

## 7. Chapter 7 Solutions

**7E1.** The criteria for this measure of “uncertainty” are:

- (1) Continuity. The measure is not discrete, but it may be bounded. So it can have a minimum and maximum, but it must be continuous in between.
- (2) Increasing with the number of events. When more distinct things can happen, and all else is equal, there is more uncertainty than when fewer things can happen.
- (3) Additivity. This is hardest one to understand, for most people. Additivity is desirable because it means we can redefine event categories without changing the amount of uncertainty.

**7E2.** This is a simple calculation in R, very much like the example on page 192:

```
# define probabilities of heads and tails
p <- c( 0.7 , 0.3 )
```

R code  
7.1

```
# compute entropy
-sum( p*log(p) )
```

```
[1] 0.6108643
```

**7E3.** Similar to the problem above, but now with four types of events:

```
# define probabilities of sides
p <- c( 0.2 , 0.25 , 0.25 , 0.3 )
```

R code  
7.2

```
# compute entropy
-sum( p*log(p) )
```

```
[1] 1.376227
```

**7E4.** When events are impossible (have probability of zero), they just fall out of the calculation:

```
# define probabilities of sides
p <- c( 1/3 , 1/3 , 1/3 )
```

R code  
7.3

```
# compute entropy
-sum( p*log(p) )
```

```
[1] 1.098612
```

**7M1.** The definition of AIC is:

$$\begin{aligned} \text{AIC} &= -2 \log \Pr(\text{data} | \text{MAP estimates}) + 2p \\ &= D_{\text{train}} + 2p \end{aligned}$$

where  $p$  is the number of parameters in the model. WAIC is a notational mess, but once you understand each part, then the notation makes sense:

$$\begin{aligned}\text{WAIC} &= -2(\text{lppd} - p_{\text{WAIC}}) \\ &= -2 \left( \sum_i \log \Pr(y_i) - \sum_i V(y_i) \right)\end{aligned}$$

where  $\Pr(y_i)$  is the likelihood of observation  $i$ , averaging over the posterior distribution, and  $V(y_i)$  is the variance in the log-likelihood of observation  $i$ , taking the variance over the posterior distribution.

Comparing these definitions, each has a term that expresses the fit to the training sample, as well as a term that expresses some penalty for flexibility in fitting to the sample. For AIC, the fit term is simply the “plug-in deviance,”  $D_{\text{train}}$ , and the penalty term is twice the number of parameters,  $p$ . For WAIC, the fit term is a fully Bayesian log-likelihood that averages over the posterior prediction for each observation separately. This is lppd. The penalty term is the sum of the variances of each log-likelihood.

From more general to less general: WAIC, AIC. When the posterior predictive mean is a good representation of the posterior predictive distribution, and the priors are effectively flat or overwhelmed by the amount of data, then WAIC and AIC will tend to agree.

**7M2.** Model selection means ranking models by information criteria (or any other criteria) and choosing the highest ranked model. Model averaging is instead weighting each model by its relative distance from the best model and creating a posterior predictive distribution comprising predictions from each model in proportion to these weights, a prediction “ensemble.” Model selection discards information about model uncertainty. Model averaging retains some of that information, but it still discards some information about model uncertainty, because it compresses the full information about the set of models into a single predictive distribution. Since the full set of ranks and information criteria values cannot be recovered from the averaged predictive distribution, some of that information has been lost. In other words, different model comparison sets with different ranks and information criteria values can produce nearly identical prediction ensembles. But those different model comparison sets may be different in other important ways. For example, inspecting the full model comparison set, with the structure of each model, often reveals why some models fit better than others.

**7M3.** When one model is fit to fewer (or different) observations, it is being judged on a different target than the other models. If fewer observations are used, the model will usually appear to perform better, because the deviance will be smaller due to having less to predict. Less to predict means less prediction error, and deviance (as well as information criteria) is a kind of accumulated error. We have to be careful about this, because most of R’s black-box regression functions like `lm`, `glm`, and `glmnet` will automatically and silently drop incomplete cases, reducing the number of observations the model is fit to. This is bad behavior for scientific software, but unfortunately it is the norm.

**7M4.** As a prior becomes more concentrated around particular parameter values, the model becomes less flexible in fitting the sample. One way to remember this is to think of the prior as representing previous learning from previous observations. So a more concentrated, or peaked, prior represents more previous data. As the model becomes less flexible, the effective number of parameters declines.

To perform some simple experiments to demonstrate this, try out this code, changing the value for `sigma` in the data list. The smaller `sigma`, the more concentrated the prior and the smaller the effective number of parameters should be.

```

y <- rnorm(10) # execute just once, to get data

# repeat this, changing sigma each time
m <- quap(
  alist(
    y ~ dnorm(mu,1),
    mu ~ dnorm(0,sigma)
  ), data=list(y=y,sigma=10) )
WAIC(m)

```

R code  
7.4

**7M5.** Informative priors reduce overfitting by reducing the sensitivity of a model to a sample. Some of the information in a sample is *irregular*, not a recurring feature of the process of interest.

**7M6.** If a prior is overly informative, then even regular features of a sample will not be learned by a model. In this case, the model may underfit the sample. In case this sounds like a terrifying balancing act, in practice there is usually a broad family of priors which achieve practically indistinguishable estimates from the same sample. The last “hard” practice problem for this chapter provides an example.

**7H1/7H2.** It’s probably obvious from the prompt that the original curve was hand-drawn and not the result of any fitting procedure. And you can see there is a high outlier point, while the rest of the points show a general increase. But let’s take this seriously and polish our modeling skills.

Lots of combinations of models will produce the same general inference. I’ll start with comparing linear to polynomial models. Then I’ll try a spline. Then I’ll consider a robust regression to cope better with the outlier.

Here is a linear fit, together with a quadratic and cubic fit:

```

library(rethinking)
data(Laffer)
d <- Laffer
d$T <- standardize( d$tax_rate )
d$R <- standardize( d$tax_revenue )

# linear model
m7H1a <- quap(
  alist(
    R ~ dnorm( mu , sigma ),
    mu <- a + b*T,
    a ~ dnorm( 0 , 0.2 ),
    b ~ dnorm( 0 , 0.5 ),
    sigma ~ dexp(1)
  ) , data=d )

# quadratic model
m7H1b <- quap(
  alist(
    R ~ dnorm( mu , sigma ),

```

R code  
7.5

```

      mu <- a + b*T + b2*T^2,
      a ~ dnorm( 0 , 0.2 ),
      c(b,b2) ~ dnorm( 0 , 0.5 ),
      sigma ~ dexp(1)
    ) , data=d )

# cubic model
m7H1c <- quap(
  alist(
    R ~ dnorm( mu , sigma ),
    mu <- a + b*T + b2*T^2 + b3*T^3,
    a ~ dnorm( 0 , 0.2 ),
    c(b,b2,b3) ~ dnorm( 0 , 0.5 ),
    sigma ~ dexp(1)
  ) , data=d )

```

As a first check, let's see how these models compare on purely predictive criteria, like PSIS:

R code 7.6

```
compare( m7H1a , m7H1b , m7H1c , func=PSIS )
```

Some Pareto  $k$  values are very high ( $>1$ )

	PSIS	SE	dPSIS	dSE	pPSIS	weight
m7H1a	93.3	26.88	0.0	NA	8.3	0.78
m7H1b	96.5	32.37	3.2	5.75	10.9	0.16
m7H1c	98.3	31.82	5.0	5.25	11.8	0.06

Not much support for the wiggly models over the linear one. But PSIS also warns about high leverage points. You can probably guess from scatterplot which point is causing the problem. If you peek at the pointwise  $k$  values, you'll see:

R code 7.7

```
PSISk(m7H1a)
```

```

[1] 0.61 0.46 0.44 0.07 -0.03 0.09 0.28 0.29 -0.07 0.04 0.43 1.85
[13] 0.34 0.34 -0.11 -0.07 -0.13 -0.08 0.31 -0.01 -0.02 -0.02 0.05 0.04
[25] 0.07 0.26 0.28 0.23 0.10

```

Point 12 is much over 1. That's the nation with the very high tax revenue value. I think it's Norway, actually, and it results from an accounting trick involving oil revenue.

