

Model choice and decision theory

Decision theory

Loss: $L(a, \theta)$ where a =action; θ =state of nature

Bayesian decision theory:

$$\min_a \left\{ \bar{L}(a) = E_{\theta|D} [L(a, \theta)] = \int L(a, \theta) p(\theta | D) d\theta \right\}$$

note separation of Loss function from
posterior/likelihood!

Profit function is the natural loss for marketing
applications!

Model Selection

We are often faced with the problem of selection from a set of models. The Bayes solution is to compute the posterior probability of each model.

For the set of models: M_1, \dots, M_k

compute:

$$p(M_i|y) = \frac{p(y|M_i)p(M_i)}{p(y)}$$

Posterior
odds
Ratio

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

= Bayes Factor \times Prior Odds

Model Probabilities cont.

For parametric models,

$$p(y|M_i) = \int p(y|\theta, M_i) p(\theta|M_i) d\theta$$

Depends on the prior! It should. One interpretation is that the model prob is the average of the “likelihood” wrt to the prior. Note the likelihood must include all normalizing constants!

$$\ell^*(y|M_i) = \mathbf{E}_{\theta|M_i} \left[\ell(\theta|y, M_i) \right]$$

Bayes Factors

$$\text{BF} = \frac{\int p(y|\theta, M_1) p(\theta|M_1) d\theta}{\int p(y|\theta, M_2) p(\theta|M_2) d\theta} = \frac{p(y|M_1)}{p(y|M_2)}$$

BF is a ratio of predictive densities. There is no distinction between nested and non-nested models!

BF depends critically on prior. Must be carefully chosen. Diffusion matters!

As prior diffusion increases for one model, BF moves away from that model.

Improper Priors are dangerous!

Model Probabilities cont.

The marginal density of the data is also the normalizing constant for the posterior.

$$p(\theta|y, M_i) = \frac{\ell(\theta|y, M_i)p(\theta|M_i)}{p(y|M_i)}$$

The numerator above is the **un-normalized posterior**. This we can always evaluate. The marginal density of the data is not always easy!

$$p(y|M_i) = \int \tilde{p}(\theta|y, M_i) d\theta = \frac{\tilde{p}(\theta|y, M_i)}{p(\theta|y, M_i)}$$

Savage-Dickey Conjugate setting

$M_0 : \phi_1 = \phi_1^h, \quad M_1 : \text{unrestricted}$ where $\phi' = (\phi'_1, \phi'_2)$.

$$\text{BF} : \frac{p(y|M_0)}{p(y|M_1)} = \frac{\int \ell(\phi_2|y)p(\phi_2)d\phi_2}{\iint \ell(\phi_1, \phi_2|y)p(\phi_1, \phi_2)d\phi_1d\phi_2}$$

$$\ell(\phi_2|y) = \ell(\phi_1, \phi_2|y) \Big|_{\phi_1=\phi_1^h}$$

a “natural” choice for prior

$$p(\phi_2 | \phi_1 = \phi_1^h) = \frac{p(\phi_1, \phi_2)}{\int p(\phi_1, \phi_2) d\phi_2} \Big|_{\phi_1=\phi_1^h}$$

Savage-Dickey Conjugate setting

$$\begin{aligned} \text{BF} &= \frac{\int \ell(\phi_1, \phi_2 | y) (p(\phi_1, \phi_2) / p(\phi_1)) d\phi_2 \Big|_{\phi_1 = \phi_1^h}}{\iint \ell(\phi_1, \phi_2 | y) p(\phi_1, \phi_2) d\phi_1 d\phi_2} \\ &= \frac{\int p(\phi_1, \phi_2 | y) d\phi_2 \Big|_{\phi_1 = \phi_1^h}}{p(\phi_1)} \\ &= \frac{p(\phi_1 | y) \Big|_{\phi_1 = \phi_1^h}}{p(\phi_1)} = \frac{\text{marginal posterior}}{\text{marginal prior}} \end{aligned}$$

