

Programming with Scalding

Krishnan Raman

Data Scientist, Twitter

<http://github.com/krishnanraman/bigdata>

What is Scalding ?

- Primary Usecase: Hadoop Map-Reduce Jobs in Scala
- Abstract away low-level Hadoop + Cascading details (but can still get to them ...)
- Identical functions as Scala – map, flatMap, filter, project, takeWhile, groupBy....
- Work well with CT primitives – Twitter’s Bijection libraries (Monoids, Groups, Rings, ...)

Powered By

[New Page](#)[Edit Page](#)[Page History](#)

Want to be added to this page? Send a tweet to [@scalding](#) or open an issue.

Company	Scalding Use Case	Code
Twitter	We use Scalding often, for everything from custom ad targeting algorithms, market insight, click prediction, traffic quality to PageRank on the Twitter graph. We hope you will use it too!	-
Etsy	We're starting to use Scalding alongside the JRuby Cascading stack described here . More to come as we use it further.	-
eBay	We use Scalding in our Search organization for ad-hoc data analysis jobs as well as more mature data pipelines that feed our production systems.	-
Gatling	We've just rebuilt our reports generation module on top of Scalding. Handy API on top of an efficient engine.	GitHub
Sonar	Our platform is built on Hadoop, Scalding, Cassandra and Storm. See Sonar's job listings .	-



Follow [@Scalding](#) on Twitter for updates.

Authors:

- Avi Bryant <http://twitter.com/avibryant>
- Oscar Boykin <http://twitter.com/posco>
- Argyris Zymnis <http://twitter.com/argyris>

Thanks for assistance and contributions:

- Sam Ritchie <http://twitter.com/sritchie>
- Aaron Siegel: <http://twitter.com/asiegel>
- Brad Greenlee: <http://twitter.com/bggreenlee>
- Edwin Chen <http://twitter.com/edchedch>
- Arkajit Dey: <http://twitter.com/arkajit>
- Krishnan Raman: http://twitter.com/dxbydt_jasq
- Flavian Vasile <http://twitter.com/flavianv>
- Chris Wensel <http://twitter.com/cwensel>
- Ning Liang <http://twitter.com/ningliang>
- Dmitriy Ryaboy <http://twitter.com/squarecog>
- Dong Wang <http://twitter.com/dongwang218>

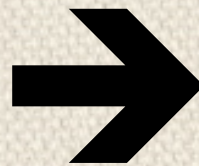
License

Copyright 2013 Twitter, Inc.

ssql : Gentle Intro to Scalding

```
open employees file1.txt
columns 3
column 1 name text
column 2 age number
column 3 income decimal
rows select age > 20
rows select age < 27
column add newincome income*1.1 + (age-20)*100
column remove age
column remove income
save file2.txt
```

```
fred 25 10000
sam 25 150000
oscar 39 200000
baker 27 234567
pascal 24 123456
ford 22 55000
haskell 56 45466
```



```
fred 11500.0
sam 165500.0
pascal 136201.6
ford 60700.0
```


ssql : Gentle Intro to Scalding

open employees file1.txt
columns 3
column 1 name text
column 2 age number
column 3 income decimal
rows select age > 20
rows select age < 27
column add newincome income*1.1 + (age-20)*100
column remove age
column remove income
save file2.txt

somescript

ssql

somescript.scala

```
import com.twitter.scalding._
class somescript(args : Args) extends Job(args){
  val employees =
    TextLine("file1.txt")
    .read
    .mapTo('line -> ('name,'age,'income) ){
      line:String =>
        val res = line.split(" ").map( _.trim).filter(_.length > 0)
        val name:String = res(0).toString
        val age:Int = res(1).toInt
        val income:Double = res(2).toDouble
        (name,age,income)
    }.filter('age) {
      age:Int =>
        (age > 20)
    }.filter('age) {
      age:Int =>
        (age < 27)
    }.map(('income,'age) -> ('newincome)) {
      columns:(Double,Int)=>
        val (income,age) = columns
        income*1.1 + (age-20)*100
    }.discard('age)
    .discard('income)
    .write(Tsv("file2.txt"))
}
```



```
1 import com.twitter.scalding._
2 class somescript(args : Args) extends Job(args){
3     val employees =
4         TextLine("file1.txt")
5         .read
6         .mapTo('line -> ('name, 'age, 'income) ){
7             line:String =>
8             val res = line.split(" ").map( _.trim).filter(_.length > 0)
9             val name:String = res(0).toString
10            val age:Int = res(1).toInt
11            val income:Double = res(2).toDouble
12            (name,age,income)
13        }.filter('age) {
14            age:Int =>
15            (age > 20)
16        }.filter('age) {
17            age:Int =>
18            (age < 27)
19        }.map(('income, 'age) -> ('newincome)) {
20            columns:(Double,Int)=>
21            val (income,age) = columns
22            income*1.1 + (age-20)*100
23        }.discard('age)
24        .discard('income)
25        .write(Tsv("file2.txt"))
26 }
```


Portfolio Management with Scalding

Divide \$1000 among 4 stocks –

1. Kroger, \$27
2. Abbott Labs, \$64
3. Dollar Tree, \$41
4. Monster Beverage, \$52

in the “best” possible way.

