**PS5841**

# Data Science in Finance & Insurance

# Bootstrap & Bagging

Yubo Wang

Autumn 2021
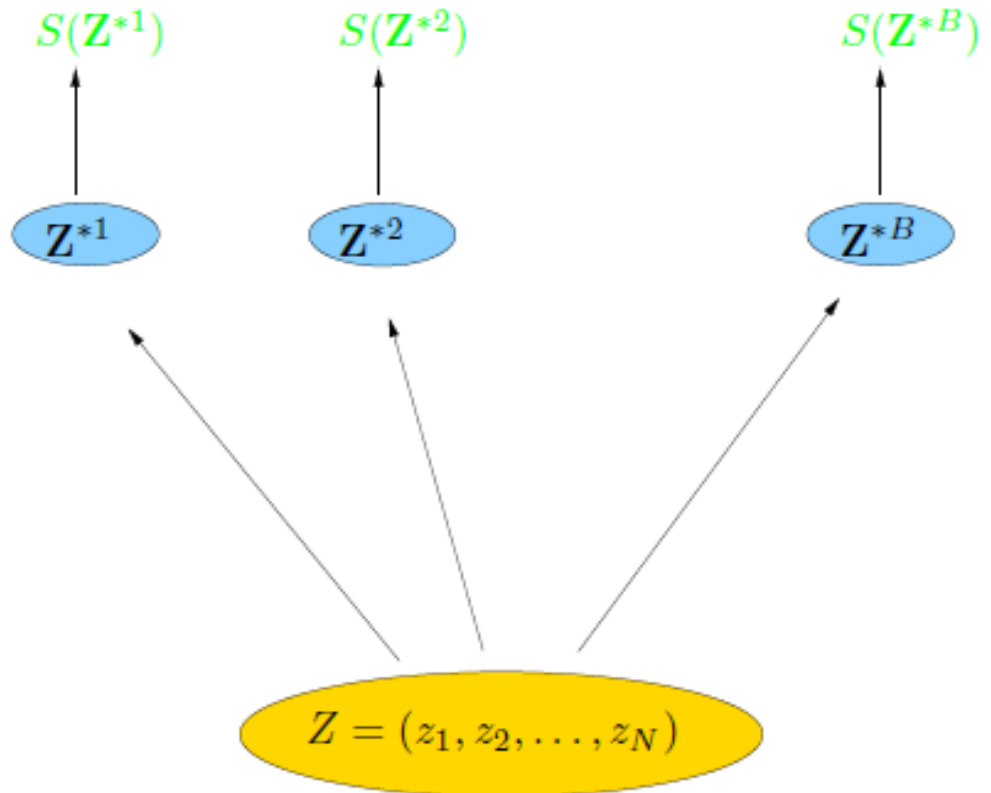
# Bootstrap

# The Bootstrap

- Emulate the process of obtaining new samples without generating additional samples
- Obtain distinct data sets by repeatedly sampling, with replacement, observations from the original data set
- Works well if the original data set is a good representation of the population
  - Risk of "garbage in, garbage out"
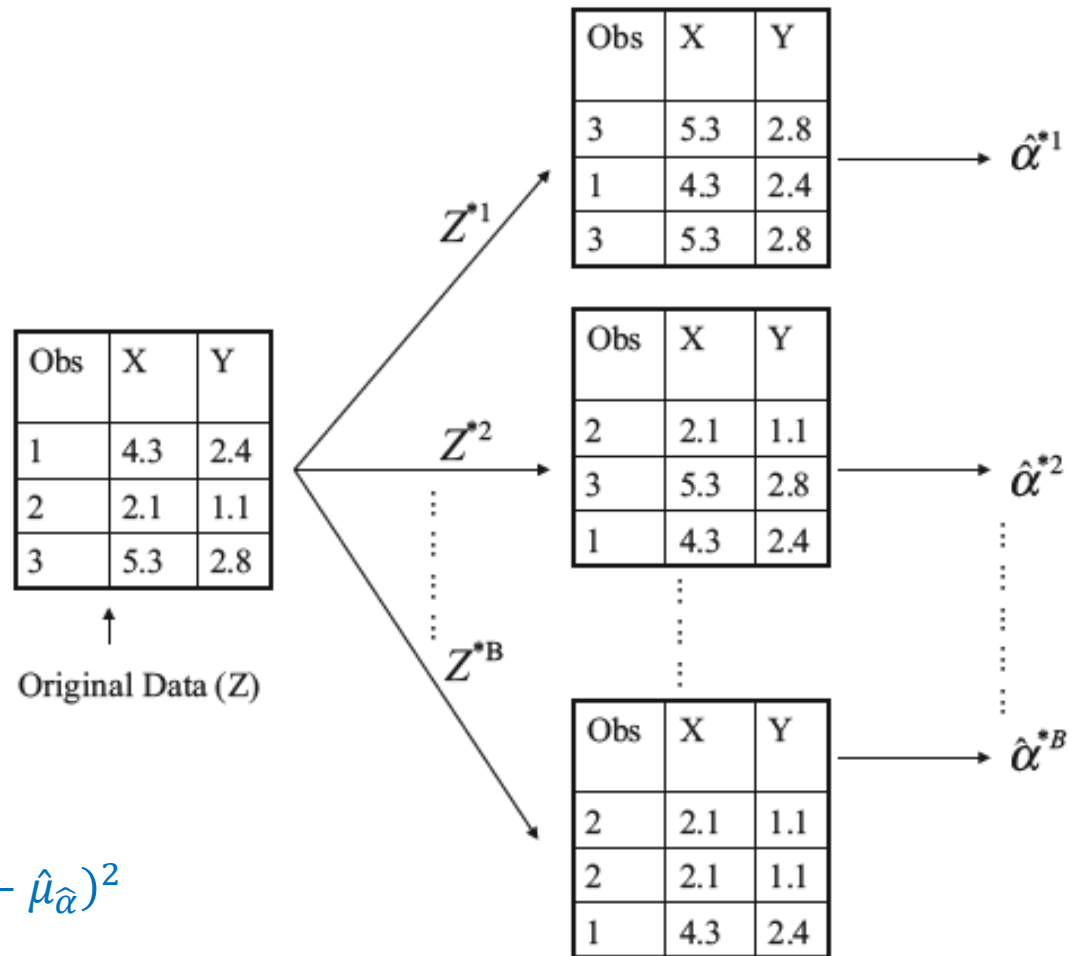
# The Bootstrap

quantity of interest $\qquad$ $S(\mathbf{Z}^{*1})$ $\qquad$ $S(\mathbf{Z}^{*2})$ $\qquad$ $S(\mathbf{Z}^{*B})$

bootstrap data sets $\qquad$ $\mathbf{Z}^{*1}$ $\qquad$ $\mathbf{Z}^{*2}$ $\qquad$ $\mathbf{Z}^{*B}$

