

Probabilistic Graphical Models

Lecture 27,28,29

Variational Inference
Mean-field algorithm
Evidence Lower Bound

Remember: Inference



K. N. Toosi
University of Technology

$$P(X_e | X_e = x_e) \Rightarrow \boxed{P(X_c)} \text{ marginalization}$$



$$\boxed{\phi_c(X_c)}$$

$$x_1^*, \dots, x_n^* = \underset{x_1, \dots, x_n}{\operatorname{arg\,max}} P(x_1, x_2, \dots, x_n)$$

$$P(X_c) = \sum_{X \setminus X_c} P(X_1, \dots, X_n)$$

How to solve this using an optimization problem?

Remember: Latent Variable Models



Latent variable models

$$P_{\theta}(X, Z) = \checkmark$$

data x^1, x^2, \dots, x^n

$$p(Z|X) = \checkmark \quad \text{Hard to compute}$$

Situations where we need normalization



K. N. Toosi
University of Technology

$$P(x) = \frac{1}{Z} \tilde{P}(x) = \frac{1}{Z} e^{F(x)}$$

Annotations: $F(x)$ is labeled "sum of MRF"; $\tilde{P}(x)$ is labeled "product of factors".

BN or MRF with ~~the~~ evidence

$$P(x) = P(x_t, x_e) \implies \text{need } P(x_t | x_e) = \frac{P(x_t, x_e)}{\sum_{x_t} P(x_t, x_e)}$$

Annotations: x_e is labeled "evidence (known)"; the denominator is circled in red.

$$P(x) = P(y, z)$$

Annotation: y is labeled "latent".

$$P(z | y) = \frac{P(z, y)}{\sum_z P(z, y)}$$

Annotations: $P(z, y)$ is circled in red and labeled $z(x_e)$; the denominator is circled in red and labeled $z(y)$.

Variational Inference on MRFs



K. N. Toosi
University of Technology

$$P(X) = \frac{1}{Z} \tilde{P}(X)$$

complex

(Hard to do inference on)

find $Q(X)$ $\left\{ \begin{array}{l} Q(X) \text{ is easy to handle} \\ Q(X) \text{ is close to } P(X) \end{array} \right.$

Kullback-Leibler (KL) Divergence



Measure distance between $Q(X), P(X)$

pgm 27 (II)

KL-divergence

$$\begin{aligned} \textcircled{Q} \text{KL}(Q \parallel P) &= \sum_x Q(x) \log \frac{Q(x)}{P(x)} = E_Q \left\{ \log \frac{Q(x)}{P(x)} \right\} \\ &= \sum_x Q(x) \left[\log Q(x) - \log P(x) \right] \end{aligned}$$

$$P(x) = Q(x) \text{ for all } x \Rightarrow \text{KL}(Q \parallel P) = 0$$

$$\text{KL}(Q \parallel P) = 0 \Rightarrow P(x) = Q(x) \text{ for } \underline{\text{all } X}$$

$$\text{KL}(Q \parallel P) \geq 0$$

almost all (continuous)

$$\text{KL}(Q \parallel P) \neq \text{KL}(P \parallel Q) \text{ in general}$$

KL Divergence and Entropy

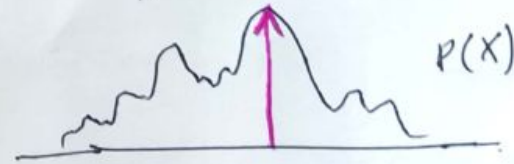


K. N. Toosi
University of Technology

$$KL(Q \parallel P) = \sum_x Q(x) \log Q(x) - \sum_x Q(x) \log P(x)$$

what choice of Q maximizes $\sum_x Q(x) \log P(x)$?

$$\sum_x Q(x) = 1$$
$$\int Q(x) = 1$$



$$\text{Entropy} = \sum_x Q(x) \log \frac{1}{Q(x)} = - \sum_x Q(x) \log Q(x) = H(Q)$$

