

chp5

Vasco Brazão

15/02/2021

```
library(tidyverse)

## -- Attaching packages ----- tidyverse 1.3.0 --

## v ggplot2 3.3.3      v purrr  0.3.4
## v tibble  3.0.6      v dplyr  1.0.4
## v tidyr   1.1.2      v stringr 1.4.0
## v readr   1.4.0      v forcats 0.5.1

## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
```

```
library(here)
```

```
## here() starts at C:/Users/admin/Documents/statistical-rethinking
```

```
library(brms)
```

```
## Loading required package: Rcpp

## Loading 'brms' package (version 2.14.4). Useful instructions
## can be found by typing help('brms'). A more detailed introduction
## to the package is available through vignette('brms_overview').

##
## Attaching package: 'brms'

## The following object is masked from 'package:stats':
##
##      ar
```

5E1

- (4) would be the standard way to write a multiple linear regression. I Suppose (2) could be valid? If we force the intercept to be 0. (3) seems plausible, but from the lack of an index on beta I would think you're forcing the beta for x to be equal to -1 * the beta for z, which is.. strange?

Van Bussel agrees https://github.com/castels/StatisticalRethinking/blob/master/Chapter%205/VanBussel_Chapter5_Questions.pdf

5E2

$\mu_{\text{latitude}_i} = \alpha + \beta_{\text{adiv}} * \text{adiv}_i + \beta_{\text{pdiv}} * \text{pdiv}_i$

5E3

$$time_i \sim \text{Normal}(\mu, \sigma)$$

$$\mu_i = \alpha + \beta_f f_i + \beta_s s_i$$

Both slope parameters should be positive.

Van Bussel agrees!

But I still can't make a stupid latex document. One day.

5E4

1, 3, 4 would be my guesses.

Van Bussel disagrees - 4 is not correct. But I still think it works?

And a latex document thing was created! I cannot believe my eyes. What fresh hell awaits me now? We shall see.. we shall see.

5M1

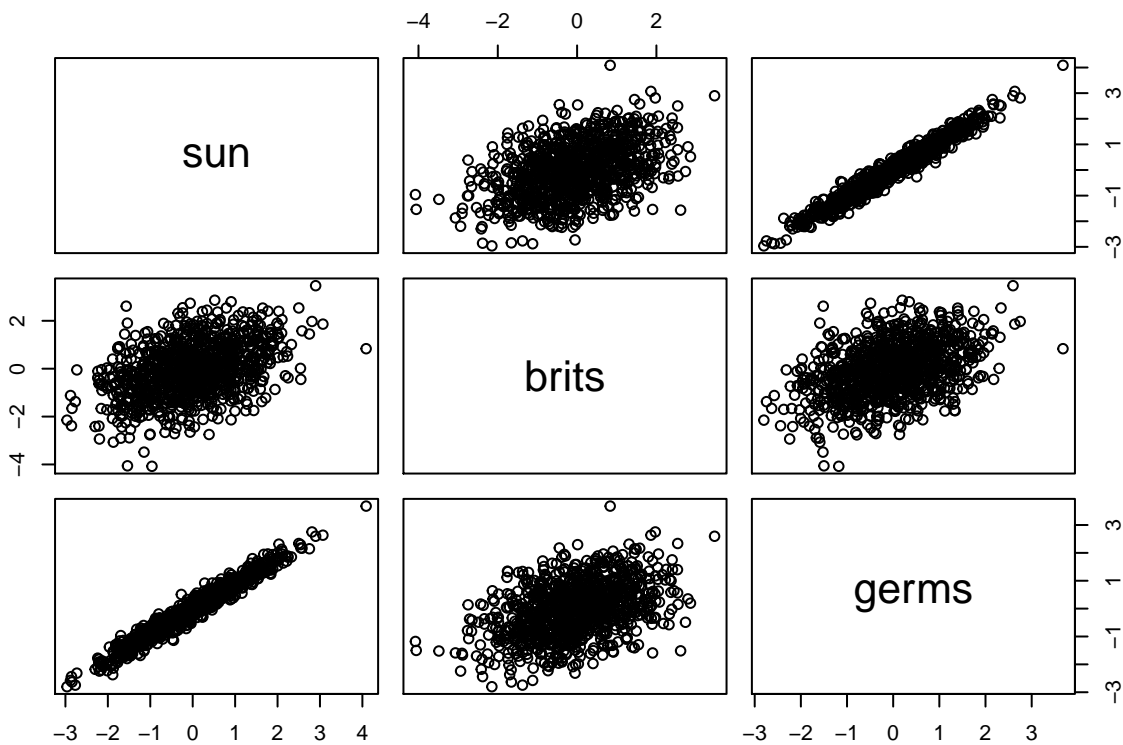
Inventing a spurious correlation.

