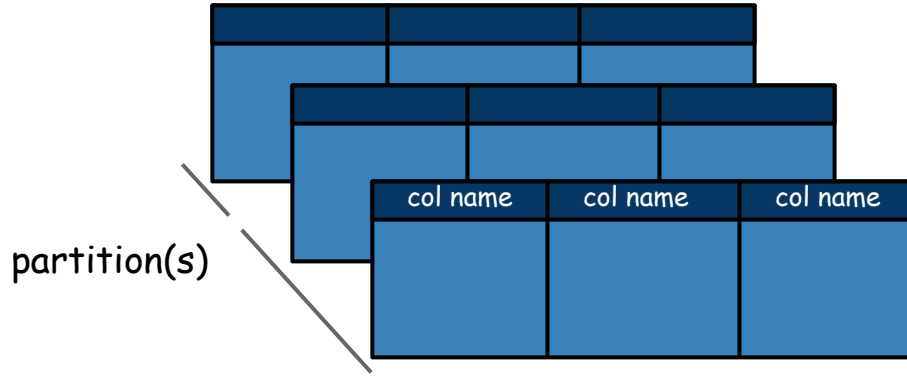# pyspark-pictures data frames
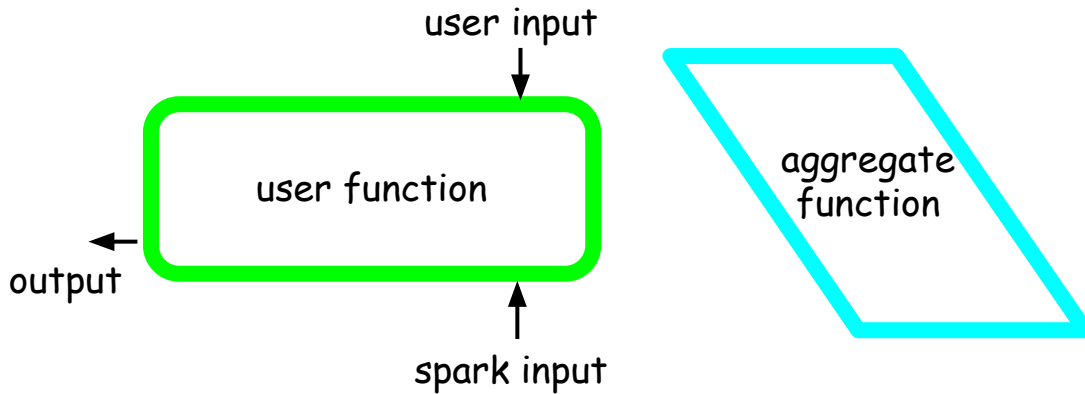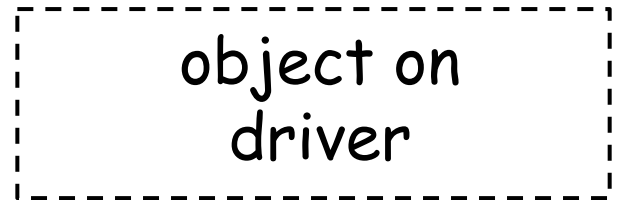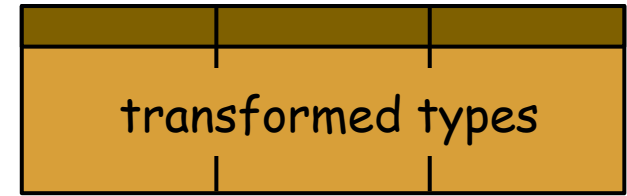
Learn the pyspark API through pictures and simple examples

https://github.com/jkthompson/pyspark-pictures

# data frame



partition(s)

col name  col name  col name

user input

user function

output

spark input

aggregate function

# data frame Row



original

transformed values

transformed types

object on driver

- 

df.B

# agg

exprs = {"B" : "aggF"}

# alias

alias = "foo"

