

# Accelerating your Python Code For GMMs with PyCUDA

Varun Nayyar

27/07/18

# Outline

Me = Math Major + Script Kiddy (Manage Expectations)

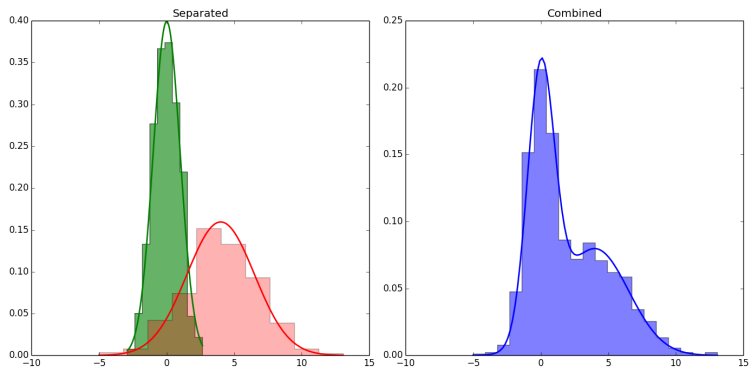
## What to Expect

- Some Math
- Iterative Process I went through
- Thinking with CUDA and Basic Syntax
- How to use PyCUDA to avoid complicated work

## What NOT to Expect

- How PyCUDA does it's magic
- Intermediate/Advanced CUDA

# Gaussian Mixture Models $d=1$ , $K=2$



K-means+=1

# Gaussian Mixture Models (GMMs)

## Density Function

For  $K$  mixtures

$$f(\mathbf{x}|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\mu_k, \Sigma_k)$$

## Log Likelihood (function of concern)

$$l(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{x}) = \sum_{i=1}^N \ln \left( \sum_{k=1}^K \pi_k \mathcal{N}(x_i|\mu_k, \Sigma_k) \right)$$

# Need for speed

## Some post-hoc realizations

- GMM likelihood formula doesn't decompose into a mathematically simple form
- However, note that the GMM likelihood has a parallelizable form in that each point of each mixture is independent (CUDA vibes)

## Computational Numbers

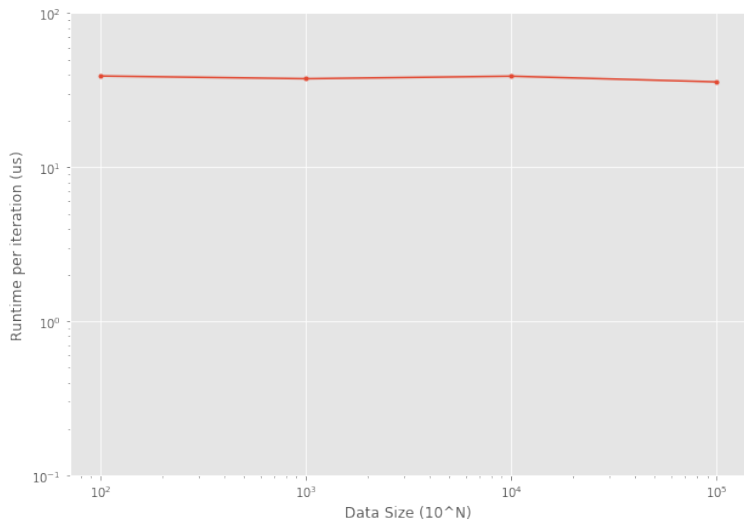
- Number of flops are of the order of  $O(NKd)$ , in my case,  $N = 10^6$ ,  $K = 8$ ,  $d = 13$ . i.e.  $O(10^8)$
- I needed to evaluate the likelihood  $10^6$  times for a fixed dataset while the parameters were varied. (Markov Chain Monte Carlo)
- i.e  $10^{14}$  floating point operations per run.

## First Attempt

- Eh, my computer is fast enough
- Pure Python (numpy)
- Took 36 us per datapoint, or 36s for whole dataset
- $10^6$  evaluations would take 1 year



# Execution speed per N



## Second Attempt

- Stand on the shoulders of giants (scikit-learn)
- Reverse engineered the likelihood evaluator
- 72x improvement!
- Took 0.5 us per datapoint, or 5s for whole dataset
- $10^6$  evaluations would take 5 days!





# Execution speed per N

