

Lesson 9 Math 341 3/9/17

11

Recall $X|D \sim \text{Bin}(n, \theta)$, n fixed, known

$D \sim \text{Beta}(\alpha, \beta)$

$D|X \sim \text{Beta}(\alpha+x, \beta+n-x)$

conjugacy

note on prior notes

$X^*|X \sim \text{Beta Binom}(n^*, \underbrace{\alpha+x+x^*}_{\alpha'}, \underbrace{\beta+n-x-x^*}_{\beta'}) := \binom{n^*}{x^*} \frac{\text{Beta}(\alpha+x+x^*, \beta+n-x-x^*)}{\text{Beta}(\alpha+x, \beta+n-x)}$

if $n^* = 1$

posterior α, β

$X^e|X \sim \text{Beta Binom}(1, \alpha+x+x^e, \beta+n-x-x^e) = \frac{\text{Beta}(\alpha+x+x^e, \beta+n-x-x^e)}{\text{Beta}(\alpha+x, \beta+n-x)} = \text{Beta}\left(\frac{\alpha+x}{\alpha+\beta+n}\right)$

we know this is Bernoulli... so set $x^e = 1$, find prob..

$$P(X^e=1|X) = \frac{\text{Beta}(\alpha+x+1, \beta+n-x)}{\text{Beta}(\alpha+x, \beta+n-x)} = \frac{\frac{\Gamma(\alpha+x+1) \Gamma(\beta+n-x)}{\Gamma(\alpha+\beta+n+1)}}{\frac{\Gamma(\alpha+x) \Gamma(\beta+n-x)}{\Gamma(\alpha+\beta+n)}} = \frac{(\alpha+x)}{(\alpha+\beta+n)}$$

$P(X^e|X)$ is dist of future X^* given data "posterior predictive dist"

$P(X)$ = dist of data observed but it could be length of n

$$= \int P(X|\theta) P(\theta) d\theta$$

\Rightarrow same form!

$$P(X|\epsilon_3)$$

$$\int P(X|\theta) P(\theta|\epsilon_3) d\theta$$

$X \sim \text{Beta Binom}(n, \alpha+x, \beta+n-x)$

AKA the "prior predictive dist"

$P(X)$... dist denominator here has another name!

(2)

If conjugate, $P(\theta)$, $P(\theta|x)$ belong to same family as well. Why?

Uniform prior: a prior which does not have a "large" effect on posterior

$$\theta \sim \text{Beta}(1, 1)$$

$$\theta|x \sim \text{Beta}\left(\overset{a}{1+x}, \overset{b}{1+n-x}\right) \Rightarrow \hat{\theta}_{\text{muse}} = \frac{x+1}{n+2} \quad \text{AKA Wilson estimate}$$

$\uparrow \quad \uparrow$
 $\text{Beta}(1,1)$ are thought "uniform" yielded probabilities of 1 success, 1 failure. That's not "no information".

What would NO info look like?

Uniform? YES

$$\theta|x \sim \text{Beta}\left(\overset{\alpha}{0}+x, \overset{\beta}{0}+n-x\right) \Rightarrow \hat{\theta}_{\text{muse}} = \frac{x}{n} = \hat{\theta}_{\text{MLE}}$$

$$\Rightarrow \theta \sim \text{Beta}(0, 0) \quad \text{what's wrong? } \alpha > 0, \beta > 0$$

this is an illegal prior! AKA "improper prior"

But posterior is proper if $x \neq 0$ and $x \neq n$

Tons of theory on this... some say improper priors okay, some say no.

Okay for us... except... you must be careful your posterior is proper!

