

# **An introduction to Bayesian regression**

Guest lecture for Psych 610/710 at UW–Madison

---

Tristan Mahr

April 13, 2017

Slides and R code that produced them are online:

[https://github.com/tjmahr/Psych710\\_BayesLecture](https://github.com/tjmahr/Psych710_BayesLecture)

I gave a similar, more code-heavy version of this talk to the R Users Group: [https://github.com/tjmahr/MadR\\_RStanARM](https://github.com/tjmahr/MadR_RStanARM)

- A little about me and how I got into Bayes
- Mathematical intuition building
- Bayesian updating
- Fitting a model with RStanARM
- Big takeaway ideas

# Background

---

- I am dissertator in Communication Sciences and Disorders
- I study word recognition in preschoolers
- For statistics, I mostly do multilevel logistic regression models
- R enthusiast

# I was once in your shoes

I learned stats and R in this course with Markus Brauer and John Curtin.

I *still* refer to the slides from this course on contrast codes.

But now I'm a “Bayesian”.

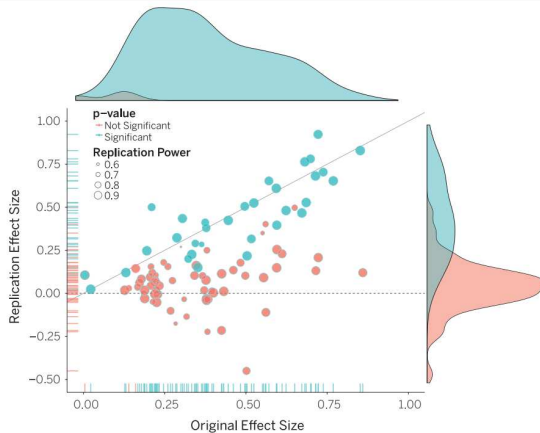
## A timeline

---

Open Science Collaboration (2015) tries to replicate 100 studies published in 3 psychology different journals in 2008.

- Boil a study down to 1 test statistic and 1 effect size.
- Replicate the study.
- Compare replication's test statistic and effect size against original.



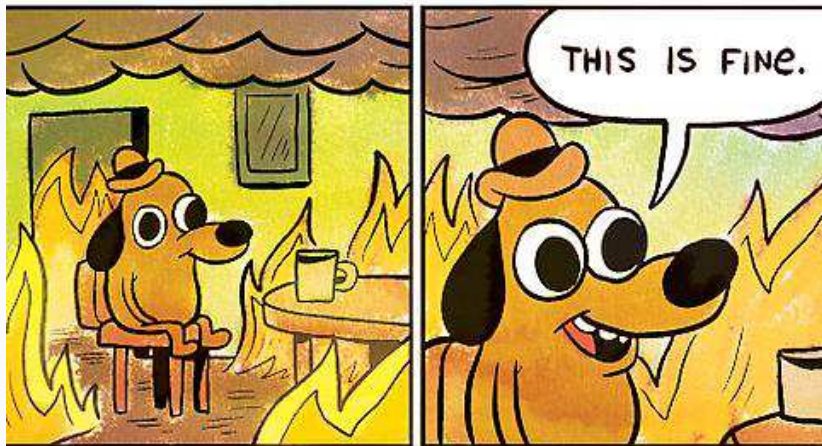


**Original study effect size versus replication effect size (correlation coefficients).** Diagonal line represents replication effect size equal to original effect size. Dotted line represents replication effect size of 0. Points below the dotted line were effects in the opposite direction of the original. Density plots are separated by significant (blue) and nonsignificant (red) effects.

**Figure 1:** Scatter plot of original vs replicated effect sizes

- Approximately 36% of the studies are replicated (same test statistic).
- On average, effect sizes in replications are half that of the original studies.

## Reactions



**Figure 2:** I don't know how to turn off the figure labeling feature

- We're doomed.

- We're doomed.
- Most findings are probably false, and we knew that already.

- We're doomed.
- Most findings are probably false, and we knew that already.
- No, this is business as usual.

- We're doomed.
- Most findings are probably false, and we knew that already.
- No, this is business as usual.
- Any credible discipline has to do this kind of house-cleaning from time to time.

# Lots of hand wringing and soul searching

Some reactionary:

- Replication creates an industry for incompetent hacks.
- Here come the methodological terrorists!

Some constructive:

- Everything is f'ed – so what else is new?
- Increased rigor and openness are a good thing.



## Crisis made me think more about questionable practices

All those unintentional acts and rituals to appease the Statistical Significance gods.

## Crisis made me think more about questionable practices

All those unintentional acts and rituals to appease the Statistical Significance gods.

**HARKing** Hypothesizing after results are known.  
Telling a story to fit the data.

## Crisis made me think more about questionable practices

All those unintentional acts and rituals to appease the Statistical Significance gods.

**HARKing** Hypothesizing after results are known.

Telling a story to fit the data.

**Garden of forking data** Conducting countless sub-tests and sub-analyses on the data.

## Crisis made me think more about questionable practices

All those unintentional acts and rituals to appease the Statistical Significance gods.

**HARKing** Hypothesizing after results are known.

Telling a story to fit the data.

**Garden of forking data** Conducting countless sub-tests and sub-analyses on the data.

***p*-hacking** Doing these tests in order to find a significant effect.

## Crisis made me think more about questionable practices

All those unintentional acts and rituals to appease the Statistical Significance gods.

**HARKing** Hypothesizing after results are known.

Telling a story to fit the data.

**Garden of forking data** Conducting countless sub-tests and sub-analyses on the data.

***p*-hacking** Doing these tests in order to find a significant effect.

**Selective reporting** Reporting only the tests that yielded a significant result.

The usual way of doing things is *insecure*.

The usual way of doing things is *insecure*.

- Perfectly fine if you know what you're doing.

The usual way of doing things is *insecure*.

- Perfectly fine if you know what you're doing.
- Works great if you pre-register analyses. Provides error control.



The usual way of doing things is *insecure*.

- Perfectly fine if you know what you're doing.
- Works great if you pre-register analyses. Provides error control.
- But vulnerable to exploitation.

The usual way of doing things is *insecure*.

- Perfectly fine if you know what you're doing.
- Works great if you pre-register analyses. Provides error control.
- But vulnerable to exploitation.
- And many people don't know what they're doing.

## My response to the crisis

I want to avoid these questionable practices.

I want to level up my stats and explore new techniques.

- Maybe more robust estimation techniques?
- Maybe machine learning techniques to complement conventional analyses?

## My response to the crisis

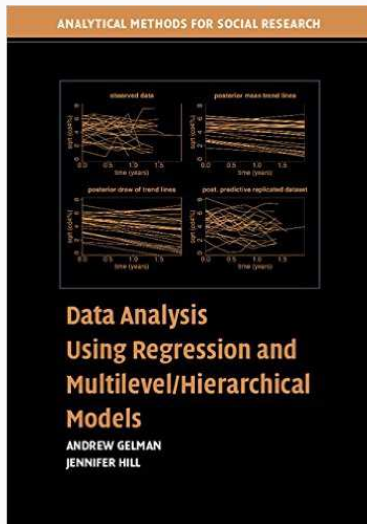
I want to avoid these questionable practices.

