

# RANDOM FORESTS

- Bagging can greatly improve the performance of CART
  - Averaging over ensemble prediction, in case of regression trees
  - Majority vote over ensemble prediction, for classification trees
- However, the 'B' bootstrapped dataset are **correlated**!

Therefore, the variance reduction due to averaging is diminished

Recall

$$\text{Var} \left[ \frac{1}{B} \sum_{b=1}^B z_b \right] = \frac{1-\rho}{B} \sigma^2 + \rho \sigma^2$$

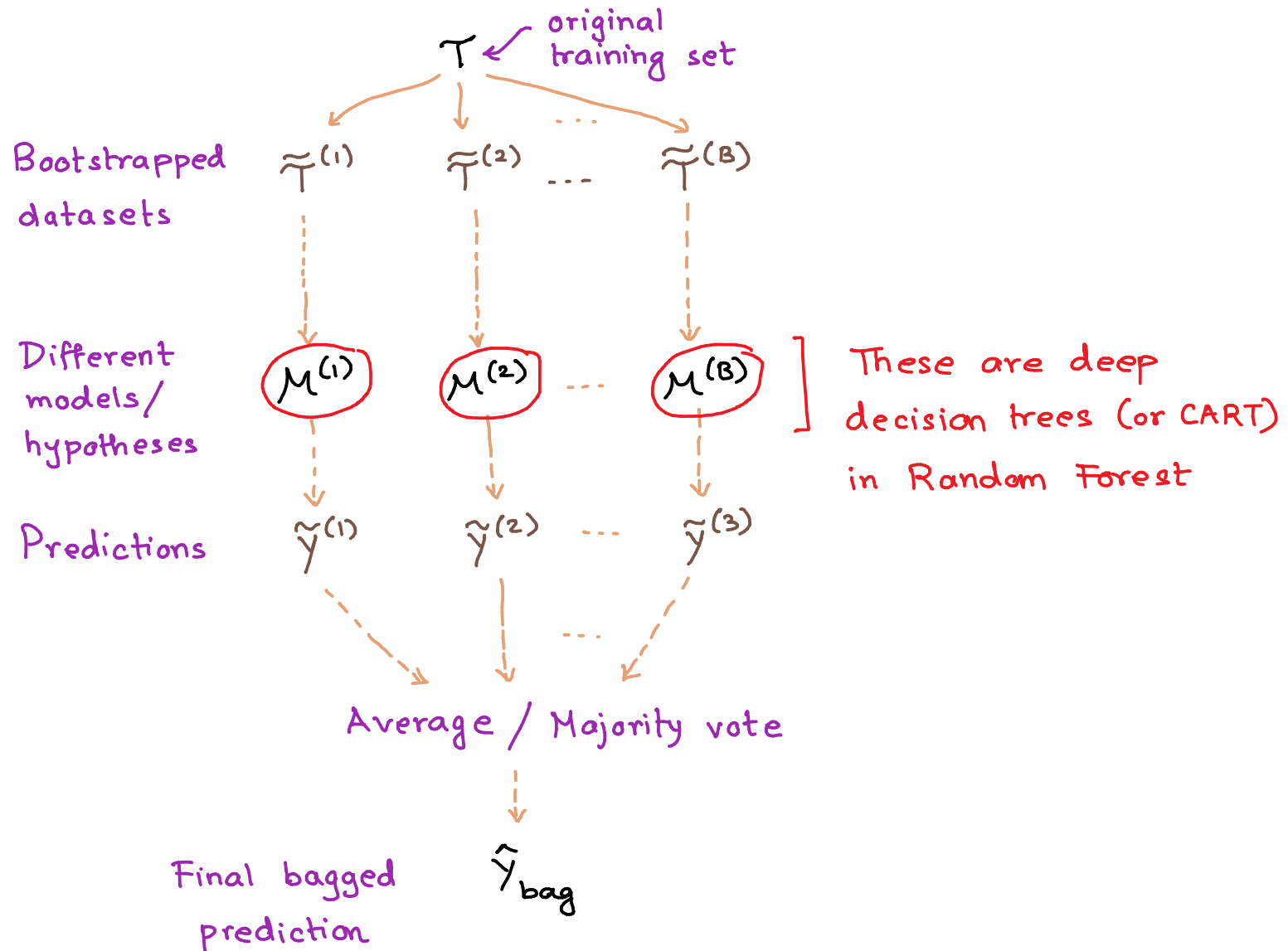
– No variance reduction when  $\rho = 1$

