```
from pyspark.sql import SparkSession
spark = SparkSession \
   .builder \
    .appName("Python Spark SQL basic example") \
    .config("spark.some.config.option", "some-value") \
    .getOrCreate()
spark
SparkSession - in-memory
     SparkContext
     Spark UI
     Version
          v3.5.5
     Master
          local[*]
     AppName
          Python Spark SQL basic example
import os
os.getcwd()
→ '/content'
df = spark.read.csv("diabetes.csv", header=True,inferSchema=True)
df.show(20)
     |Pregnancies|Glucose|BloodPressure|SkinThickness|Insulin| BMI|DiabetesPedigreeFunction| Age|Outcome|
                                                            0|33.6|
                                                                                       0.627 | 50.0 |
                                     72
                1
                       85
                                     66
                                                    29
                                                            0 26.6
                                                                                       0.351 31.0
                8
                     183
                                     64
                                                    0
                                                           0|23.3|
                                                                                      0.672|32.0|
                                                                                                         1
                      89
                                   NULL
                                                   23
                                                           94 NULL
                                                                                       NULL | NaN |
                                                                                                         0
                1
                01
                     137
                                                        168 43.1
                                                                                       2.288 | 33.0 |
                                     40
                                                   35
                                                                                                         11
                                                                                       0.201|30.0|
                5 l
                     NULL
                                     74
                                                    a١
                                                            0 25.6
                                                                                                         01
                                                   321
                                                           88| 31|
                3 |
                      78
                                     501
                                                                                       0.248 26.01
                                                                                                        1 |
               10
                      115
                                      01
                                                    01
                                                            0|35.3|
                                                                                       0.134 | 29.0 |
                                                                                                         01
                2
                      197
                                     70
                                                    45|
                                                          543|30.5|
                                                                                       0.158|53.0|
                                                                                                         1
                8
                      125
                                     96
                                                    0
                                                           0 0
                                                                                       0.232 | 54.0 |
                                                                                                         1
                4
                      110
                                      ?|
                                                    01
                                                            0|37.6|
                                                                                       0.191|30.0|
                                                                                                         0
                                                                                       0.537 34.0
               10|
                      168
                                     74 l
                                                    0
                                                            0 38
                                                                                                         1|
                                                    øj
                                                             0 27.1
                                                                                       1.441 57.0
               10
                      139
                                     80
                                                                                                         0
                      189
                                     60
                                                    23
                                                           846 | 30.1 |
                                                                                       0.398 | 59.0 |
                1
                                                                                                         1
                                                                                       0.587 | 51.0 |
                5
                      166
                                     72 l
                                                    19
                                                          175 | 25.8 |
                                                                                                         11
                7 l
                                                            01 301
                                                                                       0.484|32.0|
                      100
                                      01
                                                    01
                                                                                                         11
                                                                                       0.551|31.0|
                0
                      118
                                     84
                                                    47 l
                                                           230 | 45.8 |
                                                                                                         1|
                                     74 İ
                7
                       ##
                                                    01
                                                            0 | 29.6 |
                                                                                       0.254|31.0|
                                                                                                        1
                1
                      103
                                     30
                                                    38
                                                           83 | 43.3 |
                                                                                       0.183|33.0|
                                                                                                         0
                1|
                      115
                                     70
                                                    30
                                                           96|34.6|
                                                                                       0.529 | 32.0 |
                                                                                                         1|
     only showing top 20 rows
df.head(5)
🔁 [Row(Pregnancies=6, Glucose='148', BloodPressure='72', SkinThickness=35, Insulin=0, BMI='33.6', DiabetesPedigreeFunction=0.627,
     Age=50.0, Outcome=1),
      Row(Pregnancies=1, Glucose='85', BloodPressure='66', SkinThickness=29, Insulin=0, BMI='26.6', DiabetesPedigreeFunction=0.351,
     Age=31.0, Outcome=0),
      Row(Pregnancies=8, Glucose='183', BloodPressure='64', SkinThickness=0, Insulin=0, BMI='23.3', DiabetesPedigreeFunction=0.672,
     Age=32.0, Outcome=1),
     Row(Pregnancies=1, Glucose='89', BloodPressure=None, SkinThickness=23, Insulin=94, BMI=None, DiabetesPedigreeFunction=None,
     Age=nan, Outcome=0),
      Row(Pregnancies=0, Glucose='137', BloodPressure='40', SkinThickness=35, Insulin=168, BMI='43.1', DiabetesPedigreeFunction=2.288,
     Age=33.0, Outcome=1)]
type(df)
```

```
pyspark.sql.dataframe.DataFrame
      def __init__(jdf: JavaObject, sql_ctx: Union['SQLContext', 'SparkSession'])
         name | Benaci | avb( satar y / | max( abe / |
                                                                                                            ML
                    F|
                            150.0
       |PySpark|
                    M
                             75.0
                                         50
df.printSchema()
→ root
      |-- Pregnancies: integer (nullable = true)
      |-- Glucose: string (nullable = true)
      |-- BloodPressure: string (nullable = true)
      |-- SkinThickness: integer (nullable = true)
      |-- Insulin: integer (nullable = true)
      |-- BMI: string (nullable = true)
      |-- DiabetesPedigreeFunction: double (nullable = true)
      -- Age: double (nullable = true)
      |-- Outcome: integer (nullable = true)
df.select("Glucose").show(5)
    +----+
     |Glucose|
          148
          85
          183
          89
          137
         ----+
     only showing top 5 rows
df1 = df.select(df['Glucose'],df['Age']+1)
df1.show()
₹
     |Glucose|(Age + 1)|
          148
                   51.0
          85 l
                   32.0
          183 l
                   33.0
          89
                    NaN
          137
                   34.0
         NULL
                   31.0
          78
                   27.0
          115
                   30.0
          197
                   54.0
          125
                   55.0
          110
                   31.0
          168
                   35.01
          139
                   58.0
          189
                   60.0
          166
                   52.0
          100
                   33.0
          118
                   32.0
           ##
                   32.0
          103
                   34.0
          115
                   33.0
     only showing top 20 rows
df['Glucose']
→ Column<'Glucose'>
df[['Glucose','Age']]
DataFrame[Glucose: string, Age: double]
df.dtypes
[('Pregnancies', 'int'), ('Glucose', 'string'),
```