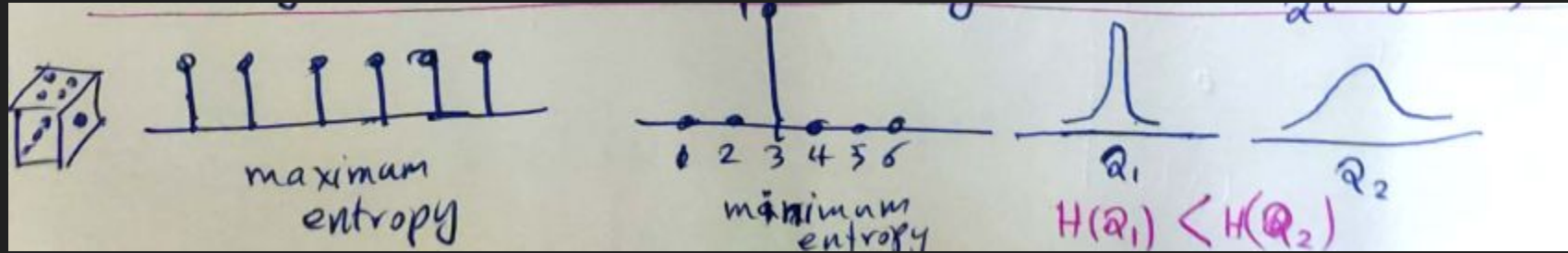
$$KL(Q \parallel P) = -H(Q) - E_Q \{ \log P(x) \}$$

minimizing $KL(Q \parallel P) \equiv$ maximizing $H(Q) + E_Q \{ \log P(x) \}$

KL Divergence and Entropy



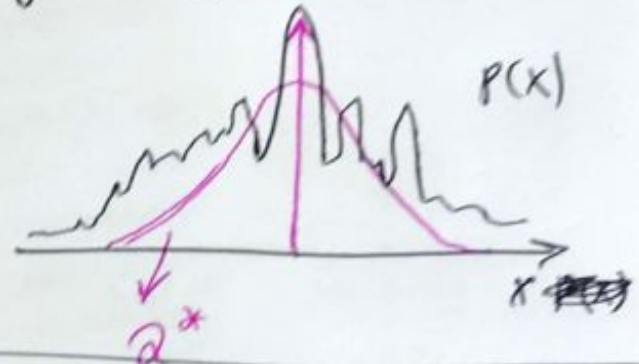
K. N. Toosi
Technology



$$KL(Q \parallel P) = -H(Q) - \mathbb{E}_Q \{ \log P(X) \}$$

pgm 27 (III)

minimize $KL(Q \parallel P)$
 $Q \in \bar{Q}_{\text{est}}$



Example: Fully Factorized Q



K. N. Toosi
University of Technology

Simple case: $Q(X) = Q(X_1, X_2, \dots, X_n) = Q_1(X_1) Q_2(X_2) \dots Q_n(X_n)$

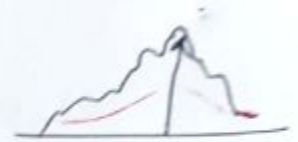
fully factorized case

KL Divergence is not Symmetric



$$P(x) = \frac{1}{Z} \tilde{P}(x) = \frac{1}{Z} e^{F(x)}$$

$$\min_Q \text{KL}(Q \parallel P) = \sum_x Q(x) \log \frac{Q(x)}{P(x)}$$



why not minimize $\text{KL}(P \parallel Q) = \sum_x \underline{P(x)} \log \frac{P(x)}{Q(x)}$?

$= \sum_x P(x) \log P(x) - \sum_x P(x) \log Q(x)$

Variational Lower Bound



$$\begin{aligned}
 KL(Q \parallel P) &= \sum_x Q(x) \left[\log Q(x) - \log P(x) \right] \\
 &= \sum_x Q(x) \left[\log Q(x) - \log \frac{1}{Z} \tilde{P}(x) \right] \\
 &= \sum_x Q(x) \log Q(x) - \sum_x Q(x) \log \tilde{P}(x) + \underbrace{\sum_x Q(x) \log Z}_{\log Z} \\
 &= \sum_x Q(x) \log Q(x) - \underbrace{\sum_x Q(x) \log \tilde{P}(x)}_{F(Q)} + \log Z
 \end{aligned}$$

$$KL(Q \parallel P) = \underbrace{\sum_x Q(x) \log \frac{Q(x)}{\tilde{P}(x)}}_{-L(Q)} + \log Z$$

$L(Q) = E_Q \left\{ \frac{\tilde{P}(x)}{Q(x)} \right\}$

$$L(Q) = +\log Z - KL(Q \parallel P) \Rightarrow L(Q) \leq \log(Z)$$

$KL > 0$ variational lower bound!

