

XGBoost and Competitors: Tree-Based Methods for Tabular Data

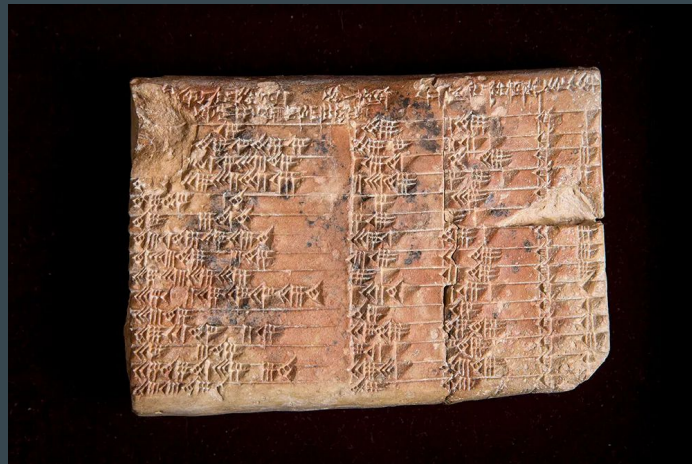
...

28.11.2025

Adéla Kubíková

What is tabular data?

- most common type
- transactional data
- heterogeneous
- no local structure
- quality of datasets
- deep learning?



Customer ID	Age	Gender	Income (\$)	Purchased
101	25	Male	40000	No
102	32	Female	55000	Yes
103	45	Male	72000	No
104	28	Female	60000	Yes
105	29	Male	50000	Yes

Contents

Motivation - why tabular data?

Decision trees

Ensemble methods - Bagging and Boosting

XGBoost, LightGBM and CatBoost

Deep neural models

Performance comparison

Summary

Evolution of Tree Algorithm



Decision
Trees

Bagging

Random
Forest

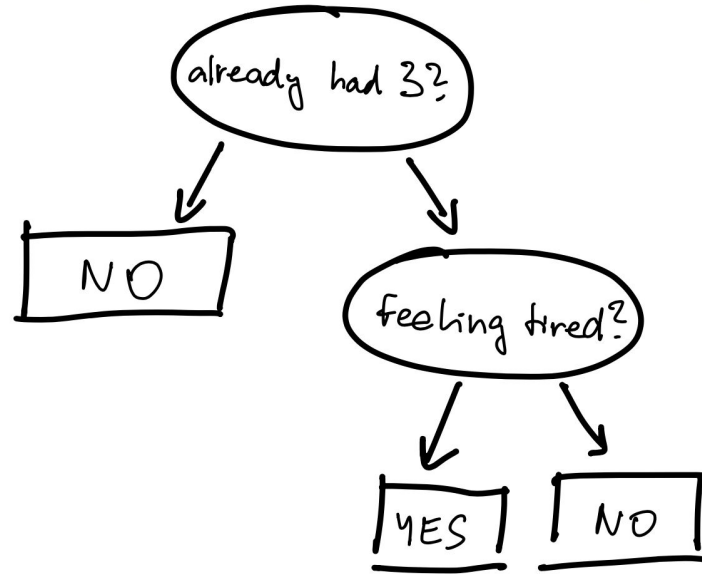
Boosting

Gradient
Boosting

XG-Boost

Base learner: Decision trees

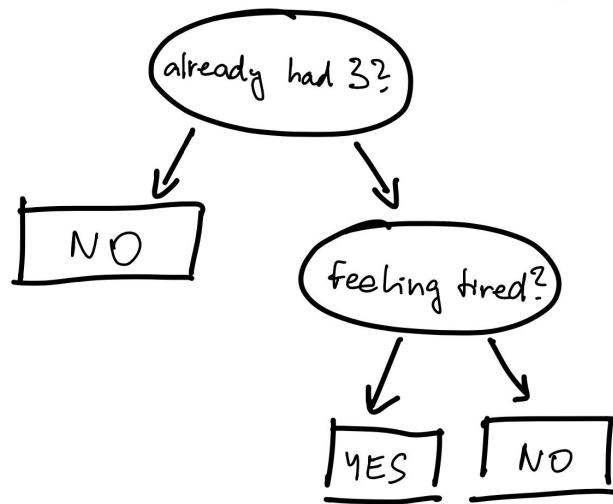
→ Should I have some coffee?