('BloodPressure', 'string'),

```
('SkinThickness', 'int'),
  ('Insulin', 'int'),
  ('BMI', 'string'),
  ('DiabetesPedigreeFunction', 'double'),
  ('Age', 'double'),
  ('Outcome', 'int')]

dist_bp = df.select('BloodPressure').distinct()
```

dist_bp.show()

df.select('Pregnancies').distinct().count()

→ 17

df.describe().show()

→	++-	+	+			+	+	
_	summary	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesF
-	+	+					+	
	count	768	767	767	768	768	767	
	mean 3	3.84505208333333335	120.91906005221932	69.07963446475196	20.5364583333333332	79.79947916666667	31.993211488250633	0.47
	stddev	3.36957806269887	32.009944976772395	19.363065113384383	15.952217567727642	115.24400235133803	7.892243623420615	0.33
	min	0	##	0	0	0	0	
	max	17	99	?	99	846	?	
	++-						+	

df.describe(['BMI','Glucose']).show()

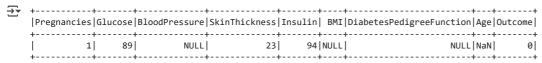
df.describe(['Age']).show()

df['Age']

```
4/3/25, 12:09 PM
                                                                 Atharva pySpark Ass.ipynb - Colab
    → Column<'Age'>
    #adding column
    df_new = df.withColumn('New Age',df.Age+2)
    df_new.show()
         |Pregnancies|Glucose|BloodPressure|SkinThickness|Insulin| BMI|DiabetesPedigreeFunction| Age|Outcome|New Age|
                         148
                                        72
                                                              0|33.6|
                                                                                        0.627|50.0|
                   1|
                         85
                                        66
                                                     29
                                                              0|26.6|
                                                                                        0.351|31.0|
                                                                                                             33.0
                   8
                         183
                                       64
                                                      0
                                                              0 23.3
                                                                                        0.672 32.0
                                                                                                             34.0
                                                                                                        1
                   11
                         89 l
                                      NULL
                                                     23
                                                             94 NULLI
                                                                                        NULL| NaN|
                                                                                                        01
                                                                                                              NaN
                                                           168 43.1
                                                                                        2.288|33.0|
                   al
                         137
                                        40
                                                     35 l
                                                                                                        11
                                                                                                             35.0
                                        74 l
                                                                                        0.201 30.01
                                                                                                        a١
                                                                                                             32.0
                   5 l
                        NULL
                                                      al
                                                             0/25.6
                   3 l
                         78 l
                                        50 l
                                                     32
                                                           88 31
                                                                                        0.248 | 26.0 |
                                                                                                        1
                                                                                                             28.0
                  10
                         115|
                                        0
                                                      0
                                                             0|35.3|
                                                                                        0.134|29.0|
                                                                                                        0|
                                                                                                             31.0
                   2
                         197
                                        70
                                                     45
                                                          543 | 30.5 |
                                                                                        0.158|53.0|
                                                                                                        1
                                                                                                             55.0
                   8
                         125
                                        96
                                                                                        0.232|54.0|
                                                                                                        1|
                                                      01
                                                             01 01
                                                                                                             56.0
                   4
                                        ? [
                                                             0 37.6
                                                                                        0.191 30.0
                                                                                                             32.0
                         110
                                                            0| 38|
0|27.1|
                   10
                         168
                                        74
                                                      0
                                                                                        0.537 34.0
                                                                                                        1|
                                                                                                             36.0
                  10
                         139
                                        80
                                                                                       1.441 57.0
                                                                                                             59.0
                                                      0
                                                                                                        0|
                                                                                        0.398 | 59.0 |
                   1
                         189 l
                                        60 l
                                                     23|
                                                            846|30.1|
                                                                                                        1
                                                                                                             61.0
                   5 İ
                         166
                                       72
                                                     19
                                                           175 25.8
                                                                                       0.587 51.0
                                                                                                        1
                                                                                                             53.0
                   7 l
                         100
                                        01
                                                      01
                                                             0| 30|
                                                                                        0.484|32.0|
                                                                                                        11
                                                                                                             34.0
                   01
                                                     47
                                                            230 | 45.8 |
                                                                                                             33.0
                         118
                                        841
                                                                                       0.551 | 31.0 |
                                                                                                        1
