

October 6<sup>th</sup>, 2022



# Introduction to Machine Learning Part #2

Yusen He, PhD - DASIL Data Scientist  
Prof. Julia Bauder – Director of DASIL  
Martin Pollack – DASIL Post-Bac Fellow

# Intro to Machine Learning Part #2

## AGENDA

Bagging & Boosting

Model Evaluation

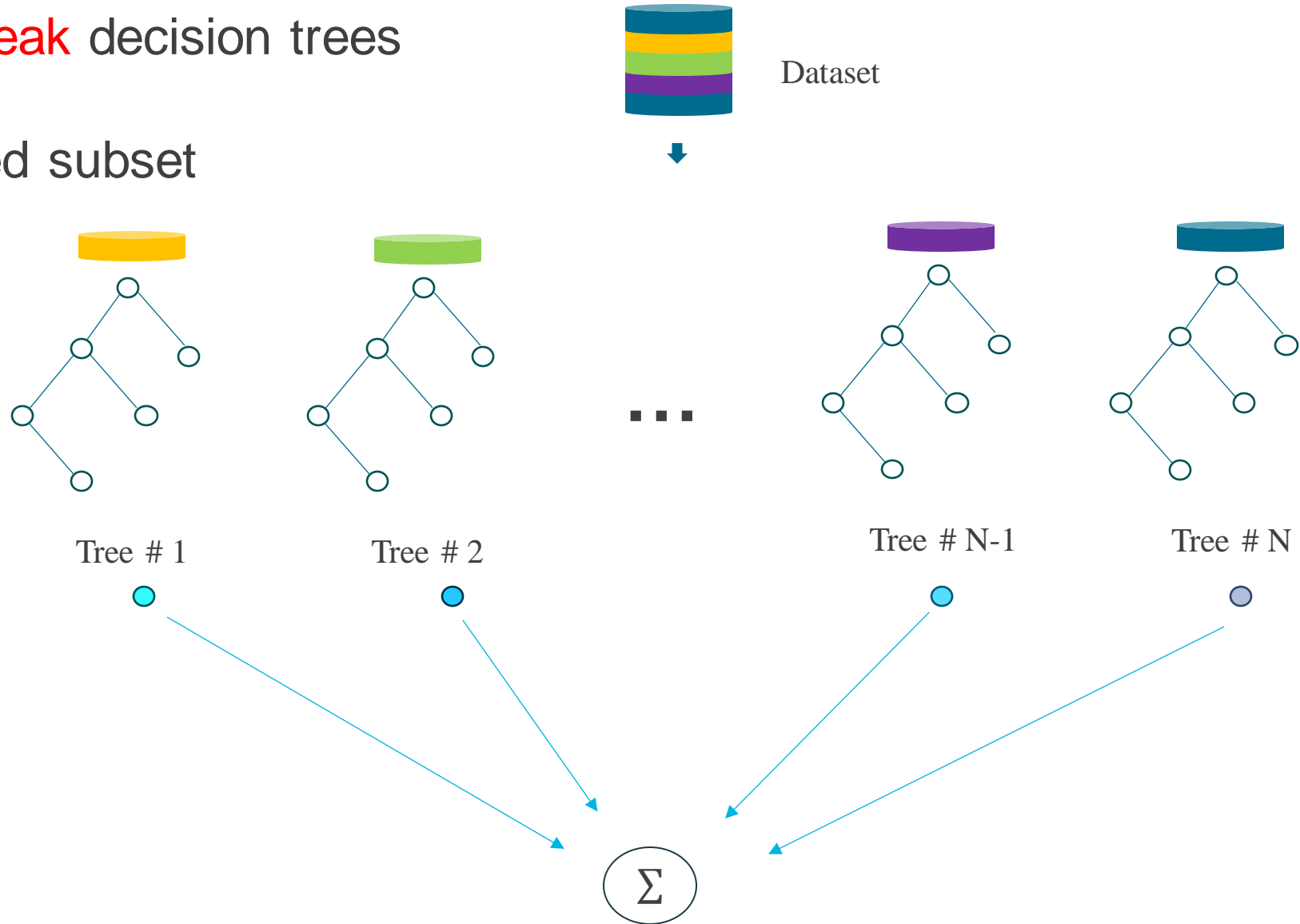
# Advanced Tree Algorithms – Random Forest

