



Streaming Algorithms You Should Know

Ted Dunning



Contact Information

Ted Dunning

Chief Applications Architect at MapR Technologies

Committer & PMC for Apache's Drill, Zookeeper & Mahout

VP of Incubator at Apache Foundation

Email tdunning@apache.org tdunning@maprtech.com

Twitter [@ted_dunning](https://twitter.com/ted_dunning)

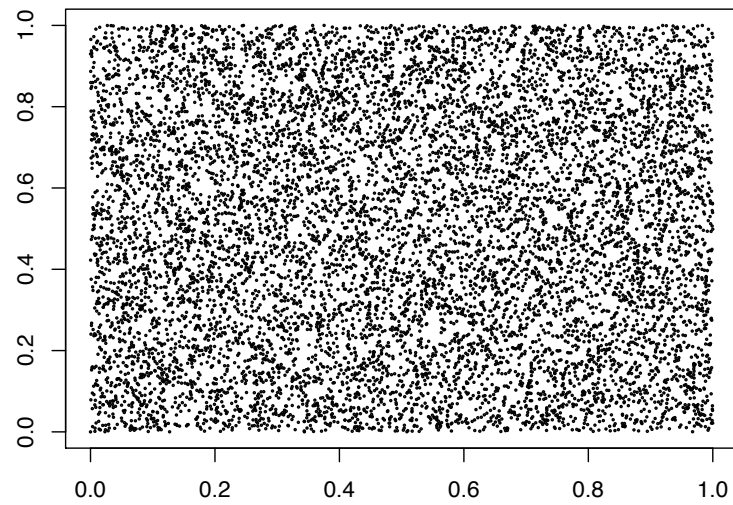
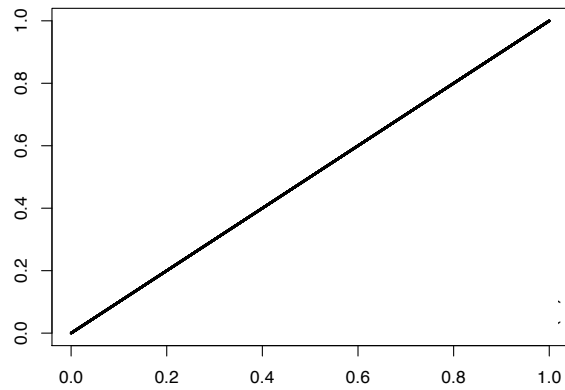


Must Approximate!

- Many important aggregates cannot be computed on-line
 - Memory cost proportional to data
 - Includes count distinct, heavy hitters, frequencies of top items, quantiles
 - Also k-means
- Most important aggregates can be approximated on-line
 - See all of the above



Key Trick - Hashing



Key Trick - Sketching

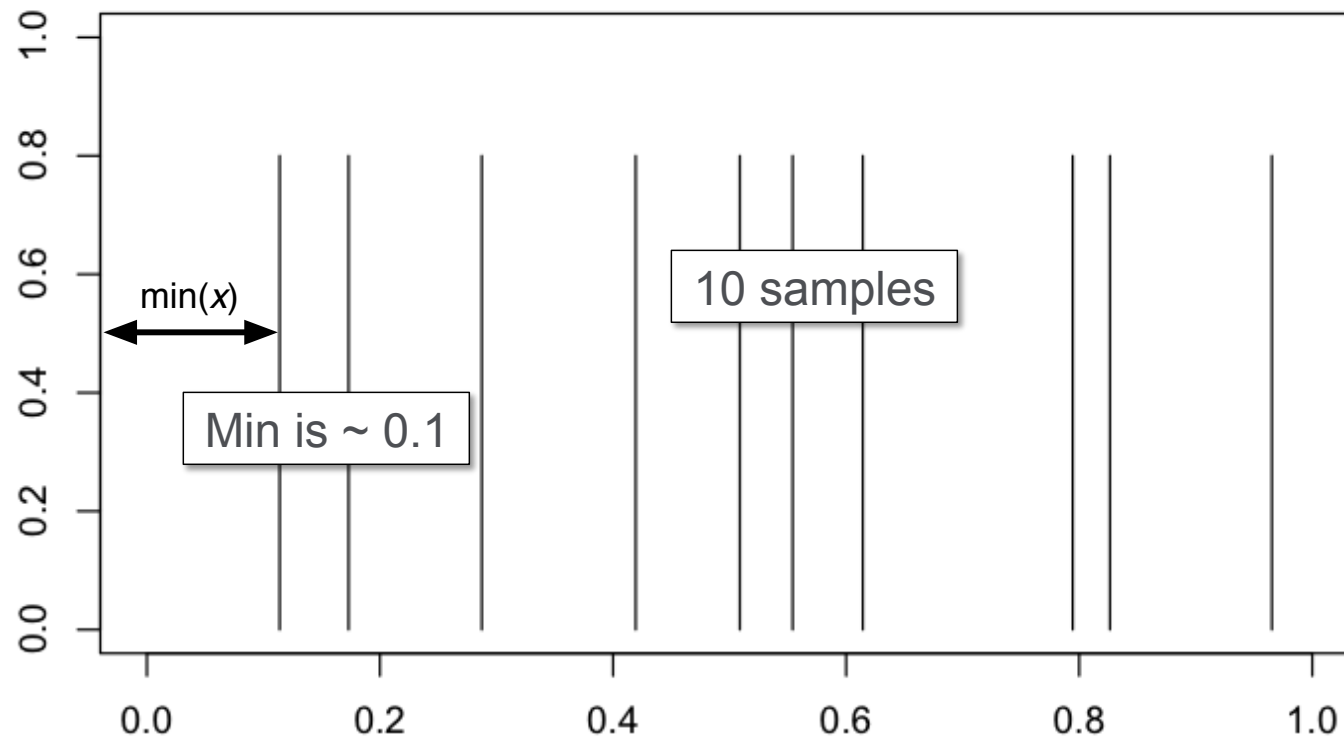
- Often we can build a very small sketch of the data
- Carefully done, we can use the sketch to get values of interest
- All of the algorithms here use sketches of some kind

