

# Bagging OOB Error



CS109A Introduction to Data Science  
Pavlos Protopapas, Kevin Rader and Chris Gumb

Wenjun Liu  
Yosemite

# Outline

---

- Motivation
- Bagging
- **Out-of-bag Error**



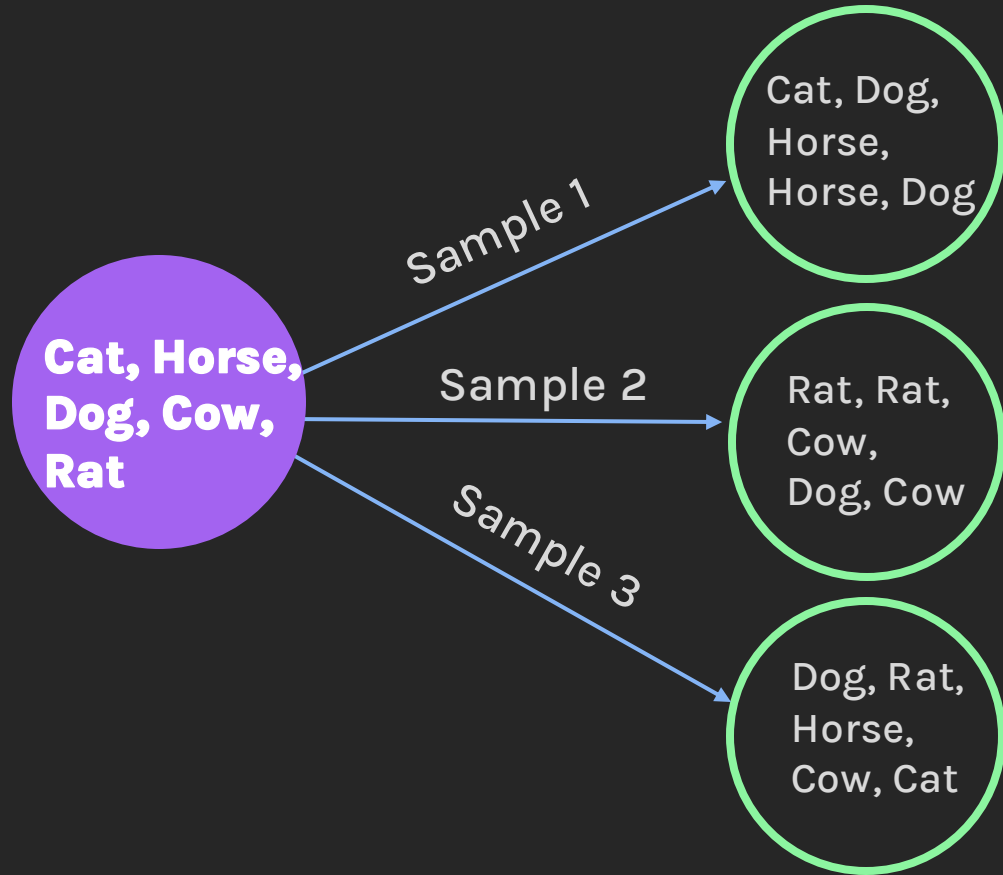
# What is OOB?

