

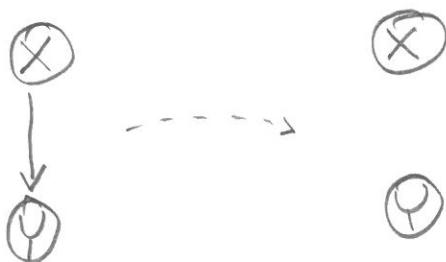
①

Variational Bayes

- find analytic approximation to intractable posterior.

$p(x|y)$ - exact/true posterior

$q(x)$ - approximate posterior



Idea: 1. can we find parameters of $q(x)$, called variational parameters, such that $q(x)$ is as similar to $p(x)$ as possible?

2. Find parameters such that the KL-divergence

$$\text{KL}(q(x) \| p(x|y)) = \int q(x) \cdot \log \frac{q(x)}{p(x|y)} dx =$$

$$= \left\{ P(x|y) = \frac{P(x,y)}{P(y)} \right\} =$$

$$= \int q(x) \cdot \log \frac{q(x) \cdot P(y)}{P(x,y)} dx =$$

$$= \int q(x) \cdot \log \frac{q(x)}{p(x,y)} dx + \underbrace{\int q(x) dx \cdot \log p(y)}_1 = ②$$

$$= \int q(x) \cdot \log \frac{q(x)}{p(x,y)} dx + \log p(y)$$

$$\Rightarrow \log p(y) = KL(q(x) \| p(x|y)) - \int q(x) \cdot \log \frac{q(x)}{p(x,y)} dx$$

$$= KL(q(x) \| p(x|y)) + \int q(x) \cdot \log \frac{p(x,y)}{q(x)} dx$$

$$\Rightarrow \log p(y) \geq \int q(x) \cdot \log \frac{p(x,y)}{q(x)} dx = \mathcal{L}(q(x))$$

\Rightarrow maximise $\mathcal{L}(q(x))$ to minimise KL
equality only when $KL(q(x) \| p(x|y)) = 0$

$$\Rightarrow q(x) = p(x|y)$$

③

Mean Field

$$q(x) = \prod_i q_i(x_i)$$

\Rightarrow we fit the marginals of the distribution

$$\mathcal{L}(q) + KL(q||P) = \log P(y)$$

$$\Rightarrow \mathcal{L}(q) \leq \log P(y)$$

$$\begin{aligned} \mathcal{L}(q) &= \int q(x) \cdot \log \frac{P(y, x)}{q(x)} dx = \left\{ q(x) = \prod_i q_i(x_i) \right\} = \\ &= \int \prod_i q_i(x_i) \cdot \log \frac{P(y, x)}{\prod_k q_k(x_k)} dx = \\ &= \int \prod_i q_i(x_i) \cdot \left(\log P(y, x) - \sum_k \log q_k(x_k) \right) dx = \end{aligned}$$

We want to have a scheme where we optimise each component in turn therefore we would like to re-write $\mathcal{L}(q)$ to single out each component as $\mathcal{L}(q) = \mathcal{L}'(q_j) + \mathcal{L}(q_{\bar{j}})$

$$\mathcal{L}(q_j) = \int \prod_{i=1}^n q_i(x_i) \cdot \left(\log p(y, x) - \sum_k \log q_k(x_k) \right) dx = \quad (4)$$

= {split integral} =

$$= \int_{\bar{j}} \int_{\bar{j}} q_j(x_j) \prod_{i=1, i \neq j}^n q_i(x_i) \left(\log p(x, y) - \sum_k \log q_k(x_k) \right) dx_{\bar{j}} dx_{\bar{j}}$$

= {move integral inside and collect terms} =

$$= \int_{\bar{j}} q_j(x_j) \underbrace{\int_{\bar{j}} \prod_{i=1, i \neq j}^n q_i(x_i) \cdot \log p(y, x) dx_{\bar{j}} dx_{\bar{j}}}_{\log f_j(x_j)} -$$

$$- \int_{\bar{j}} q_j(x_j) \underbrace{\int_{\bar{j}} \prod_{i=1, i \neq j}^n q_i(x_i) \left(\sum_{k \neq j} \log q_k(x_k) + \log q_j(x_j) \right) dx_{\bar{j}} dx_{\bar{j}}}_{\log q_j(x_j) \cdot \underbrace{\int_{\bar{j}} \prod_{i=1, i \neq j}^n q_i(x_i) dx_{\bar{j}}}_{1}} =$$

