



Apache Spark for library developers



Erik Erlandson
eje@redhat.com
[@manyaangled](https://twitter.com/manyaangled)

William Benton
willb@redhat.com
[@willb](https://twitter.com/willb)



About Will

The Silex and Isarn libraries

Reusable open-source code that works
with Spark, factored from internal apps.



We've tracked Spark releases since Spark 1.3.0.

See <https://silex.radaranalytics.io> and
<http://isarnproject.org>



**This stairway
has over
320 steps**

**Do not use
except in an
emergency**

Forecast

Basic considerations for reusable Spark code

Generic functions for parallel collections

Extending data frames with custom aggregates

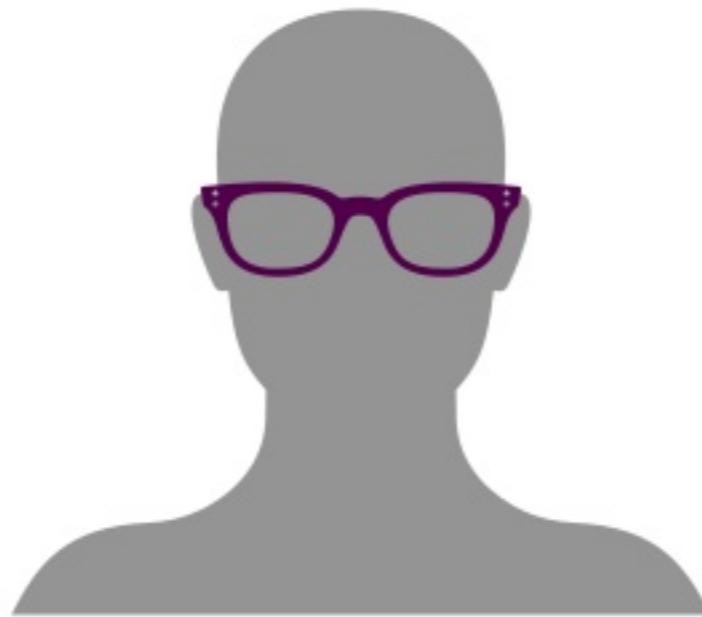
Exposing JVM libraries to Python

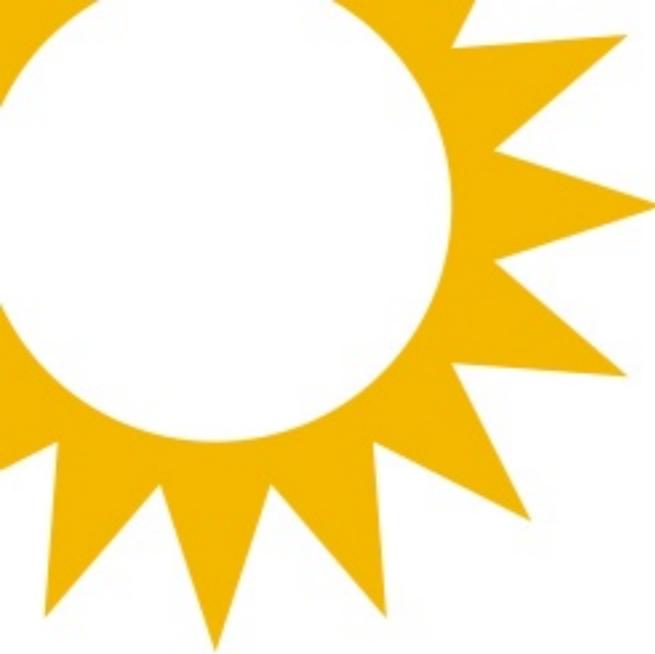
Sharing your work with the world

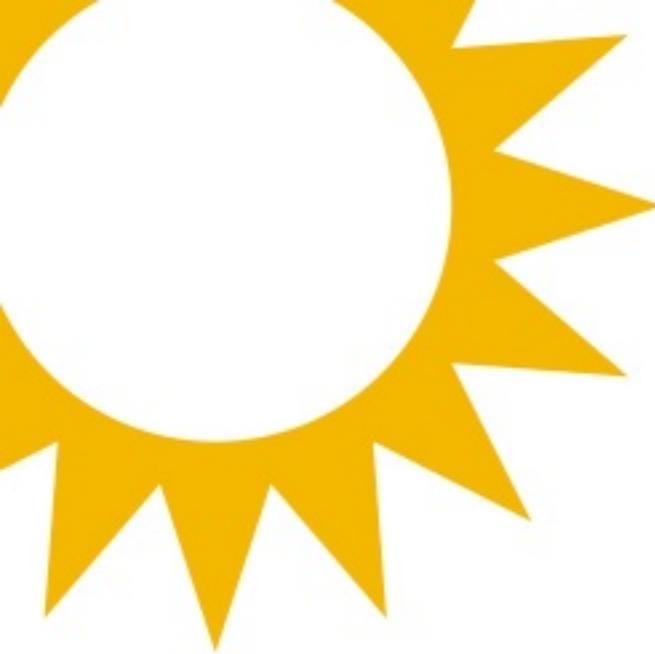


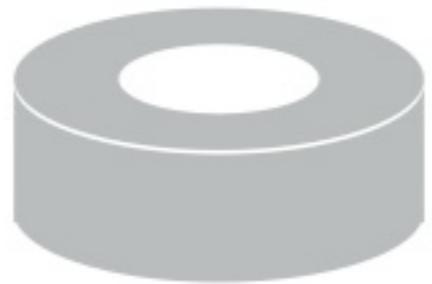
Basic considerations



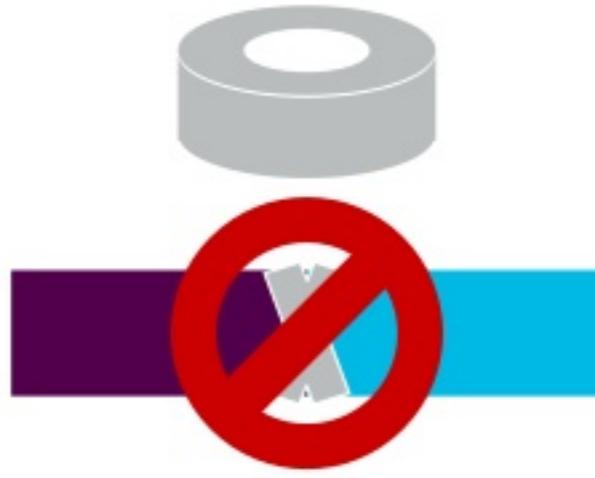












Today's main themes



Cross-building for Scala

in your SBT build definition:

```
scalaVersion := "2.11.11"
```

```
crossScalaVersions := Seq("2.10.6", "2.11.11")
```

in your shell:

```
$ sbt +compile
```

```
$ sbt "++ 2.11.11" compile
```



Cross-building for Scala

in your SBT build definition:

```
scalaVersion := "2.11.11"
```

```
crossScalaVersions := Seq("2.10.6", "2.11.11")
```

in your shell:

```
$ sbt +compile          # or test, package, publish, etc.  
$ sbt "++ 2.11.11" compile
```



Cross-building for Scala

in your SBT build definition:

```
scalaVersion := "2.11.11"
```

```
crossScalaVersions := Seq("2.10.6", "2.11.11")
```

in your shell:

```
$ sbt +compile          # or test, package, publish, etc.  
$ sbt "++ 2.11.11" compile
```



Bring-your-own Spark

in your SBT build definition:

```
libraryDependencies ++= Seq(  
    "org.apache.spark" %% "spark-core" % "2.3.0" % Provided,  
    "org.apache.spark" %% "spark-sql" % "2.3.0" % Provided,  
    "org.apache.spark" %% "spark-mllib" % "2.3.0" % Provided,  
    "org.scalatest" %% "scalatest" % "2.2.4" % Test)
```

Bring-your-own Spark

in your SBT build definition:

```
libraryDependencies ++= Seq(  
    "org.apache.spark" %% "spark-core" % "2.3.0" % Provided,  
    "org.apache.spark" %% "spark-sql" % "2.3.0" % Provided,  
    "org.apache.spark" %% "spark-mllib" % "2.3.0" % Provided,  
    "org.scalatest" %% "scalatest" % "2.2.4" % Test)
```



“Bring-your-own Spark”

in your SBT build definition:

```
libraryDependencies ++= Seq(  
    "org.apache.spark" %% "spark-core" % "2.3.0" % Provided,  
    "org.apache.spark" %% "spark-sql" % "2.3.0" % Provided,  
    "org.apache.spark" %% "spark-mllib" % "2.3.0" % Provided,  
    "joda-time" % "joda-time" % "2.7",  
    "org.scalatest" %% "scalatest" % "2.2.4" % Test)
```



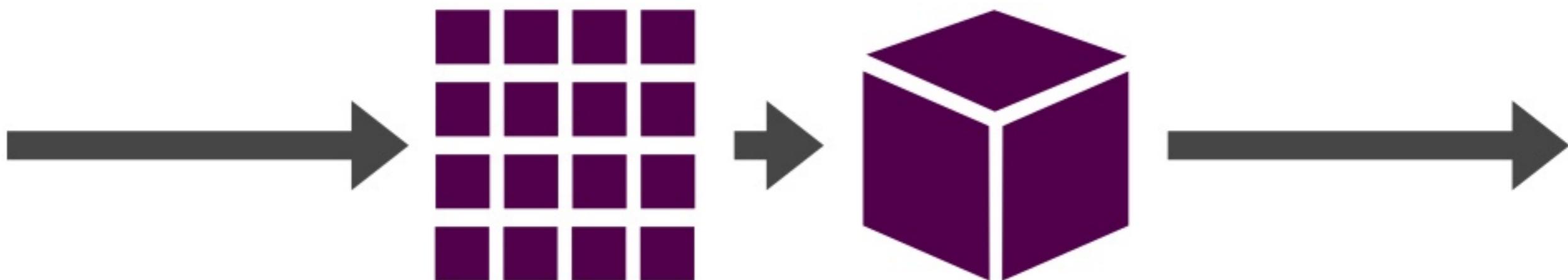
“Bring-your-own Spark”

in your SBT build definition:

```
libraryDependencies ++= Seq(  
    "org.apache.spark" %% "spark-core" % "2.3.0" % Provided,  
    "org.apache.spark" %% "spark-sql" % "2.3.0" % Provided,  
    "org.apache.spark" %% "spark-mllib" % "2.3.0" % Provided  
    "joda-time" % "joda-time" % "2.7",  
    "org.scalatest" %% "scalatest" % "2.2.4" % Test)
```



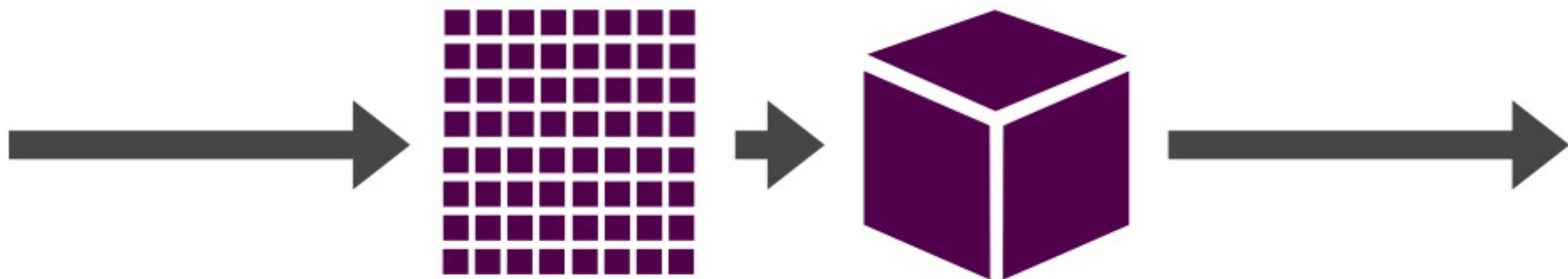
Taking care with resources



Taking care with resources



Taking care with resources



Caching when necessary

```
def step(rdd: RDD[_]) = {  
    rdd.cache()  
    result = trainModel(rdd)  
  
    result  
}
```

Caching when necessary

```
def step(rdd: RDD[_]) = {  
    rdd.cache()  
    result = trainModel(rdd)  
  
    result  
}
```

Caching when necessary

```
def step(rdd: RDD[_]) = {  
    rdd.cache()  
    result = trainModel(rdd)  
    rdd.unpersist()  
    result  
}
```

Caching when necessary

```
def step(rdd: RDD[_]) = {  
    val wasUncached = rdd.storageLevel == StorageLevel.NONE  
    if (wasUncached) { rdd.cache() }  
    result = trainModel(rdd)  
  
    result  
}
```

