

VI - variational inference - loopy belief prop.(\*) continue VI - derivation of KF for SSMSKalman filtering derivation

- recall filtering problem

- inference of  $p(x_t | y_{1:t})$  (i.e. cond. dist. of last hidden state given all observations up to now)

- use diagram as reference point for visualizing queries.

- we require:-

$p(x_t | y_{1:t})$  from a recursive procedure involving  $p(x_{t+1} | y_{1:t})$ ,  $y_{t+1}$

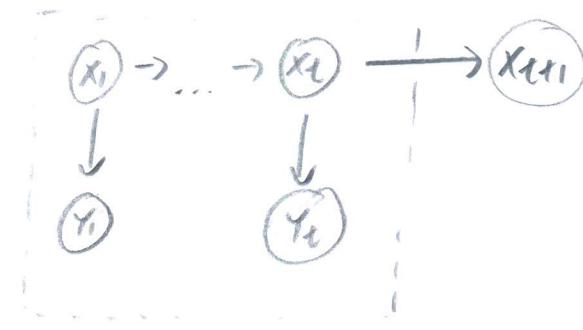
- break up problem into 2 steps:-

- 1) time update:

- compute  $p(x_{t+1} | y_{1:t})$

- from prior belief  $p(x_t | y_{1:t})$  and  
dyn. model  $p(x_{t+1} | x_t)$

(\*) solidify analogy/mental model  
to crystallise  $\rightarrow$  true understanding.



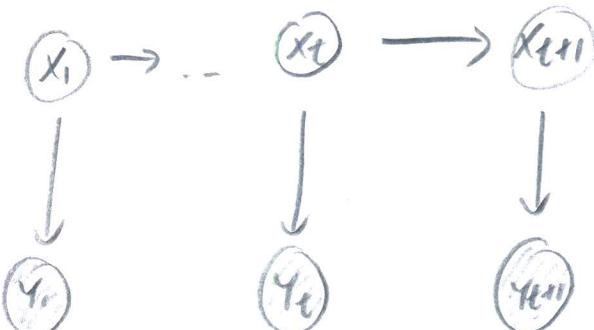
- 2) measurement update:

- compute new belief  $p(x_{t+1} | y_{1:t+1})$

- from prediction  $p(x_{t+1} | y_{1:t})$ ,

- observation  $y_{t+1}$

- obs. model  $p(y_{t+1} | x_{t+1})$



(\*) Make use of Gaussian properties everywhere

- inference problem  $\rightarrow$  dealing with <sup>new</sup> means and covariances

(\*) same technique as in FA.

(\*) Schematic of the strategy:-  
(gaussian manip.)

M.I.L.

$$\begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \sim N \left( \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \right) \quad z_1 \rightarrow z_1 z_2 \rightarrow z_1/z_2$$

corresp.

$$z_1 \sim x_t | y_{1:t} \quad x_{t+1} = Ax_t + Gw_t \rightarrow p(x_{t+1}, y_{1:t} | y_{1:t}) \rightarrow p(x_{t+1} | y_{1:t}, y_{1:t})$$

$$= f(x_t)$$

$$p(x_{t+1} | y_{1:t})$$

$$\oplus \mu_1$$

$$p(y_{1:t} | x_{t+1})$$

$$y_{1:t} = Cx_{1:t} + v_{1:t}$$

$$\mu_2 / \mu_1 \oplus$$

(A1)-tidy this up  
(makes some sense; but want  
full clarity)

- ex: conditioning is not scary; just  
some constants in eq.; focus on  
LHS of cond.

