# Bootstrap (the basics)

- Use it to get confidence intervals for a very wide variety of estimators.

- Under the bootstrap assumption:

$$\sigma^2_{\substack{\text{sampling} \\ (\text{w/replacement})}} \approx \sigma^2_{\substack{(\text{across} \\ \text{samples})}}$$

- Let $B$ (usually in the hundreds or thousands) be the number of bootstrap replicates.

- If our original sample is $(x_i, y_i)_{i=1}^{n}$ Then, for each bootstrap replicate we'll have

$$(x_{J(b)}, y_{J(b)})_{J \in I}$$

where $I$ = set of indices = $\{1, \dots, N\}$

and $\{J\}_{(b)}$ is a sequence of indices which

# of bootstrap replicate for $b = 1, \dots, B$.

are sampled with replacement from $I$.

Example. we have data sample

$(x_1, x_1), (x_2, y_2), (x_3, y_3),$

$(x_4, y_4), (x_5, y_5)$

$\Rightarrow I = \{1, 2, 3, 4, 5\}$ <span style="color:orange">(Set of indices)</span>

$(b=1) \rightarrow J_{(1)} = \{1,1, 2, 3, 4\}$ <span style="color:orange">(Sequence of indices sampled for bootstrap replicate 1).</span>

$(b=2) \rightarrow J_{(2)} = \{1, 2, 3, 3, 3\}$

So, for $b=1$ ...

$(X_{J_{(1)}}, Y_{J_{(1)}})$

For each $b$, compute $\hat{y}_{J_{(b)}} = E[Y_{J_{(b)}} | X_{J_{(b)}}]$ using our lowess estimator. Then, save $\{J_{(b)}$
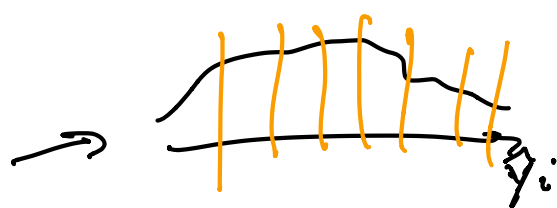
and $\hat{y}_{(b)}$

We end up with table ...

<span style="text-decoration:underline">for each observation $i$</span>

| $b \rightarrow$ | 1 | 2 | 3 | ... | B |
|---|---|---|---|---|---|
| $i \downarrow$ | | | | | |
| 1 | $\hat{y}_{1(1)}$ | — | | | |
| 2 | $\hat{y}_{2(1)}$ | $\hat{y}_{2(2)}$ | | | |
| 3 | — | $\hat{y}_{3(2)}$ | | | |
| 4 | $\hat{y}_{4(1)}$ | — | | | |
| 5 | $\hat{y}_{5(1)}$ | — | | | |
| ... | | | | | |
| N | $\hat{y}_{N(1)}$ | $\hat{y}_{N(2)}$ | | | |



Further, you can compute an estimate of standard error (observation-specific)
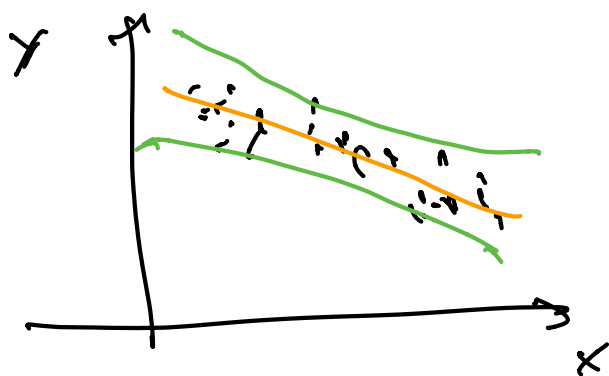
$$\sigma^{Boot}_{\hat{y}_i} = \sqrt{\frac{\sum (\hat{y}_{i(b)} - \overline{\hat{y}_i})^2}{B_i - 1}}$$