# coalesce

numPartitions = 1

# collect

# columns

# corr

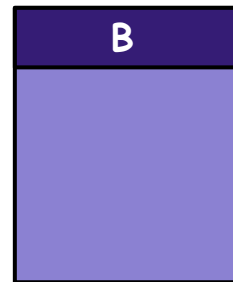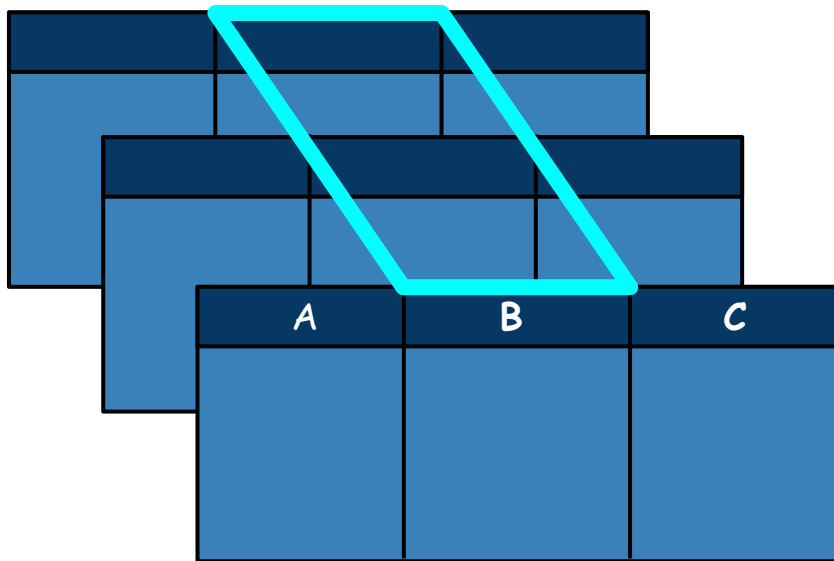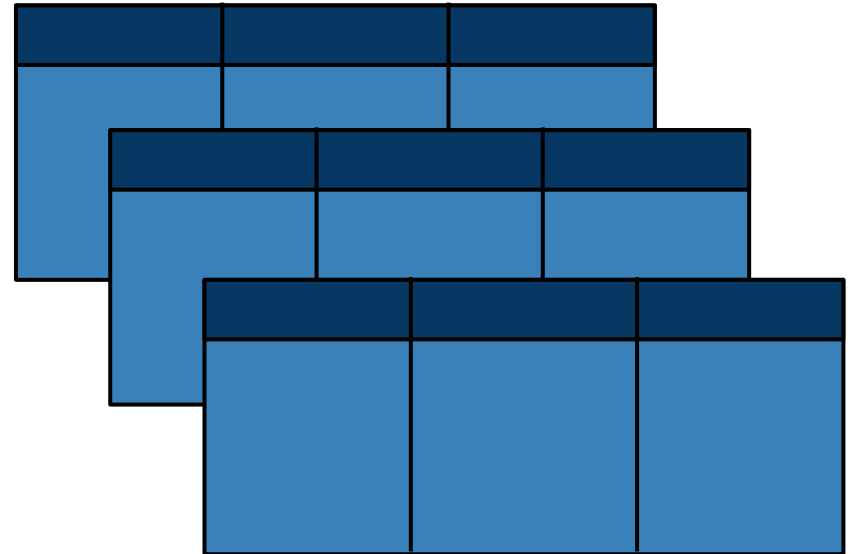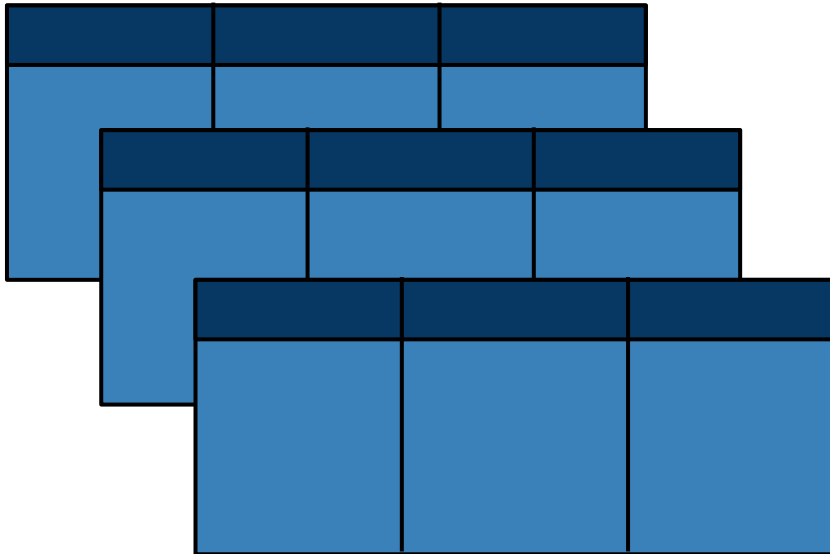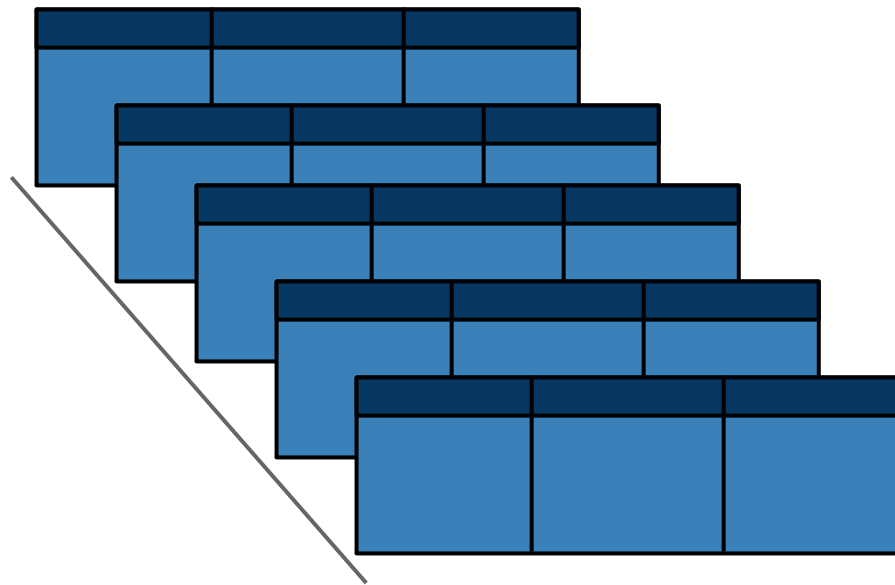col1 = A   col2 = C