I want to level up my stats and explore new techniques.

- Maybe more robust estimation techniques?
- Maybe machine learning techniques to complement conventional analyses?

I want something less finicky than statistical significance.

- $p$ -values don't mean what many people think they mean.
- Neither do confidence intervals.
- Statistical significance is not related to practical significance.



**Figure 3:** Cover of Data Analysis USING Regression and Multilevel/Hierarchical Models

I started reading the Gelman and Hill book.

- This is the book for the `arm` package.
- Still the best treatment of multilevel models in R despite being 10 years old.

It emphasizes estimation, uncertainty and simulation.

I started reading the Gelman and Hill book.

- This is the book for the `arm` package.
- Still the best treatment of multilevel models in R despite being 10 years old.

It emphasizes estimation, uncertainty and simulation.

Midway through, the book pivots to Bayesian estimation.  
(Multilevel models are kinda Bayesian because they borrow information across different clusters.)

I'm down a rabbit hole, writing Stan (Bayesian) models to fit the models from the ARM book, and there is an influx of Bayesian tools for R.

- Statistical Rethinking, a book that reteaches regression from a Bayesian perspective with R and Stan, is released.
- New version of brms is released. This package converts R model code into Stan programs.
- RStanARM is released.
- A blog post circulates: "R Users Will Now Inevitably Become Bayesians".

I eat all this up. I become a convert.



## Long story short

The replication crisis sparked my curiosity, and a wave of new tools and resources made it really easy to get started with Bayesian stats.

My goal with this approach has been to make better, more honest scientific summaries of observed data.

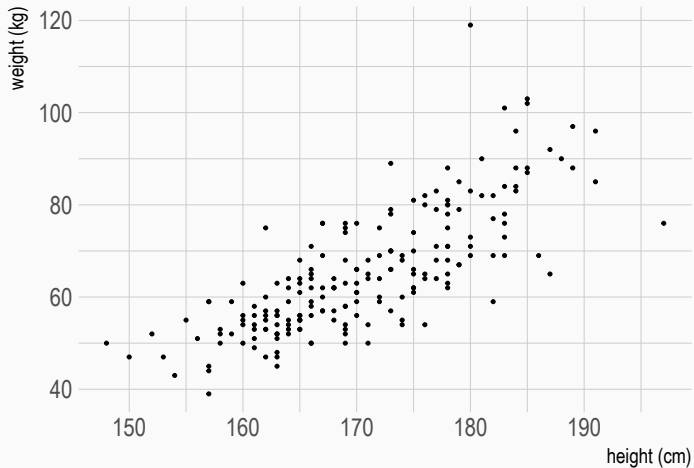
# Classical regression versus Bayesian regression in a few plots

---

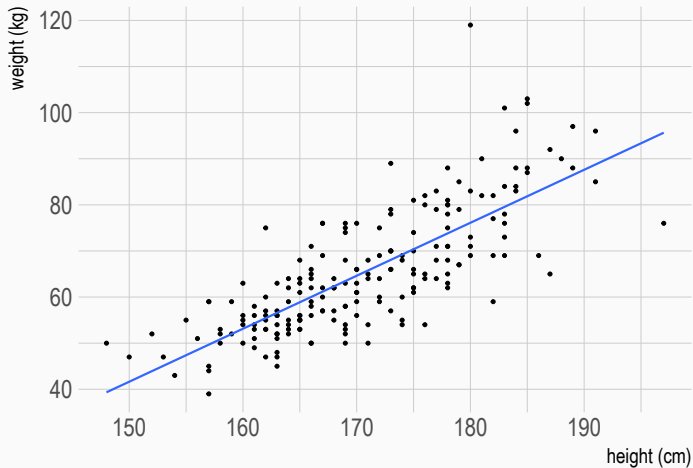
# The data

```
# Some toy data
davis <- car::Davis %>% filter(100 < height) %>% as_data_frame()
davis
#> # A tibble: 199 × 5
#>       sex weight height repwt repht
#>   <fctr>   <int>   <int> <int> <int>
#> 1      M     77    182     77    180
#> 2      F     58    161     51    159
#> 3      F     53    161     54    158
#> 4      M     68    177     70    175
#> 5      F     59    157     59    155
#> 6      M     76    170     76    165
#> 7      M     76    167     77    165
#> 8      M     69    186     73    180
#> 9      M     71    178     71    175
#> 10     M     65    171     64    170
#> # ... with 189 more rows
```

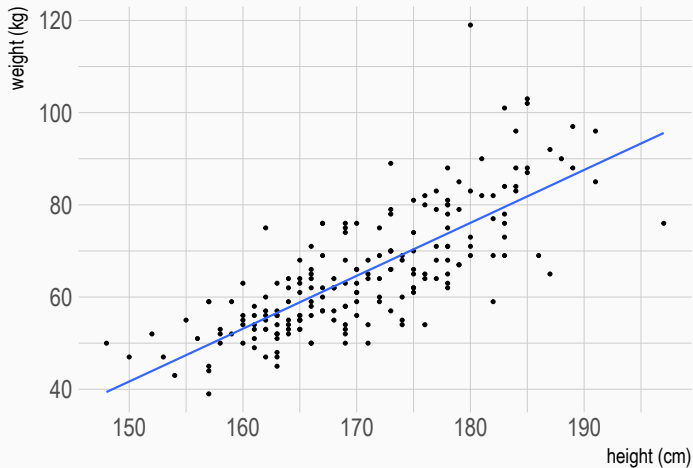
# The data



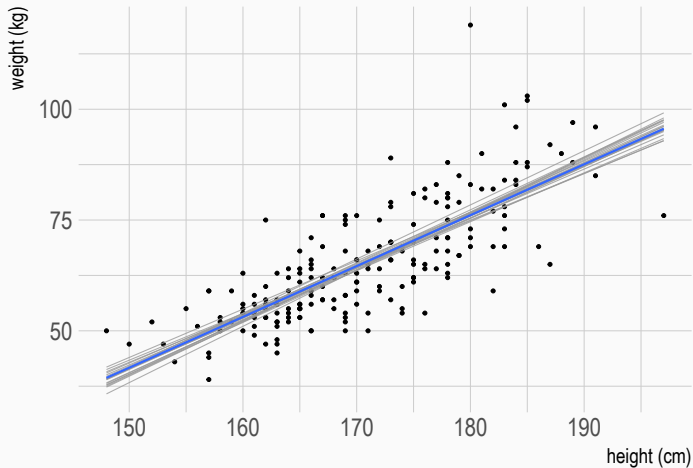
## Classical model provides the line of best fit



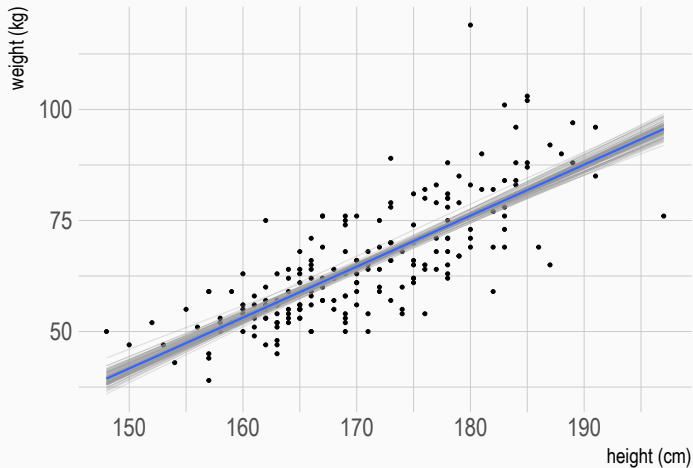
## Bayesian model's median line of fit



## Median line and 20 other lines from posterior



## Median line and 100 other lines from posterior





# Summary

- Classical: There is a single “true” line of best fit, and I’ll give my best estimate of it.
- Bayesian: There is a distribution of lines of fit—some more plausible than others—and I’ll give you samples from that distribution.

