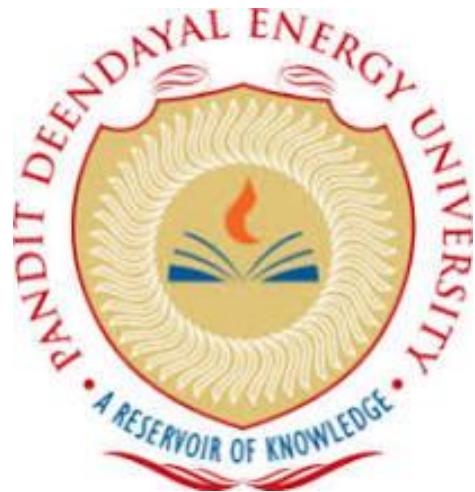


Big Data Analytics Lab (23CP309P)

Lab Report Submitted to

Pandit Deendayal Energy University, Gandhinagar

for



Bachelor of Technology
in
Computer Science & Engineering Department

Submitted by

Dhvani PATEL 21BCP116
Semester: VI
Division: 2 , Group: G4

Course Faculty

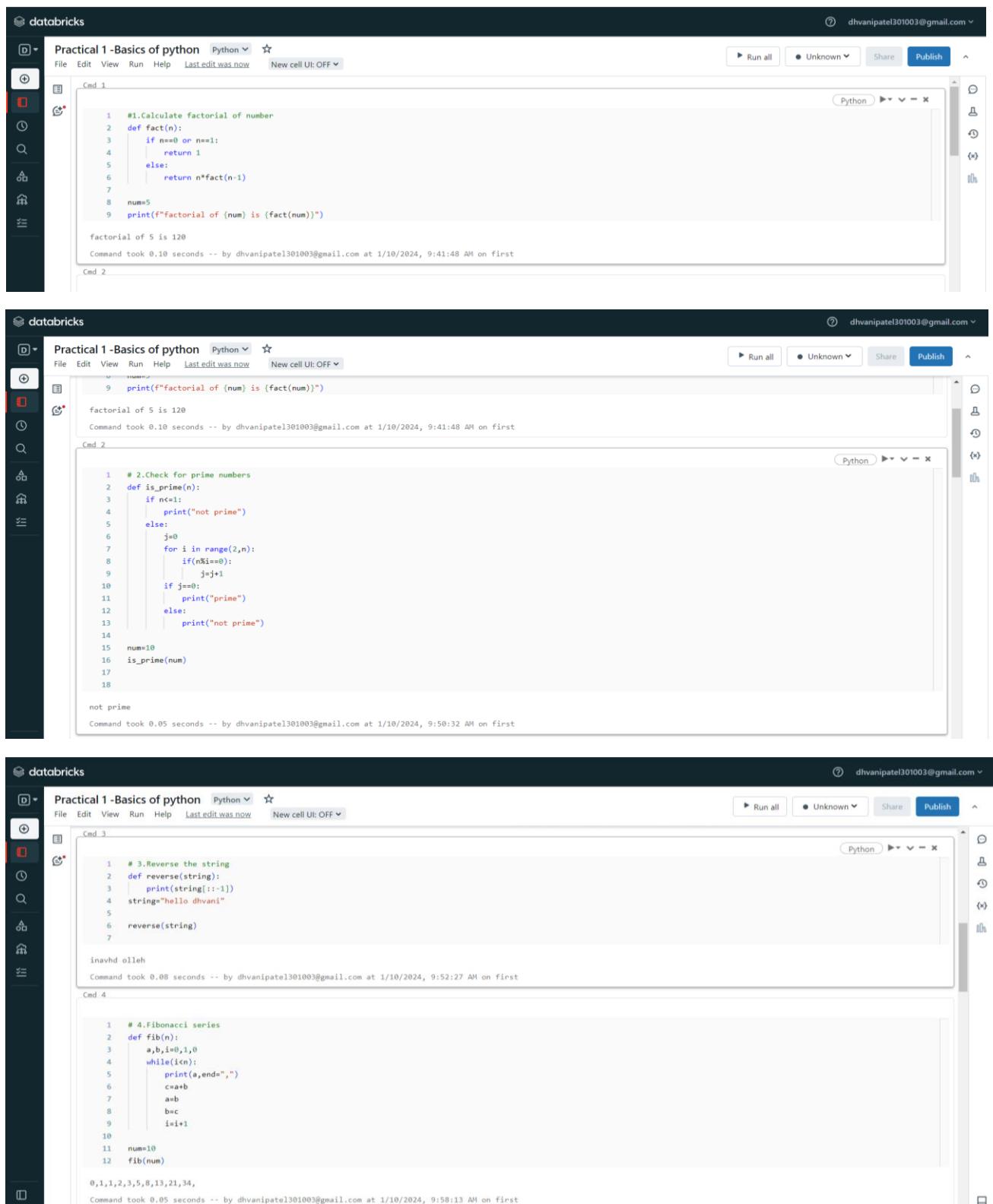
SOHAM VYAS

PANDIT DEENDAYAL ENERGY UNIVERSITY, GANDHINAGAR - 382007, India
May-2024

List of Practicals

Exp. No.	Experiment Title	Date	Signature
1	Basics of Python	10/01/24	
2	Basics of Scala	17/01/24	
3	Transformation Function	31/01/24	
4	Spark	14/02/24	
5	Spark SQL	21/02/24	
6	Data Processing using Spark	13/03/24	
7	Hadoop, Pig, Kafka, Spark and Hive installation	27/03/24	
8	Implementation of graph data structure using networkx	03/04/24	
9	Linear Regression	10/04/24	

Name: Dhvani Patel
Roll no: 21BCP116
Big Data Analytics Lab
Practical 1 : Basics of Python



The image shows three screenshots of a Databricks notebook titled "Practical 1 -Basics of python" in Python. The notebook interface includes a sidebar with icons for file operations, a command bar at the top, and a main workspace for code execution.

Screenshot 1: The code in Cmd 1 calculates the factorial of a number. The output shows the result for n=5.

```
1 #1.Calculate factorial of number
2 def fact(n):
3     if n==0 or n==1:
4         return 1
5     else:
6         return n*fact(n-1)
7
8 num=5
9 print(f"factorial of {num} is {(fact(num))}")

factorial of 5 is 120
Command took 0.10 seconds -- by dhvanipatel301003@gmail.com at 1/10/2024, 9:41:48 AM on first
```

Screenshot 2: The code in Cmd 2 checks if a number is prime. It defines a function `is_prime` and uses it to check if 10 is prime.

```
1 # 2.Check for prime numbers
2 def is_prime(n):
3     if n<=1:
4         print("not prime")
5     else:
6         j=0
7         for i in range(2,n):
8             if(n%i==0):
9                 j=j+1
10            if j==0:
11                print("prime")
12            else:
13                print("not prime")
14
15 num=10
16 is_prime(num)
17
18

not prime
Command took 0.05 seconds -- by dhvanipatel301003@gmail.com at 1/10/2024, 9:50:32 AM on first
```

Screenshot 3: The code in Cmd 3 reverses a string. It defines a function `reverse` and prints the reverse of "hello dhvani".

```
1 # 3.Reverse the string
2 def reverse(string):
3     print(string[::-1])
4 string="hello dhvani"
5
6 reverse(string)
7

olleh olleh
Command took 0.08 seconds -- by dhvanipatel301003@gmail.com at 1/10/2024, 9:52:27 AM on first
```

Screenshot 4: The code in Cmd 4 generates a Fibonacci series. It defines a function `fib` and prints the first 10 numbers of the series.

```
1 # 4.Fibonacci series
2 def fib(n):
3     a,b,i=0,1,0
4     while(i<n):
5         print(a,end=",")
6         c=a+b
7         a=b
8         b=c
9         i=i+1
10
11 num=10
12 fib(num)

0,1,1,2,3,5,8,13,21,34,
Command took 0.05 seconds -- by dhvanipatel301003@gmail.com at 1/10/2024, 9:58:13 AM on first
```

databricks

Practical 1 -Basics of python Python

File Edit View Run Help Last edit was 1 minute ago New cell UI: OFF

Run all Unknown Share Publish

Command took 0.05 seconds -- by dhvanipatel301003@gmail.com at 1/10/2024, 9:58:13 AM on first

Cmd 5

```
1 # 5.Word frequency in string
2 def word_freq(string):
3     newstring.split(" ")
4     lst=list()
5     j=0
6     for i in range(0,len(new)):
7         count=0
8         if new[i] not in lst:
9             lst.append(new[i])
10            j=j+1
11            for k in range(0,len(new)):
12                if new[i]==new[k]:
13                    count=count+1
14            print(f"(new[i]):{count}")
15
16 string="hello world hello"
17 word_freq(string)
18
```

hello:2
world:1

Command took 0.07 seconds -- by dhvanipatel301003@gmail.com at 1/10/2024, 10:13:47 AM on first

databricks

Practical 1 -Basics of python Python

File Edit View Run Help Last edit was 1 minute ago New cell UI: OFF

Run all Unknown Share Publish

hello:2
world:1

Command took 0.07 seconds -- by dhvanipatel301003@gmail.com at 1/10/2024, 10:13:47 AM on first

Cmd 6

```
1 %fs
2 ls /FileStore/tables
```

Table +

path	name	size	modificationTime
1 dbfs/FileStore/tables/iris.csv	iris.csv	3858	1704862675000

1 row | 17.25 seconds runtime

Refreshed 41 days ago

Command took 17.25 seconds -- by dhvanipatel301003@gmail.com at 1/10/2024, 10:29:55 AM on first

Cmd 7

```
1 my_df=spark.read.format("csv").option("inferSchema","true").option("header","true").load("/FileStore/tables/iris.csv")
```

(2) Spark Jobs

my_df: pyspark.sql.dataframe.DataFrame = [sepal_length: double, sepal_width: double ... 3 more fields]

Command took 1.43 seconds -- by dhvanipatel301003@gmail.com at 1/10/2024, 10:42:19 AM on first

databricks

Practical 1 -Basics of python Python

File Edit View Run Help Last edit was 1 minute ago New cell UI: OFF

Run all Unknown Share Publish

▶ (2) Spark Jobs

▶ my_df: pyspark.sql.dataframe.DataFrame = [sepal_length: double, sepal_width: double ... 3 more fields]

Command took 1.43 seconds -- by dhvanipatel301003@gmail.com at 1/10/2024, 10:42:19 AM on first

Cmd 8

```
1 display(my_df)
```

▶ (1) Spark Jobs

Table +

sepal_length	sepal_width	petal_length	petal_width	species
1 5.1	3.5	1.4	0.2	setosa
2 4.9	3	1.4	0.2	setosa
3 4.7	3.2	1.3	0.2	setosa
4 4.6	3.1	1.5	0.2	setosa
5 5	3.6	1.4	0.2	setosa
6 5.4	3.9	1.7	0.4	setosa
7 4.6	3.4	1.4	0.3	setosa

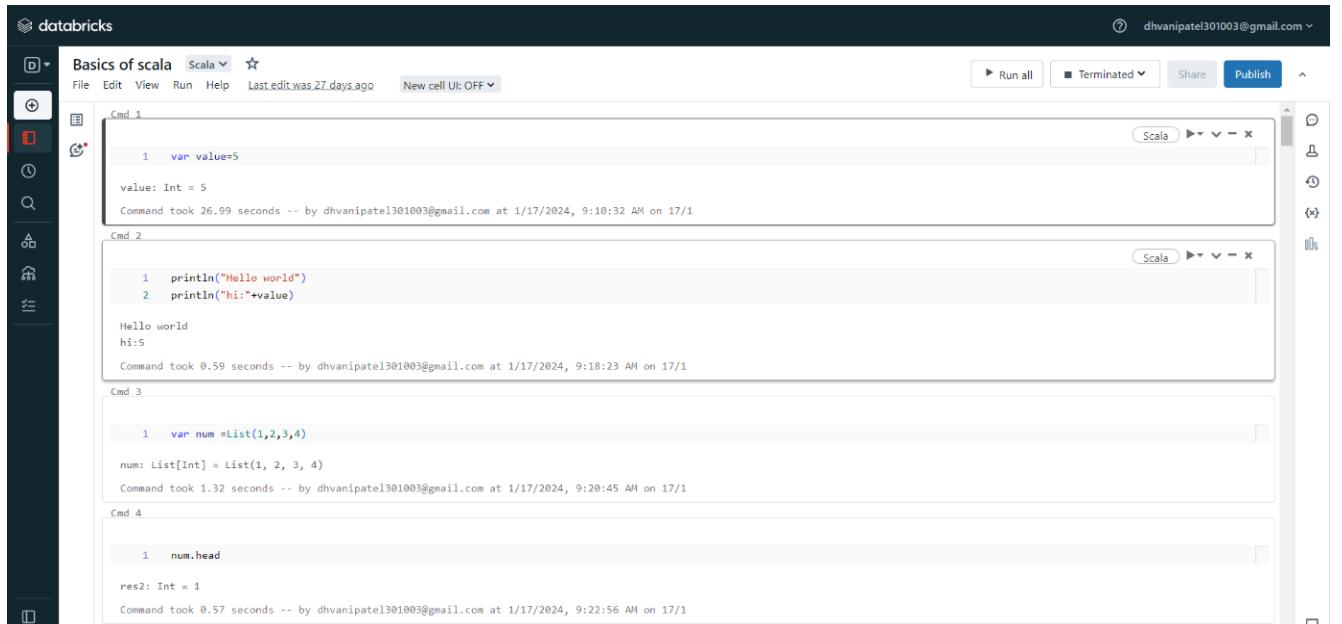
150 rows | 0.55 seconds runtime

Refreshed 41 days ago

Command took 0.55 seconds -- by dhvanipatel301003@gmail.com at 1/10/2024, 10:42:22 AM on first

[Shift+Enter] to run and move to next cell
[Esc H] to see all keyboard shortcuts

Name: Dhvani Patel
Roll no: 21BCP116
Big Data Analytics Lab
Practical 2 : Basics of Scala



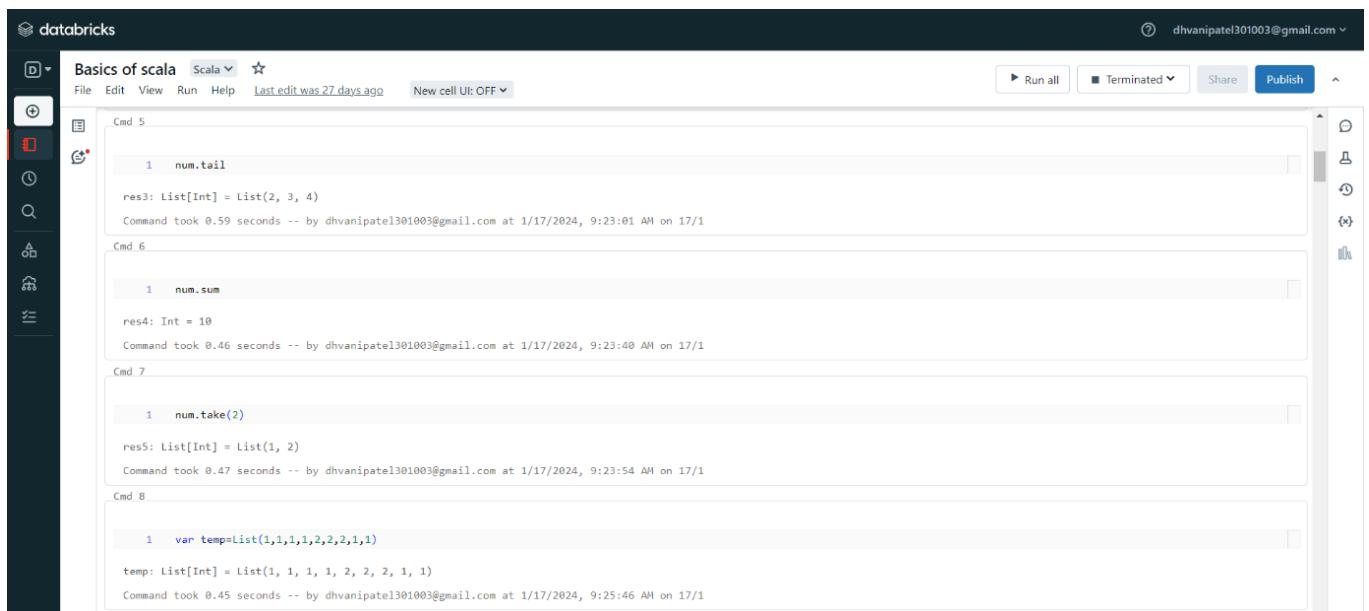
Databricks Notebook titled "Basics of scala" in Scala language. The notebook contains the following code snippets:

```
Cmd 1
1 var value=5
value: Int = 5
Command took 26.99 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:10:32 AM on 17/1

Cmd 2
1 println("Hello world")
2 println("hi:"+value)
Hello world
hi:5
Command took 0.59 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:18:23 AM on 17/1

Cmd 3
1 var num =List(1,2,3,4)
num: List[Int] = List(1, 2, 3, 4)
Command took 1.32 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:20:45 AM on 17/1

Cmd 4
1 num.head
res2: Int = 1
Command took 0.57 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:22:56 AM on 17/1
```



