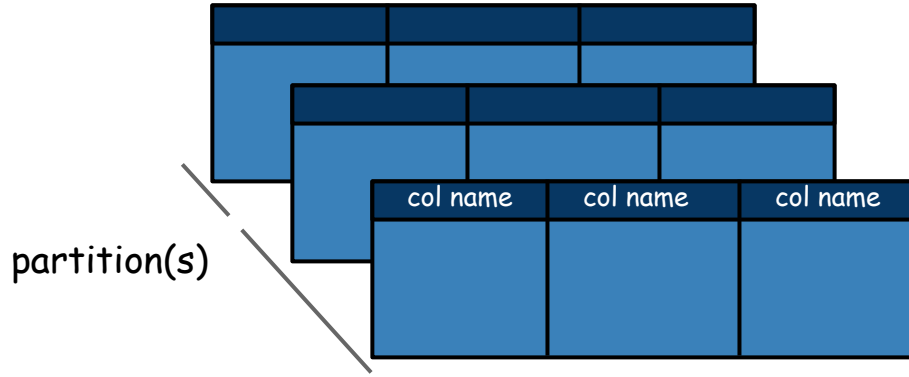# pyspark-pictures data frames
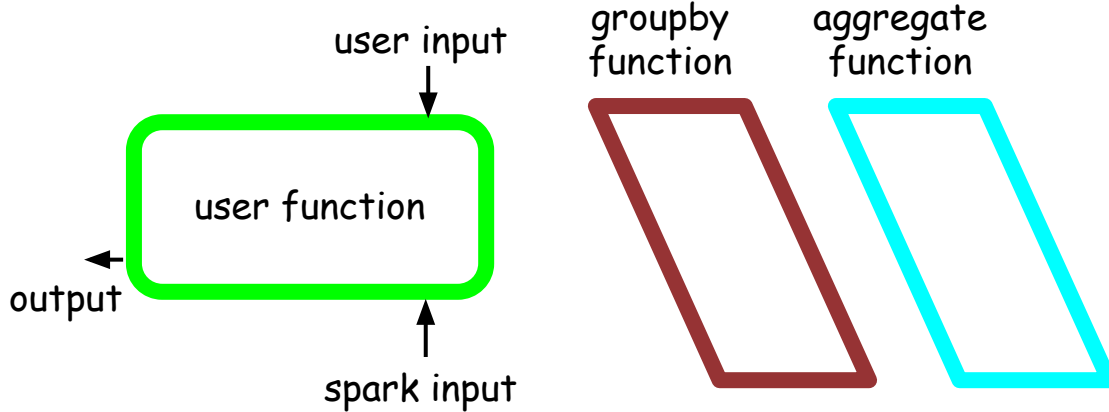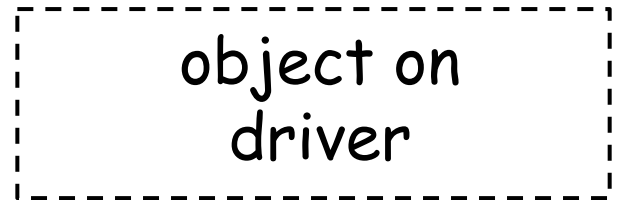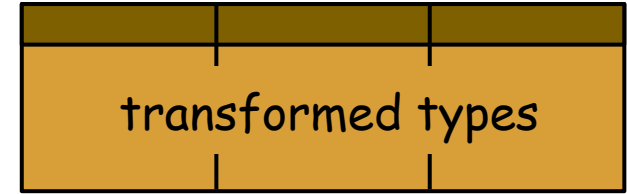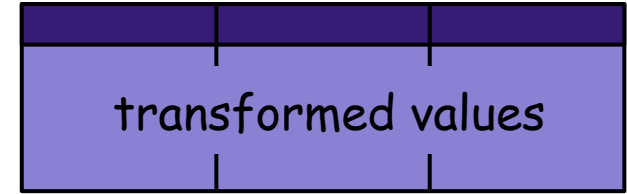
Learn the pyspark API through pictures and simple examples