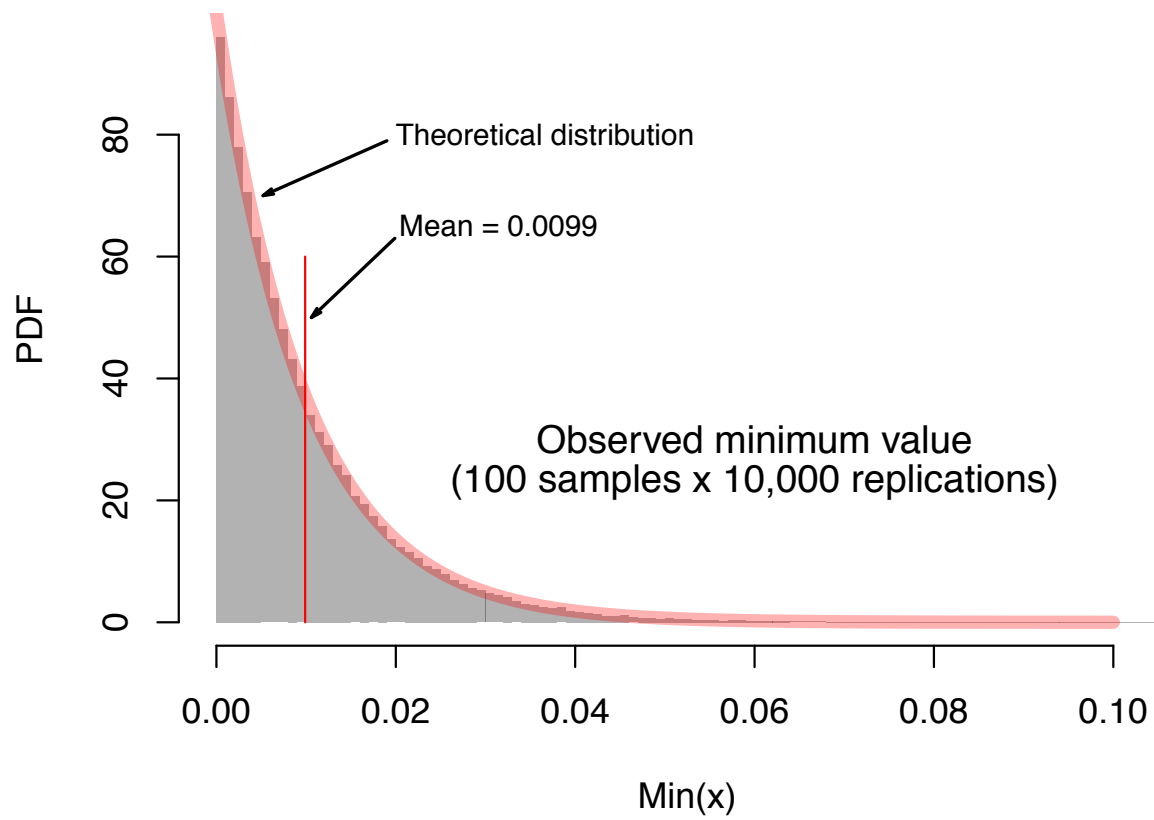


Hyper log log



Repeated Minimum

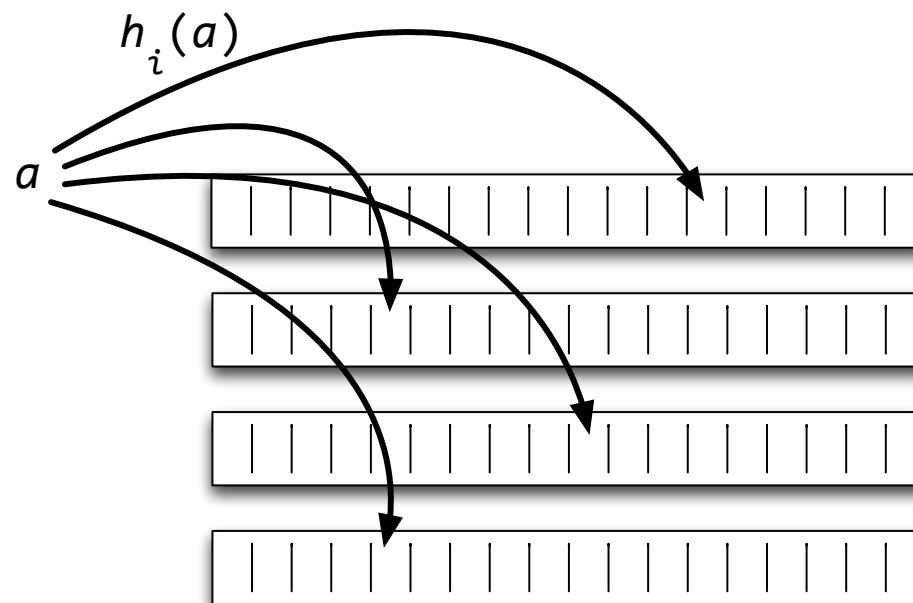




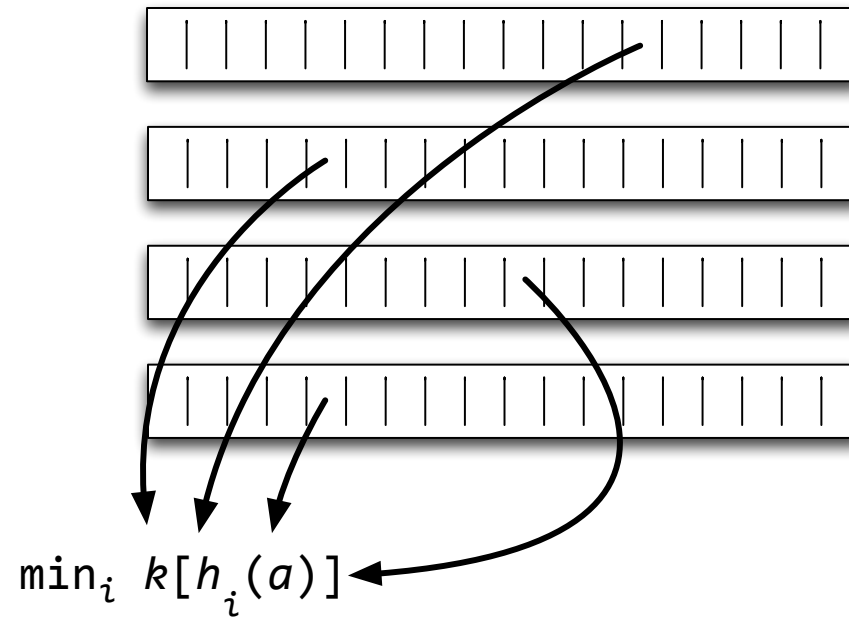
Count min sketch



Increment Hashed Locations to Insert

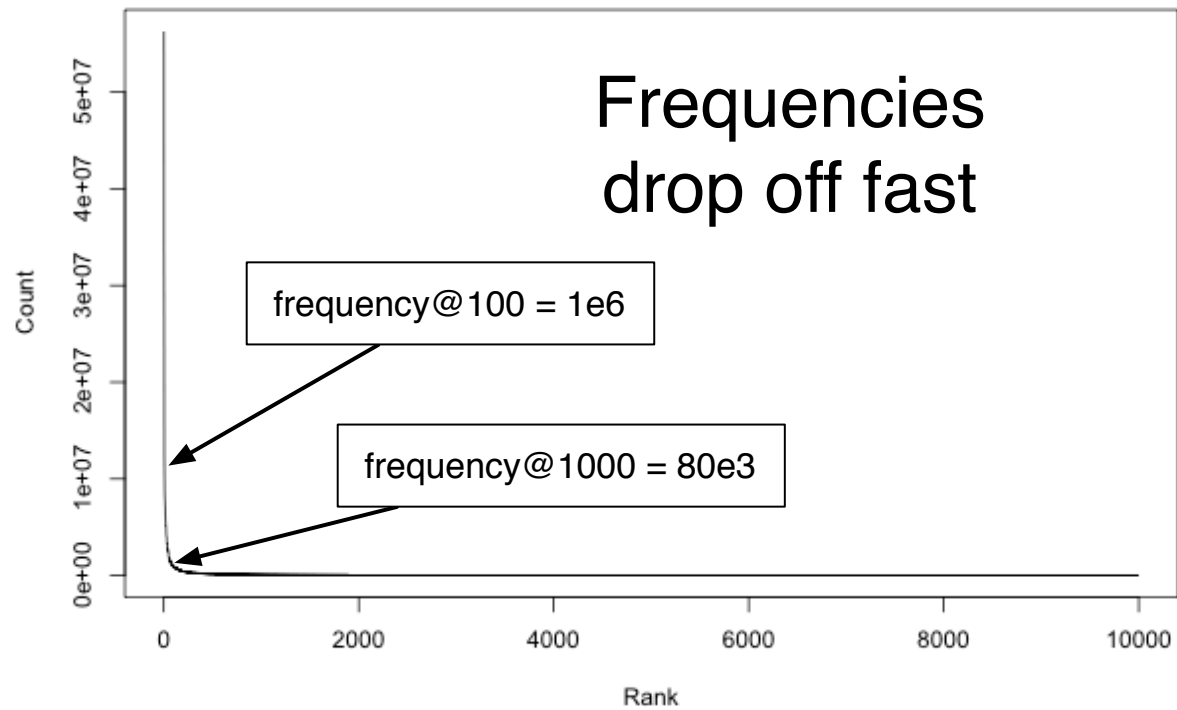


