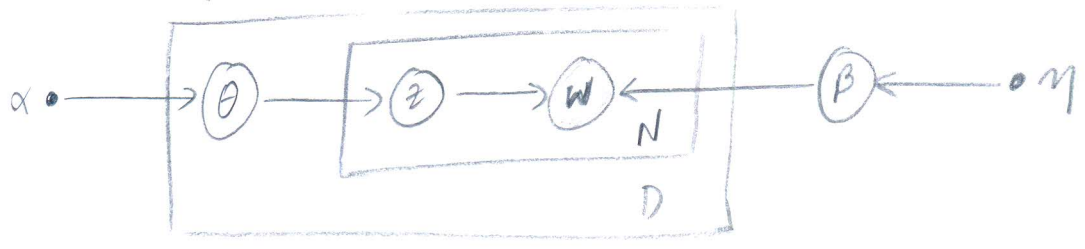


11 review - (11b) - Approximate Inference and Topic Models

- topic models
- variational inference (see 5.2020 lectures to suppl.)

(*) generalised topic model - joint likelihood (topic models / LDA covered in exercises)



- interested in posterior; data likelihood

(*) inference and learning intractable

$p(\theta_n | D) = ?$ $p(z_{n,m} | D) = ?$ (possible queries)

$p(\theta_n | D) = \frac{p(\theta_n, D)}{p(D)}$ (typo?) - 5.2020 $\rightarrow \frac{p(z_{n,m} | D)}{p(D)}$

$$= \sum_{\{z_{n,m}\}} \int \left(\prod_n \left(\prod_m p(w_{n,m} | \beta_{z_n}) p(z_{n,m} | \theta_n) \right) p(\theta_n | \alpha) \right) p(\beta | \eta) d\theta_n d\beta$$

$p(D)$

$$p(D) = \sum_{\{z_{n,m}\}} \int \dots \int \left(\prod_n \left(\prod_m p(x_{n,m} | \beta_{z_n}) p(z_{n,m} | \theta_n) \right) p(\theta_n | \alpha) \right) p(\beta | \eta) d\theta_1 \dots d\theta_n d\beta$$

(*) super exponentially complex. (compute no. terms yourself).

(*) param estimation \rightarrow intractable as inference (intractable) cannot be called as a subroutine

(*) V.I.

- introduce auxiliary/proxy distr q to $p(z)$ (prior in EX slides)
- originally got confused \rightarrow (q as proxy for true posterior $p(z|x)$)
- But it's the same $p(z|x) = \frac{p(x,z)}{p(x)}$ \leftarrow see that approx. this is indirectly approx. post.

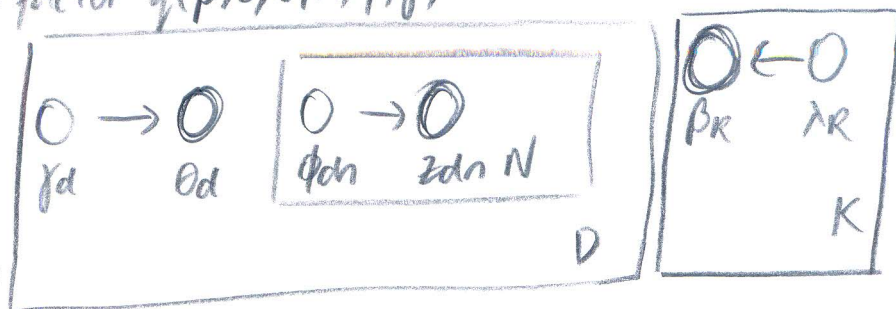
(*) MF Approximation

(*) parametric form for each marginal factor $q(\beta, z, \theta | \lambda, \phi, \gamma)$

$$q(\beta_R | \lambda_R) = \text{Dirich}(\beta_R | \lambda_R)$$

$$q(\theta_d | \gamma_d) = \text{Dirich}(\theta_d | \gamma_d)$$

$$q(z_{dn} | \phi_{dn}) = \text{Multi}(z_{dn} | \phi_{dn})$$



(*) Estimating param. of variational distri (E-step)

2. Aside: - each 'marginal factor'

$q(\beta_R), q(\theta_d), q(z_{dn})$ is itself

parametrised by variational params
(only present in $q(\cdot)$)

v.i. \longrightarrow i) find simple $q(\cdot)$ family
approx/proxy (fully factored)

ii) define parameters within this family.

MF Approx

$$q(\beta, \theta, z | \lambda, \gamma, \phi)$$

$$= \prod_{R=1}^K q(\beta_R | \lambda_R) \prod_{d=1}^D q(\theta_d | \gamma_d) \prod_{n=1}^N q(z_{dn} | \phi_{dn})$$

$$(*) \gamma^*, \lambda^*, \phi^* = \underset{\gamma, \lambda, \phi}{\text{argmin}} KL(q(\beta, \theta, z | \gamma, \phi) || p(\beta, \theta, z | \omega, \alpha, \eta))$$

(*) LDA: can compute in closed form.

(*) some algebraic hard-work now.

