# Gaussian Mixture Model + Variational Inference
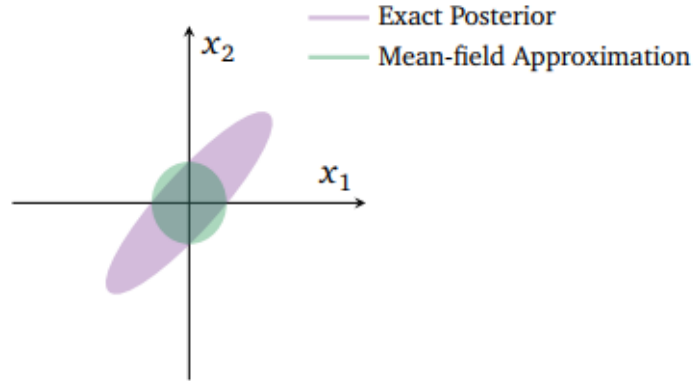
## 2.3 Visualizing the mean-field approximation.



## 2.4 Coordinate ascent mean-field variational inference algorithm

$$q_j^*(z_j) \propto \exp\left\{\mathbb{E}_{-j}\left[\log p(z_j \mid \mathbf{z}_{-j}, \mathbf{x})\right]\right\}. \tag{17}$$

$$q_j^*(z_j) \propto \exp\left\{\mathbb{E}_{-j}\left[\log p(z_j, \mathbf{z}_{-j}, \mathbf{x})\right]\right\}. \tag{18}$$

---

**Algorithm 1:** Coordinate ascent variational inference (CAVI)

**Input:** A model $p(\mathbf{x}, \mathbf{z})$, a data set $\mathbf{x}$

**Output:** A variational density $q(\mathbf{z}) = \prod_{j=1}^{m} q_j(z_j)$

**Initialize:** Variational factors $q_j(z_j)$

**while** *the ELBO has not converged* **do**

    **for** $j \in \{1, \dots, m\}$ **do**

        | Set $q_j(z_j) \propto \exp\{\mathbb{E}_{-j}[\log p(z_j \mid \mathbf{z}_{-j}, \mathbf{x})]\}$

    **end**

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

**end**

**return** $q(\mathbf{z})$

---

$$\text{ELBO}(q_j) = \mathbb{E}_j\left[\mathbb{E}_{-j}\left[\log p(z_j, \mathbf{z}_{-j}, \mathbf{x})\right]\right] - \mathbb{E}_j\left[\log q_j(z_j)\right] + \text{const}. \tag{19}$$

## 3. A complete example: Bayesian mixture Gaussian

K: quantity of clustering

n: quantity of data

$\mathbf{c} = c_{1:n}$, where $c_i$ is an indicator K-vector , $c_4 = [0,0,0,1]$

M: quantity of dimensions, and it's a scalar

$\mu = \mu_{1:K}$, real-valued mean parameters, $\mu_k = (m_k, s_k^2) = (m_k, V_k)$, k=1,2,3,...,K

where $m_k$ has a shape of $1 \times M$, and $V_k$ has a shape of $M \times M$

$X = (x_i)$, X is the dataset with a shape of $n \times M$, $i = 1,2,3,...,n$

$x_i$ is the data

### 3.1 Prior distributions:

$$p(\mu) = \text{Normal}(\mu|0, \sigma^2) = \frac{1}{\sqrt{2\pi\,\sigma^2}} \exp\left(-\frac{\mu^2}{2\sigma^2}\right)$$

$$p(c) = \frac{1}{K}, \quad \text{in details:} \quad P(c_{ik}) = \frac{1}{K}, \quad for\ k = 1,2,3,...,K$$

### 3.2 Observations:

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

$$p(X|c,\mu) = \prod_{i=1}^{n}\prod_{k=1}^{K} p(x_i|\mu_k)^{c_{ik}}$$

### 3.3 Joint distribution:

$$p(X,c,\mu) = p(\mu)\prod_{i=1}^{n}\prod_{k=1}^{K} p(x_i|\mu_k)^{c_{ik}}$$

In a more clear writing,

$$p(X,c,\mu) = p(\mu)\prod_{i=1}^{n} p(c_i)p(x_i|c_i,\mu)$$

### 3.4 From Bayesian formula we know the posterior distribution:

$$p(c,\mu|X) = \frac{p(X,c,\mu)}{p(X)}$$

## 4. variance inference

4.1 $\quad q(c; \varphi_{ik})q(\mu; m_k, S_k) \Rightarrow p(c,\mu|X)$

Now it's time to update the posteriors

### 4.2 The variational density of the mixture assignments

Mean-field VI for $q(c; \varphi_{ik})$

$$q^*(c_i; \varphi_i) \propto \exp\left\{\log p(c_i) + \mathbb{E}\left[\log p(x_i|c_i,\mu); m, s^2\right]\right\}. \tag{22}$$

