

Variational Inference

Caterina De Bacco and Isabel Valera

1 Topic modeling: the idea

Consider a collection of text corpora, for instance a set of documents.

Goal: find statistical pattern behind the data, so that we can parametrize the members of the collection by a short description. In other words, we want to find a low dimensional representation of the data.

Idea: documents are mixture of K latent variables called *topics*, and *topics* are mixtures of words. Formally, we have 3 types of latent variables:

- β_k : a distribution of words, needed to specify the topic;
- θ_d : a vector of topic proportion, needed to specify a document;
- z_{dn} : a topic assignment, needed to specify what words are seen in each document.

The data are words w , divided into documents.

The goal is to then estimate the posterior $p(\beta, \theta, z|w)$. This can then be used to perform various tasks, like classification, novelty detection, similarity or relevance judgement.

"Arts"	"Budgets"	"Children"	"Education"
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services," Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

Figure 8: An example article from the AP corpus. Each color codes a different factor from which the word is putatively generated.

Figure 1: Example of Topic Modeling. Colors represent topics, and on top we have the words' mixture denoting each topic. Figure taken from [Blei et al. \(2003\)](#).

We will see in the next Lecture how to formalize this problem and solve it. In Figure 1 you see an example of the results that we will obtain.

2 Variational Inference: the idea

We spent few Lectures discussing methods for approximating complex joint probability distributions and used the approach of statistical physics. Today we take a more pure statistics approach and focus on a particular method to address the same problem: Variational Inference (VI). VI is an inference approach that approximates probability distributions through optimization.

In the field of Bayesian statistics, one is interested in deriving the posterior distribution of model's parameters. Often though, the posterior is not analytically tractable, although the likelihood might be (recall Lecture 1).

Consider a joint density of latent variables $\mathbf{z} = (z_1, \dots, z_m)$ and data $\mathbf{x} = (x_1, \dots, x_n)$:

$$p(\mathbf{z}, \mathbf{x}) = p(\mathbf{x}|\mathbf{z})p(\mathbf{z}) \quad . \quad (1)$$

Using a model with latent variables is a common approach in many inference models, as one can use these variables to govern the data distribution. Inference in Bayesian modeling consists in posing a prior for these variables $p(\mathbf{z})$ and extracting the posterior $p(\mathbf{z}|\mathbf{x})$, given a likelihood $p(\mathbf{x}|\mathbf{z})$.

The idea behind Variational Inference, is to posit a family of *tractable* distributions \mathcal{D} over the latent variables and find one element of this set that is closest to the untractable posterior. Closeness is measured by Kullback-Leibler (KL) divergence:

$$q^*(\mathbf{z}) = \arg \min_{q_\theta(\mathbf{z}) \in \mathcal{D}} KL(q_\theta(\mathbf{z}) || p(\mathbf{z}|\mathbf{x})) \quad . \quad (2)$$

The distribution $q_\theta(\mathbf{z})$ is called *variational distribution*. Figure 2 gives a sketch of this idea.

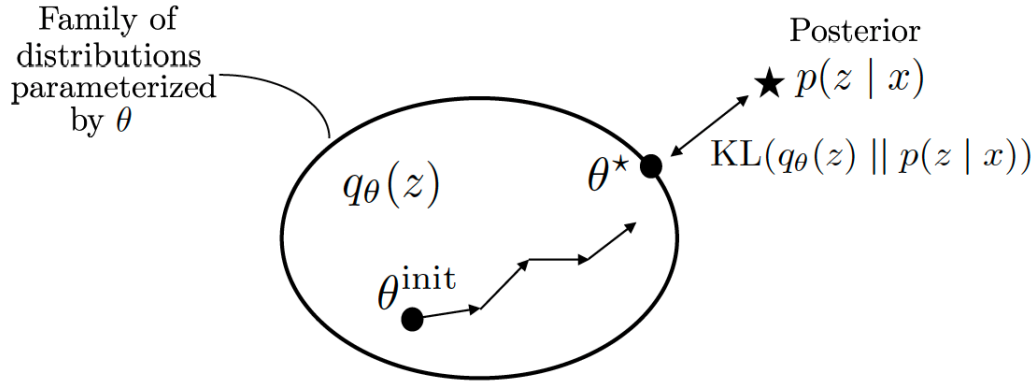


Figure 2: The idea behind Variational Inference. Figure courtesy of Francisco J. R. Ruiz.

