

today: • Bayesian approach  
 • model selection

Variational methods vs. sampling

recall for mean field

target dist.  $\left\{ \begin{array}{l} \min KL(q \parallel p) \\ q \in Q_{\text{MF}} \end{array} \right\}$  non-convex

fully factorized distributions i.e.  $q(x) = \prod q_i(x_i)$

in case where  $p(x)$  was from Ising model  
 → can coordinate descent on  $q_i$ 's  
 approximate marginal on node  $q_i(x_i=1) = \hat{\mu}_i$

say we have converged to a stationary pt.  $\hat{\mu}_i^*$   
 usually  $\hat{\mu}_i^* \neq p(x_i=1)$   
 "biased" / "inexact"

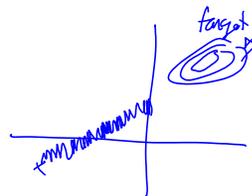
in contrast: sampling methods are usually asymptotically unbiased

example: if use Gibbs sampling to get  $\vec{x}^{(t)}$  (samples)

then  $\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T x_i^{(t)} = p(x_i=1)$

from Ergodic thm "asymptotically unbiased"

problem: "mixing time" ← how long it takes to forget initial conditions



sometimes can be very long "sticky chain"

⇒ slow convergence of Monte Carlo estimate

\* in practice, you can not use the first few samples to reduce the bias of estimate  
 (1, ..., 0, ..., 1, ...)

"burn in time"

summary:

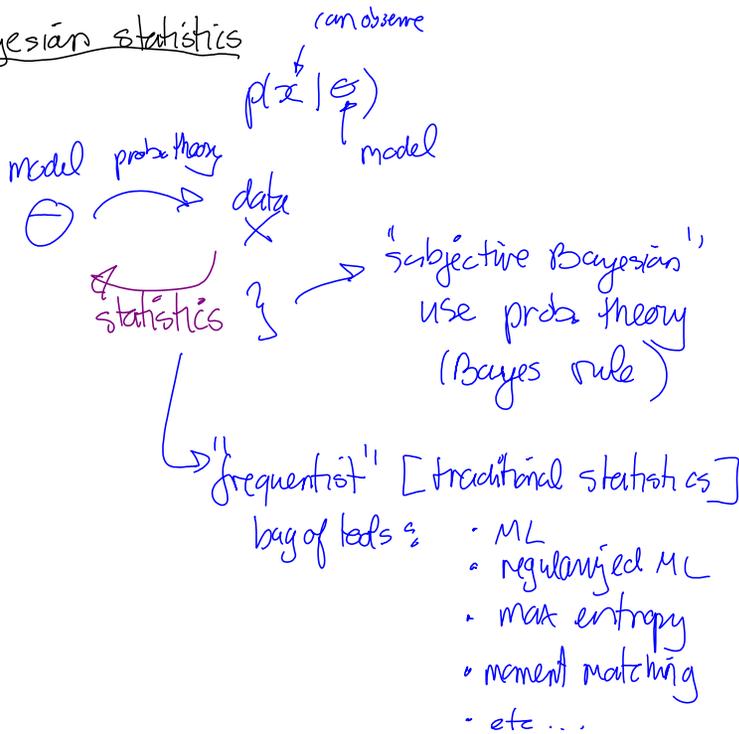
variational approach:

- often faster than sampling, but it is inexact

eg. QMR-DT network seconds for variational vs hours for sampling

- easier to debug the variational

Bayesian statistics



caricature: Bayesian is "optimist": think that can come up with good models,

$\rightarrow$  obtain a method by pulling the Bayesian crank

frequentist is more pessimistic:  $\rightarrow$  use analysis tools

Bayesian:  $p(x|\Theta)$  "likelihood model"  
 $p(\Theta)$  "prior"

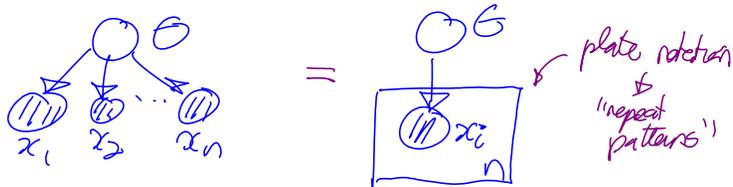
posterior:  $p(\Theta|x) = \frac{p(x|\Theta)p(\Theta)}{p(x)}$  Bayes rule  
"marginal likelihood"  $\rightarrow$  normalization