### 4.2.1

$\exp(\log p(c_i)) = \exp\left(\log\frac{1}{K}\right) = \frac{1}{K}$ is a constant and independent of $c_i$, $\mu$, $s^2$, and thus can be ignored tempor

### 4.2.2 We use this to compute the expected log probability,

$$\mathbb{E}[\log p(x_i \mid c_i, \mu)] = \sum_k c_{ik} \mathbb{E}\left[\log p(x_i \mid \mu_k); m_k, s_k^2\right] \tag{23}$$

$$= \sum_k c_{ik} \mathbb{E}\left[-(x_i - \mu_k)^2/2; m_k, s_k^2\right] + \text{const.} \tag{24}$$

$$= \sum_k c_{ik} \left(\mathbb{E}\left[\mu_k; m_k, s_k^2\right] x_i - \mathbb{E}\left[\mu_k^2; m_k, s_k^2\right]/2\right) + \text{const.} \tag{25}$$

Thus the variational update for the $i^{th}$ cluster assignment is:

$$\varphi_{ik} \propto \exp\left\{\mathbb{E}\left[\mu_k; m_k, s_k^2\right] x_i - \mathbb{E}\left[\mu_k^2; m_k, s_k^2\right]/2\right\}$$

From statistics class:

> if $w \sim \text{normal}(m_k, S_k)$,
>
> then: $E[w] = m_k$
>
> $E[ww^T] = m_k m_k^T + S_k$
>
> $E[w^T w] = m_k^T m_k + Tr(S_k)$
>
> $\mathbb{E}[\mu_k; m_k, s_k^2] = m_k$
>
> $\mathbb{E}[\mu_k^2; m_k, s_k^2] = m_k^2 + s_k^2$

$$\varphi_{ik} \propto exp\left(m_k x_i - \frac{m_k^2 + s_k^2}{2}\right)$$

1D: $\varphi_{ik} \propto \exp\left(m_k X_i - \frac{m_k^2 + Tr(S_k)}{2}\right)$

2D: $\varphi_{ik} \propto \exp\left(m_k^T X_i - \frac{m_k^T m_k + Tr(S_k)}{2}\right)$

### *4.2.3 Trick: due to mean-field assumption, S_k is a diagram with identical non-zero elements.*

$$Tr(S_k) = M \times S_k[1,1]$$

## 4.3  The variational density of the mixture-component means

### 4.3.1 Mean-field VI for $q(\mu_k)$

$$q(\mu_k) \propto \exp\left\{\log p(\mu_k) + \sum_{i=1}^n \mathbb{E}\left[\log p(x_i \mid c_i, \mu); \varphi_i, m_{-k}, s_{-k}^2\right]\right\}. \tag{27}$$

$$\log q(\mu_k) = \log p(\mu_k) + \sum_i \mathbb{E}\left[\log p(x_i \mid c_i, \boldsymbol{\mu}); \varphi_i, \mathbf{m}_{-k}, \mathbf{s}^2_{-k}\right] + \text{const.} \qquad (28)$$

$$= \log p(\mu_k) + \sum_i \mathbb{E}\left[c_{ik} \log p(x_i \mid \mu_k); \varphi_i\right] + \text{const.} \qquad (29)$$

$$= -\mu_k^2 / 2\sigma^2 + \sum_i \mathbb{E}\left[c_{ik}; \varphi_i\right] \log p(x_i \mid \mu_k) + \text{const.} \qquad (30)$$

$$= -\mu_k^2 / 2\sigma^2 + \sum_i \varphi_{ik}\left(-(x_i - \mu_k)^2/2\right) + \text{const.} \qquad (31)$$

$$= -\mu_k^2 / 2\sigma^2 + \sum_i \varphi_{ik} x_i \mu_k - \varphi_{ik}\mu_k^2/2 + \text{const.} \qquad (32)$$

$$= \left(\sum_i \varphi_{ik} x_i\right)\mu_k - \left(1/2\sigma^2 + \sum_i \varphi_{ik}/2\right)\mu_k^2 + \text{const.} \qquad (33)$$

where $p(\mu) = \text{Normal}(\mu \mid 0, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{\mu^2}{2\sigma^2}\right)$ is the prior distribution.

