

Strategy (same for both datasets):

Master node, rank==0 is

- loading the datasets
- Normalizing if needed,
- Splitting data into train and test,
- We'll concatenate y_train and x_train in one array so it is easier to send to other workers,

A fixed seed number is considered to provide the same random number every time.

We make random numbers for weights in a specific range and size depending on the dataset columns, e.g range = 0,5 or size = (1,261).

Learning rate and a fixed number of epochs and lists for the losses are considered.

- Depending on the number of epochs we choose, we loop and check if it is converged or not.
 - In each epoch, when we are in the master process, the data is splitted by the number of workers we have .
 - Each data part is sent via scatter,
 - But the weights is sent via broadcast as it is the same for all the workers.
 - In each worker the part of the data received is shuffled.
 - For each instance/row of the data received in each worker.
 - SGD is calculated.
 - Weights are updated.
 - After all the instance weights are updated.
 - Master node gathers all the weights, predicted y and true y.
 - Master node finds the mean of the weights, calculates the RMSE and append it to our list for train and it uses new mean weights to predict new values for test data and computes the RMSE of the true test target and predicted ones to append it to out test loss list.
 - If the test loss of each instance is less than a specific value then it is converged, it means it can't improve more.
 - If it is then broadcast the converge value so the program stops.
 - And finally we save the losses and show the final time.
-

Virus dataset.

The STRATEGY is the same as above.

For this dataset, we define load_virus() functions which goes into the dataset directory and list all of them into a list called files.

We go through each file and append them in another list, later we concat them to make a bigger dataframe.

The first value of each input is our target y, all the others, there is a key and a value pointing to a feature and its values exist here, we make a sparse matrix and go through each file and split them the first value is the column name and the second one is the value of that column.

We make ones as bias and add them to x and we have our final datasets.

```
20 def load_virus():
21     #freeing up the memory
22     gc.collect()
23     path = 'dataset/dataset/'
24     files = os.listdir(path)
25     li = []
26     for filename in files:
27         df = pd.read_csv(path+filename, header=None)
28         li.append(df)
29
30     dat_file = pd.concat(li, axis=0, ignore_index=True)
31
32     y = np.zeros((107856,1), float)
33     x = np.zeros((107856,479))
34     for i in dat_file.index:
35         doc = dat_file.iloc[i][0].strip().split()
36         y[i] = float(doc[0])
37
38         for element in doc[1:]:
39             k, v = element.split(":")
40             x[i,int(k)] = int(v)
41     bias = np.ones((107856,1))
42     x = np.append(bias, x, axis=1)
43
44     return x,y
```

The rest of the functions are as below.

- train_test_split() makes data into two splits by chosen fraction.
- row_normalizer(), normalize each row.
- predic(), predict new y based on x and weights received.
- sgd(), calculates the main part of our program, stochastic gradient descent.
- calculate_RMSE(), calculate RMSE between prediction and true values.

```

46 #splitting dataset
47 def train_test_split(data ,frac=0.7):
48     data = pd.DataFrame(data)
49     #train_test split
50     train = data.sample(frac=frac, random_state=10)
51     test = data.drop(train.index)
52
53     return train.to_numpy(), test.to_numpy()
54
55 #normalizing based on rows
56 def row_normalizer(x: np.ndarray):
57     return x/np.linalg.norm(x, ord=2, axis=1, keepdims=True)
58
59 #prediciton
60 def predict(x, betas):
61     return np.matmul(x, betas.T)
62
63 #SGD function
64 def sgd(x, weight, y, y_hat, lr):
65
66     #Computing Derivatives using Chain Rule
67     dv_loss = x * (-2*(y - y_hat) )
68     #Computing new weights
69     weight = weight - lr * dv_loss
70     return weight
71
72 #calculate_RMSE
73 def calculate_RMSE(true_y, pred_y):
74     return np.sqrt(mt.mean_squared_error(true_y, pred_y))
75

