

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 uses observations from a training set T to find its closest neighbours in feature space to a given an unknown sample \mathbf{x} a prediction value \bar{y} [?]. The prediction for the k-NN classifier is usually calculated as an average or consensus vote, that is

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

The notion of 'closest' implies the use of some sort of metric. More often than not, feature vectors belong to some subset \mathbb{R}^n , allowing us to apply 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

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

Other metrics such as Manhattan norm or hamming distance for discrete feature spaces [?]. 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 is presented in algorithm ???. While this method is simple, computing

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

the distance between \mathbf{x} and each $\mathbf{x}_{t_i} \in T$ can incur a large time overhead, especially for large training sets. In order to parallize this algorithm, 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 $i \neq j$ allowing us to carry out these computations on separate processes. Liang et al. (2010) takes advantage of this independence by splitting up the training examples into p partitions and send

these partitions to slave processes to compute the distances between feature vectors in our training example and a unknown sample \mathbf{x} . The k^{th} closet training examples are found within individual slaves and are collected within the master process. Combining the k^{th} closet acquired by each process locally, the master process then sorts the reduced list of training vectors to find the global k^{th} closet training examples [?]. Pseudo code for parallel KNN is shown in ??.

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

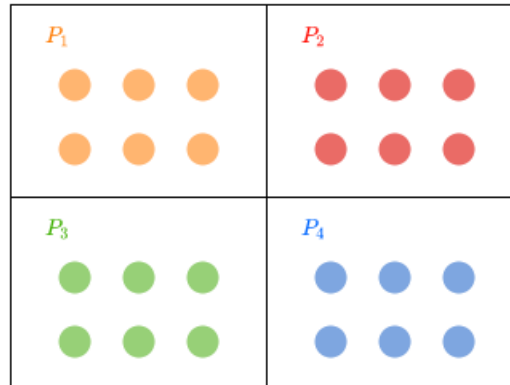
```

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

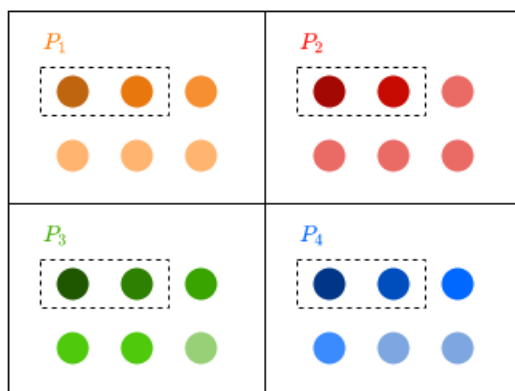
```

This algorithm is shown pictorially below

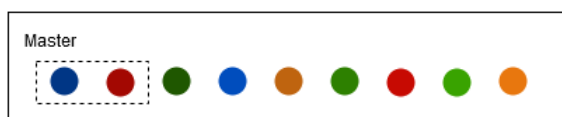
Step 1) Partition the data and send partitions to processes P_1, P_2, P_3, P_4



Step 2) Determine the closest samples within each process



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



3. Parallel Multiclass Perceptron

- Similar to a definition is class the perceptron is an algorithm that uses an affine function as a threshold boundary to classify real valued feature vectors
- The output of a perceptron is simply given as

$$f(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{w} \cdot \mathbf{x} + b > 0 \\ 0 & \text{otherwise} \end{cases}$$

where $\mathbf{w} \in \mathbb{R}^m$ with b acting as a bias

- A simple generalisation we can make to the perceptron is to allow the perceptron to predict from more than just two classes. This is known as the *multiclass perceptron*. A feature representation function maps a feature vector and its correspond label to a real valued vector. This output vector is then dot-producted with a weight vector \mathbf{w} , similar to the norm perceptron. The prediction value of an unknown feature vector for the multiclass perceptron can be computed as