---

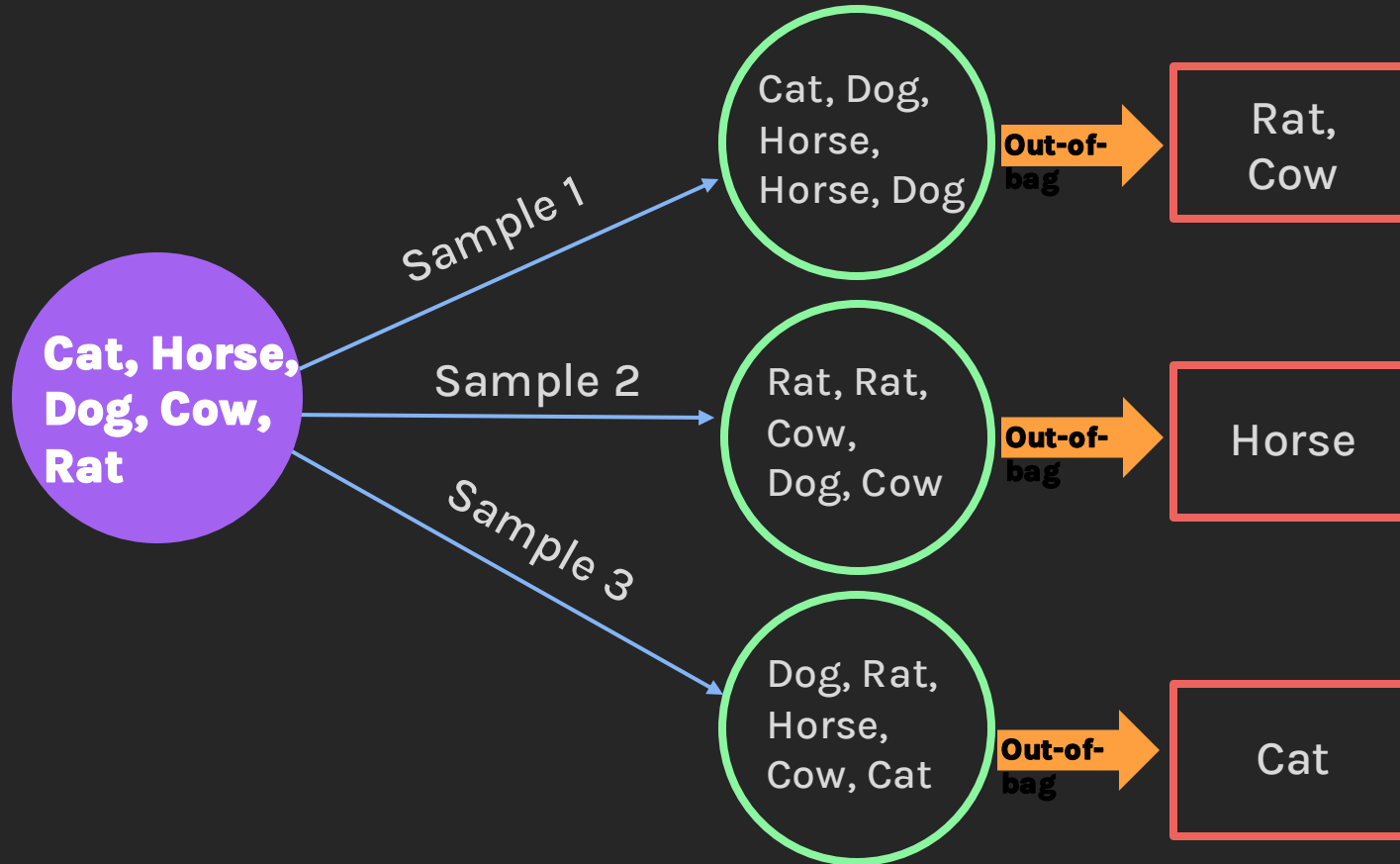


**Cat, Horse,  
Dog, Cow,  
Rat**

# What is OOB?

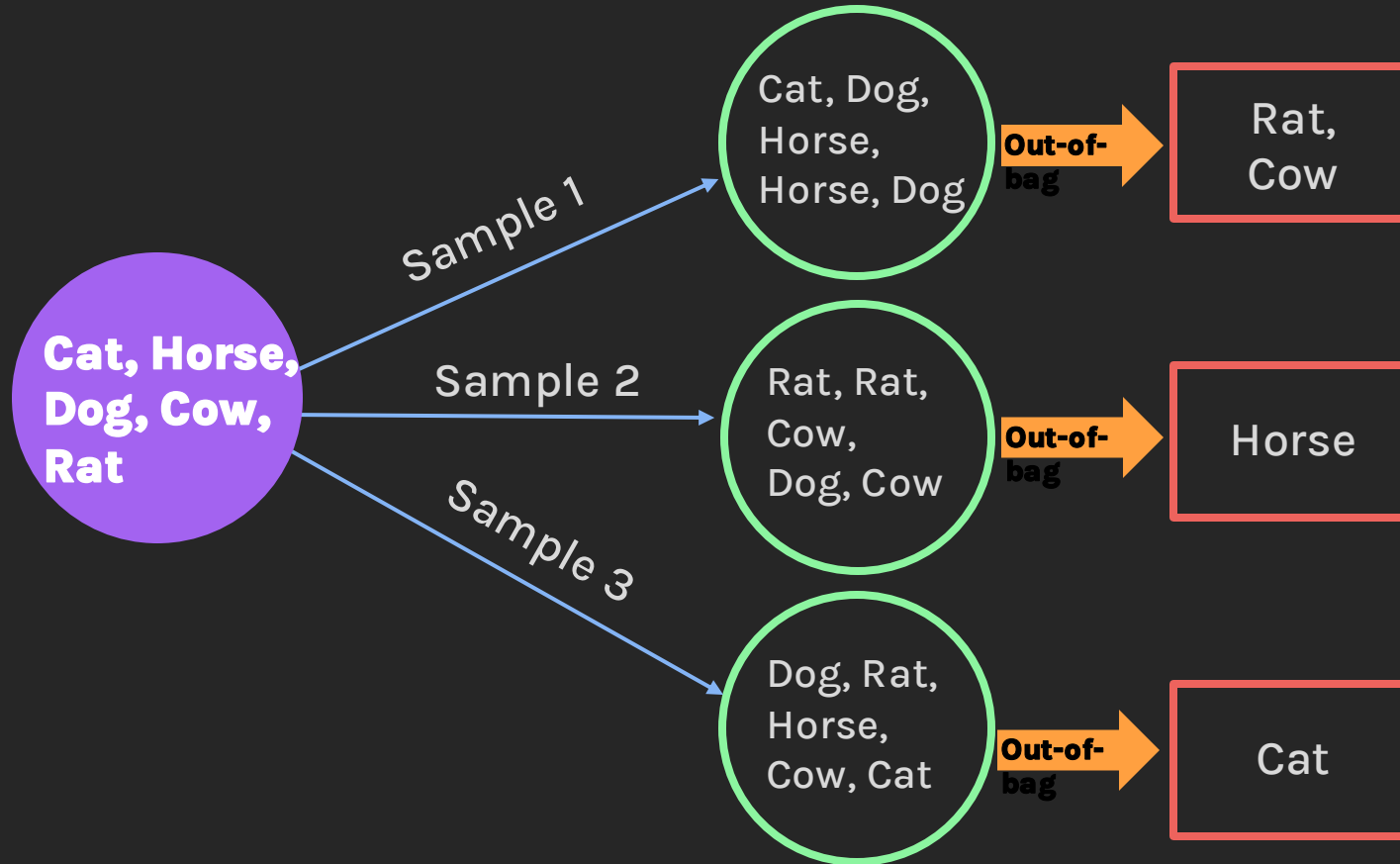


# What is OOB?



**Out-of-bag** estimate is a method of determining the prediction error whilst being trained.

# What is OOB?



**Out-of-bag** estimate is a method of determining the prediction error whilst being trained.

## Why?

- To measure generalizability.
- To replace the need for a separate measurement of performance for a validation-set performance.

Let us explore this in more details with another example

# Out-of-bag Error (OOB)

Original Data

<b>X</b>	<b>Y</b>
<b>X<sub>1</sub></b>	<b>y<sub>1</sub></b>
<b>X<sub>2</sub></b>	<b>y<sub>2</sub></b>
<b>X<sub>3</sub></b>	<b>y<sub>3</sub></b>
<b>X<sub>4</sub></b>	<b>y<sub>4</sub></b>
