

The bootstrap for sampling distributions

Assessing assumptions

- Our t -tests assume normality of variable being tested
- but, Central Limit Theorem says that normality matters less if sample is “large”
- in practice “approximate normality” is enough, but how do we assess whether what we have is normal enough?
- so far, use histogram/boxplot and make a call, allowing for sample size.

What actually has to be normal

- is: **sampling distribution of sample mean**
- the distribution of sample mean over *all possible samples*
- but we only have *one* sample!
- Idea: assume our sample is representative of the population, and draw samples from our sample (!), with replacement.
- This gives an idea of what different samples from the population might look like.
- Called *bootstrap*, after expression “to pull yourself up by your own bootstraps”.

Blue Jays attendances

```
jays$attendance
```

```
## [1] 48414 17264 15086 14433 21397 34743 44794 14184  
## [9] 15606 18581 19217 21519 21312 30430 42917 42419  
## [17] 29306 15062 16402 19014 21195 33086 37929 15168  
## [25] 17276
```

- A bootstrap sample:

```
s <- sample(jays$attendance, replace = TRUE)  
s
```

```
## [1] 21195 34743 21312 44794 16402 19014 34743 21195  
## [9] 17264 18581 19014 19217 34743 19217 14433 15062  
## [17] 16402 15062 34743 15062 15086 15168 15086 48414  
## [25] 30430
```

Getting mean of bootstrap sample

- A bootstrap sample is same size as original, but contains repeated values (eg. 15062) and missing ones (42917).
- We need the mean of our bootstrap sample:

```
mean(s)
```

```
## [1] 23055.28
```

- This is a little different from the mean of our actual sample:

```
mean(jays$attendance)
```

```
## [1] 25070.16
```

- Want a sense of how the sample mean might vary, if we were able to take repeated samples from our population.
- Idea: take lots of *bootstrap* samples, and see how *their* sample means vary.

Taking lots of bootstrap samples

- This is the same idea as simulating power, using rowwise:
 - set up dataframe with column `sim` to label the simulations
 - generate a bootstrap sample from the data for each `sim`
 - work out the mean of each sample
 - (then) plot them.

```
tibble(sim = 1:1000) %>%  
  rowwise() %>%  
  mutate(boot_sample =  
    list(sample(jays$attendance, replace = TRUE))) %>%  
  mutate(mean = mean(boot_sample)) -> boots
```

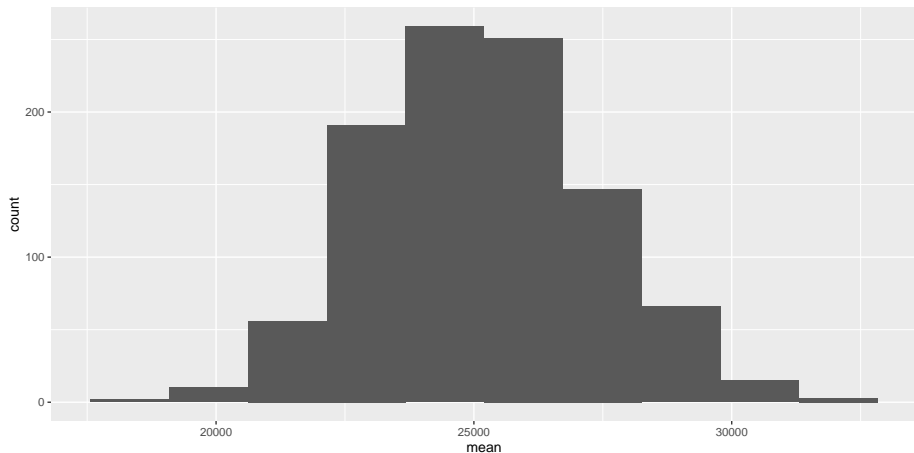
The results

```
boots
```

```
## # A tibble: 1,000 x 3
##       sim boot_sample mean
##   <int> <list>      <dbl>
## 1     1 1 <dbl [25]> 23055.
## 2     2 2 <dbl [25]> 25513.
## 3     3 3 <dbl [25]> 25563.
## 4     4 4 <dbl [25]> 29198.
## 5     5 5 <dbl [25]> 23615.
## 6     6 6 <dbl [25]> 28472.
## 7     7 7 <dbl [25]> 28648.
## 8     8 8 <dbl [25]> 23329.
## 9     9 9 <dbl [25]> 24808.
## 10    10 10 <dbl [25]> 24665.
## # ... with 990 more rows
```