$$\log q_j(x_j) \cdot \underbrace{\int_{\bar{j}} \prod_{i=1, i \neq j}^n q_i(x_i) dx_{\bar{j}}}_{1} +$$

$$+ \underbrace{\int_{\bar{j}} \prod_{i=1}^n q_i(x_i) \sum_{k \neq j} \log q_k(x_k) dx_{\bar{j}}}_{\text{const. when we fit } q_j}$$

$$= \int_{\bar{j}} q_j(x_j) \log f_j(x_j) dx_{\bar{j}} - \int_{\bar{j}} q_j(x_j) \cdot \log q_j(x_j) dx_{\bar{j}} + \text{const.} =$$

$$= \int_{\bar{j}} q_j(x_j) \cdot \log \frac{f_j(x_j)}{q_j(x_j)} dx_{\bar{j}} + \text{const.}$$

(5)

$$= - \int q_j(x_j) \cdot \log \frac{q_j(x_j)}{f_j(x_j)} dx_j + \text{const.}$$

$KL(q_j(x_j) \parallel f_j(x_j))$

$$\Rightarrow -KL(q_j(x_j) \parallel f_j(x_j)) + C$$

\Rightarrow KL - divergence between the factor $q_j(x_j)$ and the distribution when all other factors have been averaged out.

\Rightarrow maximise $L(q_j(x_j))$ by minimising $KL(q_j(x_j) \parallel f_j(x_j))$

- KL - divergence always positive

\Rightarrow minimum = \emptyset

$$KL(q_j(x_j) \parallel f_j(x_j)) = 0$$

$$\Rightarrow q_j(x_j) = f_j(x_j)$$

$$\log f_j(x_j) = \int \underbrace{\prod_{i \neq j} q_i(x_i)}_{q_j(x_j)} \log p(x, y) dx_{-j} =$$

$$= \mathbb{E}_{q_j(x_j)} [\log p(x, y)]$$

- We need to pick q so that we can compute the above expectation

6

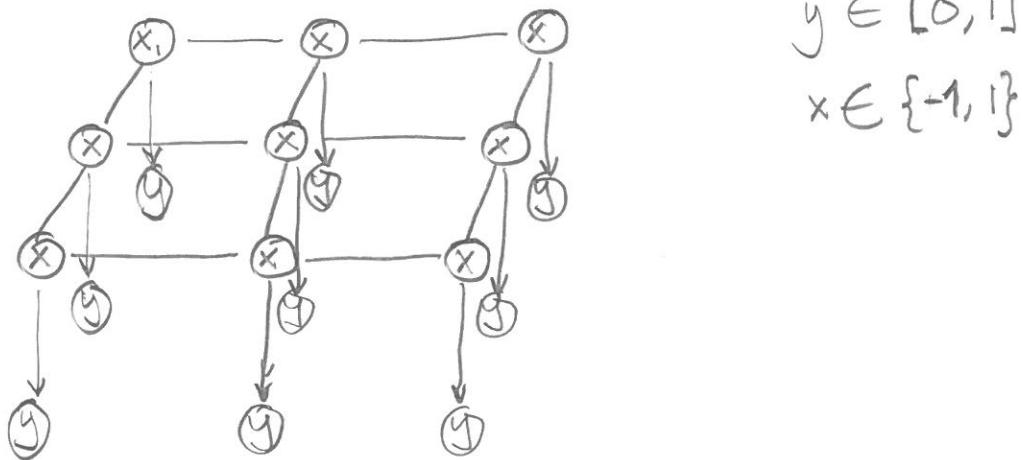
$$q(x_j) = \frac{1}{Z} \cdot \exp\left(\mathbb{E}_{q_{\bar{j}}(x_j)}[\log p(x, y)]\right)$$

as we know that $q(x_j)$ has to be a distribution we simply work in log-space

$$\Rightarrow \log(q(x_j)) = \mathbb{E}_{q_{\bar{j}}(x_j)}[\log p(x, y)] + \text{const.}$$

Example: Markow Random Field

Binary Image



x - is the binary image

y - is the noisy observations

Joint distribution: $p(x, y) = p(y|x)p(x)$

Prior: $p(x) = \frac{1}{Z_0} \cdot e^{-E_0(x)}$

$$E_0(x) = - \sum_{\tau}^D \sum_{j \in N(\tau)} w_{\tau j} x_{\tau} x_j$$

Likelihood: $p(y|x) = \prod_i p(y_i|x_i) = \sum_i e^{-L_i(x_i)}$

- Find approximate fully factorised posterior over the latent representation.