Probe Using min of Counts



Leaky counters

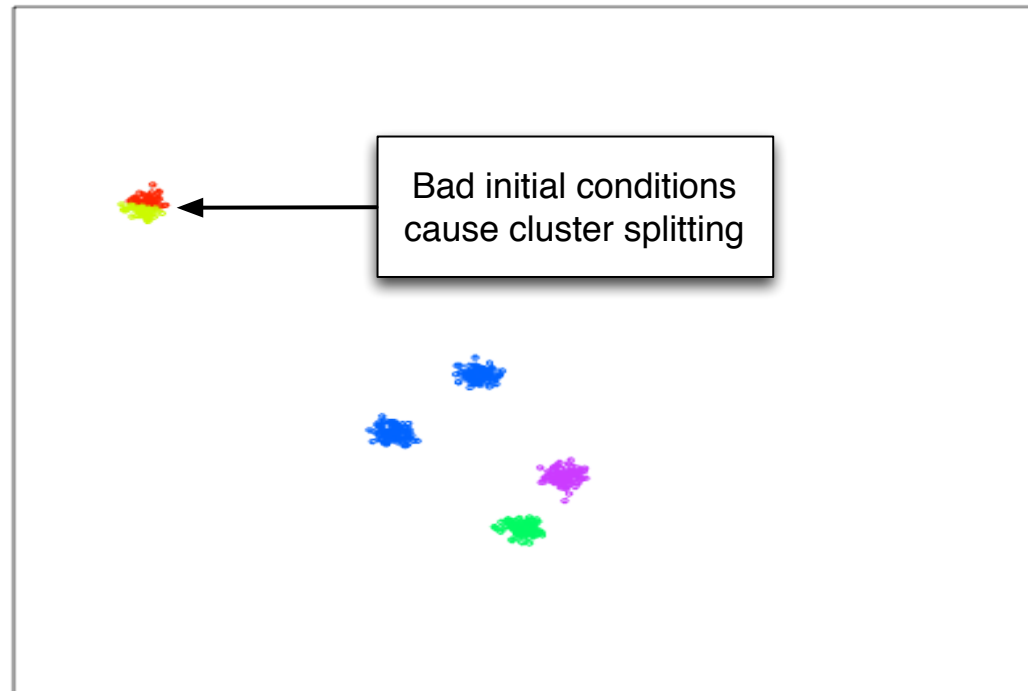




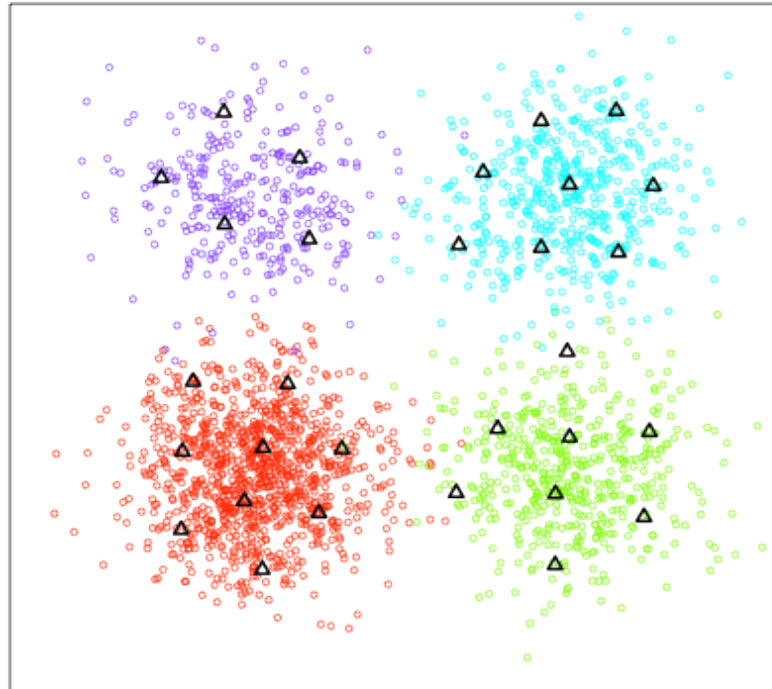
Streaming k -means



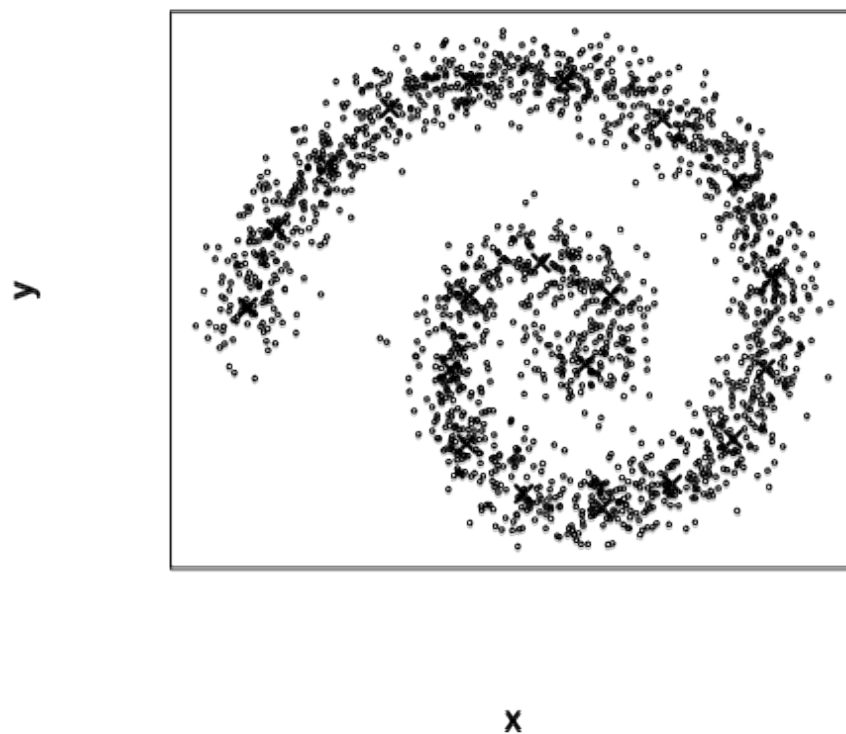
Typical k-means Failure



Many Clusters in Sketch is Fine



