

CS423 project Pitch

Group 4 - Taeyoung Kim¹, Janggun Lee², Geon Park²

April 27, 2020

¹School of Computing

²Department of Mathematical Sciences

Table of Contents

1. Gaussian Mixture Model
2. Mixture of Experts
3. Experiment
4. Advantages of using ProbProg

Gaussian Mixture Model-1

GMM(Gaussian Mixture Model) is a well known generative model in machine learning. It views data as mixture of some multivariate normal random variable.

$$p(x_i) = \sum_{k=1}^K \pi_k \mathcal{N}(x_i | \mu_k, \Sigma_k), \sum_{k=1}^K \pi_k = 1$$

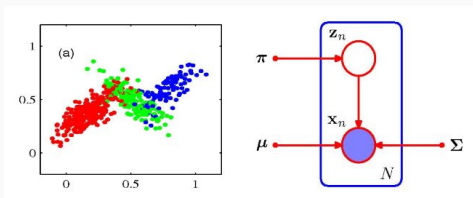


Figure 1: Visualization of GMM and model structure

Gaussian Mixture Model-2

Unlike simplicity of model itself, deriving learning algorithm is not simple. It uses **EM(Expectation-Maximization)** algorithm

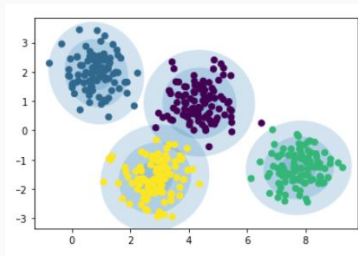
Step E : With current parameters, compute value $p(z_{ik} = 1|x_i)$, which is probability that data x_i is contained in k -th cluster.

$$p(z_{ik} = 1|x_i) = \frac{\pi_k \mathcal{N}(x_i|\mu_k, \Sigma_k)}{\sum_{l=1}^K \pi_l \mathcal{N}(x_i|\mu_l, \Sigma_l)}$$

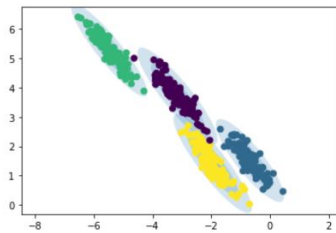
Step M : Update parameters with computed values $p(z_{ik} = 1|x_i)$.

$$\begin{aligned}\mu_k &= \sum_{i=1}^N \frac{p(z_{ik} = 1|x_i)}{\sum_{j=1}^N p(z_{jk} = 1|x_j)} x_i \\ \Sigma_k &= \sum_{i=1}^N \frac{p(z_{ik} = 1|x_i)}{\sum_{j=1}^N p(z_{jk} = 1|x_j)} (x_i - \mu_k)(x_i - \mu_k)^T \\ \pi_k &= \frac{1}{N} \sum_{i=1}^N p(z_{ik} = 1|x_i)\end{aligned}$$

Gaussian Mixture Model-3



(a) 4-component GMM with circular clusters



(b) 4-component GMM with squashed clusters

Figure 2: Examples of GMM in different dataset

Guassian Mixture Model-4

GMM in Anglican

```
1 (with-primitive-procedures [row-mean invert shape identity-matrix get-row]
2   (defquery gmm [data & [hyperparams]]
3     (let [[N D] (shape data)
4           K      (:K      hyperparams 10)
5           alpha  (:alpha  hyperparams 1.0)
6           mu-0   (:mu-0   hyperparams (row-mean data))
7           lambda-0 (:lambda-0 hyperparams (identity-matrix D))
8           nu     (:nu     hyperparams (inc D))
9           kappa  (:kappa  hyperparams 1.0)
10          pi (sample (dirichlet (repeat K alpha)))
11          lambda (into [] (map (fn [x] (sample x))
12                               (repeat K (wishart nu lambda-0)))))
13          mu (into [] (map
14                      (fn [k]
15                        (sample (mvn mu-0
16                                   (invert kappa
17                                     (get lambda k))))
18                          (range K)))
19          sigma (into [] (map invert lambda)))
20    (loop [n 0
21          z []]
22      (if (= n N)
23        z
24        (let [row (get-row data n)
25              k (sample (discrete pi))]
26          (observe (mvn (get mu k) (get sigma k)) row)
27          (recur (inc n) (conj z k))))))
```

Mixture of Experts-1

MoE(Mixture of Experts) consists of **gating model** and **experts**, where gating model divides the input into subsets, and experts learn on the subset.

