5.0	1.0	294
Silence of the Lambs, The (1991)	5.0	0.5	290
Star Wars: Episode IV - A New Hope (1977)	5.0	0.5	273
Matrix, The (1999)	5.0	0.5	261
Terminator 2: Judgment Day (1991)	5.0	0.5	253
Braveheart (1995)	5.0	0.5	248
Schindler's List (1993)	5.0	0.5	248
Fugitive, The (1993)	5.0	1.0	244
Toy Story (1995)	5.0	1.0	232
Star Wars: Episode V - The Empire Strikes Back (1980)	5.0	0.5	228
Usual Suspects, The (1995)	5.0	1.0	228
Raiders of the Lost Ark (Indiana Jones and the Raiders of the Lost Ark) (1981)	5.0	1.0	224
Star Wars: Episode VI - Return of the Jedi (1983)	5.0	0.5	222
Batman (1989)	5.0	0.5	217
American Beauty (1999)	5.0	1.0	216
Back to the Future (1985)	5.0	1.5	213
Godfather, The (1972)	5.0	1.0	210

only showing top 20 rows

userId	ct
668	5678
575	2837
458	2086
232	1421
310	1287
475	1249
128	1231
224	1182
607	1176
63	1107

userId	movieId	rating	title
668	6	5.0	Heat (1995)
668	326	4.5	To Live (Huozhe) (1994)
668	446	4.5	Farewell My Concubine (Ba wang bie ji) (1993)
668	515	4.5	Remains of the Day, The (1993)
668	593	4.5	Silence of the Lambs, The (1991)
668	594	4.5	Snow White and the Seven Dwarfs (1937)
668	608	4.5	Fargo (1996)
668	858	5.0	Godfather, The (1972)
668	898	5.0	Philadelphia Story, The (1940)
668	907	4.5	Gay Divorcee, The (1934)
668	908	5.0	North by Northwest (1959)
668	910	5.0	Some Like It Hot (1959)
668	912	5.0	Casablanca (1942)
668	913	5.0	Maltese Falcon, The (1941)
668	914	5.0	My Fair Lady (1964)
668	919	5.0	Wizard of Oz, The (1939)
668	927	4.5	Women, The (1939)
668	930	4.5	Notorious (1946)
668	945	4.5	Top Hat (1935)
668	947	4.5	My Man Godfrey (1936)

only showing top 20 rows

```
Rating:(UserID, MovieID, Rating)
-----
Rating(668,101862,4.8525842164777435)
Rating(668,5304,4.8525842164777435)
Rating(668,25961,4.8525842164777435)
Rating(668,80969,4.779325934293423)
Rating(668,93040,4.7528736838886)
Rating(668,25795,4.676957397667861)
-----
(Prediction, Rating)
(3.848087516442212,3.5)
(4.647813269020743,5.0)
(3.578002886107389,4.0)
(3.681217214985231,3.0)
(2.844685318141285,3.0)
```

Topic: 0			Topic: 1		
Terms	Index	Weight	Terms	Index	Weight
space	10665	0.046582	smile	10668	0.129227
just	10667	0.034397	just	10667	0.024922
posted	10637	0.016093	good	10663	0.022404
love	10661	0.015652	hope	10645	0.017981
photo	10639	0.013296	going	10655	0.015764
cosmic	10635	0.013212	thanks	10648	0.014945
angry	10656	0.012860	time	10662	0.014941
like	10666	0.012629	like	10666	0.014827
life	10640	0.012107	think	10659	0.014438
time	10662	0.011634	work	10649	0.012702
-----			-----		
Sum:= 0.188459750219041			Sum:= 0.28215004471848354		
Topic: 2			Topic: 3		
Terms	Index	Weight	Terms	Index	Weight
grin	10664	0.078958	like	10666	0.030890
yang	10628	0.029173	just	10667	0.020093
kita	10574	0.017318	know	10660	0.016473
disgust		10618 0.016325	good	10663	0.013343
udah	10544	0.014584	that	10651	0.012687
science		10590 0.012792	people	10658	0.012137
space	10665	0.011765	right	10654	0.012097
nggak	10501	0.011290	think	10659	0.011395
kalo	10476	0.010203	love	10661	0.010943
angry	10656	0.009313	does	10646	0.009002
-----			-----		
Sum:= 0.21172148557919923			Sum:= 0.14905966677477597		



```
+-----+
|                docs|
+-----+
|20 000 2000010 th...|
|fifty two still ...|
|please this loo...|
|gable this      ...|
|Sheridan world ...|
|  wiretap with Ca...|
|                ...|
|find plea Lamb h...|
|  find 10th aes...|
|Disney alad10 tx...|
|Alcott Gutenber...|
|  this hand dom...|
|please this loo...|
|  HORATIO breaker...|
|  HORATIO ragged ...|
|  HORATIO upward s...|
|RESEARCH Electro...|
|tradition Richard...|
+-----+
```

# Chapter 13: My Name is Bayes, Naive Bayes

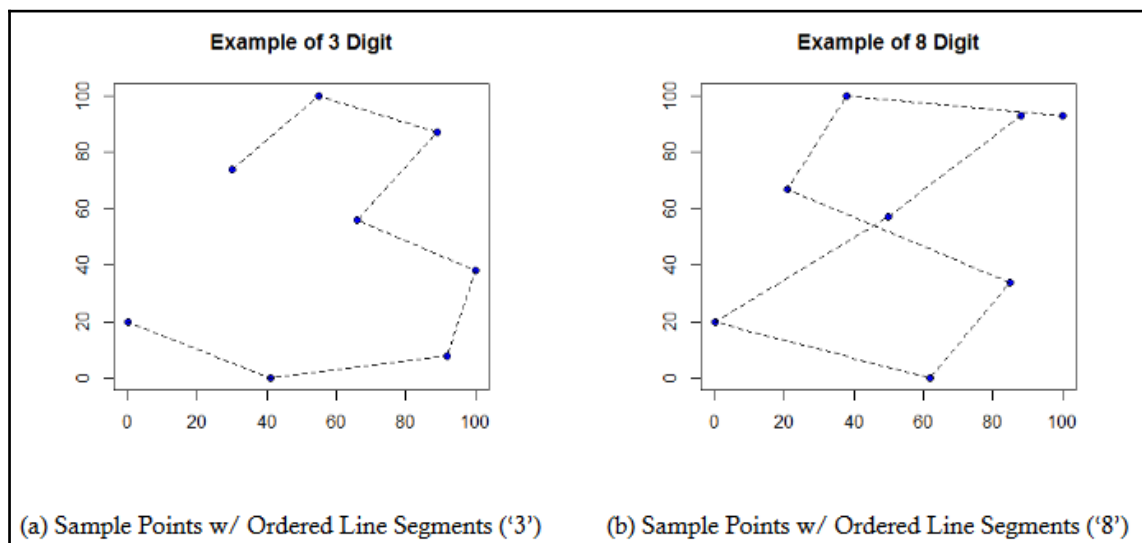


letter	xbox	ybox	width	height	onpix	xbar	ybar	x2bar	y2bar	xybar	x2ybar	xy2bar	xedge	xedgey	yedge	yedgex
T	2	8	3	5	1	8	13	0	6	6	10	8	0	8	0	8
I	5	12	3	7	2	10	5	5	4	13	3	9	2	8	4	10
D	4	11	6	8	6	10	6	2	6	10	3	7	3	7	3	9
N	7	11	6	6	3	5	9	4	6	4	4	10	6	10	2	8
G	2	1	3	1	1	8	6	6	6	6	5	9	1	7	5	10
S	4	11	5	8	3	8	8	6	9	5	6	6	0	8	9	7
B	4	2	5	4	4	8	7	6	6	7	6	6	2	8	7	10
A	1	1	3	2	1	8	2	2	2	8	2	8	1	6	2	7
J	2	2	4	4	2	10	6	2	6	12	4	8	1	6	1	7
M	11	15	13	9	7	13	2	6	2	12	1	9	8	1	1	8
X	3	9	5	7	4	8	7	3	8	5	6	8	2	8	6	7
O	6	13	4	7	4	6	7	6	3	10	7	9	5	9	5	8
G	4	9	6	7	6	7	8	6	2	6	5	11	4	8	7	8
M	6	9	8	6	9	7	8	6	5	7	5	8	8	9	8	6
R	5	9	5	7	6	6	11	7	3	7	3	9	2	7	5	11
F	6	9	5	4	3	10	6	3	5	10	5	7	3	9	6	9
O	3	4	4	3	2	8	7	7	5	7	6	8	2	8	3	8
C	7	10	5	5	2	6	8	6	8	11	7	11	2	8	5	9
T	6	11	6	8	5	6	11	5	6	11	9	4	3	12	2	4
J	2	2	3	3	1	10	6	3	6	12	4	9	0	7	1	7

only showing top 20 rows

label	features
8.0	(17, [0,1,2,3,4,5,...])
10.0	(17, [0,1,2,3,4,5,...])
9.0	(17, [0,1,2,3,4,5,...])
8.0	(17, [0,1,2,3,4,5,...])
10.0	(17, [0,1,2,3,4,5,...])
8.0	(17, [0,1,2,3,4,5,...])
5.0	(17, [0,1,2,3,4,5,...])
6.0	(17, [0,1,2,3,4,5,...])
8.0	(17, [0,1,2,3,4,5,...])
7.0	(17, [0,1,2,3,4,5,...])
6.0	(17, [0,1,2,3,4,5,...])
8.0	(17, [0,1,2,3,4,5,...])
8.0	(17, [0,1,2,3,4,5,...])
8.0	(17, [0,1,2,3,4,5,...])
9.0	(17, [0,1,2,3,4,5,...])
4.0	(17, [0,1,2,3,4,5,...])
7.0	(17, [0,1,2,3,4,5,...])
7.0	(17, [0,1,2,3,4,5,...])
8.0	(17, [0,1,2,3,4,5,...])
8.0	(17, [0,1,2,3,4,5,...])

only showing top 20 rows



label	features
8.0	(16, [0, 1, 2, 3, 4, 5, ...]
2.0	(16, [1, 2, 3, 4, 5, 6, ...]
1.0	(16, [1, 2, 3, 4, 5, 6, ...]
4.0	(16, [1, 2, 3, 4, 5, 6, ...]
1.0	(16, [1, 2, 3, 4, 5, 6, ...]
6.0	(16, [0, 1, 2, 3, 4, 5, ...]
4.0	(16, [1, 2, 3, 4, 5, 6, ...]
0.0	(16, [1, 2, 3, 4, 5, 6, ...]
5.0	(16, [0, 1, 2, 3, 4, 5, ...]
0.0	(16, [0, 1, 2, 3, 5, 6, ...]
9.0	(16, [0, 1, 2, 3, 5, 6, ...]
8.0	(16, [0, 1, 2, 3, 4, 5, ...]
5.0	(16, [0, 1, 2, 3, 4, 5, ...]
9.0	(16, [0, 1, 2, 3, 5, 6, ...]
7.0	(16, [1, 2, 3, 4, 5, 6, ...]
3.0	(16, [0, 1, 2, 3, 4, 5, ...]
3.0	(16, [0, 1, 2, 3, 4, 5, ...]
9.0	(16, [0, 1, 2, 3, 4, 5, ...]
2.0	(16, [0, 1, 2, 3, 4, 5, ...]
2.0	(16, [1, 2, 3, 4, 5, 6, ...]

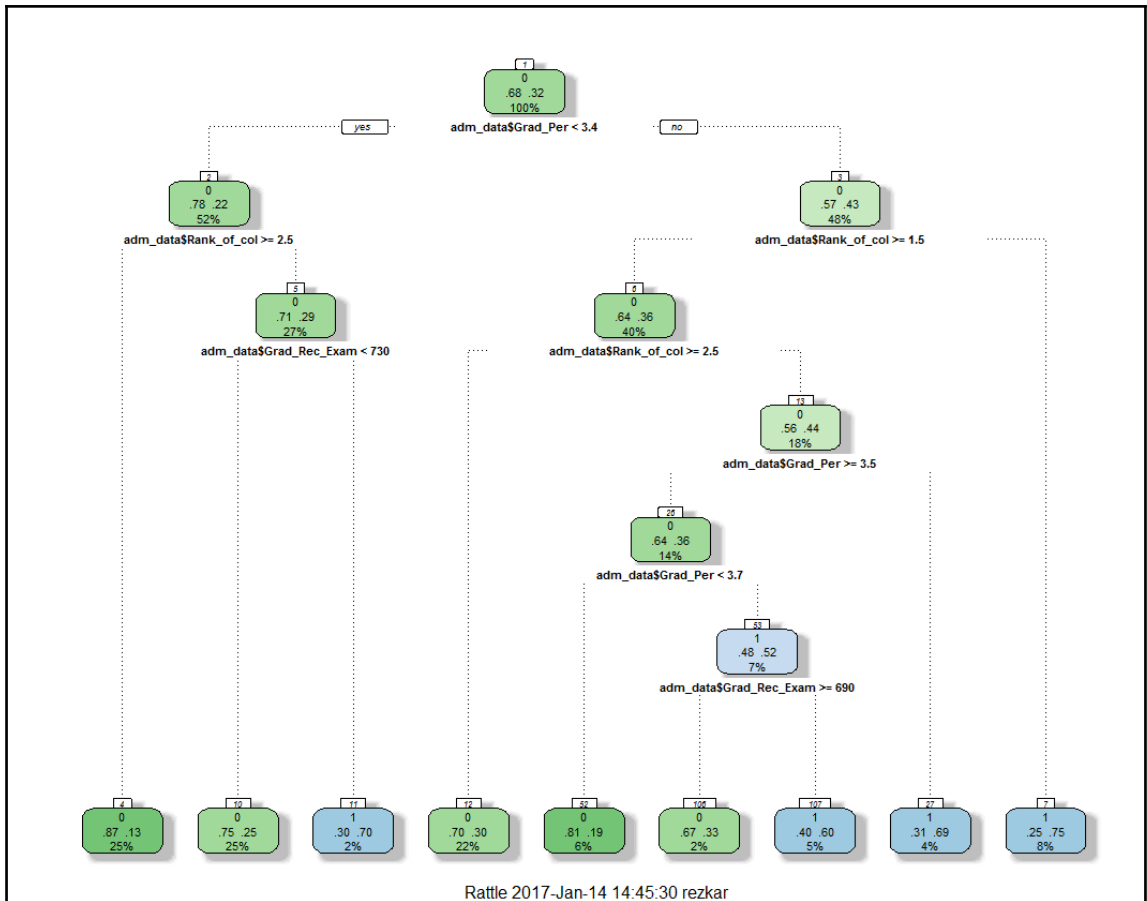
only showing top 20 rows

label	features	rawPrediction	probability	prediction
0.0	(16, [0,1,2,3,4,5,...	[-2439.0893277449...	[1.32132340702018...	4.0
0.0	(16, [0,1,2,3,4,5,...	[-1941.7868705353...	[1.0,1.5395790656...	0.0
0.0	(16, [0,1,2,3,4,5,...	[-2024.4356335162...	[1.0,1.6764090944...	0.0
0.0	(16, [0,1,2,3,4,5,...	[-1989.5775697073...	[1.0,2.2647494021...	0.0
0.0	(16, [0,1,2,3,4,5,...	[-1706.6857288506...	[1.0,5.1940219699...	0.0
0.0	(16, [0,1,2,3,4,5,...	[-1838.2628605334...	[1.0,7.2364926581...	0.0
0.0	(16, [0,1,2,3,4,5,...	[-2168.4931444350...	[1.0,6.8428584454...	0.0
0.0	(16, [0,1,2,3,4,5,...	[-2068.2067411172...	[1.0,1.1943331620...	0.0
0.0	(16, [0,1,2,3,4,5,...	[-2132.6929489447...	[1.0,1.9943684266...	0.0
0.0	(16, [0,1,2,3,4,5,...	[-1983.0451148771...	[1.0,4.9959906892...	0.0
0.0	(16, [0,1,2,3,4,5,...	[-2049.2850893323...	[1.0,1.3644883115...	0.0
0.0	(16, [0,1,2,3,4,5,...	[-1971.1755138520...	[1.0,1.6415723270...	0.0
0.0	(16, [0,1,2,3,4,5,...	[-2216.9188759036...	[1.0,1.3805417667...	0.0
0.0	(16, [0,1,2,3,4,5,...	[-2216.0583349043...	[1.0,7.7430733808...	0.0
0.0	(16, [0,1,2,3,4,5,...	[-2290.1517462265...	[1.0,1.3312677171...	0.0
0.0	(16, [0,1,2,3,4,5,...	[-2268.9492946577...	[0.01491770995335...	6.0
0.0	(16, [0,1,2,3,4,5,...	[-2377.8867352336...	[1.27336913041488...	8.0
0.0	(16, [0,1,2,3,4,5,...	[-2206.2037445466...	[1.20068275169939...	6.0
0.0	(16, [0,1,2,3,4,5,...	[-2290.1662968738...	[2.82560057752915...	8.0
0.0	(16, [0,1,2,3,4,5,...	[-2662.3029788480...	[2.38039426503477...	8.0

only showing top 20 rows

label	features
0.0	(8287348, [592160, ...
0.0	(8287348, [592137, ...
0.0	(8287348, [592137, ...
0.0	(8287348, [598864, ...
0.0	(8287348, [670767, ...
0.0	(8287348, [592137, ...
0.0	(8287348, [592137, ...
1.0	(8287348, [657980, ...
0.0	(8287348, [592208, ...
0.0	(8287348, [663584, ...
1.0	(8287348, [592137, ...
1.0	(8287348, [592137, ...
0.0	(8287348, [592137, ...
1.0	(8287348, [592137, ...
1.0	(8287348, [657930, ...
0.0	(8287348, [670767, ...
1.0	(8287348, [598111, ...
0.0	(8287348, [592137, ...
0.0	(8287348, [592188, ...
0.0	(8287348, [592137, ...

only showing top 20 rows



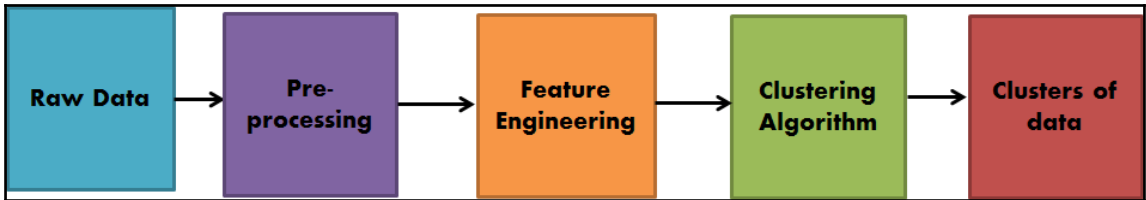
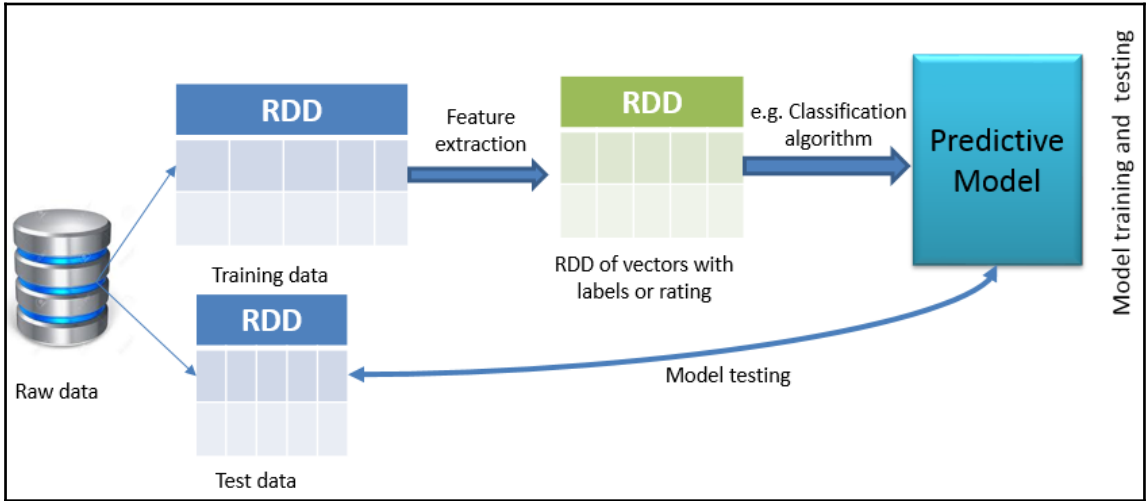
Rattle 2017-Jan-14 14:45:30 rezkar





```
Learned classification tree model:  
DecisionTreeClassificationModel (uid=dtc_fbc6a27aa70b) of depth 5 with 19 nodes  
If (feature 16 <= 7.0)  
  If (feature 16 <= 6.0)  
    If (feature 16 <= 5.0)  
      If (feature 16 <= 4.0)  
        If (feature 16 <= 3.0)  
          Predict: 9.0  
        Else (feature 16 > 3.0)  
          Predict: 7.0  
      Else (feature 16 > 4.0)  
        Predict: 5.0  
    Else (feature 16 > 5.0)  
      Predict: 3.0  
  Else (feature 16 > 6.0)  
    Predict: 1.0  
Else (feature 16 > 7.0)  
  If (feature 16 <= 8.0)  
    Predict: 0.0  
  Else (feature 16 > 8.0)  
    If (feature 16 <= 9.0)  
      Predict: 2.0  
    Else (feature 16 > 9.0)  
      If (feature 16 <= 10.0)  
        Predict: 4.0  
      Else (feature 16 > 10.0)  
        If (feature 16 <= 11.0)  
          Predict: 6.0  
        Else (feature 16 > 11.0)  
          Predict: 8.0
```

# Chapter 14: Time to Put Some Order - Cluster Your Data with Spark MLlib



Price	LotSize	Waterfront	Age	LandValue	NewConstruct	CentralAir	FuelType	HeatType	SewerType	LivingArea	PctCollege	Bedrooms	Fireplaces	Bathrooms	rooms
132500.0	0.09	0.0	42.0	50000.0	0.0	0.0	3.0	4.0	2.0	906.0	35.0	2.0	1.0	1.0	5.0
181115.0	0.92	0.0	0.0	22300.0	0.0	0.0	2.0	3.0	2.0	1953.0	51.0	3.0	0.0	2.5	6.0
109000.0	0.19	0.0	133.0	7300.0	0.0	0.0	2.0	3.0	3.0	1944.0	51.0	4.0	1.0	1.0	8.0
155000.0	0.41	0.0	13.0	18700.0	0.0	0.0	2.0	2.0	2.0	1944.0	51.0	3.0	1.0	1.5	5.0
86060.0	0.11	0.0	0.0	15000.0	1.0	1.0	2.0	2.0	3.0	840.0	51.0	2.0	0.0	1.0	3.0
120000.0	0.68	0.0	31.0	14000.0	0.0	0.0	2.0	2.0	2.0	1152.0	22.0	4.0	1.0	1.0	8.0
153000.0	0.4	0.0	33.0	23300.0	0.0	0.0	4.0	3.0	2.0	2752.0	51.0	4.0	1.0	1.5	8.0
170000.0	1.21	0.0	23.0	14600.0	0.0	0.0	4.0	2.0	2.0	1662.0	35.0	4.0	1.0	1.5	9.0
90000.0	0.83	0.0	36.0	22200.0	0.0	0.0	3.0	4.0	2.0	1632.0	51.0	3.0	0.0	1.5	8.0
122900.0	1.94	0.0	4.0	21200.0	0.0	0.0	2.0	2.0	1.0	1416.0	44.0	3.0	0.0	1.5	6.0
325000.0	2.29	0.0	123.0	12600.0	0.0	0.0	4.0	2.0	2.0	2894.0	51.0	7.0	0.0	1.0	12.0
120000.0	0.92	0.0	1.0	22300.0	0.0	0.0	2.0	2.0	2.0	1624.0	51.0	3.0	0.0	2.0	6.0
85860.0	8.97	0.0	13.0	4800.0	0.0	0.0	3.0	4.0	2.0	704.0	41.0	2.0	0.0	1.0	4.0
97000.0	0.11	0.0	153.0	3100.0	0.0	0.0	2.0	3.0	3.0	1383.0	57.0	3.0	0.0	2.0	5.0
127000.0	0.14	0.0	9.0	300.0	0.0	0.0	4.0	2.0	2.0	1300.0	41.0	3.0	0.0	1.5	8.0
89900.0	0.0	0.0	88.0	2500.0	0.0	0.0	2.0	3.0	3.0	936.0	57.0	3.0	0.0	1.0	4.0
155000.0	0.13	0.0	9.0	300.0	0.0	0.0	4.0	2.0	2.0	1300.0	41.0	3.0	0.0	1.5	7.0
253750.0	2.0	0.0	0.0	49800.0	0.0	1.0	2.0	2.0	1.0	2816.0	71.0	4.0	1.0	2.5	12.0
60000.0	0.21	0.0	82.0	8500.0	0.0	0.0	4.0	3.0	2.0	924.0	35.0	2.0	0.0	1.0	6.0
87500.0	0.88	0.0	17.0	19400.0	0.0	0.0	4.0	2.0	2.0	1092.0	35.0	3.0	0.0	1.0	6.0