Lots of Clusters Approximate Anything



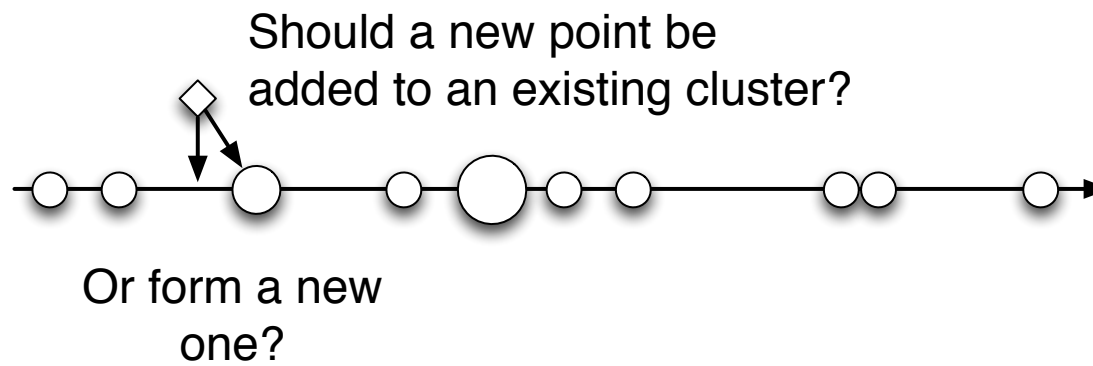
Clustering Summary

- Sketch is just lots of clusters
- Sketching can be done very approximately in one pass
- High quality clustering of the sketch is a high quality clustering of the data