# Summary belief of Bayesian

## example: biased coin

Bayesian model  $x_i \in \{0,1\}$   
 $x_i | \theta \stackrel{iid.}{\sim} \text{Bernoulli}(\theta)$   $p(x_i | \theta) = \theta^{x_i} (1-\theta)^{1-x_i}$   
 $\theta \sim \text{Unif}[0,1]$  (uniform prior)

graph model:

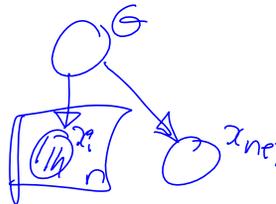


$$\begin{aligned} \text{posterior: } p(\theta | x_{1:n}) &\propto p(x_{1:n} | \theta) p(\theta) \\ &= \left( \prod_{i=1}^n p(x_i | \theta) \right) p(\theta) \\ &= \theta^{\sum_{i=1}^n x_i} (1-\theta)^{n - \sum_{i=1}^n x_i} \mathbb{1}_{[0,1]}(\theta) \\ &\stackrel{\Delta}{=} n_{\pm} \end{aligned}$$

$\nearrow$  Beta( $\alpha, \beta$ ) distribution where  $\alpha = n_{\pm} + 1$   
 $\beta = n - n_{\pm} + 1$

question: what is probability of next flip = 1?

frequentist  $\hat{\theta}_{ML} = \frac{n_{\pm}}{n}$



Bayesian integrates out the uncertainty

$$p(x_{n+1} | x_{1:n}) = \int_{\theta} p(x_{n+1} | \theta) \underbrace{p(\theta | x_{1:n})}_{\text{posterior}} d\theta$$

$\uparrow$   
 predictive distribution

$$p(x_{n+1} = 1 | x_{1:n}) = \int_{\theta} \theta p(\theta | x_{1:n}) d\theta \leftarrow \text{posterior mean}$$

mean of  $\text{Beta}(\alpha, \beta)$  is  $\frac{\alpha}{\alpha+\beta} = \frac{n_1+1}{n+2}$  here

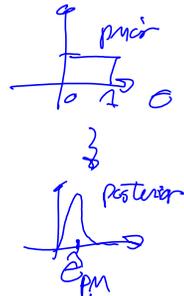
notice that for  $n=0 \rightarrow$  get  $\frac{1}{2}$  ↑ smoothed version of ML

$$\hat{\theta}_{\text{posterior mean}} = \frac{n_1}{n} \left[ \frac{n}{n+2} \right] + \frac{1}{2} \left[ \frac{2}{n+2} \right] = \frac{n_1+1}{n+2}$$

$$= \hat{\theta}_{\text{ML}} P_n + \hat{\theta}_{\text{prior}} (1-P_n)$$

convex combination of  $\hat{\theta}_{\text{ML}}$  &  $\hat{\theta}_{\text{prior}}$

$P_n \xrightarrow{n \rightarrow \infty} 1$



variance of a beta  $\frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)} = \left(\frac{n_1}{n}\right) \left(1-\frac{n_1}{n}\right) \underbrace{O\left(\frac{1}{n}\right)}_{\xrightarrow{n \rightarrow \infty} 0}$

$\hat{\theta} (1-\hat{\theta})$

posterior "contracts" around  $\hat{\theta}_{\text{PM}} \xrightarrow{n \rightarrow \infty} \hat{\theta}_{\text{ML}} \downarrow$  true parameter  
 $\hat{\theta}_{\text{ML}} \xrightarrow{n \rightarrow \infty} \theta^*$

"Bernstein von-Mises thm"

$\rightarrow$  "Bayesian central limit thm." which basically says that if prior puts non-zero mass around true model  $\theta^*$ , then posterior concentrates around  $\theta^*$  as a Gaussian asymptotically

for a Gaussian mean & mode are the same

so can approximate  $E[\theta | \text{data}]$   
 with  $\hat{\theta}_{\text{MAP}} = \underset{\theta}{\text{argmax}} p(\theta | \text{data})$   
 [fine for large  $n$ ]

revisiting the example:



$T$  coins picked randomly each flipped  $n$  times...

as a frequentist, empirical distribution on  $x_{1:n}$  will converge as  $T \rightarrow \infty$  to  $p(x_1, \dots, x_n) = \int \left( \prod_{i=1}^n p(x_i | \theta) \right) p(\theta) d\theta$

distribution of coins in jar  
↓

(continuous) mixture distribution

here,  $X_1, \dots, X_n$  are not independent  
 on the other hand,  $p(x_1, \dots, x_n) = p(x_{\pi(1)}, \dots, x_{\pi(n)})$   
 ↓ for any permutation  $\pi: 1:n \rightarrow 1:n$   
 $X_1, \dots, X_n$  are "exchangeable" weaker than independence

De Finetti's representation thm.:

$X_1, X_2, \dots$  is infinitely exchangeable

⇔ ∃ unique  $p(\theta)$  on some space  $\Theta$

s.t.  $p(x_1, \dots, x_n) = \int \left( \prod_{i=1}^n p(x_i | \theta) \right) p(\theta) d\theta$

Multinomial model

e.g. modeling words...

