

```
!apt-get install openjdk-8-jdk-headless -qq > /dev/null
```

```
!wget -q https://dlcdn.apache.org/spark/spark-3.2.1/spark-3.2.1-bin-hadoop3.2.tgz
```

```
!tar xf spark-3.2.1-bin-hadoop3.2.tgz &>/dev/null
```

```
!pip install -q findspark
```

```
import os
os.environ["JAVA_HOME"] = "/usr/lib/jvm/java-8-openjdk-amd64"
os.environ["SPARK_HOME"] = "/content/spark-3.2.1-bin-hadoop3.2"
```

```
!pip install pyspark
```

```
Collecting pyspark
  Downloading pyspark-3.2.1.tar.gz (281.4 MB)
    |████████████████████████████████████████| 281.4 MB 31 kB/s
Collecting py4j==0.10.9.3
  Downloading py4j-0.10.9.3-py2.py3-none-any.whl (198 kB)
    |████████████████████████████████████████| 198 kB 57.2 MB/s
Building wheels for collected packages: pyspark
  Building wheel for pyspark (setup.py) ... done
  Created wheel for pyspark: filename=pyspark-3.2.1-py2.py3-none-any.whl size=281853642
  Stored in directory: /root/.cache/pip/wheels/9f/f5/07/7cd8017084dce4e93e84e92efd1e1d5
Successfully built pyspark
Installing collected packages: py4j, pyspark
Successfully installed py4j-0.10.9.3 pyspark-3.2.1
```

```
!pip install pyspark[sql]
```

```
Requirement already satisfied: pyspark[sql] in /usr/local/lib/python3.7/dist-packages (
Requirement already satisfied: py4j==0.10.9.3 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: pyarrow>=1.0.0 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: pandas>=0.23.2 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: numpy>=1.17.3 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: pytz>=2017.3 in /usr/local/lib/python3.7/dist-packages (
Requirement already satisfied: python-dateutil>=2.7.3 in /usr/local/lib/python3.7/dist-
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.7/dist-packages (from
```

```
from pyspark.sql import SparkSession
```

```
spark = SparkSession.builder.getOrCreate()
```

```

from datetime import datetime, date
import pandas as pd
from pyspark.sql import Row

df = spark.createDataFrame([
    Row(a=1, b=2., c='string1', d=date(2000, 1, 1), e=datetime(2000, 1, 1, 12, 0)),
    Row(a=2, b=3., c='string2', d=date(2000, 2, 1), e=datetime(2000, 1, 2, 12, 0)),
    Row(a=4, b=5., c='string3', d=date(2000, 3, 1), e=datetime(2000, 1, 3, 12, 0))
])
df

```

```
DataFrame[a: bigint, b: double, c: string, d: date, e: timestamp]
```

```

df = spark.createDataFrame([
    (1, 2., 'string1', date(2000, 1, 1), datetime(2000, 1, 1, 12, 0)),
    (2, 3., 'string2', date(2000, 2, 1), datetime(2000, 1, 2, 12, 0)),
    (3, 4., 'string3', date(2000, 3, 1), datetime(2000, 1, 3, 12, 0))
], schema='a long, b double, c string, d date, e timestamp')
df

```

```
DataFrame[a: bigint, b: double, c: string, d: date, e: timestamp]
```

```

pandas_df = pd.DataFrame({
    'a': [1, 2, 3],
    'b': [2., 3., 4.],
    'c': ['string1', 'string2', 'string3'],
    'd': [date(2000, 1, 1), date(2000, 2, 1), date(2000, 3, 1)],
    'e': [datetime(2000, 1, 1, 12, 0), datetime(2000, 1, 2, 12, 0), datetime(2000, 1, 3, 12, 0)]
})
df = spark.createDataFrame(pandas_df)
df

```

```
DataFrame[a: bigint, b: double, c: string, d: date, e: timestamp]
```

```

rdd = spark.sparkContext.parallelize([
    (1, 2., 'string1', date(2000, 1, 1), datetime(2000, 1, 1, 12, 0)),
    (2, 3., 'string2', date(2000, 2, 1), datetime(2000, 1, 2, 12, 0)),
    (3, 4., 'string3', date(2000, 3, 1), datetime(2000, 1, 3, 12, 0))
])
df = spark.createDataFrame(rdd, schema=['a', 'b', 'c', 'd', 'e'])
df