It'll be useful to plot the posterior predictions of each model.

R code 7.8

```

T_seq <- seq( from=-3.2 , to=1.2 , length.out=30 )
la <- link( m7H1a , data=list(T=T_seq) )
lb <- link( m7H1b , data=list(T=T_seq) )
lc <- link( m7H1c , data=list(T=T_seq) )

```

```

plot( d$T , d$R , xlab="tax rate" , ylab="revenue" )
mtext( "linear model" )
lines( T_seq , colMeans(la) )
shade( apply( la , 2 , PI ) , T_seq )

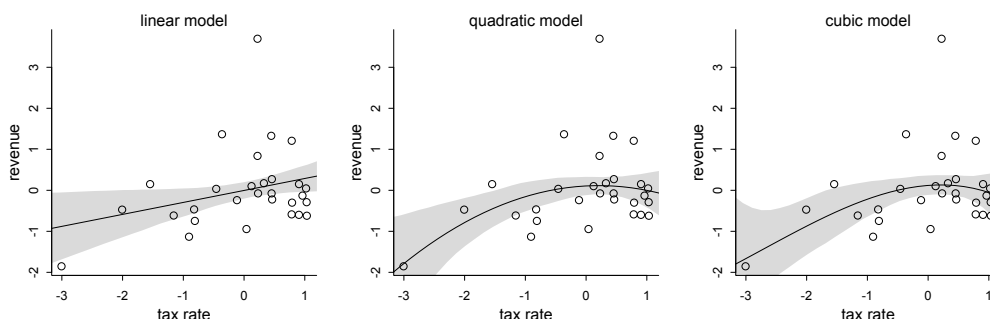
```

```

plot( d$T , d$R , xlab="tax rate" , ylab="revenue" )
mtext( "quadratic model" )
lines( T_seq , colMeans(lb) )
shade( apply( lb , 2 , PI ) , T_seq )

```

```
plot( d$T , d$R , xlab="tax rate" , ylab="revenue" )
mtext( "cubic model" )
lines( T_seq , colMeans(lc) )
shade( apply( lc , 2 , PI ) , T_seq )
```



The polynomial models do bend the right way, but they aren't bending much. If anything, they offer stronger evidence that more tax produces more revenue, not that very high tax reduces revenue.

For the sake of the exercise, let's consider a basis spline. Refer back to Chapter 4, the last section, if you've forgotten splines.

```
num_knots <- 15
knot_list <- quantile( d$T , probs=seq(0,1,length.out=num_knots) )

library(splines)
B <- bs(d$T,
  knots=knot_list[-c(1,num_knots)] ,
  degree=3 , intercept=TRUE )

m7H1s <- quap(
  alist(
    R ~ dnorm( mu , sigma ) ,
    mu <- a + B %*% w ,
    a ~ dnorm(0,1),
    w ~ dnorm(0,1),
    sigma ~ dexp(1)
  ), data=list( R=d$R , B=B ) ,
  start=list( w=rep( 0 , ncol(B) ) ) )
```

R code  
7.9

The coefficients aren't interpretable, but let's compare to the previous models:

```
compare( m7H1a , m7H1b , m7H1c , m7H1s , func=PSIS )
```

R code  
7.10

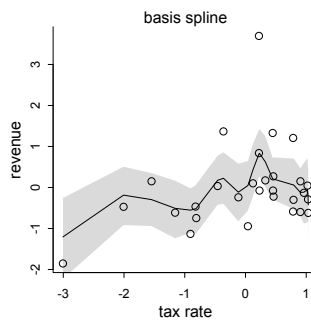
```
Some Pareto k values are very high (>1)
      PSIS    SE dPSIS  dSE pPSIS weight
m7H1a  93.1 26.81   0.0  NA   8.2   0.60
m7H1c  95.1 28.45   2.1 2.57  10.2   0.21
m7H1b  95.4 30.30   2.4 3.74  10.3   0.18
m7H1s 102.4 31.75   9.3 6.06  17.7   0.01
```

Terrible. What do the predictions look like?

```

R code
7.11 mu <- link( m7H1s )
      mu_PI <- apply( mu , 2 , PI )
      plot( d$T , d$R , xlab="tax rate" , ylab="revenue" )
      mtext( "basis spline" )
      o <- order( d$T )
      lines( d$T[o] , colMeans(mu)[o] )
      shade( mu_PI[,o] , d$T[o] )

```



If you really want to interpret that as support for the curved relationship, good luck. You might try more rigid splines, with lower polynomial order, fewer knots, and tighter weights, to see if this can be made more sensible. But clearly this is a curve-fitting exercise, not science, at this point.

For the Student-t model, we just need to replace the normal likelihood. I'll use  $\nu = 2$ , so the tails are satisfyingly thick.

```

R code
7.12 # linear model with student-t
      m7H1d <- quap(
        alist(
          R ~ dstudent( 2 , mu , sigma ),
          mu <- a + b*T,
          a ~ dnorm( 0 , 0.2 ),
          b ~ dnorm( 0 , 0.5 ),
          sigma ~ dexp(1)
        ) , data=d )

```

Comparing the previous models:

```

R code
7.13 compare( m7H1a , m7H1b , m7H1c , m7H1s , m7H1d , func=PSIS )

```

```

Some Pareto k values are very high (>1)
      PSIS    SE dPSIS  dSE pPSIS weight
m7H1d  75.0 13.70   0.0   NA   4.0     1
m7H1b  89.1 25.23  14.0 16.78   7.1     0
m7H1a  98.8 32.32  23.8 23.44  11.0     0
m7H1c 103.2 36.99  28.1 27.91  14.2     0
m7H1s 105.5 33.20  30.5 23.88  19.3     0

```

Oh PSIS likes this one. It also doesn't itself throw any Pareto- $k$  warnings. If you like, try to polynomial models with the Student-t likelihood. Overall, there is no obvious support for a strongly increasing-then-decreasing relationship between tax rate and tax revenue.

**7H3.** To compute the entropies, we just need a function to compute the entropy. Information entropy, as defined in lecture and the book, is simply:

$$H(p) = - \sum_i p_i \log(p_i)$$

where  $p$  is a vector of probabilities summing to 1. In R code this would look like:

```
H <- function(p) -sum(p*log(p))
```

R code  
7.14

I'll make a list of the birb distributions and then push each through the function above.

```
IB <- list()
IB[[1]] <- c( 0.2 , 0.2 , 0.2 , 0.2 , 0.2 )
IB[[2]] <- c( 0.8 , 0.1 , 0.05 , 0.025 , 0.025 )
IB[[3]] <- c( 0.05 , 0.15 , 0.7 , 0.05 , 0.05 )
sapply( IB , H )
```

R code  
7.15

```
[1] 1.6094379 0.7430039 0.9836003
```

The first island has the largest entropy, followed by the third, and then the second in last place. Why is this? Entropy is a measure of the evenness of a distribution. The first islands has the most even distribution of birbs. This means you wouldn't be very surprised by any particular birb. The second island, in contrast, has a very uneven distribution of birbs. If you saw any birb other than the first species, it would be surprising.

Now we need K-L distance, so let's write a function for it:

```
DKL <- function(p,q) sum( p*(log(p)-log(q)) )
```

R code  
7.16

This is the distance from  $q$  to  $p$ , regarding  $p$  as true and  $q$  as the model. Now to use each island as a model of the others, we need to consider the different ordered pairings. I'll just make a matrix and loop over rows and columns:

```
Dm <- matrix( NA , nrow=3 , ncol=3 )
for ( i in 1:3 ) for ( j in 1:3 ) Dm[i,j] <- DKL( IB[[j]] , IB[[i]] )
round( Dm , 2 )
```

R code  
7.17

```
      [,1] [,2] [,3]
[1,] 0.00 0.87 0.63
[2,] 0.97 0.00 1.84
[3,] 0.64 2.01 0.00
```

The way to read this is each row as a model and each column as a true distribution. So the first island, the first row, has the smaller distances to the other islands. This makes sense, since it has the highest entropy. Why does that give it a shorter distance to the other islands? Because it is less surprised by the other islands, due to its high entropy.

