

11 - Approximate Inference and Topic Models

(Mean Field and Loopy BP) S.2019 11b

- (*) Looks like lecture 11a. VI and Loopy B.P. not covered
- (*) Approximate inference methods started as 'tricks', ones which worked then proved.

Ex: maintain 'art of modelling' trajectory, present an example that exposes the concrete need for approx. inference

(*) Probabilistic Topic Models

- started as a class project.

(*) How to get started for a new modelling task?

(*) There is gold here: a way of thinking about the art of modelling

Ex: start with a concrete task/problem you want to solve.

- methods are invoked in service of the task!

E.g. Bird's eye view of 1 million documents.

- task \rightarrow clustering; give back cluster label.

\hookrightarrow embedding \rightarrow visualisation of cluster labels

- Representation of data \rightarrow e.g. continuous, binary, counts? (design choice)

(*) Inside each element are at a time.

(*) tasks - document embedding.

- Have each document embedded in a space.

(*) Summarising data using topics

- want the embedding to be meaningful.

(*) see how data changes over time

- Representation of topics may evolve.

(*) user interest topic modelling

- secondary task

(*) representation

- Bag of words rep.
- count article \rightarrow BOW.

(*) For each document; count no. of words, order does not matter.

(*) Each document is a vector in word-space

ex: Benefit and costs? (of representation)

- benefits \rightarrow storage, hash tables (makes data simple)
- comparing documents of different lengths (1 x 100 word vs novel)
- probability based on word-orderings \rightarrow longer docs give smaller likelihood (due to product)
 - cannot compare.
- allows comparability.
- costs \rightarrow ordering may be important for semantics

(*) BOW is a baseline representation.

(*) How to model semantics?

- without using orderings

ex: mental analogy \rightarrow certain key words give a higher prob. of a document coming from a topic.

- match a topic with keywords @ intuition.

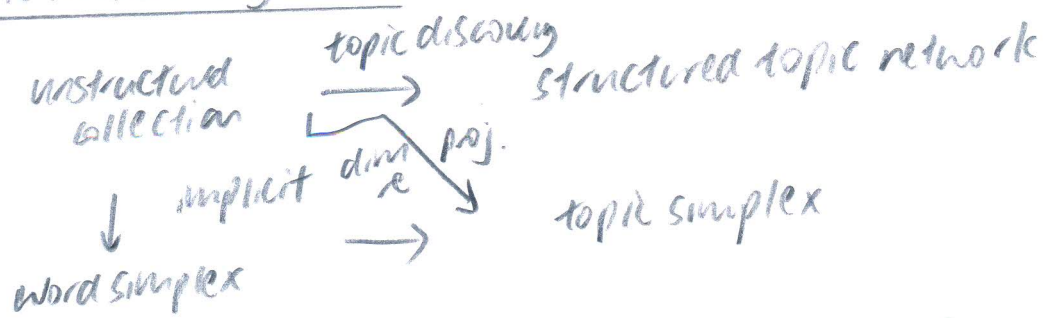
(*) A topic is a vector of words (mixture of vocab.)

A document contains topics in a p (mixing proportion)

(*) ~~How~~ Note info. compression: a document gets compressed into a weighted sum of topics (low-dimensional with some semantic meaning)

(*) Can the compare similarity of documents \rightarrow actually a vector of probabilities indicating topic mixing proper.

(*) Topic Models - Big Picture



(*) A topic corresponds to a point on the word simplex.

(*) A document — " — on the topic simplex

(*) LSI vs Topic Model (prob LSI)

- likelihood model

$$\begin{matrix} \text{words} & \text{docs} \\ \boxed{X} \end{matrix} = \begin{matrix} \text{topics} \\ \boxed{W} \end{matrix} \begin{matrix} \text{words} & \text{topic} \\ \boxed{A} \end{matrix} \begin{matrix} \text{documents} \\ \boxed{\theta^T} \end{matrix} \begin{matrix} \text{topic} \end{matrix}$$

$$\begin{matrix} \text{words} & \text{docs} \\ \boxed{p(w)} \end{matrix} = \begin{matrix} \text{words} & \text{topics} \\ \boxed{p(w|z)} \end{matrix} \begin{matrix} \text{documents} \\ \boxed{p(z)} \end{matrix}$$

- LSI \rightarrow word-document matrix decomposition (linear algebraic)

