

# Variational Inference

Loopy Belief prop 11/4/13

- HW + project progress - ambitious  
- great GS. ~~HW~~ - not here wednesday

+ Expectation propagation - HW

Last week: mean Field variational inference.

Today: LBP, in context of Ising model: Markov Random Field. Review MF <sup>in this context</sup>

Idea behind Variational inference; review for a few minutes.

let  $\{x_1, \dots, x_n\}$  let  $z_{1:m}$  be hidden vars,  $\theta$  params.  
 $z = \{z_{1:m}, \theta\}$

as always, we want the posterior distribution of latent vars:

$$p(z|x, \theta) = \frac{p(z, x, \theta)}{\int_z p(z, x, \theta) dz}$$

- posterior links data to model. used for downstream analyses: posterior predictive dist, interpretation, etc. via pt. estimates of  $z$ .

- Variational inference is one method to compute the approximate posterior dist.

- cant compute the posterior for many models.  
or - very hard. e.g. truncated gaussian, mean/var complex  
Look again at G Mixture Model. unnormalized - easy.

$$x_i \sim N(\mu_{z_i}, \Sigma) \quad z_i \sim \text{Mult}(\pi) \quad \mu_k \sim N(0, \tau^2)$$

posterior  $p(\mu, z|x) \Rightarrow \frac{\prod_{k=1}^K p(\mu_k) \prod_{i=1}^n p(z_i|\pi) p(x_i|z_i, \mu_k)}{\int \sum \prod_{k=1}^K p(\mu_k) \prod_{i=1}^n p(z_i) p(x_i|z_i, \mu_k)}$  <sup>chain rule</sup>

easy for given  $z$

difficult integral. <sup>hard.</sup> exponentially  $K^n$  terms. <sup>independence.</sup>  $d\mu_{1:k}$

We are reviewing this because we are about to launch into LDA models - hidden variables also create exponentially difficult denominator in posterior too - + hard integral.



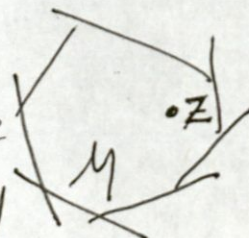
# Variational Inference

→ return to the space of parameters.  $M$

For many models (discrete, gaussian)

~~∃~~ theorem that states, for finite data, the set of parameters ( $Z$ ) consistent with these data is a convex polytope:

Optimization: - Posterior density on this space  
- approximate  $M$ , optimize within  $M$



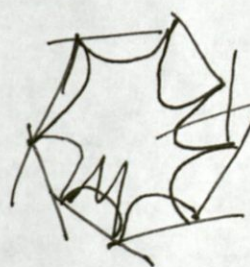
cartoon version

Mean Field: find fully factorized

Approximation to posterior:  $q(z) = \prod_{i=1}^m q(z_i)$

Compute:  $\min_z KL(q(z) || p(z|x))$

each variable is independent



→ based on  $KL$ , we know  $M_{MF}$  is a subset of the space  $M$ .

→ very possible  $z^*$  not contained within  $M_{MF}$ .

In examples we've discussed ~~we can go directly~~ ~~from KL divergence w/ unnormalized posterior.~~ ~~KL divergence is not sufficient.~~ ~~we can go directly~~

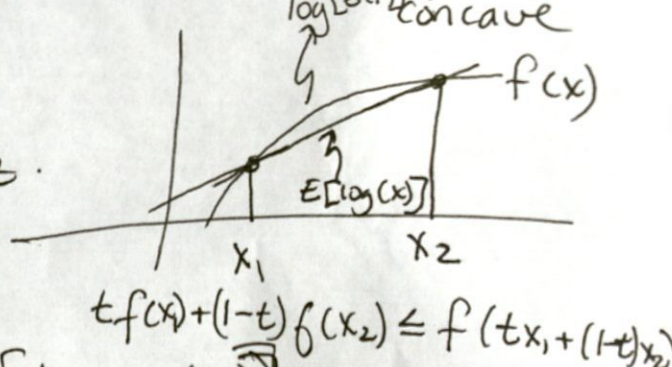
Evidence Lower Bound (ELBO)

log prob of obs:

$$\begin{aligned} \log p(x) &= \log \int_z p(x, z) dz \\ &= \log \int_z p(x, z) \frac{q(z)}{q(z)} dz \\ &= \log E_q \left[ \frac{p(x, z)}{q(z)} \right] \\ &\geq E_q [\log p(x, z)] - E_q [\log q(z)] \end{aligned}$$