Let's say that both British and German tourists like going to sunny places. For a given year, they base their decision of whether to go Portugal for vacation on the number of sunlight days of the previous year. Germans respond more strongly to this parameter because they now have an easier time going to Portugal than Brits, and they also stray less from the line because they are German.

```
n <- 1000

df <- tibble(
  sun = rnorm(n, 0, 1),
  brits = 0.5*sun + rnorm(n, 0, 1),
  germs = 0.9*sun + rnorm(n, 0, 0.2)
)
```

```
pairs(df)
```



There appears to be a positive relationship between amount of brits and amount of germans.

```
m1 <- lm(df$brits ~ df$germs)
m2 <- lm(df$brits ~ df$germs + df$sun)
```

```
summary(m1)
```

```
##
## Call:
## lm(formula = df$brits ~ df$germs)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.4552 -0.6859  0.0010  0.6245  3.3975
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -0.04160    0.03165  -1.314   0.189
## df$germs     0.49231    0.03318  14.839 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1 on 998 degrees of freedom
## Multiple R-squared:  0.1808, Adjusted R-squared:  0.1799
## F-statistic: 220.2 on 1 and 998 DF, p-value: < 2.2e-16
```

```
summary(m2)
```

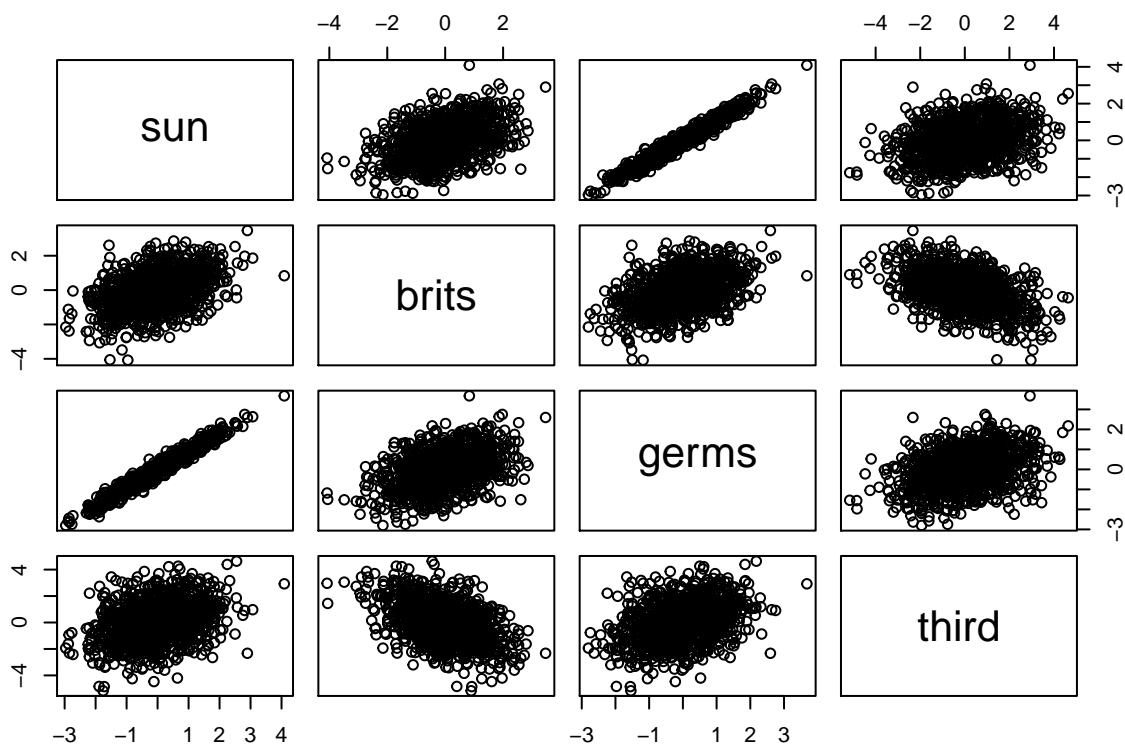
```
##
## Call:
## lm(formula = df$brits ~ df$germs + df$sun)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.5748 -0.6864  0.0192  0.6161  3.3859
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -0.04761    0.03159  -1.507  0.13210
## df$germs      0.03411    0.15585   0.219  0.82678
## df$sun        0.43180    0.14353   3.008  0.00269 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.9963 on 997 degrees of freedom
## Multiple R-squared:  0.1881, Adjusted R-squared:  0.1865
## F-statistic: 115.5 on 2 and 997 DF,  p-value: < 2.2e-16
```