**7H4.** I won't repeat the models here. They are in the text. Model `m6.9` contains both marriage status and age. Model `m6.10` contains only age. Model `m6.9` produces a confounded inference about the relationship between age and happiness, due to opening a collider path. To compare these models using WAIC:

R code  
7.18

```
compare( m6.9 , m6.10 )
```

	WAIC	pWAIC	dWAIC	weight	SE	dSE
m6.9	2714.0	3.7	0.0	1	37.54	NA
m6.10	3101.9	2.3	387.9	0	27.74	35.4

The model that produces the invalid inference, m6.9, is expected to predict much better. And it would. This is because the collider path does convey actual association. We simply end up mistaken about the causal inference. We should not use WAIC (or LOO) to choose among models, unless we have some clear sense of the causal model. These criteria will happily favor confounded models.

**7H5.** These are the models:

R code  
7.19

```
library(rethinking)
data(foxes)
d <- foxes
d$W <- standardize(d$weight)
d$A <- standardize(d$area)
d$F <- standardize(d$avgfood)
d$G <- standardize(d$groupsize)

m1 <- quap(
  alist(
    W ~ dnorm( mu , sigma ),
    mu <- a + bF*F + bG*G + bA*A,
    a ~ dnorm(0,0.2),
    c(bF,bG,bA) ~ dnorm(0,0.5),
    sigma ~ dexp(1)
  ), data=d )
m2 <- quap(
  alist(
    W ~ dnorm( mu , sigma ),
    mu <- a + bF*F + bG*G,
    a ~ dnorm(0,0.2),
    c(bF,bG) ~ dnorm(0,0.5),
    sigma ~ dexp(1)
  ), data=d )
m3 <- quap(
  alist(
    W ~ dnorm( mu , sigma ),
    mu <- a + bG*G + bA*A,
    a ~ dnorm(0,0.2),
    c(bG,bA) ~ dnorm(0,0.5),
    sigma ~ dexp(1)
  ), data=d )
m4 <- quap(
  alist(
    W ~ dnorm( mu , sigma ),
    mu <- a + bF*F,
    a ~ dnorm(0,0.2),
    bF ~ dnorm(0,0.5),
    sigma ~ dexp(1)
```



```

    ), data=d )
m5 <- quap(
  alist(
    W ~ dnorm( mu , sigma ),
    mu <- a + bA*A,
    a ~ dnorm(0,0.2),
    bA ~ dnorm(0,0.5),
    sigma ~ dexp(1)
  ), data=d )

```

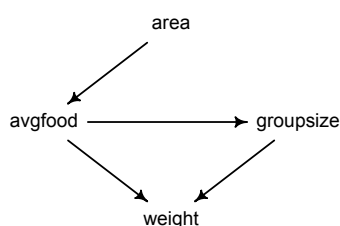
Comparing with WAIC:

```
compare( m1 , m2 , m3 , m4 , m5 )
```

R code  
7.20

	WAIC	pWAIC	dWAIC	weight	SE	dSE
m1	322.9	4.7	0.0	0.47	16.28	NA
m3	323.9	3.7	1.0	0.28	15.68	2.90
m2	324.1	3.9	1.2	0.25	16.14	3.60
m4	333.4	2.4	10.6	0.00	13.79	7.19
m5	333.7	2.7	10.8	0.00	13.79	7.24

To remind you, the DAG from the original problem is:



Notice that the top three models are m1, m3, and m2. They have very similar WAIC values. The differences are small and smaller in all cases than the standard error of the difference. WAIC sees these models are tied. This makes sense, given the DAG, because as long as a model has `groupsize` in it, we can include either `avgfood` or `area` or both and get the same inferences. Another way to think of this is that the influence of food, adjusting for group size, is (according to the DAG) the same as the influence of area, adjusting for group size, because the influence of area is routed entirely through food and group size. There are no backdoor paths.

What about the other two models, m4 and m5? These models are tied with one another, and both omit group size. Again, the influence of area passes entirely through food. So including only food or only area should produce the same inference—the total causal influence of area (or food) is just about zero. That's indeed what the posterior distributions suggest:

```
coefstab(m4,m5)
```

R code  
7.21

	m4	m5
a	0	0
bF	-0.02	NA
sigma	0.99	0.99
bA	NA	0.02
nobs	116	116

## 8. Chapter 8 Solutions

**8E1.** There are many good responses. Here are some examples.

- (1) Bread dough rises because of yeast, conditional on the presence of sugar. This is an interaction between yeast and sugar.
- (2) Education leads to higher income, conditional on field of study. Some types of degrees lead to higher income than others.
- (3) Gasoline makes a car go, unless it has no spark plugs (or belts, or wheels, etc.).

**8E2.** Only number (1) is a strict interaction. The others all imply additive influences. For (1), you could express it as “cooking with low heat leads to caramelizing, conditional on the onions not drying out.” That’s an interaction, as the effect of heat depends upon moisture. For (2), a car with many cylinders can go fast, whether or not it also has a good fuel injector. The reverse is also implied: either a fuel injector or more cylinders is sufficient to make a car go faster. For (3), the statement does not imply that the influence of parents depends upon the beliefs of friends. It just implies that one may be influenced by either parents or friends. For (4), the word “or” gives away that this is another case in which either factor, sociality or manipulative appendages (hands, tentacles), is sufficient to predict intelligence.

**8E3.**

- (1) For outcome “extent caramelized”,  $\mu_i = \alpha + \beta_H H_i + \beta_M M_i + \beta_{HM} H_i M_i$ .
- (2) For outcome “maximum speed”,  $\mu_i = \alpha + \beta_C C_i + \beta_F F_i$ .
- (3) For outcome “extent conservative”,  $\mu_i = \alpha + \beta_P P_i + \beta_F F_i$ .
- (4) For outcome “intelligence” (whatever that means when comparing an octopus to a monkey),  $\mu_i = \alpha + \beta_S S_i + \beta_M M_i$ .

**8M1.** Tulips are a winter flower in much of their natural range. High temperatures do frustrate them. This is not a *linear* interaction, because by raising temperature the effect was to prevent all blooms. A linear effect would be an additive change. But it is still correct to say that there is a three-way interaction here: the influence of water and shade depend upon one another, and both and their interaction depend upon temperature. In later chapters, you’ll see how to model non-linear responses of this kind.

**8M2.** Here is one idea. Let  $L_i$  stand for the ordinary linear model from the chapter. Then let  $H_i$  be a 0/1 indicator of whether or not the temperature was hot. Then:

$$\mu_i = L_i(1 - H_i)$$

When  $H_i = 1$ , the entire model above is zero, regardless of the value of  $L_i$ .

**8M3.** The implied relationship between ravens and wolves is one in which wolves do not need ravens, but ravens very much do benefit from wolves (at least in some places). Here's an example set of data in which this might be the case:

Region	Wolves	Ravens
1	12	43
2	15	46
3	7	28
4	30	99
5	17	60
6	70	212

Really, this “interaction” is not a static interaction effect at all, because just stating that ravens depend upon wolves implies that we can partially predict raven density with wolf density. A statistical interaction requires instead that some other third variable regulate the dependency of ravens on wolves. For example, in regions in which there is plenty of small prey for ravens to kill and consume on their own, the presence of wolves may not matter at all.

**8M4.** The direct approach is to use log-normal or exponential distributions. Both of these distributions are strictly positive. We can force the shade effect to be negative by using a minus in the linear model.