Databricks Notebook titled "Basics of scala" in Scala language. The notebook contains the following code snippets:

```
Cmd 5
1 num.tail
res3: List[Int] = List(2, 3, 4)
Command took 0.59 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:23:01 AM on 17/1

Cmd 6
1 num.sum
res4: Int = 10
Command took 0.46 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:23:40 AM on 17/1

Cmd 7
1 num.take(2)
res5: List[Int] = List(1, 2)
Command took 0.47 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:23:54 AM on 17/1

Cmd 8
1 var temp=list(1,1,1,1,2,2,2,1,1)
temp: List[Int] = List(1, 1, 1, 1, 2, 2, 2, 1, 1)
Command took 0.45 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:25:46 AM on 17/1
```

dhvanipatel301003@gmail.com

Basics of scala Scala ☆

File Edit View Run Help Last edit was 27 days ago New cell UI: OFF

Run all Terminated Share Publish

Cmd 9

```
1 temp.distinct
```

res6: List[Int] = List(1, 2)

Command took 0.37 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:26:07 AM on 17/1

Cmd 10

```
1 temp(0)
```

res7: Int = 1

Command took 0.52 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:26:44 AM on 17/1

Cmd 11

```
1 num(10)
```

IndexOutOfBoundsException: 10

Command took 0.96 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:27:31 AM on 17/1

Cmd 12

```
1 num(1)=10
```

command-44701934191313971: error: value update is not a member of List[Int]

num(1)=10

Command took 0.18 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:30:21 AM on 17/1

dhvanipatel301003@gmail.com

Basics of scala Scala ☆

File Edit View Run Help Last edit was 27 days ago New cell UI: OFF

Run all Terminated Share Publish

Cmd 13

```
1 num.size
```

res10: Int = 4

Command took 0.33 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:30:34 AM on 17/1

Cmd 14

```
1 num.length
```

res11: Int = 4

Command took 0.30 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:30:45 AM on 17/1

Cmd 15

```
1 num.reverse
```

res13: List[Int] = List(4, 3, 2, 1)

Command took 0.29 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:30:59 AM on 17/1

Cmd 16

```
1 num.min
```

res14: Int = 1

Command took 0.33 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:31:08 AM on 17/1

dhvanipatel301003@gmail.com

Basics of scala Scala ☆

File Edit View Run Help Last edit was 27 days ago New cell UI: OFF

Run all Terminated Share Publish

Cmd 17

```
1 num.max
```

res15: Int = 4

Command took 0.37 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:34:13 AM on 17/1

Cmd 18

```
1 num.isEmpty
```

res16: Boolean = false

Command took 0.40 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:34:21 AM on 17/1

Cmd 19

```
1 var empty_list>List()
```

empty_list: List[Nothing] = List()

Command took 0.91 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:36:02 AM on 17/1

Cmd 20

```
1 empty_list.isEmpty
```

res17: Boolean = true

Command took 0.36 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:36:11 AM on 17/1

databricks dhvanipatel301003@gmail.com

Basics of scala Scala File Edit View Run Help Last edit was 27 days ago New cell UI: OFF Share Publish

Cmd 21
1 var number=Array(1,2,3,4,5,6,7,8,9)
number: Array[Int] = Array(1, 2, 3, 4, 5, 6, 7, 8, 9)
Command took 0.41 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:36:53 AM on 17/1

Cmd 22
1 var lang=Array("scala","python")
lang: Array[String] = Array(scala, python)
Command took 0.38 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:37:33 AM on 17/1

Cmd 23
1 lang.tail
res18: Array[String] = Array(python)
Command took 0.47 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:37:43 AM on 17/1

Cmd 24
1 lang.head
res19: String = scala
Command took 0.42 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:37:54 AM on 17/1

databricks dhvanipatel301003@gmail.com

Basics of scala Scala File Edit View Run Help Last edit was 27 days ago New cell UI: OFF Share Publish

Cmd 25
1 number(1)=10
Command took 0.31 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:38:04 AM on 17/1

Cmd 26
1 number
res21: Array[Int] = Array(1, 10, 3, 4, 5, 6, 7, 8, 9)
Command took 0.36 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:38:11 AM on 17/1

Cmd 27
1 number(2)="hi"
command-4470193419131412:1: error: type mismatch;
found : String("hi")
required: Int
number(2)="hi"
^
Command took 0.11 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:39:05 AM on 17/1

databricks dhvanipatel301003@gmail.com

Basics of scala Scala File Edit View Run Help Last edit was 27 days ago New cell UI: OFF Share Publish

Cmd 28
1 import scala.collection.mutable.ArrayBuffercars
import scala.collection.mutable.ArrayBuffer
Command took 0.25 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:40:51 AM on 17/1

Cmd 29
1 var cars=new ArrayBuffer[String]()
cars: scala.collection.mutable.ArrayBuffer[String] = ArrayBuffer()
Command took 1.09 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:41:30 AM on 17/1

Cmd 30
1 cars += "BMW"
res23: scala.collection.mutable.ArrayBuffer[String] = ArrayBuffer(BMW)
Command took 0.36 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:41:56 AM on 17/1

Cmd 31
1 cars.append("Jaguar")
Command took 0.30 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:42:29 AM on 17/1

databricks

Basics of scala Scala

File Edit View Run Help Last edit was 27 days ago New cell UI: OFF

Run all Terminated Share Publish

Cmd. 32

```
1 cars
```

res25: scala.collection.mutable.ArrayBuffer[String] = ArrayBuffer(BMW, Jaguar)

Command took 0.40 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:42:38 AM on 17/1

Cmd. 33

```
1 cars.length
```

res26: Int = 2

Command took 0.32 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:42:50 AM on 17/1

Cmd. 34

```
1 cars += "Jaguar"
```

res27: scala.collection.mutable.ArrayBuffer[String] = ArrayBuffer(BMW, Jaguar, Jaguar)

Command took 0.43 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:43:03 AM on 17/1

Cmd. 35

```
1 cars.trimEnd(1)
```

Command took 0.24 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:43:18 AM on 17/1

databricks

Basics of scala Scala

File Edit View Run Help Last edit was 27 days ago New cell UI: OFF

Run all Terminated Share Publish

Cmd. 36

```
1 cars
```

res29: scala.collection.mutable.ArrayBuffer[String] = ArrayBuffer(BMW, Jaguar)

Command took 0.96 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:43:24 AM on 17/1

Cmd. 37

```
1 cars.trimStart(1)
```

Scala ▶ ▷ ▷ ▷ - x

Command took 0.28 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:45:56 AM on 17/1

Cmd. 38

```
1 cars
```

res31: scala.collection.mutable.ArrayBuffer[String] = ArrayBuffer(Jaguar)

Command took 0.37 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:46:00 AM on 17/1

Cmd. 39

```
1 cars.insert(1,"Bentley")
```

Command took 0.27 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:49:49 AM on 17/1

databricks

Basics of scala Scala

File Edit View Run Help Last edit was 27 days ago New cell UI: OFF

Run all Terminated Share Publish

Cmd. 40

```
1 cars
```

res34: scala.collection.mutable.ArrayBuffer[String] = ArrayBuffer(Jaguar, Bentley)

Command took 0.44 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:49:54 AM on 17/1

Cmd. 41

```
1 //Transform and map
```

2 num

res35: List[Int] = List(1, 2, 3, 4)

Command took 0.33 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:50:29 AM on 17/1

Cmd. 42

```
1 val a = num.map(x => x+1)
```

a: List[Int] = List(2, 3, 4, 5)

Command took 0.61 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:51:04 AM on 17/1

Cmd. 43

```
1 val b = num.map(y => y*y)
```

b: List[Int] = List(1, 4, 9, 16)

Command took 0.76 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:51:36 AM on 17/1

databricks

Basics of scala Scala ⚡

File Edit View Run Help Last edit was 27 days ago New cell UI: OFF

Run all Terminated Share Publish

Cmd 44

```
1 val c = b.map(y => y-1)
```

c: List[Int] = List(0, 3, 8, 15)

Command took 0.33 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:52:10 AM on 17/1

Cmd 45

```
1 var d = c.map(x => -x)
```

d: List[Int] = List(0, -3, -8, -15)

Command took 0.38 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:53:28 AM on 17/1

Cmd 46

```
1 d.map(x => -x).map(x => x+1)
```

res36: List[Int] = List(1, 4, 9, 16)

Command took 0.27 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:54:13 AM on 17/1

Cmd 47

```
1 var fruits = List("Orange", "Banana", "Apple")
```

fruits: List[String] = List(Orange, Banana, Apple)

Command took 0.29 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:54:51 AM on 17/1

databricks

Basics of scala Scala ⚡

File Edit View Run Help Last edit was 27 days ago New cell UI: OFF

Run all Terminated Share Publish

Cmd 48

```
1 fruits.map(x => (x,x.length))
```

res37: List[(String, Int)] = List((Orange,6), (Banana,6), (Apple,5))

Command took 0.74 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:55:16 AM on 17/1

Cmd 49

```
1 fruits.filter(x => (x.length > 5))
```

res38: List[String] = List(Orange, Banana)

Command took 0.45 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 9:59:34 AM on 17/1

Cmd 50

```
1 val ratings=list(2.4,5.6,8.9,7.3)
2 //Multiply by 10.filter more than 75
3 //filter and save between 60 to 75,convert back to between 1-10
```

ratings: List[Double] = List(2.4, 5.6, 8.9, 7.3)

Command took 0.38 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 10:08:05 AM on 17/1

databricks

Basics of scala Scala ⚡

File Edit View Run Help Last edit was 27 days ago New cell UI: OFF

Run all Terminated Share Publish

Cmd 51

```
1 val rate= ratings.map(x => x*10)
```

rate: List[Double] = List(24.0, 56.0, 89.0, 73.0)

Command took 0.84 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 10:15:53 AM on 17/1

Cmd 52

```
1 val filter1=rate.filter(x => (x > 75))
```

filter1: List[Double] = List(89.0)

Command took 0.24 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 10:16:29 AM on 17/1

Cmd 53

```
1 var filter2 =rate.filter(x => (x <75 && x>60))
```

filter2: List[Double] = List(73.0)

Command took 0.25 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 10:16:34 AM on 17/1

Cmd 54

```
1 var rate2=filter2.map(x => x/10)
```

rate2: List[Double] = List(7.3)

Command took 0.26 seconds -- by dhvanipatel301003@gmail.com at 1/17/2024, 10:16:37 AM on 17/1

dhvanipatel301003@gmail.com

Basics of scala Scala

```
1 def add(a:Double=100,b:Double=200):Double=
2 {
3     var sum:Double = 0
4     sum = a+b
5     return sum
6 }
```

add: (a: Double, b: Double)Double

```
1 println("Sum:"+add(3,6))
```

Sum:9.0

dhvanipatel301003@gmail.com

Basics of scala Scala

```
1 var x=1
2 if(x<3)
3 {
4     println("Less than 3")
5 }
6 else
7 {
8     println("Greater than 3")
9 }
10
```

Less than 3

x: Int = 1

dhvanipatel301003@gmail.com

Basics of scala Scala

```
1 var i=10
2 while(i>0)
3 {
4     println("value:"+i)
5     i=i-1
6 }
```

value:10
value:9
value:8
value:7
value:6
value:5
value:4
value:3
value:2
value:1
i: Int = 0

Basics of scala Scala

Last edit was 27 days ago New cell UI: OFF

Run all | Terminated | Share | Publish

```

Cmd 60
1 //matrix multiplication
2 import scala.Array
3
4 // Function to perform matrix multiplication
5 def multiplyMatrix(matrix1: Array[Array[Int]], matrix2: Array[Array[Int]]): Array[Array[Int]] = {
6   val rows1 = matrix1.length
7   val cols1 = matrix1(0).length
8   val cols2 = matrix2(0).length
9
10  val result = Array.ofDim[Int](rows1, cols2)
11
12  for (i <- 0 until rows1) {
13    for (j <- 0 until cols2) {
14      for (k <- 0 until cols1) {
15        result(i)(j) += matrix1(i)(k) * matrix2(k)(j)
16      }
17    }
18  }
19
20  result
21 }
22
23 // Example matrices
24 val matrixA = Array(Array(1, 2, 3), Array(4, 5, 6))
25 val matrixB = Array(Array(7, 8), Array(9, 10), Array(11, 12))
26
27 // Perform matrix multiplication
28 val resultMatrix = multiplyMatrix(matrixA, matrixB)
29

```

Basics of scala Scala

Last edit was 27 days ago New cell UI: OFF

Run all | Terminated | Share | Publish

```

22
23 // Example matrices
24 val matrixA = Array(Array(1, 2, 3), Array(4, 5, 6))
25 val matrixB = Array(Array(7, 8), Array(9, 10), Array(11, 12))
26
27 // Perform matrix multiplication
28 val resultMatrix = multiplyMatrix(matrixA, matrixB)
29
30 // Display the result matrix
31 println("Matrix A:")
32 matrixA.foreach(row => println(row.mkString("\t")))
33 println("\nMatrix B:")
34 matrixB.foreach(row => println(row.mkString("\t")))
35 println("\nResult Matrix:")
36 resultMatrix.foreach(row => println(row.mkString("\t")))
37

```

Matrix A:

1	2	3
4	5	6

Matrix B:

7	8
9	10
11	12

Result Matrix:

58	64
139	154

import scala.Array
multiplyMatrix: (matrix1: Array[Array[Int]], matrix2: Array[Array[Int]])Array[Array[Int]]
matrixA: Array[Array[Int]] = Array(Array(1, 2, 3), Array(4, 5, 6))
matrixB: Array[Array[Int]] = Array(Array(7, 8), Array(9, 10), Array(11, 12))
resultMatrix: Array[Array[Int]] = Array(Array(58, 64), Array(139, 154))

Command took 30.14 seconds -- by dhvanipate301003@gmail.com at 1/24/2024, 9:10:46 AM on 24/1

Name: Dhvani Patel

Roll no: 21BCP116

Big Data Analytics Lab

Practical 3 : Transformation Function

Databricks Notebook titled "Practical 3-Transformation Function" in Scala. The notebook has three cells:

- Cmd 1:**

```
1 val a=sc.parallelize(List("A","B","C","D"))
```

 Command took 0.51 seconds -- by dhvanipatel301003@gmail.com at 1/31/2024, 9:46:21 AM on 31/1
- Cmd 2:**

```
1 val b=a.map(x=>(x,1))
2 b.collect
3
```

 (1) Spark Jobs
b: org.apache.spark.rdd.RDD[(String, Int)] = MapPartitionsRDD[6] at map at command-1899821947669053:1
res5: Array[(String, Int)] = Array((A,1), (B,1), (C,1), (D,1))
Command took 0.54 seconds -- by dhvanipatel301003@gmail.com at 1/31/2024, 9:46:25 AM on 31/1
- Cmd 3:**

```
1 val b=a.map(_>1)
2 b.collect
3
```

 (1) Spark Jobs
b: org.apache.spark.rdd.RDD[(String, Int)] = MapPartitionsRDD[7] at map at command-4367539450063560:1
res6: Array[(String, Int)] = Array((A,1), (B,1), (C,1), (D,1))
Command took 0.61 seconds -- by dhvanipatel301003@gmail.com at 1/31/2024, 9:46:29 AM on 31/1

Databricks Notebook titled "Practical 3-Transformation Function" in Scala. The notebook has three cells:

- Cmd 4:**

```
1 val a=sc.parallelize(List(1,2,3,4,5)).map(x=>List(x,x+1,x+2))
2 a.collect
3
```

 (1) Spark Jobs
a: org.apache.spark.rdd.RDD[List[Int]] = MapPartitionsRDD[3] at map at command-4367539450063558:1
res3: Array[List[Int]] = Array(List(1, 2, 3), List(2, 3, 4), List(3, 4, 5), List(4, 5, 6), List(5, 6, 7))
Command took 1.47 seconds -- by dhvanipatel301003@gmail.com at 1/31/2024, 9:45:14 AM on 31/1
- Cmd 5:**

```
1 val a=sc.parallelize(List(1,2,3,4,5)).flatMap(x=>List(x,x+1,x+2))
2 a.collect
3
```

 (1) Spark Jobs
a: org.apache.spark.rdd.RDD[Int] = MapPartitionsRDD[9] at flatMap at command-4367539450063559:1
res7: Array[Int] = Array(1, 2, 3, 2, 3, 4, 3, 4, 5, 4, 5, 6, 5, 6, 7)
Command took 0.51 seconds -- by dhvanipatel301003@gmail.com at 1/31/2024, 9:47:10 AM on 31/1
- Cmd 6:**

```
1 val rdda = sc.parallelize(List("aaaa","bbbb","ccc"))
2 rdda.filter(_.equals("aaaa")).collect
3
```

