

Expectation-Maximization algorithm:
Gaussian mixture models and Sigmoid belief networks

Caio Corro

Gaussian Mixture Model

Random variables

- ▶ \mathcal{Y} taking values in $\{1\dots k\}$ — represents the choice of one (latent) cluster in k
- ▶ \mathcal{X} taking values in \mathbb{R}^d — represents the observed point

Generative story

1. $y \sim p_\theta(\mathcal{Y})$
2. $\mathbf{x} \sim p_\theta(\mathcal{X}|\mathcal{Y} = y)$

=> locally normalized models

Parameterization: $\theta = \{\boldsymbol{\lambda}, \boldsymbol{\mu}, \boldsymbol{\sigma}^2\}$

- ▶ Prior distribution: $\boldsymbol{\lambda} \in \Delta(k)$, i.e. $p_\theta(\mathcal{Y} = y) = \lambda_y$.
- ▶ Conditional distribution: $\boldsymbol{\mu} \in \mathbb{R}^{k \times d}$ and $\boldsymbol{\sigma}^2 \in \mathbb{R}_{++}^{k \times d}$,
i.e. $p_\theta(\mathbf{x}|y) = \prod_{i=1}^d p_\theta(x_i|y) = \prod_{i=1}^d f(x_i, \mu_{y,i}, \sigma_{y,i}^2)$,
where f is the PDF of univariate Gaussian distributions.

Gradient-based learning

Let $\mathcal{D} = \{\mathbf{x}\}_{i=1}^n$ be a training dataset of n datapoints.

Training objective

Maximize the log-likelihood of the training data (the evidence of the data)

$$\arg \max_{\theta \in \Theta} \sum_{\mathbf{x} \in \mathcal{D}} \log p_{\theta}(\mathbf{x}) = \arg \min_{\theta \in \Theta} - \sum_{\mathbf{x} \in \mathcal{D}} \log p_{\theta}(\mathbf{x}) = \arg \min_{\theta \in \Theta} - \sum_{\mathbf{x} \in \mathcal{D}} \log \underbrace{\sum_y p_{\theta}(y) p_{\theta}(\mathbf{x}|y)}_{\text{marginalization over latent variables}}$$

where θ is the set of all parameters the GMM and Θ the set of well-defined θ .

Training algorithm

- ▶ We can reparameterize the parameters so they are unconstrained
- ▶ We can simply use gradient descent on the objective with reparameterized variables! :)

Expectation-Maximization algorithm

Why another algorithm?

Gradient descent:

- ▶ Pros: trivial implementation via Pytorch
- ▶ Cons: you need to define a stepsize

Expectation-Maximization:

- ▶ Pros: no stepsize!
- ▶ Cons: you have to write the optimization code yourself

(there are other favors for EM, but outside the scope of this course)

Intuition of EM

- ▶ Derive a parameterized lower bound to maximization objective
- ▶ Interleave maximization of the lower bound parameters and the model parameters

EM objective

Evidence lower bound (ELBO)

$$\log \mathbb{E}_{p(\mathcal{Y})}[p(\mathbf{x}|\mathcal{Y})] \geq \mathbb{E}_{q(\mathcal{Y})}[\log p(\mathcal{Y})p(\mathbf{x}|\mathcal{Y})] + H[q(\mathcal{Y})]$$

where q is a proposal distribution and H the Shannon entropy.

EM objective

Evidence lower bound (ELBO)

$$\log \mathbb{E}_{p(\mathcal{Y})}[p(\mathbf{x}|\mathcal{Y})] \geq \mathbb{E}_{q(\mathcal{Y})}[\log p(\mathcal{Y})p(\mathbf{x}|\mathcal{Y})] + H[q(\mathcal{Y})]$$

where q is a proposal distribution and H the Shannon entropy.

Proposal distribution for GMMs

$$q_{\phi}(\mathcal{Y} = y | \mathcal{X} = \mathbf{x}) = \phi_y^{(\mathbf{x})}$$

where $\phi^{(\mathbf{x})} \in \Delta(k)$ for each \mathbf{x} in the training set.