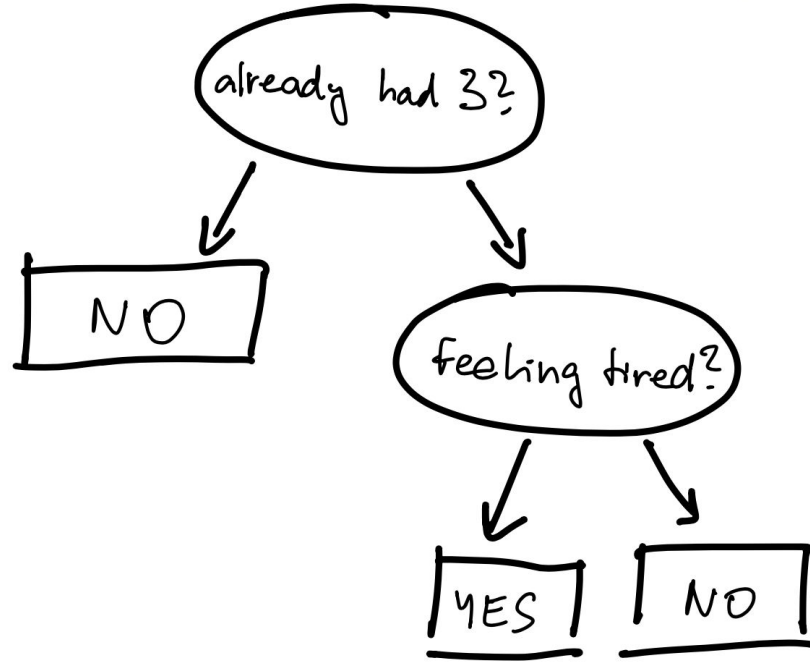
What is a decision tree?

- simple model
- set of rules leading to a decision
- feature \leq threshold
- classification trees
- regression trees

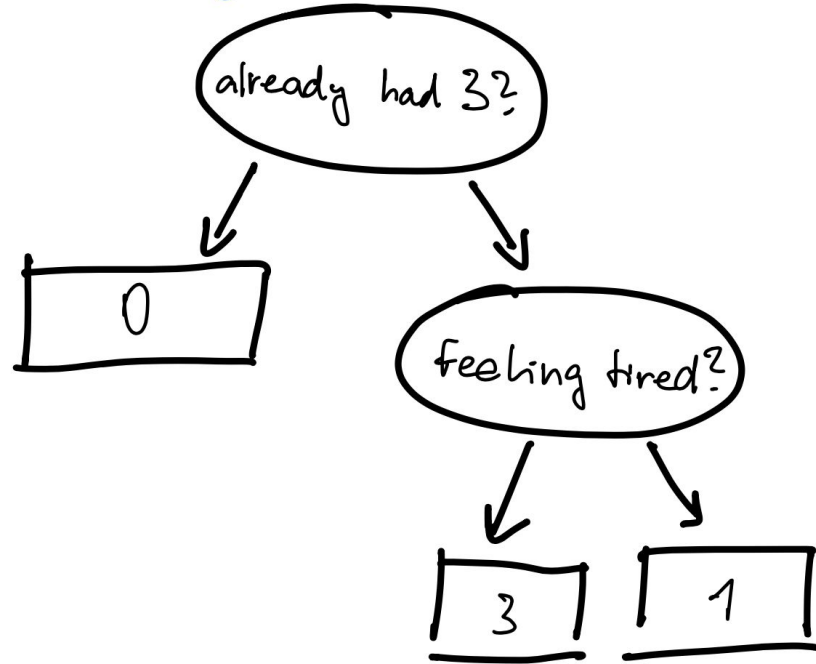
→ Should I have some coffee?



→ Should I have some coffee?



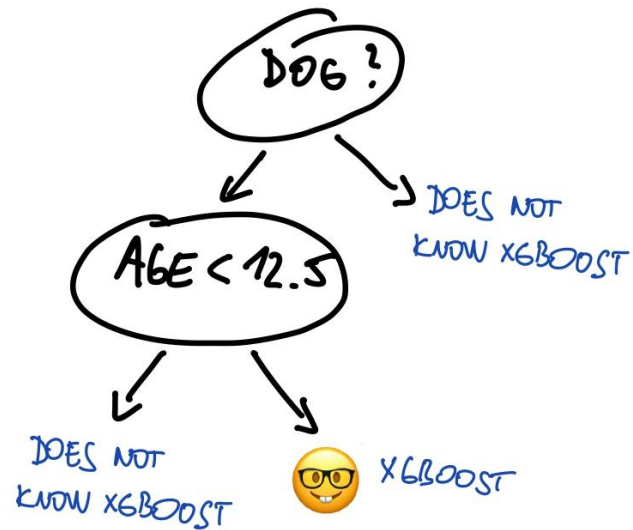
→ How many cups should I drink?



How to build a decision tree?

- classification: minimizing Gini Impurity (Entropy, Information Gain)
- regression: minimizing prediction error (MSE, MAE)

Is a woman	Has a dog	Age	Knows XGBoost
YES	YES	7	NO
YES	NO	12	NO
NO	YES	18	YES
NO	YES	35	YES
YES	YES	38	YES
YES	NO	50	NO
NO	NO	83	NO



Is a woman	Has a dog	Age	Knows XGBoost
YES	YES	7	NO
YES	NO	12	NO
NO	YES	18	YES
NO	YES	35	YES
YES	YES	38	YES
YES	NO	50	NO
NO	NO	83	NO



Is a woman	Has a dog	Age	Knows XGBoost
YES	YES	7	NO
YES	NO	12	NO
NO	YES	18	YES
NO	YES	35	YES
YES	YES	38	YES
YES	NO	50	NO
NO	NO	83	NO