Pearson's r

$$r = \frac{\sum_i (A_i - \bar{A})(C_i - \bar{C})}{\sqrt{\sum_i (A_i - \bar{A})^2}\sqrt{\sum_i (C_i - \bar{C})^2}}$$

# count



4

# cov

col1 = A    col2 = C



Sample Covariance

$$\frac{1}{N-1}\sum_i (A_i - \bar{A})(C_i - \bar{C})$$

# crosstab

# cube

col1 = A   col2 = C



| A | B | C |
|---|---|---|
| 1 | 10 | 100 |
| 2 | 20 | 200 |

| A | C | agg(A) | agg(B) | agg(C) |
|---|---|--------|--------|--------|
| null | 100 | | | |
| null | null | | | |
| 1 | null | | | |
| null | 200 | | | |
| 2 | null | | | |
| 2 | 200 | | | |
| 1 | 100 | | | |

# describe

cols = ['A','B']

# distinct

# drop

col = 'A'

# dropDuplicates

subset = ['B']

# drop_duplicates

subset = ['B']



| A | B | C |
|---|---|---|
| foo | foo | foo |
| bar | foo | baz |
| bar | foo | baz |

| A | B | C |
|---|---|---|
| foo | foo | foo |
| bar | foo | baz |

# dropna

how = 'any'  subset = ['A', 'B']

| A | B | C |
|---|---|---|
| foo | null | foo |

| A | B | C |
|---|---|---|
| null | foo | baz |

| A | B | C |
|---|---|---|
| bar | foo | null |

| A | B | C |
|---|---|---|
| bar | foo | null |

# dtypes

| A | B | C |
|---|---|---|
| foo | 3 | 0.2 |

| A | B | C |
|---|---|---|
| bar | 2 | 1.0 |

| A | B | C |
|---|---|---|
| bar | 1 | 2.2 |

[('A','string'), ('B', 'int'), ('C', 'float')]

# explain

extended = True

| A | B | C |
|---|---|---|
| foo | 3 | 0.2 |
| bar | 2 | 1.0 |
| bar | 1 | 2.2 |

== Parsed Logical Plan ==

...

