

lec 8 31.03-02 2/21/16 $\theta | x \sim \text{Beta}(\alpha+x, \beta+n-x)$ $\hat{\theta}_{MLE} = \frac{x}{n} = \frac{\alpha}{\alpha+\beta}$ $\frac{1}{\theta_{MLE}}$

How do we really be uniform?? How about $\alpha=\beta=0$? $\alpha+\beta$
i.e. $\rho=0$

$\rho = \frac{\alpha+\beta}{\alpha+\beta+1} = 0 \Rightarrow \hat{\theta}_{MLE} = \frac{x}{n} = \hat{\theta}_{MLE}$. We get frequent
(Huber & L) estimate with Bayesian interpretation!

But $\theta \sim \text{Beta}(\alpha, \beta)$ Param space $\alpha, \beta \in (0, \infty)$ $\alpha, \beta \neq 0$!

Why? $\theta \sim \text{Beta}(\alpha, \beta) = \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)} \frac{1}{\theta(1-\theta)}$

$\Gamma(t) := \int_0^\infty x^{t-1} e^{-x} dx \Rightarrow \Gamma(0) = \int_0^\infty \frac{e^{-x}}{x} dx = \infty$
Just a uv-substn

Further $\text{Beta}(0,0) \propto \frac{1}{\theta(1-\theta)}$

$\int_0^1 \frac{1}{\theta(1-\theta)} d\theta = [\ln(\theta) - \ln(1-\theta)]_0^1 = \infty$

Is an "improper density"

Since it does not integrate to 1. $\theta \sim \text{Beta}(0,0)$ is called an "improper prior"

However, $\theta | x = \text{Beta}(\alpha+x, \beta+n-x) = \text{Beta}(x, n-x)$

is a "proper posterior density"

as long as $x \neq 0$ & $x \neq n$

i.e. at least one success and at least one failure in the data.

Are improper priors "legal"?

Yes, but you need to ensure posterior is proper & inference not valid.

~~Why not just pick $\alpha \approx 0$ and $\beta \approx 0$ eg $\alpha = 1 \times 10^{-10}$, $\beta = 1 \times 10^{-10}$?~~

Why not just pick $\alpha \approx 0$ and $\beta \approx 0$ eg $\alpha = 1 \times 10^{-10}$, $\beta = 1 \times 10^{-10}$?
You can! Laplace is uglier but it makes the statistics happy.

Bayes (0) is known as the Haldane prior (1932). This is the limit $\alpha, \beta \rightarrow 0$.

Consider: if $x=0$ or $x=n \Rightarrow P(\theta|x) = \text{Bay}(0)$ or $P(\theta|x) = \text{Bay}(1)$.

Bad! And... anti-laplace in same way.

Also, it is equivalent of "seeing" -1 successes and -1 failures. In a way it's pushing you toward $\theta=0$ or $\theta=1$.

Jaynes "prior state of knowledge of complete ignorance i.e. you don't even know if successes or failures are even possible".

Bayes-Laplace-Objective prior: ignorance but you are confident both successes and failures are possible.

Not explicitly covered

Recall likelihood function

$$P(x; \theta) = \mathcal{L}(\theta; x) = \prod_{i=1}^n P(x_i; \theta)$$

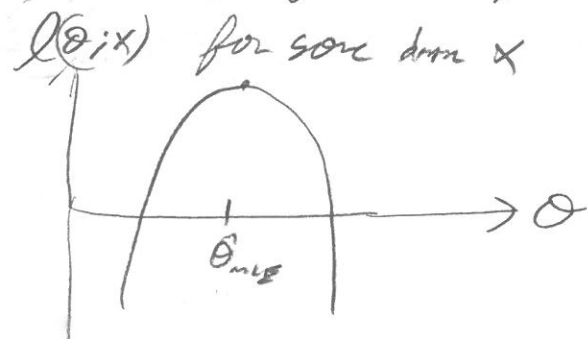
$$\mathcal{L}(\theta; x)$$

log-likelihood

to get MLE

$$\hat{\theta}_{MLE} := \underset{\theta \in \Theta}{\operatorname{argmax}} \left\{ \prod_{i=1}^n P(x_i; \theta) \right\} = \underset{\theta \in \Theta}{\operatorname{argmax}} \left\{ \sum_{i=1}^n \ln(P(x_i; \theta)) \right\}$$