Thinking about the interaction is harder. But our prior information tells us that more water increases the effect of light. So similarly more water decreases the impact of shade. So the interaction should be negative. So we can force the interaction effect to also have a log-normal or exponential, but make it negative by adding a minus in the linear model.

Here's the code with log-normal. Remember (or look up) that the mean of a log-normal is  $\exp(\mu + \sigma^2/2)$ . So we can't make these values very large at all, if we want prior mean to be reasonable. Let's load the data and try a first guess:

```
library(rethinking)
data(tulips)
d <- tulips
d$blooms_std <- d$blooms / max(d$blooms)
d$water_cent <- d$water - mean(d$water)
d$shade_cent <- d$shade - mean(d$shade)

m8M4a <- quap(
  alist(
    blooms_std ~ dnorm( mu , sigma ) ,
    mu <- a + bw*water_cent - bs*shade_cent - bws*water_cent*shade_cent ,
    a ~ dnorm( 0.5 , 0.25 ) ,
    bw ~ dlnorm( 0 , 0.25 ) ,
    bs ~ dlnorm( 0 , 0.25 ) ,
    bws ~ dlnorm( 0 , 0.25 ) ,
    sigma ~ dexp( 1 )
  ) , data=d )
```

R code  
8.1

Before looking at the fit, let's do prior simulation to see if these priors make sense.

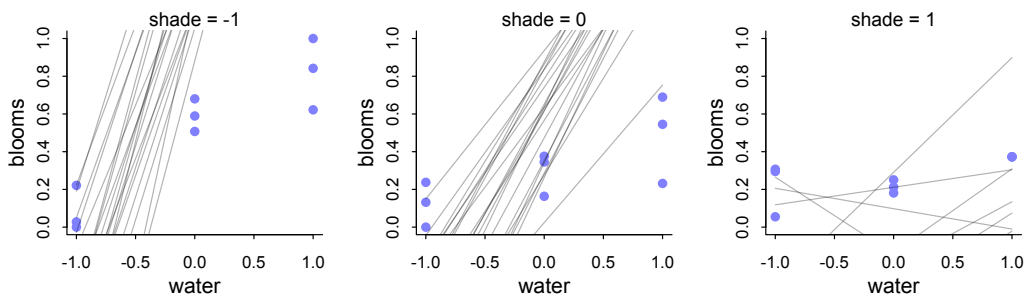
```
p <- extract.prior( m8M4a )
par(mfrow=c(1,3),cex=1.1) # 3 plots in 1 row
for ( s in -1:1 ) {
```

R code  
8.2

```

idx <- which( d$shade_cent==s )
plot( d$water_cent[idx] , d$blooms_std[idx] , xlim=c(-1,1) , ylim=c(0,1) ,
      xlab="water" , ylab="blooms" , pch=16 , col=range(2) )
mtext( concat( "shade = " , s ) )
mu <- link( m8M4 , post=p , data=data.frame( shade_cent=s , water_cent=-1:1 ) )
for ( i in 1:20 ) lines( -1:1 , mu[i,] , col=col.alpha("black",0.3) )
}

```



Wow, that's terrible. First guess was no good. We don't want to fit the data visually here, but we do want prior associations within the realm of possibility. So we need to make the prior means closer to zero. Let's try:

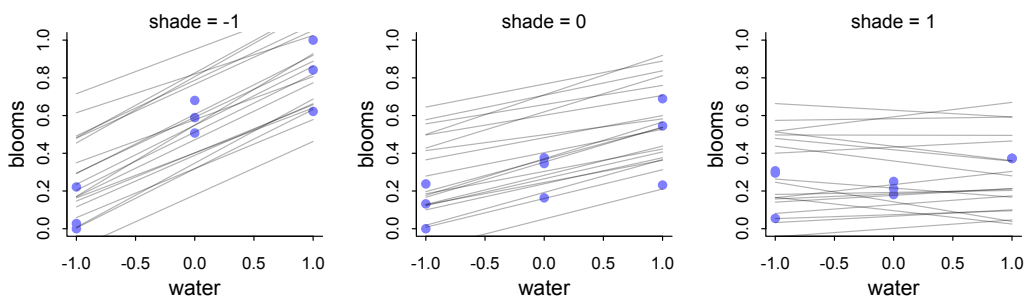
R code  
8.3

```

m8M4b <- quap(
  alist(
    blooms_std ~ dnorm( mu , sigma ) ,
    mu <- a + bw*water_cent - bs*shade_cent - bws*water_cent*shade_cent ,
    a ~ dnorm( 0.5 , 0.25 ) ,
    bw ~ dlnorm( -2 , 0.25 ) ,
    bs ~ dlnorm( -2 , 0.25 ) ,
    bws ~ dlnorm( -2 , 0.25 ) ,
    sigma ~ dexp( 1 )
  ) , data=d )

```

Extract the prior and repeat the plotting code:



That's a lot better, mostly within the empirically possible range of the outcome at least. There is still plenty of uncertainty for the data to hone down, however.

**8H1.** Let's begin by loading the data and inspecting the bed variable.

```
library(rethinking)
data(tulips)
d <- tulips
d$bed
```

R code  
8.4

```
[1] a a a a a a a a b b b b b b b b c c c c c c c c
Levels: a b c
```

This is a factor with three levels. So we could either code two dummy variables to contain the same information or rather one index variable. Both approaches were introduced in Chapter 5. I'll show both approaches here, to provide additional examples.

First, the dummy variable approach. We'll need two dummy variables, one less than the number of categories. This will construct them:

```
d$bed_b <- ifelse( d$bed=="b" , 1 , 0 )
d$bed_c <- ifelse( d$bed=="c" , 1 , 0 )
```

R code  
8.5

For bed a, it will be absorbed into the intercept, and then the coefficients for each dummy variable will contain *contrasts* with bed a.

And here is the model using these dummy variables. I'm also going to center the variables, just like in the chapter.

```
d$blooms_std <- d$blooms / max(d$blooms)
d$water_cent <- d$water - mean(d$water)
d$shade_cent <- d$shade - mean(d$shade)
m1 <- quap(
  alist(
    blooms_std ~ dnorm( mu , sigma ),
    mu <- a + bw*water_cent + bs*shade_cent +
      bws*water_cent*shade_cent +
      b_bed_b*bed_b + b_bed_c*bed_c ,
    a ~ dnorm( 0 , 0.25 ),
    c(bw,bs,bws,b_bed_b,b_bed_c) ~ dnorm(0,0.25),
    sigma ~ dexp( 1 )
  ) , data=d )
precis(m1)
```

R code  
8.6

	mean	sd	5.5%	94.5%
a	0.27	0.03	0.21	0.32
bw	0.21	0.03	0.17	0.25
bs	-0.11	0.03	-0.15	-0.07
bws	-0.14	0.03	-0.19	-0.09
b_bed_b	0.12	0.05	0.04	0.20
b_bed_c	0.14	0.05	0.06	0.22
sigma	0.11	0.01	0.08	0.13

This table is a bit monstrous to look at, but comparing to the interaction model fit in the chapter reveals that everything is basically the same, except for the presence now of the bed parameters. Both beds “b” and “c” appear to have done better than bed “a” did. How can I tell? Because both `b_bed_b` and `b_bed_c` are reliably positive.