                   7 |
                          ##|
                                        74
                                                      0
                                                             0|29.6|
                                                                                        0.254|31.0|
                                                                                                        1|
                                                                                                             33.0
                   1|
                         103
                                        30
                                                     38
                                                             83 | 43.3 |
                                                                                        0.183|33.0|
                                                                                                        0
                                                                                                             35.0
                   1
                         115
                                        70 l
                                                     30
                                                             96|34.6|
                                                                                        0.529 | 32.0 |
                                                                                                             34.0
        only showing top 20 rows
    #df.select(isnull())
    df new= df new.drop('New Age')
    df new.show()
         |Pregnancies|Glucose|BloodPressure|SkinThickness|Insulin| BMI|DiabetesPedigreeFunction| Age|Outcome|
                  72
                                                              0|33.6|
                                                                                        0.627 | 50.0 |
                   1
                          85
                                        66
                                                     29
                                                              0 | 26.6 |
                                                                                        0.351|31.0|
                        183
                                        64
                                                              0 | 23.3 |
                                                                                        0.672|32.0|
                   8
                                                      01
                                                                                                        11
                         89
                                      NULLİ
                                                     23
                                                             94 NULL
                                                                                        NULL | NaN |
                   1
                                                                                                        0
                                                          168 43.1
                   0
                        137
                                                                                        2.288 33.0
                                        40
                                                     35
                                                                                                        1
                                                             0|25.6|
                   5
                       NULL
                                        741
                                                      01
                                                                                        0.201 | 30.0 |
                                                                                                        01
                                                           88| 31|
                   3 |
                         78 l
                                        50 l
                                                     32
                                                                                        0.248 | 26.0 |
                                                                                                        1
                  10
                         115
                                        01
                                                      01
                                                             0|35.3|
                                                                                        0.134 | 29.0 |
                                                                                                        0
                   2
                         197
                                        70
                                                     45
                                                            543 | 30.5 |
                                                                                        0.158 | 53.0 |
                                                                                                        1
                                                            0 0
                   8
                         125
                                        96
                                                      0
                                                                                        0.232 | 54.0 |
                                                                                                        1
                   41
                         110
                                        ?|
                                                      0
                                                             0|37.6|
                                                                                        0.191 | 30.0 |
                                                                                                        0
                                                                                        0.537|34.0|
                   10
                         168
                                        74
                                                             0| 38|
                                                              0 27.1
                   10
                         139
                                        80
                                                      01
                                                                                        1.441 | 57.0 |
                                                            846 | 30.1 |
                                                                                       0.398 59.0
                   1
                         189
                                        60
                                                     23
                                                                                                        1
                   5
                         166
                                        72
                                                     19
                                                           175 | 25.8 |
                                                                                        0.587 | 51.0 |
                                                                                                        1
                   7 l
                         100
                                                             0| 30|
                                                                                       0.484|32.0|
                                        01
                                                      01
                                                                                                        11
                   al
                         118
                                        84 |
                                                     47
                                                            230 | 45.8 |
                                                                                        0.551|31.0|
                                                                                                        11
                   7 |
                          ## |
                                        74 l
                                                      01
                                                             0 29.6
                                                                                       0.254 31.0
                                                                                                        11
                                                             83 | 43.3 |
                   11
                         103
                                        30
                                                     38
                                                                                        0.183 | 33.0 |
                                                                                                        01
                   1|
                         115|
                                        70
                                                     30
                                                             96|34.6|
                                                                                        0.529|32.0|
                                                                                                        1|
        only showing top 20 rows
```