<b>X<sub>5</sub></b>	<b>y<sub>5</sub></b>
•	•
•	•
•	•
<b>X<sub>n</sub></b>	<b>y<sub>n</sub></b>

Predictor/Feature

PROTOPAPAS

Response/Target



# Out-of-bag Error (OOB)

Original Data

$X$	$Y$
$X_1$	$y_1$
$X_2$	$y_2$
$X_3$	$y_3$
$X_4$	$y_4$
$X_5$	$y_5$
$\cdot$	$\cdot$
$\cdot$	$\cdot$
$\cdot$	$\cdot$
$X_n$	$y_n$

*Bootstrap Sample*

1

$X$	$Y$
$X_4$	$y_4$
$X_9$	$y_9$
$X_{11}$	$y_{11}$
$X_{21}$	$y_{21}$
$X_{35}$	$y_{35}$
$\cdot$	$\cdot$
$\cdot$	$\cdot$
$\cdot$	$\cdot$
$X_k$	$y_k$



Predictor/Feature

PROTOPAPAS

Response/Target

# Out-of-bag Error (OOB)

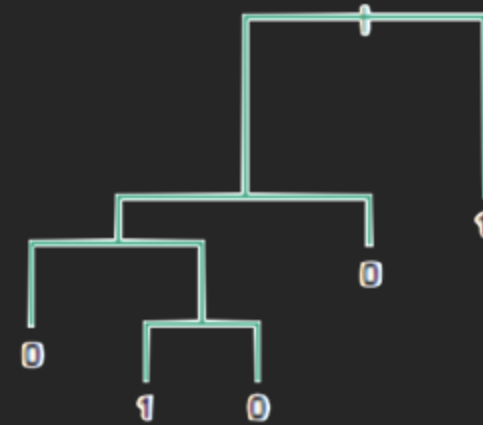
Original Data

X	Y
$x_1$	$y_1$
$x_2$	$y_2$
$x_3$	$y_3$
$x_4$	$y_4$
$x_5$	$y_5$
⋮	⋮
$x_n$	$y_n$

*Bootstrap Sample*

X	Y
$x_4$	$y_4$
$x_9$	$y_9$
$x_{11}$	$y_{11}$
$x_{21}$	$y_{21}$
$x_{35}$	$y_{35}$
⋮	⋮
$x_k$	$y_k$

