

Lecture 23 - scribbles

Tuesday, November 27, 2018 14:39

today: • Bayesian approach
• model selection

Bayesian methods



"frequentist" : bag of tools

- $\hat{\theta}$ { ML
- regularized ML
- max. entropy
- moment matching
- ERM

"subjective Bayesian"

→ use probability everywhere there is uncertainty

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

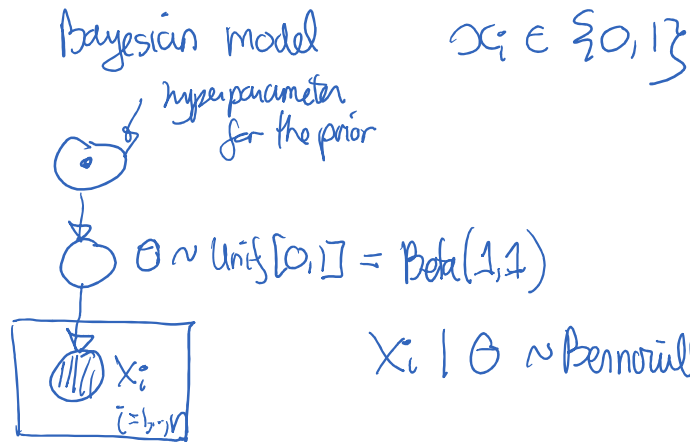
caricature: Bayesian is "optimist": they think you can get "good" models

⇒ obtain a method by doing prob. inference in model

frequentist is "pessimist" → use analysis tools

Example: biased coin:

Example: biased coin:



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

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

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

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

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

\hookrightarrow "conjugate prior" to the Bernoulli likelihood model

more generally;

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

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

if posterior $p(\theta|x, \alpha) \in F$ for any $x \sim X \in \mathcal{E}$
 i.e. \exists an α' s.t. $p(\theta|x, \alpha) = p(\theta|\alpha')$

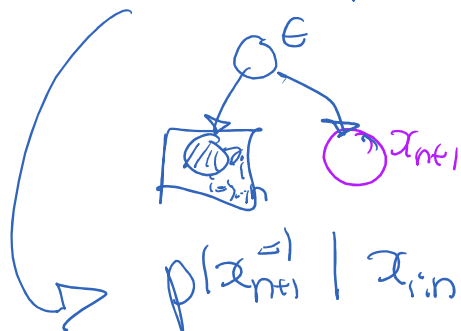
sidenote: if use conjugate priors in a DGM
 then Gibbs sampling can be easy

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

* Bayesian in action:

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

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



Bayesian: integrate out uncertainty

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

\downarrow
 θ

using cond. indep

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

$$E[\theta | \text{data}] = \frac{\alpha}{\alpha + \beta} = \frac{n_1 + 1}{n_1 + 1 + n - n_1 + 1} = \frac{n_1 + 1}{n + 2}$$

[with $\alpha_0 = \beta_0 = 1$ & uniform prior]

$$E[\theta | \text{data}] = \frac{\alpha}{\alpha + \beta} = \frac{n_1 + 1}{n_1 + 1 + n - n_1 + 1} = \frac{n_1 + 1}{n + 2} \quad [\text{with } \alpha_0 = \beta_0 = 1 \text{ (uniform prior)}]$$

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

$$p_n \xrightarrow{n \rightarrow \infty} 1 \quad \text{i.o.} \quad \hat{\theta}_{\text{posterior mean}} \xrightarrow{n \rightarrow \infty} \hat{\theta}_{\text{ML}} = \text{"true } \theta^* \text{"} = \theta^*$$

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

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

"Bernstein von-Mises thm."