# Building mathematical intuitions

---

# Caveat

- These slides and these examples are meant to illustrate the pieces of Bayes theorem.
- This is not a rigorous mathematical description of Bayesian probability or regression.

$p(A \mid B)$  : probability of A given B

## Conditional probability review

$p(A \mid B)$  : probability of A given B

Suppose that 95% of emails with the phrase “investment opportunity” are spam.

## Conditional probability review

$p(A \mid B)$  : probability of A given B

Suppose that 95% of emails with the phrase “investment opportunity” are spam.

$$p(\text{spam email} \mid \text{"investment opportunity"}) = .95$$

What would this probability express?

$$p(\text{"investment opportunity"} \mid \text{spam email})$$

## Conditional probability review

What would this probability express?

$$p(\text{"investment opportunity"} \mid \text{spam email})$$

That ordering matters.  $p(A \mid B)$  is not the same as  $p(B \mid A)$ .



# Bayes' theorem

A theorem about conditional probability.

$$p(B \mid A) = \frac{p(A \mid B) * p(B)}{p(A)}$$

I can never remember this equation with letters. Here's how I prefer to write it.

$$p(\text{hypothesis} \mid \text{data}) = \frac{p(\text{data} \mid \text{hypothesis}) * p(\text{hypothesis})}{p(\text{data})}$$

$$p(\text{hypothesis} \mid \text{data}) = \frac{p(\text{data} \mid \text{hypothesis}) * p(\text{hypothesis})}{p(\text{data})}$$

The “hypothesis” is typically something unobserved or unknown. It’s what you want to learn about using the data.

$$p(\text{hypothesis} \mid \text{data}) = \frac{p(\text{data} \mid \text{hypothesis}) * p(\text{hypothesis})}{p(\text{data})}$$

The “hypothesis” is typically something unobserved or unknown. It’s what you want to learn about using the data.

For regression models, the “hypothesis” is a parameter (intercept, slopes or error terms).

Bayes theorem tells you the probability of the hypothesis given the data.

How plausible is some hypothesis given the data?

$$p(\text{hypothesis} \mid \text{data}) = \frac{p(\text{data} \mid \text{hypothesis}) * p(\text{hypothesis})}{p(\text{data})}$$

How plausible is some hypothesis given the data?

$$p(\text{hypothesis} \mid \text{data}) = \frac{p(\text{data} \mid \text{hypothesis}) * p(\text{hypothesis})}{p(\text{data})}$$

Pieces of the equation:

$$\text{posterior} = \frac{\text{likelihood} * \text{prior}}{\text{average likelihood}}$$

## Classifying emails

I got an email with the word “cialis” in it. Is it spam?

I got an email with the word “cialis” in it. Is it spam?

- What I want to know is spam-ness (versus ham-ness).
- What I have is an email with the word “cialis”.



# Classifying emails

I got an email with the word “cialis” in it. Is it spam?

- What I want to know is spam-ness (versus ham-ness).
- What I have is an email with the word “cialis”.

$$P(\text{spam} \mid \text{"cialis"}) = \frac{P(\text{"cialis"} \mid \text{spam}) * P(\text{spam})}{P(\text{"cialis"})}$$

## Email example

The two unconditional probabilities are base rates that need to be accounted for.

## Email example

The two unconditional probabilities are base rates that need to be accounted for.

The prior is the frequency of spam in general. The average likelihood is the frequency of the word “cialis” in emails.

$$P(\text{spam} \mid \text{"cialis"}) = \frac{\text{"cialis" freq. in spam} * \text{spam rate}}{\text{"cialis" freq.}}$$

## “Bayesianism”

Some people would argue that using Bayes theorem is not “Bayesian”. After all, in this example, we’re just counting the frequency of events.

It’s kind of weird, but it is also true.

Simple event-counting is not what people usually mean by the word “Bayesian”.

## The “Bayesianism” form of Bayes’ theorem

$$\text{updated information} = \frac{\text{likelihood of data} * \text{prior information}}{\text{average likelihood of data}}$$

Bayes’ theorem provides a systematic way to update our knowledge as we encounter new data.

## The “Bayesianism” form of Bayes’ theorem

$$\text{updated information} = \frac{\text{likelihood of data} * \text{prior information}}{\text{average likelihood of data}}$$

Bayes’ theorem provides a systematic way to update our knowledge as we encounter new data.

$$\text{updated beliefs} \propto \text{likelihood of data} * \text{prior beliefs}$$

- Update your beliefs in proportion to how well the data fits those beliefs.
- Your beliefs have probabilities. You can quantify your uncertainty about what you know.

**Okay, but what is likelihood?**

---

*Sidenote: This is nifty. A lot of my stats training made more sense once I had a broader understanding of likelihood.*



# First, what are models?!

What is a statistical model?

# First, what are models?!

What is a statistical model?

It's a description of how the data could have been generated.

*IQ scores are normally distributed.*

## IQ example

*IQ scores are normally distributed.*

$$IQ_i \sim \text{Normal}(\underbrace{\mu}_{\text{mean}}, \underbrace{\sigma}_{\text{SD}})$$

(The  $\sim$  means “sampled from” or “drawn from”.)

$\mu$  and  $\sigma$  are parameters for this model that change the center and spread of the normal bell curve.

The normative IQ model has  $\mu = 100$  and  $\sigma = 15$ .

## Likelihood measures fit

How likely are the data in a given model?

## Likelihood measures fit

How likely are the data in a given model?

I never see it explained this way, but I think of likelihood as “fit”.

How the well data fits in a given model.

## An IQ example

We found some IQ scores in an old, questionable dataset.

```
library(dplyr)
iqs <- car::Burt$IQbio
iqs
#> [1] 82 80 88 108 116 117 132 71 75 93 95 88 111 63
#> [15] 77 86 83 93 97 87 94 96 112 113 106 107 98
```

## An IQ example

We found some IQ scores in an old, questionable dataset.

