05/02/24. i) Py Spark: - It is apache spark library written in Pytton to run Pytton applications cur of other spark confabilities + curring. Pysfall we can run app 1101 on distributed Clayle 2) was Apriche spouk! It is an open-source unified analytims largine word too dange Scale date processing. Spork is designed to be fast Hervide & Easy to use, making it a popular Choice for processing large scale dute sets. -> It operations on billions & tollions of det on distributed clusters too these pasts than traditional applications. -) con von on single/multi- node cluster. -) was realist to address the limitations of Mapreduce by doing in- memory prolessing -) Tt reuse date by using in meny CHON 18 3) was Ryspack - Python lenchdig Num Py, tensorphon. Malmant, Trivago, canodi Etc. -) Spydu IDE, 6 Jupyter Note book to run Pysfaux app. 1 Amalonda-also1) Pyspall features h MOY, PERMI In memory lompilation I distributed Pow cerring ushed Parallelise. Jan be used with many duster managey Espark, Yain, Mesoc Etc.)) fault-tolerance) Inhild offmitables whenevery date Haves. Emmilable. SUPPONTS ANSISAL _) Lazy Evaluation o _) Carrie & Persistance. s) Adv of Bysparc: =) It is general Purpose, in-memory, distributed Processing Engole allows you to Process data Estimently in distributed cluster. -) AFF running on Spark are 100 x faster than usual; -) will get except benitets from Pyspaise for dela Newsons Pipelines - El lan process deta from MADOUP MDIS, Aws sz, low Elk. also used to process real-thre data wong Stoloring & Katka I has men ML & graph to bratices

-) Roses Java Versions 8, 11813. -) Sala Version 2.12 & 212. beyond. -) Sala Version 2.12 & 212. beyond. -) Worles in master stare architecture. -) worles in master stare architecture. -) master-drover slaver works -) Pyspark modules & Packagust. -) Pyspark RDD (Pyspark. 200) -> Data frame & Sali (Pyspark. see 1) -> Messource (Pyspark. streams) -) Messource (Pyspark. resource) -> Resource (Pyspark. resource) -> Resource (Pyspark. resource) -> Packagus. org are & Packagus.	n) is a Version of Python Par
-) Pyspalc & S. is compatible with Pyth. 3.8 & new a -) R. 35, Jana Versions 8, 11213. -) Scala version 2.12 & 2.13. beyond. 3) Pyspalc architecture for works in master stark architecture. -) world, in master stark architecture. -) master-drover slaves works 8) Ob Pyspark modules & Porchagus. -) Pyspark RDD C Pyspark. & DD -> Data frame & SOL (Pyspark Septembry) -> Maisbil n. mil., Pyspark. streamy) -> Misbil n. mil., Pyspark. mill) -> Misbil n. mil., Pyspark. mill) -> Presponence (Pyspark. resource) -> Pospark - Prekagus. org. are 3 Party liberty	6) Star VVSPays
-) Roses Java Versons 8, 11813. -) Sala Verson 2.12 & 2.12. beyond. -) Sala Verson 2.12 & 2.12. beyond. -) Worles in master stare architection. -) master-drover slaver works e) Or Pyspark modules & Packaget -) Pyspark RDD (Pyspark. 200) Data frame 2 Soli (Pyspark. sell) MLib (n. mi, Pysparkimili) MLib (n. mi, Pysparkimili) Mongh frams (Graph Frames) Resource (Pyspark. resource) -) Resource (Pyspark. resource)	Lafter 1
-) Roses Java Versons 8, 11813. -) Sala Verson 2.12 & 2.12. beyond. -) Sala Verson 2.12 & 2.12. beyond. -) Worles in master stare architection. -) master-drover slaver works e) Or Pyspark modules & Packaget -) Pyspark RDD (Pyspark. 200) Data frame 2 Soli (Pyspark. sell) MLib (n. mi, Pysparkimili) MLib (n. mi, Pysparkimili) Mongh frams (Graph Frames) Resource (Pyspark. resource) -) Resource (Pyspark. resource)	
-) Roses Java Versons 8, 11813. -) Sala Verson 2.12 & 2.12. beyond. -) Sala Verson 2.12 & 2.12. beyond. -) Worles in master stare architection. -) master-drover slaver works e) Or Pyspark modules & Packaget -) Pyspark RDD (Pyspark. 200) Data frame 2 Soli (Pyspark. sell) MLib (n. mi, Pysparkimili) MLib (n. mi, Pysparkimili) Mongh frams (Graph Frames) Resource (Pyspark. resource) -) Resource (Pyspark. resource)	-) Py spark & S -15 companied with Python
-) Sala Veusion 2.12 & 2.13. beyond. D) Pyspall carelitecture - -) works in master stark architecture.) master-drover stares works e) Ob Pyspark modules & Packagush -) Pyspark RDD (Pyspark. 2000) Data frame & Sali (Pyspark. 521) Mesto 1 mil , Pyspark streams) Mesto 1 mil , Pyspark mill) Mestourice (Pyspark resolute) Resourice (Pyspark resolute) Brack - Packagus org are 3 pachy liberty	3.8 & new a
-) Sala Veusion 2.12 & 2.13. beyond. D) Pyspall carelitecture - -) works in master stark architecture.) master-drover stares works e) Ob Pyspark modules & Packagush -) Pyspark RDD (Pyspark. 2000) Data frame & Sal (Pyspark. 521) Messon Trans (Pyspark. 5treams) Messon Trans (Graph Frams) Resource (Pyspark. resource) Resource (Pyspark. resource) Spark - Packagus org are 3 party liberty	-) R=35, Jana Vensons 8, 11215,
-) works in master stark architecture) master-driver starks works 2) Or Pyspark modules & Packagust -) Pyspark RDD (Pyspark. 2000) Data trame 2 Sol (Pyspark. 521) 1) Straing (Pyspark. streams) Miss (Graph Frames) Miss (Pyspark. resource) Presource (Pyspark. resource) Resource (Pyspark. resource)	-) Sah verston 2-12 & 2-13. berjond.
2) O Pyspark modulus & Packagust -) Pyspark RDD (Pyspark. 2000) Data frame 2 SOL (Pyspark seg 1) MLSD (MI , Ryspark. mth) MSD (Pyspark streamy) MSD (Pyspark streamy) MSD (Pyspark streamy) MSD (Pyspark streamy) Presource (Pyspark strong) Spark - Packagus org are 3 Party libration	
2) O Pyspark modulus & Packagust -) Pyspark RDD (Pyspark. 2000) Data frame 2 SOL (Pyspark seg 1) MLSD (MI , Ryspark. mth) MSD (Pyspark streamy) MSD (Pyspark streamy) MSD (Pyspark streamy) MSD (Pyspark streamy) Presource (Pyspark strong) Spark - Packagus org are 3 Party libration	-) world in martin stark and The
Pyspaik modules & Packagust -) Pyspaik RDD (Pyspaik RDD) Data frame & Sall (Pyspaik S21) Mish (Pyspaik Streamy) Mish (Min Min Ryspaik Mill) Creaps frams (Graph Frams) Resource (Pyspaik resource) Packagus org are & Packagus org are & Packagus org are & Packagus	- drover claves works
Data frame 2 SQL (PySPark SQ1) Smaning (PySpark Streamy) MLibl M. M. Pyspark mill) Creaph frams (Graph Frames) Resource (PySpark resource) Brack - Packagus org are 3 Pachy liberty	2) @ Pyspark modules & Packagust
MLSS (Pyspark streamy) MLSS (MI , MI , Pyspark mill) Grouph france (Graph Frances) Resource (Pyspark resource) Spark - Packagus org are 3 Party libally	
MLIST M. MI, Ryspark. mlib) Greagh francs (Graph Francis) Resource (Pripark risouris) Spark - Packagus org are 3 Party libration	
Park - Parkagu org are 3 Park libalis	mana TR.
Park - Parkagu org are 3 Park libalis	1) Japan K. Stolanny)
Spark-Packagus org are 3 Party liberty	MILSD (M , M), BysPay Kimilib)
Spark-Packagus org are 3 Party liberty	(Graph Francis (Graph Francis)
85:079 are 3 Party liberty	Resource (Pyspar)c - resource
85:079 are 3 Party liberty	
	1 le Park
	Brack-Packagus org on 2 Pach libalige
and the second of the second o	35:079 are 3 Party liberty