Set $\frac{d}{d\theta} [\underbrace{\ell(\theta; x)}_{\text{score function}}] = 0$ and hope you find θ_{MLE}
 $\ell'(\theta; x)$



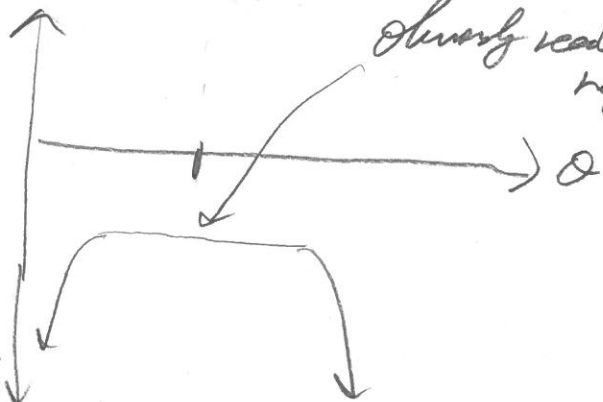
$$I(\theta) := \text{Var}[\ell'(\theta; x)]$$
$$= E[(\ell'(\theta; x))^2] - \underbrace{(E(\ell'(\theta; x)))^2}_{\text{shown to be zero}}$$



$$= E[(\ell'(\theta; x))^2]$$

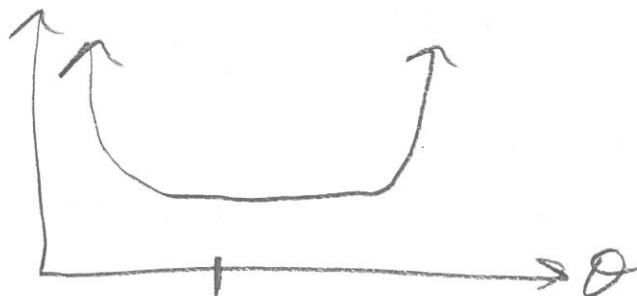
$$\ell''(\theta; x)$$

obviously needs to be negative



more interpretable

$$I(\theta) := E[-\ell''(\theta; x)]$$



On avg over all datasets,
how large is this second derivative?
If large, lots of information in data
about MLE. If small, need lots of
data to get to an MLE.

In fact $\hat{\theta}_{MLE} \sim N(\theta, (\sqrt{\frac{1}{I(\theta)}})^2)$

If information large $\hat{\theta}_{MLE} \approx \theta$. If small, large variance...

e.g. $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bern}(\theta)$

$I(\theta) = E\{l''(\theta; x)\}$

$l(\theta; x) = \sum_{i=1}^n \ln p(x_i; \theta) = \sum_{i=1}^n \ln(\theta^{x_i} (1-\theta)^{1-x_i}) = \sum_{i=1}^n x_i \ln(\theta) + (1-x_i) \ln(1-\theta)$

$l'(\theta; x) = \frac{\partial}{\partial \theta} [] = \sum_{i=1}^n \frac{x_i}{\theta} - \frac{1-x_i}{1-\theta}$

$l''(\theta; x) = \frac{\partial}{\partial \theta} [] = \sum_{i=1}^n -\frac{x_i}{\theta^2} - \frac{1-x_i}{(1-\theta)^2} = -\frac{n\bar{x}}{\theta^2} + \frac{n-n\bar{x}}{(1-\theta)^2}$

$-l''(\theta; x) = \frac{n\bar{x}}{\theta^2} - \frac{n-n\bar{x}}{(1-\theta)^2} = n \left(\frac{\bar{x}}{\theta^2} + \frac{1-\bar{x}}{(1-\theta)^2} \right)$

$I(\theta) = E[] = n \left(\frac{E(\bar{x})}{\theta^2} + \frac{1-E(\bar{x})}{(1-\theta)^2} \right) = n \left(\frac{\theta}{\theta^2} + \frac{1-\theta}{(1-\theta)^2} \right) = n \left(\frac{1}{\theta} + \frac{1}{1-\theta} \right) = n \left(\frac{1-\theta}{\theta(1-\theta)} + \frac{\theta}{\theta(1-\theta)} \right) = \boxed{\frac{n}{\theta(1-\theta)}}$

θ is a const.

