

Command No.	Description	PySpark Command	Pandas Command
1	Load a CSV file	<code>df = spark.read.csv('file.csv')</code>	<code>df = pd.read_csv('file.csv')</code>
2	Display the first 5 rows	<code>df.show(5)</code>	<code>df.head(5)</code>
3	Print the schema	<code>df.printSchema()</code>	<code>df.dtypes</code>
4	Select a column	<code>df.select('column1')</code>	<code>df['column1']</code>
5	Filter rows	<code>df.filter(df['column1'] > 50)</code>	<code>df[df['column1'] > 50]</code>
6	Group by a column and count	<code>df.groupBy('column1').count()</code>	<code>df['column1'].value_counts()</code>
7	Sort by a column in descending order	<code>df.orderBy(df['column1'].desc())</code>	<code>df.sort_values('column1', ascending=False)</code>
8	Add a new column	<code>df.withColumn('new_column', df['column1'] * 2)</code>	<code>df['new_column'] = df['column1'] * 2</code>
9	Drop a column	<code>df.drop('column1')</code>	<code>df.drop('column1', axis=1)</code>
10	Count the number of rows	<code>df.count()</code>	<code>df.shape[0]</code>
11	Count distinct rows	<code>df.distinct().count()</code>	<code>df.nunique()</code>
12	Generate descriptive statistics	<code>df.describe().show()</code>	<code>df.describe()</code>
13	Fill NA/Null values	<code>df.fillna(value)</code>	<code>df.fillna(value)</code>
14	Apply a function to each row	<code>df.rdd.map(lambda x: (x,)).toDF()</code>	<code>df.apply(lambda x: function(x))</code>
15	Join two dataframes	<code>df.join(other_df, df['column1'] == other_df['column2'])</code>	<code>df.merge(other_df, left_on='column1', right_on='column2')</code>
16	Aggregate data	<code>df.agg({'column1': 'sum'})</code>	<code>df['column1'].sum()</code>
17	Rename a column	<code>df.withColumnRenamed('old_name', 'new_name')</code>	<code>df.rename(columns={'old_name': 'new_name'})</code>
18	Create a temporary view	<code>df.createOrReplaceTempView('table')</code>	<code>df.to_sql('table', con)</code>
19	Execute SQL query	<code>spark.sql('SELECT * FROM table')</code>	<code>pd.read_sql_query('SELECT * FROM table', con)</code>
20	Get number of partitions (PySpark) or shape of dataframe (Pandas)	<code>df.rdd.getNumPartitions()</code>	<code>df.shape</code>
21	Repartition dataframe (only in PySpark)	<code>df.repartition(5)</code>	N/A
22	Write dataframe to CSV	<code>df.write.csv('path')</code>	<code>df.to_csv('path')</code>
23	Write dataframe to Parquet	<code>df.write.parquet('path')</code>	<code>df.to_parquet('path')</code>
24	Write dataframe to JSON	<code>df.write.json('path')</code>	<code>df.to_json('path')</code>
25	Save dataframe as a table	<code>df.write.saveAsTable('table')</code>	<code>df.to_sql('table', con)</code>

Command No.	Description	PySpark Command	Pandas Command
26	Drop duplicates	<code>df.dropDuplicates()</code>	<code>df.drop_duplicates()</code>
27	Get column names	<code>df.columns</code>	<code>df.columns</code>
28	Set a column as index	N/A	<code>df.set_index('column1')</code>
29	Reset index	N/A	<code>df.reset_index()</code>
30	Get column data type	<code>df.dtypes</code>	<code>df.dtypes</code>
31	Change column data type	<code>df.withColumn("column1", df["column1"].cast(IntegerType()))</code>	<code>df['column1'].astype('int')</code>
32	Count nulls in column	<code>df.filter(df['column1'].isNull()).count()</code>	<code>df['column1'].isnull().sum()</code>
33	Replace nulls in column	<code>df.na.fill({'column1': 'value'})</code>	<code>df['column1'].fillna('value')</code>
34	Apply a function to a column	<code>df.withColumn('column1', func(df['column1']))</code>	<code>df['column1'].apply(func)</code>
35	Concatenate two columns	<code>df.withColumn('new_column', F.concat(df['column1'], df['column2']))</code>	<code>df['new_column'] = df['column1'] + df['column2']</code>
36	Split a column	<code>df.withColumn('split_column', F.split(df['column1'], 'delimiter'))</code>	<code>df['column1'].str.split('delimiter')</code>
37	Extract year from date	<code>df.withColumn('year', F.year(df['date_column']))</code>	<code>df['date_column'].dt.year</code>
38	Extract month from date	<code>df.withColumn('month', F.month(df['date_column']))</code>	<code>df['date_column'].dt.month</code>
39	Extract day from date	<code>df.withColumn('day', F.dayofmonth(df['date_column']))</code>	<code>df['date_column'].dt.day</code>
40	Get the current date	<code>df.withColumn('current_date', F.current_date())</code>	<code>pd.to_datetime('today')</code>
41	Add days to date	<code>df.withColumn('new_date', F.date_add(df['date_column'], 5))</code>	<code>df['date_column'] + pd.DateOffset(days=5)</code>
42	Subtract days from date	<code>df.withColumn('new_date', F.date_sub(df['date_column'], 5))</code>	<code>df['date_column'] - pd.DateOffset(days=5)</code>
43	Get the difference between two dates	<code>df.withColumn('date_diff', F.datediff(df['date1'], df['date2']))</code>	<code>(df['date1'] - df['date2']).dt.days</code>
44	Extract hour from timestamp	<code>df.withColumn('hour', F.hour(df['timestamp_column']))</code>	<code>df['timestamp_column'].dt.hour</code>