

# The Two-Sample Bootstrap

*Math 445, Spring 2017*

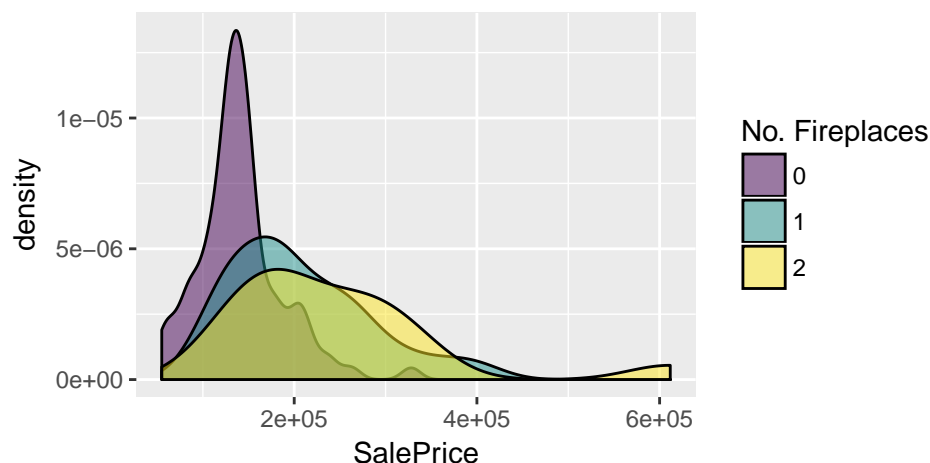
## Motivating example

Are houses with fireplaces more expensive? The file `AmesHousing2010.csv` contains a random sample of houses sold in Ames, IA during 2010. The data set contains details of the sales and the property. Suppose that you are interested in estimating the difference in average price between houses with and without a fireplace. This is a situation where the **two sample bootstrap** can be used to construct a confidence interval.

```
housing <- read.csv("../data/AmesHousing2010.csv")
```

First, let's see what the data look like

```
library(ggplot2)
library(viridis)
ggplot(data = housing, mapping = aes(x = SalePrice, fill = factor(Fireplaces))) +
  geom_density(alpha = 0.5) +
  scale_fill_viridis("No. Fireplaces", discrete = TRUE)
```



Suppose further that you are only interested in houses with either no fireplace or one fireplace

```
library(dplyr)
sub_housing <- filter(housing, Fireplaces <= 1)

# Quick summary stats
price_stats <-
  sub_housing %>%
  group_by(Fireplaces) %>%
  summarise(min = min(SalePrice),
            Q1 = quantile(SalePrice, .25),
            median = median(SalePrice),
            Q3 = quantile(SalePrice, .75),
            max = max(SalePrice),
            mean = mean(SalePrice),
            sd = sd(SalePrice),
            n = n())
price_stats
```

```
## # A tibble: 2 × 9
##   Fireplaces  min      Q1 median      Q3    max    mean      sd    n
##   <int> <int> <dbl> <dbl> <dbl> <int> <dbl> <dbl> <int>
## 1         0 55000 117750 136500 155000 328000 140676.4 46175.91  87
## 2         1 99500 144875 189000 251675 410000 205717.8 75590.35  62
```

So we see that our **point estimate** is

```
obs_diff_means <- price_stats$mean[2] - price_stats$mean[1]
obs_diff_means

## [1] 65041.4
```

## Algorithm

Given independent samples of sizes  $m$  and  $n$  from two populations,

1. Draw a resample of size  $m$  with replacement from the first sample and a separate resample of size  $n$  from the second sample. Compute a statistic that compares the two groups, such as the difference between the two sample means.
2. Repeat this resampling process many times, say 10,000.
3. Construct the bootstrap distribution of the statistic. Inspect its spread, bias, and shape.

## Example implementation

The below code snippet performs the two sample bootstrap in order to construct a confidence interval for the difference in average sales price of homes in Ames, IA with and without fireplaces.

Since we have many extra columns in our data set, it's easier to first select only the two columns of interest.

