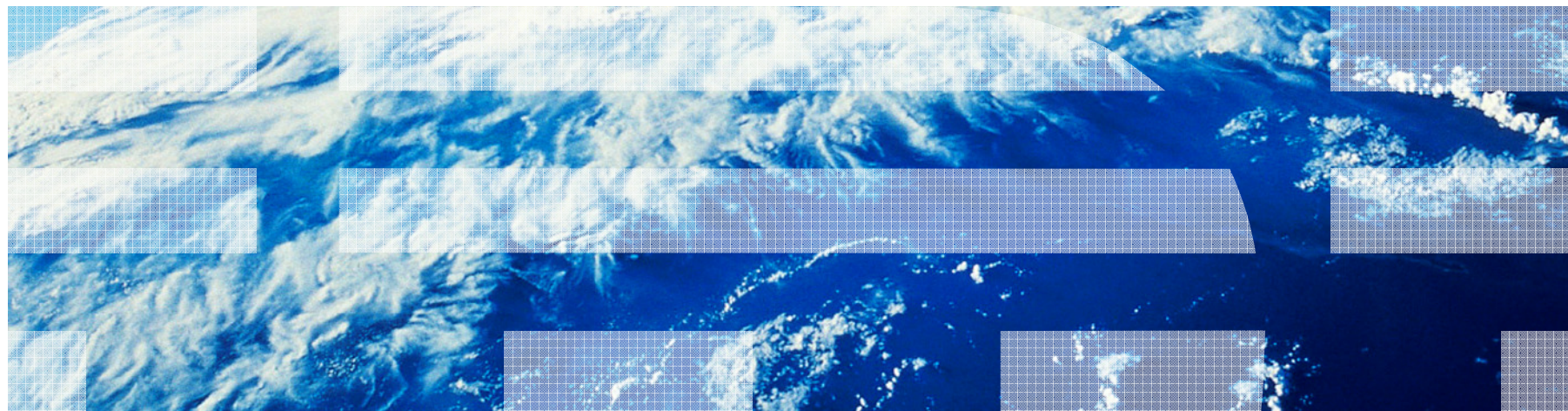

Scalable and Numerically Stable Descriptive Statistics in SystemML

Yuanyuan Tian, Shirish Tatikonda, Berthold Reinwald
IBM Almaden Research Center



SystemML

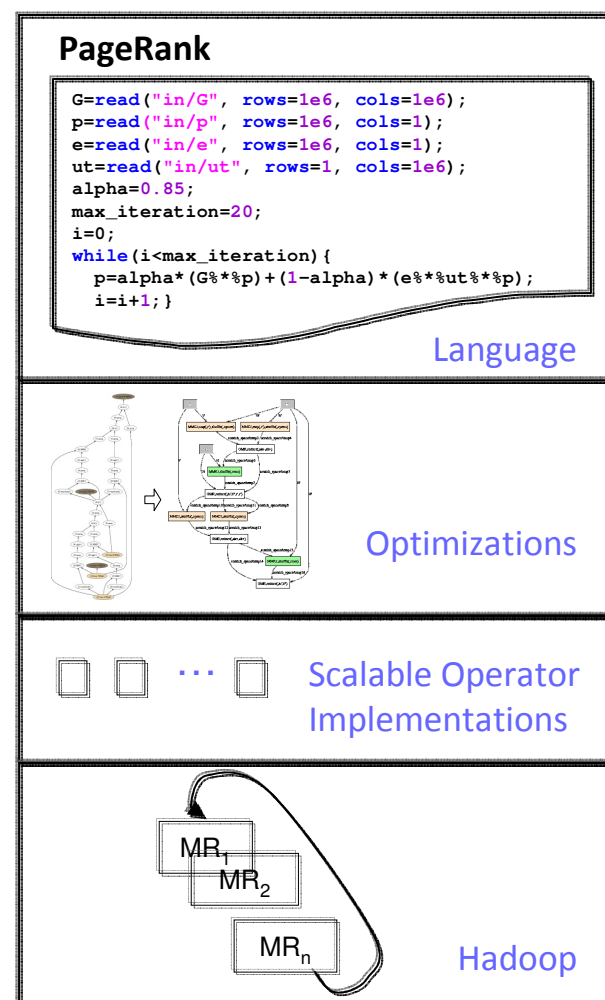
- ❑ Pervasive need to enable machine learning (ML) on massive datasets
- ❑ Increasing interest in implementing a ML algorithms on MapReduce
- ❑ Directly implementing ML algorithms on MapReduce is challenging
- ❑ Solution: **SystemML** – A Declarative Approach to Machine Learning on MapReduce