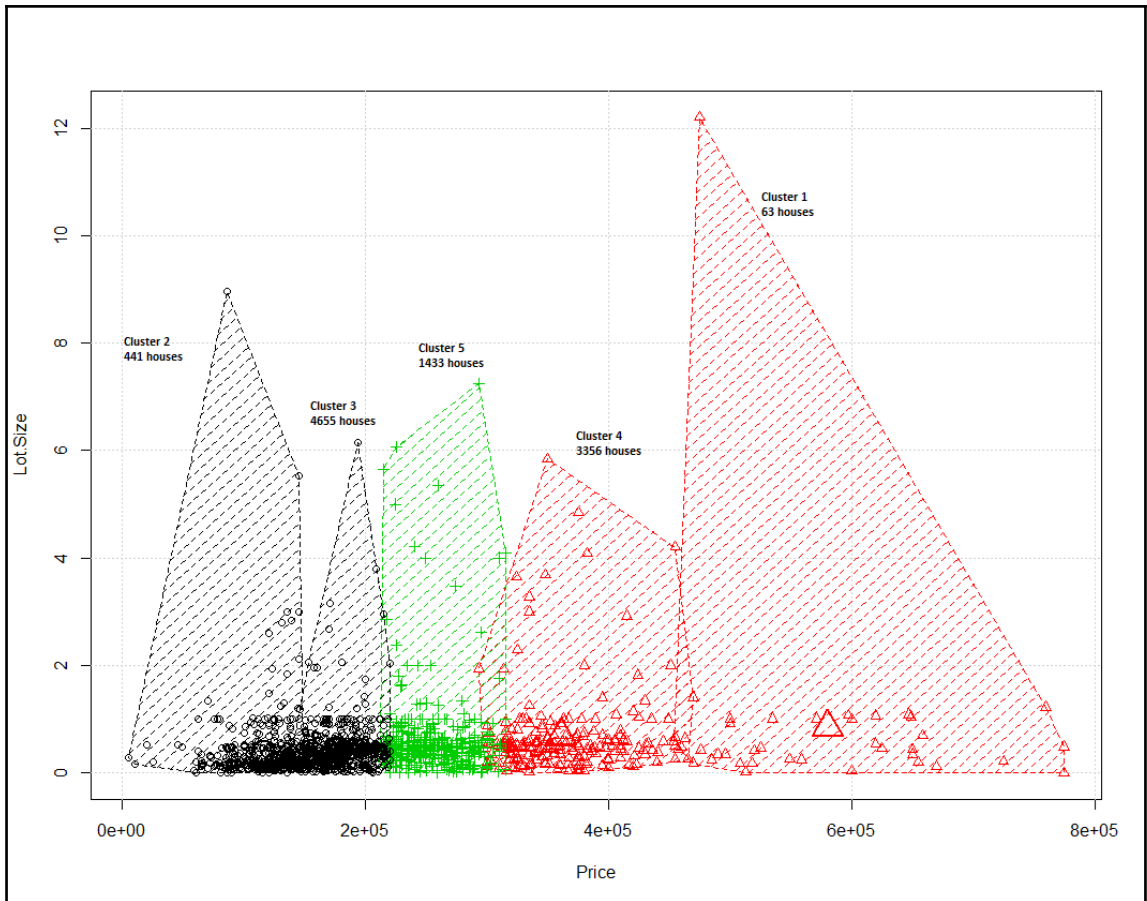
only showing top 20 rows

Price	CLUSTER
132500.0	4
181115.0	3
109000.0	0
155000.0	0
86060.0	0
120000.0	0
153000.0	3
170000.0	0
90000.0	3
122900.0	0
325000.0	0
120000.0	3
85860.0	0
97000.0	0
127000.0	0
89900.0	0
155000.0	0
253750.0	4
60000.0	0
87500.0	0

only showing top 20 rows

Price	LotSize	Waterfront	Age	LandValue	NewConstruct	CentralAir	FuelType	HeatType	SewerType	LivingArea	PctCollege	Bedrooms	Fireplaces	Bathrooms	rooms	CLUSTER
132500.0	0.21	0.0	77.0	3500.0	0.0	0.0	2.0	2.0	3.0	1379.0	36.0	3.0	0.0	1.0	7.0	4
132500.0	0.37	0.0	19.0	13000.0	0.0	0.0	3.0	4.0	3.0	1988.0	63.0	2.0	0.0	1.0	5.0	4
132500.0	0.37	0.0	19.0	13000.0	0.0	0.0	3.0	4.0	3.0	1988.0	63.0	2.0	0.0	1.0	4.0	4
132500.0	0.09	0.0	42.0	50000.0	0.0	0.0	3.0	4.0	2.0	906.0	35.0	2.0	1.0	1.0	5.0	4
253750.0	2.0	0.0	0.0	49800.0	0.0	1.0	2.0	2.0	1.0	2816.0	71.0	4.0	1.0	2.5	12.0	4
290000.0	0.66	0.0	15.0	31200.0	0.0	1.0	2.0	2.0	2.0	2305.0	51.0	4.0	1.0	2.5	11.0	4
290000.0	0.46	0.0	22.0	48000.0	0.0	1.0	2.0	2.0	3.0	2030.0	64.0	4.0	1.0	2.5	10.0	4
290000.0	0.61	0.0	34.0	32300.0	0.0	0.0	2.0	3.0	3.0	2728.0	64.0	4.0	1.0	2.5	10.0	4
290000.0	0.12	0.0	3.0	108300.0	0.0	1.0	2.0	2.0	3.0	1620.0	57.0	3.0	1.0	2.5	7.0	4
290000.0	1.0	1.0	33.0	21700.0	0.0	0.0	4.0	2.0	2.0	944.0	27.0	1.0	1.0	1.0	4.0	4
290000.0	0.15	0.0	13.0	400.0	0.0	1.0	2.0	2.0	3.0	1758.0	47.0	2.0	1.0	2.5	6.0	4
290000.0	0.51	0.0	7.0	39100.0	0.0	0.0	2.0	2.0	3.0	2362.0	64.0	4.0	1.0	2.5	8.0	4
290000.0	0.71	1.0	73.0	61000.0	0.0	0.0	4.0	2.0	2.0	1838.0	71.0	4.0	0.0	2.0	8.0	4
205980.0	0.14	0.0	1.0	45200.0	1.0	1.0	2.0	2.0	3.0	1983.0	64.0	3.0	1.0	2.5	5.0	4
275000.0	0.54	0.0	19.0	30200.0	0.0	0.0	2.0	3.0	3.0	2175.0	64.0	4.0	1.0	2.5	10.0	4
275000.0	0.47	0.0	35.0	27800.0	0.0	0.0	2.0	3.0	3.0	2588.0	64.0	4.0	1.0	2.5	10.0	4
275000.0	0.37	0.0	14.0	31200.0	0.0	1.0	2.0	2.0	2.0	2011.0	40.0	4.0	1.0	2.5	8.0	4
275000.0	0.61	0.0	21.0	16100.0	0.0	1.0	2.0	2.0	2.0	2486.0	62.0	4.0	1.0	2.5	11.0	4
275000.0	0.46	0.0	7.0	18400.0	0.0	0.0	2.0	2.0	3.0	1865.0	57.0	3.0	0.0	2.5	8.0	4
275000.0	0.03	0.0	16.0	27000.0	0.0	1.0	2.0	2.0	3.0	1812.0	57.0	2.0	1.0	2.5	7.0	4

only showing top 20 rows



summary	Price	LotSize	Waterfront	CLUSTER
count	4655	4655	4655	4655
mean	162537.34135338347	0.4691321160042959	0.003007518796992	0.0
stddev	51449.17174680274	0.6264212879059081	0.05476420278337016	0.0
min	10300.0	0.0	0.0	0
max	600000.0	8.97	1.0	0

summary	Price	LotSize	Waterfront	CLUSTER
count	3356	3356	3356	3356
mean	208313.6853396901	0.5529678188319437	0.006555423122765197	0.0
stddev	55025.18531388466	0.6481204374941402	0.08071177527503304	0.0
min	5000.0	0.01	0.0	0
max	600000.0	7.24	1.0	0

```

weight=0.062914
mu=[0.7808989073647417,0.027594804693120447,43.592594389644596,
sigma=
1.3102878656532657    -0.008781556531403927    ... (15 total)
-0.008781556531403927    0.026833331447068984    ...
-7.778315613332572      -0.22791269559972674    ...
-1994.3759086646505     953.5473866305449      ...
-0.047375946352474545   -0.004141141861757399   ...
0.037134027030446944    -8.540323116658249E-4   ...
0.18951000788031888     0.006470217824087922    ...
0.09271897737610613     0.0036900139379503023   ...
-0.18856936310448533    0.002380060277798572    ...
68.41825719674146      -4.1514710237897745     ...
0.03436783576161607     -0.1939395215895821     ...
-0.013693506665595243   -0.020101496691984813   ...
0.07766511419505753     -7.837804459298127E-4   ...
0.05849863330847962     -4.91629596665423E-4    ...
0.21074653933320067     -0.03620197853953523    ...

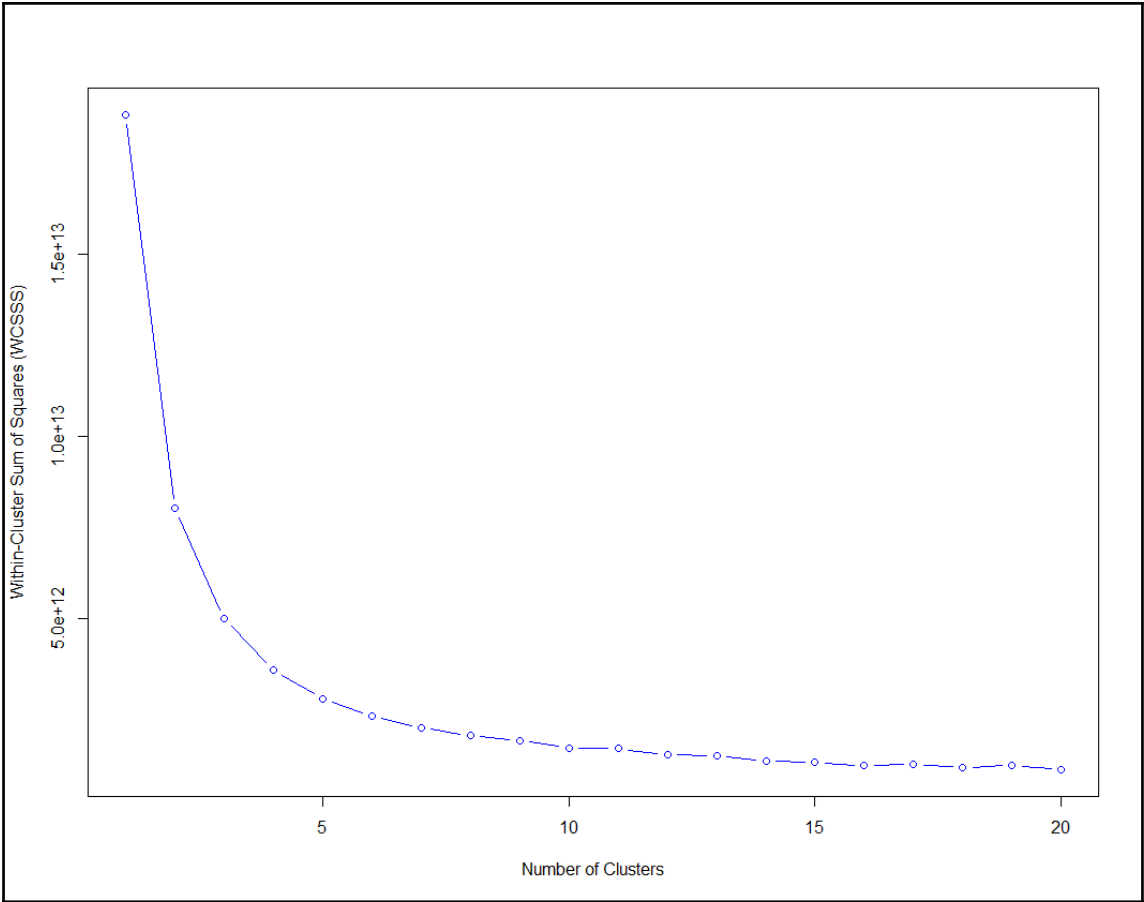
```

```
weight=0.062916
mu=[0.7808963058537262,0.027594300109707637,43.59271302953474,4
sigma=
1.310268484362083      -0.008781324169799203  ... (15 total)
-0.008781324169799203  0.026832854711163028  ...
-7.778307903368862     -0.22791180189966537  ...
-1994.312293144824     953.5408564069132     ...
-0.047374773188269     -0.004140996766905095 ...
0.03713394906123063    -8.538912936897718E-4 ...
0.18950712664774483    0.006469719690237634 ...
0.09271774632030978    0.003690014536237195 ...
-0.18856667725951917   0.002380147914518987 ...
68.41855307532745      -4.151157768413957    ...
0.03435556638123301    -0.19393666789368957 ...
-0.013692860075417968  -0.020100969638081837 ...
0.07766539257767785    -7.836426721759263E-4 ...
0.058499301884249254   -4.913759714689525E-4 ...
0.21074627737377502    -0.03620069422593031 ...
```

```
weight=0.062914
mu=[0.7808992393728132,0.027594857592962586,43.592578423292814,
sigma=
1.3102898209522704      -0.008781582527563107  ... (15 total)
-0.008781582527563107  0.026833381427386702  ...
-7.778313620309155     -0.2279126919235134  ...
-1994.3796587358433     953.5483992884145     ...
-0.04737604387775983    -0.004141154489440008 ...
0.037134040114905426    -8.540440057795036E-4 ...
0.18951016842223162    0.006470255770598361 ...
0.0927190997773798     0.003690015897543276 ...
-0.18856955926488908   0.0023800577384774702 ...
68.41823645079883      -4.151497911732236    ...
0.03436931403004916    -0.1939398262567097   ...
-0.013693555188677108  -0.020101546506840207 ...
0.07766507105646259    -7.837925362348146E-4 ...
0.05849857421583572    -4.916501535027908E-4 ...
0.21074661417458015    -0.036202094871478796 ...
```

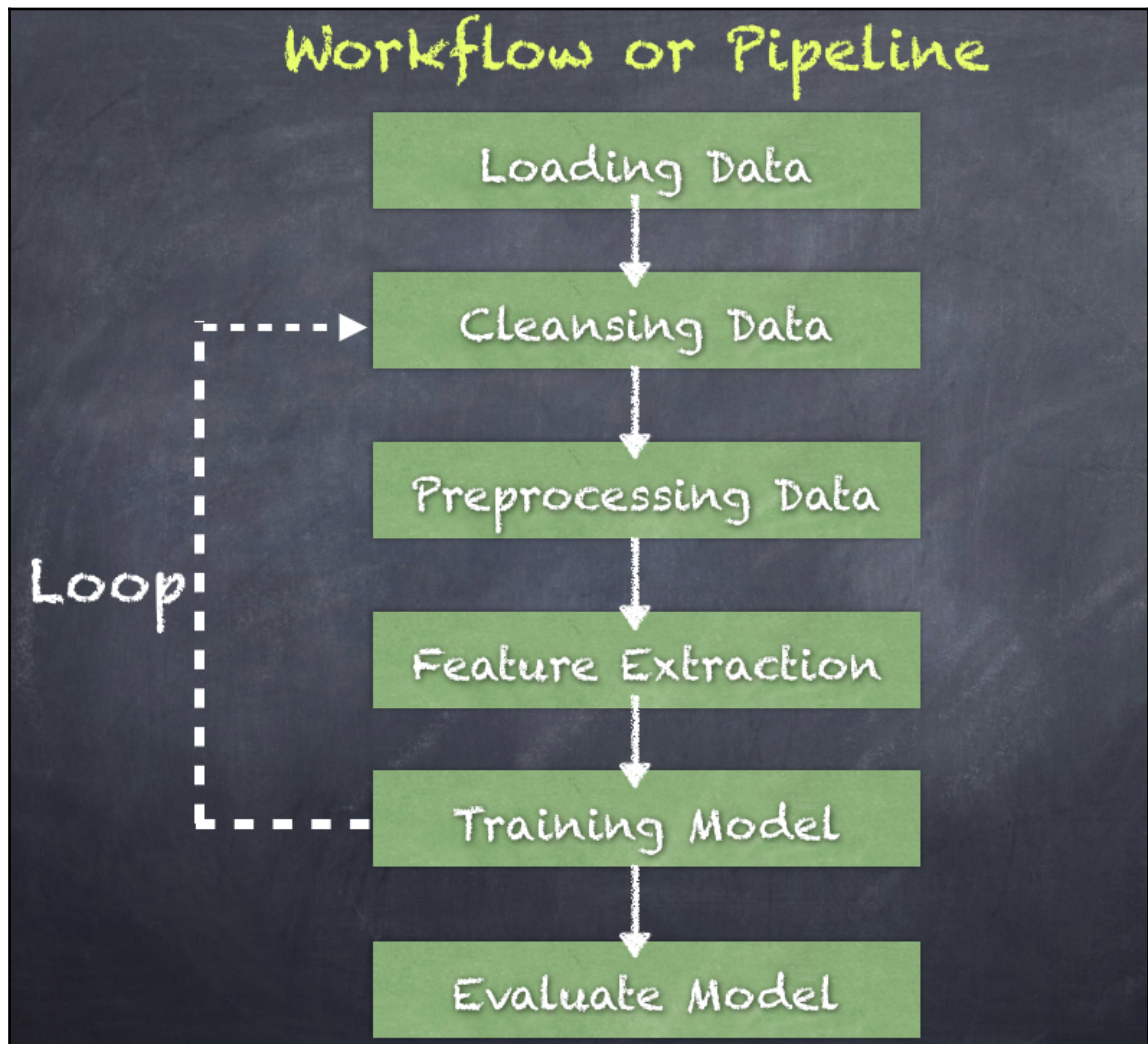
```
weight=0.062914
mu=[0.7808992393728132,0.027594857592962586,43.592578423292814,
sigma=
1.3102898209522704      -0.008781582527563107    ... (15 total)
-0.008781582527563107  0.026833381427386702    ...
-7.778313620309155     -0.2279126919235134     ...
-1994.3796587358433    953.5483992884145      ...
-0.04737604387775983   -0.004141154489440008   ...
0.037134040114905426   -8.540440057795036E-4   ...
0.18951016842223162    0.006470255770598361   ...
0.0927190997773798     0.003690015897543276   ...
-0.18856955926488908   0.0023800577384774702   ...
68.41823645079883      -4.151497911732236     ...
0.03436931403004916    -0.1939398262567097    ...
-0.013693555188677108  -0.020101546506840207   ...
0.07766507105646259    -7.837925362348146E-4   ...
0.05849857421583572    -4.916501535027908E-4   ...
0.21074661417458015    -0.036202094871478796   ...
```

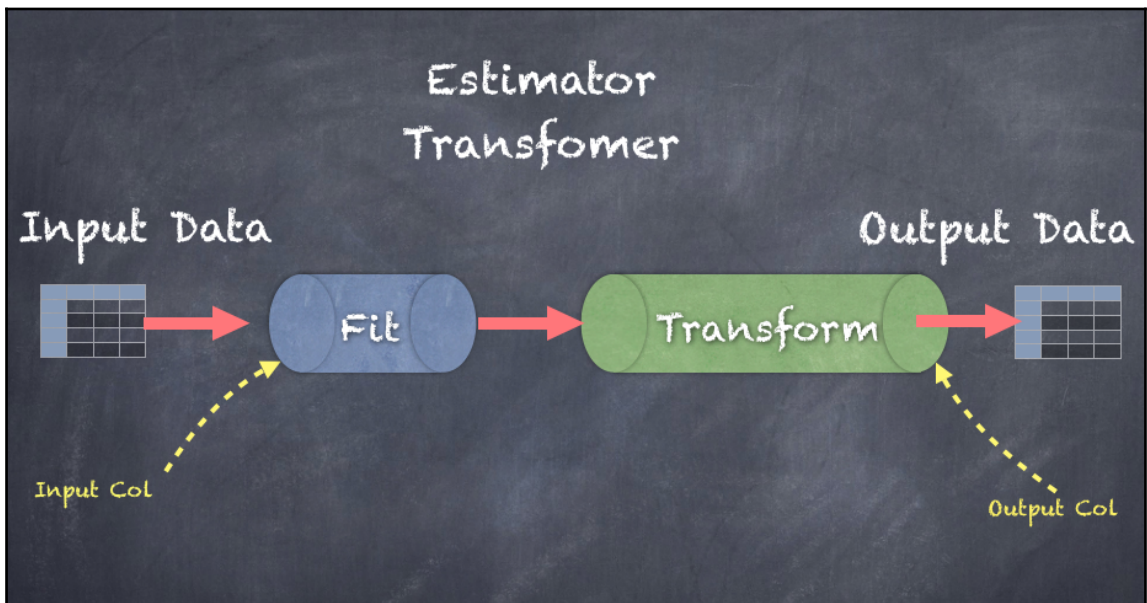
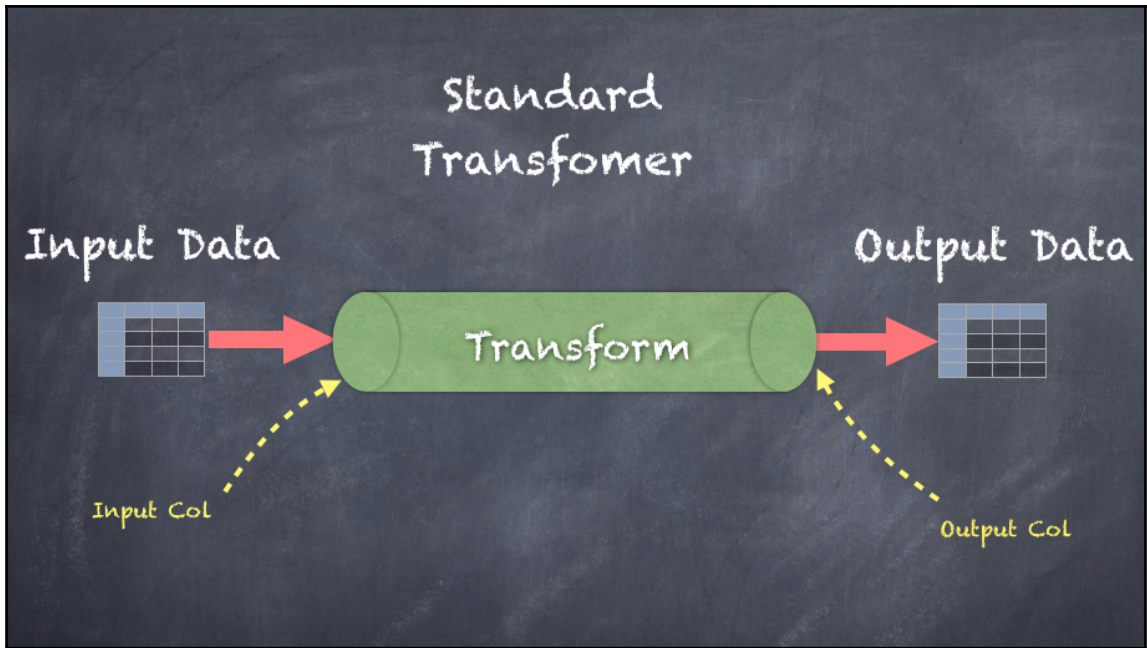
```
weight=0.748341
mu=[0.4058231162813968,0.0023199464802859684,22.644237040865264,
sigma=
0.17589935013824443    1.318739954100798E-4    ... (15 total)
1.318739954100798E-4   0.0023145643286145772   ...
0.1927311852469183     0.029438024228507845   ...
1108.2263125353527     116.97771555806634     ...
7.447161249472683E-4   -2.8238392691917215E-5  ...
-0.007840809295101704  -2.8837901025969534E-4  ...
0.023475650829693853   0.0019903868910481287   ...
-0.007865077010407383  -8.710340648340969E-5   ...
-0.06405665959152722   -5.685489003490323E-4   ...
54.81575059498132      0.06205695119789516    ...
0.08164968346291952    -0.02645130673658966    ...
0.07800513455809353    -0.001357416457442531   ...
0.023165924984143067   -1.9025716915325663E-4  ...
0.03441444210703747    -8.997788887508309E-5   ...
0.17565249306476233    -0.0017924278494132312  ...
```



