

```

1   Course: Statistical_Inference
2   Lesson: Resampling
3
4   - Class: text
5   Output: "Resampling. (Slides for this and other Data Science courses may be found at
6   github https://github.com/DataScienceSpecialization/courses/. If you care to use
7   them, they must be downloaded as a zip file and viewed locally. This lesson
8   corresponds to 06_Statistical_Inference/13_Resampling.)"
9
10  - Class: text
11  Output: In this lesson, you get a bonus! We'll talk about two topics in statistical
12  inference, bootstrapping AND permutation testing. These both fall under the broader
13  category of resampling methods. We'll start with bootstrapping.
14
15  - Class: text
16  Output: The bootstrap is a handy tool for making statistical inferences. It is used
17  in constructing confidence intervals and calculating standard errors for statistics
18  that might be difficult for some reason (e.g., lack of data or no closed form).
19  Wikipedia tells us that bootstrapping is a technique which "allows estimation of the
20  sampling distribution of almost any statistic using very simple methods." Simple is
21  good, right?
22
23  - Class: text
24  Output: The beauty of bootstrapping is that it avoids complicated mathematics and
25  instead uses simulation and computation to infer distributional properties you might
26  not otherwise be able to determine.
27
28  - Class: text
29  Output: It's relatively new, developed in 1979, by Bradley Efron, a Stanford
30  statistician. The basic bootstrap principle uses OBSERVED data to construct an
31  ESTIMATED population distribution using random sampling with replacement. From this
32  distribution (constructed from the observed data) we can estimate the distribution of
33  the statistic we're interested in.
34
35  - Class: mult_question
36  Output: So, in bootstrapping the observed data substitutes for what?
37  AnswerChoices: a population; a statistic; a hypothesis; observations
38  CorrectAnswer: a population
39  AnswerTests: omnitest(correctVal='a population')
40  Hint: We'll used the observed data and sample it with replacement, just as we would
41  do to find out information about some population.
42
43  - Class: mult_question
44  Output: So, in bootstrapping if the observed data is the population, what would the
45  random samplings correspond to?
46  AnswerChoices: a population; a statistic; a hypothesis; observations
47  CorrectAnswer: observations
48  AnswerTests: omnitest(correctVal='observations')
49  Hint: Sampling from a population would give us observations, right?
50
51  - Class: text
52  Output: In effect, the original observed sample substitutes for the population. Our
53  samplings become observations from which we estimate a statistic and get an idea
54  about its distribution. This lets us better understand the underlying population
55  (from which we didn't have enough data).
56
57  - Class: mult_question
58  Output: Here's a critical point. In constructing the estimated distribution we sample
59  the observed data WITH replacement. If the original sample is n long and we sampled n
60  times without replacement what would we get?
61  AnswerChoices: the original sample permuted; an entirely new sample; a better sample;
62  a worse sample
63  CorrectAnswer: the original sample permuted
64  AnswerTests: omnitest(correctVal='the original sample permuted')
65  Hint: Sampling without replacement n times permutes the original sample.
66
67  - Class: text
68  Output: The motivating example from the slides involves computing the average of 50
69  rolls of a die. Of course we can do this theoretically when we know that the die is

```

fair. Remember, $E(x) = \text{Sum}(x \cdot p(x))$ for $x=1,2,\dots,6$, and $p(x)=1/6$ for all values of x .

```
45
46 - Class: cmd_question
47 Output: For the heck of it, compute the expected die roll for a fair die.
48 CorrectAnswer: 3.5
49 AnswerTests: equiv_val(3.5)
50 Hint: Type sum(1:6)/6 at the command prompt.
51
52 - Class: figure
53 Output: Theoretically, the average is 3.5. Here, we've run code and plotted a
    histogram after we took 1000 such averages, each of 50 dice rolls. Note the unusual
    y-axis scale. We're displaying this as a density function so the area of the
    salmon-colored region is theoretically 1. With this scale, though, all the heights of
    the bins actually add up to 5. So you have to multiply each height by .2 and add up
    all the results to get 1.