 (1) Spark Jobs
rdda: org.apache.spark.rdd.RDD[String] = ParallelCollectionRDD[10] at parallelize at command-4367539450063561:1
res8: Array[String] = Array(aaaa)
Command took 1.09 seconds -- by dhvanipatel301003@gmail.com at 1/31/2024, 9:49:08 AM on 31/1

databricks

Practical 3-Transformation Function Scala ⚡

File Edit View Run Help Last edit was 20 days ago New cell UI: OFF

Run all Terminated Share Publish

Cmd 7

```
1 rdd.a.filter(_.contains("a")).collect

▶ (1) Spark Jobs
res9: Array[String] = Array(aaaa)
Command took 0.68 seconds -- by dhvanipatel301003@gmail.com at 1/31/2024, 9:54:54 AM on 31/1
```

Cmd 8

```
1 val a = sc.parallelize(List(("Mumbai",4000),("Delhi",2000),("Chennai",1000),("Kolkata",7000)))

a: org.apache.spark.rdd.RDD[(String, Int)] = ParallelCollectionRDD[13] at parallelize at command-4367539450063563:1
Command took 0.91 seconds -- by dhvanipatel301003@gmail.com at 1/31/2024, 9:57:14 AM on 31/1
```

Cmd 9

```
1 a.filter(_._1.contains("ai")).collect

▶ (1) Spark Jobs
res10: Array[(String, Int)] = Array((Mumbai,4000), (Chennai,1000))
Command took 0.63 seconds -- by dhvanipatel301003@gmail.com at 1/31/2024, 9:59:29 AM on 31/1
```

Cmd 10

```
1 //filter data less than 3000
2 a.filter(_._2.<(3000)).collect

▶ (1) Spark Jobs
res11: Array[(String, Int)] = Array((Delhi,2000), (Chennai,1000))
Command took 0.58 seconds -- by dhvanipatel301003@gmail.com at 1/31/2024, 10:00:25 AM on 31/1
```

The screenshot shows a Databricks notebook interface. The title bar reads "Practical 3-Transformation Function" and "Scala". The menu bar includes File, Edit, View, Run, Help, and a status message "Last edit was 20 days ago". The top right corner shows the user "dhvanipatel301003@gmail.com" and a help icon. Below the menu is a toolbar with icons for back, forward, run, and share. The main area contains two command cells:

Cmd 11

```
1 //filter data between 3000 and 6000
2 // a.filter(_._2 > (3000)).filter(_._2 <= (6000)).collect
3 a.filter(x => x._2 > (3000) && x._2 < (6000)).collect
```

(1) Spark Jobs

```
res14: Array[(String, Int)] = Array((Mumbai,4000))
```

Command took 0.48 seconds -- by dhvanipatel301003@gmail.com at 1/31/2024, 10:16:01 AM on 31/1

Cmd 12

```
1 val a = sc.parallelize(1 to 1000)
```

```
a: org.apache.spark.rdd.RDD[Int] = ParallelCollectionRDD[20] at parallelize at command-4367539450063567:1
```

Command took 0.38 seconds -- by dhvanipatel301003@gmail.com at 1/31/2024, 10:16:43 AM on 31/1

databricks

Practical 3-Transformation Function Scala

File Edit View Run Help Last edit was 20 days ago New cell UI: OFF Run all Terminated Share Publish

```
Cmd 14
1 a.sample(false,0.2).collect
(1) Spark Jobs
res16: Array[Int] = Array(1, 3, 15, 23, 28, 33, 50, 52, 54, 56, 63, 68, 69, 76, 80, 82, 83, 95, 102, 104, 105, 111, 121, 128, 130, 133, 135, 139, 147, 151, 154, 159, 160, 174, 179, 196, 201, 210, 215, 218, 220, 222, 228, 247, 250, 255, 258, 260, 262, 265, 268, 274, 276, 282, 291, 293, 294, 299, 305, 313, 314, 316, 319, 320, 325, 327, 329, 331, 333, 338, 341, 356, 360, 375, 386, 399, 405, 406, 408, 411, 424, 429, 437, 440, 445, 447, 457, 460, 471, 473, 480, 493, 494, 497, 506, 510, 513, 521, 524, 527, 530, 537, 539, 543, 544, 547, 549, 553, 554, 562, 564, 567, 572, 573, 577, 580, 581, 582, 584, 586, 587, 595, 600, 604, 607, 614, 632, 633, 635, 636, 637, 638, 650, 654, 657, 661, 668, 676, 690, 693, 694, 695, 700, 719, 724, 727, 734, 737, 738, 741, 743, 749, 752, 756, 758, 776, 777, 780, 783, 784, 787, 794, 798, 800, 811, 818, 823, 824, 826, 828, 839, 840, 843, 850, 855, 857, 863, 865, 866, 869, 880, 889, 890, 891, 898, 903, 908, 911, 917, 918, 921, 923, 924, 926, 927, 939, 943, 944, 945, 947, 948, 953, 957, 958, 968, 969, 971, 976, 978, 981, 982, 984, 994, 998, 999)
Command took 0.55 seconds -- by dhvanipatel301003@gmail.com at 1/31/2024, 10:17:33 AM on 31/1
```

```
Cmd 15
1 val temp = a.sample(false,0.2).collect
(1) Spark Jobs
temp: Array[Int] = Array(1, 8, 9, 24, 26, 27, 28, 32, 33, 41, 44, 46, 50, 55, 56, 60, 61, 62, 72, 73, 79, 93, 95, 99, 104, 105, 106, 112, 114, 116, 117, 118, 126, 145, 149, 155, 158, 161, 164, 167, 169, 170, 178, 179, 183, 185, 187, 189, 196, 198, 199, 204, 207, 210, 217, 228, 230, 232, 242, 244, 249, 250, 251, 267, 272, 277, 281, 288, 298, 304, 316, 320, 321, 323, 325, 339, 340, 342, 351, 354, 355, 356, 361, 362, 364, 365, 367, 374, 381, 383, 388, 399, 402, 406, 411, 418, 422, 432, 439, 443, 444, 447, 452, 457, 460, 475, 479, 488, 496, 498, 516, 519, 522, 528, 540, 541, 545, 572, 575, 581, 585, 590, 594, 598, 606, 616, 635, 646, 652, 659, 664, 675, 681, 686, 693, 695, 706, 712, 713, 716, 718, 720, 722, 729, 732, 738, 739, 746, 741, 743, 749, 756, 757, 766, 769, 776, 777, 780, 786, 787, 806, 815, 828, 829, 845, 846, 848, 849, 851, 854, 857, 858, 879, 885, 908, 913, 923, 924, 925, 927, 934, 937, 938, 940, 941, 942, 945, 955, 958, 967, 972, 977, 983, 984, 986, 987, 990, 995, 999)
Command took 1.18 seconds -- by dhvanipatel301003@gmail.com at 1/31/2024, 10:22:35 AM on 31/1
```

```
Cmd 16
1 temp.length
res17: Int = 203
Command took 0.40 seconds -- by dhvanipatel301003@gmail.com at 1/31/2024, 10:22:43 AM on 31/1
```

databricks

Practical 3-Transformation Function Scala

File Edit View Run Help Last edit was 20 days ago New cell UI: OFF Run all Terminated Share Publish

```
Cmd 17
1 val a = sc.parallelize(List(1,2,1,1,1,2))
a: org.apache.spark.rdd.RDD[Int] = ParallelCollectionRDD[23] at parallelize at command-4367539450063572:1
Command took 0.31 seconds -- by dhvanipatel301003@gmail.com at 1/31/2024, 10:24:32 AM on 31/1
```

```
Cmd 18
1 a.sample(true,0.5).collect
(1) Spark Jobs
res18: Array[Int] = Array(1, 1)
Command took 0.61 seconds -- by dhvanipatel301003@gmail.com at 1/31/2024, 10:24:52 AM on 31/1
```

```
Cmd 19
1 val a = sc.parallelize(1 to 7)
a: org.apache.spark.rdd.RDD[Int] = ParallelCollectionRDD[33] at parallelize at command-4367539450063578:1
Command took 0.81 seconds -- by dhvanipatel301003@gmail.com at 1/31/2024, 10:27:29 AM on 31/1
```

```
Cmd 20
1 val b = sc.parallelize(5 to 10)
b: org.apache.spark.rdd.RDD[Int] = ParallelCollectionRDD[34] at parallelize at command-4367539450063574:1
Command took 0.30 seconds -- by dhvanipatel301003@gmail.com at 1/31/2024, 10:27:32 AM on 31/1
```

databricks

Practical 3-Transformation Function Scala

File Edit View Run Help Last edit was 20 days ago New cell UI: OFF Run all Terminated Share Publish

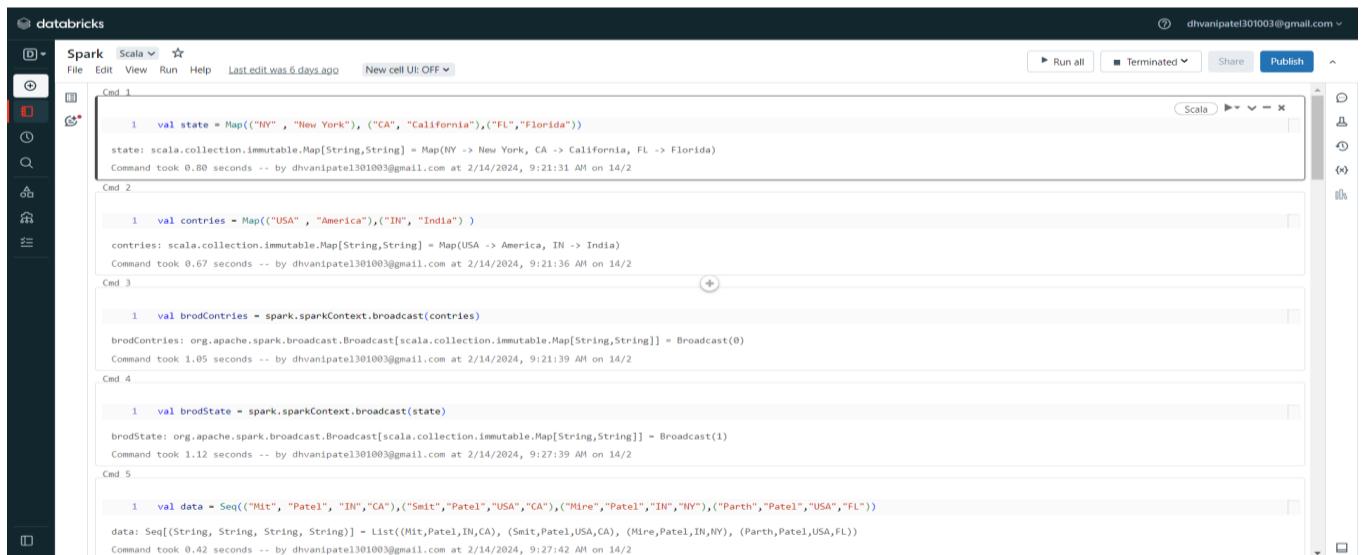
```
Cmd 21
1 a.union(b).collect
(1) Spark Jobs
res21: Array[Int] = Array(1, 2, 3, 4, 5, 6, 7, 5, 6, 7, 8, 9, 10)
Command took 0.45 seconds -- by dhvanipatel301003@gmail.com at 1/31/2024, 10:27:34 AM on 31/1
```

```
Cmd 22
1 a.intersection(b).collect
(1) Spark Jobs
res22: Array[Int] = Array(5, 6, 7)
Command took 0.68 seconds -- by dhvanipatel301003@gmail.com at 1/31/2024, 10:27:36 AM on 31/1
```

```
Cmd 23
1 a.distinct(1).collect
(1) Spark Jobs
res23: Array[Int] = Array(4, 1, 6, 3, 7, 5, 2)
Command took 1.77 seconds -- by dhvanipatel301003@gmail.com at 1/31/2024, 10:29:50 AM on 31/1
```

```
Cmd 24
1
```

Name: Dhvani Patel
Roll no: 21BCP116
Big Data Analytics Lab
Practical 4 : Spark



Databricks Notebook interface showing Scala code execution. The code defines state, countries, broadcasted countries, broadcasted state, and a sequence of data. The output shows the creation of Broadcast objects and their values.

```
1 val state = Map("NY", "New York"), ("CA", "California"), ("FL", "Florida"))
state: scala.collection.immutable.Map[String, String] = Map(NY -> New York, CA -> California, FL -> Florida)
Command took 0.80 seconds -- by dhvanipatel301003@gmail.com at 2/14/2024, 9:21:31 AM on 14/2

1 val countries = Map(("USA", "America"), ("IN", "India"))
countries: scala.collection.immutable.Map[String, String] = Map(USA -> America, IN -> India)
Command took 0.67 seconds -- by dhvanipatel301003@gmail.com at 2/14/2024, 9:21:36 AM on 14/2

1 val brodCountries = spark.sparkContext.broadcast(countries)
brodCountries: org.apache.spark.broadcast.Broadcast[scala.collection.immutable.Map[String, String]] = Broadcast(0)
Command took 1.05 seconds -- by dhvanipatel301003@gmail.com at 2/14/2024, 9:21:39 AM on 14/2

1 val brodState = spark.sparkContext.broadcast(state)
brodState: org.apache.spark.broadcast.Broadcast[scala.collection.immutable.Map[String, String]] = Broadcast(1)
Command took 1.12 seconds -- by dhvanipatel301003@gmail.com at 2/14/2024, 9:27:39 AM on 14/2

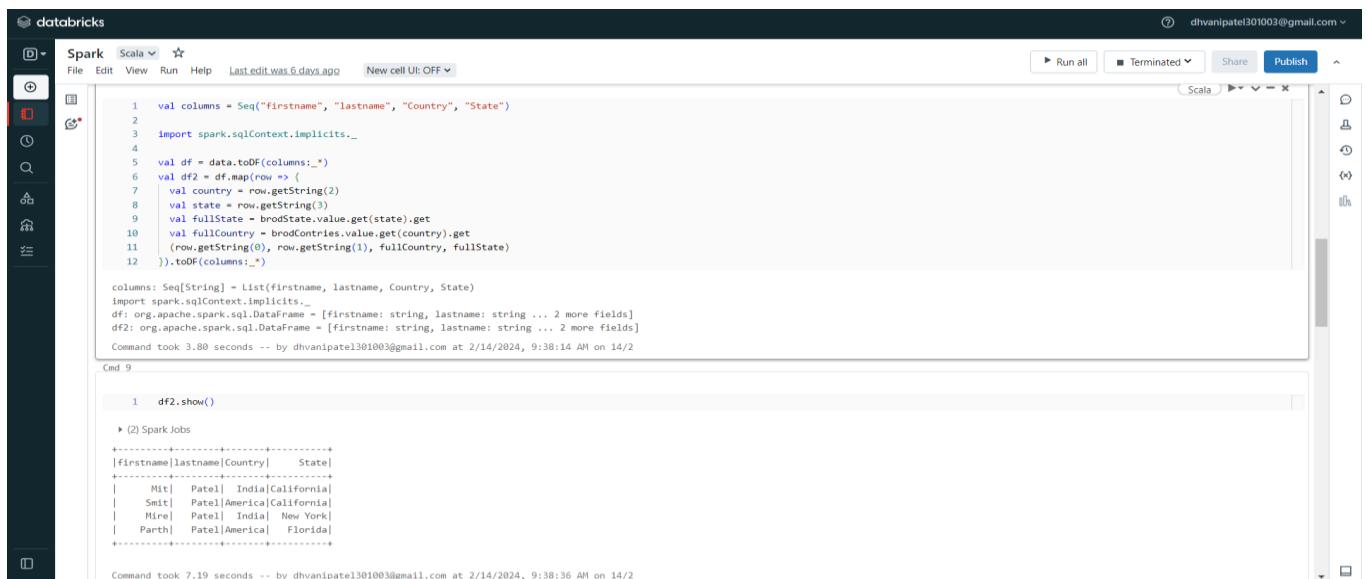
1 val data = Seq(("Mit", "Patel", "IN", "CA"), ("Smit", "Patel", "USA", "CA"), ("Mire", "Patel", "IN", "NY"), ("Parth", "Patel", "USA", "FL"))
data: Seq[(String, String, String, String)] = List((Mit, Patel, IN, CA), (Smit, Patel, USA, CA), (Mire, Patel, IN, NY), (Parth, Patel, USA, FL))
Command took 0.42 seconds -- by dhvanipatel301003@gmail.com at 2/14/2024, 9:27:42 AM on 14/2
```



Databricks Notebook interface showing Scala code execution. The code parallelizes data and maps each row with broadcasted country, state, and full state information. The output shows the creation of an RDD and its mapped version.

```
1 val rdd = spark.sparkContext.parallelize(data)
rdd: org.apache.spark.rdd.RDD[(String, String, String, String)] = ParallelCollectionRDD[1] at parallelize at command-4463440681818442:1
Command took 0.42 seconds -- by dhvanipatel301003@gmail.com at 2/14/2024, 9:27:44 AM on 14/2

1 val rdd2 = rdd.map(f =>
  2   val country = f._3
  3   val state = f._4
  4   val fullCountry = brodCountries.value.get(country).get
  5   val fullState = brodState.value.get(state).get
  6   (f._1, f._2, fullCountry, fullState)
  7 )
rdd2: org.apache.spark.rdd.RDD[(String, String, String, String)] = MapPartitionsRDD[3] at map at command-4463440681818438:1
Command took 1.12 seconds -- by dhvanipatel301003@gmail.com at 2/14/2024, 9:37:00 AM on 14/2
```