```

```
DataFrame[a: bigint, b: double, c: string, d: date, e: timestamp]
```

```

df.show()
df.printSchema()

```

```
+---+---+-----+-----+-----+-----+
```

```

|  a|  b|      c|      d|      e|
+---+---+-----+-----+-----+
|  1|2.0|string1|2000-01-01|2000-01-01 12:00:00|
|  2|3.0|string2|2000-02-01|2000-01-02 12:00:00|
|  3|4.0|string3|2000-03-01|2000-01-03 12:00:00|
+---+---+-----+-----+-----+

```

```
root
```

```

|-- a: long (nullable = true)
|-- b: double (nullable = true)
|-- c: string (nullable = true)
|-- d: date (nullable = true)
|-- e: timestamp (nullable = true)

```

```
df.show(1)
```

```

+---+---+-----+-----+-----+
|  a|  b|      c|      d|      e|
+---+---+-----+-----+-----+
|  1|2.0|string1|2000-01-01|2000-01-01 12:00:00|
+---+---+-----+-----+-----+
only showing top 1 row

```

```
spark.conf.set('spark.sql.repl.eagerEval.enabled', True)
df
```

```
df.show(1, vertical=True)
```

```

-RECORD 0-----
a  | 1
b  | 2.0
c  | string1
d  | 2000-01-01
e  | 2000-01-01 12:00:00
only showing top 1 row