$\theta \sim \text{Beta}(1, 1) \Rightarrow$ it's diffuse but successes and failures are known to be possible

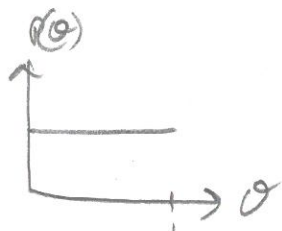
$\theta \sim \text{Beta}(0, 0) \Rightarrow$ successes and failures not known to be possible
AKA complete ignorance

Haldane prior (1932)

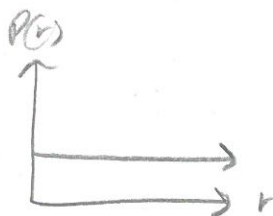
Another problem

$\theta \sim U(0,1)$ means every prob equally likely

What if I consider the odds, $r := \frac{\theta}{1-\theta}$. Am I indifferent on this scale?



$\theta \in (0,1)$



$r \in (0, \infty)$

No...
not possible!

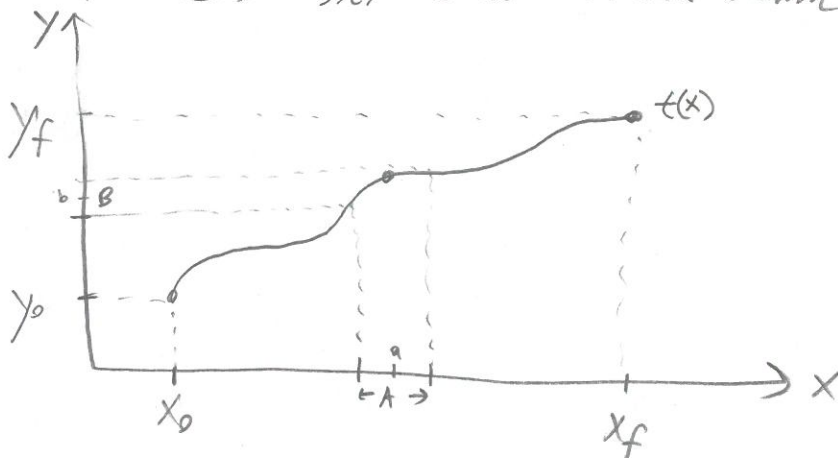
the principle of indifference has a problem!

What is PDF of R?

Math 621 covers transformation of variables.

Imagine r.v.'s X, Y with densities f_X, f_Y with f_X known, f_Y unknown

Suppose $Y = t(X)$ s.t. t is invertible function



$$P(X \in A) \approx f_X(a) A$$

$$P(Y \in B) \approx f_Y(b) B$$

$$P(X \in A) = P(Y \in B) \Rightarrow f_X(a) A = f_Y(b) B \quad \text{Since } A > 0, B > 0$$

let area be small

$$f_X(a) |dx| = f_Y(b) |dy| \quad \text{s.t. } b = t(a) \text{ or } a = t^{-1}(b)$$

$$\text{Solve for } f_Y(y) \Rightarrow f_Y(b) = f_X(a) \left| \frac{dx}{dy} \right|$$

let $b = y$ the same dummy variable

$$\Rightarrow f_Y(y) = f_X(t^{-1}(y)) \left| \frac{d}{dy} [t^{-1}(y)] \right|$$

$$\text{Supp}(X) = [x_0, x_f], \text{ Supp}(Y) = [y_0, y_f]$$

What is PDF of R ? Or $U(0,1) \Rightarrow f_\theta(\theta) = 1$ $R - RX = X \Rightarrow R = X + RX \Rightarrow R = (-X)X \Rightarrow$

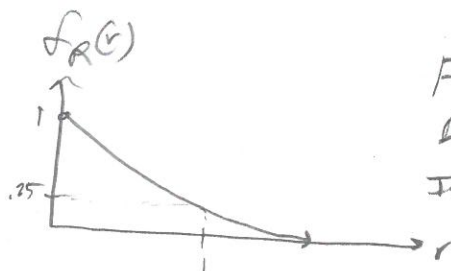
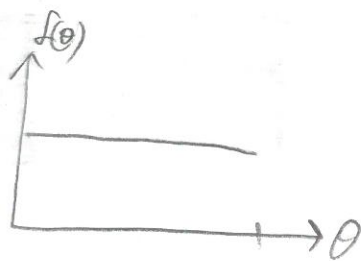
$$R = t(\theta) = \frac{\theta}{1-\theta} \Rightarrow \theta = t^{-1}(R) = \frac{R}{R+1}$$

reparametrization!

