Hogwild: A Lock-Free Approach to Parallelizing Stochastic Gradient Descent

CS 584: Big Data Analytics

### Sparse Separable Cost Function

 Minimize function that can be separated into small subsets of feature components

$$\min\sum_{e\in E} f_e(x_e)$$

- x<sub>e</sub> denotes the values of the vector x on coordinates indexed by e
- Key observation: For many machine learning problems of interest, sub-vector x<sub>e</sub> contains very small number of components (is sparse compared to n and |E|)

### Example: Sparse Functions

- Sparse SVM: Features are very sparse (only few non-zero components)
- Matrix Completion: Only provided entries from a sparse sampling of data (e.g., each user rates very few movies)
- Graph Cuts: Find the minimum number of cuts to separate entities (group similar items together)

### Two Key Issues with SGD

- Difficult to parallelize
  - Inherently serial algorithm where parameter is updated after seeing an example
  - Multi-threading may be useless as each one needs to wait for another one to proceed
- Notoriously hard to tune
  - Learning rate will drastically affect the speed of the algorithm

### An Attempt at Parallel SGD

- Each thread draws random sample from training data
  - Acquire a lock on current state of parameters
  - Read the parameter
  - Update the parameter with SGD step
  - Release lock

# Acquiring locks can take 1000x longer than to make an update!

### HOGWILD!: Asynchronous SGD

- Remove all thread locks from parallel SGD code
- Threads allowed to overwrite one another
- Gradients can be computed on stale versions of the "current solution"
- Anything goes!

### HOGWILD! Update for Individual Processors

Algorithm 1 HOGWILD! update for individual processors

- 1: **loop**
- 2: Sample e uniformly at random from E
- 3: Read current state  $x_e$  and evaluate  $G_e(x)$
- 4: for  $v \in e$  do  $x_v \leftarrow x_v \gamma b_v^T G_e(x)$
- 5: end loop

#### Wait... What? How can this work?

## HOGWILD! Convergence

- Processors overwrite each other's work, but the sparsity of the gradient helps greatly
  - Updates don't interfere too much with one another
- Updates can be old by the time they are applied, but in practice there is a certain bound on their age
- Depending on the overlap of the nonzero components, processing time, and other quantities, the algorithm can essentially recover the behavior of SGD

### HOGWILD! Performance Comparison



### HOGWILD! CPU Time vs. Threads



Round robin gets worse as more threads are added!

### Summary

- Practical lesson: Don't lock!
- HOGWILD! started trend in asynchronous algorithms for model training
- We should all run HOGWILD!