New objective

$$\begin{aligned} & \max_{\theta \in \Theta} \sum_{\mathbf{x} \in \mathcal{D}} \log \sum_y p_{\theta}(y) p_{\theta}(\mathbf{x}|y) \\ & \geq \max_{\theta \in \Theta, \phi \in \Phi} \sum_{\mathbf{x} \in \mathcal{D}} \mathbb{E}_{q_{\phi}(\mathcal{Y}|\mathbf{x})}[\log p_{\theta}(y) p_{\theta}(\mathbf{x}|y)] + H[q_{\phi}(\mathcal{Y}|\mathbf{x})] \end{aligned}$$

The two steps of EM

EM objective

$$\begin{aligned} & \max_{\theta \in \Theta, \phi \in \Phi} \sum_{\mathbf{x} \in \mathcal{D}} \mathbb{E}_{q_{\phi}(\mathcal{Y}|\mathbf{x})} [\log p_{\theta}(y)p_{\theta}(\mathbf{x}|y)] + H[q_{\phi}(\mathcal{Y}|\mathbf{x})] \\ & = \max_{\theta \in \Theta, \phi \in \Phi} \sum_{\mathbf{x} \in \mathcal{D}} \text{ELBO}(\mathbf{x}, \theta, \phi) \end{aligned}$$

EM algorithm

Compute a sequence of parameters $\phi^{(1)}, \theta^{(1)}, \phi^{(2)}, \theta^{(2)}, \phi^{(3)}, \theta^{(3)} \dots$ as follows:

- ▶ E step: $\phi^{(t+1)} = \arg \max_{\phi \in \Phi} \sum_{i=1}^n \text{ELBO}(\mathbf{x}, \theta^{(t)}, \phi)$
- ▶ M step: $\theta^{(t+1)} = \arg \max_{\theta \in \Theta} \sum_{i=1}^n \text{ELBO}(\mathbf{x}, \theta, \phi^{(t+1)})$

Comments:

- ▶ the new parameters ϕ computed in the E step are immediately used for the M step
- ▶ this is just a block coordinate ascent algorithm

Expectation step 1/4

$$\phi^{(t+1)} = \arg \max_{\phi \in \Phi} \sum_{\mathbf{x} \in \mathcal{D}} \text{ELBO}(\mathbf{x}, \theta^{(t)}, \phi^{(t)})$$

Parameter decomposition

The parameters of the proposal distribution decomposes over training examples, so we can compute each one independently:

$$\phi^{(t+1)}(\mathbf{x}) = \arg \max_{\phi \in \Delta(k)} \text{ELBO}(\mathbf{x}, \theta^{(t)}, \phi)$$

Intuition

For a given set of parameters θ , the E step computes the best proposal distribution possible (i.e. so the bound is the best possible bound)

Expectation step 2/4

KL divergence

$$D_{\text{KL}}[q(\mathcal{Y})|p(\mathcal{Y})] = \sum_y q(y) \log \frac{q(y)}{p(y)}$$

- ▶ non-negative, i.e. $D_{\text{KL}}[q(\mathcal{Y})|p(\mathcal{Y})] \geq 0$
- ▶ non-necessarily symmetric, i.e. $D_{\text{KL}}[q(\mathcal{Y})|p(\mathcal{Y})] \neq D_{\text{KL}}[p(\mathcal{Y})|q(\mathcal{Y})]$
- ▶ null iff the two distributions are equal: $D_{\text{KL}}[q(\mathcal{Y})|p(\mathcal{Y})] = 0 \iff q(y) = p(y), \forall y$

Expectation step 3/4

Evidence lower bound (ELBO)

$$\log \mathbb{E}_{p(\mathcal{Y})}[p(\mathbf{x}|\mathcal{Y})] \geq \mathbb{E}_{q(\mathcal{Y})}[\log p(\mathcal{Y})p(\mathbf{x}|\mathcal{Y})] + H[q(\mathcal{Y})]$$

where q is a proposal distribution and H the Shannon entropy.

