

Bayesian Statistics: Techniques & Models

Wk 1: Statistical Modelling & Mc

↳ Real models are complicated

MCMC: Allows sampling of post. w/
no analytical solutions

- Expands models possible for bayesians

→ Statistical Modelling

Stats - process of planning, collecting &
analyzing data to draw conclusions

↳ Many applied fields

Model \leftarrow imitate & approx data gen. process

given an estimate

- ← How will number change if measured again?
- ← What patterns are in the noise?

Objectives

- (1) Quantify uncertainty
- (2) Inference
- (3) Measure support for hypo
- (4) Prediction

• Modeling Process

↳ Stat. model process.

- ① Understand the problem
- ② Plan & collect data
- ③ Explore data
- ④ Postulate Model
- ⑤ fit model (This class: bayesian)
- ⑥ Check model
- ⑦ Iterative
- ⑧ Use model

→ Bayesian Modeling

data

↳ heights
of men, $n = 15$

Models can
have increasing
complexity

↳ Model... $y_i = \mu + \varepsilon_i$

↑ error, iid Norm

ie $\varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2) i=1\dots N$

$y_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$

↳ frequentist model

How do we turn freq. model into bayesian model?

↳ add priors to hyperparameters

i.e., μ & σ^2 are RNS

Likelihood \leftarrow describes how data is generated

$$p(y|\theta)$$

Prior \leftarrow describes probability of parameters

$$p(\theta)$$

Posterior \leftarrow beliefs of parameters given y

$$p(\theta|y) = \frac{p(y|\theta)p(\theta)}{p(y)}, p(y) = \int p(y|\theta)p(\theta)d\theta \\ = \int p(y,\theta)d\theta$$

↑
get joint dist &
marginalise over θ

→ Model specification

$$y_i | \mu, \sigma^2 \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2) \quad \text{priors}$$

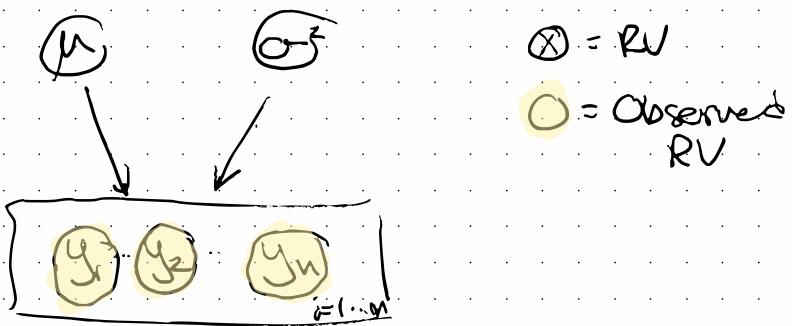
if μ or σ^2 are known

$$\text{prior}(\mu) = N(\mu_0, \sigma_0^2)$$

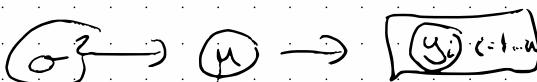
inverse
gamma

$$\text{prior}(\sigma^2) = \text{IG}(\nu_0, \beta_0)$$

convert to graphic representation



↳ you can even make μ dependent on σ



$$\text{s.t. } y_i | \mu, \sigma^2 \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$$

$$\mu | \sigma^2 \sim N(\mu_0, \frac{\sigma^2}{\omega_0})$$

$$\sigma^2 \sim \text{IG}(\nu_0, \beta_0)$$

Hence, joint dist. for posterior =

$$P(y_1, \dots, y_n, \mu, \sigma^2) = P(y_1, \dots, y_n | \mu, \sigma^2) * P(\mu | \sigma^2) * P(\sigma^2)$$

$$= \prod_{i=1}^n \left[N(y_i | \mu, \sigma^2) \right] * N(\mu | \mu_0, \frac{\sigma^2}{w_0})$$

\uparrow Likelihood of data \uparrow priors

thus, equal to numerator in Bayes theorem

↳ posterior is proportional to this.

If this value can be integrated, then work is done.

↳ problem, often non-analytical solutions

→ Non conjugate models

$n = 10$, measuring % change in employee count for a company

$$y_i | \mu \sim N(\mu, 1)$$

$$\mu \sim t(0, 1, 1)$$

$$p(\mu | y_1, \dots, y_n) \propto \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} e^{-\frac{(y_i - \mu)^2}{2}} \frac{1}{\pi(1 + \mu^2)}$$

$$\propto e^{-\frac{1}{2}\sum(y_i - \mu)^2} * \frac{1}{1 + \mu^2}$$

↓
expand & simplify

$$\propto \frac{e^{n(\bar{y}\mu - \mu^2/2)}}{1 + \mu^2}$$

↑ we can find posterior dist $p(\theta)$
but cannot be $\int d\theta$ easily.