Are these normal?

```
ggplot(boots, aes(x=mean)) + geom_histogram(bins=10)
```



Comments

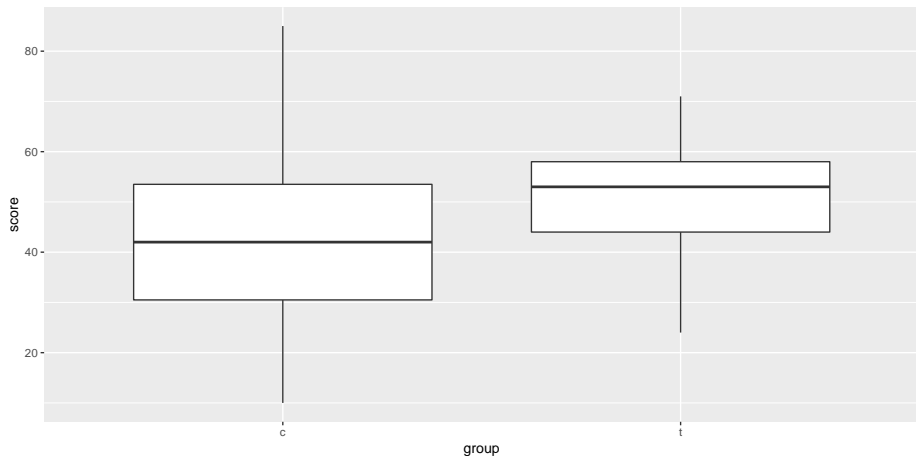
- This is very close to normal
- The bootstrap says that the sampling distribution of the sample mean is close to normal, even though the distribution of the data is not
- A sample size of 25 is big enough to overcome the skewness that we saw
- This is the Central Limit Theorem in practice
- It is surprisingly powerful.
- Thus, the t -test is actually perfectly good here.

Two samples

- Assumption: *both* samples are from a normal distribution.
- In practice, each sample is “normal enough” given its sample size, since Central Limit Theorem will help.
- Use bootstrap on each group independently, as above.

Kids learning to read

```
ggplot(kids, aes(x=group, y=score)) + geom_boxplot()
```



Getting just the control group

```
kids %>% filter(group=="c") -> controls
controls
```

```
## # A tibble: 23 x 2
```

```
##   group score
```

```
##   <chr> <dbl>
```

```
## 1 c      42
```

```
## 2 c      33
```

```
## 3 c      46
```

```
## 4 c      37
```

```
## 5 c      43
```

```
## 6 c      41
```

```
## 7 c      10
```

```
## 8 c      42
```

```
## 9 c      55
```

```
## 10 c     19
```

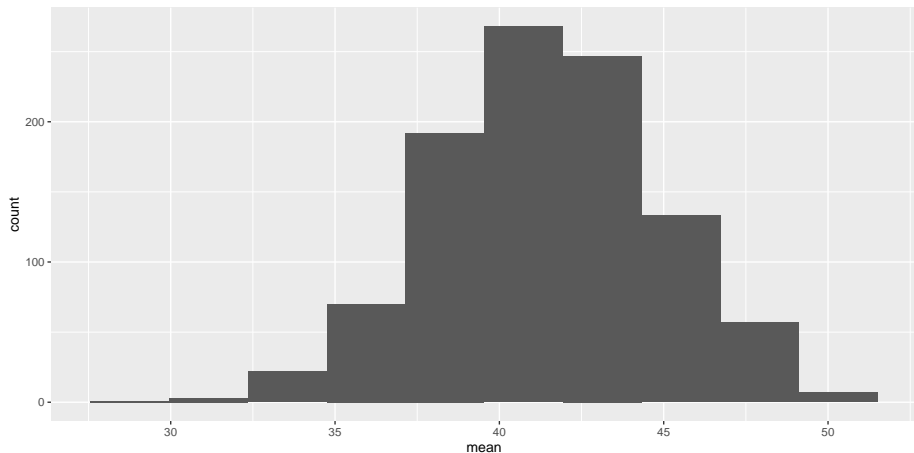
```
## # with 13 more rows
```

Bootstrap these

```
tibble(sim = 1:1000) %>%  
  rowwise() %>%  
  mutate(boot =  
    list(sample(controls$score, replace = TRUE))) %>%  
  mutate(mean = mean(boot)) -> boots
```

Plot

```
ggplot(boots, aes(x = mean)) + geom_histogram(bins=10)
```

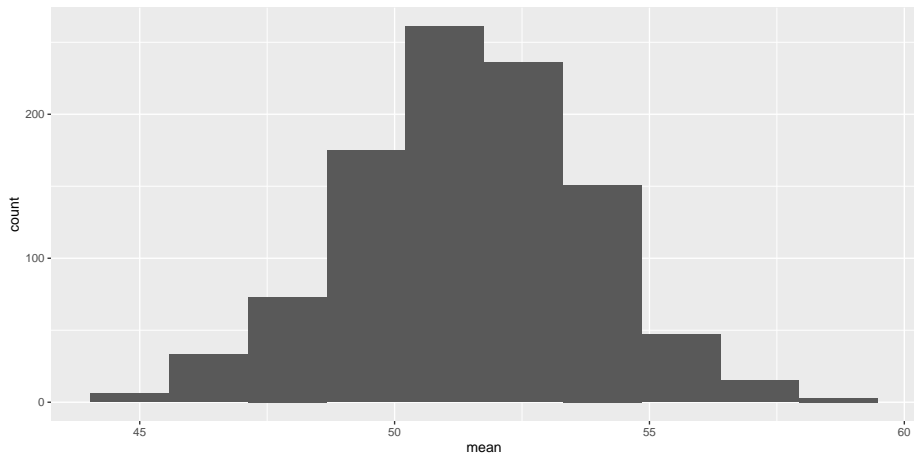


... and the treatment group:

```
kids %>% filter(group=="t") -> treats
tibble(sim = 1:1000) %>%
  rowwise() %>%
  mutate(boot =
    list(sample(treats$score, replace = TRUE))) %>%
  mutate(mean = mean(boot)) -> boots
```

Histogram

```
ggplot(boots, aes(x = mean)) + geom_histogram(bins = 10)
```



Comments

- sampling distributions of sample means both look pretty normal
- as we thought, no problems with our two-sample t at all.