# data frame



partition(s)

user input

user function

output

spark input

groupby function

aggregate function

# data frame row



original

transformed values

transformed types

object on driver

col name col name col name

agg

# alias

cache

# coalesce

numPartitions = 1

# collect

# columns

# corr



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



3

# cov



Sample Covariance

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

# crosstab

col1 = 'from'   col2 = 'to'

| from | to | amt |
|------|-----|-----|
| Carol | Dave | 0.3 |
| Bob | Carol | 0.2 |
| Alice | Bob | 0.1 |

| from_to | Bob | Carol | Dave |
|---------|-----|-------|------|
| Carol | 0 | 0 | 1 |
| Bob | 0 | 1 | 0 |
| Alice | 1 | 0 | 0 |

# cube

*cols = 'from', 'to'

| | | |
|---|---|---|
| Alice | Carol | 0.2 |

| from | **to** | **amt** |
|---|---|---|
| Alice | Bob | 0.1 |

| | | |
|---|---|---|
| null | null | |

| | | |
|---|---|---|
| Alice | null | |

| | | |
|---|---|---|
| null | Carol | |

| | | |
|---|---|---|
| null | Bob | |

| | | |
|---|---|---|
| Alice | Carol | |

| from | **to** | **agg(amt)** |
|---|---|---|
| Alice | Bob | |

# describe



| from | to | amt |
|------|-----|-----|
| Alice | Bob | 0.1 |
| Bob | Carol | 0.2 |
| Carol | Dave | 0.3 |

| summary | amt |
|---------|-----|
| count | |
| mean | |
| stdev | |
| min | |
| max | |

# distinct

| from | to | amt |
|---|---|---|
| Bob | Carol | 0.2 |
| Carol | Dave | 0.3 |
| Bob | Carol | 0.2 |
| Alice | Bob | 0.1 |

| from | to | amt |
|---|---|---|
| Carol | Dave | 0.3 |
| Bob | Carol | 0.2 |
| Alice | Bob | 0.1 |

# drop

col = 'amt'

| from | to | amt |
|---|---|---|
| Carol | Dave | 0.3 |
| Bob | Carol | 0.2 |
| Alice | Bob | 0.1 |

| from | to |
|---|---|
| Carol | Dave |
| Bob | Carol |
| Alice | Bob |

# dropDuplicates

subset = ['from','to']

| | | |
|---|---|---|
| Bob | Carol | 0.2 |

| | | |
|---|---|---|
| Bob | Carol | 0.3 |

| | | |
|---|---|---|
| Bob | Carol | 0.2 |

| from | to | amt |
|---|---|---|
| Alice | Bob | 0.1 |

| | | |
|---|---|---|
| Bob | Carol | 0.2 |

| from | to | amt |
|---|---|---|
| Alice | Bob | 0.1 |

# dropna

how = 'any'  subset = ['from', 'to']

| | | |
|---|---|---|
| Bob | Carol | 0.2 |

| | | |
|---|---|---|
| Carol | null | 0.3 |

| | | |
|---|---|---|
| Bob | Carol | Na |

| from | to | amt |
|---|---|---|
| null | Bob | 0.1 |

| | | |
|---|---|---|
| Bob | Carol | 0.2 |

| from | to | amt |
|---|---|---|
| Bob | Carol | Na |

# dtypes

| from | to | amt |
|------|------|-----|
| Carol | Dave | 0.3 |
| Bob | Carol | 0.2 |
| Alice | Bob | 0.1 |

[('from','string'), ('to', 'string'), ('amt', 'double')]

# explain

extended = True

| from | to | amt |
|------|------|-----|
| Carol | Dave | 0.3 |
| Bob | Carol | 0.2 |
| Alice | Bob | 0.1 |

== Parsed Logical Plan ==

...

== Analyzed Logical Plan ==

...

== Optimized Logical Plan ==

...

== Physical Plan ==

...

