

Variational Inference:

Aim: we want to know $P(z|X)$.

Given: Prior of $P(z)$: A prior knowledge of what $P(z)$ is like, no business with observed data X . A candidate searching family $q(z)$, which means that what kind of distribution is most likely to be used to approximate $P(z|X)$.

Theory dependence:

- Measuring how $q(z)$ approximate the true target $P(z|X) \rightarrow$ KL divergence.

$$\begin{aligned} \text{KL}(q(z) || P(z|X)) &= Z_q(\log q(z)) - \underbrace{Z_q(\log P(z|X))}_{\text{we want this to be replaced.}} \\ &= Z_q(\log q(z)) - Z_q(\log P(z, x)) + \underbrace{\log P(x)}_{\text{|| Intractable}} \end{aligned}$$

$$\begin{aligned} \text{Then we get concept of ELBO: } Z_{\text{ELBO}}(q(z)) &= Z(\log P(X|z)) - Z(\log q(z)) \\ &= Z(\log P(X|z)) + Z(\log P(z)) - Z(\log q(z)) \\ &= Z(\log P(X|z)) + \text{KL}(q || P(z)) \end{aligned}$$

- From above, we can see that we continuously searching optimal q , monitoring ELBO (indicating how near it is between q and P) to decide when we stop updating q with parameters.

Algorithms:

Initialize parameters deciding $q_z(z|\theta) \Rightarrow \theta_0$.

WHILE NOT converge:

 updating $q_z(z|\theta)$ using updating rules for θ

 Computing Z_{ELBO} to decide whether it is converged.

 where the prior of $P(z)$ engaged.

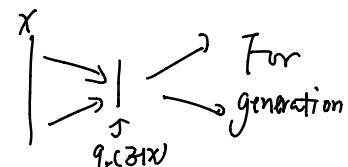
mean-field assumption

$$q_z(z) = \prod_{i=1}^m q_z(z_i)$$

Then we get full joint prob

Variational Autoencoders:

Proposal of an AEB algorithm which can optimize learning:



Inference as optimization

Here inference is proposed as an optimization problem, which means that we want to find optimal q , which can minimize $\text{KL}(q||p)$

Expectation maximization

EM can be seen as combination of two steps:

F step: Inference $q(h|v) = p(h|v)$ (Here in EM we can particularly calculate the exact form of $p(h|v)$, thus inference of $p(h|v)$ directly done, which also means that $\text{KL}(q||p)=0$, in F step.

M step: Learning algorithms using maximum likelihood, $\Theta^* = \underset{\Theta}{\operatorname{argmax}} P(v|\Theta, h)$, after updating $\Theta = \Theta^*$. $q(h|v) \neq p(h|v, \Theta^*)$, thus $\text{KL}(q||p) \neq 0$, again back to F step to do the inference.

MAP inference and Sparse coding

- MAP inference: Instead of specifying full distribution $p(h|v)$, we just use one of the maximum to stand for $h^* = \underset{h}{\operatorname{argmax}} P(h|v)$

- Sparse Coding:

$$\text{Sparse Priority: } P(h_i) = \lambda e^{-\lambda|h_i|}$$

$$\text{Data distribution: } P(X|h) = N(v | Wh + b, \beta^T I)$$

Sparse coding minimize:

$$J(H, w) = \sum_{i,j} |H_{i,j}| + \sum_{i,j} (V - HW^T)_{ij}^2 \quad \text{Optimizing this problem by iteratively maximizing wrt H and wrt w.}$$

Variational Inference and Learning

- Term Clarification: Inference, we always refer to find distribution $P(h|v)$, learning we refer to find optimal parameters Θ .
- Two KL divergence stand for:

$KL(Q||P)$: Q must be small where P is small (The divergence we need to minimize when we do variational inference)

$KL(P||Q)$: Q must be large where P is large (The divergence we want to minimize when we are learning with maximum likelihood)

- Discrete Latent Variables

- In simplest case, h_i is binary indicating which mode is engaged, we can assume $q(h_i=1|v) = \hat{h}_i$, then solving optimization problems with iteratively fixed point: $\frac{\partial}{\partial \hat{h}_i} L = 0$

- Binary sparse coding as an example:

$$P(h_i=1) = \delta(b_i) \rightarrow \text{parameters for learning}$$

$P(v|h) = \mathcal{N}(v|wh, \beta^{-1})$ Observed data consists of summation of different Gaussians
Then for maximizing likelihood:

$$\begin{aligned} \frac{\partial}{\partial b_i} \log P(v) &= \frac{\frac{\partial}{\partial b_i} P(v)}{P(v)} = \frac{\frac{\partial}{\partial b_i} \sum_h P(v|h) \cdot p(h)}{P(v)} \\ &= \sum_h P(v|h) \frac{\frac{\partial}{\partial b_i} p(h)}{P(v)} = \sum_h p(h|v) \frac{\partial}{\partial b_i} \log P(h) \end{aligned}$$

Here we need to inference $p(h|v)$

\Rightarrow When engaging mean-field approximation, and we select $q(h_i=1|v) = \hat{h}_i$, we get: $L(v, \theta, q) = Z(\log p(h) + \log P(v|h) - \log q(h|v))$ for variational inference.

$\underbrace{\quad}_{\text{prior}}$ $\underbrace{\quad}_{\text{likelihood}}$ $\underbrace{\quad}_{\text{estimation}}$

updating q , then learning b_i , we got the final solution.

- Relationship with RNN: According to the updating $q(h_i|v)$ rule, we can see that we update h_i iteratively, which means that we change $q(h_i|v)$ one "h" pertime according to last time results:



- Relationship with autoencoder: $x \rightarrow$ encode by h_i modified \rightarrow decode for optimizing $\Theta \rightarrow$ encode for h_j modified $\rightarrow \dots \rightarrow$ converge.
- Calculus of variations
A math technique used for continuous variational inference when updating q_j .
- Continuous latent variables
Same algorithms, use calculus of variations to maximize ℓ with respect to q_j .
- Interaction between learning and inferencing.
When we are learning parameters Θ^* , it is hard to judge if our variational inference works well, if the true optimal Θ^* is too complicated for our q , to learn $p(h|v, \Theta^*)$ our learning will fail or mis-leading.

Learned Approximate Inference

- What neural network can help us?
when using variational inference, it's like a mapping from all v to $q^*(h|v, \Theta)$, namely the calculation process of q , when given v (or searching process), this mapping $q^* = f(v, \Theta)$ can be modeled using neural network.
- Wake-sleep
Here comparing biological dreaming to calculating $p(w, h)$ and inference $p(h|v)$, awaking to be learning Θ , which maximizes $p(v|h, \Theta)$
- VAE: Definition of loss $P(z|x)$ and $Q(z|x)$, then use neural networks to learn it by itself.