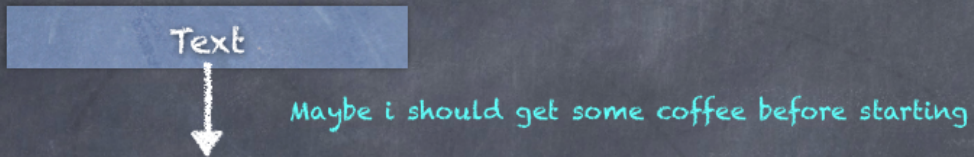


## Chapter 15: Text Analytics Using Spark ML

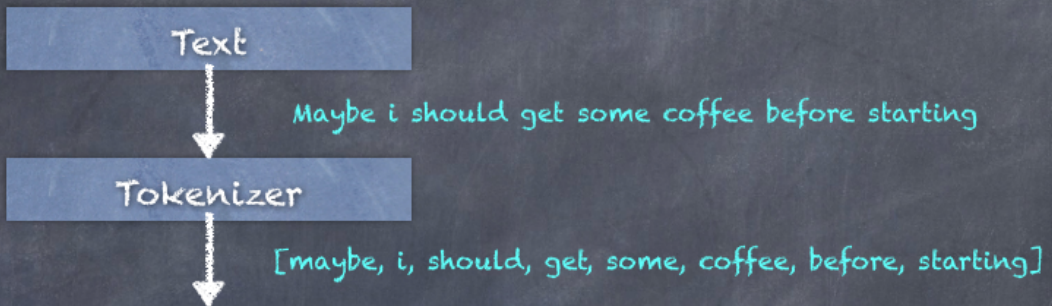


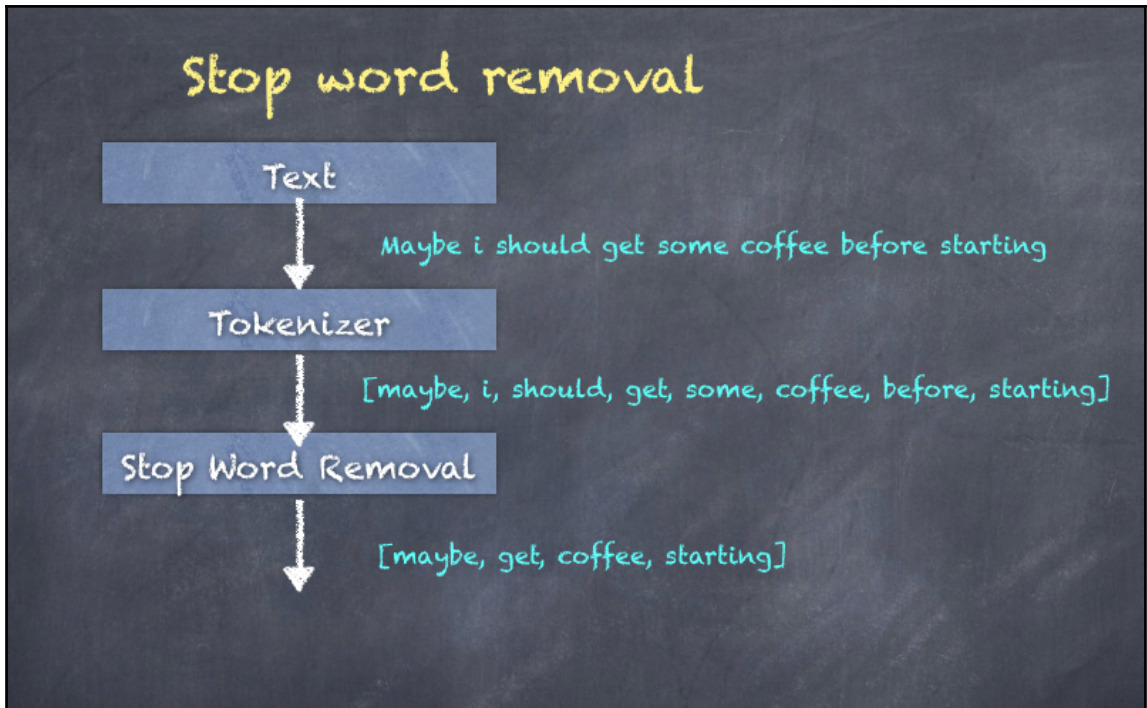


## Input Text Sentence

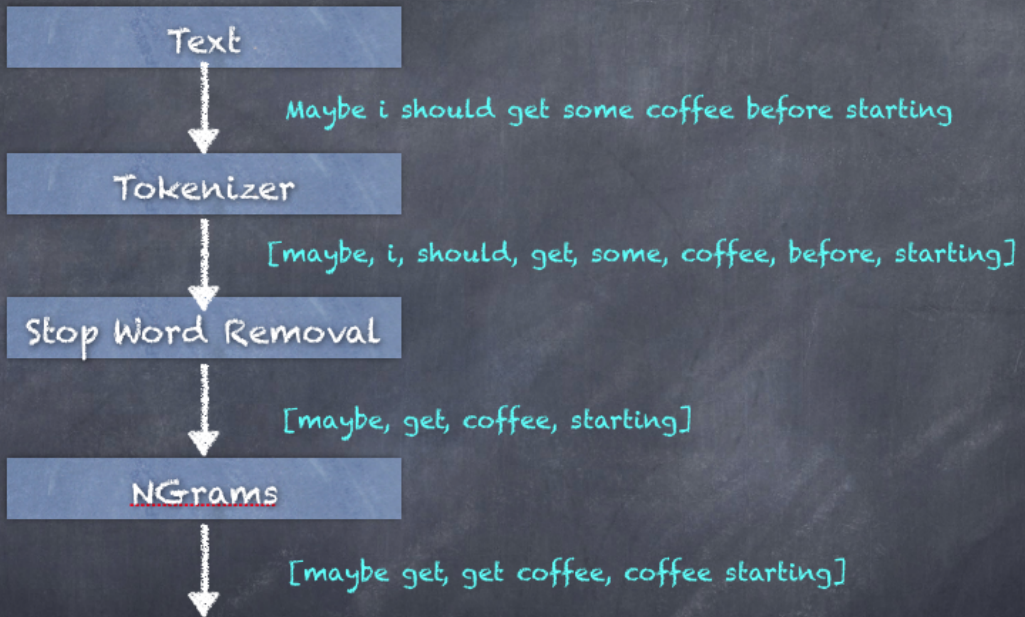


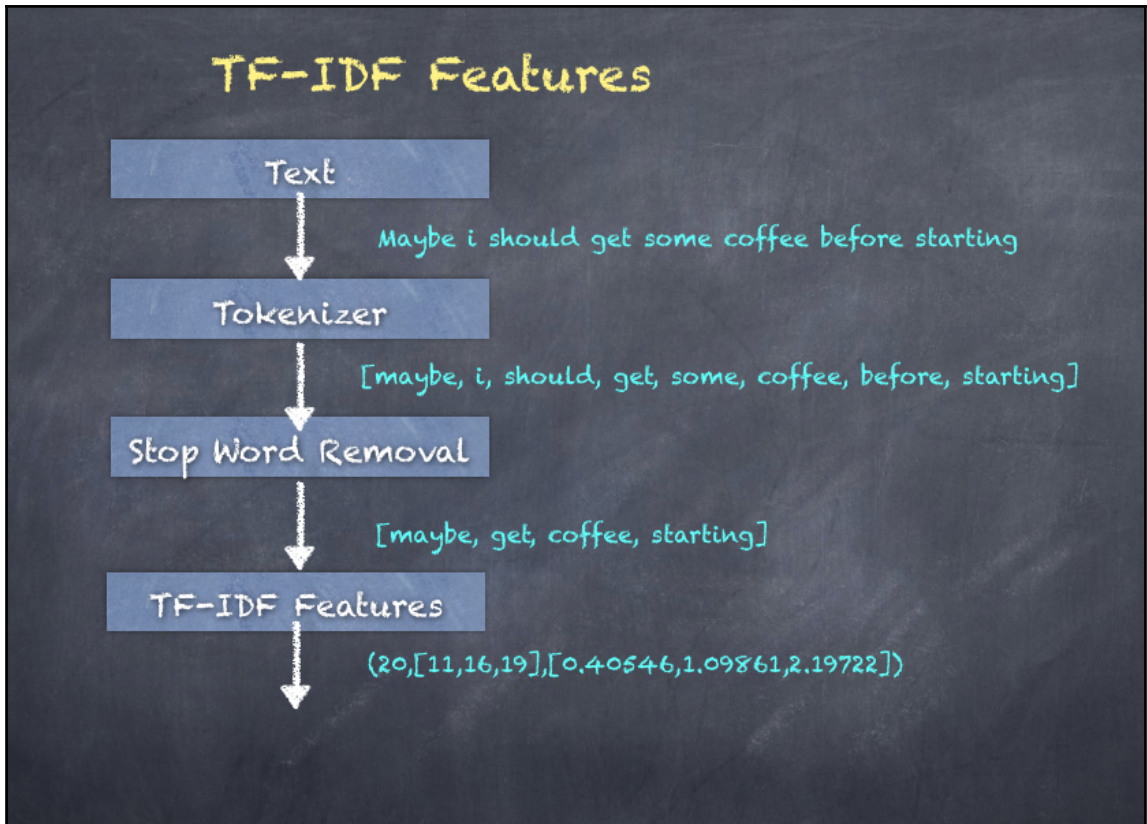
## Space delimited Tokenizer



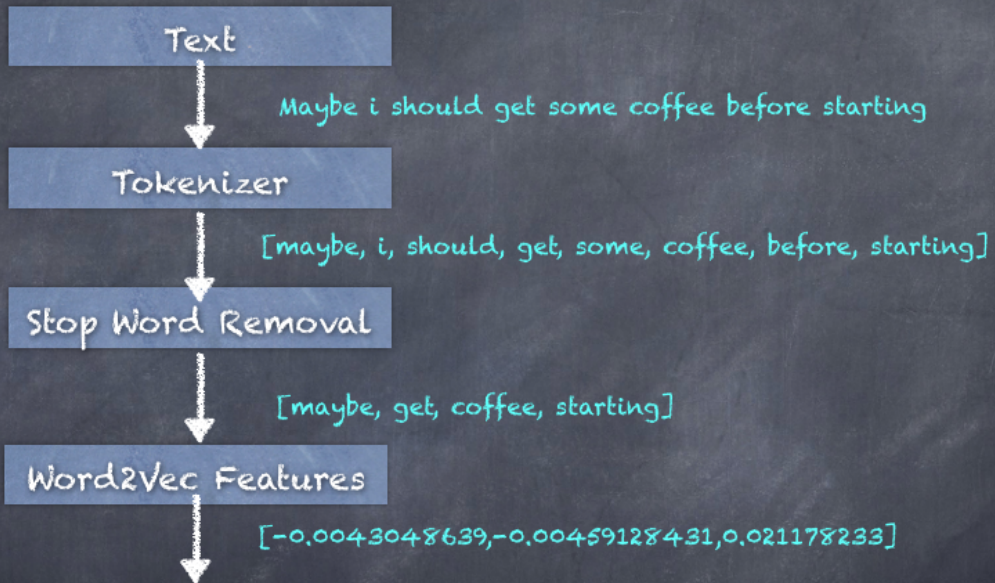


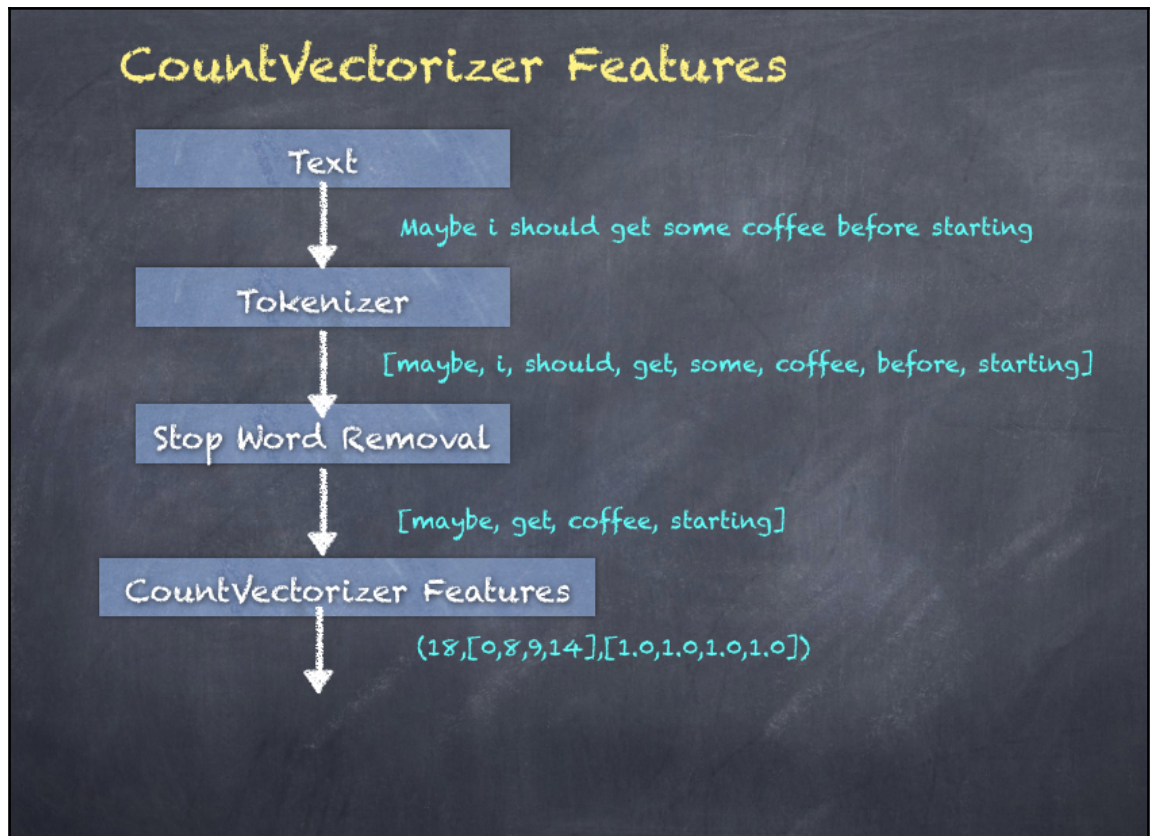
## N-Grams with $N = 2$



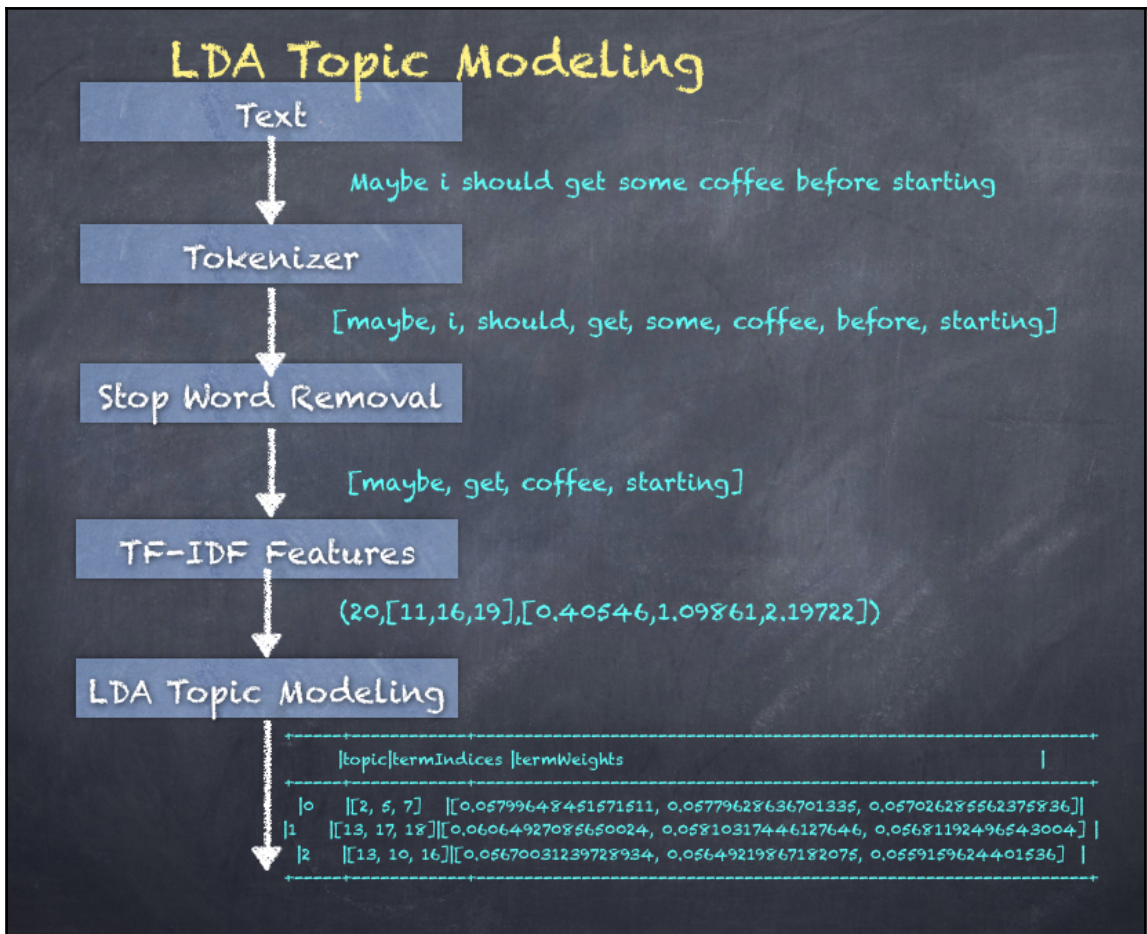


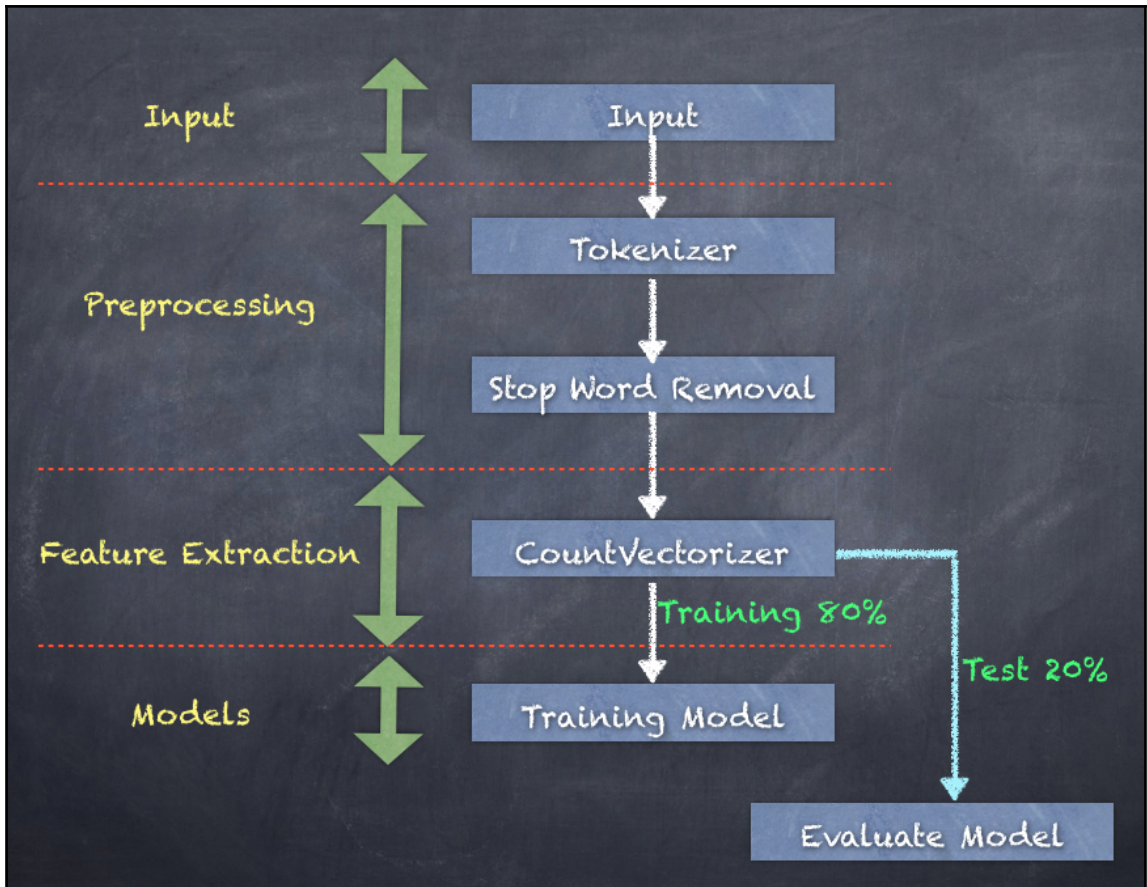
## Word2Vec Features



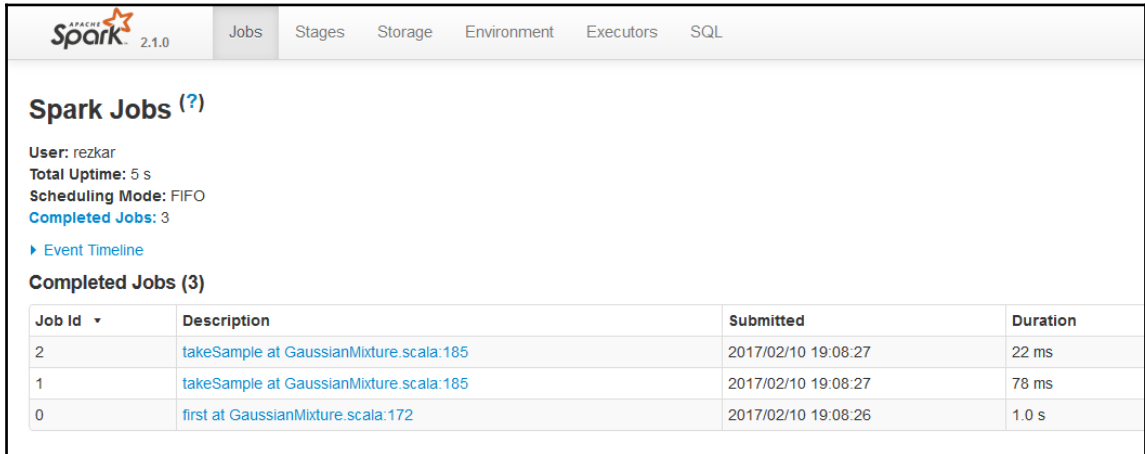








# Chapter 16: Spark Tuning



The screenshot shows the Apache Spark 2.1.0 interface. At the top, there is a navigation bar with tabs for 'Jobs', 'Stages', 'Storage', 'Environment', 'Executors', and 'SQL'. The 'Jobs' tab is selected. Below the navigation bar, the page title is 'Spark Jobs (?)'. The user is identified as 'rezkar'. The 'Total Uptime' is 5 s, and the 'Scheduling Mode' is FIFO. There are 3 completed jobs. A link for 'Event Timeline' is provided. Below this, a table titled 'Completed Jobs (3)' lists the jobs with their IDs, descriptions, submission times, and durations.

Job Id	Description	Submitted	Duration
2	<a href="#">takeSample at GaussianMixture.scala:185</a>	2017/02/10 19:08:27	22 ms
1	<a href="#">takeSample at GaussianMixture.scala:185</a>	2017/02/10 19:08:27	78 ms
0	<a href="#">first at GaussianMixture.scala:172</a>	2017/02/10 19:08:26	1.0 s

2.1.0

Jobs
Stages
Storage
Environment
Executors
SQL

## Spark Jobs (?)

**User:** rezkar  
**Total Uptime:** 6 s  
**Scheduling Mode:** FIFO  
**Completed Jobs:** 19

▼ Event Timeline  
 Enable zooming

**Executors**

- Added
- Removed

**Jobs**


- Succeeded
- Failed
- Running

400
600
800
000
200
400
600
800
000
200
400
600
800

19:09:43
19:09:44
19:09:45

### Completed Jobs (19)

Job Id ▾	Description	Submitted	Duration
18	<a href="#">show at GaussianMixtureModelDemo.scala:67</a>	2017/02/10 19:09:48	54 ms
17	<a href="#">treeAggregate at GaussianMixture.scala:201</a>	2017/02/10 19:09:48	24 ms
16	<a href="#">treeAggregate at GaussianMixture.scala:201</a>	2017/02/10 19:09:48	23 ms
15	<a href="#">treeAggregate at GaussianMixture.scala:201</a>	2017/02/10 19:09:48	21 ms
14	<a href="#">treeAggregate at GaussianMixture.scala:201</a>	2017/02/10 19:09:48	22 ms



APACHE  
**Spark**™ 2.1.0

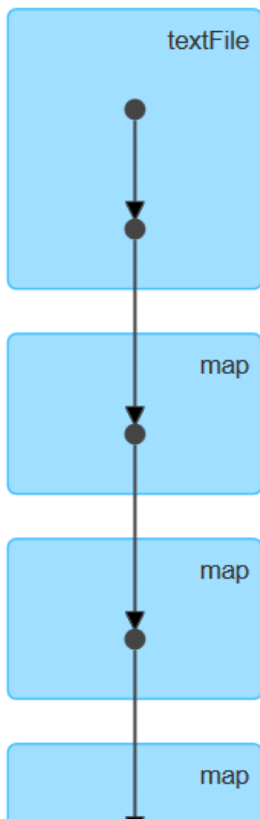
Jobs

## Details for Job 0

**Status:** RUNNING  
**Active Stages:** 1

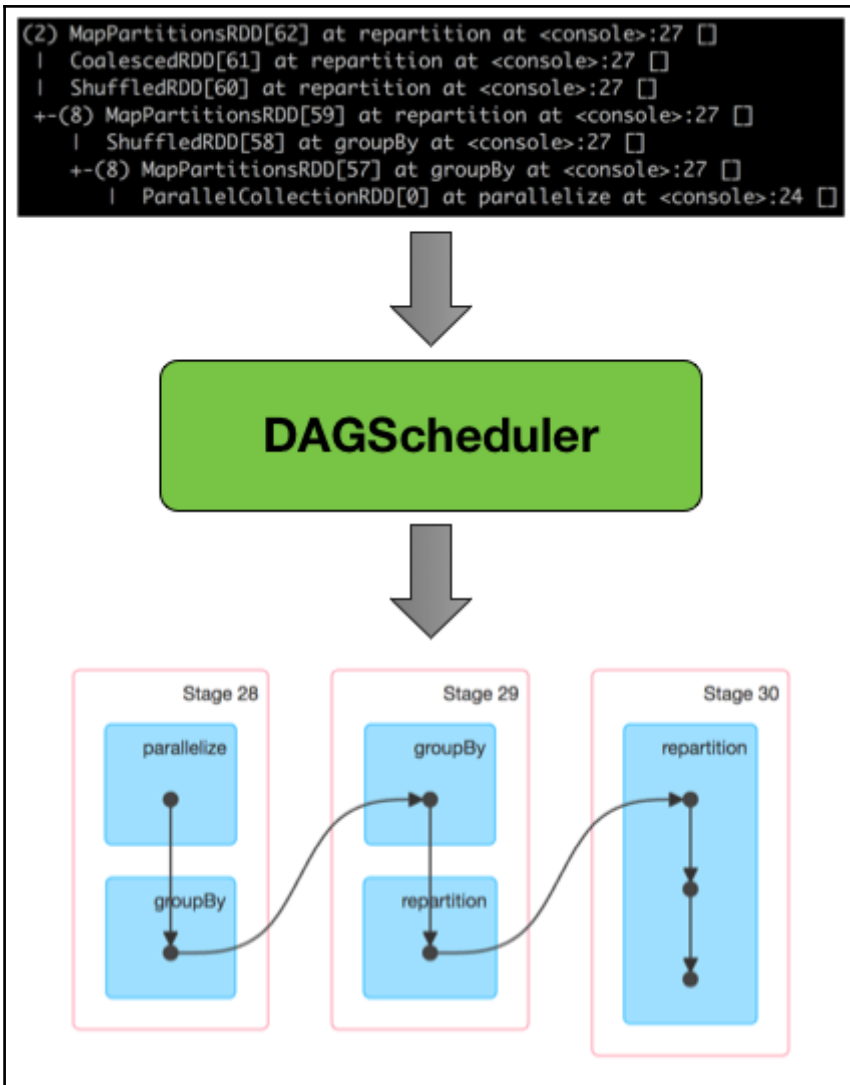
- ▶ [Event Timeline](#)
- ▼ [DAG Visualization](#)

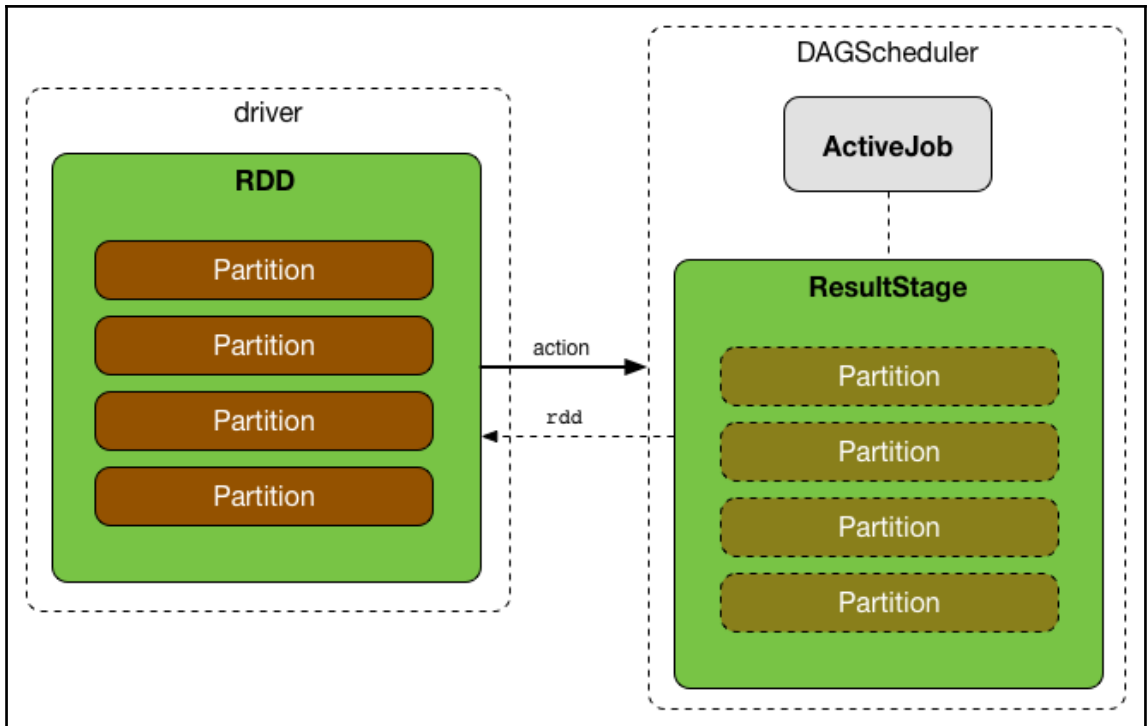
Stage 0

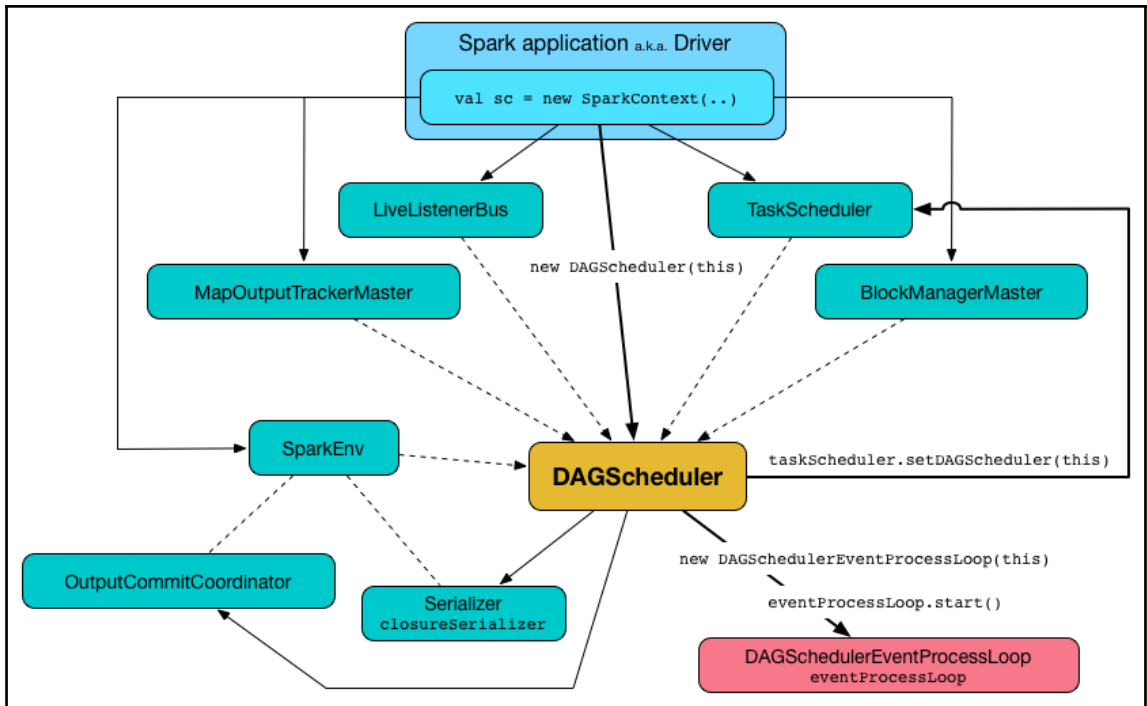


```
graph TD; A(( )) --> B(( )); B --> C(( )); C --> D(( ));
```

The diagram shows a vertical flow of four operations within a stage. The first operation is a light blue rounded rectangle labeled 'textFile'. Below it is a black dot, followed by a downward-pointing arrow, another black dot, and then a light blue rounded rectangle labeled 'map'. This sequence repeats: another downward arrow, black dot, and light blue rounded rectangle labeled 'map'. Finally, a last downward arrow, black dot, and light blue rounded rectangle labeled 'map' are shown at the bottom. The entire sequence is enclosed in a light blue rounded rectangle.







spark 2.1.0 Jobs Stages Storage Environment Executors SQL

### Stages for All Jobs

Active Stages: 1  
Completed Stages: 12

#### Active Stages (1)

Stage Id	Description	Submitted	Duration	Tasks: Succeeded/Total
12	treeAggregate at GaussianMixture.scala:201	Unknown	Unknown	0/2

#### Completed Stages (12)

Stage Id	Description	Submitted	Duration	Tasks: Succeeded/Total
11	treeAggregate at GaussianMixture.scala:201	2017/02/10 19:28:26	17 ms	2/2
10	treeAggregate at GaussianMixture.scala:201	2017/02/10 19:28:26	19 ms	2/2
9	treeAggregate at GaussianMixture.scala:201	2017/02/10 19:28:26	18 ms	2/2
8	treeAggregate at GaussianMixture.scala:201	2017/02/10 19:28:26	23 ms	2/2
7	treeAggregate at GaussianMixture.scala:201	2017/02/10 19:28:26	21 ms	2/2
6	treeAggregate at GaussianMixture.scala:201	2017/02/10 19:28:26	21 ms	2/2
5	treeAggregate at GaussianMixture.scala:201	2017/02/10 19:28:26	24 ms	2/2
4	treeAggregate at GaussianMixture.scala:201	2017/02/10 19:28:26	35 ms	2/2
3	treeAggregate at GaussianMixture.scala:201	2017/02/10 19:28:26	0.2 s	2/2
2	takeSample at GaussianMixture.scala:185	2017/02/10 19:28:26	20 ms	2/2
1	takeSample at GaussianMixture.scala:185	2017/02/10 19:28:25	50 ms	2/2
0	first at GaussianMixture.scala:172	2017/02/10 19:28:25	0.8 s	1/1



Summary Metrics for 2 Completed Tasks			
Metric	Min	25th percentile	Median
Duration	0.2 s	0.2 s	0.2 s
GC Time	0 ms	0 ms	0 ms
Input Size / Records	27.6 KB / 1	27.6 KB / 1	28.6 KB / 1

Aggregated Metrics by Executor					
Executor ID ^	Address	Task Time	Total Tasks	Failed Tasks	Killed Tasks
driver	10.2.17.13:53512	0.5 s	2	0	0