→ Prevented Bayesian methods from being main stream for a long time

→ Monte Carlo Estimation

↳ calculating parameters via MC

Say $\theta \sim \text{Gamma}(\alpha, b)$ st $\alpha = 2$

$$b = \frac{1}{3}$$

↙ LOTUS

$$\begin{aligned} E(\theta) &= \int_0^\infty \theta p(\theta) d\theta = \int_0^\alpha \theta \frac{b^\alpha}{\Gamma(\alpha)} \theta^{\alpha-1} e^{-b\theta} d\theta \\ &= \frac{\alpha}{b} \end{aligned}$$

or simulate θ_i^* , $i = 1 \dots n$

$$\text{Avg}(\theta_i^*) \rightarrow E(\theta) \text{ as } n \rightarrow \infty$$

$$= \frac{1}{n} \sum_{i=1}^n \theta_i^* = \bar{\theta}^*$$

$$\text{Var}(\theta_i^*) = \int_0^\alpha (\theta - E(\theta))^2 p(\theta) d\theta$$

continued...

Say we have $h(\theta)$ st $\int h(\theta) p(\theta) d\theta$

$$= E(h(\theta)) \approx h(\bar{\theta}^*)$$

$$\text{ie. } E[h(\theta)] = \int_0^{\infty} \sum_{0 < \theta < 5} (\theta) p(\theta) d\theta$$

$$h(\theta) = \sum_{0 < \theta < 5} (\theta) = \int_0^5 1 \cdot p(\theta) d\theta$$

$$= \Pr[0 < \theta < 5]$$

$$\approx \frac{1}{m} \sum_{i=1}^m \sum_{0 < \theta < 5} (\theta_i^*)$$

↙ people
do this all
the time

→ by CLT, $\bar{\theta}^* \sim N(E(\theta), \frac{\text{var}(\theta)}{m})$

$$\text{var}(\theta) = \frac{1}{m} \sum_{i=1}^m (\theta_i^* - \bar{\theta}_i^*)^2$$

thus

$$\sqrt{\frac{\text{var}(\theta)}{m}} = \text{standard error}$$

Example

$$y | \phi \sim \text{Binom}(10, \phi)$$

$$\phi \sim \text{Beta}(3, 2)$$

$$\text{thus } p(y, \phi) = p(\phi)p(y|\phi)$$

How do we simulate from this posterior?

- repeat
- ① ϕ_i^* from $\text{Beta}(2, 2)$
 - ② w/ given ϕ_i^* , draw $y_i^* \sim \text{Bin}(10, \phi_i^*)$
 - ③ generates (y_i^*, ϕ_i^*)
- ↑
draws from posterior distribution

→ Markov Chains

↳ Markov chains have sequences of observations w/ the property that $p(x_t | x_{t-1})$

i.e. $p(x_1, x_2, \dots, x_t) = p(x_1) p(x_2 | x_1) \dots p(x_t | x_{t-1})$

↳ i.e. you have a coin & as you flip it you increment/decrement a number

Where this number goes is called a random walk

If we have multiple hidden states, we can calculate the prob of one to another w/ a transition matrix



$$\begin{bmatrix} 1 \rightarrow 1 & 1 \rightarrow 2 & 1 \rightarrow 3 \\ 2 \rightarrow 1 & 2 \rightarrow 2 & 2 \rightarrow 3 \\ 3 \rightarrow 1 & 3 \rightarrow 2 & 3 \rightarrow 3 \end{bmatrix}$$

↖ rows of mat M
add b /

A or α

Applying A repeatedly will sometimes make it converge to steady state Z, find w/ Lin Alg

→ Metropolis Hastings

↳ Allows us to sample from prob dist
using MC w/ stationary dist = target dist

i.e. our posterior

$$\text{ie } p(\theta) \propto g(\theta) \dots p(\theta) = c g(\theta)$$

↑ c is unknown

Metropolis Hastings

① Select θ_0 .

