

STA610 Case Study 1

Emily Gentles, Weiyl Liu, Jack McCarthy, Qinzhe Wang

06 October, 2021

Qinzhe - Checker: Double-checks the work for reproducibility and errors. Also responsible for submitting the report and presentation files. Coordinator: Keeps everyone on task and makes sure everyone is involved. Also responsible for coordinating team meetings and defining the objectives for each meeting.

Emily - Presenter: Primarily responsible for organizing and putting the team presentations together.

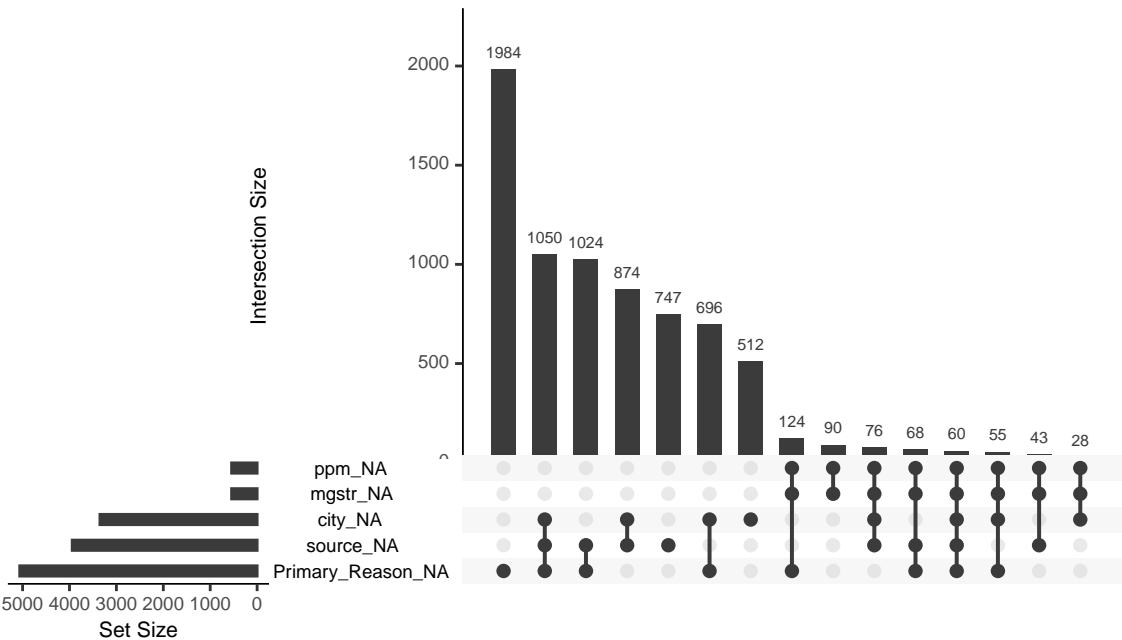
Jack - Programmer: Primarily responsible for all things coding. The programmer is responsible for putting everyone’s code together and making sure the final product is “readable”.

Wiyi - Writer: Primarily responsible for putting together the final report.

