

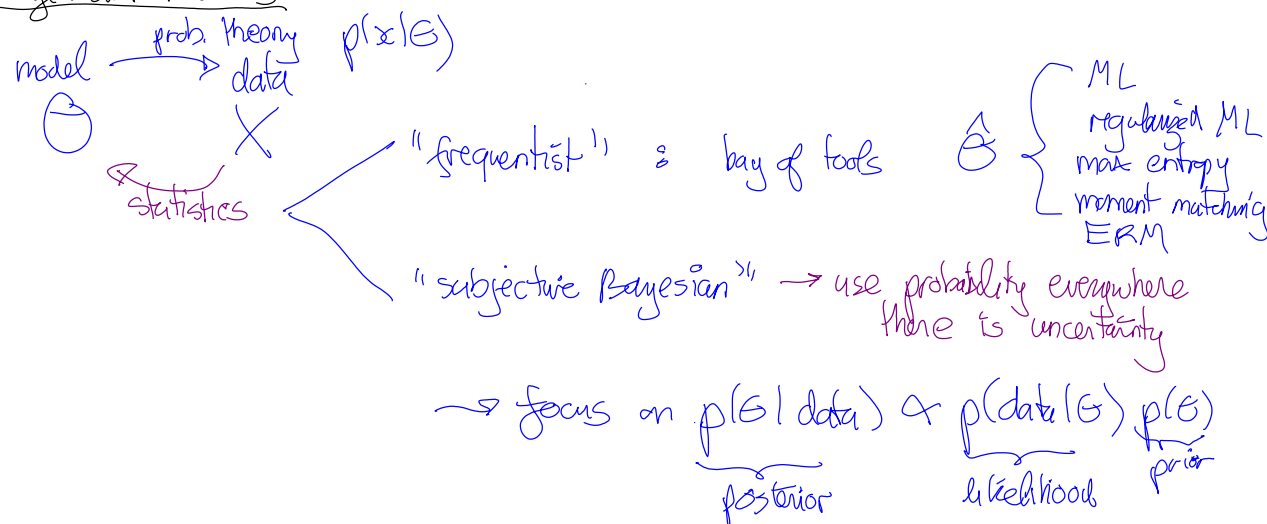
# Lecture 23 - scribbles

Tuesday, November 28, 2017

14:34

today: • Bayesian approach  
• model selection

## Bayesian methods



caricature: Bayesian is "optimist": she thinks you can get good models

⇒ obtain a method by doing probabilistic inference in model

frequentist is "pessimist" ⇒ use analysis tools

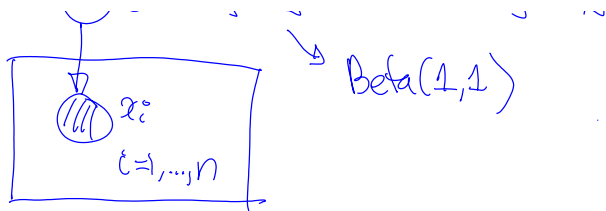
## Example: biased coin

Bayesian model  $x_i \in \{0,1\}$

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

$\theta \sim \text{unif}[0,1]$   $\theta \sim \text{unif}[0,1]$

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



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

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

$n_1$

$\uparrow$   $\text{Beta}(\theta | \alpha, \beta)$  where  $\alpha = n_1 + 1$   
 $\beta = n - n_1 + 1$

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

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

here  $\text{Beta}(\theta | \alpha_0, \beta_0)$  is a "conjugate prior" to the Bernoulli likelihood model

more generally, consider family  $F$  of dist.:  $F = \{ p(\theta | \alpha) : \alpha \in \mathcal{A} \}$

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

if posterior  $p(\theta | x, \alpha) \in F$

i.e.  $\exists \alpha' \ni \alpha. p(\theta | x, \alpha) = p(\theta | \alpha')$

$$\parallel$$

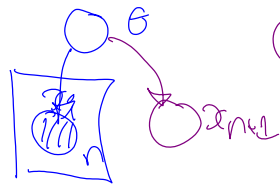
$$\frac{p(x | \theta) p(\theta | \alpha)}{\int_{\mathcal{E}} p(x | \theta') p(\theta' | \alpha) d\theta'}$$

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

→ Bayesian in action:

question: what is the probs. that next flip = 1?

frequentist:  $\hat{\theta}_{ML} = \frac{n_1}{n}$



Bayesian: integrate out uncertainty

using cond. indep.

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

"predictive dist"

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

$$E[\theta | \text{data}] = \frac{\alpha}{\alpha + \beta} = \frac{n_1 + 1}{n + 2} \quad (\text{with } \alpha = \beta = 1 \text{ i.e. uniform prior})$$

$$\hat{\theta}_{\text{posterior mean}} = \underbrace{\frac{n_1}{n}}_{\hat{\theta}_{ML}} \underbrace{\left[ \frac{1}{n+2} \right]}_{p_n} + \underbrace{\frac{1}{2}}_{\theta_{\text{prior mean}}} \underbrace{\left[ \frac{2}{n+2} \right]}_{1-p_n}$$