54 Figure: plot1.R
55 FigureType: new
56
57 - Class: text
58 Output: The point is, the empirical matches the theoretical. Yay! The highest bin is
    centered at 3.5 just as the math predicted. So what?
59
60 - Class: figure
61 Output: What if some joker wanted you to run the same experiment with a die he gave
    you and he warned you that the dice was loaded? In other words, it wasn't fair. It
    has some random distribution like this.
62 Figure: plot2.R
63 FigureType: new
64
65 - Class: text
66 Output: The outcomes aren't equally likely, are they? So when you do your 1000 runs
    of 50 rolls each, the density of the means looks different.
67
68 - Class: cmd_question
69 Output: We've done this for you and put the result in g2. Type print(g2) now to see
    the picture.
70 CorrectAnswer: print(g2)
71 AnswerTests: omnittest(correctExpr='print(g2)')
72 Hint: Type print(g2) at the command prompt.
73
74 - Class: text
75 Output: Picture's a little different, right? Although this example is a bit
    contrived, it illustrates an important concept. We really want a distribution of
    means and we have only one set of observations. (In this case it was the empirical
    distribution associated with the unfair die - the big blue picture.) We used that one
    distribution, to "create" many (1000) distributions by sampling with replacement from
    the given one. We sampled 50000 times so we created 1000 distributions of 50 rolls
    each.
76
77 - Class: text
78 Output: We then calculated the mean of each of our created distributions and got a
    distribution of means. Sampling the one distribution many times gives us some
    variability in the resulting statistics we calculate. We can then calculate the
    standard error and confidence intervals associated with the statistic.
79
80 - Class: mult_question
81 Output: Before we go on to more theory, here's another example in which we'll try to
    find a distribution of medians of a population. Do you recall what a median is?
82 AnswerChoices: 50th percentile; a point halfway between rare and well-done; the most
    frequent outcome; a person who talks to spirits
83 CorrectAnswer: 50th percentile
84 AnswerTests: omnittest(correctVal='50th percentile')
85 Hint: Half the outcomes are above it and half below it.
86
87 - Class: cmd_question
88 Output: Recall the father and son height data. Once again, we've loaded it for you.
    We've placed the height of the sons in the vector sh and the length of this vector
    is stored in the variable nh. Use the R command head to look at the first few entries
    of sh.
```

```

89   CorrectAnswer: head(sh)
90   AnswerTests: omnitest(correctExpr='head(sh)')
91   Hint: Type head(sh) at the command prompt.
92
93   - Class: cmd_question
94   Output: Now look at nh to see how long sh is.
95   CorrectAnswer: nh
96   AnswerTests: omnitest(correctExpr='nh')
97   Hint: Type nh at the command prompt.
98
99   - Class: text
100  Output: Now we'll create 1000 distributions of the same length as the original sh.
    We'll do this by sampling sh with replacement 1000*nh times and store the results in
    an array with 1000 rows, each with nh entries. Then we'll take the median of each row
    and plot the result.
101
102  - Class: text
103  Output: Note that every time we draw from the empirical distribution sh, each of its
    nh data points is equally likely to be pulled, therefore the probability of drawing
    any one is 1/nh. The 1000 samples we create will vary from the original.
104
105  - Class: figure
106  Output: Here's the resulting density curve. This estimates the distribution of
    medians. The thick vertical line shows where the median of the original, observed
    data sh lies.
107  Figure: fatherson.R
108  FigureType: new
109
110  - Class: cmd_question
111  Output: We stored the 1000 medians of the resampled sets in the vector
    resampledMedians. Use the R function median to compute the median of numbers in this
    vector.
