

Test 2 – Take Home Problems

Stat 341 — Spring 2017

Loading some data

The command below will load some objects for this test. In particular, this will load

- the model objects `m1`, `...`, `m7` (so you don't have to wait for Stan to fit them)
- the `Pallets2` data frame

```
# load some data for this test.  
load(file = url("http://www.calvin.edu/~rpruim/data/s341/Test2.Rda"))
```

Pallet repair



Figure 1:

A local company repairs pallets like the ones in the image above. They are interested to know which of their four employees are more or less efficient at this. Since they don't know much statistics, they enlist your help as a consultant to answer their question.

You don't know much about pallets, so you ask them "How many pallets can a person repair in one day?" The answer: "Oh that depends. I'm guessing it's about 100 to 150." We'll use that information below to specify some reasonable (but quite flat) priors.

The data they provide (in `Pallets2`) gives the number of pallets repaired by each of four workers on each of five days. I've already added indexed versions of `employee` and `day` to the data.

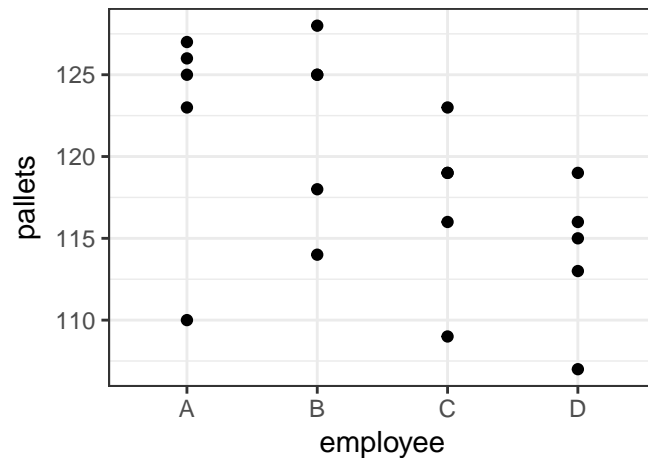
```
Pallets2 %>% head()
```

```
##   pallets employee   day emp_idx day_idx  
## 1     123        A day1      1      1  
## 2     127        A day2      1      2  
## 3     126        A day3      1      3  
## 4     125        A day4      1      4  
## 5     110        A day5      1      5
```

```
## 6      118      B day1      2      1
```

Here's a quick look at the data.

```
gf_point(pallets ~ employee, data = Pallets2)
```



Two simple models

As you have seen in class, it is often good to start with a simple model and build up from there. Here are two models. (Remember, you don't need to run this code if you load the models above.)

```
set.seed(123)
m1 <- map2stan(
  alist(
    pallets ~ dnorm(mu, sigma),
    mu <- a,
    a ~ dnorm(125, 20),
    sigma ~ dcauchy(0, 4)
  ),
  data = Pallets2
)
```

```
set.seed(123)
m2 <- map2stan(
  alist(
    pallets ~ dnorm(mu, sigma),
    mu <- b[emp_idx],
    b[emp_idx] ~ dnorm(125, 20),
    sigma ~ dcauchy(0, 4)
  ),
  data = Pallets2
)
```

1. Which model fits the data (in `Pallets2`) better? How do you know?
2. Which model would you expect to fit next week's data better? How much better? Explain.
3. If you stopped here (we're not going to do that), what would you conclude about overall employee-to-employee variation? Explain.
4. Which employee does model `m2` think is best (most efficient)? Worst?

5. How sure is model m2 that the most efficient employee is more efficient than the least efficient employee? Use posterior sampling to quantify your answer.

Day to day

Time for another model. This time we are looking at how the number of pallets repaired by one employee depends on the day.

```
set.seed(123)
m3 <- map2stan(
  alist(
    pallets ~ dnorm(mu, sigma),
    mu <- d[day_idx],
    d[day_idx] ~ dnorm(125, 20),
    sigma ~ dcauchy(0, 4)
  ),
  data = Pallets2
)
```

6. What does m3 say about day to day variation in the number of pallets repaired?
7. Comparing m2 and m3, which seems to have a greater impact on the number of pallets repaired: the employee or the day? Explain.

Employees and Days

Below are four models that include both day and employee.

8. Why is it important to include both employee and day in our model when we are really only interested in differences among the employees?

```
set.seed(44444)
m4 <- map2stan(
  alist(
    pallets ~ dnorm(mu, sigma),
    mu <- b[emp_idx] + d * day_idx,
    b[emp_idx] ~ dnorm(125, 20),
    d ~ dnorm(0, 10),
    sigma ~ dcauchy(0, 4)
  ),
  data = Pallets2, refresh = 0, iter = 5000, warmup = 1000
)
```

```
set.seed(55555)
m5 <- map2stan(
  alist(
    pallets ~ dnorm(mu, sigma),
    mu <- b[emp_idx] + d1 * day_idx + d2 * day_idx2,
    b[emp_idx] ~ dnorm(125, 20),
    c(d1, d2) ~ dnorm(0, 10),
    sigma ~ dcauchy(0, 4)
  ),
  data = Pallets2 %>% mutate(day_idx2 = day_idx^2),
  refresh = 0, iter = 7000, warmup = 1000
)
```

```

set.seed(66666)
m6 <- map2stan(
  alist(
    pallets ~ dnorm(mu, sigma),
    mu <- b[emp_idx] + d[day_idx],
    b[emp_idx] ~ dnorm(125, 20),
    d[day_idx] ~ dnorm(0, 10),
    sigma ~ dcauchy(0, 4)
  ),
  data = Pallets2, refresh = 0, iter = 7000, warmup = 1000)

```

```

set.seed(77777)
m7 <- map2stan(
  alist(
    pallets ~ dnorm(mu, sigma),
    mu <- b[emp_idx] + d[day_idx],
    b[emp_idx] ~ dnorm(0, 10),
    d[day_idx] ~ dnorm(125, 20),
    sigma ~ dcauchy(0, 4)
  ),
  data = Pallets2, refresh = 0, iter = 7000, warmup = 1000)

```

9. The only difference between the definitions `m6` and `m7` is the priors. What does this difference in priors do?
10. What would happen if we used `dnorm(125, 20)` for both (sets of) normal priors? Why is this not as good as `m6` and `m7`?
11. What would happen if we used `dnorm(0, 10)` for both (sets of) normal priors? Why is this not as good as `m6` and `m7`?

Choose Wisely

12. Time to pick your favorite model. Explain how you chose it and provide at least one plot to show how well/poorly it works.
13. Are there any indications of problems with the `map2stan()` fit for your favorite model? (This should probably have been asked earlier, but I didn't want you to have to do it for all of the models.)