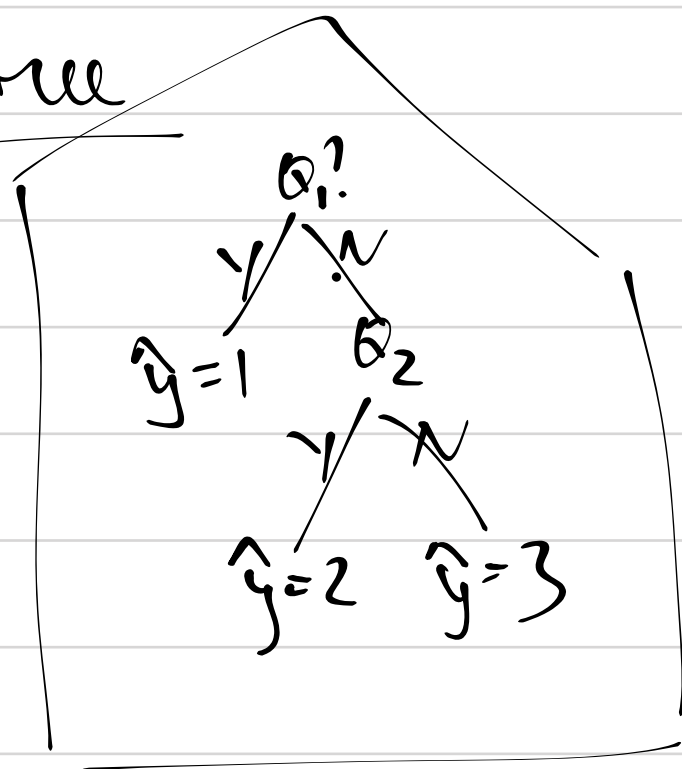


Random Forest + Data Splitting Strategies.

one tree



- ⊕ very interpretable
- ⊖ forecast quality is not very high

random forest = many trees.

Why many trees can be better?

Ex future value: Y_F ← indep of $\hat{Y}_1 \dots \hat{Y}_m$

Some independent forecasts that have the same bias; same variance
 $\hat{Y}_1, \hat{Y}_2, \hat{Y}_3 \dots \hat{Y}_m$

$$\begin{aligned} \text{a) } \underline{MSE(\hat{Y}_1)} &= ? & E((Y_F - \hat{Y}_1)^2) \\ \text{b) } MSE(\hat{Y}) &= ? & = E((Y_F - \hat{Y})^2) \\ \hat{Y} &= \frac{\hat{Y}_1 + \hat{Y}_2 + \dots + \hat{Y}_m}{m} \end{aligned}$$

$$E(W^2) = \text{Var}(W) + (E(W))^2 \quad \text{Var}(W) = E(W^2) - (E(W))^2$$

$$\text{MSE}(\hat{Y}_1) = \text{Var}(Y_F - \hat{Y}_1) + \left(\mathbb{E}(Y_F - \hat{Y}_1) \right)^2 =$$

$$= \text{Var}(Y_F) + \text{Var}(\hat{Y}_1) + \left(\mathbb{E}(Y_F) - \mathbb{E}(\hat{Y}_1) \right)^2$$

$$\text{MSE}(\hat{Y}) = \text{Var}(Y_F - \hat{Y}) + \left(\mathbb{E}(Y_F - \hat{Y}) \right)^2 =$$

$$= \text{Var}(Y_F) + \text{Var}(\hat{Y}) + \left(\mathbb{E}(Y_F) - \mathbb{E}(\hat{Y}) \right)^2$$

$$\hat{Y} = \frac{\hat{Y}_1 + \hat{Y}_2 + \dots + \hat{Y}_m}{m}$$

$$\text{Var}(\hat{Y}_1) > \text{Var}(\hat{Y}) = \frac{\text{Var}(\hat{Y}_1)}{m}$$

$$n \text{Var}(\bar{Y}) = \frac{\sigma^2}{n}$$

$$\mathbb{E}(\hat{Y}) = \mathbb{E}\left(\frac{\hat{Y}_1 + \dots + \hat{Y}_m}{m}\right) = \mathbb{E}(\hat{Y}_1)$$

It would be great to have m independent trees instead of one tree and just average the forecasts of all the trees!