$X | \theta \sim \text{Mult}(\theta, 1)$  where  $\theta \in \Delta_K$   
 i.e.  $\sum_{l=1}^K \theta_l = 1$   
 $\theta_l \geq 0$

$\hat{\theta}_l^{ML} = \frac{n_l}{n}$  if  $K > n$   
 then some  $\hat{\theta}_l^{ML} = 0$  for some  $l$   
 ⇒ overfitting

as a Bayesian, put prior on  $\Delta_K = \Theta$

a convenient property of prior family is "conjugacy"

consider family of distributions  $F = \{ p(\theta | \alpha) : \alpha \in A \}$

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

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

if posterior  $p(\theta|x, \alpha) = \frac{p(x|\theta)p(\theta|\alpha)}{p(x|\alpha)}$

↑ hyperparameter

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

for multinomial:

likelihood  $p(x_{1:n}|\theta) = \prod_{i=1}^n p(x_i|\theta)$

$$= \prod_{l=1}^k \theta_l^{\sum_{i=1}^n x_{i,l}}$$

$x_i = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix}$

if use prior  $\propto \prod_{l=1}^k \theta_l^{\alpha_l}$  stuff

Dirichlet distribution  
→ dist. over  $\Delta_k$

$$\text{Dir}(\theta|\alpha_1, \dots, \alpha_k) = \frac{1}{B(\vec{\alpha})} \prod_{l=1}^k \theta_l^{\alpha_l - 1}$$

convention

valid  $\alpha_l > 0$

$$E[\theta_l|\alpha_1, \dots, \alpha_k] = \frac{\alpha_l}{\sum_j \alpha_j}$$

$$\text{variance}(\theta_l) = O\left(\frac{1}{\sum_j \alpha_j}\right)$$

- $\alpha_l = 1 \rightarrow$  get uniform distribution
- $k = 2 \rightarrow$  get Beta distribution
- $\alpha_l < 1 \cup$  shape distribution (no mode)
- $\alpha_l \geq 1 \cap$  unimodal bump

\* for multinomial model; if <sup>prior:</sup>  $p(\theta|\alpha) = \text{Dir}(\theta|\vec{\alpha})$

posterior  $p(\theta|x_1, \dots, x_n) \propto \prod_{l=1}^k \theta_l^{n_l + \alpha_l - 1}$

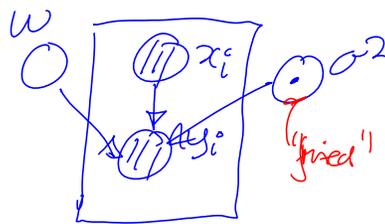
$$= \text{Dir}(\theta|(n+\alpha)_{l=1}^k)$$

posterior mean:  $E[\theta_l|\text{data}] = \frac{n_l + \alpha_l}{n + \sum \alpha_l}$  → "prior"

$$n + \sum_j \dots \text{('aints')}$$

### Bayesian linear regression

say  $y = w^T x + \epsilon$   $\epsilon \sim N(0, \sigma^2)$



likelihood model:  $p(y|x, w) = N(y | w^T x, \sigma^2)$

prior:  $p(w) = N(w | 0, \frac{I}{\lambda})$  (conjugate)  $\lambda \leftarrow \text{precision}$

posterior:  $p(w | y_{1:n}, x_{1:n})$  is also Gaussian

with covariance  $\hat{\Sigma}_n = \lambda I + \frac{X^T X}{\sigma^2}$   $n \times d$  design matrix  $\vec{y} = (y_1, \dots, y_n)$

posterior mean  $\hat{\mu}_n = \hat{\Sigma}_n^{-1} (X^T \vec{y})$

$\hookrightarrow$  same as in ridge regression with  $\hat{\lambda} = \sigma^2 \lambda$

as a Bayesian;