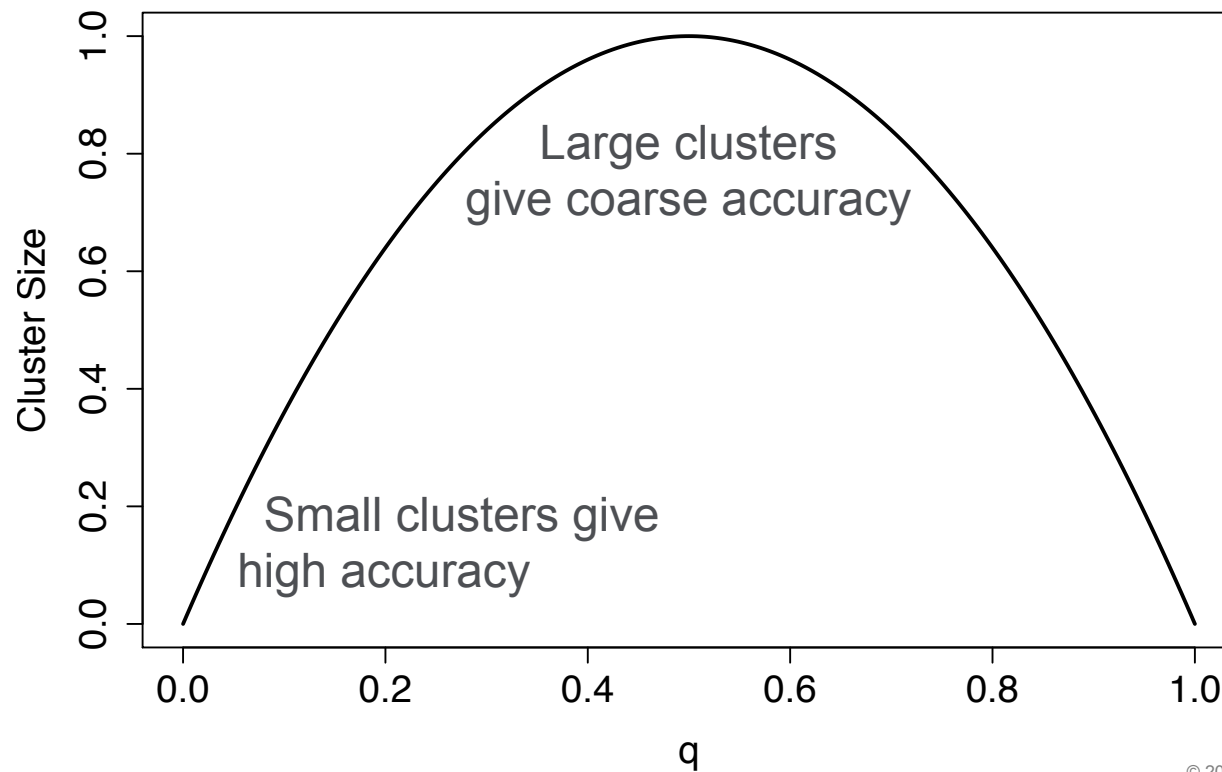


Finally, t -digest

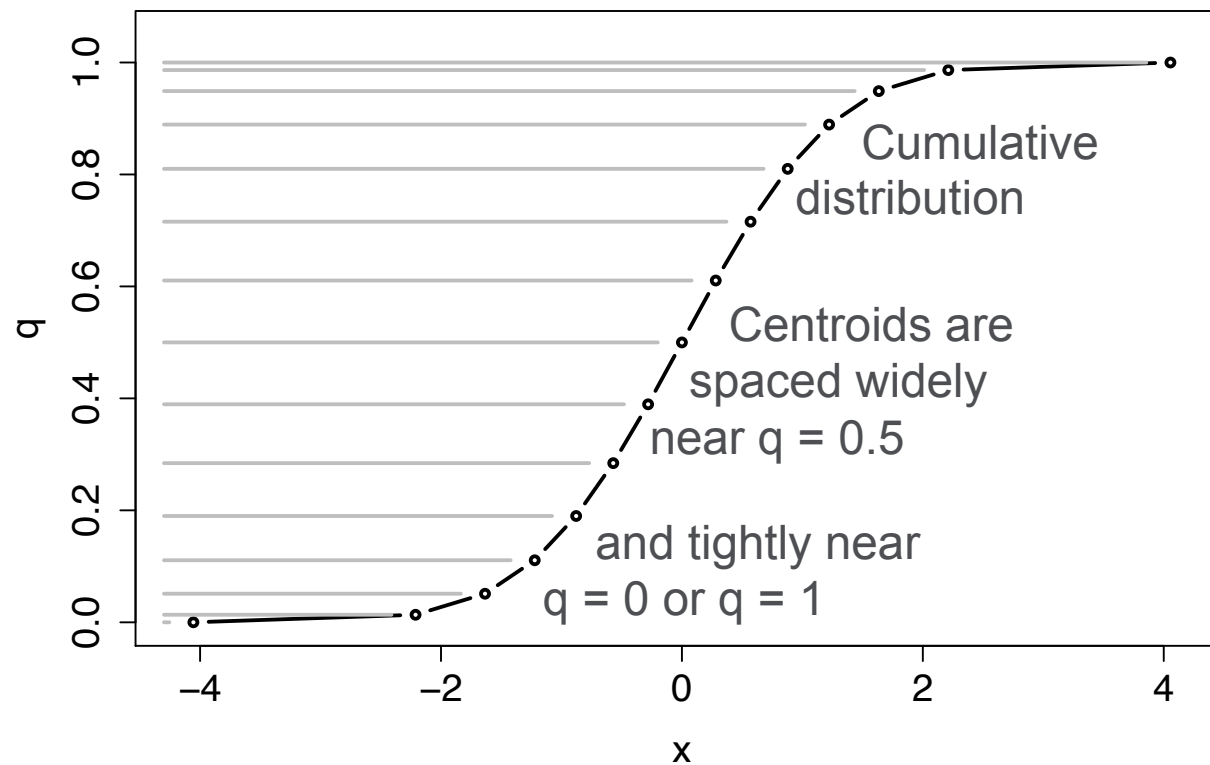