→ we need to introduce randomness in tree growth process.

randomness

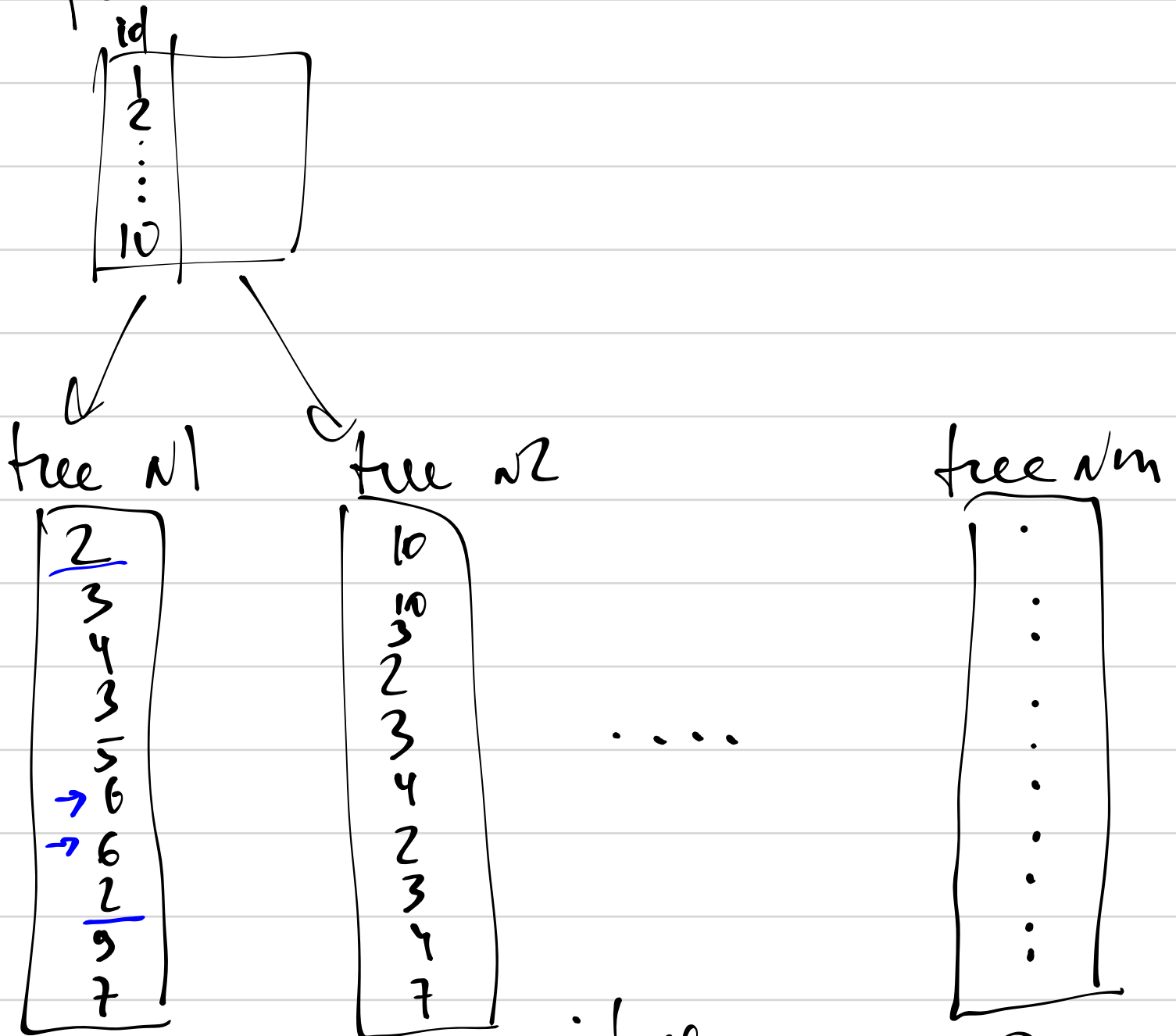
two sources of randomness

id	y_i	x_i	w_i
1	y_1	x_1	w_1
2	\vdots	\vdots	\vdots
3	\vdots	\vdots	\vdots
4	\vdots	\vdots	\vdots
\vdots	\vdots	\vdots	\vdots
n	y_n	x_n	w_n