Caching when necessary

```
def step(rdd: RDD[_]) = {  
    val wasUncached = rdd.storageLevel == StorageLevel.NONE  
    if (wasUncached) { rdd.cache() }  
    result = trainModel(rdd)  
    if (wasUncached) { rdd.unpersist() }  
  
    result  
}
```

```
var nextModel = initialModel
for (int i = 0; i < iterations; i++) {
    val current = sc.broadcast(nextModel)
    val newState = current.mapPartitions((parts, index) => parts.map((id, state) => (id, state.map((x, y) => (x, y * current.value)))))
    current.unpersist()
    nextModel = modelFromState(newState)
    current.unpersist()
}
```

```
var nextModel = initialModel
for (int i = 0; i < iterations; i++) {
    val current = sc.broadcast(nextModel)
    val newState = mapPartitionsWithIndex((i, iter) => {
            val state = iterState[0].value
            val model = iterModel[0].value
            val newModel = state.map((x, y) => (x, y * model))
            (i, newModel)
    })
    nextModel = modelFromState(newState)
    current.unpersist
}
```

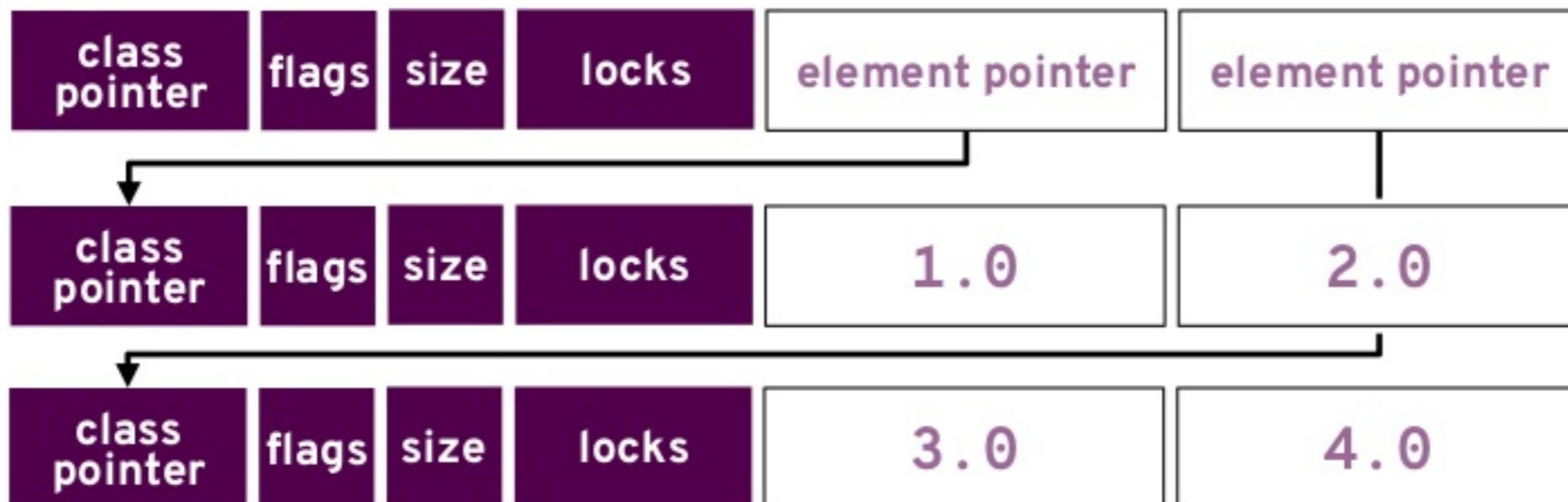
```
var nextModel = initialModel
for (int i = 0; i < iterations; i++) {
    val current = sc.broadcast(nextModel)
    val newState = current.value.map((state, example) =>
        state.map((currentValue, id, example) =>
            if (id == 0) currentValue else id * currentValue
        ).reduce(_ + _)
    ).reduce(_ + _)
    nextModel = modelFromState(newState)
    current.unpersist
}
```

Minding the JVM heap

```
val mat = Array(Array(1.0, 2.0), Array(3.0, 4.0))
```

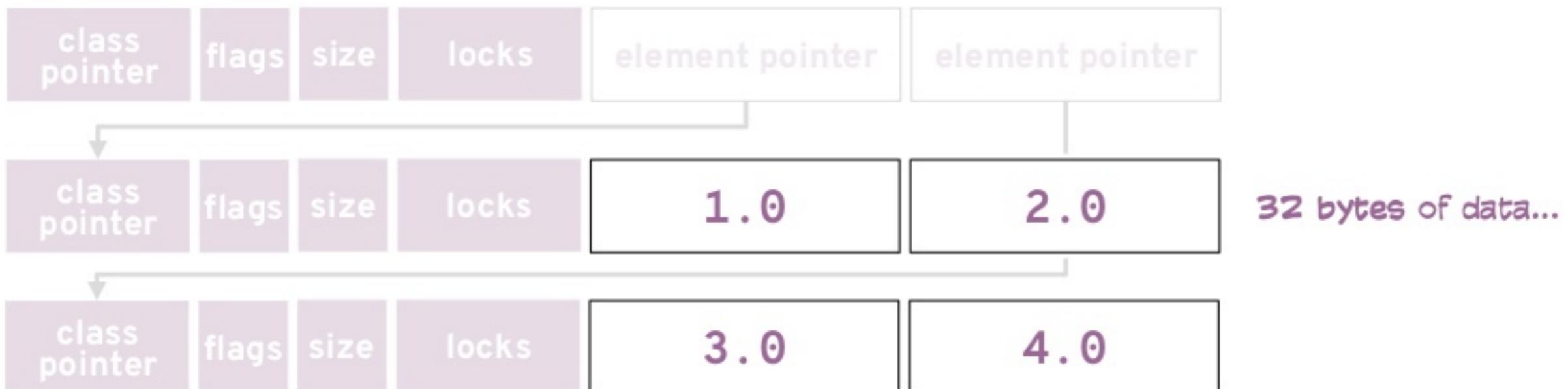
Minding the JVM heap

```
val mat = Array(Array(1.0, 2.0), Array(3.0, 4.0))
```



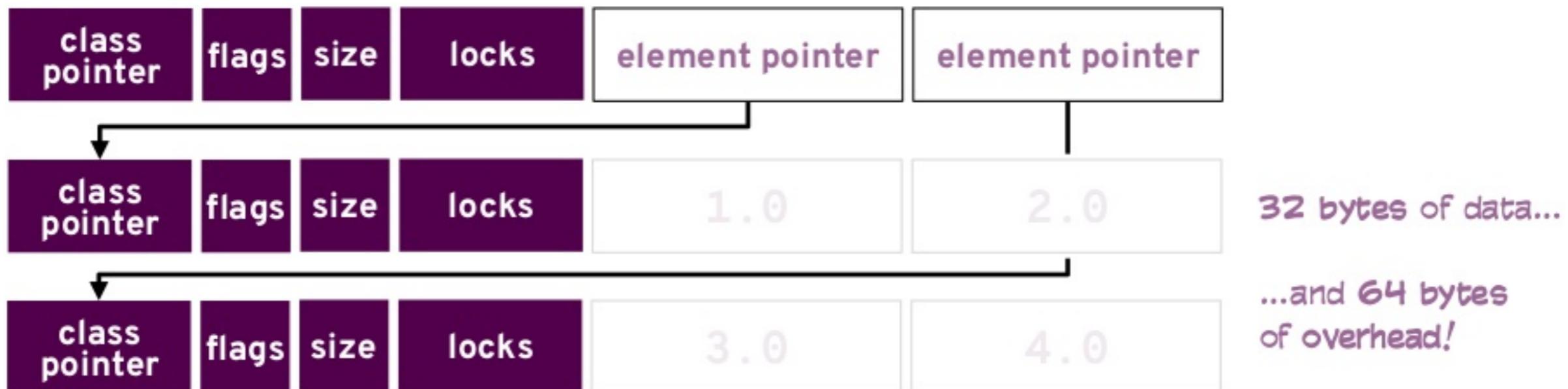
Minding the JVM heap

```
val mat = Array(Array(1.0, 2.0), Array(3.0, 4.0))
```



Minding the JVM heap

```
val mat = Array(Array(1.0, 2.0), Array(3.0, 4.0))
```





Continuous integration for Spark libraries and apps

local[*]