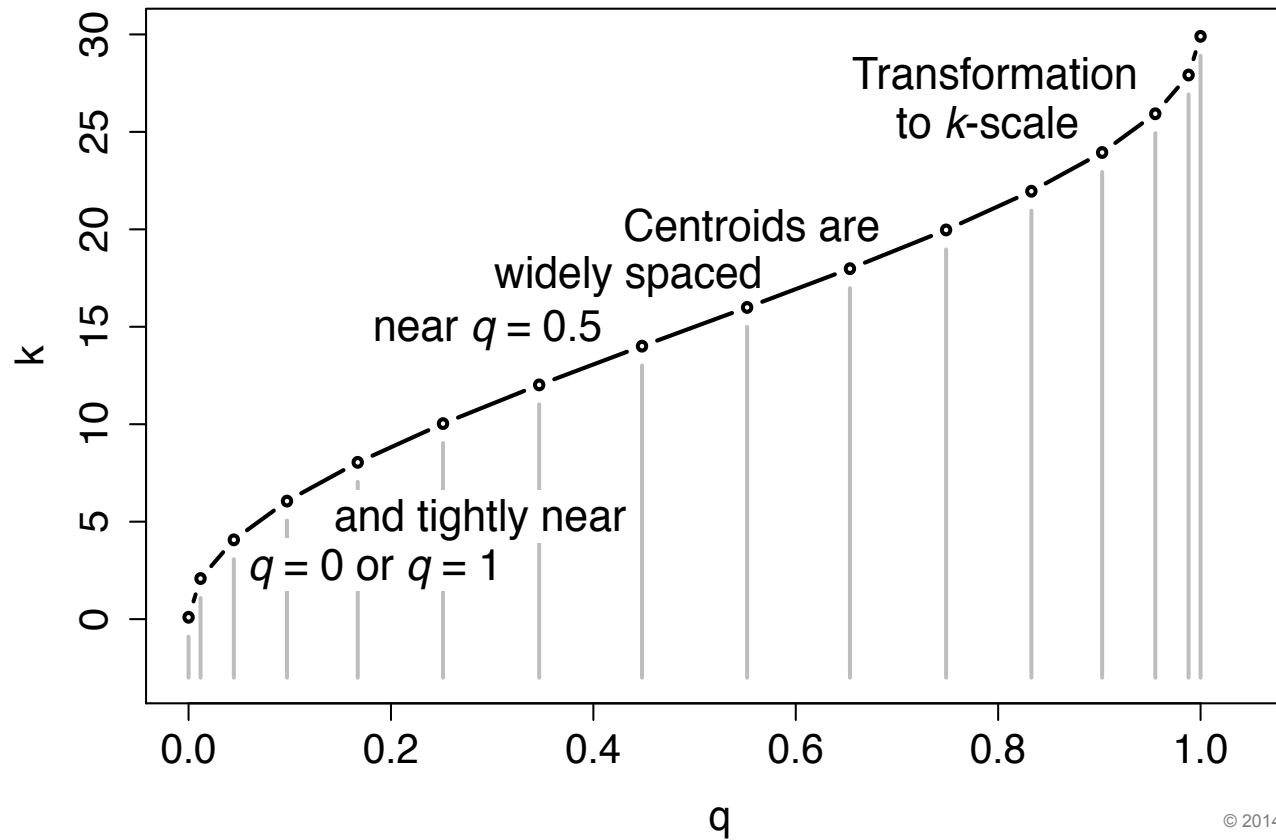
Variable Cluster Size for Constant Relative Accuracy