↪ "Bayesian CLT" : basically says that if prior put non-zero mass on true parameter θ^* [i.e. $x_i \sim p(x|\theta^*)$]

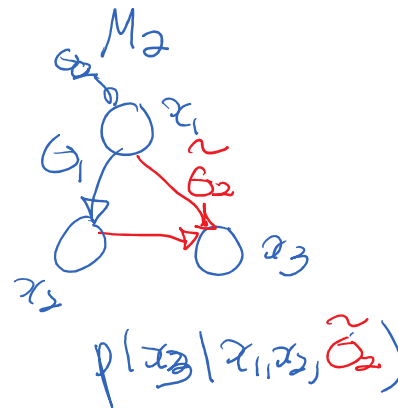
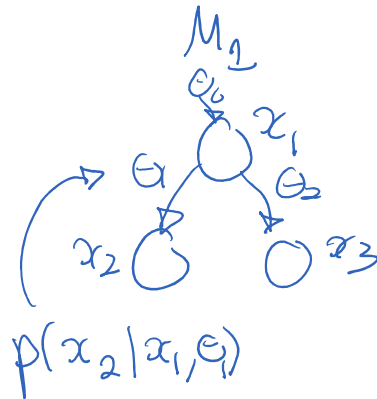
then posterior concentrates around θ^* as a Gaussian asymptotically

recall from hwk 1 : multinomial model

Pitichlet dist. is conjugate to multinomial model

Model selection:

say we want to choose between 2 DOM



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

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

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

different space (pointing to $\tilde{\theta}_2$)

cautious notation (pointing to the quotes around "model")

how to choose between models?

can't 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 $M_1 \subseteq M_2$

(ie you would always choose "bigger model")
 → as frequentist, use cross-validation i.e. $\log p(\text{test data} | \hat{\theta}_{ML}(\text{train data}), M=M_1)$

Bayesian alternatives

true Bayesian → sum over models (integrate out uncertainty)
 introduce prior over models $p(M)$

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

marginal model prob.
 posterior on Θ given data D & model M
 Bayesian standard predictive dist. for one model
 doing model averaging
 [note: $p(x_{\text{new}} | \theta, M, D) = p(x_{\text{new}} | \theta, M)$]
 $p(x_{\text{new}} | \text{data}, M)$

15h35

⊗ in model selection, forced to pick one model

\Rightarrow pick model that maximizes $p(M|\text{data}) \propto \underbrace{p(\text{data}|M)}_{\text{Likelihood}} p(M)$
 $p(\text{data}|M) = \int p(\text{data}|\theta, M) p(\theta|M) d\theta$
"marginal likelihood"

to compare two models, look at

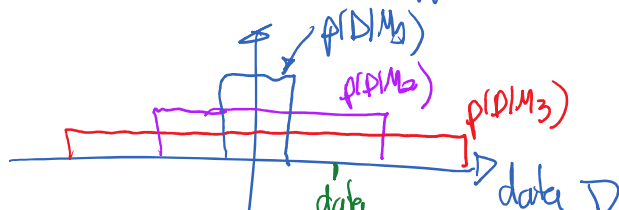
$$\frac{p(M=M_1|\mathcal{D})}{p(M=M_2|\mathcal{D})} = \frac{\underbrace{p(\mathcal{D}|M_1)}_{\text{Bayes factor}} p(M_1)}{\underbrace{p(\mathcal{D}|M_2)}_{\text{Bayes factor}} p(M_2)}$$

\downarrow
 prior ratio

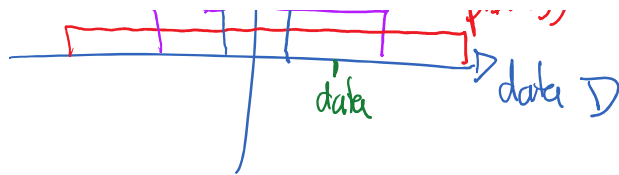
"uniform prior over models"; then we can pick among k models M_1, \dots, M_k
 by maximizing $p(\text{data}|M=M_i)$
"empirical Bayes"
"type II ML"

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

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



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

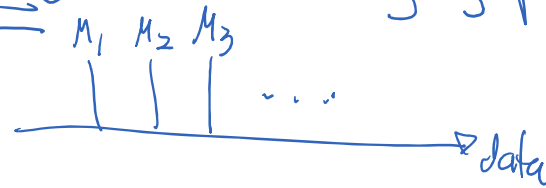


vs.