Why do we care?

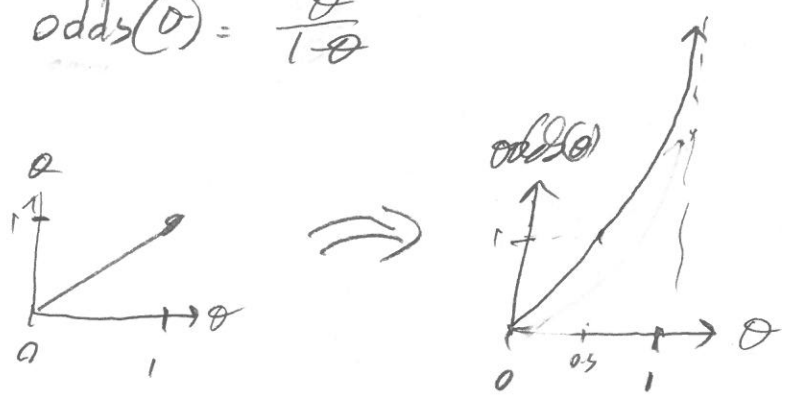
Recall $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bern}(\theta)$, $\theta \sim U(0,1)$

Why is principle of indifference bad?

principle of indifference

It doesn't work on a diffuse scale

e.g. $\text{odds}(\theta) = \frac{\theta}{1-\theta}$

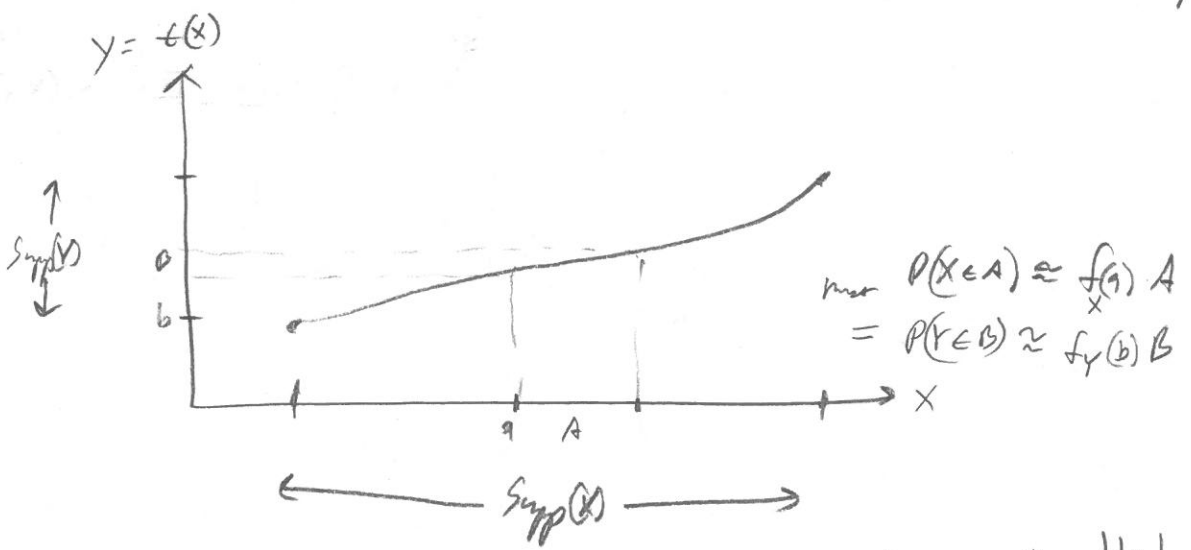


If I'm different about θ , should I be different about the odds(θ)?

$f(\theta) = 1 \mathbb{1}_{\theta \in [0,1]}$

$g(\theta) = \frac{\theta}{1-\theta}$ What is PDF of $g(\theta)$, a monotonic \Rightarrow 1:1 function of θ ?