↑
response
variable

predictors / regressors / features

source n1: before growing the next tree
take a sample (with replacement) of
n observations out of n from original
sample.



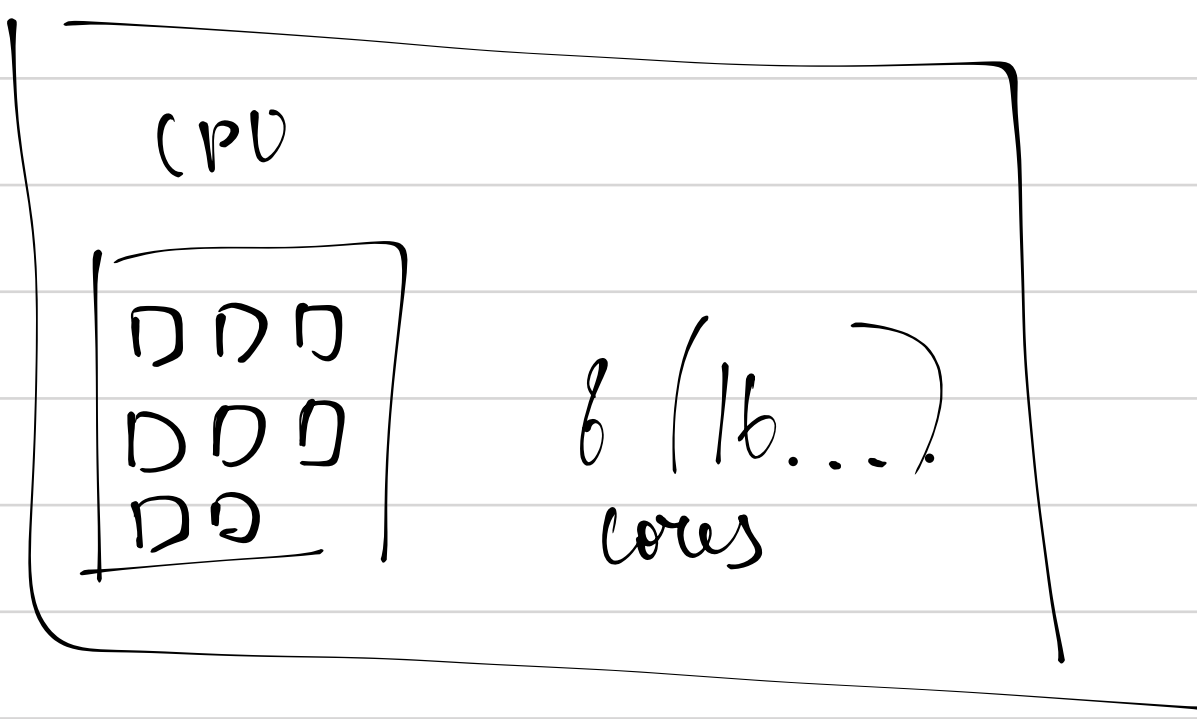
