HBP: *Hotness Balanced Partition* for Prioritized Iterative Graph Computations

Shufeng Gong, Yanfeng Zhang, and Ge Yu Northeastern University, China

Distributed Graph Computation



Distributed Graph Computation





Different graph partition schemes may result in different graph processing performance.





Workload Balance



one pair of vertices communicate between different partitions

Workload Balance



Workload Balance

one pair of vertices communicate between different partitions Less Communication 3

Graph Partition

each partition has the same number of vertices (four vertices) Workload Balance

 $\widehat{2}$

three pairs of vertices communicate between different partitions *More Communication* 3

Graph Partition



wait time

(2)

one pair of vertices communicate between different partitions Less Communication

3



A good partition can

Reduce wait time
 Minimize communication cost



Goals: 1. Balance the number of vertices 2. Minimize the number of edge cuts

Related Works

Metis¹, Leopard², Fennel³, HDRF⁴, Spinner⁵, Sheep⁶, NE⁷

They assume the workload depends on vertices (vertex centric) or edges (edge centric), and the communication cost depends on the number of edge cuts or vertex replications.

- 1. Karypis, George, and Vipin Kumar. "A fast and high quality multilevel scheme for partitioning irregular graphs." SISC 20.1 (1998): 359-392.
- 2. Huang, Jiewen, and Daniel J. Abadi. "Leopard: Lightweight edge-oriented partitioning and replication for dynamic graphs." Proceedings of the VLDB 2016.
- 3. Tsourakakis, Charalampos, et al. "Fennel: Streaming graph partitioning for massive scale graphs." WSDM 2014.
- 4. Petroni, Fabio, et al. "HDRF: Stream-based partitioning for power-law graphs." CIKM 2015.
- 5. Martella, Claudio, et al. "Spinner: Scalable graph partitioning in the cloud." ICDE 2017.
- 6. Margo, Daniel, and Margo Seltzer. "A scalable distributed graph partitioner." VLDB 2015.
- 7. Zhang, Chenzi, et al. "Graph edge partitioning via neighborhood heuristic." KDD 2017.

Synchronous Frameworks



Prior graph partition methods are designed based on synchronous framework.

Sync & Async

Some work pays attention to asynchronous frameworks since they can accelerate a class of iterative algorithms.



Sync & Async





Vertices can be processed at any time

Sync & Async



In each super step, each vertex is only processed once, and each edge only delivers one message· The number of updates on each vertex and the number of messages passed through each edge are not consistent.

Hotness of Vertices



Hotness of Vertices



Hotness of Vertices





Partition graph into G1 and G2

Worker1 processes G1 Worker2 processes G2



The number of vertices in G1 and G2 is equal. The sum of hotness in G1 and G2 is not equal.

> Prior vertex balanced graph partition methods do not take the vertex hotness into account when partitioning graphs.



Two workers are running.

There may be a lot of invalid computation in woker2 since vertices in G2 need fewer iterations.



Two workers are running.

There may be a lot of invalid computation in woker2 since vertices in G2 need fewer iterations.

Sum of hotness of G2 and G1 is equal







Priority Scheduling

Accelerate graph computing by selecting hot vertices to process preferentially.







Convergence Progress Bar

Worker1 process a Worker2 process m





Convergence Progress Bar

Worker1 process g Worker2 process n

A hot vertex can make more contributions to approaching the fixed point \cdot





Convergence Progress Bar

Worker1 process h Worker2 process l

A hot vertex can make more contributions to approaching the fixed point.



The worker2 is still not as efficient as we expected, because no matter how to schedule, worker2 always processes vertices with low hotness.

The hotness distribution of each partition should be the same as that of the original graph.

prioritized iterative graph computations

- Assign the same amount of vertex hotness to workers
- Minimize the variance between hotness distributions of each partition and the original graph.
- Minimize the communication cost between partitions.
 both synchronous and asynchronous frameworks

- Assign the same amount of vertex hotness to workers
- *Minimize the variance between hotness distributions of each partition and the original graph.*
- Minimize the communication cost between partitions.

$$\begin{split} HJS(P||P_i) = &\frac{1}{2} \Biggl[\sum_{j=1}^{z} P(\mathcal{H}_j) log \Biggl(\frac{P(\mathcal{H}_j)}{\frac{P(\mathcal{H}_j) + P_i(\mathcal{H}_j)}{2}} \Biggr) + \\ &\sum_{j=1}^{z} P_i(\mathcal{H}_j) log \Biggl(\frac{P_i(\mathcal{H}_j)}{\frac{P_i(\mathcal{H}_j) + P(\mathcal{H}_j)}{2}} \Biggr) \Biggr] \\ P(\mathcal{H}_j) = &\frac{\sum_{v \in \mathcal{H}_j} h_v}{\sum_{v \in V} h_v} \end{split}$$
 HJS distance is to rest between hotness dist partition and the order

HJS distance is to measure the variance between hotness distributions of each partition and the original graph.

- Assign the same amount of vertex hotness to partitions
- *Minimize the HJS distance between each partition and the original graph.*
- Minimize the communication cost between partitions.