Theory of transformation of variables, given $Y = t(X)$ and $f_X(x)$ find $f_Y(y)$



then $P(X \in A) \approx f_X(a) A$
 $= P(Y \in B) \approx f_Y(b) B$

$X = t^{-1}(Y)$

$\Rightarrow f_X(x) |dx| = f_Y(y) |dy| \Rightarrow f_Y(y) = f_X(x) \left| \frac{dx}{dy} \right|$

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

BAD if non-monotonic

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

$$f_{\text{odds}}(0) = f_{\theta}(\tau^{-1}(0)) \left| \frac{d}{d\theta} [\tau^{-1}(0)] \right| = (1) \left| \frac{1}{(0+1)^2} \right| = \frac{1}{(0+1)^2} \mathbb{1}_{0 \in (0, \infty)}$$

$$0 = \frac{\theta}{1-\theta} \Rightarrow 0 - 0\theta = \theta \Rightarrow 0 = \theta + \theta \Rightarrow 0 = \theta(2) \Rightarrow \theta = \frac{0}{2} = \tau^{-1}(0)$$

Is this a PDF? $\int_0^{\infty} \frac{1}{(0+1)^2} d0 = 1$ ✓

$0 \notin U(0,1)$. Is this a contradiction??? If indifferent about $\theta \Rightarrow$ indifferent odds(θ)

Consider $l := \text{logit}(\theta) := \ln(\text{odds}(\theta)) = \ln\left(\frac{\theta}{1-\theta}\right)$ $\text{Supp}(l) = \mathbb{R}$

which is a very well-used "link function" and is "logistic regression" or "logit regression".

$$f_l(l) = f_x\left(\frac{e^l}{1+e^l}\right) \left| \frac{d}{dl} \left[\frac{e^l}{1+e^l} \right] \right| = \frac{e^l}{(1+e^l)^2} \neq U(0,1)!$$

$$l = \ln \frac{\theta}{1-\theta} \Rightarrow e^l = \frac{\theta}{1-\theta} \Rightarrow e^l - e^l \theta = \theta \Rightarrow e^l = \theta + e^l \theta \Rightarrow e^l = \theta(1+e^l) \Rightarrow \theta = \frac{e^l}{1+e^l} = \tau^+(l)$$

The second you go to quota scale, principle of indifference goes bye-bye!
It's as if going rewriting logits near 0 more times than logit elements...

Note: $\theta \sim \text{Beta}(0,0) \propto \frac{1}{\theta(1-\theta)}$ Haldane prior

$$l = \text{logit}(\theta) := \ln\left(\frac{\theta}{1-\theta}\right)$$

$$f_l(l) = f_x\left(\frac{e^l}{1+e^l}\right) \frac{e^l}{(1+e^l)^2} = \frac{1}{\left(\frac{e^l}{1+e^l}\right)\left(1-\frac{e^l}{1+e^l}\right)} \frac{e^l}{(1+e^l)^2} = \frac{e^l}{\left(\frac{e^l}{1+e^l}\right)\left(\frac{1}{1+e^l}\right)(1+e^l)^2} = 1 \in \mathbb{R}$$

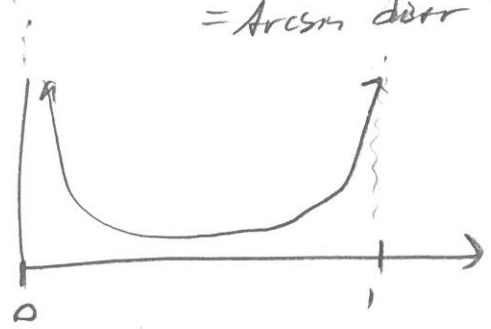
\Rightarrow Uniform prior on logit scale (but improper).

But is the $\theta \sim p(\theta)$ s.t. it is the same on all transformed scales? So you don't have to worry if it changes a little bit? YES

Let $p(\theta) \propto \sqrt{I(\theta)}$ \leftarrow Fisher Information of the likelihood

e.g. in the binomial likelihood model

$$p(\theta) \propto \sqrt{\frac{n}{\theta(1-\theta)}} = n^{\frac{1}{2}} \theta^{\frac{1}{2}} (1-\theta)^{\frac{1}{2}} \propto \theta^{\frac{1}{2}} (1-\theta)^{\frac{1}{2}} \propto \text{Beta}\left(\frac{1}{2}, \frac{1}{2}\right) = \text{"Arcsin distr"}$$



$\text{Beta}(0,0)$	$<$	$\text{Beta}(\frac{1}{2}, \frac{1}{2})$	$<$	$\text{Beta}(1,1)$
improper				$= U(0,1)$
Haldane		Jeffreys prior		objective on the θ scale
complete ignorance		"half ignorance"		ignorance

So is $\theta \sim \text{Beta}(\frac{1}{2}, \frac{1}{2})$ really more arbitrary than $\theta \sim \text{Beta}(\frac{1}{2}, \frac{1}{2})$?