compute predictive dist.  $p(y_{new} | x_{new}, x_{1:n}, y_{1:n})$

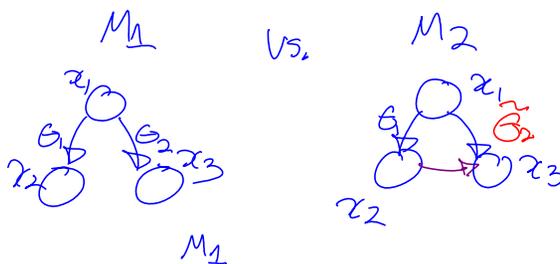
$$= \int_w \underbrace{p(y_{new} | x_{new}, w)}_{\text{Gaussian dis. model}} \underbrace{p(w | \text{data})}_{\text{Gaussian posterior}} dw$$

$$= N(y_{new} | \hat{\mu}_n^T x_{new}, \sigma_{\text{predictive}}^2)$$

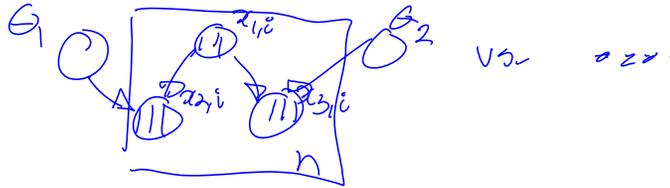
$$\sigma_{\text{predictive}}^2(x_{new}) = \underbrace{\sigma^2}_{\text{noise model}} + \underbrace{x_{new}^T \hat{\Sigma}_n^{-1} x_{new}}_{\text{from posterior covariance}}$$

### Model selection

say want to choose between



as a Bayesian



as a frequentist,

$$\hat{\Theta}_{M_1}^{ML} = \underset{\Theta_1, \Theta_2}{\text{argmax}} \log p(\text{data} | \Theta_1, \Theta_2, M_1)$$

$$\hat{\Theta}_{M_2}^{ML} = \underset{\substack{\Theta_1, \hat{\Theta}_2 \\ \text{different space}}}{\text{argmax}} \log p(\text{data} | \Theta_1, \hat{\Theta}_2, M_2)$$

compare  $\log p(\text{data} | \hat{\Theta}_{M_1}^{ML}, M_1)$  vs.  $\log p(\text{data} | \hat{\Theta}_{M_2}^{ML}, M_2)$ ?

here  $M_1 \ll M_2 \Rightarrow \uparrow \leq \uparrow$

here, likelihood is useless...

instead, use cross-validation

\* here, Bayesian alternative:

true Bayesian, sum over models  
(integrate out the uncertainty)

$$p(x_{\text{new}} | \underset{\text{data}}{D}) = \sum_M \int_{\Theta} p(x_{\text{new}} | \Theta, M) \underbrace{p(M, \Theta | D)}_{p(\Theta | D, M) p(M | \text{data})} d\Theta$$

$$= \sum_M \underbrace{p(M | \text{data})}_{\text{model averaging}} \left[ \int_{\Theta} p(x_{\text{new}} | \Theta, M) \underbrace{p(\Theta | \text{data}, M)}_{\text{standard predictive dist. for one model}} d\Theta \right]$$

\* if force to pick model  $j$

pick model which maximize  $p(M | \text{data}) \underbrace{p(\text{data} | M)}_{\text{"marginal likelihood"}}$

$$p(D | M) = \int_{\Theta} p(D | \Theta, M) p(\Theta | M) d\Theta$$

here, stopping mean  $\frac{p(M=M_1 | D)}{p(M=M_2 | D)}$

Bayes factor

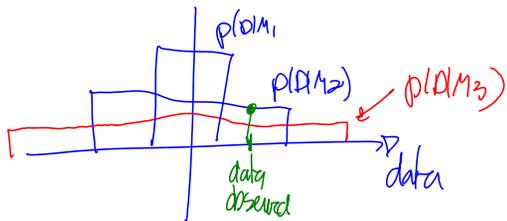
$$\frac{p(M_1 | D)}{p(M_2 | D)} = \frac{p(D | M_1) p(M_1)}{p(D | M_2) p(M_2)}$$

JG

pick  $\hat{M} = \underset{M}{\text{argmax}} p(D | M) \leftarrow \begin{array}{l} \text{"empirical Bayes"} \\ \text{"type 2 ML"} \end{array}$

too many models  $\rightarrow$  can <sup>still</sup> overfit

Caroon why marginal likelihood works vs. ML:

$$M_1 \subseteq M_2 \subseteq M_3$$


BIC criterion  $\rightarrow$  approximation  $p(D|M)$