112  CorrectAnswer: median(resampledMedians)
113  AnswerTests: obliterate("resamples"); omnitest(correctExpr='median(resampledMedians)')
114  Hint: Type median(resampledMedians) at the command prompt.
115
116  - Class: cmd_question
117  Output: Now compute the median of the original sample sh.
118  CorrectAnswer: median(sh)
119  AnswerTests: omnitest(correctExpr='median(sh)')
120  Hint: Type median(sh) at the command prompt.
121
122
123  - Class: text
124  Output: Pretty close, right? Now back to theory. Suppose you have a statistic that
    estimates some population parameter, but you don't know its sampling distribution.
    The bootstrap principle uses the distribution defined by the observed data to
    approximate the sampling distribution of that statistic.
125
126  - Class: text
127  Output: The nice thing about bootstrapping is that you can always do it with
    simulation. The general procedure follows by first simulating B complete data sets
    from the observed data by sampling with replacement. Make sure B is large and that
    you're sampling WITH replacement to create data sets the same size as the original.
128
129  - Class: text
130  Output: This approximates drawing from the sampling distribution of that statistic,
    at least as well as the data approximates the true population distribution. By
    calculating the statistic for each simulated data set and using these simulated
    statistics we can either define a confidence interval (e.g. find the 2.5 and the 97.5
    percentiles) or take the standard deviation to estimate a standard error of that
    statistic.
131
132  - Class: text
133  Output: Notice that this process doesn't use any fancy math or asymptotics. The only
    assumption behind it is that the observed sample is representative of the underlying
    population.
134
135  - Class: text

```

```

136 Output: We've created the vector fh for you which contains the fathers' heights from
the father son data we've been working with. It's the same length as the sons' data
(1078) which is stored in nh. B, the number of bootstraps we want has been set to
1000. We'll do an example now in small steps.

137
138 - Class: cmd_question
139 Output: Our one sample of observed data is in the vector fh. Use the R function
sample to sample fh nh*B times. Set the argument replace to TRUE. Put the result in
the variable sam.
140 CorrectAnswer: sam <- sample(fh,nh*B,replace=TRUE)
141 AnswerTests: expr_creates_var('sam'); omnitest(correctExpr='sam <-
sample(fh,nh*B,replace=TRUE)')
142 Hint: Type sam <- sample(fh,nh*B,replace=TRUE) at the command prompt.
143
144 - Class: cmd_question
145 Output: Now form sam into a matrix with B rows and nh columns. Use the R function
matrix and put the result in resam
146 CorrectAnswer: resam <- matrix(sam,B,nh)
147 AnswerTests: expr_creates_var('resam'); omnitest(correctExpr='resam <-
matrix(sam,B,nh)')
148 Hint: Type resam <- matrix(sam,B,nh) at the command prompt.
149
150 - Class: cmd_question
151 Output: Now use the R function apply to take the median (third argument) of each row
of resam (first argument). Put the result in meds. The second argument, the number 1,
specifies that the application of the function is to the rows of the first argument.
152 CorrectAnswer: meds <- apply(resam,1,median)
153 AnswerTests: obliterate("sam"); expr_creates_var('meds'); omnitest(correctExpr='meds
<- apply(resam,1,median)')
154 Hint: Type meds <- apply(resam,1,median) at the command prompt.
155
156 - Class: cmd_question
157 Output: Now look at the difference between the median of fh and the median of meds.
158 CorrectAnswer: median(meds)-median(fh)
159 AnswerTests: ANY_of_exprs("median(meds)-median(fh)","median(fh)-median(meds)")
160 Hint: Type median(meds)-median(fh) or median(fh)-median(meds) at the command prompt.
161
162 - Class: cmd_question
163 Output: Pretty close, right? Now use the R function sd to estimate the standard
error of the vector meds.
164 CorrectAnswer: sd(meds)
165 AnswerTests: obliterate("resam"); omnitest(correctExpr='sd(meds)')
166 Hint: Type sd(meds) at the command prompt.
167
168 - Class: cmd_question
169 Output: We previously did this same process for the sons' data and stored the
resampled medians in the 1000-long vector resampledMedians. Find the standard error
of resampledMedians.