② for i ... n:

a) draw candidate $\theta^* \sim q(\theta^* | \theta_{i-1})$

$$b) \alpha = \frac{q(\theta^*) / q(\theta^* | \theta_{i-1})}{q(\theta_{i-1}) / q(\theta_{i-1} | \theta^*)}$$

$$= \frac{q(\theta^*) q(\theta_{i-1} | \theta^*)}{q(\theta_{i-1}) q(\theta^* | \theta_{i-1})}$$

$$q(\theta_{i-1}) q(\theta^* | \theta_{i-1})$$

c) if $\alpha \geq 1$, accept θ^* , $\theta_i := \theta^*$

if $0 < \alpha < 1$, " ", $\theta_i := \theta^*$ w/ prob α

else reject θ^* & set $\theta_i := \theta_{i-1}$ w/ prob 1 - α

a markov chain

Notes on candidate generating distribution

$$q(\theta^* | \theta_t)$$

↳ generates candidates to step toward

we can use different distributions.

→ Random Walk Met-Hastings

uses normal distribution

$$\therefore q(\theta_{t+1} | \theta^*) = q(\theta^* | \theta_{t+1})$$

$$\kappa = \frac{q(\theta^*)}{q(\theta_{t+1})} \quad \begin{matrix} \leftarrow \text{evaluating backwards} \\ \text{since } q \text{ evaluates to } p \end{matrix}$$

This should only be done if $q \sim N$ approximates p well.

Consider this problem,

(lets also say

Θ is fair

1 is loaded

$\Theta \in \{\text{fair, loaded}\}$ (coin example)

Prior, $P(\Theta = \text{loaded}) = .6$

Likelihood, $f(x|\Theta) = \begin{cases} \binom{5}{x} \left(\frac{1}{2}\right)^5 & \{\Theta = \text{fair}\} \\ + \binom{5}{x} (.7)^x (.3)^{5-x} & \{\Theta = \text{loaded}\} \end{cases}$

post, $P(\Theta = \text{loaded} | X=2) = .388$

↑ can we simulate
this w/ mcmc?

- 1) Start @ either $\Theta_0 = \text{fair}$ or $\Theta_0 = \text{loaded}$
- 2) for $i = 1, \dots, m$

a) Propose Θ^* to be the other state
as Θ_{t-1}

b) $\alpha = \frac{g(\Theta^*) / g(\Theta^* | \Theta_{i-1})}{g(\Theta_{i-1}) / g(\Theta_{i-1} | \Theta^*)}$

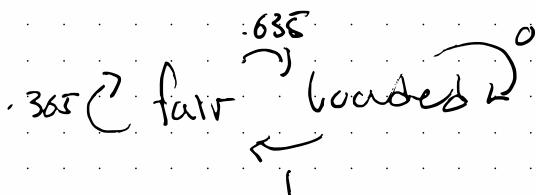
$$= \frac{f(x=2 | \Theta^*) / 1}{f(x=2 | \Theta_{i-1}) / 1} \quad \begin{array}{l} \leftarrow \text{simple case} \\ \text{prob 1 propose } \Theta \\ \text{as other state.} \end{array}$$

$$\text{If } \theta^* = \text{loaded}, \alpha = \frac{0.0744}{0.125} = 0.635 \quad \begin{matrix} L' \text{ accept} \\ \text{w/prob} \end{matrix}$$

$$= 0.635$$

$$\theta^* = \text{fair}, \alpha = \frac{0.125}{0.0744} = 1.574$$

\uparrow always accept



$$\left(\begin{array}{cc} 0.365 & 0.635 \\ 1 & 0 \end{array} \right), \pi = [0.612 \ 0.388]$$

Q

$$\hookrightarrow \text{Stationary dist} = [0.612 \ 0.388]$$

So, it works

\hookrightarrow we only need to maximise post via est w/ prior & likelihood
no need for norm. const.

Side bar: Chain rule of probability

So far we've used the rule of pr.,
so, what is it?

$$P(A, B) = P(A|B)P(B) \text{ by definition}$$

↳ gives bayes form: $P(A|B) = \frac{P(B|A)P(A)}{P(B)}$ ✓

Chain rule of prob.