Now let's do the model over again, using an index variable instead. This approach estimates a unique intercept for each category. To construct the index variable, you can use the convenient `coerce_index` function:

```
R code
8.7 ( d$bed_idx <- coerce_index( d$bed ) )
```

```
[1] 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2 2 3 3 3 3 3 3 3 3
```

So now bed “a” has index 1, bed “b” index 2, and bed “c” index 3. To fit the model:

```
R code
8.8 m2 <- quap(
  alist(
    blooms_std ~ dnorm( mu , sigma ),
    mu <- a[bed_idx] + bw*water_cent + bs*shade_cent +
      bws*water_cent*shade_cent ,
    a[bed_idx] ~ dnorm( 0 , 0.25 ),
    c(bw,bs,bws) ~ dnorm( 0 , 0.25 ),
    sigma ~ dexp( 1 )
  ) , data=d )
precis(m2,depth=2)
```

	mean	sd	5.5%	94.5%
a[1]	0.26	0.04	0.21	0.32
a[2]	0.39	0.04	0.33	0.44
a[3]	0.40	0.04	0.34	0.46
bw	0.21	0.03	0.17	0.25
bs	-0.11	0.03	-0.15	-0.07
bws	-0.14	0.03	-0.19	-0.09
sigma	0.11	0.01	0.08	0.13

These are (nearly) the same estimates. Again we see that beds “b” and “c” did better than bed “a”. Bed “c” did appear to grow a little better than bed “b”. But what’s the posterior distribution of that difference? We can calculate it with samples, just like the examples in Chapter 5:

```
R code
8.9 post <- extract.samples(m2)
diff_b_c <- post$a[,2] - post$a[,3]
PI( diff_b_c )
```

	5%	94%
	-0.09631392	0.06708758

So while the expected difference is there, the posterior distribution of the difference has a lot of probability on both sides of zero.

So what to make of all of this? Including bed in the analysis doesn’t change any qualitative inference about the experiment, even though there probably were differences between the beds.

**8H2.** Now let’s compare the model from the previous problem with the interaction model from the chapter.

```
R code
8.10 m3 <- quap(
  alist(
    blooms_std ~ dnorm( mu , sigma ),
    mu <- a + bw*water_cent + bs*shade_cent +
      bws*water_cent*shade_cent ,
    a ~ dnorm( 0 , 0.25 ),
    c(bw,bs,bws) ~ dnorm( 0 , 0.25 ),
    sigma ~ dexp( 1 )
```

```
) , data=d )
compare(m2,m3)
```

```
      WAIC      SE dWAIC    dSE pWAIC weight
m2 -22.8 10.16    0.0   NA  10.0    0.57
m3 -22.3 10.18    0.6  8.35   6.4    0.43
```

The model with bed, m2, does a little bit better in the WAIC comparison. But the difference is very small. Why? Because while there is some evidence that bed mattered, the treatments mattered a lot more. So including bed in the model doesn't help prediction much, as least as far as WAIC can estimate. This is all as it should be: this was a factorial experiment. In the experimental design, there is no correlation between bed and treatment. So even when there are effects of bed, we can still get good measures of the treatments. In observational studies, trying to control for common confounds like bed is much more important, typically.

**8H3.** (a) Let's load the data and fit the model, as in the chapter:

```
library(rethinking)
data(rugged)
d <- rugged
d$log_gdp <- log(d$rgdppc_2000)
dd <- d[ complete.cases(d$rgdppc_2000) , ]
dd$log_gdp_std <- dd$log_gdp / mean(dd$log_gdp)
dd$rugged_std <- dd$rugged / max(dd$rugged)
dd$cid <- ifelse( dd$cont_africa==1 , 1 , 2 )

m8.3 <- quap(
  alist(
    log_gdp_std ~ dnorm( mu , sigma ) ,
    mu <- a[cid] + b[cid]*( rugged_std - 0.215 ) ,
    a[cid] ~ dnorm( 1 , 0.1 ) ,
    b[cid] ~ dnorm( 0 , 0.3 ) ,
    sigma ~ dexp( 1 )
  ) , data=dd )
```

R code  
8.11

Now let's look at the Pareto  $k$  values, sorted and displayed with the country names:

```
k <- PSISk( m8.3 )
o <- order( k , decreasing=TRUE )
data.frame( country=as.character( dd$isocode ) ,
  k=k , dd$rugged_std , dd$log_gdp_std )[o,]
```

R code  
8.12

	country	k	dd.rugged_std	dd.log_gdp_std
145	SYC	0.64	0.7876491454	1.1501264
150	TJK	0.49	0.8547242825	0.7826922
24	BWA	0.40	0.0291841342	1.0507430
118	NPL	0.35	0.8131247985	0.8438960
62	GNQ	0.34	0.0901322154	1.1304719
84	KGZ	0.34	0.6912286359	0.8632546
93	LSO	0.34	1.0000000000	0.8994088
27	CHE	0.28	0.7676555950	1.2110100
167	YEM	0.26	0.3745565946	0.7830325

```

55      GAB  0.25  0.0351499516      1.0237147
57      GEO  0.23  0.5899709771      0.8851899
108     MUS  0.22  0.1530151564      1.0768738
144     SWZ  0.22  0.4938729442      0.9922796
...

```

Note that the Pareto  $k$  values depend upon samples, so the values you see will be a little different. But the basic order should be similar. Seychelles (SYC) does have the largest Pareto  $k$  value. But several other countries also have pretty high values. What do these countries have in common? That are poorly fit the trend, but in different ways. Seychelles has an unusually high GDP. Tajikistan (TJK) has an unusually low GDP. Botswana (BWA) has an unusually high GDP for such a flat African country.

(b) Simply change the likelihood to Student-t:

R code  
8.13

```

m8.3t <- quap(
  alist(
    log_gdp_std ~ dstudent( 2 , mu , sigma ) ,
    mu <- a[cid] + b[cid]*( rugged_std - 0.215 ) ,
    a[cid] ~ dnorm( 1 , 0.1 ) ,
    b[cid] ~ dnorm( 0 , 0.3 ) ,
    sigma ~ dexp( 1 )
  ) , data=dd )

```

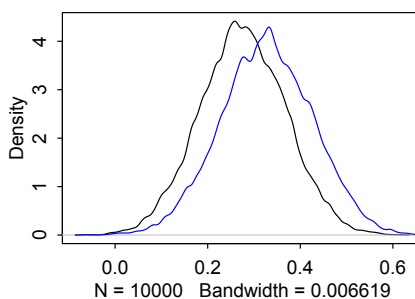
We want to know how the posterior distribution of the difference between  $b[1]$  and  $b[2]$  has changed by using the Student-t distribution. So let's do that comparison directly:

R code  
8.14

```

postN <- extract.samples( m8.3 )
postT <- extract.samples( m8.3t )
diffN <- postN$b[,1] - postN$b[,2]
diffT <- postT$b[,1] - postT$b[,2]
dens( diffN )
dens( diffT , add=TRUE , col="blue" )

```



Blue is the Student-t model. The difference has actually grown a little. The thing about the Student-t model is that it treats all the of the extreme points, not just the ones you notice with your eyes. So it can sometimes surprise you.

**8H4.** (a) You could have chosen to fit a number of different models. I'm going to fit three models that represent a basic analysis of the hypothesis: (1) A model predicting log-lang-per-capita with only a constant, (2) a model with only mean growing season as a predictor, and (3) a model that includes log(area) as a covariate. Then I'll compare them, using WAIC.



But keep in mind that WAIC is not choosing a model for us. The causal theory chooses the model. In this case that means both area and growing season influence the number of languages. These two variables may be associated with one another through other processes. For example, as a consequence of world history, larger national territories (area) may be found in particular places and therefore be associated with particular growing seasons. Since larger nations (area) tend to have fewer languages, there are clearly historical processes at work here. Can you make a DAG that expresses this?

