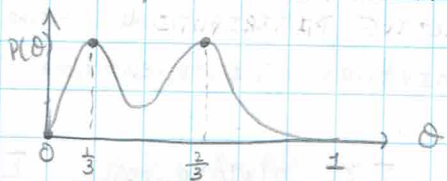Suppose $x|\theta \sim Bin(n, \theta)$ & $\theta \sim Beta(\alpha, \beta)$. We know that this implies

① $\theta|x \sim Beta(x+\alpha, n-x+\beta)$ &

② $X \sim BetaBin(n, \alpha, \beta)$

What if instead, we had $x|\theta \sim Bin(n, \theta)$ & $\theta \not\sim Beta(\alpha, \beta)$. I.e. if



Clearly $\theta \not\sim Beta(\alpha, \beta)$ since Beta is unimodal & this is bimodal.

What would be our posterior? Well

$$p(\theta|x) = \frac{P(x|\theta)p(\theta)}{P(x)} \quad \alpha \quad P(x|\theta)p(\theta)$$

So $k(\theta|x) = p(x|\theta)p(\theta)$. Hence

① we can sample from $\theta \in \theta_G$ OR

② get $k(\theta_g|x)$ via grid sampling

③ use a "TYPE OF CONJUGACY"

Let us explore ③.

Assume $p(\theta) = \gamma_1 Beta(\alpha_1, \beta_1) + \gamma_2 Beta(\alpha_2, \beta_2)$ where
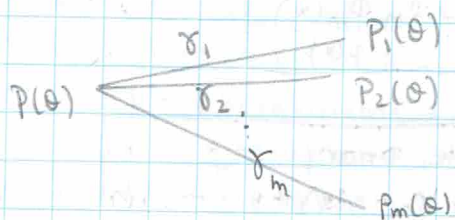
$$\gamma_1 + \gamma_2 = 1$$

Suppose, in particular,

$$p(\theta) = \tfrac{1}{2} Beta(3,3) + \tfrac{1}{2} Beta(2,4)$$

$\left.\begin{array}{c} \\ \\ \end{array}\right\}$ CONJUGATE MIXTURE PRIOR

In general a CONJUGATE MIXTURE PRIOR looks like

$$p(\theta) = \sum_{m=1}^{M} \gamma_m P_m(\theta) \quad s.t. \quad \sum \gamma_m = 1 \quad \& \quad P_m(\theta) \text{ is a } \underline{density}.$$

Visually we see $p(\theta)$ is broken down into M parts:



We can think of each $\gamma_i$ as our PRIOR WEIGHTS.

If we have a CONJUGATE MIXTURE PRIOR then

$$p(x) = \int_\theta p(x,\theta)\, d\theta = \int_\theta p(x|\theta)\, p(\theta)\, d\theta$$

$$= \int_\theta p(x|\theta) \sum_m \gamma_m\, p_m(\theta)\, d\theta$$

$$= \int_\theta p(x|\theta) \left( \gamma_1 p_1(\theta) + \gamma_2 p_2(\theta) + \cdots + \gamma_m p_m(\theta) \right)^{d\theta} d\theta = \sum_m \gamma_m \int_\theta p(x|\theta)\, p_m(\theta)\, d\theta$$

$$= \sum_m \gamma_m\, \underset{p_m(x)}{\underbrace{p_m(x)}} \rightarrow \text{this is the PRIOR PREDICTIVE DISTRIBUTION UNDER THAT SINGLE}$$

PRIOR IN THE MIXTURE DISTRIBUTION.

Hence,

$$p(\theta|x) = \frac{p(x|\theta)p(\theta)}{p(x)} = \frac{p(x|\theta)\sum_m \gamma_m p_m(\theta)}{p(x)} = \frac{\sum_m \gamma_m\, p(x|\theta)\, p_m(\theta)}{p(x)} = \frac{\sum_m \gamma_m\, p(x|\theta)\,(1)\, p_m(\theta)}{p(x)}$$

$$= \frac{\sum_m \gamma_m\, p(x|\theta)\, \frac{p_m(x)}{p_m(x)}\, p_m(\theta)}{p(x)} = \frac{\sum_m \gamma_m\, \frac{p(x|\theta)\, p_m(\theta)}{p_m(x)}\, p_m(x)}{p(x)}$$

$$= \frac{\sum_m \gamma_m\, p_m(\theta|x)\, p_m(x)}{p(x)} = \frac{\sum_m \gamma_m\, p_m(x)}{p(x)} \cdot p_m(\theta|x)$$

Now $p(x)$ & $p_m(x)$ is a constant so

$$\boxed{= \sum \gamma_m'\, p_m(\theta|x)}$$

**Thm:** If $p(\theta) = \sum \gamma_m\, p_m(\theta)$ then

$$p(\theta|x) = \sum \gamma_m'\, p_m(\theta|x)$$

where $\gamma_i$ is your prior weight on mixture model $i$, $\gamma_i'$ is your posterior weight on model $i$, $p_m(\theta|x)$ is your posterior model per single mixture & where $\sum_{m=1}^{M} \frac{\gamma_m\, p_m(x)}{p(x)}$

Hence this looks CONJUGATE. The prior & posterior are of the same form.