== RDD ==

# fillna

value = 'unknown''  subset = ['from', 'to']

| Carol | null | 0.3 |
|---|---|---|

| | Bob | Carol | nan |
|---|---|---|---|

| from | **to** | **amt** |
|---|---|---|
| null | Bob | 0.1 |

| Carol | unknown | 0.3 |
|---|---|---|

| | Bob | Carol | nan |
|---|---|---|---|

| from | **to** | **amt** |
|---|---|---|
| unknown | Bob | 0.1 |

# filter

condition = "amt > 0.1"

| | | |
|---|---|---|
| Carol | Dave | 0.3 |

| | | |
|---|---|---|
| Bob | Carol | 0.2 |

| from | **to** | **amt** |
|---|---|---|
| Alice | Bob | 0.1 |

| | | |
|---|---|---|
| Carol | Dave | 0.3 |

| **from** | **to** | **amt** |
|---|---|---|
| Bob | Carol | 0.2 |

# first

| from | to | amt |
|------|-----|-----|
| Carol | Dave | 0.3 |
| Bob | Carol | 0.2 |
| Alice | Bob | 0.1 |

Row(from='Alice', to='Bob', amt=0.1)

# flatMap

# foreach

| from | to | amt |
|------|-----|-----|
| Carol | Dave | 0.3 |
| Bob | Carol | 0.2 |
| Alice | Bob | 0.1 |

side effects
(e.g print) ←

*no return value,
original DataFrame
unchanged

# foreachPartition



| from | to | amt |
|------|-----|-----|
| Carol | Dave | 0.3 |
| Bob | Carol | 0.2 |
| Alice | Bob | 0.1 |

side effects (e.g print)

*no return value, original DataFrame unchanged

# freqItems

cols = ['from', 'amt']  support = 0.8

| from | to | amt |
|------|------|------|
| Carol | Bob | 0.1 |
| Alice | Bob | 0.5 |
| Alice | Bob | 0.1 |
| Alice | Dave | 0.1 |
| Bob | Carol | 0.1 |

| from_freqItems | amt_freqItems |
|----------------|---------------|
| [Alice] | [0.1] |

# groupBy (groupby)



GroupedData Object
with methods: agg, avg, count, max, mean, min, pivot, sum

# groupBy(col1).avg(col2)

col1 = 'from'  col2 = 'amt'

| from | to | amt |
|------|------|------|
| Carol | Dave | 0.3 |
| Alice | Carol | 0.2 |
| Alice | Bob | 0.1 |

| from | avg(amt) |
|------|----------|
| Carol | 0.3 |
| Alice | 0.15 |

# head

n = 2

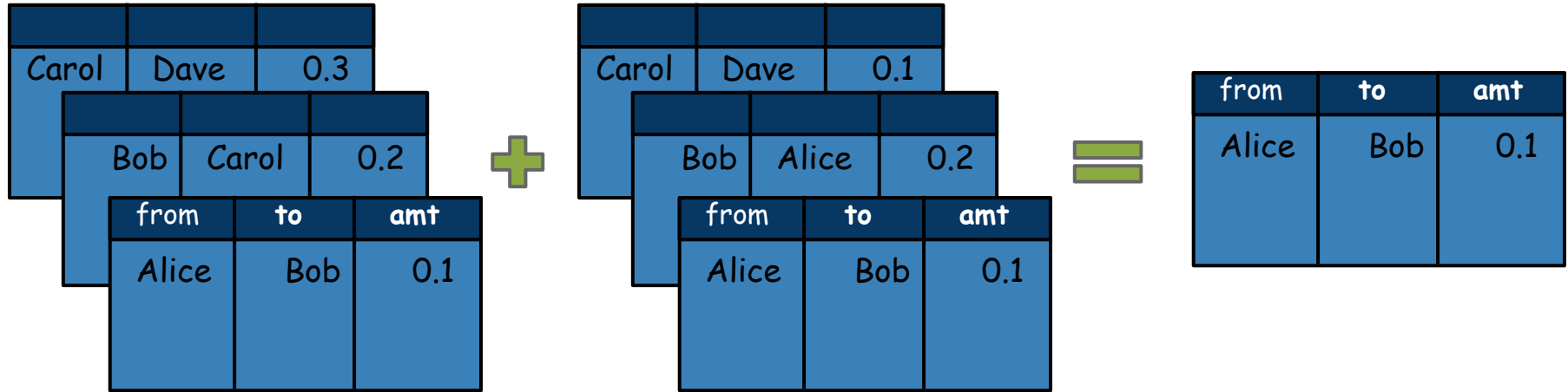| from | to | amt |
|------|-----|-----|
| Carol | Dave | 0.3 |
| Bob | Carol | 0.2 |
| Alice | Bob | 0.1 |

