

251 Project

Derek Walton & Bryce Martin

Project idea: we want to compare the quality of the two closest malls to BYU: Provo City Center & University Parkway. We will be gathering sample ratings by randomly picking stores from each and getting their ratings off of google maps. Our parameters of interest are the average rating for the stores in each location (our ‘populations’), as well as the standard deviation of the ratings. Another thing we could look into is any correlation between number of ratings and overall rating.

We’ll likely model the ratings with a beta distribution (adjusting the 5 star to be ‘100%’ and 1 star to be ‘0%’) and model the count data with a normal distribution (taking the log to transform the data).

The rating data is easy to (manually) get from google maps, and we’ll randomly pick stores based on the malls’ respective websites with all their stores listed.

We hope that at the end of this analysis we will be able to determine the overall quality of both malls, how close the stores tend to be of similar quality, and possibly even if the the more ‘niche’ stores (stores with lower counts of ratings) have a different rating on average than the high rating count stores. This can help local residents make informed decisions on which mall to go to when they are looking for consistent quality or fewer people.

Data Prep

Given the output of the confidence intervals, we can not confidently say that the $\log(n)$ has any significant correlative effect on the ratings.

We are therefore skip trying to estimate that as a parameter and estimate only two things for each population: the true ratings for each and the standard deviation of those ratings.

Because there isn’t a known posterior distribution for our rating data, we use a Monte Carlo approximation to analyze our data. Our prior is an uninformative uniform distribution because we wanted to ensure a prior that would least effect our posterior distribution. For our likelihood, we assume a good approximation may be the beta distribution with $\alpha = 3$ and $\beta = 1.5$, as we imagine the data is more left skewed. We acknowledge that the beta distribution doesn’t allow for the endpoints, so we slightly change the data on the endpoints (5 star to

4.99, 1 star to 1.01). Since we aren't predicting how any individual would rate the store but the stores overall rating, we can confidently say that a given true store's rating is not equal to exactly 1 or 5 stars. This allows us to use the beta distribution for our likelihood.

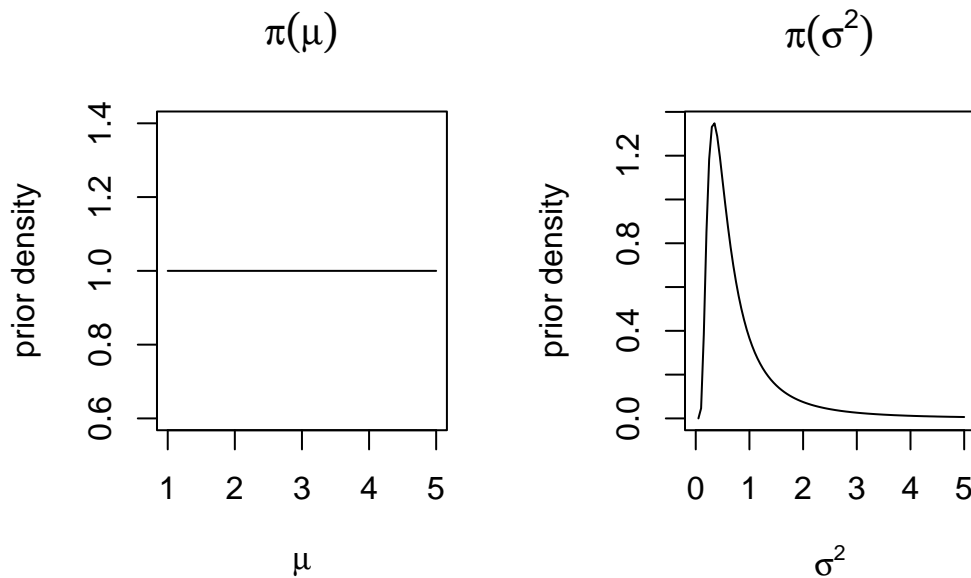
Our model for a single store would be:

x_i = Average rating for the i th store at a specific mall

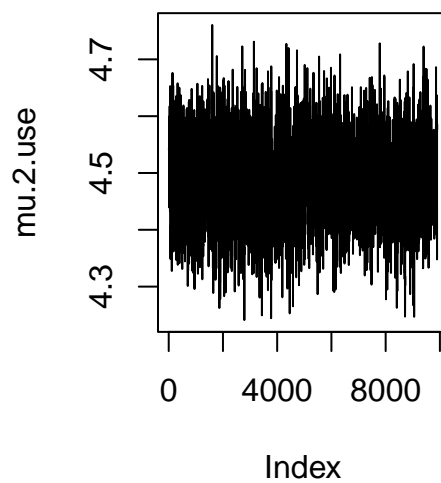
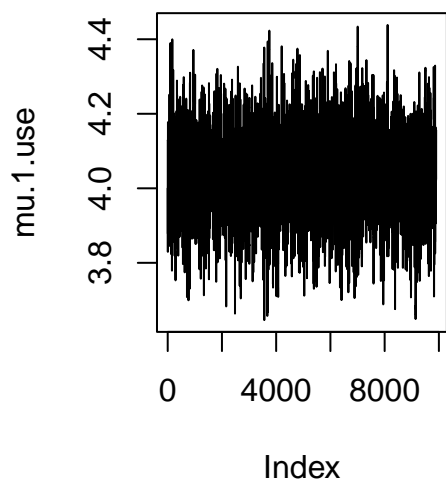
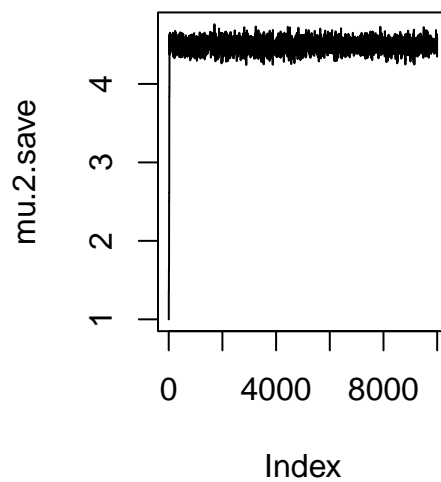
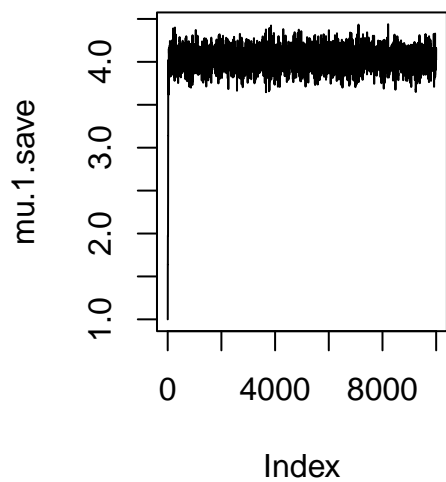
Data = x_1, x_2, \dots, x_n