SystemML Overview

- ❑ **High-level language** with ML specific constructs
 - Syntax is similar to R and Matlab
 - Matrix, vector and scalar data types
 - Linear algebra and mathematics operators

- ❑ **Optimizations** based on data and system characteristics
 - Cost-based and rule-based optimization
 - Job generation heuristics

- ❑ **Scalable** implementations on Hadoop
 - Handle large sparse data sets
 - Parallelization is *transparent* to end users



ICDE 2011

Example ML Algorithms Supported in SystemML

- ❑ Classification: linear SVMs, logistic regression
- ❑ Regression: linear regression
- ❑ Matrix Factorization: NMF, SVD, PCA
- ❑ Clustering: k-means
- ❑ Ranking: PageRank, HITS
- ❑ Data Exploration: **Descriptive Statistics** ← **focus of this talk**
 - **Univariate Statistics:**
 - **Scale:** Sum, Mean, Harmonic mean, Geometric mean, Min, Max, Range, Median, Quantile, Inter-quartile-mean, Variance, Standard deviation, Coefficient of variance, Central moment, Skewness, Kurtosis, Standard error of mean, Standard error of skewness, Standard error of kurtosis
 - **Categorical:** Mode, Per-category frequencies
 - **Bivariate Statistics:**
 - **Scale-scale:** Covariance, Pearson correlation
 - **Scale-categorical:** Eta, ANOVA F measure
 - **Categorical-categorical:** Chi-square coefficient, Cramer's V, Spearman correlation

How hard is it?

- ❑ Seemingly trivial to implement in MapReduce
 - Most descriptive statistics can be written in certain summation form

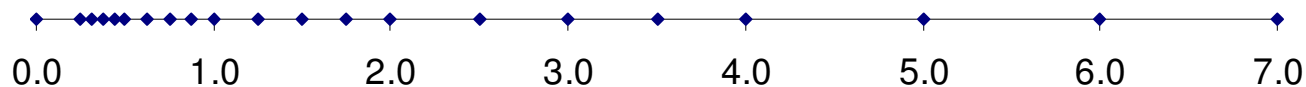
Sum	$s = \sum x_i$
Mean	$\bar{x} = \frac{1}{n} \sum x_i$
Variance	$s^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2$ $= \frac{1}{n-1} \sum x_i^2 - \frac{1}{n(n-1)} (\sum x_i)^2$

- ❑ **Pitfall:** these straight forward implementations can lead to disasters in *numerical accuracy*
- ❑ Problem gets **worse** with increasing volumes of data

Background: Floating Point Numbers

- ❑ Source of Inaccuracy: finite precision arithmetic for floating point numbers
- ❑ Floating point number system $F \subset \mathbb{R}$: $y = \pm\beta^e \times .d_1d_2 \dots d_t$ ($0 \leq d_i \leq \beta-1$)
 - *base* β ,
 - *precision* t ,
 - *exponent range* $e_{\min} \leq e \leq e_{\max}$
 - IEEE double: $\beta=2$, $t=53$, $e_{\min}=-1021$, $e_{\max}=1024$
- ❑ Floating point numbers are **not** equally spaced.
 - Example number system: $\beta = 2$, $t = 3$, $e_{\min} = -1$, and $e_{\max} = 3$

