

```
In [24]: import pandas as pd
        from pyspark.sql import SparkSession
        import findspark
        from pyspark.sql import Row

In [2]: findspark.init()
        findspark.find()

Out[2]: '/Users/youngjinseo/anaconda3/lib/python3.10/site-packages/pyspark'

In [3]: spark = SparkSession.builder.appName('practice1').config("spark.driver.bindAddress", "127.0.0.1").getOrCreate()
        spark

Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
24/09/19 10:46:10 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable

Out[3]: SparkSession - in-memory
```

SparkContext

Spark UI

Version	v3.5.2
Master	local[*]
AppName	practice1

```
In [4]: df = spark.read.csv('/Users/youngjinseo/Desktop/미이생/tested.csv', header = True )
        df.show()

+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|PassengerId|Survived|Pclass|      Name| Sex| Age|SibSp|Parch|      Ticket|  Fare|Cabin|Embarked|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|      892|      0|      3| Kelly, Mr. James| male|34.5|  0|  0| 330911| 7.8292| NULL| Q|
|      893|      1|      3| Wilkes, Mrs. Jame...|female| 47|  1|  0| 363272|  7| NULL| S|
|      894|      0|      2| Myles, Mr. Thomas...| male| 62|  0|  0| 240276| 9.6875| NULL| Q|
|      895|      0|      3| Wirz, Mr. Albert| male| 27|  0|  0| 315154| 8.6625| NULL| S|
|      896|      1|      3| Hirvonen, Mrs. Al...|female| 22|  1|  1| 3101298|12.2875| NULL| S|
|      897|      0|      3| Svensson, Mr. Joh...| male| 14|  0|  0| 7538| 9.225| NULL| S|
|      898|      1|      3| Connolly, Miss. Kate|female| 30|  0|  0| 330972| 7.6292| NULL| Q|
|      899|      0|      2| Caldwell, Mr. Alb...| male| 26|  1|  1| 248738|  29| NULL| S|
|      900|      1|      3| Abraham, Mrs. Jos...|female| 18|  0|  0| 2657| 7.2292| NULL| C|
|      901|      0|      3| Davies, Mr. John...| male| 21|  2|  0| A/4 48871| 24.15| NULL| S|
|      902|      0|      3| Ilieff, Mr. Ylio| male|NULL|  0|  0| 349220| 7.8958| NULL| S|
|      903|      0|      1| Jones, Mr. Charle...| male| 46|  0|  0| 694|  26| NULL| S|
|      904|      1|      1| Snyder, Mrs. John...|female| 23|  1|  0| 21228|82.2667| B45| S|
|      905|      0|      2| Howard, Mr. Benjamin| male| 63|  1|  0| 24065|  26| NULL| S|
|      906|      1|      1| Chaffee, Mrs. Her...|female| 47|  1|  0| W.E.P. 5734| 61.175| E31| S|
|      907|      1|      2| del Carlo, Mrs. S...|female| 24|  1|  0| SC/PARIS 2167|27.7208| NULL| C|
|      908|      0|      3| Keane, Mr. Daniel| male| 35|  0|  0| 233734| 12.35| NULL| Q|
|      909|      0|      3| Assaf, Mr. Gerios| male| 21|  0|  0| 2692| 7.225| NULL| C|
|      910|      1|      3| Ilmakangas, Miss...|female| 27|  1|  0| STON/O2. 3101270| 7.925| NULL| S|
|      911|      1|      3| Assaf Khalil, Mr...|female| 45|  0|  0| 2696| 7.225| NULL| C|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
only showing top 20 rows
```

## GROUPBY

```
In [5]: group_df = df.groupBy('Sex').count()
        group_df.show()

+-----+-----+
| Sex|count|
+-----+-----+
|female| 152|
| male| 266|
+-----+-----+
```