– Highest variance reduction when  $\rho = 0$

- **Idea of Random Forest**: De-correlate the 'B' trees by injecting additional randomness when constructing each tree

Random Forest = Bagging + Decision Trees (with random feature subset selection)

### Bagging



## Random Feature Subsets

- While growing a decision tree, one selects the best input feature  $x_j$  from all 'p' input variables  $x_1, x_2, \dots, x_p$  for splitting a node
- In random forest, we pick a random subset consisting of  $q \leq p$  features and only consider these 'q' input features for possible splits

Example

A bootstrapped dataset

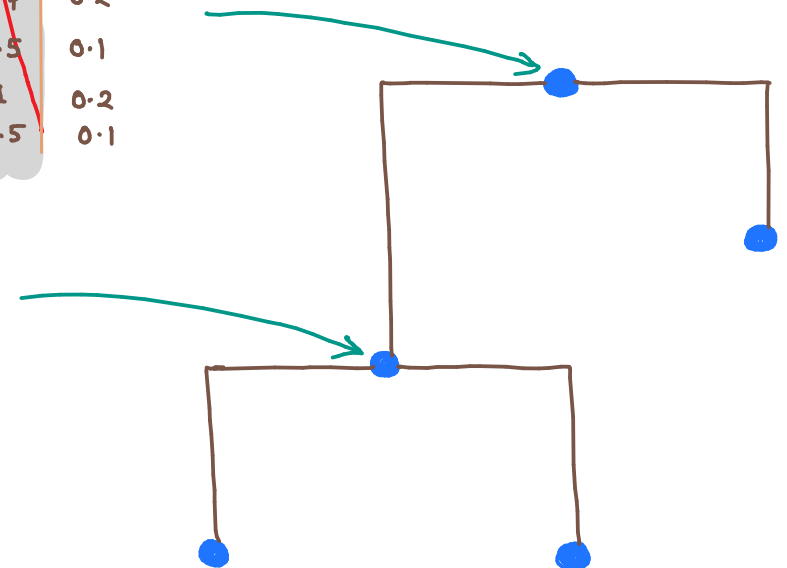
$x_1$	$x_2$	$x_3$	$x_4$
1.1	3.5	1.7	0.2
1.1	3.5	1.7	0.2
-0.9	3.3	0.5	0.1
1.2	3.2	-1	0.2
-0.9	3.3	0.5	0.1

$p = 4$  (# of inputs)

<del><math>x_1</math></del>	$x_2$	<del><math>x_3</math></del>	$x_4$
<del>1.1</del>	3.5	<del>1.7</del>	0.2
<del>1.1</del>	3.5	<del>1.7</del>	0.2
<del>-0.9</del>	3.3	<del>0.5</del>	0.1
<del>1.2</del>	3.2	<del>-1</del>	0.2
<del>-0.9</del>	3.3	<del>0.5</del>	0.1

<del><math>x_1</math></del>	<del><math>x_2</math></del>	$x_3$	$x_4$
<del>1.1</del>	<del>3.5</del>	1.7	0.2
<del>1.1</del>	<del>3.5</del>	1.7	0.2
<del>-0.9</del>	<del>3.3</del>	0.5	0.1
<del>1.2</del>	<del>3.2</del>	-1	0.2
<del>-0.9</del>	<del>3.3</del>	0.5	0.1

$q = 2$  random subsets



## Random forest algorithm

Inputs:  $\mathcal{T} = \{ \mathbf{x}^{(i)}, y^{(i)} \}_{i=1}^N$ ;  $\mathbf{x} \in \mathbb{R}^P$

for  $b=1$  to  $B$ , do (can run in parallel)

(a) Draw a bootstrap dataset  $\tilde{\mathcal{T}}^{(b)}$  of size  $N$  from  $\mathcal{T}$

(b) Grow a regression (or classification) tree by repeating the steps below, until a minimum node size is reached:

- Select a random subset consisting of  $q \leq P$  inputs
- Find the best splitting variable  $x_j$  among the 'q' selected inputs
- Split the node into two children with  $\{x_j \leq s\}$  and  $\{x_j > s\}$