Let's load the data and build some models:

```
library(rethinking)
data(nettle)
d <- nettle
d$L <- standardize( log( d$num.lang / d$k.pop ) )
d$A <- standardize( log(d$area) )
d$G <- standardize( d$mean.growing.season )

m0 <- quap(
  alist(
    L ~ dnorm(mu,sigma),
    mu <- a + 0,
    a ~ dnorm(0,0.2),
    sigma ~ dexp(1)
  ) , data=d )

m1 <- quap(
  alist(
    L ~ dnorm(mu,sigma),
    mu <- a + bG*G,
    a ~ dnorm(0,0.2),
    bG ~ dnorm(0,0.5),
    sigma ~ dexp(1)
  ) , data=d )

m2 <- quap(
  alist(
    L ~ dnorm(mu,sigma),
    mu <- a + bG*G + bA*A,
    a ~ dnorm(0,0.2),
    c(bG,bA) ~ dnorm(0,0.5),
    sigma ~ dexp(1)
  ) , data=d )

compare(m0,m1,m2)
```

R code  
8.15

	WAIC	SE	dWAIC	dSE	pWAIC	weight
m2	205.6	16.10	0.0	NA	4.7	0.53
m1	206.0	15.57	0.3	3.90	3.7	0.46
m0	213.7	17.22	8.1	8.09	2.7	0.01

This is very strong support for using mean growing season as a predictor of language diversity. Controlling for log(area) has little effect, as can be seen by comparing the estimates from models m1 and m2:

```
coefstab(m1,m2)
```

R code  
8.16

	m1	m2

```

a          0          0
bG         0.34       0.29
sigma      0.92       0.91
bA         NA        -0.17
nobs       74        74

```

You can check the standard errors, too, to ensure that the estimate for `mean.growing.season` is reliably positive in both cases.

Plotting the predicted relationship for `m2`:

R code  
8.17

```

G_seq <- seq( from=-2.5 , to=2 , length.out=30 )
new_dat <- data.frame( A=0 , G=G_seq )

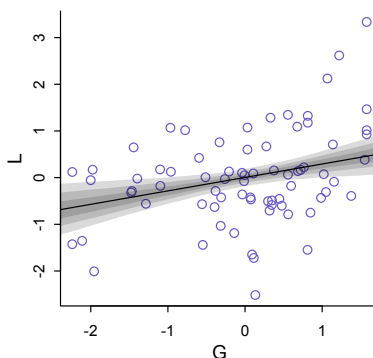
mu <- link( m2 , data=new_dat )

plot( L ~ G , data=d , col="slateblue" )
lines( G_seq , colMeans(mu) )

# plot several intervals with shading
for ( p in c(0.5,0.79,0.95) ) {
  mu_PI <- apply( mu , 2 , PI , prob=p )
  shade( mu_PI , G_seq )
}

```

And here's the plot:



Hm, that relationship doesn't actually look very linear. Most of the error is above the line, at long growing season values on the right side. In any event, there does seem to be some positive relationship between the two variables, as predicted. Maybe we shouldn't be surprised by how poor the linear relationship represents the data. After all, there are a ton of unmeasured variables that influence language density. For example, imperial expansions have mainly reduced language diversity, but have done so unequally in different regions.

(b) Fitting analogous models and comparing them (`m0` is same as before):

R code  
8.18

```

d$S <- standardize( d$sd.growing.season )
m4 <- quap(
  alist(
    L ~ dnorm(mu,sigma),
    mu <- a + bS*S,
    a ~ dnorm(0,0.2),

```

```

      bS ~ dnorm(0,0.5),
      sigma ~ dexp(1)
    ) , data=d )
m5 <- quap(
  alist(
    L ~ dnorm(mu,sigma),
    mu <- a + bS*S + bA*A,
    a ~ dnorm(0,0.2),
    c(bS,bA) ~ dnorm(0,0.5),
    sigma ~ dexp(1)
  ) , data=d )
compare(m0,m4,m5)

```

	WAIC	SE	dWAIC	dSE	pWAIC	weight
m4	211.0	17.36	0.0	NA	3.7	0.51
m5	211.8	17.21	0.8	3.86	5.4	0.34
m0	213.5	17.11	2.6	4.37	2.5	0.14

A very similar story, although not as strong a relationship. Take a look at the marginal posterior from m5:

```
precis(m5)
```

R code  
8.19

	mean	sd	5.5%	94.5%
a	0.00	0.10	-0.15	0.15
bS	-0.14	0.12	-0.34	0.06
bA	-0.19	0.12	-0.39	0.01
sigma	0.94	0.08	0.82	1.06

Including `log(area)` generates considerable uncertainty about the direction of the effect of `sd.growing.season`. Still, the MAP is negative, as predicted. Of course, there is no reason to think the effect of `S` should be linear, as this model assumes.

(c) Here's the interaction model:

```

m6 <- quap(
  alist(
    L ~ dnorm(mu,sigma),
    mu <- a + bG*G + bS*S + bGS*G*S,
    a ~ dnorm(0,0.2),
    c(bG,bS,bGS) ~ dnorm(0,0.5),
    sigma ~ dexp(1)
  ) , data=d )

```

R code  
8.20

Okay, so let's inspect the marginal posterior:

```
precis(m6)
```

R code  
8.21

	mean	sd	5.5%	94.5%
a	0.00	0.09	-0.14	0.15
bG	0.23	0.11	0.05	0.41
bS	-0.23	0.10	-0.39	-0.07
bGS	-0.24	0.10	-0.40	-0.08
sigma	0.85	0.07	0.74	0.97

You can see above that the interaction coefficient is reliably negative. What that means exactly will require some more work.

Let's see what the predictions look like. Let's show the interaction in a panel of six plots. The top row will show the relationship between `log-lang-per-capita` and `mean.growing.season`, across values of `sd.growing.season`. The bottom row will show `log-lang-per-capita` on `sd.growing.season`, across values of `mean.growing.season`. This is just to show the two-way interaction from both perspectives.

I'm also going to use transparency, as a function of distance from the value on the top of each plot, to show how the data change through the 3rd, un-plotted, dimension. The function `col.dist` in the `rethinking` package handles this for you. I didn't use it in the book, so here is an example of how it can be helpful. You can get more details about how it works from the help `?col.dist`. The key issue is the standard deviation value, which determines how quickly color fades as individual points move away from some reference value. The reference value in these plots will be the value of the third variable displayed on the top margin. You'll want to play with the standard deviation value to get a sense of how it works. Larger values mean less fading.

This code will draw a triptych of interactions, varying `mean.growing.season` along the horizontal axis and `sd.growing.season` across the plots.

R code  
8.22

```
# pull out 10%, 50%, and 95% quantiles of sd.growing.season
# these values will be used to make the three plots
S_seq <- quantile(d$S,c(0.1,0.5,0.95))

# now loop over the three plots
# draw languages against mean.growing.season in each
G_seq <- seq(from=-2.5,to=2,length.out=30)
par(mfrow=c(1,3),cex=1.1) # set up plot window for row of 3 plots
for ( i in 1:3 ) {
  S_val <- S_seq[i] # select out value for this plot
  new.dat <- data.frame(
    G = G_seq,
    S = S_val )
  mu <- link( m6 , data=new.dat )
  mu.mean <- apply( mu , 2 , mean )
  mu.PI <- apply( mu , 2 , PI )

  # fade point color as function of distance from sd.val
  cols <- col.dist( d$S , S_val , 2 , "slateblue" )

  plot( L ~ G , data=d , col=cols , lwd=2 )
  mtext( paste("S =",round(S_val,2)) , 3 )
  lines( G_seq , mu.mean )
  shade( mu.PI , G_seq )
}
```

And this code will produce the analogous triptych in which `sd.growing.season` is varied on the horizontal axis.