## sort & orderBy

```
In [6]: sort_df = df.sort('PassengerId')
        sort_df.show()

+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|PassengerId|Survived|Pclass|      Name| Sex| Age|SibSp|Parch|      Ticket|  Fare|Cabin|Embarked|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|      1000|      0|      3| Willer, Mr. Aaro...| male|NULL|  0|  0| 3410| 8.7125| NULL| S| |
|      1001|      0|      2| Swane, Mr. George| male|18.5|  0|  0| 248734|  13| F| S|
|      1002|      0|      2| Stanton, Mr. Samu...| male| 41|  0|  0| 237734| 15.0458| NULL| C|
|      1003|      1|      3| Shine, Miss. Elie...|female|NULL|  0|  0| 330968| 7.7792| NULL| Q|
|      1004|      1|      1| Evans, Miss. Edit...|female| 36|  0|  0| PC 17531| 31.6792| A29| C|
|      1005|      1|      3| Buckley, Miss. Ka...|female|18.5|  0|  0| 329944| 7.2833| NULL| Q|
|      1006|      1|      1| Straus, Mrs. Isid...|female| 63|  1|  0| PC 17483|221.7792|C55 C57| S|
|      1007|      0|      3| Chronopoulos, Mr...| male| 18|  1|  0| 2680| 14.4542| NULL| C|
|      1008|      0|      3| Thomas, Mr. John| male|NULL|  0|  0| 2681| 6.4375| NULL| C|
|      1009|      1|      3| Sandstrom, Miss. ...|female|  1|  1|  1| PP 9549|  16.7| G6| S|
|      1010|      0|      1| Beattie, Mr. Thomso| male| 36|  0|  0| 13050| 75.2417| C6| C|
|      1011|      1|      2| Chapman, Mrs. Joh...|female| 29|  1|  0| SC/AH 29037|  26| NULL| S|
|      1012|      1|      2| Watt, Miss. Bertha J|female| 12|  0|  0| C.A. 33595| 15.75| NULL| S|
|      1013|      0|      3| Kiernan, Mr. John| male|NULL|  1|  0|  0| 367227|  7.75| NULL| Q|
|      1014|      1|      1| Schabert, Mrs. Pa...|female| 35|  1|  0| 13236| 57.75| C28| C|
|      1015|      0|      3| Carver, Mr. Alfre...| male| 28|  0|  0| 392095| 7.25| NULL| S|
|      1016|      0|      3| Kennedy, Mr. John| male|NULL|  0|  0|  0| 368783| 7.75| NULL| Q|
|      1017|      1|      3| Cribb, Miss. Laur...|female| 17|  0|  1| 371362| 16.1| NULL| S|
|      1018|      0|      3| Brobeck, Mr. Karl...| male| 22|  0|  0| 350045| 7.7958| NULL| S|
|      1019|      1|      3| McCoy, Miss. Alicia|female|NULL|  2|  0|  0| 367226| 23.25| NULL| Q|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
only showing top 20 rows
```

```
In [7]: order_df = df.orderBy('PassengerId')
        order_df.show()

+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|PassengerId|Survived|Pclass|      Name| Sex| Age|SibSp|Parch|      Ticket|  Fare|Cabin|Embarked|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|      1000|      0|      3| Willer, Mr. Aaro...| male|NULL|  0|  0| 3410| 8.7125| NULL| S| |
|      1001|      0|      2| Swane, Mr. George| male|18.5|  0|  0| 248734|  13| F| S|
|      1002|      0|      2| Stanton, Mr. Samu...| male| 41|  0|  0| 237734| 15.0458| NULL| C|
|      1003|      1|      3| Shine, Miss. Elie...|female|NULL|  0|  0| 330968| 7.7792| NULL| Q|
|      1004|      1|      1| Evans, Miss. Edit...|female| 36|  0|  0| PC 17531| 31.6792| A29| C|
|      1005|      1|      3| Buckley, Miss. Ka...|female|18.5|  0|  0| 329944| 7.2833| NULL| Q|
|      1006|      1|      1| Straus, Mrs. Isid...|female| 63|  1|  0| PC 17483|221.7792|C55 C57| S|
|      1007|      0|      3| Chronopoulos, Mr...| male| 18|  1|  0| 2680| 14.4542| NULL| C|
|      1008|      0|      3| Thomas, Mr. John| male|NULL|  0|  0| 2681| 6.4375| NULL| C|
|      1009|      1|      3| Sandstrom, Miss. ...|female|  1|  1|  1| PP 9549|  16.7| G6| S|
|      1010|      0|      1| Beattie, Mr. Thomso| male| 36|  0|  0| 13050| 75.2417| C6| C|
|      1011|      1|      2| Chapman, Mrs. Joh...|female| 29|  1|  0| SC/AH 29037|  26| NULL| S|
|      1012|      1|      2| Watt, Miss. Bertha J|female| 12|  0|  0| C.A. 33595| 15.75| NULL| S|
|      1013|      0|      3| Kiernan, Mr. John| male|NULL|  1|  0|  0| 367227|  7.75| NULL| Q|
|      1014|      1|      1| Schabert, Mrs. Pa...|female| 35|  1|  0| 13236| 57.75| C28| C|
|      1015|      0|      3| Carver, Mr. Alfre...| male| 28|  0|  0| 392095| 7.25| NULL| S|
|      1016|      0|      3| Kennedy, Mr. John| male|NULL|  0|  0|  0| 368783| 7.75| NULL| Q|
|      1017|      1|      3| Cribb, Miss. Laur...|female| 17|  0|  1| 371362| 16.1| NULL| S|
|      1018|      0|      3| Brobeck, Mr. Karl...| male| 22|  0|  0| 350045| 7.7958| NULL| S|
|      1019|      1|      3| McCoy, Miss. Alicia|female|NULL|  2|  0|  0| 367226| 23.25| NULL| Q|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
only showing top 20 rows
```