original data set $\qquad$ $Z = (z_1, z_2, \ldots, z_N)$

# Bootstrap Data Sets

$$\hat{\mu}_{\hat{\alpha}} = \frac{1}{B} \sum_{r=1}^{B} \hat{\alpha}^{*r}$$

| Obs | X | Y |
|-----|-----|-----|
| 3 | 5.3 | 2.8 |
| 1 | 4.3 | 2.4 |
| 3 | 5.3 | 2.8 |

$Z^{*1}$ → $\hat{\alpha}^{*1}$

| Obs | X | Y |
|-----|-----|-----|
| 1 | 4.3 | 2.4 |
| 2 | 2.1 | 1.1 |
| 3 | 5.3 | 2.8 |

Original Data (Z)

$Z^{*2}$

| Obs | X | Y |
|-----|-----|-----|
| 2 | 2.1 | 1.1 |
| 3 | 5.3 | 2.8 |
| 1 | 4.3 | 2.4 |

→ $\hat{\alpha}^{*2}$

$Z^{*B}$

| Obs | X | Y |
|-----|-----|-----|
| 2 | 2.1 | 1.1 |
| 2 | 2.1 | 1.1 |
| 1 | 4.3 | 2.4 |

→ $\hat{\alpha}^{*B}$

$$\hat{\sigma}_{\hat{\alpha}}^2 = \frac{1}{B-1} \sum_{r=1}^{B} (\hat{\alpha}^{*r} - \hat{\mu}_{\hat{\alpha}})^2$$

