

reading data from a table or creating data frame from the external files

```
# Read CSV file into table
spark.read.option("header",true).csv("/Users/admin/simple-zipcodes.csv")
spark.read.csv("/Users/admin/simple-zipcodes.csv",header=True,inferschema=True)
# reading data from text file
df2 = spark.read.text("/src/resources/file.txt")
#Creating from JSON file
df2 = spark.read.json("/src/resources/file.json")

#Using append save mode, we can append a dataframe to an existing parquet file. Incase to overwrite use overwrite save mode.

df.write.mode('append').parquet("/tmp/output/people.parquet")
df.write.mode('overwrite').parquet("/tmp/output/people.parquet")
#writing with partition
df.write.partitionBy("gender","salary").mode("overwrite").parquet("/tmp/output/people2.parquet")

# writing to csv
# Saving modes
df2.write.mode('overwrite').csv("/tmp/spark_output/zipcodes")
# You can also use this
df2.write.format("csv").mode('overwrite').save("/tmp/spark_output/zipcodes")
# Read CSV file into table
spark.read.option("header",true).csv("/Users/admin/simple-zipcodes.csv")
spark.read.csv("/Users/admin/simple-zipcodes.csv",header=True,inferschema=True)

"""
overwrite - mode is used to overwrite the existing file.

append - To add the data to the existing file.

ignore - Ignores write operation when the file already exists.

error - This is a default option when the file already exists, it returns an error."""
```

```
|| rdd = spark.sparkContext.parallelize([(i, i**2, i**3) for i in range(10)])
```

```
|| df=rdd.toDF()
```

```
|| df.show()
+---+---+---+
| _1| _2| _3|
+---+---+---+
| 0| 0| 0|
| 1| 1| 1|
| 2| 4| 8|
| 3| 9| 27|
| 4| 16| 64|
| 5| 25| 125|
| 6| 36| 216|
| 7| 49| 343|
| 8| 64| 512|
| 9| 81| 729|
+---+---+---+
```

```
|| df=df.withColumnRenamed("_1","ID").withColumnRenamed("_2","squareID").withColumnRenamed("_3","cubeID")
```

```
|| from pyspark.sql.functions import *
```

```
|| df.withColumn("New",df.squareID*3).show()
```

```
+---+---+---+---+
| ID|squareID|cubeID|New|
+---+---+---+---+
| 0| 0| 0| 0|
| 1| 1| 1| 3|
| 2| 4| 8| 12|
| 3| 9| 27| 27|
| 4| 16| 64| 48|
| 5| 25| 125| 75|
| 6| 36| 216| 108|
| 7| 49| 343| 147|
| 8| 64| 512| 192|
| 9| 81| 729| 243|
+---+---+---+---+
```

```
data = [('James','Smith','M',3000), ('Anna','Rose','F',4100),
        ('Robert','Williams','M',6200)]

columns = ["firstname","lastname","gender","salary"]

df1=spark.createDataFrame(data=data,schema=columns)
```

```
|| df1.createOrReplaceTempView("person")
```

```
|| dfT=spark.sql("select salary,avg(salary) as Average_salary from person where firstname in ('Anna','James') group by salary")
```

```
|| dfT.show()
+---+---+
|salary|Average_salary|
+---+---+
| 4100| 4100.0|
| 3000| 3000.0|
+---+---+
```

```
|| df1.show()
+---+---+---+---+
|firstname|lastname|gender|salary|
+---+---+---+---+
| James| Smith| M| 3000|
| Anna| Rose| F| 4100|
```

```
data = [{"first_name": "James", "last_name": "Smith", "dob": "1991-04-01", "gender": "M", "salary": 3800},
        {"first_name": "Michael", "last_name": "Rose", "dob": "2000-05-19", "gender": "M", "salary": 4000},
        {"first_name": "Robert", "last_name": "Williams", "dob": "1978-09-05", "gender": "M", "salary": 4000},
        {"first_name": "Maria", "last_name": "Anne Jones", "dob": "1967-12-01", "gender": "F", "salary": 4000},
        {"first_name": "Jen", "last_name": "Mary", "dob": "1980-02-17", "gender": "F", "salary": 4100}]

columns = ["first_name", "last_name", "dob", "gender", "salary"]
df2 = spark.createDataFrame(data=data, schema = columns)
```

```
table=spark.sql("select * From Table").show()
```

firstname	middlename	lastname	dob	gender	salary
James		Smith	1991-04-01	M	3000
Michael	Rose		2000-05-19	M	4000
Robert		Williams	1978-09-05	M	4000
Maria	Anne	Jones	1967-12-01	F	4000
Jen	Mary	Brown	1980-02-17	F	-1

```
spark.sql("select count(dob),salary from table group by salary").show()
```

count(dob)	salary
3	4000
1	3000
1	-1

```
spark.sql("select * from tier1_edw_core.email").printSchema()
spark.sql("select * from tier1_edw_core.email").show()
```