Woman

Dog

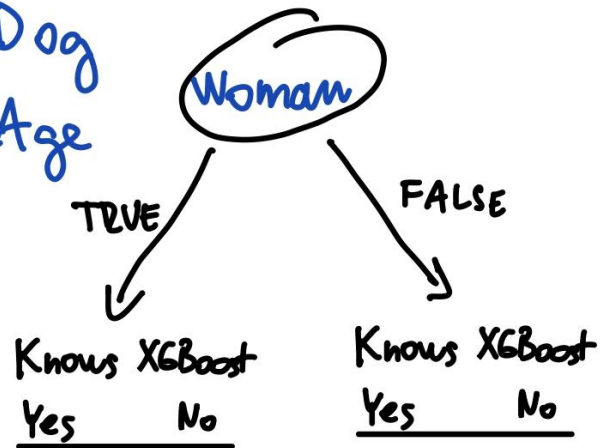
Age



Is a woman	Has a dog	Age	Knows XGBoost
YES	YES	7	NO
YES	NO	12	NO
NO	YES	18	YES
NO	YES	35	YES
YES	YES	38	YES
YES	NO	50	NO
NO	NO	83	NO

~~Woman~~

Dog
Age

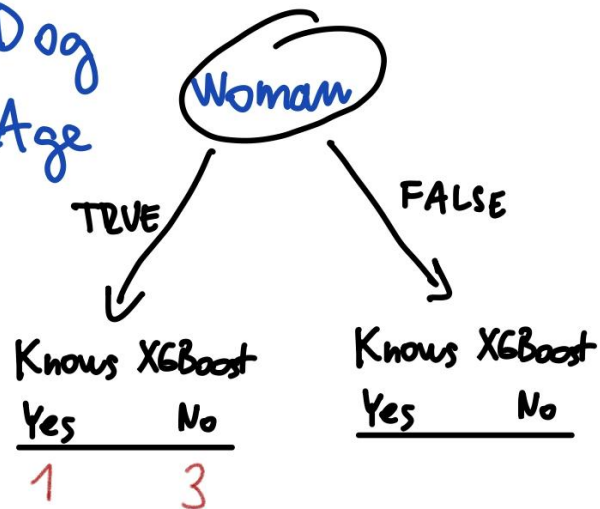


↓ ↓

Is a woman	Has a dog	Age	Knows XGBoost
YES	YES	7	NO
YES	NO	12	NO
NO	YES	18	YES
NO	YES	35	YES
YES	YES	38	YES
YES	NO	50	NO
NO	NO	83	NO

~~Woman~~

Dog
Age

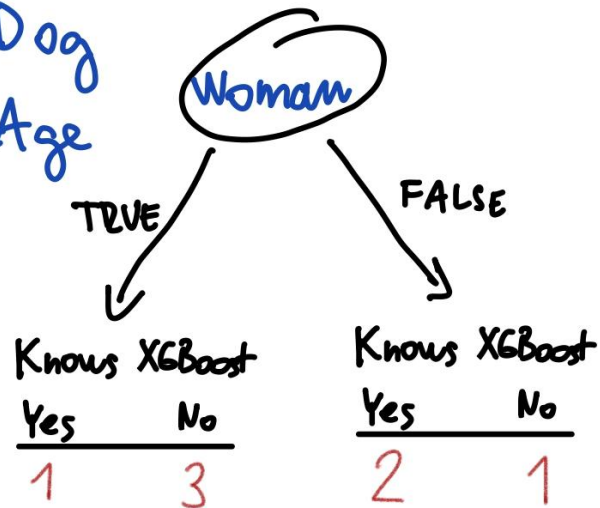


↓ ↓

Is a woman	Has a dog	Age	Knows XGBoost
YES	YES	7	NO
YES	NO	12	NO
NO	YES	18	YES
NO	YES	35	YES
YES	YES	38	YES
YES	NO	50	NO
NO	NO	83	NO

~~Woman~~

Dog
Age

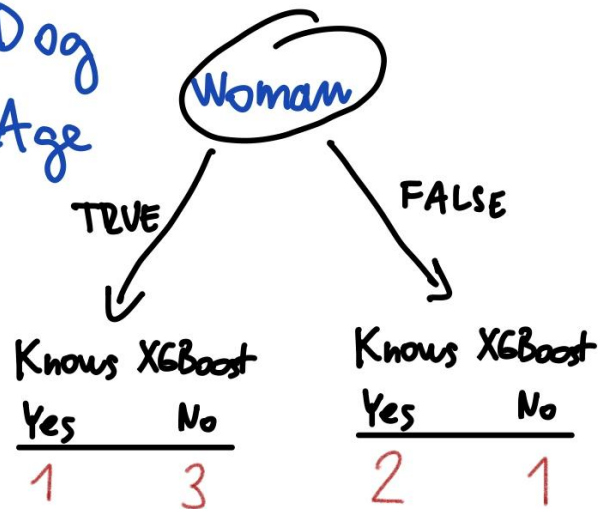


↓ ↓


Is a woman	Has a dog	Age	Knows XGBoost
YES	YES	7	NO
YES	NO	12	NO
NO	YES	18	YES
NO	YES	35	YES
YES	YES	38	YES
YES	NO	50	NO
NO	NO	83	NO