- Topic Models \rightarrow conceptually similar (probabilistic inference)

- matrix decomposition techniques (e.g. LSI) make it difficult to bypass computationally expensive matrix inversions (especially for log-matrices) (batch)

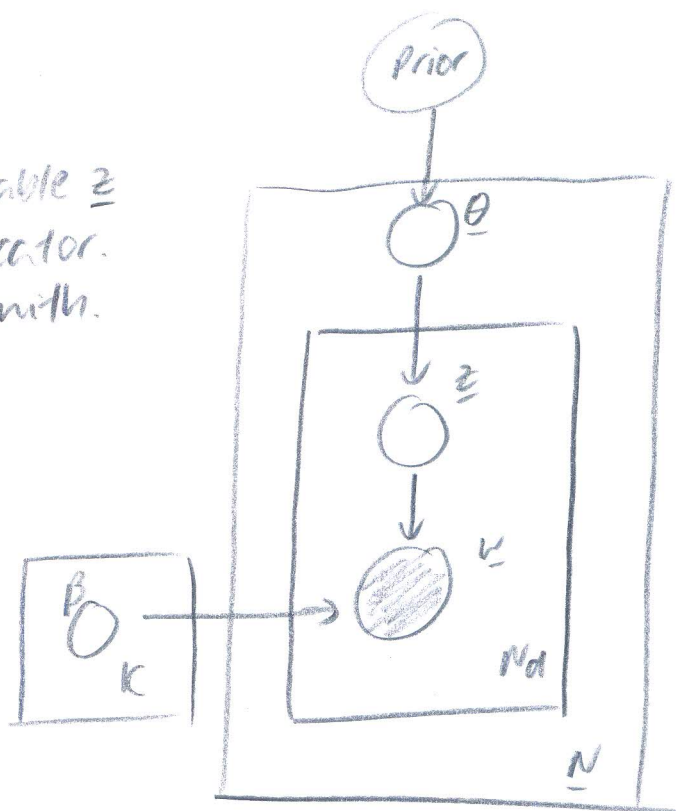
(*) Probabilistic methods allow ways of dealing with this issue in a piece-wise iterative manner (e.g. LASSO vs batch LR)
i.e. local operations to

(*) Admixture Models

- skip

(*) Topic Models

- generative process for a document.
- for every word there is a latent variable z indicating its topic assignment/indicator.
- i.e. what topic a word is affiliated with.
- topic indicator z comes from a weight vector θ
- every document is a vector of topical weights
- prior: distn of different topics in corpus.



generating:

- 1) draw a topical weight from a prior distribution vector/document.

Given this, n th out of
for every, N_d words:

- draw a topic indicator $z_n \sim \text{multinomial}(\theta)$
- conditioning on the topic indicator of the n th word z_n ;
sample the word (using a coll. of word-frequency distns.)
- draw $w_n | z_n, \{\beta_k\}$ from $\text{multinom}(\beta_{z_n})$

(v): presentation is not the clearest

(*) Choices of priors

- Dirichlet (LDA), Blei et al. 2003
- Logistic Normal, Blei & Lafferty (2005), Ahmed & Xing (2006).
- (see the facets of these w/ modelling intuition)

(*) Generative semantics of LONTAM

- captures intuition of topics being highly correlated
e.g. sports and topics
- use a covariance matrix?
- only in MVGS (at simplest)

- OK, so use MVGS \rightarrow to yield χ

ISSUE: we cannot use χ to sample \mathbf{z} $\rightarrow \chi$ is Gaussian vector with -ve components; NOT Multinomial

- Apply a transformation, apply exponential, then normalise to get appropriate prior.

$\chi \rightarrow \theta$ (logistic normal)

EX: Model design can be arbitrary, flexible, to a degree.

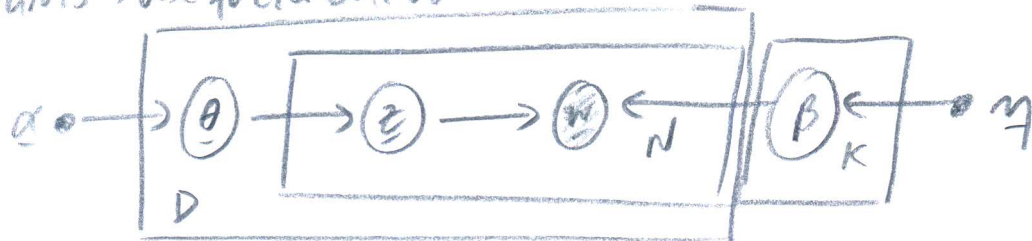
(*) Posterior inference results

- Familiar \rightarrow recap (A1) - review.
- See blowed slides

(*) Joint likelihood

(A2) - check dim.