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

type II ML can still overfit when have many models:

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



how to compute marginal likelihood:

use approximations \leftarrow variational inference
sampling

Bayesian information criterion

BIC is a (rough) approximation of $\log p(\text{data}|M) \approx \log p(\text{data} | \hat{\Theta}_{ML}, M) - \frac{d \log n}{2}$

$\dim(\Theta_M)$

complexity penalty

use Laplace approximation

$$p(D|M) = \int_{\Theta} \underbrace{\left[\prod_{i=1}^n p(x_i | \Theta, M) p(\Theta|M) \right]}_{\exp(-n h(\Theta))} d\Theta$$

where $h(\Theta) = -\frac{1}{n} \sum_{i=1}^n \log p(x_i | \Theta, M) + \log \frac{p(\Theta|M)}{n}$

\rightarrow do Taylor expansion of this around $\hat{\Theta}_{ML}$

↳ do Taylor expansion of this around $\hat{\Theta}_{MAP}$

2 approximations:
 - keep only terms which grows with n
 - replace $\hat{\Theta}_{MAP}$ by $\hat{\Theta}_{MLE}$
 } get BIC

BIC is "consistent"

Gaussian networks:

$$X \sim N(\mu, \Sigma) \quad \mu \in \mathbb{R}^p \quad \Sigma \in \mathbb{R}^{p \times p} \quad \Sigma \succ 0$$

$$p(x; \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^p |\Sigma|}} \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

$$(x-\mu)^T \Sigma^{-1} (x-\mu) = \underbrace{x^T \Sigma^{-1} x}_{\text{linear form on matrices}} - 2\mu^T \Sigma^{-1} x + \mu^T \Sigma^{-1} \mu$$

$$\text{tr}(x^T \Sigma^{-1} x) = \text{tr}(\Sigma^{-1} x x^T)$$

$$= \langle \Sigma^{-1}, x x^T \rangle$$

canonical parameter

$$\Lambda \triangleq \Sigma^{-1}$$

linear form on matrices

$$-2 \langle \underbrace{\Sigma^{-1} \mu}_{\Delta} , x \rangle$$

$$\rightarrow \eta = \Sigma \mu$$

parameter
 $\Lambda \triangleq \Sigma^{-1}$
 precision matrix

$$\begin{aligned} \Rightarrow \mu &= \Sigma \eta \\ &= \Lambda^{-1} \eta \end{aligned}$$

sufficient statistics $T(x) = \begin{pmatrix} x \\ -\frac{1}{2} x x^T \end{pmatrix}$

canonical parameter $\tilde{\eta}(\theta) = \begin{pmatrix} \eta \\ \Lambda \end{pmatrix} = \begin{pmatrix} \Sigma^{-1} \mu \\ \Sigma^{-1} \end{pmatrix}$

$$p(x; \eta, \Lambda) = \exp(\eta^T x + \langle \Lambda, -\frac{1}{2} x x^T \rangle - \underbrace{\left[\frac{1}{2} \eta^T \Lambda^{-1} \eta + \frac{N}{2} \log 2\pi - \frac{1}{2} \log |\Lambda| \right]}_{A(\eta, \Lambda)})$$

$$\Omega = \left\{ (\eta, \Lambda) : \begin{array}{l} \eta \in \mathbb{R}^p, \Lambda \succ 0, \\ \Lambda \in \mathbb{R}^{p \times p}, \Lambda = \Lambda^T \end{array} \right\} \quad A(\eta, \Lambda) \text{ valid for } (\eta, \Lambda) \in \Omega$$

useful exercise: $\nabla_{\eta} A(\eta, \Lambda) = \mathbb{E}[x] = \mu = \Lambda^{-1} \eta$

$$\nabla_{\Lambda} A(\eta, \Lambda) = \mathbb{E}\left[-\frac{1}{2} x x^T\right]$$

UGM viewpoint:
$$p(x; \eta, \Lambda) = \exp\left(-\frac{1}{2} \sum_{i,j} \Lambda_{ij} x_i x_j + \sum_i \eta_i x_i - A(\eta, \Lambda)\right)$$

$p \in \mathcal{J}(G)$ where $E \triangleq \{ \xi_{ij} \}$ s.t. $\xi_{ij} \neq 0$

zeros in precision matrix \Rightarrow cond. indep. properties $\textcircled{*}$

"Gaussian network"