CPU



Memory

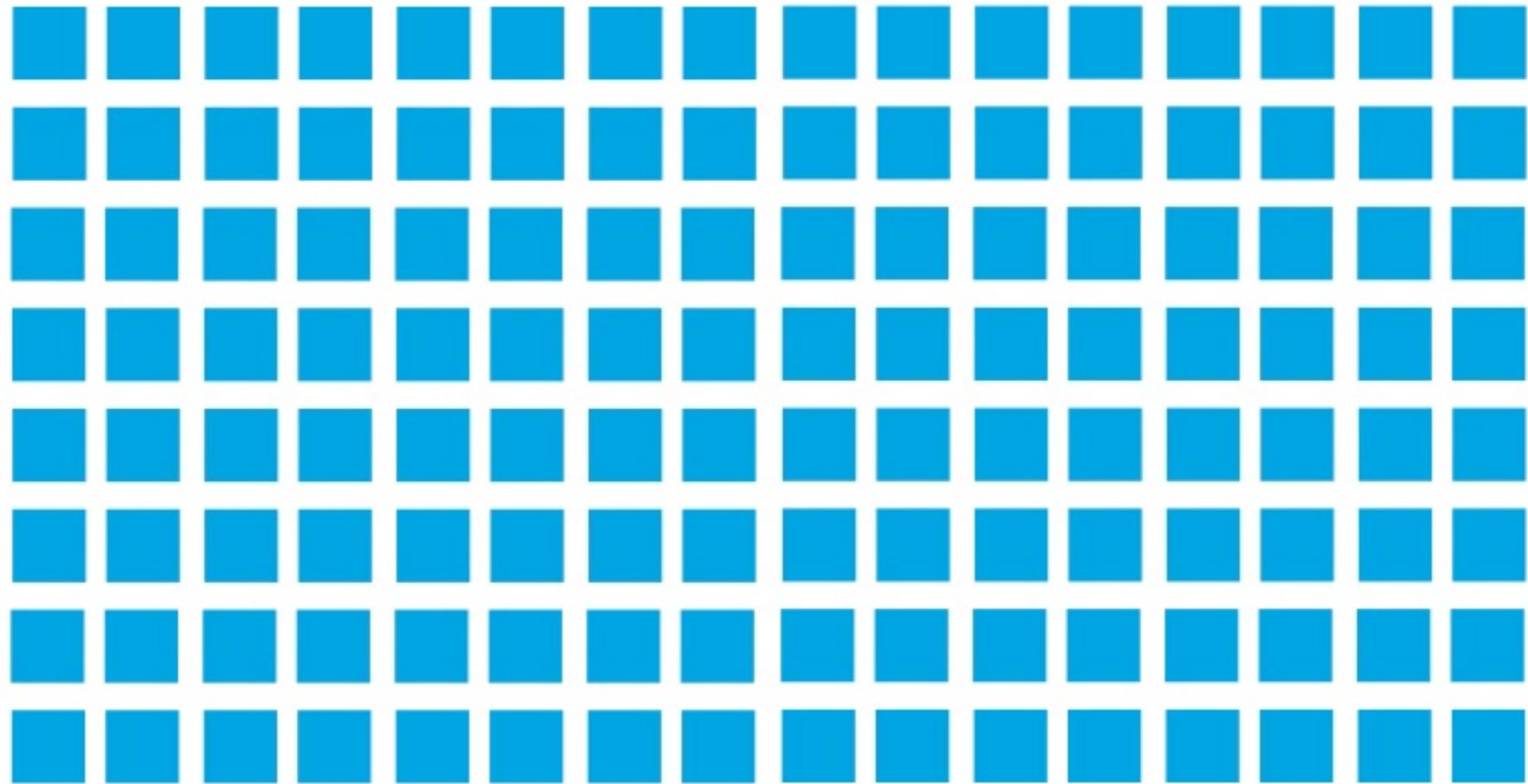


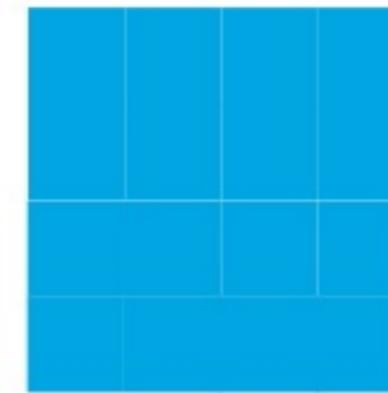
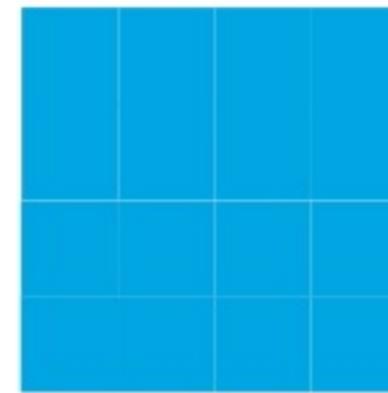
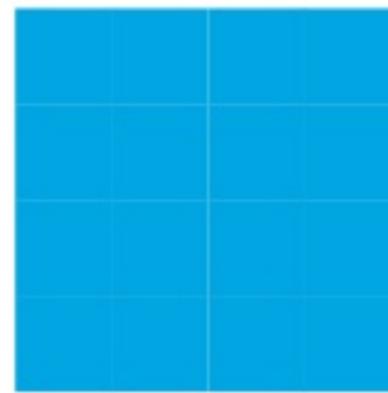
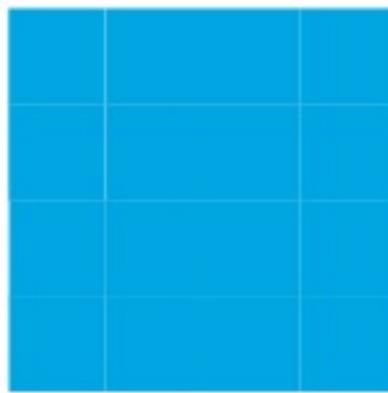
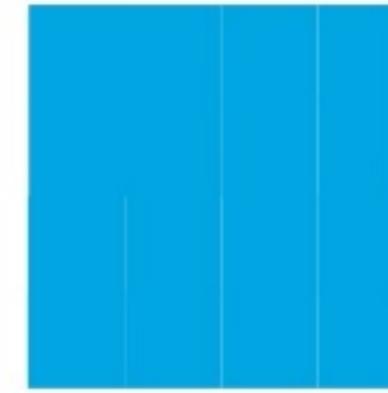
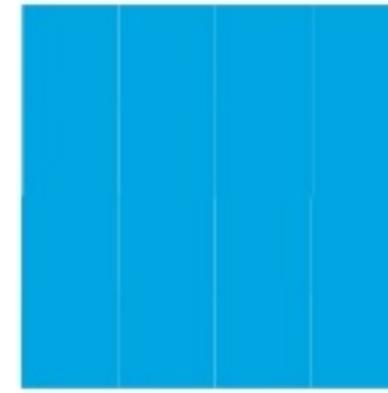
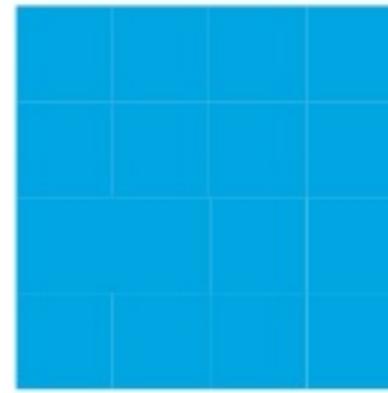
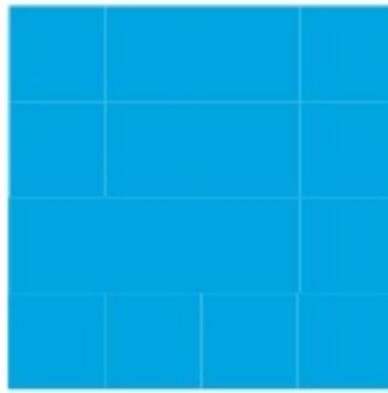






local[2]







Writing generic code for Spark's parallel collections

The RDD is invariant

$T \leq U \neq \text{RDD}[T] \leq \text{RDD}[U]$

The RDD is invariant

$T \leq U \neq \text{RDD}[T] \leq \text{RDD}[U]$

dog *animal*

$T <: U \quad \cancel{RDD[T] <: RDD[U]}$

```
trait HasUserId { val userid: Int }
case class Transaction(override val userid: Int,
                      timestamp: Int,
                      amount: Double)
extends HasUserId {}

def badKeyByUserId(r: RDD[HasUserId]) = r.map(x => (x.userid, x))
```

$T <: U \quad \cancel{RDD[T] <: RDD[U]}$

```
trait HasUserId { val userid: Int }
case class Transaction(override val userid: Int,
                      timestamp: Int,
                      amount: Double)
extends HasUserId {}

def badKeyByUserId(r: RDD[HasUserId]) = r.map(x => (x.userid, x))
```

```
val xacts = spark.parallelize(Array(  
    Transaction(1, 1, 1.0),  
    Transaction(2, 2, 1.0)  
))
```

```
badKeyByUserId(xacts)
```

```
<console>: error: type mismatch;  
  found   : org.apache.spark.rdd.RDD[Transaction]  
  required: org.apache.spark.rdd.RDD[HasUserId]
```

```
Note: Transaction <: HasUserID, but class RDD is invariant in type T.  
You may wish to define T as +T instead. (SLS 4.5)
```

```
    badKeyByUserId(xacts)
```



```
val xacts = spark.parallelize(Array(  
    Transaction(1, 1, 1.0),  
    Transaction(2, 2, 1.0)  
))
```

```
badKeyByUserId(xacts)
```

```
<console>: error: type mismatch;  
  found    : org.apache.spark.rdd.RDD[Transaction]  
  required: org.apache.spark.rdd.RDD[HasUserId]
```

Note: Transaction <: HasUserID, but class RDD is invariant in type T.

You may wish to define T as +T instead. (SLS 4.5)

```
    badKeyByUserId(xacts)
```



3 JUNE 2018

Spark's RDD API, variance, and typeclasses

This brief post is based on material that [Erik](#) and I didn't have time to cover in our [Spark+AI Summit talk](#); it will show you how to use Scala's *implicit parameter* mechanism to work around an aspect of the RDD API that can make it difficult to write generic functions. This post will be especially useful for experienced Spark users who are relatively new to Scala.

If you've written reusable code that uses Spark's RDD API, you might have run into headaches related to *variance*. The RDD is an *invariant API*, meaning that `RDD[T]` and `RDD[U]` are *unrelated types* if `T` and `U` are *different types* – even if there is a subtyping relation between `T` and `U`.

Let's say you had a Scala trait and some concrete class extending that trait, like these:

```
trait HasUserId { val userid: Int }

case class Transaction(override val userid: Int,
                      timestamp: Int,
                      amount: Double)
  extends HasUserId {}
```

You might then want to write a function operating on an RDD of any type that is a subtype of your



An example: natural join

A	B	C	D	E
■■■■■	■■■■■	■■■■■	■■■■■	■■■■■
■■■■■	■■■■■	■■■■■	■■■■■	■■■■■
■■■■■	■■■■■	■■■■■	■■■■■	■■■■■
■■■■■	■■■■■	■■■■■	■■■■■	■■■■■
■■■■■	■■■■■	■■■■■	■■■■■	■■■■■
■■■■■	■■■■■	■■■■■	■■■■■	■■■■■
■■■■■	■■■■■	■■■■■	■■■■■	■■■■■

A	B	E	X	Y
■■■■■	■■■■■	■■■■■	■■■■■	■■■■■
■■■■■	■■■■■	■■■■■	■■■■■	■■■■■
■■■■■	■■■■■	■■■■■	■■■■■	■■■■■
■■■■■	■■■■■	■■■■■	■■■■■	■■■■■
■■■■■	■■■■■	■■■■■	■■■■■	■■■■■
■■■■■	■■■■■	■■■■■	■■■■■	■■■■■
■■■■■	■■■■■	■■■■■	■■■■■	■■■■■

An example: natural join

A	B	C	D	E
Light Gray				
Dark Gray				
Light Gray				
Light Gray				
Dark Gray				

A	B	E	X	Y
Light Gray				
Light Gray				
Light Gray				
Light Gray				
Dark Gray				

An example: natural join

A	B	C	D	E	X	Y
■	■	■	■	■	■	■
■	■	■	■	■	■	■

Ad-hoc natural join

```
df1.join(df2, df1("a") === df2("a") &&  
         df1("b") === df2("b") &&  
         df1("e") === df2("e"))
```

```
def natjoin(left: DataFrame, right: DataFrame): DataFrame = {
    val lcols = left.columns
    val rcols = right.columns
    val ccols = lcols.toSet intersect rcols.toSet

    if(ccols.isEmpty)
        left.limit(0).crossJoin(right.limit(0))
    else
        left
            .join(right, ccols.map {col => left(col) === right(col)}.reduce(_ && _))
            .select(lcols.collect { case c if ccols.contains(c) => left(c) } ++
                    lcols.collect { case c if !ccols.contains(c) => left(c) } ++
                    rcols.collect { case c if !ccols.contains(c) => right(c) } :_*)
}
```



```
def natjoin(left: DataFrame, right: DataFrame): DataFrame = {
    val lcols = left.columns
    val rcols = right.columns
    val ccols = lcols.toSet intersect rcols.toSet

    if(ccols.isEmpty)
        left.limit(0).crossJoin(right.limit(0))
    else
        left
            .join(right, ccols.map {col => left(col) === right(col)}.reduce(_ && _))
            .select(lcols.collect { case c if ccols.contains(c) => left(c) } ++
                    lcols.collect { case c if !ccols.contains(c) => left(c) } ++
                    rcols.collect { case c if !ccols.contains(c) => right(c) } :_*)
}
```



```
def natjoin(left: DataFrame, right: DataFrame): DataFrame = {  
    val lcols = left.columns  
    val rcols = right.columns  
    val ccols = lcols.toSet intersect rcols.toSet  
  
    if(ccols.isEmpty)  
        left.limit(0).crossJoin(right.limit(0))  
    else  
        left  
            .join(right, ccols.map {col => left(col) === right(col)}.reduce(_ && _))  
            .select(lcols.collect { case c if ccols.contains(c) => left(c) } ++  
                    lcols.collect { case c if !ccols.contains(c) => left(c) } ++  
                    rcols.collect { case c if !ccols.contains(c) => right(c) } :_*)  
}
```

```
def natjoin(left: DataFrame, right: DataFrame): DataFrame = {
    val lcols = left.columns
    val rcols = right.columns
    val ccols = lcols.toSet intersect rcols.toSet

    if(ccols.isEmpty)
        left.limit(0).crossJoin(right.limit(0))
    else
        left
            .join(right, ccols.map {col => left(col) === right(col)}.reduce(_ && _))
            .select(lcols.collect { case c if ccols.contains(c) => left(c) } ++
                    lcols.collect { case c if !ccols.contains(c) => left(c) } ++
                    rcols.collect { case c if !ccols.contains(c) => right(c) } :_*)
}
```



```
def natjoin(left: DataFrame, right: DataFrame): DataFrame = {
    val lcols = left.columns
    val rcols = right.columns
    val ccols = lcols.toSet intersect rcols.toSet

    if(ccols.isEmpty)
        left.limit(0).crossJoin(right.limit(0))
    else
        left
            .join(right, ccols.map {col => left(col) === right(col)} .reduce(_ && _))
            .select(lcols.collect { case c if ccols.contains(c) => left(c) } ++
                    lcols.collect { case c if !ccols.contains(c) => left(c) } ++
                    rcols.collect { case c if !ccols.contains(c) => right(c) } : _*)
}
```

dynamically constructing expressions



```
def natjoin(left: DataFrame, right: DataFrame): DataFrame = {
    val lcols = left.columns
    val rcols = right.columns
    val ccols = lcols.toSet intersect rcols.toSet

    if(ccols.isEmpty)
        left.limit(0).crossJoin(right.limit(0))
    else
        left
            .join(right, ccols.map {col => left(col) === right(col)}.reduce(_ && _))
            .select(lcols.collect { case c if ccols.contains(c) => left(c) } ++
                    lcols.collect { case c if !ccols.contains(c) => left(c) } ++
                    rcols.collect { case c if !ccols.contains(c) => right(c) } :_*)
}
```

dynamically constructing expressions



```
def natjoin(left: DataFrame, right: DataFrame): DataFrame = {
    val lcols = left.columns
    val rcols = right.columns
    val ccols = lcols.toSet intersect rcols.toSet

    if(ccols.isEmpty)
        left.limit(0).crossJoin(right.limit(0))
    else
        left
            .join(right, ccols.map {col => left(col) === right(col)} .reduce(_ && _))
            .select(lcols.collect { case c if ccols.contains(c) => left(c) } ++
                    lcols.collect { case c if !ccols.contains(c) => left(c) } ++
                    rcols.collect { case c if !ccols.contains(c) => right(c) } : _*)
}
```

dynamically constructing expressions



```
def natjoin(left: DataFrame, right: DataFrame): DataFrame = {
    val lcols = left.columns
    val rcols = right.columns
    val ccols = lcols.toSet intersect rcols.toSet

    if(ccols.isEmpty)
        left.limit(0).crossJoin(right.limit(0))
    else
        left
            .join(right, ccols.map {col => left(col) === right(col)} .reduce(_ && _))
            .select(lcols.collect { case c if ccols.contains(c) => left(c) } ++
                    lcols.collect { case c if !ccols.contains(c) => left(c) } ++
                    rcols.collect { case c if !ccols.contains(c) => right(c) } : _*)
}
```