### 4.3.2 For 1D GMM:

这里做一个小小的修正, 我习惯把括号写完：

$$\log q(\mu_k) = -\frac{\mu_k^2}{2\sigma^2} + \sum_i^n \left(\varphi_{ik} x_i \mu_k - \frac{\varphi_{ik}\mu_k^2}{2}\right) + const$$

$$\Rightarrow q(\mu_k; m_k, s_k^2) = normal(\mu_k \mid m_k, s_k^2)$$

$$m_k = \frac{\sum_i \varphi_{ik} x_i}{1/\sigma^2 + \sum_i \varphi_{ik}}, \qquad s_k^2 = \frac{1}{1/\sigma^2 + \sum_i \varphi_{ik}}. \qquad (34)$$

### 4.3.2 For Multi-Dimension GMM:

$$\log q(\mu_k) = -\frac{\mu_k V^{-1} \mu_k^T}{2} + \sum_i^n \left(\varphi_{ik} x_i \mu_k^T - \frac{\varphi_{ik}\mu_k V_0^{-1} \mu_k^T}{2}\right) + const$$

$$\Rightarrow q(\mu_k; m_k, V_k) = normal(\mu_k \mid m_k, V_k)$$

$$m_k = \left\{\sum_i^n \varphi_{ik} x_i\right\}\left(V^{-1} + V_0^{-1}\sum_i^n \varphi_{ik}\right)^{-1}$$

$$V_k = \left(V^{-1} + V_0^{-1}\sum_i^n \varphi_{ik}\right)^{-1}$$

where $V_0 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$, $V = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$

## 4.4 Stochastic Variational Inference for GMM
### 4.4.1 Algorithm for Stochastic Variational Inference

## Stochastic Variational Inference

**Input:** data $\mathbf{x}$, model $p(\beta, \mathbf{z}, \mathbf{x})$.

Initialize $\lambda$ randomly.   Set $\rho_t$ appropriately.

**repeat**

Sample $j \sim \text{Unif}(1, \ldots, n)$.

Set local parameter $\phi \leftarrow \mathbb{E}_\lambda\big[\eta_\ell(\beta, x_j)\big]$.

Set intermediate global parameter

$$\hat{\lambda} = \alpha + n\mathbb{E}_\phi[t(Z_j, x_j)].$$

Set global parameter

$$\lambda = (1 - \rho_t)\lambda + \rho_t\hat{\lambda}.$$

**until** *forever*

$q(\boldsymbol{c}; \varphi_{ik})$ 不变,   $q(\boldsymbol{\mu_k})$ 采用随机优化 (参考随机梯度下降)

Deriving from Mean-field Variational Inference:

$$q(\mu_k) \propto \exp\left\{ \log p(\mu_k) + \sum_{i=1}^{n} \mathbb{E}\big[\log p(x_i \mid c_i, \boldsymbol{\mu}); \varphi_i, \mathbf{m}_{-k}, \mathbf{s}_{-k}^2\big] \right\}. \tag{27}$$

In Stochastic Variational Inference we may update  for global variables $\lambda$

Randomly sample $j \sim Unif(1,2,3, \ldots, n)$.