- Bagging: A set of **weak** decision trees

- Randomly sampled subset

- With **replacement**

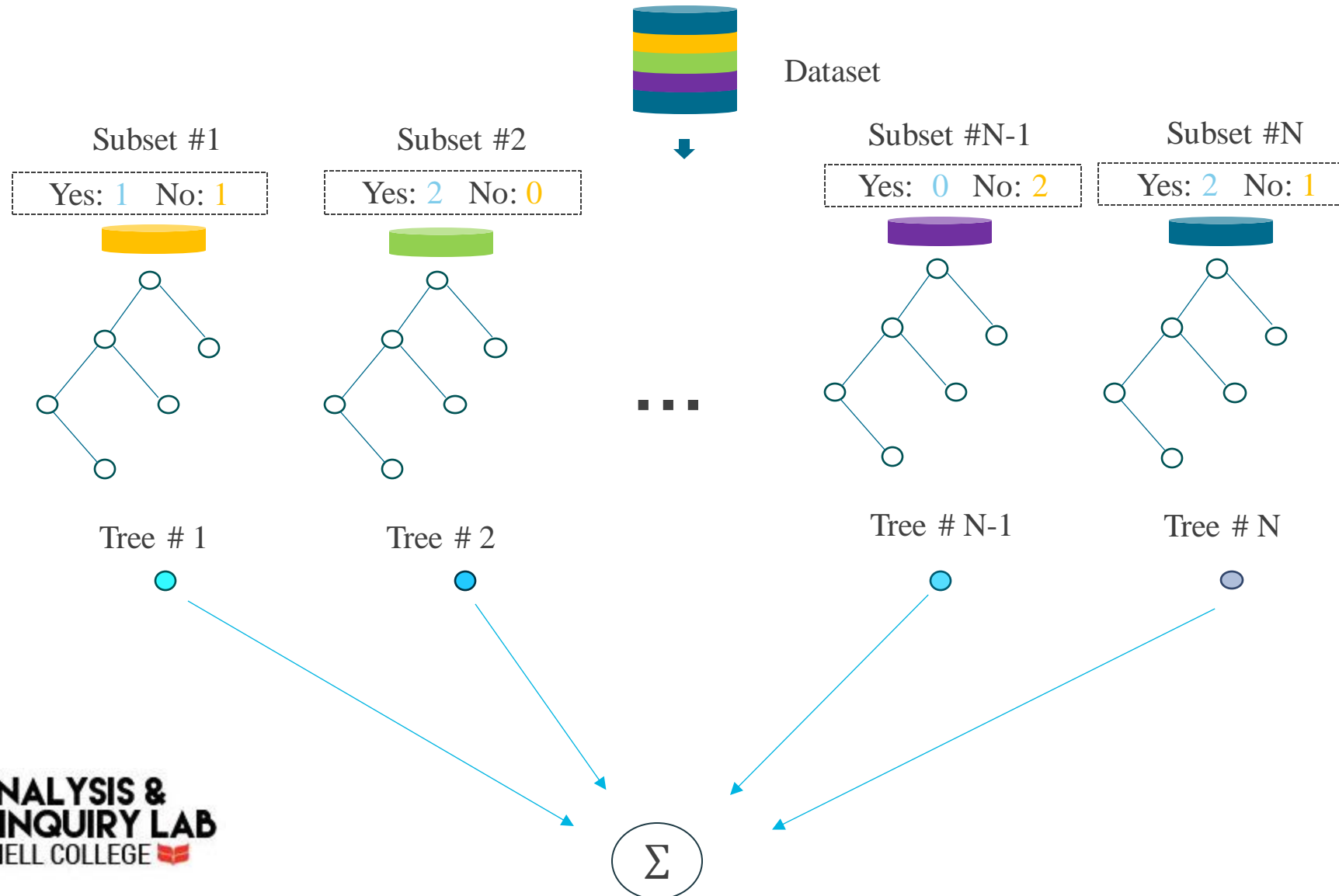
- One tree per set



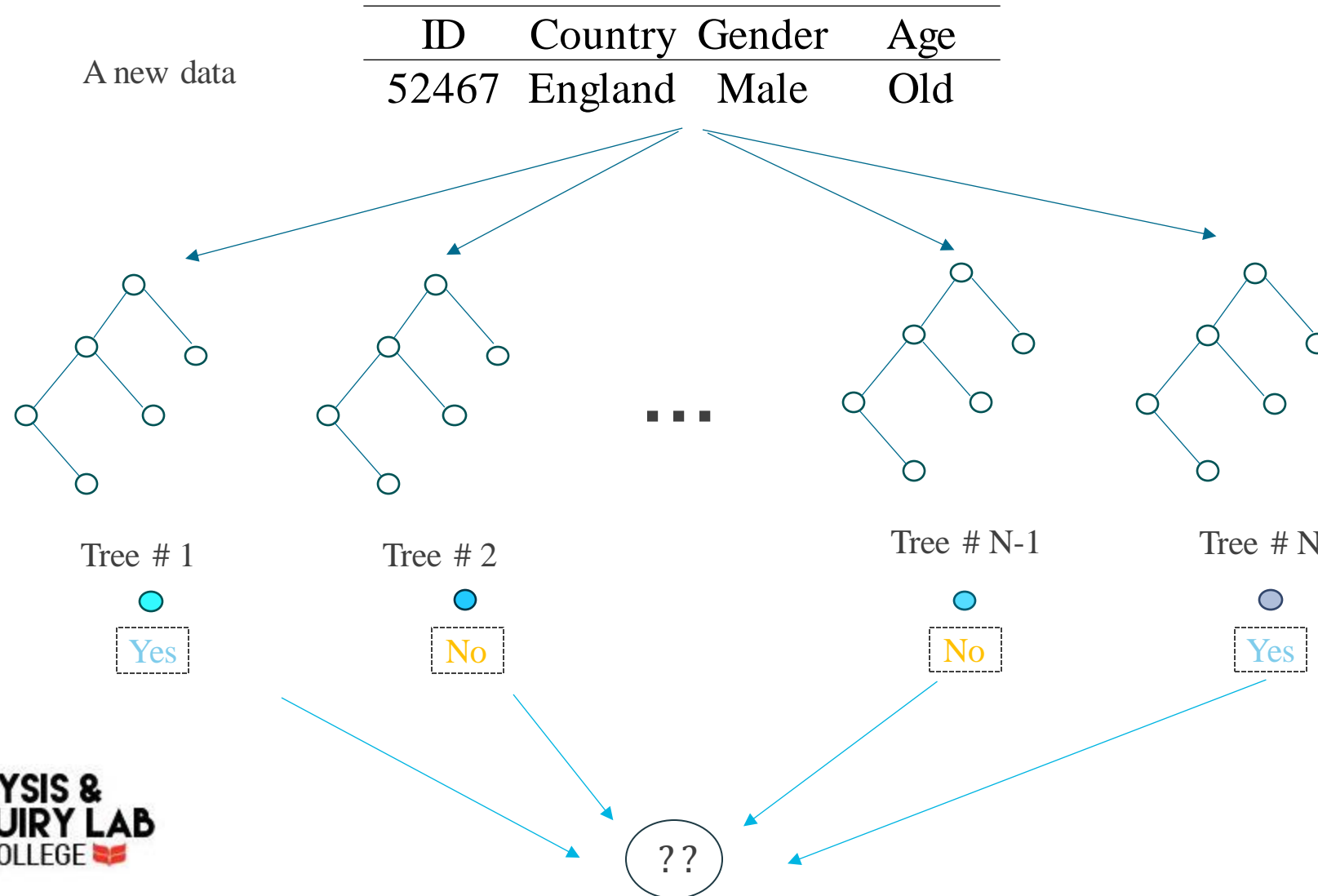
# Advanced Tree Algorithms – Random Forest

Predictor variable					Target
	ID	Country	Gender	Age	Leave
Subset #1	1	Scotland	Male	Old	No
	2	Scotland	Male	Old	No
Subset #2	3	England	Male	Young	Yes
	4	Wales	Male	Old	Yes
...	5	Wales	Female	Old	Yes
Subset #N	6	Wales	Female	Young	No