[**left.a === right.a, left.b === right.b, ...**]

```
def natjoin(left: DataFrame, right: DataFrame): DataFrame = {
    val lcols = left.columns
    val rcols = right.columns
    val ccols = lcols.toSet intersect rcols.toSet

    if(ccols.isEmpty)
        left.limit(0).crossJoin(right.limit(0))
    else
        left
            .join(right, ccols.map {col => left(col) === right(col)}.reduce(_ && _))
            .select(lcols.collect { case c if ccols.contains(c) => left(c) } ++
                    lcols.collect { case c if !ccols.contains(c) => left(c) } ++
                    rcols.collect { case c if !ccols.contains(c) => right(c) } :_*)
}
```

left.a === right.a && left.b === right.b && ...

```
def natjoin(left: DataFrame, right: DataFrame): DataFrame = {
    val lcols = left.columns
    val rcols = right.columns
    val ccols = lcols.toSet intersect rcols.toSet

    if(ccols.isEmpty)
        left.limit(0).crossJoin(right.limit(0))
    else
        left
            .join(right, ccols.map {col => left(col) === right(col)} .reduce(_ && _))
            .select(lcols.collect { case c if ccols.contains(c) => left(c) } ++
                    lcols.collect { case c if !ccols.contains(c) => left(c) } ++
                    rcols.collect { case c if !ccols.contains(c) => right(c) } : _*)
}
```

left.a === right.a && left.b === right.b && ...

```
def natjoin(left: DataFrame, right: DataFrame): DataFrame = {
    val lcols = left.columns
    val rcols = right.columns
    val ccols = lcols.toSet intersect rcols.toSet

    if(ccols.isEmpty)
        left.limit(0).crossJoin(right.limit(0))
    else
        left
            .join(right, ccols.map {col => left(col) === right(col)}.reduce(_ && _))
            .select(lcols.collect { case c if ccols.contains(c) => left(c) } ++
                    lcols.collect { case c if !ccols.contains(c) => left(c) } ++
                    rcols.collect { case c if !ccols.contains(c) => right(c) } : _*)
}
```



```
def natjoin(left: DataFrame, right: DataFrame): DataFrame = {
    val lcols = left.columns
    val rcols = right.columns
    val ccols = lcols.toSet intersect rcols.toSet

    if(ccols.isEmpty)
        left.limit(0).crossJoin(right.limit(0))
    else
        left
            .join(right, ccols.map {col => left(col) === right(col)}.reduce(_ && _))
            .select(lcols.collect { case c if ccols.contains(c) => left(c) } ++
                    lcols.collect { case c if !ccols.contains(c) => left(c) } ++
                    rcols.collect { case c if !ccols.contains(c) => right(c) } : _*)
}
```

dynamically constructing column lists



```
def natjoin(left: DataFrame, right: DataFrame): DataFrame = {
    val lcols = left.columns
    val rcols = right.columns
    val ccols = lcols.toSet intersect rcols.toSet

    if(ccols.isEmpty)
        left.limit(0).crossJoin(right.limit(0))
    else
        left
            .join(right, ccols.map {col => left(col) === right(col)}.reduce(_ && _))
            .select(lcols.collect { case c if ccols.contains(c) => left(c) } ++
                    lcols.collect { case c if !ccols.contains(c) => left(c) } ++
                    rcols.collect { case c if !ccols.contains(c) => right(c) } : _*)
}
```

dynamically constructing column lists



```
case class DFWithNatJoin(df: DataFrame)
  extends NaturalJoining {
  def natjoin(other: DataFrame): DataFrame = super.natjoin(df, other)
}

object NaturalJoin extends NaturalJoining {
  object implicits {
    implicit def dfWithNatJoin(df: DataFrame) = DFWithNatJoin(df)
  }
}
```

```
case class DFWithNatJoin(df: DataFrame)
    extends NaturalJoining {
  def natjoin(other: DataFrame): DataFrame = super.natjoin(df, other)
}

object NaturalJoin extends NaturalJoining {
  object implicits {
    implicit def dfWithNatJoin(df: DataFrame) = DFWithNatJoin(df)
  }
}
```

```
case class DFWithNatJoin(df: DataFrame)
    extends NaturalJoining {
  def natjoin(other: DataFrame): DataFrame = super.natjoin(df, other)
}

object NaturalJoin extends NaturalJoining {
  object implicits {
    implicit def dfWithNatJoin(df: DataFrame) = DFWithNatJoin(df)
  }
}

import NaturalJoin.implicits._
df.natjoin(otherdf)
```

User-defined functions

```
{"a": 1, "b": "wilma", ..., "x": "club"}  
{"a": 2, "b": "betty", ..., "x": "diamond"}  
{"a": 3, "b": "fred", ..., "x": "heart"}  
{"a": 4, "b": "barney", ..., "x": "spade"}
```

User-defined functions

```
{"a": 1, "b": "wilma", ..., "x": "club"}
```

```
{"a": 2, "b": "betty", ..., "x": "diamond"}
```

```
{"a": 3, "b": "fred", ..., "x": "heart"}
```

```
{"a": 4, "b": "barney", ..., "x": "spade"}
```



wilma	club
betty	diamond
fred	heart
barney	spade

```
import json
from pyspark.sql.types import *
from pyspark.sql.functions import udf

def selectively_structure(fields):
    resultType = StructType([StructField(f, StringType(), nullable=True)
        for f in fields])
    def impl(js):
        try:
            d = json.loads(js)
            return [str(d.get(f)) for f in fields]
        except:
            return [None] * len(fields)
    return udf(impl, resultType)
```

```
import json
from pyspark.sql.types import *
from pyspark.sql.functions import udf

def selectively_structure(fields):
    resultType = StructType([StructField(f, StringType(), nullable=True)
        for f in fields])
    def impl(js):
        try:
            d = json.loads(js)
            return [str(d.get(f)) for f in fields]
        except:
            return [None] * len(fields)
    return udf(impl, resultType)
```

```
import json
from pyspark.sql.types import *
from pyspark.sql.functions import udf

def selectively_structure(fields):
    resultType = StructType([StructField(f, StringType(), nullable=True)
        for f in fields])
    def impl(js):
        try:
            d = json.loads(js)
            return [str(d.get(f)) for f in fields]
        except:
            return [None] * len(fields)
    return udf(impl, resultType)
```

```
import json
from pyspark.sql.types import *
from pyspark.sql.functions import udf

def selectively_structure(fields):
    resultType = StructType([StructField(f, StringType(), nullable=True)
        for f in fields])
    def impl(js):
        try:
            d = json.loads(js)
            return [str(d.get(f)) for f in fields]
        except:
            return [None] * len(fields)
    return udf(impl, resultType)
```

```
import json
from pyspark.sql.types import *
from pyspark.sql.functions import udf

def selectively_structure(fields):
    resultType = StructType([StructField(f, StringType(), nullable=True)
        for f in fields])
    def impl(js):
        try:
            d = json.loads(js)
            return [str(d.get(f)) for f in fields]
        except:
            return [None] * len(fields)
    return udf(impl, resultType)
```

```
import json
from pyspark.sql.types import *
from pyspark.sql.functions import udf

def selectively_structure(fields):
    resultType = StructType([StructField(f, StringType(), nullable=True)
        for f in fields])
    def impl(js):
        try:
            d = json.loads(js)
            return [str(d.get(f)) for f in fields]
        except:
            return [None] * len(fields)
    return udf(impl, resultType)
```

```
extract_bx = selectively_structure(["b", "x"])
```

```
structured_df = df.withColumn("result", extract_bx("json"))
```

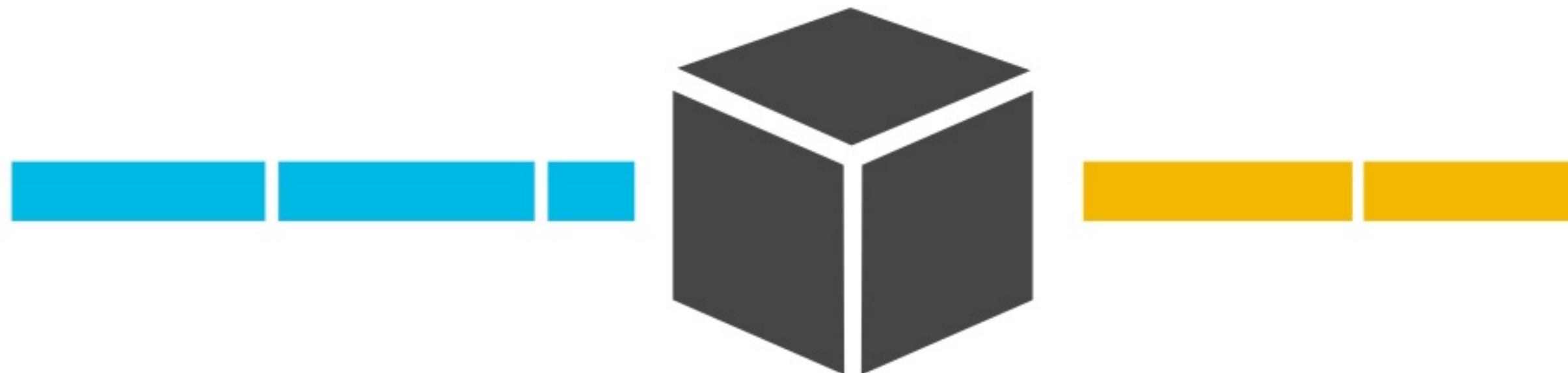
Spark's ML pipelines



model.**transform**(df)

#SAISDD6

Spark's ML pipelines



model.**transform**(df)

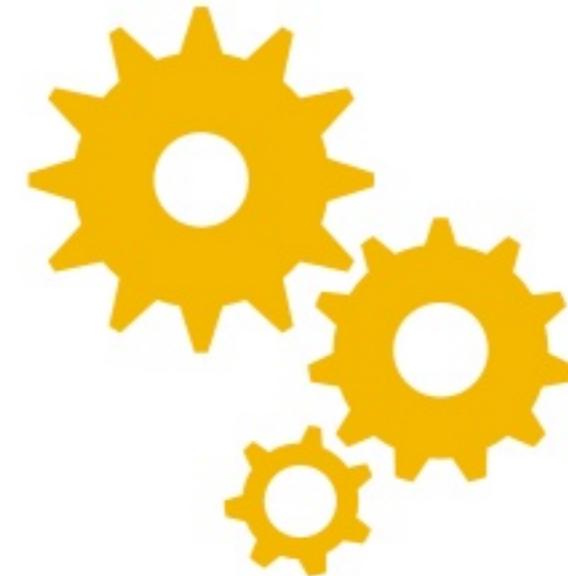
#SAISDD6

Spark's ML pipelines



`estimator.fit(df)`

Spark's ML pipelines



`estimator.fit(df)`

`model.transform(df)`

Working with ML pipelines



model.**transform**(df)

