

Parallel Principle Component Analysis using Covariance Matrix and SVD in CUDA

Outlines:

1. Aim

To reduce the dimensionality of the input data such that we don't lose valuable information and also reduce the computation overhead. This would be achieved using standard PCA algorithm on top on which we apply parallelism (as explained in the sections below).

2. SVD

The main objective in the PCA is to calculate the eigen vector and eigen values of the covariance matrix (of the input data). This vector matrix is then used to project the original data with higher dimension to a lower dimension data.

We will decompose our covariance matrix using the svd as implemented in cuSolver library.

<https://docs.nvidia.com/cuda/cusolver/index.html> (<https://docs.nvidia.com/cuda/cusolver/index.html>).

3. Analyse the optimal dimension

After we have calculated eigen values and eigen vectors, we aim to find the best reduced dimension which can represent our data without losing useful information. This is done by using a threshold value and using the inequality as below.

$$\frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^n \lambda_i} > \theta$$

Where k is the reduced number of dimensions and n is the total number of dimensions and theta is a threshold value. Note that the numerator picks up the eigen values one by one in descending order.

4. Projected data

The projected data is constructed by using the k eigen vectors (eigen vectors corresponding to the eigen values used in the inequality above).

$$ProjectedData = XV[:, : k]$$

Here X is the original data and V is the eigen vector matrix.

5. Comparing performance

Applying some Machine Learning algorithms (Classification, Clustering, etc.) and comparing the results between Projected Data and Original Data. This will give us a measure of how good we can perform even with reduced dimensions and in turn saving much precious computation time.

6. Algorithm

1. Read the input data(intrusion.csv) into a 2d- matrix.
2. Normalize the data by calculating mean and standard deviation. Let's call this data matrix A.
3. Calculating Covariance Matrix.($\text{matmul}(A_transpose, A)$). Let's call this matrix B.
4. Compute the SVD of B using cuSolver.
5. Find the optimal dimension using method explained in section 3. (using threshold of 90%)
6. Project the data with reduced dimension as explained in section 4.
7. Store this data from step 6 into a csv file.

Apply appropriate computations on this reduced dimension csv file and compare the results with original csv file.

6.Results

Serialized code took around 1.2s.

CUDA parallelized code took around 2.3s.

This is down to the multiple memory copy instructions from device to host and vice versa.