If A_1, \dots, A_n & $P(A_1 \cap \dots \cap A_{n-1}) > 0$

$$P(A_1 \cap \dots \cap A_n) = P(A_1) P(A_2 | A_1) P(A_3 | A_1 \cap A_2) \dots P(A_n | A_1 \cap \dots \cap A_{n-1})$$

proof by induction:

If $n=2$ then chain rule holds

$$P(A_1 \cap A_2) = P(A_1) P(A_2 | A_1)$$

Let $B = A_1 \cap \dots \cap A_{n-1}$

$$P(B) = P(A_1) P(A_2 | A_1) \dots P(A_{n-1} | A_1 \cap \dots \cap A_{n-2})$$

$$P(B \cap A_n) = P(A_n | B) P(B) \quad \blacksquare$$

→ Gibbs Sampling

↳ We saw how to draw post w/ 1 variable parameter

What about multiple parameters

↳ say $p(\theta, \phi | y) \propto q(\theta, \phi | y)$

↳ update values 1 at a time using new values for each draw

JAGS → Just Another Gibbs Sampler

Note that

$$P(\theta, \phi | y) = \underbrace{P(\phi | y)}_{\text{no } \theta} P(\theta | \phi, y)$$

by chain rule of probability

↳ $p(\theta | \phi, y) \propto p(\theta, \phi | y) \propto q(\theta, \phi)$

$$p(\phi | \theta, y) \propto p(\theta, \phi | y) \propto q(\theta, \phi)$$

Basically: for each unknown param, start w/ each one & update one at a time

→ Algorithm

① Init. Θ_0, ϕ_0

② for i in $1 \dots n$

a) using Θ_{i-1} , draw $\Theta_i \sim P(\Theta | \phi_{i-1}, y)$

b) using Θ_i , draw $\phi_i \sim P(\phi | \Theta_i, y)$

(Θ_i, ϕ_i)

Example

↳ Similar var-cov dist from before

$y_i | \mu, \sigma^2 \sim N(\mu, \sigma^2), i = 1, \dots, n$

$\mu \sim N(\mu, \sigma_0^2)$ ← we have multiple
 $\sigma^2 \sim \text{Ga}(\nu_0, \beta_0)$ parameters, which
is why we need Gibbs

$P(\mu, \sigma^2 | y_1, \dots, y_n)$

$\propto p(y_1, \dots, y_n | \mu, \sigma^2) p(\mu) p(\sigma^2)$

↑ What does this density look like?

$$\text{post} \propto \prod_{i=1}^n [N(y_i | \mu, \sigma^2)] N(\mu | \mu_0, \sigma_0^2) \text{Ga}(\sigma^2 | \nu_0, \beta_0)$$

↪ insert actual densities & drop proportionality
consts.

$$\propto (\sigma^2)^{-\frac{n}{2}} e^{-\sum_{i=1}^n \frac{(y_i - \mu)^2}{2\sigma^2}} e^{-\frac{(\mu - \mu_0)^2}{2\sigma^2}} e^{-\frac{-(\nu_0 + 1)}{\sigma^2}} e^{-\frac{-\beta_0}{\sigma^2}}$$

full posterior up to proportionality

Now that we have post to prop. we can change
& as necessary for Gibbs sampling

ie

$$P(\mu | \sigma^2, y_1, \dots, y_n)$$

$$\propto e^{-\sum_{i=1}^n \frac{(y_i - \mu)^2}{2\sigma^2}} e^{-\frac{(\mu - \mu_0)^2}{2\sigma^2}} \quad \checkmark$$

&

$$P(\sigma^2 | \mu, y_1, \dots, y_n)$$

$$\propto (\sigma^2)^{-\frac{n}{2}} e^{-\sum_{i=1}^n \frac{(y_i - \mu)^2}{2\sigma^2}} (\sigma^2)^{-\frac{-(\nu_0 + 1)}{2}} e^{-\frac{-\beta_0}{\sigma^2}} \quad \checkmark$$

↪ Both can be simplified further
but not so important.

↪ they have recognizable
distributions to sample
from

Note: What happens if we can't draw from dist directly like we did for the previous example?

↳ Place mcmc algorithm inside of Gibbs sampling.

→ Assessing Convergence

↳ trace plot → Shows where sampler is spending its time

- Shouldn't show long term trends

- No wandering

↳ though, you can still get convergence even when wandering

this can be assessed using **auto correlation**

AC will be high when chain is searching the space

↳ When good value is found, autocorr will be low. (ie Randomness dictates search)

Effective Sample Size

↳ How many indep samples from stationary dist you would need to draw to have equiv info in your chain.

coda has function **effectiveSize**

↳ check my effective sample size goes a lot of M.F., so it is worth doing.

you can also use **raftery diagnostic**

↳ Shows how much more samples you'd need for some accuracy w/ some prob for a quantile est

↳ vs how many from completely indep chain.

↳ Make sure to check documentation

What happens if it takes long time to sample from stationary dist?