Mixture models like GMM are good for generative modelings. By MoE, we can use them on discriminative models, like regression or classification problem.



Figure 3: Mixture of Experts

Mixture of Experts-2

Definition (Inverse Problem)

Given data (x_i, y_i) such that $y_i = f(x_i) + \epsilon$, inverse problem is finding some function g such that $x \cong g(f(x))$.

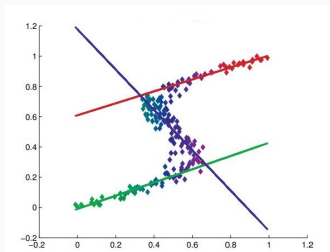


Figure 4: Example of GMM based MoE on inverse problem

MoE in inverse problems:

If f is not invertible, the model needs to choose one of x that $g(x) \cong y$. Choosing one of them is done by the gates of MoE, and approximation of subspace is done by each experts.

Mixture of Experts-3

We may have experts of MoE as MoEs. This creates **hierarchical MoE**, which is a finite tree where every leaf nodes are experts, and internal nodes are gates.

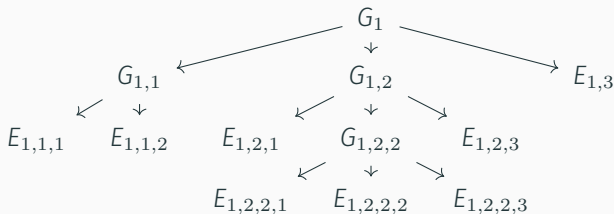


Figure 5: Example of hierarchical mixture of experts model

Experiment

We will use hierarchical mixture model, where gates are given by GMM. Target is inverse problem mentioned above, but specific dataset is not yet chosen.

Our algorithm is composed of 3 parts, GMM Generator, Expanding Decider, Expert Generator.

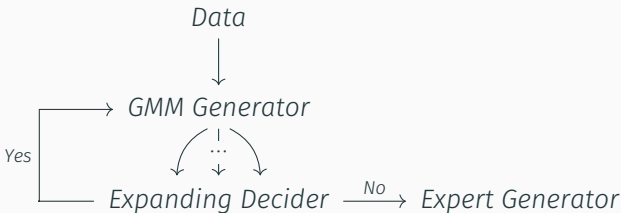


Figure 6: Data flow diagram for our model

Advantage of using ProbProg-1

EM algorithm is hard to implement, hard to understand, and even prone to local minima. Though EM algorithm for GMM is simpler, its derivation includes some matrix calculus, and EM algorithm for mixture of experts is more complex.

Using probabilistic programming will solve this computation problem.

Advantage of using ProbProg-2

Hyperparameter - for example, number of clusters - tuning in GMM is not simple. Usual approach is trying for various choices, then evaluate them, and this is time-consuming, as well as difficult to automate.

We expect that including hyperparameters as distribution, then inferring posterior distribution on themselves will solve this.

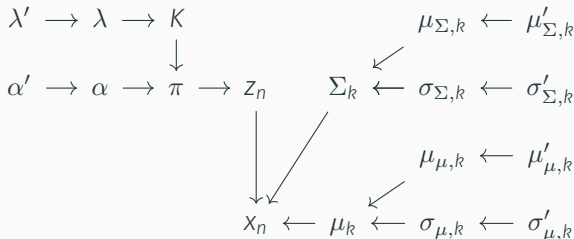


Figure 7: Independence structure of auto tuning GMM

Advantage of using ProbProg-3

Major application in using hierarchical MoE is that splitting condition can be included in inference. This allows regularization on depth of MoE while learning.

Using recursion, hierarchical MoE can be expressed as a finite tree, a structure that classical modeling in statistics can not simply express.

References

Figure 1, 2 : Mathematical Foundation for Artificial Intelligence (2019, Fall), Department of Mathematical Sciences, Prof. Ganguk Hwang.

Slide 5 : Anglican page, examples 'GAUSSIAN MIXTURE MODEL: IRIS DATA'

Figure 4 : Machine Learning: a Probabilistic Perspective by Kevin Patrick Murphy, Ch 11.

Thanks!