Obs1: other traditional models for inference that go under the Monte Carlo family are based on *sampling*, rather than optimization.

Obs2: other types of cost functions to be optimized could also be considered. Here we focus only on $KL(q||p)$.

3 The problem

Objective: estimate the posterior $p(\mathbf{z}|\mathbf{x})$ given the data.

We have that:

$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{x})} \quad , \quad (3)$$

where $p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z}$ is called the *evidence*. This integral is often unavailable in closed form, but this is needed in order to calculate the posterior. This is why inference in these cases is hard.

3.1 Example: Bayesian Mixture of Gaussians.

We already saw this example in the previous Lectures, but will recap here the problem. Consider a mixture of unit-variance univariate Gaussians. There are K mixture components, one for each Gaussian, with mean $\boldsymbol{\mu} = (\mu_1, \dots, \mu_K)$. In Figure 3 we plot an example of such dataset.

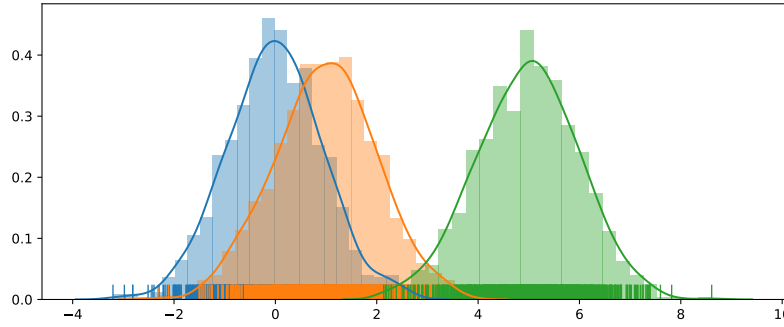


Figure 3: Example of a GMM dataset with 3 clusters and their Gaussians components. Here the means of the Gaussians are $\mu_1 = 0, \mu_2 = 1, \mu_3 = 5$.

The means in turn are also random variables, drawn from a prior $p(\mu_k)$ which we assume to be Gaussian $\mathcal{N}(0, \sigma^2)$; σ^2 is an hyper-parameter. The model works as follow: we assume that each data point x_i is extracted from one of the Gaussians, but we do not know which one. We thus need to introduce a cluster assignment latent variable c_i that tells which Gaussian x_i was drawn from. The $c_i = \{1, \dots, K\}$ is assumed to be categorical, but is encoded in an indicator K -vector, with all zeros except one equal to 1 entry corresponding to the cluster. Formally, the model is:

$$\mu_k \sim \mathcal{N}(0, \sigma^2) \quad k = 1, \dots, K \quad (4)$$

$$c_i \sim \text{Categorical}(1/K, \dots, 1/K) \quad i = 1, \dots, n \quad (5)$$

$$x_i | c_i, \boldsymbol{\mu} \sim \mathcal{N}(c_i^T \boldsymbol{\mu}, 1) \quad i = 1, \dots, n \quad (6)$$

Obs1: notice that each variable is drawn from *one single* Gaussian, and *not* from a mixture of Gaussians! The *whole* set of variables is a mixture of Gaussian.

The joint density of data and parameters is then:

$$p(\boldsymbol{\mu}, \mathbf{c}, \mathbf{x}) = p(\boldsymbol{\mu}) \prod_{i=1}^n p(c_i) p(x_i | c_i, \boldsymbol{\mu}) \quad . \quad (7)$$

The latent variables in this case are $\mathbf{z} = (\boldsymbol{\mu}, \mathbf{c})$. The evidence is then:

$$p(\mathbf{x}) = \int p(\boldsymbol{\mu}) \prod_{i=1}^n \sum_{c_i} p(c_i) p(x_i | c_i, \boldsymbol{\mu}) d\boldsymbol{\mu} \quad . \quad (8)$$

Question: can you calculate that integral?

Unfortunately no. This is because the integrand does contain a separate factor for each μ_k , thus preventing the reduction onto one-dimensional integrals. The sum over the cluster assignments runs over K^n configurations, i.e. exponential in K , which is only doable for very small values of n .