ELBO gap

$$\log \mathbb{E}_{p(\mathcal{Y})}[p(\mathbf{x}|\mathcal{Y})] - \mathbb{E}_{q(\mathcal{Y})}[\log p(\mathcal{Y})p(\mathbf{x}|\mathcal{Y})] + H[q(\mathcal{Y})] = D_{\text{KL}}[q(\mathcal{Y})|p(\mathcal{Y}|\mathbf{x})]$$

where D_{KL} is the Kullback-Leibler divergence.

Expectation step 4/4

$$\phi^{(t+1)}(\mathbf{x}) = \arg \max_{\phi \in \Delta(k)} \text{ELBO}(\mathbf{x}, \theta^{(t)}, \phi)$$

How to solve the E step?

By the ELBO gap, we know that the objective is maximized if the proposal distribution is equal to the posterior distribution!

$$q(y|\mathbf{x}) = p(y|\mathbf{x}) = \frac{p(y)p(\mathbf{x}|y)}{\sum_{y'} p(y')p(\mathbf{x}|y')}$$

=> very easy to compute the optimal parameters ϕ in the E step

Maximization step

$$\theta^{(t+1)} = \arg \max_{\theta \in \Theta} \sum_{\mathbf{x} \in \mathcal{D}} \text{ELBO}(\mathbf{x}, \theta, \phi^{(t+1)})$$

How to solve the M step?

Ignore the constraints on the variance parameters, and simply compute the closed form expression using first order optimality methods!

Solution: looks like "weighted" means and variances.

Exercises

1. Compute the ELBO gap
2. Derive the E step solution using KKT conditions instead of the ELBO gap — is the result expected?
3. Compute the closed form solution for the M step (for the E step it's too trivial)

Sigmoid belief network (two layers only)

Random variables

- ▶ \mathcal{Y} taking values in $[0, 1]^k$ (latent)
- ▶ \mathcal{X} taking values in $[0, 1]^d$ (observed)

Parameterization $\theta = \{\mathbf{a}, \mathbf{B}, \mathbf{c}\}$

$$p_{\theta}(\mathbf{x}, \mathbf{y}) = p_{\theta}(\mathbf{y})p_{\theta}(\mathbf{x}|\mathbf{y}) = \prod_{i=1}^k p_{\theta}(y_i) \prod_{i=1}^d p_{\theta}(x_i|\mathbf{y}) = \prod_{i=1}^k \frac{\exp(y_i a_i)}{1 + \exp(a_i)} \prod_{i=1}^d \frac{\exp(x_i(\mathbf{B}_i \mathbf{y} + c_i))}{1 + \exp(\mathbf{B}_i \mathbf{y} + c_i)}$$

where $\mathbf{a} \in \mathbb{R}^k$, $\mathbf{B} \in \mathbb{R}^{d \times k}$ and $\mathbf{c} \in \mathbb{R}^d$.

Generative story

1. $\mathbf{y} \sim p_{\theta}(\mathcal{Y})$
2. $\mathbf{x} \sim p_{\theta}(\mathcal{X}|\mathcal{Y} = \mathbf{y})$

=> sampling from independant Bernoullis

SBN training

Can we directly use gradient ascent to?

Evidence of a training datapoint $\mathbf{x} \in \mathcal{D}$:

$$\log p_{\theta}(\mathbf{x}) = \log \sum_{\mathbf{y}} p_{\theta}(\mathbf{y}) p_{\theta}(\mathbf{x}|\mathbf{y})$$

\Rightarrow sum over 2^k values for \mathbf{y} , **intractable!!!**

We can't even compute the objective, not even mentioning the gradient

SBN training

Can we directly use gradient ascent to?