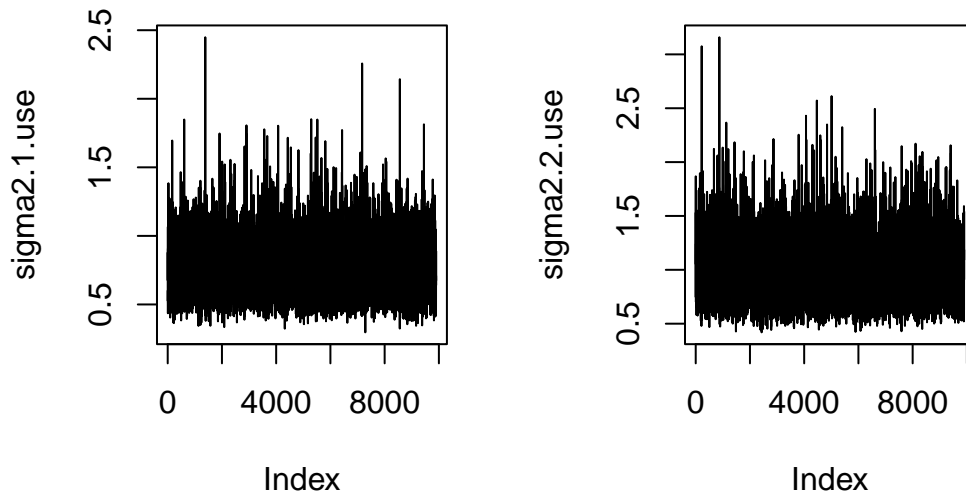
Prior: $X \sim \text{Unif}(0, 1)$

Likelihood: $f(\text{Data}|\mu, \sigma^2) \sim \text{Beta}(3, 1.5)$

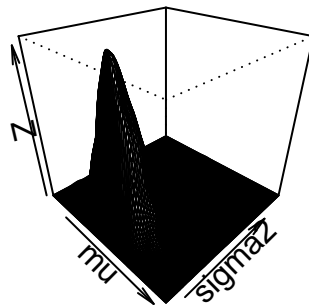


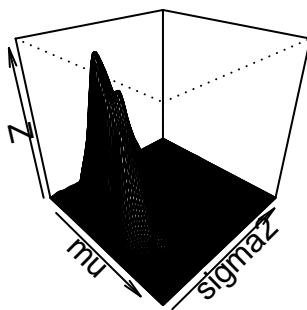
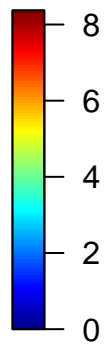
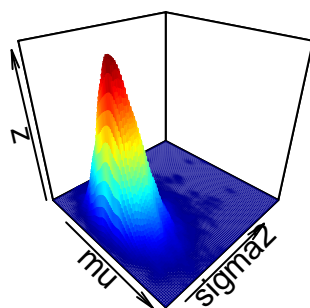
Because we're trying to estimate two parameters of interest, we will use Gibbs Sampling to estimate both the average rating as well as the variance of the ratings.

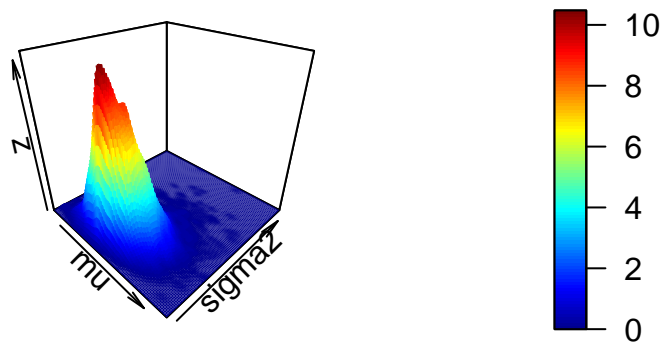




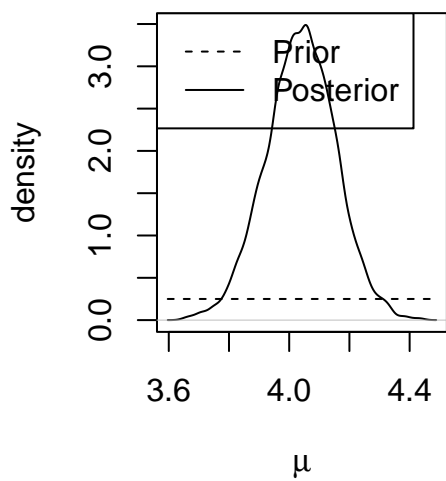
We are not able to derive the posterior/full conditional distributions for the parameters we are estimating since there is no known distribution for a normal prior/beta likelihood or a double inverse gamma distribution. However we are able to plot them as a 3D graph, and not that there doesn't appear to be any significant correlation between the two for either populations.