170 CorrectAnswer: sd(resampledMedians)
171 AnswerTests: omnitest(correctExpr='sd(resampledMedians)')
172 Hint: Type sd(resampledMedians) at the command prompt.
173
174 - Class: cmd_question
175 Output: Now we'll find a 95% confidence interval for the sons' data with the R
function quantile. The first argument is the vector of resampledMedians and the
second is the expression c(.025,.975). Do this now.
176 CorrectAnswer: quantile(resampledMedians,c(.025,.975))
177 AnswerTests: omnitest(correctExpr='quantile(resampledMedians,c(.025,.975))')
178 Hint: Type quantile(resampledMedians,c(.025,.975)) at the command prompt.
179
180 - Class: cmd_question
181 Output: Pretty close quantiles, right? Now do the same thing for the fathers' data.
Recall that it's stored in the vector meds.
182 CorrectAnswer: quantile(meds,c(.025,.975))
183 AnswerTests: omnitest(correctExpr='quantile(meds,c(.025,.975))')
184 Hint: Type quantile(meds,c(.025,.975)) at the command prompt.
185
186 - Class: text
187 Output: Another pair of close quantiles, but notice that these quantiles of the

```

```

188     fathers' medians differ from those of the sons.
189 - Class: text
190 Output: Bootstrapping is a very diverse and complicated topic and we just skimmed the
    surface here. The technique we showed you is nonparametric, that is, it's not based
    on any parameterized family of probability distributions. We used only one set of
    observations that we assumed to be representative of the population.
191
192 - Class: text
193 Output: Finally, the confidence intervals we calculated might not perform very well
    because of biases but the R package bootstrap provides an easy fix for this problem.
194
195 - Class: text
196 Output: Now, to permutation testing, another handy tool used in group comparisons. As
    bootstrapping did, permutation testing samples a single dataset a zillion times and
    calculates a statistic based on these samplings.
197
198 - Class: text
199 Output: Permutation testing, however, is based on the idea of exchangeability of group
    labels. It measures whether or not outcomes are independent of group identity. Our
    zillion samples simply permute group labels associated with outcomes. We'll see an
    example of this.
200
201 - Class: figure
202 Output: Here's a picture from the dataset InsectSprays which contains counts of the
    number of bugs killed by six different sprays.
203 Figure: insectSprays.R
204 FigureType: new
205
206 - Class: text
207 Output: We'll use permutation testing to compare Spray B with Spray C.
208
209 - Class: cmd_question
210 Output: Use the R command dim to find the dimensions of InsectSprays.
211 CorrectAnswer: dim(InsectSprays)
212 AnswerTests: omnitest(correctExpr='dim(InsectSprays)')
213 Hint: Type dim(InsectSprays) at the R prompt.
214
215 - Class: cmd_question
216 Output: Now use the R command names to find what the two columns of InsectSprays
    contain.
217 CorrectAnswer: names(InsectSprays)
218 AnswerTests: omnitest(correctExpr='names(InsectSprays)')
219 Hint: Type names(InsectSprays) at the R prompt.
220
221 - Class: text
222 Output: We'll use permutation testing to compare Spray B with Spray C. We subsetting
    data for these two sprays into a data frame subdata. Moreover, the two data frames
    Bdata and Cdata contain the data for their respective sprays.
223
224 - Class: cmd_question
225 Output: Now use the R command range on Bdata$count to find the minimum and maximum
    counts for Spray B.
226 CorrectAnswer: range(Bdata$count)
227 AnswerTests: omnitest(correctExpr='range(Bdata$count)')
228 Hint: Type range(Bdata$count) at the R prompt.
229
230 - Class: cmd_question
231 Output: The picture makes more sense now, right? Now do the same for Spray C. Its
    data is in Cdata.