Decision Tree



Response/Target

Predictor/Feature

# Out-of-bag Error (OOB)

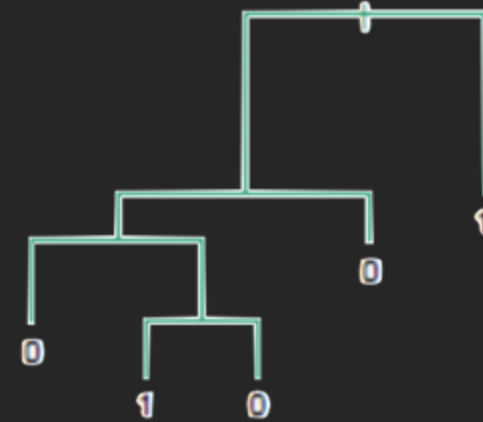
Original Data

X	Y
$x_1$	$y_1$
$x_2$	$y_2$
$x_3$	$y_3$
$x_4$	$y_4$
$x_5$	$y_5$
⋮	⋮
$x_n$	$y_n$

*Bootstrap Sample*

X	Y
$x_1$	$y_1$
$x_3$	$y_3$
$x_5$	$y_5$
$x_{21}$	$y_{21}$
$x_{35}$	$y_{35}$
⋮	⋮
$x_k$	$y_k$

Decision Tree



Used and unused data

X	Y
$x_1$	$y_1$
$x_2$	$y_2$
$x_3$	$y_3$
$x_4$	$y_4$
$x_5$	$y_5$
⋮	⋮
$x_n$	$y_n$

Response/Target

Predictor/Feature

# Out-of-bag Error (OOB)

Original Data

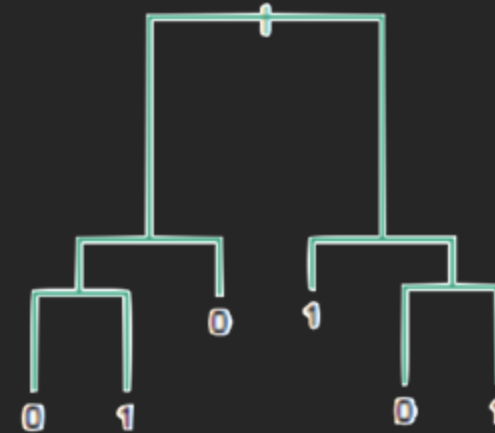
X	Y
$X_1$	$y_1$
$X_2$	$y_2$
$X_3$	$y_3$
$X_4$	$y_4$
$X_5$	$y_5$
⋮	⋮
$X_n$	$y_n$

*Bootstrap Sample*

2

X	Y
$X_5$	$y_5$
$X_7$	$y_7$
$X_{13}$	$y_{13}$
$X_{27}$	$y_{27}$
$X_{32}$	$y_{32}$
⋮	⋮
$X_k$	$y_k$

