

Today: Notes 3, get as far as we can

Monday: Lab, wrap up Notes 3 → Lab due Wed night
HW due Wed. night,
graded on completion

Wednesday: Quiz

I'll Post solutions on Thurs

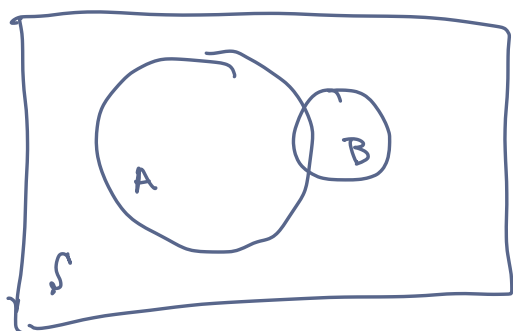
Tuesday: Stat speaker @ 4:15 / 4:30
→ EC opportunity!

03: BAYESIAN ESTIMATION

Larsen & Marx 5.8

Prof Amanda Luby

1 Bayes Theorem

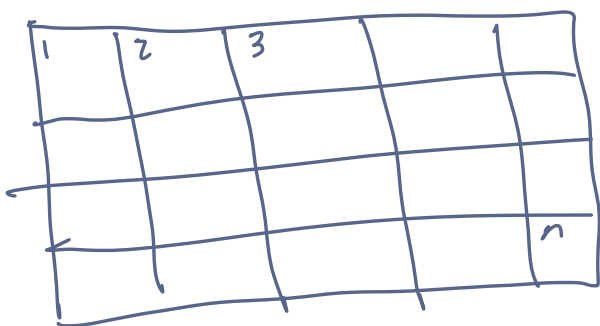


Idea: if you know $P(A|B)$, how can you find $P(B|A)$?
"inverse probability"

$$P(B|A) = \frac{P(A|B)P(B)}{P(A)} = \frac{P(A|B)P(B)}{P(A|B)P(B) + P(A|B^c)P(B^c)}$$

Bayesian statistics is a set of techniques that are based on inverse probabilities calculated using Bayes' theorem. Relative to "classical techniques" (MoM and MLE), Bayesian estimation provides a way to incorporate "prior knowledge" into the estimation of parameters.

Example: 1968 submarine went missing USS Scorpion



A_1 : Sub sunk in sec 1

\vdots

A_n : Sub sunk in sec n

Solicited $P(A_1) \dots P(A_n)$ from experts

① Idea: Pick largest A_i call A_k and search that one first

B_k : Sub would be found in k if k was searched
- function of water depth

B_k^c : k was searched & sub not found

2 key pieces:

① incorporate "prior knowledge"

$$\textcircled{2} P(A_k | B_k^c) = \frac{P(B_k^c | A_k) P(A_k)}{P(B_k^c | A_k) P(A_k) + P(B_k^c | A_k^c) P(A_k^c)}$$

② Mechanism to update $P(A_k)$ with new information

③ $P(\text{Sub sunk in } k \mid \text{not found in } k)$ becomes "updated" $P(A_k) \rightarrow P^*(A_k)$

④ renormalize $P(A_j)$ for $j \neq k \rightarrow P^*(A_j)$
Search largest $P^*(A_j)$ and repeat $P^{**}(A_k)$ etc...

Classical Statistics	Bayesian Statistics
<i>Probability</i> refers to limiting relative frequencies. Probabilities are objective properties of the real world.	<i>Probability</i> describes a degree of belief, not a limiting frequency. As such, we can make probability statements about lots of things, not just data which are subject to random variation. For example, I might say that "the probability that Albert Einstein drank a cup of tea on August 1, 1948 is .35". This does not refer to any limiting frequency. It reflects my strength of belief that the proposition is true.
<i>Parameters</i> are fixed, unknown constants, and the data we observe is random. Because they are constant, no useful probability statements can be made about parameters.	<i>Parameters</i> are random, and the data that we observe are fixed. We can therefore make probability statements about parameters.
Statistical procedures should be designed to have well-defined long-run frequency properties. For example, a 95% confidence interval should capture the true value of the parameter at least 95% of the time.	We make inferences about a parameter θ by producing a probability distribution for θ . Inferences, such as point estimates and interval estimates, may then be extracted from this distribution.

Bayesian inference is a controversial approach because it inherently embraces a subjective notion of probability. The field of statistics generally puts more emphasis on frequentist methods although Bayesian methods definitely have a presence.

