

CS395T
Computational Statistics with
Application to Bioinformatics

Prof. William H. Press
Spring Term, 2011
The University of Texas at Austin

Lecture 12

What is the uncertainty in quantities other than the fitted coefficients:

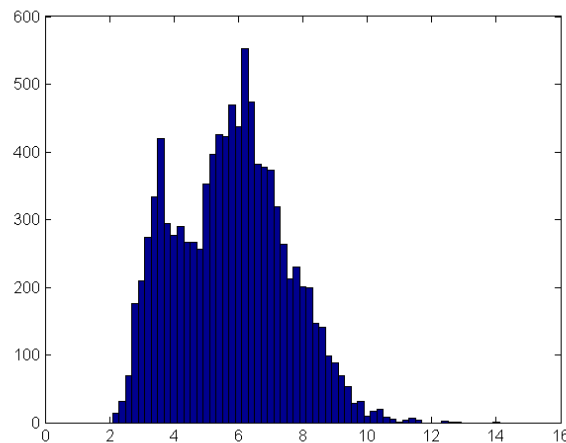
Method 2: Sample from the posterior distribution

1. Generate a large number of (vector) \mathbf{b} 's

$$\mathbf{b} \sim \text{MVNormal}(\mathbf{b}_0, \Sigma_b)$$

2. Compute your $f(\mathbf{b})$ separately for each \mathbf{b}

3. Histogram

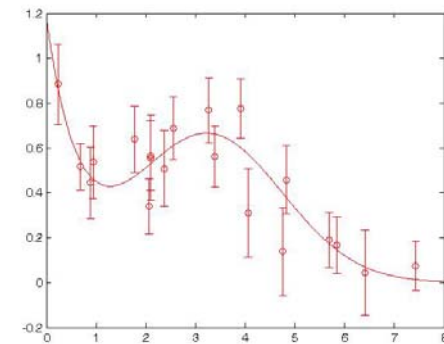
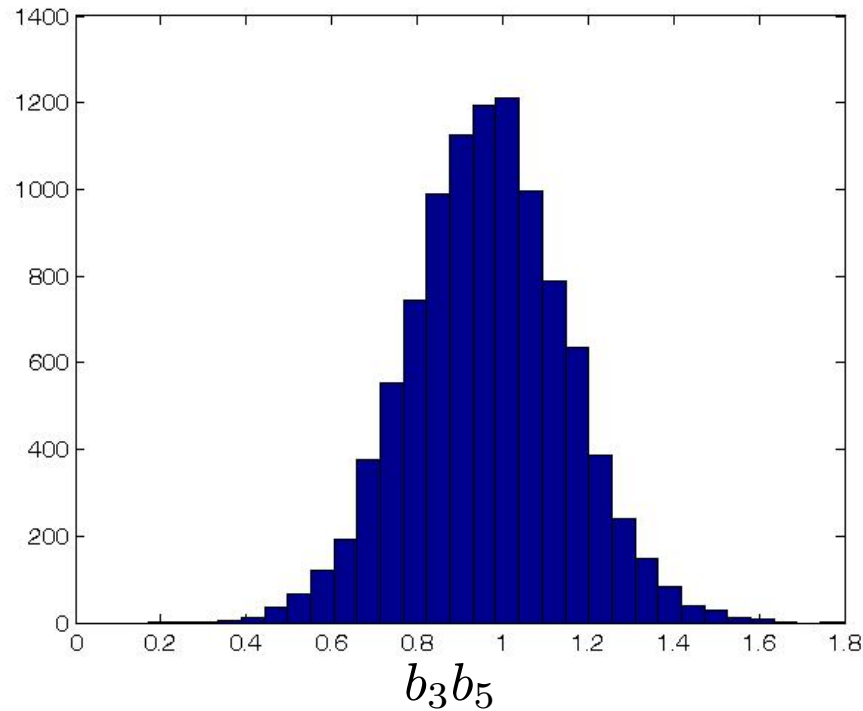


Note again that \mathbf{b} is typically (close to) m.v. normal because of the CLT, but your (nonlinear) f may not, in general, be anything even close to normal!