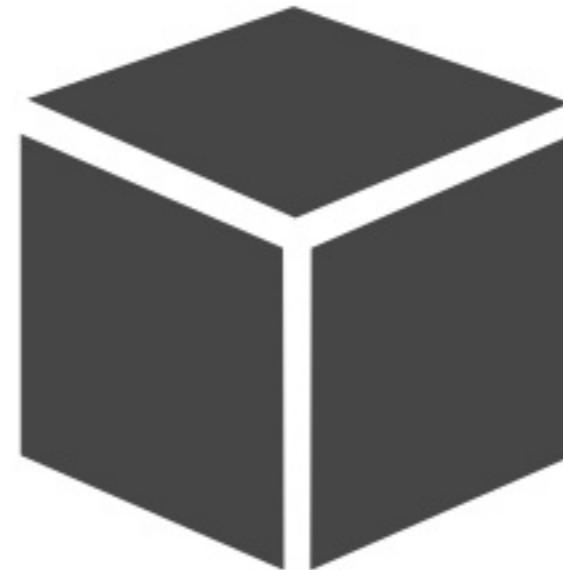
Working with ML pipelines



`model.transform(df)`

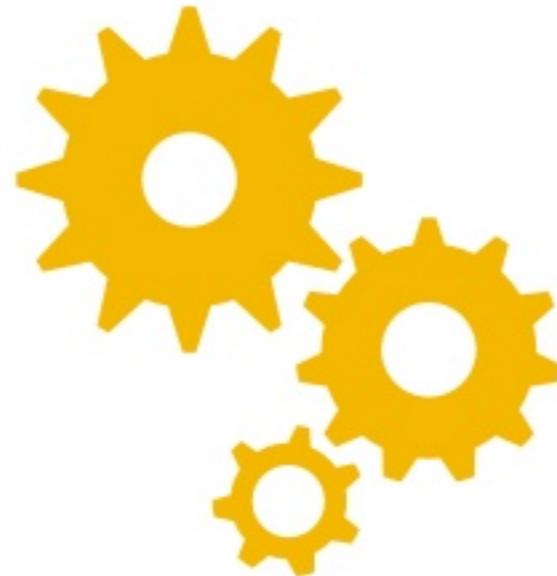
Spark's ML pipelines

Spark's ML pipelines

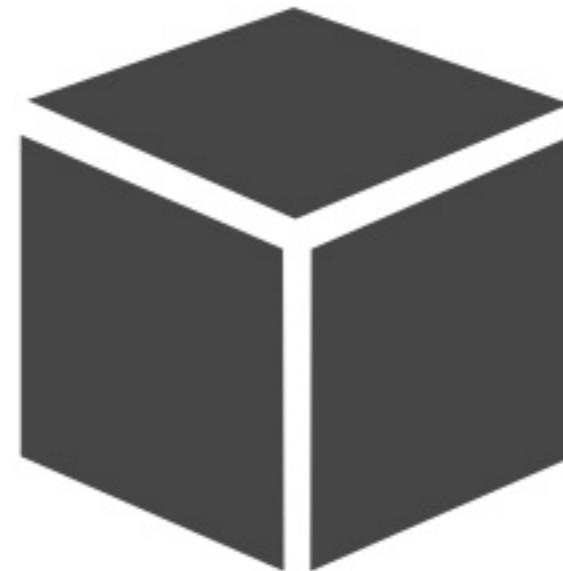


`model.transform(df)`

Spark's ML pipelines

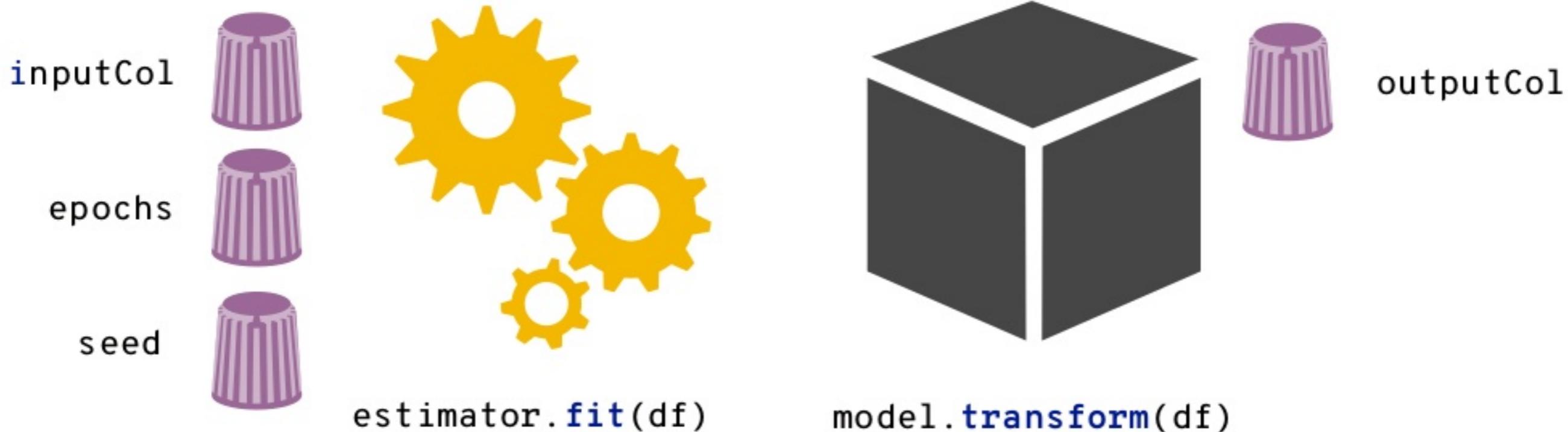


`estimator.fit(df)`



`model.transform(df)`

Spark's ML pipelines



Building Machine Learning Algorithms on Apache Spark: Scaling Out and Up

There are lots of reasons why you might want to implement your own machine learning algorithms on Spark: you might want to experiment with a new idea, try and reproduce results from a recent research paper, or simply to use an existing technique that isn't implemented in MLlib.

In this talk, we'll walk through the process of developing a new machine learning algorithm for Spark. We'll start with the basics, by considering how we'd design a scale-out parallel implementation of our unsupervised learning technique. The bulk of the talk will focus on the details you need to know to turn an algorithm design into an efficient parallel implementation on Spark.

We'll start by reviewing a simple RDD-based implementation, show some improvements, point out some pitfalls to avoid, and iteratively extend our implementation to support contemporary Spark features like ML Pipelines and structured query processing. We'll conclude by briefly examining some useful techniques to complement scale-out performance by scaling our code up, taking advantage of specialized hardware to accelerate single-worker performance.

You'll leave this talk with everything you need to build a new machine learning technique that runs on Spark.

Session hashtag: #DS4SALS



Forecast

Basic considerations for reusable Spark code

Generic functions for parallel collections

Extending data frames with custom aggregates

Exposing JVM libraries to Python

Sharing your work with the world

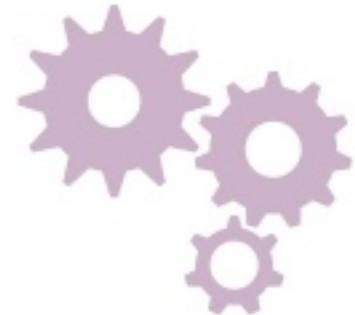


About Erik

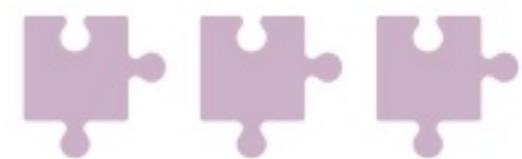


User-defined aggregates: the fundamentals

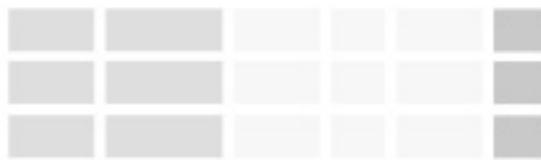
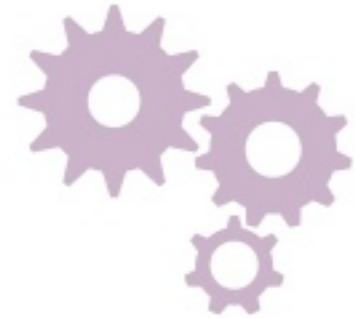
Three components



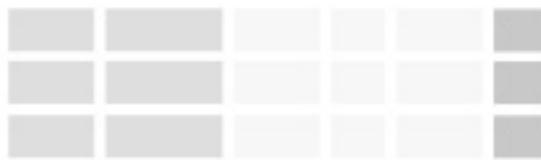
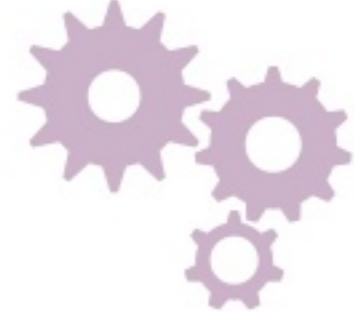
Three components



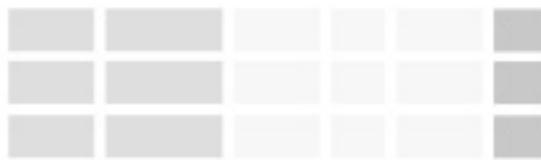
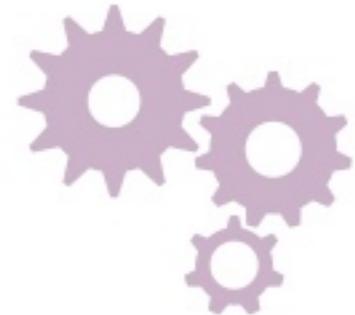
Three components



Three components



Three components





User-defined aggregates: the implementation

```
case class TDigestUDAF[N](deltaV: Double, maxDiscreteV: Int)
  (implicit num: Numeric[N], dataTpe: TDigestUDADataType[N])
extends UserDefinedAggregateFunction {

  def deterministic: Boolean = false

  def inputSchema: StructType =
    StructType(StructField("x", dataTpe.tpe) :: Nil)

  def bufferSchema: StructType =
    StructType(StructField("tdigest", TDigestUDT) :: Nil)

  def dataType: DataType = TDigestUDT
}
```

```
case class TDigestUDAF[N](deltaV: Double, maxDiscreteV: Int)
  (implicit num: Numeric[N], dataTpe: TDigestUDADataType[N])
extends UserDefinedAggregateFunction {

  def deterministic: Boolean = false

  def inputSchema: StructType =
    StructType(StructField("x", dataTpe.tpe) :: Nil)

  def bufferSchema: StructType =
    StructType(StructField("tdigest", TDigestUDT) :: Nil)

  def dataType: DataType = TDigestUDT
```

```
case class TDigestUDAF[N](deltaV: Double, maxDiscreteV: Int)
  (implicit num: Numeric[N], dataTpe: TDigestUDADataType[N])
extends UserDefinedAggregateFunction {

  def deterministic: Boolean = false

  def inputSchema: StructType =
    StructType(StructField("x", dataTpe.tpe) :: Nil)

  def bufferSchema: StructType =
    StructType(StructField("tdigest", TDigestUDT) :: Nil)

  def dataType: DataType = TDigestUDT
```

```
case class TDigestUDAF[N](deltaV: Double, maxDiscreteV: Int)
  (implicit num: Numeric[N], dataTpe: TDigestUDAFDataType[N])
extends UserDefinedAggregateFunction {

  def deterministic: Boolean = false

  def inputSchema: StructType =
    StructType(StructField("x", dataTpe.tpe) :: Nil)

  def bufferSchema: StructType =
    StructType(StructField("tdigest", TDigestUDT) :: Nil)

  def dataType: DataType = TDigestUDT
```

```
case class TDigestUDAF[N](deltaV: Double, maxDiscreteV: Int)
  (implicit num: Numeric[N], dataTpe: TDigestUDADataType[N])
extends UserDefinedAggregateFunction {

  def deterministic: Boolean = false

  def inputSchema: StructType =
    StructType(StructField("x", dataTpe.tpe) :: Nil)

  def bufferSchema: StructType =
    StructType(StructField("tdigest", TDigestUDT) :: Nil)

  def dataType: DataType = TDigestUDT
```



```
case class TDigestUDAF[N](deltaV: Double, maxDiscreteV: Int)
  (implicit num: Numeric[N], dataTpe: TDigestUDADataType[N])
extends UserDefinedAggregateFunction {

  def deterministic: Boolean = false

  def inputSchema: StructType =
    StructType(StructField("x", dataTpe.tpe) :: Nil)

  def bufferSchema: StructType =
    StructType(StructField("tdigest", TDigestUDT) :: Nil)

  def dataType: DataType = TDigestUDT
```

```
case class TDigestUDAF[N](deltaV: Double, maxDiscreteV: Int)
  (implicit num: Numeric[N], dataTpe: TDigestUDAFDataType[N])
extends UserDefinedAggregateFunction {

  def deterministic: Boolean = false

  def inputSchema: StructType =
    StructType(StructField("x", dataTpe.tpe) :: Nil)

  def bufferSchema: StructType =
    StructType(StructField("tdigest", TDigestUDT) :: Nil)

  def dataType: DataType = TDigestUDT
```

```
case class TDigestUDAF[N](deltaV: Double, maxDiscreteV: Int)
  (implicit num: Numeric[N], dataTpe: TDigestUDAFDataType[N])
extends UserDefinedAggregateFunction {

  def deterministic: Boolean = false

  def inputSchema: StructType =
    StructType(StructField("x", dataTpe.tpe) :: Nil)

  def bufferSchema: StructType =
    StructType(StructField("tdigest", TDigestUDT) :: Nil)

  def dataType: DataType = TDigestUDT
```

Four main functions: initialize



initialize

Four main functions: initialize



initialize

```
def initialize(buf: MutableAggregationBuffer): Unit = {
    buf(0) = TDigestSQL(TDigest.empty(deltaV, maxDiscreteV))
}

def evaluate(buf: Row): Any = buf.getAs[TDigestSQL](0)
```

```
def initialize(buf: MutableAggregationBuffer): Unit = {
    buf(0) = TDigestSQL(TDigest.empty(deltaV, maxDiscreteV))
}

def evaluate(buf: Row): Any = buf.getAs[TDigestSQL](0)
```

Four main functions: evaluate



evaluate

Four main functions: evaluate



evaluate

```
def initialize(buf: MutableAggregationBuffer): Unit = {
    buf(0) = TDigestSQL(TDigest.empty(deltaV, maxDiscreteV))
}

def evaluate(buf: Row): Any = buf.getAs[TDigestSQL](0)
```

```
def initialize(buf: MutableAggregationBuffer): Unit = {
    buf(0) = TDigestSQL(TDigest.empty(deltaV, maxDiscreteV))
}

def evaluate(buf: Row): Any = buf.getAs[TDigestSQL](0)
```