4 The Evidence Lower Bound (ELBO).

Recall that the goal of VI is to optimize (we omit the explicit dependence on θ):

$$q^*(\mathbf{z}) = \arg \min_{q(\mathbf{z}) \in \mathcal{Q}} KL(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x})) \quad . \quad (9)$$

However, this objective is not computable because we have to compute $\log p(\mathbf{x})$, which is not feasible as we saw before. To see this, let's unpack the KL divergence:

$$KL(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x})) = \mathbb{E}_q[\log q(\mathbf{z})] - \mathbb{E}_q[\log p(\mathbf{z}|\mathbf{x})] \quad (10)$$

$$= \mathbb{E}_q[\log q(\mathbf{z})] - \mathbb{E}_q[\log p(\mathbf{z}, \mathbf{x})] + \log p(\mathbf{x}) \quad , \quad (11)$$

which requires calculating $\log p(\mathbf{x})$ as said before.

To skip computing the exact KL, we propose an *alternative* optimization objective.

Question:How do we choose this?

We pick an objective equivalent to the same KL as before up to a constant (in q):

$$\text{ELBO}(q) := \mathbb{E}_q[\log p(\mathbf{z}, \mathbf{x})] - \mathbb{E}_q[\log q(\mathbf{z})] \quad (12)$$

$$= \log p(\mathbf{x}) - KL(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x})) \quad . \quad (13)$$

Obs1: maximizing $\text{ELBO}(q)$ is equivalent to minimizing $KL(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x}))$.

4.1 Properties of the ELBO

If we rewrite the ELBO unpacking the first term yields:

$$\text{ELBO}(q) = \mathbb{E}_q[\log p(\mathbf{x}|\mathbf{z})] + \mathbb{E}_q[\log p(\mathbf{z})] - \mathbb{E}_q[\log q(\mathbf{z})] \quad (14)$$

$$= \mathbb{E}_q[\log p(\mathbf{x}|\mathbf{z})] - KL(q(\mathbf{z})||p(\mathbf{z})) \quad . \quad (15)$$

In other words, the ELBO is the sum of the expected log likelihood of the data and (minus) the KL divergence between the variational distribution $q(\mathbf{z})$ and the prior $p(\mathbf{z})$.

Property 1: maximizing the ELBO is then encouraging a combination of increasing the expected log likelihood, i.e. finding a $q(\mathbf{z})$ that explains the data, and choosing $q(\mathbf{z})$ close to the prior.

Because $KL(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x})) \geq 0$, we have that:

$$\log p(\mathbf{x}) = \text{ELBO}(q) + KL(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x})) \quad (16)$$

$$\geq \text{ELBO}(q) \quad . \quad (17)$$

Property 2: the ELBO is a lower bound to the evidence (this is the reason of its name).

4.2 The Mean-Field variational family

One relevant aspect of Variational Inference, is that we can pick a variational family at our choice. Of course, this does not guarantee that we will make a good choice. We expect that the more complex the q the better the approximation, at the cost of losing the tractability.

If we want to go for making things analytically tractable, which includes for instance being able to compute $\mathbb{E}_q[\cdot]$, then the best choice is to consider a fully factorized family as we already saw in Lecture 7. This is the Mean-Field variational family:

$$q(\mathbf{z}) = \prod_{j=1}^m q_j(z_j) \quad . \quad (18)$$

Obs1: each latent variable z_j is governed by its own variational factor and these are all *independent*.

Obs2: notice that the variational distribution $q(\mathbf{z})$ does *not* depend on the data \mathbf{x} . The dependence on the data is taken care of inside the ELBO, through the expected log likelihood term.

Obs3: the subindex j of $q_j(\cdot)$ is there so that one can specify a different type of distribution for each latent variable. For instance, for binary latent variables $q_j(\cdot)$ can be Bernoulli, for continuous numbers $q_j(\cdot)$ can be Gaussian, and so on.

5 Coordinate Ascent Mean-Field Variational Inference (CAVI)

So far we have seen the theory being VI and a class of variational family, the Mean-Field distribution. We haven't said anything yet about how to solve the optimization problem. Here we show how to tackle this task by describing the CAVI algorithm.

Idea: iteratively optimize each single variational factor holding the others fixed, until we reach a local maximum of the ELBO.

We will use the following result. Consider z_j . The *complete conditional* $p(z_j | \mathbf{z}_{\setminus j}, \mathbf{x})$ of z_j is a conditional density given all the other latent variables and data.

