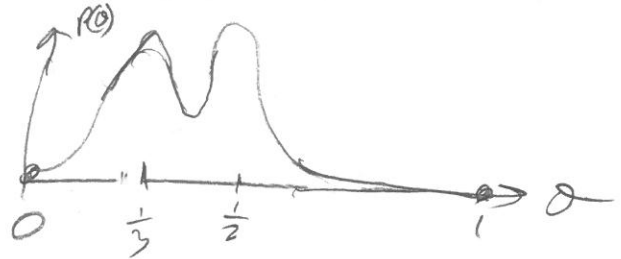


lecture 19 Mon 30.03-02 5/2/16

Recall ... $X|\theta \sim \text{bin}(n, \theta)$, $\theta \sim \text{beta}(\alpha, \beta)$

$$\Rightarrow \theta|X \sim \text{beta}(x+\alpha, n-x+\beta), \quad X \sim \text{Bin}(n, \alpha, \beta)$$

What if $\theta \neq \text{beta}(\alpha, \beta)$ e.g.



Can you still use conjugacy?

No, but Yes... what if $P(\theta) \approx \sum_{m=1}^M \delta_m P_m(\theta)$

$M=1$
prior mixing prop.

$$P_m(\theta) = \text{beta}(\alpha_m, \beta_m) \quad \text{st} \quad \sum_{m=1}^M \delta_m = 1$$

ie

a mixture of

betas... is this conjugate?

No .. but the math is nice...

Conjugate mixture priors

$$\sum_{m=1}^M \delta_m P_m(\theta)$$

$$P(\theta|x) = \frac{P(x|\theta) \sum_{m=1}^M \delta_m P_m(\theta)}{P(x)}$$

$$= \frac{\sum \delta_m P(x|\theta) P_m(\theta)}{P(x)}$$

$$= \frac{\sum \delta_m \frac{P(x|\theta) P_m(\theta)}{P_m(x)} P_m(x)}{P(x)}$$

$$\int P(x|\theta) \sum_{m=1}^M \delta_m P_m(\theta) d\theta = \sum_{m=1}^M \delta_m \int P(x|\theta) P_m(\theta) d\theta$$

$$\text{if } \delta_m = \frac{1}{M} \quad \sum_{m=1}^M \frac{P_m(x)}{P(x) + P_m(x)} P_m(\theta|x) \propto \sum_{m=1}^M P_m(x) P_m(\theta|x)$$

$$= \sum_{m=1}^M \frac{\delta_m P_m(x)}{P(x)} P_m(\theta|x) \propto \sum_{m=1}^M \delta_m P_m(x) P_m(\theta|x)$$

posterior then conjugate prior

e.g. two betas $\theta_1 \sim \text{beta}(\alpha_1, \beta_1)$, $\theta_2 \sim \text{beta}(\alpha_2, \beta_2)$, $\delta_1 = \delta_2 = \frac{1}{2}$

$$P(\theta|x) = \frac{1}{\text{Bin}(n, \alpha_1, \beta_1) + \text{Bin}(n, \alpha_2, \beta_2)} \left(\text{Bin}(n, \alpha_1, \beta_1) \text{beta}(x+\alpha_1, n-x+\beta_1) + \text{Bin}(n, \alpha_2, \beta_2) \text{beta}(x+\alpha_2, n-x+\beta_2) \right)$$

not a conj. distr!

Real life ex.

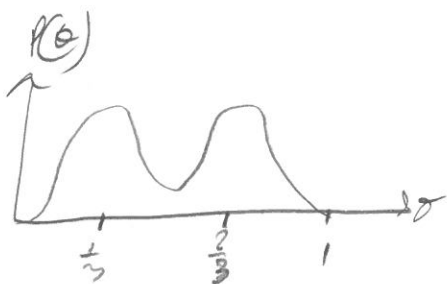
magnum "Slip" coin by spinning down. $P(H) = \frac{1}{2}$

BUT... if the side is shaved, the $P(H)$ is biased... assume bias is $\frac{2}{3}$ or $\frac{1}{3}$ which is not

Assume shaved sounds H or sounds T is equal prob... *detectable with higher eye*

$$\text{So } P(\theta) \approx \frac{1}{2} \text{Bern}(10, 20) + \frac{1}{2} \text{Bern}(20, 10)$$

the strengths come from prior knowledge that the estimates of the strong



coin "spin" $h=10, x=3$ (10 trials)