Our example:

```
bees = mvnrnd(bfit, covar, 10000);  
humps = bees(:, 3) .* bees(:, 5);  
hist(humps, 30);  
std(humps)
```

std = 0.1833



Does it matter that I use the full covar, not just the 2x2 piece for parameters 3 and 5?

Compare linear propagation of errors to sampling the posterior

- Note that even with lots of data, so that the distribution of the b 's really \rightarrow multivariate normal, a derived quantity might be very non-Normal.
 - In this case, sampling the posterior is a good idea!
- For example, the ratio of two normals of zero mean is Cauchy
 - which is very non-Normal!
- So, sampling the posterior is a more powerful method than linear propagation of errors.
 - even when optimistically (or in ignorance) assuming multivariate Gaussian for the fitted parameters
- In fact, sampling the posterior distribution of large Bayesian models whose parameters are not at all Gaussian is, under the name MCMC, the most powerful technique in modern computational statistics.
 - we'll come back to this!

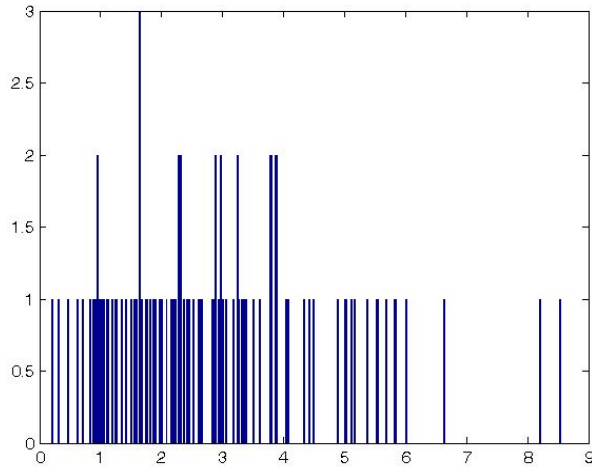
Method 3: Bootstrap resampling of the data

- We applied some end-to-end process to a data set (set of data points) and got a number f out
- The data set was drawn from a population of data points in repetitions of the identical experiment
 - which we don't get to see, unfortunately
 - we see only a sample of the population
- We'd like to draw new data sets from the population, reapply the process, and see the distribution of answers
 - this would tell us how accurate the original answer, on average, was
 - but we can't: we don't have access to the population
- **However, the data set itself is an estimate of the population pdf!**
 - **in fact, it's the only estimate we've got!**
- So we draw from the data set – with replacement – many “fake” data sets of equal size, and carry out the proposed program
 - does this sound crazy? for a long time many people thought so!
 - Bootstrap theorem [glossing over technical assumptions]: **The distribution of any resampled quantity around its full-data-set value estimates (naively: “asymptotically has the same histogram as”) the distribution of the data set value around the population value.**



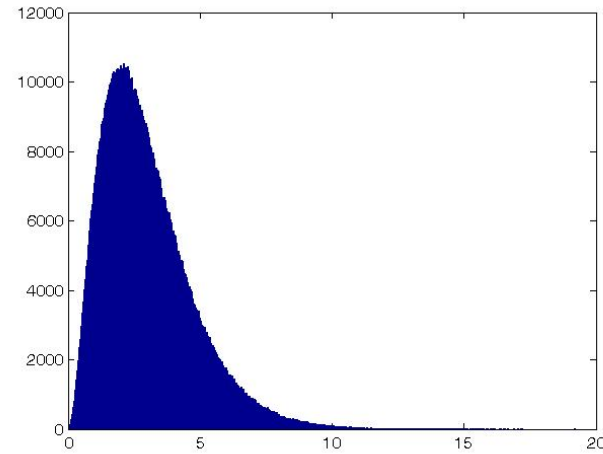
Let's try a simple example where we can see the "hidden" side of things, too.