# Advanced Tree Algorithms – Random Forest



# Advanced Tree Algorithms – Random Forest



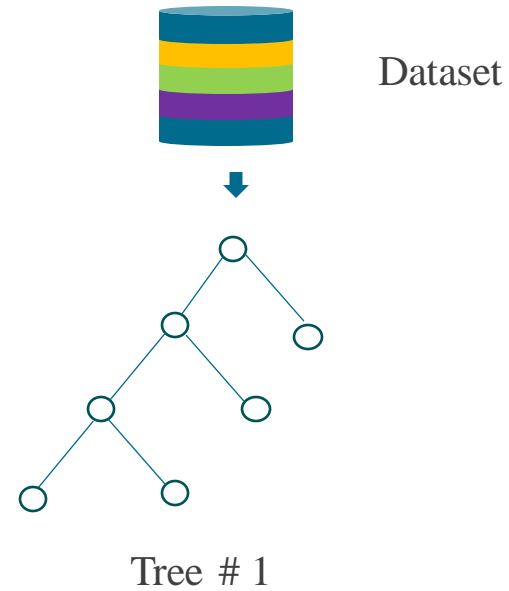
# Advanced Tree Algorithms – XGBoosting

---

- Boosting: Stacked **strong** decision trees
  - Focus on one dataset
  - Focus more on errors
  - Overfitting

