

STAT4402 Tutorial

Michael Ciccotosto-Camp

University of Queensland

1. Introduction

Most computational work done within science is impractical to perform on commercial laptops and desktops, typically due to the extremely high memory and processing demands. Hence, almost every university and industry has its own high-performance computer to carry out such strenuous problems. A lot of research within the realm of Machine Learning benefits from access to high performance machinery as most of them require matrix computation (which can be efficiently carried out on GPU clusters) and can be processed in parallel.

In this tutorial we will revisit two models covered in the first few weeks of lectures, these being the SDG linear regressor and KNN classifier. The naive implementation of both these algorithms are fairly inefficient, so we shall look at some ways in which these two methods can be decomposed and parallelised.

2. Parallel KNN

- The k^{th} Nearest-Neighbor (k-NN) methods use observations in the training set T closest in feature space to a given unknown sample \mathbf{x} to directly find its corresponding prediction \bar{y}
- The prediction for the k-NN classifier is usually calculated as

$$\bar{y}(\mathbf{x}) = \sum_{\mathbf{x}_i \in N_k(\mathbf{x})} y_i$$

- The notion of 'closest' implies the use of some sort of metric. More often than not, feature vectors belong to \mathbb{R}^n allowing us to use commonly used metrics to define distance between vectors in our feature space. For our purposes, we shall use the

Euclidean norm as a measurement of determining how close two feature values are to each other. The Euclidean norm is simply defined as

$$d(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^n (x_i - y_i)^2 \right)^{\frac{1}{2}}$$

- When a unknown sample \mathbf{x} is to be classified, a k-NN classifier computes the distance between \mathbf{x} and the other points within the training set T . The training data is then sorted by distance and the k^{th} closests training samples are then used to predict \mathbf{x} .
- A simple K-NN algorithm works as follows

Algorithm 1: Serial k-NN

input : Training data T , an unlabelled sample \mathbf{x} and a value k

output: Predicted class $\bar{y}(\mathbf{x})$

- 1 Computes distance $d(\mathbf{x}, \mathbf{x}_{t_i})$ for each $\mathbf{x}_{t_i} \in T$;
- 2 $N_k(\mathbf{x}) \leftarrow$ the k^{th} closest \mathbf{x}_{t_i} determined by $d(\mathbf{x}, \mathbf{x}_{t_i})$;
- 3 $\bar{y}(\mathbf{x}) \leftarrow \sum_{\mathbf{x}_i \in N_k(\mathbf{x})} y_i$;

Result: $\bar{y}(\mathbf{x})$

- While this method is simple, computing the distance between \mathbf{x} and each $\mathbf{x}_{t_i} \in T$ can incur a large time overhead, especially for large training sets
- One important observation is that computing the distances $d(\mathbf{x}, \mathbf{x}_{t_i})$ and $d(\mathbf{x}, \mathbf{x}_{t_j})$ where $\mathbf{x}_{t_i}, \mathbf{x}_{t_j} \in T$ and $i \neq j$ may be done completely independently of each other meaning these computations may be carried out on separate processes.
- CITATION takes advantage of this independence to have distance computation carried out on different processes.

This algorithm is shown pictorially below

3. Parallel SGD

- As before let T a set of training samples $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^m = \{z_i\}_{i=1}^m$
- Let $C(\mathbf{w})$ be the cost function

$$C(\mathbf{w}) = \frac{1}{N} \sum_{i=0}^m C_{z_i}(\mathbf{w})$$

Algorithm 2: Parallel k-NN

input : Training data T , an unlabelled sample \mathbf{x} , a value k and the number of processes to perform the algorithm p