The Mean-field Algorithm



Meanfield inference $q(x) = q(x_1, \dots, x_n) = \prod_{i=1}^n q_i(x_i)$
fully factorized

$$q^* = \operatorname{argmin}_q \sum_x q(x) \log \frac{q(x)}{\tilde{P}(x)} + \log Z$$

$$= \operatorname{argmin}_q \sum_x q(x) \log \frac{q(x)}{\tilde{P}(x)}$$

$$\sum_x q(x) \log \frac{q(x)}{\tilde{P}(x)} = \boxed{\sum_x q(x) \log q(x)} - \underbrace{\sum_x q(x) \log \tilde{P}(x)}_{F(x)}$$

$$\sum_x q(x) \log q(x) = \sum_{x_1} \sum_{x_2} \dots \sum_{x_n} q_1(x_1) q_2(x_2) \dots q_n(x_n) \left[\sum_{i=1}^n \log q_i(x_i) \right]$$

$$= \sum_{i=1}^n \sum_{x_i} q_i(x_i) \log q_i(x_i)$$

The Mean-field Algorithm



K. N. Toosi
University of Technology

$$\begin{aligned}\sum_x q(x) \log \tilde{P}(x) &= \sum_x q(x) F(x) = \sum_x q(x) \sum_c F_c(x_c) \\ &= \sum_c \sum_x q(x) F_c(x_c) = \sum_c \sum_{x_c} \left[\sum_{x \setminus x_c} q(x_c, x \setminus x_c) \right] F_c(x_c) \\ &= \sum_c \sum_{x_c} q(x_c) F_c(x_c)\end{aligned}$$

$$\sum_i q_i(x_i) \log q_i(x_i) - \sum_c \sum_{x_c} q(x_c) F_c(x_c)$$

The Mean-field Algorithm for Pairwise MRF



$$P(x) = \frac{1}{Z} e^{F(x)} \quad , \quad F(x) = \sum_{i=1}^n F_i(x_i) + \sum_{(i,j) \in \mathcal{E}} F_{ij}(x_i, x_j)$$

$$\begin{aligned} -L(q, \tilde{P}) &= \sum_i \sum_{x_i} q_i(x_i) \log q_i(x_i) - \sum_{i=1}^n q_i(x_i) F_i(x_i) \\ &\quad - \sum_{(i,j) \in \mathcal{E}} q_i(x_i) q_j(x_j) F_{ij}(x_i, x_j) \end{aligned}$$

$$\begin{aligned} L(q, \tilde{P}) &= - \sum_i \sum_{x_i} q_i(x_i) \log q_i(x_i) + \sum_{i=1}^n \sum_{x_i} q_i(x_i) F_i(x_i) \\ &\quad + \sum_{(i,j) \in \mathcal{E}} \sum_{x_i} \sum_{x_j} q_i(x_i) q_j(x_j) F_{ij}(x_i, x_j) \end{aligned}$$

The Mean-field Algorithm for Pairwise MRF



K. N. Toosi
University of Technology

$$Q(X) = Q(X_1, X_2, \dots, X_n) = \prod_{i=1}^n q_i(X_i)$$

pgm 28 (II)

mean-field

$$P(X) = \frac{1}{Z} \tilde{P}(X) = \frac{1}{Z} e^{F(X)}$$

$$\sum_{i=1}^n F_i(X_i) + \sum_{(i,j) \in \mathcal{E}} F_{ij}(X_i, X_j)$$

$$\mathcal{L}(Q, \tilde{P}) = - \sum_{i=1}^n \sum_{X_i} \underline{q_i(X_i)} \log q_i(X_i) + \sum_{i=1}^n \sum_{X_i} q_i(X_i) F_i(X_i)$$

$$+ \sum_{(i,j) \in \mathcal{E}} \sum_{X_i} \sum_{X_j} q_i(X_i) q_j(X_j) F_{ij}(X_i, X_j)$$

The Mean-field Algorithm for Pairwise MRF



$\max_{\theta, p} L(\theta, p)$ subject to $\sum_{k=1}^L q_{i,k} = 1$ for $i=1 \dots n$

coordinate ascent
 Loop for $i=1 \dots n$

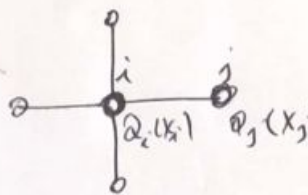
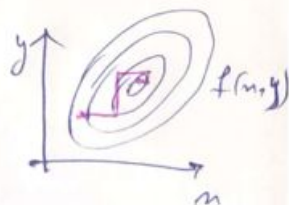
$\max_{q_{i1}, q_{i2}, \dots, q_{iL}}$