(\*) Compute means and covariances within the schematic

(A2) review

(\*) predict step

- asymmetrical model  $x_{t+1} = Ax_t + Gw_t \quad w_t \sim N(0, Q)$

$$\hat{x}_{t+1|t} = E[x_{t+1} | y_1, \dots, y_t] = A\hat{x}_{t|t}$$

$$P_{t+1|t} = E[(x_{t+1} - \hat{x}_{t+1|t})(x_{t+1} - \hat{x}_{t+1|t})^T | y_1, \dots, y_t]$$

$$= E[(Ax_t + Gw_t - \hat{x}_{t+1|t})(Ax_t + Gw_t - \hat{x}_{t+1|t})^T | y_1, \dots, y_t]$$

$$= AP_{t|t}A^T + GPG^T$$

$\hat{x}_{t+1|t}$  and  $P_{t+1|t}$  are conditional means and covariances (\*)

(A3) - notation

$$\text{observation model} \quad y_{t+1} = Cx_t + v_t \quad v_t \sim N(0, R)$$

$$E[y_{t+1}|y_1, \dots, y_t] = E[(x_{t+1} + v_{t+1})|y_1, \dots, y_t] = \hat{x}_{t+1|t}$$

$$E[(y_{t+1} - \hat{x}_{t+1|t})(y_{t+1} - \hat{x}_{t+1|t})^T | y_1, \dots, y_t] = P_{t+1|t} C^T + R$$

### (\*) update step

- From previous slide:-

we have joint  $p(x_{t+1}, y_{t+1}|y_1, \dots, y_t) = N(\underline{m}_{t+1}, \underline{V}_{t+1})$  where:-

$$\underline{m}_{t+1} = \begin{pmatrix} \hat{x}_{t+1|t} \\ \hat{C}_{t+1|t} \end{pmatrix} \quad \underline{V}_{t+1} = \begin{pmatrix} P_{t+1|t} & P_{t+1|t} C^T \\ C P_{t+1|t} & C P_{t+1|t} C^T + R \end{pmatrix}$$

- Now use conditional Gaussian formulae for partitioned matrices to get  $p(x_{t+1}|y_{t+1}, y_{1:t})$

### (\*) Kalman Filter

- Has closed form measured and time updates:-

measurement update :-  $\hat{x}_{t+1|t+1} = \hat{x}_{t+1|t} + K_{t+1}(y_{t+1} - \hat{x}_{t+1|t})$

$$P_{t+1|t+1} = P_{t+1|t} - K_{t+1} P_{t+1|t} \quad ?$$

-  $K_{t+1}$  - Kalman gain matrix

time updates:  $\hat{x}_{t+1|t} = A\hat{x}_{t|t}$

$$P_{t+1|t} = AP_{t|t}A^T + GQG^T$$

$$K_{t+1} = P_{t+1|t} C^T (C P_{t+1|t} C^T + R)^{-1}$$

(\*) Review notation and intuition.

Ex. Be sensitive to conditioning in terms of info being used at each point.

## (\*) Example of KF in 1D

- just apply KF equations

## (\*\*) Intuition behind these

### (\*) KF intuition (\*\*)

- KF update of mean:-  $\hat{x}_{\text{true|t+1}} = \hat{x}_{\text{true|t}} + K_{t+1} (\underline{z}_{t+1} - \hat{x}_{\text{true|t}})$

$$= \frac{(\sigma_t + \sigma_x) \underline{z}_t + \sigma_z \hat{x}_{\text{true|t}}}{\sigma_t + \sigma_x + \sigma_z}$$

- innov.

• new belief is convex combo of updates from prior and observation, weighted by Kalman gain matrix

$$K_{t+1} = P_{\text{true|t}} C^T (C P_{\text{true|t}} C^T + R)^{-1}$$

- observation unreliable,  $\sigma_z$  (i.e.  $R$ ) large so  $K_{t+1}$  is small  $\rightarrow$  more att. to prediction

- old prior is reliable, large  $\sigma_t$ , or process is unpredictable, large  $\sigma_x$ , we pay more attention to observation.

### (\*) Where do A, G, C matrices come from?

- A parameter estimation (as opposed to filtering, which is an inference problem)
  - ↳ i.e. compact p.d. of latent var, given data in this context.
- inference assumes learning / param est. has already occurred.

