

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

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

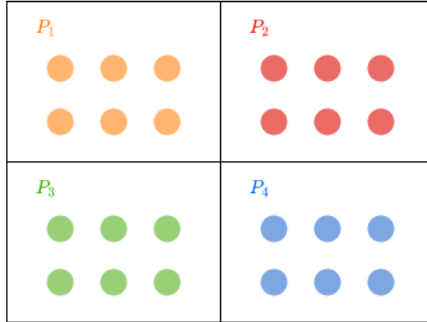
Algorithm 2: Parallel k-NN

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

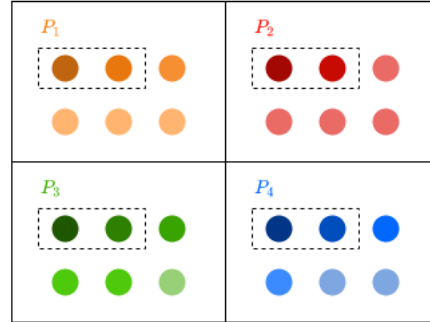
output: Predicted class $\bar{y}(\mathbf{x})$

- 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