Asymptotic methods (Laplace)

$$p(y | M_i) = \int \exp(\Gamma(\theta)) d\theta \quad \Gamma(\theta) = \log(\tilde{p}(\theta|y))$$

$$\approx \int \exp\left(\Gamma(\tilde{\theta}) - \frac{1}{2}(\theta - \tilde{\theta})' H(\tilde{\theta})(\theta - \tilde{\theta})\right) d\theta$$

$$= \exp(\Gamma(\tilde{\theta})) (2\pi)^{p_i/2} |H(\tilde{\theta})|^{-1/2} \quad \text{where } H(\tilde{\theta}) = -\frac{\partial^2 \Gamma(\theta)}{\partial \theta \partial \theta'} \Big|_{\theta=\tilde{\theta}}$$

$$\frac{\partial^2 \Gamma(\theta)}{\partial \theta \partial \theta'} = \frac{\partial^2 \log(p(\theta))}{\partial \theta \partial \theta'} + \frac{\partial^2 \log(p(y | \theta))}{\partial \theta \partial \theta'}$$

Asymptotic methods (Laplace)

Laplace approximation depends on:

prior evaluated @ posterior mode

prior curvature at mode

“highest” quality asymptotic approximation

Note: posterior can be easier to maximize than likelihood.

Tails on informative priors can deal with unidentified parameters and non-existence of MLEs.

Asymptotic methods (Laplace)

Laplace :

$$\text{BF} \approx (2\pi)^{(p_1-p_2)/2} \frac{p_1(\tilde{\theta})p_1(y|\tilde{\theta}_1)|H_1(\tilde{\theta})|^{-1/2}}{p_2(\tilde{\theta})p_2(y|\tilde{\theta}_2)|H_2(\tilde{\theta})|^{-1/2}}$$

number of parms in each
model

Asymptotic methods (BIC)

Alternative(dumber): expand about the MLE and throw out terms that don't depend on n.

$$\begin{aligned} p(y | M_i) &\approx \exp\left(\Gamma\left(\hat{\theta}_{\text{mle}}\right)\right) (2\pi)^{p_i/2} \left| \text{Inf}\left(\hat{\theta}_{\text{mle}}\right) \right|^{-1/2} \\ &= \exp\left(\Gamma\left(\hat{\theta}_{\text{mle}}\right)\right) (2\pi)^{p_i/2} n^{-p_i/2} \left| \text{Inf}_i\left(\hat{\theta}_{\text{mle}}\right) / n \right|^{-1/2} \\ &\approx \tilde{p}\left(y | \hat{\theta}_{\text{MLE}}, M_i\right) n^{-p_i/2} \quad \text{as } n \rightarrow \infty \end{aligned}$$

goes to constant matrix as n goes to infinity

Asymptotic methods (BIC)

$$\text{BIC} : \log(p(M_i | y)) \approx \log(\ell_i(\hat{\theta})) - \frac{p_i}{2} \log(n)$$

or

$$\text{BF} \approx \log(\text{LR}_{1,2}) - \frac{\Delta p}{2} \log(n)$$

all that remains
of the prior!

ridiculously inaccurate!

useful in two ways:

- 1). pick model with highest prob
- 2). establishes intuition that BF “automatically penalizes for “overfitting”

Computing Model Probs

For non-conjugate problems, there are three approaches:

1. Importance Sampling
2. Use of MCMC draws (NR)
3. Chib's Method (useful for latent var models)

BF using MCMC draws

$$\begin{aligned}\int \frac{q(\theta)}{\tilde{p}(\theta | y, M_i)} p(\theta | y, M_i) d\theta &= \int \frac{q(\theta)}{p(\theta | M_i) p(y | \theta, M_i)} p(\theta | y, M_i) d\theta \\ &= \frac{1}{p(y | M_i)} \int q(\theta) d\theta \\ &= \frac{1}{p(y | M_i)}\end{aligned}$$

$$E_{\theta | y, M_i} \left[\frac{q(\theta)}{\tilde{p}(\theta | y, M_i)} \right] = \frac{1}{p(y | M_i)}$$

Special case (Newton-Raftery)

If $q(\theta) = p(\theta | M_i)$,

$$p(y | M_i) = \frac{1}{\mathbb{E}_{\theta|y, M_i} \left[\frac{1}{\ell(\theta | M_i)} \right]}$$

LogMargDenNR

$$\hat{p}(y | M_i) = \frac{1}{\frac{1}{R} \sum_{r=1}^R \frac{1}{\ell(\theta^r | M_i)}}$$

Appeal: uses MCMC draws and likelihood

Problems: extreme sensitivity to outliers.