$\Rightarrow \frac{\partial}{\partial q_{i,l}} L(\theta, p) + \lambda \left(\sum_{k=1}^L q_{i,k} - 1 \right) = 0$

$= -\log q_{i,l} - 1 + F_i(l) + \sum_{j \in N_i} \sum_{k=1}^L q_{j,k'} F_{ij}(l, k') + \lambda = 0$

for $l=1, 2, \dots, L$

$\sum_{k=1}^L q_{i,k} = 1$

The Mean-field Algorithm for Pairwise MRF



pgm 28

$$\Rightarrow \log q_{il} = F_i(l) + \sum_{j \in N_i} \sum_{k=1}^L q_{jk} F_{ij}(l, k) + \lambda - 1$$

$$q_{il} = \exp(F_i(l) + \sum_{j \in N_i} \sum_{k=1}^L q_{jk} F_{ij}(l, k)) e^{\lambda - 1}$$

$l = 1, 2, \dots, L$ $\sum_{k=1}^L q_{ikl} = 1$

$q'_{il} = \exp(F_i(l) + \sum_{j \in N_i} \sum_{k=1}^L q_{jk} F_{ij}(l, k))$

$$q_{il} \leftarrow \frac{q'_{il}}{\sum_{k=1}^L q'_{ik}}$$

mean-field message passing

Meanfield message passing



$q'_{il} = \exp (F_i(l) + \sum_{j=N_i} \sum_{k=1}^L q_{jk} F_{ij}(l,k))$

$q_{il} \leftarrow \frac{q'_{il}}{\sum_{k=1}^L q'_{ik}}$ mean-field message passing

CRF

Variational Inference: Continuous Case



K. N. Toosi
University of Technology

Case 2: Continuous $Q_i(X_i)$



Among all $Q_i(X)$ with $Q_i(X) \geq 0$, $\int Q_i(X) d\mu = 1$
 $i=1, \dots, n$

find the one that maximizes $L(Q, \tilde{P})$
minimizes $KL(Q \| P)$

Variational Calculus $f: \mathbb{R}^m \rightarrow \mathbb{R}^n$ is a
optimize $J(f)$ s.t. so some constraint function

Latent Variable Models and Variational Learning



Variational Learning

$$P_{\theta}(x) = \frac{1}{Z(\theta)} \tilde{P}_{\theta}(x) = \frac{1}{Z(\theta)} e^{F_{\theta}(x)}$$

Latent Variable Models

$$P_{\theta}(x) = \int P_{\theta}(x, z) dz$$

Data x^1, x^2, \dots, x^m

$$\max_{\theta} \sum_i \log P_{\theta}(x^i)$$

$$\max_{\theta} \ell(\theta) = \max_{\theta} \sum_{i=1}^m \log P_{\theta}(x^i)$$

$$= \max_{\theta} \sum_{i=1}^m \log \int P_{\theta}(x^i, z) dz$$

Classic Method fail!



Solution 1

$$\frac{\partial \ell(\theta)}{\partial \theta_k} = \frac{\partial}{\partial \theta_k} \sum_{i=1}^m \log \int p_{\theta}(x^i, z) dz$$

$$= \sum_{i=1}^m \frac{\int \frac{\partial}{\partial \theta_k} p_{\theta}(x^i, z) dz}{\int p_{\theta}(x^i, z) dz}$$