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

i.e.  $\hat{\theta}_{\text{posterior mean}} \xrightarrow{n \rightarrow \infty} \hat{\theta}_{ML} = \text{"true } \theta \text{"}$

$$\xrightarrow{n \rightarrow \infty} 0 \text{ as } n \rightarrow \infty$$

$$\begin{aligned} \text{Variance of a Beta} : \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)} &= \left(\frac{n_1}{n}\right)\left(1-\frac{n_1}{n}\right)\tilde{O}\left(\frac{1}{n}\right) \\ &= \hat{\theta}_{ML}(1-\hat{\theta}_{ML})O\left(\frac{1}{n}\right) \end{aligned}$$

posterior "contracts" around  $\hat{\theta}_{\text{post. mean}} \xrightarrow{n \rightarrow \infty} \hat{\theta}_{ML} = \theta^*$

"Bernstein von-Mises thm"

→ "Bayesian CLT": basically says that if prior put non-zero mass on true parameter  $\theta^*$  (i.e.  $x_i \sim p(x|\theta^*)$ ) then posterior concentrates around  $\theta^*$  as a Gaussian asymptotically

recall from hwk 1: multinoulli model

$X|\theta \sim \text{Mult}(\theta)$  where  $\theta \in \Delta_K$

$$\hat{\theta}_{ML} = \left(\frac{n\theta}{n}\right)_{k=1}^K$$

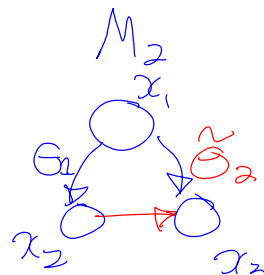
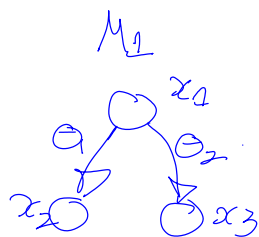
putting Dirichlet prior over  $\theta$   $\theta \sim \text{Dir}(\alpha)$

we got Dirichlet posterior  $\theta | \text{data} \sim \text{Dir}((\alpha_k + n_k)_{k=1}^K)$

thus Dirichlet is conjugate prior to Multinoulli likelihood

Model selection:

say we want to choose between 2 DGM



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

as a frequentist:  $\hat{\Theta}_{M_1}^{ML} = \arg \max_{\Theta, \Theta_2} \log p(\text{data} | \Theta, \Theta_2, \text{model} = M_1)$  (causal notation)

$\hat{\Theta}_{M_2}^{ML} = \arg \max_{\Theta, \Theta_2} \log p(\text{data} | \Theta, \Theta_2, \text{model} = M_2)$

$\Theta_1, \Theta_2$  different space

how to choose between models?

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

because LHS  $\leq$  RHS since  $M_1 \subseteq M_2$

(i.e. you would always choose "bigger model")

→ as frequentist, use cross-validation i.e.  $\log p(\text{test data} | \hat{\Theta}_{M_2}^{ML}(\text{train data}), M = M_1)$

Bayesian alternative:

true Bayesian → sum over models  
(integrate out uncertainty)

introduce prior over models  $p(M)$

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

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

$p(M | D) \propto p(M) p(D | M)$

$$\underbrace{p(\theta | D, M)}_{\text{posterior on } \theta \text{ given data } D \text{ and model } M} \underbrace{p(M | D)}_{\text{marginal model prob.}}$$

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

doing Model averaging

Standard Bayesian predictive for one model

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

⊛ in model selection, forced to pick one model  
 $\Rightarrow$  pick model that maximizes

$$p(M | \text{data}) \propto p(\text{data} | M) p(M)$$

$$p(\text{data} | M) = \underbrace{\int_{\theta \in \Theta_M} p(\text{data} | \theta, M) p(\theta | M) d\theta}_{\text{"marginal likelihood"}}$$

to compare two models, look at

$$\frac{p(M=M_1 | D)}{p(M=M_2 | D)} = \underbrace{\frac{p(D | M_1)}{p(D | M_2)}}_{\text{Bayes factor}} \underbrace{\frac{p(M_1)}{p(M_2)}}_{\text{prior ratio}}$$

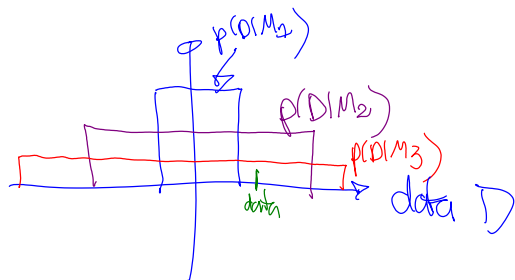
"uniform prior over models"; then we can pick among  $K$  models  $M_1, \dots, M_K$

by maximizing  $p(\text{data} | M = M_i)$   
 "empirical Bayes"  
 "type II ML"

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

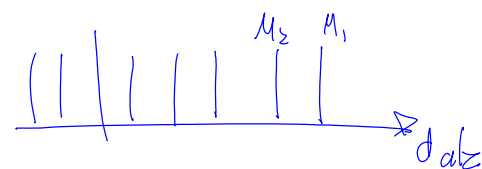
Zoubin's cartoon: suppose  $M_1 \leq M_2 \leq M_3$

$p(\text{data} | \mathcal{G})$



$p(D|M)$  is normalized over  $D$   
 vs  
 $p(D | \hat{E}_{ML}(D), M)$

type II ML can still overfit when have many models say e.g.  $p(D|M) = \mathcal{S}(D, M)$



$$\mathcal{S}(D, M) = \begin{cases} 1 & \text{if } D=M \\ 0 & \text{o.w.} \end{cases}$$

how to compute marginal likelihood:

use approximations  $\leftarrow$  variational inference  
 sampling