Jensen's Inequality

$$\log E(x) \geq E[\log(x)]$$





③ Recall from last week, variational objective

$$KL(q \| \tilde{p}) = \underbrace{-E_q[\log p(z, x)] + E_q[\log q(z)]}_{\text{negative elbo}} + \log p(x)$$

Thus: ELBO is KL divergence w/ unnormalized posterior.  
Since  $p(x)$  indep of  $q$ , this is equiv. to minimizing ELBO.

Now: choose  $q$  st. these expectations are computable  
• maximize  $q(z)$  w.r.t. elbo to get as tight an approximation to  $p(z, x)$  as possible.

Ising Model (MRF)

Denoising image:  $y = \{-1, 1\}$  pixels  
 $x = \{-1, 1\}$  denoiser pixels

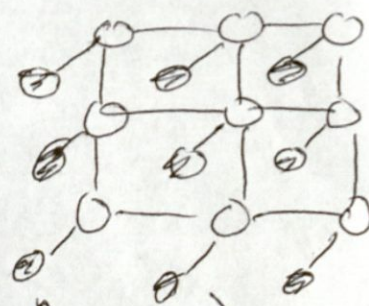
$$\psi_s(x) = p(y_i | x_i) \triangleq L_i(x_i)$$

$$\psi_{st}(x_s, x_t) = x_s x_t \quad w_{st} = 1 \text{ - neighbors want to have similar states}$$

$$p(y|x) = \prod_{i=1}^n \exp(-L_i(x_i)) \text{ likelihood.}$$

$$p(x) = \frac{1}{Z_0} \exp\left(-\sum_{i=1}^n \sum_{j \in \text{nei}(i)} x_i x_j\right) \text{ prior}$$

$$\text{Posterior: } p(x|y) = \frac{1}{Z} \exp\left(-\sum_{i=1}^n \sum_{j \in \text{nei}(i)} x_i x_j - \sum_{i=1}^n L_i(x_i)\right)$$



Mean Field  $q(x) = \prod_{i=1}^n q(x_i)$  let  $\mu_i$  be "mean value" or variational param

$$\log(q_i(x_i)) = E_{\tilde{p}}(\log \tilde{p}(x))$$

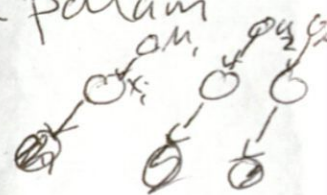
↳ basic coordinate ascent from last class.

$$\text{ELBO: for } x_i \quad \log(q_i(x_i)) = E_{\tilde{p}} \left[ x_i \sum_{j \in \text{nei}(i)} x_j + L_i(x_i) + c \right] \quad \text{21.43 book}$$

$$q_i(x_i) \propto \exp \left\{ x_i \sum_{j \in \text{nei}(i)} \mu_j + L_i(x_i) \right\}$$

$$\text{so } q_i(x_i = +1) = \frac{\exp \left\{ \sum_{j \in \text{nei}(i)} x_j + L_i(+1) \right\}}{\sum_{x_i \in \{-1, +1\}} \exp \left\{ \sum_{j \in \text{nei}(i)} x_j + L_i(x_i) \right\}}$$