### (\*) Complexity of one KF step., Inf prob 2, RTS smoother, RTS deriv.

- Review  $\approx$  in OM time

## (\*) learning SSMS

- (\*) use EM. inference from ICF or RTS filtering forms E-step!
- (\*) Form complete log-like.
- (\*) from E-step  $\rightarrow$  compute suffi. statistics using inference results to give you a posterior.

M-step  $\rightarrow$  MLE

Ex: You can and should derive m your spare time.

- (\*) very principle  $\rightarrow$  inference can be treated as a subroutine for learning in a partially obs. setting

## (\*) nonlinear systems

- use derivatives, Taylor exp., approxim., linearization.

- (\*) L.S thinking  $\rightarrow$  you can approximate functions to the order you can afford
- (\*) Ex A lot of seemingly difficult, obtuse algebra can be viewed/reduced through an understanding of principles motivating them (often new).

## UI - Approximate Inference and Topic Models

(Mean Field and Loopy BP) S. 2019 IIb

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

(\*) consider each element one 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.
- just 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 ( $1 \times 100$  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  $\oplus$  intuition.

(\*) A topic is a vector of words (in term of vocab.)

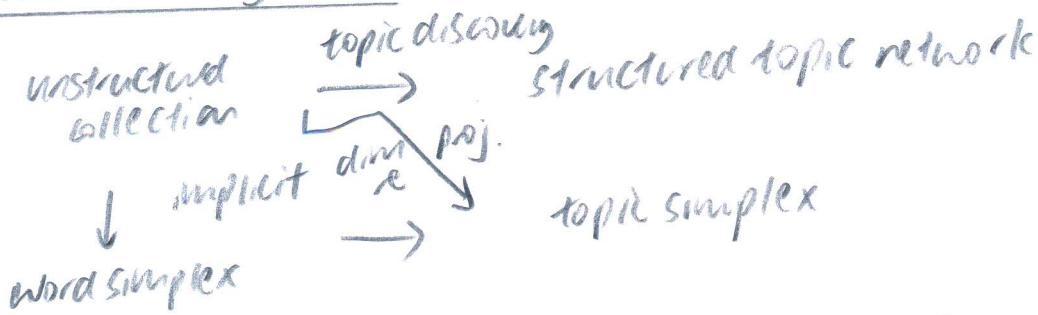
A document contains topics in a p (mixing proportion)