232 CorrectAnswer: range(Cdata$count)
233 AnswerTests: omnitest(correctExpr='range(Cdata$count)')
234 Hint: Type range(Cdata$count) at the R prompt.
235
236 - Class: text
237 Output: From the ranges (as well as the picture), the sprays look a lot different.
    We'll test the (obviously false) null hypothesis that their means are the same.
238
239 - Class: cmd_question

```

```

240 Output: To make the analysis easier we've defined two arrays for you, one holding
the counts for sprays B and C. It's call BCcounts. Look at it now.
241 CorrectAnswer: BCcounts
242 AnswerTests: omnitest(correctExpr='BCcounts')
243 Hint: Type BCcounts at the command prompt.
244
245 - Class: cmd_question
246 Output: The second array we've defined holds the spray identification and it's
called group. These two arrays line up with each other, that is, the first 12 entries
of counts are associated with spray B and the last 12 with spray C. Look at group
now.
247 CorrectAnswer: group
248 AnswerTests: omnitest(correctExpr='group')
249 Hint: Type group at the command prompt.
250
251 - Class: cmd_question
252 Output: We've also defined for you a one-line function testStat which takes two
parameters, an array of counts and an array of associated identifiers. It assumes all
the counts come from group B or group C. It subtracts the mean of the counts from
group C from the mean of the counts of group B. Type testStat with no parentheses and
no arguments to see how it's defined.
253 CorrectAnswer: testStat
254 AnswerTests: omnitest(correctExpr='testStat')
255 Hint: Type testStat at the command prompt.
256
257 - Class: cmd_question
258 Output: Now set a variable obs by invoking testStat with the arguments BCcounts and
group and assigning the result to obs.
259 CorrectAnswer: obs <- testStat(BCcounts,group)
260 AnswerTests: expr_creates_var('obs'); omnitest(correctExpr='obs <-
testStat(BCcounts,group)')
261 Hint: Type obs <- testStat(BCcounts,group) at the command prompt.
262
263 - Class: cmd_question
264 Output: Take a peek at obs now.
265 CorrectAnswer: obs
266 AnswerTests: omnitest(correctExpr='obs')
267 Hint: Type obs at the command prompt.
268
269 - Class: cmd_question
270 Output: Pretty big difference, right? You can check this by using mean on
Bdata$count and on Cdata$count and subtracting the latter from the former.
Equivalently, you can just apply mean to Bdata$count-Cdata$count. Do either one now.
271 CorrectAnswer: mean(Bdata$count)-mean(Cdata$count)
272 AnswerTests:
ANY_of_exprs("mean(Bdata$count)-mean(Cdata$count)","mean(Bdata$count-Cdata$count)")
273 Hint: Type mean(Bdata$count)-mean(Cdata$count) at the command prompt.
274
275 - Class: mult_question
276 Output: So, mean(Bdata$count)-mean(Cdata$count) equals mean(Bdata$count-Cdata$count)
because ?
277 AnswerChoices: the data is special; mean is linear; mathemagic;
278 CorrectAnswer: mean is linear
279 AnswerTests: omnitest(correctVal='mean is linear')
280 Hint: Recall the mean or average is linear so the mean of a sum is the sum of the
means.
281
282 - Class: text
283 Output: Now this is where the permutation testing starts to involve resampling. We're
going to test whether or not the particular group association of the counts affects
the difference of the means.
284
285 - Class: text
286 Output: We'll keep the same array of counts, just randomly relabel them, by permuting
the group array. R makes this process very easy. Calling the function sample (which
we've used several times in this lesson) with one argument, an array name, will
simply permute the elements of that array.
287
288

```



```

289 - Class: cmd_question
290 Output: Call sample now on the array group to see what happens.
291 CorrectAnswer: sample(group)
292 AnswerTests: omnitest(correctExpr='sample(group)')
293 Hint: Type sample(group) at the command prompt.
294
295 - Class: text
296 Output: "The labels are all mixed up now. We'll do this permuting of labels and then
we'll recalculate the difference of the means of the two \"new\" (really newly
labelled) groups."