$$q(x) = \prod_i q(x_i, \mu_i)$$

\uparrow variational parameter

Write down joint distribution

(8)

$$\log p(x, y) = \log P(y|x)p(x) =$$

$$= \log \left(\prod_i e^{-L_i(x_i)} \cdot \frac{1}{Z_0} e^{\sum_j w_{ij} x_i x_j} \right) =$$

$$= \sum_i L_i(x_i) + \sum_i \sum_{j \in N(i)} w_{ij} x_i x_j + \text{const.} =$$

= $\begin{cases} \text{If we are only seeing each term} \\ \text{In turn we can take the non } x_i \\ \text{terms into the const} \end{cases} \quad] =$

$$= L_i(x_i) + x_i \sum_{j \in N(i)} w_{ij} x_j + \text{const.}$$

Lets compute the expectation to get the update.

$$\log q_i(x_i) = \log f_i(x_i) = \int \prod_{j \neq i} q_j(x_j) \cdot \log p(x, D) dx_{\neq i} =$$

$$= \left\{ \mathbb{E}_{q_j}[x_j] = \mu_j \right\} =$$

$$= \int_{j \neq i} q_j(x_j) \left(L_i(x_i) + x_i \sum_{j \in N(i)} w_{ij} x_j + \text{const.} \right) dx_{\neq i} =$$

$$= \left\{ \mathbb{E}[ax] = a \cdot \mathbb{E}[x] \right\} = L_i(x_i) + x_i \cdot \underbrace{\sum_{j \in N(i)} w_{ij} \mu_j}_{m_i} + \text{const.}$$

$$\Rightarrow q_i(x_i) \propto e^{x_i \cdot m_i + L_i(x_i)}$$

⑨

$q_i(x_i)$ is a distribution which means it needs to sum to 1. However it is a really simple task as x_i can only take two values -1 or 1.

$$\begin{aligned}
 q_i(x_i=1) &= \frac{1}{q(x_i=1) + q(x_i=-1)} \cdot q(x_i=1) = \\
 &= \frac{e^{m_i + l_i(1)}}{e^{m_i + l_i(1)} + e^{-m_i + l_i(-1)}} = \\
 &= \left\{ \frac{e^a}{e^a + e^b} = \frac{1}{e^a(e^a + e^b)} = \frac{1}{1 + e^{b-a}} \right\} = \\
 &= \frac{1}{1 + e^{-2m_i - l_i(1) + l_i(-1)}} = \frac{1}{1 + e^{-2(m_i + \frac{1}{2}(l_i(1) - l_i(-1))}}} = \\
 &= \frac{1}{1 + e^{-2a_i}} = \text{Sigm}(2a_i)
 \end{aligned}$$

$$\Rightarrow q_i(x_i=-1) = \text{Sigm}(-2a_i)$$

A sigmoid posterior makes perfect sense for a classification task.

⑩

We want to update the variational parameter of each proposal distribution.

$$\begin{aligned}
 \mu_i &= E_{q_i(x_i)}[x_i] = (+1) \cdot q_i(x_i=1) + (-1) \cdot q_i(x_i=-1) = \\
 &= \frac{1}{1+e^{-2\alpha_i}} - \frac{1}{1+e^{2\alpha_i}} = \frac{e^{\alpha_i}}{e^{\alpha_i} + e^{-\alpha_i}} - \frac{e^{-\alpha_i}}{e^{-\alpha_i} + e^{\alpha_i}} = \frac{e^{\alpha_i} - e^{-\alpha_i}}{e^{\alpha_i} + e^{-\alpha_i}} \\
 &= \tanh(\alpha_i) = \tanh\left(m_i + \frac{1}{2}(L_i(1) - L_i(-1))\right) = \\
 &= \tanh\left(\sum_{j \in \text{EN}(i)} w_{ij} \mu_j + \frac{1}{2}(L_i(1) - L_i(-1))\right)
 \end{aligned}$$

Iteratively update these in turn

$$\mu_i^t = \tanh\left(\sum_{j \in \text{EN}(i)} w_{ij} \mu_j^{t-1} + \frac{1}{2}(L_i(1) - L_i(-1))\right)$$

Or a clamped update

$$\mu_i^t = (1-\lambda)\mu_i^{t-1} + \lambda \cdot \tanh\left(\sum_{j \in \text{EN}(i)} w_{ij} \mu_j^{t-1} + \frac{1}{2}(L_i(1) - L_i(-1))\right)$$