Databricks Notebook interface showing Scala code execution. The code reads a CSV file, imports implicits, and performs a join operation using broadcast variables. The output shows the creation of a DataFrame and its display.

```
1 val columns = Seq("firstname", "lastname", "Country", "State")
2
3 import spark.sqlContext.implicits._
4
5 val df = data.toDF(columns:_*)
6 val df2 = df.map(row => {
  7   val country = row.getString(2)
  8   val state = row.getString(3)
  9   val fullState = brodState.value.get(state).get
10   val fullCountry = brodCountries.value.get(country).get
11   (row.getString(0), row.getString(1), fullCountry, fullState)
12 }).toDF(columns:_*)
columns: Seq[String] = List(firstname, lastname, Country, State)
import spark.sqlContext.implicits._
df: org.apache.spark.sql.DataFrame = [firstname: string, lastname: string ... 2 more fields]
df2: org.apache.spark.sql.DataFrame = [firstname: string, lastname: string ... 2 more fields]
Command took 3.80 seconds -- by dhvanipatel301003@gmail.com at 2/14/2024, 9:38:14 AM on 14/2

1 df2.show()
#(2) Spark Jobs
+-----+-----+-----+-----+
|firstname|lastname|Country|State|
+-----+-----+-----+-----+
| Mit | Patel | India | California |
| Smit | Patel | America | California |
| Mire | Patel | India | New York |
| Parth | Patel | America | Florida |
+-----+-----+-----+-----+
Command took 7.19 seconds -- by dhvanipatel301003@gmail.com at 2/14/2024, 9:38:36 AM on 14/2
```

dhvanipatel301003@gmail.com

Spark Scala

File Edit View Run Help Last edit was 6 days ago New cell UI: OFF

Run all Terminated Share Publish

Cmd 10

```
1 val longAcc = spark.sparkContext.longAccumulator("SUM")  
longAcc: org.apache.spark.util.LongAccumulator = LongAccumulator(id: 111, name: Some(SUM), value: 0)  
Command took 0.94 seconds -- by dhvanipatel301003@gmail.com at 2/14/2024, 9:39:49 AM on 14/2
```

Cmd 11

```
1 val rdd = spark.sparkContext.parallelize(Array(1, 2, 3, 4, 5))  
2 rdd.foreach(x => longAcc.add(x))  
3 rdd.collect  
(2) Spark Jobs  
rdd: org.apache.spark.rdd.RDD[Int] = ParallelCollectionRDD[6] at parallelize at command-4463440681818450:1  
res3: Array[Int] = Array(1, 2, 3, 4, 5)  
Command took 1.57 seconds -- by dhvanipatel301003@gmail.com at 2/14/2024, 9:44:38 AM on 14/2
```

Cmd 12

```
1 longAcc.value  
res4: Long = 15  
Command took 0.29 seconds -- by dhvanipatel301003@gmail.com at 2/14/2024, 9:44:42 AM on 14/2
```

Cmd 13

```
1 spark.sparkContext.setLogLevel("Error")  
Command took 0.42 seconds -- by dhvanipatel301003@gmail.com at 2/14/2024, 9:44:45 AM on 14/2
```

dhvanipatel301003@gmail.com

Spark Scala

File Edit View Run Help Last edit was 6 days ago New cell UI: OFF

Run all Terminated Share Publish

Cmd 14

```
1 val inputRDD = spark.sparkContext.parallelize(List(("2", 1), ("B", 30), ("A", 20), ("B", 30), ("C", 40), ("B", 60))  
inputRDD: org.apache.spark.rdd.RDD[(String, Int)] = ParallelCollectionRDD[9] at parallelize at command-4463440681818449:1  
Command took 0.33 seconds -- by dhvanipatel301003@gmail.com at 2/14/2024, 9:47:24 AM on 14/2
```

Cmd 15

```
1 val listRDD = spark.sparkContext.parallelize(List(1, 2, 3, 4, 5, 2, 3))  
listRDD: org.apache.spark.rdd.RDD[Int] = ParallelCollectionRDD[10] at parallelize at command-4463440681818451:1  
Command took 0.31 seconds -- by dhvanipatel301003@gmail.com at 2/14/2024, 9:47:27 AM on 14/2
```

Cmd 16

```
1 def param0 = (acc:Int, v:Int) => acc + v  
2 def param1 = (acc1:Int, acc2:Int) => acc1+acc2  
3 println("Aggregate: " + listRDD.aggregate(0) (param0, param1))  
(1) Spark Jobs  
Aggregate: 20  
param0: (Int, Int) => Int  
param1: (Int, Int) => Int  
Command took 1.19 seconds -- by dhvanipatel301003@gmail.com at 2/14/2024, 9:50:22 AM on 14/2
```

dhvanipatel301003@gmail.com

Spark Scala

File Edit View Run Help Last edit was 6 days ago New cell UI: OFF

Run all Terminated Share Publish

Command took 1.19 seconds -- by dhvanipatel301003@gmail.com at 2/14/2024, 9:50:22 AM on 14/2

Cmd 17

```
1 def param3 = (acc:Int, v:(String, Int)) => acc + v._2  
2 def param2 = (acc1: Int, v2: Int) => acc1 + v2  
3 println("Aggregate: " + inputRDD.aggregate(0) (param3, param2))  
(1) Spark Jobs  
Aggregate: 181  
param3: (Int, (String, Int)) => Int  
param2: (Int, Int) => Int  
Command took 0.53 seconds -- by dhvanipatel301003@gmail.com at 2/14/2024, 9:50:34 AM on 14/2
```

Cmd 18

```
1
```

[Shift+Enter] to run and move to next cell
[Esc H] to see all keyboard shortcuts

Name: Dhvani Patel
Roll no: 21BCP116
Big Data Analytics Lab
Practical 5: Spark SQL



Databricks Notebook titled "Practical 5 Spark sql" in Scala. The code creates an RDD "a" from 1 to 10, then collects it into an array "res2".

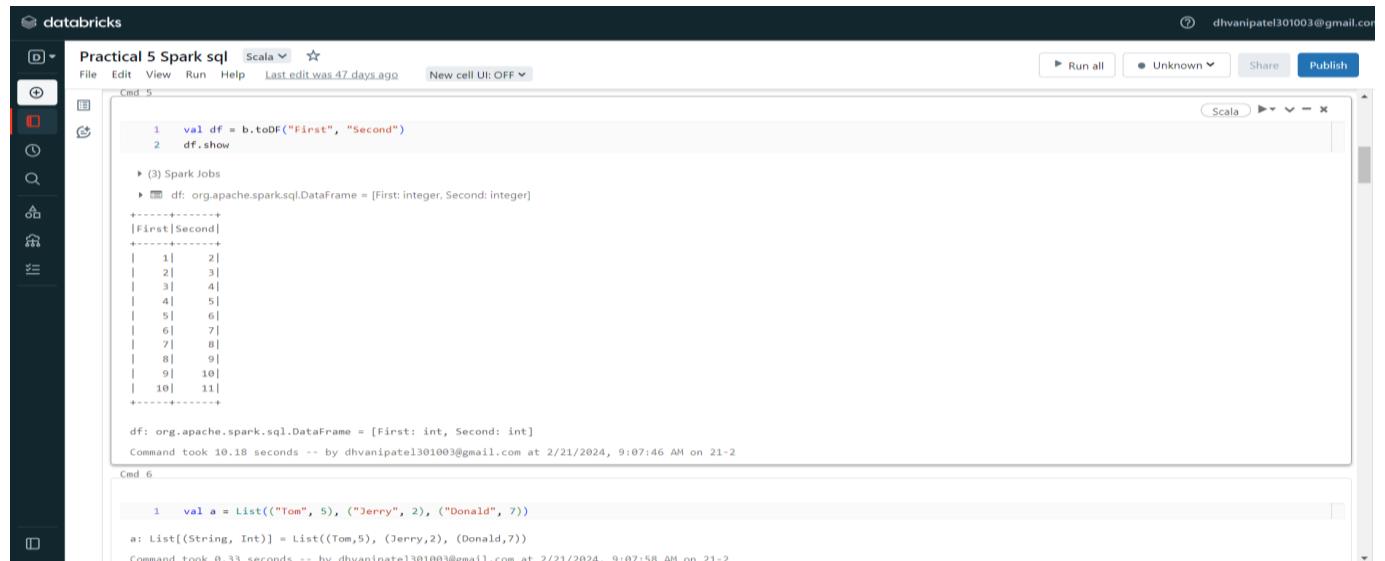
```
1 val a = sc.parallelize(1 to 10)
a: org.apache.spark.rdd.RDD[Int] = ParallelCollectionRDD[0] at parallelize at command-4463440681818456:1
Command took 0.79 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:07:13 AM on 21-2

Cmd. 3

1 a.collect
▶ (1) Spark Jobs
res2: Array[Int] = Array(1, 2, 3, 4, 5, 6, 7, 8, 9, 10)
Command took 2.52 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:07:37 AM on 21-2

Cmd. 4

1 val b = a.map(x=>(x,x+1))
2 b.collect
▶ (1) Spark Jobs
b: org.apache.spark.rdd.RDD[(Int, Int)] = MapPartitionsRDD[1] at map at command-4463440681818459:1
res3: Array[(Int, Int)] = Array((1,2), (2,3), (3,4), (4,5), (5,6), (6,7), (7,8), (8,9), (9,10), (10,11))
Command took 0.75 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:07:43 AM on 21-2
```

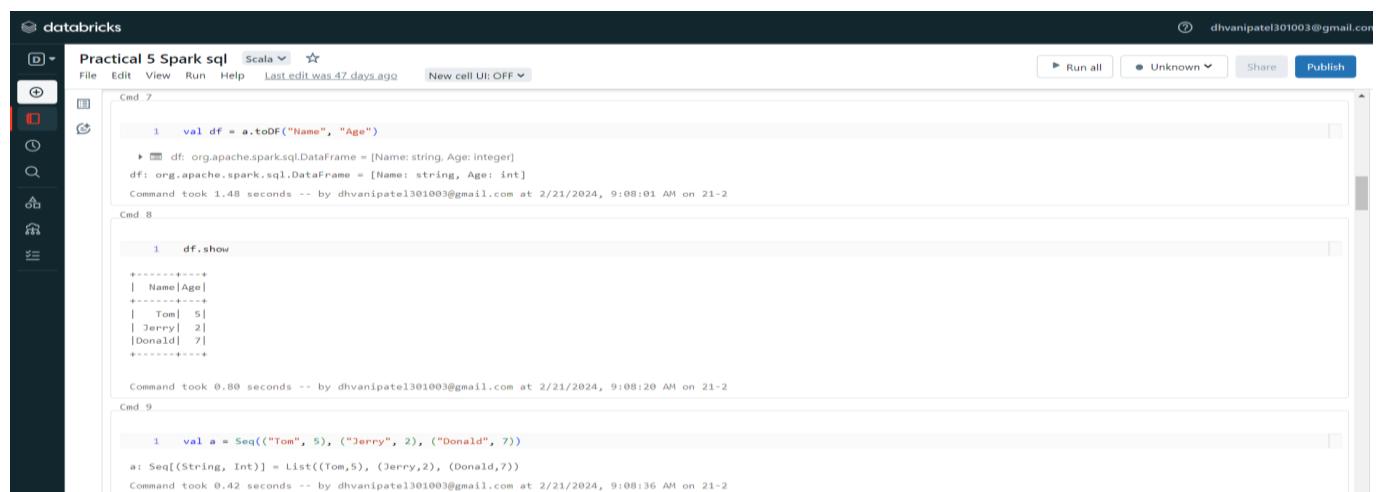


Databricks Notebook titled "Practical 5 Spark sql" in Scala. The code converts a list of tuples into a DataFrame "df" and shows its contents.

```
1 val df = b.toDF("First", "Second")
2 df.show
▶ (3) Spark Jobs
df: org.apache.spark.sql.DataFrame = [First: integer, Second: integer]
+-----+-----+
|First|Second|
+-----+-----+
| 1 | 2 |
| 2 | 3 |
| 3 | 4 |
| 4 | 5 |
| 5 | 6 |
| 6 | 7 |
| 7 | 8 |
| 8 | 9 |
| 9 | 10 |
| 10 | 11 |
+-----+-----+
df: org.apache.spark.sql.DataFrame = [First: int, Second: int]
Command took 10.18 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:07:46 AM on 21-2

Cmd. 6

1 val a = List(("Tom", 5), ("Jerry", 2), ("Donald", 7))
a: List[(String, Int)] = List((Tom,5), (Jerry,2), (Donald,7))
Command took 0.33 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:07:58 AM on 21-2
```



Databricks Notebook titled "Practical 5 Spark sql" in Scala. The code converts a list of tuples into a DataFrame "df" and shows its contents.

```
1 val df = a.toDF("Name", "Age")
2 df.show
▶ (3) Spark Jobs
df: org.apache.spark.sql.DataFrame = [Name: string, Age: integer]
df: org.apache.spark.sql.DataFrame = [Name: string, Age: int]
Command took 1.48 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:08:01 AM on 21-2

Cmd. 8

1 df.show
+-----+-----+
| Name|Age|
+-----+-----+
| Tom| 5 |
| Jerry| 2 |
| Donald| 7 |
+-----+-----+
Command took 0.80 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:08:20 AM on 21-2

Cmd. 9

1 val a = Seq(("Tom", 5), ("Jerry", 2), ("Donald", 7))
a: Seq[(String, Int)] = List((Tom,5), (Jerry,2), (Donald,7))
Command took 0.42 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:08:36 AM on 21-2
```

dataricks Practical 5 Spark sql Scala ⚡

File Edit View Run Help Last edit was 47 days ago New cell UI: OFF

Cmd 10

```
1 val df = a.toDF("Name", "Age")  
2 df: org.apache.spark.sql.DataFrame = [Name: string, Age: integer]  
df: org.apache.spark.sql.DataFrame = [Name: string, Age: int]  
Command took 0.85 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:08:46 AM on 21-2
```

Cmd 11

```
1 df.show  
  
+---+---+  
| Name|Age|  
+---+---+  
| Tom | 5 |  
| Jerry| 2 |  
| Donald| 7 |  
+---+---+  
  
Command took 0.50 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:08:54 AM on 21-2
```

Cmd 12

```
1 df.registerTempTable("Cartoon")  
command-23613877692441:1: warning: method registerTempTable in class Dataset is deprecated (since 2.0.0): Use createOrReplaceTempView(viewName) instead.  
df.registerTempTable("Cartoon")  
^  
  
Command took 0.48 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:09:01 AM on 21-2
```

dataricks Practical 5 Spark sql Scala ⚡

File Edit View Run Help Last edit was 47 days ago New cell UI: OFF

Cmd 13

```
1 df.createOrReplaceTempView("Cartoon")  
Command took 0.36 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:09:09 AM on 21-2
```

Cmd 14

```
1 sqlContext.sql("select * from Cartoon where Name='Tom'").show  
  
+---+---+  
| Name|Age|  
+---+---+  
| Tom | 5 |  
+---+---+  
  
Command took 0.56 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:09:22 AM on 21-2
```

Cmd 15

```
1 sqlContext.sql("select count(*) from Cartoon").show  
  (2) Spark Jobs  
  +-----+  
  | count(1)|  
  +-----+  
  | 3|  
  +-----+  
  
Command took 2.62 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:09:31 AM on 21-2
```

dataricks Practical 5 Spark sql Scala ⚡

File Edit View Run Help Last edit was 47 days ago New cell UI: OFF

Cmd 16

```
1 // Question: To create a JSON File, upload to DBFS,  
2 // printSchema() select query with all names filter and identify age > 23 groupBy Age count it and show it  
Command took 0.79 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:10:23 AM on 21-2
```

Cmd 17

```
1 val df1 = spark.read.format("json").load("dbfs:/FileStore/shared_uploads/dhvanipatel301003@gmail.com/data.json")  
  (1) Spark Jobs  
  +-----+  
  | df1: org.apache.spark.sql.DataFrame = [Age: long, id: long ... 1 more field]|  
  df1: org.apache.spark.sql.DataFrame = [Age: bigint, id: bigint ... 1 more field]  
Command took 4.80 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:16:20 AM on 21-2
```

Cmd 18

```
1 df1.show  
  (1) Spark Jobs  
  +---+---+  
  | Age| id | name|  
  +---+---+  
  | 27 |1201| Kriti|  
  | 22 |1202| Utkarsh|  
  | 25 |1203| Sneha|  
  | 20 |1204| Harsh|  
  | 35 |1205| Sona|  
  +---+---+
```

dhvanipatel301003@gmail.com

Practical 5 Spark sql Scala

File Edit View Run Help Last edit was 47 days ago New cell UI: OFF

Cmd 19

```
1 df1.printSchema()
```

root
|-- Age: long (nullable = true)
|-- id: long (nullable = true)
|-- name: string (nullable = true)