→ logistic function to get  $\mu_i$ .  
write out for  $-1$ , take average





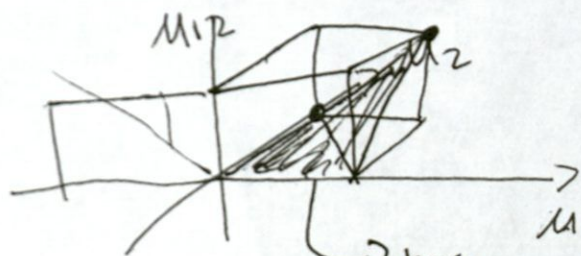
④ Then

11/4/13

$$\mu_i = +1 (q_i(x_i=+1)) + -1 (q_i(x_i=-1))$$

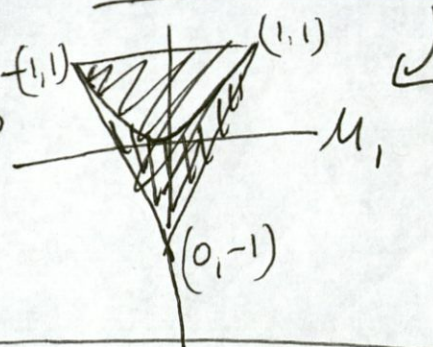
→ MF = GIBBS:  $z \propto \exp \left[ \sum_i \mu_i \log P(z_i | x) \right]$  - looks like conditional

return to drawing



→ marginal polytope  $P(z)$   
 flip over to -1, -1 on each axis.  $\mu_1, \mu_2 \in \{-1, 1\}$   
 Fully factorized  
 $\mu_1^2 = \mu_2^2$

slice at  $\mu_1 = \mu_2$ .



Loopy belief Propagation: LBP.

~~now~~. Now:  $q$  only locally consistent  
 not a valid joint probability.

LC:  $q(x_i) = \sum_{x_j} q(x_i, x_j) \quad \forall \text{ pairs } (i, j) \in E$

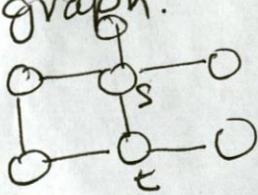
Recall Belief ~~propagation~~ propagation in a tree is exact. Not true in a graphical model with loops. E.g. Forward ~~backward~~ peeling (in tree) → final answer after inside-outside pass.

Idea: apply Belief Propagation to loopy graph.

Same Ising model.

$$\psi_s(x_s) = L_s(x_s)$$

$$\psi_{st}(x_s, x_t) = x_s x_t$$



Algorithm

- messages  $m_{s \rightarrow t}(x_t) = 1 \quad \forall \text{ edges.}$
- beliefs  $\mu_s = 1 \quad \forall \text{ nodes } s.$

repeat:  $m_{s \rightarrow t}(x_t) = \sum_{x_s} L_s(x_s) x_s x_t \prod_{u \in \text{neigh}(s) \setminus t} m_{u \rightarrow s}(x_s)$



5

cont.

11/4/13.

$$M_s \propto L_s(x_s) \prod_{t \in \text{neighbor}(s)} m_{t \rightarrow s}(x_s).$$

return  $M_s$ .

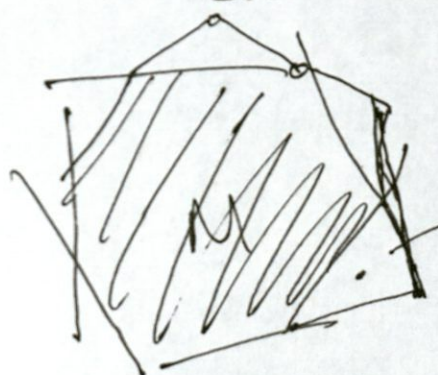
→ Synchronous updates.

Will LBP converge? maybe not - no guarantees.

but methods to help - dampening to avoid oscillation

- asynchronous steps.

- Scheduling (tree based reparametrization)



$$M_s = M_{LBP}$$

When graph is tree

o.w. ~~MC~~  $M \subset M_{LBP}$  outer polytope

$O(|V| + |E|)$  constraints (local consistency)

why:

E Propagation:

In HW.

Twiston MF.

MF vs LBP.

- MF not exact in Trees

- MF - node marginals / LBP node + edge marginals.

- MF - more local optima.