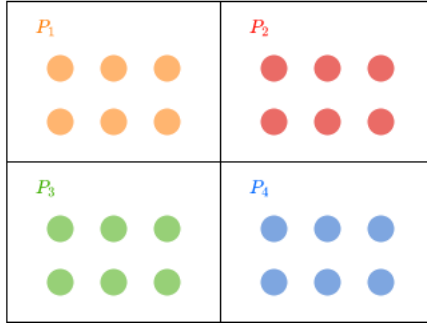
output: Predicted class $\bar{y}(\mathbf{x})$

```

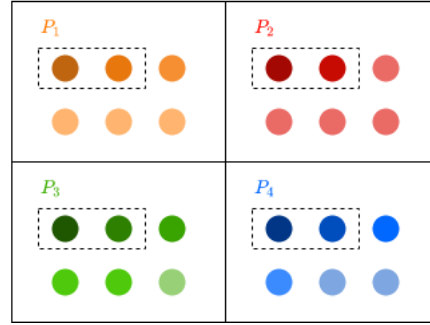
1  $\{T_1, T_2, \dots, T_p\} \leftarrow$  an equal partition of  $T$ ;
2 for  $T_i \in \{T_1, T_2, \dots, T_p\}$  concurrently do
3    $N_{k_i}(\mathbf{x}) \leftarrow$  the  $k^{th}$  nearest neighbors from  $T_i$ ;
4 end
5  $N_k(\mathbf{x}) \leftarrow$  the  $k^{th}$  closest neighbors from  $N_{k_1}(\mathbf{x}), N_{k_2}(\mathbf{x}), \dots, N_{k_p}(\mathbf{x})$ ;
6  $\bar{y}(\mathbf{x}) \leftarrow \sum_{\mathbf{x}_i \in N_k(\mathbf{x})} y_i$ ;
Result:  $\bar{y}(\mathbf{x})$ 

```

1) Partition the data and send partitions to processes P_1, P_2, P_3, P_4



2) Determine the closest samples within each process



3) Collect the k^{th} closest samples and order them in the master node get the global k^{th} closest samples



- For gradient descent algorithms, we wish to find a weight vector \mathbf{w}^* that minimizes our cost function, that is

$$\mathbf{w}^* = \arg \min_{\mathbf{w} \in \mathbb{R}^n} \sum_{i=0}^m C_{z_i}(\mathbf{w})$$

- We shall also introduct the notation $G \triangleq \frac{\partial C}{\partial \mathbf{w}}$ and $G_z \triangleq \frac{\partial C_z}{\partial \mathbf{w}}$ to simplify gradient notation as well as $H \triangleq \frac{\partial G}{\partial \mathbf{w}}$ and $H_z \triangleq \frac{\partial G_z}{\partial \mathbf{w}}$ to simplify Hessian notation.
- At each step of the SGD, a sample $z_j = (\mathbf{x}_j, y_j)$ is uniformly selected from the training

set to update the existing weight vector as

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t G_{z_j}(\mathbf{w}_t)$$

- where η_t is just the learning rate at iteration t .
- Say a process performs SGD of a data set T_1 to get from a weight \mathbf{w}_g to weight \mathbf{w}_1 . When processing another training set T_2 , a sequential SGD algorithm would have started at weight \mathbf{w}_1 to reach a possibly different weight vector \mathbf{w}_h .
- To parallelize the SDG algorithm we wish to start computing of training set T_2 on weight vector \mathbf{w}_1 while simultaneously running the training set T_1 on weight vector \mathbf{w}_g , but \mathbf{w}_1 is not know until SGD is finished with T_1 . So how do we ge around this?
- One method is to soundly combine models from different processes, in the hopes of achieving a weight vector if SGD was simply run sequentially.
- This requires adjusting the computation of T_2 to account for the staleness $\mathbf{w}_1 - \mathbf{w}_g$ in the initial model.
- To do so, the second model performs its computations instead on $\mathbf{w}_g - \Delta\mathbf{w}$ where $\Delta\mathbf{w}$ is an unknown symbolic vector
- This allows the second model to run in parallel and not have to wait until \mathbf{w}_1 is produced from the succeeding model.
- Once the first process is done, the second process takes $\Delta\mathbf{w}$ to be $\mathbf{w}_1 - \mathbf{w}_g$
- This technique can be easily extended to an arbitrary number of processors.
- Let $S_T(\mathbf{w})$ represent the SGD computation of a training T from an initial weight vector \mathbf{w} , for example $S_{T_1}(\mathbf{w}_g) = \mathbf{w}_1$
- To come up with a model combiner we need to think about how we can calculate