Tasks (2)						
Index ^	ID	Attempt	Status	Locality Level	Executor ID / Host	Launch Time
0	4	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2017/02/04 12:57:01
1	5	0	SUCCESS	PROCESS_LOCAL	driver / localhost	2017/02/04 12:57:01

Storage				
RDD Name	Storage Level	Cached Partitions	Fraction Cached	Size in Memory
MapPartitionsRDD	Memory Deserialized 1x Replicated	2	100%	281.8 KB
*SerializeFromObject [assertNotNull(input[0, com.chapter13.Clustering.GaussianMixtureModelDemo\$Land, true], top level Product input object) Price AS Price#17, assertNotNull(input[0, com.chapter13.Clustering.GaussianMixtureModelDemo\$Land, true], top level Product input object) LotSize AS LotSize#18, assertNotNull(input[0, com.chapter13.Clustering.GaussianMixtureModelDemo\$Land, true], top level Product input object) Waterfront AS Waterfront#19, assertNotNull(input[0, com.chapter13.Clustering.GaussianMixtureModelDemo\$Land, true], top level Product input object) Age AS Age#20, assertNotNull(input[0, com.chapter13.Clustering.GaussianMixtureModelDemo\$Land, true], top level Product input object) LandValue AS LandValue#21, assertNotNull(input[0, com.chapter13.Clustering.GaussianMixtureModelDemo\$Land, true], top level Product input object) NewConstruct AS NewConstruct#22, assertNotNull(input[0, com.chapter13.Clustering.GaussianMixtureModelDemo\$Land, true], top level Product input object) CentralAir AS CentralAir#23...	Memory Deserialized 1x Replicated	2	100%	220.5 KB
MapPartitionsRDD	Memory Deserialized 1x Replicated	2	100%	244.3 KB

APACHE Spark 2.1.0 Jobs Stages Storage Environment Executors SQL

### RDD Storage Info for MapPartitionsRDD

Storage Level: Memory Deserialized 1x Replicated  
 Cached Partitions: 1  
 Total Partitions: 2  
 Memory Size: 5.3 MB  
 Disk Size: 0.0 B

Data Distribution on 1 Executors

Host	Memory Usage
192.168.192.1:52341	5.3 MB (4.0 GB Remaining)

1 Partitions

Block Name	Storage Level	Size in Memory
rdd_4_0	Memory Deserialized 1x Replicated	5.3 MB

APACHE Spark 2.1.0 Jobs Stages Storage Environment Executors SQL

### Environment

Runtime Information

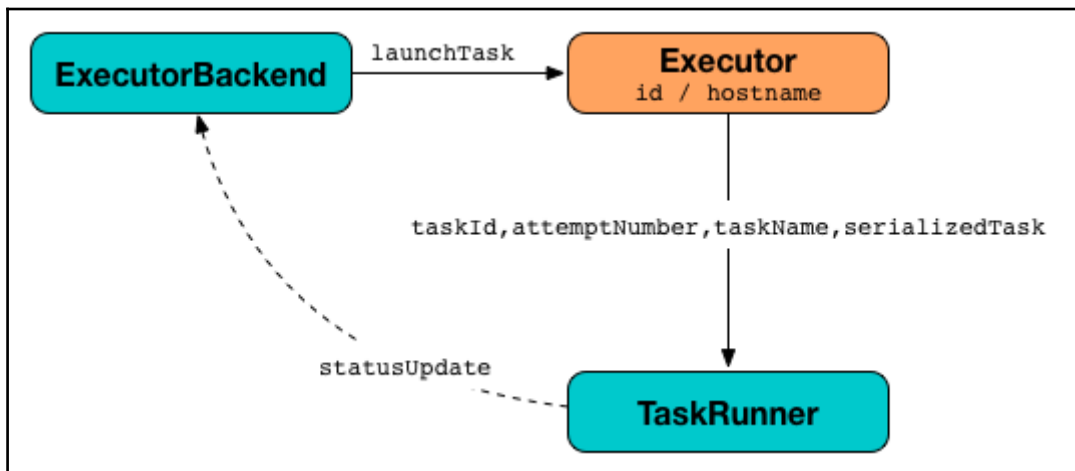
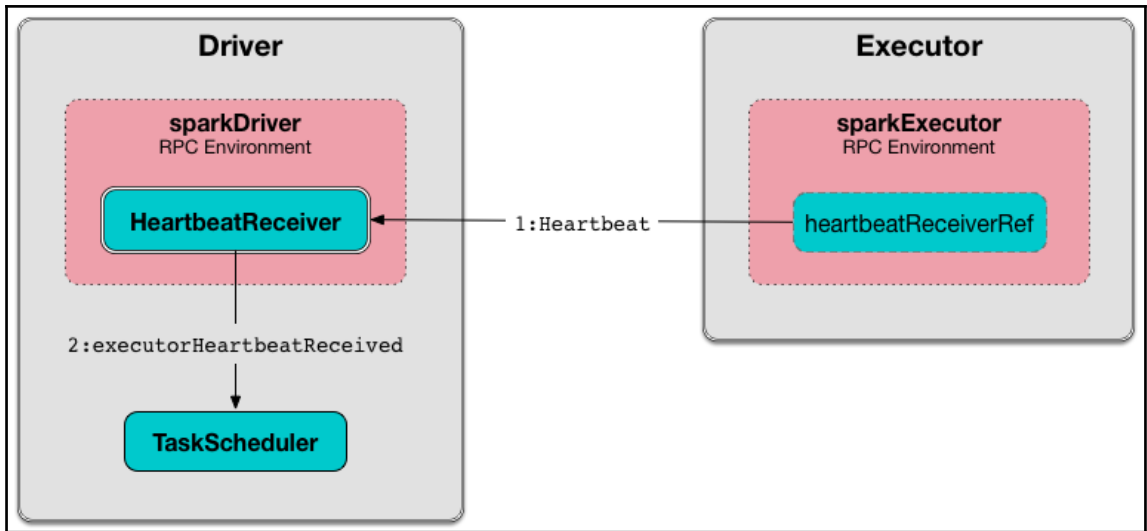
Name	Value
Java Home	C:\Program Files\Java\jdk1.8.0_102\jre
Java Version	1.8.0_102 (Oracle Corporation)
Scala Version	version 2.11.8

Spark Properties

Name	Value
spark.app.id	local-1485627820379
spark.app.name	SparkDFebay
spark.driver.host	10.2.16.255
spark.driver.port	64649
spark.executor.id	driver
spark.master	local[*]
spark.scheduler.mode	FIFO

System Properties

Name	Value
awt.toolkit	sun.awt.windows.WToolkit
file.encoding	UTF-8
file.encoding.pkg	sun.io



APACHE **spark** 2.1.0    Jobs   Stages   Storage   Environment   **Executors**

### Executors

**Summary**

	RDD Blocks	Storage Memory	Disk Used	Cores	Active Tasks	Failed Tasks	Complete Tasks	Total Tasks	Task Time (GC Time)
Active(1)	4	263.4 KB / 4.3 GB	0.0 B	8	1	0	0	1	0 ms (0 ms)
Dead(0)	0	0.0 B / 0.0 B	0.0 B	0	0	0	0	0	0 ms (0 ms)
<b>Total(1)</b>	<b>4</b>	<b>263.4 KB / 4.3 GB</b>	<b>0.0 B</b>	<b>8</b>	<b>1</b>	<b>0</b>	<b>0</b>	<b>1</b>	<b>0 ms (0 ms)</b>

**Executors**

Show  entries

Executor ID	Address	Status	RDD Blocks	Storage Memory	Disk Used	Cores	Active Tasks	Failed Tasks	Complete Tasks	Total Tasks	Task Time (GC Time)
driver	10.2.16.255:64829	Active	4	263.4 KB / 4.3 GB	0.0 B	8	1	0	0	1	0 ms (0 ms)

Showing 1 to 1 of 1 entries

APACHE **spark** 2.1.0    **Jobs**   Stages   Storage   Environment   Executors   SQL

### Spark Jobs (?)


User: rezkar  
 Total Uptime: 7 s  
 Scheduling Mode: FIFO  
 Completed Jobs: 25

▼ Event Timeline  
 Enable zooming

The chart displays two events on a timeline from 19:19:37 to 19:19:39. The first event, 'Executor driver added', is marked with a blue dot at approximately 19:19:37.5. The second event, 'first at GaussianMixt', is marked with a blue dot at approximately 19:19:39.5. The legend indicates that blue dots represent 'Added' executors and 'Succeeded' jobs.

**Completed Jobs (25)**

Job Id	Description	Submitted	Duration
24	<a href="#">show at GaussianMixtureModelDemo.scala:75</a>	2017/02/10 19:19:42	14 ms
23	<a href="#">run at ThreadPoolExecutor.java:1142</a>	2017/02/10 19:19:42	12 ms
22	<a href="#">show at GaussianMixtureModelDemo.scala:74</a>	2017/02/10 19:19:42	15 ms
21	<a href="#">run at ThreadPoolExecutor.java:1142</a>	2017/02/10 19:19:42	13 ms
20	<a href="#">show at GaussianMixtureModelDemo.scala:71</a>	2017/02/10 19:19:42	13 ms
19	<a href="#">run at ThreadPoolExecutor.java:1142</a>	2017/02/10 19:19:42	23 ms


2.1.0

Jobs
Stages
Storage
Environment
Executors
SQL

## Spark Jobs (?)

User: rezkar  
 Total Uptime: 6 s  
 Scheduling Mode: FIFO  
 Active Jobs: 1  
 Completed Jobs: 12  
[Event Timeline](#)

### Active Jobs (1)

Job Id	Description	Submitted	Duration	Stages: Succeeded/Total
12	treeAggregate at GaussianMixture.scala:201	(kill) 2017/02/10 19:32:16	30 ms	0/1

### Completed Jobs (12)

Job Id	Description	Submitted	Duration	Stages: Succeeded/Total
11	treeAggregate at GaussianMixture.scala:201	2017/02/10 19:32:16	57 ms	1/1
10	treeAggregate at GaussianMixture.scala:201	2017/02/10 19:32:16	35 ms	1/1
9	treeAggregate at GaussianMixture.scala:201	2017/02/10 19:32:16	68 ms	1/1
8	treeAggregate at GaussianMixture.scala:201	2017/02/10 19:32:16	28 ms	1/1
7	treeAggregate at GaussianMixture.scala:201	2017/02/10 19:32:16	33 ms	1/1
6	treeAggregate at GaussianMixture.scala:201	2017/02/10 19:32:16	34 ms	1/1
5	treeAggregate at GaussianMixture.scala:201	2017/02/10 19:32:16	46 ms	1/1
4	treeAggregate at GaussianMixture.scala:201	2017/02/10 19:32:16	44 ms	1/1
3	treeAggregate at GaussianMixture.scala:201	2017/02/10 19:32:15	0.2 s	1/1
2	takeSample at GaussianMixture.scala:185	2017/02/10 19:32:15	29 ms	1/1
1	takeSample at GaussianMixture.scala:185	2017/02/10 19:32:15	69 ms	1/1

```

# Set everything to be logged to the console
log4j.rootCategory=INFO, console
log4j.appender.console=org.apache.log4j.ConsoleAppender
log4j.appender.console.target=System.err
log4j.appender.console.layout=org.apache.log4j.PatternLayout
log4j.appender.console.layout.ConversionPattern=%d{yy/MM/dd HH:mm:ss} %p %c{1}: %m%n

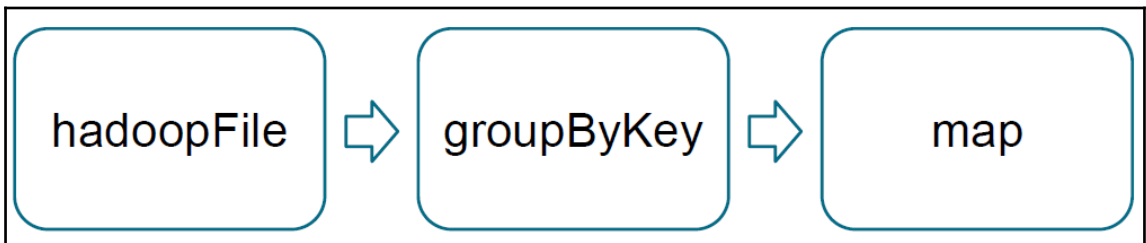
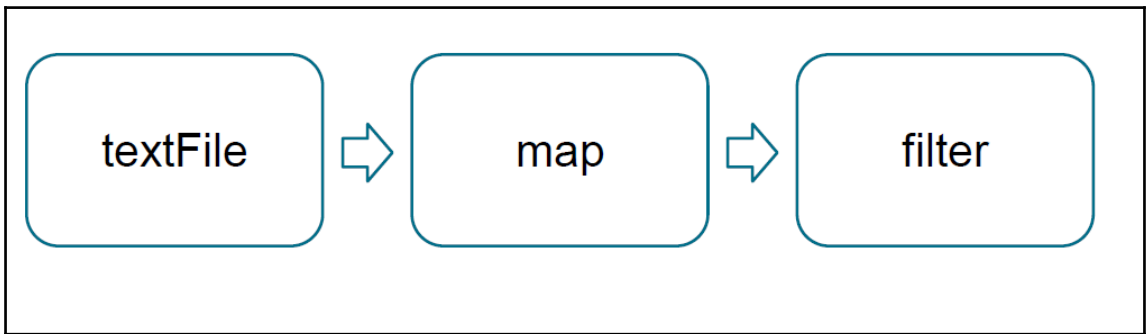
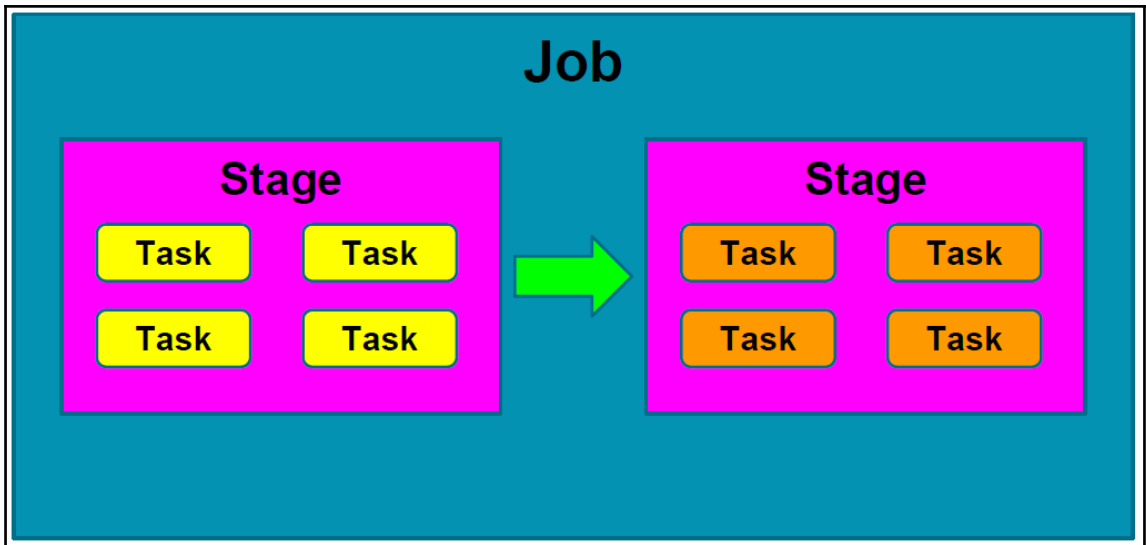
# Set the default spark-shell log level to WARN. When running the spark-shell, the
# log level for this class is used to overwrite the root logger's log level, so that
# the user can have different defaults for the shell and regular Spark apps.
log4j.logger.org.apache.spark.repl.Main=WARN

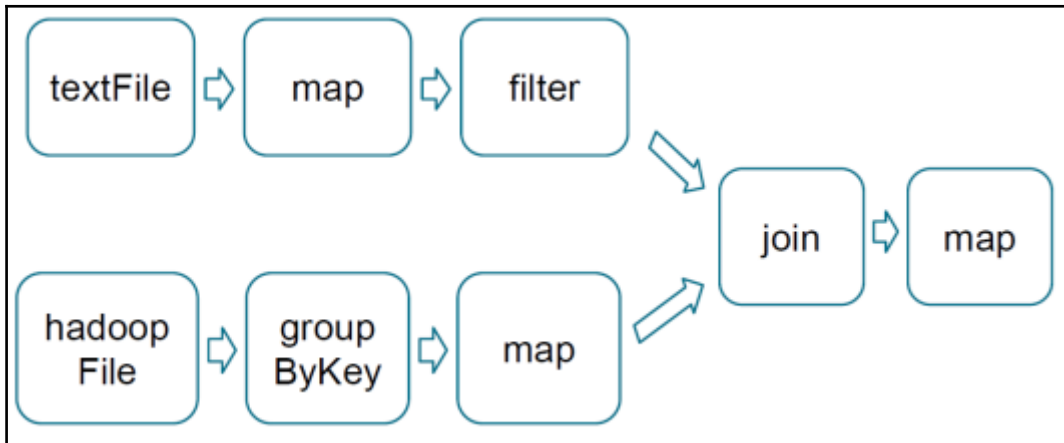
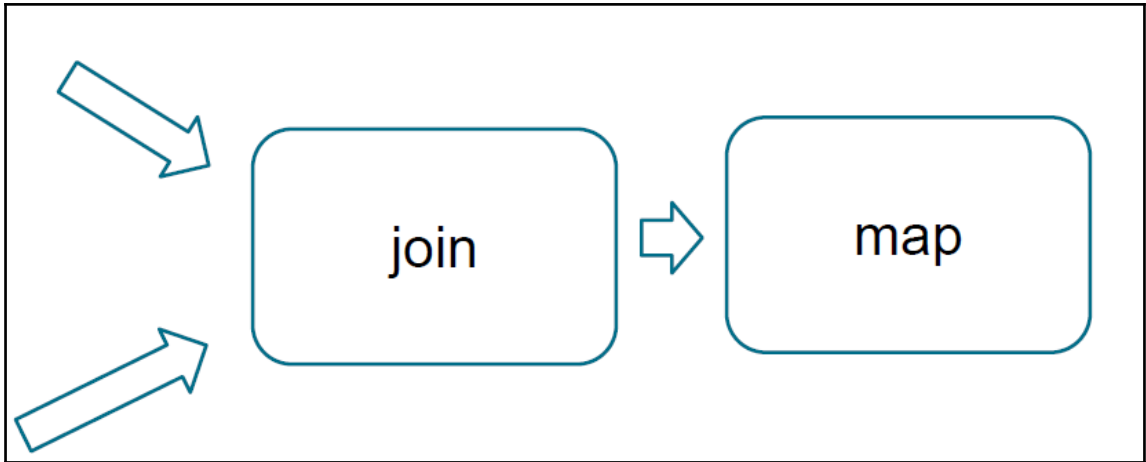
# Settings to quiet third party logs that are too verbose
log4j.logger.org.spark_project.jetty=WARN
log4j.logger.org.spark_project.jetty.util.component.AbstractLifeCycle=ERROR
log4j.logger.org.apache.spark.repl.SparkIMain$exprTyper=INFO
log4j.logger.org.apache.spark.repl.SparkILoop$SparkILoopInterpreter=INFO
log4j.logger.org.apache.parquet=ERROR
log4j.logger.parquet=ERROR

# SPARK-9183: Settings to avoid annoying messages when looking up nonexistent UDFs in SparkSQL with Hive support
log4j.logger.org.apache.hadoop.hive.metastore.RetryingHMSHandler=FATAL
log4j.logger.org.apache.hadoop.hive.ql.exec.FunctionRegistry=ERROR

```

Environment Variable	Meaning
SPARK_MASTER_HOST	Bind the master to a specific hostname or IP address, for example a public one.
SPARK_MASTER_PORT	Start the master on a different port (default: 7077).
SPARK_MASTER_WEBUI_PORT	Port for the master web UI (default: 8080).
SPARK_MASTER_OPTS	Configuration properties that apply only to the master in the form "-Dx=y" (default: none). See below for a list of possible options.
SPARK_LOCAL_DIRS	Directory to use for "scratch" space in Spark, including map output files and RDDs that get stored on disk. This should be on a fast, local disk in your system. It can also be a comma-separated list of multiple directories on different disks.
SPARK_WORKER_CORES	Total number of cores to allow Spark applications to use on the machine (default: all available cores).
SPARK_WORKER_MEMORY	Total amount of memory to allow Spark applications to use on the machine, e.g. 1000m, 2g (default: total memory minus 1 GB); note that each application's <i>individual</i> memory is configured using its <code>spark.executor.memory</code> property.
SPARK_WORKER_PORT	Start the Spark worker on a specific port (default: random).
SPARK_WORKER_WEBUI_PORT	Port for the worker web UI (default: 8081).
SPARK_WORKER_DIR	Directory to run applications in, which will include both logs and scratch space (default: SPARK_HOME/work).
SPARK_WORKER_OPTS	Configuration properties that apply only to the worker in the form "-Dx=y" (default: none). See below for a list of possible options.
SPARK_DAEMON_MEMORY	Memory to allocate to the Spark master and worker daemons themselves (default: 1g).
SPARK_DAEMON_JAVA_OPTS	JVM options for the Spark master and worker daemons themselves in the form "-Dx=y" (default: none).
SPARK_PUBLIC_DNS	The public DNS name of the Spark master and workers (default: none).








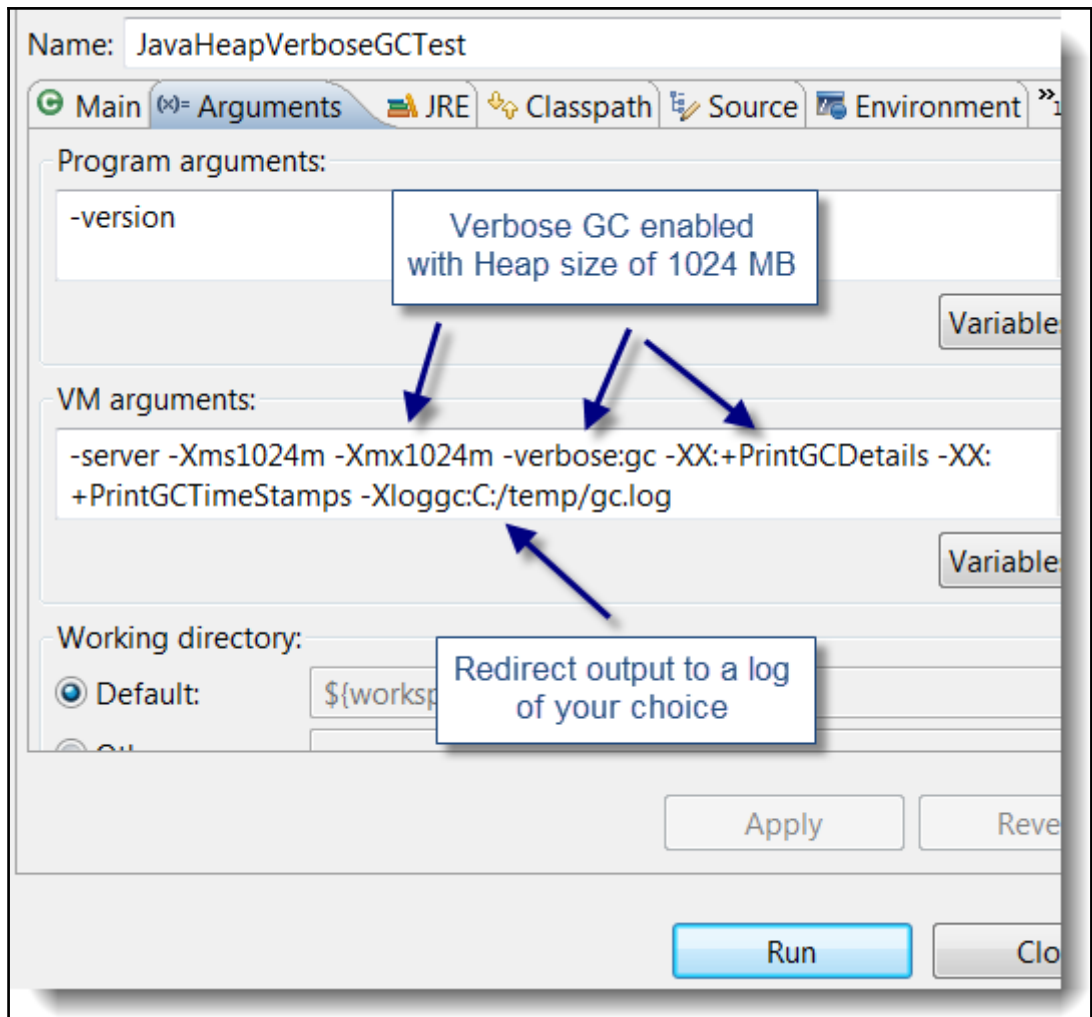
```

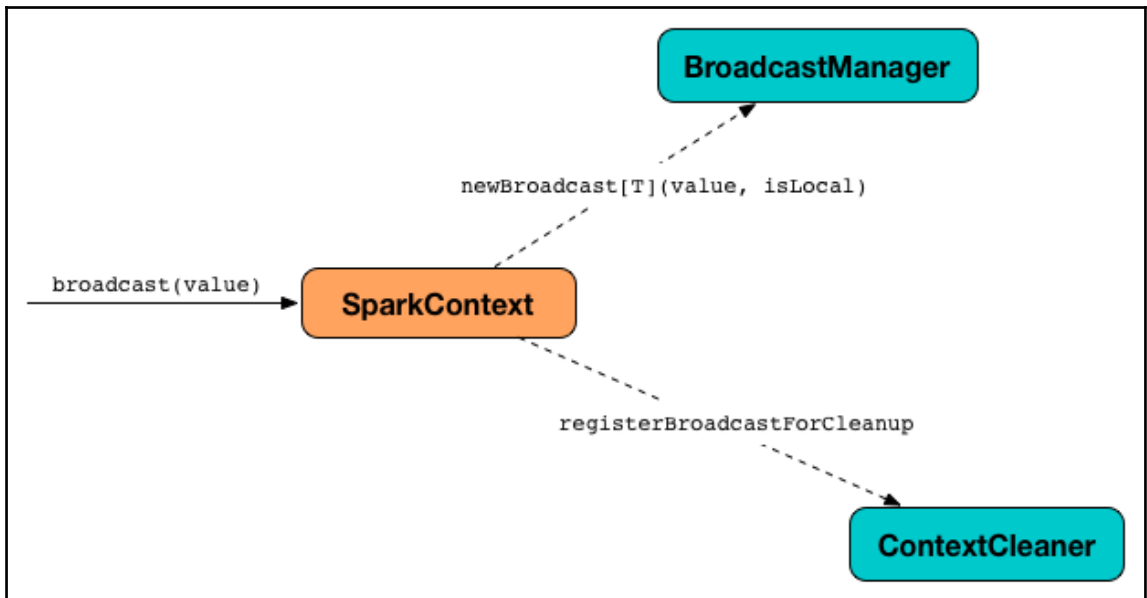
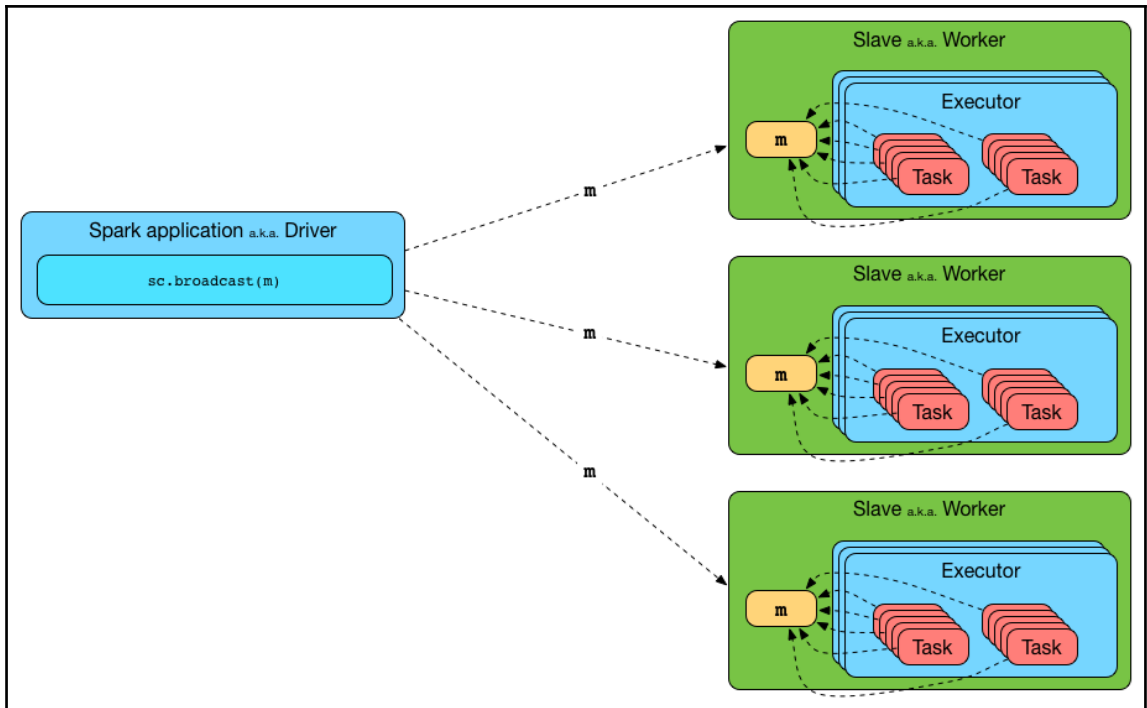
/11/20 17:20:58 INFO TaskSchedulerImpl: Removed TaskSet 1.0, whose tasks have all completed, from pool
/11/20 17:20:58 INFO TaskSchedulerImpl: Removed TaskSet 2.0, whose tasks have all completed, from pool
/11/20 17:20:58 INFO DAGScheduler: Failed to run collect at ReceiverTracker.scala:270
/11/20 17:20:58 INFO TaskSchedulerImpl: Cancelling stage 1
Exception in thread "Thread-53" org.apache.spark.SparkException: Job aborted due to stage failure: All masters are unresponsive!
Giving up.
  at org.apache.spark.scheduler.DAGScheduler.org$apache$spark$scheduler$DAGScheduler$$failJobAndIndependentStages
(DAGScheduler.scala:1033)
  at org.apache.spark.scheduler.DAGScheduler$$anonfun$abortStage$1.apply(DAGScheduler.scala:1017)
  at org.apache.spark.scheduler.DAGScheduler$$anonfun$abortStage$1.apply(DAGScheduler.scala:1015)
  at scala.collection.mutable.ResizableArray$class.foreach(ResizableArray.scala:59)
  at scala.collection.mutable.ArrayBuffer.foreach(ArrayBuffer.scala:47)
  at org.apache.spark.scheduler.DAGScheduler.abortStage(DAGScheduler.scala:1015)
  at org.apache.spark.scheduler.DAGScheduler$$anonfun$handleTaskSetFailed$1.apply(DAGScheduler.scala:633)
  at org.apache.spark.scheduler.DAGScheduler$$anonfun$handleTaskSetFailed$1.apply(DAGScheduler.scala:633)
  at scala.Option.foreach(Option.scala:236)
  at org.apache.spark.scheduler.DAGScheduler.handleTaskSetFailed(DAGScheduler.scala:633)
  at org.apache.spark.scheduler.DAGSchedulerEventProcessActor$$anonfun$receive$2.applyOrElse(DAGScheduler.scala:1207)
  at akka.actor.ActorCell.receiveMessage(ActorCell.scala:498)
  at akka.actor.ActorCell.invoke(ActorCell.scala:456)
  at akka.dispatch.Mailbox.processMailbox(Mailbox.scala:237)
  at akka.dispatch.Mailbox.run(Mailbox.scala:219)
  at akka.dispatch.ForkJoinExecutorConfigurator$AkkaForkJoinTask.exec(AbstractDispatcher.scala:386)
  at scala.concurrent.forkjoin.ForkJoinTask.doExec(ForkJoinTask.java:260)
  at scala.concurrent.forkjoin.ForkJoinPool$WorkQueue.runTask(ForkJoinPool.java:1339)
  at scala.concurrent.forkjoin.ForkJoinPool.runWorker(ForkJoinPool.java:1979)
  at scala.concurrent.forkjoin.ForkJoinWorkerThread.run(ForkJoinWorkerThread.java:107)
/11/20 17:20:58 INFO DAGScheduler: Failed to run take at DStream.scala:593
/11/20 17:20:58 INFO TaskSchedulerImpl: Cancelling stage 2
/11/20 17:20:58 INFO JobsScheduler: Starting job streaming job 1416484202000 ms.0 from job set of time 1416484202000 ms
/11/20 17:20:58 INFO SparkContext: Starting job: take at DStream.scala:593
/11/20 17:20:58 ERROR JobsScheduler: Error running job streaming job-1416484202000-ms-0
org.apache.spark.SparkException: Job aborted due to stage failure: All masters are unresponsive! Giving up.
  at org.apache.spark.scheduler.DAGScheduler.org$apache$spark$scheduler$DAGScheduler$$failJobAndIndependentStages
(DAGScheduler.scala:1033)
  at org.apache.spark.scheduler.DAGScheduler$$anonfun$abortStage$1.apply(DAGScheduler.scala:1017)
  at org.apache.spark.scheduler.DAGScheduler$$anonfun$abortStage$1.apply(DAGScheduler.scala:1015)
  at scala.collection.mutable.ResizableArray$class.foreach(ResizableArray.scala:59)
  at scala.collection.mutable.ArrayBuffer.foreach(ArrayBuffer.scala:47)
  at org.apache.spark.scheduler.DAGScheduler.abortStage(DAGScheduler.scala:1015)

