

data generating process (DGP)

$$P(X, Y) = P(Y | X = x) P(X)$$

matrix r.v.      vector r.v.      constant matrix

Imagine you get to see  $D_1, D_2, \dots, D_m$  where  $m$  is large. iid from  $P(X, Y)$

Now you can fit  $g_1 = A(D_1), \dots, g_m = A(D_m)$ . Then you can use the avg:

$$g_{\text{avg}} = \frac{g_1 + \dots + g_m}{m} \approx E[G]$$

$$\begin{aligned} \text{MSE} &= \sigma^2 + E_x [\text{Bias}[g_{\text{avg}}]^2] + E_x [\text{Var}[g_{\text{avg}}]] \\ &= \sigma^2 + E_x \left[ E \left[ \left( \frac{g_1 + \dots + g_m}{m} - f \right)^2 \right] \right] + E_x \left[ \text{Var} \left[ \frac{g_1 + \dots + g_m}{m} \right] \right] \end{aligned}$$

$g_i = A(D_i)$  are iid

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

$$= \sigma^2 + E_x [\text{Bias}[g_i]] + \frac{1}{m} \underbrace{E_x [\text{Var}[g_i]]}_{\text{constant}}$$

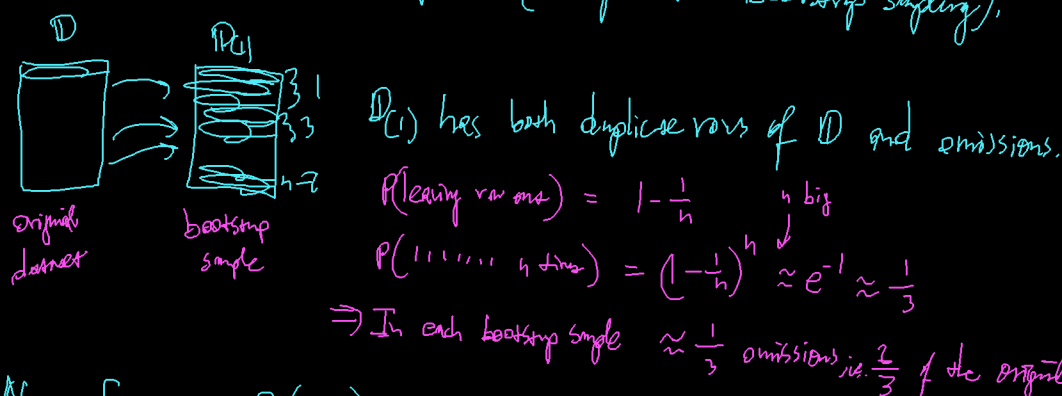
if  $m$  is large  $\downarrow$

$$\approx \sigma^2 + E_x [\text{Bias}[g_i]] \approx \sigma^2 \quad \text{pick } \mathcal{H} \text{ very expressive}$$

homotetic minimum MSE.

Leo Breiman invented "bagging" (Bootstrap AG-reg+ING) in 1994.

with one  $D$ : what can we do?  
 Why not consider a sampling of the  $D$ :  $D_{(1)}, D_{(2)}, \dots, D_{(m)}$  where  $m$  is large  
 How to sample? Sample  $n$  rows/units with replacement. (Non-parametric bootstrap sampling).



Now, fit  $g_{(1)} = A(D_{(1)}), \dots, g_{(m)} = A(D_{(m)})$

$g_{(1)}, \dots, g_{(m)}$  are not the same but they will be dependent.  $g_{\text{avg}} = \frac{g_{(1)} + \dots + g_{(m)}}{m}$

Back to Math 241... let  $X_1, \dots, X_n$  iid ident. distr

$$\text{Var}[\bar{X}] = \text{Var} \left[ \frac{1}{n} \sum X_i \right] = \frac{1}{n^2} \sum \text{Var}[X_i] = \frac{1}{n^2} n \sigma^2 = \frac{\sigma^2}{n}$$

let  $X_1, \dots, X_n$  be dependent but identically distr. let  $\sigma^2 = \text{Var}[X_i]$  and have the same dependence. let  $\rho = \text{Corr}[X_i, X_j] \quad i \neq j$

$$\text{Corr}[X_i, X_j] = \frac{\text{Cov}[X_i, X_j]}{\text{SE}[X_i] \text{SE}[X_j]} \Rightarrow \rho = \frac{\text{Cov}[X_i, X_j]}{\sigma \cdot \sigma} \Rightarrow \text{Cov}[X_i, X_j] = \sigma^2 \rho$$

$$\begin{aligned} \text{Var}[\bar{X}] &= \frac{1}{n^2} \text{Var} \left[ \sum X_i \right] = \frac{1}{n^2} \left( \underbrace{\text{Var}[X_1] + \dots + \text{Var}[X_n]}_{\text{if } n \text{ large}} + \sum_{i \neq j} \text{Cov}[X_i, X_j] \right) \\ &= \frac{1}{n^2} (n \sigma^2 + (n^2 - n) \sigma^2 \rho) = \frac{\sigma^2}{n} + \frac{n-1}{n} \sigma^2 \rho = \sigma^2 \left( \frac{1}{n} + \frac{n-1}{n} \rho \right) \\ &= \frac{1}{n} (\sigma^2 + (n-1) \sigma^2 \rho) \\ &= \rho \sigma^2 + \frac{1-\rho}{n} \sigma^2 \end{aligned}$$

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

If  $g_{(1)}, g_{(2)}$  have same correlation  $\rho$  and we use  $g_{\text{avg}}$  then...

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

let  $\mathcal{H}$  be very expressive  $\Rightarrow \text{Bias}[g_i] = 0 \Rightarrow \text{Bias}[g_{\text{avg}}] = 0$

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

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

gain from bagging.

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  $\rho$  as much as possible.