R code  
8.23

```
# pull out 10%, 50%, and 95% quantiles of mean.growing.season
G_seq <- quantile(d$G,c(0.1,0.5,0.95))

# now loop over the three plots
x.seq <- seq(from=-1.7,to=4,length.out=30)
par(mfrow=c(1,3),cex=1.1) # set up plot window for row of 3 plots
for ( i in 1:3 ) {
```

```

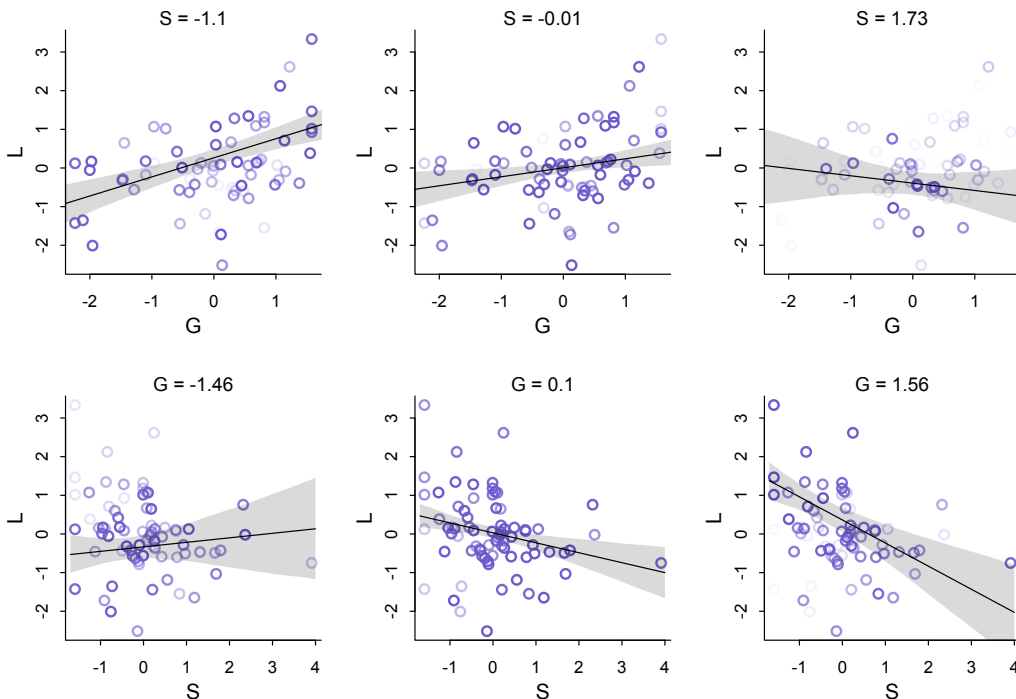
G_val <- G_seq[i] # select out value for this plot
new.dat <- data.frame(
  G = G_val ,
  S = x.seq )
mu <- link( m6 , data=new.dat )
mu.mean <- apply( mu , 2 , mean )
mu.PI <- apply( mu , 2 , PI )

# fade point color as function of distance from sd.val
cols <- col.dist( d$G , G_val , 5 , "slateblue" )

plot( L ~ S , data=d , col=cols , lwd=2 )
mtext( paste("G =",round(G_val,2)) , 3 )
lines( x.seq , mu.mean )
shade( mu.PI , x.seq )
}

```

And here are both triptych constructions:



So the models suggest that mean growing season increases language diversity, unless the variance in growing season is also high (top row). Simultaneously, variance in growing season decreases language diversity, unless the mean growing season is very short (bottom row).

**8H5.** Load the data and construct the variables we'll need:

```

library(rethinking)
data(Wines2012)

```

R code  
8.24

```
d <- Wines2012

dat_list <- list(
  S = standardize(d$score),
  jid = as.integer(d$judge),
  wid = as.integer(d$wine)
)
```

The model is straightforward. The only issue is the priors. Since I've standardized the outcome, we can use the ordinary  $N(0,0.5)$  prior from the examples in the text with standardized outcomes. Then the prior outcomes will stay largely within the possible outcome space. A bit more regularization than that wouldn't be a bad idea either.

R code  
8.25

```
m1 <- quap(
  alist(
    S ~ dnorm( mu , sigma ),
    mu <- a[jid] + w[wid],
    a[jid] ~ dnorm(0,0.5),
    w[wid] ~ dnorm(0,0.5),
    sigma ~ dexp(1)
  ), data=dat_list )
```

Since this is your first MCMC homework, we'll spend some time inspecting the chains to ensure they worked. First, the diagnostics that `precis` provides:

R code  
8.26

```
precis( m1 , 2 )
```

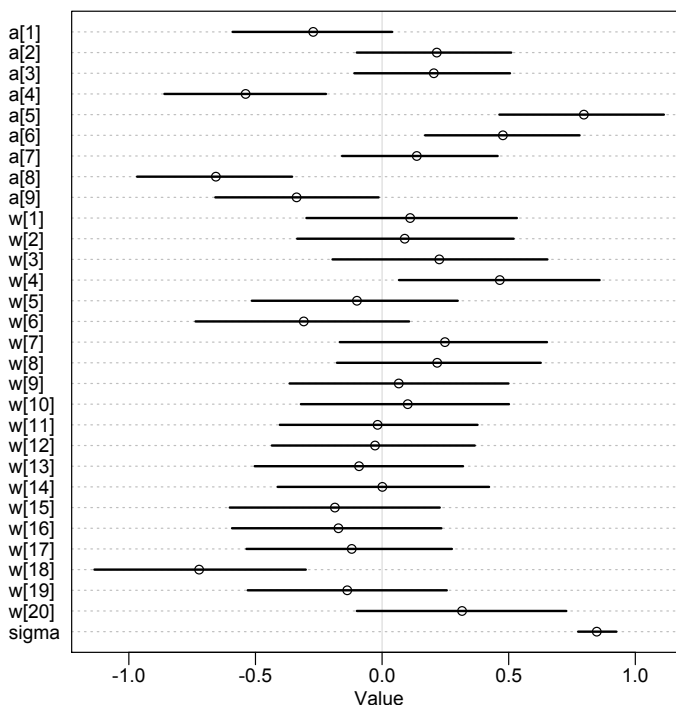
	mean	sd	5.5%	94.5%
a[1]	-0.28	0.19	-0.58	0.01
a[2]	0.22	0.19	-0.08	0.51
a[3]	0.21	0.19	-0.09	0.51
a[4]	-0.55	0.19	-0.85	-0.25
a[5]	0.81	0.19	0.51	1.11
a[6]	0.48	0.19	0.19	0.78
a[7]	0.13	0.19	-0.16	0.43
a[8]	-0.67	0.19	-0.97	-0.37
a[9]	-0.35	0.19	-0.65	-0.05
w[1]	0.12	0.25	-0.27	0.51
w[2]	0.09	0.25	-0.30	0.48
w[3]	0.24	0.25	-0.16	0.63
w[4]	0.48	0.25	0.09	0.87
w[5]	-0.11	0.25	-0.50	0.28
w[6]	-0.32	0.25	-0.71	0.07
w[7]	0.25	0.25	-0.14	0.64
w[8]	0.24	0.25	-0.16	0.63
w[9]	0.07	0.25	-0.32	0.46
w[10]	0.10	0.25	-0.29	0.50
w[11]	-0.01	0.25	-0.40	0.38
w[12]	-0.03	0.25	-0.42	0.37
w[13]	-0.09	0.25	-0.48	0.30
w[14]	0.01	0.25	-0.39	0.40
w[15]	-0.19	0.25	-0.58	0.20
w[16]	-0.17	0.25	-0.57	0.22

```
w[17] -0.12 0.25 -0.52 0.27
w[18] -0.75 0.25 -1.14 -0.35
w[19] -0.14 0.25 -0.53 0.25
w[20] 0.33 0.25 -0.06 0.73
sigma 0.79 0.04 0.72 0.85
```

Now let's plot these parameters so they are easier to interpret:

```
plot( precis( m1 , 2 ) )
```

R code  
8.27



The  $a$  parameters are the judges. Each represents an average deviation of the scores. So judges with lower values are harsher on average. Judges with higher values liked the wines more on average. There is some noticeable variation here. It is fairly easy to tell the judges apart.

The  $w$  parameters are the wines. Each represents an average score across all judges. Except for wine 18 (a New Jersey red I think), there isn't that much variation. These are good wines, after all. Overall, there is more variation from judge than from wine.