Command took 0.71 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:17:16 AM on 21-2

Cmd 20

```
1 df1.select("Name", "Age").show()
```

(1) Spark Jobs

	Name	Age
1	Kriti	27
2	Utkarsh	22
3	Sneha	25
4	Harsh	20
5	Sonal	35

Command took 1.21 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:17:25 AM on 21-2

dhvanipatel301003@gmail.com

Practical 5 Spark sql Scala

File Edit View Run Help Last edit was 47 days ago New cell UI: OFF

Cmd 21

```
1 df1.createOrReplaceTempView("Employee")
```

Command took 0.46 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:17:35 AM on 21-2

Cmd 22

```
1 df1.filter("age > '23'").show()
```

(1) Spark Jobs

Age	id	name
27	1201	Kriti
25	1203	Sneha
35	1205	Sonal

Command took 1.19 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:17:46 AM on 21-2

dhvanipatel301003@gmail.com

Practical 5 Spark sql Scala

File Edit View Run Help Last edit was 47 days ago New cell UI: OFF

Cmd 23

```
1 df1.filter(df1("age") > 23).show()
```

(1) Spark Jobs

Age	id	name
27	1201	Kriti
25	1203	Sneha
35	1205	Sonal

Command took 0.93 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:17:58 AM on 21-2

Cmd 24

```
1 df1.groupBy("age").count().show
```

(2) Spark Jobs

age	count
22	1
25	1
27	1
35	1
20	1

Command took 2.58 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:18:05 AM on 21-2

databricks dhvanipatel301003@gmail.com

Practical 5 Spark sql Scala ⚡

File Edit View Run Help Last edit was 47 days ago New cell UI: OFF

Cmd 25

```
1 val rdda = sc.parallelize(1 to 1000)
2 rdda.collect()
3 val rddb = sc.parallelize(List("BMW", "Mercedes", "Toyota", "Audi"))
4 rddb.collect()
```

▶ (2) Spark Jobs

```
rdda: org.apache.spark.rdd.RDD[Int] = ParallelCollectionRDD[31] at parallelize at command-236138776692454:1
rddb: org.apache.spark.rdd.RDD[String] = ParallelCollectionRDD[32] at parallelize at command-236138776692454:3
res23: Array[String] = Array(BMW, Mercedes, Toyota, Audi)
```

Command took 0.97 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:18:20 AM on 21-2

Cmd 26

```
1 rdda.partitions.length
```

res24: Int = 8

Command took 0.42 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:18:28 AM on 21-2

Cmd 27

```
1 rddb.partitions.length
```

res25: Int = 8

Command took 0.37 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:18:33 AM on 21-2

databricks dhvanipatel301003@gmail.com

Practical 5 Spark sql Scala ⚡

File Edit View Run Help Last edit was 47 days ago New cell UI: OFF

Cmd 28

```
1 val rdda = sc.parallelize(1 to 1000, 10)
2 rdda.collect()
3 rdda.partitions.length
```

▶ (1) Spark Jobs

```
rdda: org.apache.spark.rdd.RDD[Int] = ParallelCollectionRDD[33] at parallelize at command-236138776692457:1
res26: Int = 10
```

Command took 0.47 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:18:40 AM on 21-2

Cmd 29

```
1 rdda.take(10)
```

▶ (1) Spark Jobs

```
res27: Array[Int] = Array(1, 2, 3, 4, 5, 6, 7, 8, 9, 10)
```

Command took 0.41 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:18:51 AM on 21-2

Cmd 30

```
1 rdda.count()
```

▶ (1) Spark Jobs

```
res38: Long = 1000
```

Command took 0.40 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:31:16 AM on 21-2

databricks dhvanipatel301003@gmail.com

Practical 5 Spark sql Scala ⚡

File Edit View Run Help Last edit was 47 days ago New cell UI: OFF

Cmd 31

```
1 rdda.saveAsTextFile("dbfs:/FileStore/shared_uploads/dhvanipatel301003@gmail.com/random.txt")
```

FileAlreadyExistsException: Output directory dbfs:/FileStore/shared_uploads/dhvanipatel301003@gmail.com/random.txt already exists

Command took 0.34 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:31:53 AM on 21-2

Cmd 32

```
1 val rddRead = sc.textFile("dbfs:/FileStore/shared_uploads/dhvanipatel301003@gmail.com/random.txt")
```

rddRead: org.apache.spark.rdd.RDD[String] = dbfs:/FileStore/shared_uploads/dhvanipatel301003@gmail.com/random.txt MapPartitionsRDD[54] at textFile at command-236138776692461:1

Command took 0.29 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:31:54 AM on 21-2

Cmd 33

```
1 rddRead.count()
```

res41: Long = 0

Command took 0.29 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:31:58 AM on 21-2

Cmd 34

```
1 rddRead.take(10)
```

res42: Array[String] = Array()

Command took 0.23 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:32:01 AM on 21-2

dhvanipatel301003@gmail.com

spark sql continue Scala

```
1 print("Dhvani patel")
Dhvani patel
Command took 29.05 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:04:23 AM on 21-2
```

```
1 val book = sc.textFile("dbfs:/FileStore/shared_uploads/dhvanipatel301003@gmail.com/book.txt")
book: org.apache.spark.rdd.RDD[String] = dbfs:/FileStore/shared_uploads/dhvanipatel301003@gmail.com/book.txt MapPartitionsRDD[45] at textFile at command-236138776692434:1
Command took 0.32 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:30:42 AM on 21-2
```

```
1 book.collect
▶ (1) Spark Jobs
res2: Array[String] = Array(hello, dhvani, patel)
Command took 0.87 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:30:46 AM on 21-2
```

```
1 val a = book.flatMap(x=>x.split(" ")).collect
▶ (1) Spark Jobs
a: Array[String] = Array(hello, dhvani, patel)
Command took 1.53 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:35:41 AM on 21-2
```

dhvanipatel301003@gmail.com

spark sql continue Scala

```
1 val b = a.map(y=>(y,1))
b: Array[(String, Int)] = Array((hello,1), (dhvani,1), (patel,1))
Command took 0.42 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:35:47 AM on 21-2
```

```
1 val a = book.flatMap(x=>x.split(" ")).map(y=>(y,1)).reduceByKey((x,y) => (x+y)).collect
▶ (1) Spark Jobs
a: Array[(String, Int)] = Array((patel,1), (dhvani,1), (hello,1))
Command took 1.99 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:37:22 AM on 21-2
```

```
1 val a = book.flatMap(x=>x.split(" ")).map(y=>(y,1)).reduceByKey((x,y) => (x+y)).sortBy(_.value,false).collect
▶ (2) Spark Jobs
a: Array[(String, Int)] = Array((patel,1), (hello,1), (dhvani,1))
Command took 1.25 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:38:22 AM on 21-2
```

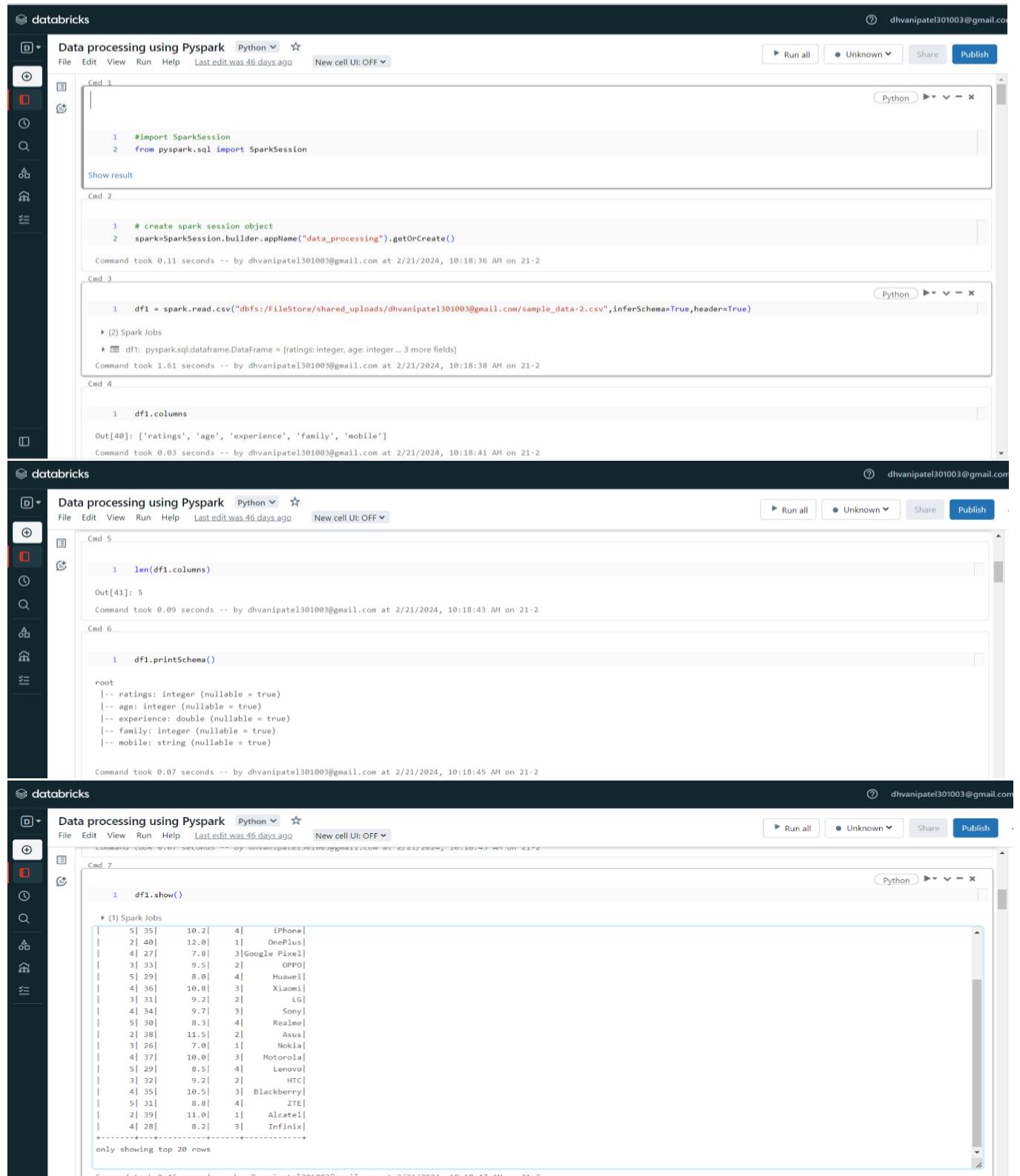
dhvanipatel301003@gmail.com

spark sql continue Scala

```
1 val a = book.flatMap(x=>x.split(" ")).map(y=>(y,1)).reduceByKey((x,y) => (x+y)).sortBy(_.value,false).take(1)
▶ (3) Spark Jobs
a: Array[(String, Int)] = Array((patel,1))
Command took 1.61 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:40:56 AM on 21-2
```

```
1 val a = book.flatMap(x=>x.split(" ")).map(y=>(y,1)).reduceByKey((x,y) => (x+y)).sortBy(_.value,false).filter(_.value > 1).collect
▶ (2) Spark Jobs
a: Array[(String, Int)] = Array()
Command took 1.03 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 9:41:32 AM on 21-2
```

Name: Dhvani Patel
Roll no: 21BCP116
Big Data Analytics Lab
Practical 6 : Data Preprocessing using Pyspark



The image shows three screenshots of Databricks notebooks titled "Data processing using Pyspark" in Python. The first notebook (Cmd 1) imports SparkSession and creates a spark session object. The second notebook (Cmd 2) reads a CSV file named "sample_data-2.csv" into a DataFrame df1. The third notebook (Cmd 3) prints the columns of df1 and its schema. The fourth notebook (Cmd 4) shows the top 20 rows of df1.

```
1 #import SparkSession
2 from pyspark.sql import SparkSession

1 # create spark session object
2 spark=SparkSession.builder.appName("data_processing").getOrCreate()

df1 = spark.read.csv("dbfs:/FileStore/shared_uploads/dhvani Patel301003@gmail.com/sample_data-2.csv",inferSchema=True,header=True)

df1.columns

len(df1.columns)

df1.printSchema()

df1.show()
```

Index	Age	Experience	Family	Mobile
1	35	10.2	4	iPhone
2	40	12.0	1	OnePlus
3	27	7.8	3	Google Pixel
4	33	9.5	2	OPPO
5	29	8.0	4	Huawei
6	36	10.8	3	Xiaomi
7	31	9.2	2	LG
8	34	9.7	3	Sony
9	30	8.3	4	Realme
10	38	11.5	2	Asus
11	26	7.0	1	Nokia
12	37	10.0	3	Motorola
13	29	8.5	4	Lenovo
14	32	9.2	2	HTC
15	35	10.5	3	Blackberry
16	31	8.8	4	ZTE
17	39	11.0	1	Alcatel
18	28	8.2	3	Infinix

databricks

Data processing using Pyspark Python

File Edit View Run Help Last edit was 46 days ago New cell UI: OFF

Run all Unknown Share Publish

Cmd 8

```
1 df1.select('ratings','age','mobile').show(5)
```

▶ (1) Spark Jobs

	ratings	age	mobile
1	3 32	Vivo	
2	4 28	Samsung	
3	5 35	iPhone	
4	2 40	OnePlus	
5	4 27	Google Pixel	

only showing top 5 rows

Command took 0.38 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:18:49 AM on 21-2

databricks

Data processing using Pyspark Python

File Edit View Run Help Last edit was 46 days ago New cell UI: OFF

Run all Unknown Share Publish

Cmd 9

```
1 df1.describe().show()
```

▶ (2) Spark Jobs

summary	ratings	age	experience	family	mobile
count	47	47	47	47	47
mean	3.6808510638297873	32.40425531914894	9.295744680851064	2.6808510638297873	null
stddev	1.023769313747283	[3.848543408548855]	1.1760208231618	1.023769313747283	null
min	2	26	7.0	1 Alcatel	
max	5	40	12.0	4 iPhone	

Command took 1.70 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:18:50 AM on 21-2

Cmd 10

```
1 from pyspark.sql.types import StringType,DoubleType,IntegerType
```

Command took 0.06 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:18:54 AM on 21-2

databricks

Data processing using Pyspark Python

File Edit View Run Help Last edit was 46 days ago New cell UI: OFF

Run all Unknown Share Publish

Cmd 11

```
1 df1.withColumn("age_after_10_yrs", (df1["age"]+10)).show(10,False)
```

▶ (1) Spark Jobs

	ratings	age	experience	family	mobile	age_after_10_yrs
1	3 32	9.0	3 Vivo	42		
2	4 28	8.5	2 Samsung	38		
3	5 35	10.2	4 iPhone	45		
4	2 40	12.0	1 OnePlus	50		
5	4 27	7.8	3 Google Pixel	37		
3	3 33	9.5	2 OPPO	43		
5	2 29	8.0	4 Huawei	39		
4	3 36	10.8	3 Xiaomi	46		
3	3 31	9.2	2 LG	41		
4	3 34	9.7	3 Sony	44		

only showing top 10 rows

Command took 0.52 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:19:30 AM on 21-2

databricks

Data processing using Pyspark Python ⭐

File Edit View Run Help Last edit was 46 days ago New cell UI: OFF

Cmd 12

```
1 df1.withColumn("age_double", (df1["age"]).cast(DoubleType())).show(10, False)
```

▶ (1) Spark Jobs

	ratings	age	experience	family	mobile	age_double
3	32	9.0	3	Vivo	32.0	
4	28	8.5	2	Samsung	28.0	
5	35	10.2	4	iPhone	35.0	
2	40	12.0	1	OnePlus	40.0	
4	27	7.8	3	Google Pixel	27.0	
3	33	9.5	2	OPPO	33.0	
5	29	8.0	4	Huawei	29.0	
4	36	10.8	3	Xiaomi	36.0	
3	31	9.2	2	LG	31.0	
4	34	9.7	3	Sony	34.0	

only showing top 10 rows

Command took 0.55 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:20:04 AM on 21-2

databricks

Data processing using Pyspark Python ⭐

File Edit View Run Help Last edit was 46 days ago New cell UI: OFF

Cmd 13

```
1 df.filter(df["mobile"]=="OnePlus").show()
```

▶ (1) Spark Jobs

	ratings	age	experience	family	mobile
2 40 12.0 1 OnePlus					
5 29 8.7 4 OnePlus					
4 37 10.8 3 OnePlus					

Command took 0.51 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:20:06 AM on 21-2

Cmd 14

```
1 # filter the records
2 df.filter(df["mobile"]=="Oneplus").select("age","ratings","mobile").show()
```

▶ (1) Spark Jobs

	age	ratings	mobile

Command took 0.42 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:20:09 AM on 21-2

databricks

Data processing using Pyspark Python ⭐

File Edit View Run Help Last edit was 46 days ago New cell UI: OFF

Cmd 15

```
1 # filter multiple condition
2 df1.filter(df1["mobile"]=="OnePlus").filter(df1["experience"]>10).show()
```