## na.drop & na.fill

```
In [8]: nadrop_df = df.na.drop()
        nadrop_df.show(10)

+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|PassengerId|Survived|Pclass|      Name| Sex| Age|SibSp|Parch|      Ticket|  Fare|Cabin|Embarked|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|      904|      1|      1| Snyder, Mrs. John...|female| 23|  1|  0| 21228|82.2667| B45| S|
|      906|      1|      1| Chaffee, Mrs. Her...|female| 47|  1|  0| W.E.P. 5734| 61.175| E31| S|
|      916|      1|      1| Ryerson, Mrs. Art...|female| 48|  1|  3| PC 17608|262.375|B57 B59 B63 B66| C|
|      918|      1|      1| Ostby, Miss. Hele...|female| 22|  0|  1| 113509|61.9792| B36| C|
|      920|      0|      1| Brady, Mr. John B...| male| 41|  0|  1| 240276| 30.5| A21| S|
|      926|      1|      1| Mook, Mr. Philip...| male| 30|  1|  0| 13236| 57.75| C78| C|
|      936|      0|      1| Kimball, Mrs. Edw...|female| 45|  1|  0| 11753|52.5542| D19| S|
|      938|      0|      1| Chevre, Mr. Paul ...| male| 45|  0|  0| PC 17594| 29.77| A9| C|
|      940|      1|      1| Bucknell, Mrs. Wi...|female| 60|  0|  0| 11813|76.2917| D15| C|
|      942|      0|      1| Smith, Mr. Lucien...| male| 24|  1|  0| 13695|  60| C31| S|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
only showing top 10 rows
```

```
In [9]: fill_df = df.na.fill({'Cabin':None'})
        fill_df.show(10)

+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|PassengerId|Survived|Pclass|      Name| Sex| Age|SibSp|Parch|      Ticket|  Fare|Cabin|Embarked|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|      892|      0|      3| Kelly, Mr. James| male|34.5|  0|  0| 330911| 7.8292| None| Q|
|      893|      1|      3| Wilkes, Mrs. Jame...|female| 47|  1|  0| 363272|  7| None| S|
|      894|      0|      2| Myles, Mr. Thomas...| male| 62|  0|  0| 240276| 9.6875| None| Q|
|      895|      0|      3| Wirz, Mr. Albert| male| 27|  0|  0| 315154| 8.6625| None| S|
|      896|      1|      3| Hirvonen, Mrs. Al...|female| 22|  1|  1| 3101298|12.2875| None| S|
|      897|      0|      3| Svensson, Mr. Joh...| male| 14|  0|  0| 7538| 9.225| None| S|
|      898|      1|      3| Connolly, Miss. Kate|female| 30|  0|  0| 330972| 7.6292| None| Q|
|      899|      0|      2| Caldwell, Mr. Alb...| male| 26|  1|  1| 248738|  29| None| S|
|      900|      1|      3| Abraham, Mrs. Jos...|female| 18|  0|  0| 2657| 7.2292| None| C|
|      901|      0|      3| Davies, Mr. John ...| male| 21|  2|  0| A/A 4 48871| 24.15| None| S|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
only showing top 10 rows
```