Introduction

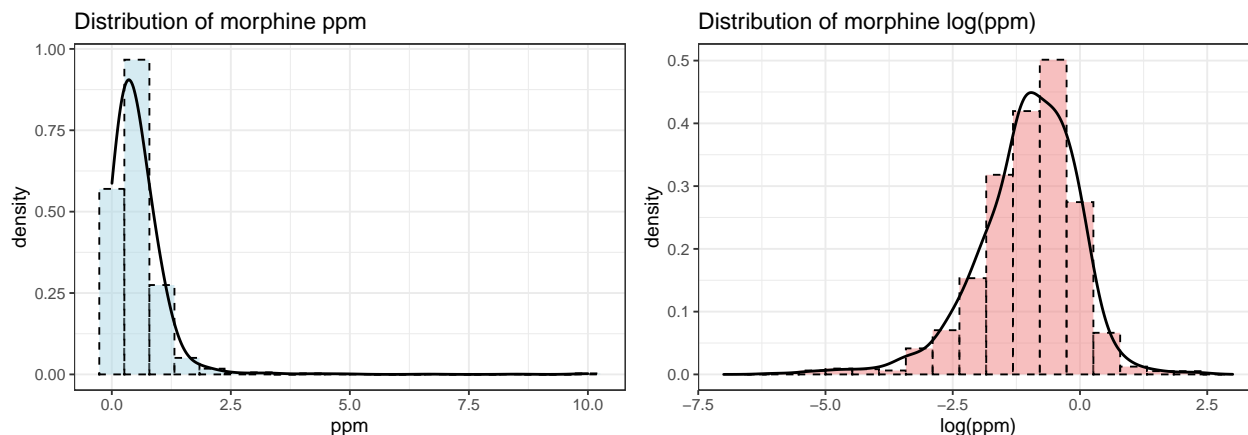
EDA

Missing Values



Response Distribution

First, a look at the distributions of the response variable “ppm”. Observations with ppm between the 0.1 and 99.9 percentiles were considered so as to avoid the influence of extreme outliers on the analysis of the ppm distribution.



The distribution of ppm is clearly right-skewed, and it is strictly nonnegative in value, so a log transformation may be appropriate. The distribution of $\log(\text{ppm})$ is given above, and appears closer to the desired normal.

state vs. $\log(\text{ppm})$

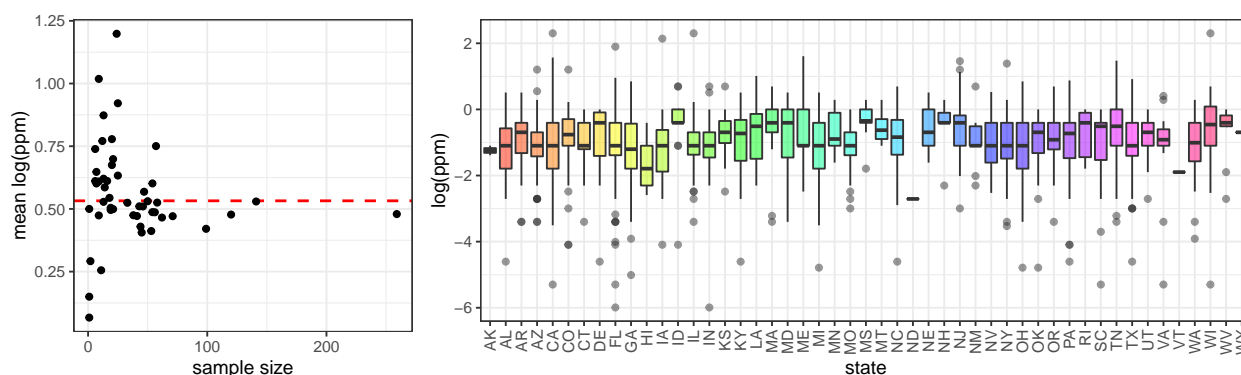
We see that there are 4 states that have a sample size of 1, North Dakota, Vermont, Washington DC, and Wyoming, as well as 1 state that has a sample size of 2, Alaska. Due to the extremely small sample sizes we decided to remove these states from our dataset to avoid computational instability.

Table 1: 7 states with smallest sample size

North Dakota	Vermont	Washington, DC	Wyoming	Alaska
1	1	1	1	2

Table 2: 7 states with largest sample size

Arizona	Michigan	Texas	Florida	California
71	99	120	141	259

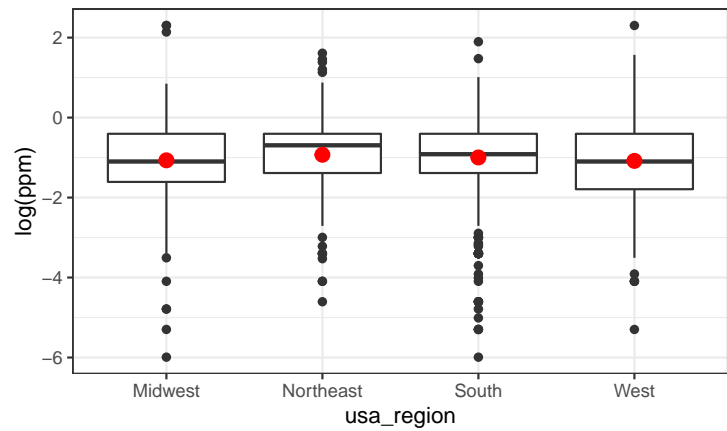


We observe that the within-state means for states with higher sample sizes in general adhere more closely to the grand mean. It is also evident that the $\log(\text{ppm})$ distributions differ little as compared to the within-state variance. This is conducive to the borrowing of information between states.

