

# Analysis of amenity

2019-04-03

## Data Prep

1. Dropped records with missing amenity values at t1 and t2. Call the resulting data set `df`.
2. Created long-format version:
  - `df_long`: long-format version of the full set `df`

## Analyze the full set `df`

Table 1: Sample Summary Statistics of amenity

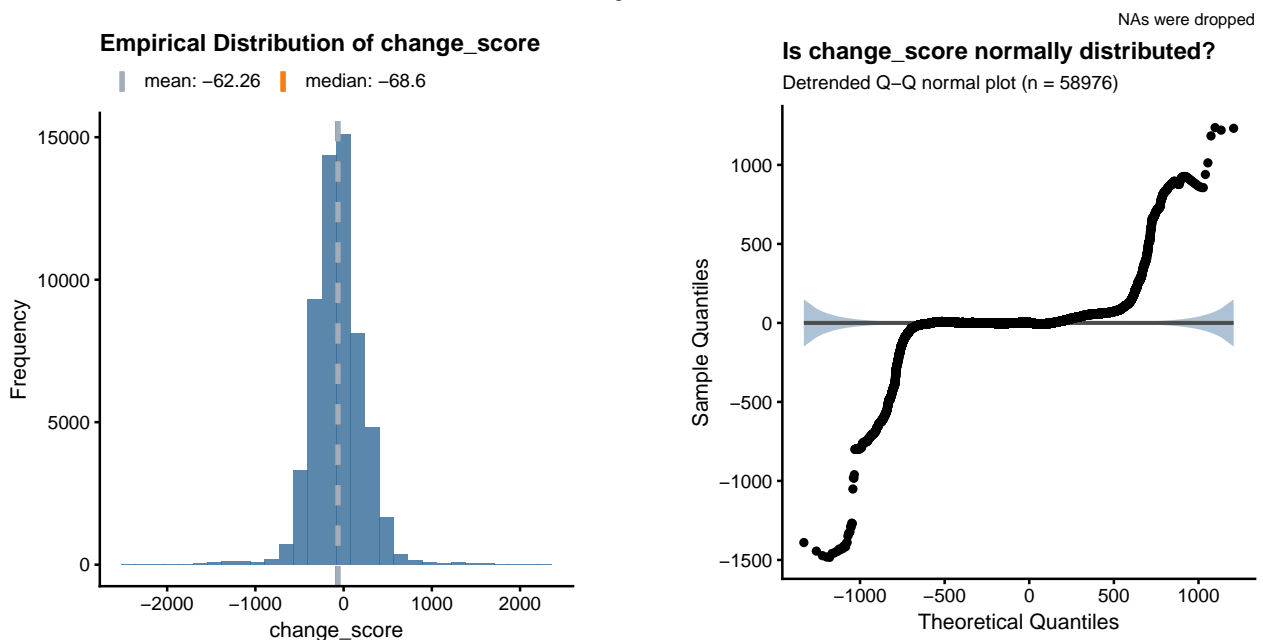
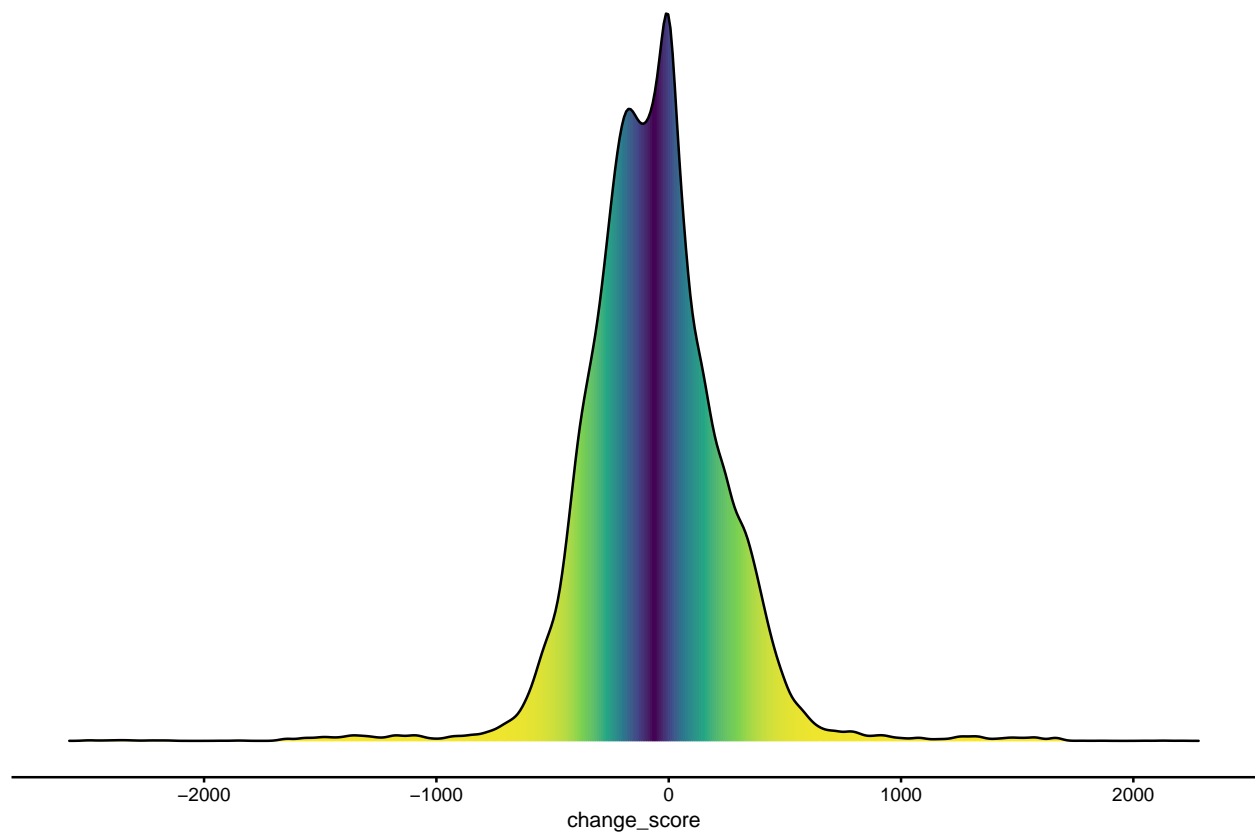
time	n_tribes	n	mean	SEM
t1	184	58976	839.7657	1.434148
t2	184	58976	777.5040	1.338598