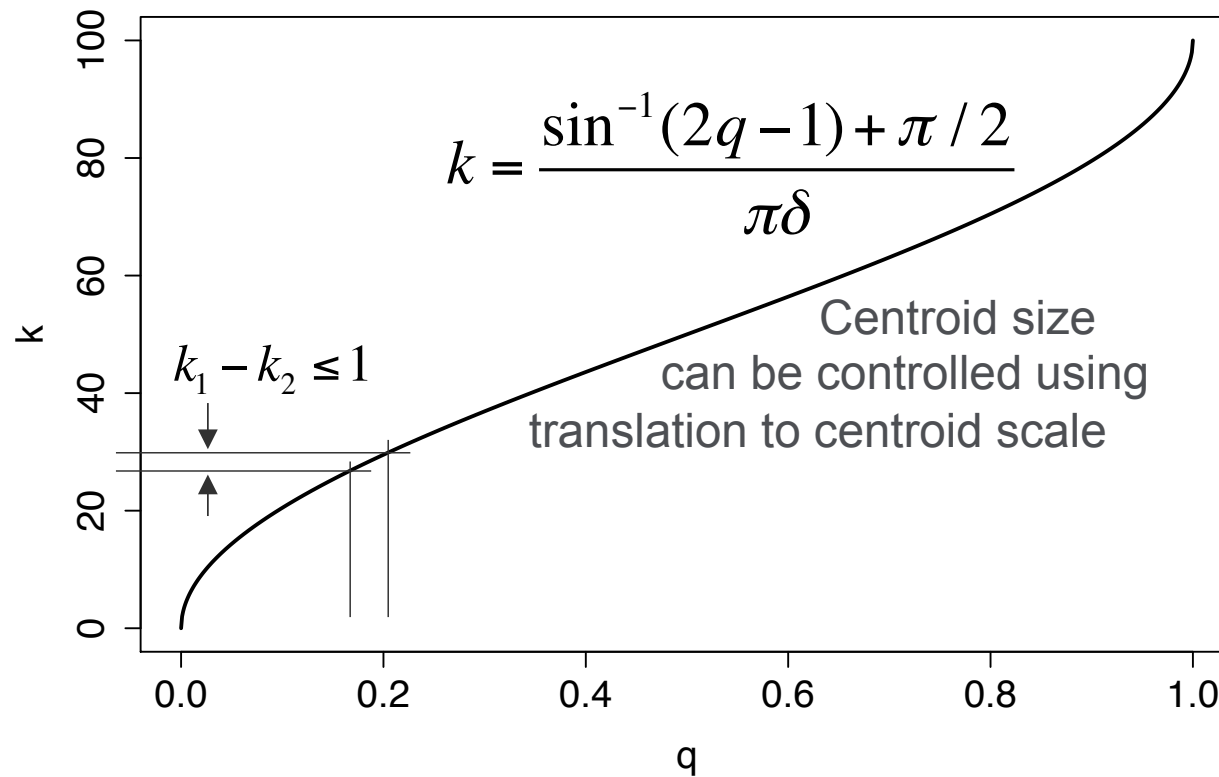
Second-Order Accuracy via Interpolation



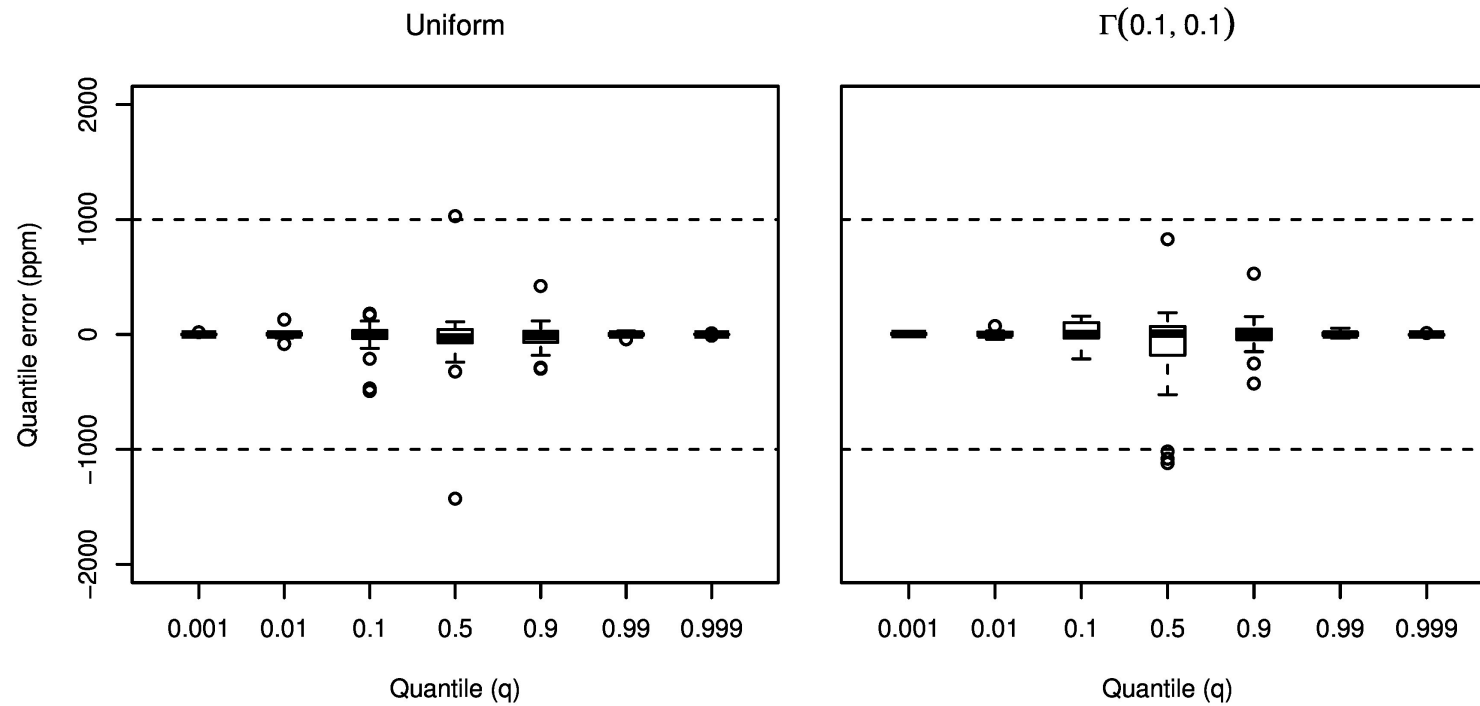
Size Limits via Constant Steps in k -scale