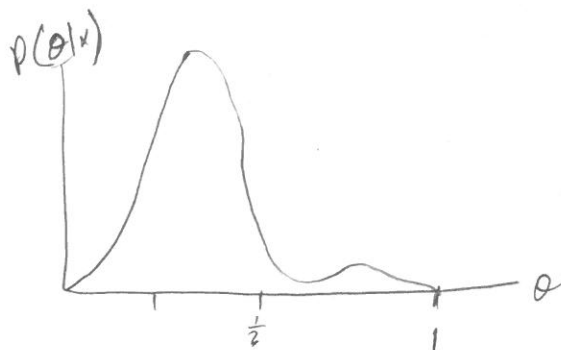
$$P(\theta|x) = \frac{1}{\underbrace{\text{dbb}(3,10,10,20)}_{0.228} + \underbrace{\text{dbb}(3,10,20,10)}_{0.0277}} \left(\begin{array}{l} \text{dbb}(3,10,10,20) \text{Bern}(13, 27) \\ + \text{dbb}(3,10,20,10) \text{Bern}(23, 17) \end{array} \right)$$

How to sample??
draw and compare "Monte Carlo"

Recall
 = rdbinomial(n, α, β)
 = pbinomial(x, n, α, β)
 = {dbinomial(p, n, α, β)
 = dbinomial(x, n, α, β)

↑
value at $P(X=x)$
if $X \sim \text{Bin}(n, \alpha, \beta)$

$$= 0.892 \text{Bern}(13, 27) + 0.108 \text{Bern}(23, 17)$$



$$\hat{\theta}_{\text{MMSE}} = 0.892 \frac{13}{40} + 0.108 \frac{23}{40} = 0.352$$

$$\hat{\theta}_{\text{MAP}} = \arg\max \{ P(\theta|x) \}$$

not so simple...

Gently,

Review,

13

How to get $\hat{\theta}_{MAP}$? ① If you know unrolled, use Norm approx
② Use grid sampling. make be good at compare all vals. Take max.

$$\hat{\theta}_{MAP} := \argmax \{P(\theta|x)\} = \argmax \{L(\theta|x)\} = \argmax \left\{ \sum_{m=1}^M \gamma_m \ln(x) \ln(\theta(x)) \right\}$$

for the beta-binomial case,

$$\frac{d}{d\theta} \left[\sum \gamma_m \text{BetaBin}(n, \alpha_m, \beta_m) \text{Beta}(x+\alpha_m, n-x+\beta_m) \right] \stackrel{\text{set}}{=} 0$$

$$\frac{d}{d\theta} \left[\sum \gamma_m \left(\binom{n}{x} \frac{b(x+\alpha_m, n-x+\beta_m)}{b(\alpha_m, \beta_m)} \right) \left(\frac{1}{b(x+\alpha_m, n-x+\beta_m)} \theta^{x+\alpha_m-1} (1-\theta)^{n-x+\beta_m-1} \right) \right] \stackrel{\text{set}}{=} 0$$

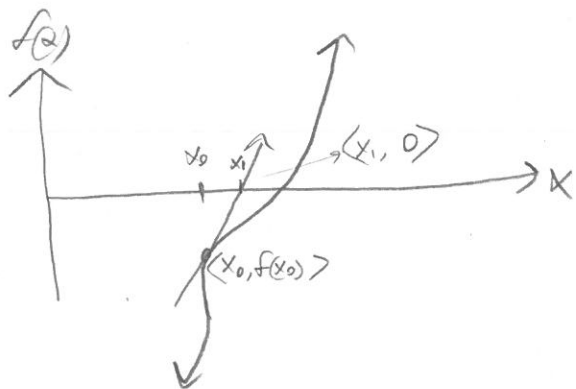
$$\sum \frac{\gamma_m}{b(\alpha_m, \beta_m)} \frac{d}{d\theta} \left[\right] = \sum \frac{\gamma_m}{b(\alpha_m, \beta_m)} \left((x+\alpha_m-1) \theta^{x+\alpha_m-2} (1-\theta)^{n-x+\beta_m-1} - (n-x+\beta_m-1) \theta^{x+\alpha_m-1} (1-\theta)^{n-x+\beta_m-2} \right)$$

= 0

No closed form sol for θ !

What to do? Numerical method...