# Example   (1)

- $X \sim \mathcal{N}(\mu_X = 0, \sigma_X^2 = 10^2)$
- Training Set   $Z = (x_1, x_2, \ldots, x_{N=100})$
- Quantity of interest $\bar{X} = \frac{1}{N} \sum_{i=1}^{N} x_i$

$$\texttt{true\_mean}(\bar{X}) = \mu_X = 0$$

$$\texttt{true\_se}(\bar{X}) = \frac{\sigma_X}{\sqrt{N}} = 1$$

# Example   (2)

- $X \sim \mathcal{N}(\mu_X = 0, \sigma_X^2 = 10^2)$
- Training Set   $Z = (x_1, x_2, \ldots, x_{N=100})$
- Quantity of interest $\bar{X} = \frac{1}{N} \sum_{i=1}^{N} x_i$
- Generate B=1000 <span style="color:red">true</span> training sets

$$\texttt{true\_sampling\_mean}(\bar{X}) = \frac{1}{B} \sum_{r=1}^{B} \overline{x}^{(r)}$$

$$\texttt{true\_sampling\_se}(\bar{X})$$

$$= \left[ \frac{1}{B-1} \sum_{r=1}^{B} \left( \bar{x}^{(r)} - \texttt{true\_sampling\_mean}(\bar{X}) \right)^2 \right]^{\frac{1}{2}}$$

# Example   (3)

- $X \sim \mathcal{N}(\mu_X = 0, \sigma_X^2 = 10^2)$
- Training Set   $Z = (x_1, x_2, \ldots, x_{N=100})$
- Quantity of interest $\bar{X} = \frac{1}{N} \sum_{i=1}^{N} x_i$
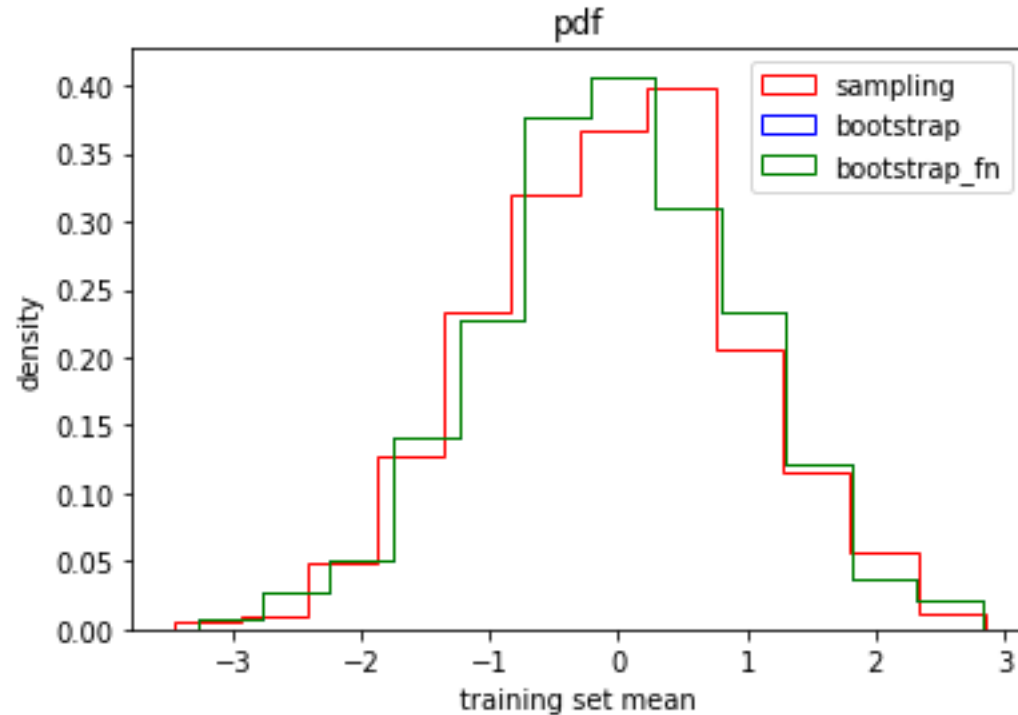- Generate B=1000 <span style="color:red">bootstrap</span> training sets

$$\texttt{boot\_sampling\_mean}(\bar{X}) = \frac{1}{B} \sum_{r=1}^{B} \bar{x}^{*(r)}$$

$$\texttt{boot\_sampling\_se}(\bar{X})$$

$$= \left[ \frac{1}{B-1} \sum_{r=1}^{B} \left( \bar{x}^{*(r)} - \texttt{boot\_sampling\_mean}(\bar{X}) \right)^2 \right]^{\frac{1}{2}}$$