1) PYSPORK RDD - Resilient Distributed. data let: Create RDDF ways. distributes a set of Collection looding an Enternal dataset of objects. L'vit We have to import of i) by wong parallelize () tunction. from Pyspark - Sql . Import Sparksersion Space = spacession \ · builder) · approance (" Byther spack Create & DD Enerph") · Configl Spak · Some · Con fig - option, "some-value") df = SPark . Spark Content . Parallelize (([[1213, 'abc')] (4, (16 def'), (7,8,9, 'ghf'))), to DF (['(011', '(012', '(013', '(014')) df. show() = Prils ROD data

Syntan;

from Pyspark Session.

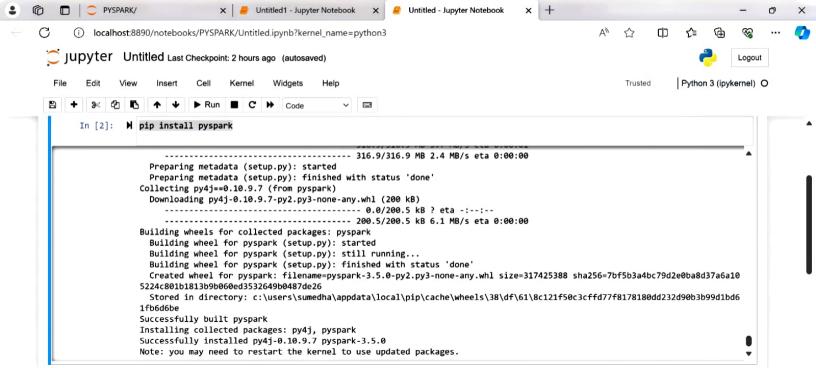
Forcate Sparksession.

Spark = Sparksession.

Spark = Spark session.

master ("local (1)")

app Mame



```
In [3]: | import pyspark
            from pyspark.sql import SparkSession
            # Create SparkSession
            spark = SparkSession.builder.appName("SparkByExamples.com") .getOrCreate()
            dataList = [("Java", 20000), ("Python", 100000), ("Scala", 3000)]
            rdd=spark.sparkContext.parallelize(dataList)
            result = rdd.collect()
            ("RDD Contents:", result)
   Out[3]: ('RDD Contents:', [('Java', 20000), ('Python', 100000), ('Scala', 3000)])
```