**8H6.** The easiest way to code the data is to use indicator variables. Let's look at that approach first. I'll do an index variable version next. I'll use the three indicator variables  $w$  (NJ wine),  $j$  (American NJ), and  $r$  (red wine).

```
dat_list2 <- list(
  S = standardize(d$score),
  W = d$wine.amer,
  J = d$judge.amer,
  R = ifelse(d$flight=="red",1L,0L)
)
```

R code  
8.28

The model structure is just a linear model with an ordinary intercept. I'll put a relatively tight prior on the intercept, since it must be near zero (centered outcome). What about the coefficients for the indicator variables? Let's pretend we haven't already seen the results from Problem 1—there aren't any big wine differences to find there. Without that cheating foresight, we should consider what the most extreme effect could be. How big could the difference between NJ and French wines be? Could it be a full standard deviation? If so, then maybe a  $\text{Normal}(0,0.5)$  prior makes sense, since they place a full standard deviation difference out in the tails of the prior. I'd personally be inclined to something even tighter, so that it regularizes more. But let's go with these wide priors, which nevertheless stay within the outcome space. It would make even more sense to put a tighter prior on the difference between red and white wines—on average they should be the no different, because judges only compare within flights. Here's the model:

R code  
8.29

```
m2a <- quap(
  alist(
    S ~ dnorm( mu , sigma ),
    mu <- a + bW*W + bJ*J + bR*R,
    a ~ dnorm( 0 , 0.2 ),
    c(bW,bJ,bR) ~ dnorm( 0 , 0.5 ),
    sigma ~ dexp(1)
  ), data=dat_list2 )
precis( m2a )
```

	mean	sd	5.5%	94.5%
a	-0.01	0.12	-0.21	0.18
bW	-0.18	0.13	-0.40	0.04
bJ	0.23	0.13	0.01	0.44
bR	-0.01	0.14	-0.22	0.21
sigma	0.98	0.05	0.90	1.06

As expected, red and wines are on average the same—bR is right on top of zero. American judges seem to be more on average slightly more generous with ratings—bJ is slightly but reliably above zero. American wines have slightly lower average ratings than French wines—bW is mostly below zero, but not very large in absolute size.

Okay, now for an index variable version. The thing about index variables is that you can easily end up with more parameters than in an equivalent indicator variable model. But it's still the same posterior distribution. You can convert from one to the other (if the priors are also equivalent). We'll need three index variables:

R code  
8.30

```
dat_list2b <- list(
  S = standardize(d$score),
  wid = d$wine.amer + 1L,
  jid = d$judge.amer + 1L,
  fid = ifelse(d$flight=="red",1L,2L)
)
```

Now wid is 1 for a French wine and 2 for a NJ wine, jid is 1 for a French judge and 2 for an American judge, and fid is 1 for red and 2 for white. Those 1L numbers are just the R way to type the number as an integer—"1L" is the integer 1, while "1" is the real number 1. We want integers for an index variable.

Now let's think about priors for the parameters that correspond to each index value. Now the question isn't how big the difference could be, but rather how far from the mean an indexed category could be. If we use  $\text{Normal}(0,0.5)$  priors, that would make a full standard deviation difference from the global mean rare. It will also match what we had above, in a crude sense. Again, I'd be tempted



to something narrow, for the sake of regularization. But certainly something like  $\text{Normal}(0,10)$  is flat out silly, because it makes impossible values routine. Let's see what we get:

```
m2b <- quap(
  alist(
    S ~ dnorm( mu , sigma ),
    mu <- w[wid] + j[jid] + f[fid],
    w[wid] ~ dnorm( 0 , 0.5 ),
    j[wid] ~ dnorm( 0 , 0.5 ),
    f[wid] ~ dnorm( 0 , 0.5 ),
    sigma ~ dexp(1)
  ), data=dat_list2b )
precis( m2b , 2 )
```


R code  
8.31

```
      mean   sd  5.5% 94.5%
w[1]   0.09 0.30 -0.39  0.57
w[2]  -0.09 0.30 -0.57  0.39
j[1]  -0.12 0.30 -0.60  0.36
j[2]   0.12 0.30 -0.36  0.60
f[1]   0.00 0.30 -0.48  0.48
f[2]   0.00 0.30 -0.47  0.48
sigma  0.98 0.05  0.90  1.06
```

To see that this model is the same as the previous, let's compute contrasts. The contrast between American and French wines is:

```
post <- extract.samples(m2b)
diff_w <- post$w[,2] - post$w[,1]
precis( diff_w )
```

R code  
8.32

```
'data.frame': 10000 obs. of 1 variables:
      mean   sd  5.5% 94.5%      histogram
diff_w -0.18 0.15 -0.42  0.05 
```

That's almost exactly the same mean and standard deviation as *bw* in the first model. The other contrasts match as well.

**8H7.** I'll use the indicator variable approach here, because it'll be much easier. Once you start using MCMC in the next chapter, it'll be possible to define very flexible parameter structures. Then the index approach will be easy again.

For the indicator approach, we can use the same predictor variables as before:

```
dat_list2 <- list(
  S = standardize(d$score),
  W = d$wine.amer,
  J = d$judge.amer,
  R = ifelse(d$flight=="red",1L,0L)
)
```

R code  
8.33

It's the model that is different.

```
m3 <- quap(
  alist(
```

R code  
8.34

```

S ~ dnorm( mu , sigma ),
mu <- a + bW*W + bJ*J + bR*R +
      bWJ*W*J + bWR*W*R + bJR*J*R,
a ~ dnorm(0,0.2),
c(bW,bJ,bR) ~ dnorm(0,0.5),
c(bWJ,bWR,bJR) ~ dnorm(0,0.25),
sigma ~ dexp(1)
), data=dat_list2 )

```

I used the same priors as before for the main effects. I used tighter priors for the interactions. Why? Because interactions represent sub-categories of data, and if we keep slicing up the sample, differences can't keep getting bigger. Again, the most important thing is not to use flat priors like Normal(0,10) that produce impossible outcomes.

R code  
8.35

```
precis(m3)
```

	mean	sd	5.5%	94.5%
a	-0.05	0.13	-0.25	0.16
bW	-0.07	0.17	-0.35	0.20
bJ	0.21	0.18	-0.08	0.49
bR	0.09	0.18	-0.21	0.38
bWJ	-0.02	0.18	-0.31	0.27
bWR	-0.23	0.18	-0.52	0.06
bJR	0.05	0.18	-0.24	0.34
sigma	0.98	0.05	0.89	1.06

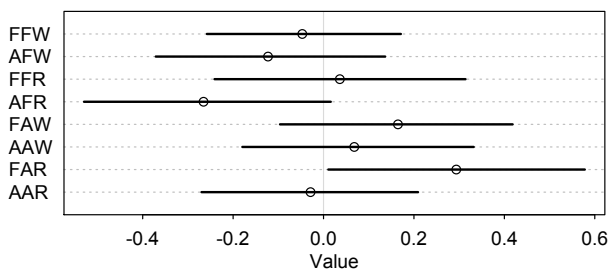
Reading the parameters this way is not easy. But right away you might notice that bW is now close to zero and overlaps it a lot on both sides. NJ wines are no longer on average worse. So the interactions did something. Glancing at the interaction parameters, you can see that only one of them has much mass away from zero, bWR, the interaction between NJ wines and red flight, so red NJ wines. To get the predicted scores for red and white wines from both NJ and France, for both types of judges, we can use link:

R code  
8.36

```

pred_dat <- data.frame(
  W = rep( 0:1 , times=4 ),
  J = rep( 0:1 , each=4 ),
  R = rep( c(0,0,1,1) , times=2 )
)
mu <- link( m3 , data=pred_dat )
row_labels <- paste( ifelse(pred_dat$W==1,"A","F") ,
                     ifelse(pred_dat$J==1,"A","F") ,
                     ifelse(pred_dat$R==1,"R","W") , sep="" )
plot( precis( list(mu=mu) , 2 ) , labels=row_labels )

```



I've added informative labels. FFW means: French wine, French judge, White wine. So the first four rows are as judged by French judges. The last four are as judged by American judges. The two rows that jump out are the 4th and the 2nd-to-last, AFR and FAR. Those are NJ red wines as judged by French judges and French red wines as judged by American judges. French judges didn't like NJ reds so much (really only one NJ red, if you look back at Problem 1). And American judges liked French reds more. Besides these two interactions, notice that it is very hard to figure this out from the table of coefficients.