1

7

1

```
|1925e75e-5e1c-4f7...| 1|
|2bfe2159-d5d1-448...| 1|
|e7b9e725-f267-48b...| 1|
|e92733e1-c551-491...| 1|
|f39a6677-a793-43f...| 1|
```

```
spark.sql("select * from tier1_dap_marketing_email.visitor_preference").printSchema()
spark.sql("select * from tier1_dap_marketing_email.visitor_preference").show()
```

```
root
 |-- uuid: string (nullable = true)
 |-- emailhash: string (nullable = true)
 |-- userrole: string (nullable = true)
 |-- category: string (nullable = true)
 |-- subscribed: boolean (nullable = true)
 |-- source: string (nullable = true)
 |-- updatetime: long (nullable = true)
 |-- min_subscribe_date: string (nullable = true)
 |-- brandid: integer (nullable = true)

+-----+-----+-----+-----+-----+-----+-----+-----+
|      uuid|      emailhash|userrole| category|subscribed|  source|  updatetime| min_subscribe_date|brandid|
+-----+-----+-----+-----+-----+-----+-----+-----+
|67fbd95f-eef0-44a...|a91d6307dcb076746...|TRAVELER|PROMOTION|   false|    null|1594314054450|2014-04-16 19:38:...|   126|
|f1f7ad7c-09b7-443...|e5eeaid3cee5c143b...|TRAVELER|PROMOTION|    true| booking|1675886322603|2014-04-16 19:38:...|   126|
|1374188a-3b4f-44f...|d71bd207ad2bf3493...|TRAVELER|PROMOTION|   false|    null|1594314054450|2014-04-16 19:38:...|   126|
|16b8b911-3945-414...|3114e7e0345ff6073...|TRAVELER|PROMOTION|   false| booking|1675885554109|2014-04-16 19:38:...|   126|
|3b49450a-6c1d-445...|f57hfa5f6f32eefha...|TRAVELER|PROMOTION|   false|BACKFILL|1594314054450|2014-04-16 19:38:...|   126|
```

```
spark.sql("select * from tier1_dap_marketing_email.visitor_preference v left join tier1_edw_core.email e on v.uuid=e.user_uuid").show()
```

uuid	emailhash	userrole	category	subscribed	source	updatetime	min_subscribe_date	brandid	source_id	first_name_parallax	last_name_parallax
67fbd95f-eef0-44a...	a91d6307dcb076746...	TRAVELER	PROMOTION	false	BACKFILL	1594314054450	2014-04-16 19:38:...	126	615 00000839-6988-4f1...	\$parallax\$Mi7TaWe...	\$parallax\$z8Macsv...
f1f7ad7c-09b7-443...	e5eeaid3cee5c143b...	TRAVELER	PROMOTION	true	booking	1675886322603	2014-04-16 19:38:...	126	00000839-6988-4f1...	00000839-6988-4f1...	00000839-6988-4f1...
1374188a-3b4f-44f...	d71bd207ad2bf3493...	TRAVELER	PROMOTION	false	BACKFILL	1594314054450	2014-04-16 19:38:...	126	00000839-6988-4f1...	00000839-6988-4f1...	00000839-6988-4f1...
16b8b911-3945-414...	3114e7e0345ff6073...	TRAVELER	PROMOTION	false	booking	1675885554109	2014-04-16 19:38:...	126	00000839-6988-4f1...	00000839-6988-4f1...	00000839-6988-4f1...
3b49450a-6c1d-445...	f57hfa5f6f32eefha...	TRAVELER	PROMOTION	false	BACKFILL	1594314054450	2014-04-16 19:38:...	126	00000839-6988-4f1...	00000839-6988-4f1...	00000839-6988-4f1...

1