Sampling-based Randomized Algorithms for Big Data Analytics

Matteo Riondato (matteo@cs.brown.edu, http://cs.brown.edu/~matteo) – Advisor: Prof. Eli Upfal, Dept. of Computer Science, Brown University, Providence, RI



Data Analytics Task	Contributions & Comparison with prev. work
Frequent Itemsets and Association Rules	Sampling algorithm – <mark>smaller sample size</mark>
	(MapReduce+sampling) alg. <mark>– more scalable</mark>
	Stat. test for false positives – more power
Betwenness Centrality	Sampling algorithm – <mark>smaller sample size</mark>
	Tighter analysis of existing algorithm
Database Query Selectivity	Sampling algorithm – <mark>smaller sample size</mark>

3. Approximations and Limitations of Classic Approach

Tradeoff between sample size and quality of approximation is well studied

Given $(u_i)_{1 \le i \le n} u_i \in [0, 1]$, $(\tilde{u}_i)_{1 \le i \le n}$ is (ε, δ) -approximation of $(u_i)_{1 \le i \le n}$ if

 $\left| \Pr(\exists i \text{ s.t. } |u_i - \tilde{u}_i| > \varepsilon) < \delta \right|$

Our goal: compute sample size |S| to obtain (ε, δ) -approximation

Classic bounds can not handle Big Data Variety (i.e., high values of n) E.g., with Chernoff bound + Union bound, sample size is

$$|\mathcal{S}| = O\left(\frac{1}{\varepsilon^2}\left(\ln n + \ln \frac{1}{\delta}\right)\right)$$

Dependency on $\ln n$ is too much for typical data analytics tasks E.g., in Frequent Itemsets mining, $\ln n$ is number of items, can be O(10⁴).

VC-dimension overcomes this issue: $|\mathcal{S}| = O\left(\frac{1}{\varepsilon^2}\left(\mathsf{VC}((u_i)_{1 \le i \le n}) + \ln \frac{1}{\delta}\right)\right)$

7. Estimating Betweenness Centrality

R. Kornaropoulos. "Fast Approximation of Betweenness Centrality through Sampling", ACM WSDM'14 Betweenness centrality: measure of vertex importance in graphs Settings: Graph G = (V, E) |V| = n, |E| = m

: fraction of Shortest Paths in G that pass through v



Exact algorithm for $(b(v))_{v \in V}$ takes time $O(nm + n^2 \log n)$ [Brandes01] Our goal: fast computation of (ε, δ) -approximation using sampling Algorithm

 $\mathsf{VD}(G) \leftarrow \max\{|p|, p \in \mathbb{S}_G\}$ $r \leftarrow (1/2\varepsilon^2)(\lfloor \log_2(\mathsf{VD}(G) - 2) \rfloor + 1 + \ln(1/\delta))$ $\mathsf{b}(v) \leftarrow 0, \forall v \in V$ for $i \leftarrow 1, \ldots, r$ $(u, v) \leftarrow random pair of vertices$ $S_{uv} \leftarrow \text{all SPs from } u \text{ to } v$ (BFS, Dijkstra, bidirectional search) $p \leftarrow \mathsf{random} \mathsf{SP} \mathsf{from} \mathcal{S}_{uv}$ $\tilde{\mathsf{b}}(w) \leftarrow \tilde{\mathsf{b}}(w) + 1/r, \forall w \in \mathsf{Int}(p)$ return $b(v), \forall v \in V$

Theorem: $(\tilde{b}(v))_{v \in V}$ is a (ε, δ) -approximation for $(b(v))_{v \in V}$

Rangeset for betweenness centrality: $D = \mathbb{S}_G, \mathcal{F} = \{\mathcal{T}_v, v \in V\}$

Theorem: $VC(\mathbb{S}_G, \mathcal{F}) \leq \lfloor \log_2(VD(G) - 2) \rfloor + 1$

Evaluation: C implementation, patch for igraph, on SNAP graphs