$$\bar{y}(\mathbf{x}) = \arg \min_y \mathbf{w} \cdot f(\mathbf{x}, y)$$

- The iterative procedure to create our weight vector from the training set is also similar, although this time if we make an incorrect prediction we move the weight vector in the direction of $f(\mathbf{x}, y) - f(\mathbf{x}, \bar{y})$.
- Algorithm for the serial multiclass perceptron is shown in algorithm ???. Note, for the sake of simplicity a learning rate parameter has been omitted with in the training algorithm
- One obvious way to parallize the algorithm from ??? is to split the training set T into p partitions. One can then train an initial weight on each of these partitions independently on the p processes and then take a weighted mixture of the weights from each of the different processes. The mixture coefficients will be given as the vector $\boldsymbol{\mu}$ where $\mu_i \geq 0$ and $\sum_i \mu_i = 1$. Algorithm for this naive parallel multiclass perceptron is shown in algorithm ???.

Algorithm 3: Serial Multiclass Perceptron

input : Training data T , an initial weight vector $\mathbf{w}^{(0)}$

output: Updated weight vector

```
1  $k \leftarrow 0$ ;
2 while until convergence do
3   for  $(\mathbf{x}_i, y_i) \in T$  do
4      $\bar{y} \leftarrow \arg \min_y \mathbf{w}^{(k)} \cdot f(\mathbf{x}, y)$ ;
5      $k \leftarrow k + 1$ ;
6     if  $\bar{y} \neq y$  then
7        $\mathbf{w}^{(k+1)} \leftarrow \mathbf{w}^{(k)} + f(\mathbf{x}, y) - f(\mathbf{x}, \bar{y})$ ;
8   end
9 end
Result:  $\mathbf{w}^{(k)}$ 
```

Algorithm 4: Naive Parallel Multiclass Perceptron

input : Training data T , an initial weight vector $\mathbf{w}^{(0)}$

output: Updated weight vector

```
1  $k \leftarrow 0$ ;
2  $\{T_1, T_2, \dots, T_p\} \leftarrow$  an equal partition of  $T$ ;
3 for  $T_i \in \{T_1, T_2, \dots, T_p\}$  concurrently do
4    $\mathbf{w}^{(i)} = \text{SerialMulticlassPerceptron}(T_i, \mathbf{w}^{(0)})$ ;
5 end
6  $\mathbf{w} = \sum_i \mu_i \mathbf{w}^{(i)}$ ;
Result:  $\mathbf{w}$ 
```

- While ?? is easily understood and makes effective use of parallel resources without incurring too much over-head in inter-process communication, there is one large short-coming of this algorithm. Even for linearly separable datasets T , algorithm ?? does not always converge to a separating hyperplane.
- In fact in Question ??? you yourself will be asked to provide a linearly separable dataset for which algorithm ?? will not converge.
- To salvage this parallel algorithm we can make sure each processors is communicating

its results before finalization.

- That is we can collect weights from processors after each single epoch, mix these the weights together and then redistribute this new weight vector to our processes again for a new epoch.
- This new iterative parameter mixing idea is shown in algorithm.

Algorithm 5: Parallel Multiclass Perceptron

input : Training data T , an initial weight vector $\mathbf{w}^{(0)}$

output: Updated weight vector

```

1  $\mathbf{w} \leftarrow \mathbf{w}^{(0)}$ ;
2 for  $n \in \{1, 2, \dots, N\}$  do
3   for  $T_i \in \{T_1, T_2, \dots, T_p\}$  concurrently do
4      $\mathbf{w}^{(i,n)} = \text{OnePerceptronEpoch}(T_i, \mathbf{w})$ ;
5   end
6    $\mathbf{w} \leftarrow \sum_i \mu_i \mathbf{w}^{(i,n)}$ ;
7 end
Result:  $\mathbf{w}$ 

