

- today:
- finish variational inference
  - Bayesian
  - model selection & causality

mean field continuation

$$\min_{q \in Q_{MF}} KL(q \parallel p)$$

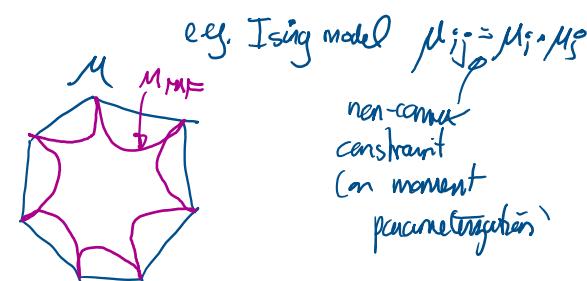
$$\hookrightarrow \{q : q(x) = \prod_i q_i(x_i)\}$$

[see lecture 22 in 2017 for "marginal polytope perspective"]

lecture 22 Fall 2017 link

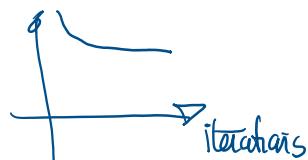
•  $KL(\cdot \parallel p)$  is convex fct. of  $q$

but  $Q_{MF}$  is non-convex constraint set



but can monitor progress

$$KL(q^{(t)} \parallel p) + \text{cst.}$$



pros & cons of variational methods

⊕ optimization based  
→ often faster to run  
→ easier to debug

⊖ biased estimate

$$\mathbb{E}_{q^{(t)}} [f(z)] \neq \mathbb{E}_p [f(z)]$$

Vs sampling

⊖ noisy → harder to debug  
mixing problem for chains

⊕ unbiased estimate

$$\mathbb{E} [\mathbb{E}_{q^{(t)}} [f(z)]] = \mathbb{E}_p [f(z)]$$

with respect to random samples

to make sure chain has mixed

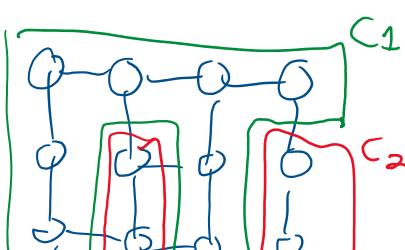
$\frac{1}{N} \sum_{i=1}^N f(z^{(i)})$

$t = 10$  in a chain

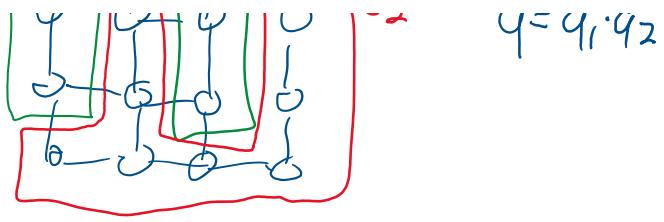
structured mean field

idea  $q(z) = \prod_{j=1}^K q_j(z_j)$  where  $C_1, \dots, C_K$  is a partition of  $V$

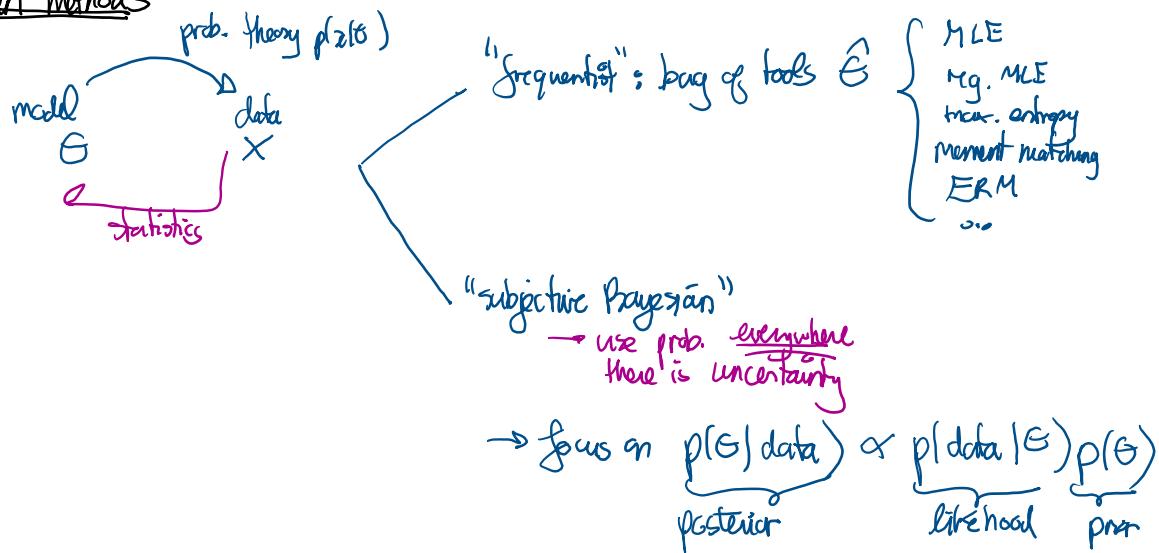
and  $q_j$ 's are tractable distribution  
(for example free UGM for  $q_j$ )