df_null = df.filter(df['BloodPressure'].isNull())

df_null.show()



df.select('Pregnancies').distinct().show()

```
|Pregnancies|
```

 +
12
1
13
6
3
5
15
9
17
4
8
7
10
11
14
2
0

df.show(20)

→ *	+	+	·		·				++
ت	Pregnancies	Glucose	BloodPressure	SkinThickness	 Insulin	BMI	 DiabetesPedigreeFunction	Age	Outcome
	6	148	72	35	0	33.6	0.627	50.0	1
	1	85	66	29	0	26.6	0.351	31.0	0
	8	183	64	0	0	23.3	0.672	32.0	1
	1	89	NULL	23	94	NULL	NULL	NaN	0
	0	137	40	35	168	43.1	2.288	33.0	1
	5	NULL	74	0	0	25.6	0.201	30.0	0
	3	78	50	32	88	31	0.248	26.0	1
	10	115	0	0	0	35.3	0.134	29.0	0
	2	197	70	45	543	30.5	0.158	53.0	1
	8	125	96	0	0	0	0.232	54.0	1
	4	110		0	0	37.6	0.191	30.0	j 0j
	10	168	74	0	0	38	0.537	34.0	1
	10	139	80	0	0	27.1	1.441	57.0	0
	1	189	60	23	846	30.1	0.398	59.0	1
	j 5	166	72	19	175	25.8	0.587	51.0	1
	7	100	0	0	0	30	0.484	32.0	1
	j 0	118	84	47	230	45.8	0.551	31.0	j 1
	j 7	##	:	0	0	29.6	0.254	31.0	j 1
	j 1	103	30	38	83	43.3	•		
	j 1	115	:	30		34.6	•		
	+	+	++		+		+	+	++

only showing top 20 rows

df1 = df.na.drop(how='any')

df1.show(20)

Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outc
6l	+ 148	72	 35	 0	 33.6	+ 0.627	+ 50.0	+
1	85	66	29		26.6		31.0	İ
8	183	64	0	0	23.3	0.672	32.0	İ
0	137	40	35	168	43.1	2.288	33.0	ĺ
3	78	50	32	88	31	0.248	26.0	İ
10	115	0	0	0	35.3	0.134	29.0	ĺ
2	197	70	45	543	30.5	0.158	53.0	
8	125	96	0	0	0	0.232	54.0	
4	110	?	0	0	37.6	0.191	30.0	
10	168	74	0	0	38	0.537	34.0	
10	139	80	0	0	27.1	1.441	57.0	
1	189	60	23	846	30.1	0.398	59.0	
5	166	72	19	175	25.8	0.587	51.0	
7	100	0	0	0	30	0.484	32.0	
0	118	84	47	230	45.8	0.551	31.0	
7	##	74	0	0	29.6	0.254	31.0	
1	103	30	38	83	43.3	0.183	33.0	
1	115	70	30	96	34.6	0.529	32.0	
3	126	88		235	39.3			
8	99	84	0	0	?	0.388	50.0	

only showing top 20 rows

df1 = df.na.drop(how='any',thresh=7)

df1.show(20)