Decision Tree 2



Used and unused

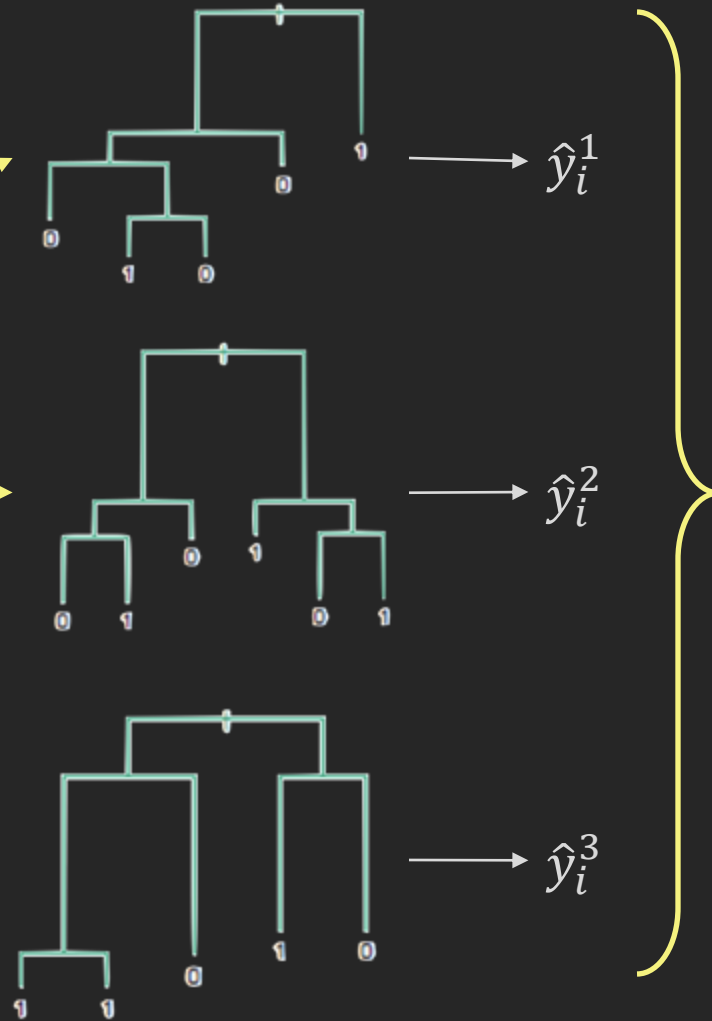
data

X	Y
$X_1$	$y_1$
$X_2$	$y_2$
$X_3$	$y_3$
$X_4$	$y_4$
$X_5$	$y_5$
⋮	⋮
$X_n$	$y_n$

# Point-wise out-of-bag error

B Trees that did not see  $\{X_i, y_i\}$

X	Y
$X_1$	$y_1$
$X_2$	$y_2$
$X_3$	$y_3$
$X_4$	$y_4$
$X_5$	$y_5$
..	..
$X_i$	$y_i$
..	..
$X_n$	$y_n$



- Identify observations the trained models have not seen
- Get the predictions for these observations from the models

# Point-wise out-of-bag error

Take **majority** for **classification** and **average** for **regression** tasks as the validation prediction for that observation

Point-wise  
prediction

$$\hat{y}_{i,pw} = \text{majority}(\hat{y}_i^j)$$

Classification

Point-wise  
out-of-bag  
error

$$e_i = \mathbb{I}(\hat{y}_{i,pw} \neq y_i)$$

Regression

$$\hat{y}_{i,pw} = \frac{1}{B} \sum_{j \in B} \hat{y}_{i,j}$$

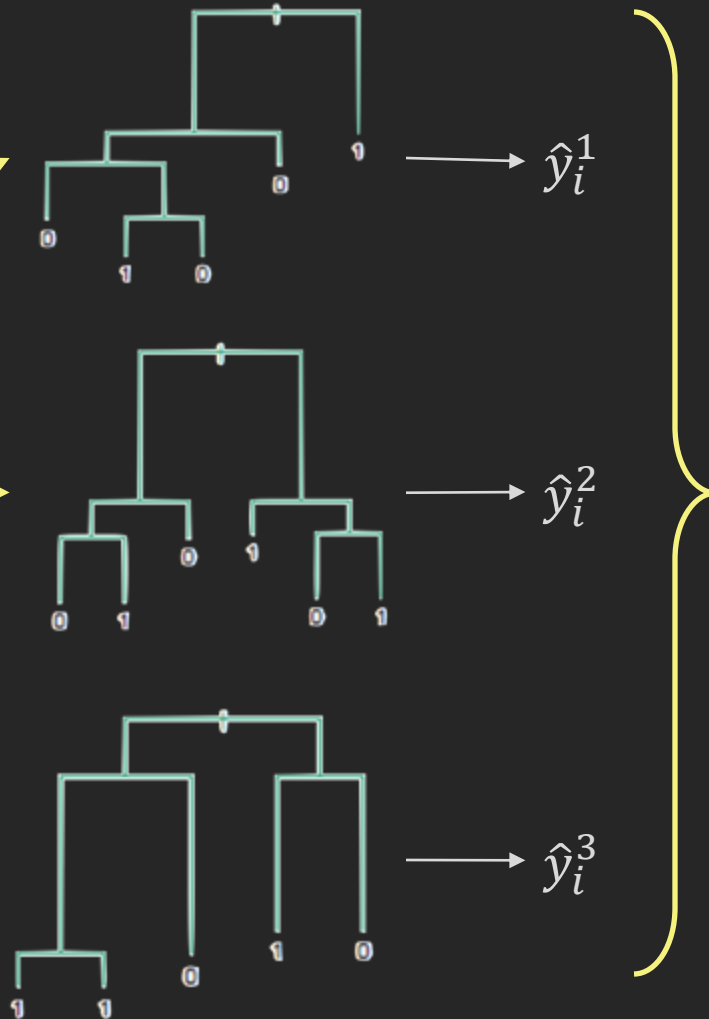
$$e_i = (y_i - \hat{y}_{i,pw})^2$$



# Point-wise out-of-bag error

X	Y
$X_1$	$y_1$
$X_2$	$y_2$
$X_3$	$y_3$
$X_4$	$y_4$
$X_5$	$y_5$
..	..
$X_i$	$y_i$
..	..
$X_n$	$y_n$