~~Woman~~

Dog
Age

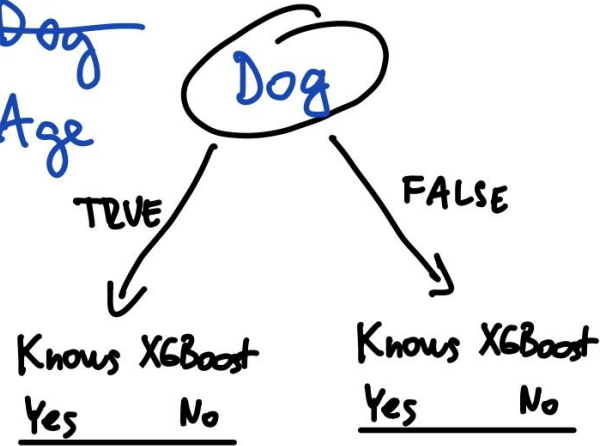



Impure leaves



Is a woman	Has a dog	Age	Knows XGBoost
YES	YES	7	NO
YES	NO	12	NO
NO	YES	18	YES
NO	YES	35	YES
YES	YES	38	YES
YES	NO	50	NO
NO	NO	83	NO

~~Woman~~
~~Dog~~
 Age



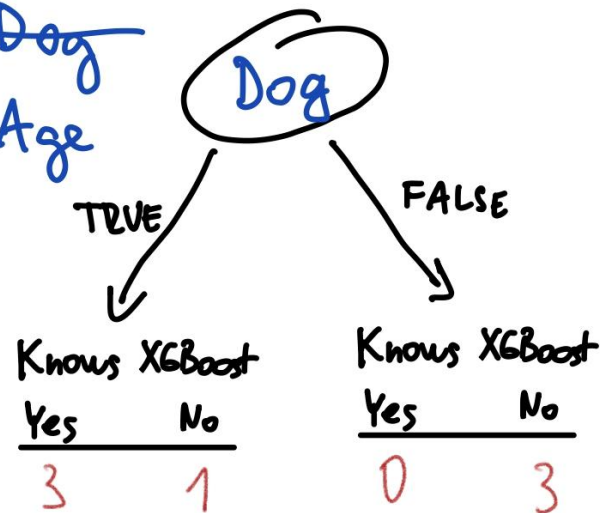



Is a woman	Has a dog	Age	Knows XGBoost
YES	YES	7	NO
YES	NO	12	NO
NO	YES	18	YES
NO	YES	35	YES
YES	YES	38	YES
YES	NO	50	NO
NO	NO	83	NO

~~Woman~~

~~Dog~~

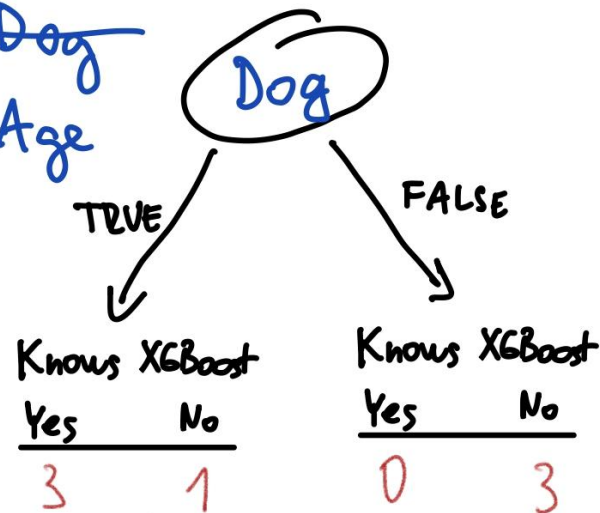
Age






Is a woman	Has a dog	Age	Knows XGBoost
YES	YES	7	NO
YES	NO	12	NO
NO	YES	18	YES
NO	YES	35	YES
YES	YES	38	YES
YES	NO	50	NO
NO	NO	83	NO

~~Woman~~
Dog
Age



Only this leaf is impure

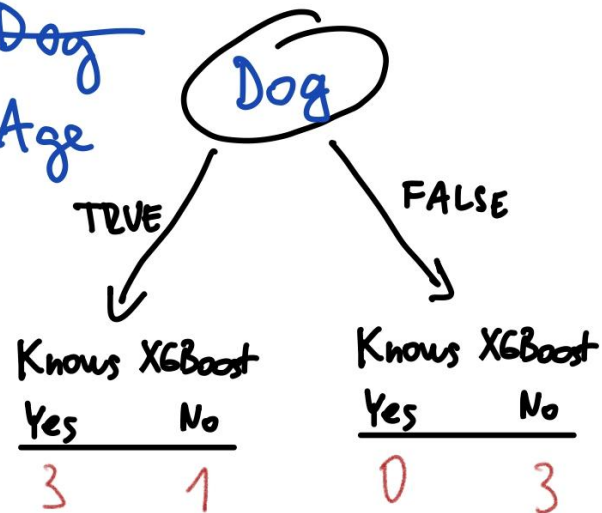


Is a woman	Has a dog	Age	Knows XGBoost
YES	YES	7	NO
YES	NO	12	NO
NO	YES	18	YES
NO	YES	35	YES
YES	YES	38	YES
YES	NO	50	NO
NO	NO	83	NO