$$f_R(r) = \underbrace{f_\theta(t^{-1}(r))}_{=1} \left| \frac{d}{dr} [t^{-1}(r)] \right| = \left| \frac{(r+1) - (r+1)^2}{(r+1)^2} \right| = \frac{1}{(r+1)^2}$$

is this a density? $\int_{f_R(R)} f_R(r) dr = 1 \Rightarrow \int_0^\infty \frac{1}{(r+1)^2} dr = \left[\frac{r}{r+1} \right]_0^\infty = 1 \checkmark$

Now $f_R(r) \neq f_\theta(\theta) \Rightarrow$ If you are indifferent about θ , you are not indifferent about r 's

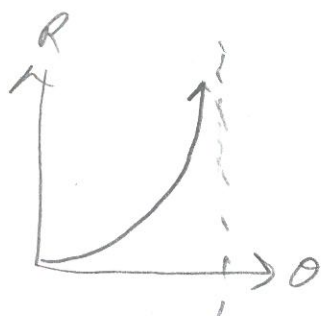


or any other monotone transformation of θ

Fisher used this to show Bayes is stupid!

Is it a problem if

the prior is discrete?



$$\theta \in [0, 0.5] \Rightarrow R \in [0, 1]$$

$$\theta \in [0.5, 1] \Rightarrow R \in [1, \infty)$$

$$r = \frac{0.1}{1-0.1} = \frac{1}{9}$$

$$\begin{aligned} f_\theta(\theta=0.1) &= 1 \\ f_R(\theta=0.1) &= \frac{1}{(1/9+1)^2} = 0.81 \neq 1 \end{aligned}$$

Is there a prior $\theta \sim p(\theta)$ s.t. any reparametrization will yield the same prior density?

$\theta \sim \text{Bern}(1,1)$
 $\theta \sim U(0,1)$ Laplace prior
 $\theta \sim \text{Beta}(0,0)$ Haldane prior

}

Uniform, Vague, Weak,
 Diffuse

⇓

don't affect inference too much.

Strategy
protocol

Is there a way to choose an uniform prior that would be the same under reparameterizations?

Likelihood Model

$p(x|\theta) \xrightarrow{\text{strategy}} \text{pick } p(\theta)$ and make a reparameterization:

$\phi = t(\theta)$: s.t. t is 1:1 and monotonic

$p(x|\phi) \xrightarrow{\text{strategy}} \text{pick } p(\phi)$

wouldn't it be nice if

$$p(\phi) = p(t^{-1}(\phi)) \left| \frac{d}{d\phi} [t^{-1}(\phi)] \right|$$

As in the strategy itself would ensure this "invariance" wouldn't break?

This is the strategy Jeffreys found ≈ 1930 's.

Before we get there, we need two pieces.

① "Kernels"

② Fisher Information

16

Keywords) Recall ...

$$P(\theta|x) = \frac{P(x|\theta) P(\theta)}{P(x)} \propto P(x|\theta) P(\theta) \quad \text{w/ ?} \quad P(x) \text{ not known of } \theta!$$

$$f(x; \theta) \propto g(x; \theta)$$

this can be diff shw $\exists c \in \mathbb{R}$ s.t.

$f(x; \theta) = \frac{1}{c} g(x; \theta)$ where $\frac{1}{c}$ is called $h(\theta)$

How to find c ?