[Row(from=u'Alice', to=u'Bob', amt=0.1),
Row(from=u'Bob', to=u'Carol', amt=0.2)]

# intersect

| | | |
|---|---|---|
| Carol | Dave | 0.3 |

| | | |
|---|---|---|
| Bob | Carol | 0.2 |

| from | to | amt |
|---|---|---|
| Alice | Bob | 0.1 |

**+**

| | | |
|---|---|---|
| Carol | Dave | 0.1 |

| | | |
|---|---|---|
| Bob | Alice | 0.2 |

| from | to | amt |
|---|---|---|
| Alice | Bob | 0.1 |

**=**

| from | to | amt |
|---|---|---|
| Alice | Bob | 0.1 |

# isLocal

False

# join

joinExprs = x.to==y.name
joinType = 'inner'

| | | |
|---|---|---|
| Carol | Dave | 0.3 |

| | | |
|---|---|---|
| Bob | Carol | 0.2 |

| from | to | amt |
|---|---|---|
| Alice | Bob | 0.1 |

**+**

| | |
|---|---|
| Dave | 80 |

| | |
|---|---|
| Bob | 40 |

| name | age |
|---|---|
| Alice | 20 |

**=**

| | | | | |
|---|---|---|---|---|
| Alice | Bob | 0.1 | Bob | 40 |

| from | to | amt | name | age |
|---|---|---|---|---|
| Carol | Dave | 0.3 | Dave | 80 |

# limit

num = 2

map

# mapPartitions

# na

| | | |
|---|---|---|
| Bob | Carol | 0.2 |

| | | |
|---|---|---|
| Carol | null | 0.3 |

| | | |
|---|---|---|
| Bob | Carol | Na |

| from | to | amt |
|---|---|---|
| null | Bob | 0.1 |

**DataFrameNaFunctions**

| | | |
|---|---|---|
| Bob | Carol | 0.2 |

| | | |
|---|---|---|
| Carol | null | 0.3 |

| | | |
|---|---|---|
| Bob | Carol | Na |

| from | to | amt |
|---|---|---|
| null | Bob | 0.1 |

# orderBy

cols = ['from'], ascending = [False]

| from | to | amt |
|---|---|---|
| Carol | Dave | 0.3 |
| Bob | Carol | 0.2 |
| Alice | Bob | 0.1 |

| from | to | amt |
|---|---|---|
| Alice | Bob | 0.1 |
| Bob | Carol | 0.2 |
| Carol | Dave | 0.3 |

# persist

strorageLevel =
StorageLevel(MEMORY_ONLY_SER)

# printSchema

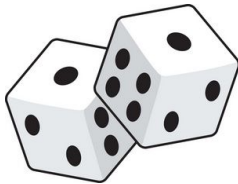| from | to | amt |
|------|-----|-----|
| Carol | Dave | 0.3 |
| Bob | Carol | 0.2 |
| Alice | Bob | 0.1 |