How about $t(\theta) = \text{odds}(\theta) = \frac{\theta}{1-\theta}$?

$$f_{\text{odds}}(\theta) = f_{\theta}\left(\frac{\theta}{\theta+1}\right) \left| \frac{d}{d\theta} \left[\frac{\theta}{\theta+1} \right] \right| = \frac{1}{\text{Beta}(\frac{1}{2}, \frac{1}{2})} \left(\frac{\theta}{\theta+1}\right)^{-\frac{1}{2}} \left(\frac{\theta}{\theta+1}\right)^{-\frac{1}{2}} \left| \frac{1}{(\theta+1)^2} \right|$$

\uparrow
 $t^{-1}(\theta)$

$\propto \frac{1}{\theta} \frac{1}{(\theta+1)^2} = \frac{1}{\theta(\theta+1)}$

$$P(\theta) \propto \sqrt{f_{\text{odds}}(\theta)} = \sqrt{\frac{1}{\theta(\theta+1)}} = \theta^{-\frac{1}{2}} (\theta+1)^{-\frac{1}{2}} \propto \text{Beta}\left(\frac{1}{2}, \frac{1}{2}\right)$$

Invariance of

MAGIC!!!

Proof of Jeffreys Prior

Consider $\theta \sim P(\theta)$ a prior s.t. $P(\theta) \propto \sqrt{I(\theta)}$
 i.e. the Jeffreys prior.

WTS $P(\phi) \propto \sqrt{I(\phi)}$ for arbitrary $\phi = t(\theta)$

$$\begin{aligned} \text{We know } P_{\phi}(\phi) &= P_{\theta}(t^{-1}(\phi)) \left| \frac{d\theta}{d\phi} \right| \\ &= P_{\theta}(\theta) \left| \frac{d\theta}{d\phi} \right| \propto \sqrt{I(\theta)} \left| \frac{d\theta}{d\phi} \right| \\ &= \sqrt{I(\theta) \left(\frac{d\theta}{d\phi} \right)^2} \end{aligned}$$

$$\begin{aligned}
&= \sqrt{E\left[\left(\ell(\theta; y)\right)^2\right] \left(\frac{d\theta}{d\phi}\right)^2} \\
&= \sqrt{E\left[\left(\frac{d}{d\theta} \left[\ln(\mathcal{L})\right]\right)^2 \left(\frac{d\theta}{d\phi}\right)^2\right]} \\
&= \sqrt{E\left[\left(\frac{d \ln(\mathcal{L})}{d\theta} \cdot \frac{d\theta}{d\phi}\right)^2\right]} = \sqrt{E\left[\left(\frac{d \ln \mathcal{L}}{d\phi}\right)^2\right]} = \sqrt{I(\phi)} \checkmark
\end{aligned}$$

So $\text{Beta}\left(\frac{1}{2}, \frac{1}{2}\right)$ is the Jeffreys/Invariant prior
and this makes it a "natural" choice.

Another proof

$$p(\theta) \propto \sqrt{I(\theta)} \Rightarrow$$

WTS $p(\phi) \propto \sqrt{I(\phi)}$ for arbitary $\phi := h(\theta)$

We know $p(\phi) = p(\theta) \left| \frac{d\theta}{d\phi} \right| \propto \sqrt{I(\theta)} \left| \frac{d\theta}{d\phi} \right|$

If this is $\propto \sqrt{I(\phi)}$ we're done

$$= \sqrt{E\left[-\left[\frac{d^2}{d\phi^2} \left[\underbrace{\ln p(x|\phi)}_{g(\phi)}\right]\right]\right]} = \sqrt{E\left[-\frac{1}{d\phi^2} [g(\phi)]\right]} \quad \text{let } \theta = h^{-1}(\phi)$$