## na.replace

```
In [11]: replace_df = df.na.replace("female","woman")
        replace_df.show()

+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|PassengerId|Survived|Pclass|      Name| Sex| Age|SibSp|Parch|      Ticket|  Fare|Cabin|Embarked|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|      892|      0|      3| Kelly, Mr. James| male|34.5|  0|  0| 330911| 7.8292| NULL| Q|
|      893|      1|      3| Wilkes, Mrs. Jame...|woman| 47|  1|  0| 363272|  7| NULL| S|
|      894|      0|      2| Myles, Mr. Thomas...| male| 62|  0|  0| 240276| 9.6875| NULL| Q|
|      895|      0|      3| Wirz, Mr. Albert| male| 27|  0|  0| 315154| 8.6625| NULL| S|
|      896|      1|      3| Hirvonen, Mrs. Al...|woman| 22|  1|  1| 3101298|12.2875| NULL| S|
|      897|      0|      3| Svensson, Mr. Joh...| male| 14|  0|  0| 7538| 9.225| NULL| S|
|      898|      1|      3| Connolly, Miss. Kate|woman| 30|  0|  0| 330972| 7.6292| NULL| Q|
|      899|      0|      2| Caldwell, Mr. Alb...| male| 26|  1|  1| 248738|  29| NULL| S|
|      900|      1|      3| Abraham, Mrs. Jos...|woman| 18|  0|  0| 2657| 7.2292| NULL| C|
|      902|      0|      3| Davies, Mr. John ...| male| 21|  2|  0| A/4 48871| 24.15| NULL| S|
|      903|      0|      3| Jones, Mr. Charle...| male| 46|  0|  0| 694|  26| NULL| S|
|      904|      1|      1| Snyder, Mrs. John...|woman| 23|  1|  0| 21228|82.2667| B45| S|
|      905|      0|      2| Howard, Mr. Benjamin| male| 63|  1|  0| 24065|  26| NULL| S|
|      906|      1|      1| Chaffee, Mrs. Her...|woman| 47|  1|  0| W.E.P. 5734| 61.175| E31| S|
|      907|      1|      2| del Carlo, Mrs. S...|woman| 24|  1|  0| SC/PARIS 2167|27.7208| NULL| C|
|      908|      0|      3| Keane, Mr. Daniel| male| 35|  0|  0| 233734| 12.35| NULL| Q|
|      909|      0|      3| Assaf, Mr. Gerios| male| 21|  0|  0| 2692| 7.225| NULL| C|
|      910|      1|      3| Ilmakangas, Miss...|woman| 27|  1|  0| STON/O2. 3101270| 7.925| NULL| S|
|      911|      1|      3| Assaf Khalil, Mr...|woman| 45|  0|  0| 2696| 7.225| NULL| C|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
only showing top 20 rows
```

## describe & summary

```
In [12]: df.describe(['Age']).show()

+-----+-----+
|summary|      Age|
+-----+-----+
| count|      332|
| mean|30.272590361445783|
| stddev|14.181209235624424|
| min|  0.17|
| max|  9|
+-----+-----+

In [13]: df.select('Ticket','Age','Fare').summary().show()

+-----+-----+-----+-----+
|summary|      Ticket|      Age|      Fare|
+-----+-----+-----+-----+
| count|      418|      332|      417|
| mean|223850.98986486485|30.272590361445783| 35.6271884892086|
| stddev| 369523.7764694362|14.181209235624424|55.907576179973844|
| min|      110469|  0.17|  0|
| 25%|      17464.0|  21.0|  7.8958|
| 50%|      230136.0|  27.0| 14.4542|
| 75%|      347080.0|  39.0| 31.5|
| max| W.E.P. 5734|  9| 93.5|
+-----+-----+-----+-----+

24/09/19 10:47:36 WARN SparkStringUtils: Truncated the string representation of a plan since it was too large. This behavior can be adjusted by setting 'spark.sql.debug.maxToStringFields'.
```

