

# STATS305A - Lecture 17

John Duchi  
Scribed by Michael Howes

11/18/21

## Contents

<b>1</b>	<b>The bootstrap</b>	<b>1</b>
1.1	Motivation	1
1.2	Resampling bootstrap	1
1.3	Theory of the bootstrap	2
1.4	The bootstrap in linear models	3
1.5	Weighted bootstrap	3
<b>2</b>	<b>M-estimators and robust regression</b>	<b>4</b>

## 1 The bootstrap

Today we will discuss two new techniques, the bootstrap and M-estimators.

### 1.1 Motivation

Suppose we have an estimate  $\hat{\beta}$  from some procedure. How can we get a handle on the sampling variability of  $\hat{\beta}$ ? Two approaches are:

- Make modelling assumptions such as  $y = X\beta + \varepsilon$  where  $\varepsilon \stackrel{\text{iid}}{\sim} (0, \sigma^2)$ .
- Use the variability in the sample  $(X_i, Y_i)_{i=1}^n \stackrel{\text{iid}}{\sim} P$ .

The bootstrap is an example of the second approach. The bootstrap is due to our own Brad Efron (1979). The idea is if we define the *empirical distribution* of a sample  $(x_i, y_i)_{i=1}^n$  to be:

$$\hat{P} = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{(x_i, y_i)}.$$

Then the idea of the bootstrap is that sampling from  $\hat{P}$  should approximate sampling from the true distribution  $P$ . There are many ways to implement this idea.

### 1.2 Resampling bootstrap

The resampling bootstrap is the most commonly used form of the bootstrap. It is often what people mean when they refer to the bootstrap. The procedure is as follow. First pick a large number  $B$  which will be the number of bootstrap resamples. For  $b = 1, 2, \dots, B$ , do the following:

- Draw a sample of size  $n$   $(X^{*b}, Y^{*b})$  from  $\hat{P}$  where  $n$  is the size of the original sample. This is done by drawing indices  $i_1, i_2, \dots, i_n$  from  $[n] = \{1, \dots, n\}$  uniformly at random with replacement and then setting  $X^{*b} = (X_{i_1}, \dots, X_{i_n})$  and  $Y^{*b} = (Y_{i_1}, \dots, Y_{i_n})$ .
- Then calculate  $\hat{\beta}^{*b} = \hat{\beta}(X^{*b}, Y^{*b})$  from the bootstrap resample.

We can then use  $(\hat{\beta}^{*b})_{b=1}^B$  to estimate the sampling distribution of  $\hat{\beta}$  when we draw  $(X, Y) \stackrel{\text{iid}}{\sim} P$ . There are two forms of error in this approximation:

- There is error from the Monte-Carlo sampling. This can be reduced by increasing  $B$  which should also be taken as large as computationally possible.
- There is also error because  $\hat{P} \neq P$ . This error cannot be eliminated.

### 1.3 Theory of the bootstrap

How can one mathematically describe the “success” of the bootstrap? Suppose we have a parameter of interest  $\theta \in \mathbb{R}$  and we are interested in an estimator  $\hat{\theta}$ . For example we may have  $\hat{\theta} = c^T \hat{\beta}$  or  $\hat{\theta} = \hat{\beta}_j$  or  $\hat{\theta} = \|\hat{\beta}\|_2^2$ . Define

$$J_n(t, P) = \mathbb{P}(\sqrt{n}(\hat{\theta}_n - \theta^*) \leq t),$$

where  $\hat{\theta}_n = \hat{\theta}(X_1, \dots, X_n; Y_1, \dots, Y_n)$  is an estimate calculated on an i.i.d. sample drawn from  $P$  and  $\theta^*$  is the population parameter. The function  $J_n$  is the (normalized) CDF of the sampling distribution of  $\hat{\theta}_n$  about  $\theta^*$ . We can similarly define

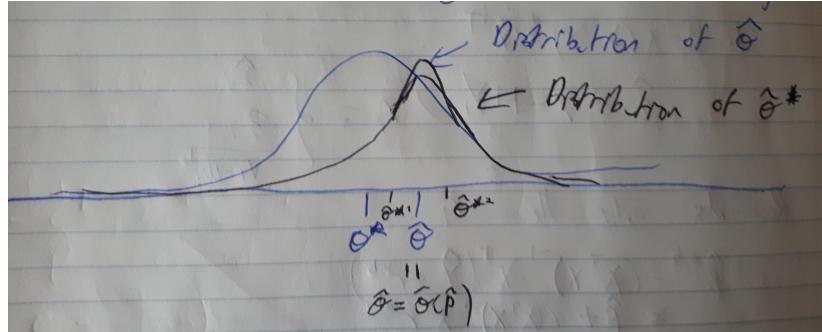
$$J_n(t, \hat{P}) = \mathbb{P}(\sqrt{n}(\hat{\theta}_n^* - \hat{\theta}_n) \leq t),$$