```

```
df.columns
```

```
['a', 'b', 'c', 'd', 'e']
```

```
df.printSchema()
```

```
root
|-- a: long (nullable = true)
|-- b: double (nullable = true)
|-- c: string (nullable = true)
|-- d: date (nullable = true)
|-- e: timestamp (nullable = true)
```

```
df.select("a", "b", "c").describe().show()
```

```
+-----+-----+-----+
|summary|  a|  b|      c|
+-----+-----+-----+
|  count|  3|  3|      3|
|   mean|2.0|3.0|   null|
| stddev|1.0|1.0|   null|
|    min|  1|2.0|string1|
|    max|  3|4.0|string3|
+-----+-----+-----+
```

```
df.collect()
```

```
[Row(a=1, b=2.0, c='string1', d=datetime.date(2000, 1, 1), e=datetime.datetime(2000, 1,
Row(a=2, b=3.0, c='string2', d=datetime.date(2000, 2, 1), e=datetime.datetime(2000, 1,
Row(a=3, b=4.0, c='string3', d=datetime.date(2000, 3, 1), e=datetime.datetime(2000, 1,
```



```
df.take(1)
```

```
[Row(a=1, b=2.0, c='string1', d=datetime.date(2000, 1, 1), e=datetime.datetime(2000, 1,
```



```
df.toPandas()
```

```
df.a
```

```
Column<'a'>
```

```
from pyspark.sql import Column
from pyspark.sql.functions import upper
```

```
type(df.c) == type(upper(df.c)) == type(df.c.isNull())
```

```
True
```

```
df.select(df.c).show()
```

```
+-----+
|      c|
+-----+
|string1|
|string2|
|string3|
+-----+
```

```
df.withColumn('upper_c', upper(df.c)).show()
```

```
+---+---+-----+-----+-----+-----+
| a| b|      c|      d|      e|upper_c|
+---+---+-----+-----+-----+-----+
|  1|2.0|string1|2000-01-01|2000-01-01 12:00:00|STRING1|
|  2|3.0|string2|2000-02-01|2000-01-02 12:00:00|STRING2|
|  3|4.0|string3|2000-03-01|2000-01-03 12:00:00|STRING3|
+---+---+-----+-----+-----+-----+
```

```
df.filter(df.a == 1).show()
```

```
+---+---+-----+-----+-----+
| a| b|      c|      d|      e|
+---+---+-----+-----+-----+
|  1|2.0|string1|2000-01-01|2000-01-01 12:00:00|
+---+---+-----+-----+-----+
```

```
df = spark.createDataFrame([
    ['red', 'banana', 1, 10], ['blue', 'banana', 2, 20], ['red', 'carrot', 3, 30],
    ['blue', 'grape', 4, 40], ['red', 'carrot', 5, 50], ['black', 'carrot', 6, 60],
```

```
[ 'red', 'banana', 7, 70], [ 'red', 'grape', 8, 80]], schema=[ 'color', 'fruit', 'v1', 'v2']
df.show()
```

```
+-----+-----+-----+
|color| fruit| v1| v2|
+-----+-----+-----+
|  red|banana|  1| 10|
| blue|banana|  2| 20|
|  red|carrot|  3| 30|
| blue| grape|  4| 40|
|  red|carrot|  5| 50|
|black|carrot|  6| 60|
|  red|banana|  7| 70|
|  red| grape|  8| 80|
+-----+-----+-----+
```

```
df.groupby('color').avg().show()
```

```
+-----+-----+-----+
|color|avg(v1)|avg(v2)|
+-----+-----+-----+
|  red|    4.8|   48.0|
| blue|    3.0|   30.0|
|black|    6.0|   60.0|
+-----+-----+-----+
```

```
def plus_mean(pandas_df):
    return pandas_df.assign(v1=pandas_df.v1 - pandas_df.v1.mean())
```

```
df.groupby('color').applyInPandas(plus_mean, schema=df.schema).show()
```

```
+-----+-----+-----+
|color| fruit| v1| v2|
+-----+-----+-----+
|black|carrot|  0| 60|
| blue|banana| -1| 20|
| blue| grape|  1| 40|
|  red|banana| -3| 10|
|  red|carrot| -1| 30|
|  red|carrot|  0| 50|
|  red|banana|  2| 70|
|  red| grape|  3| 80|
+-----+-----+-----+
```

```
df1 = spark.createDataFrame(
    [(20000101, 1, 1.0), (20000101, 2, 2.0), (20000102, 1, 3.0), (20000102, 2, 4.0)],
    ('time', 'id', 'v1'))
```

```
df2 = spark.createDataFrame(
    [(20000101, 1, 'x'), (20000101, 2, 'y')],
    ('time', 'id', 'v2'))

def asof_join(l, r):
    return pd.merge_asof(l, r, on='time', by='id')

df1.groupby('id').cogroup(df2.groupby('id')).applyInPandas(
    asof_join, schema='time int, id int, v1 double, v2 string').show()
```

```
+-----+-----+-----+
|   time| id| v1| v2|
+-----+-----+-----+
|20000101| 1|1.0| x|
|20000102| 1|3.0| x|
|20000101| 2|2.0| y|
|20000102| 2|4.0| y|
+-----+-----+-----+
```

```
df.write.csv('foo.csv', header=True)
spark.read.csv('foo.csv', header=True).show()
```

```
+-----+-----+-----+
|color| fruit| v1| v2|
+-----+-----+-----+
|  red|carrot| 5| 50|
|black|carrot| 6| 60|
|  red|banana| 7| 70|
|  red| grape| 8| 80|
|  red|banana| 1| 10|
| blue|banana| 2| 20|
|  red|carrot| 3| 30|
| blue| grape| 4| 40|
+-----+-----+-----+
```

```
df.write.parquet('bar.parquet')
spark.read.parquet('bar.parquet').show()
```

```
+-----+-----+-----+
|color| fruit| v1| v2|
+-----+-----+-----+
|  red|carrot| 5| 50|
|black|carrot| 6| 60|
|  red|banana| 7| 70|
|  red| grape| 8| 80|
|  red|banana| 1| 10|
| blue|banana| 2| 20|
|  red|carrot| 3| 30|
| blue| grape| 4| 40|
```

```
+-----+-----+-----+
```

```
df.write.orc('zoo.orc')
spark.read.orc('zoo.orc').show()
```

```
+-----+-----+-----+
|color| fruit| v1| v2|
+-----+-----+-----+
|  red|carrot|  5| 50|
|black|carrot|  6| 60|
|  red|banana|  7| 70|
|  red| grape|  8| 80|
|  red|banana|  1| 10|
| blue|banana|  2| 20|
|  red|carrot|  3| 30|
| blue| grape|  4| 40|
+-----+-----+-----+
```

```
df.createOrReplaceTempView("tableA")
spark.sql("SELECT count(*) from tableA").show()
```

```
+-----+
|count(1)|
+-----+
|      8|
+-----+
```


✓ 0s completed at 3:23 PM