Distribution of likelihood values is non-uniform!

Marketing decisions and BDT

$$\max_x E[\pi(x | \Omega)]$$

There exists a range of marketing actions, x .

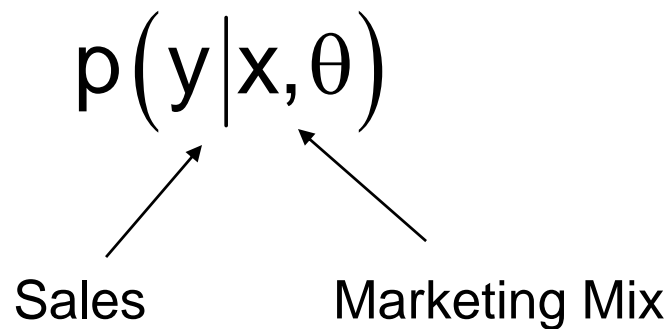
Actions are taken with respect to various information sets, Ω , reflecting what is known about model parameters.

Best action maximizes profits, π .

Loss/profit function

Canonical Example:

Sales Response Function



$$\max_x \mathbb{E}_{y|x} \left[\pi(y|x, \theta) \right]$$

Full Bayes

Sources of Uncertainty:

Distribution of sales given x , θ

Uncertainty in θ (summarized in posterior)

note: posterior is a means to the end of computing predictive distribution of sales

$$\begin{aligned}\pi^*(\mathbf{x}) &= \mathbf{E}_{\theta} \left[\mathbf{E}_{y|\mathbf{x},\theta} \left[\pi(y|\mathbf{x},\theta) \right] \right] \\ &\equiv \mathbf{E}_{\theta} \left[\bar{\pi}(\mathbf{x}|\theta) \right]\end{aligned}$$

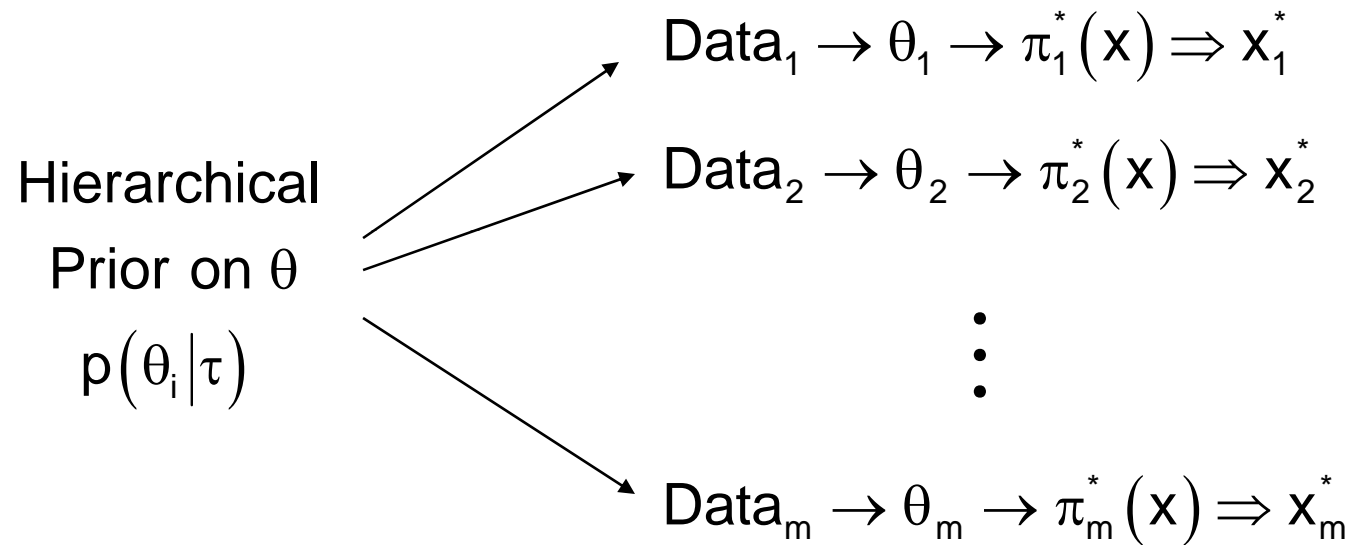
Why full Bayes?