ideal : $m \rightarrow \infty$

$m = 10000$

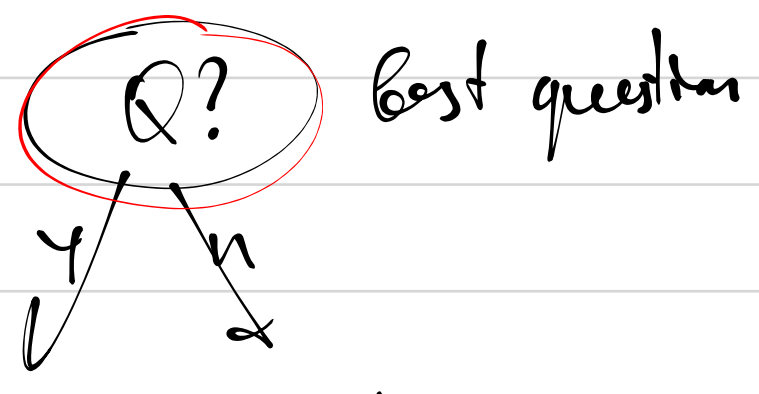
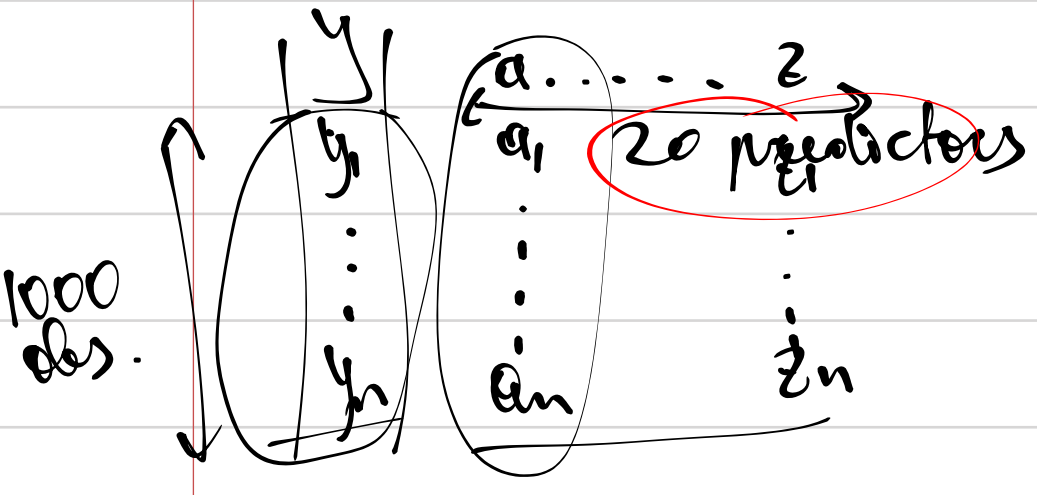
source n2

Source NZ

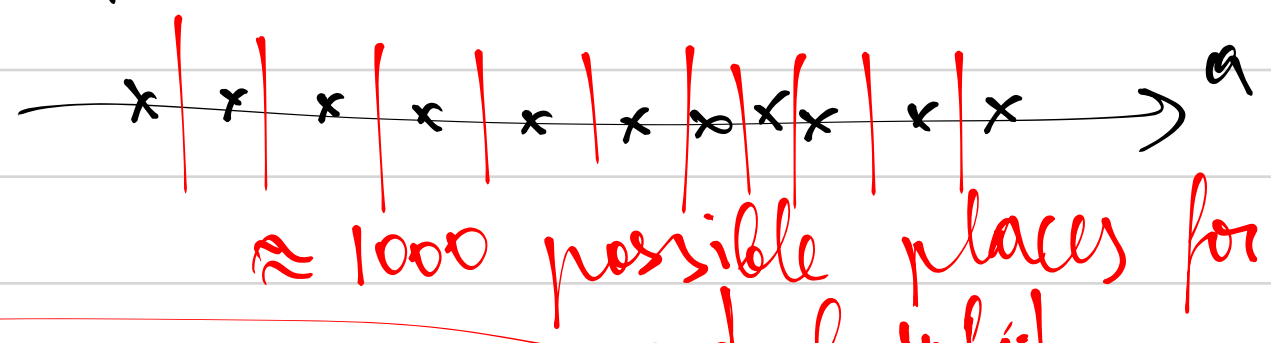
random
preslection
of variables before
each split.