== Analyzed Logical Plan ==

...

== Optimized Logical Plan ==

...

== Physical Plan ==

...

== RDD ==

# fillna

value = 'unknown"  subset = ['A', 'B']

| A | B | C |
|---|---|---|
| foo | null | foo |
| null | foo | baz |
| bar | foo | null |

| A | B | C |
|---|---|---|
| foo | unknown | foo |
| unknown | foo | baz |
| bar | foo | null |

# filter

condition = "A = foo"

| A | B | C |
|---|---|---|
| foo | foo | foo |

| A | B | C |
|---|---|---|
| bar | foo | baz |

| A | B | C |
|---|---|---|
| bar | foo | baz |

| A | B | C |
|---|---|---|
| foo | foo | foo |

# first

| A | B | C |
|---|---|---|
| bar | foo | baz |

Row(*A*='bar', *B*='foo', *C*='baz')

flatMap

# foreach



side effects
(e.g print)

*no return value,
original DataFrame
unchanged

# foreachPartition



side effects
(e.g print)

*no return value,
original DataFrame
unchanged

# freqItems

cols = ['A','C']  support = 0.5

| | | |
|---|---|---|
| foo | foo | foo |

| | | |
|---|---|---|
| bar | foo | baz |

| A | B | C |
|---|---|---|
| bar | foo | baz |

| A_freqItems | C_freqItems |
|---|---|
| [bar,foo] | [baz,foo] |

# groupBy

['A', 'C']

| A | B | C |
|---|---|---|
| 1 | 10 | 100 |
| 2 | 20 | 200 |

| A | C | agg(A) | agg(B) | agg(C) |
|---|---|--------|--------|--------|
| 1 | 100 | | | |
| 2 | 200 | | | |

# groupby

['A', 'C']

# head

n = 2

| A | B | C |
|---|---|---|
| bar | foo | baz |

[Row(A='bar', B='foo', C='baz'),
Row(A='bar', B='bar', C='baz')]

# intersect

# isLocal

| | | |
|---|---|---|
| foo | foo | foo |

| | | |
|---|---|---|
| bar | bar | baz |

| A | B | C |
|---|---|---|
| bar | foo | baz |

True

# join

joinExprs = 'A', joinType = 'inner'

# limit

num = 2



| | | |
|---|---|---|
| foo | foo | foo |

| | | |
|---|---|---|
| bar | foo | baz |

| A | B | C |
|---|---|---|
| bar | foo | baz |

| | | |
|---|---|---|
| bar | foo | baz |

| A | B | C |
|---|---|---|
| bar | foo | baz |

# map

# mapPartitions

# na

| | | |
|---|---|---|
| foo | null | foo |

| | | |
|---|---|---|
| null | foo | baz |

| A | B | C |
|---|---|---|
| bar | foo | null |

# orderBy

cols = ['A', 'C'],  ascending = [True, False]

| A | B | C |
|---|---|---|
| baz | baz | 3 |

| A | B | C |
|---|---|---|
| bar | foo | 1 |

| A | B | C |
|---|---|---|
| baz | foo | 2 |

| A | B | C |
|---|---|---|
| baz | foo | 2 |

| A | B | C |
|---|---|---|
| baz | baz | 3 |

| A | B | C |
|---|---|---|
| bar | foo | 1 |

# persist

| A | B | C |
|---|---|---|
| foo | null | foo |
| null | foo | baz |
| bar | foo | null |

# printSchema

| A | B | C |
|---|---|---|
| foo | null | foo |

| A | B | C |
|---|---|---|
| null | foo | baz |

| A | B | C |
|---|---|---|
| bar | foo | null |

stdout

```
root
 |-- A: string (nullable = true)
 |-- B: string (nullable = true)
 |-- C: string (nullable = true)
```

# randomSplit

weights = [0.8,0.2]

rdd

# registerTempTable

name = "myTable"