▶ (1) Spark Jobs

	ratings	age	experience	family	mobile
2 40 12.0 1 OnePlus					

Command took 0.58 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:20:12 AM on 21-2

Cmd 16

```
1 df1.filter((df1['mobile']=='OnePlus')&(df1['experience']>10)).show()
```

▶ (1) Spark Jobs

	ratings	age	experience	family	mobile
2 40 12.0 1 OnePlus					

Command took 0.35 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:20:15 AM on 21-2

databricks

Data processing using Pyspark Python ⭐

File Edit View Run Help Last edit was 46 days ago New cell UI: OFF

Cmd 17

```
1 df1.select("mobile").distinct().count()
```

▶ (3) Spark Jobs

Out[55]: 25

Command took 1.24 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:20:20 AM on 21-2

Cmd 18

```
1 df1.groupBy("mobile").count().show(5,False)
```

▶ (2) Spark Jobs

mobile	count
Infinix	2
Nokia	2
Sony	2
Alcatel	2
Motorola	2

only showing top 5 rows

Command took 0.94 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:20:23 AM on 21-2

databricks

Data processing using Pyspark Python ⭐

File Edit View Run Help Last edit was 46 days ago New cell UI: OFF

Cmd 19

```
1 df1.groupBy("mobile").mean().show(10,False)
```

▶ (2) Spark Jobs

mobile	avg(ratings)	avg(age)	avg(experience)	avg(family)
Infinix	3.0	32.5	9.25	2.0
Nokia	3.5	27.0	7.6	2.0
Sony	3.0	36.5	10.6	2.0
Alcatel	3.5	34.0	9.85	2.5
Motorola	3.5	35.0	9.9	2.5
OPPO	2.5	35.0	9.9	1.5
Realme	4.0	31.0	8.9	3.0
iPhone	5.0	32.5	9.35	4.0
Huawei	4.0	27.5	7.9	3.0
Xiaomi	4.5	32.5	9.55	3.5

only showing top 10 rows

Command took 1.07 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:20:26 AM on 21-2

databricks

Data processing using Pyspark Python ⭐

File Edit View Run Help Last edit was 46 days ago New cell UI: OFF

Cmd 20

```
1 df1.groupBy("mobile").sum().show(5,False)
```

▶ (2) Spark Jobs

mobile	sum(ratings)	sum(age)	sum(experience)	sum(family)
Infinix	6	65	18.5	4
Nokia	7	54	15.2	4
Sony	6	73	21.2	4
Alcatel	7	68	19.7	5
Motorola	7	70	19.8	5

only showing top 5 rows

Command took 0.99 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:20:31 AM on 21-2

databricks

Data processing using Pyspark Python ⭐

File Edit View Run Help Last edit was 46 days ago New cell UI: OFF

Cmd 21

```
1 df1.groupBy("mobile").min().show(5, False)
```

▶ (2) Spark Jobs

mobile	min(ratings)	min(age)	min(experience)	min(family)
Infinix	28	8.2	1	
Nokia	3	26	7.0	1
Sony	2	34	9.7	1
Alcatel	2	29	8.7	1
Motorola	3	33	9.8	2

only showing top 5 rows

Command took 1.12 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:20:34 AM on 21-2

databricks

Data processing using Pyspark Python ⭐

File Edit View Run Help Last edit was 46 days ago New cell UI: OFF

Cmd 22

```
1 df1.groupBy("mobile").max().show(5, False)
```

▶ (2) Spark Jobs

mobile	max(ratings)	max(age)	max(experience)	max(family)
Infinix	4	37	10.3	3
Nokia	4	28	8.2	3
Sony	4	39	11.5	3
Alcatel	5	39	11.0	4
Motorola	4	37	10.0	3

only showing top 5 rows

Command took 0.97 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:20:37 AM on 21-2

databricks

Data processing using Pyspark Python ⭐

File Edit View Run Help Last edit was 46 days ago New cell UI: OFF

Cmd 23

```
1 df.groupby("mobile").agg({"experience": "sum"}).show(5, False)
```

▶ (2) Spark Jobs

mobile	sum(experience)
Infinix	18.5
Nokia	15.2
Sony	21.2
Alcatel	19.7
Motorola	19.8

only showing top 5 rows

Command took 0.94 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:21:29 AM on 21-2

Cmd 24

```
1 from pyspark.sql.functions import udf
```

Command took 0.17 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:28:46 AM on 21-2

databricks

Data processing using Pyspark Python

File Edit View Run Help Last edit was 46 days ago New cell UI: OFF

Cmd 25

```
1 def price_range(brand):
2     if brand in ["Samsung", "iPhone"]:
3         return "High Price"
4     elif brand=="Xiaomi":
5         return "Mid price"
6     else:
7         return "Low price"
```

Command took 0.09 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:28:49 AM on 21-2

databricks

Data processing using Pyspark Python

File Edit View Run Help Last edit was 46 days ago New cell UI: OFF

Cmd 26

```
1 brand_udf=udf(price_range, StringType())
2 df.withColumn("price-range",brand_udf(df["mobile"])).show(10,False)
```

▶ (1) Spark Jobs

rating	age	experience	family	mobile	price-range
3	32	9.0	3	Vivo	Low price
4	28	8.5	2	Samsung	High Price
5	35	10.2	4	iPhone	High Price
2	40	12.0	1	OnePlus	Low price
4	27	7.8	3	Google Pixel	Low price
3	33	9.5	2	OPPO	Low price
5	29	8.0	4	Huawei	Low price
4	36	10.8	3	Xiaomi	Mid price
3	31	9.2	2	LG	Low price
4	34	9.7	3	Sony	Low price

only showing top 10 rows

Command took 1.77 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:28:55 AM on 21-2

Cmd 27

```
1 from pyspark.sql.functions import pandas_udf , PandasUDFType
```

Command took 0.06 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:36:51 AM on 21-2

databricks

Data processing using Pyspark Python

File Edit View Run Help Last edit was 46 days ago New cell UI: OFF

Cmd 28

```
1 def remaining_yrs(age):
2     yrs_left = 100-age
3     return yrs_left
```

Command took 0.12 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:36:54 AM on 21-2

Cmd 29

```
1 length_udf = pandas_udf(remaining_yrs, IntegerType())
```

Command took 1.79 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:36:57 AM on 21-2

databricks

Data processing using Pyspark Python ⭐

File Edit View Run Help Last edit was 46 days ago New cell UI: OFF

Run all Unknown Share Publish

Cmd_30

```
1 df1.withColumn("yrs_left", length_udf(df1['age'])).show(10,False)
```

▶ (1) Spark Jobs

	ratings	age	experience	family	mobile	yrs_left
3	32	9.0	3	Vivo	68	
4	28	8.5	2	Samsung	72	
5	35	10.2	4	iPhone	65	
2	40	12.0	1	OnePlus	60	
4	27	7.8	3	Google Pixel	73	
3	33	9.5	2	OPPO	67	
5	29	8.0	4	Huawei	71	
4	36	10.8	3	Xiaomi	64	
3	31	9.2	2	LG	69	
4	34	9.7	3	Sony	66	

only showing top 10 rows

Command took 2.51 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:37:31 AM on 21-2

databricks

Data processing using Pyspark Python ⭐

File Edit View Run Help Last edit was 46 days ago New cell UI: OFF

Run all Unknown Share Publish

Cmd_32

```
1 prod_udf=pandas_udf(prod[DoubleType]())
2 df1.withColumn("product",prod_udf(df1["ratings"] ,df1['experience'])).show(10,False)
```

▶ (1) Spark Jobs

	ratings	age	experience	family	mobile	product
3	32	9.0	3	Vivo	27.0	
4	28	8.5	2	Samsung	34.0	
5	35	10.2	4	iPhone	51.0	
2	40	12.0	1	OnePlus	24.0	
4	27	7.8	3	Google Pixel	31.2	
3	33	9.5	2	OPPO	28.5	
5	29	8.0	4	Huawei	40.0	
4	36	10.8	3	Xiaomi	43.2	
3	31	9.2	2	LG	27.599999999999998	
4	34	9.7	3	Sony	38.8	

only showing top 10 rows

Command took 1.59 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:49:24 AM on 21-2

Cmd_33

```
1 df1.dropDuplicates()
```

Out[78]: DataFrame[ratings: int, age: int, experience: double, family: int, mobile: string]

Command took 0.14 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:49:42 AM on 21-2

databricks

Data processing using Pyspark Python ⭐

File Edit View Run Help Last edit was 46 days ago New cell UI: OFF

Run all Unknown Share Publish

Cmd_34

```
1 df.count()
```

▶ (2) Spark Jobs

Out[79]: 4

Command took 0.48 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:49:46 AM on 21-2

Cmd_35

```
1 df_new = df1.drop('mobile')
```

Out[80]: pyspark.sql.dataframe.DataFrame = [ratings: integer, age: integer ... 2 more fields]

Command took 0.05 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:49:54 AM on 21-2

databricks

Data processing using Pyspark Python ⭐

File Edit View Run Help Last edit was 46 days ago New cell UI: OFF

Run all Unknown Share Publish

Cmd 36

```
1 df_new.show()
```

▶ (1) Spark Jobs

5	35	10.2	4
2	40	12.0	1
4	27	7.8	3
3	33	9.5	2
5	29	8.0	4
4	36	10.8	3
3	31	9.2	2
4	34	9.7	3
5	30	8.3	4
2	38	11.5	2
3	26	7.0	1
4	37	10.0	3
5	29	8.5	4
3	32	9.2	2
4	35	10.5	3
5	31	8.8	4
2	39	11.0	1
4	28	8.2	3

+-----+-----+-----+

only showing top 20 rows

Command took 0.47 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:50:15 AM on 21-2

databricks

Data processing using Pyspark Python ⭐

File Edit View Run Help Last edit was 46 days ago New cell UI: OFF

Run all Unknown Share Publish

Cmd 38

```
1 pwd
```

Out[82]: '/databricks/driver'

Command took 0.03 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:50:32 AM on 21-2

Cmd 38

```
1 write_uri = '/home/jovyan/work/df_csv'
```

Command took 0.10 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:50:38 AM on 21-2

Cmd 39

```
1 df1.coalesce(1).write.format("csv").option("header","true").save(write_uri)
```

▶ (1) Spark Jobs

Command took 2.36 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:50:57 AM on 21-2

Cmd 40

```
1 parquet_uri = '/home/jovyan/work/df_parquet'
```

Command took 0.04 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:51:11 AM on 21-2

Cmd 41

```
1 df1.write.format('parquet').save(parquet_uri)
```

▶ (1) Spark Jobs

Command took 3.78 seconds -- by dhvanipatel301003@gmail.com at 2/21/2024, 10:51:24 AM on 21-2

Name: Dhvani Patel

Rollno: 21BCP116

Big Data Analytics Lab

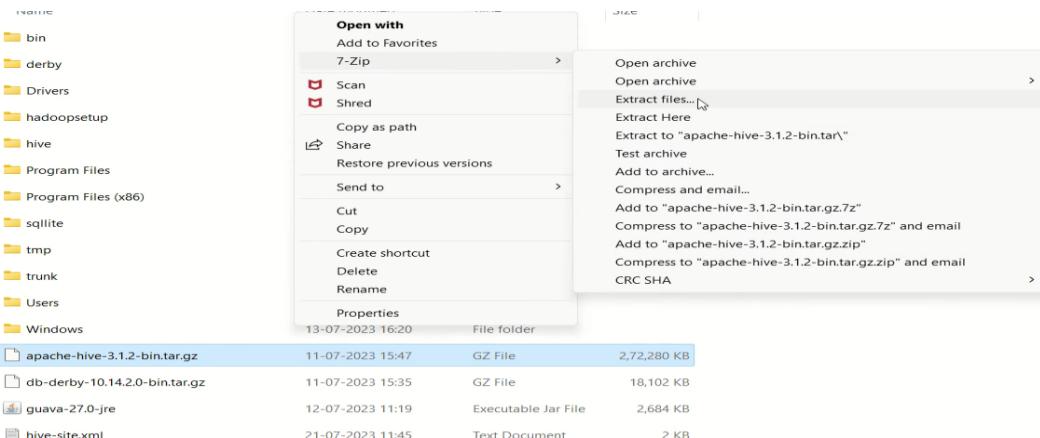
Practical 7 : Hadoop, Pig, Kafka, Spark and Hive installation

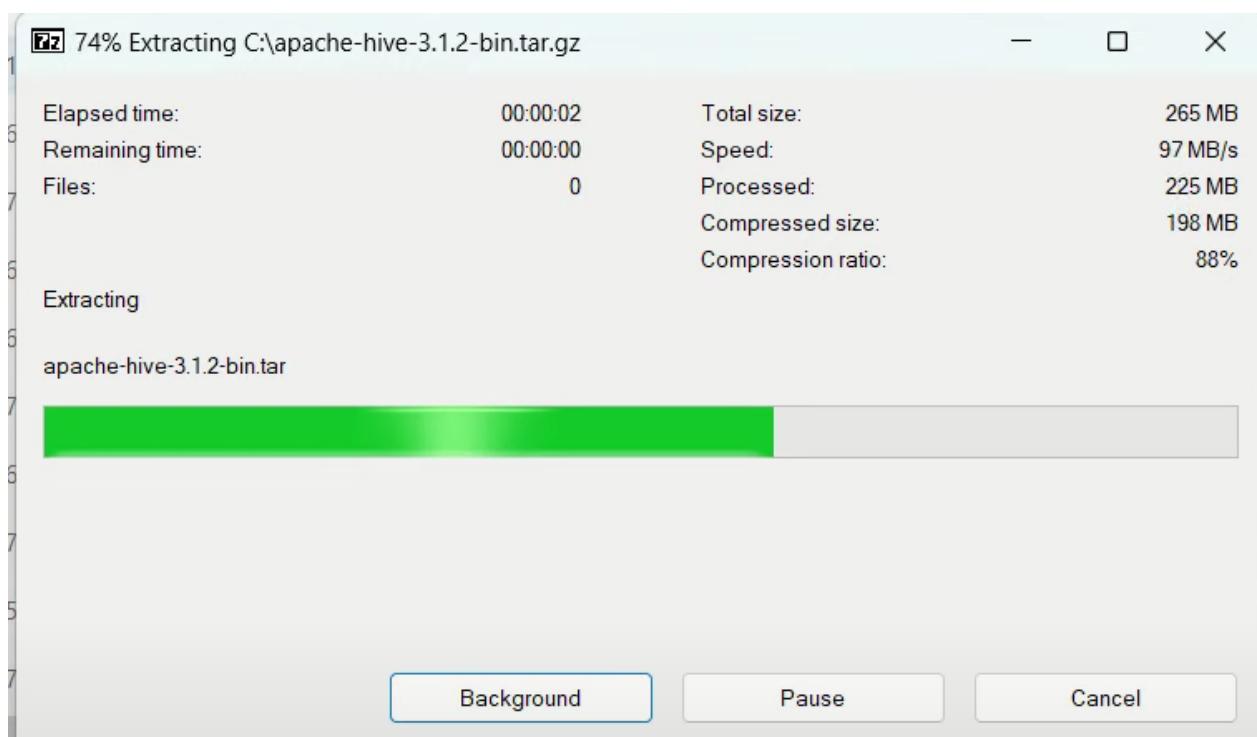
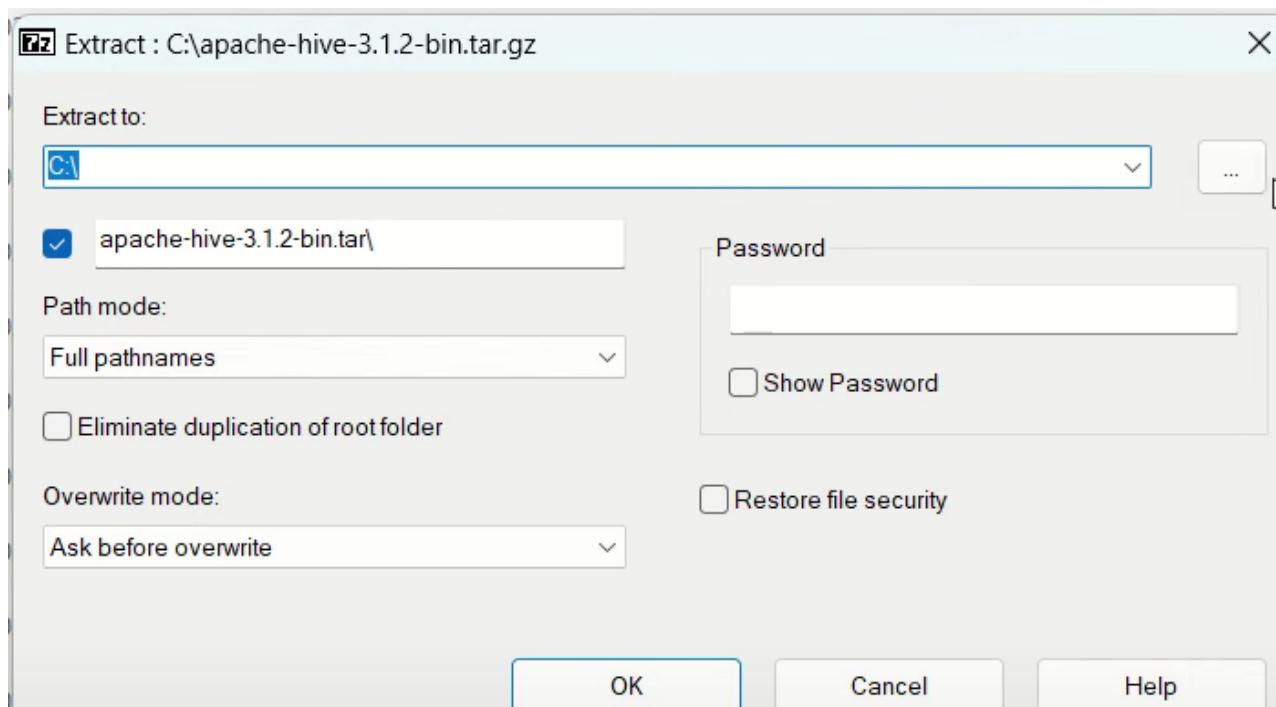
Apache Hive:

The screenshot shows the Apache Derby Downloads page. The top navigation bar includes Home, Quick Start, Download, Community, Documentation, and Resources. The Download menu is expanded, showing Overview, The Apache Software Foundation, and a search bar. The main content area is titled "Apache Derby: Downloads" and lists releases categorized by Java version: "For Java 17 and Higher", "For Java 9 and Higher", "For Java 8 and Higher", "For Java 6 and Higher", and "For Java 4 and Higher". Each category contains a list of specific release links.