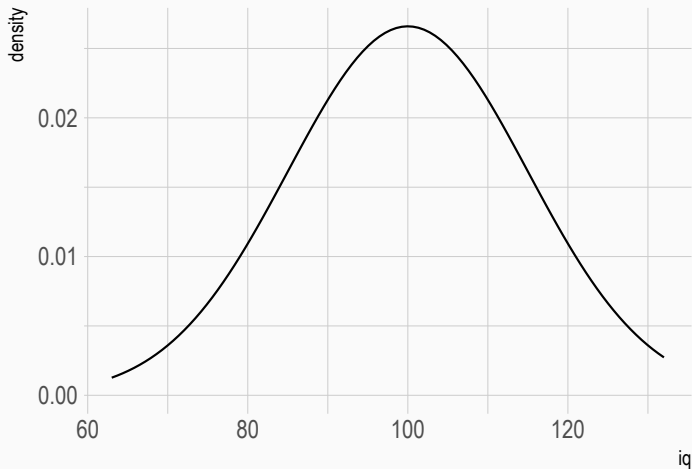
```
library(dplyr)
iqs <- car::Burt$IQbio
iqs
#> [1] 82 80 88 108 116 117 132 71 75 93 95 88 111 63
#> [15] 77 86 83 93 97 87 94 96 112 113 106 107 98
```

IQs are designed to have a normal distribution with a population mean of 100 and an SD of 15.

How well do these data *fit* in that kind of bell curve?

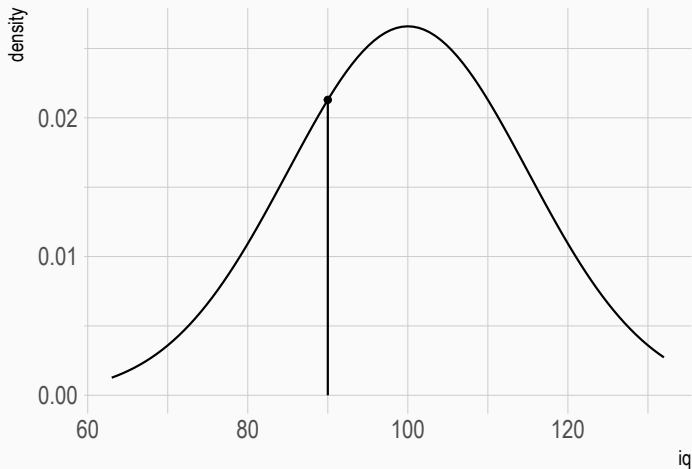


## Density as height on a bell curve



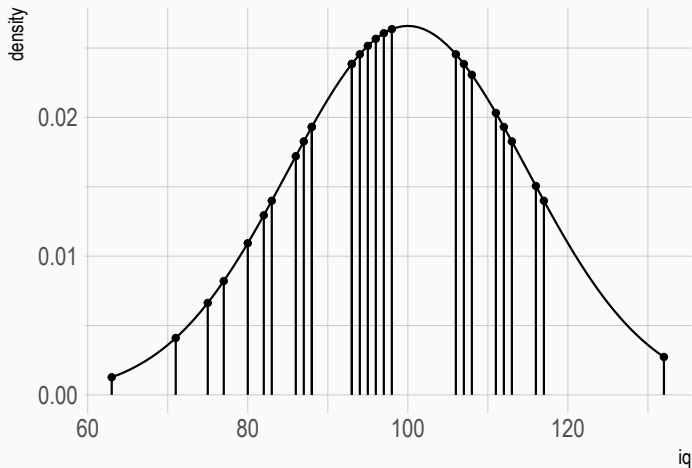
**Figure 4:** A hypothetical bell curve with a mean of 100 and SD of 15.

## Density measures likelihood

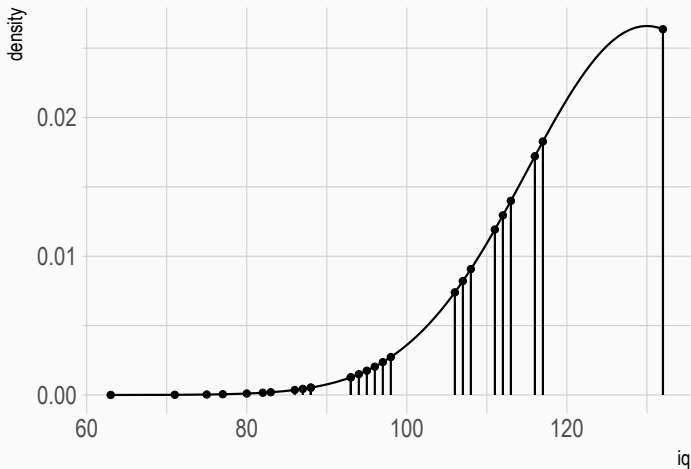


**Figure 5:** Likelihood of an IQ of 90

- Height of each point on curve is density around that point.
- Higher density regions are more likely.
- Data farther from peak density is less likely.



**Figure 6:** Density of IQ scores drawn a bell curve with mean 100.



**Figure 7:** Density of IQ scores drawn a bell curve with mean 130. The fit is terrible.

Density function `dnorm(xs, mean = 100, sd = 15)` tells us the height of each value in `xs` when drawn on a normal bell curve.

```
# likelihood (density) of each point
```

```
dnorm(iqs, 100, 15) %>% round(3)
```

```
#> [1] 0.013 0.011 0.019 0.023 0.015 0.014 0.003 0.004 0.007
```

```
#> [10] 0.024 0.025 0.019 0.020 0.001 0.008 0.017 0.014 0.024
```

```
#> [19] 0.026 0.018 0.025 0.026 0.019 0.018 0.025 0.024 0.026
```

Density function `dnorm(xs, mean = 100, sd = 15)` tells us the height of each value in `xs` when drawn on a normal bell curve.

```
# likelihood (density) of each point
```

```
dnorm(iqs, 100, 15) %>% round(3)
```

```
#> [1] 0.013 0.011 0.019 0.023 0.015 0.014 0.003 0.004 0.007
```

```
#> [10] 0.024 0.025 0.019 0.020 0.001 0.008 0.017 0.014 0.024
```

```
#> [19] 0.026 0.018 0.025 0.026 0.019 0.018 0.025 0.024 0.026
```

Likelihood of all points is the product. These quantities get vanishingly small so we sum their logs instead. (Hence, **log-likelihoods.**)

```
# 2 * 10-50 is vaaaaaaanishingly small!
```

```
prod(dnorm(iqs, 100, 15))
```

```
#> [1] 2.276823e-50
```

```
# log scale
```

```
sum(dnorm(iqs, 100, 15, log = TRUE))
```

```
#> [1] -114.3065
```

Log-likelihoods provide a measure of how well the data fit a given normal distribution.

Which mean best fits the data? Below average IQ (85), average IQ (100), or above average IQ (115)? (Higher is better.)



Log-likelihoods provide a measure of how well the data fit a given normal distribution.

Which mean best fits the data? Below average IQ (85), average IQ (100), or above average IQ (115)? (Higher is better.)

```
sum(dnorm(iqs, 85, 15, log = TRUE))  
#> [1] -119.0065  
sum(dnorm(iqs, 100, 15, log = TRUE))  
#> [1] -114.3065  
sum(dnorm(iqs, 115, 15, log = TRUE))  
#> [1] -136.6065
```

```
sum(dnorm(iqs, 85, 15, log = TRUE))  
#> [1] -119.0065  
sum(dnorm(iqs, 100, 15, log = TRUE))  
#> [1] -114.3065  
sum(dnorm(iqs, 115, 15, log = TRUE))  
#> [1] -136.6065
```

Of these three, the data fit best with the “population average” mean (100).