Et voila. when sun is in the model, knowing how many germans went to Portugal does not really tell us much more about how many brits went.

5M2

Now imagine a third country. Those citizens go to Portugal in higher numbers when they know their German friends will be going, but they prefer to avoid Brits on the beaches, for whatever reason.

```
df <- df %>%
  mutate(
    third = rnorm(n, mean = germs - brits)
  )
pairs(df)
```



```
lm(third ~ germs, data = df)
```

```
##
## Call:
## lm(formula = third ~ germs, data = df)
##
## Coefficients:
## (Intercept)      germs
##    0.07003      0.55148
```

```
lm(third ~ brits, data = df)
```

```
##
## Call:
## lm(formula = third ~ brits, data = df)
##
## Coefficients:
## (Intercept)      brits
##    0.01362    -0.64645
```

```
lm(third ~ germs + brits, data = df)
```

```
##
## Call:
```

```
## lm(formula = third ~ germs + brits, data = df)
##
## Coefficients:
## (Intercept)      germs      brits
##      0.02692      1.06163     -1.03624
```

Indeed! When both are included in the model, the coefficients are much greater. When only one is included, its relationship is masked.

5M3

If, in one year, a person marries, divorces, and marries again, they will have added two counts of marriage for that one year – marriage rate is affected by remarriages. Thus, in a society where divorce is accepted but being married is highly desirable, we could reasonably assume a casual link wherein the divorce rate influences the marriage rate. In years with (say) 0 divorces, marriages would mostly be first marriages. In years with a very high number of divorces, we might see that same baseline of new marriages PLUS a surge of remarriages driving up the marriage rate.

You could do a simple linear regression, though that would not inform you about the direction of causality. Alternatively, if you have reason to believe that the time between divorce and remarriage is about a year, you should see an effect of divorce at year 0 on marriage at year 1, but not so much the other way around.

5M4

Found the LDS data on wikipedia.

```
data(WaffleDivorce, package = "rethinking")

d <- WaffleDivorce %>% as_tibble()
lds <- readxl::read_xlsx(path = here("chp5/ldsdata.xlsx"),
                        col_names = TRUE)

d2 <- full_join(d, lds, by = c("Location" = "State")) %>%
  select(Loc, MedianAgeMarriage, Marriage, Divorce, LDS) %>%
  filter(!is.na(Loc)) %>%
  mutate(
    across(.cols = MedianAgeMarriage:LDS,
           .fns = ~ rethinking::standardize(.))
  )
```