```

Algorithm 6: One Perceptron Epoch

input : Training data T_i , an initial weight vector $\mathbf{w}^{(0)}$

output: Updated weight vector

```

1  $k \leftarrow 0$ ;
2  $\mathbf{w} \leftarrow \mathbf{w}^{(0)}$ ;
3 for  $(\mathbf{x}_i, y_i) \in T_i$  do
4    $\bar{y} \leftarrow \arg \max_y \mathbf{w}^{(k)} \cdot f(\mathbf{x}, y)$ ;
5   if  $\bar{y} \neq y$  then
6      $\mathbf{w}^{(k+1)} \leftarrow \mathbf{w}^{(k)} + f(\mathbf{x}, y) - f(\mathbf{x}, \bar{y})$ ;
7      $k \leftarrow k + 1$ ;
8 end
Result:  $\mathbf{w}^{(k)}$ 

```

- The obvious and main drawback of algorithm ?? over ?? is that there is much more inter-process communication overhead in reducing $\mathbf{w}^{(i,n)}$ and sending new values of \mathbf{w} .
- However, we will see that algorithm ?? will in fact converge for linearly separable training sets.
- We can show this in the following Theorem.

Theorem Assuming a training set T is linearly separable by a margin λ , let $k_{i,n}$ be the number of mistakes that occurred on partition i of the training set during the n^{th} epoch of training. For any N , when training a perceptron with fixed parameter mixing in algorithm ??

$$\sum_{n=1}^N \sum_{i=1}^S \mu_{i,n} k_{i,n} \leq \frac{R^2}{\lambda^2}$$

Using some notation from before, let $\mathbf{w}^{(i,n)}$ be the weight acquired on the n^{th} epoch of training by process i and let $\mathbf{w}^{(i,n)-k}$ be the weight vector for process i k iterations back. Let $\mathbf{w}^{(avg,n)}$ be the weighted mixture weight for the n^{th} epoch, that is

$$\mathbf{w}^{(avg,n)} = \sum_{i=1}^S \mu_{i,n} \mathbf{w}^{(i,n)}$$

we can easily enough show that

$$\mathbf{u} \cdot \mathbf{w}^{(i,n)} \leq \mathbf{u} \cdot \mathbf{w}^{(avg,n-1)} + k_{i,n} \lambda \quad (\text{R1})$$

and

$$\|\mathbf{w}^{(i,n)}\|^2 \leq \|\mathbf{w}^{(avg,n-1)}\|^2 + k_{i,n} R^2 \quad (\text{R2})$$

Using ?? and ?? we can inductively show the following results

$$\mathbf{u} \cdot \mathbf{w}^{(i,n)} \geq \sum_{n=1}^N \sum_{i=1}^S \mu_{i,n} k_{i,n} \lambda \quad (\text{H1})$$

and

$$\|\mathbf{w}^{(i,n)}\|^2 \leq \sum_{n=1}^N \sum_{i=1}^S \mu_{i,n} k_{i,n} R^2 \quad (\text{H2})$$

we can see that ?? implies $\|\mathbf{w}^{(i,n)}\| \geq \sum_{n=1}^N \sum_{i=1}^S \mu_{i,n} k_{i,n} \lambda$ since $\mathbf{v} \cdot \mathbf{w} \leq \|\mathbf{u}\| \|\mathbf{w}\|$ by the Cauchy Schwarz inequality and $\|\mathbf{u}\| = 1$. For the base case $\mathbf{w}^{(avg,1)}$. So using ?? and the

fact that $\mathbf{w}^{(avg,0)} = 0$ and

$$\mathbf{u} \cdot \mathbf{w}^{(avg,1)} = \sum_{i=1}^S \mu_{i,1} \mathbf{u} \cdot \mathbf{w}^{(i,1)} \geq \sum_{i=1}^S \mu_{i,1} k_{i,n} \lambda$$

To show ??, the base case can be written as

$$\begin{aligned} \|\mathbf{w}^{(avg,1)}\|^2 &= \left\| \sum_{i=1}^S \mu_{i,1} \mathbf{w}^{(i,1)} \right\|^2 \\ &= \sum_{i=1}^S \mu_{i,1} \|\mathbf{w}^{(i,1)}\|^2 \\ &= \sum_{i=1}^S \mu_{i,1} k_{i,n} R^2 \end{aligned}$$