These are 4 stocks with a low beta (risk) and high average volume.

$$\text{Solve: } 27x + 64y + 41z + 52w = 1000$$


```
1  import com.twitter.scalding._
2  import com.twitter.scalding.mathematics.Combinatorics
3
4  class Portfolios(args : Args) extends Job(args) {
5
6      val cash = 1000.0 // money at hand
7      val error = 1 // its ok if we cannot invest the last dollar
8      val (kr,abt,dltr,mnst) = (27.0,64.0,41.0,52.0) // share prices
9      val stocks = IndexedSeq( kr,abt,dltr,mnst)
10
11      Combinatorics.weightedSum( stocks, cash,error).write( Tsv("invest.txt"))
12  }
```



```
1 import cern.colt.matrix.{DoubleFactory2D, DoubleFactory1D }
2 import cern.colt.matrix.linalg.Algebra
3 import java.util.StringTokenizer
4
5 object BestPortfolio {
6   def main(args:Array[String]) = {
7     val alg = Algebra.DEFAULT
8
9     // define the correlation matrix
10    val data = Array(Array(0.448, 0.177, 0.0, 0.017),
11                      Array(0.177, 0.393, 0.177, 0.237),
12                      Array(0.0, 0.177, 0.237, 0.06),
13                      Array(0.017, 0.237, 0.06, 0.19))
14    val corr = DoubleFactory2D.dense.make(data)
15    val file = scala.io.Source.fromFile("invest.txt").getLines.toList
16
17    // convert the tab-delimited weights in the file to a row vector
18    def getWeights(s:String) = {
19      val weights = s.split("\t").map( x=> x.toDouble)
20      DoubleFactory1D.dense.make(weights)
21    }
22    // compute risks per tuple and sort by risk
23    file.map(line=> {
24      val w = getWeights(line)
25      ( line, alg.mult( alg.mult(corr,w), w))
26    }).sortBy(x=>x._2)
27    .foreach( println )
28  }
29 }
```


weighted sums...

```
1 $ cat invest.txt
2 0 0 13 9
3 0 1 0 18
4 0 1 19 3
5 0 2 1 16
6 0 2 20 1
7 0 3 7 10
8 0 4 8 8
9 0 5 9 6
10 0 6 15 0
11 0 8 3 7
12 0 9 4 5
13 0 10 5 3
14 0 14 0 2
15 1 0 6 14
16 1 1 12 8
17 1 2 13 6
18 1 3 0 15
19 1 3 14 4
20 1 4 1 13
21 1 5 2 11
22 1 6 8 5
23 1 7 9 3
24 1 8 10 1
25 1 11 4 2
26 1 12 5 0
27 2 0 18 4
28 2 1 5 13
29 2 1 19 2
30 2 2 6 11
31 2 3 7 9
32 2 4 13 3
33 2 5 14 1
34 2 6 1 10
35 2 7 2 8
36 2 8 3 6
37 2 9 9 0
38 3 0 11 9
39 3 1 12 7
40 3 2 18 1
41 3 3 0 14
42 3 4 6 8
43 3 5 7 6
44 3 6 8 4
45 3 9 2 5
46 3 10 3 3
47 3 11 4 1
48 4 0 4 14
```

```
50 4 2 11 6
51 4 3 12 4
52 4 4 13 2
53 4 5 0 11
54 4 6 1 9
55 4 7 7 3
56 4 8 8 1
57 4 12 3 0
58 5 0 16 4
59 5 1 17 2
60 5 2 4 11
61 5 3 5 9
62 5 4 6 7
63 5 5 12 1
64 5 7 0 8
65 5 8 1 6
66 5 9 2 4
67 6 0 9 9
68 6 1 10 7
69 6 2 11 5
70 6 5 5 6
71 6 6 6 4
72 6 7 7 2
73 6 10 1 3
74 6 11 2 1
75 7 0 2 14
76 7 0 16 3
77 7 1 3 12
78 7 2 4 10
79 7 3 10 4
80 7 4 11 2
81 7 7 0 7
82 7 8 6 1
83 8 0 9 8
84 8 1 15 2
85 8 2 16 0
86 8 3 3 9
87 8 4 4 7
88 8 5 5 5
89 8 9 0 4
90 8 10 1 2
91 9 0 2 13
92 9 1 8 7
93 9 2 9 5
94 9 3 10 3
95 9 6 4 4
96 9 7 5 2
97 9 11 0 1
98 10 0 14 3
```

```
100 10 1 15 1
101 10 2 2 10
102 10 3 3 8
103 10 4 9 2
104 10 5 10 0
105 11 0 7 8
106 11 1 8 6
107 11 2 14 0
108 11 4 2 7
109 11 5 3 5
110 11 6 4 3
111 12 0 0 13
112 12 1 1 11
113 12 2 7 5
114 12 3 8 3
115 12 4 9 1
116 12 7 3 2
117 12 8 4 0
118 13 0 12 3
119 13 1 13 1
120 13 2 0 10
121 13 3 1 8
122 13 4 2 6
123 13 5 8 0
124 14 0 5 8
125 14 1 6 6
126 14 2 7 4
127 14 5 1 5
128 14 6 2 3
129 14 7 3 1
130 15 0 12 2
131 15 2 0 9
132 15 3 6 3
133 15 4 7 1
134 15 8 2 0
135 16 0 5 7
136 16 1 11 1
137 16 4 0 6
138 16 5 1 4
139 17 1 4 6
140 17 2 5 4
141 17 3 6 2
142 17 6 0 3
143 17 7 1 1
144 18 0 10 2
145 18 4 5 1
146 19 0 3 7
147 19 1 4 5
```

```
148 19 6 0 2
149 20 2 3 4
150 20 3 4 2
151 21 0 8 2
152 21 1 9 0
153 22 0 1 7
154 22 1 2 5
155 22 2 3 3
156 23 0 8 1
157 23 3 2 2
158 23 4 3 0
159 24 0 1 6
160 24 1 7 0
161 25 1 0 5
162 25 2 1 3
163 25 3 2 1
164 26 0 6 1
165 26 4 1 0
166 27 1 0 4
167 28 3 0 1
168 29 0 4 1
169 34 0 2 0
170 37 0 0 0
```