297
298 - Class: cmd_question
299 Output: "We'll relabel and calculate the difference of means 10000 times and store
the differences (of means) in the array perms. Here's what the code looks like perms
<- sapply(1 : 10000, function(i) testStat(BCcounts, sample(group))). Try it now."
300 CorrectAnswer: "perms <- sapply(1 : 10000, function(i) testStat(BCcounts,
sample(group)))"
301 AnswerTests: "expr_creates_var('perms'); omnitest(correctExpr='perms <- sapply(1 :
10000, function(i) testStat(BCcounts, sample(group)))'"
302 Hint: "Type perms <- sapply(1 : 10000, function(i) testStat(BCcounts, sample(group)))
at the command prompt."
303
304 - Class: cmd_question
305 Output: We can take the mean of the virtual array of the boolean expression perms >
obs. Do this now.
306 CorrectAnswer: mean(perms>obs)
307 AnswerTests: omnitest(correctExpr='mean(perms>obs)')
308 Hint: Type mean(perms>obs) at the command prompt.
309
310 - Class: text
311 Output: So on average 0 of the permutations had a difference greater than the
observed. That means we would reject the null hypothesis that the means of the two
sprays were equal.
312
313
314 - Class: figure
315 Output: Here's a histogram of the difference of the means. Looks pretty normal,
right? We can see that the distribution runs roughly between -10 and +10 and it's
centered around 0. The vertical line shows where the observed difference of means was
and we see that it's pretty far away from the distribution of the resampled
permutations. This means that group identification did matter and sprays B and C were
quite different.
316 Figure: insectHisto.R
317 FigureType: new
318
319 - Class: figure
320 Output: Here's the picture of the InsectSprays again. Suppose we run the same
experiment, this time comparing sprays D and E, which look more alike. We've
redefined testStat to look at these sprays and subtract the mean of spray E from the
mean of spray D.
321 Figure: insectSprays2.R
322 FigureType: new
323
324 - Class: cmd_question
325 Output: We've also stored off the D and E data in DEcounts and the group labels in
group. Run testStat now with DEcounts and group.
326 CorrectAnswer: testStat(DEcounts,group)
327 AnswerTests: omnitest(correctExpr='testStat(DEcounts,group)')
328 Hint: Type testStat(DEcounts,group) at the command prompt.
329
330 - Class: cmd_question
331 Output: "We've stored off this value, 1.416667, in the variable obs for you. Now run
the permutation command, with DEcounts. Here it is, perms <- sapply(1 : 10000,
function(i) testStat(DEcounts, sample(group)))"
332 CorrectAnswer: "perms <- sapply(1 : 10000, function(i) testStat(DEcounts,
sample(group)))"
333 AnswerTests: "expr_creates_var('perms'); omnitest(correctExpr='perms <- sapply(1 :
10000, function(i) testStat(DEcounts, sample(group)))'"
334 Hint: "Type perms <- sapply(1 : 10000, function(i) testStat(DEcounts, sample(group)))

```

```
335         at the command prompt."
336 - Class: figure
337 Output: Finally, we can plot the histogram of the distribution of the difference of
the means. We see that with these sprays the observed difference of means (the
vertical line) is closer to the mean of the permuted labels. This indicates that
sprays D and E are quite similar and we fail to reject the null hypothesis that the
means were equal.
338 Figure: insectHisto.R
339 FigureType: new
340
341 - Class: text
342 Output: Congrats! We hope you weren't bugged too much by this lesson and feel like
you've pulled yourself up by your bootstraps.
343
344 - Class: mult_question
345 Output: "Would you like to receive credit for completing this course on
346 Coursera.org?"
347 CorrectAnswer: NULL
348 AnswerChoices: Yes;No
349 AnswerTests: coursera_on_demand()
350 Hint: ""
351
```