- PLMS \rightarrow use factorization law!



$$p(\beta, \theta, \mathbf{z}, \mathbf{w}) = \prod_{k=1}^K p(\beta_k | \eta) \prod_{d=1}^D p(\theta_d | \alpha) \prod_{n=1}^N p(z_n | \theta_d) p(w_n | z_n, \beta)$$

(*) we want posterior of any of latent variables given observed words \mathbf{w} and likelihood of the word.

(*) Inference and learning both intractable.

(P3) - review the intractability to get a real sense of how powerful approx inference.

e.g. $p(\theta_n | D)$ and $p(D)$

- no known technique for performing these techniques exactly.

(*) Approximate Inference

i) variational inference

- turns solution of an inference problem \rightarrow sol. of an optimization problem.

ii) MCMC

(*) Variational Inference (VI)

(P4) - review logic

- Consider generative model $p_\theta(x|z)$, prior $p(z)$

- Joint distri: $p_\theta(x, z) = p_\theta(x|z)p(z)$

- Assume variational distri $q_\phi(z|x)$

- objective: maximise lower bound for log-likelihood: -

$$\log p(x) = \text{KL}(q_\phi(z|x) \| p_\theta(z|x)) + \int_z q_\phi(z|x) \log \frac{p_\theta(x, z)}{q_\phi(z|x)} dz \quad ?$$

$$\geq \int_z q_\phi(z|x) \log \frac{p_\theta(x, z)}{q_\phi(z|x)}$$

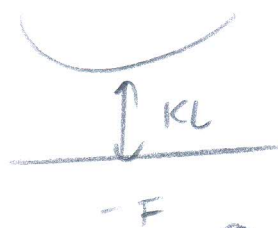
$$:= \mathcal{L}(\theta; \phi; x)$$

equivalently; minimise free-energy (upper bound on log-like.)
(surrogate or target)

$$\mathcal{F}(\theta; \phi; x) = -\log p(x) + \text{KL}(q_\phi(z|x) \| p_\theta(z|x))$$

(*) distance between free energy and ...
(-F)

$$(*) q_{\phi}(z|x) = p_{\theta}(z|x) \Rightarrow KL(q||p) = 0$$



(*) Variational Inference

(AS): review
↓

(*) Maximize variational lower bound:-

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

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

(W): A little abstract.

(*) E-step: Maximize \mathcal{L} wrt ϕ with θ fixed:-

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

If closed form sol. exist:-

$$q_{\phi}^*(z|x) \propto \exp[\log p_{\theta}(x, z)] \quad (\text{do not set this to } p(z|x) \text{ mode v. t.})$$

(*) M-step: Maximize \mathcal{L} wrt θ ; with ϕ fixed.

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

ex use $q_{\phi}(z|x)$ i.e. an inference step on z (latent) given data x

- make sure q_{ϕ} is 'good' in sense of being 'close' to $p_{\theta}(z|x)$

using KL-divergence as a measure of closeness

- Why not just set $q_{\theta}(z|x) = p_{\theta}(z|x)$? $\rightarrow p_{\theta}$ intractable

(*) mean-field assumption (in topic models)

$$\text{- True posterior: } p(\beta, \theta, z | w) = \frac{p(\beta, \theta, z, w)}{p(w)}$$

(*) Break dependency using fully factorised distri:-

$$q(\beta, \theta, z) = \prod_k q(\beta_k) \prod_d q(\theta_d) \prod_n q(z_{dn})$$

- Product of non-coupled, individual distri (marginals)

Ex: Why does this make it things easier?

W:- Can optimise individually by breaking into subproblems

(AB): Review - not entirely sure of
i) logic
ii) dimensionality

(*) Mean-Field Approx

(AB): Pick up during review. (Ex gone through quickly)

- Read/glance/skim papers

(*) Co-ordinate ascent algorithm for LDA

- Get on iterative program.

- (AB): Review pseudocode.

Ex: A lot of material to digest

- next lectures \rightarrow more examples of approx. inference

- Give a 'grand theory' to unify/better understood.