Fact: the *optimal* $q_j(z_j)$ is proportional to the exponentiated expected log of the complete conditional of z_j :

$$q_j^*(z_j) \propto \exp \left\{ \mathbb{E}_{\setminus j} [\log p(z_j | \mathbf{z}_{\setminus j}, \mathbf{x})] \right\} \quad , \quad (19)$$

where $\mathbb{E}_{\setminus j}[\cdot]$ is over $\prod_{l \neq j} q_l(z_l)$ (recall that $q_l(z_l)$ are currently being held fixed).

An equivalent result is:

$$q_j^*(z_j) \propto \exp \left\{ \mathbb{E}_{\setminus j} [\log p(z_j, \mathbf{z}_{\setminus j}, \mathbf{x})] \right\} \quad , \quad (20)$$

where we considered instead the log of the joint distribution.

Proof. Let's rewrite the ELBO by isolating a variational factor $q_j(z_j)$ absorbing into a constant terms that do not depend on it, and use the Mean-Field approximation:

$$\text{ELBO}(q_j) = \mathbb{E}_j [\mathbb{E}_{\setminus j} [\log p(z_j, \mathbf{z}_{\setminus j}, \mathbf{x})]] - \mathbb{E}_j [\log q_j(z_j)] + \text{const} \quad (21)$$

$$= -KL(q_j || q_j^*) + \text{const} \quad . \quad (22)$$

Thus we maximize the ELBO when we minimize the $KL(q_j || q_j^*)$; this is minimized when $q_j(z_j) \equiv q_j^*(z_j)$.

5.1 Example: Bayesian Mixture of Gaussians (continued).

Let's go back to the GMM and apply what just learned. The MF family for the latent parameters μ_k and c_i is:

$$q(\boldsymbol{\mu}, \mathbf{c}) = \prod_{k=1}^K q(\mu_k; m_k, s_k^2) \prod_{i=1}^n q(c_i; \rho_i) \quad (23)$$

where we choose:

$$q(\mu_k; m_k, s_k^2) = \mathcal{N}(m_k, s_k^2) \quad (24)$$

$$q(c_i; \rho_i) = \text{Categorical}(\rho_i) \quad , \quad (25)$$

and ρ_i is a K -dimensional vector.

Obs1: we chose Gaussian and Categorical distributions. In principle we could have chosen other types. However, the choice impacts the goodness of the approximation. For this case, Gaussian and Categorical are indeed the optimal choice for the Mean-Field family.

Variational update of the cluster assignment c_i . Using equation (20) yields:

$$\begin{aligned} q^*(c_i; \rho_i) &\propto \exp \left\{ \mathbb{E}_{\setminus c_i} \left[\log \left(p(c_i) \prod_{l \neq i} p(c_l) p(\boldsymbol{\mu}) p(\mathbf{x}|c_i, \mathbf{c}_{\setminus i}, \boldsymbol{\mu}) \right) \right] \right\} \\ &= \exp \left\{ \log p(c_i) + \mathbb{E}_{\setminus c_i} \left[\sum_{l \neq i} \log p(c_l) + \log p(\boldsymbol{\mu}) \right] + \mathbb{E}_{\setminus c_i} [\log p(\mathbf{x}|c_i, \mathbf{c}_{\setminus i}, \boldsymbol{\mu})] \right\} \\ &= \exp \left\{ \log p(c_i) + \sum_{l \neq i} \mathbb{E}_{q(c_l)} [\log p(c_l)] + \sum_k \mathbb{E}_{q(\mu_k)} [\log p(\mu_k)] + \mathbb{E}_{\setminus c_i} \left[\sum_i \log p(x_i|c_i, \boldsymbol{\mu}) \right] \right\} \end{aligned} \quad (26)$$

The second and third terms are constant in c_i , so they can be neglected. The likelihood term can be further unpacked by keeping only terms containing c_i . Recall that using c_i as an indicator we have:

$$p(x_i|c_i, \boldsymbol{\mu}) = \prod_{k=1}^K p(x_i|\mu_k)^{c_{ik}} \quad . \quad (28)$$