Now it's all standardized and in d2.

```
lm(Divorce ~ MedianAgeMarriage + Marriage + LDS, data = d2) %>% summary

##
## Call:
## lm(formula = Divorce ~ MedianAgeMarriage + Marriage + LDS, data = d2)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.05378 -0.47978 -0.05498  0.53054  1.95093
##
```

```
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    5.659e-16  1.084e-01   0.000   1.0000
## MedianAgeMarriage -7.555e-01  1.604e-01  -4.709 2.32e-05 ***
## Marriage         6.338e-03  1.649e-01   0.038  0.9695
## LDS             -3.446e-01  1.297e-01  -2.658  0.0108 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.7667 on 46 degrees of freedom
## Multiple R-squared:  0.4482, Adjusted R-squared:  0.4122
## F-statistic: 12.45 on 3 and 46 DF,  p-value: 4.347e-06
```

```
b5m4 <-
  brm(data = d2,
    family = gaussian,
    Divorce ~ 1 + MedianAgeMarriage + Marriage + LDS,
    prior = c(prior(normal(0, 0.2), class = Intercept),
      prior(normal(0, 0.5), class = b),
      prior(exponential(1), class = sigma)),
    iter = 2000, warmup = 1000, chains = 4, cores = 4,
    seed = 5,
    sample_prior = T,
    file = "b5m4")

print(b5m4)
```

```
## Family: gaussian
## Links: mu = identity; sigma = identity
## Formula: Divorce ~ 1 + MedianAgeMarriage + Marriage + LDS
## Data: d2 (Number of observations: 50)
## Samples: 4 chains, each with iter = 2000; warmup = 1000; thin = 1;
##          total post-warmup samples = 4000
##
## Population-Level Effects:
##              Estimate Est.Error l-95% CI u-95% CI Rhat Bulk_ESS Tail_ESS
## Intercept         -0.00      0.10   -0.20    0.20 1.00     4078     2604
## MedianAgeMarriage  -0.68      0.15   -0.97   -0.38 1.00     3129     2736
## Marriage            0.04      0.16   -0.25    0.35 1.00     2856     2894
## LDS                -0.31      0.13   -0.57   -0.04 1.00     3312     2811
##
## Family Specific Parameters:
##              Estimate Est.Error l-95% CI u-95% CI Rhat Bulk_ESS Tail_ESS
## sigma         0.78      0.08    0.64    0.96 1.00     3418     2103
##
## Samples were drawn using sampling(NUTS). For each parameter, Bulk_ESS
## and Tail_ESS are effective sample size measures, and Rhat is the potential
## scale reduction factor on split chains (at convergence, Rhat = 1).
```

```
summary(b5m4)
```

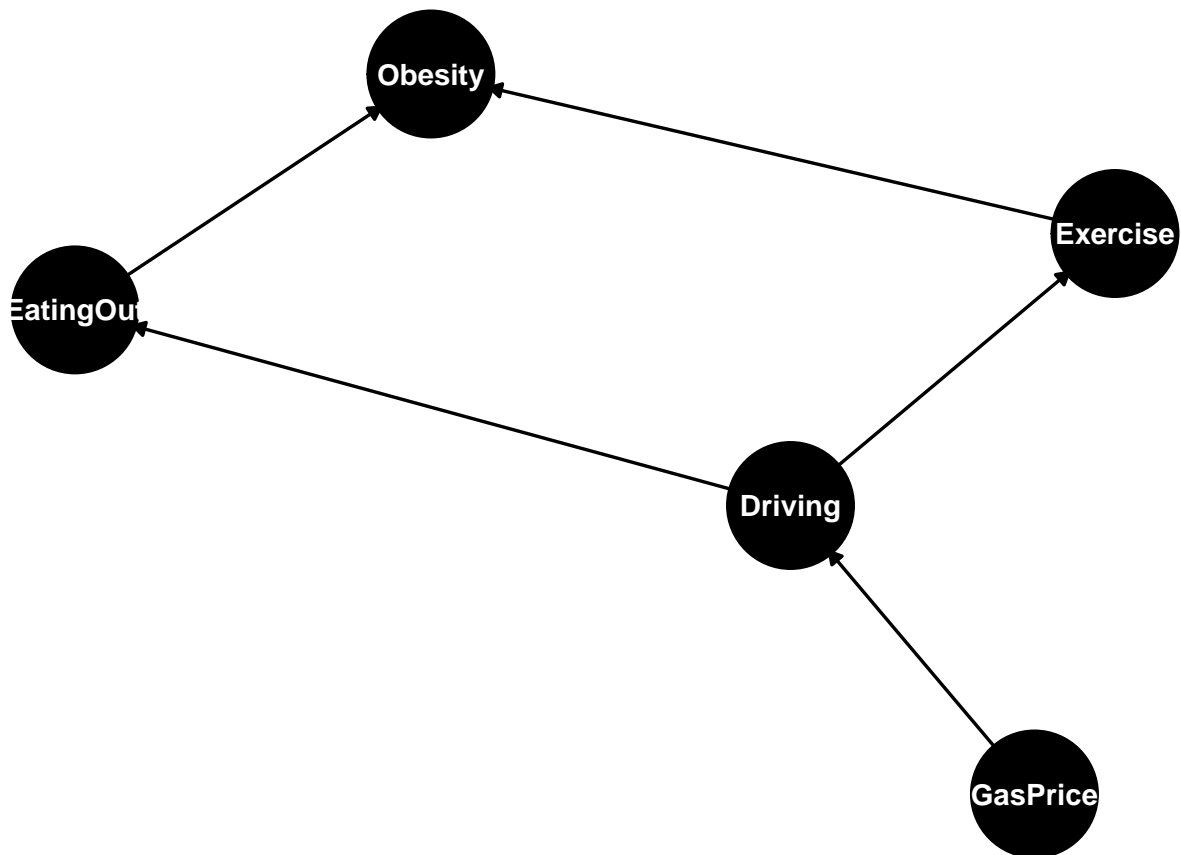
```
## Family: gaussian
## Links: mu = identity; sigma = identity
```

```
## Formula: Divorce ~ 1 + MedianAgeMarriage + Marriage + LDS
## Data: d2 (Number of observations: 50)
## Samples: 4 chains, each with iter = 2000; warmup = 1000; thin = 1;
## total post-warmup samples = 4000
##
## Population-Level Effects:
##      Estimate Est.Error l-95% CI u-95% CI Rhat Bulk_ESS Tail_ESS
## Intercept      -0.00      0.10   -0.20    0.20 1.00     4078     2604
## MedianAgeMarriage -0.68      0.15   -0.97   -0.38 1.00     3129     2736
## Marriage          0.04      0.16   -0.25    0.35 1.00     2856     2894
## LDS              -0.31      0.13   -0.57   -0.04 1.00     3312     2811
##
## Family Specific Parameters:
##      Estimate Est.Error l-95% CI u-95% CI Rhat Bulk_ESS Tail_ESS
## sigma      0.78      0.08      0.64      0.96 1.00     3418     2103
##
## Samples were drawn using sampling(NUTS). For each parameter, Bulk_ESS
## and Tail_ESS are effective sample size measures, and Rhat is the potential
## scale reduction factor on split chains (at convergence, Rhat = 1).
```