↳ remove samples as "burn in" period

We can find this visually on a trace plot

Gelman-Rubin diagnostic

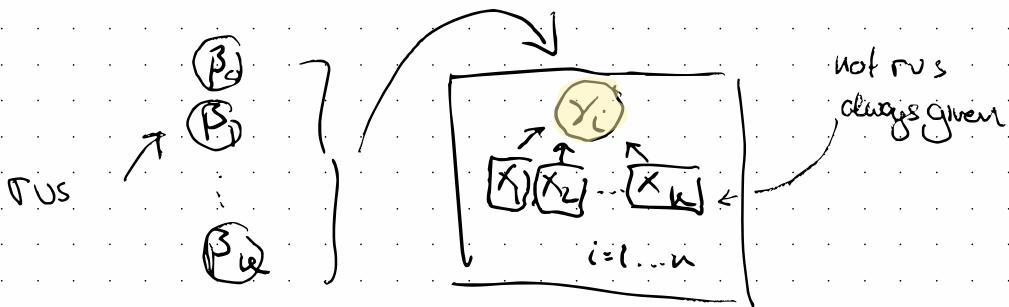
↳ checks if multiple chains have converged. Should give value close to 1.

↳ function `gelman.diag`

\rightarrow Bayesian Linear Regression

$$y_i = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + \epsilon_i, \epsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2), i=1 \dots n$$

$$y_i | x_i, \beta, \sigma^2 \sim N(\beta x_i, \sigma^2)$$

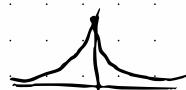


β ← while holding all other x const. a 1-fold
rise in x will change the mean of y_i by β

Note that all β & σ have their own priors

Note that adding Laplace prior to
 β is Lasso regression

$$\hookrightarrow p(\beta) = \frac{1}{Z} e^{-\|\beta\|}$$



→ ANOVA (Analysis of Variance)

"Factors"

Categorical vars, obs belong to groups

- ↳ comparing responses inside vs betw. groups
- ↳ If var betw > ins, conclude grouping effect

re	Sand	Food size	(for website)
	music	small	
	non-music	med.	$2 \times 3 = 6$ total combes
		large	

→ model: $y_{ij} | g_i, \mu, \sigma^2 \sim N(\mu_{g_i}, \sigma^2)$

$$g_i \in \{1, \dots, G\}$$

$$j = 1, \dots, n$$

Alternative

$$E(y_{ij}) = \beta_0 + \beta_1 X_{1,i} + \dots + \beta_{G-1} X_{G-1,i}$$

↑
dummy vars (indicates that group $i = \text{true}$)

Note: Alternative formulation shows how ANOVA
is related to linear regression

consider {sound: A, fast: b?}

		B		
		1	2	3
A	1	μ_{11}	μ_{12}	μ_{13}
	2	μ_{21}	μ_{22}	μ_{23}

Cell means model

$$E(y_i) = \mu + \alpha_2 I_{(a_i=2)} + \beta_2 I_{(b_i=2)} \text{ Additive Model}$$

$$+ \beta_3 I_{(b_i=3)}$$

Appropriate when no interaction between independent variables

Logistic Regression

What about classification tasks

↳ linear regression doesn't quite work

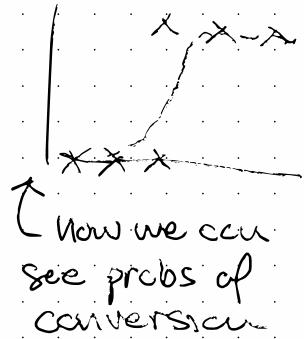
likelihood here is bernoulli

$$y_i | (\phi_i \text{ ind Bern}(\phi_i)) \quad i=1, \dots, n$$

$$\text{Odds reg: } E(y_i) = \beta_0 + \beta_1 x_i \quad \begin{array}{l} \leftarrow \text{use CMC function} \\ \text{to bound } E(y_i) \text{ between 0 & 1} \end{array}$$

$$\frac{1}{1 + e^{-(\beta_0 + \beta_1 x_i)}}$$

Note that this is the result of using
 $\log\left(\frac{y}{1-y}\right)$ as the link (logit)



Poisson Regression

What about count data?