```

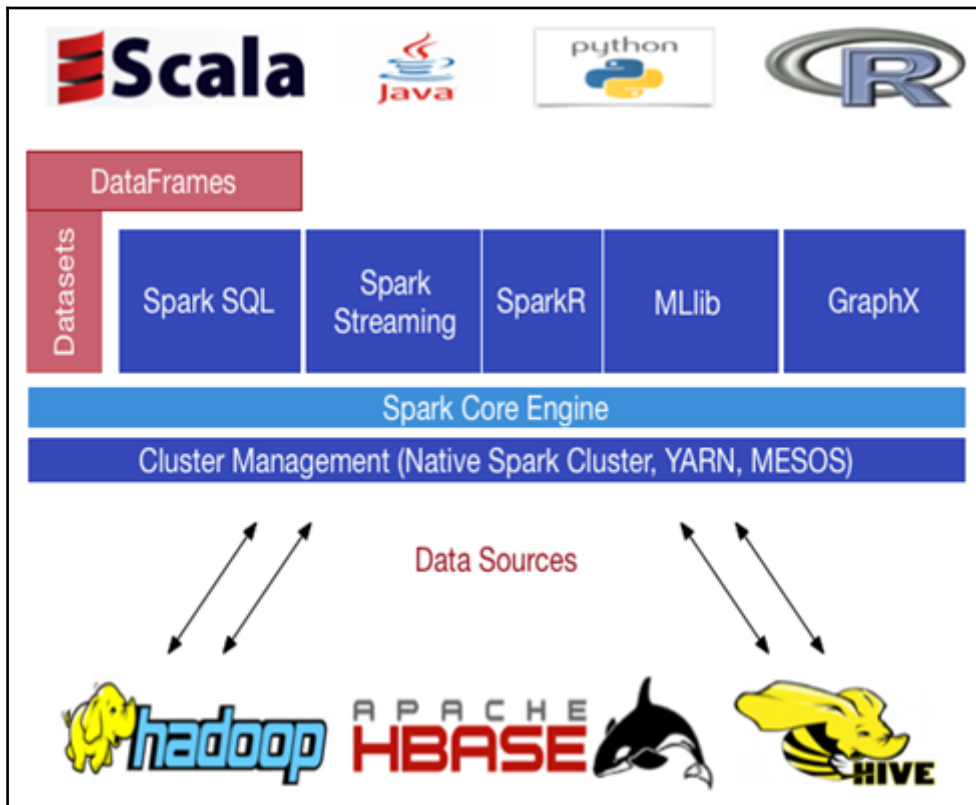
Job Scheduling Information	Diagnostic Info
NA	<p>Job initialization failed: java.io.IOException: Filesystem closed at org.apache.hadoop.hdfs.DFSCClient.checkOpen(DFSCClient.java:241) at org.apache.hadoop.hdfs.DFSCClient.access\$800(DFSCClient.java:74) at org.apache.hadoop.hdfs.DFSCClient\$DFSOutputStream.closeInternal(DFSCClient.java:3667) at org.apache.hadoop.hdfs.DFSCClient\$DFSOutputStream.close(DFSCClient.java:3626) at org.apache.hadoop.fs.FSDataOutputStream\$PositionCache.close(FSDataOutputStream.java:61) at org.apache.hadoop.fs.FSDataOutputStream.close(FSDataOutputStream.java:86) at org.apache.hadoop.security.Credentials.writeTokenStorageFile(Credentials.java:171) at org.apache.hadoop.mapred.JobInProgress.generateAndStoreTokens(JobInProgress.java:3528) at org.apache.hadoop.mapred.JobInProgress.initTasks(JobInProgress.java:696) at org.apache.hadoop.mapred.JobTracker.initJob(JobTracker.java:4207) at org.apache.hadoop.mapred.FairScheduler\$JobInitializer\$InitJob.run(FairScheduler.java:291) at java.util.concurrent.ThreadPoolExecutor\$Worker.runTask(ThreadPoolExecutor.java:886) at java.util.concurrent.ThreadPoolExecutor\$Worker.run(ThreadPoolExecutor.java:908) at java.lang.Thread.run(Thread.java:662)</p>

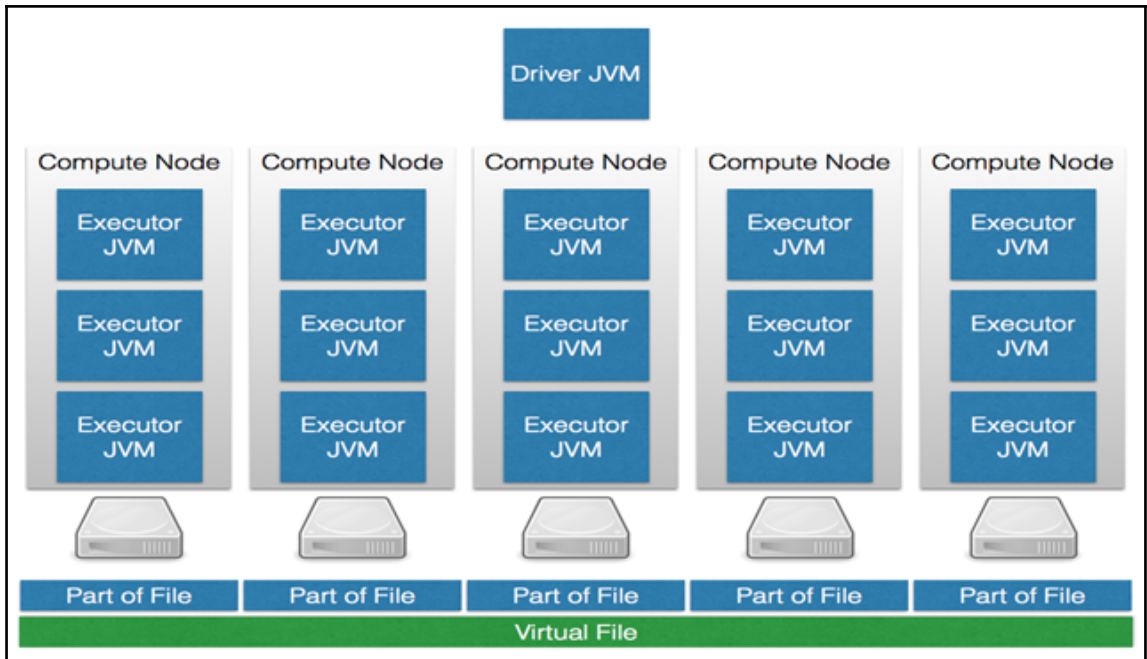
Task Index	Task ID	Status	Locality Level	Executor	Launch Time	Duration	GC Time
1	0	SUCCESS	NODE_LOCAL		2014/06/13 13:14:16	12.82 h	9.59 h
2	1	SUCCESS	NODE_LOCAL		2014/06/13 13:14:16	12.00 h	8.97 h
3	2	SUCCESS	NODE_LOCAL		2014/06/13 13:14:16	12.39 h	9.16 h
0	3	SUCCESS	NODE_LOCAL		2014/06/13 13:14:16	12.09 h	8.88 h
6	4	SUCCESS	NODE_LOCAL		2014/06/13 13:14:16	11.65 h	8.54 h
4	5	SUCCESS	NODE_LOCAL		2014/06/13 13:14:16	11.68 h	8.62 h
7	6	SUCCESS	NODE_LOCAL		2014/06/13 13:14:16	12.19 h	9.12 h
12	7	SUCCESS	NODE_LOCAL		2014/06/13 13:14:16	11.62 h	8.50 h
8	8	SUCCESS	NODE_LOCAL		2014/06/13 13:14:16	12.57 h	9.40 h
9	9	SUCCESS	NODE_LOCAL		2014/06/13 13:14:16	12.02 h	8.98 h
5	10	SUCCESS	NODE_LOCAL		2014/06/13 13:14:16	12.24 h	9.04 h
11	11	SUCCESS	NODE_LOCAL		2014/06/13 13:14:16	11.11 h	8.15 h
10	12	SUCCESS	NODE_LOCAL		2014/06/13 13:14:16	11.84 h	8.68 h
13	13	SUCCESS	NODE_LOCAL		2014/06/13 13:14:16	11.85 h	8.74 h
18	14	SUCCESS	NODE_LOCAL		2014/06/13 13:14:16	12.26 h	9.17 h

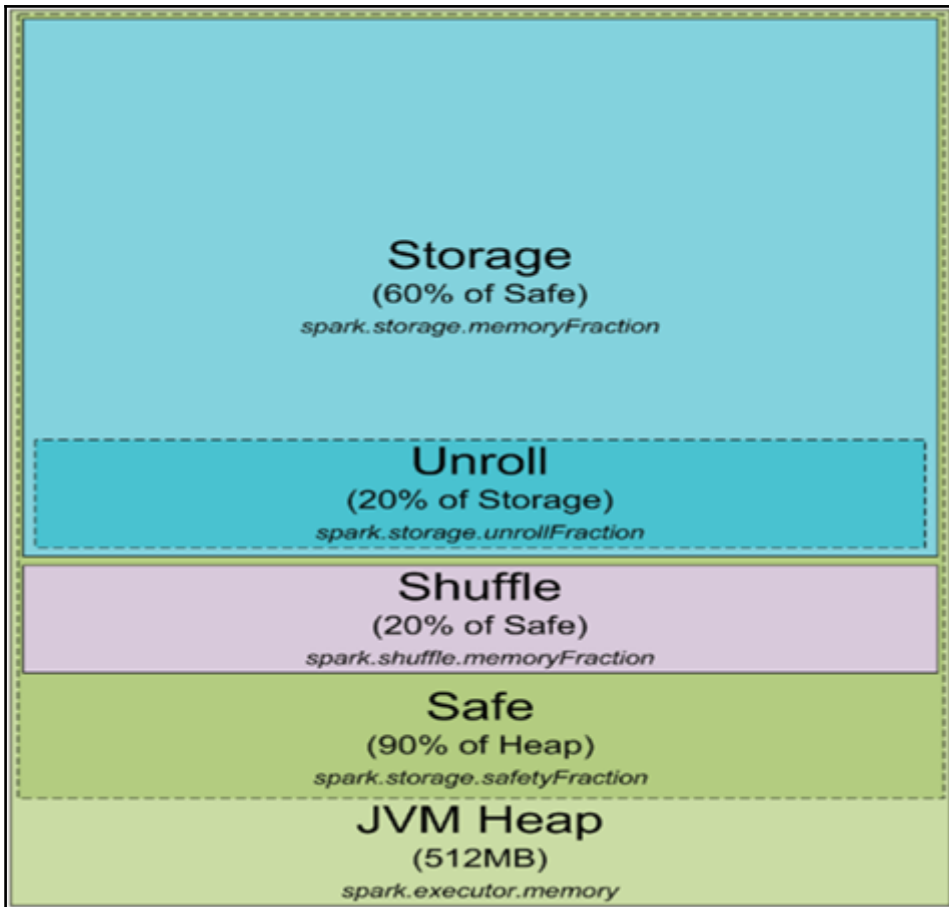


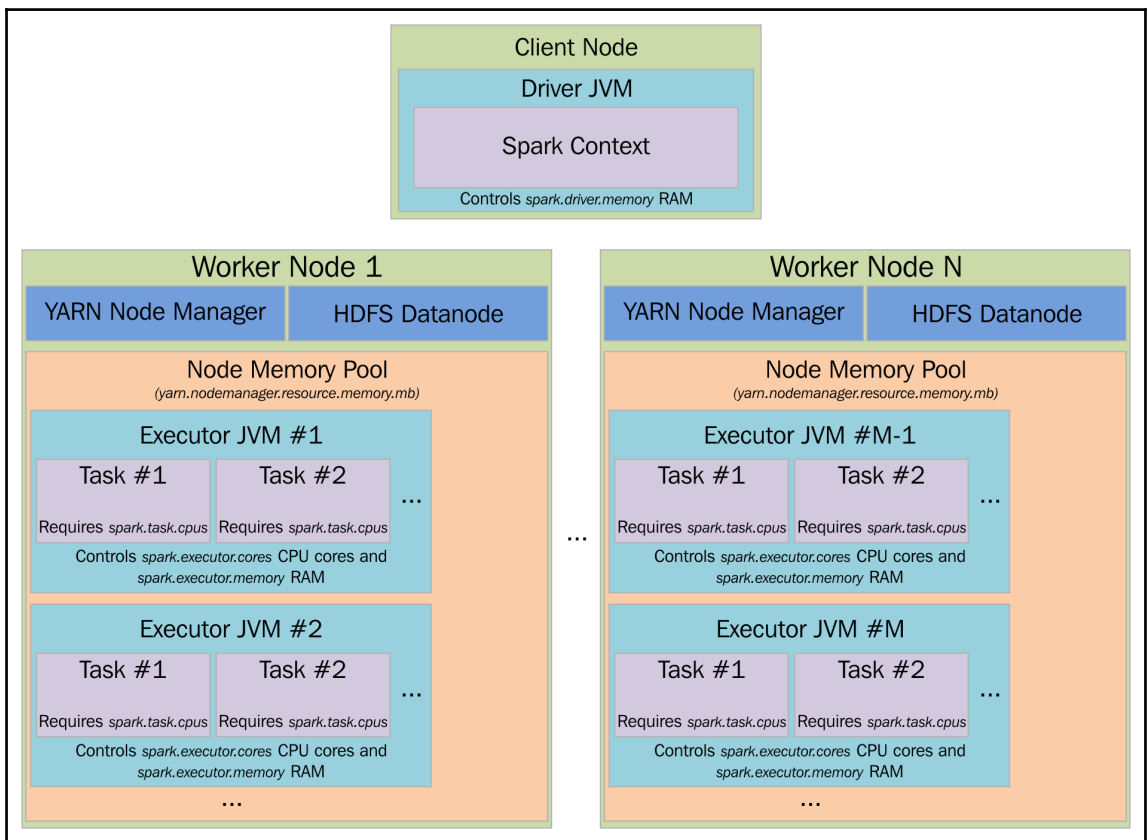
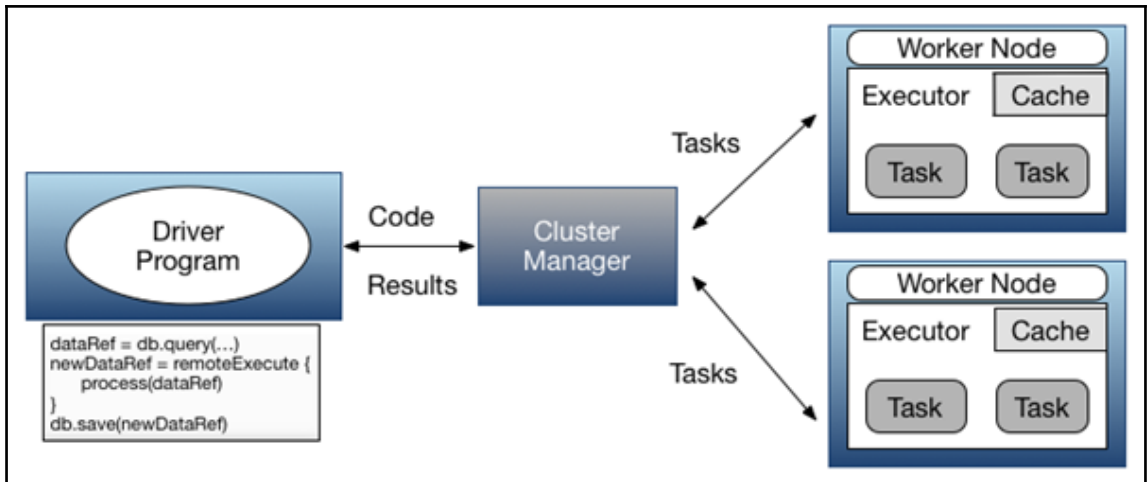


# Chapter 17: Time to Go to ClusterLand - Deploying Spark on a Cluster

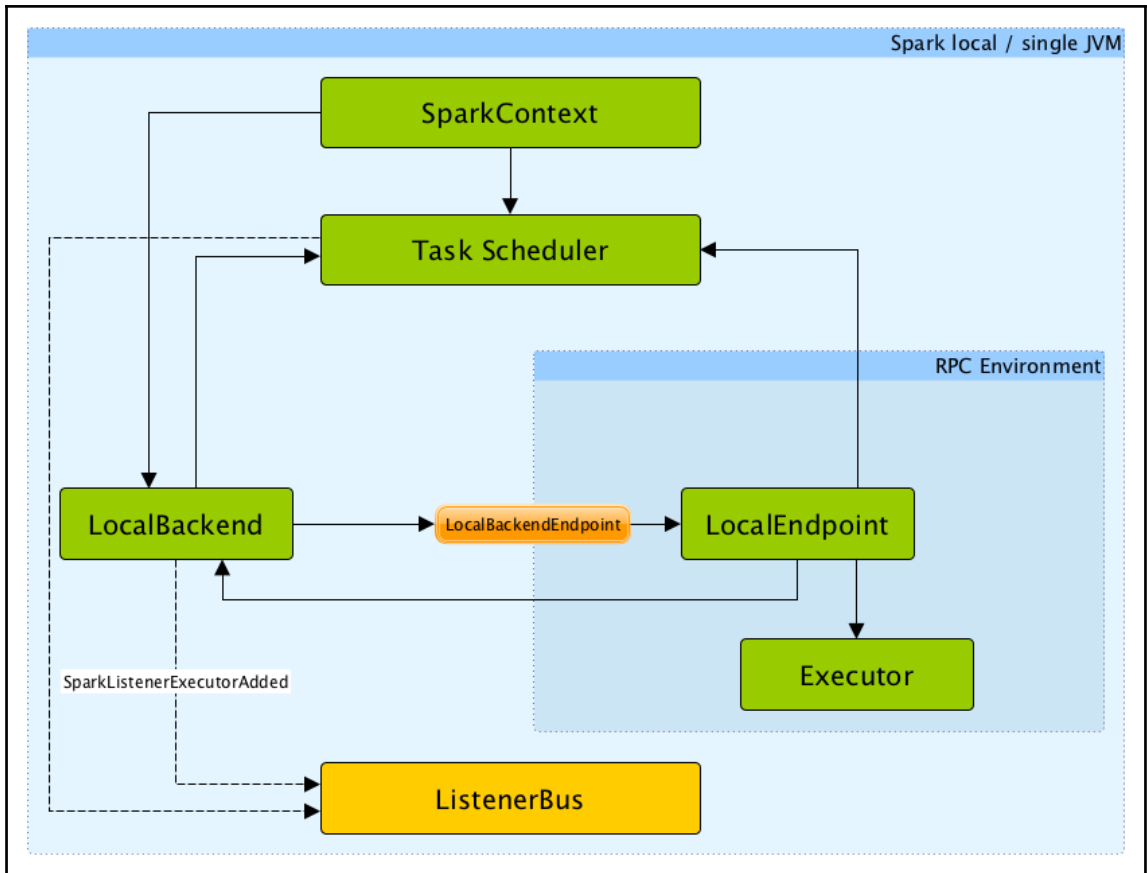


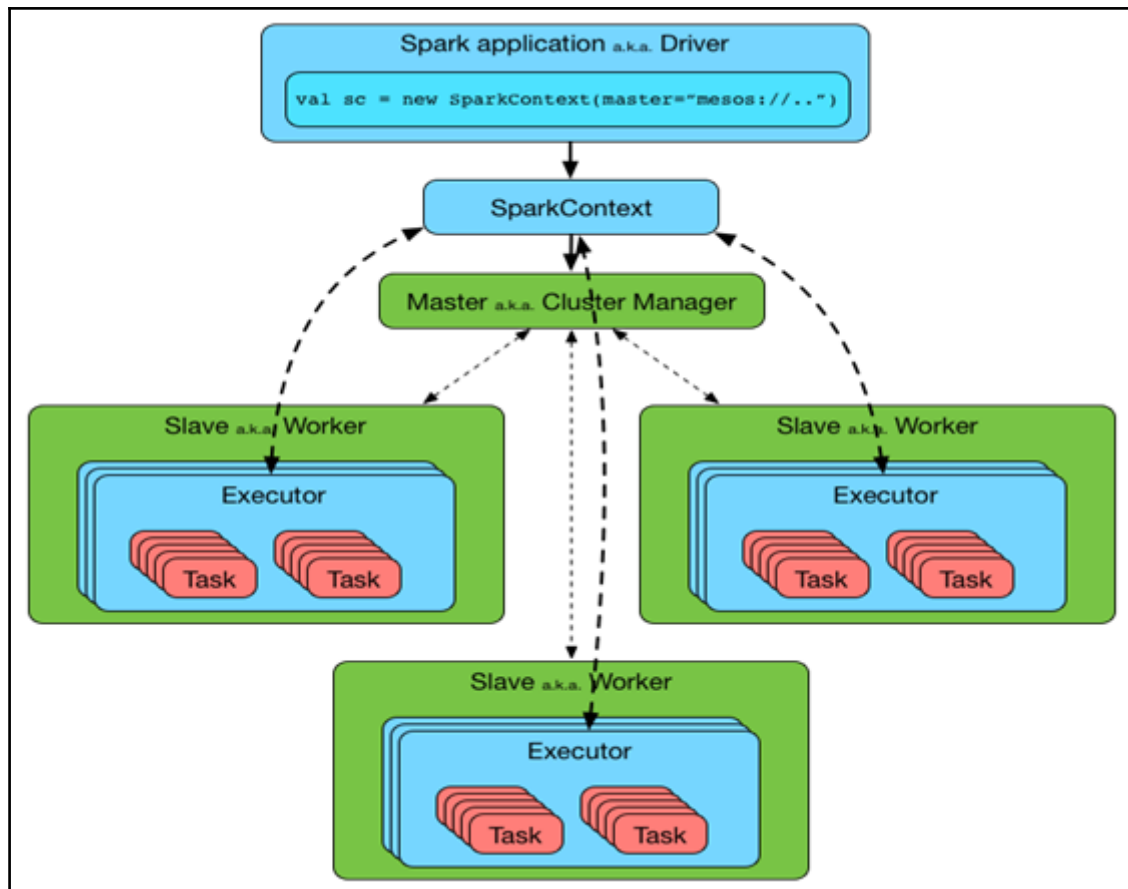













Term	Meaning
Application	User program built on Spark. Consists of a <i>driver program</i> and <i>executors</i> on the cluster.
Application jar	A jar containing the user's Spark application. In some cases users will want to create an "uber jar" containing their application along with its dependencies. The user's jar should never include Hadoop or Spark libraries, however, these will be added at runtime.
Driver program	The process running the main() function of the application and creating the SparkContext
Cluster manager	An external service for acquiring resources on the cluster (e.g. standalone manager, Mesos, YARN)
Deploy mode	Distinguishes where the driver process runs. In "cluster" mode, the framework launches the driver inside of the cluster. In "client" mode, the submitter launches the driver outside of the cluster.
Worker node	Any node that can run application code in the cluster
Executor	A process launched for an application on a worker node, that runs tasks and keeps data in memory or disk storage across them. Each application has its own executors.
Task	A unit of work that will be sent to one executor
Job	A parallel computation consisting of multiple tasks that gets spawned in response to a Spark action (e.g. save, collect); you'll see this term used in the driver's logs.
Stage	Each job gets divided into smaller sets of tasks called <i>stages</i> that depend on each other (similar to the map and reduce stages in MapReduce); you'll see this term used in the driver's logs.