$$S_T(\mathbf{w} + \Delta\mathbf{w})$$

- Assuming S_T is differentiable at $\mathbf{w} + \Delta\mathbf{w}$, we get following by consider the Taylor series of S_T about the point $\mathbf{w} + \Delta\mathbf{w}$

$$S_T(\mathbf{w} + \Delta\mathbf{w}) = S_T(\mathbf{w}) + S'_T(\mathbf{w}) \cdot \Delta\mathbf{w} + \mathcal{O}(|\Delta\mathbf{w}|^2)$$

- We will introduce the notation $M_D \triangleq S'_T$ as the model combiner. In the equation above the model combiner captures the first order information of how a change in $\Delta \mathbf{w}$ will effect the SGD
- When $\Delta \mathbf{w}$ is sufficiently small, one can neglect higher order terms and only use the model combiner to combine models from different processes.
- From CITATION we can show that for a sequence of input examples z_1, z_2, \dots, z_n the model combiner can be computed as

$$M_D(\mathbf{w}) = \prod_{i=1}^n (\mathbf{I} - \eta_i \cdot H_{z_i}(S_{T_{i-1}}(\mathbf{w})))$$

where $S_{T_0}(\mathbf{w}) = \mathbf{w}$. This result can be easily shown by applying the chain rule to $S_T(\mathbf{w}) = S_{T_n}(S_{T_{n-1}}(\dots(S_1(\mathbf{w}))))$.

- Thus to create a parallelized SGD, each of the p processors start with the same initial global weight vector \mathbf{w}_g to compute its local model $S_{T_i}(\mathbf{w}_g)$ and model combiner $M_{T_i}(\mathbf{w}_g)$ in parallel. A subsequent reduction phase computes \mathbf{w}_i by adjusting the input of processor i by adjusting by the staleness introduced in the $i - 1$ processor

$$\mathbf{w}_i = S_{T_i}(\mathbf{w}_g) + M_{T_i}(\mathbf{w}_g) \cdot (\mathbf{w}_{i-1} - \mathbf{w}_g)$$

- The algorithm for parallel SGD is summaries below

Algorithm 3: Parallel SGD

input : Training data T , an initial weight vector \mathbf{w}_g and the number of processes to perform the algorithm p

output: Updated weight vector

```

1  $\{T_1, T_2, \dots, T_p\} \leftarrow$  an equal partition of  $T$ ;
2 for  $T_i \in \{T_1, T_2, \dots, T_p\}$  concurrently do
3   | Compute  $S_{T_i}(\mathbf{w}_g)$  and  $M_{T_i}(\mathbf{w}_g)$ ;
4 end
5 for  $i \in \{1, 2, \dots, p\}$  do
6   |  $\mathbf{w}_i \leftarrow S_{T_i}(\mathbf{w}_g) + M_{T_i}(\mathbf{w}_g) \cdot (\mathbf{w}_{i-1} - \mathbf{w}_g)$ ;
7 end
```

Result: \mathbf{w}_p

4. Introduction to High Performance Computing

- So all these parallel algorithms seem great and all, but they won't actually benefit us anything if we don't have the computing power to run them!
- While commercial bought laptops usually have more than one core in them, their usage is taken up by processes running in the background, so if you ran anyone of these algorithms on your own machine, chances are you won't see much of an improvement in performance
- You might be asking now, *what sort of computers can run these algorithms to actually see performance improvements?* Well I'm glad you asked! The answer High Performance Computers. High Performance Computers (or just HPCs) are very large machines consisting of hundreds or even thousands of cores to run scientific or analytic programs on. The processes on HPCs are monitored by a special operating system so that any jobs that you submit to be run of a HPC gets exactly the number of processes and amount of memory you've asked for (provided that HPC system is able to provide those resources).
- As a student taking STAT4402 you should have access to UQ's `getafix` HPC.