Translation Between Quantile and Cluster

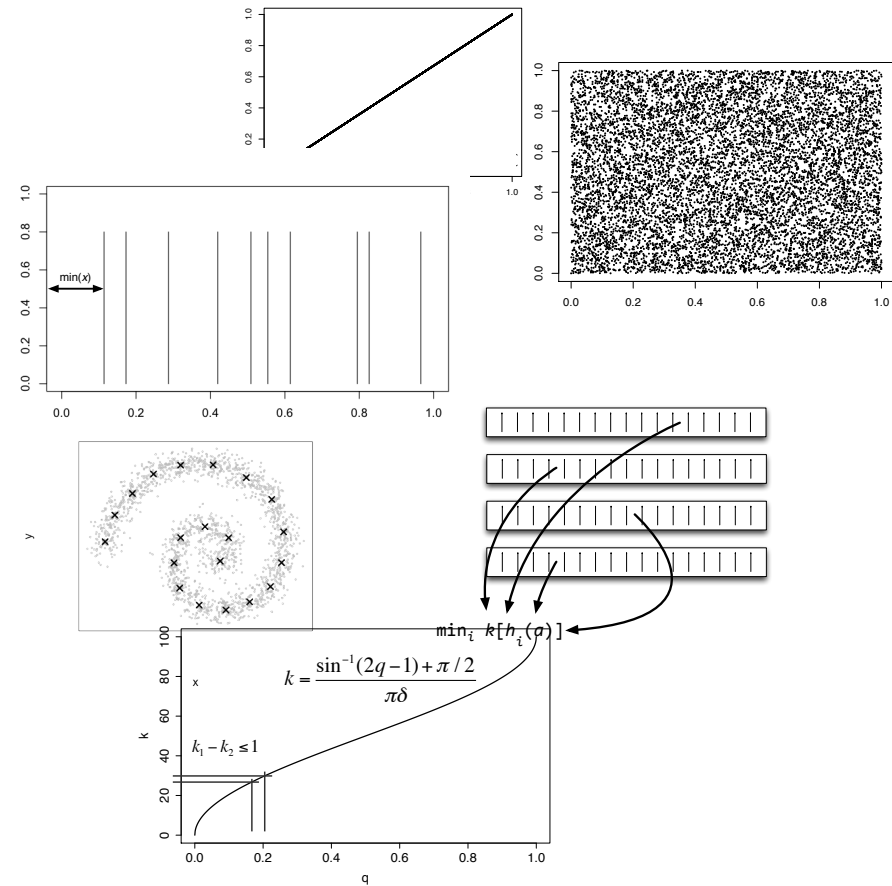


Actual Results



Summary

- Hashing and sketching
- Hyper log log = count distinct
- Count min = count(s)
- Streaming k-means
- Quantiles via t-digest



MAPR®



Free on-demand Hadoop training
leading to certification

Start becoming an expert now
mapr.com/training

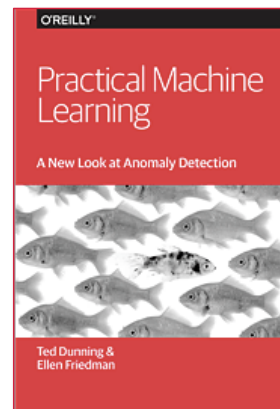


Short Books by Ted Dunning & Ellen Friedman

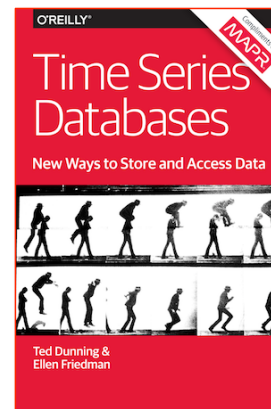
- Published by O'Reilly in 2014 and 2015
- For sale from Amazon or O'Reilly
- Free e-books currently available courtesy of MapR



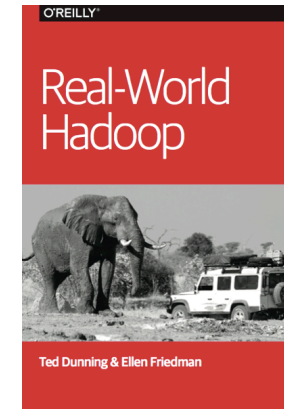
<http://bit.ly/recommendation-ebook>



<http://bit.ly/ebook-anomaly>



<http://bit.ly/mapr-tsdb-ebook>

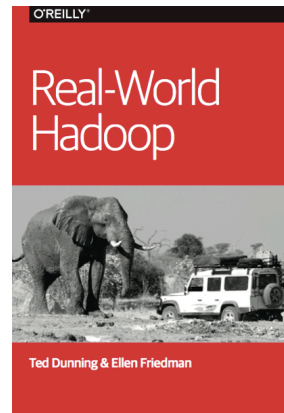


<http://bit.ly/ebook-real-world-hadoop>



Real World Hadoop

by Ted Dunning and Ellen Friedman © Feb 2015 (published by O'Reilly)



Free copies at book signing today



Thank You!



Q&A

Engage with us!

@mapr



maprtech

mapr-technologies



MapR

tdunning@mapr.tech.com



maprtech