stdout

```
root
 |-- from: string (nullable = true)
 |-- to: string (nullable = true)
 |-- amt: double (nullable = true)
```
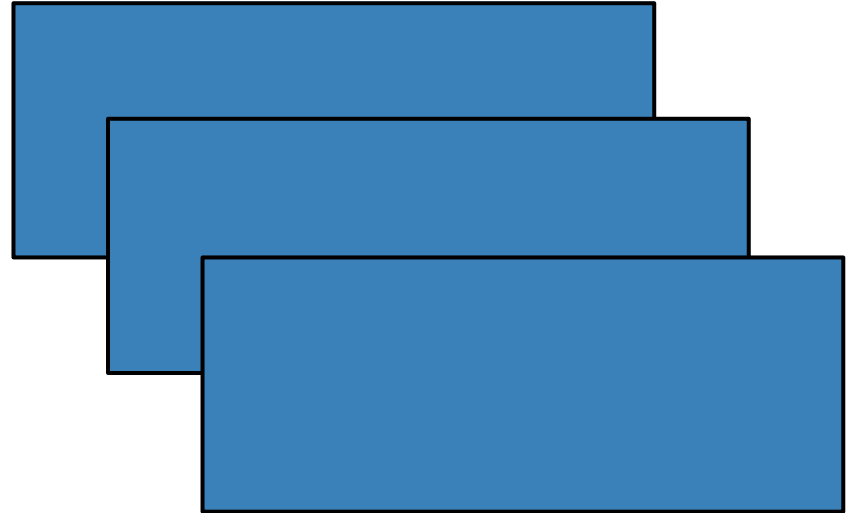
# randomSplit

weights = [0.5,0.5]

| Carol | Dave | 0.3 |
|-------|------|-----|

| Bob | Carol | 0.2 |
|-----|-------|-----|

| from | to | amt |
|-------|------|------|
| Alice | Bob | 0.1 |

| Carol | Dave | 0.3 |
|-------|------|-----|

| from | to | amt |
|------|-------|------|
| Bob | Carol | 0.2 |

| from | to | amt |
|-------|-----|------|
| Alice | Bob | 0.1 |

rdd

# registerTempTable
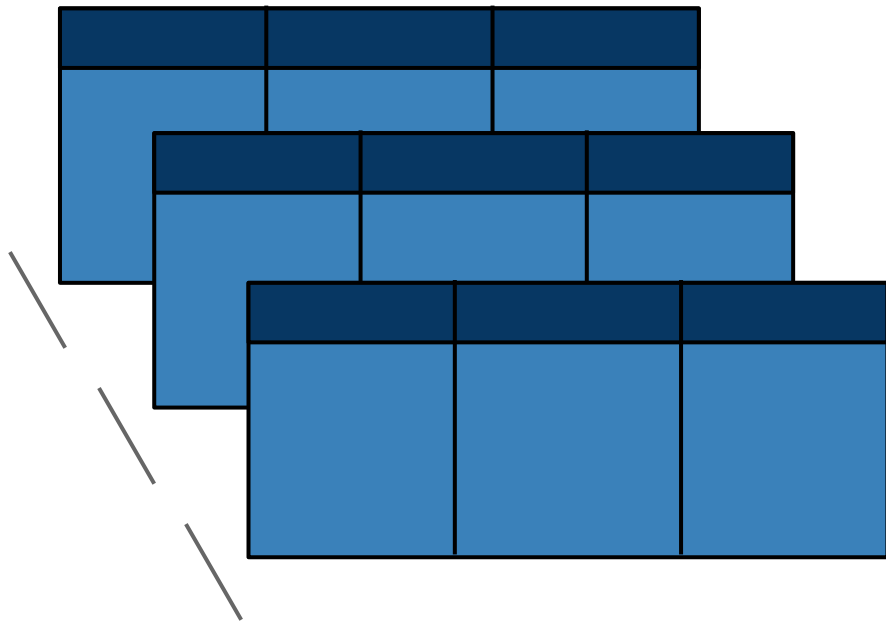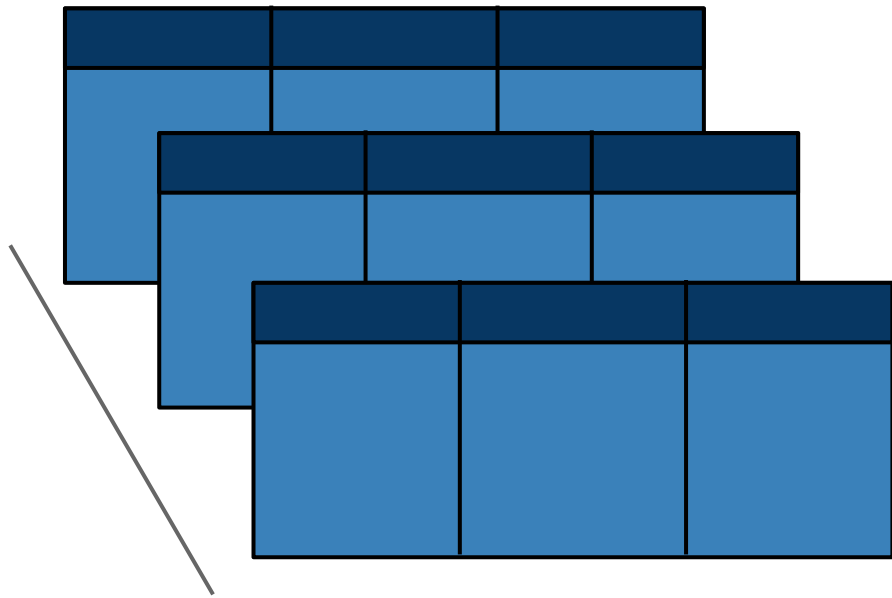
name = "TRANSACTIONS"