(*) update each marginal.

$$\text{update} - q(\theta_d) \propto \exp \left\{ \underset{(i)}{\mathbb{E} \pi_n} q(z_{dn}) \left[\underset{(ii)}{\log p(\theta_d | \alpha)} + \sum_n \log p(z_{dn} | \theta_d) \right] \right\}$$

(i) Dirichlet topic prop. distri.

$$p(\theta_d | \alpha) \propto \exp \left\{ \sum_{R=1}^K (\alpha_R - 1) \log \theta_{dR} \right\}$$

(ii) Multinomial word-label frequency distri

$$p(z_{dn} | \theta_d) \propto \exp \left\{ \sum_{R=1}^K \mathbb{1}(z_{dn} = R) \log \theta_{dR} \right\}$$

(*) To obtain:-

$$q(\theta_d) \propto \exp \left\{ \sum_{k=1}^K \left(\sum_{n=1}^N q(z_{dn}=k) + \alpha_k - 1 \right) \log \theta_{dk} \right\} \leftarrow \text{O/S!} \quad \text{- get this result.}$$

- relies on Dirichlet (conjugate) prior

(*) update each msg.

- similar to $q(\theta_d | y_d)$; obtain optimal param ϕ_{dn}^* for $q(z_{dn} | \phi_{dn})$

$$q(z_{dn}=k | \phi_{dn}) = \phi_{dn}(k) = \beta_k(w_{dn}) \exp \left\{ \psi(y_d(k)) - \psi\left(\sum_{j=1}^K y_d(j)\right) \right\}$$

- Represented multinomial

- And optimal parameters λ_k^* for $q(\beta_k | \lambda_k)$

$$\lambda_k(j) = \eta(j) + \sum_{d=1}^D \sum_{n=1}^{N_d} \phi_{dn}^* \mathbb{1}\{w_{dn}=j\}$$

(*) Iterating to convergence yields MF approx. to posterior
(as subroutine?)

- $\lambda_k(j)$ - frequency count of j^{th} word frequency in topic k ; which is clear from summation over y doc. at every word and delta fn.

(*) Coordinate ascent for LDA

- Have 3 equations for E-step (inference subroutine)

ex: key message:-

$$p(\text{""}) \quad q(\text{""}) = \prod q_i(y_i) \quad \eta^* = \argmin KL(q || p)$$

(*) Variational est. / learning (print S.2020 add. slides as suppl.)

- maximise variational lower bound:-

$$\mathcal{L}(\theta, \phi, x) = \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x|z) + KL(q_\phi(z|x) || p(z))]$$

$$= \log p(z) - KL(q_\phi(z|x) \| p_\theta(z|x))$$

- previous: =

- variational E-step

- maximise \mathcal{L} w.r.t ϕ with θ fixed

$$\max_{\phi} \mathcal{L}(\theta, \phi; x)$$

- keep original model params θ fixed;

optimise ϕ
(variational param.)

- variat. est: -

$$\gamma^*, \lambda^*, \phi^* = \underset{\gamma, \lambda, \phi}{\operatorname{argmin}} KL(q(\beta, \theta, z | \gamma, \phi) \| p(\beta, \theta, z | w, \alpha, \eta))$$

variational params

(*) Additional params in M-step: We want to know θ and β
(not present in V.E.) - why?

- computationally not significant (as Bayesian?) (~)

- θ and β are integrated out

$$\alpha \rightarrow \theta \rightarrow z \quad \begin{array}{l} \text{write/specify} \\ p(z|\alpha) \text{ directly} \\ \text{Dir}(n+\alpha) \end{array}$$

↑
integrated out

- Hence, no need for θ .

- this is why M-step often ignored

- M-step: estimate θ with ϕ fixed

$$p(\theta|z, x)$$

- below is doable;

- ppl. don't bother

- fix variational param.; combine counts.

$$\text{MAP}_{\theta, \beta} \mathcal{L}(\theta, \phi; x)$$

- or even hyperparam

$$\max_{x, \eta} \mathcal{L}(\theta, \phi; x)$$

(*) E-step already provides

ϕ_{dn} γ_{dn} λ_R
 \downarrow \downarrow \downarrow
 z_{dn} θ_d β_R

(sufficient statistics for parameters below)

via exp family;
GLIMS

(?) 0/52
 - clarity

- topical context avg.

(*) e.g. $\lambda_R(j) = \eta(j) + \sum_{d=1}^D \sum_{n=1}^{N_d} \phi_{dn}^*(k) \mathbb{I}(w_{dn}=j)$

prior pseudocount of that word over any topic

indic. weighted by n^{th} word being in topic k

indicator that n^{th} word in document d is " j ".

(*) Estimating $\phi^*, \gamma^*, \lambda^* \rightarrow$ variational E-step; gives you 'everything else' as they're sufficient statistics

(*) Evaluation of gold truth in topic models

- difficult to run on 'synthetic' gold truth
- what is 'ground truth' of a document topics?

(*) V.I. - v. important

- MF - approx - v. imp. - can 'solve' arbitrarily complex models (vague)
- amounts to removal of edges
 - KL between true and proxy

(*) (0/53): see Mohamed slides for more on WPA, VI. (2016) NIPS.