Note: for each update it need a new $j$,  where $j$ = 1,2,3,...,n

$$\hat{\lambda} = \log p(\mu_k) + n \times E\big[\log p(\boldsymbol{x_j} \mid c_j, \boldsymbol{\mu}); \varphi_j, \boldsymbol{m}_{-k}, s_{-k}^2\big]$$
$$\lambda = (1 - \rho_t)\lambda + \rho_t \hat{\lambda}$$

$\Rightarrow$

### 4.4.2 For 1D GMM:

Randomly sample $j \sim Unif(1,2,3, \ldots, n)$.

$$\log q(\mu_k) = -\frac{\mu_k^2}{2\sigma^2} + n \times (\varphi_{jk}x_j\mu_k - \frac{\varphi_{jk}\mu_k^2}{2}) + const$$
$$\Rightarrow q(\mu_k; m_k, s_k^2) = normal(\mu_k \mid m_k, s_k^2)$$

$$m_k = \frac{n \times \varphi_{jk}x_j}{\frac{1}{\sigma^2} + n \times \varphi_{jk}}, \qquad s_k^2 = \frac{1}{\frac{1}{\sigma^2} + n \times \varphi_{jk}}$$

### 4.4.3 For Multi-Dimension GMM:

Randomly sample $j \sim Unif(1,2,3, \ldots, n)$.

$$\log q(\mu_k) = -\frac{\mu_k V^{-1}\mu_k^T}{2} + n \times (\varphi_{jk}x_j\mu_k^T - \frac{\varphi_{jk}\mu_k V_0^{-1}\mu_k^T}{2}) + const$$

$$\Rightarrow q(\mu_k; m_k, V_k) = normal(\mu_k | m_k, V_k)$$

$$m_k = (n \times \varphi_{jk} x_j)\left(V^{-1} + V_0^{-1} \, n \times \varphi_{jk}\right)^{-1}$$

$$V_k = \left(V^{-1} + V_0^{-1} \, n \times \varphi_{jk}\right)^{-1}$$

where n is the quantity of all data,

$V_0 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$ is from the prior distribution for $p(\mu) = normal(\mu | 0, V_0)$

$V = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$ is the observation for $p(x_i | c_i, \mu)$

## 5.Comparison of Stochastic VI and Mean-field VI
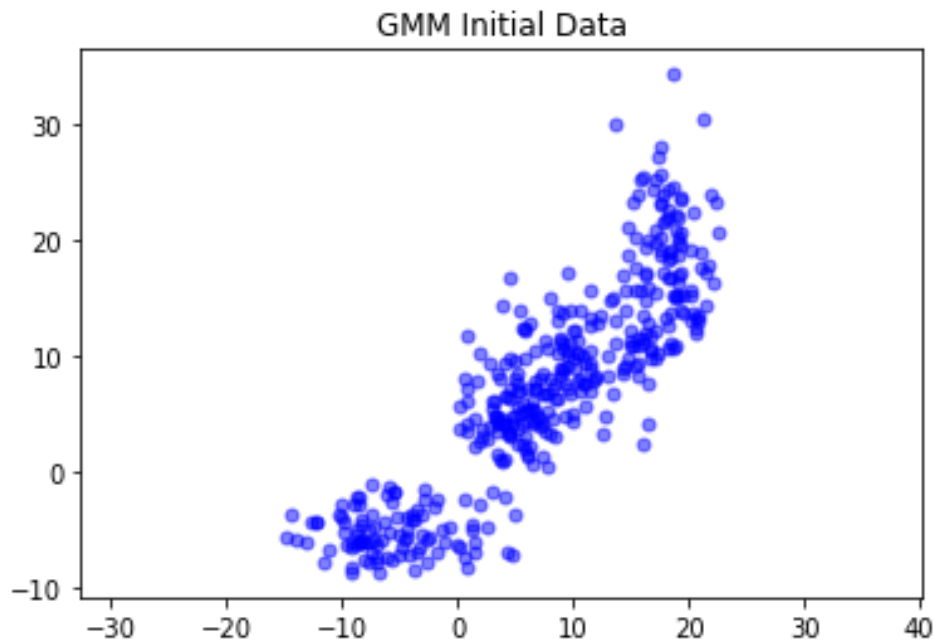
### 5.1 Classical VI is inefficient:
- Do some local computation *for each data point.*
- Aggregate these computations to re-estimate global structure.
- Repeat.
- This cannot handle massive data.