This expression is found by first applying Jensen's inequality followed by ?? as well as the fact that $\|\mathbf{w}^{(avg,0)}\|^2 = 0$. Now considering the general case, $\|\mathbf{w}^{(avg,N)}\|^2$

$$\begin{aligned} \mathbf{u} \cdot \mathbf{w}^{(avg,N)} &= \sum_{i=1}^S \mu_{i,N} (\mathbf{u} \cdot \mathbf{w}^{(i,N)}) \\ &\geq \sum_{i=1}^S \mu_{i,N} (\mathbf{u} \cdot \mathbf{w}^{(i,N-1)} + k_{i,N} \lambda) \quad (\text{Using ??}) \\ &= \mathbf{u} \cdot \mathbf{w}^{(avg,N-1)} + \sum_{i=1}^S \mu_{i,N} \mathbf{u} \cdot \mathbf{w}^{(i,N-1)} k_{i,N} \lambda \\ &\geq \left[\sum_{n=1}^{N-1} \sum_{i=1}^S \mu_{i,n} k_{i,n} \gamma \right] + \sum_{i=1}^S \mu_{i,N} k_{i,N} \\ &= \sum_{n=1}^N \sum_{i=1}^S \mu_{i,n} k_{i,n} \gamma \end{aligned}$$

Which thusly proves ???. To show the inductive case for ?? we have

$$\begin{aligned}
\|\mathbf{w}^{(\text{avg}, N)}\|^2 &\leq \sum_{i=1}^S \mu_{i,N} \|\mathbf{w}^{(i,N)}\|^2 \\
&= \sum_{i=1}^S \mu_{i,N} \left(\|\mathbf{w}^{(\text{avg}, N-1)}\|^2 + k_{i,N} R^2 \right) \\
&= \|\mathbf{w}^{(\text{avg}, N-1)}\|^2 + \sum_{i=1}^S \mu_{i,N} k_{i,N} R^2 \\
&\leq \left[\sum_{n=1}^{N-1} \sum_{i=1}^S \mu_{i,n} k_{i,n} R^2 \right] + \sum_{i=1}^S \mu_{i,N} k_{i,N} R^2 \\
&= \sum_{n=1}^N \sum_{i=1}^S \mu_{i,n} k_{i,n} R^2
\end{aligned}$$

This shows ???. Combining the results of ?? and ?? as well as the fact $\|\mathbf{w}^{(\text{avg}, N)}\| > \mathbf{u} \cdot \mathbf{w}^{(\text{avg}, N)}$

$$\begin{aligned}
\left[\sum_{n=1}^N \sum_{i=1}^S \mu_{i,n} k_{i,n} \right]^2 \gamma^2 &\leq \left[\sum_{n=1}^N \sum_{i=1}^S \mu_{i,n} k_{i,n} \right] R^2 \\
\sum_{n=1}^N \sum_{i=1}^S \mu_{i,n} k_{i,n} &\leq \frac{R^2}{\lambda^2}
\end{aligned}$$

This inequality tells us that if the weight are distributed uniformly then the number of mistakes is bounded and that there is enough convergence between iterations to guarantee convergence.

4. 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 for our model to learn off

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

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

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

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

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

Let's introduce 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 get around this? One method is to soundly combine models from different processes in the hopes of achieving a weight vector had the SGD was instead 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 the 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 [?] 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 summarised in algorithm ??