Substituting into the previous formula yields:

$$\mathbb{E}_{\setminus c_i} [\log p(x_i|c_i, \boldsymbol{\mu})] = \mathbb{E}_{q(\boldsymbol{\mu})} [\log p(x_i|c_i, \boldsymbol{\mu})] = \sum_k \mathbb{E}_{q(\mu_k)} [c_{ik} \log p(x_i|\mu_k)] \quad (29)$$

$$= \sum_k c_{ik} \mathbb{E}_{q(\mu_k)} [\log p(x_i|\mu_k)] \quad (30)$$

$$= \sum_k c_{ik} \mathbb{E}_{q(\mu_k)} [-(x_i - \mu_k)^2 / 2] + \text{const} \quad (31)$$

$$= -\sum_k c_{ik} x_i^2 / 2 - \sum_k c_{ik} \mathbb{E}_{q(\mu_k)} [\mu_k^2 / 2] + x_i \sum_k c_{ik} \mathbb{E}_{q(\mu_k)} [\mu_k] + \text{const} \quad (32)$$

$$= x_i \sum_k c_{ik} \mathbb{E}_{q(\mu_k)} [\mu_k] - \frac{1}{2} \sum_k c_{ik} \mathbb{E}_{q(\mu_k)} [\mu_k^2] + \text{const} \quad (33)$$

$$= x_i \sum_k c_{ik} m_k - \frac{1}{2} \sum_k c_{ik} (s_k^2 + m_k^2) + \text{const} \quad . \quad (34)$$

We can finally substitute into Eq. (27):

$$q^*(c_i; \rho_i) = \prod_k \rho_{ik}^{*c_{ik}} \propto \exp \left\{ \log p(c_i) + x_i \sum_k c_{ik} m_k - \frac{1}{2} \sum_k c_{ik} (s_k^2 + m_k^2) \right\} \quad (35)$$

$$= p(c_i) \exp \left\{ x_i \sum_k c_{ik} m_k - \frac{1}{2} \sum_k c_{ik} (s_k^2 + m_k^2) \right\} \quad (36)$$

$$\propto \prod_k \exp \left\{ \left[x_i m_k - \frac{1}{2} (s_k^2 + m_k^2) \right] c_{ik} \right\} \quad . \quad (37)$$

Which means that the optimal parameter for the categorical variational distribution is:

$$\rho_{ik}^* \propto \exp \left[x_i m_k - \frac{1}{2} (s_k^2 + m_k^2) \right] . \quad (38)$$

Obs1: these are a function of both data and the variational parameters of the mixture component, but *not* of the other ρ_j^* .

Variational update of the mixture components' means μ_k . We can repeat similar calculations for the variational distributions $q(\mu_k; m_k, s_k^2)$. Again, use equation (20) to get:

$$q(\mu_k) \propto \exp \left\{ \log p(\mu_k) + \sum_i \mathbb{E}_{q(\mu_k)} [\log p(x_i | c_i, \mu)] \right\} . \quad (39)$$

Using again Eq. (28) and keeping only terms containing μ_k yields:

$$\log q(\mu_k) = \log p(\mu_k) + \sum_i \mathbb{E}_{q(c_i)} [c_{ik}] \log p(x_i | \mu_k) + \text{const} \quad (40)$$

$$= -\frac{\mu_k^2}{2\sigma^2} - \sum_i \rho_{ik} \frac{(x_i - \mu_k)^2}{2} + \text{const} \quad (41)$$

$$= \mu_k \sum_i \rho_{ik} x_i - \frac{\mu_k^2}{2} \left(\frac{1}{\sigma^2} + \sum_i \rho_{ik} \right) + \text{const} . \quad (42)$$

For those of you familiar with exponential families, this means that $q(\mu_k)$ is a member of it with sufficient statistics $\{\mu_k, \mu_k^2\}$ and natural parameters $\{\sum_i \rho_{ik} x_i, -\frac{1}{2} (\frac{1}{\sigma^2} + \sum_i \rho_{ik})\}$. This means that it is a Gaussian distribution. For a Gaussian of mean μ and variance σ^2 , the natural parameters are generally $\eta \equiv (\frac{\mu}{\sigma^2}, -\frac{1}{2\sigma^2})$. Doing the mapping, we obtain:

$$m_k = \frac{\sum_i \rho_{ik} x_i}{\frac{1}{\sigma^2} + \sum_i \rho_{ik}} \quad (43)$$

$$s_k^2 = \frac{1}{\frac{1}{\sigma^2} + \sum_i \rho_{ik}} . \quad (44)$$

Obs1: to derive these results, we never used the assumption that $q(\mu_k)$ is a Gaussian. This was coming as a result looking at the shape in terms of exponential families. Thus, for this case, the Gaussian variational distribution is actually the optimal one for the mixing components (as anticipated before).