~~Woman~~

~~Dog~~

Age

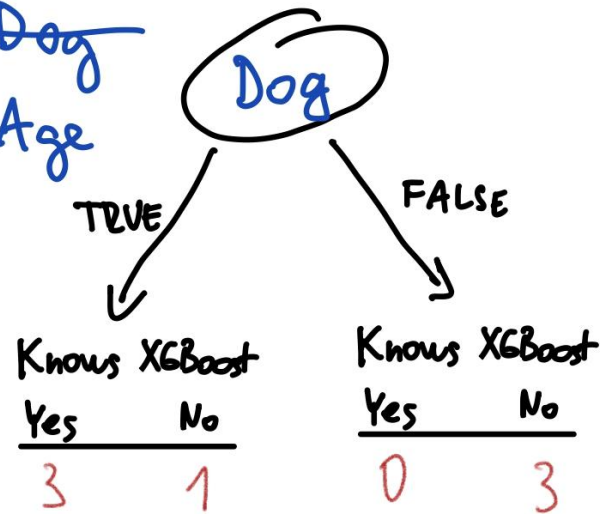


How to quantify which predictor is better?



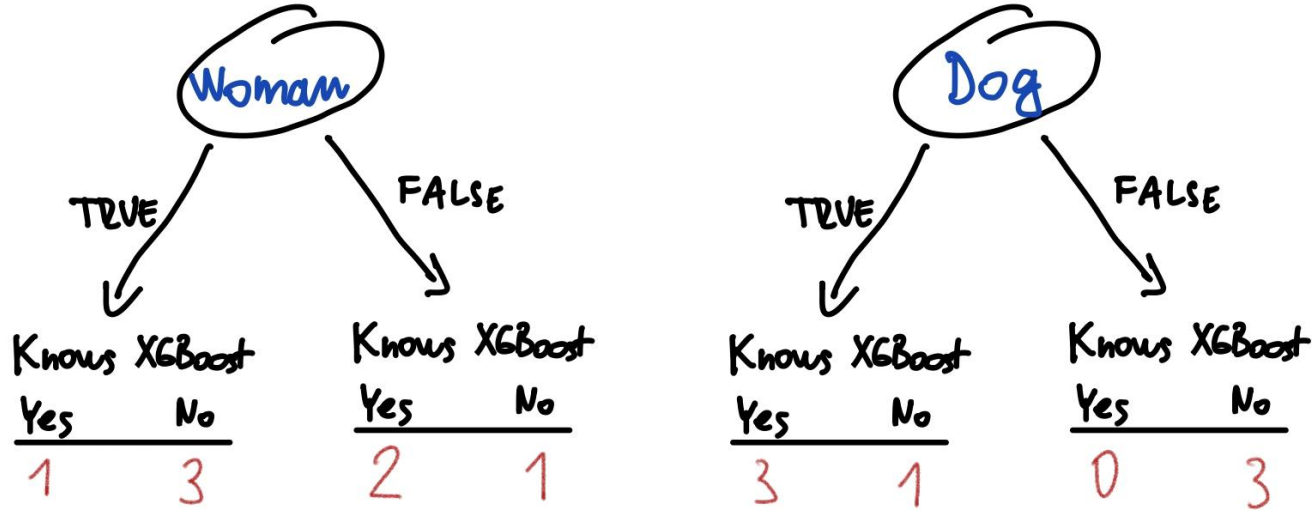
Is a woman	Has a dog	Age	Knows XGBoost
YES	YES	7	NO
YES	NO	12	NO
NO	YES	18	YES
NO	YES	35	YES
YES	YES	38	YES
YES	NO	50	NO
NO	NO	83	NO