↳ bin reg? Problem, counts aren't -ve

Variance might not be constant

$$y_i \sim \text{Pois}(\lambda_i) \quad (i=1, \dots, n)$$

Direct reg: $E(y_i) = \beta_0 + \beta_1 x_{1,i}$ ↪ same problem w/ logistic regression

use log link

$$\log(\lambda_i) = \beta_0 + \beta_1 x_1$$

Note

$$\Rightarrow E(y_i) = \lambda_i = e^{\beta_0 + \beta_1 x_{1,i}} \quad \leftarrow E(\log(y)) \neq \log(E(y))$$

Prior Sensitivity Analysis

↳ Detect dependence of result on prior

- Why did you choose the prior you did

- How do different priors change things?

One thing you can try is to create a skeptical prior that disfavours your hypothesis.

↳ if you can still support your own hypothesis then you have good justifications to think so.

Hierarchical Models

So far all observations have been independent

Not always true

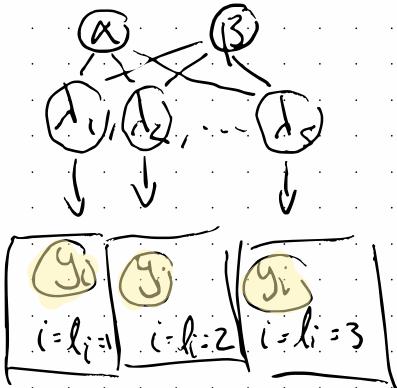
Example: 15C cookies. 3C from 5 bcs.

fully indep. $y_{il} \sim \text{Pois}(\lambda_l)$ ($i = 1, \dots, 50$)

loc dep. $y_{il} | \lambda_i, \lambda_l \sim \text{Pois}(\lambda_{li})$ $\lambda_i \in \{1, \dots, 5\}$

$\lambda_l | \alpha, \beta \sim \text{Ga}(\alpha, \beta)$ $\lambda_l = 1, \dots, 5$

Graphical Representation



Alternatives

- Make all lectures independent
- Ignore differences between lectures

- fit 5 models for all.

ignores data from other lectures

Linear Regression As Hierarchical Model

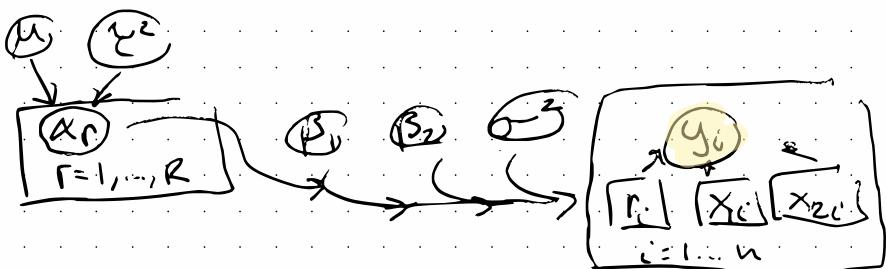
$$c = 1, \dots, n$$

Linear model: $y_i | \alpha_c, \beta, \sigma^2 \text{ iid } N(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2}, \sigma^2)$

Random intercept model, change β_c

$y_{ir} | \beta_1, \beta_2, \alpha_r, \sigma^2 \text{ iid } N(\alpha_r + \beta_1 x_{r1} + \beta_2 x_{r2}, \sigma^2) \quad r = 1, \dots, R$

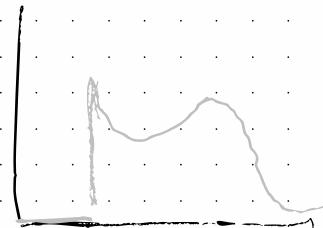
$\alpha_r | \mu, \tau^2 \sim N(\mu, \tau^2) \quad r = 1, \dots, R$



Mixture Models

Problem, what if data doesn't follow a standard distribution?

- ↳ Mixture dists are a combo of multiple distributions



i.e. a PDF for the graph above might look like this:

$$p(y) = .4 e^{-y} I_{y \geq 0} + .6 \frac{1}{2\pi} e^{-\frac{(y-3)^2}{2}}$$

- ↳ we can see how these dists combine here →



The general form for an MM is

$$\hookrightarrow p(y) = \sum_{j=1}^J \omega_j \cdot f_j(y) \quad j = 1, \dots, J \text{ pdfs}$$

$$\text{s.t. } \sum_{j=1}^J \omega_j = 1, \quad f_j(y) \in \{\text{PDF}\}$$

we can simulate from ~~this~~^{this} PM by selecting
from either $f_j(y)$ w/ w_j probability

↳ i.e. $[j=1, j=2]$ w/ prob [.4, .6]

when j is selected, take rv from $f_j(y)$

Issue: When we have real data, we
can't see j , latent variable

↳ we can use bayesian hierarchical
models to infer them