## withColumnRenamed

```
In [14]: renamed_df = df.withColumnRenamed("Sex","Gender")
        renamed_df.show()

+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|PassengerId|Survived|Pclass|      Name|Gender| Age|SibSp|Parch|      Ticket|  Fare|Cabin|Embarked|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|      892|      0|      3| Kelly, Mr. James| male|34.5|  0|  0| 330911| 7.8292| NULL| Q|
|      893|      1|      3| Wilkes, Mrs. Jame...|female| 47|  1|  0| 363272|  7| NULL| S|
|      894|      0|      2| Myles, Mr. Thomas...| male| 62|  0|  0| 240276| 9.6875| NULL| Q|
|      895|      0|      3| Wirz, Mr. Albert| male| 27|  0|  0| 315154| 8.6625| NULL| S|
|      896|      1|      3| Hirvonen, Mrs. Al...|female| 22|  1|  1| 3101298|12.2875| NULL| S|
|      897|      0|      3| Svensson, Mr. Joh...| male| 14|  0|  0| 7538| 9.225| NULL| S|
|      898|      1|      3| Connolly, Miss. Kate|female| 30|  0|  0| 330972| 7.6292| NULL| Q|
|      899|      0|      2| Caldwell, Mr. Alb...| male| 26|  1|  1| 248738|  29| NULL| S|
|      900|      1|      3| Abraham, Mrs. Jos...|woman| 18|  0|  0| 2657| 7.2292| NULL| C|
|      901|      0|      3| Davies, Mr. John ...| female| 21|  2|  0| A/4 48871| 24.15| NULL| S|
|      902|      0|      3| Davies, Mr. John ...| male|NULL|  0|  0| 349220| 7.8958| NULL| S|
|      903|      0|      1| Jones, Mr. Charle...| male| 46|  0|  0| 694|  26| NULL| S|
|      904|      1|      1| Snyder, Mrs. John...|female| 23|  1|  0| 21228|82.2667| B45| S|
|      905|      0|      2| Howard, Mr. Benjamin| male| 63|  1|  0| 24065|  26| NULL| S|
|      906|      1|      1| Chaffee, Mrs. Her...|female| 47|  1|  0| W.E.P. 5734| 61.175| E31| S|
|      907|      1|      2| del Carlo, Mrs. S...|female| 24|  1|  0| SC/PARIS 2167|27.7208| NULL| C|
|      908|      0|      3| Keane, Mr. Daniel| male| 35|  0|  0| 233734| 12.35| NULL| Q|
|      909|      0|      3| Assaf, Mr. Gerios| male| 21|  0|  0| 2692| 7.225| NULL| C|
|      910|      1|      3| Ilmakangas, Miss...|female| 27|  1|  0| STON/O2. 3101270| 7.925| NULL| S|
|      911|      1|      3| Assaf Khalil, Mr...|female| 45|  0|  0| 2696| 7.225| NULL| C|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
only showing top 20 rows
```

## Head

```
In [15]: head_row = df.head()
        print(head_row)

Row(PassengerId='892', Survived=0, Pclass=3, Name='Kelly, Mr. James', Sex='male', Age=34.5, SibSp=0, Parch=0, Ticket='330911', Fare='7.8292', Cabin=None, Embarked='Q')
```

## toPandas

```
In [17]: pandas_df = df.toPandas()
        pandas_df.head(5)

Out[17]:
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	892	0	3	Kelly, Mr. James	male	34.5	0	0	330911	7.8292	None	Q
1	893	1	3	Wilkes, Mrs. James (Ellen Needs)	female	47	1	0	363272	7	None	S
2	894	0	2	Myles, Mr. Thomas Francis	male	62	0	0	240276	9.6875	None	Q
3	895	0	3	Wirz, Mr. Albert	male	27	0	0	315154	8.6625	None	S
4	896	1	3	Hirvonen, Mrs. Alexander (Helga E Lindqvist)	female	22	1	1	3101298	12.2875	None	S

```

In [27]: from pyspark.sql import Row

#예제 데이터프레임
data = {
    'Name': ['Alice', 'Bob', 'Charlie'],
    'Age': [24, 27, 22],
    'City': ['New York', 'Los Angeles', 'Chicago']}

In [28]: #DataFrame
data = pd.DataFrame(data)
df2 = spark.createDataFrame(data)
df2.show()

+-----+-----+-----+
| Name|Age|      City|
+-----+-----+-----+
| Alice| 24| New York|
| Bob| 27| Los Angeles|
| Charlie| 22| Chicago|
+-----+-----+-----+

In [29]: #Row
df3= spark.createDataFrame([
    Row(Name = 'Alice',Age = 24, City = 'New York'),
    Row(Name = 'Bob',Age = 27, City = 'Los Angeles'),
    Row(Name = 'Charlie',Age = 22, City = 'Chicago'),
])

df3.show()

+-----+-----+-----+
| Name|Age|      City|
+-----+-----+-----+
| Alice| 24| New York|
| Bob| 27| Los Angeles|
| Charlie| 22| Chicago|
+-----+-----+-----+
```