**Thm:** For the mixture model conjugate prior,

$$p(\theta|x) \propto \sum \gamma_m\, p_m(x) \cdot p_m(\theta|x)$$

**Proof:** follows straight from $p(\theta|x) = \sum \gamma_m'\, p_m(\theta|x)$

**Thm:** In the special case all prior weights are equal, i.e. $\gamma_m = \frac{1}{M}$ then $p(\theta|x) = \sum \frac{p_m(x)}{p_1(x) + \cdots + p_M(x)} \cdot p_m(\theta|x)$

$$\propto \sum p_m(x)\, p_m(\theta|x)$$

With these theorems, let us do an example: Suppose
$X \sim Bin(n, \theta)$ for a fixed $n$ & suppose
$\theta_1 \sim Beta(\alpha_1, \beta_1)$, $\theta_2 \sim Beta(\alpha_2, \beta_2)$, $\gamma_1 = \frac{1}{2} = \gamma_2$.

Then $p(\theta) = \sum \gamma_m p_m(\theta)$
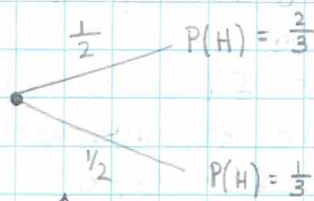$$= \frac{1}{2} Beta(\alpha_1, \beta_1) + \frac{1}{2} Beta(\alpha_1, \beta_1) \quad \&$$

$P_1(x) = BetaBin(n, \alpha_1, \beta_1)$, $P_1(\theta|x) = Beta(\alpha_1 + x, \beta_1 + (n-x))$

$P_2(x) = BetaBin(n, \alpha_2, \beta_2)$, $P_2(\theta|x) = Beta(\alpha_2 + x, \beta_2 + (n-x))$

Thus $p(\theta|x) = \dfrac{\sum P_m(x) \cdot P_m(\theta|x)}{P_1(x) + P_2(x)} = \dfrac{P_1(x) P_1(\theta|x) + P_2(x) P_2(\theta|x)}{P_1(x) + P_2(x)}$

$$= \frac{BetaBin(n, \alpha_1, \beta_1)\, Beta(\alpha_1 + x, \beta_1 + (n-x)) + BetaBin(n, \alpha_2, \beta_2)\, Beta(\alpha_2 + x, \beta_2 + (n-x))}{BetaBin(n, \alpha_1, \beta_1) + BetaBin(n, \alpha_2, \beta_2)}$$
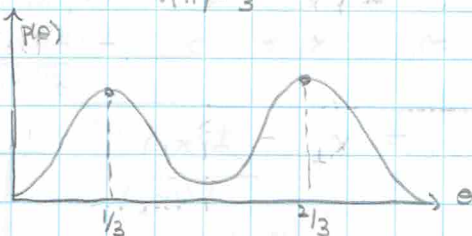
Now suppose a magician shaves off any side of a coin w.p. $\frac{1}{2}$. He then spins the coin! So we have



AND SUPPOSE FROM OUR PRIOR KNOWLEDGE,
$P(\theta) \approx \frac{1}{2} Beta(10,20) + \frac{1}{2} Beta(20,10)$

Then $p(\theta)$ looks like



Suppose we spin a coin 10 times & get 3 heads. What is the posterior?

Recall in $R$,
```
x = rbetabin (n, α, β)
p = pbetabinom (x, n, α, β)
x = qbetabinom (p, n, α, β)
d = betabin (x, n, α, β)  ← ~ for discrete distributions this
                                    is p(x)
```

using these commands & the formula presented above becomes

$p(\theta|x) = p(\theta|x=3) = \dfrac{dbb(3,10,10,20)\, Beta(13,27) + dbb(3,10,20,10)\, Beta(23,17)}{dbb(3,10,10,20) + dbb(3,10,20,10)}$

$$= .892\, Beta(13,27) + .108\, Beta(23,17)$$

<u>Plotting yields</u>



\* which is STILL BIMODAL b/c the data has yet to
overpower the prior \*

: We now take a small detour: Suppose you have a function
$f(x)$ s.t. $f$ is differentiable & you want to find $\tilde{x}$ s.t.
_i.e. you want to find a ROOT_
$f(\tilde{x}) = 0.$ How do you do this? Here is an algorithm:



① Guess where the solution is & call
your guess $x_0$
② Draw tangent line from $x_0$
③ Find the x-intercept of $x_0$
④ Repeat ①-③ until
$$\left| x_{t+1} - x_t \right| < \varepsilon \quad \text{for } \varepsilon \text{ a}$$
predefined tolerance level

NOTE IN STEP ② WE ARE DRAWING A LINE through the point
$(x_0, f(x_0))$ (our notation here doesn't <u>agree</u>) By the
point slope formula, this is the line
$$y - f(x_0) = f'(x_0)(x - x_0)$$

Now in step ③ we are getting the $x$-intercept, i.e. the value
for where $y = 0$. Hence
$$-f(x_0) = f'(x_0)(x_1 - x_0)$$
$$\Rightarrow \quad x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}$$

Generally we have
$$\boxed{x_{t+1} = x_t - \frac{f(x_t)}{f'(x_t)}}$$
NEWTON-RAPHSON
ALGORITHM. This is
an example of an
ITERATIVE
NUMERICAL
ALGORITHM

☆ NOTE: YOU MUST CHOOSE
$x_0$ & CHOOSE AN $\varepsilon$ ☆ (SO YOU ENTER, $x_0, \varepsilon$, & YOUR
FUNCTION ☆
We will soon use this result.

We know by the Theorem which is boxed in on PAGE 2 that
if $p(\theta) = \sum \gamma_m p_m(\theta)$ then $p(\theta|x) = \sum \gamma'_m p_m(\theta|x)$.
What is $\hat{\theta}_{mmse}$?

<u>Thm</u>: $\hat{\theta}_{mmse} = \int_\theta \theta p(\theta|x) d\theta = \int_\theta \theta \cdot \sum \gamma'_m p_m(\theta|x) d\theta$

$= \sum \gamma'_m \int_\theta \theta p_m(\theta|x) d\theta = \sum \gamma'_m E_m[\theta|x]$ SO

$$\boxed{\hat{\theta}_{mmse} = E[\theta|x] = \sum \gamma'_m E_m[\theta|x]}$$

This is a very nice result which matches our intuition! The expected value of $\theta|x$, or the expected value of the sum $\sum \gamma_m' P_m(\theta|x)$, is the sum of $\hat{\theta}_p^{mmse}$ times the posterior weight $\gamma_m'$ where $\hat{\theta}_p^{mmse} = E_p[\theta|x]$, i.e. the expected value of $P_m(\theta|x)$.

As an example consider the $n=10$, $x=3$ heads example. We found
$$p(\theta|x) = p(\theta|x=3) = .892 \text{ Beta}(13,27) + .108 \text{ Beta}(23,17).$$
Recall for $Y \sim \text{Beta}(\alpha,\beta)$, $E[Y] = \frac{\alpha}{\alpha+\beta}$
$$\hat{\theta}_{mmse} = E[\theta|x]$$
$$= .892\left(\frac{13}{40}\right) + .108\left(\frac{23}{40}\right)$$
$$= .352!$$

A more difficult question to answer is what is $\hat{\theta}_{MAP}$? Well,
$$\hat{\theta}_{MAP} = \underline{\text{argmax}}\{p(\theta|x)\} = \text{argmax}\{k(\theta|x)\}.$$
So to find $\hat{\theta}_{MAP}$, it suffices to maximize $k(\theta|x)$. To find the max of $k(\theta|x)$ we now have 3 ways:

1) From Lecture 15, $k(\theta|x) \approx \alpha$ to $N\left(\hat{\theta}_{MAP}, \left(\sqrt{\frac{1}{g''(\hat{\theta}_{MAP}|x)}}\right)^2\right)$

2) From Lecture 15, $k(\theta|x)$ can grid sample.

3) We can do this the old fashioned way, take the derivative of $k(\theta|x)$ & set it equal to 0.

Let us do (3): $p(\theta|x) = \sum \gamma_m \binom{n}{x} \frac{B(x+\alpha_m, n-x+\beta_m)}{B(\alpha_m, \beta_m)} \cdot \frac{1}{B(x+\alpha_m, n-x+\beta_m)} \cdot \theta^{x+\alpha_m-1}(1-\theta)^{n-x+\beta_m-1}$

$$\Rightarrow k(\theta|x) = \sum \gamma_m \cdot \frac{1}{B(\alpha_m, \beta_m)} \cdot \theta^{x+\alpha_m-1}(1-\theta)^{n-x+\beta_m-1}$$

$$\Rightarrow k'(\theta|x) = \sum \frac{\gamma_m}{B(\alpha_m, \beta_m)}\left[(x+\alpha_m-1)\theta^{x+\alpha_m-2}(1-\theta)^{n-x+\beta_m-1} - (n-x+\beta_m-1)\theta^{x+\alpha_m-1}(1-\theta)^{n-x+\beta_m-2}\right]$$