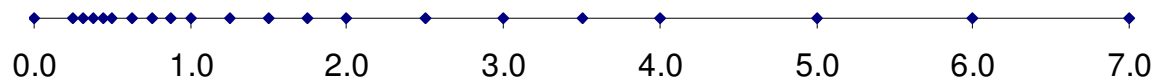
$F = \{ 0, \pm 0.25, \pm 0.3125, \pm 0.375, \pm 0.4375, \pm 0.5, \pm 0.625, \pm 0.75, \pm 0.875, \pm 1.0, \pm 1.25, \pm 1.5, \pm 1.75, \pm 2.0, \pm 2.5, \pm 3.0, \pm 3.5, \pm 4.0, \pm 5.0, \pm 6.0, \pm 7.0 \}$



Background: Numerical Accuracy

- Example number system: $\beta = 2$, $t = 3$, $e_{\min} = -1$, and $e_{\max} = 3$

$F = \{0, \pm 0.25, \pm 0.3125, \pm 0.375, \pm 0.4375, \pm 0.5, \pm 0.625, \pm 0.75, \pm 0.875, \pm 1.0, \pm 1.25, \pm 1.5, \pm 1.75, \pm 2.0, \pm 2.5, \pm 3.0, \pm 3.5, \pm 4.0, \pm 5.0, \pm 6.0, \pm 7.0\}$



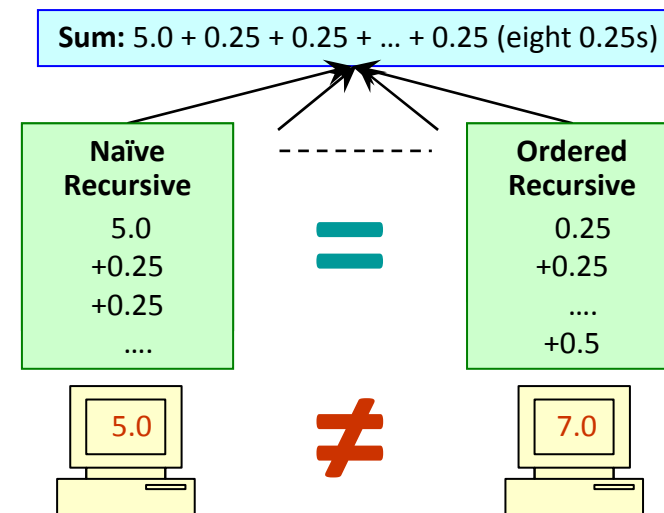
relative error = $|x-x'|/|x|$

- x : true value
- x' : value in float point

<ul style="list-style-type: none"> round off: 6.375 and 5.625 			
<i>true value:</i> 6.375 and 5.625	<i>computed:</i> 6.0 and 6.0	<i>relative error:</i> 0.059 and 0.067	
<ul style="list-style-type: none"> big number + many small numbers: $5.0 + 0.25 + 0.25 + \dots + 0.25$ (eight 0.25) \rightarrow <i>pitfall for summation</i> 			
<i>true value:</i> 7.0	<i>computed:</i> 5.0	<i>relative error:</i> 0.286	
	$5.0 + 0.25 = 5.25$ round to 5.0		
<ul style="list-style-type: none"> subtraction of 2 big and similar numbers: $6.375 - 5.625$ \rightarrow <i>catastrophic cancellation for subtraction</i> 			
<i>true value:</i> 0.75	<i>computed:</i> 0	<i>relative error:</i> 1.0	
	6.375 round to 6.0		
	5.625 round to 6.0		

Background: Numerical Stability

- Algebraically equivalent algorithms for the same calculation produce different results on digital computers
 - Some damp out errors
 - Some magnify errors



- *Numerical stability* is a desirable property of numerical algorithms
 - Algorithms that can be proven not to magnify approximation errors are called *numerically stable*

Importance of Numerical Stability

- ❑ Numerical stability issue has been largely ignored in big data processing
 - e.g. PIG and HIVE, are using well-known unstable algorithms for computing some basic statistics

- ❑ How about software floating point packages, e.g. BigDecimal?
 - Arbitrary precision, but *very* slow
 - +, -, *: **2 orders of magnitude** slower
 - /: **5 orders of magnitude** slower

- ❑ **Goal of this Talk:** share our experience on descriptive statistics algorithms for big data
 - **Scalable** – database people already understand
 - **Numerically Stable** – need more attention!!!!

