

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

mean field continuation

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

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

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

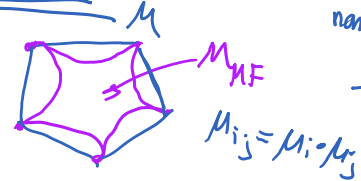
• $KL(\cdot \| p)$ is a convex of q

but Q_{MF} is a non-convex constraint set \Rightarrow can get stuck in local minima

Ising model

$$M_{ij} = M_i \cdot M_j$$

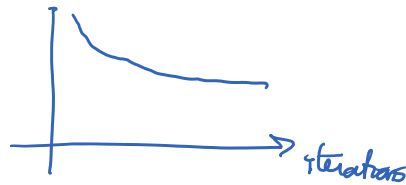
non-convex constraint



[lecture 22 Fall 2017 link](#)

but can monitor progress

$$KL(q^{(t)} \| p) + \underline{cost}$$



pros & cons of variational methods

vs. sampling

- ⊕ optimization based \Rightarrow often faster to run & easier to debug

- ⊖ noisy \Rightarrow harder to debug
- missing problem for chains

⊖ biased estimate

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

⊕ unbiased estimate

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

with respect to random sample

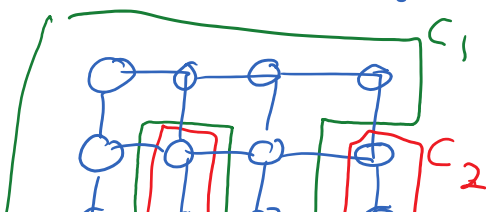
structured mean field :

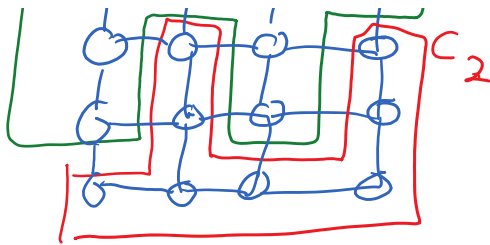
[lecture 22 Fall 2017 link](#)

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

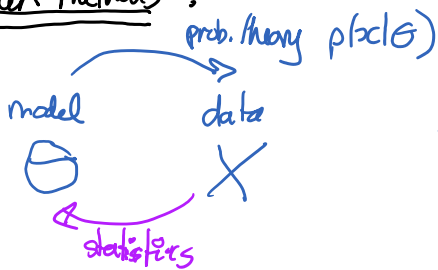
where C_1, \dots, C_K is a partition of V

and q_j 's are tractable distributions (for example tree UGM)





Bayesian methods :



"frequentist": bag of tools Θ { MLE, reg. ML, max entropy, moment matching, ERM }

"subjective Bayesian"
 → use proba everywhere
 there is uncertainty

→ focus on $\underbrace{p(\theta | \text{data})}_{\text{posterior}} \propto \underbrace{p(\text{data} | \theta)}_{\text{likelihood}} \underbrace{p(\theta)}_{\text{prior}}$

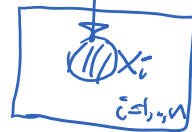
critique: Bayesian is "optimist" & "good"
 they think you can get good models

⇒ obtain a method by doing inference in model

frequentist is "pessimist" → use analysis tools
 $\alpha_0, \beta_0 \leftarrow$ hyperparameters for the prior

Example: biased coin:

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



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 $\propto p(\theta | x_{1:n}) \propto \left(\prod_i p(x_i | \theta) \right) p(\theta)$
 $= \theta^{\sum_i x_i} (1-\theta)^{n - \sum_i x_i} \theta^{\alpha_0 - 1} (1-\theta)^{\beta_0 - 1} \mathbb{1}_{[0, 1]}(\theta)$
 $\stackrel{\sum_i x_i \triangleq n_1}{=} \theta^{n_1} (1-\theta)^{n - n_1} \theta^{\alpha_0 - 1} (1-\theta)^{\beta_0 - 1} \mathbb{1}_{[0, 1]}(\theta)$

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

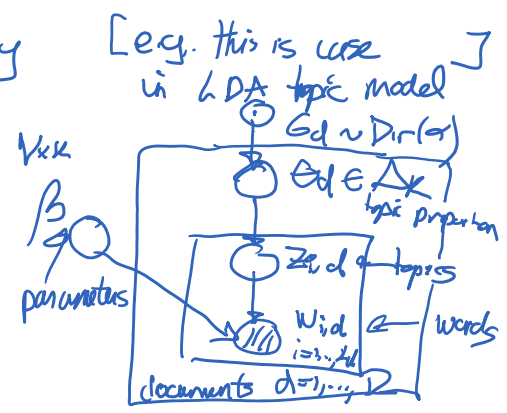
↳ "conjugate prior" to the Bernoulli likelihood model

more generally,

consider a family F of dist. $F = \{p(\theta|\alpha) : \alpha \in \mathcal{H}\}$
 say that F is a "conjugate family" to observation model $p(x|\theta)$
 if posterior $p(\theta|x, \alpha) \in F$ for any $x \in \mathcal{X} \times \mathcal{O}$
 i.e. \exists an $\alpha'(x, \alpha)$ s.t. $p(\theta|x, \alpha) = p(\theta|\alpha')$