Visible side (sample):



These happen to be drawn from a Gamma distribution.

Hidden side (population):



Statistic we are interested in happens to be (it could be anything):

$$\frac{\text{mean of distribution}}{\text{median of distribution}}$$

```
sammedian = median(sample)
sammean = mean(sample)
samstatistic = sammean/sammedian
sammedian =
    2.6505
sammean =
    2.9112
samstatistic =
    1.0984
```

How accurate is this?

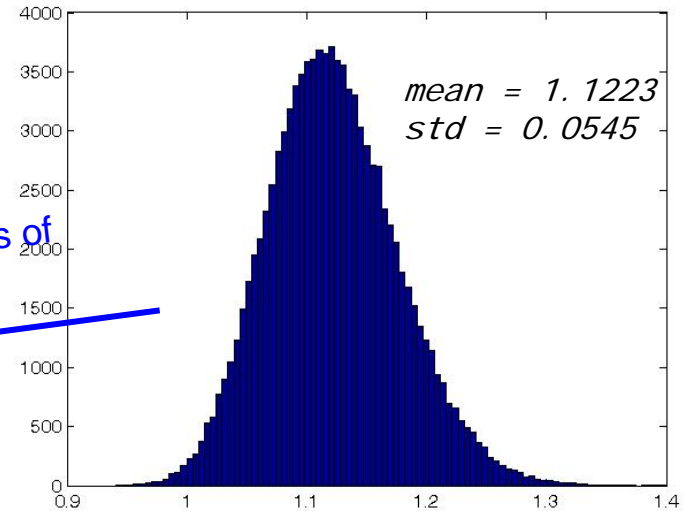
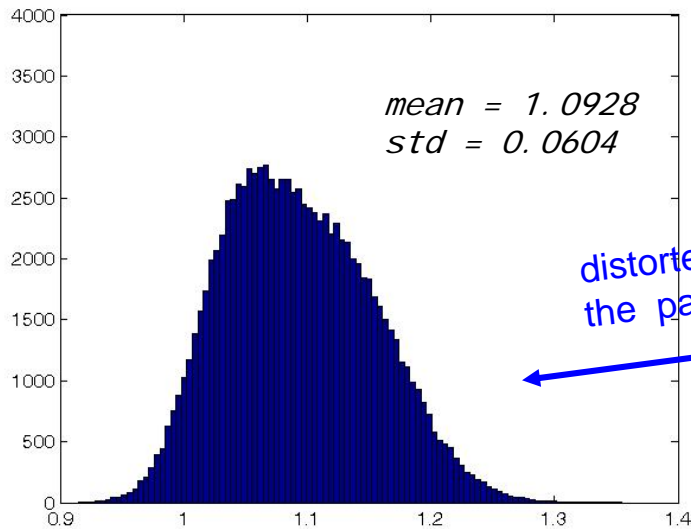
```
themedian = median(bigsample)
themean = mean(bigsample)
thestatistic = themean/themedian
themedian =
    2.6730
themean =
    2.9997
thestatistics =
    1.1222
```

To estimate the accuracy of our statistic, we bootstrap...

```
ndata = 100;
nboot = 100000;
val s = zeros(nboot, 1);
for j=1:nboot,
    choose = randsample(ndata, ndata, true);
    val s(j) = mean(sample(choose))
              /median(sample(choose));
end
hist(val s, 100)
```

new sample of integers in 1:ndata, with replacement

```
ndata = 100;
nboot = 100000;
val s = zeros(nboot, 1);
for j=1:nboot,
    sam = randg(3, [ndata 1]);
    val s(j) = mean(sam)/median(sam);
end
hist(val s, 100)
```



distorted by peculiarities of the particular data set

Things to notice:

The mean of resamplings does not improve the original estimate! (Same data!)

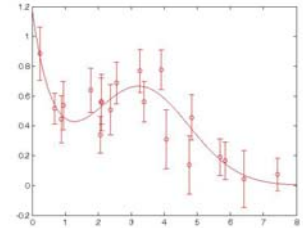
The distribution around the mean is not identical to that of the population. But it is close and would become identical asymptotically for large *ndata* (not *nboot*!).