So it looks like we have 169 unique ways to divvy up the cash.

Lets take one of the 169 tuples. How about (5,4,6,7) ?

So 5 shares of Kroger @ \$27 = \$135

4 shares of Abbott Labs @ \$64 = \$256

6 shares of Dollar Tree @ \$41 = \$246

7 shares of Monster @ \$52 = \$468

$135+256+246+468 = \$1001$

But is (5,4,6,7) the best-possible tuple among the 169 choices ?

Correlations among equities...

4x4 correlation matrix

KR	ABT	DLTR	MNST
0.448	0.177	0	0.017
0.177	0.393	0.177	0.237
0	0.177	0.237	0.06
0.017	0.237	0.06	0.19

Given a vector of weights W , the net risk imposed by this matrix = $W'CW$

Goal: Want the 4-tuple with the least risk.

The 4-tuple with the least risk...

```
// compute risks per tuple and sort by risk
```

```
file.map(line=> {  
    val w = getWeights(line)  
    ( line, alg.mult( alg.mult(corr,w), w))  
}).sortBy(x=>x._2)
```

```
res11: List[(String, Double)] = List(  
  (1 0 6 14,56.776),  
  (4 0 4 14,56.824),  
  (2 1 5 13,57.544),  
  (4 1 5 12,58.552),  
  (2 2 6 11,59.646),  
  .....
```



"Best" tuple

"Worst" tuple



Under the hood...weightedSum(..)

```
1 def weightedSum( weights:IndexedSeq[Double], result:Double, error:Double)(implicit flowDef:FlowDef):Pipe = {
2   val numWeights = weights.size
3   val allColumns = (1 to numWeights).map( x=> Symbol("k"+x))
4
5   // create as many single-column pipes as the number of weights
6   val pipes = allColumns.zip(weights).map( x=> {
7     val (name,wt) = x
8     IterableSource( (0.0 to result by wt), name).read
9   }).zip( allColumns )
10
11   val first = pipes.head
12   val accum = (first._1, List[Symbol](first._2))
13   val rest = pipes.tail
14
15   val res = rest.foldLeft(accum)((a,b)=>{
16
17     val (apipe, aname) = a
18     val (bpipe, bname) = b
19     val allc = (List(aname)).flatten ++ List[Symbol](bname)
20     // Cross two pipes, Create a temp column that stores intermediate results
21     // Apply progressive filtering on the temp column
22     // Once all pipes are crossed, test for temp column within error bounds of result
23     ( apipe.crossWithSmaller(bpipe)
24       .map(allc->'temp){
25         x:TupleEntry =>
26         val values = (0 until allc.size).map( i=> x.getDouble( i.asInstanceOf[java.lang.Integer]))
27         values.sum
28       }.filter('temp){
29         x:Double => if( allc.size == numWeights) (math.abs(x-result)<= error) else (x <= result)
30       }.discard('temp), allc )
31   })._1.unique(allColumns)
32
33   (1 to numWeights).zip(weights).foldLeft( res) ((a,b) => {
34     val (num,wt) = b
35     val myname = Symbol("k"+num)
36     a.map( myname->myname){ x:Int => (x/wt).toInt }
37   })
38 }
```