hard to compute

Solution 2: Expectation - Maximization

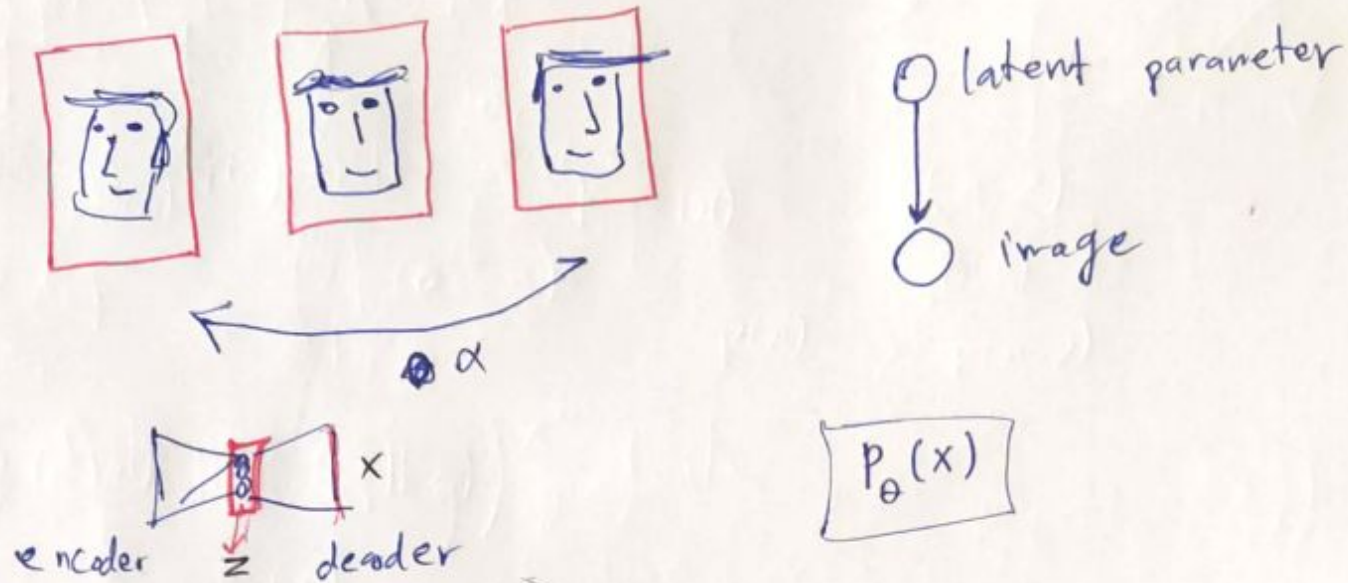
⇒ Need posterior $p_{\theta}(z|x^i)$ for $i=1, \dots, m$
posterior *hard to compute*

Variational Learning



K. N. Toosi
University of Technology

Variational Learning for latent variable models pgm 29



Remember: Latent Variable Models



model $P_{\theta}(x) = \sum_z P_{\theta}(x, z)$ latent variable model

Data = $\{x^1, x^2, \dots, x^m\} = \mathcal{D}$

$$\max_{\theta} \ell(\theta, \mathcal{D}) = \sum_{i=1}^m \log P_{\theta}(x^i) = \sum_{i=1}^m \log \sum_z P_{\theta}(x^i, z)$$

Usually $P_{\theta}(x, z) = \boxed{P_{\theta}(x|z) P(z)}$

or $\boxed{P_{\theta_1}(x|z) P_{\theta_2}(z)}$

Example

Remember: Latent Variable Models



Example

$$P_{\theta}(x, z) = P_{\theta}(x|z) P(z)$$

$$P_{\theta}(x|z) \cong \mathcal{N}(x; \mu(z), I\sigma(z))$$

- 1 $z \sim P(z)$
- 2 compute $\mu(z), \sigma(z)$
- 3 $x \sim P_{\theta}(x|z) = \mathcal{N}(x; \mu(z), I\sigma(z))$

a deep neural net

How to train the neural network

Need to compute $\frac{\partial}{\partial \theta_i} \sum_{i=1}^m \log \sum_z P_{\theta}(x^i, z)$

Intractable \leftarrow

$$= \frac{\partial}{\partial \theta_i} \sum_{i=1}^m \log \left(\sum_z P_{\theta}(x^i|z) P(z) \right)$$

Scanned with CamScanner

for EM need $P(z|x)$ \rightarrow posterior \rightarrow intractable!

VI on MRFs vs Latent Variables Models



K. N. Toosi
University of Technology

MRF	latent variable model
$P_{\theta}(y) = \frac{1}{Z(\theta)} \tilde{P}(y)$	$\underline{P_{\theta}(z x)} = \frac{1}{P_{\theta}(x)} P_{\theta}(x, z)$
$Z(\theta) = \sum_y \tilde{P}(y)$	$P_{\theta}(x) = \sum_z P_{\theta}(x, z)$
minimize over q $KL(q \parallel P_{\theta}(y))$	minimize over $q(z x)$ $KL(\underline{q(z x)} \parallel P_{\theta}(z x))$ ↓ approximate posterior Example: $q(z x) = \prod q_i(z_i x)$