```

Starting from the master node.

```

10
79 if rank == 0:
80     print('Loading data...', flush=True)
81     x,y = load_virus()
82     print(f"X shape: {x.shape}")
83     x = normalize_rows(x)
84     print('Splitting into train and test sets...', flush=True)
85     x_train, x_test = train_test_split(x ,frac=0.7)
86     y_train, y_test = train_test_split(y ,frac=0.7)
87
88     data = np.concatenate((y_train.reshape(-1,1), x_train), axis=1)
89
90
91 #Initializing weights
92
93 w = np.random.randint(0,2,size=((1,480)))
94
95 #learning rate
96 lr = 0.00001
97
98 #Setting the number of epochs
99 N = 100
100 train_loss = []
101 test_loss = []
102
103 converged = False #For checking convergence
104 epoch = 0
105
106 if rank == 0:
107     #Starting time
108     t0 = MPI.Wtime()
109

```

```

110 #Training for N Epochs
111 while (not converged) and (epoch < N):
112
113     if rank == 0:
114         #print onlt every 10 steps
115         if epoch%10==0:
116             print("Epoch : ", epoch, flush=True)
117
118         #splitting
119         data_part = np.array_split(data,size)
120     else:
121         data_part = None
122     #Sending data using scatter and weights via bcast
123     recvd_dt = comm.scatter(data_part,root=0)
124     weights = comm.bcast(w,root=0)
125     row_, col_ = np.shape(recvd_dt)
126     y_hat = np.zeros((row_,1))
127
128     #Shuffling the received dataset at each worker
129     np.take(recvd_dt,np.random.permutation(recvd_dt.shape[0]),axis=0,out=recvd_dt)
130     #Epoch
131     for i in range(0, row_):
132         X_part = recvd_dt[i,1:]
133         y_part = recvd_dt[i,0]
134         #Prediction of instance i, forward pass
135         y_hat[i,0] = predict(X_part, weights.T)
136
137         #Computing new weights
138         weights = sgd(X_part, weights, y_part, y_hat[i,0], lr)
139
140     #Gathering weights at Worker 0
141     w_r = comm.gather(weights, root=0)
142     #Gathering predicted output (yy) and original data (true_dt) at Worker 0
143     y_hat_recvd = comm.gather(y_hat, root=0)
144     true_dt = comm.gather(recvd_dt[:,0], root=0)
145
146     if rank == 0 and converged == False:
147         #Worker 0 computing mean of the weights
148         w = np.mean(w_r, axis=0)
149
150         pred_y = np.vstack(y_hat_recvd)
151         true_y = np.hstack(true_dt)
152
153         #Worker 0 computing calculate RMSE
154         calculate_RMSE_train = calculate_RMSE(true_y, pred_y)
155         train_loss.append(calculate_RMSE_train)
156
157         #test prediction
158         pred_y_test = predict(x_test, w.T)
159         #Computing RME Testing
160         calculate_RMSE_test = calculate_RMSE(y_test, pred_y_test)
161         test_loss.append(calculate_RMSE_test)
162
163         #Checking for convergence in RMSE test
164         if epoch > 0 and abs((test_loss[epoch-1] - test_loss[epoch])) < 10**-6:
165             print(test_loss[epoch-1]- test_loss[epoch], flush=True)
166             converged = True
167     #the converged flag is sent back for stopping epochs
168     converged = comm.bcast(converged,root=0)
169     epoch += 1
170
171 if rank == 0:
172     #Ending time
173     t1 = MPI.Wtime()
174     df = pd.DataFrame(list(zip(train_loss, test_loss)),columns=['train_loss', 'test_loss'])
175     df.to_csv(f"loss_virus{size}.csv", index=False)
176     #t1 = MPI.Wtime()
177     time_final = t1 - t0
178     print("With P = {} Workers, the process took {} seconds".format(size, time_final), flush=True)

```

Here is the screenshot of running with 4 workers.

```
1 !mpirun -n 4 python virus.py
```

Loading data...

X shape: (107856, 480)

Splitting into train and test sets...