Uncertainty in θ -- NO PLUG-In

Especially important with low amount of information –
disaggregate data

$$\pi^*(\mathbf{x}) = \mathbf{E}_{\theta} \left[\bar{\pi}(\mathbf{x}|\theta) \right] \neq \bar{\pi}(\mathbf{x}|\theta) \Big|_{\theta=\hat{\theta}}$$

Disaggregate decisions



Information: 1. “Borrowed” via Prior 2. Data on Individual Units

Information/customization

Decisions customized to the individual units.

- Product design

- Price/promotion

Extent of Customization *should* depend on:

- Differences in units

- Uncertainty in θ at individual level

 - e.g. with very noisy data, you should be back at uniform decision!


Value of disaggregate information

Compare Profits based on disaggregate with aggregate decisions

Disaggregate Profits:

$$\Pi_{\text{disagg}} = \sum_i \pi_i(x_i^*)$$

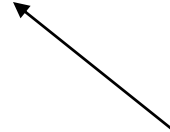
Optimal
Decision for
ith unit



Aggregate Profits:

$$\Pi_{\text{agg}} = \sum_i \pi_i(x^*)$$

Optimal
Uniform
Decision



Value of disaggregate information

Valuation through comparison of:

“Aggregate Profits” -- profits from “uniform” rule or same action for every unit

“Disaggregate Profits” – profits from full customization allowing for optimal Bayes decisions at each unit

Shows the applied orientation of marketing as a decision and information oriented discipline

Ex: Value of HH Purchase Info

Rossi et al (1995).

What is the value of household level purchase information?

Requires a metric.

Use profits from customized couponing. Print coupons whose face value is customized to the consumer's preferences – to the extent we can infer these using limited data.

Unit-Level Model and Heterogeneity

Diagonal MNP

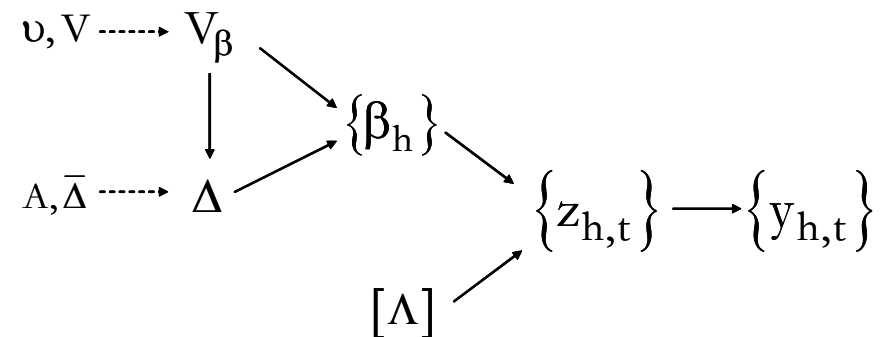
$$y_{h,t} = j \text{ if } \max(z_{h,t} = X_{h,t}\beta_h + \varepsilon_{h,t}) = z_{j,h,t}$$

$$\varepsilon_{h,t} \sim N(0, \Lambda)$$

Heterogeneity (observed demos and unobs)

$$\beta_h = \Delta' z_h + v_h$$

$$v_h \sim N(0, V_\beta)$$



Information Sets and Decisions

1. Full – purchase history information and demos
customized coupons – a function of hh data
2. Demos Only – only demo info
customized coupons – function of demos only
3. “Blanket” – no information about household
but info on distribution of parameters across
hhs.
“blanket” or uniform coupon