Numerically Stable Summation

- ❑ **Naïve Recursive:** unstable
- ❑ **Sorted Recursive:** better for nonnegative but needs an expensive sort
- ❑ **Kahan:** efficient and stable [*Kahan 1965*]

- Recursive summation with a **correction term** to compensate rounding error

$$(s', c') = \text{KahanAdd}(s, c, a)$$

s/s': old/new sum

$$a' = a + c$$

c/c': old/new correct

$$s' = s + a'$$

a: number to add

$$c' = a' - (s' - s)$$

Stability Property: relative error bound independent of problem size n , when n is less than $O(10^{16})$ for IEEE doubles

- ❑ **MR Kahan in SystemML:**

- Mapper: apply Kahan to compute partial sum and correction
 - Reducer: apply Kahan on partial results to compute sum

$$(s, c) = \text{KahanAdd}(s_1, c_1 + c_2, s_2)$$

s₁/s₂: partial sums

c₁/c₂: partial corrects

s: total sum

c: total correct

Our Proof: relative error bound independent of problem size n , as long as each mapper processes less than $O(10^{16})$ numbers

Kahan: $5.0 + 0.25 + 0.25 + \dots + 0.25$

	sum	correction
5.0	5.0	0
+ 0.25		
5.0	5.0	0.25
+ 0.25		
6.0	6.0	-0.5
+ 0.25		
6.0	6.0	-0.25
+ 0.25		
6.0	6.0	0

.....

10 **MR Kahan scales to larger problem size with numerical accuracy.**

Numerically Stable Mean

❑ **Naïve sum/count:** unstable

❑ **Incremental:** stable [*Chan et al 1979*]

n_1, n_2 : partial counts, μ_1, μ_2 : partial means

$$n = n_1 + n_2$$

$$\delta = \mu_2 - \mu_1$$

$$\mu = \mu_1 + \delta n_2 / n$$

❑ **MR Incremental:** adapt Incremental to MapReduce

$$n = n_1 + n_2$$

$$\delta = \mu_2 - \mu_1$$

$$\mu = \mu_1 \boxplus \delta n_2 / n \quad (\boxplus \text{ -- KahanAdd, maintain a correction term for } \mu)$$

Numerically Stable Higher-Order Statistics

- ❑ **Higher order statistics:** central moment, variance, standard deviation, skewness, kurtosis (core: central moment $m_p = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^p$)
- ❑ **2-Pass:** 1st pass to compute mean, 2nd pass to compute m_p
 - Stable, but needs 2 scans of data
- ❑ **Textbook 1-Pass:** textbook rewrites $m_2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \frac{1}{n^2} (\sum_{i=1}^n x_i)^2$
 - **Notoriously unstable** (due to *catastrophic cancellation*), can even produce negative result for m_p when $p\%2=0$
 - Unfortunately, widely used in practice
- ❑ **Incremental:** stable* & 1 pass [*Bennett et al 2009*]

$$n = n_a + n_b, \delta = \mu_b - \mu_a, \mu = \mu_a \boxplus n_b \frac{\delta}{n}$$

$$M_p = M_{p,a} \boxplus M_{p,b} \boxplus \left\{ \sum_{j=1}^{p-2} \binom{p}{j} \left[\left(-\frac{n_b}{n}\right)^j M_{p-j,a} + \left(\frac{n_a}{n}\right)^j M_{p-j,b} \right] \delta^j + \left(\frac{n_a n_b}{n}\right)^p \left[\frac{1}{n_b^{p-1}} - \left(\frac{-1}{n_a}\right)^{p-1} \right] \right\}$$
- ❑ **MR Incremental:** adapt Incremental to MapReduce
 - Use *KahanAdd*

Numerically Stable Covariance

- ❑ **2-Pass:** 1st pass to compute means, 2nd pass to compute covariance
 - Stable, but needs 2 scans of data