Four main functions: update



update

Four main functions: update



update

```
def update(buf: MutableAggregationBuffer, input: Row): Unit = {
    if (!input.isNullAt(0)) {
        buf(0) = TDigestSQL(buf.getAs[TDigestSQL](0).tdigest +
            num.toDouble(input.getAs[N](0)))
    }
}

def merge(buf1: MutableAggregationBuffer, buf2: Row): Unit = {
    buf1(0) = TDigestSQL(buf1.getAs[TDigestSQL](0).tdigest ++
        buf2.getAs[TDigestSQL](0).tdigest)
}
```

```
def update(buf: MutableAggregationBuffer, input: Row): Unit = {
    if (!input.isNullAt(0)) {
        buf(0) = TDigestSQL(buf.getAs[TDigestSQL](0).tdigest +
            num.toDouble(input.getAs[N](0)))
    }
}

def merge(buf1: MutableAggregationBuffer, buf2: Row): Unit = {
    buf1(0) = TDigestSQL(buf1.getAs[TDigestSQL](0).tdigest ++
        buf2.getAs[TDigestSQL](0).tdigest)
}
```

```
def update(buf: MutableAggregationBuffer, input: Row): Unit = {
    if (!input.isNullAt(0)) {
        buf(0) = TDigestSQL(buf.getAs[TDigestSQL](0).tdigest +
            num.toDouble(input.getAs[N](0)))
    }
}

def merge(buf1: MutableAggregationBuffer, buf2: Row): Unit = {
    buf1(0) = TDigestSQL(buf1.getAs[TDigestSQL](0).tdigest ++
        buf2.getAs[TDigestSQL](0).tdigest)
}
```

```
def update(buf: MutableAggregationBuffer, input: Row): Unit = {
  if (!input.isNullAt(0)) {
    buf(0) = TDigestSQL(buf.getAs[TDigestSQL](0).tdigest +
      num.toDouble(input.getAs[N](0)))
  }
}

def merge(buf1: MutableAggregationBuffer, buf2: Row): Unit = {
  buf1(0) = TDigestSQL(buf1.getAs[TDigestSQL](0).tdigest ++
    buf2.getAs[TDigestSQL](0).tdigest)
}
```

```
def update(buf: MutableAggregationBuffer, input: Row): Unit = {
  if (!input.isNullAt(0)) {
    buf(0) = TDigestSQL(buf.getAs[TDigestSQL](0).tdigest +
      num.toDouble(input.getAs[N](0)))
  }
}

def merge(buf1: MutableAggregationBuffer, buf2: Row): Unit = {
  buf1(0) = TDigestSQL(buf1.getAs[TDigestSQL](0).tdigest ++
    buf2.getAs[TDigestSQL](0).tdigest)
}
```

```
def update(buf: MutableAggregationBuffer, input: Row): Unit = {
  if (!input.isNullAt(0)) {
    buf(0) = TDigestSQL(buf.getAs[TDigestSQL](0).tdigest +
      num.toDouble(input.getAs[N](0)))
  }
}

def merge(buf1: MutableAggregationBuffer, buf2: Row): Unit = {
  buf1(0) = TDigestSQL(buf1.getAs[TDigestSQL](0).tdigest ++
    buf2.getAs[TDigestSQL](0).tdigest)
}
```

```
def update(buf: MutableAggregationBuffer, input: Row): Unit = {
  if (!input.isNullAt(0)) {
    buf(0) = TDigestSQL(buf.getAs[TDigestSQL](0).tdigest +
      num.toDouble(input.getAs[N](0)))
  }
}
```

```
def merge(buf1: MutableAggregationBuffer, buf2: Row): Unit = {
  buf1(0) = TDigestSQL(buf1.getAs[TDigestSQL](0).tdigest ++
    buf2.getAs[TDigestSQL](0).tdigest)
}
```

Four main functions: merge

1

2

merge



Four main functions: merge

1 + 2

merge

```
def update(buf: MutableAggregationBuffer, input: Row): Unit = {
    if (!input.isNullAt(0)) {
        buf(0) = TDigestSQL(buf.getAs[TDigestSQL](0).tdigest +
            num.toDouble(input.getAs[N](0)))
    }
}

def merge(buf1: MutableAggregationBuffer, buf2: Row): Unit = {
    buf1(0) = TDigestSQL(buf1.getAs[TDigestSQL](0).tdigest ++
        buf2.getAs[TDigestSQL](0).tdigest)
}
```

```
def update(buf: MutableAggregationBuffer, input: Row): Unit = {
    if (!input.isNullAt(0)) {
        buf(0) = TDigestSQL(buf.getAs[TDigestSQL](0).tdigest +
            num.toDouble(input.getAs[N](0)))
    }
}

def merge(buf1: MutableAggregationBuffer, buf2: Row): Unit = {
    buf1(0) = TDigestSQL(buf1.getAs[TDigestSQL](0).tdigest ++
        buf2.getAs[TDigestSQL](0).tdigest)
}
```



User-defined aggregates: User-defined types

User-defined types

```
package org.apache.spark.isarnproject.sketches.udt

@SQLUserDefinedType(udt = classOf[TDigestUDT])
case class TDigestSQL(tdigest: TDigest)

class TDigestUDT extends UserDefinedType[TDigestSQL] {
    def userClass: Class[TDigestSQL] = classOf[TDigestSQL]

    // ...
}
```

User-defined types

```
package org.apache.spark.isarnproject.sketches.udt

@SQLUserDefinedType(udt = classOf[TDigestUDT])
case class TDigestSQL(tdigest: TDigest)

class TDigestUDT extends UserDefinedType[TDigestSQL] {
    def userClass: Class[TDigestSQL] = classOf[TDigestSQL]

    // ...
}
```

User-defined types

```
package org.apache.spark.isarnproject.sketches.udt

@SQLUserDefinedType(udt = classOf[TDigestUDT])
case class TDigestSQL(tdigest: TDigest)

class TDigestUDT extends UserDefinedType[TDigestSQL] {
    def userClass: Class[TDigestSQL] = classOf[TDigestSQL]

    // ...
}
```

Implementing custom types

```
class TDigestUDT extends UserDefinedType[TDigestSQL] {  
    def userClass: Class[TDigestSQL] = classOf[TDigestSQL]  
  
    override def pyUDT: String =  
        "isarnproject.sketches.udt.tdigest.TDigestUDT"  
  
    override def typeName: String = "tdigest"  
  
    def sqlType: DataType = StructType(  
        StructField("delta", DoubleType, false) ::  
        /* ... */  
        StructField("clustM", ArrayType(DoubleType, false), false) ::  
        Nil)
```



```
class TDigestUDT extends UserDefinedType[TDigestSQL] {  
    def userClass: Class[TDigestSQL] = classOf[TDigestSQL]  
  
    override def pyUDT: String =  
        "isarnproject.sketches.udt.tdigest.TDigestUDT"  
  
    override def typeName: String = "tdigest"  
  
    def sqlType: DataType = StructType(  
        StructField("delta", DoubleType, false) ::  
        /* ... */  
        StructField("clustM", ArrayType(DoubleType, false), false) ::  
        Nil)
```



```
class TDigestUDT extends UserDefinedType[TDigestSQL] {  
    def userClass: Class[TDigestSQL] = classOf[TDigestSQL]  
  
    override def pyUDT: String =  
        "isarnproject.sketches.udt.tdigest.TDigestUDT"  
  
override def typeName: String = "tdigest"  
  
    def sqlType: DataType = StructType(  
        StructField("delta", DoubleType, false) ::  
        /* ... */  
        StructField("clustM", ArrayType(DoubleType, false), false) ::  
        Nil)
```



```
class TDigestUDT extends UserDefinedType[TDigestSQL] {  
    def userClass: Class[TDigestSQL] = classOf[TDigestSQL]  
  
    override def pyUDT: String =  
        "isarnproject.sketches.udt.tdigest.TDigestUDT"  
  
    override def typeName: String = "tdigest"  
  
    def sqlType: DataType = StructType(  
        StructField("delta", DoubleType, false) ::  
        /* ... */  
        StructField("clustM", ArrayType(DoubleType, false), false) ::  
        Nil)
```



```
def serialize(tdsql: TDigestSQL): Any = serializeTD(tdsql.tdigest)

private[sketches] def serializeTD(td: TDigest): InternalRow = {
    val TDigest(delta, maxDiscrete, nclusters, clusters) = td
    val row = new GenericInternalRow(5)
    row.setDouble(0, delta)
    row.setInt(1, maxDiscrete)
    row.setInt(2, nclusters)
    val clustX = clusters.keys.toArray
    val clustM = clusters.values.toArray
    row.update(3, UnsafeArrayData.fromPrimitiveArray(clustX))
    row.update(4, UnsafeArrayData.fromPrimitiveArray(clustM))
    row
}
```

```
def serialize(tdsql: TDigestSQL): Any = serializeTD(tdsql.tdigest)

private[sketches] def serializeTD(td: TDigest): InternalRow = {
    val TDigest(delta, maxDiscrete, nclusters, clusters) = td
    val row = new GenericInternalRow(5)
    row.setDouble(0, delta)
    row.setInt(1, maxDiscrete)
    row.setInt(2, nclusters)
    val clustX = clusters.keys.toArray
    val clustM = clusters.values.toArray
    row.update(3, UnsafeArrayData.fromPrimitiveArray(clustX))
    row.update(4, UnsafeArrayData.fromPrimitiveArray(clustM))
    row
}
```

```
def serialize(tdsql: TDigestSQL): Any = serializeTD(tdsql.tdigest)

private[sketches] def serializeTD(td: TDigest): InternalRow = {
    val TDigest(delta, maxDiscrete, nclusters, clusters) = td
    val row = new GenericInternalRow(5)
    row.setDouble(0, delta)
    row.setInt(1, maxDiscrete)
    row.setInt(2, nclusters)
    val clustX = clusters.keys.toArray
    val clustM = clusters.values.toArray
    row.update(3, UnsafeArrayData.fromPrimitiveArray(clustX))
    row.update(4, UnsafeArrayData.fromPrimitiveArray(clustM))
    row
}
```

```
def serialize(tdsql: TDigestSQL): Any = serializeTD(tdsql.tdigest)

private[sketches] def serializeTD(td: TDigest): InternalRow = {
    val TDigest(delta, maxDiscrete, nclusters, clusters) = td
    val row = new GenericInternalRow(5)
    row.setDouble(0, delta)
    row.setInt(1, maxDiscrete)
    row.setInt(2, nclusters)
    val clustX = clusters.keys.toArray
    val clustM = clusters.values.toArray
    row.update(3, UnsafeArrayData.fromPrimitiveArray(clustX))
    row.update(4, UnsafeArrayData.fromPrimitiveArray(clustM))
    row
}
```

```
def deserialize(td: Any): TDigestSQL = TDigestSQL(deserializeTD(td))

private[sketches] def deserializeTD(datum: Any): TDigest =
  datum match { case row: InternalRow =>
    val delta = row.getDouble(0)
    val maxDiscrete = row.getInt(1)
    val nclusters = row.getInt(2)
    val clustX = row.getArray(3).toDoubleArray()
    val clustM = row.getArray(4).toDoubleArray()
    val clusters = clustX.zip(clustM)
      .foldLeft(TDigestMap.empty) { case (td, e) => td + e }
    TDigest(delta, maxDiscrete, nclusters, clusters)
  }
```



```
def deserialize(td: Any): TDigestSQL = TDigestSQL(deserializeTD(td))

private[sketches] def deserializeTD(datum: Any): TDigest =
  datum match { case row: InternalRow =>
    val delta = row.getDouble(0)
    val maxDiscrete = row.getInt(1)
    val nclusters = row.getInt(2)
    val clustX = row.getArray(3).toDoubleArray()
    val clustM = row.getArray(4).toDoubleArray()
    val clusters = clustX.zip(clustM)
      .foldLeft(TDigestMap.empty) { case (td, e) => td + e }
    TDigest(delta, maxDiscrete, nclusters, clusters)
  }
```



```
def deserialize(td: Any): TDigestSQL = TDigestSQL(deserializeTD(td))

private[sketches] def deserializeTD(datum: Any): TDigest =
  datum match { case row: InternalRow =>
    val delta = row.getDouble(0)
    val maxDiscrete = row.getInt(1)
    val nclusters = row.getInt(2)
    val clustX = row.getArray(3).toDoubleArray()
    val clustM = row.getArray(4).toDoubleArray()
    val clusters = clustX.zip(clustM)
      .foldLeft(TDigestMap.empty) { case (td, e) => td + e }
    TDigest(delta, maxDiscrete, nclusters, clusters)
  }
```



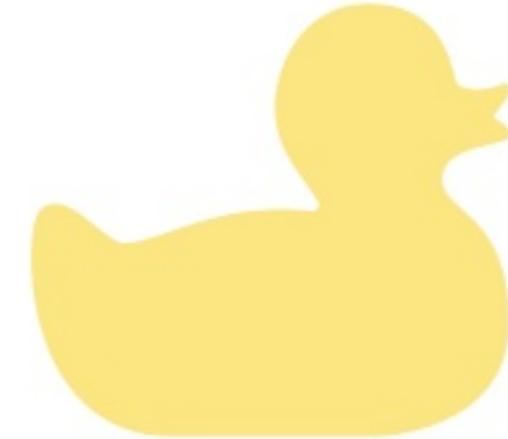
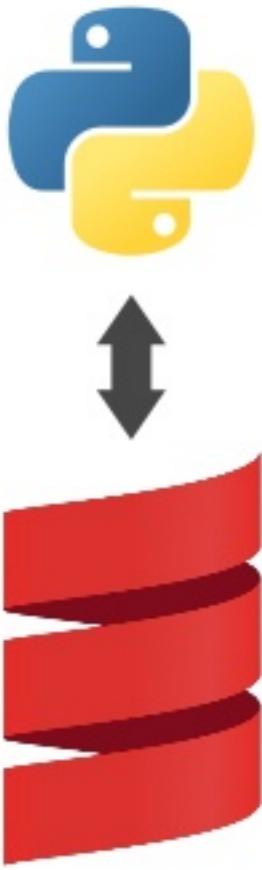