TRANSACTIONS

# repartition

# replace

to_replace = 'Dave'   value = 'David'

| from | to | amt |
|------|------|------|
| Carol | Dave | 0.3 |
| Bob | Carol | 0.2 |
| Alice | Bob | 0.1 |

| from | to | amt |
|------|------|------|
| Carol | David | 0.3 |
| Bob | Carol | 0.2 |
| Alice | Bob | 0.1 |

# rollup

cols = ['from','to']



| from | to | amt |
|------|------|-----|
| Carol | Dave | 0.3 |
| Bob | Carol | 0.2 |
| Alice | Bob | 0.1 |

| from | to | agg(amt) |
|------|------|----------|
| null | null | |
| Carol | null | |
| Bob | null | |
| Carol | Dave | |
| Alice | null | |
| Bob | Carol | |
| Alice | Bob | |

sample

# sampleBy

| from | to | amt |
|------|------|------|
| Bob | Carol | 0.6 |
| Bob | Bob | 0.5 |
| Alice | Dave | 0.4 |
| Alice | Alice | 0.3 |
| Alice | Carol | 0.2 |
| Alice | Bob | 0.1 |

| from | to | amt |
|------|------|------|
| Bob | Carol | 0.6 |
| Bob | Bob | 0.5 |
| Alice | Bob | 0.3 |

# schema



|  |  |  |
|---|---|---|
| from | to | amt |

# select

cols = ['from', 'amt']

| from | to | amt |
|------|-----|-----|
|      |     |     |

| from | amt |
|------|-----|
|      |     |

# selectExpr

expr = ["substr(from,1,1)", "amt + 10"]

| from | to | amt |
|------|------|------|
| Carol | Dave | 0.3 |
| Bob | Carol | 0.2 |
| Alice | Bob | 0.1 |

| from | amt |
|------|------|
| C | 10.3 |
| B | 10.2 |
| A | 10.1 |

# show

| from | to | amt |
|---|---|---|
| Alice | Bob | 0.1 |
| Bob | Carol | 0.2 |
| Carol | Dave | 0.3 |

stdout

```
+-----+-----+---+
| from|   to|amt|
+-----+-----+---+
|Alice|  Bob|0.1|
|  Bob|Carol|0.2|
|Carol| Dave|0.3|
+-----+-----+---+
```

# sort

cols = ['to']

| from | to | amt |
|------|-----|-----|
| Carol | Alice | 0.3 |
| Bob | Carol | 0.2 |
| Alice | Bob | 0.1 |

| from | to | amt |
|------|-----|-----|
| Bob | Carol | 0.2 |
| Alice | Bob | 0.1 |
| Carol | Alice | 0.3 |

