

**Imperial College  
London**

COURSEWORK

PROBABILISTIC INFERENCE (CO-493)

---

**Mean Field Variational Inference**

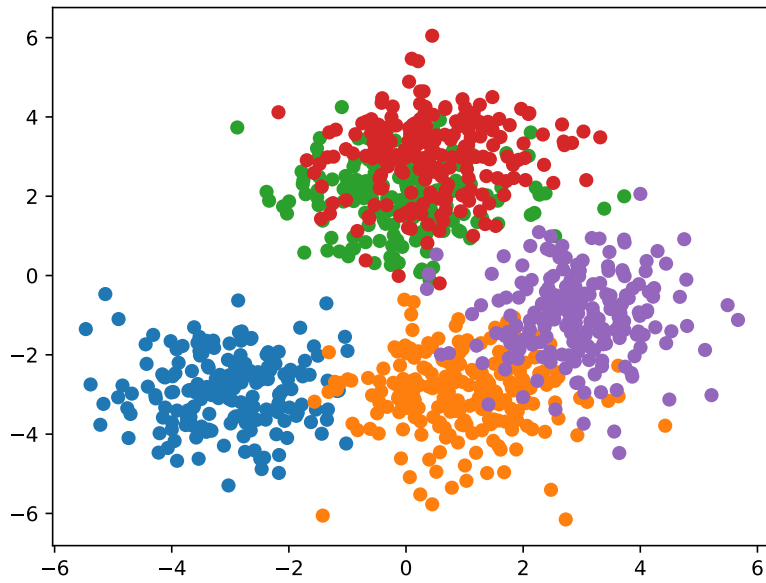
---

---

# Mean Field Variational Inference

## Coursework Description

Variational inference is a powerful tool in modern machine learning that frames the problem of approximating intractable densities as an optimization problem. The goal of this coursework is to build an intuitive understanding of variational inference by implementing the necessary steps for approximate inference in a Gaussian Mixture Model (GMM).



**Figure 1:** Samples from a 2D mixture of Gaussians with five components.

You are given exemplary 2D data from a mixture of Gaussians with five components in Figure 1. The aim is to model the dataset of a  $D$ -dimensional GMM with  $K$  components. The generative process is specified as:

$$\mu_k \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}), k = 1, \dots, K, \quad (1)$$

$$c_i \sim \text{Cat}(1/K), i = 1, \dots, N, \quad (2)$$

$$\mathbf{x}_i | c_i, \mu \sim \prod_k \mathcal{N}(\mu_k, \mathbf{I})^{c_{ik}}, i = 1, \dots, N, k = 1, \dots, K, \quad (3)$$

where  $\mathbf{x}_i \in \mathbb{R}^D$  are the observations,  $c_i \in \{0, 1\}^K$  are one-hot vectors indicating mixture component assignment,  $\mu_k \in \mathbb{R}^D$  are the means of the mixture components,  $\sigma^2 \in \mathbb{R}$  is the prior variance on the means which is assumed to be fixed across dimensions and  $\mathbf{I}$  denotes the identity matrix.

---

Note that we have only specified a prior distribution over the means of the mixture components and the variance of the observations is assumed to be known and fixed. This is a (partially) Bayesian GMM where the latent variables are the means of the mixture components  $\mu_k$  and the cluster assignment variables  $c_i$ .

Given the dataset  $X = \{\mathbf{x}_n\}_{n=1}^N$  the objective is to approximate the posterior distribution  $p(\mu, c|X)$ . To do so, we use a **mean field** variational approximation (all the latent variables factorize):

$$q(\mu, c) = \prod_{k=1}^K q(\mu_k) \prod_{i=1}^N q(c_i) \quad (4)$$

where  $q(\mu_k)$  is a Gaussian distribution on the mean of the  $k$ -th mixture component,  $q(c_i)$  is a categorical distribution on the mixture assignment of the  $i$ -th observation. Formally, we have

$$\mu_k \sim \mathcal{N}(\mathbf{m}_k, \text{diag}(\mathbf{s}_k^2)), \quad k = 1, \dots, K \quad (5)$$

$$c_i \sim \text{Cat}(\pi_i), \quad i = 1, \dots, N \quad (6)$$

where  $\mathbf{m}_k \in \mathbb{R}^D$ ,  $\mathbf{s}_k^2 \in \mathbb{R}^D$  and  $\pi_i \in [0, 1]^K$  are variational parameters that we are going to optimize. Note that  $\text{diag}(\mathbf{s}_k^2)$  is a diagonal  $D \times D$  matrix, i.e. we have assumed that the means of the mixture components factorize across dimensions.