- ❑ **Textbook 1-Pass:** textbook rewrites $c_2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \frac{1}{n-1} \sum_{i=1}^n x_i y_i - \frac{1}{n(n-1)} \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right)$
 - **Notoriously unstable** (catastrophic cancellation)
 - Unfortunately, widely used in practice

- ❑ **Incremental:** practically stable & 1 pass [Bennett et al 2009]

$$n = n_a + n_b, \quad \delta_x = \mu_{x,b} - \mu_{x,a}, \quad \mu_x = \mu_{x,a} \boxplus n_b \frac{\delta_x}{n}, \quad \delta_y = \mu_{y,b} - \mu_{y,a}, \quad \mu_y = \mu_{y,a} \boxplus n_b \frac{\delta_y}{n}$$

$$C = C_a \boxplus C_b \boxplus \frac{n_a n_b}{n} \delta_x \delta_y$$

- ❑ **MR Incremental:** adapt Incremental to MapReduce
 - Use *KahanAdd*

Experiment Results

Example Univariate Statistics

Ranges : R1= [1.0 – 1.5), R2= [1000.0 – 1000.5), R3= [1000000.0 – 1000000.5)

Range	Size (million)	Sum		Mean		Variance		Standard Deviation	
		SML	Naïve	SML	Naïve	SML	Textbook	SML	Textbook
R2	10	16.8	14.4	16.5	14.4	15.4	5.9	15.9	6.2
	100	16.1	13.4	16.9	13.4	15.6	5.3	15.8	5.6
	1000	16.6	13.1	16.4	13.1	16.2	4.9	16.4	5.2
R3	10	15.9	14.0	16.3	14.0	14.4	0	14.7	0
	100	16.0	13.1	16.9	13.1	12.9	Negative	13.2	NA
	1000	16.3	12.9	16.5	12.9	13.2	Negative	13.5	NA

Example Bivariate Statistics

Range	Size (million)	CoVariance		Pearson-R	
		SML	Textbook	SML	Textbook
R1 vs R2	10	15.0	8.4	15.1	6.2
	100	15.6	8.5	15.4	6.4
	1000	16.0	8.7	15.7	6.2
R2 vs R3	10	13.5	3.0	13.5	3.0
	100	12.8	2.8	12.7	NA
	1000	13.6	3.9	13.8	NA

Significantly better accuracy!
No sacrifice to performance!

- **Data Sets:** uniform distribution in 3 ranges
 - R1= [1.0 – 1.5), R2= [1000.0 – 1000.5), R3= [1000000.0 – 1000000.5)
 - modeled after NIST StRD benchmark
- **Accuracy Measure:**
 - *LRE* = $\log(\text{relative error})$
 - # significant digits that match between the computed value and the true value.
 - true value: produced by BigDecimal with precision 1000.

Lessons Learned (1/2)

- ❑ Many existing numerical stable techniques can be adapted to the distributed environment
- ❑ Divide-and-conquer design helps in scaling to larger data sets while achieving good numerical accuracy
 - e.g. MR Kahan can handle more data than Kahan with numerical accuracy
- ❑ Kahan technique is useful beyond simple summation

R1 10 million points

Variance		Std	
⌘	+	⌘	+
16.0	13.5	15.9	13.8

R1 vs R2 10 million points

Covariance		Pearson-R	
⌘	+	⌘	+
15.0	14.2	15.1	13.0

Lessons Learned (2/2)

- ❑ Shifting can be used for improved accuracy
 - Accuracy degrades as magnitude of values increases
 - Achieve better accuracy by shifting the data points by a constant prior to computation

10 million points

Range	Sum		Mean		Variance		Standard Deviation	
	SML	Naïve	SML	Naïve	SML	Textbook	SML	Textbook
1000 – 1000.5	16.8	14.4	16.5	14.4	15.4	5.9	15.9	6.2
1000,000 – 1000,000.5	15.9	14.0	16.3	14.0	14.4	0	14.7	0

- ❑ Performance need not be sacrificed for accuracy

❑ Recommended Reading:

- Nicholas J. Higham. *Accuracy and Stability of Numerical Algorithms*. SIAM, 2nd edition, 2002.

❑ **Thanks! and Questions?**