[seer.cancer.gov/popdata/download.html](#)

scala scalding scalding hit cascading jdk email colt cal science cgkit cost of production

Datasets & Software

Datasets

- + [SEER Data 1973-2009](#)
- + [Standard Population Data](#)
- [US Mortality Data](#)
- [US Population Data](#)
 - [Download US Population Data](#)
 - [Single Year of Age County Population Estimates](#)
 - [Estimates Used in SEER*Stat Software](#)
 - [Intercensal Population Impact on Rates](#)
 - [End of Decade Corrections](#)
 - [Previous Data Release](#)
- [County Attributes](#)
- [SEER Linked Databases](#)
- + [Specialized SEER*Stat Datasets](#)

Statistical Software

- + [SEER*Stat](#)
- + [SEER*Prep](#)
- + [HD*Calc](#)

Download US Population Data - 1969-2011

Note: All of the populations change annually going back to the last decennial. [Used in NCI's SEER*Stat Software.](#)

Files containing US population data are provided for analysis using your own [SEER*Stat](#) (you must use the [SEER*Prep](#) software to create a database for analysis).

Populations are available for the following time periods and "races":

- 1969-2011 County-level: White, Black, Other;
- 1990-2011 County-level: Expanded Races (White, Black, American Indian or Alaska Native, Hispanic or Latino, Non-Hispanic) by Origin (Hispanic, Non-Hispanic);
- 1981-2011 State-level: Expanded Races (White, Black, American Indian or Alaska Native, Hispanic or Latino, Non-Hispanic). See [State-Level Race and Hispanic Origin Population Estimates](#).
- 1981-2011 State-level: White and Non-White by Hispanic Origin

County- and state-level population files with 19 age groups (<1, 1-4, ..., 80+) and 85+ age groups (<1, 1, 2, ..., 84, 85+) are provided below.

For the All States Combined and four states (Alabama, Louisiana, Mississippi, Texas), population estimates are available for 2005: the standard set based on July 1, 2005 has been adjusted for the population shifts due to hurricanes Katrina (August 29) and Rita (September 24). For more information, see [Adjusted Populations for the Counties/Parishes Affected by Hurricanes Katrina and Rita](#).

Data Files

The data are stored in text files and provided here as Windows self-extracting archives.

- ☐ [SEER Data 1973-2009](#)
- ☐ [Standard Population Data](#)
 - [US Mortality Data](#)
- ☐ [US Population Data](#)
 - [Download US Population Data](#)
 - [Single Year of Age County Population Estimates](#)
 - [Estimates Used in SEER*Stat Software](#)
 - [Intercensal Population Impact on Rates](#)
 - [End of Decade Corrections](#)
 - [Previous Data Release](#)

- [SEER*Stat](#)
- [SEER*Prep](#)
- [HD*Calc](#)

Note: All of the populations change annually going back to the last decennial
Used in NCI's SEER*Stat Software.

Files containing US population data are provided for analysis using your own [SEER*Stat](#) (you must use the [SEER*Prep](#) software to create a database for

Populations are available for the following time periods and "races":

- 1969-2011 County-level: White, Black, Other;
- 1990-2011 County-level: Expanded Races (White, Black, American Indian or Alaska Native, Asian, Pacific Islander) by Origin (Hispanic, Non-Hispanic);
- 1981-2011 State-level: Expanded Races (White, Black, American Indian or Alaska Native, Asian, Pacific Islander). See [State-Level Race and Hispanic Origin Population Estimates](#)
- 1981-2011 State-level: White and Non-White by Hispanic Origin

County- and state-level population files with 19 age groups (<1, 1-4, ..., 80- age groups (<1, 1, 2, ..., 84, 85+) are provided below.

For the All States Combined and four states (Alabama, Louisiana, Mississippi) population estimates are available for 2005: the standard set based on July 2005 has been adjusted for the population shifts due to hurricanes Katrina (August 29, 2005) and Rita (September 24, 2005). For more information, see [Adjusted Populations for the Counties/Parishes Affected by Hurricanes Katrina and Rita](#).