We just used a **maximum likelihood** criterion to choose among these alternatives!

## Likelihood summary

We have some model of how the data could be generated. This model has tuneable parameters.

*The IQs are drawn from a normal distribution with an SD of 15 and some unknown mean.*

Likelihood is how well the observed data fit in a particular data-generating model.

Classical regression's "line of best fit" finds model parameters that maximize the likelihood of the data.

# Bayesian models

---

## Bayesian updating

Let's consider all integer values from 70 to 130 as equally probable means for the IQs. This is a flat or uniform prior.

Here's our model.

$$IQ_i \sim \text{Normal}(\mu, \sigma = 15) \quad [\text{likelihood}]$$

$$\mu \sim \{\text{integers from 70 to 130}\} \quad [\text{prior for } \mu]$$

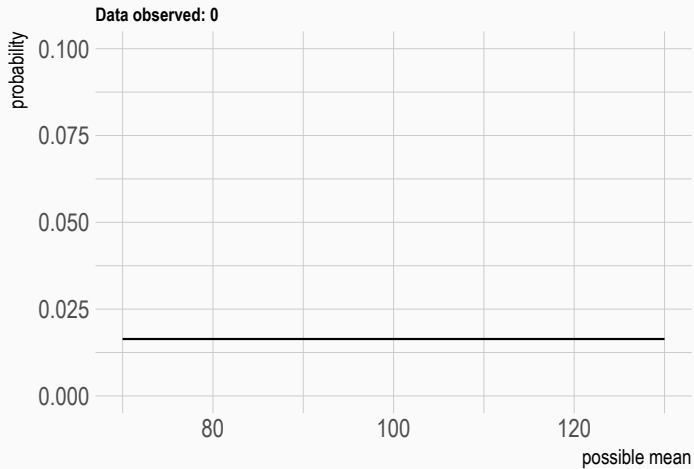
We are going to use **grid approximation** for this example. That means systematically exploring about a bunch of parameter values. (It's mostly useful for illustrating how Bayes' theorem works.)

```
df_iq_model <- data_frame(  
  # Candidate mean value  
  mean = 70:130,  
  # Probability of each candidate mean right now  
  prob = 1 / length(mean),  
  # Probability of each candidate mean during the last update  
  previous = NA_real_)
```

```

# Probabilities sum to 1
sum(df_iq_model$prob)
#> [1] 1
df_iq_model
#> # A tibble: 61 × 3
#>   mean      prob previous
#>   <int>    <dbl>    <dbl>
#> 1     70 0.01639344      NA
#> 2     71 0.01639344      NA
#> 3     72 0.01639344      NA
#> 4     73 0.01639344      NA
#> 5     74 0.01639344      NA
#> 6     75 0.01639344      NA
#> 7     76 0.01639344      NA
#> 8     77 0.01639344      NA
#> 9     78 0.01639344      NA
#> 10    79 0.01639344      NA
#> # ... with 51 more rows

```





We observe one data-point,  $y = 82$ , and update our prior information using the likelihood of the data at each possible mean.

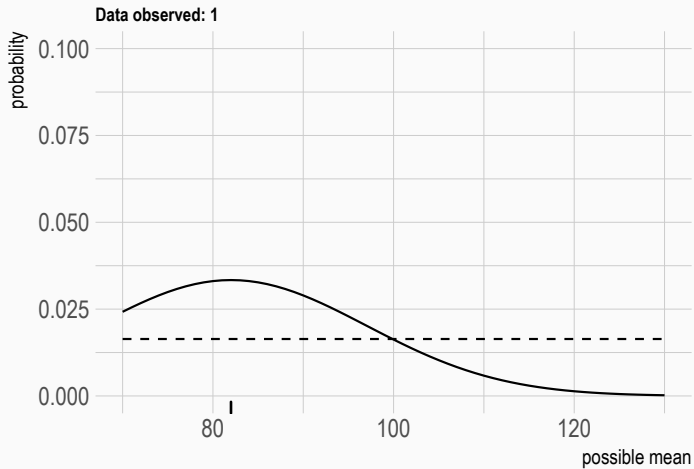
```
df_iq_model$previous <- df_iq_model$prob  
likelihoods <- dnorm(iqs[1], df_iq_model$mean, 15)  
# numerator of bayes theorem  
df_iq_model$prob <- likelihoods * df_iq_model$prob  
sum(df_iq_model$prob)  
#> [1] 0.01306729
```

We observe one data-point,  $y = 82$ , and update our prior information using the likelihood of the data at each possible mean.

```
df_iq_model$previous <- df_iq_model$prob  
likelihoods <- dnorm(iqs[1], df_iq_model$mean, 15)  
# numerator of bayes theorem  
df_iq_model$prob <- likelihoods * df_iq_model$prob  
sum(df_iq_model$prob)  
#> [1] 0.01306729
```

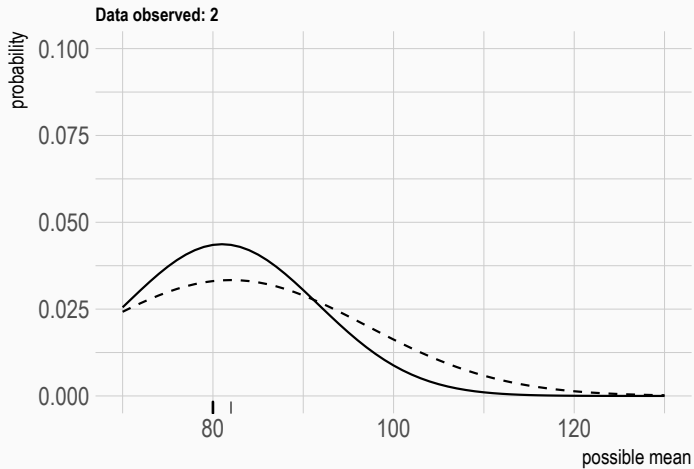
That's not right! We need the *average likelihood* to ensure that the probabilities add up to 1. This is why it's sometimes called a *normalizing constant*.

```
# include denominator of bayes theorem  
df_iq_model$prob <- df_iq_model$prob / sum(df_iq_model$prob)  
sum(df_iq_model$prob)  
#> [1] 1
```