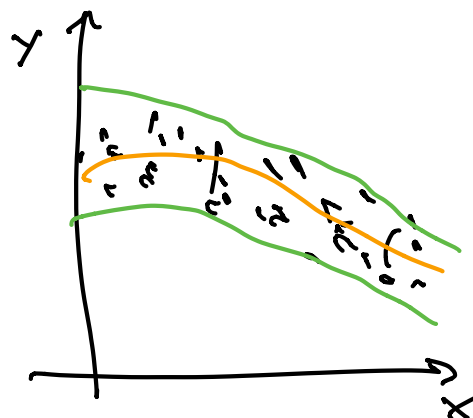<span style="color:orange">Average across bootstrap replicates</span>

$B$ is fixed, $B_i$ is variable

Comparing the width and shape of CI for OLS
vs. Bootstrap CI's, we observe:



OLS

① OLS CI's are tighter ( considers an infinite
# of samples).

② On the other hand, our CI's for Lowess
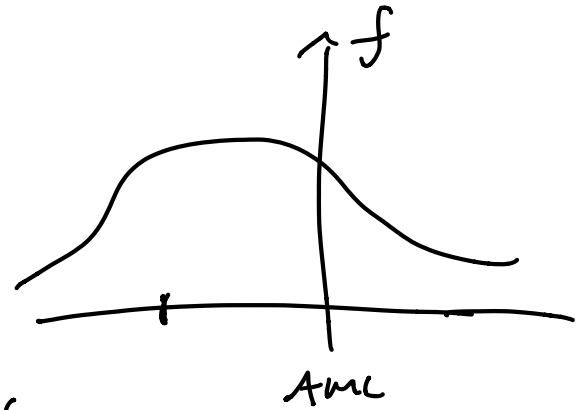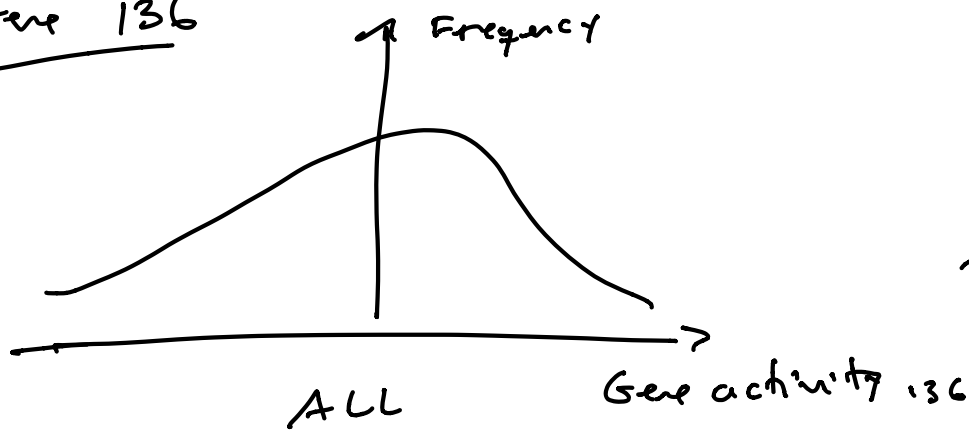just consider $\underline{B}$ samples.

bootstrap

---

Big data implications on hypothesis testing.

- 72 Leukemia Patients
  - 47 with ALL ( Acute Lymphoblastic Leukemia)
  - 25 with AML ( Acute Myeloid Leukemia)

— Measurements available for 7128 genes
for each patient !
→ Let's focus on 1 of them (Gene 136)

Gene 136



ALL                  Gene activity 136                AML

- Classic way to check if means are equal
is using a t-test for hypothesis

$$H_0 : \mu_{AML} = \mu_{ALL}$$

$$t = \frac{\overline{Y}_{136}^{AML} - \overline{Y}_{136}^{ALL}}{\widehat{se}}$$

- $\widehat{se}$ = standard error estimate.

- For simplicity, let's assume **equal** variances.

$$\widehat{se} = \hat{\sigma} \left( \frac{1}{n_1} + \frac{1}{n_2} \right)^{1/2}$$

where $n_2 = 47$
$n_1 = 25$

- $V = n_1 + n_2 - 2 = 70$ (degrees of freedom)