Data Files

The data are stored in text files and provided here as Windows self-extracting

Population Data Dictionary

Population Data Dictionary - 1969-2011

File Format

- Fixed length ASCII text records (28 bytes)
- One population per record/line
- Windows (CR/LF) line delimiters
- All numeric data is zero filled to the left

Variable Name and Values	Start Column	Length	Data Type
Year (1969, 1970, 1971...)	1	4	numeric
State postal abbreviation "KR" is used for the dummy state created to represent hurricane Katrina/Rita evacuees	5	2	character
State FIPS code Field is 9-filled for dummy state created to represent hurricane Katrina/Rita evacuees	7	2	numeric
County FIPS code Field is 9-filled for dummy state created to represent hurricane Katrina/Rita evacuees	9	3	numeric
Registry 01 = San Francisco-Oakland SMSA 02 = Connecticut 20 = Detroit (Metropolitan) 21 = Hawaii 22 = Iowa 23 = New Mexico 25 = Seattle (Puget Sound) 26 = Utah 27 = Atlanta (Metropolitan) 29 = Alaska Natives	12	2	numeric

Population data – 12 million records

1	1969AL01001991910000000159
2	1969AL01001991910100000657
3	1969AL01001991910200001137
4	1969AL01001991910300000956
5	1969AL01001991910400000721
6	1969AL01001991910500000424
7	1969AL01001991910600000585
8	1969AL01001991910700000637
9	1969AL01001991910800000607
10	1969AL01001991910900000523
11	1969AL01001991911000000466
12	1969AL01001991911100000386
13	1969AL01001991911200000357
14	1969AL01001991911300000268
15	1969AL01001991911400000179
16	1969AL01001991911500000126
17	1969AL01001991911600000076
18	1969AL01001991911700000049
19	1969AL01001991911800000031
20	1969AL01001991920000000151
21	1969AL01001991920100000566
22	1969AL01001991920200001037
23	1969AL01001991920300000979
24	1969AL01001991920400000683
25	1969AL01001991920500000557
26	1969AL01001991920600000705
27	1969AL01001991920700000669
28	1969AL01001991920800000614
29	1969AL01001991920900000496
30	1969AL01001991921000000395
31	1969AL01001991921100000372
32	1969AL01001991921200000377
33	1969AL01001991921300000320
34	1969AL01001991921400000240

US FIPS Codes.csv (3000+ records)

1	Alabama,Autauga,01,001
2	Alabama,Baldwin,01,003
3	Alabama,Barbour,01,005
4	Alabama,Bibb,01,007
5	Alabama,Blount,01,009
6	Alabama,Bullock,01,011
7	Alabama,Butler,01,013
8	Alabama,Calhoun,01,015
9	Alabama,Chambers,01,017
10	Alabama,Cherokee,01,019
11	Alabama,Chilton,01,021
12	Alabama,Choctaw,01,023
13	Alabama,Clarke,01,025
14	Alabama,Clay,01,027
15	Alabama,Cleburne,01,029
16	Alabama,Coffee,01,031
17	Alabama,Colbert,01,033
18	Alabama,Conecuh,01,035
19	Alabama,Coosa,01,037
20	Alabama,Covington,01,039
21	Alabama,Crenshaw,01,041
22	Alabama,Cullman,01,043
23	Alabama,Dale,01,045
24	Alabama,Dallas,01,047
25	Alabama,De Kalb,01,049
26	Alabama,Elmore,01,051
27	Alabama,Escambia,01,053
28	Alabama,Etowah,01,055
29	Alabama,Fayette,01,057
30	Alabama,Franklin,01,059
31	Alabama,Geneva,01,061
32	Alabama,Greene,01,063