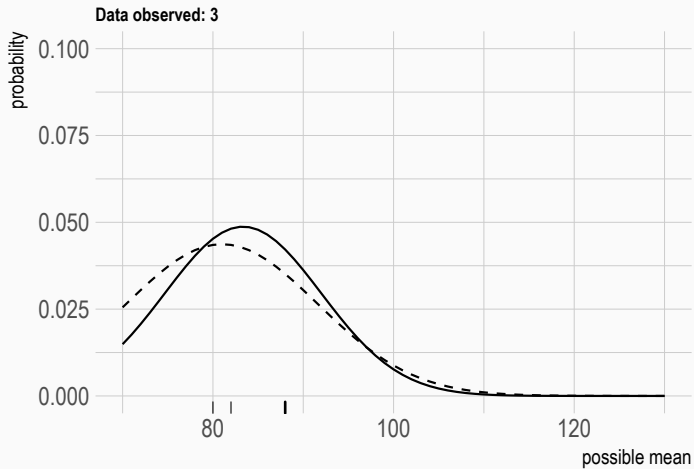
We observe another data-point and update the probability with the likelihood again.

```
df_iq_model$previous <- df_iq_model$prob
likelihoods <- dnorm(iqs[2], df_iq_model$mean, 15)
df_iq_model$prob <- likelihoods * df_iq_model$prob
# normalize
df_iq_model$prob <- df_iq_model$prob / sum(df_iq_model$prob)
df_iq_model
#> # A tibble: 61 × 3
#>   mean      prob previous
#>   <int>    <dbl>    <dbl>
#> 1    70 0.02551519 0.02422865
#> 2    71 0.02801128 0.02549920
#> 3    72 0.03047942 0.02671736
#> 4    73 0.03287154 0.02786958
#> 5    74 0.03513768 0.02894257
#> 6    75 0.03722765 0.02992359
#> 7    76 0.03909289 0.03080065
#> 8    77 0.04068830 0.03156283
#> 9    78 0.04197406 0.03220045
```



And one more...

```
df_iq_model$previous <- df_iq_model$prob
likelihoods <- dnorm(iqs[3], df_iq_model$mean, 15)
df_iq_model$prob <- likelihoods * df_iq_model$prob
# normalize
df_iq_model$prob <- df_iq_model$prob / sum(df_iq_model$prob)
df_iq_model
#> # A tibble: 61 × 3
#>   mean      prob previous
#>   <int>    <dbl>    <dbl>
#> 1     70 0.01490139 0.02551519
#> 2     71 0.01768232 0.02801128
#> 3     72 0.02070434 0.03047942
#> 4     73 0.02392174 0.03287154
#> 5     74 0.02727304 0.03513768
#> 6     75 0.03068201 0.03722765
#> 7     76 0.03405991 0.03909289
#> 8     77 0.03730890 0.04068830
#> 9     78 0.04032654 0.04197406
#> 10    79 0.04301093 0.04291726
#> # with 51 more rows
```



## An animation of these steps

[https://github.com/tjmahr/MadR\\_RStanARM/blob/master/assets/simple-updating.gif](https://github.com/tjmahr/MadR_RStanARM/blob/master/assets/simple-updating.gif)



# Connecting Bayes' theorem to linear regression

---

# Linear models

I learned stats in this course, so I bet you probably write regression models as a one-liner like:

$$\underbrace{y_i}_{\text{observation}} = \underbrace{\alpha + \beta_1 x_{1i}}_{\text{predicted mean given } x} + \underbrace{\epsilon_i}_{\text{random error}}$$

# Linear models

I learned stats in this course, so I bet you probably write regression models as a one-liner like:

$$\underbrace{y_i}_{\text{observation}} = \underbrace{\alpha + \beta_1 x_{1i}}_{\text{predicted mean given } x} + \underbrace{\epsilon_i}_{\text{random error}}$$

Data generating model: Observation  $y_i$  is a draw from a normal distribution centered around a mean.

We estimate the mean with a constant “intercept” term  $\alpha$  plus a linear combination of predictor variables (just  $x_1$  for now).

Let's re-write the model to make the normal-distribution part clearer. No more one-liner.

$$y_i \sim \text{Normal}(\text{mean} = \mu_i, \text{SD} = \sigma) \quad [\text{likelihood}]$$

$$\mu_i = \alpha + \beta_1 * x_{1i} \quad [\text{linear model}]$$

Let's re-write the model to make the normal-distribution part clearer. No more one-liner.

$$y_i \sim \text{Normal}(\text{mean} = \mu_i, \text{SD} = \sigma) \quad [\text{likelihood}]$$

$$\mu_i = \alpha + \beta_1 * x_{1i} \quad [\text{linear model}]$$

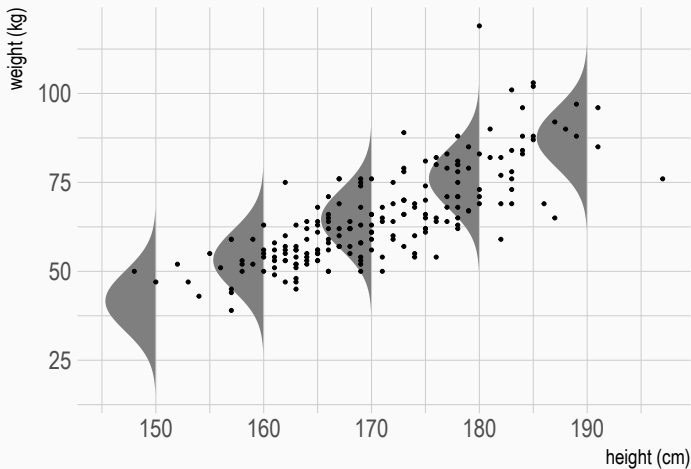
Observation  $y_i$  is a draw from a normal distribution centered around a mean  $\mu_i$  with a standard deviation of  $\sigma$ .

The mean is a constant term  $\alpha$  plus a linear combination of predictor variables (just  $x_1$  for now).

(These equations describe the same models. It's just a different kind of notation.)

## Weight by height model

Consider a model of weight predicted by height. . .



**Figure 8:** It's like a tunnel of bell curves. The center of it moves with  $x$ .



## Bayesian stats

To make the model Bayesian, we need to give prior distributions to parameters.

The parameters we need to estimate for regression:  $\alpha, \beta_1, \sigma$ .

## Bayesian stats

To make the model Bayesian, we need to give prior distributions to parameters.

The parameters we need to estimate for regression:  $\alpha, \beta_1, \sigma$ .

$$y_i \sim \text{Normal}(\mu_i, \sigma) \quad [\text{likelihood}]$$

$$\mu_i = \alpha + \beta_1 * x_{1i} \quad [\text{linear model}]$$

$$\alpha \sim \text{Normal}(0, 10) \quad [\text{prior for } \alpha]$$

$$\beta_1 \sim \text{Normal}(0, 5) \quad [\text{prior for } \beta_1]$$

$$\sigma \sim \text{HalfCauchy}(0, 5) \quad [\text{prior for } \sigma]$$

# What's the point?

- A classical model provides one model of many plausible models of the data. It'll find the parameters that maximize likelihood.
- A Bayesian model is a model of models. We get a *distribution of models* that are consistent with the data.

## But this is where things get difficult!

Parameters we need to estimate:  $\alpha, \beta_1, \sigma$

$$\text{posterior} = \frac{\text{likelihood} * \text{prior}}{\text{average likelihood}}$$

## But this is where things get difficult!

Parameters we need to estimate:  $\alpha, \beta_1, \sigma$

$$\text{posterior} = \frac{\text{likelihood} * \text{prior}}{\text{average likelihood}}$$

$$P(\alpha, \beta, \sigma | x) = \frac{P(x | \alpha, \beta, \sigma) P(\alpha, \beta, \sigma)}{\iiint P(x | \alpha, \beta, \sigma) P(\alpha, \beta, \sigma) d\alpha d\beta d\sigma}$$

Things get gnarly. This is the black-box step.