Norm: $\int_{\text{supp}(x)} f(x) dx = 1$

$$\text{Now: } \frac{f(x_1; \theta)}{f(x_2; \theta)} = \frac{\cancel{h(\theta)}}{\cancel{h(\theta)}} \frac{g(x_1; \theta)}{g(x_2; \theta)}$$

$$\int_{\text{supp}(f)} g(x) dx = \int c f(x) dx = c \underbrace{\int f(x) dx}_1 = c = \int g(x) dx$$

Note: $\int g(x) dx < \infty$ and $\int g(x) dx > 0$ Note: $g(x), f(x)$ are 1:1

So little is

r.v.b can be changed by other
kernels

Kenil

$$p(x) \propto \binom{n}{x} \theta^x (1-\theta)^{n-x} \frac{1}{\beta \alpha \Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} \propto \theta^x (1-\theta)^{n-x} \theta^{\alpha-1} (1-\theta)^{\beta-1} = \theta^{x+\alpha-1} (1-\theta)^{n-x+\beta-1} \propto \text{Beta}(x+\alpha, n-x+\beta)$$

Wz? $\mathcal{O}_{\text{Ben}}(2\beta) := \frac{1}{(2\beta)!} \mathcal{O}^{\alpha,1} (1-\theta)^{b-1} \propto \mathcal{O}^{\alpha,1} (1-\theta)^{b-1} = \underbrace{\theta^9 (1-\theta)^6}_{\text{Ben}}$

the kernel of
the lion

$$X| \theta \sim \text{Bin}(n, \theta) = \binom{n}{x} \theta^x (1-\theta)^{n-x} = \frac{n!}{x!(n-x)!} \theta^x (1-\theta)^n (1-\theta)^{-x}$$

$$\propto \frac{1}{x!(n-x)!} \left(\frac{\theta}{1-\theta}\right)^x \leftarrow \text{kernel of the binomial!!}$$

If you have a score function where $P(\theta|x) \propto \text{kernel} \dots$ you're done! EASIER than solving explicitly.

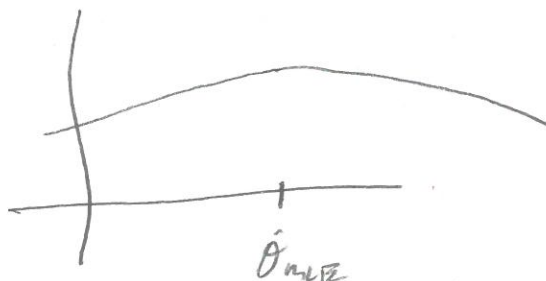
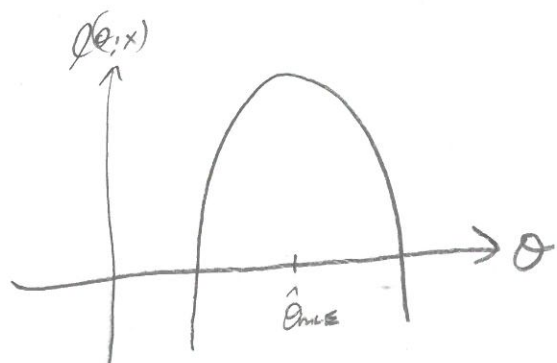
7

Fisher Info

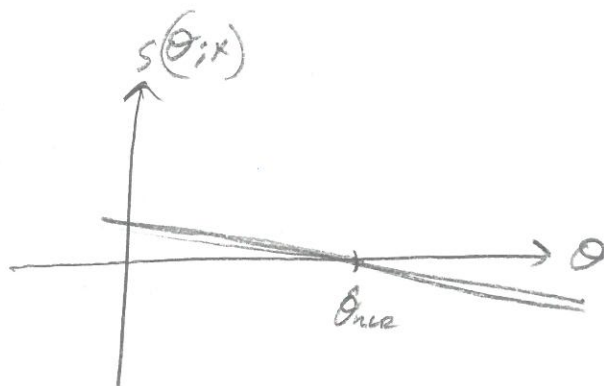
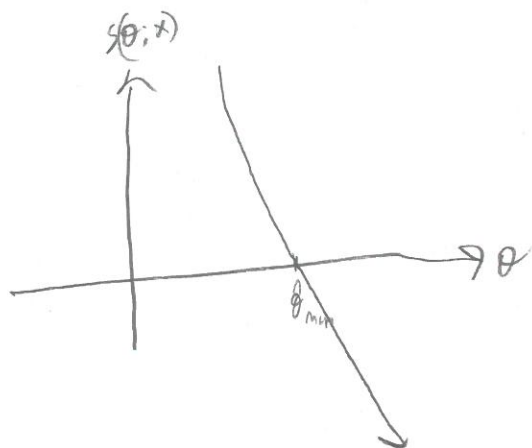
Real likelihood

