

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

mean field approximation

$$\min_{q \in \mathcal{Q}_{MF}} KL(q || p)$$

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

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

[lecture 22 Fall 2017 link](#)

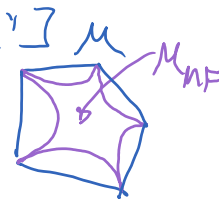
• $KL(\cdot || p)$ is a convex fct. of q

but \mathcal{Q}_{MF} is a non-convex constraint set

Ising model

$$\mu_{ij} = \mu_i \cdot \mu_j$$

non-convex constraint



but can monitor progress



pros & cons of variational methods

⊕ optimization based
 ⇒ often faster to run
 & easier to debug

⊖ biased estimate
 $E_{q(z)} [f(z)] \neq E_p [f(z)]$

vs. Sampling

⊖ noisy ⇒ harder to debug
 mixing problem for chains

⊕ unbiased estimate
 $E [E_{q(z)} [f(z)]] = E_p [f(z)]$
 with respect to random sample

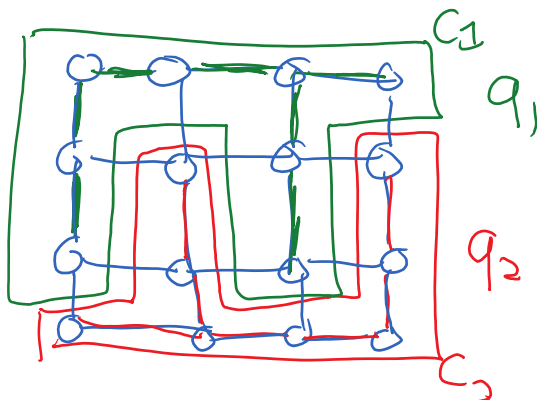
structured mean field:

[lecture 22 Fall 2017 link](#)

$$\text{idea } q(z) = \prod_{j=1}^k q_j(z_{C_j})$$

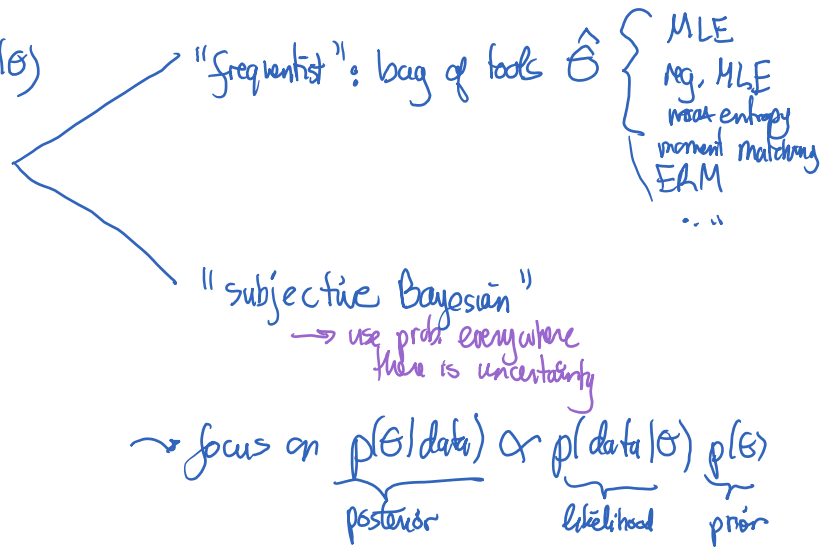
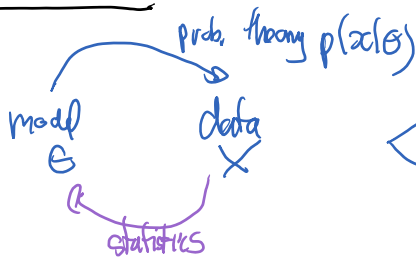
where C_1, \dots, C_k is a partition of V
 and q_j 's are tractable distributions

(for example tree UGM)



$$q = q_1 \cdot q_2$$

Bayesian methods:



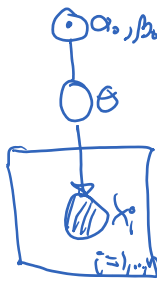
caricature: Bayesian is "optimist"
they think you can get "good" models

⇒ obtain a method by doing inference in model

frequentist is "pessimist" → use analysis tools

Example: biased coin

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



α_0, β_0 ← hyperparameters for prior

e.g.

$$\theta \sim \text{Unif}[0,1] = \text{Beta}(1,1)$$

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

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

posterior: $p(\theta | x_{1:n}) \propto \left(\prod_{i=1}^n p(x_i | \theta) \right) p(\theta)$

$$= \theta^{\sum_{i=1}^n x_i} (1-\theta)^{n - \sum_{i=1}^n x_i} \theta^{\alpha_0 - 1} (1-\theta)^{\beta_0 - 1} \mathbb{1}_{[0,1]}(\theta)$$

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

↳ "conjugate prior" to the Bernoulli likelihood model

conjugate priors

consider a family F of dist. $F = \{ p(\theta | \alpha) : \alpha \in \Omega \}$ on θ