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

(\*) can then compute similarity of documents

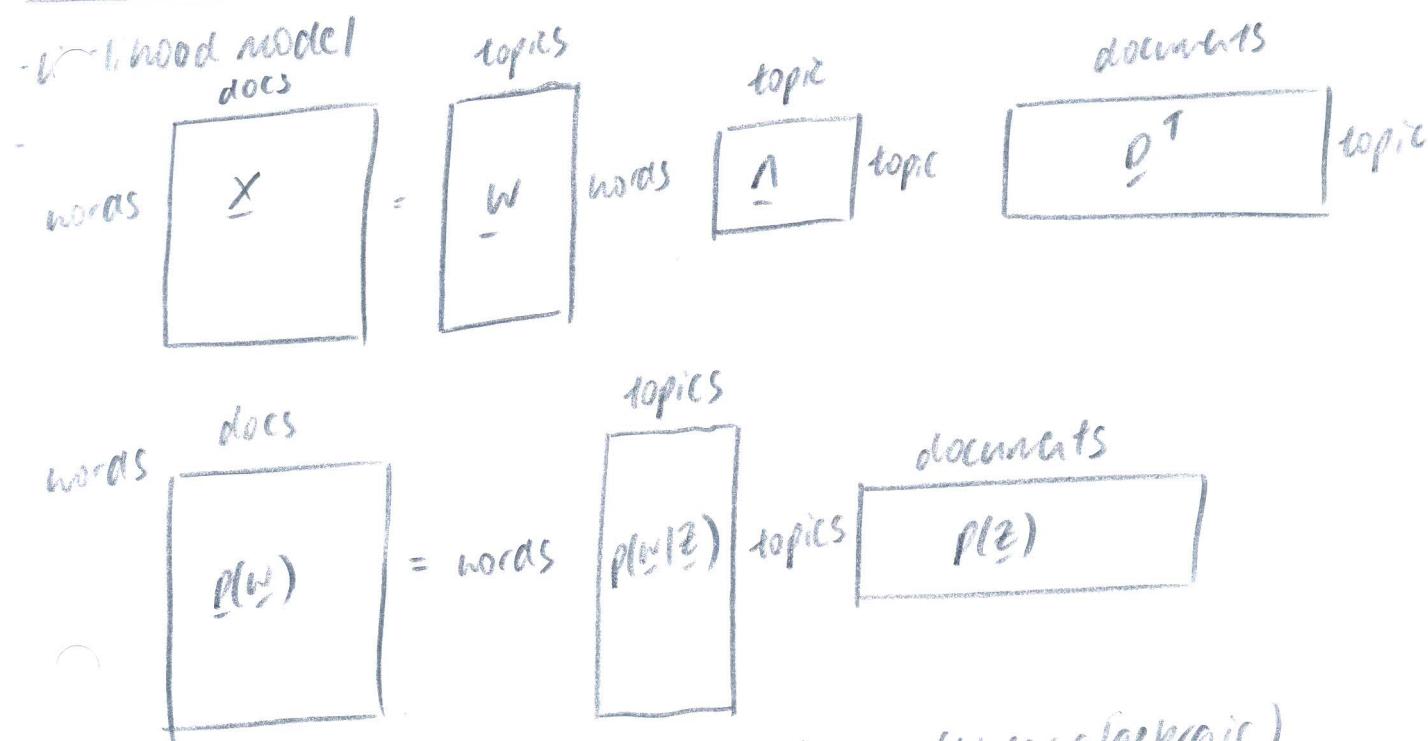
actually a vector of probabilities indicating topic mixing propor.

## (\*) Topic Models - Big Picture



- (\*) A topic corresponds to a point on the word simplex.
- (\*) A document is on the topic simplex

## (\*) LSI vs Topic Model (probLSI)



- 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

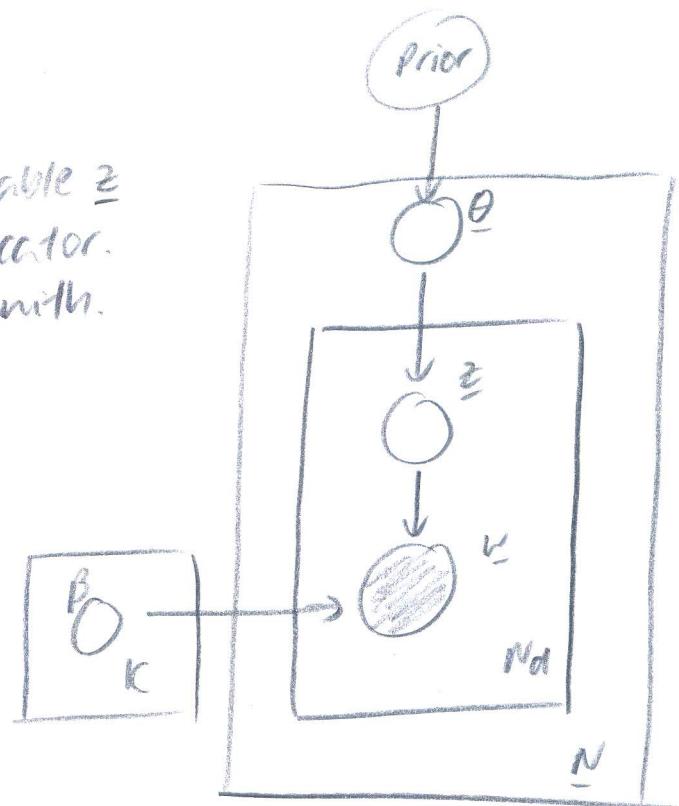
- generative process for a document.