5M5

A DAG might be useful here.

```
ggdag::dagify(
  Obesity ~ Exercise + EatingOut,
  Exercise ~ Driving,
  EatingOut ~ Driving,
  Driving ~ GasPrice
) %>%
  ggdag(node_size = 22) +
  theme_void()
```

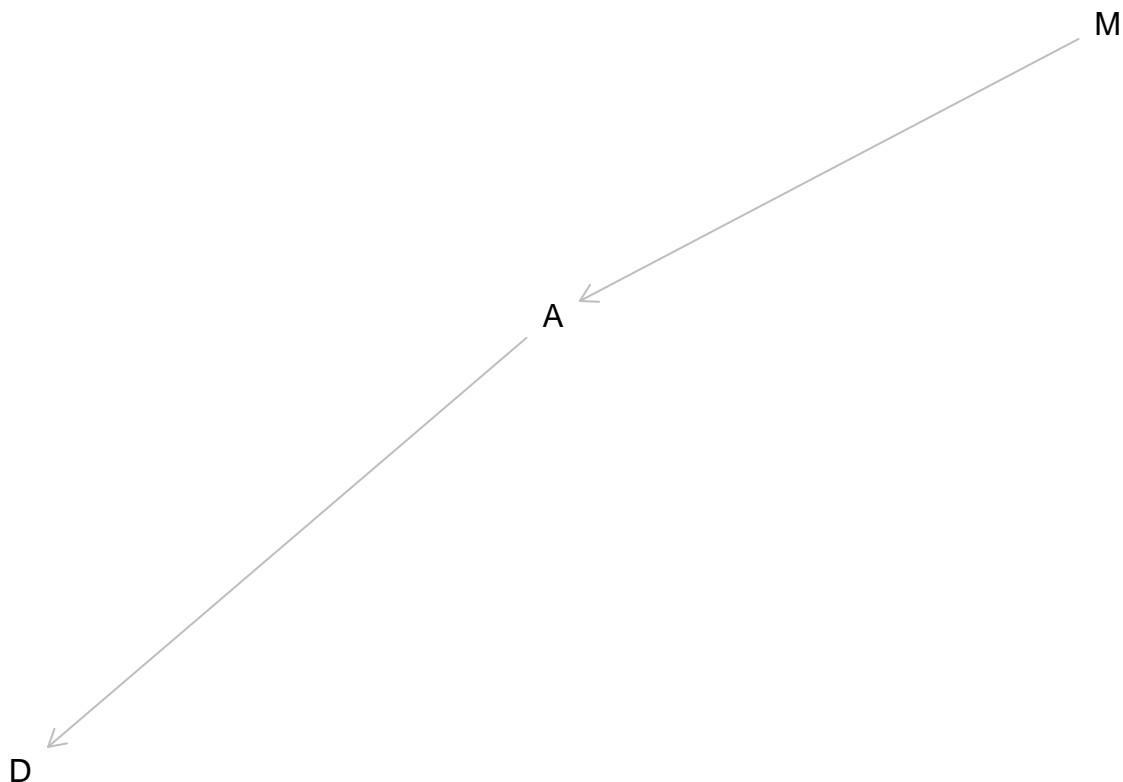
How ugly! And it changes every time I run it. I embrace the chaos.