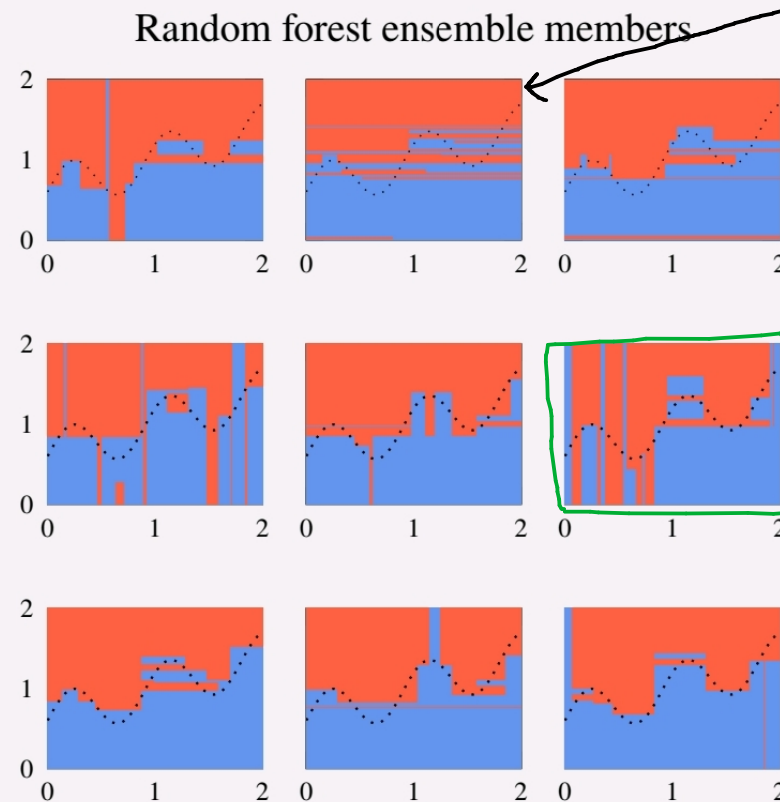
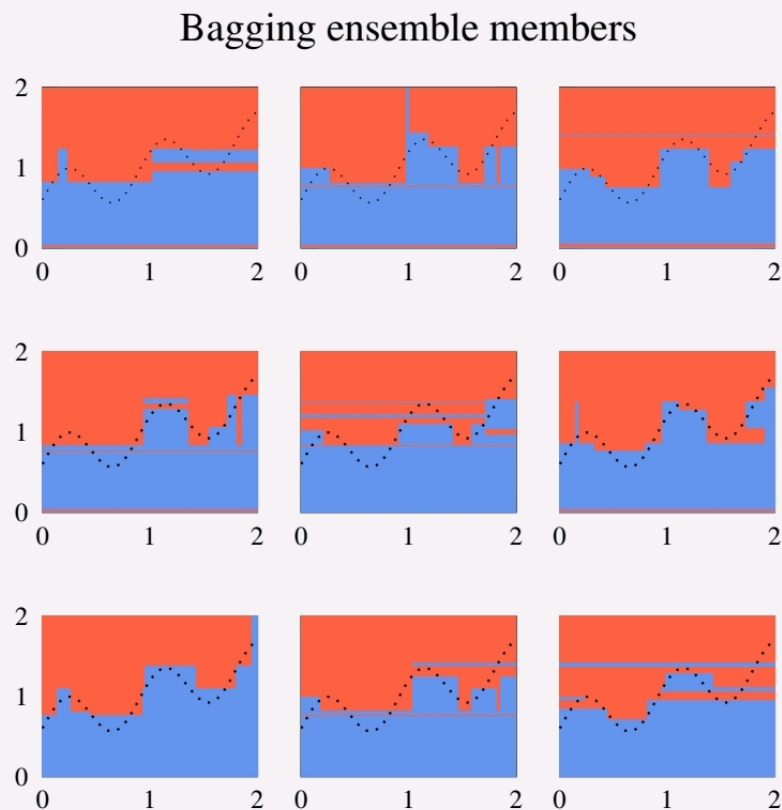