You often have to customize your own bootstrap code, but it is not a big deal

```

ndata = 20;
nboot = 1000;
vals = zeros(nboot, 1);
ymodel = @(x, b) b(1)*exp(-b(2)*x)+b(3)*exp(-(1/2)*((x-b(4))/b(5)). ^2);
for j=1:nboot,
    samp = randsample(ndata, ndata, true); new sample of integers in 1:ndata, with replacement
    xx = x(samp);
    yy = y(samp);
    ssi g = sig(samp);
    chi sqfun = @(b) sum(((ymodel (xx, b)-yy). /ssi g). ^2);
    bguess = [1 2 .7 3.14 1.5];
    options = optimset(' MaxFunEval s' , 10000, ' MaxIter' ,
        10000, ' Tol Fun' , 0.001);
    [b fval fl ag] = fminsearch(chi sqfun, bguess, options);
    if (fl ag == 1), vals(j) = b(3)*b(5);
    else vals(j) = 100; end
end
hist(vals(vals < 2), 30);
std(vals(vals < 2))

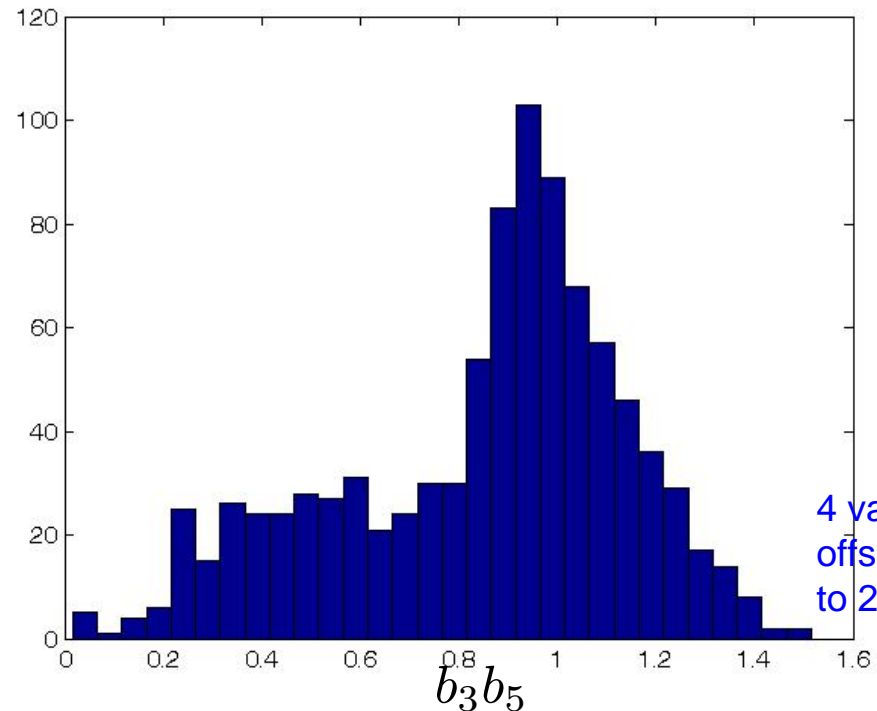
```