- for every word there is a latent variable  $z$  indicating its topic assignment/indicator.
- i.e. what topic a word is affiliated with.

- topic marker for  $z$  comes from a weight vector  $\theta$

- every document is a vector of topical weights

prior: distri of different topics in corpus.



## generating:

- 1) draw a topical weight from a prior distribution  $\theta$

- given this,  $n$  times

for every  $N_d$  words:

- draw a topic indicator  $z_n \sim \text{multinomial}(\theta)$

- conditioning on the topic indicator of the  $n^{\text{th}}$  word  $z_n$ ;

- sample the word (using a coll. of word-frequency distris.)

- draw  $w|z_n, \beta_{z_n}$  from  $\text{multinom}(\beta_{z_n})$

⑥: 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 LDA

- captures intuition of topics being highly correlated  
e.g. sports and topics

- use a covariance matrix?

- only in MFGS (not simplest)

- DTC, so use MFGS  $\rightarrow$  to yield  $\gamma$

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

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

$$\gamma \rightarrow \theta \text{ (logistic normal)}$$

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

## (\*) Posterior inference results

- Familiar  $\rightarrow$  recap

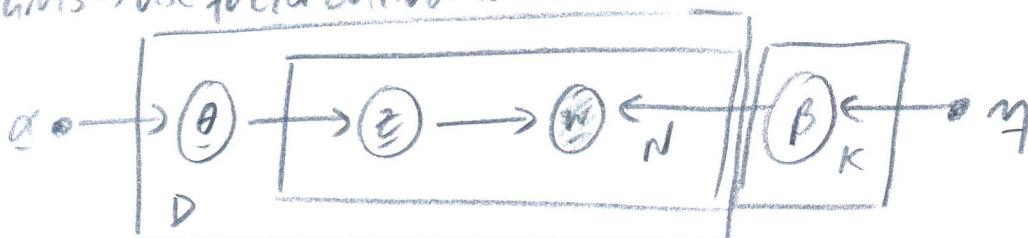
Ⓐ - review.

- Ⓑ see allowed slides

## (\*) Joint likelihood

Ⓑ - check dims.

- PMS  $\rightarrow$  use factorisation law!



$$p(\beta, \theta, z, w) = \prod_{n=1}^N p(\beta_n | \eta) \prod_{d=1}^D p(\theta_{dn} | \alpha) \prod_{n=1}^N p(z_{dn} | \theta_{dn}) p(w_{dn} | z_{dn}, \beta)$$

- (\*) we want posterior of any of latent variables given observed words  $\gamma^w$  and likelihood of the word.

(\*) Inference and learning both intractable.

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

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

- 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 (cont.)

(A4) - review logic

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

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

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

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

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

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

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

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

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

(\*) Distance between free energy and ...  
(-F)

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



### (\*) variational inference

(\*) Maximize variational lower bound:-

$$L(\theta; \phi; x) = E_{q_\phi(z|x)} [\log p_\theta(z|x)] + KL(q_\phi(z|x) || p_\theta(z))$$

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

Ⓐ: review  
↓

Ⓑ: A little abstract.

(\*) C-step: Maximise L wrt  $\phi$  with  $\theta$  fixed:-

$$\max_\phi L(\theta; \phi; x)$$

If closed form sol. exist:-

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

(\*) M-step: Maximise L wrt  $\theta$ ; with  $\phi$  fixed.

$$\max_\theta L(\theta; \phi; x)$$

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

- make sure  $q_\phi$  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)

- 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

(\*\*) Review - not entirely see of i) logic  
ii) dimensionality

## (\*) Mean-Field Approx

(\*) Pick up during review. (Ex gone through quickly)

- Read / glance / skim papers

## (\*) Co-ordinate ascent algorithm for LDA

- Get an iterative program.

- (\*\*) Review pseudocode.

Ex: A lot of material to digest

- next lectures  $\rightarrow$  more examples of approx. inference

- give a 'grand theory' to unify / better understand.