and log likelihood

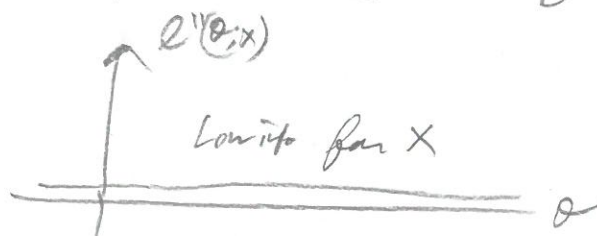
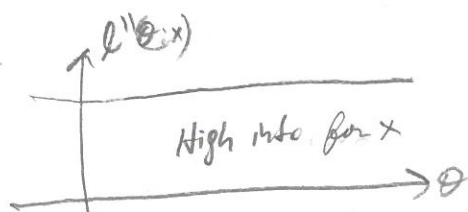
$$L(\theta; x) = P(x; \theta) \quad \ell(\theta; x) := \ln(L(\theta; x))$$



Define $s(\theta; x) := \ell'(\theta; x)$



$$I(\theta) := \text{Var}_x[s(\theta; x)] = \overset{\text{Stroby}}{\downarrow} \mathbb{E}_x[s(\theta; x)^2] = \overset{\text{Stroby}}{\downarrow} \mathbb{E}[\ell''(\theta; x)]$$



$I(\theta)$ measures how much information is in X for a r.v.

lets see this for $X \sim \text{Binom}(n, \theta)$ for fixed n

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

$$l(\theta; x) = \ln \binom{n}{x} + x \ln(\theta) + (n-x) \ln(1-\theta)$$

$$l'(\theta; x) = \frac{x}{\theta} - \frac{n-x}{1-\theta}$$

$$l''(\theta; x) = -\frac{x}{\theta^2} - \frac{n-x}{(1-\theta)^2} \quad (-1) \cdot (-1) \cdot (-1)$$

recall $E[aX+c] = aE(X) + c$

$$I(\theta) = E_x[-l''(\theta; x)] = E\left[\frac{x}{\theta^2} + \frac{n-x}{(1-\theta)^2}\right] = \frac{E(X)}{\theta^2} + \frac{n-E(X)}{(1-\theta)^2} = \frac{n\theta}{\theta^2} + \frac{n-n\theta}{(1-\theta)^2} = n\left(\frac{1}{\theta} + \frac{1}{1-\theta}\right)$$

$$= n\left(\frac{1}{\theta(1-\theta)}\right)$$

Not a function of X ! X is averaged out...

If $\theta = \frac{1}{2}$, $n=1$, How much info? $I(\frac{1}{2}) = 4$ ← the r.v. does not have too much info about θ on average

If $\theta = \frac{1}{100}$, $n=1$ $I(\frac{1}{100}) = 101.01$ ← the r.v. has a ton of info

Why?

Why should there be a multi. factor of n ?
 more data \Rightarrow more info. Because binomial is n bernoullis... more bernoulli data \Rightarrow more info

Back to the issue... CONSIDER:

What if $p(\theta) \propto \sqrt{I(\theta)}$ AKA the Jeffreys prior