Epoch : 0

Epoch : 10

Epoch : 20

Epoch : 30

Epoch : 40

Epoch : 50

Epoch : 60

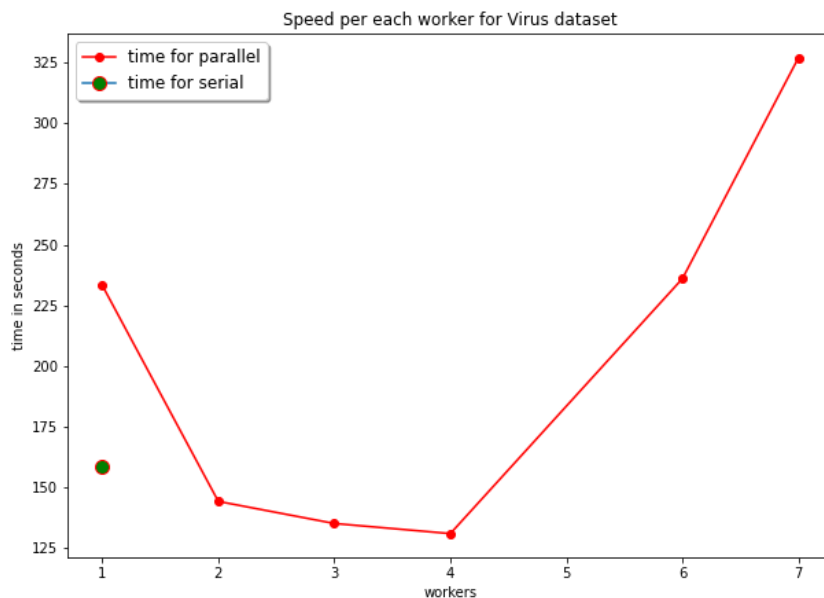
Epoch : 70

Epoch : 80

Epoch : 90

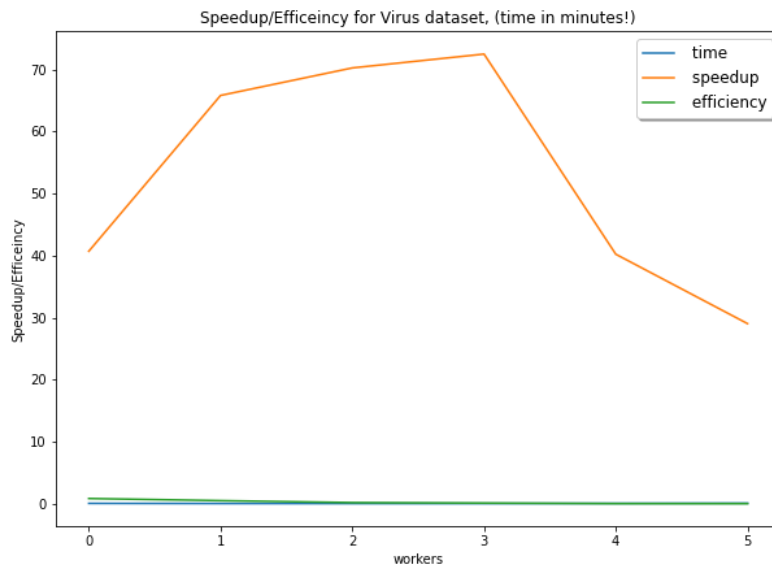
With P = 4 Workers, the process took 130.96906805038452 seconds

Let's look at the time elapsed using different workers.



Looking at the graph we see that the sequential code is faster than parallel code using only one worker. Using 2 workers is definitely faster than one and so until 4th workers. Using 6th and more is not fast anymore as my machine has only 4 cores.

Let's look at the speedup/efficiency graph.



Looking at the above graph we see that we have a huge speed up until using 3 workers and then the speed goes down. The efficiency overall goes down. But more details are observed in this table.

```

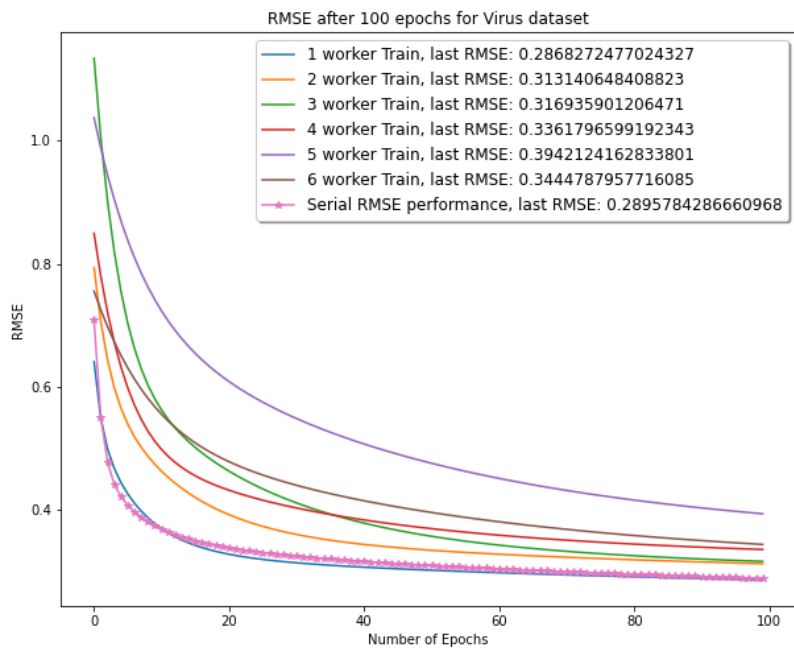
: 1 time_virus
:

```

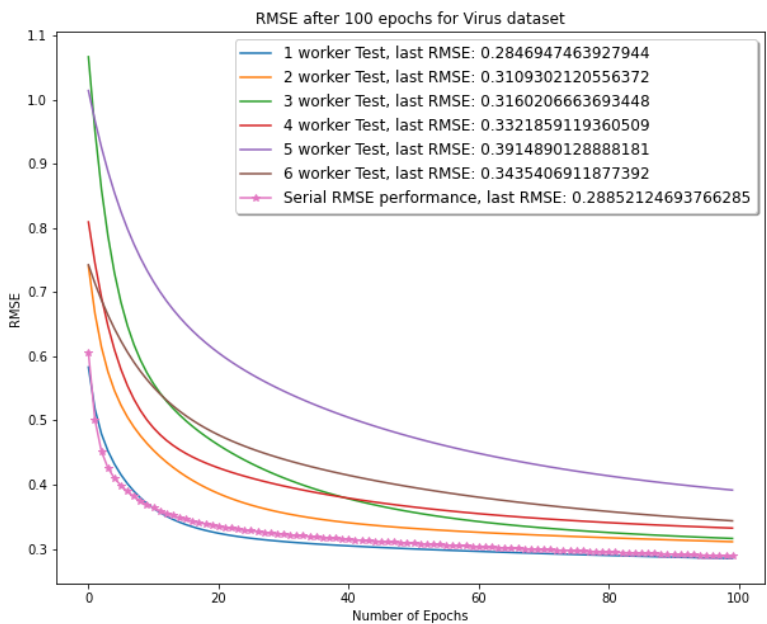
	workers	time	speedup	efficiency
0	1	0.064786	40.718604	0.848725
1	2	0.040067	65.840266	0.504409
2	3	0.037536	70.278991	0.186993
3	4	0.036380	72.511816	0.113861
4	6	0.065559	40.238802	0.027713
5	7	0.090825	29.044867	0.019868

RMSE Performance.

Looking at the graph the RMSE is going down and we have a good accuracy. But still sequential performance and using one worker is giving us better accuracy than more workers.



Test RMSE Observation:



Even in the test set we see better performance in serial code and using one core, also 3 workers give us a good RMSE score.

- Using one worker it didn't converge but we have a good RMSE score.

- 2 workers did not converge, RMSE is worse than using one worker.
 - 3 workers did not converge, RMSE is worse than using 1/2 workers.
 - 4 workers did not converge, RMSE is worse than before.
 - 6 workers did not converge, RMSE is worse than before.
 - 7 workers did not converge, RMSE is worse than before.
-

KDD dataset.

For this data set to save time I first concat the three datasets for learning and validation and validation target in one data frame.

Replaced “ ” with 0 and applied one hot encoding for features which are categorical plus the columns that have less than 5 unique values, they are converted to categorical features and one hot encoding applied.

Next I checked the Pearson correlation and those which had positive correlation were chosen. Finally we have a data set with shape (191779,261) bias added.