# Advanced Tree Algorithms – XGBoosting

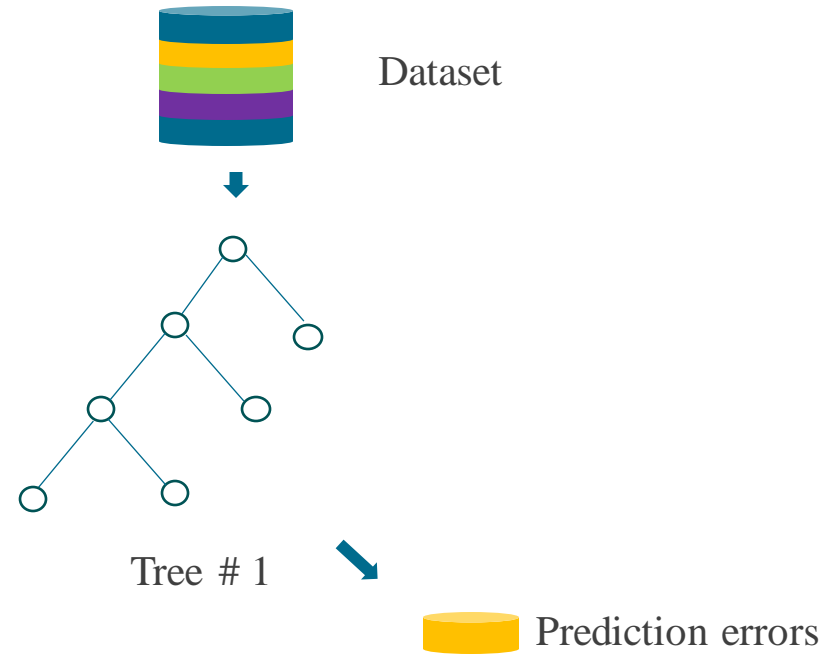
- Step 1: Fit a decision tree





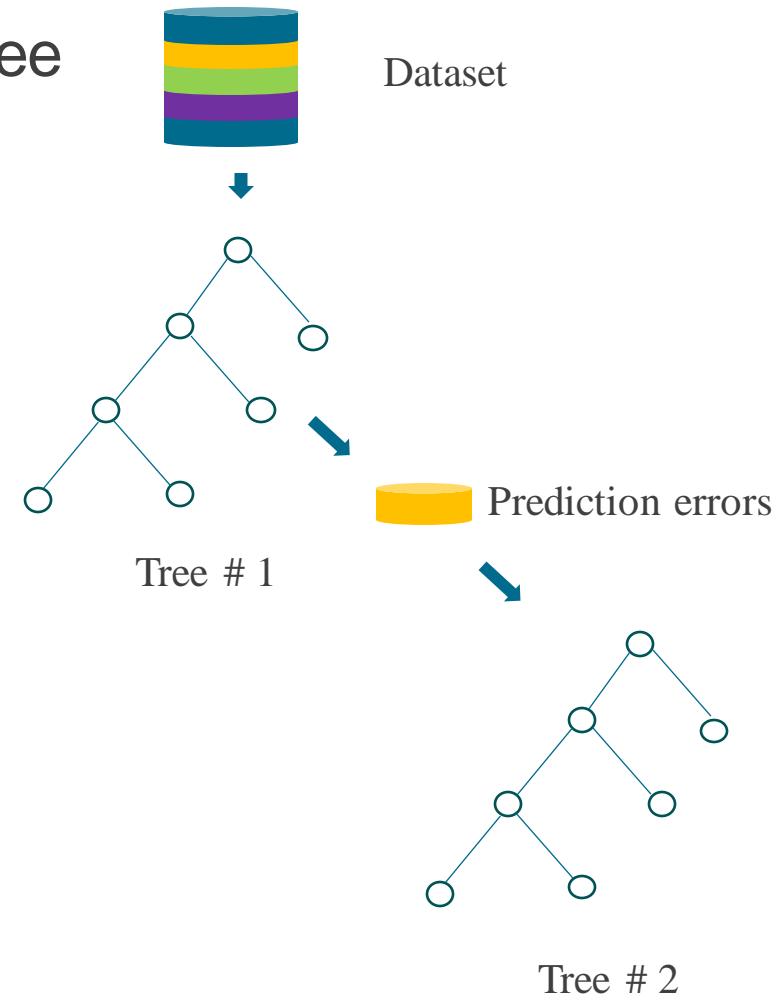
# Advanced Tree Algorithms – XGBoosting

- Step 2: Compute prediction errors



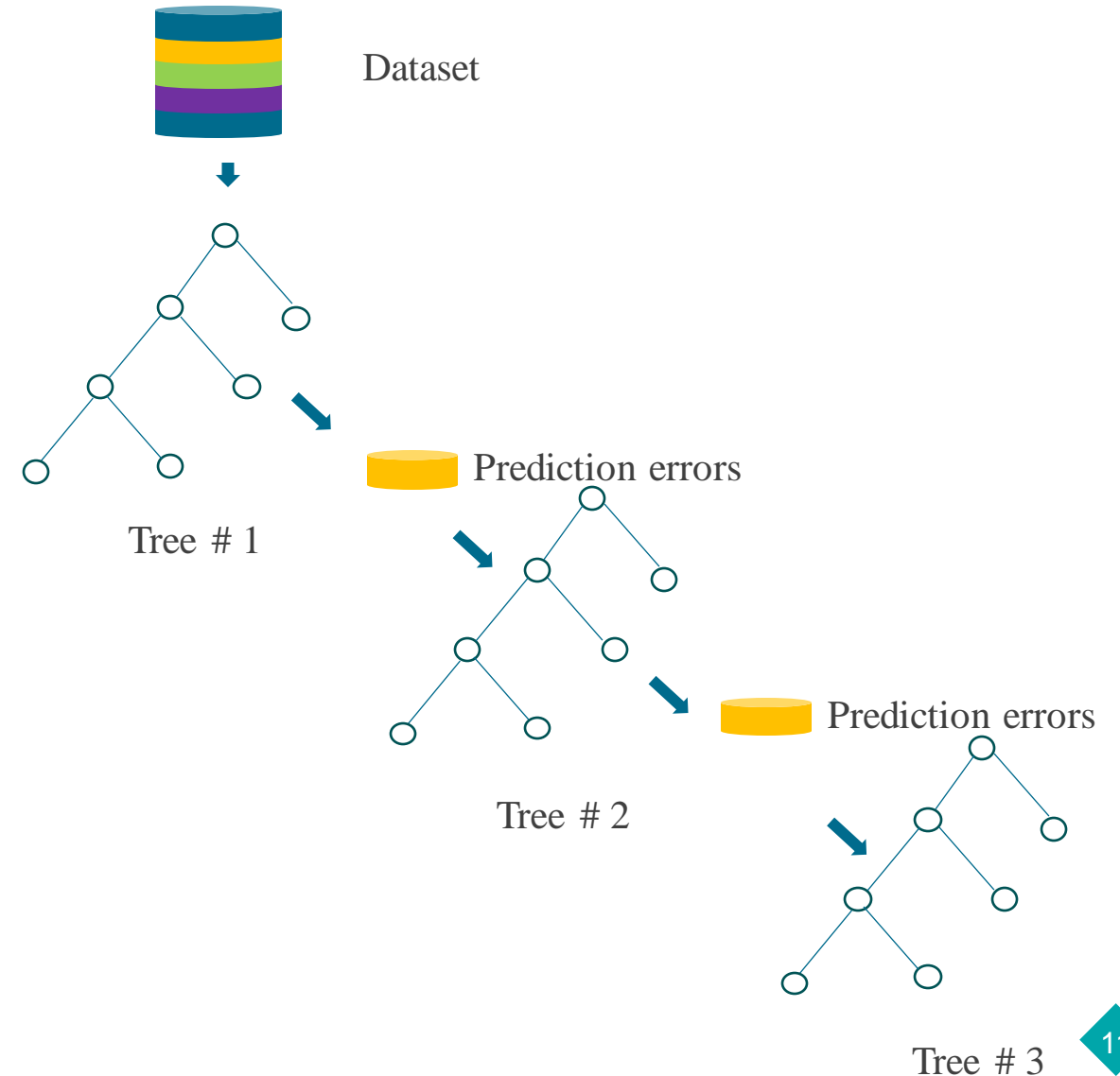
# Advanced Tree Algorithms – XGBoosting

- Step 3: Stacked with another decision tree
  - Focus on error **only**



# Advanced Tree Algorithms – XGBoosting

- Step N: Stack trees iteratively



## Pop Quiz!

---

- Talk to your neighbor about the difference between bagging and boosting.

Which one is more likely to overfitting?

Which one is likely to underfitting?

Which one has higher computational complexity?

