

Lecture 22 - scribbles

Friday, November 18, 2016
13:23

today: - ML in graph model
- Bayesian methods
↳ model selection

Estimation in GM

DGM: parametric family $P_{\Theta} = \left\{ p_{\Theta}(x) = \prod_i p(x_i | x_{\pi_i}, \Theta_i) \right\}$
 DGM $G=(V,E)$ $\Theta = (\Theta_1, \dots, \Theta_M)$
 $\Theta_i \in \mathcal{C}_i = \mathcal{C}_1 \times \dots \times \mathcal{C}_M$
 $M=p$
 independent parameters

i.e. no tying of parameters

⇒ MLE decouples in p independent problems

$$p(\text{data} | \Theta) = \prod_{i=1}^n p(x^{(i)} | \Theta) = \prod_{i=1}^n \prod_{j=1}^p p(x_j^{(i)} | x_{\pi_j}^{(i)}, \Theta_j)$$

$$\log L = \sum_{j=1}^p \left(\sum_{i=1}^n \log p(x_j^{(i)} | x_{\pi_j}^{(i)}, \Theta_j) \right) = \sum_{j=1}^p f_j(\Theta_j)$$

example: for discrete h.v. ⇒ $\hat{\Theta}_j^{(MLE)} = \frac{\text{proportion of observations } x_j = k, x_{\pi_j} = \text{something}}{\text{# } (x_{\pi_j} = \text{something})}$

⊗ latent variable (unobserved)

⇒ use EM

UGM: example for exponential family:

$$p(x|m) = \exp\left(\sum_c m_c^T T_c(x_c) - A(m)\right)$$

$$\downarrow$$

$$\exp(m_c^T T_c(x_c)) = \psi_c(x_c)$$

gradient ascent on log-likelihood

$$\frac{1}{n} \sum_{i=1}^n \log p(x^{(i)} | m) = \sum_c m_c^T \left(\frac{1}{n} \sum_{i=1}^n T_c(x_c^{(i)}) \right) - \frac{1}{n} A(m)$$

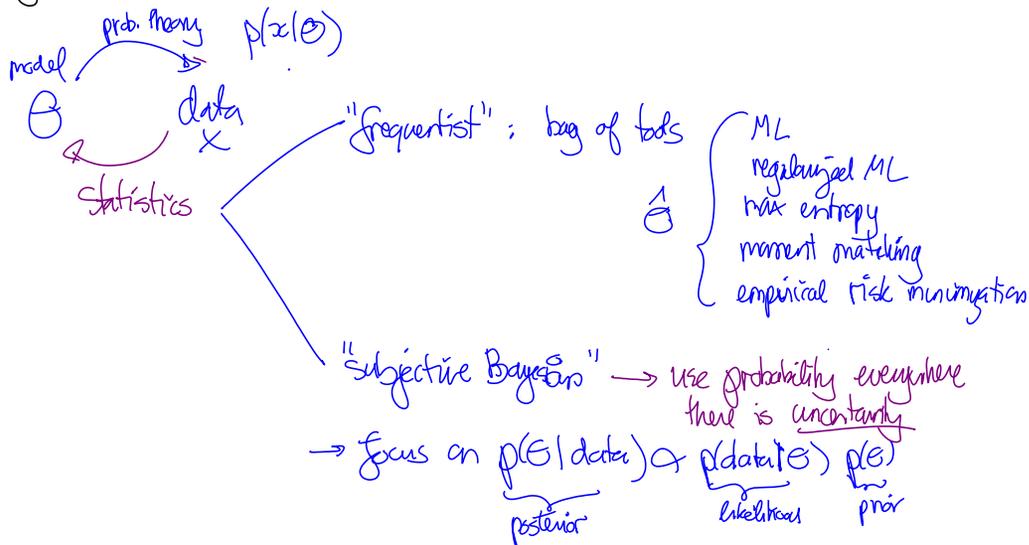
$$\nabla_{m_c} (\quad) = \hat{m}_c - \underbrace{E_{p(x|m)} [T_c(x_c)]}_{\text{(empirical moment)}}$$

to compute this, need inference $M_C(n)$

perhaps use approximate inference \leftarrow variational sampling

Ising model
 e.g. $T_{ij}(x_i, x_j) \in \{0, 1\}$
 $= x_i x_j$
 $M_{ij} = P(x_i=1, x_j=1 | n)$

Bayesian methods:



conclusion: Bayesian is "optimist": think you can get good models \Rightarrow obtain a method by doing probabilistic inference in model

frequentist is "pessimist": \rightarrow use analysis tools

Example: biased coin:

Bayesian model $x_i \in \{0, 1\}$
 $x_i | \theta \sim \text{Bernoulli}(\theta)$ $p(x_i | \theta) = \theta^{x_i} (1-\theta)^{1-x_i}$

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

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

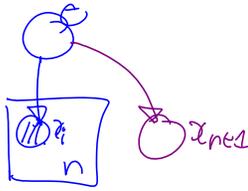
$= \theta^{\sum_{i=1}^n x_i} (1-\theta)^{n - \sum_{i=1}^n x_i} \mathbb{1}_{[0,1]}(\theta)$

\uparrow Beta(α, β) where $\alpha = n_{\uparrow} + 1$
 $\beta = n - n_{\uparrow} + 1$

questions: what is prob. of next flip = 1?