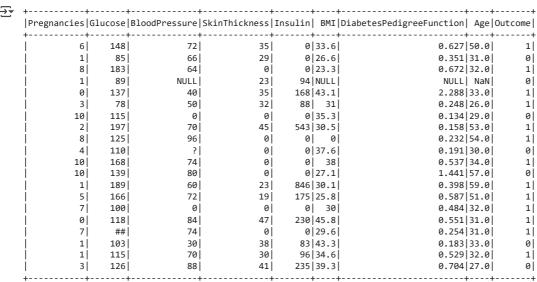
₹

egnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcom
6	148	72	35	0	33.6	0.627	50.0	
1	85	66	29	0	26.6	0.351	31.0	
8	183	64			23.3	•		
0	137	40			43.1	•		
5	NULL	74			25.6	•		
3	78	50			31	•		
10		0			35.3	•		
2		70			30.5			
8		96		0				
4		}			37.6	•		
10								
10		80			27.1	•		
1	189	60			30.1	•		
5	166				25.8	•		
7	100	0		0				
0	118	84			45.8	•		
7	##	74			29.6	•		
1		30			43.3	•		
1		70			34.6			
3	126	88	41	235	39.3	0.704	27.0	1

only showing top 20 rows

df1 = df.na.drop(how='any',subset=['Glucose'])

df1.show(20)



only showing top 20 rows

df1 = df.na.fill("Missing",subset=["BloodPressure","Glucose"])

df1.show(20)

∓

+		·				++	+
Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction Age	Outcome
+		·				++	+
6	148	72	35	0	33.6	0.627 50.0	1
1	85	66	29	0	26.6	0.351 31.0	0
8	183	64	0	0	23.3	0.672 32.0	1
1	89	Missing	23	94	NULL	NULL NaN	0
0	137	40	35	168	43.1	2.288 33.0	1
5	Missing	74	0	0	25.6	0.201 30.0	0
3	78	50	32	88	31	0.248 26.0	1
10	115	0	0	0	35.3	0.134 29.0	0
2	197	70	45	543	30.5	0.158 53.0	1
8	125	96	0	0	0	0.232 54.0	1
4	110	?	0	0	37.6	0.191 30.0	0
10	168	74	0	0	38	0.537 34.0	1
10	139	80	0	0	27.1	1.441 57.0	0
1	189	60	23	846	30.1	0.398 59.0	1
5	166	72	19	175	25.8	0.587 51.0	1
7	100	0	0	0	30	0.484 32.0	1
0	118	84	47	230	45.8	0.551 31.0	1
7	##	74	0	0	29.6	0.254 31.0	1
1	103	30	38	83	43.3	0.183 33.0	0
1	115	70	30	96	34.6	0.529 32.0	1

₹

only showing top 20 rows

from pyspark.ml.feature import Imputer

imputer = Imputer(inputCols=['Age'], outputCols=['Imputed Age'],strategy="mean")

imputer.fit(df).transform(df).show(20)