# Intro to Machine Learning Part #2

## AGENDA

Bagging & Boosting

**Model Evaluation**

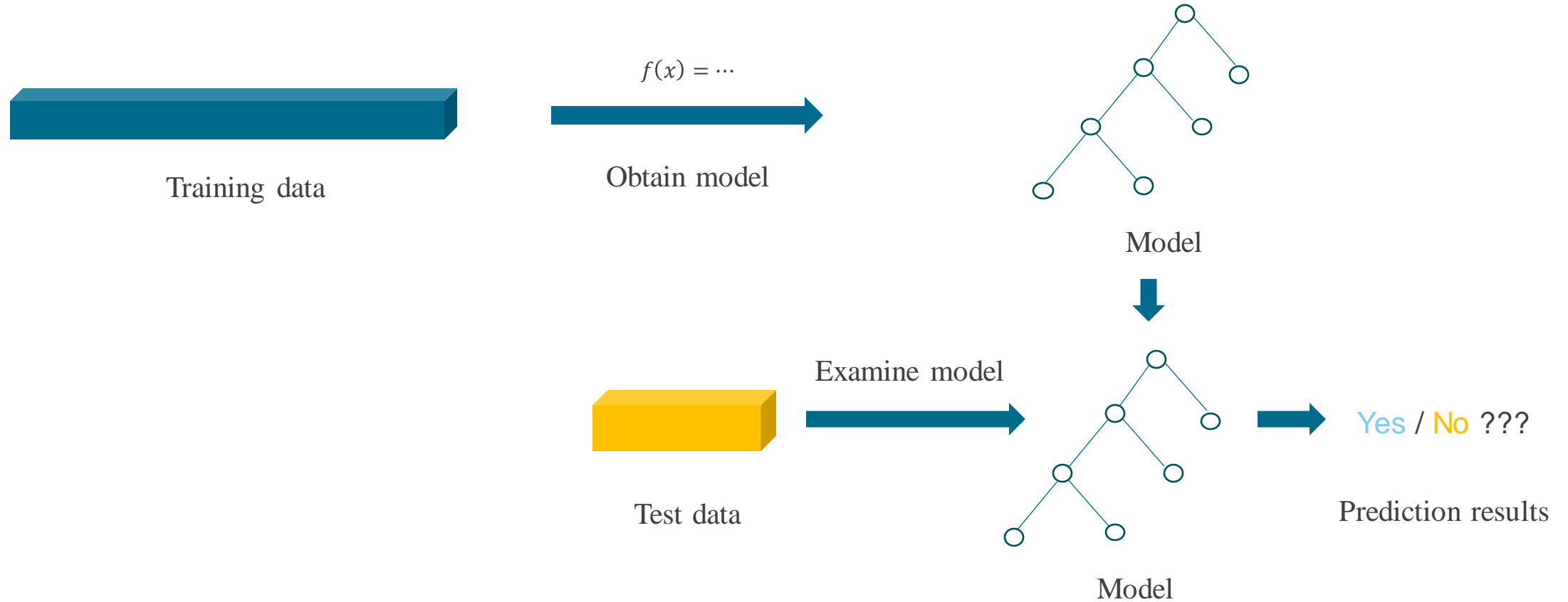
# Training / Testing Dataset Split

- Training dataset versus test dataset



# Training / Testing Dataset Split

- Training dataset versus test dataset



# Training / Testing Dataset Split

---

- Single train-test split



Training data



Test data

- Single train-test split with validation dataset



Training data



Validation data

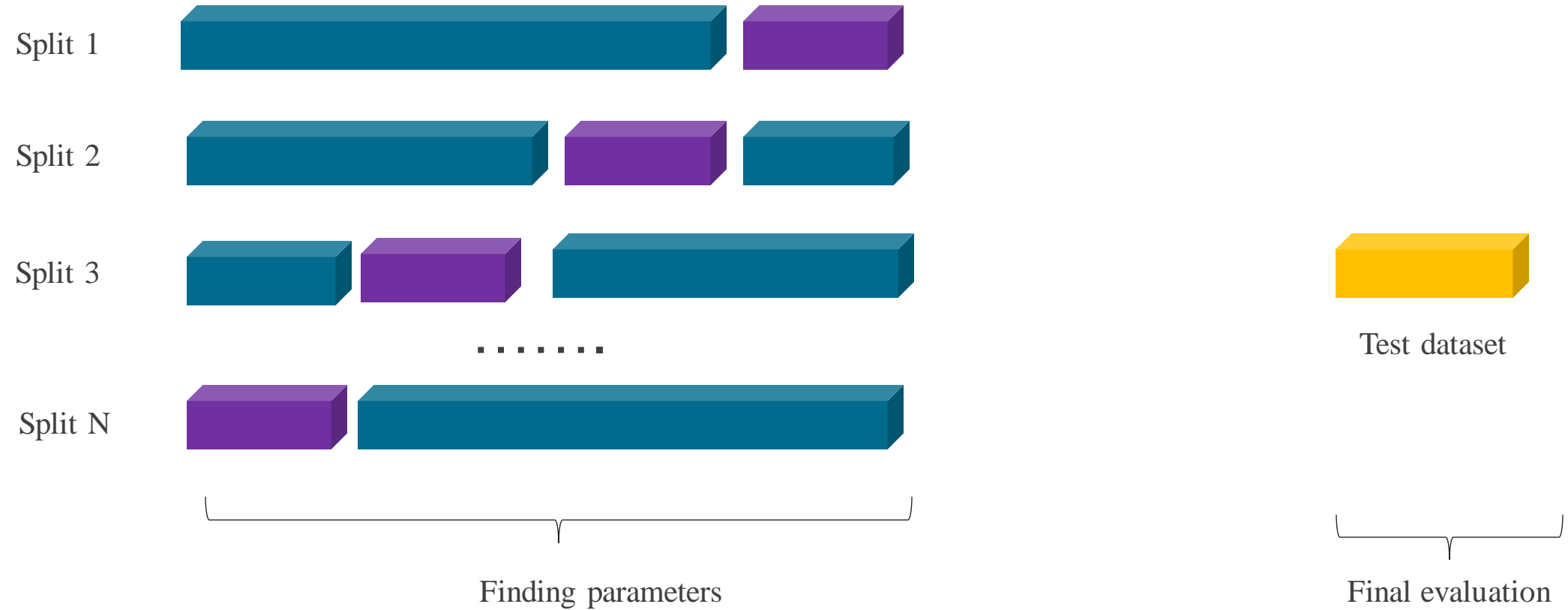


Test data



# Training / Testing Dataset Split

- Cross-validation



# Measurement Metrics - Classification

- Confusion Matrix:

		actual value		
		$p$	$n$	total
prediction outcome	$p'$	True Positive	False Positive	$P'$
	$n'$	False Negative	True Negative	$N'$
total		$P$	$N$	

# Measurement Metrics - Classification

- Confusion Matrix:

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN}$$

$$Sensitivity = \frac{TP}{\text{Actual positive}} = \frac{TP}{TP + FN}$$

$$Specificity = \frac{TN}{\text{Actual negative}} = \frac{TN}{TN + FP}$$

		actual value		
		<i>p</i>	<i>n</i>	total
prediction outcome	<i>p'</i>	True Positive	False Positive	P'
	<i>n'</i>	False Negative	True Negative	N'
total		P	N	

# Measurement Metrics - Classification

---

- Discuss:

- If I am designing a COVID-19 test, is having a higher sensitivity or specificity more important?