Solve $f(x) = 0$ for x , where f is diff. val function



Step 1 guess an x_0 as the root

Step 2 draw tangent line at $f(x_0)$

Step 3: find its x intercept = x_1

Step 4: Repeat steps 1-3 with $x_0 = x_1$,
may this converge

$$|x_{t+1} - x_t| < \epsilon \text{ (threshold)}$$

where ϵ is small

(this algorithm is an iterative method)

Exp 2 ~~can~~ can be translated as follows. Back to 7th grade..

at slope small (a,b)

$$y - b = m(x - a)$$

$$y - f(x_0) = f'(x_0)(x - x_0)$$

x_1 is solution for x -intercept i.e. $y=0$

$$-f(x_0) = f'(x_0)(x_1 - x_0)$$

$$\Rightarrow -\frac{f(x_0)}{f'(x_0)} = x_1 - x_0 \Rightarrow x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}$$

1669 - Newton for polynomial

1690 - Raphson

1740 - Simpson (John's version)

$$x_2 = x_1 - \frac{f(x_1)}{f'(x_1)}$$

$$x_{t+1} = x_t - \frac{f(x_t)}{f'(x_t)}$$

Call it the Newton-Raphson Method

Let's get $\hat{\theta}_{MAP}$. $f(x)$

$$\theta_0 \stackrel{?}{=} 0.33$$

$$\theta_1 = 0.37$$

$$k(\theta|x) = \underbrace{\frac{1}{B(1,2)}}_{C_1} \theta^{12} (1-\theta)^{26} + \underbrace{\frac{1}{B(23,17)}}_{C_2} \theta^{23} (1-\theta)^{16}$$

$$k'(\theta|x=3) = C_1 (12 \theta^{11} (1-\theta)^{26} + 26 \theta^{12} (1-\theta)^{25}) + C_2 (23 \theta^{22} (1-\theta)^{16} + 16 \theta^{23} (1-\theta)^{15})$$

$$f'(x) = k''(\theta|x=3) = \checkmark$$

$$N-R \text{ gives } \hat{\theta}_{MAP} = 0.31577 \neq 0.352 = \hat{\theta}_{MLE}$$

What did we do?

$$p(\theta|x_1, \dots, x_m) \propto p(\theta) \prod_{i=1}^m p(x_i|\theta)$$

"Hierarchical Model"

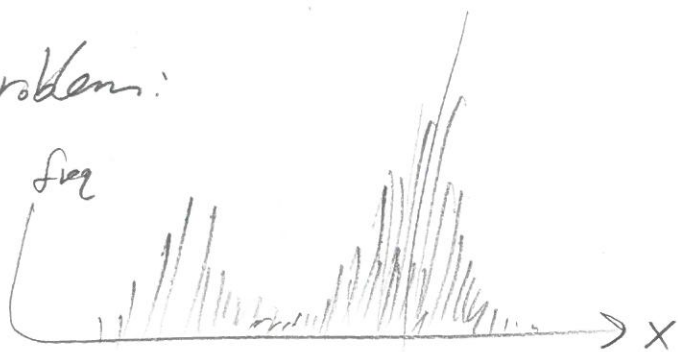
$$p(\theta|x) = \frac{p(\theta) \sum_{i=1}^m p(x_i|\theta)}{p(x)}$$

prior (since it is the prior probab of θ) but we depend on other quantities

hyperprior w/ hyperhyperparameters

Not correct

New problem:



this looks like...

Mixture model

$$x_1, \dots, x_n \stackrel{\text{each}}{\sim} p(x | \underbrace{\theta_0, \sigma_0^2, \theta_1, \sigma_1^2, \epsilon}_{\text{param vector is 5-dim}}) = \epsilon N(\theta_0, \sigma_0^2) + (1-\epsilon) N(\theta_1, \sigma_1^2)$$

What if we just use about a best guess of params? MLE

$$\text{Set } \nabla \ln L(\theta_0, \sigma_0^2, \theta_1, \sigma_1^2, \epsilon) = 0$$

$$\nabla \ln \prod_{i=1}^n \epsilon \frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{1}{2\sigma_0^2}(x_i - \theta_0)^2} + (1-\epsilon) \frac{1}{\sqrt{2\pi\sigma_1^2}} e^{-\frac{1}{2\sigma_1^2}(x_i - \theta_1)^2}$$

good luck!

What to do?

Grid search $\theta_0 \in [\dots]$, $\theta_1 \in [\dots]$, $\epsilon \in (0,1)$, $\sigma_0^2 \in [\dots]$, $\sigma_1^2 \in [\dots]$

In 5 dimensions $1000^5 = 1e15 \dots$ still okay... but what if

$$x_1, \dots, x_n \stackrel{\text{each}}{\sim} \sum_{m=1}^M \epsilon_m N(\theta_m, \sigma_m^2), \quad M \text{ large} \dots \text{grid search fails} \dots$$