B Trees that did not see  $\{X_i, y_i\}$



Point-wise prediction

$$\hat{y}_{i,pw} = \text{majority}(\hat{y}_i^j)$$

Classification

$$e_i = \mathbb{I}(\hat{y}_{i,pw} \neq y_i)$$

Point-wise out-of-bag error

Regression

$$\hat{y}_{i,pw} = \frac{1}{B} \sum_{j \in B} \hat{y}_{i,j}$$

$$e_i = (y_i - \hat{y}_{i,pw})^2$$

# OOB Error

We average the point-wise out-of-bag errors over the full training set.

Classification

n

$$Error_{OOB} = \frac{1}{N} \sum_i^N e_i = \frac{1}{N} \sum_i^N \mathbb{I}(\hat{y}_{i,pw} \neq y_i)$$

Regression

$$Error_{OOB} = \frac{1}{N} \sum_i^N e_i = \frac{1}{N} \sum_i^N (y_i - \hat{y}_{i,pw})^2$$

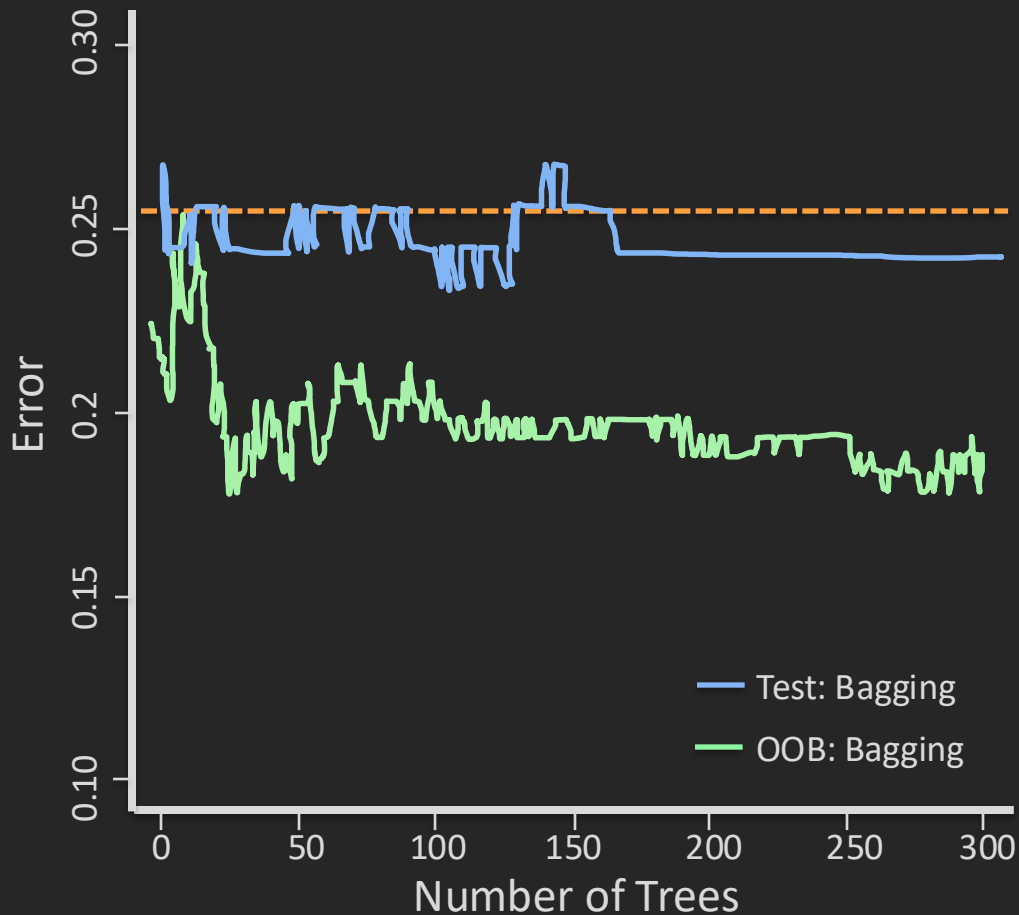
# Out-of-Bag Error: Summary

With ensemble methods, we get a new metric for assessing the predictive performance of the model, the *out-of-bag error*.