~~Woman~~
~~Dog~~
Age



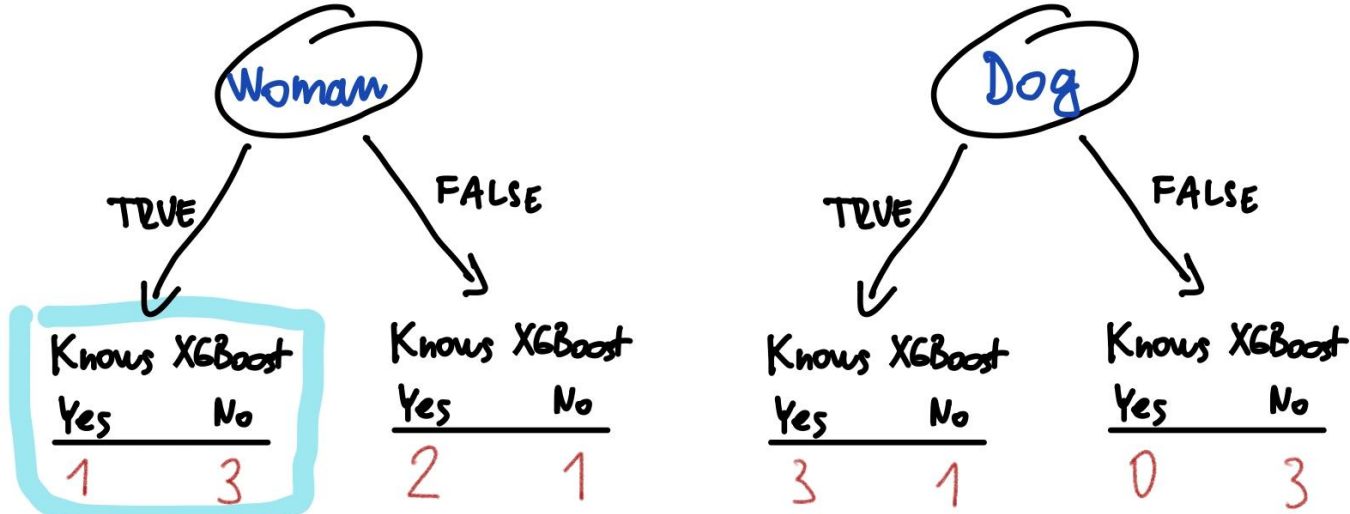
→ Gini impurity
(Entropy, Information Gain)

CALCULATING GINI IMPURITY



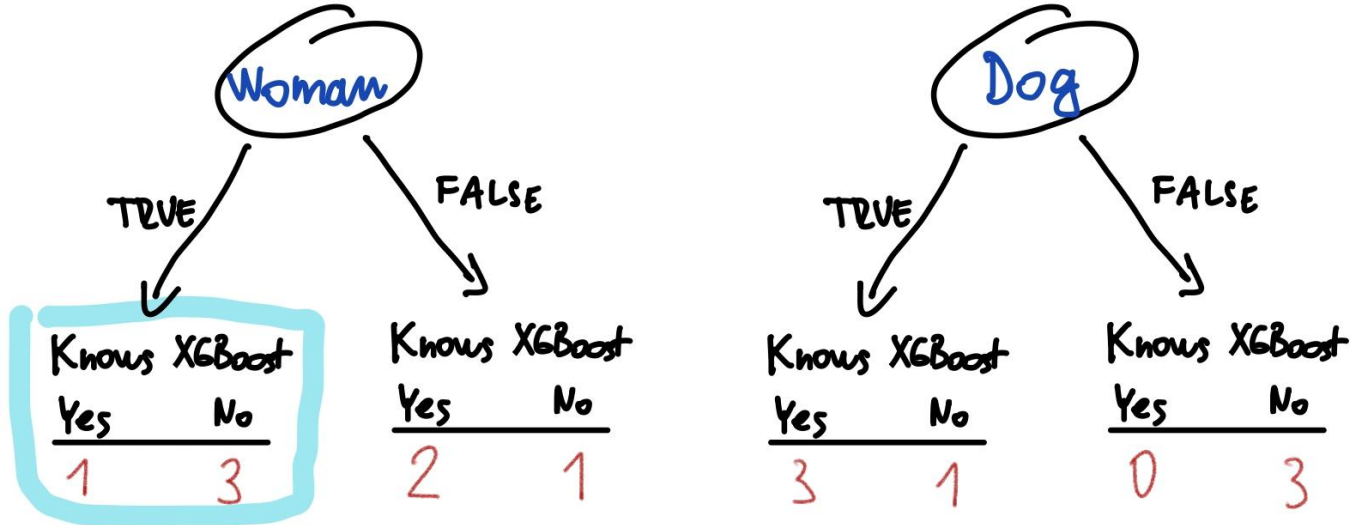
$$\text{G.I. for a Leaf} = 1 - (P[\text{Yes}])^2 - (P[\text{No}])^2$$

CALCULATING GINI IMPURITY



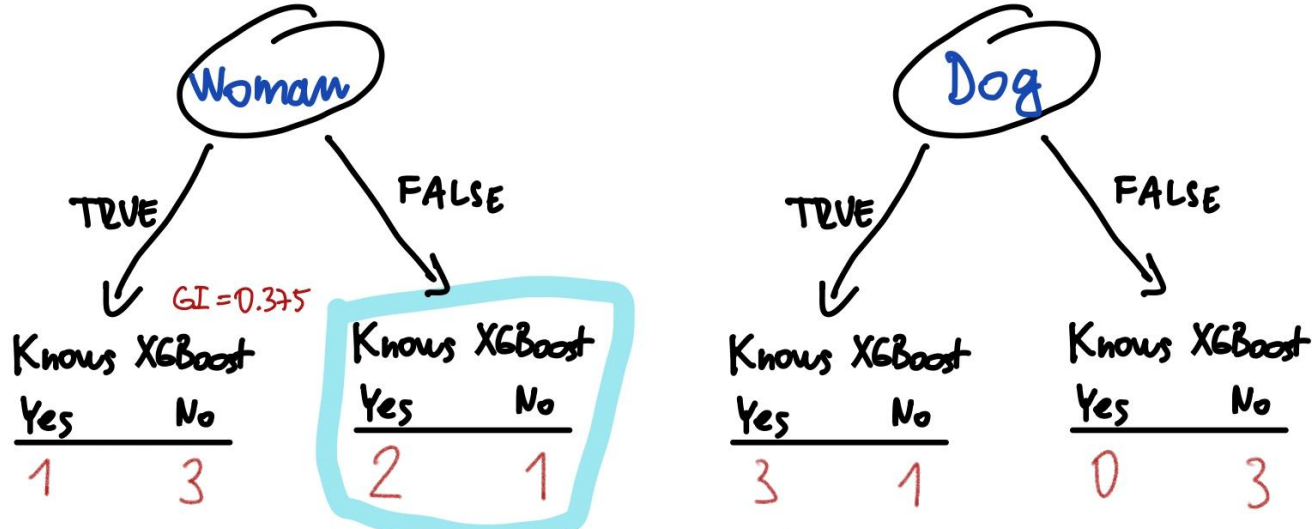
$$\begin{aligned}\text{G.I. for a Leaf} &= 1 - (P[\text{Yes}])^2 - (P[\text{No}])^2 \\ &= 1 - \left(\frac{1}{1+3}\right)^2 - \left(\frac{3}{1+3}\right)^2\end{aligned}$$