Imputed Ag	Outcome	Age	DiabetesPedigreeFunction	BMI	Insulin	SkinThickness	BloodPressure	ilucose	Pregnancies
50.	1	50.0	0.627	33.6	0	35	72	148	6
31.	0	31.0	0.351	26.6	0	29	66	85	1
32.	1	32.0	0.672	23.3	0	0	64	183	8
33.2568448500651	0	NaN	NULL	NULL	94	23	NULL	89	1
33.	1	33.0	2.288	43.1	168	35	40	137	0
30.	0	30.0	0.201	25.6	0	0	74	NULL	5
26.	1	26.0	0.248	31	88	32	50	78	3
29.	0	29.0	0.134	35.3	0	0	0	115	10
53.	1	53.0	0.158	30.5	543	45	70	197	2
54.	1	54.0	0.232	0	0	0	96	125	8
30.	0	30.0	0.191	37.6	0	0	?	110	4
34.	1	34.0	0.537	38	0	0	74	168	10
57.	0	57.0	1.441	27.1	0	0	80	139	10
59.	1	59.0	0.398	30.1	846	23	60	189	1
51.	1	51.0	0.587	25.8	175	19	72	166	5
32.	1	32.0	0.484	30	0	0	0	100	7
31.	1	31.0	0.551	45.8	230	47	84	118	0
31.	1	31.0	0.254	29.6	0	0	74	##	7
33.	0	33.0	0.183	43.3	83	38	30	103	1
32.	1	32.0	0.529	34.6	96	30	70	115	1

only showing top 20 rows

df.filter("Age>50").show()

Pregnancies G	lucose Blo	odPressure S	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Out
++- 1	89	 NULL	23	+ l 94	NULL	+ NULL	++ NaN	,
j 2	197	70	45		30.5		53.0	
8	125	96	0	0	0	0.232	54.0	
10	139	80	0	0	27.1	1.441	57.0	
1	189	60	23	846	30.1	0.398	59.0	
5	166	72	19		25.8		51.0	
11	143	94	33		36.6			
13	145	82	19		22.2	•		
5	109	75	26					
4	111	72	47		37.1	!	56.0	
9	171	110	24		45.4	•		
8	176	90	34		33.7	•		
2	109	92	0		42.7	•		
4	134	72	0		23.8			
4	146	92	0		31.2	•		
5	132	80	0		26.8	•		
0	105	84	0		27.9	•		
3	128	78	0		21.1	•		
5 8	147 181	78 68	0 36		33.7 30.1	•		

only showing top 20 rows

df.filter("Age>50").select(['Age','Insulin','Glucose','Outcome']).show()

→	++	+	+	+
	Age Ir	nsulin Gl	.ucose 0ut	tcome
	+			+
	NaN	94	89	0
	53.0	543	197	1
	54.0	0	125	1
	57.0	0	139	0
	59.0	846	189	1
	51.0	175	166	1
	51.0	146	143	1
	57.0	110	145	0
	60.0	0	109	0
	56.0	207	111	1
	54.0	240	171	1
	58.0	300	176	1
	54.0	0	109	0
	60.0	0	134	1
	61.0	0	146	1

```
69.0
          0
                132
                         0
162.0
          0
                105
                         1
|55.0|
          0
                128
                          0
60.0
         495
                181
                         1
only showing top 20 rows
```

 $\label{linear_def} $$ df.filter(\sim(df['Age']>50) \& (df['Glucose']>175)).select(['Age','Insulin','Glucose','Outcome']).show() $$ $$ df.filter(\sim(df['Age']>50) \& (df['Glucose']>175)).select(['Age','Insulin','Glucose','Outcome']).show() $$ $$ df.filter(\sim(df['Age']>50) \& (df['Glucose']>175)).select(['Age','Insulin','Glucose','Outcome']).show() $$ $$ df.filter(\sim(df['Age']>50) \& (df['Glucose']>175)).select(['Age','Insulin','Glucose','Outcome']).show() $$ $$ df.filter(\sim(df['Age']>50) \& (df['Glucose']>175)).select(['Age','Insulin','Glucose','Outcome']).show() $$ $$ df.filter(\sim(df['Age']>50) \& (df['Glucose']>175)).select(['Age','Insulin','Glucose','Outcome']).show() $$ $$ df.filter(\sim(df['Age']>50) \& (df['Glucose']>175)).select(['Age']>50) \& (df['Glucose']>50) \& (df.filter(\sim(df['Age']>50) \& (df.fil$

```
| Age|Insulin|Glucose|Outcome|
|32.0|
           0
                          1|
41.0
           0
                 196
                          1
          70
                 180
                          0
26.0
125.0
           0
                 180
                          1
         304
41.0
                 187
                          1
143.0
          al
                 188
                          1
         130
36.0
                 179
                          1
41.0
           0
                 194
                          1
|41.0|
           0
                 184
                          1|
21.0
         478
                 177
                          1
|31.0|
         744
                 197
23.0
           0
                 179
                          1|
49.0
           0
                 184
                          1
24.0
         375 l
                 193
                          0
134.0
                          0
         130
                 191
129.0
           0
                 182
                          1
37.0
           0
                 179
                          0
41.0
           0
                 180
                          1
41.0
           al
                 178
                          1
29.0
         249
                 196
                          1|
```