First we define some functions:

- `load_data_kdd()` to load csv file we already cleaned and preprocessed, it adds bias to it and create x and y matrix.
- `train_test_split()` makes data into two splits by chosen fraction.
- `row_normalizer()`, normalize each row.
- `predic()`, predict new y based on x and weights received.
- `sgd()`, calculates the main part of our program, stochastic gradient descent.
- `calculate_RMSE()`, calculate RMSE between prediction and true values.


```

20 #reading dataset
21 def load_data_kdd(data = "kdd_final_filtered.csv"):
22     #freeing up the memory
23     gc.collect()
24     #reading the saved dataset
25     kdd = pd.read_csv(data, low_memory=False)
26     #target
27     y = kdd["TARGET_D"].to_numpy()
28     x = kdd.drop(columns=["TARGET_D"], axis=1)
29
30     bias = pd.DataFrame(np.ones(x.shape[0]).reshape(-1,1), columns = ["bias"])
31     x = pd.concat([bias, x], axis=1).to_numpy()
32
33     return x, y
34
35 #splitting dataset
36 def train_test_split(data ,frac=0.7):
37     data = pd.DataFrame(data)
38     #train test split
39     train = data.sample(frac=frac, random_state=10)
40     test = data.drop(train.index)
41
42     return train.to_numpy(), test.to_numpy()
43
44 #normalizing based on rows
45 def row_normalizer(x: np.ndarray):
46     return x/np.linalg.norm(x, ord=2, axis=1, keepdims=True)
47
48 #prediciton
49 def predict(x, betas):
50     return np.matmul(x, betas.T)
51
52 #SGD function
53 def sgd(x, weight, y, y_hat, lr):
54
55     #Computing Derivatives using Chain Rule
56     dv_loss = x * (-2*(y - y_hat) )
57     #Computing new weights
58     weight = weight - lr * dv_loss
59
60     return weight
61
62 #calculate RMSE
63 def calculate_RMSE(true_y, pred_y):
64     return np.sqrt(mt.mean_squared_error(true_y, pred_y))

```

The rest of the code is the same as explained in the STRATEGY.

```

69 if rank == 0:
70     print('Loading data...', flush=True)
71     x,y = load_data_kdd()
72     print(f"X shape: {x.shape}")
73     #print("Normalizing")
74
75     x = row_normalizer(x)
76
77     print('Splitting into train and test sets...', flush=True)
78     x_train, x_test = train_test_split(x ,frac=0.7)
79     y_train, y_test = train_test_split(y ,frac=0.7)
80
81     data = np.concatenate((y_train.reshape(-1,1), x_train), axis=1)
82
83
84     #fixing seed
85     np.random.seed(2021)
86     #random weights
87     w = np.random.randint(0,5,size=((1,261)))
88
89     #learning rate
90     lr = 0.00001
91
92     #Setting the number of epochs
93     N = 100
94     train_loss = []
95     test_loss = []
96
97     converged = False #For checking convergence
98     epoch = 0
99
100 if rank == 0:
101     #Starting time
102     t0 = MPI.Wtime()
103
104
105 while (not converged) and (epoch < N):
106
107     if rank == 0:
108         #print onlt every 10 steps
109         if epoch%10==0:
110             print("Epoch : ", epoch, flush=True)
111
112         #splitting
113         data_part = np.array_split(data,size)
114     else:
115         data_part = None
116     #Sending data using scatter and weights via bcast
117     recvd_dt = comm.scatter(data_part,root=0)
118     weights = comm.bcast(w,root=0)
119     row_, col_ = np.shape(recvd_dt)
120     y_hat = np.zeros((row_,1))
121
122     #Shuffling the received dataset at each worker
123     np.take(recvd_dt,np.random.permutation(recvd_dt.shape[0]),axis=0,out=recvd_dt)
124     #Epoch
125     for i in range(0, row_):
126         X_part = recvd_dt[i,1:]
127         y_part = recvd_dt[i,0]
128         #Prediction of instance i, forward pass
129         y_hat[i,0] = predict(X_part, weights.T)
130
131         #Computing new weights
132         weights = sgd(X_part, weights, y_part, y_hat[i,0], lr)
133
134     #Gathering weights at Worker 0
135     w_r = comm.gather(weights, root=0)
136     #Gathering predicted output (yy) and original data (true_dt) at Worker 0
137     y_hat_recvd = comm.gather(y_hat, root=0)
138     true_dt = comm.gather(recvd_dt[:,0], root=0)
139