Obs2: you can obtain the same results by taking a different approach. Which one?

Obs3: in general, the complete conditional can be a complicated distribution. However, in many cases this is a member of the exponential family. If this is the case, then Eq. (19) simplifies.

The algorithm then works as in Algorithm 1.

By applying this algorithm for the example of Figure 3 you can see how it performs in Figure 4.

Variational Inference: summary

- VI is an inference technique to approximate a complicated distribution with a tractable one.
- It relies on optimizing the ELBO, this has the same optimum as that of the KL divergence between the exact and the variational distribution.
- The Mean-Field family is a popular choice for the variational distribution.
- It allows for efficient Coordinate Ascent updates (CAVI).

A main reference for this lecture is [Blei et al. \(2017\)](#).

Algorithm 1: CAVI for a Gaussian mixture model

Input: Data \mathbf{x} , number of components K , prior variance of component means σ^2
Output: Variational densities $q(\mu_k; m_k, s_k^2)$ (Gaussian) and $q(c_i; \rho_i)$ (K -categorical)
Initialize: Variational parameters $\mathbf{m} = m_{1:K}$, $\mathbf{s}^2 = s_{1:K}^2$, and $\boldsymbol{\varphi} = \rho_{1:n}$
while the ELBO has not converged **do**
 for $i \in \{1, \dots, n\}$ **do**
 Set $\rho_{ik} \propto \exp\left[x_i m_k - \frac{1}{2}(s_k^2 - m_k^2)\right]$
 end
 for $k \in \{1, \dots, K\}$ **do**
 Set $m_k \leftarrow \frac{\sum_i \rho_{ik} x_i}{\frac{1}{\sigma^2} + \sum_i \rho_{ik}}$
 Set $s_k^2 \leftarrow \frac{1}{\frac{1}{\sigma^2} + \sum_i \rho_{ik}}$
 end
 Compute ELBO($\mathbf{m}, \mathbf{s}^2, \boldsymbol{\varphi}$)
end
return $q(\mathbf{m}, \mathbf{s}^2, \boldsymbol{\varphi})$

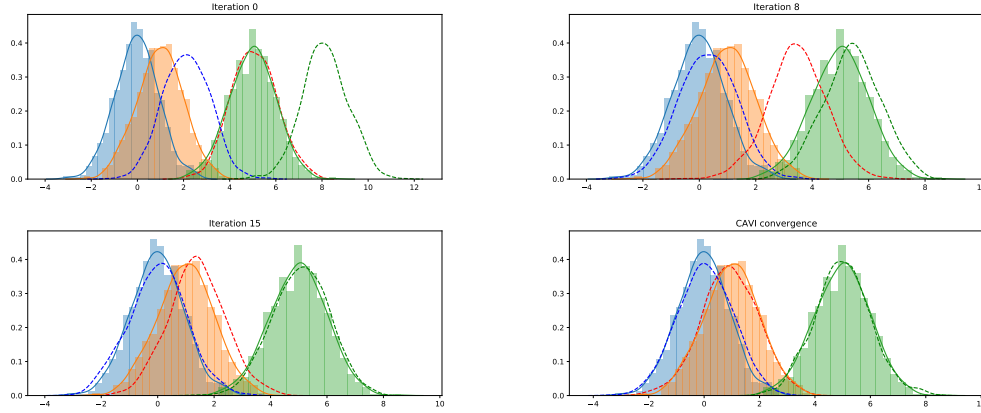


Figure 4: Result of CAVI on GMM of Figure 3. The plots are at 4 different iteration steps: at the beginning, at two changing points of the ELBO and at convergence. Dashed lines are corresponding to the variational estimates, regular lines are the exact Gaussians (those used to generate the data).

References

- D. M. Blei, A. Y. Ng, and M. I. Jordan, Journal of machine Learning research **3**, 993 (2003).
D. M. Blei, A. Kucukelbir, and J. D. McAuliffe, Journal of the American Statistical Association **112**, 859 (2017).