Master URL	Meaning
local	Run Spark locally with one worker thread (i.e. no parallelism at all).
local [K]	Run Spark locally with K worker threads (ideally, set this to the number of cores on your machine).
local [*]	Run Spark locally with as many worker threads as logical cores on your machine.
spark://HOST:PORT	Connect to the given <a href="#">Spark standalone cluster</a> master. The port must be whichever one your master is configured to use, which is 7077 by default.
mesos://HOST:PORT	Connect to the given <a href="#">Mesos</a> cluster. The port must be whichever one your is configured to use, which is 5050 by default. Or, for a Mesos cluster using ZooKeeper, use mesos://zk://... To submit with --deploy-mode cluster, the HOST:PORT should be configured to connect to the <a href="#">MesosClusterDispatcher</a> .
yarn	Connect to a <a href="#">YARN</a> cluster in client or cluster mode depending on the value of --deploy-mode. The cluster location will be found based on the HADOOP_CONF_DIR or YARN_CONF_DIR variable.

```

17/02/14 12:31:02 INFO Executor: Finished task 0.0 in stage 0.0 (TID 0). 3543 bytes result sent to driver
17/02/14 12:31:02 INFO TaskSetManager: Finished task 0.0 in stage 0.0 (TID 0) in 215 ms on localhost (executor driver) (1/1)
17/02/14 12:31:02 INFO TaskSchedulerImpl: Removed TaskSet 0.0, whose tasks have all completed, from pool
17/02/14 12:31:02 INFO DAGScheduler: ResultStage 0 (show at KMeansDemo.scala:56) finished in 0.225 s
17/02/14 12:31:02 INFO DAGScheduler: Job 0 finished: show at KMeansDemo.scala:56, took 0.322031 s
17/02/14 12:31:02 INFO CodeGenerator: Code generated in 19.812394 ms
-----
| Price|LotSize|Waterfront| Age|LandValue|NewConstruct|CentralAir|FuelType|HeatType|SewerType|LivingArea|PctCollege|Bedrooms|Fireplaces|Bathrooms|rooms|
-----
132500.0| 0.09| 0.0| 42.0| 50000.0| 0.0| 0.0| 3.0| 4.0| 2.0| 906.0| 35.0| 2.0| 1.0| 1.0| 5.0|
181115.0| 0.92| 0.0| 0.0| 22300.0| 0.0| 0.0| 2.0| 3.0| 2.0| 1953.0| 51.0| 3.0| 0.0| 2.5| 6.0|
109000.0| 0.19| 0.0| 133.0| 7300.0| 0.0| 0.0| 2.0| 3.0| 3.0| 1944.0| 51.0| 4.0| 1.0| 1.0| 8.0|
155000.0| 0.41| 0.0| 13.0| 18700.0| 0.0| 0.0| 2.0| 2.0| 2.0| 1944.0| 51.0| 3.0| 1.0| 1.5| 5.0|
86060.0| 0.11| 0.0| 0.0| 15000.0| 1.0| 1.0| 2.0| 2.0| 3.0| 840.0| 51.0| 2.0| 0.0| 1.0| 3.0|
120000.0| 0.68| 0.0| 31.0| 14000.0| 0.0| 0.0| 2.0| 2.0| 2.0| 1152.0| 22.0| 4.0| 1.0| 1.0| 8.0|
153000.0| 0.4| 0.0| 33.0| 23300.0| 0.0| 0.0| 4.0| 3.0| 2.0| 2752.0| 51.0| 4.0| 1.0| 1.5| 8.0|
170000.0| 1.21| 0.0| 23.0| 14600.0| 0.0| 0.0| 4.0| 2.0| 2.0| 1662.0| 35.0| 4.0| 1.0| 1.5| 9.0|
90000.0| 0.83| 0.0| 36.0| 22200.0| 0.0| 0.0| 3.0| 4.0| 2.0| 1632.0| 51.0| 3.0| 0.0| 1.5| 8.0|
122900.0| 1.94| 0.0| 4.0| 21200.0| 0.0| 0.0| 2.0| 2.0| 1.0| 1416.0| 44.0| 3.0| 0.0| 1.5| 6.0|
325000.0| 2.29| 0.0| 123.0| 12600.0| 0.0| 0.0| 4.0| 2.0| 2.0| 2894.0| 51.0| 7.0| 0.0| 1.0| 12.0|
120000.0| 0.92| 0.0| 1.0| 22300.0| 0.0| 0.0| 2.0| 2.0| 2.0| 1624.0| 51.0| 3.0| 0.0| 2.0| 6.0|
85860.0| 8.97| 0.0| 13.0| 4800.0| 0.0| 0.0| 3.0| 4.0| 2.0| 704.0| 41.0| 2.0| 0.0| 1.0| 4.0|
97000.0| 0.11| 0.0| 153.0| 3100.0| 0.0| 0.0| 2.0| 3.0| 3.0| 1383.0| 57.0| 3.0| 0.0| 2.0| 5.0|
127000.0| 0.14| 0.0| 9.0| 300.0| 0.0| 0.0| 4.0| 2.0| 2.0| 1300.0| 41.0| 3.0| 0.0| 1.5| 8.0|
89900.0| 0.0| 0.0| 88.0| 2500.0| 0.0| 0.0| 2.0| 3.0| 3.0| 936.0| 57.0| 3.0| 0.0| 1.0| 4.0|
155000.0| 0.13| 0.0| 9.0| 300.0| 0.0| 0.0| 4.0| 2.0| 2.0| 1300.0| 41.0| 3.0| 0.0| 1.5| 7.0|
253750.0| 2.0| 0.0| 0.0| 49900.0| 0.0| 1.0| 2.0| 2.0| 1.0| 2816.0| 71.0| 4.0| 1.0| 2.5| 12.0|
60000.0| 0.21| 0.0| 82.0| 8500.0| 0.0| 0.0| 4.0| 3.0| 2.0| 924.0| 35.0| 2.0| 0.0| 1.0| 6.0|
87500.0| 0.88| 0.0| 17.0| 19400.0| 0.0| 0.0| 4.0| 2.0| 2.0| 1092.0| 35.0| 3.0| 0.0| 1.0| 6.0|
-----
only showing top 20 rows

17/02/14 12:31:02 INFO ContextCleaner: Cleaned accumulator 3
17/02/14 12:31:03 INFO BlockManagerInfo: Removed broadcast_1_piece0 on 10.2.16.255:53581 in memory (size: 9.4 KB, free: 4.0 GB)
17/02/14 12:31:03 INFO SparkContext: Starting job: takeSample at KMeans.scala:353
17/02/14 12:31:03 INFO DAGScheduler: Got job 1 (takeSample at KMeans.scala:353) with 2 output partitions
17/02/14 12:31:03 INFO DAGScheduler: Final stage: ResultStage 1 (takeSample at KMeans.scala:353)
17/02/14 12:31:03 INFO DAGScheduler: Parents of final stage: List()
    
```



## Spark Master at spark://ubuntu:7077

**URL:** spark://ubuntu:7077  
**REST URL:** spark://ubuntu:6066 (cluster mode)  
**Alive Workers:** 0  
**Cores in use:** 0 Total, 0 Used  
**Memory in use:** 0.0 B Total, 0.0 B Used  
**Applications:** 0 Running, 0 Completed  
**Drivers:** 0 Running, 0 Completed  
**Status:** ALIVE

**Workers**


Worker Id	Address	State

**Running Applications**

Application ID	Name	Cores	Memory per Node	Submitted Time

**Completed Applications**

Application ID	Name	Cores	Memory per Node	Submitted Time

 **Spark Worker at 192.168.12.129:35079**

ID: worker-20170214044222-192.168.12.129-35079  
Master URL: spark://ubuntu:7077  
Cores: 1 (0 Used)  
Memory: 1024.0 MB (0.0 B Used)

[Back to Master](#)

**Running Executors (0)**

ExecutorID	Cores	State	Memory
------------	-------	-------	--------

 **Spark Master at spark://ubuntu:7077**

URL: spark://ubuntu:7077  
REST URL: spark://ubuntu:6066 (cluster mode)  
Alive Workers: 1  
Cores in use: 1 Total, 0 Used  
Memory in use: 1024.0 MB Total, 0.0 B Used  
Applications: 0 [Running](#), 0 [Completed](#)  
Drivers: 0 [Running](#), 0 [Completed](#)  
Status: ALIVE

**Workers**

Worker Id	Address	State
<a href="#">worker-20170214044222-192.168.12.129-35079</a>	192.168.12.129:35079	ALIVE

**Running Applications**

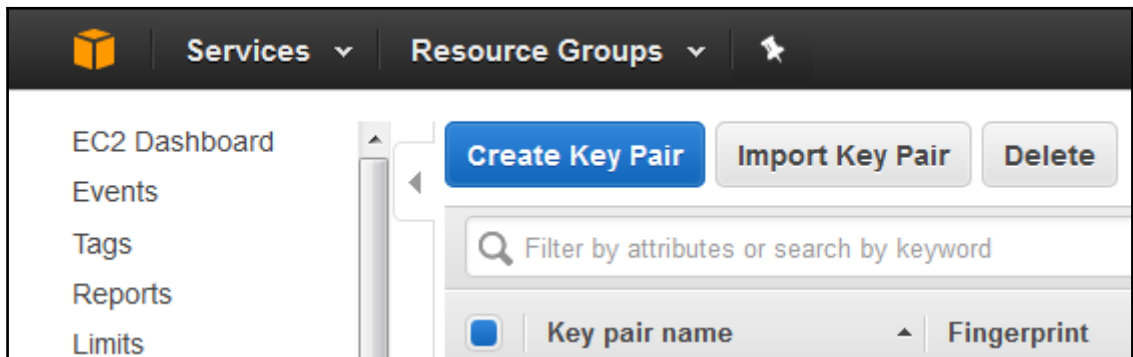
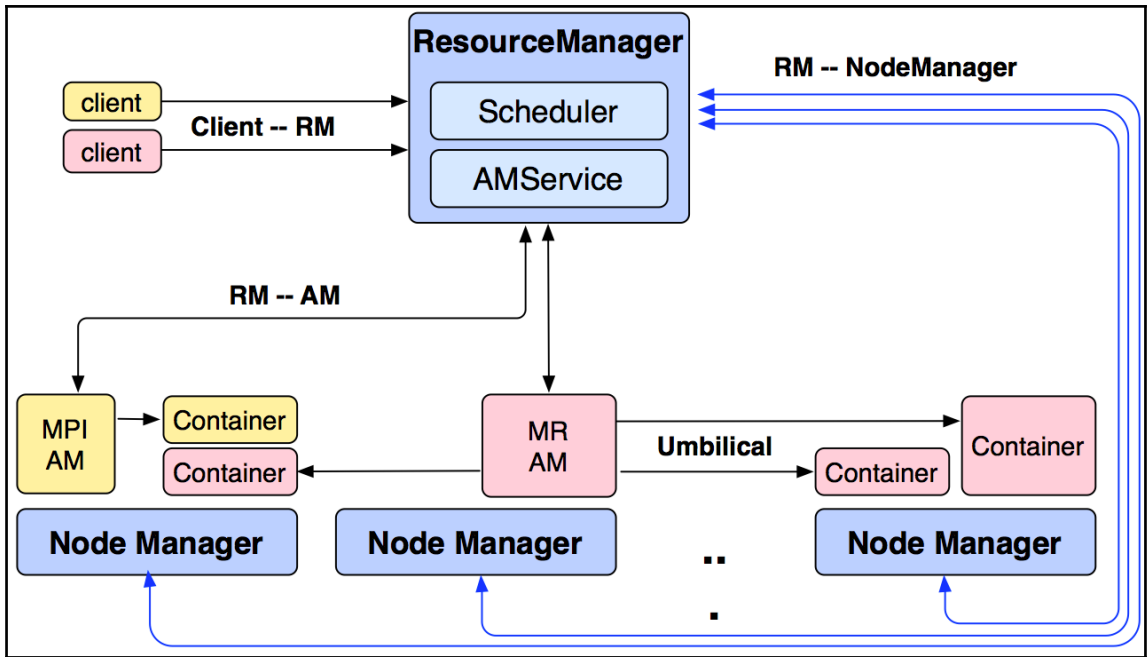
Application ID	Name	Cores	Memory per Node	Submitted Time
----------------	------	-------	-----------------	----------------

**Completed Applications**

Application ID	Name	Cores	Memory per Node	Submitted Time
----------------	------	-------	-----------------	----------------

Argument	Meaning
-h HOST, --host HOST	Hostname to listen on
-i HOST, --ip HOST	Hostname to listen on (deprecated, use -h or --host)
-p PORT, --port PORT	Port for service to listen on (default: 7077 for master, random for worker)
--webui-port PORT	Port for web UI (default: 8080 for master, 8081 for worker)
-c CORES, --cores CORES	Total CPU cores to allow Spark applications to use on the machine (default: all available); only on worker
-m MEM, --memory MEM	Total amount of memory to allow Spark applications to use on the machine, in a format like 1000M or 2G (default: your machine's total RAM minus 1 GB); only on worker
-d DIR, --work-dir DIR	Directory to use for scratch space and job output logs (default: SPARK_HOME/work); only on worker
--properties-file FILE	Path to a custom Spark properties file to load (default: conf/spark-defaults.conf)

- `sbin/start-master.sh` - Starts a master instance on the machine the script is executed on.
- `sbin/start-slaves.sh` - Starts a slave instance on each machine specified in the `conf/slaves` file.
- `sbin/start-slave.sh` - Starts a slave instance on the machine the script is executed on.
- `sbin/start-all.sh` - Starts both a master and a number of slaves as described above.
- `sbin/stop-master.sh` - Stops the master that was started via the `bin/start-master.sh` script.
- `sbin/stop-slaves.sh` - Stops all slave instances on the machines specified in the `conf/slaves` file.
- `sbin/stop-all.sh` - Stops both the master and the slaves as described above.



The screenshot displays the Spark Master web interface for a cluster. At the top, the browser address bar shows the URL: `ec2-52-19-229-38.eu-west-1.compute.amazonaws.com:8080`. The page title is "Spark Master at spark://ec2-52-19-229-38.eu-west-1.compute.amazonaws.com:7077".

Key information displayed includes:

- URL:** `spark://ec2-52-19-229-38.eu-west-1.compute.amazonaws.com:7077`
- REST URL:** `spark://ec2-52-19-229-38.eu-west-1.compute.amazonaws.com:6066 (cluster mode)`
- Alive Workers:** 2
- Cores in use:** 16 Total, 0 Used
- Memory in use:** 56.4 GB Total, 0.0 B Used
- Applications:** 0 Running, 0 Completed
- Drivers:** 0 Running, 0 Completed
- Status:** ALIVE

**Workers**

Worker Id	Address	State	Cores	Memory
<a href="#">worker-20160325211558-172.31.8.198-45995</a>	172.31.8.198:45995	ALIVE	8 (0 Used)	28.2 GB (0.0 B Used)
<a href="#">worker-20160325211601-172.31.8.197-44556</a>	172.31.8.197:44556	ALIVE	8 (0 Used)	28.2 GB (0.0 B Used)

**Running Applications**

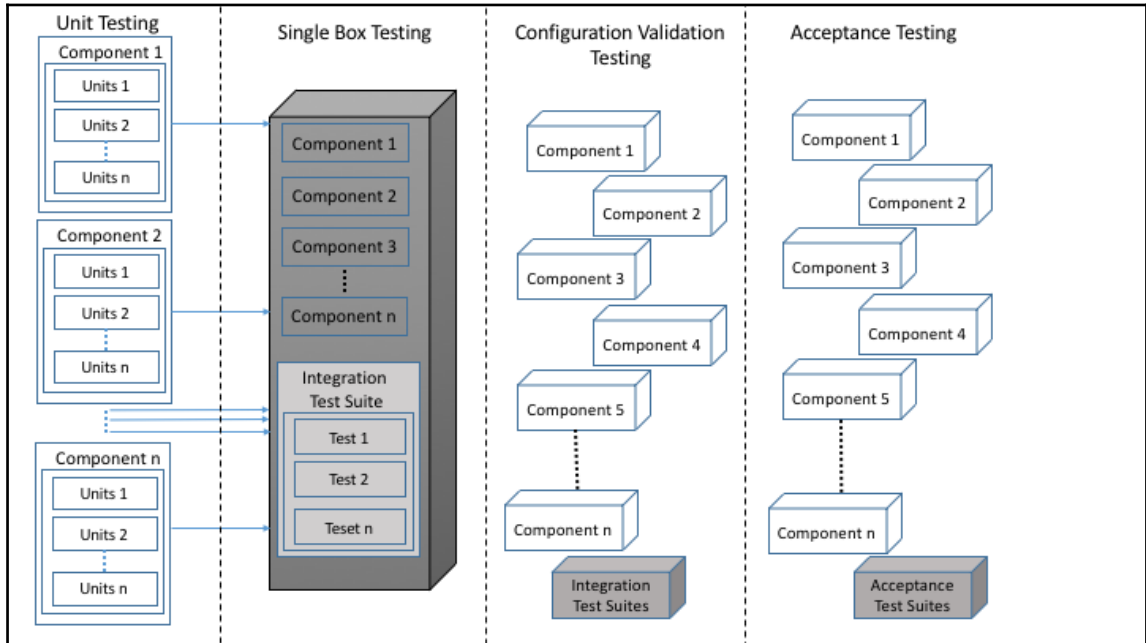
Application ID	Name	Cores	Memory per Node	Submitted Time	User	State	Duration
----------------	------	-------	-----------------	----------------	------	-------	----------

**Completed Applications**

Application ID	Name	Cores	Memory per Node	Submitted Time	User	State	Duration
----------------	------	-------	-----------------	----------------	------	-------	----------



# Chapter 18: Testing and Debugging Spark



```
Exception in thread "main" java.lang.AssertionError: assertion failed
    at scala.Predef$.assert(Predef.scala:156)
    at com.chapter16.SparkTesting.SimpleScalaTest$.main(SimpleScalaTest.scala:7)
    at com.chapter16.SparkTesting.SimpleScalaTest.main(SimpleScalaTest.scala)
```

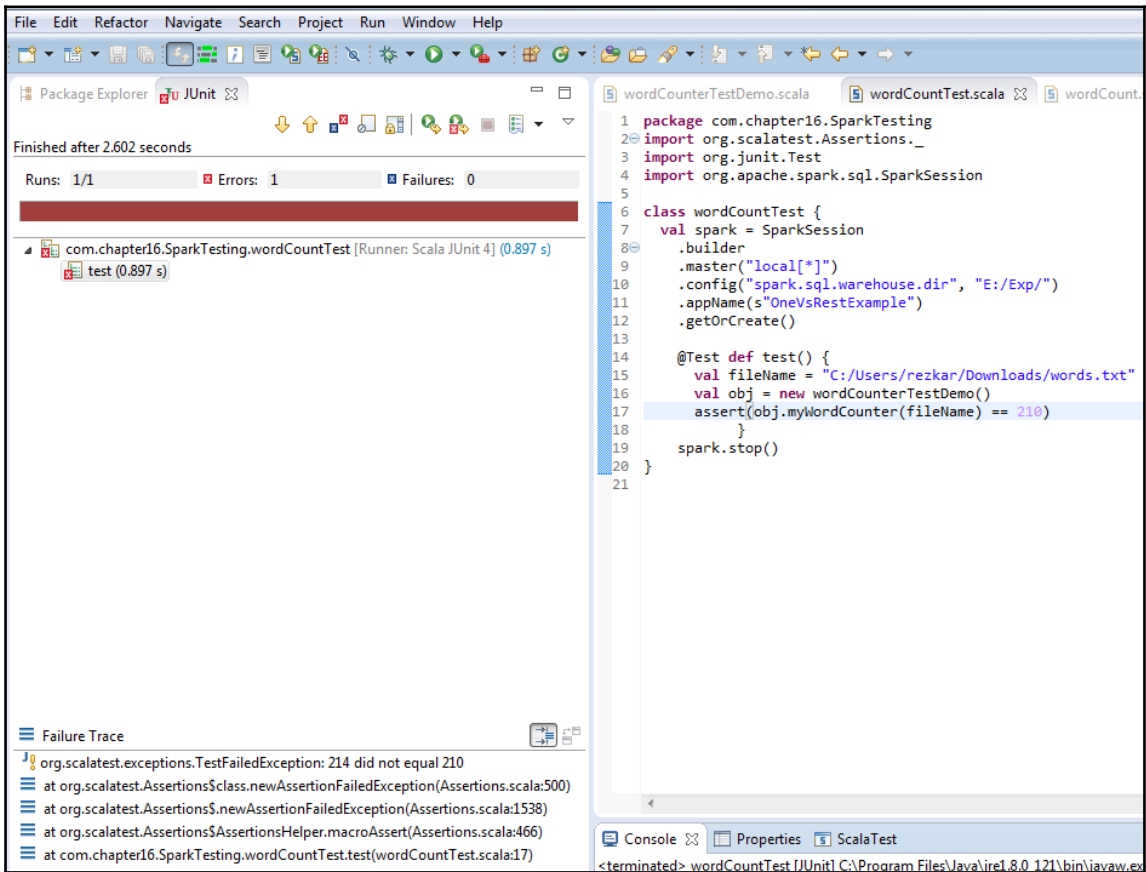
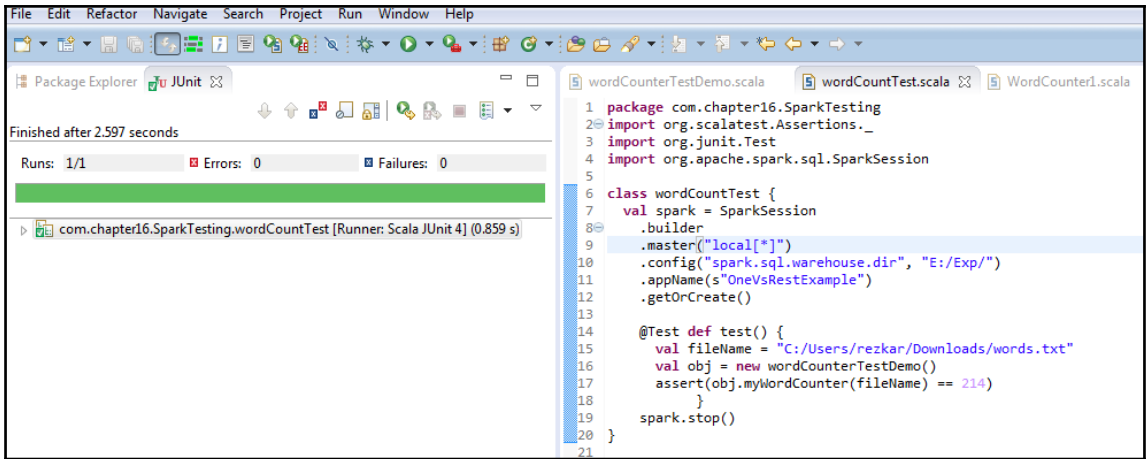
```
Exception in thread "main" org.scalatest.exceptions.TestFailedException: 2 did not equal 1
    at org.scalatest.Assertions$class.newAssertionFailedException(Assertions.scala:500)
    at org.scalatest.Assertions$.newAssertionFailedException(Assertions.scala:1538)
    at org.scalatest.Assertions$AssertionsHelper.macroAssert(Assertions.scala:466)
    at com.chapter16.SparkTesting.SimpleScalaTest$.main(SimpleScalaTest.scala:8)
    at com.chapter16.SparkTesting.SimpleScalaTest.main(SimpleScalaTest.scala)
```

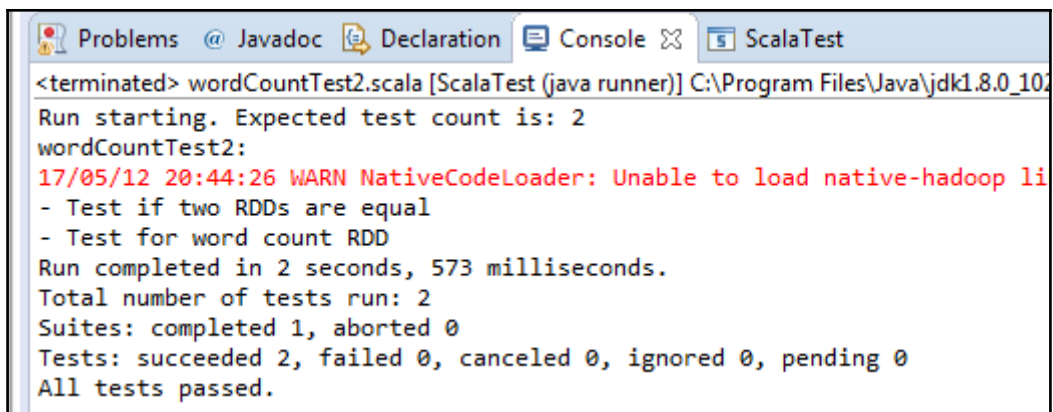
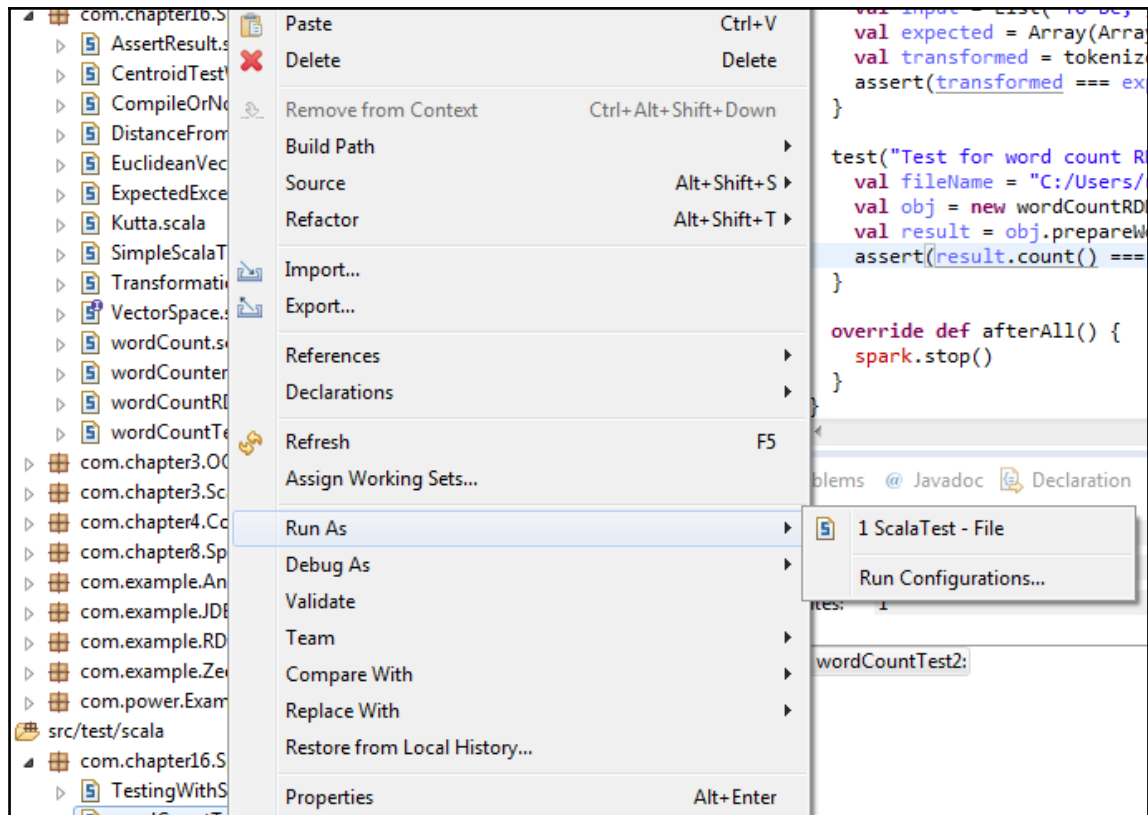
```
Exception in thread "main" org.scalatest.exceptions.TestFailedException: Expected 3, but
got 4
  at org.scalatest.Assertions$class.newAssertionFailedException(Assertions.scala:495)
  at org.scalatest.Assertions$.newAssertionFailedException(Assertions.scala:1538)
  at org.scalatest.Assertions$class.assertResult(Assertions.scala:1226)
  at org.scalatest.Assertions$.assertResult(Assertions.scala:1538)
  at com.chapter16.SparkTesting.AssertResult$.main(AssertResult.scala:8)
  at com.chapter16.SparkTesting.AssertResult.main(AssertResult.scala)
```

```
Exception in thread "main" org.scalatest.exceptions.TestFailedException
  at org.scalatest.Assertions$class.newAssertionFailedException(Assertions.scala:493)
  at org.scalatest.Assertions$.newAssertionFailedException(Assertions.scala:1538)
  at org.scalatest.Assertions$class.fail(Assertions.scala:1313)
  at org.scalatest.Assertions$.fail(Assertions.scala:1538)
  at com.chapter16.SparkTesting.ExpectedException$.main(ExpectedException.scala:9)
  at com.chapter16.SparkTesting.ExpectedException.main(ExpectedException.scala)
```

```
AssertDoesNotCompile True
AssertTypeError True
AssertCompiles True
Exception in thread "main" org.scalatest.exceptions.TestFailedException: Expected a
compiler error, but got none for code: val a: Int = 1
  at com.chapter16.SparkTesting.CompileOrNot$.main(CompileOrNot.scala:15)
  at com.chapter16.SparkTesting.CompileOrNot.main(CompileOrNot.scala)
```







```
Problems @ Javadoc Declaration Console ScalaTest
<terminated> wordCountTest2.scala [ScalaTest (java runner)] C:\Program Files\Java\jdk1.8.0_102\bin\javaw.exe (12 May 2017, 20:47:04)
Run starting. Expected test count is: 2
wordCountTest2:
17/05/12 20:47:05 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable

[Stage 0:>] (0 + 0) / 8]
- Test if two RDDs are equal *** FAILED ***
  Array(Array("To", "be:"), Array("on", "not", "to", "be:"), Array("that", "is", "the", "question-"), Array("William", "Shakespeare")) did not equal Array(Array("To", "be")
- Test for word count RDD *** FAILED ***
  214 did not equal 210 (wordCountTest2.scala:33)
Run completed in 2 seconds, 749 milliseconds.
Total number of tests run: 2
Suites: completed 1, aborted 0
Tests: succeeded 0, failed 2, canceled 0, ignored 0, pending 0
*** 2 TESTS FAILED ***
```