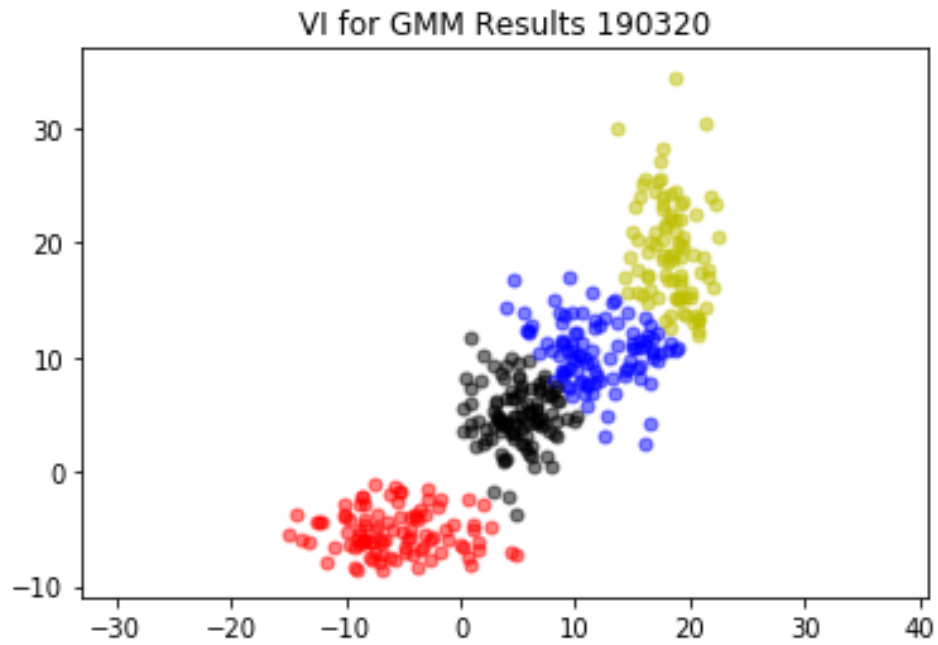### 5.2  Stochastic VI more efficient:
- Stochastic variational inference (SVI) scales VI to massive data.
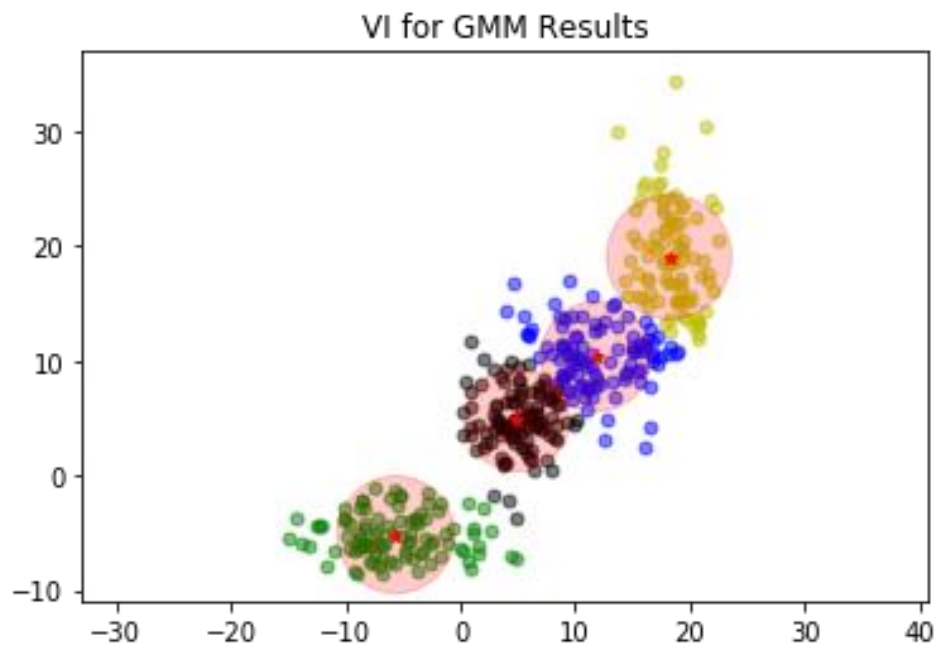
## 6. simulation results

### 6.1 Initial data:



GMM Initial Data

### 6.2 Visualization of results 1

VI for GMM Results 190320

## 6.3 Visualization of results 2



VI for GMM Results

7. references

1. Variational Inference: A Review for Statisticians, https://arxiv.org/abs/1601.00670
2. VARIATIONAL INFERENCE: FOUNDATIONS AND INNOVATIONS
   http://www.cs.columbia.edu/~blei/talks/Blei_VI_tutorial.pdf