# sortWithinPartitions



Left partition stack:

| from | to | amt | p_id |
|------|------|------|------|
| Carol | Alice | 0.3 | 2 |
| Bob | Carol | 0.2 | 2 |
| Alice | Bob | 0.1 | 1 |

Right partition stack:

| from | to | amt | p_id |
|------|------|------|------|
| Bob | Carol | 0.2 | 2 |
| Carol | Alice | 0.3 | 2 |
| Alice | Bob | 0.1 | 1 |

stat

stat

# subtract

| | | |
|---|---|---|
| Carol | Dave | 0.3 |

| | | |
|---|---|---|
| Bob | Carol | 0.2 |

| from | to | amt |
|---|---|---|
| Alice | Bob | 0.1 |

−

| | | |
|---|---|---|
| Carol | Dave | 0.1 |

| | | |
|---|---|---|
| Bob | Carol | 0.2 |

| from | to | amt |
|---|---|---|
| Alice | Bob | 0.1 |

=

| from | to | amt |
|---|---|---|
| Carol | Dave | 0.3 |

# take

num = 2

# toDF

cols = ["seller", "buyer"]

| from | to | amt |
|---|---|---|
| | | |

| seller | buyer | amt |
|---|---|---|
| | | |

# toJSON

cols = ['to']

# toPandas

| from | to | amt |
|------|------|------|
| Carol | Alice | 0.3 |

| | | |
|------|------|------|
| Bob | Carol | 0.2 |

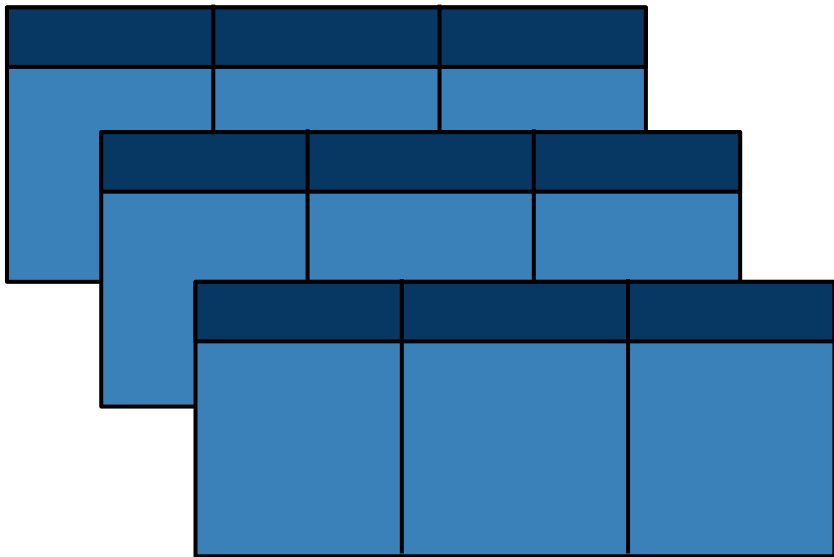| from | to | amt |
|------|------|------|
| Alice | Bob | 0.1 |

```
   from  to   amt
0  Alice Bob  0.1
1  Bob   Carol 0.2
2  Carol Alice 0.3
```

# unpersist

# where (filter)

condition = "amt > 0.1"

| | | |
|---|---|---|
| Carol | Dave | 0.3 |

| | | |
|---|---|---|
| Bob | Carol | 0.2 |

| from | to | amt |
|---|---|---|
| Alice | Bob | 0.1 |

| | | |
|---|---|---|
| Carol | Dave | 0.3 |

| from | to | amt |
|---|---|---|
| Bob | Carol | 0.2 |

# withColumn

colName = 'conf'

# withColumnRenamed

existing = 'amt'  col = 'amount'



| from | to | amt |
|------|-----|-----|
|      |     |     |

| from | to | amount |
|------|-----|--------|
|      |     |        |

write