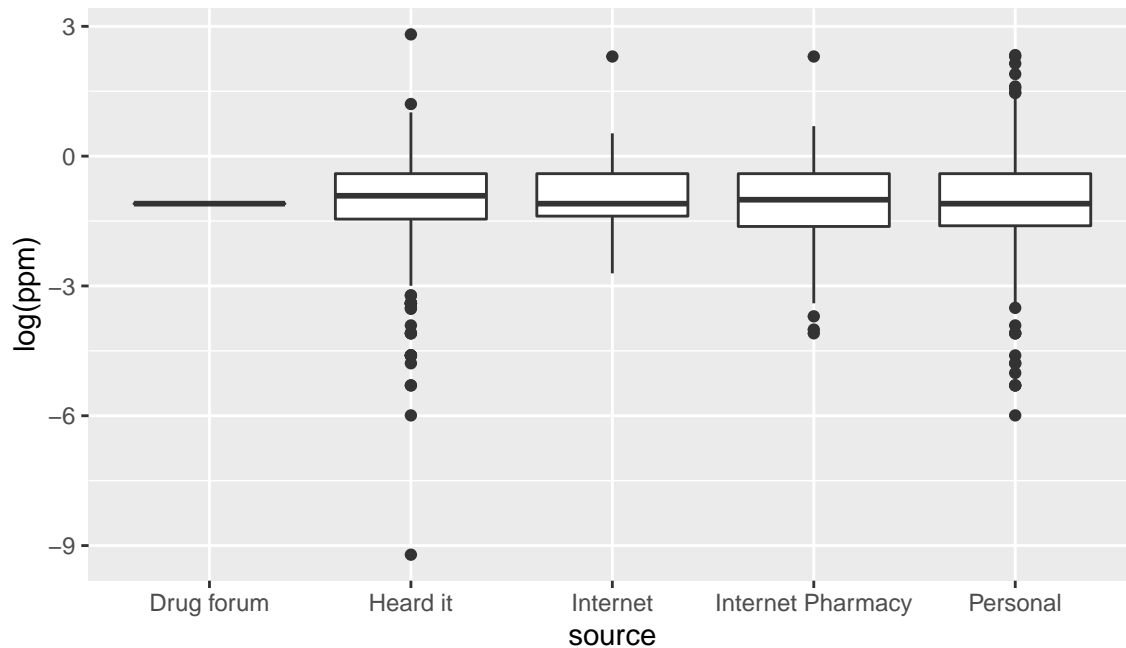
region vs. $\log(\text{ppm})$

We also have access to the broader region in which a purchase is made. This could be useful if we wanted to develop a simpler model that still captured variation by purchase location.

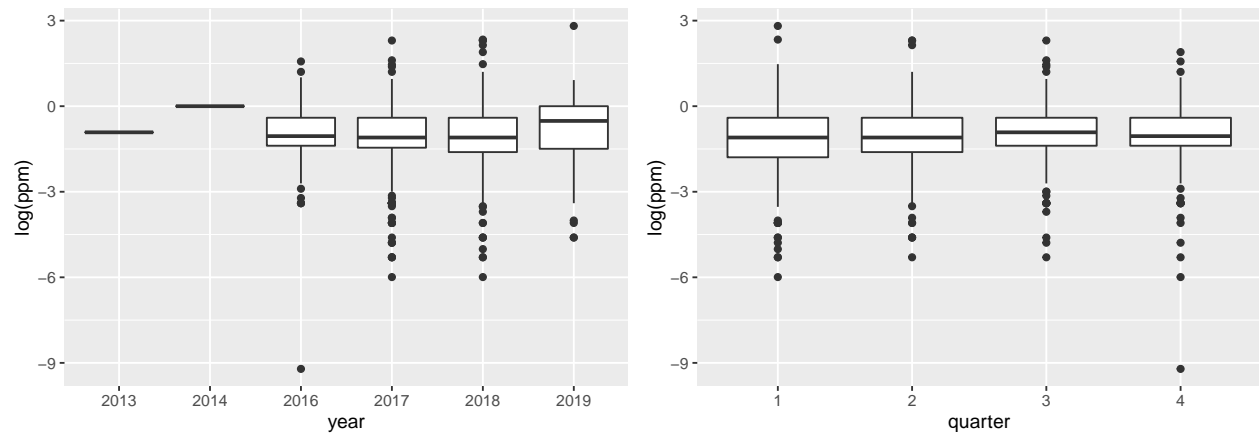
	usa_region	n	mean
1	Midwest	386	-1.069
2	Northeast	191	-0.930
3	South	673	-0.998
4	West	583	-1.083