```

```

40 if rank == 0 and converged == False:
41     #Worker 0 computing mean of the weights
42     w = np.mean(w_r, axis=0)
43
44     pred_y = np.vstack(y_hat_recvd)
45     true_y = np.hstack(true_dt)
46
47     #Worker 0 computing calculate_RMSE
48     calculate_RMSE_train = calculate_RMSE(true_y, pred_y)
49     train_loss.append(calculate_RMSE_train)
50
51     #test prediction
52     pred_y_test = predict(x_test, w.T)
53     #Computing RME Testing
54     calculate_RMSE_test = calculate_RMSE(y_test, pred_y_test)
55     test_loss.append(calculate_RMSE_test)
56
57     #Checking for convergence in RMSE test
58     if epoch > 0 and abs((test_loss[epoch-1] - test_loss[epoch])) < 10**-6:
59         print(test_loss[epoch-1] - test_loss[epoch], flush=True)
60         converged = True
61     #the converged flag is sent back for stopping epochs
62     converged = comm.bcast(converged, root=0)
63     epoch += 1
64
65 if rank == 0:
66     #End time
67     t1 = MPI.Wtime()
68     df = pd.DataFrame(list(zip(train_loss, test_loss)), columns=['train_loss', 'test_loss'])
69     df.to_csv(f"loss{size}.csv", index=False)
70     #t1 = MPI.Wtime()
71     time_final = t1 - t0
72     print("With P = {} Workers, the process took {} seconds".format(size, time_final), flush=True)
73

```

Running the code.

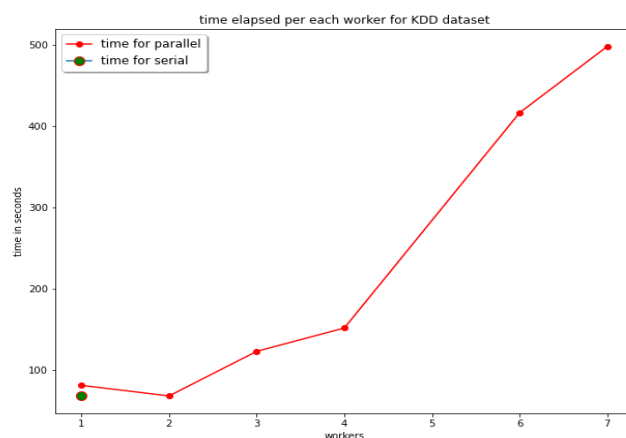
```
1 !mpiexec -n 2 python kdd_.py
```

```

Loading data...
X shape: (191779, 261)
Splitting into train and test sets...
Epoch : 0
Epoch : 10
Epoch : 20
8.017752826106062e-07
With P = 2 Workers, the process took 68.61729502677917 seconds

```

Looking at the time elapsed for each time with different workers.

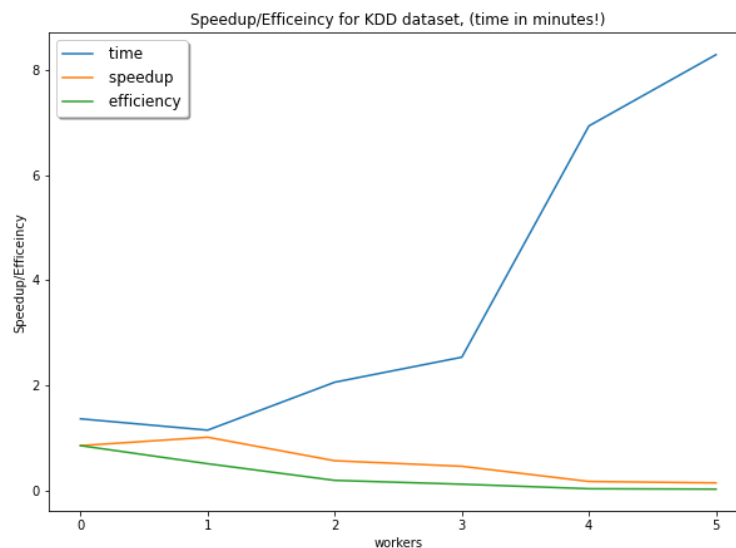


- Looking at this graph we see that using 2 workers gave us the quickest run time, as my machine has only 4 cores it is understandable why it takes longer time for more workers.
- Below is the total time, speedup and efficiency for this program.

