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
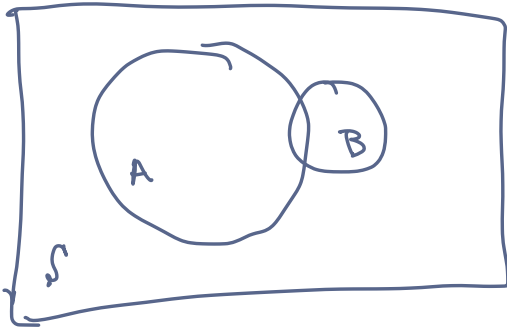Wednesday: Quiz                          I'll Post solutions on Tues

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

# 03: BAYESIAN ESTIMATION

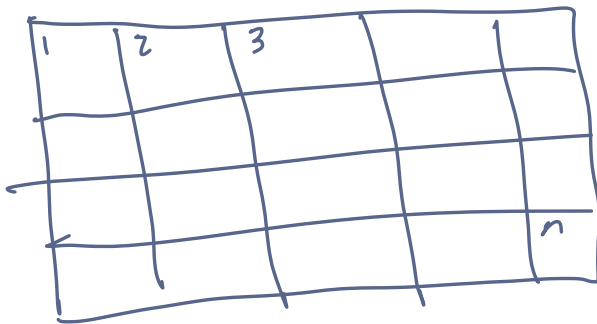Larsen & Marx 5.8

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 $P(A_1) \dots \sim P(A_n)$ 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"

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

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

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

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

| Classical Statistics | Bayesian Statistics |
| --- | --- |
| *Probability* refers to limiting relative frequencies. Probabilities are objective properties of the real world. | *Probability* 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. |
| *Parameters* are fixed, unknown constants, and the data we observe is random. Because they are constant, no useful probability statements can be made about parameters. | *Parameters* 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 $\theta$ by producing a probability distribution for $\theta$. 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:**

$$f_\theta(\theta) \qquad P_\theta(\theta) \text{ if discrete}$$

degree of belief about $\theta$ <u>before</u> we see any data

sub example : $P(A_k)$'s

2. **Statistical model for data:**

$$f_x(x|\theta) : \text{belief about the data given a parameter } \theta$$

NOTE: $f_x(x;\theta)$ is different than $f_x(x|\theta)$

sub example : $P(B_k)$'s

3. **Posterior distribution:**

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

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

If we see $w_1, ..., w_n$ replace $P_W(w|\theta)$ with $\prod_{i=1}^{n} P_W(w_i|\theta) = \mathcal{L}(\theta, w)$

---

**Posterior distribution**

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

$$f_{\theta|w} = \begin{cases} \dfrac{P_W(w|\theta) f_\theta(\theta)}{\int_{-\infty}^{\infty} P_W(w|\theta) f_\theta(\theta) d\theta} & w \text{ discrete} \\[4mm] \dfrac{f_W(w|\theta) f_\theta(\theta)}{\int_{-\infty}^{\infty} f_W(w|\theta) f_\theta(\theta) d\theta} & w \text{ continuous} \end{cases}$$

If $\theta$ is discrete, replace integrals w/ sums and $f_\theta$ with $P_\theta$

---

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

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

---

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

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

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

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

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

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

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

$= \underbrace{\binom{n}{x} \dfrac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)}}_{\substack{\text{factor constant} \\ \text{out of top \& bottom}}} \underbrace{\theta^{x+\alpha-1} (1-\theta)^{n-x+\beta-1}}_{\text{"kernel" of a Beta pdf}}$

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

$f_{\theta|x} = \dfrac{\dfrac{\binom{n}{x}\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)}}{\dfrac{\binom{n}{x}\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)}} \cdot \dfrac{\theta^{x+\alpha-1}(1-\theta)^{n-x+\beta-1}}{\int \theta^{x+\alpha-1}(1-\theta)^{n-x+\beta-1} d\theta}$   3

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

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

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

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

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.

$\Rightarrow$ we don't have to go through the trouble if we recognize the kernel

Shortcut : $f_{\theta|X} \propto f_{X|\theta}(X|\theta) f_{\theta}(\theta)$

$$\propto \theta^{X+\alpha-1}(1-\theta)^{n-X+\beta-1}$$

$\uparrow$

"proportional to"

"up to a normalization constant"

Our Bayes Estimator is theoretical mean of
$\underline{\qquad\qquad\qquad\qquad\qquad\qquad\qquad}$
posterior

$E(\theta|X)$

for Beta: $E(\theta|X) \cdot \dfrac{X+\alpha}{n-X+\beta+X+\alpha} = \dfrac{X+\alpha}{n+\beta+\alpha}$

Today:
- Wrap up NOTES 3
- Overview of Quiz Expectations
- Lab 02

- Quiz 1 on Wed
- HW 3 & lab 2 due Wed night
- Post solutions on Tues
- Colloquium speaker @ 4:15
- Moving Wed OH to Tuesday 2:30-4

## Recap

posterior distribution:

$$f_{\theta|x} \propto f_{x|\theta}(x|\theta) f_\theta(\theta)$$

If $X_1, \ldots, X_n$ is our data, replace $f(x|\theta)$
with $f(x_1, \ldots, x_n) = \prod f(x_i|\theta) = L(x_i|\theta)$