Extending PySpark with your Scala library











```
# class to access the active Spark context for Python
from pyspark.context import SparkContext

# gateway to the JVM from py4j
sparkJVM = SparkContext._active_spark_context._jvm

# use the gateway to access JVM objects and classes
thisThing = sparkJVM.com.path.to.this.thing
```



```
# class to access the active Spark context for Python
from pyspark.context import SparkContext

# gateway to the JVM from py4j
sparkJVM = SparkContext._active_spark_context._jvm

# use the gateway to access JVM objects and classes
thisThing = sparkJVM.com.path.to.this.thing
```



```
# class to access the active Spark context for Python
from pyspark.context import SparkContext

# gateway to the JVM from py4j
sparkJVM = SparkContext._active_spark_context._jvm

# use the gateway to access JVM objects and classes
thisThing = sparkJVM.com.path.to.this.thing
```

A Python-friendly wrapper

```
package org.isarnproject.sketches.udaf

object pythonBindings {
    def tdigestDoubleUDAF(delta: Double, maxDiscrete: Int) =
        TDigestUDAF[Double](delta, maxDiscrete)
}
```

```
package org.isarnproject.sketches.udaf

object pythonBindings {
    def tdigestDoubleUDAF(delta: Double, maxDiscrete: Int) =
        TDigestUDAF[Double](delta, maxDiscrete)
}
```

```
package org.isarnproject.sketches.udaf

object pythonBindings {
    def tdigestDoubleUDAF(delta: Double, maxDiscrete: Int) =
        TDigestUDAF[Double](delta, maxDiscrete)
}
```

```
from pyspark.sql.column import Column, _to_java_column, _to_seq
from pyspark.context import SparkContext

# one of these for each type parameter Double, Int, Long, etc
def tdigestDoubleUDAF(col, delta=0.5, maxDiscrete=0):
    sc = SparkContext._active_spark_context
    pb = sc._jvm.org.isarnproject.sketches.udaf.pythonBindings
    tdapply = pb.tdigestDoubleUDAF(delta, maxDiscrete).apply
    return Column(tdapply(_to_seq(sc, [col], _to_java_column)))
```

```
from pyspark.sql.column import Column, _to_java_column, _to_seq
from pyspark.context import SparkContext

# one of these for each type parameter Double, Int, Long, etc
def tdigestDoubleUDAF(col, delta=0.5, maxDiscrete=0):
    sc = SparkContext._active_spark_context
    pb = sc._jvm.org.isarnproject.sketches.udaf.pythonBindings
    tdapply = pb.tdigestDoubleUDAF(delta, maxDiscrete).apply
    return Column(tdapply(_to_seq(sc, [col], _to_java_column)))
```

```
from pyspark.sql.column import Column, _to_java_column, _to_seq
from pyspark.context import SparkContext

# one of these for each type parameter Double, Int, Long, etc
def tdigestDoubleUDAF(col, delta=0.5, maxDiscrete=0):
    sc = SparkContext._active_spark_context
    pb = sc._jvm.org.isarnproject.sketches.udaf.pythonBindings
    tdapply = pb.tdigestDoubleUDAF(delta, maxDiscrete).apply
    return Column(tdapply(_to_seq(sc, [col], _to_java_column)))
```

```
class TDigestUDT(UserDefinedType):  
    @classmethod  
    def sqlType(cls):  
        return StructType([  
            StructField("delta", DoubleType(), False),  
            StructField("maxDiscrete", IntegerType(), False),  
            StructField("nclusters", IntegerType(), False),  
            StructField("clustX", ArrayType(DoubleType(), False), False),  
            StructField("clustM", ArrayType(DoubleType(), False))])  
  
    # ...
```

```
class TDigestUDT(UserDefinedType):  
    @classmethod  
    def sqlType(cls):  
        return StructType([  
            StructField("delta", DoubleType(), False),  
            StructField("maxDiscrete", IntegerType(), False),  
            StructField("nclusters", IntegerType(), False),  
            StructField("clustX", ArrayType(DoubleType(), False), False),  
            StructField("clustM", ArrayType(DoubleType(), False), False)])  
  
    # ...
```

```
class TDigestUDT(UserDefinedType):  
    # ...  
    @classmethod  
    def module(cls):  
        return "isarnproject.sketches.udt.tdigest"  
  
    @classmethod  
    def scalaUDT(cls):  
        return "org.apache.spark.isarnproject.sketches.udt.TDigestUDT"  
  
    def simpleString(self):  
        return "tdigest"
```

```
class TDigestUDT(UserDefinedType):  
    # ...  
    @classmethod  
    def module(cls):  
        return "isarnproject.sketches.udt.tdigest"  
  
    @classmethod  
    def scalaUDT(cls):  
        return "org.apache.spark.isarnproject.sketches.udt.TDigestUDT"  
  
    def simpleString(self):  
        return "tdigest"
```

```
class TDigestUDT(UserDefinedType):  
    # ...  
    @classmethod  
    def module(cls):  
        return "isarnproject.sketches.udt.tdigest"  
  
    @classmethod  
    def scalaUDT(cls):  
        return "org.apache.spark.isarnproject.sketches.udt.TDigestUDT"  
  
    def simpleString(self):  
        return "tdigest"
```

```
class TDigestUDT(UserDefinedType):
    # ...
    def serialize(self, obj):
        return (obj.delta, obj.maxDiscrete, obj.nclusters, \
                [float(v) for v in obj.clustX], \
                [float(v) for v in obj.clustM])

    def deserialize(self, datum):
        return TDigest(datum[0], datum[1], datum[2], datum[3], datum[4])
```

```
class TDigestUDT(UserDefinedType):
    # ...
    def serialize(self, obj):
        return (obj.delta, obj.maxDiscrete, obj.nclusters, \
                [float(v) for v in obj.clustX], \
                [float(v) for v in obj.clustM])

    def deserialize(self, datum):
        return TDigest(datum[0], datum[1], datum[2], datum[3], datum[4])
```

```
class TDigestUDT extends UserDefinedType[TDigestSQL] {  
    // ...  
    override def pyUDT: String =  
        "isarnproject.sketches.udt.tdigest.TDigestUDT"  
}
```

Python code in JAR files

```
mappings in (Compile, packageBin) += Seq(  
  (baseDirectory.value / "python" / "isarnproject" / "__init__.pyc") ->  
    "isarnproject/__init__.pyc",  
  (baseDirectory.value / "python" / "isarnproject" / "sketches" / "__init__.pyc") ->  
    "isarnproject.sketches/__init__.pyc",  
  (baseDirectory.value / "python" / "isarnproject" / "sketches" / "udaf" / "__init__.pyc") ->  
    "isarnproject.sketches.udaf/__init__.pyc",  
  (baseDirectory.value / "python" / "isarnproject" / "sketches" / "udaf" / "tdigest.py") ->  
    "isarnproject.sketches.udaf/tdigest.py",  
  (baseDirectory.value / "python" / "isarnproject" / "sketches" / "udt" / "__init__.pyc") ->  
    "isarnproject.sketches.udt/__init__.pyc",  
  (baseDirectory.value / "python" / "isarnproject" / "sketches" / "udt" / "tdigest.py") ->  
    "isarnproject.sketches.udt/tdigest.py"  
)
```



```
mappings in (Compile, packageBin) ++= Seq(  
  (baseDirectory.value / "python" / "isarnproject" / "__init__.pyc") ->  
    "isarnproject/__init__.pyc",  
  (baseDirectory.value / "python" / "isarnproject" / "sketches" / "__init__.pyc") ->  
    "isarnproject/sketches/__init__.pyc",  
  (baseDirectory.value / "python" / "isarnproject" / "sketches" / "udaf" / "__init__.pyc") ->  
    "isarnproject/sketches/udaf/__init__.pyc",  
  (baseDirectory.value / "python" / "isarnproject" / "sketches" / "udaf" / "tdigest.py") ->  
    "isarnproject/sketches/udaf/tdigest.py",  
  (baseDirectory.value / "python" / "isarnproject" / "sketches" / "udt" / "__init__.pyc") ->  
    "isarnproject/sketches/udt/__init__.pyc",  
  (baseDirectory.value / "python" / "isarnproject" / "sketches" / "udt" / "tdigest.py") ->  
    "isarnproject/sketches/udt/tdigest.py"  
)
```

```
mappings in (Compile, packageBin) +++= Seq(  
  (baseDirectory.value / "python" / "isarnproject" / "__init__.pyc") ->  
    "isarnproject/__init__.pyc",  
  (baseDirectory.value / "python" / "isarnproject" / "sketches" / "__init__.pyc") ->  
    "isarnproject.sketches/__init__.pyc",  
  (baseDirectory.value / "python" / "isarnproject" / "sketches" / "udaf" / "__init__.pyc") ->  
    "isarnproject.sketches.udaf/__init__.pyc",  
  (baseDirectory.value / "python" / "isarnproject" / "sketches" / "udaf" / "tdigest.py") ->  
    "isarnproject.sketches.udaf/tdigest.py",  
  (baseDirectory.value / "python" / "isarnproject" / "sketches" / "udt" / "__init__.pyc") ->  
    "isarnproject.sketches.udt/__init__.pyc",  
  (baseDirectory.value / "python" / "isarnproject" / "sketches" / "udt" / "tdigest.py") ->  
    "isarnproject.sketches.udt/tdigest.py"  
)
```



```
mappings in (Compile, packageBin) +++= Seq(  
  (baseDirectory.value / "python" / "isarnproject" / "__init__.pyc") ->  
    "isarnproject/__init__.pyc",  
  (baseDirectory.value / "python" / "isarnproject" / "sketches" / "__init__.pyc") ->  
    "isarnproject/sketches/__init__.pyc",  
  (baseDirectory.value / "python" / "isarnproject" / "sketches" / "udaf" / "__init__.pyc") ->  
    "isarnproject/sketches/udaf/__init__.pyc",  
  (baseDirectory.value / "python" / "isarnproject" / "sketches" / "udaf" / "tdigest.pyc) ->  
    "isarnproject/sketches/udaf/tdigest.py",  
  (baseDirectory.value / "python" / "isarnproject" / "sketches" / "udt" / "__init__.pyc") ->  
    "isarnproject/sketches/udt/__init__.pyc",  
  (baseDirectory.value / "python" / "isarnproject" / "sketches" / "udt" / "tdigest.py") ->  
    "isarnproject/sketches/udt/tdigest.py"  
)
```



Cross-building for Python

```
lazy val compilePython = taskKey[Unit]("Compile python files")

compilePython := {
    val s: TaskStreams = streams.value
    s.log.info("compiling python...")
    val stat = (Seq(pythonCMD, "-m", "compileall", "python/") !)
    if (stat != 0) {
        throw new IllegalStateException("python compile failed")
    }
}

(packageBin in Compile) <=>
    (packageBin in Compile).dependsOn(compilePython)
```

```
lazy val compilePython = taskKey[Unit]("Compile python files")

compilePython := {
    val s: TaskStreams = streams.value
    s.log.info("compiling python...")
    val stat = (Seq(pythonCMD, "-m", "compileall", "python/") !)
    if (stat != 0) {
        throw new IllegalStateException("python compile failed")
    }
}

(packageBin in Compile) <=>
    (packageBin in Compile).dependsOn(compilePython)
```

```
lazy val compilePython = taskKey[Unit]("Compile python files")

compilePython := {
    val s: TaskStreams = streams.value
    s.log.info("compiling python...")
    val stat = (Seq(pythonCMD, "-m", "compileall", "python/") !)
    if (stat != 0) {
        throw new IllegalStateException("python compile failed")
    }
}

(packageBin in Compile) <=>
    (packageBin in Compile).dependsOn(compilePython)
```

Using versioned JAR files

```
$ pyspark --packages \
'org.isarnproject:isarn-sketches-spark_2.11:0.3.0-sp2.2-py2.7'
```

Using versioned JAR files

```
$ pyspark --packages \
'org.isarnproject:isarn-sketches-spark_2.11:0.3.0-sp2.2-py2.7'
```

Using versioned JAR files

```
$ pyspark --packages \
'org.isarnproject:isarn-sketches-spark_2.11:0.3.0-sp2.2-py2.7'
```



**Show your work:
publishing results**

Developing with git-flow

```
$ brew install git-flow      # macOS  
$ dnf install git-flow      # Fedora  
$ yum install git-flow      # CentOS  
$ apt-get install git-flow  # Debian and friends
```

(Search the internet for “git flow” to learn more!)

```
# Set up git-flow in this repository
$ git flow init

# Start work on my-awesome-feature; create
# and switch to a feature branch
$ git flow feature start my-awesome-feature
$ ...

# Finish work on my-awesome-feature; merge
# feature/my-awesome-feature to develop
$ git flow feature finish my-awesome-feature
```

```
# Start work on a release branch
$ git flow release start 0.1.0
# Hack and bump version numbers
$ ...
# Finish work on v0.1.0; merge
# release/0.1.0 to develop and master;
# tag v0.1.0
$ git flow release finish 0.1.0
```

```
echo ~/devel/silex on develop(v0.2.0-8-gc258e1a) tracking origin/develop
2280 silex:develop? % 2018-10-03 14:35:41 willb ttys002
```