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

Bayesian, integrates out uncertainty



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

↑
predictive distribution

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

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

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

$n \rightarrow \infty$
 $p_n \rightarrow 1$

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

$$= \hat{\theta}_{ML} (1 - \hat{\theta}_{ML}) \xrightarrow{n \rightarrow \infty} 0$$

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

"Bernstein von-Mises thm"
→ "Bayesian CLT"

basically says that if prior put non-zero mass around the true model θ^* then posterior concentrates around θ^* as a Gaussian asymptotically

recall from hwk 1: multinomial model

$$X | \theta \sim \text{Mult}(\theta, 1) \quad \text{where } \theta \in \Delta_K$$

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

putting Dirichlet prior over θ ,

$$\theta \sim \text{Dir}(\alpha)$$

we also get Dirichlet posterior $\theta | \text{data} \sim \text{Dir}\left(\sum x_i + n \mathbf{1}_{K+1}\right)$

→ Dirichlet family is a "conjugate prior" for Multinomial likelihood

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

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

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

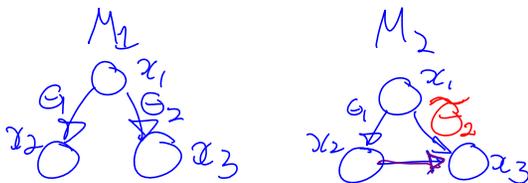
i.e. $\exists \alpha'$ s.t. $p(\epsilon|x, \alpha) = p(\epsilon|\alpha')$

$$\frac{p(z|\epsilon) p(\alpha)}{\int_{\epsilon} p(z|\epsilon) p(\alpha) d\epsilon}$$

side note: if use conjugate family pairs in a DBM; then Gibbs sampling can be easy

Model selection:

say want to choose between 2 DBM



(note here $M_1 \subseteq M_2$)

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

$\hat{\theta}_{M_2}^1 = \text{argmax}_{\theta_1, \tilde{\theta}_2} \log p(\text{data} | \theta_1, \tilde{\theta}_2, \text{model} = M_2)$
different space

if compare $\log p(\text{data} | \hat{\theta}_{M_1}, M = M_1)$ vs $\log p(\text{data} | \hat{\theta}_{M_2}, M = M_2)$?

since here $M_1 \subseteq M_2 \Rightarrow \text{LHS} \leq \text{RHS}$

i.e. cannot use likelihood directly to choose model

\rightarrow instead, use cross-validation (i.e. $\log p(\text{test data} | \hat{\theta}_{M_i}^{ML}(\text{train data}), M = M_i)$...)

Bayesian alternative \Rightarrow

true Bayesian \rightarrow sum over models (integrate out uncertainty)

introduce prior over models $p(M)$

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

data

$$p(\theta | D, M) p(M | D)$$

posterior on parameter gives data and model

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

↑ model averaging
 ↓ standard Bayesian predictive dist. for one model

⊗ in model selection forced to pick one model ; pick model which maximizes $p(M | \text{data}) \propto p(\text{data} | M) p(M)$

"marginal likelihood"

$$p(\text{data} | M) = \int_{\theta} p(\text{data}(\theta, M)) p(\theta | M) d\theta$$

to compare two models, look at

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

prior ratio

Bayes factor

choosing between k models M_1, \dots, M_k

pick M_i which maximizes $p(\text{data} | M=M_i)$

"empirical Bayes"

"type II ML"

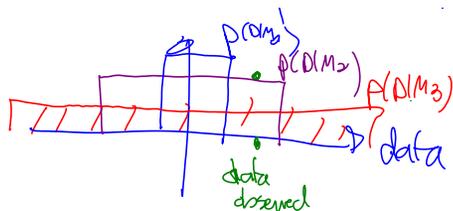
is this a good thing?

it is better than ML because $p(\text{data} | M=M_i)$ is normalized over the possible datasets

[vs. $p(\text{data} | \hat{\theta}_{ML}(\text{data}), M=M_i)$ which could be high on all datasets]

counter: suppose $M_1 \subseteq M_2 \subseteq M_3$ (increasing complexity)

$$p(D | M=M_i) = \text{uniform on same support}$$

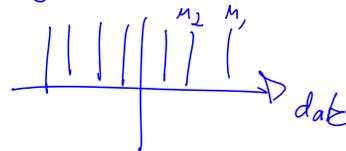


(pick intermediate M_2)

depending on M_0

"type II ML" is less prone to overfitting than standard ML

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



how to compute Marginal likelihood?

usually need approximations ← variational inf. sampling

Bayesian information criterion

is a (rough) approximation of $\log p(\text{data} | M)$

$$\approx \log p(\text{data} | \hat{\Theta}_{ML}, M) - \frac{d}{2} \log(n)$$

$\nearrow \text{dim}(\Theta_M)$
 complexity penalty

use Laplace approximation

$$p(D|M) = \int_{\Theta} \prod_{i=1}^n p(x_i | \Theta, M) p(\Theta|M) d\Theta$$

$$\approx \int_{\Theta} \exp(-nh(\Theta)) d\Theta$$

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

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

2 approximations → keep only term which grows with n } get BIC
 replace $\hat{\Theta}_{MAP}$ by $\hat{\Theta}_{ML}$

BIC is "consistent"