Then

$$f(\theta|x^n) = \frac{f(x^n|\theta) f(\theta)}{\int f(x^n|\theta) f(\theta)\, d\theta} = \frac{L_n(\theta) f(\theta)}{\int L_n(\theta) f(\theta)\, d\theta} \overset{\text{\propto}}{\propto} L_n(\theta) f(\theta)$$

constant that does not depend on $\theta$

**Conjugate prior**

When the prior and posterior are in the same family of distributions (same name) we say the prior is conjugate for that likelihood

Ex: Beta is the conjugate prior for binomial likelihood

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

parameters        hyper parameters

Goal: find posterior distribution $\theta | X^n$

$$f_{\theta|X}(\theta|X) \propto f_{X^n}(X|\theta) f_\theta(\theta) = \left[\prod \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(X_i - \theta)^2}{2\sigma^2}\right)\right] \cdot \frac{1}{\sqrt{2\pi b^2}} \exp\left[\frac{(\theta - a)^2}{2b^2}\right]$$

$$\propto \exp\left(-\frac{1}{2\sigma^2}\sum(X_i - \theta)^2\right) \exp\left(-\frac{1}{2b^2}(\theta - a)^2\right)$$

$$= \exp\left(-\frac{1}{2\sigma^2}\sum(X_i - \theta)^2 - \frac{1}{2b^2}(\theta - a)^2\right)$$

$$= \exp\left(-\frac{\sum X_i^2}{2\sigma^2} + \frac{2\theta\sum X_i}{2\sigma^2} - \frac{n\theta^2}{2\sigma^2} - \frac{\theta^2}{2b^2} + \frac{2a\theta}{2b^2} - \frac{a^2}{2b^2}\right)$$

$$= \exp\left(\theta^2\left(\frac{-n}{2\sigma^2} - \frac{1}{2b^2}\right) + \theta\left(\frac{2\sum X_i}{2\sigma^2} + \frac{2a}{2b^2}\right) + \left(\frac{\sum X_i^2}{2\sigma^2} - \frac{a^2}{2b^2}\right)\right)$$

From here, want to show the above expression can be written as

$$\exp\left(-\frac{1}{2\sigma_*^2}(\theta - \theta_*)^2\right) = \exp\left(-\frac{1}{2\sigma_*^2}(\theta^2 - 2\theta\theta_* + \theta_*^2)\right)$$

Equate like terms:

$$-\frac{\theta^2}{2\sigma_*^2} = -\theta^2\left(\frac{n}{2\sigma^2} + \frac{1}{2b^2}\right)$$

$$\frac{1}{\sigma_*^2} = \left(\frac{nb^2 + \sigma^2}{\sigma^2 b^2}\right)$$

$$\boxed{\sigma_*^2 = \frac{\sigma^2 b^2}{nb^2 + \sigma^2}}$$

$$\frac{2\theta\theta_*}{2\sigma_*^2} = \theta\left(\frac{\sum X_i}{\sigma^2} + \frac{a}{b^2}\right)$$

$$\frac{\theta_*}{\sigma_*^2} = \frac{\sum X_i}{\sigma^2} + \frac{a}{b^2}$$

$$\theta_* = \sigma_*^2\left(\frac{\sum X_i}{\sigma^2} + \frac{a}{b^2}\right)$$

$$\theta_* = \frac{\sigma^2 b^2 \sum X_i}{\sigma^2} + \frac{\sigma^2 b^2 a}{(nb^2 + \sigma^2)b^2}$$

$$= \frac{b^2 \sum X_i}{nb^2 + \sigma^2} + \frac{\sigma^2 a}{nb^2 + \sigma^2} = \boxed{\frac{b^2 \sum X_i + \sigma^2 a}{nb^2 + \sigma^2}}$$

$$\Rightarrow \theta|X \sim N(\theta_*, \sigma_*^2) = N\left(\boxed{\frac{b^2 \sum X_i + \sigma^2 a}{nb^2 + \sigma^2}, \frac{\sigma^2 b^2}{nb^2 + \sigma^2}}\right)$$

bayes estimator: $E(\theta|X) = \frac{b^2}{b^2 + \sigma^2/n}\bar{X} + \frac{\sigma^2/n}{b^2 + \sigma^2/n} \cdot a$

weighted average of sample mean (also MLE) $\bar{X}$ and prior mean $a$. As $n \to \infty$, $E(\theta|X) \to \bar{X}$

# Quiz 1

- Terminology
  - parameter us estimator
  - estimator vs estimate
  - pdf vs likelihood
  - sample vs population
  - prior vs posterior

- Basic Integration
  - polynomial: $x^2 + x + c$
  - $e^x$
  - $\ln(k)$
  - simple chain rules $\int \sin(2x)$

- Given a pdf, find Expected value & variance
  - recognize named distributions
  - use properties of expected value + variance
  - simple integrations

- Given $E(\hat{\theta})$ and $V(\hat{\theta})$
  - comment on bias + efficiency

- Estimation
  - MLE - set up + find
  - MOM
  - Bayes Estimator : multiply $L_n(\theta) \cdot f_\theta(\theta)$
    use kernel to recognize named posterior

Quiz :
~30 minutes

mix of concepts
and mechanics

OH
today 11:30-12:30
tues 2:30-4