CALCULATING GINI IMPURITY



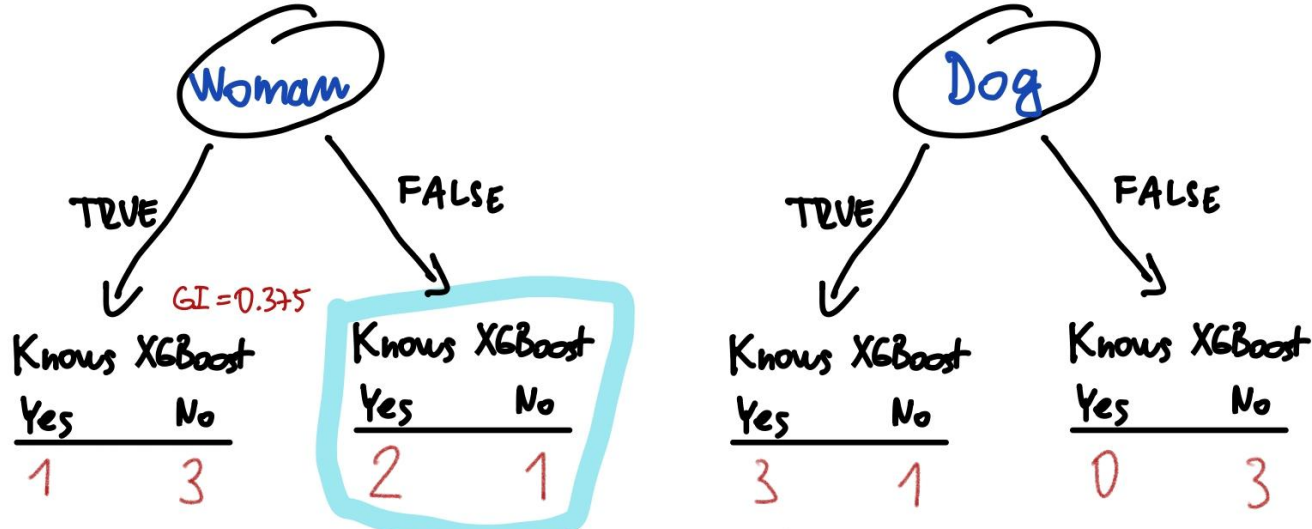
$$\begin{aligned}\text{G.I. for a Leaf} &= 1 - (P[\text{Yes}])^2 - (P[\text{No}])^2 \\ &= 1 - \left(\frac{1}{1+3}\right)^2 - \left(\frac{3}{1+3}\right)^2 \\ &= \underline{0.375}\end{aligned}$$

CALCULATING GINI IMPURITY



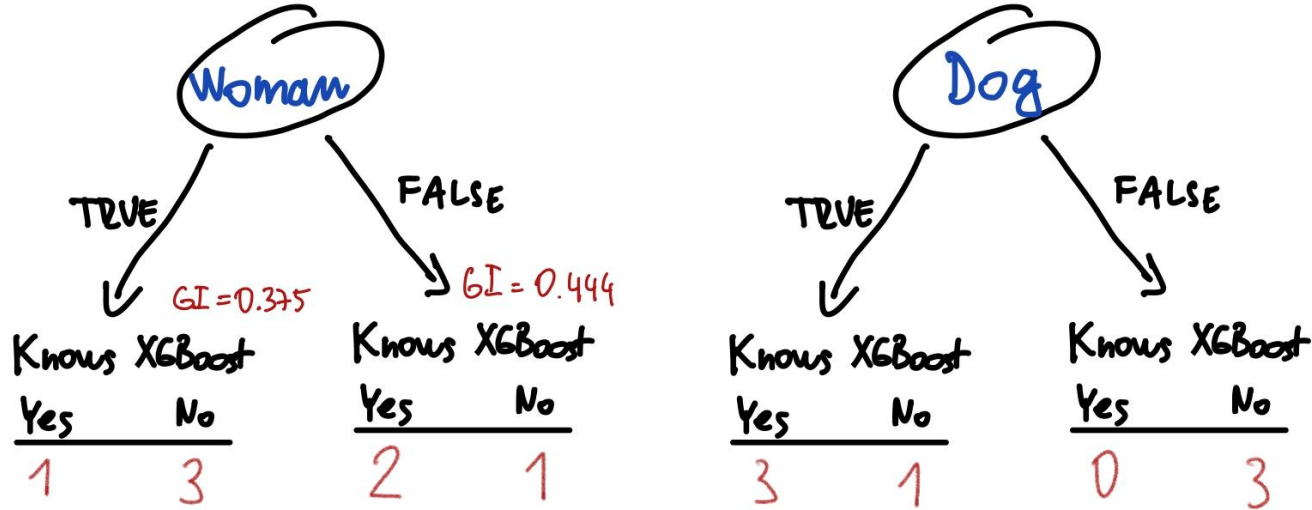
$$G.I. \text{ for a Leaf} = 1 - (P[\text{Yes}])^2 - (P[\text{No}])^2$$