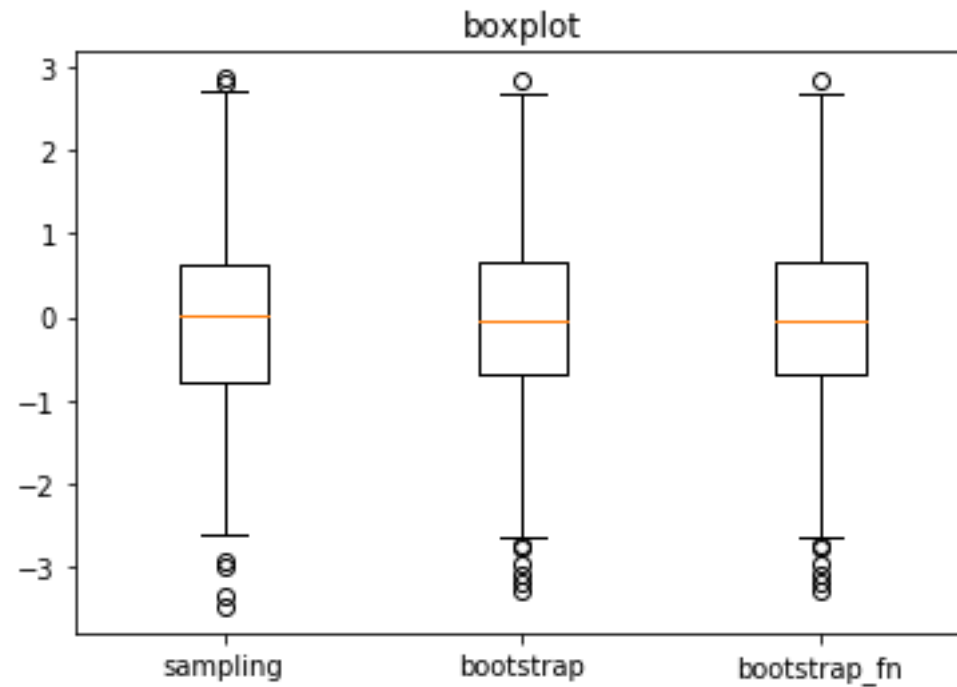# Example   (4)



```
theoretical mean & se            =    0.000000,  1.000000
"true" sampling mean & se        =  -0.045906,  1.009381
"bootstrap" sampling mean & se   =  -0.033890,  1.005676
```

# Example (5)



boxplot

# Bagging

# Averaging to reduce variance in SL

Recall that given a set of $n$ independent observations $Z_1, \ldots, Z_n$, each with variance $\sigma^2$, the variance of the mean $\bar{Z}$ of the observations is given by $\sigma^2/n$. In other words, *averaging a set of observations reduces variance.* Hence a natural way to reduce the variance and hence increase the prediction accuracy of a statistical learning method is to take many training sets from the population, build a separate prediction model using each training set, and average the resulting predictions. In other words, we could calculate $\hat{f}^1(x), \hat{f}^2(x), \ldots, \hat{f}^B(x)$ using $B$ separate training sets, and average them in order to obtain a single low-variance statistical learning model, given by

$$\hat{f}_{\text{avg}}(x) = \frac{1}{B} \sum_{b=1}^{B} \hat{f}^b(x).$$

# Bagging

Of course, this is not practical because we generally do not have access to multiple training sets. Instead, we can bootstrap, by taking repeated samples from the (single) training data set. In this approach we generate $B$ different bootstrapped training data sets. We then train our method on the $b$th bootstrapped training set in order to get $\hat{f}^{*b}(x)$, and finally average all the predictions, to obtain

$$\hat{f}_{\text{bag}}(x) = \frac{1}{B} \sum_{b=1}^{B} \hat{f}^{*b}(x).$$

This is called bagging.

# That was