$$q = q_1 \cdot q_2$$



## Bayesian methods



caricature:

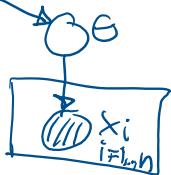
- Bayesian is "optimist"  
they think you can get "good" models
- ⇒ obtain a method by doing inference in model

hyperparameters  
for prior  $\alpha_0, \beta_0$

• frequentist is "pessimist" → use analysis tools

Example: biased coin

$$X_i | \Theta \sim \text{Bernoulli}(\Theta)$$



$$\text{e.g. } \Theta \sim \text{Unif}[\alpha_0] = \text{Beta}(1,1)$$

$$p(\Theta) = \text{Beta}(\Theta | \alpha_0, \beta_0)$$

$$p(x_i | \Theta) = \Theta^{x_i} (1-\Theta)^{1-x_i}$$

$$\begin{aligned} \text{posterior} &\propto \left( \prod_{i=1}^n p(x_i | \Theta) \right) p(\Theta) \\ &= \Theta^{\sum x_i} (1-\Theta)^{n-\sum x_i} \Theta^{\alpha_0-1} (1-\Theta)^{\beta_0-1} \text{Beta}(\Theta | \alpha_0, \beta_0) \end{aligned}$$

$$\Rightarrow p(\Theta | \text{data}) = \text{Beta}(\Theta | \alpha_0 + n_1, \beta_0 + n - n_1)$$

↳ "conjugate prior" to the Bernoulli likelihood model

Conjugate priors:

consider a family  $F$  of dist. on  $G$   $F = \{p(\Theta|\alpha) : \alpha \in A\}$

say that  $F$  is a "conjugate family" to observation model  $p(x|\Theta)$

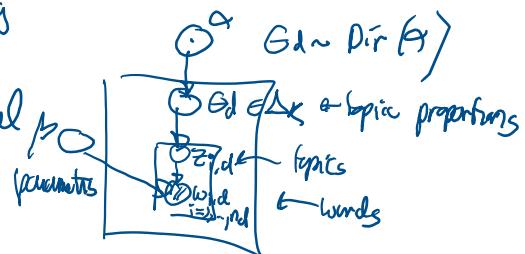
If posterior  $p(\Theta|x, \alpha) \in F$  for any  $x \sim X | \Theta$

If posterior  $p(\theta|x, \alpha) \in F$  for any  $x \sim X | \epsilon$   
 i.e.  $\exists$  some  $\alpha'(x, \alpha)$  s.t.  $p(\epsilon|x, \alpha) = p(\epsilon|\alpha')$

Side note: if use conjugate priors in a DGM  
 then Gibbs sampling might be easy

e.g. LDA topic model  $p(\epsilon|z, \alpha)$

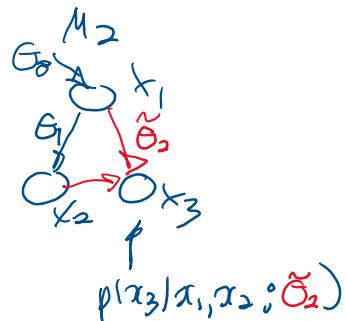
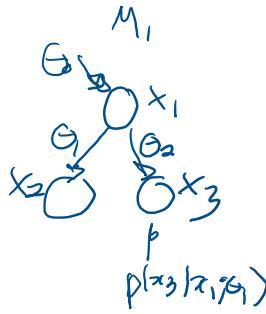
example hwk 1  
 Dirichlet prior  
 is conjugate to  
 multinomial likelihood  
 model



15h34

### Model selection

say we want to choose between 2 DGM



[note here that " $M_1 \subseteq M_2$ "]

"corner relation"

as a frequentist  $\hat{\Theta}_{M_1}^{MLE} = \underset{G_0, G_1, G_2}{\text{argmax}} \log p(\text{data} | G_0, G_1, G_2, \text{"model } M_1")$