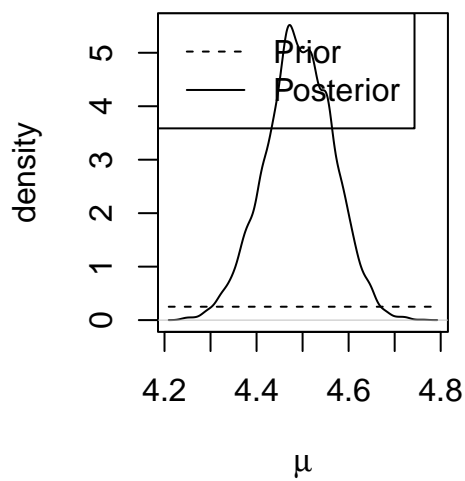




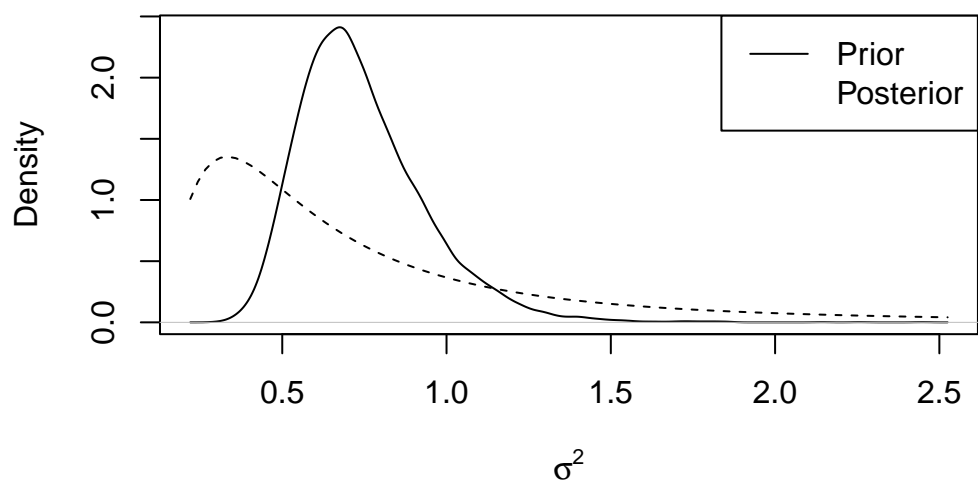
$\pi(\mu \text{ (University Mall)} \mid \text{data})$



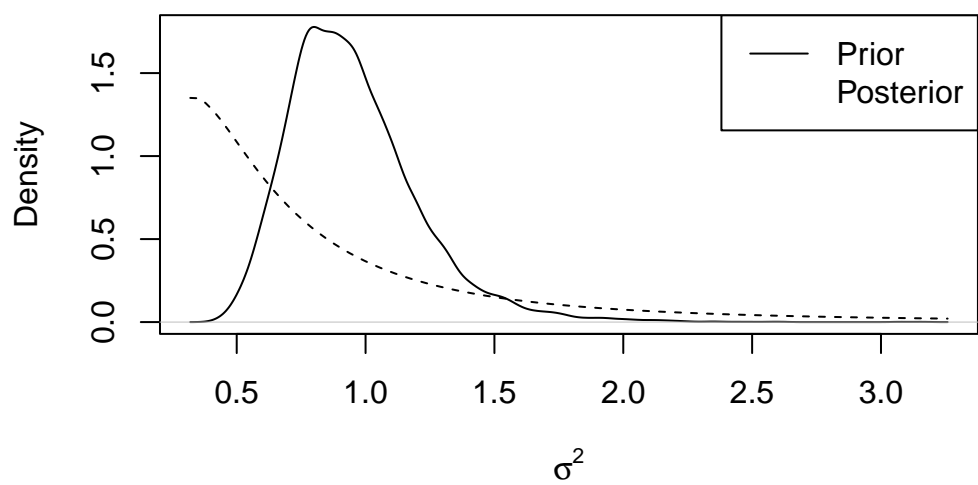
$\pi(\mu \text{ (Provo Mall)} \mid \text{data})$