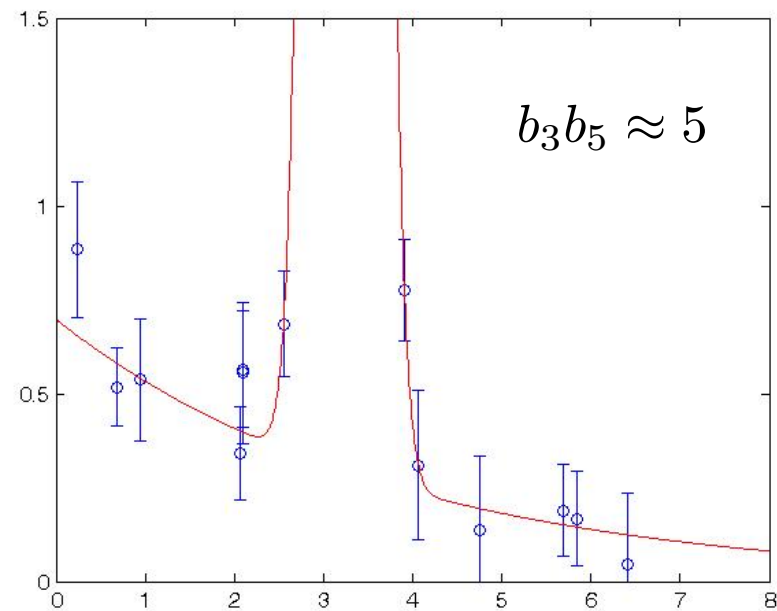
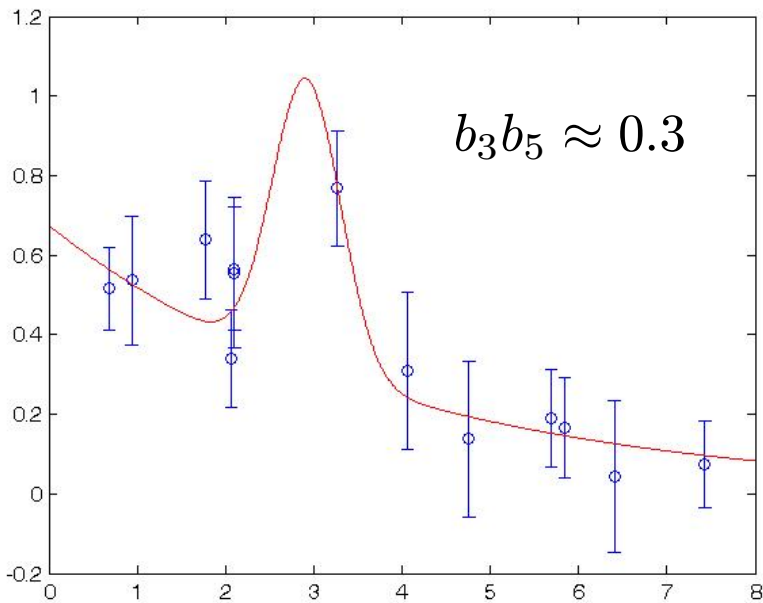
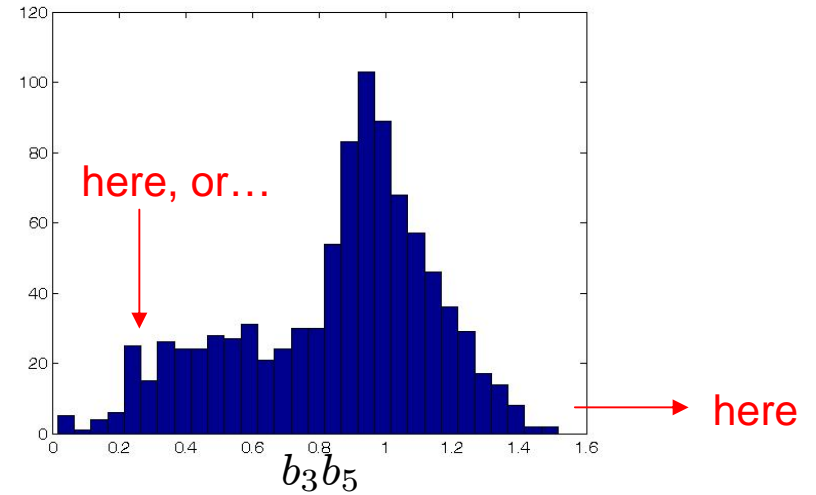
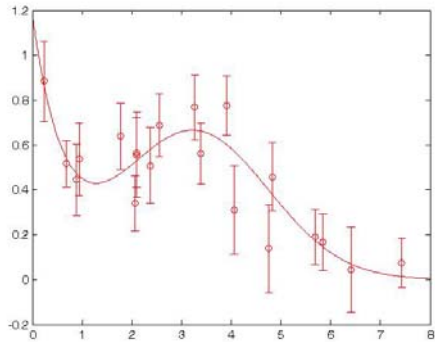
here is the embedded "whole statistical analysis of a data set" inside the bootstrap loop

0.2924

So we get the peak around 1, as before, but a much broader distribution.