Predictive Distributions

Full Info Set:

$$p(\beta_h | \{y_1, \dots, y_h, \dots, y_H, X_1, \dots, X_h, \dots, X_H\}, Z) = \int p(\beta_h | y_h, X_h, \Delta, V_\beta) p(\Delta, V_\beta | \{y_1, \dots, y_h, \dots, y_H, X_1, \dots, X_h, \dots, X_H\}, Z) d\Delta dV_\beta$$

Demos Only Set:

$$\int p(\beta_h | z_h, \Delta, V_\beta) p(\Delta, V_\beta | \text{Info}) d\Delta dV_\beta$$

Blanket:

$$\int p(\beta_h | z_h, \Delta, V_\beta) p(z_h) p(\Delta, V_\beta | \text{Info}) dz_h d\Delta dV_\beta$$

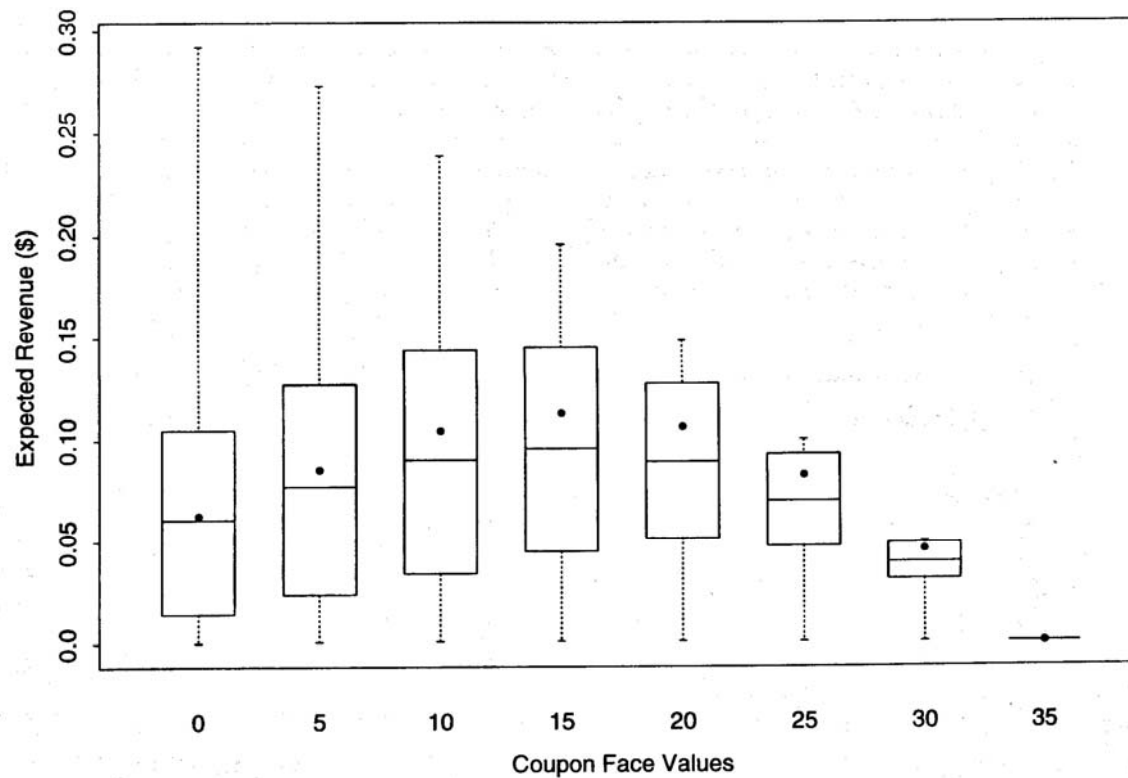
from a sample
of households

integrate out demos!

Couponing Problem

$$\max_F \pi(F) = \int \Pr[i|\beta_h, \Lambda, X(F)] (M - F) p(\beta_h, \Lambda | \Omega^*) d\beta_h d\Lambda$$

info set



Profit Results

Table 5 Relative Value of the Information Sets

Information Set	Net Revenue	Gain Relative to Blanket
Full	0.1570	2.55
Choices-Only	0.1529	1.93
One Obs	0.1500	1.56
Demos-Only	0.1467	1.12
Blanket	0.1459	1.0
No Coupon	0.1380	