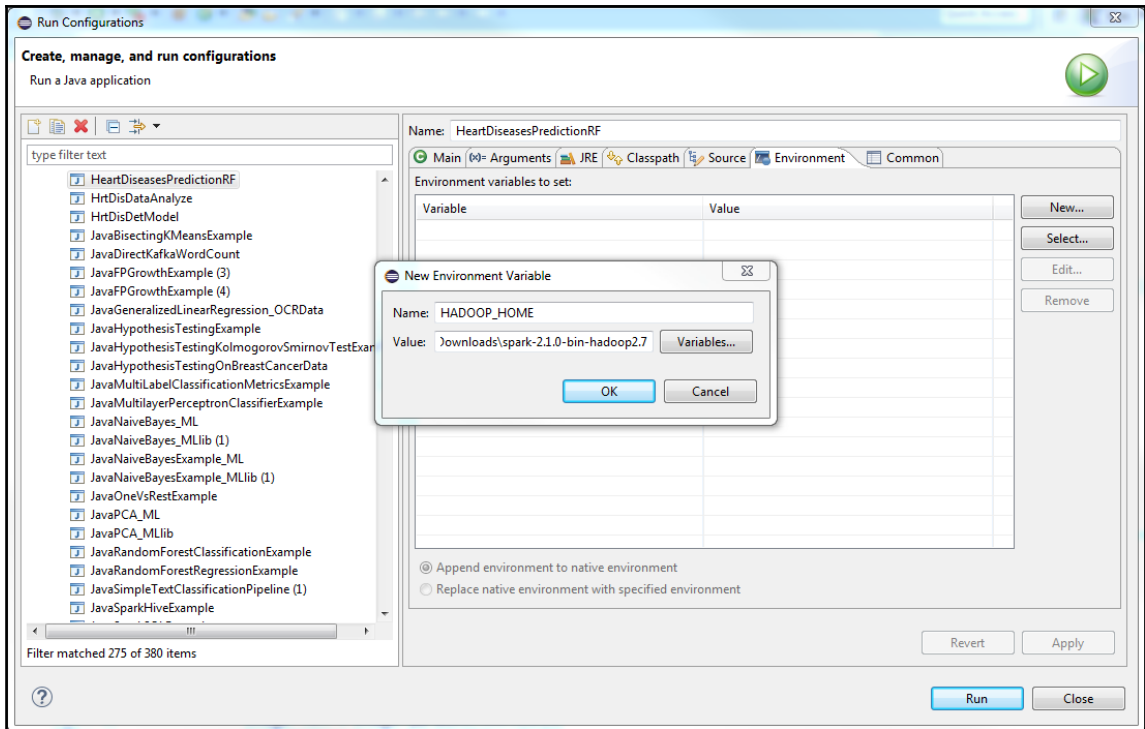
```
Console Problems
<terminated> TransformationTestWithSparkTestingBase.scala [ScalaTest (java runner)] C:\Program Files\Java\jdk1.8.0_102\bin\javaw.exe (19 Feb 2017, 00:49:48)
Run starting. Expected test count is: 3
TransformationTestWithSparkTestingBase:
17/02/19 00:49:49 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
- works, obviously!

[Stage 0:>] (0 + 0) / 8] - broken
- Testing RDD transformations using a shared Spark Context
Run completed in 2 seconds, 534 milliseconds.
Total number of tests run: 3
Suites: completed 1, aborted 0
Tests: succeeded 3, failed 0, canceled 0, ignored 0, pending 0
All tests passed.
```

```
Console Problems ScalaTest
Tests: 3/3 Succeeded: 3 Failed: 0
Ignored: 0 Pending: 0 Canceled: 0
Suites: 1 Aborted: 0

> TransformationTestWithSparkTestingBase (2.179 s) Stack Trace
```

```
17/02/26 13:22:00 ERROR Shell: Failed to locate the winutils binary in the hadoop binary path
java.io.IOException: Could not locate executable null\bin\winutils.exe in the Hadoop binaries.
    at org.apache.hadoop.util.Shell.getQualifiedBinPath(Shell.java:278)
    at org.apache.hadoop.util.Shell.getWinUtilsPath(Shell.java:300)
    at org.apache.hadoop.util.Shell.<clinit>(Shell.java:293)
    at org.apache.hadoop.util.StringUtils.<clinit>(StringUtils.java:76)
    at org.apache.hadoop.mapred.FileInputFormat.setInputPaths(FileInputFormat.java:362)
    at org.apache.spark.SparkContext$$anonfun$hadoopFile$1$$anonfun$30.apply(SparkContext.scala:1014)
    at org.apache.spark.SparkContext$$anonfun$hadoopFile$1$$anonfun$30.apply(SparkContext.scala:1014)
    at org.apache.spark.rdd.HadoopRDD$$anonfun$getJobConf$6.apply(HadoopRDD.scala:179)
    at org.apache.spark.rdd.HadoopRDD$$anonfun$getJobConf$6.apply(HadoopRDD.scala:179)
    at scala.Option.foreach(Option.scala:257)
    at org.apache.spark.rdd.HadoopRDD.getJobConf(HadoopRDD.scala:179)
    at org.apache.spark.rdd.HadoopRDD.getPartitions(HadoopRDD.scala:198)
    at org.apache.spark.rdd.RDD$$anonfun$partitions$2.apply(RDD.scala:252)
    at org.apache.spark.rdd.RDD$$anonfun$partitions$2.apply(RDD.scala:250)
    at scala.Option.getOrElse(Option.scala:121)
    at org.apache.spark.rdd.RDD.partitions(RDD.scala:250)
    at org.apache.spark.rdd.MapPartitionsRDD.getPartitions(MapPartitionsRDD.scala:35)
    at org.apache.spark.rdd.RDD$$anonfun$partitions$2.apply(RDD.scala:252)
    at org.apache.spark.rdd.RDD$$anonfun$partitions$2.apply(RDD.scala:250)
    at scala.Option.getOrElse(Option.scala:121)
    at org.apache.spark.rdd.RDD.partitions(RDD.scala:250)
    at org.apache.spark.rdd.MapPartitionsRDD.getPartitions(MapPartitionsRDD.scala:35)
    at org.apache.spark.rdd.RDD$$anonfun$partitions$2.apply(RDD.scala:252)
    at org.apache.spark.rdd.RDD$$anonfun$partitions$2.apply(RDD.scala:250)
    at scala.Option.getOrElse(Option.scala:121)
    at org.apache.spark.rdd.RDD.partitions(RDD.scala:250)
    at org.apache.spark.rdd.MapPartitionsRDD.getPartitions(MapPartitionsRDD.scala:35)
    at org.apache.spark.rdd.RDD$$anonfun$partitions$2.apply(RDD.scala:252)
    at org.apache.spark.rdd.RDD$$anonfun$partitions$2.apply(RDD.scala:250)
    at scala.Option.getOrElse(Option.scala:121)
    at org.apache.spark.rdd.RDD.partitions(RDD.scala:250)
    at org.apache.spark.rdd.MapPartitionsRDD.getPartitions(MapPartitionsRDD.scala:35)
    at org.apache.spark.rdd.RDD$$anonfun$partitions$2.apply(RDD.scala:252)
    at org.apache.spark.rdd.RDD$$anonfun$partitions$2.apply(RDD.scala:250)
    at scala.Option.getOrElse(Option.scala:121)
```



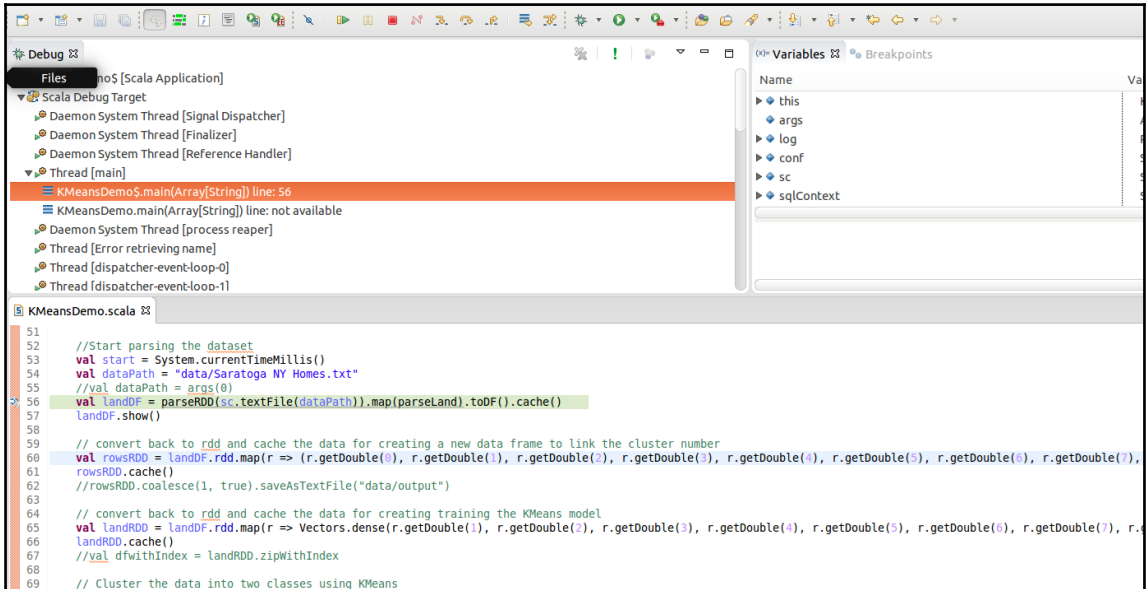
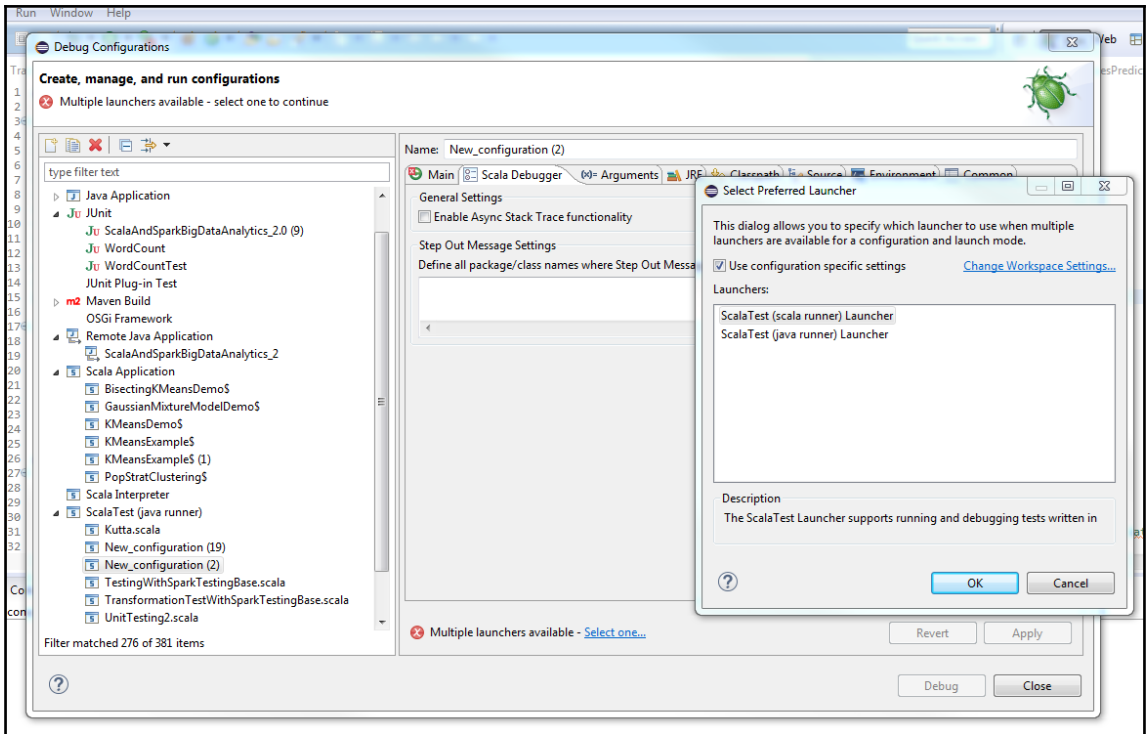
```
# Set everything to be logged to the console
log4j.rootCategory=INFO, console
log4j.appender.console=org.apache.log4j.ConsoleAppender
log4j.appender.console.target=System.err
log4j.appender.console.layout=org.apache.log4j.PatternLayout
log4j.appender.console.layout.ConversionPattern=%d{yy/MM/dd HH:mm:ss} %p %c{1}: %m%n

# Set the default spark-shell log level to WARN. When running the spark-shell, the
# log level for this class is used to overwrite the root logger's log level, so that
# the user can have different defaults for the shell and regular Spark apps.
log4j.logger.org.apache.spark.repl.Main=WARN

# Settings to quiet third party logs that are too verbose
log4j.logger.org.spark_project.jetty=WARN
log4j.logger.org.spark_project.jetty.util.component.AbstractLifeCycle=ERROR
log4j.logger.org.apache.spark.repl.SparkIMain$exprTyper=INFO
log4j.logger.org.apache.spark.repl.SparkILoop$SparkILoopInterpreter=INFO
log4j.logger.org.apache.parquet=ERROR
log4j.logger.parquet=ERROR

# SPARK-9183: Settings to avoid annoying messages when looking up nonexistent UDFs in SparkSQL with Hive support
log4j.logger.org.apache.hadoop.hive.metastore.RetryingHMSHandler=FATAL
log4j.logger.org.apache.hadoop.hive.ql.exec.FunctionRegistry=ERROR
```

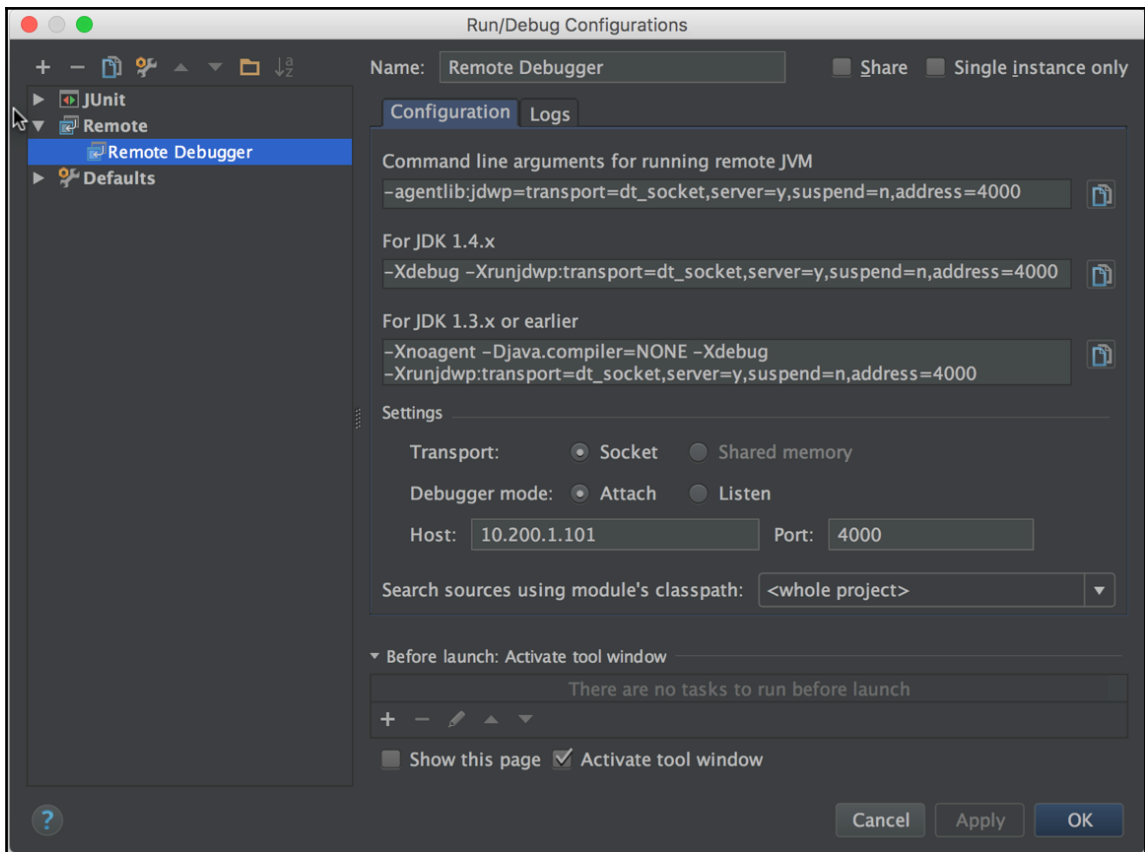


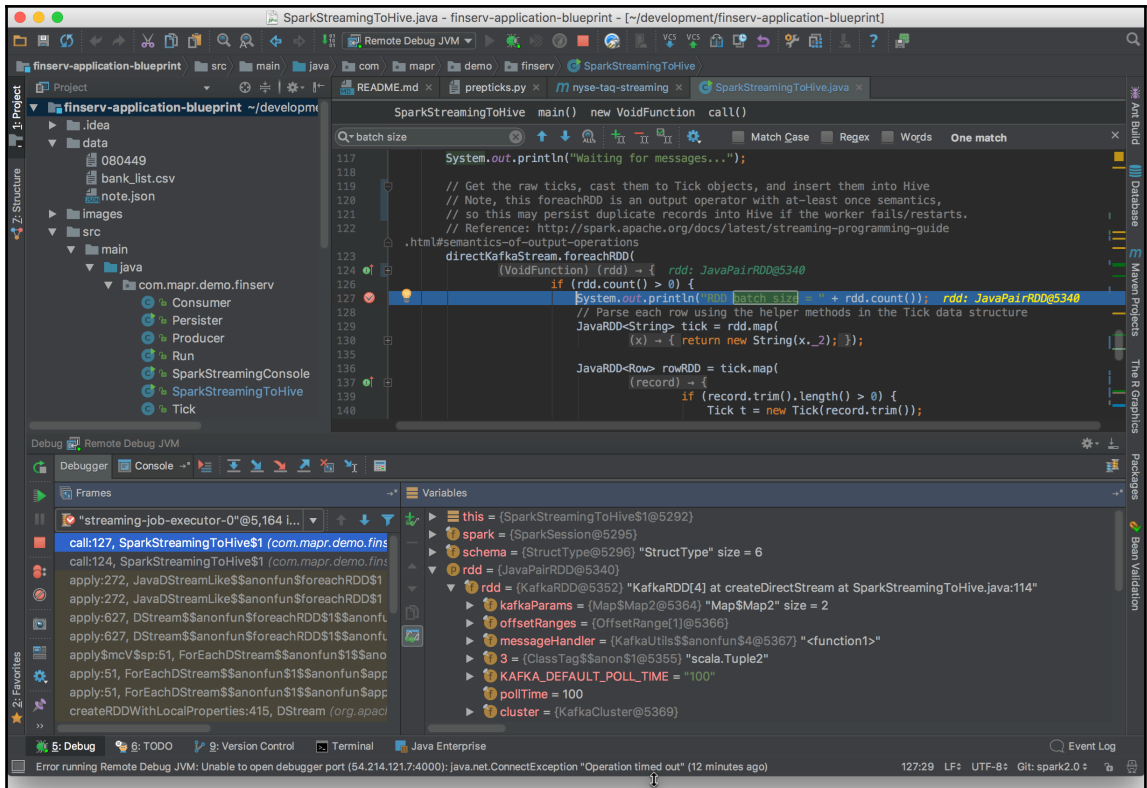


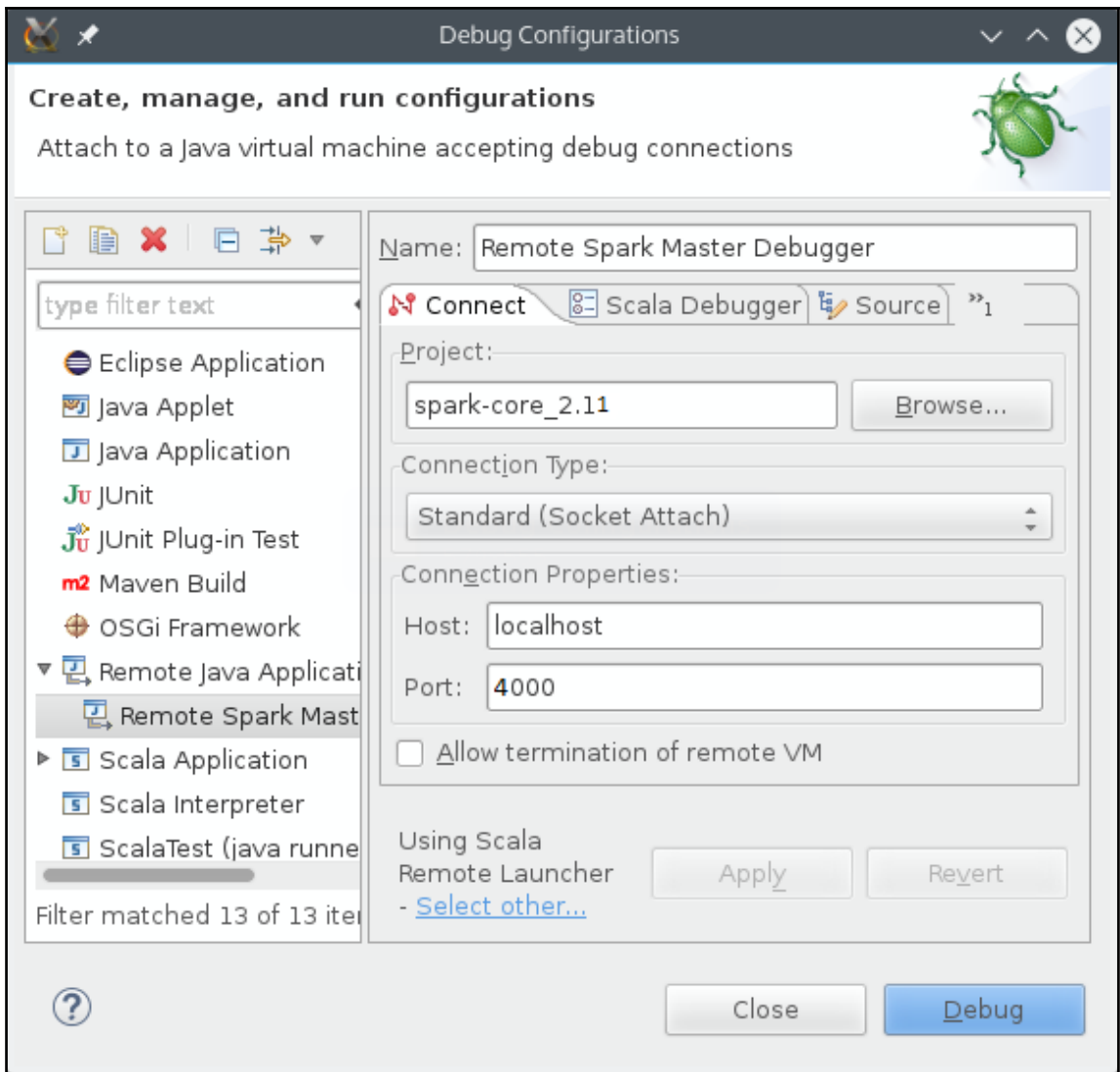
```
<terminated>KMeansDemo$ [Scala Application]
<terminated, exit value: 0>/usr/lib/jvm/java-8-oracle/bin/java (Feb 19, 2017, 12:34:48 PM)

KMeansDemo.scala
81 val end = System.currentTimeMillis()
82 println("Model building and prediction time: " + {end - start} + "ms")
83
84 // Compute and print the prediction accuracy for each house
85 model.predict(landRDD).foreach(println)
86 landDF.show()
87
88 // Get the prediction from the model with the ID so we can link them back to other information
89 val predictions = rowsRDD.map{r => (r._1, model.predict(Vectors.dense(r._2, r._3, r._4, r._5, r._6, r._7, r._8))}
90 val conMat = predictions.collect().toMap.values
91 println(conMat)
92
93
94 // convert the rdd to a dataframe
95 val predCluster = predictions.toDF("Price", "CLUSTER")
96 predCluster.show()
97
98
99 // Join the prediction DataFrame with the old dataframe to know the individual cluster number for each house
100 val newDF = landDF.join(predCluster, "Price")

Console
<terminated> KMeansDemo$ [Scala Application] /usr/lib/jvm/java-8-oracle/bin/java (Feb 19, 2017, 12:34:48 PM)
|253750.0|      3|
| 60000.0|      2|
| 87500.0|      2|
+-----+
only showing top 20 rows
MapLike(2, 2, 2, 1, 0, 2, 0, 0, 0, 0, 2, 1, 3, 3, 2, 3, 0, 2, 2, 0, 3, 2, 2, 2, 1, 0, 0, 0, 3, 2, 3, 3, 3, 2, 0, 1, 3,
17/02/19 12:35:09 WARN root: Finished
```