Ok, so...

- We could regress obesity on each of the predictors separately, just to confirm our intuitions.
- We could regress obesity on gas price and driving. We would expect that the coefficient for GasPrice would shrink, since Driving would be holding most of the information
- We could regress obesity on EatingOut and Exercise. The coefficients should shrink because they are collinear.
- We could regress obesity on EatingOut, Exercise, and Driving. The coefficient for Driving should shrink?

5H1

```
dag <- dagitty::dagitty(  
  "dag{M -> A; A -> D}",  
  layout = TRUE  
)  
  
plot(dag)
```



```
dagitty::impliedConditionalIndependencies(dag)
```

```
## D _||_ M | A
```

Without conditioning, all pairs of variables should be associated.

D is independent of M, conditioned on A. (A here is the mediator)

This is consistent with the pattern of results from the three models (m5.1, m5.2, and m5.3).

5H2

Taking the previous DAG, I am to fit a model and use it to estimate the effect of halving a State's marriage rate (M).

Side-note: I recently read this blogpost: <https://elevanth.org/blog/2018/07/14/statistical-rethinking-edition-2-eta-2020/>

In which Richard says:

“First, I force the reader to explicitly specify every assumption of the model. Some readers of the first edition lobbied me to use simplified formula tools like brms or rstanarm. Those are fantastic packages, and graduating to use them after this book is recommended. But I don't see how a person can come to understand the model when using those tools. The priors being hidden isn't the most limiting part. Instead, since linear model formulas like $y \sim (1|x) + z$ don't show the parameters, nor even all of the terms, it is not easy to see how the mathematical model relates

to the code. It is ultimately kinder to be a bit cruel and require more work. So the formula lists remain. In this book, you are programming the log-posterior, down to the exact relationship between each variable and coefficient. You'll thank me later."

So I just decided that I will use the rethinking package to fit the models after all, and then do them in brms as well. This is becoming a giant mess. I embrace the chaos.

```
# we will want these for later: mean and sd of MedianAgeMarriage
M_mean <- mean(d$Marriage)
M_sd <- sd(d$Marriage)

d2 <- d2 %>%
  mutate(
    M = Marriage,
    A = MedianAgeMarriage,
    D = Divorce
  )
```

Now the models we need to estimate would be $M \rightarrow A$, and then $A \rightarrow D$, right? No need for a model that has M and A?

wjakethompson disagrees (see <https://github.com/wjakethompson/sr2-solutions/blob/main/03-more-linear-models.Rmd>). I will include both then, but I'm not very confident.

```
m5h2 <- rethinking::quap(
  alist(
    ## M -> A -> D
    D ~ dnorm(mu, sigma),
    mu <- a + bA*A + bM*M,
    a ~ dnorm(0, 0.2),
    bA ~ dnorm(0, 0.5),
    bM ~ dnorm(0, 0.5),
    sigma ~ dexp(1),

    ## M -> A
    A ~ dnorm(mu_A, sigma_A),
    mu_A <- aA + bAA*M,
    aA ~ dnorm(0, 0.2),
    bAA ~ dnorm(0, 0.5),
    sigma_A ~ dexp(1)
  ), data = d2
)

rethinking::precis(m5h2)
```