where  $\hat{\theta}_n^* = \hat{\theta}(X^*, Y^*)$  where  $(X^*, Y^*)$  is an i.i.d. sample of size  $n$  from  $\hat{P}$  (that is  $(X^*, Y^*)$  is a bootstrap sample). Note that for the distribution  $\hat{P}$ ,  $\hat{\theta}_n$  is the population parameter. This is why we are now centring at  $\hat{\theta}_n$ . We call  $J_n(t, \hat{P})$  the bootstrap CDF of the statistic  $\hat{\theta}_n$ .

While calculating or approximating  $J_n(t, P)$  maybe be very hard, we can always get a handle on  $J_n(t, \hat{P})$  by drawing more bootstrap samples. Under some regularity conditions we have that with probability 1,

$$\sup_{t \in \mathbb{R}} |J_n(t, P) - J_n(t, \hat{P})| \xrightarrow{n} 0.$$

So for large  $n$ ,  $J_n(t, \hat{P})$  is close to  $J_n(t, P)$ . It can be shown that under some assumptions the rate of convergence is  $\sqrt{\frac{\log n}{n}}$ . We thus have a picture that looks like the one below. The bootstrap density of  $\hat{\theta}^*$  is close to the true density of  $\hat{\theta}$ .



We can use the bootstrap to define  $1 - \alpha$  confidence intervals for parameters  $\theta \in \mathbb{R}$ . One option is

$$CI_\alpha = \left[ \text{Quantile}_{\alpha/2}(\hat{\theta}^{*b}), \text{Quantile}_{1-\alpha/2}(\hat{\theta}^{*b}) \right].$$

Another option is

$$CI'_\alpha = \hat{\theta} \pm z_{1-\frac{\alpha}{2}} \sqrt{\text{Var}(\hat{\theta}^{*b} | \hat{P})},$$

where  $z_{1-\frac{\alpha}{2}}$  is the  $1 - \frac{\alpha}{2}$  quartile of the standard normal. One can show that

$$\mathbb{P}(\theta^* \in CI_\alpha), \mathbb{P}(\theta^* \in CI'_\alpha) \xrightarrow{n} 1 - \alpha.$$

## 1.4 The bootstrap in linear models

One variant of the bootstrap is called the *parametric bootstrap* which can be used if one has faith in their model. The procedure is as follows. In the linear model  $y = X\beta + \varepsilon$ ,

- Fit model  $\hat{\beta}$  based on original sample.
- Resample from the distribution  $\hat{\beta}$  defines in the linear model ie  $y^* = X\hat{\beta} + \varepsilon^*$  but only resample the residuals.
- To resample the residuals, let  $H = X(X^T X)^{-1}X^T$ ,  $\hat{\varepsilon} = Y - HY$  and  $r_i = \hat{\varepsilon}_i$  or  $r_i = \frac{\hat{\varepsilon}_i}{\sqrt{1-H_{ii}}}$  (both are used in practice).
- Generate bootstrap resamples by resampling *only*  $r_i$  to get  $r_i^{*b}$ . Our bootstrap resample of  $\hat{\beta}$  is then

$$\hat{\beta}^{*b} = (X^T X)^{-1} X^T (y + r^{*b}) = \hat{\beta} + (X^T X)^{-1} X^T r^{*b}.$$

We can then use the bootstrap resamples  $\hat{\beta}^{*b}$  as described above. Some pros of this method are:

- It works well when we have a fixed design matrix  $X$ .
- It is more efficient provided the model is true.

A major drawback is:

- It really relies on the model assumptions. It especially relies on the noise being homoskedastic. Problems occur if  $\mathbb{E}[\varepsilon_i^2] = \sigma_i^2 \neq \sigma_j^2 = \mathbb{E}[\varepsilon_j^2]$ .

Thus people often use the *non-parametric bootstrap* for the linear model which does not have the above drawback. Here we resample from  $\hat{P}$  to get bootstrap resamples  $(X^{*b}, y^{*b})$  and then we calculate

$$\hat{\beta}^{*b} = ((X^{*b})^T X^{*b})^{-1} (X^{*b})^T y^{*b}.$$

This doesn't rely on the model being true and it allows for non-homoskedastic noise. There is a drawback that for small sample sizes, the matrix  $X^{*b}$  may be rank deficient but for reasonably large  $n$ , this does not happen often.