Coordinate ascent variational inference (CAVI) is an iterative algorithm that finds a local optimum of the evidence lower bound (ELBO). CAVI takes advantage of the mean field assumption and updates each factor according to:

$$q^*(z_j) \propto \exp\left(\mathbb{E}_{-j}[\log p(z_j, \mathbf{z}_{-j}, X)]\right) \quad (7)$$

where the subscript  $-j$  denotes all variables except the  $j$ -th one, so that the expectation is with respect to all variational factors except  $j$ .

## File Descriptions / Submission

- `vi_assignment.py`: Skeleton file that will be graded. Put all your code in this file.
- `vi_extra.ipynb`: Jupyter notebook where you can plot the data, run CAVI and visualize the variational distribution after implementing the code.

Submit your final version to CATE via the LabTS system. **Use Python 3.6 for your implementation.**

---

## The Evidence Lower Bound

The ELBO is a lower bound on the log-marginal likelihood  $\log p(X) = \log \int_{\mu, c} p(X, \mu, c)$  which is generally intractable. The ELBO is a key feature of variational inference and is given by:

$$\text{ELBO} = \mathbb{E}_{q(\mu, c)} \left[ \frac{p(X, \mu, c)}{q(\mu, c)} \right] \quad (8)$$

### Task 1: 20 Marks

Complete the function `log_joint()` which returns the expectation of the log-density of the model under the variational distribution:  $\mathbb{E}_{q(\mu, c)} [\log p(X, \mu, c)]$ .

**NOTE:** Make sure the function returns the value after throwing away constant terms, meaning terms that do not depend on any variational parameters.

**NOTE:** An  $n$ -dimensional Gaussian with mean  $\mu$  and diagonal covariance matrix  $\text{diag}(\sigma^2) = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2)$  is the same as a collection of  $n$  independent Gaussian random variables with mean  $\mu_i$  and variance  $\sigma_i^2$ , respectively.

### Task 2: 15 Marks

Complete the function `log_var()` which returns the expectation of the log-variational-density of the model under the variational distribution:  $\mathbb{E}_{q(\mu, c)} [\log q(\mu, c)]$ .

**NOTE:** Make sure the function returns the value after throwing away constant terms, meaning terms that do not depend on any model parameters or variational parameters.

### Task 3: 5 Marks

Complete the `elbo()` function using the `log_joint()` and `log_var()` functions from above.

## Variational Parameter Updates

Equation (7) is the general form for the optimal updates. In order to derive the updates you can plug in the model definition, throw away constant terms and perform the relevant expectations.

### Task 4: 20 Marks

Complete the function `update_pi()` which returns the normalized probabilities of cluster assignments.

---

**NOTE:** After throwing away constant terms you are left with an expression  $\pi_i \propto [\exp(p_1), \dots, \exp(p_K)]$ . Normalize by dividing by the sum across  $K$  for each  $i$ .

**Task 5: 20 Marks**

Complete the function `update_m()` which returns the updated mean parameters for  $q(\mu)$ .

**NOTE:** You might find the Gaussian distribution in its canonical form useful for the derivation:

$$p(x|\eta, \Lambda) = \alpha \exp\{\eta^T x + \frac{1}{2} x^T \Lambda x\} \quad (9)$$

Coverting the canonical parameters  $\eta, \Lambda$  to moment parameters, we have  $\mu = \Lambda^{-1} \eta$  and  $\sigma^2 = \Lambda^{-1}$ .  $\alpha$  is a normalizing constant.

**NOTE:** Remember to throw away constant terms.

**Task 6: 20 Marks**

Complete the function `update_s2()` which returns the updated variance parameters for  $q(\mu)$ .

**NOTE:** Remember to throw away constant terms.