Given a training set and an ensemble of models, we compute the *out-of-bag error* by

1. For each point  $x_i$  in the training set, we average the predicted outputs  $\hat{y}_i'$ s. To do so we only use the  $B$  trees whose bootstrap training set excludes this point.
2. We compute the error of this averaged prediction. We call this the **point-wise out-of-bag error**.
3. We average the point-wise out-of-bag error over the full training set  $N$ .

# Why OOB Error? COMPARING OOB AND CROSS VALIDATION



- While using the cross-validation technique, every validation set has already been seen in training by a few decision trees and hence there is a **leakage of data**.
- OOB Error prevents leakage and yields a better model with lower variance or less overfitting.
- There is also **lesser computational cost** for OOB as compared to CV for bagging.

# Bagging

---

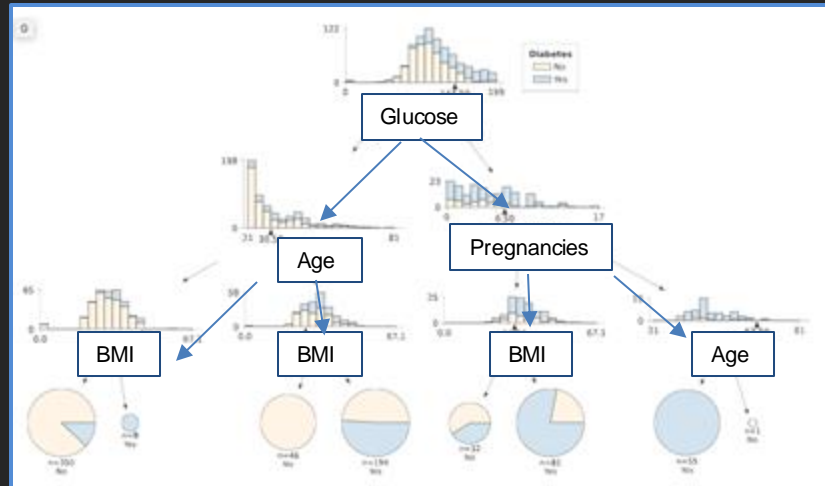
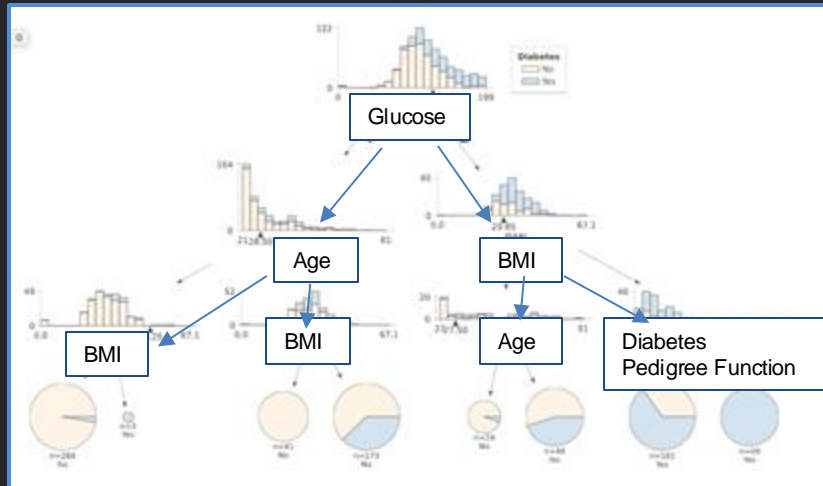
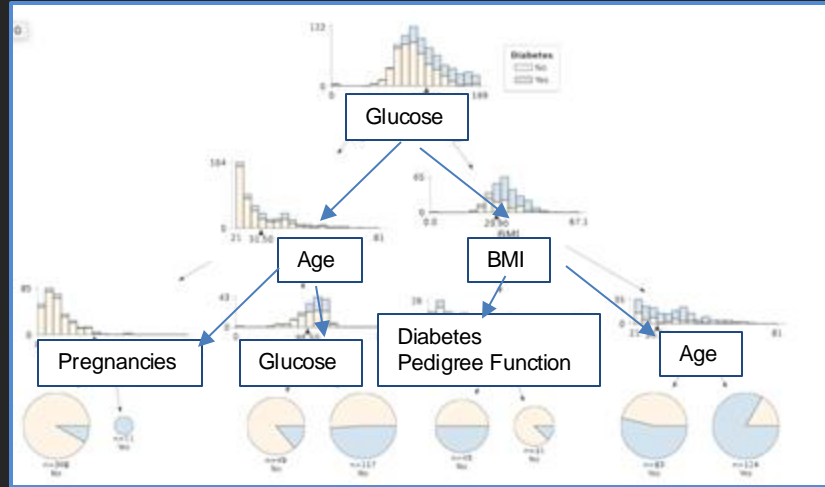
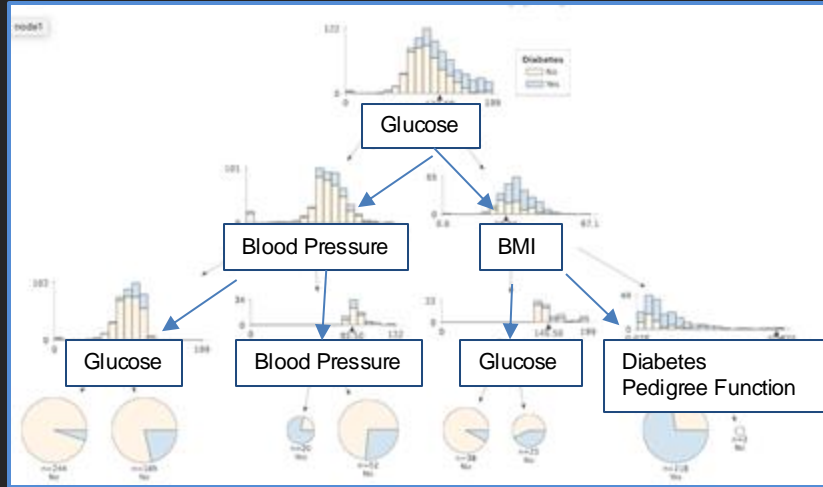
## Interpretability:

Still an issue and we will address this later.

If the individual trees are too shallow, the ensembled model can still **underfit**. Even if we combine many underfitting trees we will still underfit.

If the individual trees are too large, the ensembled model could still overfit.

# Drawbacks of Bagging



For each bootstrap, we build a decision tree.

Created by: Dr. Rahul Dave



# Improving on Bagging

In practice, the trees in Bagging tend to be **highly correlated**.

- Suppose we have an **extremely strong predictor**,  $x_j$ , in the training set amongst **moderate predictors**. Then the greedy learning algorithm ensures that most of the models in the ensemble will choose to split on  $x_j$  in early iterations.
- However, we assumed (or hope) that each tree in the ensemble is **independently and identically distributed**.

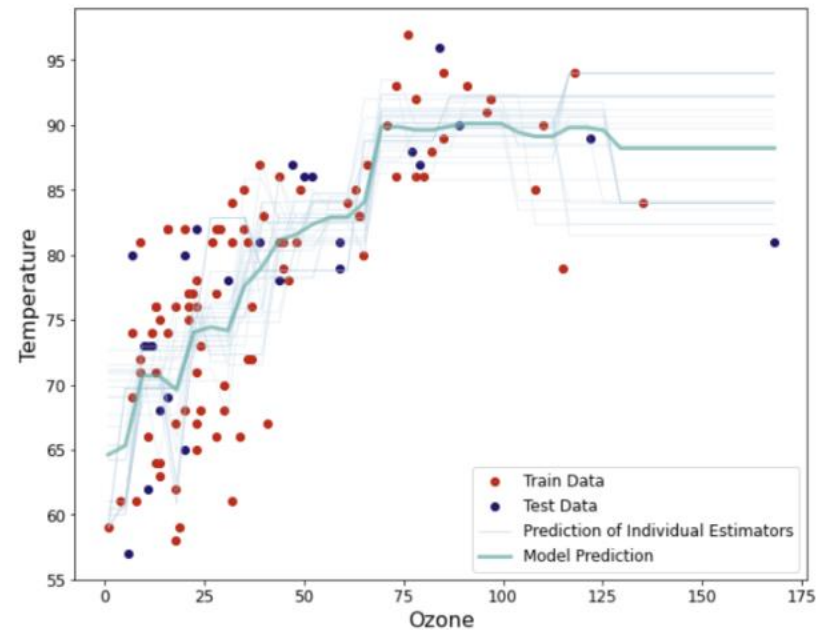
Next Monday, on 'Tree Mysteries Unveiled': Can trees ever truly be independent? 🌳 The secrets unraveled! Tune in and unlock the enigma... Only at the Monday Lecture!"





# Exercise: Regression with Bagging

The aim of this exercise is to understand regression using Bagging.



Thank you