Convert pop.txt to a legit Scalding Source

```
44 case class USPopulationSource(override val p:String) extends TextLine(p) {
45   override def transformForRead(pipe : Pipe) = {
46     import Dsl._
47     RichPipe(pipe).mapTo('line->( 'year,'state,'fips,'isWhite,'isBlack,
48                                   'isHispanic,'isMale,'age,'population)) {
49       record:String =>
50         val year:Int          = record.slice(0,0+4).toInt
51         val state:String      = record.slice(4,4+2)
52         val fips:String       = record.slice(6,6+5)
53         val isWhite:Boolean   = record.slice(13,13+1).toInt == 1
54         val isBlack:Boolean   = record.slice(13,13+1).toInt == 2
55         val isHispanic:Boolean = record.slice(14,14+1).toInt == 1
56         val isMale:Boolean    = record.slice(15,15+1).toInt == 1
57         val age:Int           = 5*(record.slice(16,16+2).toInt -1)
58         val population:Int     = record.slice(18,18+8).toInt
59
60         (year,state,fips,isWhite,isBlack, isHispanic,isMale,age,population)
61     }
62   }
63 }
```



```

65 Find THE FASTEST GROWING COUNTY IN THE UNITED STATES over the 1969-2011 timeframe
66
67 class PopulationStats(args:Args) extends Job(args) {
68   val people = USPopulationSource("pop.txt").read
69   val fipspipe = TextLine("US_FIPS_Codes.csv").read.mapTo('line->('state,'county,'fips)) {
70     line:String =>
71     var arr = line.split(",")
72     (arr(0),arr(1),(arr(2)+arr(3)))
73   }
74   people.groupBy('year, 'fips){
75     group => group.plus[Int]('population->'population)
76   }.groupBy('fips) {
77     val init = (0,0.0d)
78     type X = (Int,Double)
79     type T = (Int,Int)
80
81     // foldLeft[X,T](fieldDef : (Fields,Fields))(init : X)(fn : (X,T) => X)
82     group => group.foldLeft[X,T]( ('population,'year) -> ('dummy,'growth))(init:X) {
83
84       (x:X, t:T) =>
85       val (population,year) = t
86       val (dummy, growth ) = x
87       year match {
88         case 1969 => (population, 0.0d)
89         case 2011 => (population,(population-dummy)/(dummy+0.0d))
90         case _ => if (dummy==0) (population,0.0d) else (dummy,0.0d)
91       }
92     }
93   }.project('fips,'growth)
94   .joinWithSmaller(('fips-> 'fips), fipspipe)
95   .groupAll(_.sortBy('growth))
96   .write(Tsv("growth.txt"))
97 }

```


RESULTS:

13053	-0.5959766162310867	Georgia	Chattahoochee
38083	-0.5825892857142857	North Dakota	Sheridan
38013	-0.575839766325892	North Dakota	Burke
54047	-0.5703608502224419	West Virginia	McDowell
48101	-0.5474033816425121	Texas	Cottle
38047	-0.5436781609195402	North Dakota	Logan
38023	-0.5362287210824968	North Dakota	Divide
38087	-0.5334632878492528	North Dakota	Slope
38037	-0.5310054184226369	North Dakota	Grant
48301	-0.5276381909547738	Texas	Loving

---- some 3000 counties here ----

13117	10.637015231025215	Georgia	Forsyth
12097	10.880533448053345	Florida	Osceola
48397	10.907133440749963	Texas	Rockwall
13135	11.063892016788289	Georgia	Gwinnett
48157	11.168507788849015	Texas	Fort Bend
48491	11.323119312014695	Texas	Williamson
08117	11.39344262295082	Colorado	Summit
48085	11.507137247655564	Texas	Collin
12035	21.135939986360537	Florida	Flagler
08035	36.49576488706366	Colorado	Douglas

CONCLUSION:

Over the 1969-2011 timeframe, Douglas County experienced the highest population growth (3600%)
 The top-10 counties with the most pop growth are shown above - they are in Texas, Florida & Colorado.


```

99 Goal: You are a male between the age of (20,40).
100 Where should you go to maximize your chances of hooking up with the opposite sex ?
101
102 Scalding code: sexratio(2011,(20,40),people,fipspipe)
103
104 where sexratio as below:
105
106 def sexratio(baseYear:Int, agerange:(Int,Int), people:Pipe, fipspipe:Pipe) = {
107     people.filter('year,'age) {
108         yearage:(Int,Int) =>
109         val (year,age) = yearage
110         (year == baseYear) && (age>=agerange._1) && (age<=agerange._2) // marriageable
111     }.groupBy('fips, 'isMale){
112         group => group.plus[Int]('population->'population)
113     }.groupBy('fips) {
114         val init = 0.0d
115         type X = Double
116         type T = (Int,Boolean)
117         // foldLeft[X,T](fieldDef : (Fields,Fields))(init : X)(fn : (X,T) => X)
118         group => group.foldLeft[X,T]( ('population,'isMale) -> ('sexratio))(init:X) {
119             (x:X, t:T) =>
120             val ratio = x
121             val (population,isMale) = t
122             if(ratio==0.0) (population+0.0d)
123             else {
124                 val sexratio = ratio/(population+0.0d)
125                 if (isMale) 1.0/sexratio else sexratio
126             }
127         }
128     }.joinWithSmaller(('fips-> 'fips), fipspipe)
129     .project('sexratio, 'state,'county)
130     .groupAll(_.sortBy('sexratio))
131     .write(Tsv("sexratio"+baseYear+".txt"))
132 }

```


SORTED RESULTS:

SEX RATIO (M/F), State, County

0.5809725158562368	Georgia	Pulaski
0.6842327150084317	West Virginia	Summers
0.6924882629107981	Wyoming	Niobrara
0.7	Missouri	Audrain
0.7070063694267515	Indiana	Parke
0.7194008037997808	Missouri	Livingston
0.7434895833333334	Virginia	Fluvanna
0.764525993883792	North Dakota	Hettinger
0.7655172413793103	Texas	Oldham
0.7692307692307692	Texas	Loving
0.7831548198636806	Ohio	Union

....3000 more counties here

3.4205729166666665	Georgia	Wheeler
3.56390977443609	Colorado	Bent
4.472380952380952	Pennsylvania	Forest
4.957894736842105	Texas	Concho
5.52020202020202	Colorado	Crowley

TL:DR;

Pulaski County, Georgia has the most advantageous sex ratio of **0.58** (**Male/Female** in the **(20,40)** agegroup)

~ **Two** women **for** every man!


```

134 Goal: Race Analysis: find the county with the highest ratio of whites
135 to the total population. Do the same for African Americans.
136
137 /*
138  Want: white people in an agegroup
139 */
140 race(2011,(20,40),people,fipspipe,true)
141
142 /*
143  Want: black people in an agegroup
144 */
145 race(2011,(20,40),people,fipspipe,false)
146
147
148 def race(baseYear:Int, agerange:(Int,Int), people:Pipe, fipspipe:Pipe, white:Boolean) = {
149   val filename = (if (white) "white" else "black")+baseYear+".txt"
150
151   // first get total population ( ie. all races ) of agegroup for baseYear
152   val pipe1 = people.filter('year,'age){
153     yearage:(Int,Int) =>
154     val (year,age) = yearage
155     (year == baseYear) && (age>=agerange._1) && (age<=agerange._2)
156   }.groupBy('isMale,'fips){
157     // add (fe)male blacks, (fe)male hispanics, (fe)male whites etc.
158     group => group.plus[Int]('population->'population)
159   }.groupBy('fips){
160     group => group.plus[Int]('population->'population) // add males + females
161   }

```



```

163 // now count people of certain race in certain age group
164 val pipe2 = people.filter('year','age, 'isWhite','isBlack) {
165   yearagewb:(Int,Int,Boolean,Boolean) =>
166   val (year,age,isWhite,isBlack) = yearagewb
167   (year == baseYear) && (age>=agerange._1) && (age<=agerange._2) &&
168   (if (white) isWhite else isBlack)
169 }.groupBy('isMale','fips){
170   // add (fe)male blacks, (fe)male hispanics, (fe)male whites etc.
171   group => group.plus[Int]('population->'population)
172 }.groupBy('fips){
173   group => group.plus[Int]('population->'mypop) // add males + females
174 }
175
176 // now compute ratio & sort by ratio
177 val pipe3 = pipe1.joinWithSmaller(('fips-> 'fips), pipe2)
178 .map( ('mypop, 'population)->'raceratio){
179   mypop_total:(Int,Int) =>
180   val( mypop, population) = mypop_total
181   mypop/(population+0.0d)
182 }
183 .project('population, 'raceratio,'fips)
184
185 // now join with verbose fips
186 pipe3.joinWithSmaller(('fips-> 'fips), fipspipe)
187 .project('raceratio,'population, 'state,'county)
188 .groupAll(_.sortBy('raceratio))
189 .write(Tsv(filename))
190 }

```



```
194 Black:
195 .....
196
197 0.7854832914121451 2783 Mississippi Humphreys
198 0.7856462369006033 3149 Alabama Wilcox
199 0.7889369592088998 3236 Alabama Lowndes
200 0.7913015254787407 3081 Georgia Hancock
201 0.794999045619393 10478 Virginia Petersburg City
202 0.799015275849009 7921 Mississippi Coahoma
203 0.8423519386450788 2347 Alabama Greene
204 0.8453242514545196 7047 Alabama Macon
205 0.87166636835511 5587 Mississippi Holmes
206 0.8917171717171717 2475 Mississippi Jefferson
207
208 White
209 .....
210
211 0.9939713639788997 2654 West Virginia Clay
212 0.9940298507462687 670 Nebraska Frontier
213 0.9946686976389947 1313 Tennessee Pickett
214 0.9950510361892979 6466 West Virginia Lincoln
215 0.9955357142857143 672 Nebraska Sherman
216 0.9961240310077519 1032 Nebraska Chase
217 0.9964850615114236 569 Nebraska Greeley
218 1.0 93 Montana Petroleum
219 1.0 400 Nebraska Garfield
220 1.0 155 Nebraska Grant
221 1.0 152 Nebraska Keya Paha
222 1.0 116 Nebraska Loup
223
224
225 TL:DR; Racial makeup of youth in the 20-40 agegroup is:
226
227 100% white, in several counties of Nebraska - Loup, Keya Paha, Grant, Garfoeld etc.
228 In general, Nebraska has an overwhelming white population ( >99%) in several counties.
229
230 The racial makeup of Jefferson county in Mississippi is 90% African American.
231 In general, Mississippi & Alabama dominate AA rankings as seen above.
```


