

Bootstrapping with real data

The goal of this lab is to load some real data (from an MEG experiment), estimate several biological response parameters, and then use bootstrap to estimate the validity of those parameters. This is a long lab—you may need an extra session, so make sure your Matlab code gets saved if you need to use it next time.

- 1) Download the Matlab file containing the data from the Lab page of the course website, or from <http://terpconnect.umd.edu/~jzsimon/biol708L/lab/megExData.mat> , and load it into Matlab. It has MEG data from a single channel, where a subject repeatedly listened to a 2 second 200 Hz pure tone.
 - a) Once you have loaded the data, use `whos` or the Matlab workspace to see what the new variables are that you have loaded. There should be 6 total. The most important are `NS`, `data`, `fs`, and `time`. `NS` is the number of trials that the stimulus was presented and response recorded. `fs` is the sampling frequency. `time` is a vector of time steps in ms. `data`, the most important, contains the time-series data: the recorded magnetic field at each times step, for each of `NS` trials, with values $\sim 10^{-13}$ Tesla.
 - b) Plot the first 4 trials, e.g. `plot(time, data(:,m))` for 4 different `m`. Do they look similar? They should look more different than similar. (FYI, the stimulus turns on at $t = 0$ ms, which is not the earliest moment of `time`).
 - c) Below, you will need the indices of `time` that correspond to $t = 0$ ms and $t = 400$ ms, since we will only search in that interval. Use

```
t0index = find(diff(sign(time)));
t400index = find(diff(sign(time-400)));
```

What is `time(t0index)`?—it should be near 0. What is `time(t400index)`? Why does this work?
 - d) Calculate the following “plug-in” (i.e. re-usable in bootstrap) estimates:
 - i) The mean response over all trials. Call it `meandata`. Hint: This should be as long as `time`. Plot it; it should have a coherent response for the first few hundred ms after $t = 0$ ms (the “onset response”) and after $t = 2000$ ms (“offset response”).
 - ii) The minimum (most negative) value of the mean response (call it `minval`) and the time at which it occurs (call it `mintime`), between 0 and 400 ms. You must use an automated method, not “by eye”, because we will repeat this measurement hundreds of times in bootstrap. I use

```
[minval,minindex] = min(meandata(t0index:t400index));
mintime = time(minindex+t0index-1)
```

Do the results look agree with your eye’s estimates? For this channel, it should be between 100 and 300 ms (e.g. “near” 200 ms).
 - iii) Repeat (ii) for the maximum response between 0 and 400 ms, called `maxval` and `maxtime`. Do the results look agree with your eye’s estimates? For this channel, it should be near 300 ms.
 - e) Plotting with context:
 - i) Plot `data` and `meandata` overlapping in the same figure. Plot `data` first, in magenta, so that it doesn’t overwhelm the plot, then plot `meandata` in black. You can do this in one line with

```
plot(time, data, 'm', time, meandata, 'k')
```

Clearly, the individual signals have a lot of noise, in addition to signal (seen only in the mean).

- ii) Plot `meandata` and put crosshairs where the minimum and maximum values occur. I use:

```
plot(time, meandata, 'k')
line([0 0]+mintime, 2.5e-13*[1 -1], 'color', 'r')
line([time(1) time(end)], minval*[1 1], 'color', 'r')
line([0 0]+maxtime, 2.5e-13*[1 -1], 'color', 'b')
line([time(1) time(end)], maxval*[1 1], 'color', 'b')
```

Alternate a few times between `axis tight`, to get the big picture, and

```
axis([0 600 -2.5e-13 2.5e-13]),
```

 to see where the action is.

- 2) Prepare the bootstrap. We'll use $B = 500$ bootstrap samplings, and the reshufflings of the data in `p`:

```
B = 500;
p = repmat(1:NS, B, 1);
p = p(reshape(randperm(B*NS), B, NS));
```

Just as the actual data contains `NS` samples, each of which is a time series, each bootstrap sampling will also contain `NS` samples, each of which is a time series. Verify to your own satisfaction that each row of `p` contains a valid index reshuffling. What should be the maximum value in each row, and why? How many rows should there be, and why? How many columns, and why?