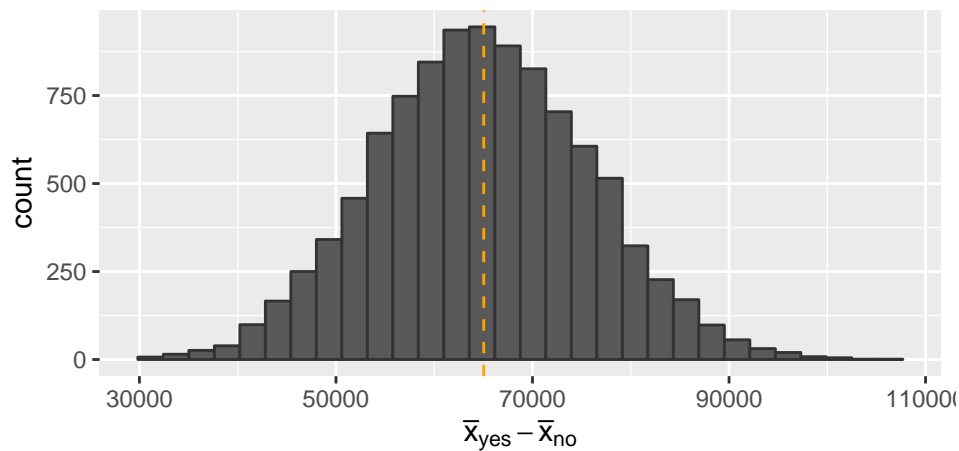
```
boot_df <- select(sub_housing, SalePrice, Fireplaces)

N <- 10^4

no_fp <- subset(boot_df, select = SalePrice, subset = Fireplaces == 0, drop = TRUE)
fp <- subset(boot_df, select = SalePrice, subset = Fireplaces == 1, drop = TRUE)

price_diff_mean <- numeric(N)
for (i in 1:N) {
  no_fp_sample <- sample(no_fp, replace = TRUE)
  fp_sample <- sample(fp, replace = TRUE)
  price_diff_mean[i] <- mean(fp_sample) - mean(no_fp_sample)
}
```

## Bootstrap distribution



```
# bootstrap mean
mean(price_diff_mean)

## [1] 64896.79

# bias
mean(price_diff_mean) - obs_diff_means

## [1] -144.6105

# standard error
sd(price_diff_mean)

## [1] 10855.58

# bias/se
(mean(price_diff_mean) - obs_diff_means) / sd(price_diff_mean)

## [1] -0.01332131
```

### Bootstrap percentile confidence intervals

```
quantile(price_diff_mean, probs = c(.05, .95))

##      5%      95%
## 47071.32 82947.21
```

### Other statistics

Instead of focusing on the difference in means, we could instead focus on the ratio of means.

### Implementation

```
N <- 10^4

no_fp <- subset(boot_df, select = SalePrice, subset = Fireplaces == 0, drop = TRUE)
fp <- subset(boot_df, select = SalePrice, subset = Fireplaces == 1, drop = TRUE)

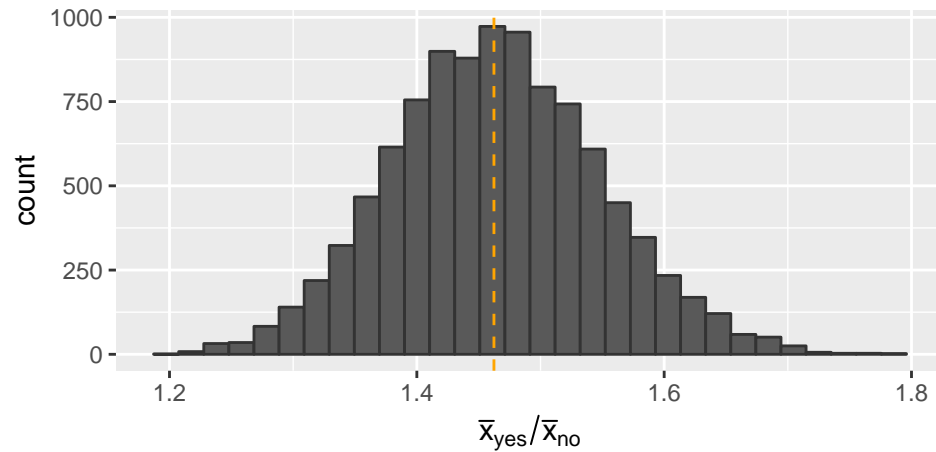
price_ratio_mean <- numeric(N)
for (i in 1:N) {
  no_fp_sample <- sample(no_fp, replace = TRUE)
```

```

fp_sample <- sample(fp, replace = TRUE)
price_ratio_mean[i] <- mean(fp_sample) / mean(no_fp_sample)
}

```

### Bootstrap distribution



```

# bootstrap mean
mean(obs_ratio_mean)

```

```
## [1] 1.462348
```

```

# bias
mean(price_ratio_mean) - obs_ratio_mean

```

```
## [1] 0.001524396
```

```

# standard error
sd(price_ratio_mean)

```

```
## [1] 0.08445088
```

```

# bias/se
(mean(price_ratio_mean) - obs_ratio_mean) / sd(price_ratio_mean)

```

```
## [1] 0.01805069
```

### Bootstrap percentile confidence intervals

```
quantile(price_ratio_mean, probs = c(.05, .95))
```

```
##          5%          95%
## 1.327712 1.607014
```