I

```
echo ~/devel/silex on develop(v0.2.0-8-gc258e1a) tracking origin/develop
2280 silex:develop? % 2018-10-03 14:35:41 willb ttys002
```

I

```
echo ~/devel/silex on develop(v0.2.0-11-gd7c12e8) tracking origin/develop
2299 silex:develop? % cat project/plugins.sbt    2018-10-03 15:54:41 willb ttys002
```

```
echo ~/devel/silex on develop(v0.2.0-11-gd7c12e8) tracking origin/develop
2299 silex:develop? % cat project/plugins.sbt    2018-10-03 15:54:41 willb ttys002
```

Maven Central

Bintray

not really

easy to set up for library developers

trivial

trivial

easy to set up for library users

mostly

yes, via sbt

easy to publish

yes, via sbt + plugins

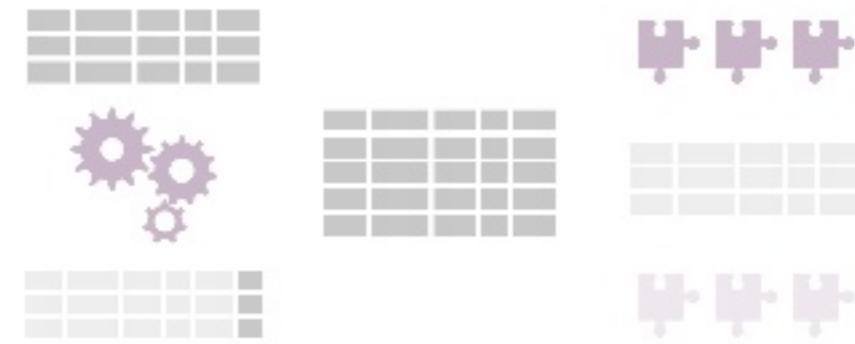
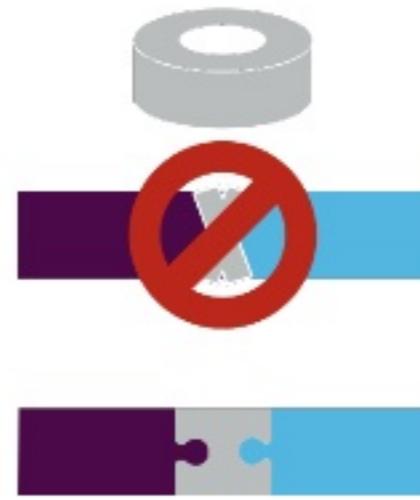
yes

easy to resolve artifacts

mostly

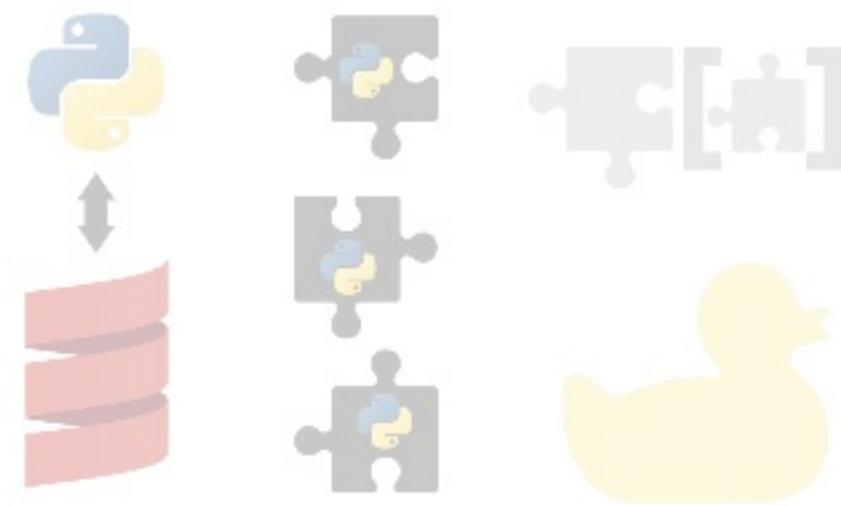


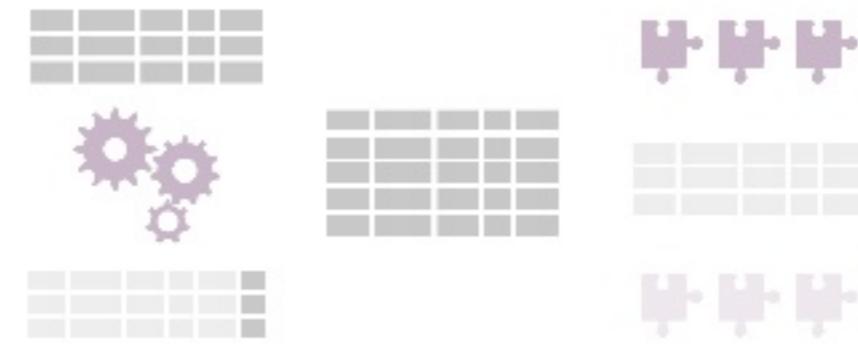
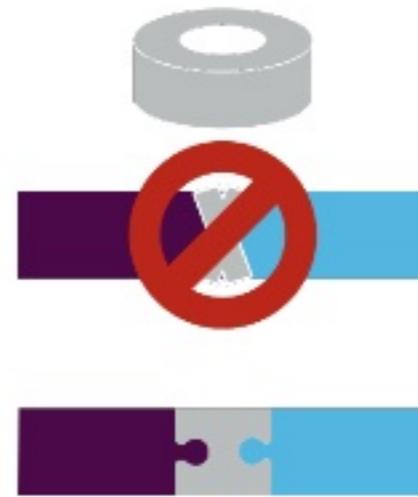
Conclusions and takeaways



estimator.fit(df)

model.transform(df)

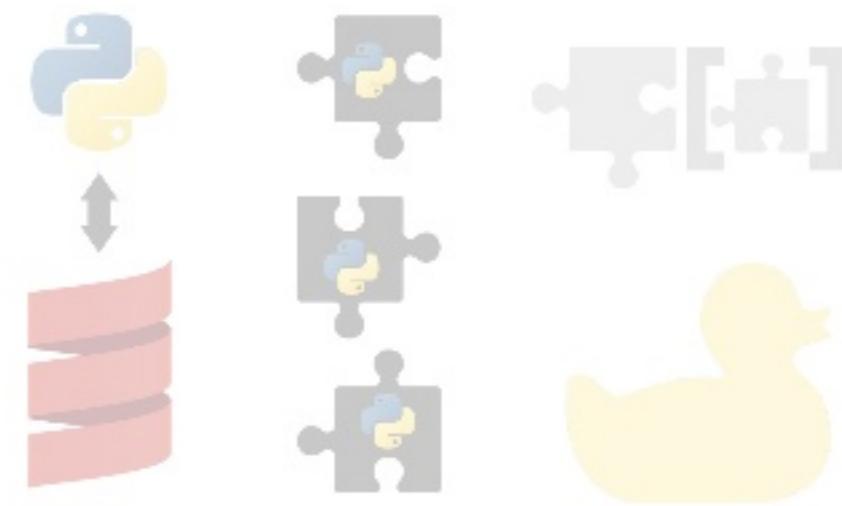


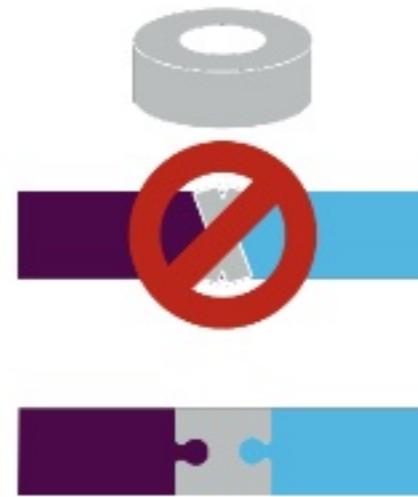


`estimator.fit(df)`



`model.transform(df)`

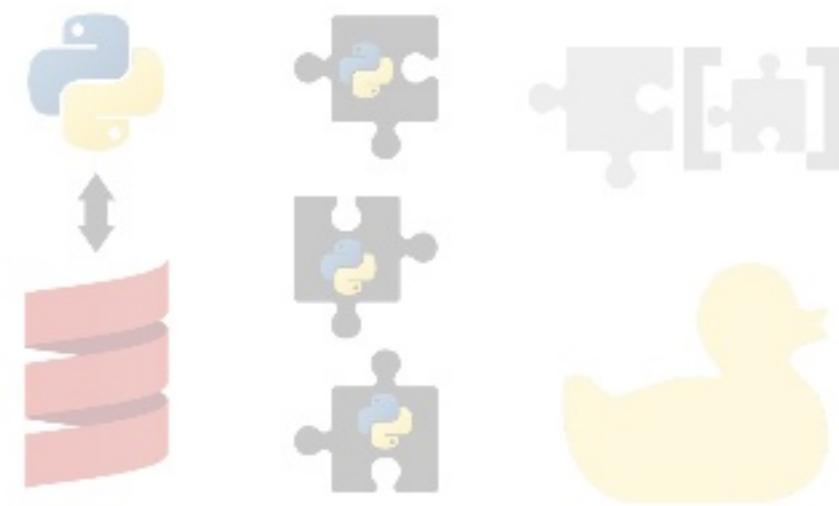


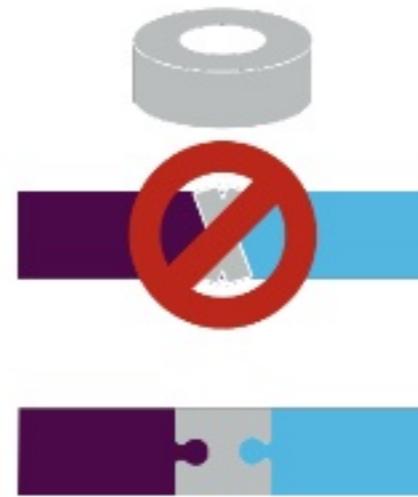


`estimator.fit(df)`



`model.transform(df)`

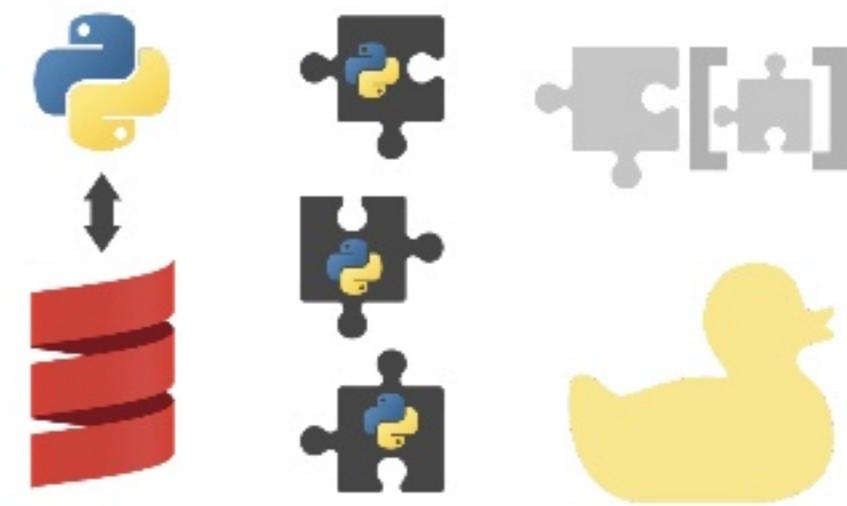


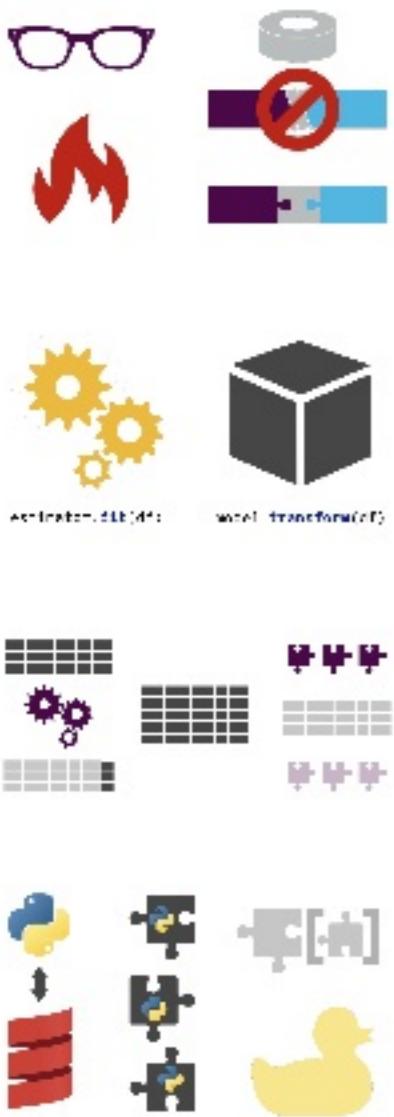


`estimator.fit(df)`



`model.transform(df)`





KEEP IN TOUCH

[https://radalytics.io](https://radanalytics.io)

eje@redhat.com • @manyangled

willb@redhat.com • @willb