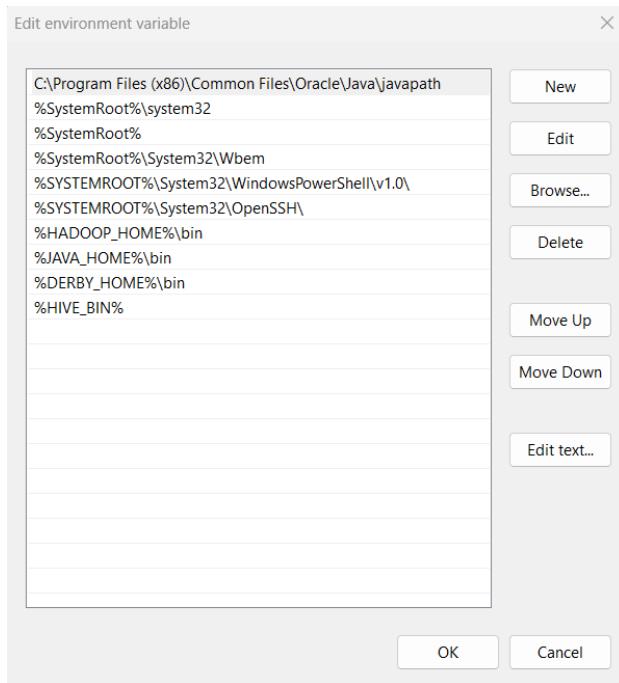
Index of /hive/hive-3.1.2

Name	Last modified	Size	Description
Parent Directory		-	
apache-hive-3.1.2-bin.tar.gz	2019-08-26 20:20	266M	
apache-hive-3.1.2-bin.tar.gz.asc	2019-08-26 20:20	833	
apache-hive-3.1.2-bin.tar.gz.sha256	2019-08-26 20:20	95	
apache-hive-3.1.2-src.tar.gz	2019-08-26 20:20	24M	
apache-hive-3.1.2-src.tar.gz.asc	2019-08-26 20:20	833	
apache-hive-3.1.2-src.tar.gz.sha256	2019-08-26 20:20	95	





HIVE_BIN	%HIVE_HOME%\bin
HIVE_HOME	C:\hive\apache-hive-3.1.2-bin\apache-hive-3.1.2-bin
HIVE_LIB	%HIVE_HOME%\lib
DERBY_HOME	C:\derby\db-derby-10.14.2.0-bin\db-derby-10.14.2.0-bin



Apache spark:

APACHE Spark

Download Libraries Documentation Examples Community Developers Apache Software Foundation ▾

Download Apache Spark™

- Choose a Spark release: 3.5.0 (Sep 13 2023) ▾
- Choose a package type: Pre-built for Apache Hadoop 3.3 and later
- Download Spark: spark-3.5.0-bin-hadoop3.tgz
- Verify this release using the 3.5.0 signatures, checksums and project release KEYS by following these procedures.

Note that Spark 3 is pre-built with Scala 2.12 in general and Spark 3.2+ provides additional pre-built distribution with Scala 2.13.

Link with Spark
Spark artifacts are hosted in Maven Central. You can add a Maven dependency with the following coordinates:

```
groupId: org.apache.spark
artifactId: spark-core_2.12
version: 3.5.0
```

Installing with PyPi
PySpark is now available in pypi. To install just run `pip install pyspark`.

Convenience Docker Container Images
Spark Docker Container Images are available from DockerHub; these images contain non-ASF software and may be subject to different license terms.

Latest News
Spark 3.3.4 released (Dec 16, 2023)
Spark 3.4.2 released (Nov 30, 2023)
Spark 3.5.0 released (Sep 13, 2023)
Spark 3.3.3 released (Aug 21, 2023)

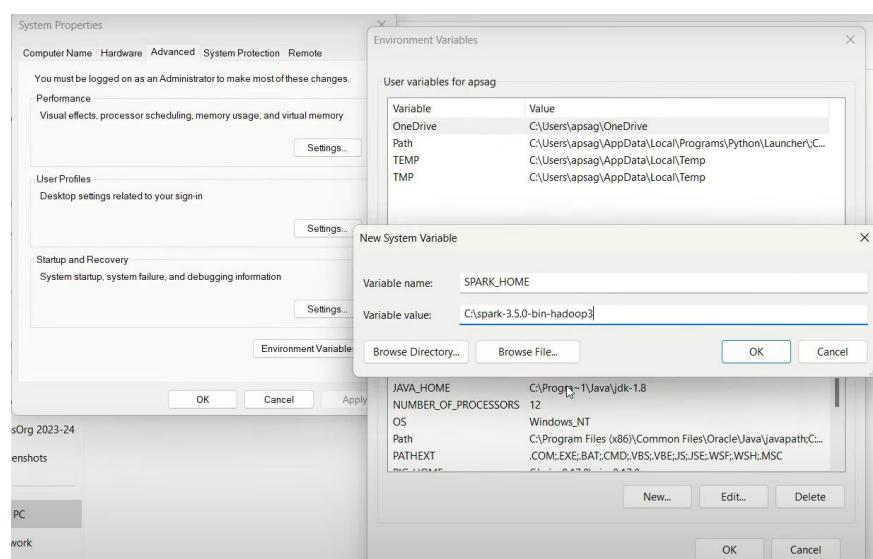
Archive

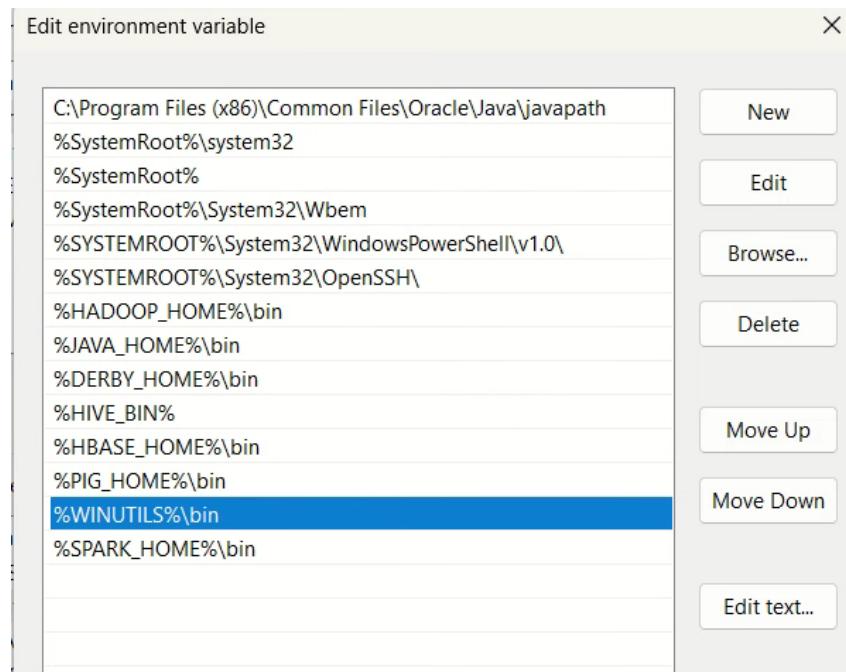
COMMUNITY CODE

DOWNLOAD SPARK

Built-in Libraries:
SQL and DataFrames
Spark Streaming
MLlib (machine learning)
GraphX (graph)

Third-Party Projects





```

Windows PowerShell
Copyright (C) Microsoft Corporation. All rights reserved.

Install the latest PowerShell for new features and improvements! https://aka.ms/PSWindows

PS C:\Windows\system32> cd C:\spark-3.5.0-bin-hadoop3\bin\
PS C:\spark-3.5.0-bin-hadoop3\bin> spark-shell
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
Spark context Web UI available at http://amrit:4040
Spark context available as 'sc' (master = local[*], app id = local-1703438479261).
Spark session available as 'spark'.
Welcome to

    _/\_ _/\_ _/\_ _/\_
   / \ / \ / \ / \ / \
  /   \ / \ / \ / \ / \
 /     \ / \ / \ / \ / \
/       \ / \ / \ / \ / \
 \         \ / \ / \ / \
  \           \ / \ / \
   \             \ / \
    \               \ /
     \                 \
      \                   \
       \                     \
        \                       \
         \                         \
          \                           \
           \                             \
            \                               \
             \                                 \
              \                                   \
               \                                     \
                \                                       \
                 \                                         \
                  \                                           \
                   \                                             \
                    \                                               \
                     \                                                 \
                      \                                                 \
                       \                                                 \
                        \                                                 \
                         \                                                 \
                          \                                                 \
                           \                                                 \
                            \                                                 \
                             \                                                 \
                              \                                                 \
                               \                                                 \
                                \                                                 \
                                 \                                                 \
                                  \                                                 \
                                   \                                                 \
                                    \                                                 \
                                     \                                                 \
                                      \                                                 \
                                       \                                                 \
                                        \                                                 \
                                         \                                                 \
                                          \                                                 \
                                           \                                                 \
                                            \                                                 \
                                             \                                                 \
                                              \                                                 \
                                               \                                                 \
                                                \                                                 \
                                                 \                                                 \
                                                  \                                                 \
                                                   \                                                 \
                                                    \                                                 \
                                                     \                                                 \
                                                      \                                                 \
                                                       \                                                 \
                                                        \                                                 \
                                                         \                                                 \
                                                          \                                                 \
                                                           \                                                 \
                                                            \                                                 \
                                                             \                                                 \
                                                              \                                                 \
                                                               \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
                                                                 \                                                 \
                                                                \                                                 \
................................................................

```

Hadoop installation

Java 8 installation:

Solaris x64 (SRV4 package)	124.37 MB	jdk-8u202-solaris-x64.tar.Z
Solaris x64	85.38 MB	jdk-8u202-solaris-x64.tar.gz
Windows x86	201.64 MB	jdk-8u202-windows-i586.exe
Windows x64	211.58 MB	jdk-8u202-windows-x64.exe



Apache Hadoop Download Documentation Community Development Help Apache Software Foundation

Download

Hadoop is released as source code tarballs with corresponding binary tarballs for convenience. The downloads are distributed via mirror sites and should be checked for tampering using GPG or SHA-512.

Version	Release date	Source download	Binary download	Release notes
3.4.0	2024 Mar 17	source (checksum signature)	binary (checksum signature) binary-aarch64 (checksum signature)	Announcement
3.3.6	2023 Jun 23	source (checksum signature)	binary (checksum signature) binary-aarch64 (checksum signature)	Announcement
2.10.2	2022 May 31	source (checksum signature)	binary (checksum signature)	Announcement

Sponsor the ASF

We suggest the following location for your download:
<https://dlcdn.apache.org/hadoop/common/hadoop-3.4.0/hadoop-3.4.0-src.tar.gz>

Alternate download locations are suggested below.

It is essential that you [verify the integrity](#) of the downloaded file using the PGP signature ([.asc](#) file) or a hash ([.ad5](#) or [.sha*](#) file).

HTTP

<https://dlcdn.apache.org/hadoop/common/hadoop-3.4.0/hadoop-3.4.0-src.tar.gz>

BACKUP SITES

<https://dlcdn.apache.org/hadoop/common/hadoop-3.4.0/hadoop-3.4.0-src.tar.gz>

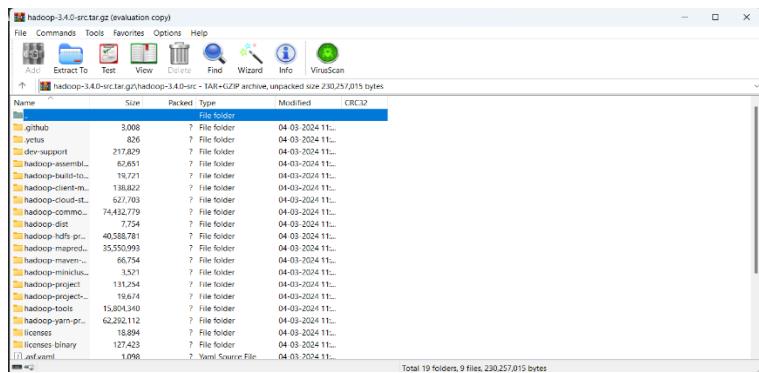
Downloads

Today

[hadoop-3.4.0-src.tar.gz](https://dlcdn.apache.org/hadoop/common/hadoop-3.4.0/hadoop-3.4.0-src.tar.gz)
https://dlcdn.apache.org

This PC > New Volume (D:) > hadoop-3.4.0-src >

Name	Date modified	Type	Size
hadoop-common-project	08-04-2024 01:22 PM	File folder	
hadoop-yarn-project	08-04-2024 01:22 PM	File folder	
licenses	08-04-2024 01:21 PM	File folder	
NOTICE	04-03-2024 11:14 AM	Text Document	2 KB
README	04-03-2024 11:14 AM	Text Document	1 KB



Pig Installation :

Apache > Pig >

APACHE hadoop

Project **Wiki**

Apache Pig Releases

Download News

- 19 June, 2017: release 0.17.0 available
- 8 June, 2016: release 0.16.0 available
- 6 June, 2015: release 0.15.0 available
- 20 November, 2014: release 0.14.0 available
- 14 March, 2014: release 0.13.0 available
- 14 April, 2014: release 0.12.1 available
- 14 October, 2013: release 0.12.0 available
- 1 April, 2013: release 0.11.1 available
- 21 February, 2013: release 0.11.0 available
- 03 January, 2013: release 0.10.1 available
- 25 August, 2012: release 0.10.0 available
- 22 January, 2012: release 0.9.2 available

Search the site with google Search Last Published: 02/23/2021 03:15:42 PDF Incognito

apache.org/dyn/closer.cgi/pig

Sponsor the ASF

Community Projects Downloads Learn Resources & Tools About

**THE APACHE®
SOFTWARE FOUNDATION
ESTABLISHED 1999**

We suggest the following location for your download:

<https://dlcdn.apache.org/pig>

Alternate download locations are suggested below.

It is essential that you [verify the integrity](#) of the downloaded file using the PGP signature ([.asc](#) file) or a hash ([.md5](#) or [.sha256](#) file).