Evidence of a training datapoint $\mathbf{x} \in \mathcal{D}$:

$$\log p_{\theta}(\mathbf{x}) = \log \sum_{\mathbf{y}} p_{\theta}(\mathbf{y}) p_{\theta}(\mathbf{x}|\mathbf{y})$$

=> sum over 2^k values for \mathbf{y} , **intractable!!!**

We can't even compute the objective, not even mentioning the gradient

Can we apply EM?

- ▶ The proposal distribution $q_{\phi}(\mathbf{y}|\mathbf{x})$ would require 2^k parameters per training point
- ▶ The E step closed-form expression requires summing over 2^k values (denominator in Bayes rule)

=> **intractable again :(**

Mean Field Theory (MFT)

Mean Field assumption

Assume independence between dimension of latent space in the proposal distribution:

$$q_{\phi}(\mathbf{z}|\mathbf{x}) = \prod_{i=1}^k (\phi_i^{(\mathbf{x})})^{z_i} (1 - \phi_i^{(\mathbf{x})})^{1-z_i}$$

where $\phi^{(\mathbf{x})} \in [0, 1]^k$ are the parameters of the proposal distribution associated with observation \mathbf{x} .

What does it changes?

- ▶ Cons: (probably) not possible to have a tight ELBO anymore (gap = 0)
- ▶ Pros: can make computation tractable

ELBO for SBNs

$$\log p_{\theta}(\mathbf{x}) = \log \sum_{\mathbf{y} \in \mathcal{Y}} p_{\theta}(\mathbf{y}) p_{\theta}(\mathbf{x}|\mathbf{y})$$

$$\geq \underbrace{\mathbb{E}_{q_{\phi}(\mathcal{Y}|\mathbf{x})}[\log p_{\theta}(\mathcal{Y})]}_{\text{(a)}} + \underbrace{\mathbb{E}_{q_{\phi}(\mathcal{Y}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathcal{Y})]}_{\text{(b)}} + \underbrace{H^{\text{S}}[q_{\phi}(\mathcal{Y}|\mathbf{x})]}_{\text{(c)}}$$

ELBO for SBNs

$$\begin{aligned}\log p_{\theta}(\mathbf{x}) &= \log \sum_{\mathbf{y} \in \mathcal{Y}} p_{\theta}(\mathbf{y}) p_{\theta}(\mathbf{x}|\mathbf{y}) \\ &\geq \underbrace{\mathbb{E}_{q_{\phi}(\mathcal{Y}|\mathbf{x})}[\log p_{\theta}(\mathcal{Y})]}_{\text{(a)}} + \underbrace{\mathbb{E}_{q_{\phi}(\mathcal{Y}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathcal{Y})]}_{\text{(b)}} + \underbrace{H^{\text{S}}[q_{\phi}(\mathcal{Y}|\mathbf{x})]}_{\text{(c)}} \\ &\geq \underbrace{\langle \mathbf{a}, \phi^{(\mathbf{x})} \rangle - \sum_{i=1}^k \log(1 + \exp(a_i))}_{\text{(a)}}\end{aligned}$$

ELBO for SBNs

$$\log p_{\theta}(\mathbf{x}) = \log \sum_{\mathbf{y} \in \mathcal{Y}} p_{\theta}(\mathbf{y}) p_{\theta}(\mathbf{x} | \mathbf{y})$$

$$\geq \underbrace{\mathbb{E}_{q_{\phi}(\mathcal{Y} | \mathbf{x})}[\log p_{\theta}(\mathcal{Y})]}_{\text{(a)}} + \underbrace{\mathbb{E}_{q_{\phi}(\mathcal{Y} | \mathbf{x})}[\log p_{\theta}(\mathbf{x} | \mathcal{Y})]}_{\text{(b)}} + \underbrace{H^S[q_{\phi}(\mathcal{Y} | \mathbf{x})]}_{\text{(c)}}$$

$$\geq \underbrace{\langle \mathbf{a}, \boldsymbol{\phi}^{(\mathbf{x})} \rangle - \sum_{i=1}^k \log(1 + \exp(a_i))}_{\text{(a)}}$$

$$\underbrace{\langle \mathbf{x}, \mathbf{B}\boldsymbol{\phi}^{(\mathbf{x})} \rangle + \langle \mathbf{x}, \mathbf{c} \rangle - \sum_{i=1}^d \log \left(1 + \exp(c_i) \prod_{j=1}^k (1 - \phi_j^{(\mathbf{x})} + \phi_j^{(\mathbf{x})} \exp(B_{i,j})) \right)}_{\text{lower bound on (b)}}$$