# Chapter 19: PySpark and SparkR

```
asif@ubuntu:~$ cd $SPARK_HOME
asif@ubuntu:~/Spark$ ./bin/pyspark
Python 2.7.6 (default, Oct 26 2016, 20:30:19)
[GCC 4.8.4] on linux2
Type "help", "copyright", "credits" or "license" for more information.
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
17/04/24 09:49:02 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using
17/04/24 09:49:02 WARN Utils: Your hostname, ubuntu resolves to a loopback address: 127.0.1.1; using 19:
17/04/24 09:49:02 WARN Utils: Set SPARK_LOCAL_IP if you need to bind to another address
17/04/24 09:49:06 WARN ObjectStore: Failed to get database global_temp, returning NoSuchObjectException
Welcome to

  ____      _
 / ___|    / \
| |  | |  / _ \
| |  | | / ___ \
| |  | | \___/ \
| |  | |
| |  | |
|_|  |_|

 version 2.1.0

Using Python version 2.7.6 (default, Oct 26 2016 20:30:19)
SparkSession available as 'spark'.
>>> █
```

```
asif@ubuntu:~$ cd $SPARK_HOME
asif@ubuntu:~/Spark$ ./bin/pyspark
Python 2.7.6 (default, Oct 26 2016, 20:30:19)
[GCC 4.8.4] on linux2
Type "help", "copyright", "credits" or "license" for more information.
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
17/04/24 09:49:02 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using
17/04/24 09:49:02 WARN Utils: Your hostname, ubuntu resolves to a loopback address: 127.0.1.1; using 19:
17/04/24 09:49:02 WARN Utils: Set SPARK_LOCAL_IP if you need to bind to another address
17/04/24 09:49:06 WARN ObjectStore: Failed to get database global_temp, returning NoSuchObjectException
Welcome to

  ____      _
 / ___|    / \
| |  | |  / _ \
| |  | | / ___ \
| |  | | \___/ \
| |  | |
| |  | |
|_|  |_|

 version 2.1.0

Using Python version 2.7.6 (default, Oct 26 2016 20:30:19)
SparkSession available as 'spark'.
>>> █
```

```
MyPySpark.py x  pca_example.py x  pySparkDemo.py x  Query.py x  svm_with_sgd_example.py x  PySparkDemo2.py x
if __name__ == ...
1  from __future__ import print_function
2  import os
3  import sys
4
5  # Path for spark source folder
6  os.environ['SPARK_HOME'] = "C:/Users/rezkar/Downloads/spark-2.1.0-bin-hadoop2.7/"
7  os.environ['HADOOP_HOME'] = "C:/Users/rezkar/Downloads/spark-2.1.0-bin-hadoop2.7/"
8
9  # Append pyspark to Python Path
10 sys.path.append("C:/Users/rezkar/Downloads/spark-2.1.0-bin-hadoop2.7/python/")
11 sys.path.append("C:/Users/rezkar/Downloads/spark-2.1.0-bin-hadoop2.7/python/lib/py4j-0.10.4-src.zip")
12
13 try:
14     from pyspark.ml.feature import PCA
15     from pyspark.ml.linalg import Vectors
16     from pyspark.sql import SparkSession
17     print ("Successfully imported Spark Modules")
18
19 except ImportError as e:
20     print ("Can not import Spark Modules", e)
21     sys.exit(1)
22
```

```
Run  pca_example
C:\Python27\python.exe C:/Users/rezkar/PycharmProjects/SparkWithPython/pca_example.py
Successfully imported Spark Modules
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
17/04/24 18:11:15 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using built-in
[Stage 6:>
(0 + 8) / 8]17/04/24 18:11:25 WARN BLAS: Failed to load implementation from: com.github.fommil.netlib.NativeRefBLAS
17/04/24 18:11:25 WARN LAPACK: Failed to load implementation from: com.github.fommil.netlib.NativeSystemLAPACK
17/04/24 18:11:25 WARN LAPACK: Failed to load implementation from: com.github.fommil.netlib.NativeRefLAPACK
-----
```

```
Run pca_example
C:\Python27\python.exe C:/Users/rezkar/PycharmProjects/SparkWithPython/pca_example.py
Successfully imported Spark Modules
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
17/04/26 14:50:59 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform..
(0 + 8) / 8]17/04/26 14:51:10
[Stage 6:>
17/04/26 14:51:10 WARN BLAS: Failed to load implementation from: com.github.fommil.netlib.NativeR
17/04/26 14:51:10 WARN LAPACK: Failed to load implementation from: com.github.fommil.netlib.Nativ
17/04/26 14:51:10 WARN LAPACK: Failed to load implementation from: com.github.fommil.netlib.Nativ
+-----+
|pcaFeatures                                     |
+-----+
|[1.6485728230883807,-4.013282700516296,-5.524543751369388] |
|[-4.645104331781534,-1.1167972663619026,-5.524543751369387]|
|[-6.428880535676489,-5.337951427775355,-5.524543751369389] |
+-----+

SUCCESS: The process with PID 7416 (child process of PID 4100) has been terminated.
SUCCESS: The process with PID 4100 (child process of PID 6488) has been terminated.
SUCCESS: The process with PID 6488 (child process of PID 1200) has been terminated.

Process finished with exit code 0
```



```
+-----+-----+
|label|          features|
+-----+-----+
| 8.0|(17, [0,1,2,3,4,5,...|
|10.0|(17, [0,1,2,3,4,5,...|
| 9.0|(17, [0,1,2,3,4,5,...|
| 8.0|(17, [0,1,2,3,4,5,...|
|10.0|(17, [0,1,2,3,4,5,...|
| 8.0|(17, [0,1,2,3,4,5,...|
| 5.0|(17, [0,1,2,3,4,5,...|
| 6.0|(17, [0,1,2,3,4,5,...|
| 8.0|(17, [0,1,2,3,4,5,...|
| 7.0|(17, [0,1,2,3,4,5,...|
| 6.0|(17, [0,1,2,3,4,5,...|
| 8.0|(17, [0,1,2,3,4,5,...|
| 8.0|(17, [0,1,2,3,4,5,...|
| 8.0|(17, [0,1,2,3,4,5,...|
| 9.0|(17, [0,1,2,3,4,5,...|
| 4.0|(17, [0,1,2,3,4,5,...|
| 7.0|(17, [0,1,2,3,4,5,...|
| 7.0|(17, [0,1,2,3,4,5,...|
| 8.0|(17, [0,1,2,3,4,5,...|
| 8.0|(17, [0,1,2,3,4,5,...|
+-----+-----+
only showing top 20 rows
```

```

root
|-- year: string (nullable = true)
|-- month: string (nullable = true)
|-- day: string (nullable = true)
|-- dep_time: string (nullable = true)
|-- dep_delay: string (nullable = true)
|-- arr_time: string (nullable = true)
|-- arr_delay: string (nullable = true)
|-- carrier: string (nullable = true)
|-- tailnum: string (nullable = true)
|-- flight: string (nullable = true)
|-- origin: string (nullable = true)
|-- dest: string (nullable = true)
|-- air_time: string (nullable = true)
|-- distance: string (nullable = true)
|-- hour: string (nullable = true)
|-- minute: string (nullable = true)

```

year	month	day	dep_time	dep_delay	arr_time	arr_delay	carrier	tailnum	flight	origin	dest	air_time	distance	hour	minute
2013	1	1	517	2	830	11	UA	N14228	1545	EWR	IAH	227	1400	5	17
2013	1	1	533	4	850	20	UA	N24211	1714	LGA	IAH	227	1416	5	33
2013	1	1	542	2	923	33	AA	N619AA	1141	JFK	MIA	160	1089	5	42
2013	1	1	544	-1	1004	-18	B6	N804JB	725	JFK	BQN	183	1576	5	44
2013	1	1	554	-6	812	-25	DL	N668DN	461	LGA	ATL	116	762	5	54
2013	1	1	554	-4	740	12	UA	N39463	1696	EWR	ORD	150	719	5	54
2013	1	1	555	-5	913	19	B6	N516JB	507	EWR	FLL	158	1065	5	55
2013	1	1	557	-3	709	-14	EV	N829AS	5708	LGA	IAD	53	229	5	57
2013	1	1	557	-3	838	-8	B6	N593JB	79	JFK	MCO	140	944	5	57
2013	1	1	558	-2	753	8	AA	N3ALAA	301	LGA	ORD	138	733	5	58
2013	1	1	558	-2	849	-2	B6	N793JB	49	JFK	PBI	149	1028	5	58
2013	1	1	558	-2	853	-3	B6	N657JB	71	JFK	TPA	158	1005	5	58
2013	1	1	558	-2	924	7	UA	N29129	194	JFK	LAX	345	2475	5	58
2013	1	1	558	-2	923	-14	UA	N53441	1124	EWR	SFO	361	2565	5	58
2013	1	1	559	-1	941	31	AA	N3DUAA	707	LGA	DFW	257	1389	5	59
2013	1	1	559	0	702	-4	B6	N708JB	1806	JFK	BOS	44	187	5	59
2013	1	1	559	-1	854	-8	UA	N76515	1187	EWR	LAS	337	2227	5	59
2013	1	1	600	0	851	-7	B6	N595JB	371	LGA	FLL	152	1076	6	0
2013	1	1	600	0	837	12	MQ	N542MQ	4650	LGA	ATL	134	762	6	0
2013	1	1	601	1	844	-6	B6	N644JB	343	EWR	PBI	147	1023	6	1

only showing top 20 rows

```
+-----+-----+-----+-----+-----+
|number|product_name| ransaction_id|   website|price|   date|
+-----+-----+-----+-----+-----+
|   0|   jeans|30160906182001|   ebay.com| 100|12-02-2016|
|   1|  camera|70151231120504| amazon.com| 450|09-08-2017|
|   2|  laptop|90151231120504|   ebay.ie| 1500|07--5-2016|
|   3|   book|80151231120506|  packt.com|   45|03-12-2016|
|   4|  drone| 8876531120508|alibaba.com| 120|01-05-2017|
+-----+-----+-----+-----+-----+
```

```
+-----+-----+-----+-----+-----+
|number|product_name| ransaction_id|   website|price|   date|
+-----+-----+-----+-----+-----+
|   0|   jeans|30160906182001|   ebay.com| 100|12-02-2016|
|   1|  camera|70151231120504| amazon.com| 450|09-08-2017|
|   2|  laptop|90151231120504|   ebay.ie| 1500|07--5-2016|
|   3|   book|80151231120506|  packt.com|   45|03-12-2016|
|   4|  drone| 8876531120508|alibaba.com| 120|01-05-2017|
+-----+-----+-----+-----+-----+
```

```
+-----+-----+
|product_name|price|
+-----+-----+
|   laptop| 1500|
+-----+-----+
```

```
+-----+-----+
|product_name|price|
+-----+-----+
|   camera| 450|
|   book| 45|
|   laptop| 1500|
|   drone| 120|
|   jeans| 100|
+-----+-----+
```

Student	Course	Score
Jason	Math	87
Jason	Science	32
Jason	Geography	126
Jason	History	12
Jason	IT	17
Jason	Statistics	37
John	Math	143
John	Science	54
John	Geography	146
John	History	54
John	IT	26
John	Statistics	171
Geroge	Math	102
Geroge	Science	146
Geroge	Geography	5
Geroge	History	112
Geroge	IT	163
Geroge	Statistics	175
David	Math	27
David	Science	4

only showing top 20 rows

Student	Course	Score	Grade
Jason	Math	87	B
Jason	Science	32	D
Jason	Geography	126	A
Jason	History	12	D
Jason	IT	17	D
Jason	Statistics	37	D
John	Math	143	A
John	Science	54	D
John	Geography	146	A
John	History	54	D
John	IT	26	D
John	Statistics	171	A
Geroge	Math	102	A
Geroge	Science	146	A
Geroge	Geography	5	D
Geroge	History	112	A
Geroge	IT	163	A
Geroge	Statistics	175	A
David	Math	27	D
David	Science	4	D
David	Geography	1	D
David	History	13	D
David	IT	60	C
David	Statistics	19	D

```
+-----+-----+-----+
|Student|Score|Grade|
+-----+-----+-----+
| Jason|  42|   D|
| Jason| 153|   A|
| Jason| 120|   A|
| Jason|  99|   A|
| Jason| 110|   A|
| Jason| 150|   A|
|  John|  21|   D|
|  John|  45|   D|
|  John|   1|   D|
|  John| 138|   A|
|  John| 168|   A|
|  John|  90|   A|
|Geroge|  84|   B|
|Geroge|  84|   B|
|Geroge| 192|   A|
|Geroge| 192|   A|
|Geroge|  10|   D|
|Geroge| 132|   A|
| David|  93|   A|
| David| 127|   A|
+-----+-----+-----+
only showing top 20 rows
```

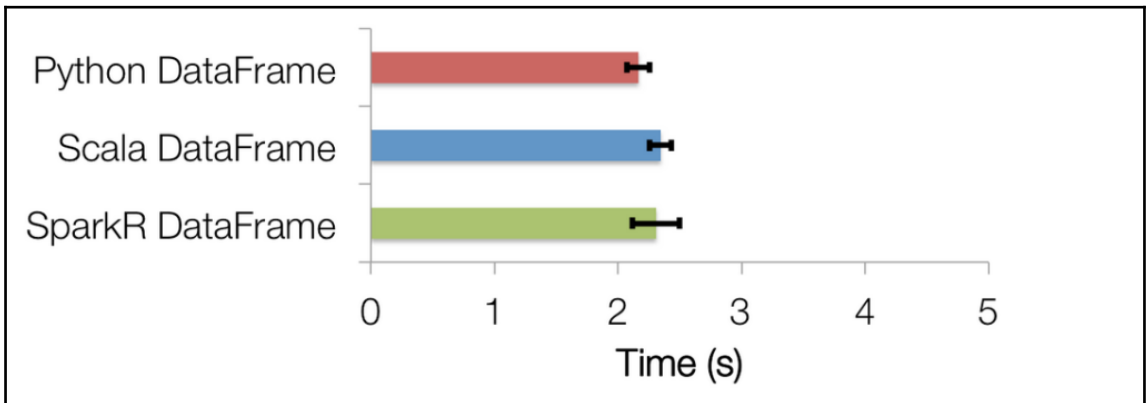
```
+---+---+---+---+---+
|_1|_2|_3|_4|_42|
+---+---+---+---+---+
| 0|tcp|http|SF|normal.|
| 0|tcp|http|SF|normal.|
| 0|tcp|http|SF|normal.|
| 0|tcp|http|SF|normal.|
| 0|tcp|http|SF|normal.|
+---+---+---+---+---+
only showing top 5 rows
```

```
smurf. 2807886
neptune. 1072017
normal. 972781
satan. 15892
ipsweep. 12481
portsweep. 10413
nmap. 2316
back. 2203
warezclient. 1020
teardrop. 979
pod. 264
guess_passwd. 53
buffer_overflow. 30
land. 21
warezmaster. 20
imap. 12
rootkit. 10
loadmodule. 9
ftp_write. 8
multihop. 7
phf. 4
perl. 3
spy. 2
```

```

Final centers: [array([ 4.10612163e+00,  6.36522840e-02,  4.85948958e-02,
-2.21319176e-03, -1.51849176e-02,  1.59666681e-02,
-1.37464150e-02,  4.63552710e-03, -2.80722691e-01,
 1.01178785e-01,  7.90818282e-02,  1.62820689e-01,
 1.08778945e-01,  3.21998554e-01, -8.41384069e-03,
 6.05393588e-02,  0.00000000e+00,  3.30078588e-02,
-2.46237569e-02, -1.14832651e+00, -1.19575475e+00,
-3.71645499e-01, -3.67973482e-01,  8.19357206e-01,
 8.14955084e-01, -3.26320418e-01,  4.33755203e+00,
-1.82859395e-01,  1.79392516e-01, -1.71925941e+00,
-1.75521881e+00,  6.82285609e+00,  2.23215018e-01,
-1.16133090e-01, -3.68177485e-01, -3.66477378e-01,
 8.07658804e-01,  8.18438116e-01]), array([-6.69802290e-02, -1.36283222e-03, -1.65369293e-03,
-2.21319176e-03, -1.51849176e-02, -1.64391576e-03,
-2.65266109e-02, -4.38631465e-03, -4.09296131e-01,
-2.00370428e-03, -8.21527723e-03, -4.60861589e-03,
-3.04988915e-03, -9.62851412e-03, -8.41384069e-03,
-2.85810713e-02,  0.00000000e+00, -5.21653093e-04,
-2.88684412e-02,  6.87674624e-01,  7.54010775e-01,
-4.65800760e-01, -4.65512939e-01, -2.48364764e-01,
-2.48177638e-01,  5.39551929e-01, -2.55781037e-01,
-2.01125081e-01,  3.42806366e-01,  6.19909484e-01,
 5.98368428e-01, -2.82739959e-01,  8.20664819e-01,
-1.56479158e-01, -4.66075407e-01, -4.65194517e-01,
-2.50690649e-01, -2.49676723e-01]), array([-6.69767578e-02, -1.86749297e-03, -1.65012194e-03,
-2.21319176e-03, -1.51849176e-02, -1.64391576e-03,
-2.64973873e-02, -4.38631465e-03, -4.09177709e-01,
-1.99486560e-03, -8.21527723e-03, -4.60861589e-03,

```





Property Name	Property group	spark-submit equivalent
spark.master	Application Properties	--master
spark.yarn.keytab	Application Properties	--keytab
spark.yarn.principal	Application Properties	--principal
spark.driver.memory	Application Properties	--driver-memory
spark.driver.extraClassPath	Runtime Environment	--driver-class-path
spark.driver.extraJavaOptions	Runtime Environment	--driver-java-options
spark.driver.extraLibraryPath	Runtime Environment	--driver-library-path

	year	month	day	dep_time	dep_delay	arr_time	arr_delay	carrier	tailnum	flight	origin	dest	air_time	distance	hour	minute
1	2013	1	1	517	2	830	11	UA	N14228	1545	EWB	IAH	227	1400	5	17
2	2013	1	1	533	4	850	20	UA	N24211	1714	LGA	IAH	227	1416	5	33
3	2013	1	1	542	2	923	33	AA	N619AA	1141	JFK	MIA	160	1089	5	42
4	2013	1	1	544	-1	1004	-18	B6	N804JB	725	JFK	BQN	183	1576	5	44
5	2013	1	1	554	-6	812	-25	DL	N668DN	461	LGA	ATL	116	762	5	54
6	2013	1	1	554	-4	740	12	UA	N39463	1696	EWB	ORD	150	719	5	54
7	2013	1	1	555	-5	913	19	B6	N516JB	507	EWB	FLL	158	1065	5	55
8	2013	1	1	557	-3	709	-14	EV	N829AS	5708	LGA	IAD	53	229	5	57
9	2013	1	1	557	-3	838	-8	B6	N593JB	79	JFK	MCO	140	944	5	57
10	2013	1	1	558	-2	753	8	AA	N3ALAA	301	LGA	ORD	138	733	5	58
11	2013	1	1	558	-2	849	-2	B6	N793JB	49	JFK	PBI	149	1028	5	58
12	2013	1	1	558	-2	853	-3	B6	N657JB	71	JFK	TPA	158	1005	5	58
13	2013	1	1	558	-2	924	7	UA	N29129	194	JFK	LAX	345	2475	5	58
14	2013	1	1	558	-2	923	-14	UA	N53441	1124	EWB	SFO	361	2565	5	58
15	2013	1	1	559	-1	941	31	AA	N3DUAA	707	LGA	DFW	257	1389	5	59
16	2013	1	1	559	0	702	-4	B6	N708JB	1806	JFK	BOS	44	187	5	59
17	2013	1	1	559	-1	854	-8	UA	N76515	1187	EWB	LAS	337	2227	5	59
18	2013	1	1	600	0	851	-7	B6	N595JB	371	LGA	FLL	152	1076	6	0
19	2013	1	1	600	0	837	12	MQ	N542MQ	4650	LGA	ATL	134	762	6	0
20	2013	1	1	601	1	844	-6	B6	N644JB	343	EWB	PBI	147	1023	6	1

```

root
|-- year: integer (nullable = true)
|-- month: integer (nullable = true)
|-- day: integer (nullable = true)
|-- dep_time: string (nullable = true)
|-- dep_delay: string (nullable = true)
|-- arr_time: string (nullable = true)
|-- arr_delay: string (nullable = true)
|-- carrier: string (nullable = true)
|-- tailnum: string (nullable = true)
|-- flight: integer (nullable = true)
|-- origin: string (nullable = true)
|-- dest: string (nullable = true)
|-- air_time: string (nullable = true)
|-- distance: integer (nullable = true)
|-- hour: string (nullable = true)
|-- minute: string (nullable = true)

```

year	month	day	dep_time	dep_delay	arr_time	arr_delay	carrier	tailnum	flight	origin	dest	air_time	distance	hour	minute
2013	1	1	517	2	830	11	UA	N14228	1545	EWB	IAH	227	1400	5	17
2013	1	1	533	4	850	20	UA	N24211	1714	LGA	IAH	227	1416	5	33
2013	1	1	542	2	923	33	AA	N619AA	1141	JFK	MIA	160	1089	5	42
2013	1	1	544	-1	1004	-18	B6	N804JB	725	JFK	BQN	183	1576	5	44
2013	1	1	554	-6	812	-25	DL	N668DN	461	LGA	ATL	116	762	5	54
2013	1	1	554	-4	740	12	UA	N39463	1696	EWB	ORD	150	719	5	54
2013	1	1	555	-5	913	19	B6	N516JB	507	EWB	FLL	158	1065	5	55
2013	1	1	557	-3	709	-14	EV	N829AS	5708	LGA	IAD	53	229	5	57
2013	1	1	557	-3	838	-8	B6	N593JB	79	JFK	MCO	140	944	5	57
2013	1	1	558	-2	753	8	AA	N3ALAA	301	LGA	ORD	138	733	5	58

only showing top 10 rows

```

root
|-- year: integer (nullable = true)
|-- month: integer (nullable = true)
|-- day: integer (nullable = true)
|-- dep_time: string (nullable = true)
|-- dep_delay: string (nullable = true)
|-- arr_time: string (nullable = true)
|-- arr_delay: string (nullable = true)
|-- carrier: string (nullable = true)
|-- tailnum: string (nullable = true)
|-- flight: integer (nullable = true)
|-- origin: string (nullable = true)
|-- dest: string (nullable = true)
|-- air_time: string (nullable = true)
|-- distance: integer (nullable = true)
|-- hour: string (nullable = true)
|-- minute: string (nullable = true)

```

year	month	day	dep_time	dep_delay	arr_time	arr_delay	carrier	tailnum	flight	origin	dest	air_time	distance	hour	minute
2013	1	1	517	2	830	11	UA	N14228	1545	EWB	IAH	227	1400	5	17
2013	1	1	533	4	850	20	UA	N24211	1714	LGA	IAH	227	1416	5	33
2013	1	1	542	2	923	33	AA	N619AA	1141	JFK	MIA	160	1089	5	42
2013	1	1	544	-1	1004	-18	B6	N804JB	725	JFK	BQN	183	1576	5	44
2013	1	1	554	-6	812	-25	DL	N668DN	461	LGA	ATL	116	762	5	54
2013	1	1	554	-4	740	12	UA	N39463	1696	EWB	ORD	150	719	5	54
2013	1	1	555	-5	913	19	B6	N516JB	507	EWB	FLL	158	1065	5	55
2013	1	1	557	-3	709	-14	EV	N829AS	5708	LGA	IAD	53	229	5	57
2013	1	1	557	-3	838	-8	B6	N593JB	79	JFK	MCO	140	944	5	57
2013	1	1	558	-2	753	8	AA	N3ALAA	301	LGA	ORD	138	733	5	58

only showing top 10 rows

year	month	day	dep_time	dep_delay	arr_time	arr_delay	carrier	tailnum	flight	origin	dest	air_time	distance	hour	minute
2013	1	1	542	2	923	33	AA	N619AA	1141	JFK	MIA	160	1089	5	42
2013	1	1	606	-4	858	-12	AA	N633AA	1895	EWB	MIA	152	1085	6	6
2013	1	1	607	0	858	-17	UA	N53442	1077	EWB	MIA	157	1085	6	7
2013	1	1	623	13	920	5	AA	N3EMAA	1837	LGA	MIA	153	1096	6	23
2013	1	1	655	-5	1002	-18	DL	N997DL	2003	LGA	MIA	161	1096	6	55
2013	1	1	659	-1	1008	-7	AA	N3EKAA	2279	LGA	MIA	159	1096	6	59
2013	1	1	753	-2	1056	-14	AA	N3HMAA	2267	LGA	MIA	157	1096	7	53
2013	1	1	759	-1	1057	-30	DL	N955DL	1843	JFK	MIA	158	1089	7	59
2013	1	1	826	71	1136	51	AA	N3GVAA	443	JFK	MIA	160	1089	8	26
2013	1	1	856	-4	1222	-10	DL	N970DL	2143	LGA	MIA	158	1096	8	56

only showing top 10 rows

flight	dep_delay	origin	dest
1545	2	EWR	IAH
1714	4	LGA	IAH
496	-4	LGA	IAH
473	-4	LGA	IAH
1479	0	EWR	IAH
1220	0	EWR	IAH
1004	2	LGA	IAH
455	-1	EWR	IAH
1086	134	LGA	IAH
1461	5	EWR	IAH

only showing top 10 rows

day	avg(dep_delay)	avg(arr_delay)
1 31	9.506521	3.359225
2 28	15.743213	8.183567
3 26	9.748002	3.656098
4 27	12.083969	3.331213
5 12	15.177765	11.138973
6 22	18.712073	17.404916

	dest	NUM_FLIGHTS	AVG_DELAY	MAX_DELAY	MIN_DELAY
1	PSE	365	7.871508	NA	-1
2	MSY	3799	6.490175	NA	-1
3	BUR	371	8.175676	NA	-1
4	SNA	825	-7.868227	NA	-1
5	GRR	765	18.189560	NA	-1
6	GSO	1606	14.112601	NA	-1

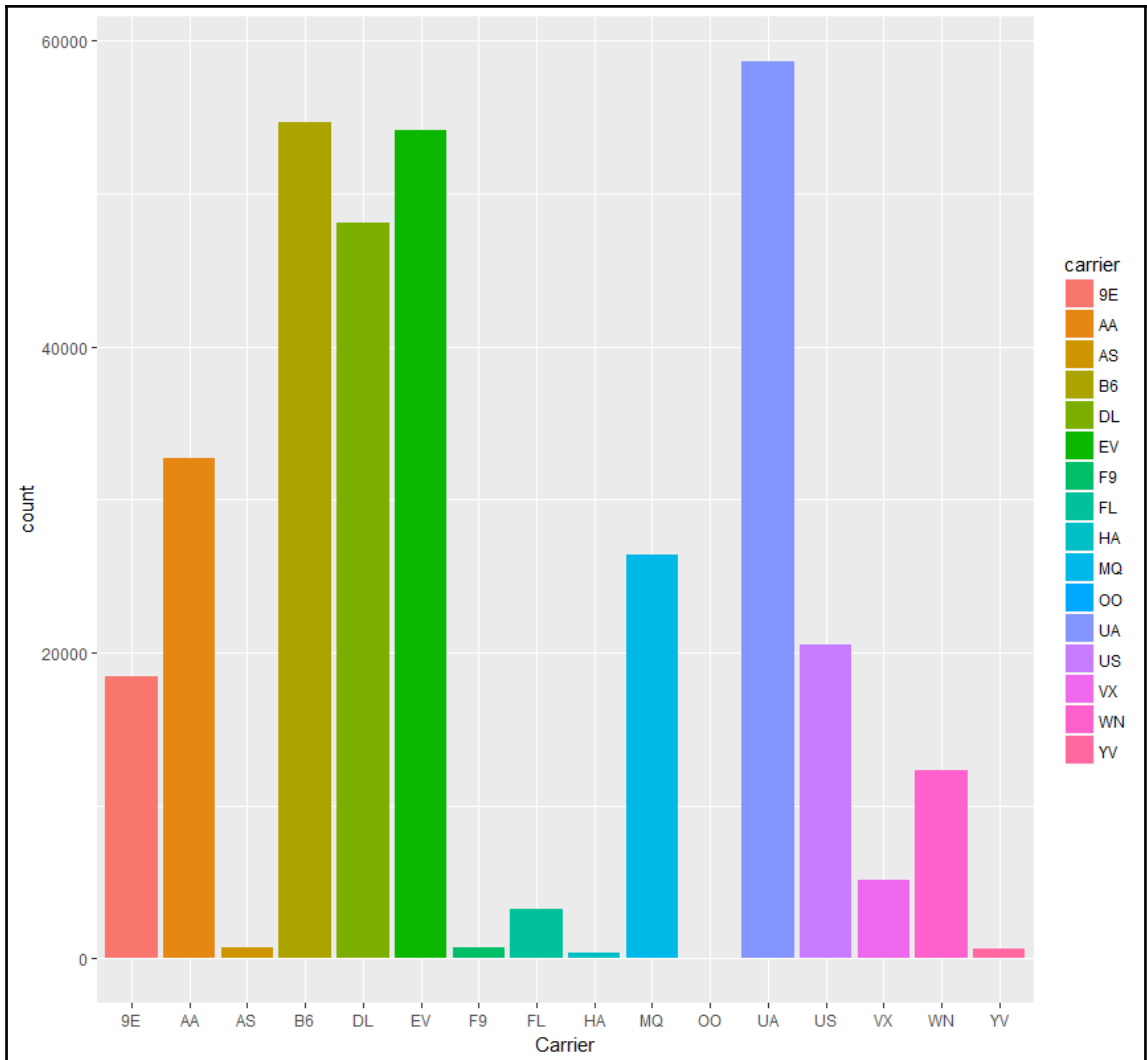
dest	origin	carrier
IAH	EWR	UA
IAH	LGA	UA
MIA	JFK	AA
BQN	JFK	B6
ATL	LGA	DL
ORD	EWR	UA
FLL	EWR	B6
IAD	LGA	EV
MCO	JFK	B6
ORD	LGA	AA

only showing top 10 rows

dest	origin	arr_delay
CLT	LGA	137
BWI	JFK	851
BOS	EWR	123
IAH	LGA	145
RIC	EWR	127
MCO	EWR	125
MCI	EWR	136
IAD	JFK	123
DAY	EWR	123
BNA	LGA	138

only showing top 10 rows

origin	dest	arr_delay
JFK	IAH	783
LGA	IAH	435
LGA	IAH	390
EWB	IAH	374
EWB	IAH	373
LGA	IAH	370
LGA	IAH	363
EWB	IAH	338
LGA	IAH	324
LGA	IAH	321
LGA	IAH	312
LGA	IAH	309
EWB	IAH	302
LGA	IAH	301
EWB	IAH	297
LGA	IAH	294
EWB	IAH	292
EWB	IAH	288
EWB	IAH	283
LGA	IAH	278



carrier	cnt
UA	58665
B6	54635
EV	54173
DL	48110
AA	32729
MQ	26397
US	20536
9E	18460
WN	12275
VX	5162
FL	3260
AS	714
F9	685
YV	601
HA	342
OO	32