

Data generating process (DGP)

$$P(x, y) = P(y|x=x) P(x)$$

\uparrow matrix r.v. \uparrow vertex r.v. \uparrow constant matrix \uparrow $\mathbb{D} = \left(\begin{array}{|c|} \hline x \\ \hline y \\ \hline \end{array} \right)$

Imagine you got to see $\mathbb{D}_1, \mathbb{D}_2, \dots, \mathbb{D}_m$ where m is large
 Now you can fit $g_1 = A(\mathbb{D}_1), \dots, g_m = A(\mathbb{D}_m)$; then you can use the avg:

$$g_{\text{avg}} = \frac{g_1 + \dots + g_m}{m} \approx E[G] \quad g_i = A(\mathbb{D}_i) \text{ are iid}$$

$$\text{MSE} = \sigma^2 + E_x [\text{Bias}[G_{\text{avg}}]^2] + E_x [\text{Var}[G_{\text{avg}}]]$$

$$= \sigma^2 + E_x \left[E \left[\left(\frac{g_1 + \dots + g_m}{m} - f \right)^2 \right] \right] + E_x \left[\text{Var} \left[\frac{g_1 + \dots + g_m}{m} \right] \right]$$

$$= \sigma^2 + E_x \left[\left(\frac{1}{m} ((E[g_1] - f) + \dots + (E[g_m] - f)) \right)^2 \right] + E_x \left[\frac{1}{m^2} (\text{Var}[g_1] + \dots + \text{Var}[g_m]) \right]$$

constant

$$\approx \sigma^2 + E_x [\text{Bias}[g_i]] + \frac{1}{m} E_x [\text{Var}[g_i]]$$

if m is large

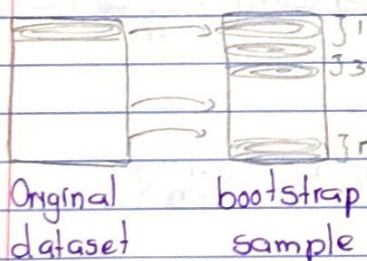
$$\approx \sigma^2 + E_x [\text{Bias}[g_i]] \approx \sigma^2 \text{ theoretical minimum MSE}$$

Pick H very expressive

Leo Breiman invented "bagging" (Bootstrap Aggregating) in 1994, with one \mathbb{D} : What now?

Why not consider a sampling of the \mathbb{D} : $\mathbb{D}_{(1)}, \mathbb{D}_{(2)}, \dots, \mathbb{D}_{(m)}$ when m is large

How to sample? Sample n rows/unit with replacement (Non-parametric bootstrap sampling).



$\mathbb{D}_{(1)}$ has both duplicate rows of \mathbb{D} and omissions.

$$P(\text{learning row one}) = 1 - \frac{1}{n}$$

$$P(\text{learning row } n \text{ times}) = (1 - \frac{1}{n})^n \approx e^{-1} \approx \frac{1}{3}$$

n big

\Rightarrow In each bootstrap sample $\approx 1/3$ omissions, i.e. $2/3$ of the original \mathbb{D} .

Now, $[i] \ g_{(i)} = A(\mathbb{D}_{(i)}), \dots, g_{(m)} = A(\mathbb{D}_{(m)})$

$g_{(1)}, \dots, g_{(m)}$ are not the same, but they will be dependent

$$g_{\text{avg}} = \frac{g_{(1)} + \dots + g_{(m)}}{m}$$

Back to math 241,

Let x_1, \dots, x_m i.i.d $\text{Var}[\bar{x}] = \text{Var}[\frac{1}{n} \sum x_i]$
 $\stackrel{\text{i.i.d}}{=} \frac{1}{n^2} \sum \text{Var}[x_i] = 0 \stackrel{\text{ident. distribution}}{=} \frac{1}{n^2} n \sigma^2 = \sigma^2/n$

let x_1, \dots, x_n be dependent but identically distribution,
 let $\sigma^2 = \text{Var}[x_i]$ and have the same dependence
 let $\rho = \text{Corr}[x_i, x_j] \ i \neq j$

$$\text{Corr}[x_i, x_j] = \frac{\text{Cov}[x_i, x_j]}{\text{SE}[x_i] \text{SE}[x_j]}$$

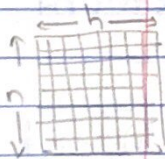
$$\Rightarrow \rho = \frac{\text{Cov}[x_i, x_j]}{\sigma \cdot \sigma} \Rightarrow \text{Cov}[x_i, x_j] = \sigma^2 \rho$$

$$\text{Var}[\bar{x}] = \frac{1}{n^2} \cdot \text{Var}[\sum x_i]$$

$$= \frac{1}{n^2} (\text{Var}[x_1] + \dots + \text{Var}[x_n] + \sum_{i \neq j} \text{Cov}[x_i, x_j])$$

if n large, $= \frac{1}{n^2} (n\sigma^2 + (n^2 - n)\sigma^2 \rho)$

$$\approx \sigma^2 \rho = \frac{\sigma^2}{n} + \frac{(n-1)}{n} \sigma^2 \rho = \sigma^2 (\frac{1}{n} - \frac{n-1}{n} \rho)$$



$$= \frac{1}{n} (\sigma^2 + n\sigma^2\rho - \sigma^2\rho)$$

$$= \frac{\rho\sigma^2 + 1-\rho}{n}\sigma^2$$

if $\rho \approx 0$ $\approx \frac{\sigma^2}{n}$

If $g_{(1)}, g_{(2)}$ have same correlation ρ and we use g_{avg} then

$$MSE = \sigma^2 + E_x [\text{Bias}[g_{avg}]^2] + E_x [\text{Var}[g_{avg}]]$$

Let H be very expressive $\Rightarrow \text{Bias}[g_{(1)}] = 0 \Rightarrow \text{Bias}[g_{avg}] = 0$

$$= \sigma^2 + E_x [\text{Var}[g_{avg}]] = \sigma^2 + E_x [\underbrace{\rho \text{Var}[g_{(1)}]}_{\text{constant}} + \frac{1-\rho}{m} \text{Var}[g_{(1)}]]$$

if $m \rightarrow \infty$ \swarrow gain from bagging.

$$= \sigma^2 + E_x [\rho \text{Var}[g_{(1)}]]$$

Since $D_{(1)}, D_{(2)}$ share observations $\Rightarrow g_{(1)}, g_{(2)}$ will be positively correlated i.e. $\rho \in (0, 1)$

We want to minimize ρ as much as possible.