HTTP

<https://dlcdn.apache.org/pig>

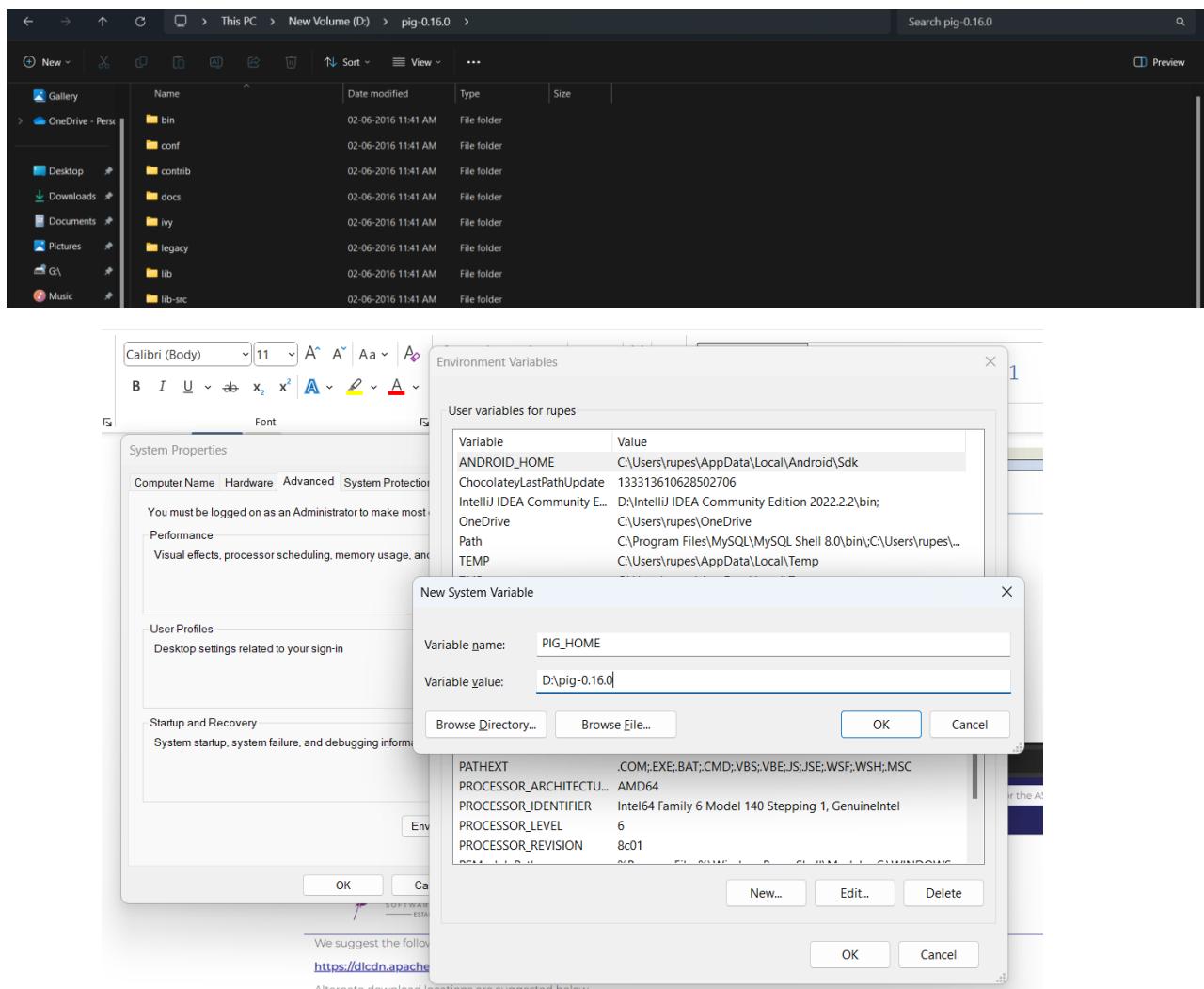
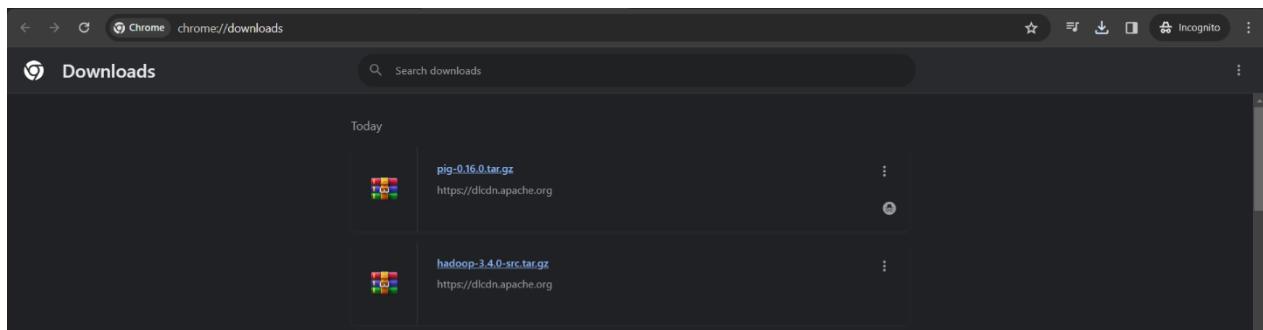
BACKUP SITES

<https://dlcdn.apache.org/pig>

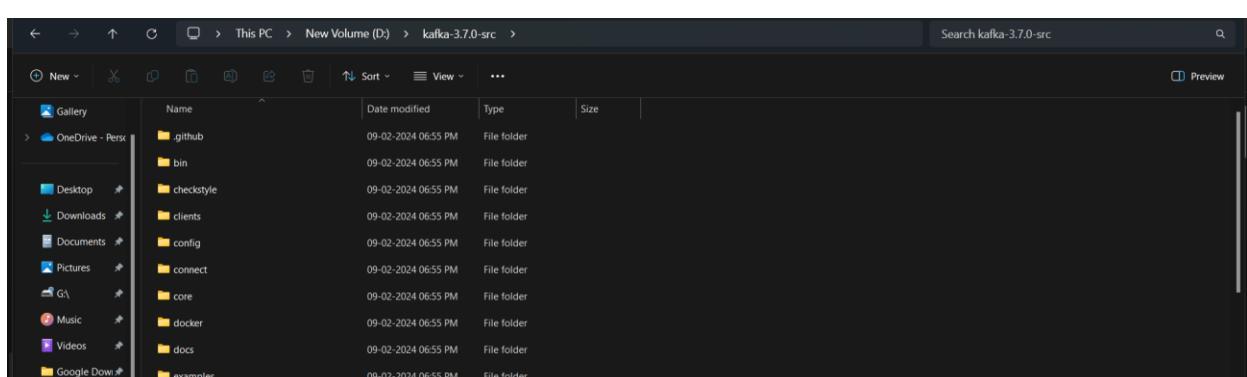
dcldn.apache.org/pig/pig-0.16.0/

Index of /pig/pig-0.16.0

Name	Last modified	Size	Description
Parent Directory		-	
README.txt	2016-06-07 17:08	1.4K	
pig-0.16.0-src.tar.gz	2016-06-07 17:08	14M	
pig-0.16.0-src.tar.gz.asc	2016-06-07 17:08	195	
pig-0.16.0-src.tar.gz.md5	2016-06-07 17:08	56	
pig-0.16.0.tar.gz	2016-06-07 17:08	169M	
pig-0.16.0.tar.gz.asc	2016-06-07 17:08	195	
pig-0.16.0.tar.gz.md5	2016-06-07 17:08	52	



Kafka installation:



kafka.apache.org/downloads

kafka

GET STARTED DOCS POWERED BY COMMUNITY APACHE DOWNLOAD KAFKA

DOWNLOAD

3.7.0 is the latest release. The current stable version is 3.7.0

You can verify your download by following these [procedures](#) and using these [KEYS](#).

3.7.0

- Released Feb 27, 2024
- [Release Notes](#)
- Docker image: [apache/kafka:3.7.0](#).
- Source download: [kafka-3.7.0-src.tgz](#) (asc, sha512)
- Binary downloads:
 - Scala 2.12 - [kafka_2.12-3.7.0.tgz](#) (asc, sha512)
 - Scala 2.13 - [kafka_2.13-3.7.0.tgz](#) (asc, sha512)

We build for multiple versions of Scala. This only matters if you are using Scala and you want a version built for the same Scala version you use. Otherwise any version should work (2.13 is recommended).

Kafka 3.7.0 includes a significant number of new features and fixes. For more information, please read our [blog post](#) and the detailed [Release Notes](#).

3.6.2

D:\kafka-3.7.0-src\config\zookeeper.properties - Notepad++

```

File Edit Search View Encoding Language Settings Tools Macro Run Plugins Window ?
zookeeper.properties
1   # Licensed to the Apache Software Foundation (ASF) under one or more
2   # contributor license agreements. See the NOTICE file distributed with
3   # this work for additional information regarding copyright ownership.
4   # The ASF licenses this file to You under the Apache License, Version 2.0
5   # (the "License"); you may not use this file except in compliance with
6   # the License. You may obtain a copy of the License at
7   #
8   #     http://www.apache.org/licenses/LICENSE-2.0
9   #
10  # Unless required by applicable law or agreed to in writing, software
11  # distributed under the License is distributed on an "AS IS" BASIS,
12  # WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
13  # See the License for the specific language governing permissions and
14  # limitations under the License.
15  # the directory where the snapshot is stored.
16  dataDir=D:/kafka-3.7.0-src/zookeeper
17  # the port at which the clients will connect
18  clientPort=2181
19  # disable the per-ip limit on the number of connections since this is a non-production config
20  maxClientCnxns=0
21  # Disable the adminserver by default to avoid port conflicts.
22  # Set the port to something non-conflicting if choosing to enable this
23  admin.enableServer=false
24  # admin.serverPort=8080
25

```

D:\kafka-3.7.0-src\config\server.properties - Notepad++

```

File Edit Search View Encoding Language Settings Tools Macro Run Plugins Window ?
zookeeper.properties server.properties
37  # If not set, it uses the value for "listeners".
38  #advertised.listeners=PLAINTEXT://your.host.name:9092
39
40  # Maps listener names to security protocols, the default is for them to be the same. See the config documentation for more details
41  #listener.security.protocol.map=PLAINTEXT:PLAINTEXT,SSL:SSL,SASL_PLAINTEXT:SASL_PLAINTEXT,SASL_SSL:SASL_SSL
42
43  # The number of threads that the server uses for receiving requests from the network and sending responses to the network
44  num.network.threads=3
45
46  # The number of threads that the server uses for processing requests, which may include disk I/O
47  num.io.threads=8
48
49  # The send buffer (SO_SNDBUF) used by the socket server
50  socket.send.buffer.bytes=102400
51
52  # The receive buffer (SO_RCVBUF) used by the socket server
53  socket.receive.buffer.bytes=102400
54
55  # The maximum size of a request that the socket server will accept (protection against OOM)
56  socket.request.max.bytes=104857600
57
58  ##### Log Basics #####
59
60  # A comma separated list of directories under which to store log files
61  log.dirs=D:/kafka-3.7.0-src/kafka-logs
62
63  # The default number of log partitions per topic. More partitions allow greater
64  # parallelism for consumption, but this will also result in more files across
65  # the brokers.
66  num.partitions=1
67
68

```

Name: Dhvani Patel

Roll no: 21BCP116

Big Data Analytics Lab

Practical 8: Implementation of graph data structure using networkx

jupyter bdalab_27_3_24 Last Checkpoint: 5 days ago

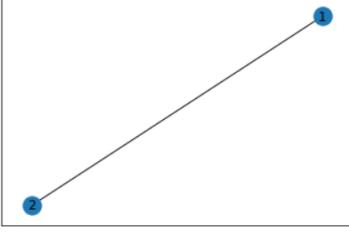
File Edit View Run Kernel Settings Help Trusted JupyterLab Python 3 (ipykernel)

```
[1]: import networkx as nx
import matplotlib.pyplot as plt

[2]: G = nx.Graph()

G.add_edge(1,2)

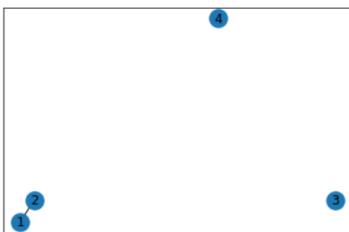
nx.draw_networkx(G)
plt.show()
```



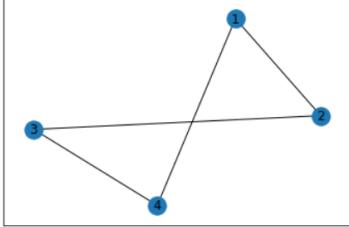
jupyter bdalab_27_3_24 Last Checkpoint: 5 days ago

File Edit View Run Kernel Settings Help Trusted JupyterLab Python 3 (ipykernel)

```
[3]: G.add_nodes_from([3,4])
nx.draw_networkx(G)
plt.show()
```



```
[4]: G.add_edge(3,4)
G.add_edges_from([(2,3),(4,1)])
nx.draw_networkx(G)
plt.show()
```



jupyter bdalab_27_3_24 Last Checkpoint: 5 days ago

File Edit View Run Kernel Settings Help Trusted

[5]: G.nodes

[5]: NodeView((1, 2, 3, 4))

[6]: G.edges

[6]: EdgeView([(1, 2), (1, 4), (2, 3), (3, 4)])

[7]: list(nx.generate_adjlist(G))

[7]: ['1 2 4', '2 3', '3 4', '4']

[8]: nx.to_dict_of_lists(G)

[8]: {1: [2, 4], 2: [1, 3], 3: [4, 2], 4: [3, 1]}

[9]: nx.to_edgelist(G)

[9]: EdgeDataView([(1, 2, {}), (1, 4, {}), (2, 3, {}), (3, 4, {})])

jupyter bdalab_27_3_24 Last Checkpoint: 5 days ago

File Edit View Run Kernel Settings Help Trusted

[10]: G.add_edge(1,3)
nx.draw_networkx(G)
plt.show()

[11]: G.degree

[11]: DegreeView({1: 3, 2: 2, 3: 3, 4: 2})

[12]: kg=nx.fast_gnp_random_graph(10000,0.01)
k = kg.degree()
plt.hist(list(dict(k).values()))
plt.show()

jupyter bdalab_27_3_24 Last Checkpoint: 5 days ago

File Edit View Run Kernel Settings Help Trusted

[12]: kg=nx.fast_gnp_random_graph(10000,0.01)
k = kg.degree()
plt.hist(list(dict(k).values()))
plt.show()

[13]: nx.draw_networkx(kg)
plt.show()

[13]: <function matplotlib.pyplot.show(close=None, block=None)>

jupyter bdalab_27_3_24 Last Checkpoint: 5 days ago

File Edit View Run Kernel Settings Help Trusted JupyterLab Python 3 (ipykernel)

```
[14]: G=nx.krackhardt_kite_graph()
nx.draw_networkx(G)
plt.show()
```

```
[15]: print(nx.has_path(G,source=1, target=9))
print(nx.shortest_path(G,source=1, target=9))
print(nx.shortest_path_length(G,source=1, target=9))
paths=list(nx.all_pairs_shortest_path(G))
paths[5][1]
```

```
True
[1, 6, 7, 8, 9]
4
```

```
[15]: {5: [5],
 0: [5, 0],
 2: [5, 2],
 3: [5, 3],
 6: [5, 6],
 7: [5, 7],
 1: [5, 0, 1],
 4: [5, 3, 4],
 8: [5, 7, 8],
 9: [5, 7, 8, 9]}
```

jupyter bdalab_27_3_24 Last Checkpoint: 5 days ago

File Edit View Run Kernel Settings Help Trusted JupyterLab Python 3 (ipykernel)

```
[16]: nx.betweenness_centrality(G)
```

```
{0: 0.023148148148148143,
 1: 0.023148148148148143,
 2: 0.0,
 3: 0.10185185185185183,
 4: 0.0,
 5: 0.23148148148148148,
 6: 0.23148148148148148,
 7: 0.38888888888888884,
 8: 0.2222222222222222,
 9: 0.0}
```

```
[17]: nx.degree_centrality(G)
```

```
{0: 0.4444444444444444,
 1: 0.4444444444444444,
 2: 0.3333333333333333,
 3: 0.6666666666666666,
 4: 0.3333333333333333,
 5: 0.5555555555555555,
 6: 0.5555555555555555,
 7: 0.3333333333333333,
 8: 0.2222222222222222,
 9: 0.1111111111111111}
```

```
[18]: nx.closeness_centrality(G)
```

```
{0: 0.5294117647058824,
 1: 0.5294117647058824,
 2: 0.5,
 3: 0.6,
 4: 0.5,
 5: 0.6428571428571429,
 6: 0.6428571428571429,
 7: 0.6,
 8: 0.42857142857142855,
 9: 0.3103448275862069}
```

jupyter bdalab_27_3_24 Last Checkpoint: 5 days ago

File Edit View Run Kernel Settings Help Trusted

[19]: nx.harmonic_centrality(G)

```
[19]: {0: 6.083333333333333, 1: 6.083333333333333, 2: 5.583333333333333, 3: 7.083333333333333, 4: 5.583333333333333, 5: 6.833333333333333, 6: 6.833333333333333, 7: 6.0, 8: 4.666666666666666, 9: 3.416666666666665}
```

[20]: nx.eigenvector_centrality(G)

```
[20]: {0: 0.35220898139203594, 1: 0.35220898139203594, 2: 0.28583473531632414, 3: 0.48102048812210046, 4: 0.28583473531632414, 5: 0.3976910106255469, 6: 0.3976910106255469, 7: 0.19586185175360382, 8: 0.048074775014202945, 9: 0.01116405857582424}
```

[21]: nx.clustering(G)

```
[21]: {0: 0.6666666666666666, 1: 0.6666666666666666, 2: 1.0, 3: 0.5333333333333333, 4: 1.0, 5: 0.5, 6: 0.5, 7: 0.3333333333333333, 8: 0, 9: 0}
```

jupyter bdalab_27_3_24 Last Checkpoint: 5 days ago

File Edit View Run Kernel Settings Help Trusted

[22]: G=nx.cubical_graph(G)

```
[22]: nx.draw(G, pos=nx.circular_layout(G),node_color='r',edge_color='b')
```

The image shows a 3D plot of a cubical graph. It consists of 8 red nodes arranged in a cube-like structure. The edges connecting the nodes are drawn in blue. The plot is set against a white background with a light gray grid.