## 1.5 Weighted bootstrap

This is sometimes called Bayesian bootstrap although it isn't really Bayesian at all. We first fix a distribution  $W$  which is supported on  $[0, \infty)$  and satisfies  $\mathbb{E}[W] = 1$  and  $\text{Var}[W] < \infty$ . Often people use the exponential distribution with scale 1.

At every bootstrap iteration  $b$ , we sample weights  $w_1^{*b}, \dots, w_n^{*b} \stackrel{\text{iid}}{\sim} W$ . We then define a distribution based on these weights and our observed sample  $(x_i, y_i)_{i=1}^n$ ,

$$\hat{P}^{*b} = \frac{\sum_{i=1}^n w_i^{*b} \mathbf{1}_{(x_i, y_i)}}{\sum_{i=1}^n w_i^{*b}}.$$

We then define

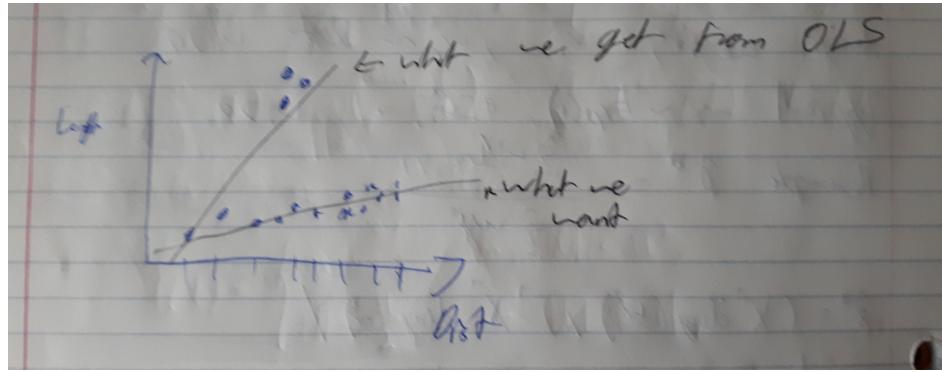
$$\begin{aligned}\hat{\beta}^{*b} &= \underset{c \in \mathbb{R}^d}{\operatorname{argmin}} \mathbb{E}_{\hat{P}^{*b}} \left[ \|y - X^T c\|_2^2 \right] \\ &= \underset{c \in \mathbb{R}^d}{\operatorname{argmin}} \frac{\sum_{i=1}^n w_i^{*b} (y_i - x_i^T c)^2}{\sum_{i=1}^n w_i^{*b}} \\ &= (X^T \operatorname{diag}(w_i^{*b}) X)^{-1} X^T \operatorname{diag}(w_i^{*b}) y.\end{aligned}$$

There are similar asymptotic results for the weighted bootstrap as there are for the resampling bootstrap.

## 2 M-estimators and robust regression

Suppose we want to fit a regression  $y_i = x_i^T \beta + \varepsilon_i$  but maybe we no longer have  $\varepsilon_i \sim N(0, \sigma^2)$ . Maybe there are some outlying  $y_i$ 's that are much larger than we'd expect for normally distributed errors.

For a historical example, imagine trying to predict phone call length from distance to family. We might expect that the further apart two people are, the less often they would catch up and so their phone calls would be longer. The resulting would look something like this:



There is a problem because there are some people who just really liked talking on the phone for a long time. These data points really mess up the linear model. Some solutions are:

- Remove outliers. This can be hard as it is not always clear what make something an outlier.
- Choose an alternative to squared error to use when fitting (Huber, 70s-90s).

We will discuss the second approach.

**Definition 1.** A *loss function* is a function  $l : \mathbb{R} \rightarrow \mathbb{R}^+$ . We will sometimes require that  $l(0) = 0$ .

**Definition 2.** Given data  $(x_i, y_i)_{i=1}^n$  and a loss function  $l$ , the *M-estimator* for predicting  $Y$  from  $X$  is the linear estimator  $\hat{f}(x) = x^T \hat{\beta}$  where

$$\hat{\beta} = \underset{b}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n l(y_i - x_i^T b).$$

The rough recipe is to choose a loss to capture properties we want the predictor to have. For example if  $Y = f(X) + \varepsilon$ , then we might want

$$\operatorname{argmin}_t \mathbb{E}[l(Y - t) | X = x] = f(x).$$

For example for when  $l$  is squared error we get the conditional mean above and when  $l$  is absolute error we get the conditional median. Some examples of loss functions are:

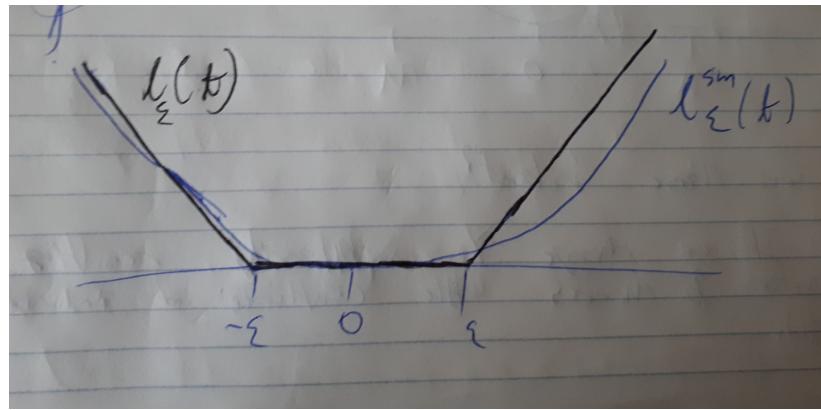
- *Absolute error loss*:  $l(t) = |t|$ . As stated above this has the property that

$$\operatorname{argmin}_t \mathbb{E}[|Y - t|] = \operatorname{median}(Y).$$

- *$\varepsilon$ -insensitive loss*: If we don't care about if the predicted value is off by less than  $\varepsilon$  we can use the loss

$$l^\varepsilon(t) = (t - \varepsilon)_+ + (-t - \varepsilon)_+.$$

This loss looks like the black function



- *Smoothed  $\varepsilon$ -insensitive loss* is a smooth version of  $\varepsilon$ -insensitive loss. It is defined by:

$$l_\varepsilon^{\text{sm}}(t) = \log(1 + \exp(t - \varepsilon)) + \log(1 + \exp(-t - \varepsilon)).$$

This is the blue function in the picture above. We also have rescaled versions of the smoothed  $\varepsilon$ -insensitive loss. For  $a > 0$  define,

$$l_{\varepsilon,a}^{\text{sm}}(t) = a \cdot l_\varepsilon^{\text{sm}}\left(\frac{t}{a}\right),$$

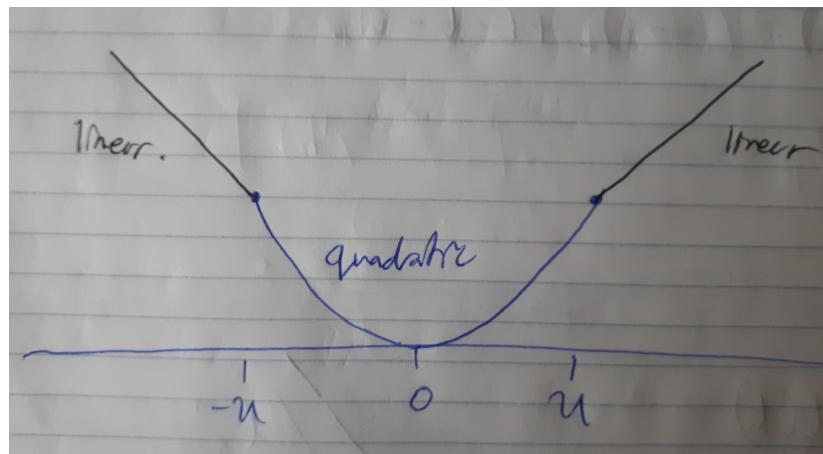
which satisfies

$$\lim_{a \searrow 0} l_{\varepsilon,a}^{\text{sm}}(t) = l_\varepsilon(t).$$

- *Huber loss*: for  $u > 0$ , define

$$l(t) = \begin{cases} \frac{1}{2u}t^2 & \text{if } |t| \leq u, \\ |t| - \frac{u}{2} & \text{if } |t| \geq u. \end{cases}$$

Huber loss looks like this:



Why are these M-estimators more “robust”? The rough idea is that if we change a single  $y_i$  arbitrarily, then the resulting estimator  $\hat{\beta}$  changes little. Thus the influence of individual observations is small.

Some loss functions give estimators with this property, others do not. For example, we have seen that under squared error, points with high leverage can have a large influence on the model.

It turns out that two things are key for an M-estimator to be robust:

- The features  $X_i$  are bounded.
- The loss function is (1)-Lipschitz. That is for all  $t, s \in \mathbb{R}$ ,  $|l(t) - l(s)| \leq t - s$  or equivalently  $|l'(t)| \leq 1$  for all  $t \in \mathbb{R}$ .

This will be discussed more later and will be explored on upcoming homework.