### Q1. Is there a difference between t1 and t2?

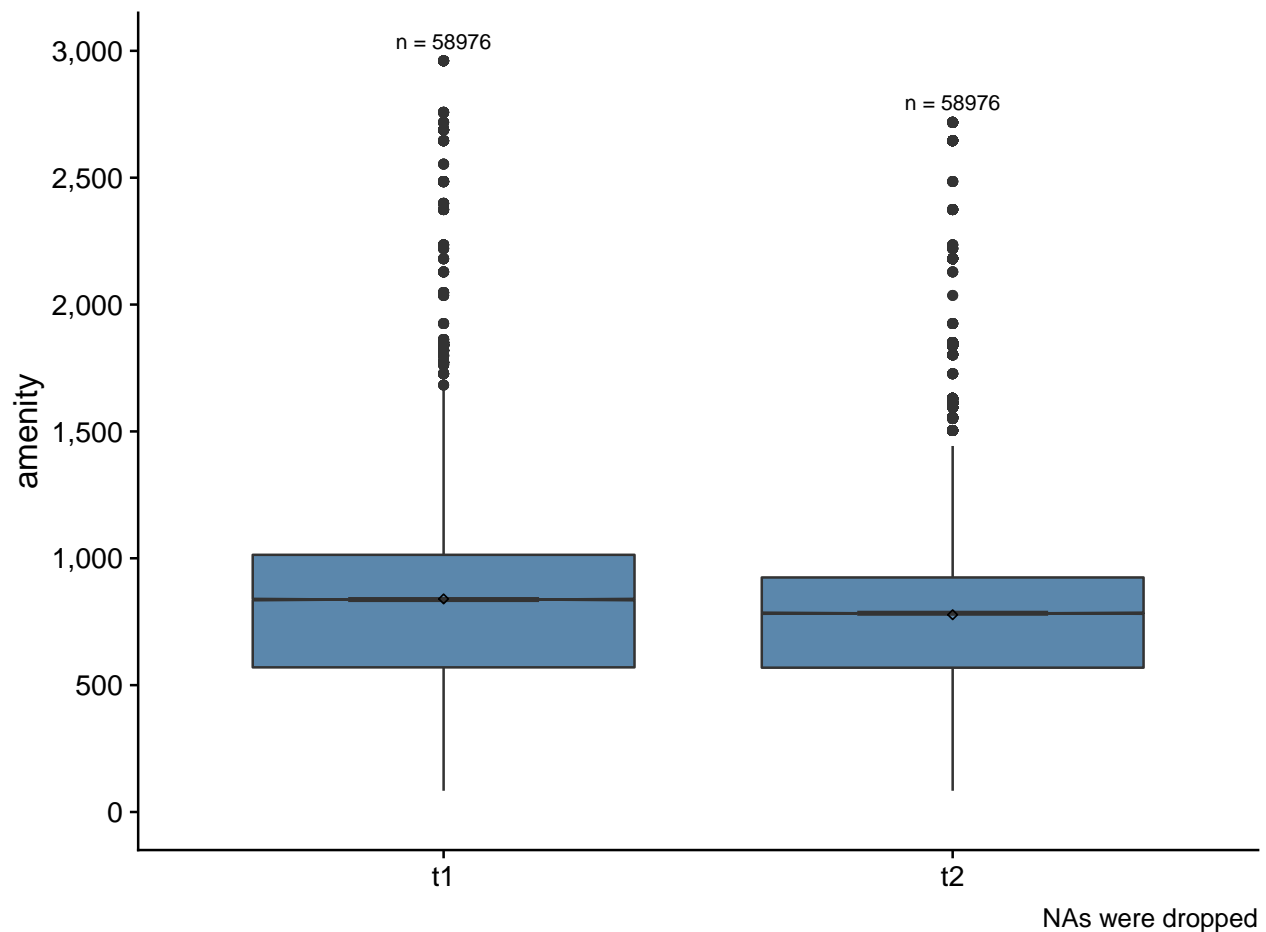
#### Descriptive Analysis

First we look at the distribution of the change scores between t2 and t1 ( $t2 - t1$ ). The following density plot, histogram and detrended qqnormal plot show that its distribution is symmetric but not normal, with long tails extending in both the positive and negative directions. A normal distribution would produce a detrended qqnormal plot with most of the data points randomly scattered around the line  $y=0$  and within the grayish blue confidence band. The left tail is slightly longer than the right tail. Both the sample mean and median are negative. It is not immediately clear from these plots that if the population mean amenity change from t1 to t2 is significantly different from zero.

amenity



Next we look at the distributions of amenity values at t1 and t2 side by side with notched boxplots. We see the two box bodies overlap, with t2 mean slightly below t1 mean. We also observe there are many outliers at t1 or t2, and these outliers inflate the sample averages. We cannot tell from these boxplots if there's a significant difference in the population mean amenity values between t2 and t1.



## Statistical Analysis

To test if there's a difference between the population mean amenity values at t1 and t2, we ran an one-way repeated measure ANOVA, as well as a linear mixed model with tribe as random effect. We chose these methods because at each time point (t1 or t2), there are multiple amenity values for each tribe. These methods account for the within-tribe correlations. According to the ANOVA output, **time** is significant with a tiny p-value of 0.0000000108. According to the linear mixed model output, **time** is significant with a p-value  $< .0001$ . Both methods show consistent results. So we conclude there's a significant difference between the mean amenity values at t2 and t1.

When reading the output from ANOVA and linear mixed model, you want to focus on the reported p-value of the term **time**. It tells you the probability of observing a difference between t1 and t2 as extreme as in the sample data due to chance or randomness. If it's small, it's more likely that the observed difference is not due to chance. To decide how small is "small", the convention is to compare the p-value with 0.05. (But you don't have to use 0.05, it's really your choice. For