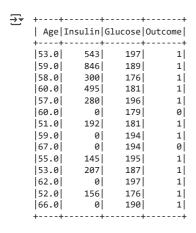
only showing top 20 rows

Start coding or generate with AI.

Start coding or generate with AI.

Start coding or generate with AI.

df.filter((df['Age']>50) & (df['Glucose']>175)).select(['Age','Insulin','Glucose','Outcome']).show()



df.groupBy(['Age']).count().show()

```
| Age|count|
|67.0|
70.0
         1|
69.0
         2
NaN
         1
149.01
         5 l
129.01
        29
64.0
         1
|47.0|
         6
42.0
        18
44.0
         8
|35.0|
        10
62.0
39.0
```

df.groupBy('BloodPressure').max().show()

→ ▼	+	+	+			++	+
ت	BloodPressure r	max(Pregnancies)	max(SkinThickness) r	max(Insulin)	max(DiabetesPedigreeFunction)	max(Age)	max(Outcome)
	++	+			·	++	+
	54	7	32	175	0.748	62.0	1
	64	8	46	415	1.731	41.0	1
	30	1	42	99	0.496	33.0	1
	85	12	54	100	1.213	56.0	1
	52	9	43	540	1.699	42.0	1
	0	13	30	0	0.933	72.0	1
	98	10	41	58	1.321	34.0	1
	110	9	46	240	0.721	54.0	1
	NULL	1	23	94	NULL	NaN	0
	96	8	0	0	0.268	54.0	1
	100	8	39	240	1.021	45.0	1
	70	15	99	744	2.329	63.0	1
	61	6	0	0	0.151	55.0	0
	75	5	32	0	0.572	60.0	1
	46	2	21	335	0.654	22.0	0
	78	14	63	293	2.42	67.0	1
	60	13	46	846	1.072	67.0	1
	90	13	51	680	0.805	66.0	1
	68	11	49	579	1.391	60.0	1
	104	5	25	0	0.435	52.0	1
	+	+	+		h	++	+

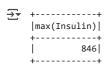
only showing top 20 rows

df.groupBy('BloodPressure').mean().show()

a	avg(Age)	avg(DiabetesPedigreeFunction)	avg(Insulin)	avg(SkinThickness)	avg(Pregnancies)	BloodPressure
0.181818	29.181818181818183	0.43290909090909097	76.81818181818181	16.363636363636363	2.727272727272727	54
0.30232	26.13953488372093	0.48055813953488374	90.13953488372093	23.348837209302324	2.4186046511627906	64
	29.5	0.3395	91.0	40.0	1.0	30
	36.833333333333336	•	22.6666666666668			85
0.27272	25.181818181818183	0.487000000000000004	91.45454545454545	16.454545454545453	2.4545454545454546	52
		0.38842857142857146		1.5142857142857142		0
0.66666	31.6666666666668	0.85066666666666	19.333333333333333	13.66666666666666	5.3333333333333333	98
0.66666	39.0	0.5733333333333334	123.33333333333333	33.66666666666664	4.333333333333333	110
	NaN	•	94.0	23.0	1.0	NULL
	•	•	0.0	0.0	4.0	96
:				36.6666666666664		100
:	31.982456140350877	1		24.05263157894737	4.087719298245614	70
:		0.151	0.0	0.0	6.0	61
:		0.382625			2.625	75
!		,	209.0	20.0	1.5	46
	•	•	•	23.733333333333334		78
	29.135135135135137	•	•	20.18918918918919		
	•	•	•	23.545454545454547		90
:	•	•		23.133333333333333		68
	46.5	0.293	0.0	12.5	2.5	104

only showing top 20 rows

df.agg({"Insulin":'max'}).show()



Start coding or generate with AI.