# repartition

numPartitions = 3

# replace

to_replace = ['foo', 'null']  value = ['foo-bar', 'unknown']

| foo | null | foo |
|-----|------|-----|

| null | foo | baz |
|------|-----|-----|

| A | B | C |
|---|---|---|
| bar | foo | null |

| foo-bar | unknown | foo-bar |
|---------|---------|---------|

| unknown | foo-bar | baz |
|---------|---------|-----|

| A | B | C |
|---|---|---|
| bar | foo-bar | unknown |

# rollup

cols = [A,C]



| | | |
|---|---|---|
| | 1 | 10 | 100 |
| A | B | C |
| 2 | 20 | 200 |

| A | C | agg(A) | agg(B) | agg(C) |
|---|---|---|---|---|
| null | null | | | |
| 1 | null | | | |
| 2 | null | | | |
| 2 | 200 | | | |
| 1 | 100 | | | |

sample

# schema

| | | |
|---|---|---|
| foo | null | foo |

| | | |
|---|---|---|
| null | foo | baz |

| A | B | C |
|---|---|---|
| bar | foo | null |

| A | B | C |
|---|---|---|

# select

cols = ['B', 'C']

| A | B | C |
|---|---|---|
| foo | foo | 1 |
| bar | foo | 2 |
| bar | foo | 3 |

| B | C |
|---|---|
| foo | 1 |
| foo | 2 |
| foo | 3 |

# selectExpr

expr = ["substr(A,1,1)", "C + 10"]

| A | B | C |
|---|---|---|
| foo | foo | 1 |
| bar | foo | 2 |
| bar | foo | 3 |

| A | C |
|---|---|
| f | 11 |
| b | 12 |
| b | 13 |

# show



| A | B | C |
|---|---|---|
| foo | 3 | 0.2 |
| bar | 2 | 1.0 |
| bar | 1 | 2.2 |

stdout

```
+---+-+---+
| A|B| C|
+---+-+---+
|foo|3|0.2|
|bar|2|1.0|
|bar|1|2.2|
+---+-+---+
```

# sort

cols = ['A', 'C'],  ascending = [True, False]



| | | |
|---|---|---|
| baz | baz | 3 |

| | | |
|---|---|---|
| bar | foo | 1 |

| A | B | C |
|---|---|---|
| baz | foo | 2 |

| | | |
|---|---|---|
| baz | foo | 2 |

| | | |
|---|---|---|
| baz | baz | 3 |

| A | B | C |
|---|---|---|
| bar | foo | 1 |

# stat

| | | |
|---|---|---|
| foo | null | foo |

| | | |
|---|---|---|
| null | foo | baz |

| A | B | C |
|---|---|---|
| bar | foo | null |

# subtract

# take

n = 2

| | | |
|---|---|---|
| foo | foo | foo |

| | | |
|---|---|---|
| bar | bar | baz |

| A | B | C |
|---|---|---|
| bar | foo | baz |

[Row(A='bar', B='foo', C='baz'),
Row(A='bar', B='bar', C='baz')]

# toJSON

| A | B | C |
|---|---|---|
| foo | 3 | 0.2 |
| bar | 2 | 1.0 |
| bar | 1 | 2.2 |

u'{"A":"foo","B":3,"C":0.2}'

u'{"A":"bar","B":2,"C":1.0}'

u'{"A":"bar","B":1,"C":2.2}'

# toPandas

| | | |
|---|---|---|
| foo | 3 | 0.2 |

| | | |
|---|---|---|
| bar | 2 | 1.0 |

| A | B | C |
|---|---|---|
| bar | 1 | 2.2 |

```
    A   B   C
0  foo  3  0.2
1  bar  2  1.0
2  bar  1  2.2
```

unionAll

# unpersist

| | | |
|---|---|---|
| foo | null | foo |

| | | |
|---|---|---|
| null | foo | baz |

| A | B | C |
|---|---|---|
| bar | foo | null |

# where

condition = "A = foo"

| | | |
|---|---|---|
| foo | foo | foo |

| | | |
|---|---|---|
| bar | foo | baz |

| A | B | C |
|---|---|---|
| bar | foo | baz |

| A | B | C |
|---|---|---|
| foo | foo | foo |

# withColumn

colName = 'D'

# withColumnRenamed

existing = 'C'  col = 'D'

# write

| | | |
|---|---|---|
| foo | null | foo |

| | | |
|---|---|---|
| null | foo | baz |

| A | B | C |
|---|---|---|
| bar | foo | null |