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

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

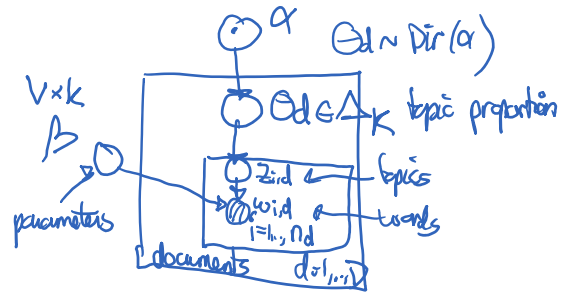
i.e. \exists some $\alpha'(x, \alpha)$ s.t. $p(\theta | x, \alpha) = p(\theta | \alpha')$

Sidonsti.

Sidenote: if use conjugate prior in a DGM then Gibbs sampling can be easy

[e.g. this is case in LDA topic model]

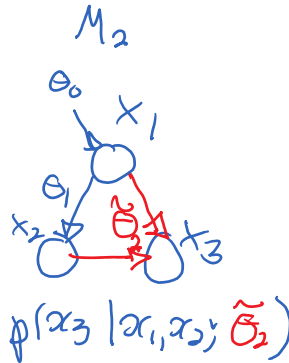
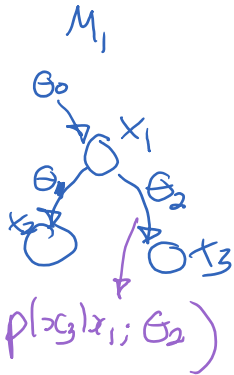
example hwk 1:
Dirichlet prior is conjugate for multinomial likelihood model



16h03

Model selection:

say we want to choose between 2 DGM



[note here that "M1 subset M2"]

"caution notation" ↓

as a frequentist: $\hat{\theta}_{M_1}^{MLE} = \text{arg max}_{\theta_0, \theta_1, \theta_2} \log p(\text{data} | \theta_0, \theta_1, \theta_2, \text{"model"}=M_1)$

$\hat{\theta}_{M_2}^{MLE} = \text{arg max}_{\theta_0, \theta_1, \tilde{\theta}_2} \log p(\text{data} | \theta_0, \theta_1, \tilde{\theta}_2, \text{"model"}=M_2)$
 ↳ different space than θ_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 $M_1 \subseteq M_2$

(ie. you would always choose "bigger model")

→ as frequentist, use cross-validation or validation set

ie. $\log p(\text{test data} | \hat{\theta}_{M_1}^{MLE}(\text{train data}))$
 $M=M_1$

Bayesian alternatives:

true Bayesian → sum over models (integrate out uncertainty about M)

introduce a prior over models $p(M)$

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

data

$$\sum_M \left[\int_{\Theta \in \Theta} p(x_{\text{new}} | \theta, M) p(\theta | D, M) d\theta \right] p(M | D)$$

standard Bayesian prediction dist. for one model

posterior on θ given data D , model M

sum over posterior over models

doing model averaging

⊛ in model selection forced to pick model

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

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

$$\int_{\Theta \in \Theta_M} p(\text{data} | \theta, M) p(\theta | M) d\theta$$

likelihood

to compare two models, look at

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

Bayes factor

prior ratio

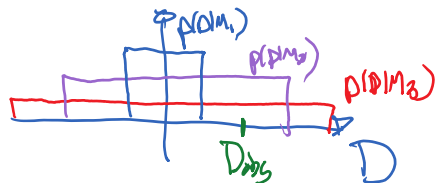
"Uniform prior over models" ⇒ then pick among k models M_1, \dots, M_k

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

"empirical Bayes" "type II ML"

When # of models is "small", then this approach is "good" (i.e. want overfit)

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



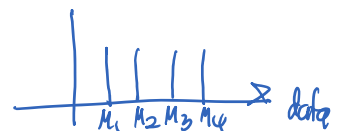
$p(D|M)$ is normalized over D

Vs.

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

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

say e.g. $p(D|M) = \delta(D|M)$



how to compute marginal likelihood:

use variational inference

how to compute marginal likelihood:

use approximations $\left\{ \begin{array}{l} \text{variational inference} \\ \text{sampling} \end{array} \right.$

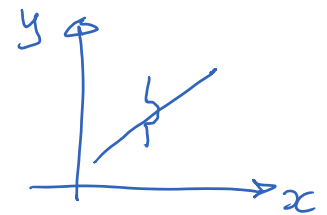
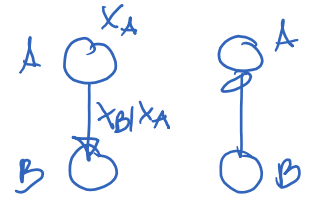
simple approximation \rightarrow Bayesian information criterion

Causality:

structural causal model: graph model + intervention model

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

identify causal direction $\left\{ \begin{array}{l} \text{via parametric assumptions} \\ \text{via interventions} \end{array} \right.$



semantic of intervening on node J

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

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)