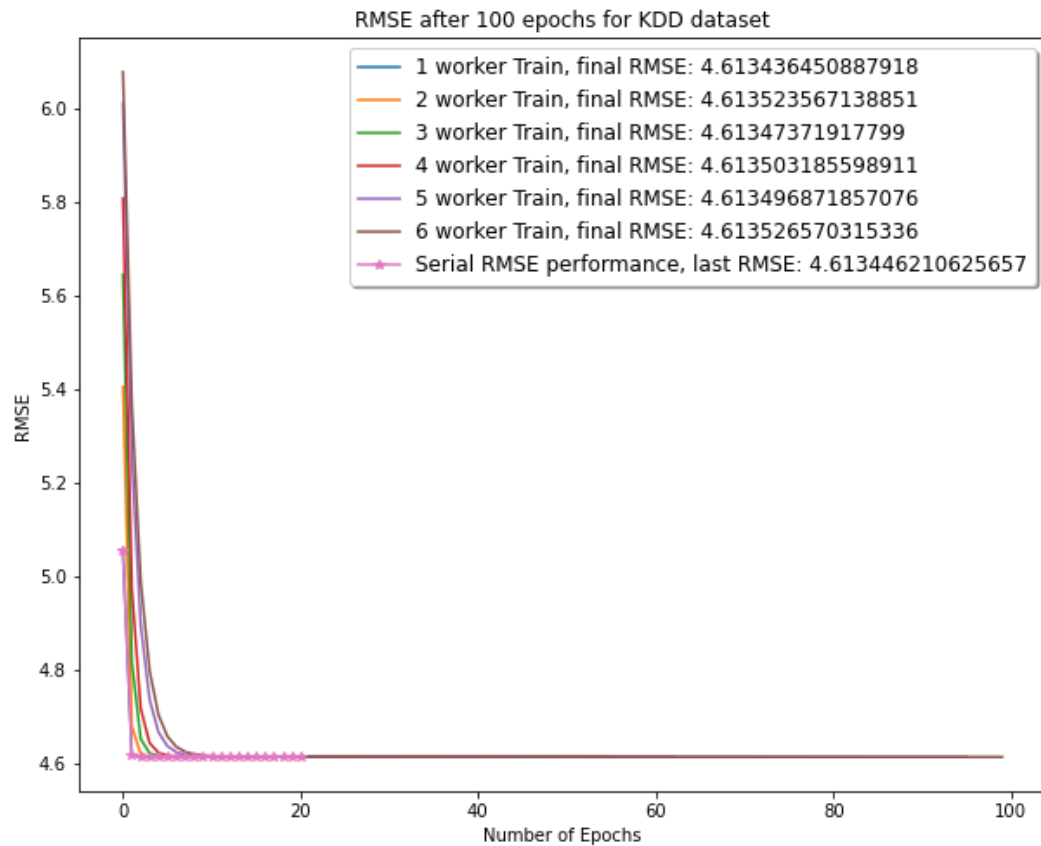
	workers	time	speedup	efficiency
0	1	1.359333	0.848725	0.848725
1	2	1.143617	1.008817	0.504409
2	3	2.056583	0.560979	0.186993
3	4	2.533133	0.455444	0.113861
4	6	6.938417	0.166277	0.027713
5	7	8.295667	0.139073	0.019868

-
- Below is the graph for speedup and efficiency , time is divided by 60 seconds so we can see more details.
- The efficiency overall goes down with the number of processes increasing! I don't know why this happens, the time for the serial program is 69.222 here how I calculated the speedup and efficiency.
- The speedup increase and later it decreases.

```
1 time["speedup"] = (69.222)/time_.time
2 time["efficiency"] = time["speedup"]/time_.workers
```



Looking at RMSE result for train and test sets.



- Using one worker it converged at the 23rd epoch.
- Using 2 workers it converged at the 30th epoch.
- Using 3 workers it converged at the 56th epoch
- Using 4 workers it converged at the 63rd epoch.
- Using 6 workers it converged at the 96th epoch
- Using 7 workers it did not converge.

Overall the RMSE is going lower but i could not get a lower RMSE even training for longer epochs and different learning rates.

Below is the testing RMSE values for KDD dataset. The behavior is as the same as train set.