| ## | | mean | sd | 5.5% | 94.5% |
|------------|--|---------------|------------|------------|------------|
| ## a | | -2.469709e-06 | 0.09707479 | -0.1551467 | 0.1551418 |
| ## bA | | -6.135226e-01 | 0.15098129 | -0.8548199 | -0.3722254 |
| ## bM | | -6.539464e-02 | 0.15077074 | -0.3063554 | 0.1755661 |
| ## sigma | | 7.851049e-01 | 0.07784015 | 0.6607013 | 0.9095085 |
| ## aA | | -5.179614e-06 | 0.08684849 | -0.1388058 | 0.1387955 |
| ## bAA | | -6.947250e-01 | 0.09572792 | -0.8477167 | -0.5417333 |
| ## sigma_A | | 6.817432e-01 | 0.06758164 | 0.5737347 | 0.7897518 |

```

M_seq <- seq(from = -3, to = 3, length.out = 50)

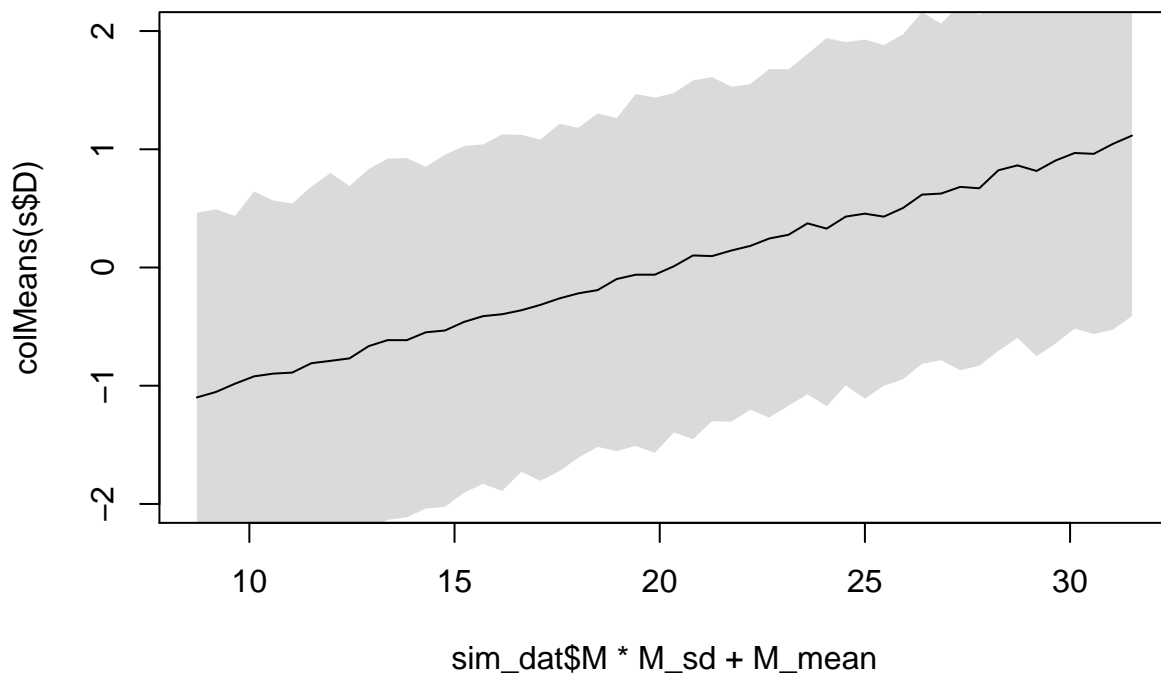
# now to simulate and plot the counterfactual

sim_dat <- data.frame(M = M_seq)

s <- rethinking::sim(m5h2,
  data = sim_dat,
  vars = c("A", "D"))

plot(sim_dat$M * M_sd + M_mean, colMeans(s$D), ylim = c(-2, 2), type = "l")
rethinking::shade(apply(s$D, 2, rethinking::PI), sim_dat$M * M_sd + M_mean)

```



The effect of halving the marriage rate depends on the rate itself. From $M = 30$ to $M = 15$, Divorce falls about 1SD; from $M = 20$ to $M = 10$, the fall is a bit smaller than 1SD.