KL-divergence



$$\begin{aligned} \text{KL}(q(z|x) \parallel P_{\theta}(z|x)) &= \sum_z q(z|x) \log \frac{q(z|x)}{P_{\theta}(z|x)} \\ &= \sum_z q(z|x) \log \frac{q(z|x)}{\frac{P_{\theta}(x,z)}{P_{\theta}(x)}} \rightarrow \text{easy} \\ &= \sum_z q(z|x) \left(\log \frac{q(z|x)}{P_{\theta}(x,z)} + \log P_{\theta}(x) \right) \rightarrow \text{hard} \\ &= \sum_z q(z|x) \log \frac{q(z|x)}{P_{\theta}(x,z)} + \sum_z q(z|x) \log P_{\theta}(x) \rightarrow \text{independent of } z \\ &= -\sum_z q(z|x) \log \frac{P_{\theta}(x,z)}{q(z|x)} + \log P_{\theta}(x) \end{aligned}$$

Evidence Lower Bound (ELBO)



K. N. Toosi
University of Technology

$$\sum_z q(z|x) \log \frac{P_\theta(x, z)}{q(z|x)} = \log P_\theta(x) - \text{KL}(q(z|x) \parallel P_\theta(z|x))$$

$\mathcal{L}(q, P_\theta)$: variational lower bound

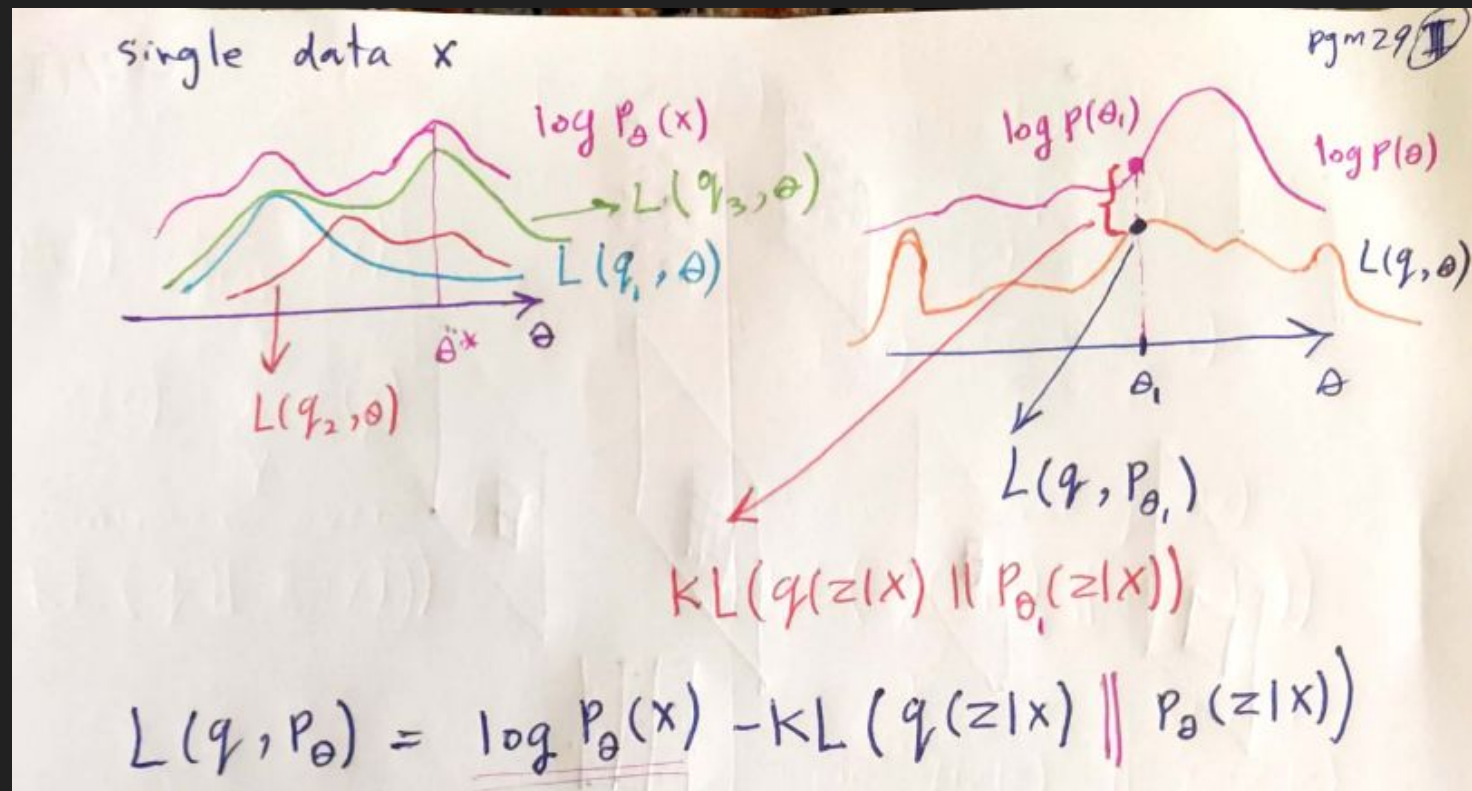
Evidence Lower Bound (ELBO)