example, 0.01 or 0.1 are also commonly used in different applications). If it's less than 0.05, we say the observed difference between t1 and t2 in the sample is likely not due to chance and hence can be generalized to the entire population. In other words, the difference is (statistically) significant. Otherwise when the p-value is greater than 0.05, we say the observed difference between t1 and t2 in the sample is likely a fluke and cannot be generalized to the entire population. In other words, the difference is not (statistically) significant.

1-way Repeated Measure ANOVA Output:

Error: tribe

```

      Df      Sum Sq  Mean Sq F value Pr(>F)
Residuals 183 8830587959 48254579

Error: tribe:time
      Df      Sum Sq  Mean Sq F value      Pr(>F)
time      1 114310567 114310567   35.89 0.0000000108 ***
Residuals 183 582798938   3184694
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Error: Within
      Df      Sum Sq  Mean Sq F value Pr(>F)
Residuals 117584 3972555800   33785

Linear Mixed Model Output:
      numDF  denDF  F-value p-value
(Intercept)      1 117767  495.1337 <.0001
time              1 117767 2955.1897 <.0001

```

## Q2. How are t1 and t2 related?

We first calculated the mean and median amenity values of each tribe at t1 and t2. The reason why we also looked at the median is because the median is extreme-value resistant while the mean is heavily influenced by outliers. We then made a scatterplot of the t2 means vs. t1 means, and another scatterplot of the t2 medians vs. t1 medians. These scatterplots showed strong positive linear relationships:

- **strong:** the tighter the dots, the stronger the correlation.
- **positive:** upward slanted trend from bottom left corner to upper right corner. Or y tends to increase as x increases.

Finally, we ran linear regressions to quantify these relationships. For the mean values, we obtained a r-squared value of 0.89, which translates to a correlation of 0.94 (the squared root of 0.89). This big positive correlation indicates the mean values at t1 and t2 tend to rise and fall together, and they are almost perfectly correlated (perfect positive correlation is 1). The slope of the line is 1.04, meaning that for every 100-unit increase in precipitation at t1, we can expect an increase of 104-unit at t2. And this is statistically significant by the tiny p-value. A similar interpretation can be done for the median values.