$\hat{\Theta}_{M_2}^{MLE} = \underset{G_0, G_1, G_2}{\text{argmax}} \log p(\text{data} | G_0, G_1, G_2, \text{"model } M_2")$   
 ↩ different space than  $\Theta_2$

how to choose between models?

can't compare  $\log p(\text{data} | \hat{\Theta}_{M_1}^{MLE}, M=M_1)$  vs  $\log p(\text{data} | \hat{\Theta}_{M_2}^{MLE}, M=M_2)$

because LHS  $\leq$  RHS since  $M_1 \subseteq M_2$   
 (i.e. you would always choose "bigger model")

→ as a frequentist, use cross-validation  
 or validation set

i.e.  $\log p(\text{test data} | \hat{\Theta}_{M_i}^{MLE}(\text{train data}), M_i)$

### Bayesian alternatives

↑ true Bayesian  $\Rightarrow$  sum over models (integrate out uncertainty about  $M$ )

True Bayesian  $\Rightarrow$  sum over models (integrate out uncertainty about  $M$ )

introduce a prior over models  $p(M)$

$$p(x_{\text{new}} | D) = \sum_M p(x_{\text{new}} | D, M) p(M | D)$$

standard Bayesian predictive dist.  
for one model

posterior over models  
 given data  $D$  & model  $M$

$$\sum_M \left[ \left( \sum_{E_n \in \mathcal{E}_M} p(x_{\text{new}} | E_n, M) p(E_n | D, M) dE_n \right) p(M | D) \right]$$

$\otimes$  in model selection, forced to pick model

$\Rightarrow$  pick model that maximizes  $p(M | \text{data}) \propto p(\text{data} | M) p(M)$

$p(\text{data} | M) = \text{"marginal likelihood"}$

$$\sum_{E_n} p(\text{data} | E_n, M) p(E_n | M) dE_n$$

to compare two models, look at

$$\frac{p(M=M_1 | D)}{p(M=M_2 | D)} = \frac{p(D | M_1)}{p(D | M_2)} \frac{p(M_1)}{p(M_2)}$$

Bayes factor      prior ratio

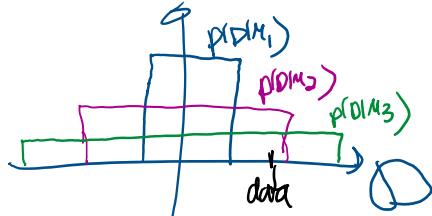
"uniform prior over models"  $\Rightarrow$  then pick among  $K$  models  $M_1, \dots, M_K$

by max  $p(\text{data} | M=M_i)$

"empirical Bayes"      "type II ML"

when # of models is "small", then this approach is "fine"  
(i.e. won't overfit)

Zoubin's cartoon: suppose  $M_1 \subseteq M_2 \subseteq M_3$



$p(D | M)$  is normalized over  $D$

$p(D | \hat{M}_{\text{MLE}}(D), M)$       [can overfit badly]

but type II ML can still overfit if have too many models

$$\text{say, e.g. } p(D|M) = S(D, M)$$



how to compute marginal likelihood:

use approximations ↘ variational inference  
sampling

simple approximation → Bayesian information criterion (BIC)

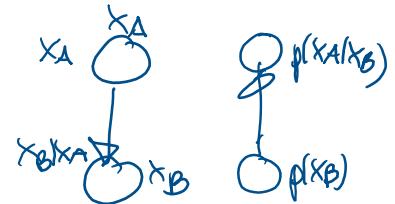
Causality:

structural causal model: graph model + intervention model

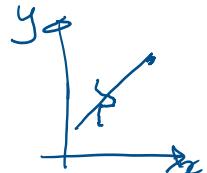
$$p(x) = \prod_{i=1}^n p(z_i | x_{\pi_i}, \theta_i)$$

semantic of intervening on node J

$$p(z | \text{intervention}_{\text{on } J}) = \left( \prod_{i \neq J} p(x_i | x_{\pi_i}, \theta_i) \right) p(x_J | \text{intervention})$$



identify causal direction ↘ via parametric assumptions  
via interventions



see thoughts of Bernhard Schölkopf on causality:

<https://arxiv.org/abs/1911.10500>

(and references therein, e.g. his book: )

Elements of Causal Inference, 2017

By Jonas Peters, Dominik Janzing and Bernhard Schölkopf

<https://mitpress.mit.edu/books/elements-causal-inference>

(available for free online)