Algebird...Twitter's Abstract Algebra library

- Model a wide class of "reductions" as a sum on some iterator
- Average, moving average, max/min, set union, approximate set size (in much less memory with HyperLogLog), approximate item counting (using CountMinSketch).
- Implementations of Monoids for interesting approximation algorithms, such as Bloom filter, HyperLogLog and CountMinSketch.
- Use CMS, HLL within Storm (realtime aggregation) and Scalding (batchmode aggregation

Trivial Monoid: Employee Monoid

257 **USECASE:** Find the oldest (max age) employee with max tenure making the least salary

258

```
259 case class Monoid(e:Set[Employee]) {  
260   def plus(a:Employee,b:Employee)= {  
261     if(a.age > b.age) a else  
262     if(b.age > a.age) b else  
263     if(a.tenure > b.tenure) a else  
264     if(b.tenure > a.tenure) b else  
265     if (a.salary > b.salary) b else  
266     if(b.salary > a.salary) a else  
267     if( math.random > 0.5) a else b  
268   }  
269   val identity = Employee(0,0,0,"id")  
270   val theUnluckyOne = e.foldLeft(identity)((a,b)=>plus(a,b))  
271 }
```

272

```
273 scala> Monoid(empDB.toSet).theUnluckyOne  
274 res2: Employee = Employee(35,6,85079,emp7435)
```

275

Nontrivial example: Abelian Group

```
234 trait Abelian
235
236 case class Zn(order:Int, zero:Int) extends Abelian {
237     def identity = zero
238     def size = order
239     def elements = (1 to order).toSeq
240     def cayley = {
241         Vector.tabulate(order,order)((x,y)=> {
242             val idx = math.abs(zero - (x+1)) // difference between x & zero
243             val timesToShift = if ((x+1)>=zero) idx else (order-idx)
244             val row = shift(elements, timesToShift )
245             row(y)
246         })
247     }
248     private def shift(s:Seq[Int], pos:Int) = {
249         val n = 1+s.indexOf(pos)
250         s.slice(n,s.size)++s.slice(0,n)
251     }
252
253     def plus(x:Int, y:Int) = cayley(x-1)(y-1)
254     def inverse(x:Int) = cayley(x-1).indexOf(zero)+1
255 }
```


CMS Monoid for counting tweets

```
1 import com.twitter.algebird.{CountMinSketchMonoid, CMS}
2 object CMSTest {
3   // sane CMS builder
4   def mkCMS( probability:Byte = 90, EPS:Double = 0.01 ) = {
5     assert( probability > 0 && probability < 100 )
6     val DELTA = 1 - probability/100.0
7     val SEED = 1
8     new CountMinSketchMonoid(EPS, DELTA, SEED)
9   }
10
11   def updateCMS( monoid:CountMinSketchMonoid, cms:Option[CMS], x:Long ):Option[CMS] = {
12     val item = monoid.create(x)
13     if( cms.isDefined ) Some(item ++ cms.get) else Some(item)
14   }
15
16   def main(args:Array[String]) = {
17
18     // make items from 1 to n, assume k has been seen n-k times
19     val n = 100
20     val stream = List.tabulate[List[Long]](n)( k => {
21       val freq = n - k
22       List.fill[Long](freq)( k )
23     }).flatten
24
25     // assume this data comes in one at a time from some stream
26     // update the CMS as the data arrives
27     val monoid = mkCMS( args(0).toByte, args(1).toDouble )
28     val cmsitem:Option[CMS] = None
29     val result = stream.foldLeft(monoid, cmsitem)( (a,b) => {
30       val (mymonoid, myitem) = a
31       val number = b
32       val newitem = updateCMS( mymonoid, myitem, number )
33       val retval = (mymonoid, newitem)
34       retval
35     })._2.get
36
37     // compare estimated frequencies vs actuals
38     val estimate = result.frequency ( 34 ).estimate
39     val actual = n - 34
40     println( "Estimate: %d, Actual : %d" + estimate, actual)
41   }
42 }
```


More goodies...

- <https://github.com/twitter/scalding>
- <https://github.com/twitter/algebird>
- <https://github.com/twitter/bijections>

Thank you! (Krishnan Raman, kraman@twitter.com)