We don't perform this integral calculus.

Instead, we rely on Markov-chain Monte Carlo simulation to get samples from the posterior.

Those samples will provide a detailed picture of the posterior.

## Finally, let's fit a model

---

## An example: Height and Weight by Sex

```
davis
```

```
#> # A tibble: 199 × 5
```

```
#>       sex weight height repwt repht
```

```
#>   <fctr>   <int>   <int> <int> <int>
```

```
#> 1      M      77     182     77    180
```

```
#> 2      F      58     161     51    159
```

```
#> 3      F      53     161     54    158
```

```
#> 4      M      68     177     70    175
```

```
#> 5      F      59     157     59    155
```

```
#> 6      M      76     170     76    165
```

```
#> 7      M      76     167     77    165
```

```
#> 8      M      69     186     73    180
```

```
#> 9      M      71     178     71    175
```

```
#> 10     M      65     171     64    170
```

```
#> # ... with 189 more rows
```



# Classical linear model

```
# Mean-center height
mean(davis$height)
#> [1] 170.5879

davis$heightC <- davis$height - mean(davis$height)

m <- glm(weight ~ heightC * sex, davis, family = gaussian())
m %>% summary() %>% coef() %>% round(3)

#>               Estimate Std. Error t value Pr(>|t|)
#> (Intercept)    60.558      1.099   55.081   0.000
#> heightC         0.623      0.135    4.626   0.000
#> sexM            7.949      1.710    4.648   0.000
#> heightC:sexM    0.373      0.190    1.964   0.051
```

# What is RStanARM?

Stan: a probabilistic programming language / MCMC sampler

RStanARM: RStan Applied Regression Modeling

- Batteries-included versions of common regression models.
- `glm -> stan_glm`, `glmer -> stan_glmer`.
- CRAN page is very good! They have lots of detailed vignettes!
- Proper successor to the `arm` package.

```
library(rstanarm)
```

```
#> Loading required package: Rcpp
```

```
#> Warning: package 'Rcpp' was built under R version 3.3.3
```

```
#> rstanarm (Version 2.14.1, packaged: 2017-01-16 18:47:11 UTC)
```

```
#> - Do not expect the default priors to remain the same in future rstanarm
```

```
#> Thus, R scripts should specify priors explicitly, even if they are just
```

```
#> - For execution on a local, multicore CPU with excess RAM we recommend
```

```
#> options(mc.cores = parallel::detectCores())
```

- So... hard-code the priors.

## Fit the model

We have to use `stan_glm()`.

- `stan_lm()` uses a different specification of the prior.

By default, it does sampling with 4 MCMC chains. Each “chain” explores the posterior distribution from random starting locations.

- Each chain is 2000 samples, but the first half are warm-up samples.
- Warm-up samples are ignored

```

stan_model <- stan_glm(
  weight ~ heightC * sex,
  data = davis,
  family = gaussian,
  # RStanARM rescales predictor variables and priors use that scaling
  prior = normal(0, 5),
  prior_intercept = normal(0, 10)
)

#>
#> SAMPLING FOR MODEL 'continuous' NOW (CHAIN 1).
#>
#> Chain 1, Iteration:    1 / 2000 [ 0%] (Warmup)
#> Chain 1, Iteration:  200 / 2000 [10%] (Warmup)
#> Chain 1, Iteration:  400 / 2000 [20%] (Warmup)
#> Chain 1, Iteration:  600 / 2000 [30%] (Warmup)
#> Chain 1, Iteration:  800 / 2000 [40%] (Warmup)
#> Chain 1, Iteration: 1000 / 2000 [50%] (Warmup)
#> Chain 1, Iteration: 1001 / 2000 [50%] (Sampling)
#> Chain 1, Iteration: 1200 / 2000 [60%] (Sampling)
#> Chain 1, Iteration: 1400 / 2000 [70%] (Sampling)
#> Chain 1, Iteration: 1600 / 2000 [80%] (Sampling)

```

# Printing the model

```
stan_model
#> stan_glm(formula = weight ~ heightC * sex, family = gaussian,
#>          data = davis, prior = normal(0, 5), prior_intercept = normal(0,
#>          10))
#>
#> Estimates:
#>               Median MAD_SD
#> (Intercept)  60.6      1.1
#> heightC      0.6      0.1
#> sexM         7.9      1.7
#> heightC:sexM 0.4      0.2
#> sigma        8.0      0.4
#>
#> Sample avg. posterior predictive
#> distribution of y (X = xbar):
#>               Median MAD_SD
#> mean_PPD 65.3      0.8
```

## One note

Predictors are centered and rescaled internally by `rstanarm`, so our priors are on the standardized scale.

- `normal(0, 5)` is a distribution of effect sizes with mean 0 and SD 5

See `?rstanarm::priors`, esp. the `autoscale` argument.

```
prior_summary(stan_model)
#> Priors for model 'stan_model'
#> -----
#> Intercept (after predictors centered)
#> ~ normal(location = 0, scale = 10)
#>      **adjusted scale = 266.87
#>
#> Coefficients
#> ~ normal(location = [0,0,0], scale = [5,5,5])
#>      **adjusted scale = [ 7.5,133.4, 11.8]
#>
```

## Getting a summary from the model

```
#> stan_glm(formula = weight ~ heightC * sex, family = gaussian,
#>           data = davis, prior = normal(0, 5), prior_intercept = normal(0,
#>           10))
#>
#> Family: gaussian (identity)
#> Algorithm: sampling
#> Posterior sample size: 4000
#> Observations: 199
#>
#> Estimates:
```

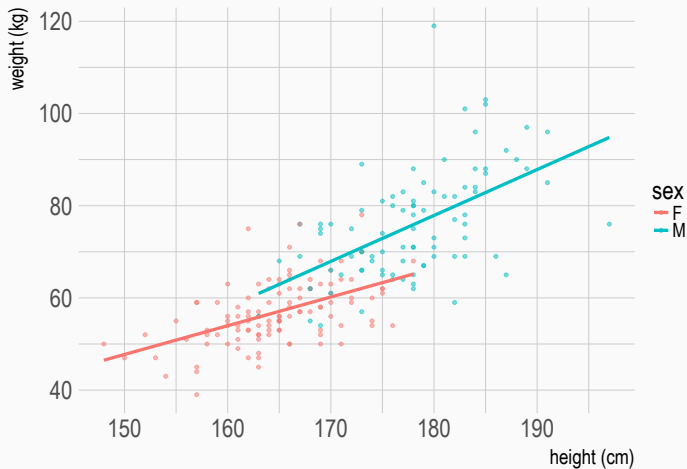
	mean	sd	2.5%	25%	50%	75%
#> (Intercept)	60.6	1.1	58.4	59.8	60.6	61.3
#> heightC	0.6	0.1	0.4	0.5	0.6	0.7
#> sexM	7.9	1.7	4.4	6.8	7.9	9.1
#> heightC:sexM	0.4	0.2	0.0	0.2	0.4	0.5
#> sigma	8.1	0.4	7.3	7.8	8.0	8.3
#> mean_PPD	65.3	0.8	63.7	64.8	65.3	65.9
#> log-posterior	-708.6	1.6	-712.6	-709.4	-708.3	-707.4



