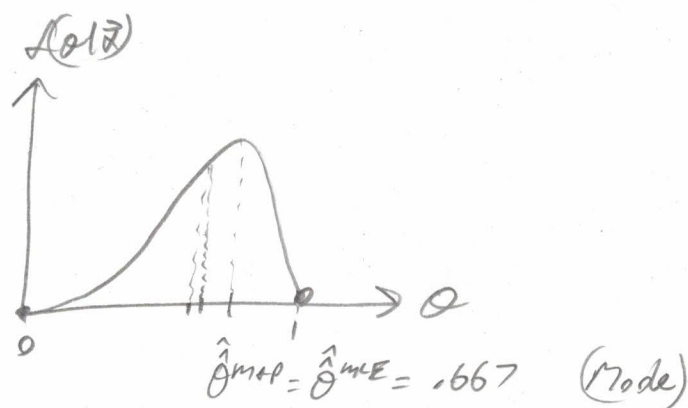


Lee 17 MATH 341/641

OBP: $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bern}(\theta)$, $f(\theta) = \mathcal{U}(\theta, 1)$, Lyphric Prior of ~~Ind~~ ~~Henne~~

$$\vec{x} = \langle 0, 1, 1 \rangle \Rightarrow f(\theta | \vec{x}) = \text{Beta}(3, 2)$$



$$\hat{\theta}_{MME} = 0.6 \text{ (mean)}$$

$$\hat{\theta}_{MME} = .619 \text{ (median)}$$

Is there another reasonable point estimate?

let $\hat{\theta}_{MMAE} := \text{MED}[\theta | \vec{x}]$ which ^{is the θ which} ~~minimizes~~ $E[|0 - \theta| | \vec{x}]$

MMAE: "minimal mean absolute error"

see next page

There is no closed form formula for median of the beta dist so...

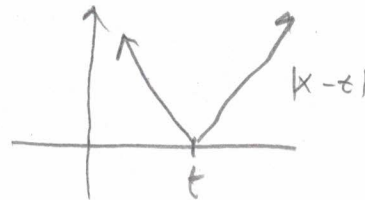
$$q_{\text{beta}}(\alpha, \beta) := \left\{ x : \int_0^x \frac{1}{B(\alpha, \beta)} t^{\alpha-1} (1-t)^{\beta-1} dt = I_x(\alpha, \beta) = \alpha \right\}$$

$$\hat{\theta}_{MMAE} = q_{\text{beta}}(0.5, \alpha, \beta) = q_{\text{beta}}(0.5, 3, 2) = .619$$

If $\vec{x} = \langle 0, 0, 0 \rangle$ $\hat{\theta}_{MMAE} = q_{\text{beta}}(0.5, 1, 1) = .159$, also reasonable for that scenario

This completes the discussion of point estimation in Bayesian Inference:
 $\hat{\theta}_{MAP}$, $\hat{\theta}_{MLE}$, $\hat{\theta}_{MMAE}$ are the 3 widely used estimators

Continuous r.v. X , continuous



$$\text{let } g(t) = E[|X - t|]$$

$$g'(t) = \frac{d}{dt}[E[|X - t|]] = E\left[\frac{d}{dt}[|X - t|]\right] = E[\mathbb{1}_{X \geq t} - \mathbb{1}_{X < t}]$$

$$= E[\mathbb{1}_{X \geq t}] - E[\mathbb{1}_{X < t}] = P(X \geq t) - P(X < t) \stackrel{\text{set}}{=} 0$$

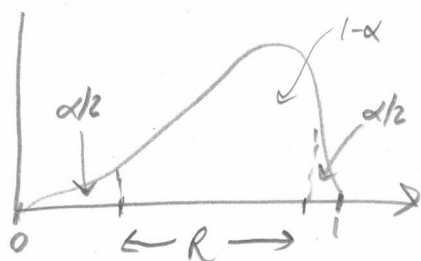
$$\Rightarrow P(X \geq t) = P(X < t)$$

$$\Rightarrow P(X \geq t) = 0.5 \text{ \& } P(X < t) = 0.5$$

$$\Rightarrow t = \text{MED}[X]$$

Let's now do goal #2: confidence sets. Now, we get what we want. let $R \subset \Theta$, a "region".

get $P(\theta \in R | \vec{x}) = 1 - \alpha$! How?