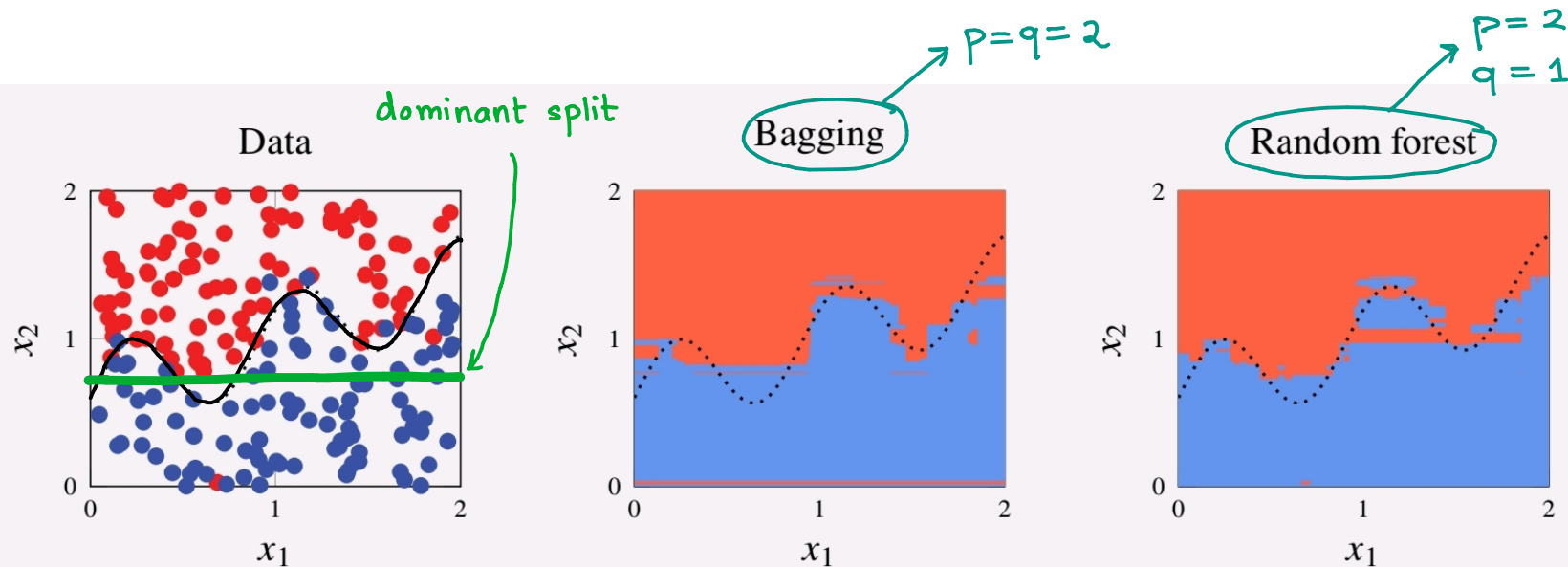
Thumb rule  
 $q = \sqrt{P}$  (for CT)  
 $q \approx P/3$  (for RT)