2 Bayesian Inference

1. Prior distribution:

$f_{\theta}(\theta)$ $P_{\theta}(\theta)$ if discrete
degree of belief about θ before we
see any data
sub example: $P(A_k)$'s

2. Statistical model for data:

$f_x(x|\theta)$: belief about the data given a parameter θ
NOTE: $f_x(x;\theta)$ is different than $f_x(x|\theta)$

sub example: $P(B_k)$'s

3. Posterior distribution:

$f_{\theta|x}(\theta|x)$: updated belief about θ after
seeing our data

ex: $P(A_k|B_k^c) \rightarrow p^*$

if we see w_1, \dots, w_n replace $P_w(w|\theta)$ with $\prod_{i=1}^n P_w(w_i|\theta) = L(\theta, w)$

Posterior distribution

Let W be a statistic dependent on parameter θ . Call its pdf $f_w(w|\theta)$. Assume that θ is the value of a random variable Θ , whose prior distribution is denoted p_Θ if discrete and f_Θ if continuous. The posterior distribution of Θ given $W = w$ is:

$$f_{\theta|w} = \begin{cases} \frac{P_w(w|\theta) f_\theta(\theta)}{\sum_{-\infty}^{\infty} P_w(w|\theta) f_\theta(\theta)} & w \text{ discrete} \\ \frac{f_w(w|\theta) f_\theta(\theta)}{\int_{-\infty}^{\infty} f_w(w|\theta) f_\theta(\theta) d\theta} & w \text{ continuous} \end{cases}$$

If θ is discrete, replace integrals w/ sums and f_θ with p_θ

4. Posterior mean: Estimator: $\hat{\theta} = E(\theta|w)$

$$= \int_{-\infty}^{\infty} \theta \cdot f_{\theta|w}(\theta|w) d\theta$$

Example: Let $X_1, \dots, X_n \sim \text{Bernoulli}(\theta)$ and suppose that θ has the prior distribution $\theta \sim \text{Beta}(\alpha, \beta)$.

$$P(X_i = x) = \theta^x (1-\theta)^{1-x} \quad x = \{0, 1\}$$

$$\text{let } X = \sum X_i \quad X \sim \text{Bin}(n, \theta)$$

$$P(X=x) = \binom{n}{x} \theta^x (1-\theta)^{n-x}$$

$$f_\theta = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} \quad 0 \leq \theta \leq 1$$

Goal: find posterior distribution of $\theta|X$: $\frac{P_X(X|\theta) f_\theta(\theta)}{\int P_X(X|\theta) f_\theta(\theta) d\theta}$

$$\text{numerator: } P_X(X|\theta) f_\theta(\theta) = \binom{n}{x} \theta^x (1-\theta)^{n-x} \cdot \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

$$= \binom{n}{x} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \underbrace{\theta^{x+\alpha-1} (1-\theta)^{n-x+\beta-1}}_{\text{"kernel" of a Beta pdf}}$$

factor constant out of top & bottom

make denominator pdf of Beta($x+\alpha, n-x+\beta$)

$$f_{\theta|X} = \frac{\binom{n}{x} \Gamma(\alpha+\beta)}{\Gamma(\alpha) \Gamma(\beta)} \cdot \frac{\theta^{x+\alpha-1} (1-\theta)^{n-x+\beta-1}}{\int \theta^{x+\alpha-1} (1-\theta)^{n-x+\beta-1} d\theta}$$

$$= \frac{\Gamma(n+\alpha+\beta)}{\Gamma(x+\alpha)\Gamma(n-x+\beta)} \theta^{x+\alpha-1} (1-\theta)^{n-x+\beta-1}$$

\Rightarrow Posterior is a Beta($x+\alpha, n-x+\beta$)

$$f_{\theta|X} = \frac{\Gamma(n+\alpha+\beta)}{\Gamma(x+\alpha)\Gamma(n-x+\beta)} \theta^{x+\alpha-1} (1-\theta)^{n-x+\beta-1}$$

What we did after recognizing the kernel of a Beta RV is mess w/ normalizing constant. If we recognize kernel, can always factor constant out of numerator & denominator, then multiply by a more useful constant in both numerator and denominator.

⇒ we don't have to go through the trouble if we recognize the kernel

$$\text{Shortcut: } f_{\theta|X} \propto f_{X|\theta}(X|\theta) f_{\theta}(\theta) \\ \propto \theta^{x+\alpha-1} (1-\theta)^{n-x+\beta-1}$$

↑

"proportional to"

"up to a normalization constant"

Our Bayes Estimator is theoretical mean of posterior

$$E(\theta|X)$$

$$\text{for Beta: } E(\theta|X) = \frac{x+\alpha}{n-x+\beta+x+\alpha} = \frac{x+\alpha}{n+\beta+\alpha}$$

Conjugate prior

Example: Let $X_1, \dots, X_n \sim N(\theta, \sigma^2)$ and suppose we take $\theta \sim N(a, b^2)$. For simplicity, let's assume σ^2 is known.