This is known as a 2-sided credible region:

$$CR_{\theta, 1-\alpha} := [Q[\theta | \vec{x}, \frac{\alpha}{2}], Q[\theta | \vec{x}, 1 - \frac{\alpha}{2}]]$$

A 95% CR would then be

$$CR_{\theta, 95\%} = [Q[\theta | \vec{x}, 25\%], Q[\theta | \vec{x}, 97.5\%]]$$

In our example $\vec{x} = \langle 9, 1 \rangle$, $p(\theta) = U(\theta, 1)$,

$$CR_{\theta, 95\%} = [\text{qbeta}(.025, 3, 2), \text{qbeta}(.975, 3, 2)]$$

$$= [0.194, 0.932] \quad \text{Makes sense! Why so wide? } n=3$$

For $\vec{x} = \langle 0, 0 \rangle$



$$CR_{\theta, 95\%} = [\text{qbeta}(.025, 1, 1), \text{qbeta}(.975, 1, 1)] = [0.006, 0.602]$$

Now let's do goal #3: hypothesis testing. Again we get what we want, we can compute $P(H_0|\vec{x})$, $P(H_1|\vec{x})$.

Define Bayesian p-value as $p_{\text{bc}} := P(H_0|\vec{x})$

Let's do one-sided tests first

$H_1: \theta > 0.5 \Rightarrow H_0: \theta \leq 0.5$ right-tail test

$$p_{\text{bc}} = P(H_0|\vec{x}) = P(\theta \leq 0.5|\vec{x}) = \int_0^{0.5} \text{beta}(3,2) d\theta = I_{0.5}(3,2)$$

regularized, incomplete beta function

In this class, use the R code notation $\text{pbeta}(x, \alpha, \beta) = I_x(\alpha, \beta)$

$p_{\text{bc}} = \text{pbeta}(0.5, 3, 2) = .3125 \Rightarrow$ Retain H_0 . Possible Type II error at $\alpha = 5\%$. Since underpowered at $n=3$.

Likewise the left-tailed test is...

$H_1: \theta < 0.5 \Rightarrow H_0: \theta \geq 0.5$

$$p_{\text{bc}} = P(H_0|\vec{x}) = P(\theta \geq 0.5|\vec{x}) = 1 - P(\theta \leq 0.5|\vec{x}) = 1 - \text{pbeta}(0.5, 3, 2) = .6875$$

\Rightarrow Retain H_0 .

at $\alpha = 5\%$

Okay let's do 2-sided tests:

$H_1: \theta \neq 0.5 \Rightarrow H_0: \theta = 0.5$

$p_{\text{bc}} = P(H_0|\vec{x}) = P(\theta = 0.5|\vec{x}) = 0!$ Since $\theta|\vec{x}$ is a continuous distn!

In fact 0.5 is a point in the support of $H_0: \theta = 0.5$ will have $p_{\text{bc}} = 0$

What happened here? Exactly what should happen. Recall is frequent hypothesis testing of $H_0: \theta = \theta_0$. If $n \rightarrow \infty \Rightarrow p \rightarrow 0$ and H_0 is always rejected? But the rejection may not be "practically significant" or "clinically significant". Who cares if $\theta = 0.500001$ for a coin? It's still fair for all "practical purposes". The Bayesian Framework forces you to make clear what practical significance means by forcing you to define L , a margin of equivalence. The hypotheses then become:

$$H_a: \theta \notin [\theta_0 \pm L] \Rightarrow H_0: \theta \in [\theta_0 \pm L]$$

$$\text{let } L = 0.02$$

$$\begin{aligned} p_{\text{val}} &= P(H_0 | \vec{x}) = P(\theta \in [0.5 \pm 0.02] | \vec{x}) = P(\theta \leq 0.52 | \vec{x}) - P(\theta \leq 0.48 | \vec{x}) \\ &= p_{\text{beta}}(0.52, 3, 2) - p_{\text{beta}}(0.48, 3, 2) = 0.060 \Rightarrow \text{Reject } H_0 \text{ at } \alpha = 5\% \end{aligned}$$

Be careful with L . If too small, you always reject!

For this reason, there is an alternate procedure for 2-sided testing to arrive at a decision

Retain H_0 ^{but} at α if $\theta_0 \in CR_{\theta, 1-\alpha}$ otherwise Reject H_0 .

In our case $0.5 \in [0.48, 0.52] \Rightarrow \text{Retain } H_0$.

Done with hypothesis testing!!

You can get a pvalue out of this by taking the maximum α s.t. θ_0 is in the region $CR_{\theta, 1-\alpha}$

Obp: $X_1, \dots, X_n \overset{iid}{\sim} \text{Bern}(\theta)$, $f(\theta) = U(0,1)$.