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

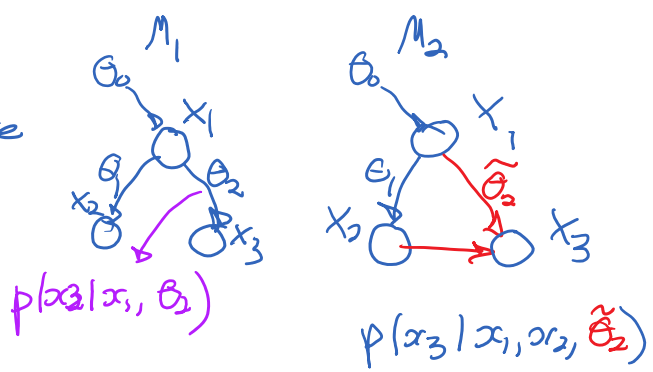
example hwk 1:
 Dirichlet prior
 is conjugate
 for multinomial likelihood model



4h33

model selection:

say we want to choose
 between 2
 DGM



(note here that " $M_1 \subseteq M_2$ ")

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

$\hat{\theta}_{M_2}^{ML} = \text{argmax}_{\theta_0, \theta_1, \hat{\theta}_2} \log p(\text{data} | \theta_0, \theta_1, \hat{\theta}_2, \text{"model"} = M_2)$
 (different space than θ_2)

how to choose between models?

can't compare $\log p(\text{data} | \hat{\theta}_{M_1}^{ML}, M=M_1)$ vs. $\log p(\text{data} | \hat{\theta}_{M_2}^{ML}, M=M_2)$
 because LHS \neq RHS since $M_1 \subseteq M_2$

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

→ as frequentist, use cross-validation
 ie. $\log p(\text{test data} | \hat{\theta}_{M_i}^{ML}(\text{train data}), M = M_i)$

Bayesian alternatives:

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

introduce a prior over models $p(M)$

$$\begin{aligned}
 p(x_{\text{new}} | D) &= \sum_M p(x_{\text{new}} | D, M) p(M | D) \\
 &= \sum_M \left[\int_{\theta \in \Theta_M} p(x_{\text{new}} | \theta, M) p(\theta | D, M) d\theta \right] p(M | D)
 \end{aligned}$$

data (pointing to x_{new})
standard Bayesian predictive distribution for one model (pointing to the integral)
sum over posterior over models (pointing to the outer sum)
posterior on θ given data D & model M (pointing to $p(\theta | D, M)$)
during model averaging (pointing to the overall expression)

(*) 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

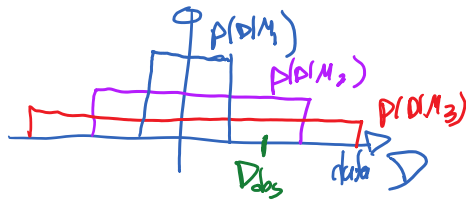
"uniform prior over models"; ⇒ then pick among k models M_1, \dots, M_k
 by max. $p(\text{data} | M = M_i)$

"empirical Bayes" "type II ML"

when # of models is "small", then this approach is fine

11.2. can't overfit)

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

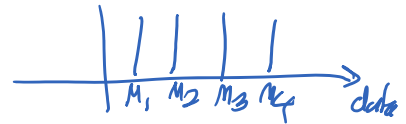


$p(D|M)$ is normalized over \mathcal{D}

vs.

$p(D|\hat{\Theta}_{ML}(\mathcal{D}), M)$ [can overfit badly]

but M_k can still overfit if have many nodes



say eg. $p(D|M) = \delta(D, M)$

how to compute marginal likelihood:

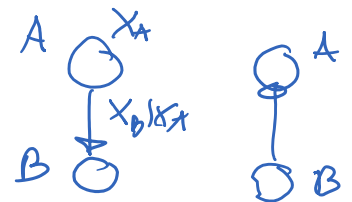
use approximations \leftarrow variational inference sampling

simple approximation \rightarrow Bayesian information criteria

Causality:

structural causal model: graph model + intervention model

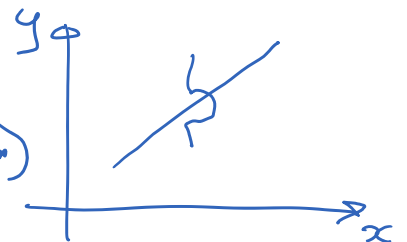
$$p(\mathcal{X}) = \prod_{i=1}^n p(x_i | x_{\pi_i}, \epsilon_i)$$



identify causal direction \leftarrow via intervention
 \leftarrow via parametric assumptions

Semantic of intervening on node J

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



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)