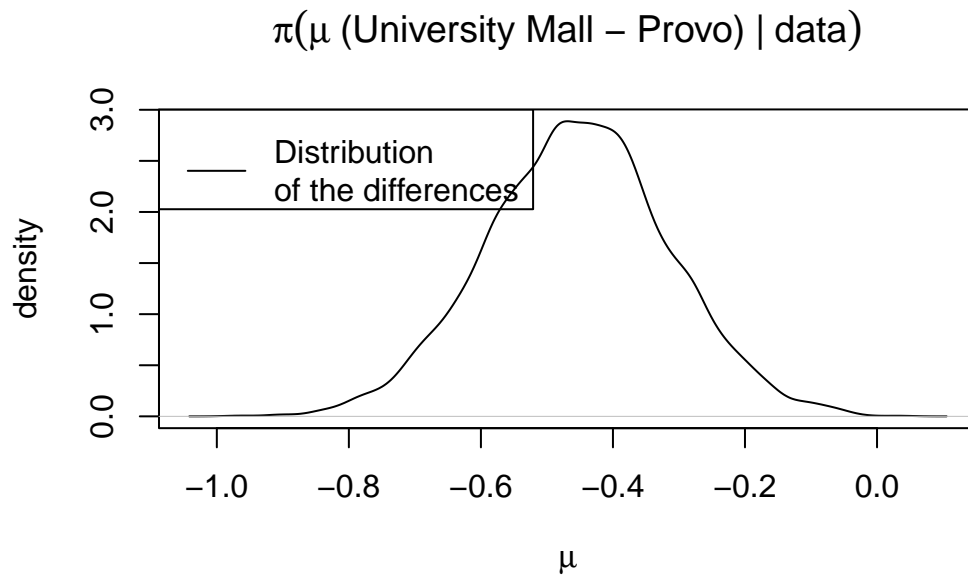


$\pi(\sigma^2 \text{ (University Mall)} | \text{ data})$



$\pi(\sigma^2 \text{ (Provo Mall)} | \text{ data})$





```
#95% credible interval
print("Mean and 95% CI for average rating of University Mall")
```

```
[1] "Mean and 95% CI for average rating of University Mall"
```

```
mean(mu.1.use)
```

```
[1] 4.038684
```

```
quantile(mu.1.use, c(.025, .975))
```

```
      2.5%      97.5%
3.808151 4.264046
```

```
print("Mean and 95% CI for average rating of Provo Towne Centre Mall")
```

```
[1] "Mean and 95% CI for average rating of Provo Towne Centre Mall"
```



```
mean(mu.2.use)
```

```
[1] 4.492282
```

```
quantile(mu.2.use, c(.025, .975))
```

```
      2.5%      97.5%  
4.343450 4.635386
```

```
# Given our data and prior knowledge, there is a 95% chance that  
# the true average rating for University Mall is between 3.81 and 4.25 and  
# the true average rating for Provo Towne Centre Mall is between 4.34 and 4.63
```

We are able to conclude that, since our 95% CIs of our average ratings don't overlap, they are statistically significant. The true average rating for University Mall is between 3.81 and 4.25, and the true average rating for Provo Towne Centre Mall is between 4.34 and 4.63.

```
#posterior mean of the average variance in mall ratings  
#95% credible interval for sigma2 for University Mall  
print("Mean and 95% CI for variance of ratings at University Mall")
```

```
[1] "Mean and 95% CI for variance of ratings at University Mall"
```

```
mean(sigma2.1.use)
```

```
[1] 0.7436179
```

```
quantile(sigma2.1.use, c(.025, .975))
```

```
      2.5%      97.5%  
0.4571692 1.1970401
```

```
#95% credible interval for sigma2 for Provo Towne Centre Mall  
print("Mean and 95% CI for variance of ratings at Provo Towne Centre Mall")
```

```
[1] "Mean and 95% CI for variance of ratings at Provo Towne Centre Mall"
```

```
mean(sigma2.2.use)
```

```
[1] 0.9529601
```

```
quantile(sigma2.2.use, c(.025, .975))
```

```
      2.5%      97.5%  
0.5721731 1.5648278
```

```
#While the variation in the plots may look different, their 95% confidence intervals overlap
```

However, we aren't able to conclude the same for our variance, since our 95% CIs of our average ratings overlap. The true average rating for University Mall is between .455 and 1.20, and the true average rating for Provo Towne Centre Mall is between 0.581 and 1.60.

Appendix (all code)