$\mathcal{L}(q, P_\theta) \leq \log P_\theta(x)$ for any choice of q

Evidence Lower Bound (ELBO)



K. N. Toosi
University of Technology



Evidence Lower Bound (ELBO)



$$L(q, p_{\theta}) = \log p_{\theta}(x) - \text{KL}(q(z|x) \parallel p_{\theta}(z|x))$$

- 1: Maximizing $L(q, p_{\theta})$ with respect to θ pushes $\log p_{\theta}(x)$ up.
- 2: Maximizing $L(q, p_{\theta})$ w.r.t θ works better when $q(z|x)$ is closer to $p_{\theta}(z|x)$ [$\text{KL}(q(z|x) \parallel p_{\theta}(z|x))$ is smaller]
- 3: Maximizing $L(q, p_{\theta})$ w.r.t. q minimizes $\text{KL}(q(z|x) \parallel p_{\theta}(z|x))$

Variational Learning



K. N. Toosi
University of Technology

$$\max_{q, \theta} \mathcal{L}(q, P_{\theta}) = \sum_z q(z|x) \log \frac{P_{\theta}(z, x)}{q(z|x)}$$

How to compute (approximate) $\mathcal{L}(q, P_{\theta}) = \sum_z \dots$?

$$\mathcal{L}(q, P_{\theta}) = \mathbb{E}_{q(z|x)} \left\{ \log \frac{P_{\theta}(z, x)}{q(z|x)} \right\}$$

$$\approx \frac{1}{P} \sum_{i=1}^P \log \frac{P_{\theta}(z^i, x)}{q(z^i|x)} \text{ where } z^1, z^2, \dots, z^P \sim q(z|x)$$

Variational Learning



Data $x^1, x^2, x^3, \dots, x^m$

$$\max_{\theta} \ell(\theta) = \sum_{i=1}^m \log P_{\theta}(x^i) = \sum_{i=1}^m \log \sum_z \log P_{\theta}(x^i, z)$$

instead

$$\max_{\theta, q} \mathcal{L}(q, P_{\theta}) = \sum_{i=1}^m \sum_z q_i(z|x^i) \log \frac{P_{\theta}(x^i, z)}{q_i(z|x^i)}$$

$$= \mathcal{L}(q, \theta) = \sum_{i=1}^m \mathbb{E}_{q_i(z|x^i)} \{ \log P_{\theta}(x^i, z) \} + H \{ q_i(z|x^i) \}$$

For each x^i a different q might be optimum.

\Rightarrow A different q_i for each x^i

Variational Learning Algorithm



Start from some $\theta \leftarrow \theta_0, q_1, q_2, \dots, q_m$
optimize w.r.t. θ

$$\frac{\partial L(q, P_\theta)}{\partial \theta} = \sum_{i=1}^m E_{q_i} \left\{ \frac{\partial}{\partial \theta} \log P_\theta(x^i; z) \right\}$$

$$= \sum_{i=1}^m E_{q_i} \left\{ \frac{\partial}{\partial \theta} \log P_\theta(x^i | z) \right\}$$

$$\sim \sum_{i=1}^m \sum_{j=1}^{m'} \frac{\partial}{\partial \theta} \log P_\theta(x^i | z_j^i)$$

$$\theta \leftarrow \theta + \lambda \frac{\partial}{\partial \theta} L(q, P_\theta)$$

$z_1^i, z_2^i, \dots, z_{m'}^i \sim q_i(z | x^i)$
usually $m'=1$

maximize e w.r.t. q_1, q_2, \dots, q_m

can do independently for each q_i