ELBO for SBNs

$$\log p_{\theta}(\mathbf{x}) = \log \sum_{\mathbf{y} \in \mathcal{Y}} p_{\theta}(\mathbf{y}) p_{\theta}(\mathbf{x} | \mathbf{y})$$

$$\geq \underbrace{\mathbb{E}_{q_{\phi}(\mathcal{Y} | \mathbf{x})}[\log p_{\theta}(\mathcal{Y})]}_{\text{(a)}} + \underbrace{\mathbb{E}_{q_{\phi}(\mathcal{Y} | \mathbf{x})}[\log p_{\theta}(\mathbf{x} | \mathcal{Y})]}_{\text{(b)}} + \underbrace{H^{\text{S}}[q_{\phi}(\mathcal{Y} | \mathbf{x})]}_{\text{(c)}}$$

$$\geq \underbrace{\langle \mathbf{a}, \boldsymbol{\phi}^{(\mathbf{x})} \rangle - \sum_{i=1}^k \log(1 + \exp(a_i))}_{\text{(a)}}$$

$$\underbrace{\langle \mathbf{x}, \mathbf{B}\boldsymbol{\phi}^{(\mathbf{x})} \rangle + \langle \mathbf{x}, \mathbf{c} \rangle - \sum_{i=1}^d \log \left(1 + \exp(c_i) \prod_{j=1}^k (1 - \phi_j^{(\mathbf{x})} + \phi_j^{(\mathbf{x})} \exp(B_{i,j})) \right)}_{\text{lower bound on (b)}}$$

$$\underbrace{\sum_{i=1}^k \left(\phi_i^{(\mathbf{x})} \log \phi_i^{(\mathbf{x})} + (1 - \phi_i^{(\mathbf{x})}) \log(1 - \phi_i^{(\mathbf{x})}) \right)}_{\text{(c)}}$$

Mean Field EM for SBN

Algorithm

- ▶ E step: maximize the ELBO:
 1. write down the first-order optimality condition for ϕ
 2. solve the resulting problem with iterative equation solving method
- ▶ M step: one step of gradient ascent on the ELBO wrt model parameters θ

..

Difference with GMMs

- ▶ Cannot use the ELBO gap and Bayes rule in EM
- ▶ No closed-form expression for step E or M

Upper bound on the evidence

Question

As we cannot close the Elbo gap in the E step,
how to evaluate the quality of the resulting bound?

Upper bound on the evidence

Question

As we cannot close the Elbo gap in the E step,
how to evaluate the quality of the resulting bound?

Variational formulation of the sigmoid

$$\sigma(u) = \inf_{\epsilon \in [0,1]} \exp(\epsilon u - H^{\text{FD}}[\epsilon])$$

where $H^{\text{FD}}[\epsilon] = -\epsilon \log \epsilon - (1 - \epsilon) \log(1 - \epsilon)$ is the Fermi-Dirac entropy.

Upper bound

For any $\epsilon \in [0, 1]^d$:

$$\log p_{\theta}(\mathbf{x}) \leq \left(\begin{array}{l} -\sum_{i=1}^d H^{\text{FD}}[\epsilon_i] + \sum_{i=1}^d \epsilon_i c_i (2x_i - 1) \\ + \sum_{j=1}^k \left(1 - \sigma(a_j) + \sigma(a_j) \exp(\sum_{i=1}^d \epsilon_i \mathbf{B}_{i,j} (2x_i - 1)) \right) \end{array} \right)$$

(not trivial to find the best variational parameters for this bound!)