But we can do a bit better: if

$$D = \alpha + \beta M$$

Then

$$D_{\frac{1}{2}M} = \alpha + \beta * \left(\frac{1}{2}M\right)$$

And we can use these to say exactly how we predict D to change when we halve M . The *difference* between the new value and the original value will be

$$D_{\frac{1}{2}M} - D = \alpha + \beta * (\frac{1}{2}M) - (\alpha + \beta M)$$

Which we simplify to

$$D_{\frac{1}{2}M} - D = -\frac{1}{2}\beta M$$

```
unscaled <- NULL
unscaled$M <- sim_dat$M * M_sd + M_mean

lm(colMeans(s$D) ~ (unscaled$M) )

##
## Call:
## lm(formula = colMeans(s$D) ~ (unscaled$M))
##
## Coefficients:
## (Intercept)      unscaled$M
##      -1.91946         0.09508
```

From this simple linear model we estimate $\beta = 0.095$. As a sanity check, we calculate that the difference in Divorce rates when we go from $M = 30$ to $M = 15$ should equal

```
-0.5 * 0.095 * 30
```

```
## [1] -1.425
```

This is a bigger jump than I predicted. Maybe it's hard to see from the graph?

Next, doing it in brms and the tidyverse.

```
# first we specify each model separately
d_model <- bf(D ~ 1 + A + M)
a_model <- bf(A ~ 1 + M)

b5h2 <-
  brm(data = d2,
      family = gaussian,
      d_model + a_model + set_rescor(FALSE),
      prior = c(prior(normal(0, 0.2), class = Intercept, resp = D),
                prior(normal(0, 0.5), class = b, resp = D),
                prior(exponential(1), class = sigma, resp = D),

                prior(normal(0, 0.2), class = Intercept, resp = A),
                prior(normal(0, 0.5), class = b, resp = A),
                prior(exponential(1), class = sigma, resp = A)),
      iter = 2000, warmup = 1000, chains = 4, cores = 4,
      seed = 5,
      file = "fits/b5h2")

print(b5h2)
```

```
## Family: MV(gaussian, gaussian)
## Links: mu = identity; sigma = identity
##      mu = identity; sigma = identity
## Formula: D ~ 1 + A + M
##      A ~ 1 + M
## Data: d2 (Number of observations: 50)
## Samples: 4 chains, each with iter = 2000; warmup = 1000; thin = 1;
##      total post-warmup samples = 4000
##
## Population-Level Effects:
##      Estimate Est.Error l-95% CI u-95% CI Rhat Bulk_ESS Tail_ESS
## D_Intercept    -0.00     0.10   -0.20    0.19 1.00    5531    3016
## A_Intercept     -0.00     0.09   -0.18    0.18 1.00    5709    2875
## D_A             -0.61     0.16   -0.92   -0.28 1.00    3595    3196
## D_M             -0.06     0.16   -0.36    0.26 1.00    3446    2991
## A_M             -0.69     0.10   -0.88   -0.50 1.00    5626    2546
##
## Family Specific Parameters:
##      Estimate Est.Error l-95% CI u-95% CI Rhat Bulk_ESS Tail_ESS
## sigma_D       0.83     0.09    0.67    1.01 1.00    4704    2532
## sigma_A       0.71     0.07    0.58    0.87 1.00    5202    2737
##
## Samples were drawn using sampling(NUTS). For each parameter, Bulk_ESS
## and Tail_ESS are effective sample size measures, and Rhat is the potential
## scale reduction factor on split chains (at convergence, Rhat = 1).
```

```
nd <- tibble(
  M = seq(from = -3, to = 3, length.out = 50)
)

# now to simulate and plot the counterfactual

nd_2 <- predict(b5h2,
  resp = "A",
  newdata = nd
) %>%
  data.frame() %>%
  bind_cols(nd) %>%
  rename(
    A = Estimate
  ) %>%
  select(A, M)

# If I don't first predict A based on M and only then predict D based on A and
# M, the relationship is negative. (If I set A = 0 like in the brms example here
# ): https://bookdown.org/content/4857/the-many-variables-the-spurious-waffles.html#counterfactual-plot.
predict(b5h2,
  resp = "D",
  newdata = nd_2
) %>%
  data.frame() %>%
  bind_cols(nd_2) %>%
  mutate(
    M = M * M_sd + M_mean
```

```
) %>%  
ggplot(aes(x = M, y = Estimate, ymin = Q2.5, ymax = Q97.5)) +  
geom_smooth(stat = "identity")
```

