DOT-K: Distributed Online Top-K Elements JOHNS HOPKINS Algorithm using Extreme Value Statistics



WHITING SCHOOL of ENGINEERING

Nick Carey, Tamás Budavári, Yanif Ahmad, Alexander Szalay Department of Computer Science

ABSTRACT

DOT-K is a one pass, communications-efficient technique useful for both estimating upper order quantiles and selecting the largest 'k' elements in the context of an extremely large (petascale) and highly distributed data environment. Our novel approach draws its foundations from Extreme Value Statistics (EVS) to reason about the statistical relationships between the tail distributions of dataset partitions while minimizing communications costs.

DOT-K ALGORITHM

- Assuming a dataset row-partitioned across many nodes, our goal is to estimate the k'th largest element and subsequently retrieve all elements greater than the estimate.
 - 1. Each node collects its largest 'k' local values and using a parameter estimator such as MLE, calculates the GPD parameters that best fit the local data
 - 2. By relating the GPD parameters collected from each node, the query issuer estimates the global k'th

EXTREME VALUE STATISTICS

- Characterizes the tail distributions, or extreme values, of random variables
- Traditionally concerned with modeling frequency of extreme environmental phenomena

PICKANDS, BALKEMA, DE HAAN THEOREM

The distribution of threshold exceedances of a sequence of I.I.D. random variables with a common continuous underlying distribution function is approximated by the Generalized Pareto Distribution (GPD) and that the approximation converges as the tail threshold rises

$$p(x|\xi,\sigma,\mu) = \frac{1}{\sigma} \left[1 + \xi \cdot \left(\frac{x-\mu}{\sigma}\right) \right]^{-\frac{1}{\sigma}}$$

Equation 1. Generalized Pareto Distribution probability density function including parameters $\boldsymbol{\varepsilon}$ (shape) $\boldsymbol{\sigma}$ (scale) and $\boldsymbol{\mu}$ (location, or threshold)

For a large class of common data distributions, the 'k' largest values in a dataset may be approximated by a GPD provided the

- largest element by numerically solving Equation 3
- 3. The k'th order statistic estimate is communicated to the distributed nodes and the exceedances are relayed back to the query issuer

EXPERIMENTS

- DOT-K was executed within the Apache Storm stream processing framework on an Amazon AWS cluster to test communication costs and query latency
- A pseudo-distributed Matlab implementation was used to evaluate DOT-K query accuracy at larger scales



- With 'P' as parallelization factor, an ideal implementation communicates 4*P total messages between workers with $6^*P + \sim k$ total real values transmitted
- Slightly more than 4*P total messages were transmitted in our empirical evaluation due to Storm overhead **Distributed Worker Count** N = 50 million K = 2.5 million Naive Top-K versus POT-K Query result latency quickly overtakes a naïve top-k implementation as Naive Top-K parallelization rises. DOT-K retrieves only the elements estimated to be in POT-K the final query result from the distributed workers, and at very large scale this property reduces the query result sort cost Parallelization Factor or Number of Distributed Nodes Relative Error vs Partition Count on IID-Partitioned PageRank Data 0.065 _[0.06 $\delta k = \left| \frac{k - k}{k} \right|$ Relative error is given as: 0.055 0.05 Berkeley Big Data Benchmark 0.045 0.045 PageRank dataset was used to test DOT-K accuracy on a highly 0.04 partitioned real dataset 0.035 Each data point is the result of 0.03 the average of ten trials with the

k'th order statistic is an appropriately high threshold

Bias-variance trade-off: lower threshold results in theoretically worse GPD approximation while higher threshold limits amount of available threshold exceedances leading to greater parameter estimation uncertainty

M-OBSERVATION RETURN LEVEL

For a given GPD, one may calculate the threshold x_m that is exceeded on average once every m observations

$$\int_{-\frac{1}{\xi}} \left[1 + \xi \cdot \left(\frac{x_m - \mu}{\sigma} \right) \right]^{-\frac{1}{\xi}} = \frac{1}{m}$$

Equation 2. Coles' M-Observation Return Level equation. $\zeta_{\rm p}$ is a constant estimated by the number of observations exceeding μ divided by total observations



Equation 3. Our modification of Coles' M-Observation Return Level. Numerically