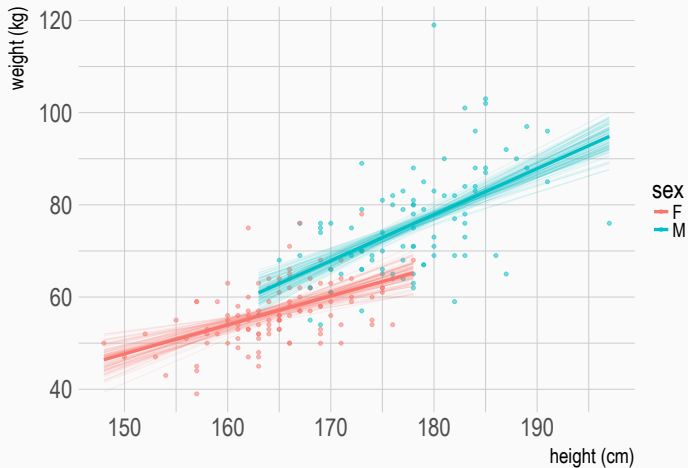
- Split into estimation and diagnostic information
- `mean_PPD` is the predicted value for a completely average observation

# The cut to the chase plot

Here is what classical linear regression does.



# Here is what Bayesian linear regression does



## Inspecting posterior samples

---

## Looking at the posterior parameter samples

Coerce to a data-frame. Columns are parameters. One row per posterior sample.

```
samples <- stan_model %>% as.data.frame() %>% tbl_df()
```

```
samples
```

```
#> # A tibble: 4,000 × 5
```

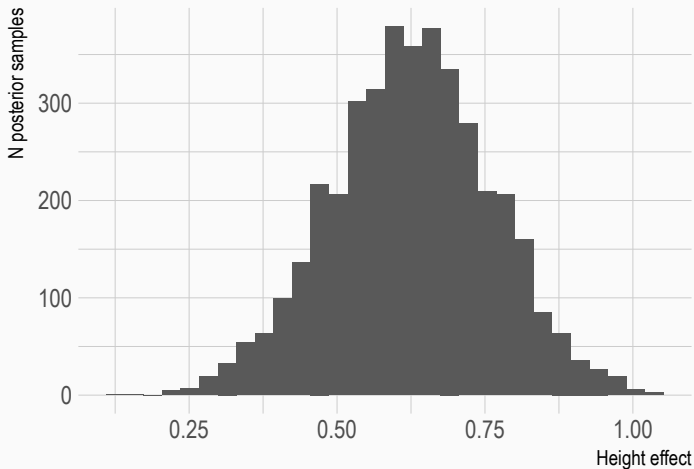
```
#>   `(Intercept)` heightC sexM `heightC:sexM` sigma
#>   <dbl>      <dbl>   <dbl>      <dbl>      <dbl>
#> 1    59.95677 0.5506608 9.517326    0.5270369 8.440430
#> 2    60.02179 0.5898664 8.466223    0.5210752 8.342935
#> 3    59.55219 0.5749104 8.013311    0.5420474 7.766122
#> 4    60.26689 0.6200884 8.822269    0.3888026 8.420006
#> 5    61.14721 0.6210841 7.166503    0.4254251 8.429904
#> 6    60.98669 0.7016688 7.178697    0.4527272 8.795731
#> 7    59.72132 0.5395581 9.325594    0.5446949 7.245603
#> 8    60.49143 0.6660482 9.162610    0.3693488 7.797503
#> 9    62.34006 0.7063126 4.022597    0.5697835 8.107349
#> 10   62.06608 0.6769759 4.352102    0.5321317 8.109990
```

```
#> # ... with 3,990 more rows
```

## We have a distribution

Any stats that can describe a distribution can describe the model's parameters now. Mean, median, skew, quantiles, etc.

## Looking at the posterior parameter samples



**Figure 9:** Histogram of height effect.

## Quantiles are post-data probabilities

If we believe there is a “true” value for a parameter, there is 90% probability that this “true” value is in the 90% interval, given our model, prior information, and the data.

The 90% interval contains the middle 90% of the parameter values.

There is a 5% chance, says the model, the height parameter that generated the data is below the 5% quantile.

```
posterior_interval(stan_model)
#>                                5%          95%
#> (Intercept) 58.72978002 62.4138859
#> heightC      0.40117962 0.8455608
#> sexM         5.01139335 10.7527402
#> heightC:sexM 0.05973589 0.6872728
#> sigma       7.42828134 8.7304518
```



```
# This is where I toured launch_shinystan(stan_model)  
# and did some other stuff.
```

## My experience with this framework

---

The models provide intuitive results.

- When we misinterpret  $p$ -values or confidence intervals, we usually are interpreting them in a Bayesian way.
- Bayesian uncertainty intervals are what we want from confidence intervals.

Bayesian models quantify uncertainty.

- Basically, if a classical model can estimate or predict something about the data, the Bayesian model can estimate a distribution for that thing too.
- Bayesian models are generative, and the posterior predictive distribution (which simulates fake data using the model) is a useful tool.

Bayesian models incorporate prior information.

- That information can be weak, moderate or strong.
- I don't say "prior beliefs" because that sounds too subjective.
- All models make assumptions and prior information, and priors make that information explicit.

Bayesian models are flexible.

- This course captures a bag of tricks ( $t$ -tests, ANOVA, ANCOVA, mixed effects) under a general framework (it's all regression).
- Bayesian regression incorporates even more tricks (missing data imputation, measurement error models, robust error models) into the framework.

Bayesian models have computational benefits.

- Multilevel models with lots of random effects probably won't converge.
- But some weak prior information will nudge the models in the right direction and make the models work.

It's different.

- People are really used to significance testing and  $p$ -values, so you have to do more hand-holding when explaining results.
- You don't get to say *significant* anymore. (I use *plausible* and *credible*.)
- People have misconceptions about subjectivism and bias.



More work before and after modeling.

- You need to specify priors for your parameters.
- Your model is a distribution, so you have to do a bit more work wrangling the data.

It takes longer.

- Classical models solve an optimization problem and provide a single set of parameter estimates.
- MCMC sampling explores the space of parameter values and provides thousands of parameter estimates.
- It can take a few hours to fit a complicated multilevel model.

It's not a cure-all. There are still insecurities.

- It's statistics and people can still misunderstand the methods and models.
- A motivated  $p$ -hacker can still exploit Bayes factors, which is why I won't discuss them.

These are some older slides on good resources for learning about Bayesian statistics.

[https://cdn.rawgit.com/tjmahr/MadR\\_RStanARM/master/04-learning-more-rpubs.html](https://cdn.rawgit.com/tjmahr/MadR_RStanARM/master/04-learning-more-rpubs.html)