Imagine $n=3$ $\vec{x} = (x_1, x_2, x_3)$ and do the inference one data point at a time using the posterior from the previous posterior as the prior:

$$f(\theta|x_1) = \frac{P(x_1|\theta)f(\theta)}{P(x_1)} \propto P(x_1|\theta)f(\theta) = \theta^{x_1}(1-\theta)^{1-x_1} \mathbb{1}_{\theta \in (0,1)} \propto \text{Beta}(x_1+1, 1-x_1+1)$$

$$f(\theta|x_2) \propto P(x_2|\theta)f(\theta|x_1) \propto \theta^{x_2}(1-\theta)^{1-x_2} \theta^{x_1}(1-\theta)^{1-x_1} \mathbb{1}_{\theta \in (0,1)} \propto \text{Beta}(x_1+x_2+1, 2-(x_1+x_2)+1)$$

$$f(\theta|x_3) \propto P(x_3|\theta)f(\theta|x_2) \propto \theta^{x_3}(1-\theta)^{1-x_3} \theta^{x_1+x_2}(1-\theta)^{2-(x_1+x_2)} \mathbb{1}_{\theta \in (0,1)} \propto \text{Beta}(x_1+x_2+x_3+1, 3-(x_1+x_2+x_3)+1)$$

This process of iterative updating is not specific to this Obp with Laplace's Prior, but holds always (Hv).

What else did we learn? Beta prior \rightarrow Beta posterior.

Let's prove this in general:

if $\alpha = \beta = 1$

Obp: $X_1, \dots, X_n \overset{iid}{\sim} \text{Bern}(\theta)$, $f(\theta) = \text{Beta}(\alpha, \beta) = U(0,1)$ special case

The parameters on the prior are called "hyperparameters"

$$f(\theta|\vec{x}) \propto P(\vec{x}|\theta)f(\theta) \propto \left(\theta^{\sum x_i}(1-\theta)^{n-\sum x_i}\right) \left(\theta^{\alpha-1}(1-\theta)^{\beta-1} \mathbb{1}_{\theta \in (0,1)}\right) \\ \propto \text{Beta}(\alpha + \sum x_i, \beta + n - \sum x_i)$$

$$\Rightarrow \hat{\theta}_{MLE} = \frac{\alpha + \sum x_i - 1}{\alpha + \beta + n - 2}$$

$$\hat{\theta}_{MMSE} = \frac{\alpha + \sum x_i}{\alpha + \beta + n}$$

$$\hat{\theta}_{JMMSE} = \text{Beta}(0.5, \alpha + \sum x_i, \beta + n - \sum x_i)$$

We say the Beta distr is the "conjugate prior" for the iid Bernoulli DGP. Conjugacy means prior and posterior are same rv (but different parameter values updated by the data).

$$\text{Beta}(\alpha, \beta) \xrightarrow{\vec{x}} \text{Beta}\left(\underbrace{\alpha + \sum x_i}_{\# \text{ successes}}, \underbrace{\beta + n - \sum x_i}_{\# \text{ failures}}\right)$$

What is the interpretation of the values of the hyperparameters?

The prior's hyperparameters are like observing fake data, $n_0 = \alpha + \beta$ ^{where} # pseudodata

If we employ Laplace's prior of indifference $\mathcal{U}(0,1) = \text{Beta}(\alpha=1, \beta=1)$

$\Rightarrow n_0 = 2$ pseudodata. This is weird. Further, consider

$$\hat{\theta}_{\text{MISE}} = \frac{\alpha + \sum x_i}{\alpha + \beta + n} = \frac{\alpha}{\alpha + \beta} \cdot \frac{\alpha + \beta}{\alpha + \beta} + \frac{\sum x_i}{\alpha + \beta + n} \cdot \frac{n}{n}$$

$$= \frac{\alpha + \beta}{\alpha + \beta + n} \frac{\alpha}{\alpha + \beta} + \frac{n}{\alpha + \beta + n} \frac{\sum x_i}{n}$$

$$= \underbrace{\rho}_{\text{prior gain}} E[\theta] + (1 - \rho) \underbrace{\hat{\theta}^{\text{MLE}}}_{\text{pure data estimate}}$$

$$\lim_{n \rightarrow \infty} \hat{\theta}_{\text{MISE}} = \hat{\theta}^{\text{MLE}} \quad \text{since} \quad \lim_{n \rightarrow \infty} \rho = 0$$

$$\text{let } \rho := \frac{\alpha + \beta}{\alpha + \beta + n},$$

the shrinkage metric,
measures strength of
prior $\in [0,1]$ on the
most popular point estimate.

$$\text{eg. } \alpha = \beta = 1, n = 3 \Rightarrow \rho = \frac{2}{5}$$

\Rightarrow 40% weight on prior,
60% weight on data

$\hat{\theta}_{\text{MISE}}$ is called a "shrinkage estimator"

since it shrinks toward $E(\theta)$