- 3) Perform the bootstrap calculations. In a for loop, e.g. with `for b = 1:B` and `end` surrounding the calculations, recalculate the same parameters as in (1), but with the bootstrap samplings. Save the bootstrap values of `mintime`, `minval`, `maxtime` and `maxval` in arrays called `mintimeB`, `minvalB`, `maxtimeB` and `maxvalB`. Recommended techniques:

- a) To keep the time information correct and only shuffle the trial numbers, replace `data` with

```
data(:, p(b, :)).
```

 For example, if in (1) you had an equation of the form

```
meandata = mean(...data...),
```

 replace it with `meandataTmp = mean(...data(:, p(b, :))...)`

The temporary variable `meandataTmp` is explained next.

- b) (This is not crucial, but it keeps things tidy so there will be less chance for mistakes and fewer index gymnastics.) Put all calculations in temporary variables, as in (a), and as in

```
[minvalTmp, minindexTmp] = min(meandataTmp(t0index:t400index));
```

```
mintimeTmp = time(minindexTmp+t0index-1);
```

...

and the after all temporary variables are created, only then put them into the bootstrap vectors:

...

```
mintimeB(b) = mintimeTmp;
```

```
minvalB(b) = minvalTmp;
```

```
maxtimeB(b) = maxtimeTmp;
```

```
maxvalB(b) = maxvalTmp;
```

- 4) Analyze the statistics of `minvalB`, `maxvalB`, `mintimeB` and `maxtimeB`.

- a) What is the mean of `minvalB`? How does it compare to `minval`? What is the standard deviation of `minvalB`? Is the bias large enough to that you should use the mean of `minvalB` instead of `minval`?

(E.g. is $|bias/\sigma| > 0.25$? ask me about bias when you get here). Use `hist()` to see the distribution of values of `minvalB`. Is it relatively smooth (in this case it should be).

- b) repeat (a) but for `mintimeB`. In this case the distribution may *not* be smooth. You may need to fiddle with the number of bins used in `hist()`, since 1) you may have outliers, and 2) you want the width of the bins to be roughly two or three times the smallest time step (so they get filled more nicely). (What is the smallest time step?)
 - c) repeat (a) but for `maxvalB`.
 - d) repeat (b) for `maxtimeB`. Same precautions apply as in (b).
- 5) Create and examine the bootstrap mean time series `meandataB`. Since all of the previous estimates (e.g. `mintime`, `minval`) are single numbers, the bootstrap collection of them is a $1 \times B$ vector. `meandata`, though, is a time series, so its bootstrap collection is a $length(time) \times B$ matrix, and in the bootstrap loop you can fill it up like this: `meandataB(:,b) = meandataTmp;`

- a) Plot the first 4 bootstrap mean time series, e.g. `plot(time, meandataB(:,b))` for different `b`. Do they look similar? They should look more similar than different (unlike 1b).
- b) Compute the mean and standard deviation of the bootstrap mean time series (across bootstrap samplings):

```
meandataBmean = mean(meandataB,2);
```

```
meandataBstd = std(meandataB,1,2);
```

Because of the way we implemented bootstrap, `meandataBmean` and `meandata` should be almost identical. Calculate `sum(abs((meandataBmean - meandata)))` to verify (Why is this so, and why does this show it?)

- c) Plotting with context: Plot `meandata` and `meandataB` overlapping in the same figure. Plot `meandataB` first, in yellow, so that it doesn't overwhelm the plot, then plot `meandata` in black. Plot `meandata + 2*meandataBstd` and `meandata - 2*meandataBstd`, to act as error bars. You can do this in one line with:

```
plot(time, meandataB, 'y', time, meandata, 'k', time, ...
```

```
meandata-2*meandataBstd, 'b', time, meandata+2*meandataBstd, 'b')
```

The bootstrap versions of the response have much less noise than the original version (in `data`) plotted above in (1ei). Also the error bars (in blue) follow the yellow bootstrap cloud well. You may wish to restrict the axes to the important time region so the graph looks less busy:

```
axis([0 600 -2.5e-13 2.5e-13])
```

- d) Add estimates of the minimum and maximum values, with their own error bars, both in amplitude and time, to the same graph:

```
line([0 0]+mintimeBmean, minvalBmean+2*[-minvalBstd minvalBstd], 'color', 'r')
```

```
line(mintimeBmean+2*[-mintimeBstd mintimeBstd], minvalBmean*[1 1], 'color', 'r')
```

```
line([0 0]+maxtimeBmean, maxvalBmean+2*[-maxvalBstd maxvalBstd], 'color', 'r')
```

```
line(maxtimeBmean+2*[-maxtimeBstd maxtimeBstd], maxvalBmean*[1 1], 'color', 'r',)
```

Now, your data should be nicely presentable.