Algorithm 7: 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:  $\mathbf{w}_p$ 
```

5. 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 cores and 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. To log into **getafix** we can simply **ssh** via terminal (you will need to be on campus to do this or use a VPN).

If you're working from home you'll need to run a few commands before we can start coding. **getafix** does not have many python libraries installed for users, so we will need to install them ourselves

We can start by creating a folder to put all of our code in, changing to that directory and creating our first python file using vim.

If you're unsure of what any of these commands do, you can simply look at the manual page for that command by typing `man` followed by the command into your terminal. For example `man mkdir` will bring up the manual page for the `mkdir` command used above. As a simple parallel code example, let's populate our newly created `parallel_sum.py` with a python script that will sum together the contents of a numpy array in parallel. You can find this code in the following git repository under the `src` folder https://github.com/Michael1CC/STAT4402_TUTORIAL. Exit out of vim by pressing `ctrl+c`, typing `:wq` into the prompt at the bottom and then hitting enter. This is great but how do we submit this script to be executed by the HPC? It wouldn't make sense to execute it in the login node we are current in since we wouldn't be enjoying the immense computing power offered by the HPC. To submit jobs to the HPC we need to create a job file, commonly referred to as a batch file, which specifies a whole bunch of parameters and conditions that the HPC will use to run our code. These parameters include things such as how many cores to use, how much memory to use, what programs to execute, where to put the output should go etc. An example batch file to run our `parallel_python.py` file is shown below. Saving this batch file as `parallel_python_batch.sh` we can submit our job to the HPC, and then checkout or standard out and standard error files generated by the submitted job. The standard out and error files filenames and paths are specified by the `#SBATCH --output` and `#SBATCH --error` parameters in the batch file, respectively. As you can see from our `parallel_python.out`, we get an array sum of 4928.8167 which is what we would expect for an array of 10000 with elements that are generated using $U[0, 1]$. If you get any errors from this example, try checking your `parallel_python.err` to see what sort of errors are being thrown.

6. Questions

Q1) a) Suppose the training set

$$T = \{([0, 1], 0), ([1, 1], 0), ([2, 1], 0), ([2, 2], 0), ([5, 8], 1), ([6, 7], 1), ([5, 7], 1), ([4, 5], 1)\}$$

was partitioned into

$$T_1 = \{([2, 1], 0), ([2, 2], 0), ([5, 8], 1), ([6, 7], 1)\}$$

$$T_2 = \{([2, 1], 0), ([2, 2], 0), ([5, 8], 1), ([6, 7], 1)\}$$

for a k-NN parallel algorithm using the 2 closest neighbours and the Manhattan distance metric to classify the an unknown sample $\mathbf{x} = [0, 0]$ use two processes.

What 2 closest neighbours would be returned by each of the processors?

- b) Would the parallel k-NN algorithm still work if only $k - 1$ closest neighbours was instead returned by each process?
- c) i. Implement your own parallel k-NN algorithm from scratch in python, only using the `threading` and `queue` libraries to perform parallelism. Use `sklearn`s iris dataset to test your model, you can access this data set using the import `from sklearn import datasets` and then creating the dataset by using `datasets.load_iris()`.
- ii. Assuming the completion of the question above, implementing these parallel algorithms yourself can be tedious and error prone. Fortunately many parallel Machine Learning algoirthms and models have been implemented for us already in `python`s `sklearn` library. Like `pytorch`, `sklearn` is free to use python library that implements many of the Machine Learning tools seen in this course. Use `sklearn`s `KNeighborsClassifier` class to train on the *breast cancer* data. Leave $\frac{1}{5}$ of the data aside to provide a test accuracy. Submit this script to `golith` using a job batch file. Use repeat training with $1 - 4$ threads and record your results. *Hint:* you will need to adjust the `n_jobs` parameter in the `KNeighborsClassifier` to chnage the number of threads used for running the k-NN algorithm.

Number of threads	1	2	3	4
Time (sec)	0.03754734	0.02993416	0.0252390	0.01833604

7. Answers

Q1)

a)

The neighbours to be returned by the first and second process would be $\{([2, 1], 0), ([2, 2], 0)\}$ and $\{([2, 1], 0), ([2, 2], 0)\}$ respectively.

7.0.1. b)

No, all k closest neighbours may randomly be given to a single process. If some of these k closest neighbours were neglected by this process then not all the k closest neighbours will be used in the final reduction step which may give differing results to the serial k-NN algorithm.

7.0.2. c)

7.0.3. i)

The code for the k-NN algorithm implemented from scratch is shown below

7.0.4. ii)

An example batch file for job submission is shown below

Code to training a parallel k-NN algorithm using `sklearn` is given below

The times for running the parallel algorithm with different numbers of processes is shown below