Final model is the average of the 'B' ensemble members

$$\hat{y}_{rf} = \frac{1}{B} \sum_{b=1}^B \tilde{y}^{(b)}$$

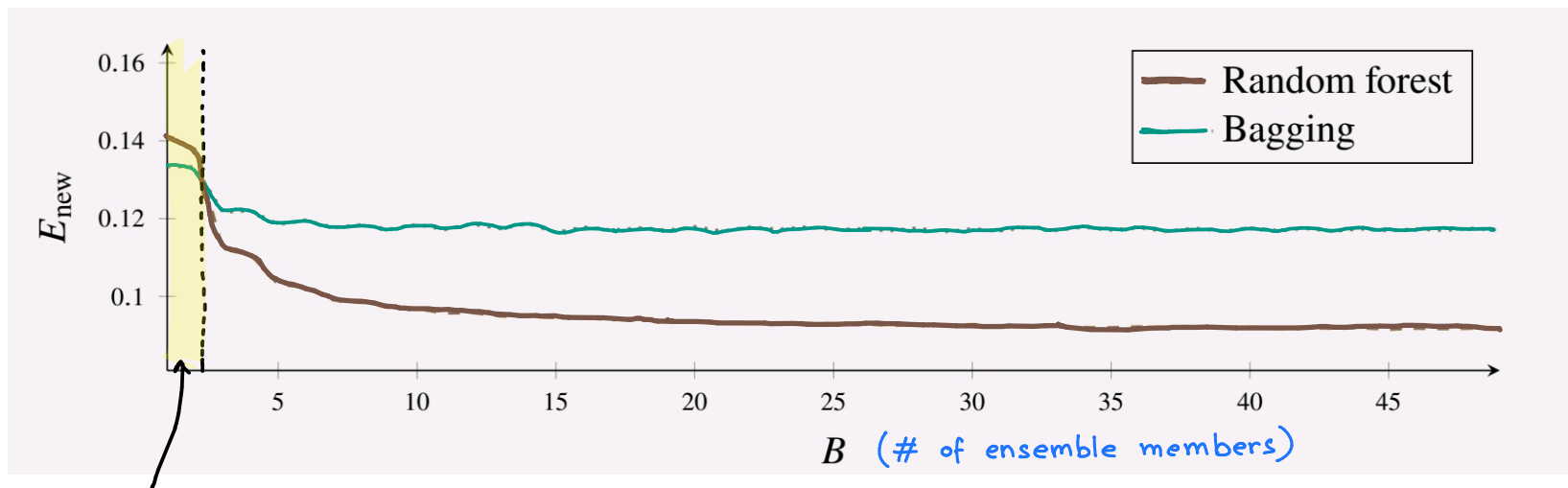
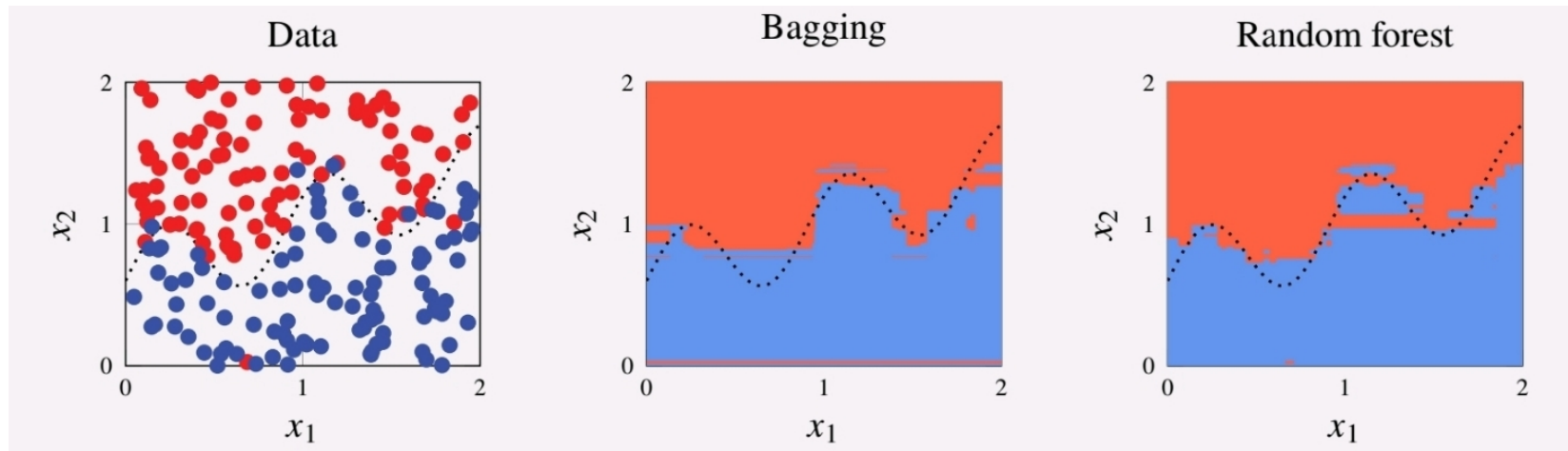
# Example of binary classification

$B = 9$  ensemble members



Random Forest Ensemble members have more individual variation than bagging

dominant splits along  $x_1$



For very small  $B$ ,  
bagging performs better  
than random forests

However, as the number of ensemble  
member increases, test error decreases  
more for random forests

- For identically distributed random variables  $\{z_b\}_{b=1}^B$

$$\text{Var} \left[ \frac{1}{B} \sum_{b=1}^B z_b \right] = \frac{1-\rho}{B} \sigma^2 + \rho \sigma^2$$

- The random input selection used in random forests:
    - increases the bias, but often very slowly ↓
    - adds to the variance ( $\sigma^2$ ) of each tree ↓
    - reduces the correlation ( $\rho$ ) between member trees ↑↑↑
  - The reduction in correlation typically has a dominant effect  
 ⇒ leads to an overall reduction in error
  - Bagging is a general technique → can be used with any base model
- Random forests consider base models as classification or regression trees