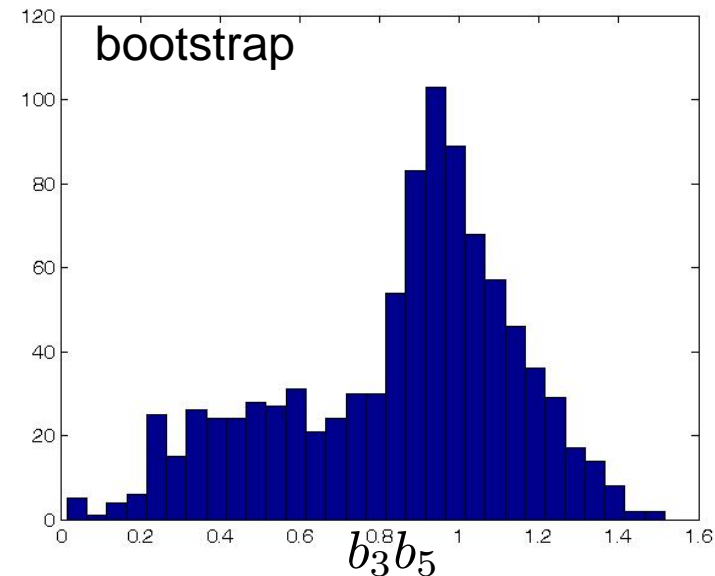
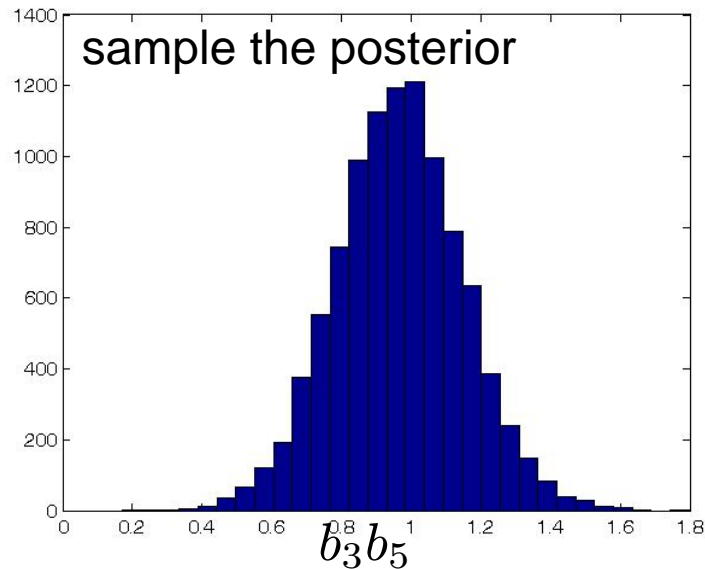


Can you guess what the extreme bootstrap cases look like, compared to the full data?



Is it “fair” for our estimate of b_3b_5 to have its accuracy “impuned” by data sets that “don’t look like” the full data? **Deep frequentist philosophical question!**

We previously compared bootstrap-from-sample to bootstrap-from-population.
More relevant, let's compare bootstrap-from-sample to sample-the-posterior:



- We could increase number of samples of posterior, and of bootstrap, to make both curves very smooth.
 - the histograms would not converge to each other!
- We could increase the size of the underlying data sample
 - from 20 (x,y) values to infinity (x,y) values
 - the histograms would converge to each other (modulo technical assumptions)
- For finite size samples, each technique is a valid answer to a different question
 - Frequentist: Imagining repetitions of the experiment, what would be the range of values obtained?
 - **And, conservatively, I shouldn't expect my experiment to be better than that, should I?**
 - Bayesian: For exactly the data that I see, what is the probability distribution of the parameters?
 - **Because maybe I got lucky and my data set really nails the parameters!**

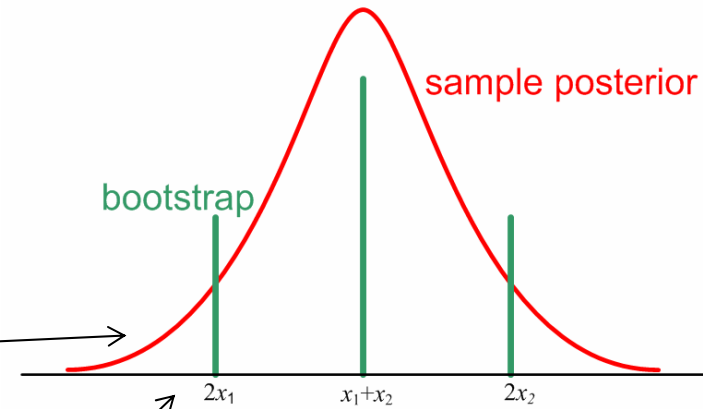
Note that sampling the posterior “honors” the stated measurement errors. Bootstrap doesn’t. That can be good!

Suppose (very toy example) the “statistic” is

$$s = x_1 + x_2$$

then the posterior probability is

$$P(s) \propto \exp \left[-\frac{1}{2} \frac{(s - x_1 - x_2)^2}{\sigma_1^2 + \sigma_2^2} \right]$$

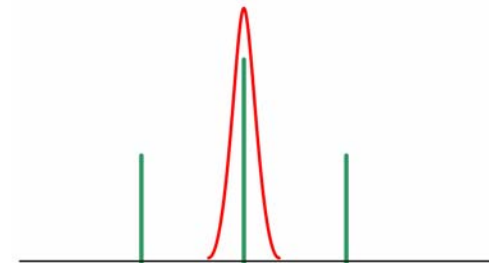


Note that this depends on the σ 's!

The bootstrap (here noticeably discrete) doesn't depend on the σ 's. In some sense it estimates them, too.

So, if the errors were badly underestimated, sampling the posterior would give too small an uncertainty, while bootstrap would still give a valid estimate.

If the errors are right, both estimates are valid. Notice that the model need not be correct. Both procedures give estimates of the statistical uncertainty of parameters of even a wrong (badly fitting) model. *But for a wrong model, your interpretation of the parameters may not mean anything!*



Compare and contrast bootstrap resampling and sampling from the posterior

Both have same goal: Estimate the accuracy of fitted parameters.

- **Bootstrap** is frequentist in outlook
 - draw from “the population”
 - even if we have only an estimate of it (the data set)
- Easy to code but computationally intensive
 - great for getting your bearings
 - must repeat your basic fitting calculation over all the data 100 or 1000 times
- Applies to both model fitting and descriptive statistics
- Fails completely for some statistics
 - e.g. (extreme example) “harmonic mean of distance between consecutive points”
 - how can you be sure that your statistic is OK (without proving theorems)?
- Doesn't generalize much
 - take it or leave it!
- It is not always obvious what you should resample over
 - things that are independent draws from a population
 - “patients not polyps”
- **Sampling from the posterior** is Bayesian in outlook
 - there is only one data set and it is never varied
 - what varies from sample to sample is the goodness of fit of the parameters
 - we don't just sit on the (frequentist's) MLE, we explore around
- In general harder to implement
 - we haven't learned how yet, except in the simple case of an assumed multivariate normal posterior
 - will come back to this later, when we do Markov Chain Monte Carlo (MCMC)
 - may or may not be computationally intensive (depending on whether there are shortcuts possible in computing the posterior)
- Rich set of variations and generalizations are possible