Response



continuous predictor \Rightarrow almost all values are different

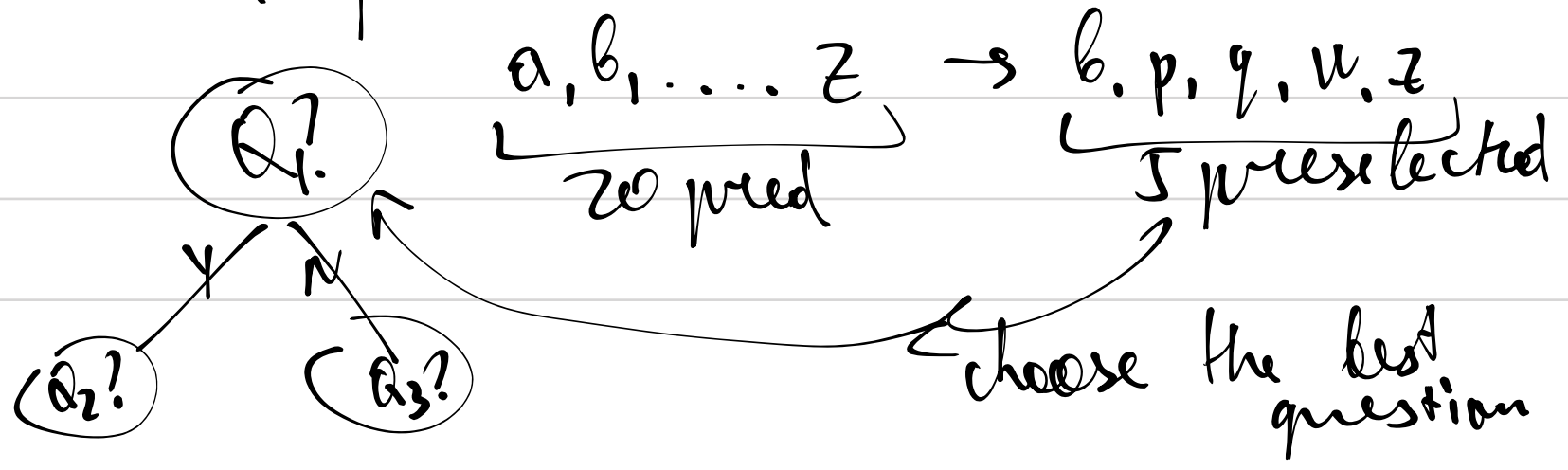


$b > 7$ $a > 6$... $z > 20$

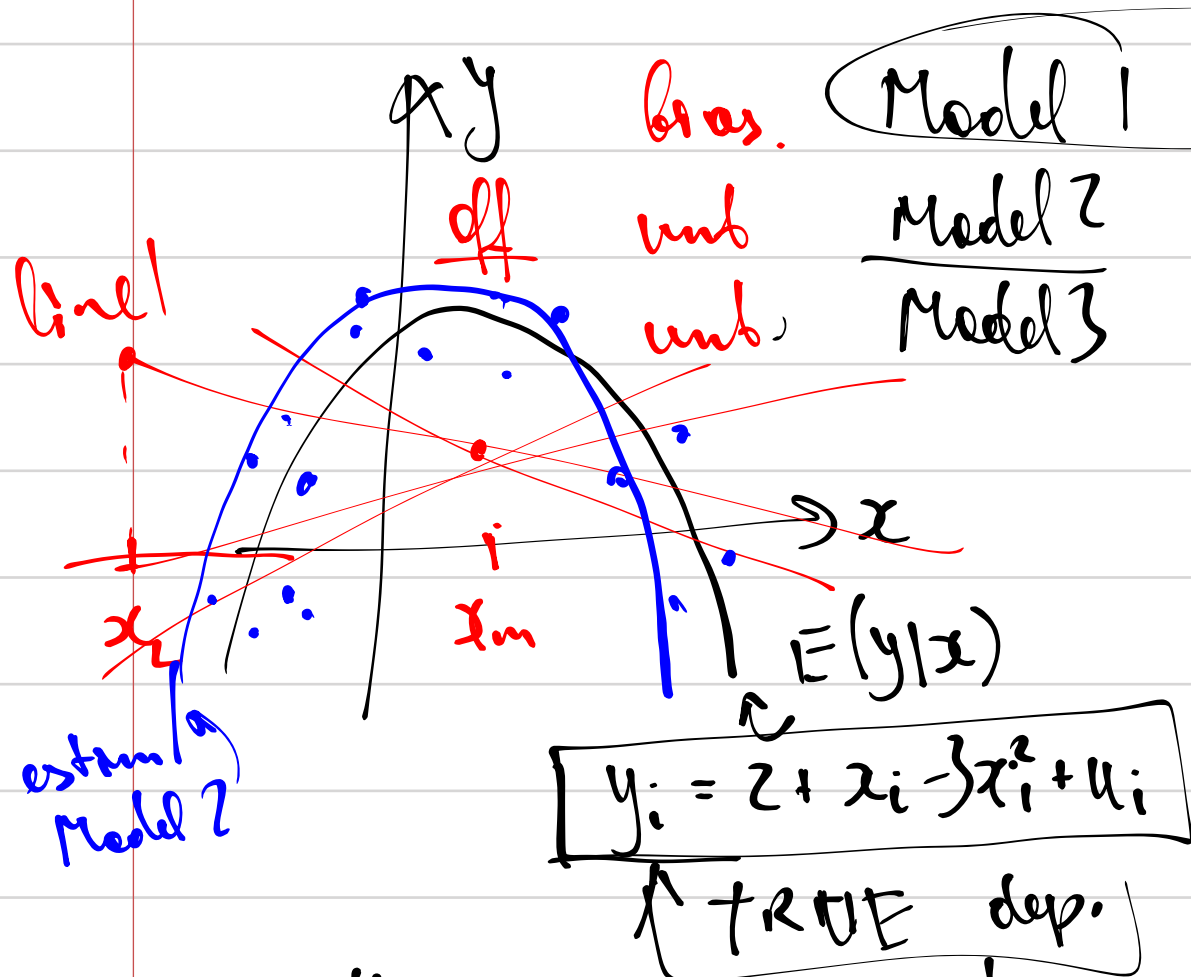
10000 trees \times 100 questions \times 20 predictors \times 1000 possible splitting points

\rightarrow introduce randomness
 \rightarrow decrease computation time

\rightarrow before each question we randomly select 5 predictors out of 20. We choose optimal question only by considering these 5 preselected predictors.



Can we do something to decrease bias?



$$\hat{y}_i = \hat{\beta}_1 + \hat{\beta}_2 \cdot x_i$$

$$\hat{y}_i = \hat{\beta}_1 + \hat{\beta}_2 \cdot x_i + \hat{\beta}_3 \cdot x_i^2$$

$$\hat{y}_i = \hat{\beta}_1 + \hat{\beta}_2 \cdot x_i + \hat{\beta}_3 \cdot x_i^2 + \hat{\beta}_4 \cdot x_i^3$$

a) which model will produce unbiased $\hat{\beta}_i$?

b) which model has lowest variance

moral : to decrease bias grow the biggest possible trees.

normally stopping criterion for random forest looks like: grow the tree until it is possible until in each node there are no more than 5 observations / .

$$\downarrow \left(E(Y_F - \hat{Y}) \right)^2 \Rightarrow \text{long-long trees}$$

$$\downarrow \text{Var}(\hat{Y}) \Rightarrow \text{big number of trees.}$$

- random forest
- ④ good quality of forecasts
(even without tuning)
 - ⑤ almost no interpretation

How to measure importance of variables?

permutation test.

idea: if y does not depend on a particular predictor x then changing the order of observations in x should not affect the forecast quality.

y_1	x_1	\longrightarrow	y_1	x_3
y_2	\vdots		\vdots	x_1
\vdots	\vdots		\vdots	x_n
\vdots	\vdots		\vdots	x_j
y_n	x_n		y_n	x_j

\rightarrow shuffle the predictor x 100 times
for each shuffle calculate the quality of the forecasts, calculate the average drop of the forecast quality.
drop of MSE due to shuffling.