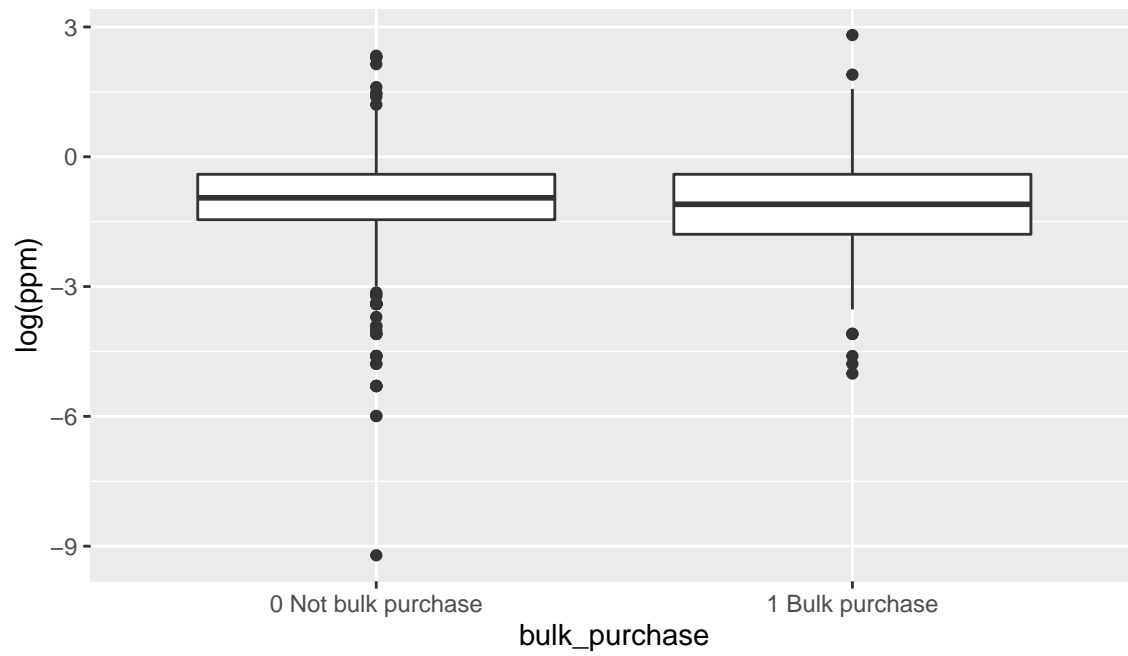
source vs. $\log(\text{ppm})$



year & quarter vs.log(ppm)



bulk_purchase vs.log(ppm)



Primary_Reason vs.log(ppm)



Model

choose grouping variable

Grouping	BIC
State	9528.336
City	9539.977
Region	9539.042

Choose **State** as our grouping variable

```
## [1] 6.192383e-07
```

```
## [1] NA
```

```
## [1] 0.01677248
```

```
## [1] 0.0009368742
```

```
## [1] 1.567514e-57
```

```
## [1] 0.2103971
```

We now have **price_date**, **bulk_purchase**, **primary_reason** and **mgstr** in our model, regarding **state** as the grouping variable.

Interaction

final model

```
log(ppm) ~ bulk_purchase + price_date + primary_reason + mgstr + (1 | state)
```

Influence

Heterogeneity across states

I removed the original model1, instead, we have final model now (no interaction version)

so the the chunk below, I changed the model to be tested to the final model.

Interclass correlation is 0.0159, very small so very little correlation across states. Including bulk purchases, the interclass correlation is 0.016, so bulk purchase actually increases the heterogeneity across states by a very small amount.

Make table with results for all models tested in ANOVA