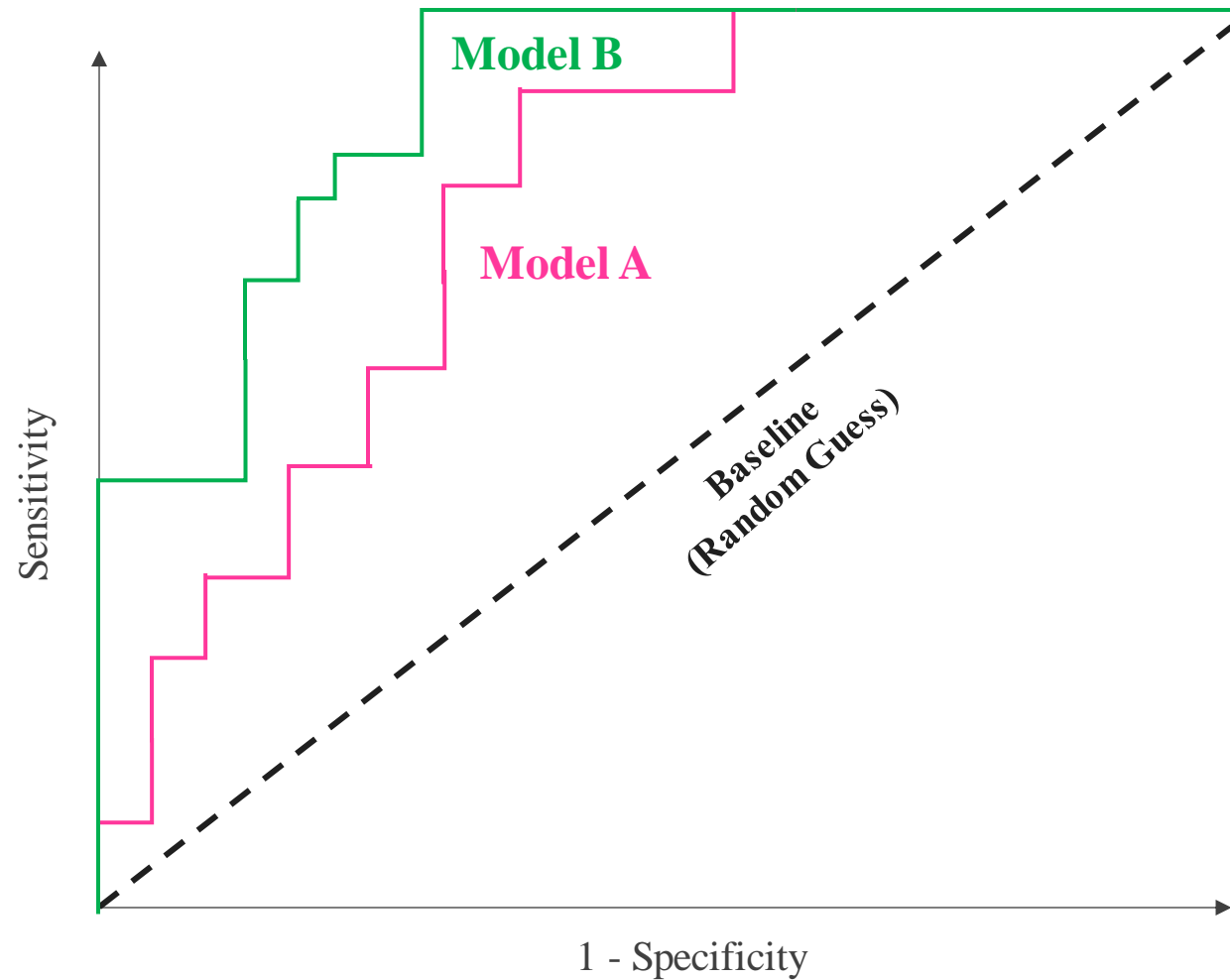
$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN}$$

$$Sensitivity = \frac{TP}{\text{Actual positive}} = \frac{TP}{TP + FN}$$

$$Specificity = \frac{TN}{\text{Actual negative}} = \frac{TN}{TN + FP}$$

# Measurement Metrics - Classification

- Area Under the ROC Curve (AUC)



# Measurement Metrics - Classification

---

- Area Under the ROC Curve
  - Measures how well the model can separate the groups
  - Shows sensitivity/specificity for each decision boundary between 0 and 1
  - Higher AUC -> better model

# Q & A

---