Setting $k'(\theta|x) = 0$ & solving for $\theta$ is EXTREMELY DIFFICULT.
If only we had an algorithm to do this ..... oh wait, we do!
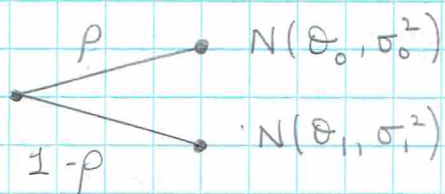Let us use the Newton-Raphson Algorithm.
Remember, to use the NRA we must specify our initial guess & specify some $\varepsilon > 0$. BASED on our graph on the top of the last page, it seems $\hat{\theta}_{MAP}$ is NEAR $\hat{\theta}_{mmse}$. So using the $\varepsilon$ from Professor Kapeler, we enter in the computer
"$\theta_0 = \hat{\theta}_{mmse} = .352$" & enter "$\theta_{t+1} = \theta_t + \frac{k'(\theta|x=3)}{k''(\theta|x=3)}$"
Note the computer also calculates the next derivative of $k'$.
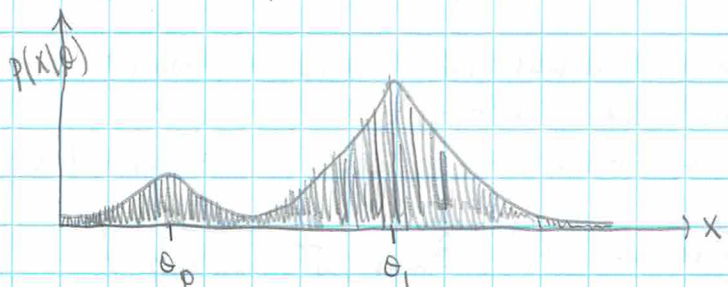We get: $\hat{\theta}_{MAP} = .31577 \neq \hat{\theta}_{mmse} = .352$.

We now begin our next unit: Let us generalize our
coin flipping example: Suppose we have the following mixture

$$\rho \quad \bullet \quad N(\theta_0, \sigma_0^2)$$
$$1-\rho \quad \bullet \quad N(\theta_1, \sigma_1^2)$$

Where $X_1, ..., X_n \overset{\text{exch}}{\sim} \rho N(\theta_0, \sigma_0^2) + (1-\rho) N(\theta_1, \sigma_1^2)$
$= p(x|\theta)$

So we have a bunch of r.v.'s each being a
mixture r.v.

So we have $p(x|\theta) = \rho N(\theta_0, \sigma_0^2) + (1-\rho) N(\theta_1, \sigma_1^2)$ looks like this:



Now $p(x|\theta)$ is really $p(x| \rho, \theta_0, \theta_1, \sigma_0^2, \sigma_1^2)$ since given $\theta$ means
you are given $\rho N(\theta_0, \sigma_0^2) + (1-\rho) N(\theta_1, \sigma_1^2)$. Thus we have
5 PARAMETERS!

The Bayesian Question is what is $p(\theta|x)$, i.e. what is
$p(\rho, \theta_0, \theta_1, \sigma_0^2, \sigma_1^2 | x)$ which is VERY, VERY, DIFFICULT.
We saw how messy things were for the normal distribution
i.e. $p(\theta, \sigma^2 | x)$. This now is even messier! Hence to do
things like compute credible regions for $\theta_0$ & such are
difficult. A more reasonable question is what is $\hat{\theta}_{MLE}$?

That is, what is the solution to $\nabla \ln \ell_0(\theta_0, \theta_1, \sigma_0^2, \sigma_1^2, \rho : X) = 0$
That is, what is the solution to

$$\nabla \ln \left( \prod_{i=1}^{n} \left( \rho \frac{1}{\sqrt{2\pi \sigma_0^2}} e^{-\frac{1}{2\sigma_0^2}(x_i - \theta_0)^2} + (1-\rho) \frac{1}{\sqrt{2\pi \sigma_1^2}} e^{-\frac{1}{2\sigma_1^2}(x_i - \theta_1)^2} \right) \right) = 0$$

We could grid search... BUT the amount of points to
plug in can lead to a quadrillion plus samples
which takes forever. Why? For
$\rho \in [0,1], \theta_0 \in [\_\_], \theta_1 \in [\_\_], \sigma_0^2 \in [\_\_], \sigma_1^2 \in [\_\_]$
$\underbrace{\phantom{\rho}}_{1000} \quad \underbrace{\phantom{\theta_0}}_{1000} \quad \underbrace{\phantom{\theta_1}}_{1000} \quad \underbrace{\phantom{\sigma_0^2}}_{1000} \quad \underbrace{\phantom{\sigma_1^2}}_{1000}$

$1000^5$ samples... impossible for a computer to sample! So how do we sample?
STAY TUNED FOR THE NEXT LECTURE.