- Assign the same amount of vertex hotness to workers
- *Minimize the HJS distance between each partition and the original graph.*
- Minimize the communication cost between partitions.

When balancing the hotness of each bin partition.

The amount of vertex hotness of each partition is equal.
 HJS distance between each partition and the original graph is 0.

The first two partition goals can be combined into one:

- Balance the hotness of each bin partition.
- Minimize the communication cost between partitions.

Per-Bin Hotness Balanced Partition



Fig. 1: hotness balance

Fig. 2: per-bin hotness balance

Per-Bin Hotness Balanced Partition

The hotness of each partition is balanced and the distribution of hotness values of each partition is the same as that of the original graph.



Both worker1 and worker2 are efficient

per-bin hotness balance

Per-Bin Hotness Balanced Partition



Partition each bin into k hotness balanced partitions.

Maximize the communication between bins in the same partitions

Minimize the communication between bins in different partitions.

An illustration of Per-Bin hotness balanced partition



The vertex with high indegrees always has higher execution priority and is likely to be hot.



The vertex with high indegrees always has higher execution priority and is likely to be hot.

Weights of edges may affect the hotness value of target vertices.





$$h_{v} = \sum_{u \in IN(v)} \frac{w_{u,v}}{\sum_{w \in OUT(u)} w_{u,w}}$$

the hotness of vertex v





$$h_{v} = \sum_{u \in IN(v)} \frac{w_{u,v}}{\sum_{w \in OUT(u)} w_{u,w}}$$

$$com_{u_v^v} = hu$$

the hotness of vertex v

the communication of edge $e_{u,v}$

Streamed Per-Bin Hotness Balanced Partition We assign each vertex to

belong to bin3

the hotness values are separated into 3 bins We assign each vertex to the partition with low hotness on the corresponding bin and low communication cost (edge information is omitted).

The first partition's third bin has smallest hotness, so this vertex is assigned into the first partition.

Bin 1

Bin 2

Bin 3





Partition 3

the hotness values are separated into 3 bins

> Bin 1 Bin 2

Bin 3



Partition 1



Partition 3

the hotness values are separated into 3 bins

Bin 1 Bin 2

Bin 3







Partition 1

Partition 2

Partition 3

the hotness values are separated into 3 bins

> Bin 1 Bin 2

Bin 3



Partition 1



Partition 2

Partition 3

 $i = \underset{i:\{\boldsymbol{h}_{ji} \leq \tau \frac{\sum_{v \in \mathcal{H}_j} h_v}{k}\}}{\arg \min} \alpha \cdot \left((\boldsymbol{h}_{ji} + h_v)^{\gamma} - \boldsymbol{h}_{ji}^{\gamma} \right) + (1 - \alpha) \cdot \left(\sum_{u \notin V_i} com_{u,v} + \sum_{w \notin V_i} com_{v,w} \right)$

where $0 \le \alpha \le 1$, $\tau > 1$, $\gamma > 1$, and h_{ji} is the sum of hotness values of i-th partition in the j-th bin.

Streamed Per-Bin Hotness Balanced Partition *imbalance cost*

$$i = \underset{i:\{\boldsymbol{h}_{ji} \leq \tau \frac{\sum_{v \in \mathcal{H}_{j}} h_{v}}{k}}{\operatorname{arg\,min}_{k}} \left[\alpha \cdot \left((\boldsymbol{h}_{ji} + h_{v})^{\gamma} - \boldsymbol{h}_{ji}^{\gamma} \right) \right] + (1 - \alpha) \cdot \left(\sum_{u \notin V_{i}} com_{u,v} + \sum_{w \notin V_{i}} com_{v,w} \right) \right]$$

Communication cost

where $0 \le \alpha \le 1$, $\tau > 1$, $\gamma > 1$, and h_{ji} is the sum of hotness values of i-th partition in the j-th bin.

Experiment: Preparation

Platform	Alibaba Cloud
System	Maiter
Cluster	1 master (ecs.cs.large)
	4 slaves (ecs.cs.large)
Algorithm	PageRank, PHP
Competitors	Hash, Fennel ¹ , HotGraph ²
Data sets	Twitter (TW), LiveJournal (LJ) Hollywood (HW)

- 1. Tsourakakis, Charalampos, et al. "Fennel: Streaming graph partitioning for massive scale graphs." Proceedings of the WSDM 2014.
- 2. Zhang, Yu, et al. "HotGraph: Efficient asynchronous processing for real-world graphs." *IEEE TOC 2016*.

Experiments: Runtime Comparison



SPb-HBP is our stream-based per-bin hotness balanced partition

Experiments: Communication Cost Comparison



Experiments: Scalability



Experiments: Scalability



Conclusion

- The analysis of the existing graph partition methods finds that they are not suitable for asynchronous graph processing systems.
- A new graph partition idea, hotness balanced partition, is proposed.
- A stream-based per-bin hotness balanced partition algorithm is proposed.



Any questions please email gongsf@stumail·neu·edu·cn