CALCULATING GINI IMPURITY



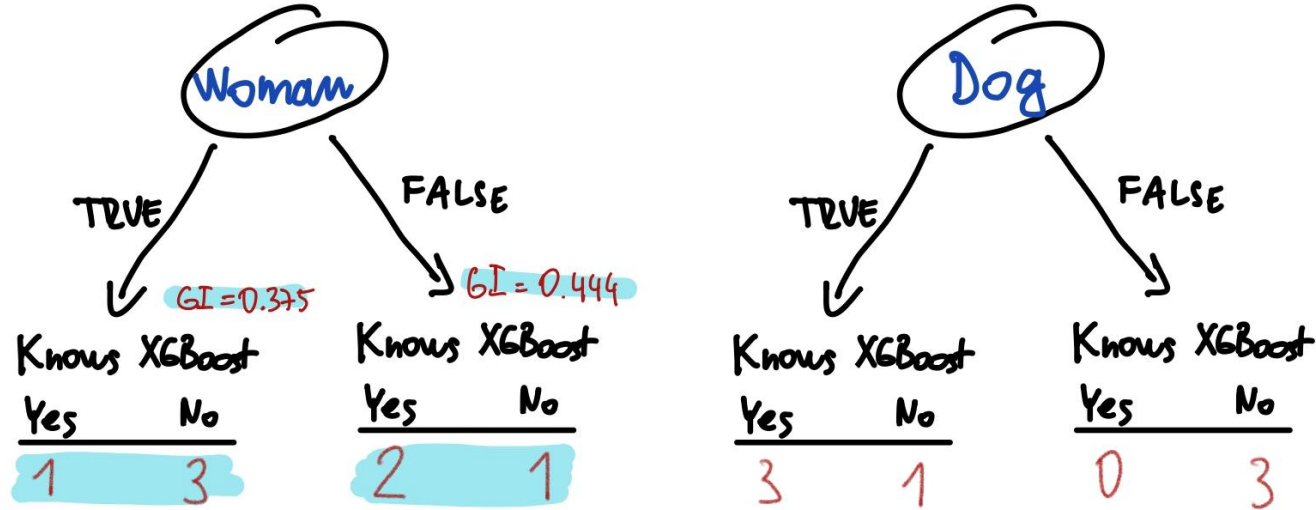
$$\begin{aligned} \text{G.I. for a Leaf} &= 1 - (P[\text{Yes}])^2 - (P[\text{No}])^2 \\ &= 1 - \left(\frac{2}{2+1}\right)^2 - \left(\frac{1}{2+1}\right)^2 \\ &= 0.444 \end{aligned}$$

CALCULATING GINI IMPURITY



Total Gini Impurity = weighted average of Gini Impurities for the Leaves

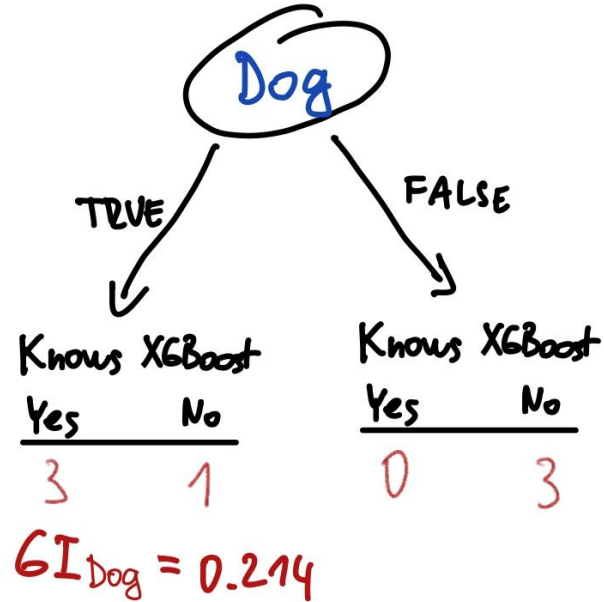
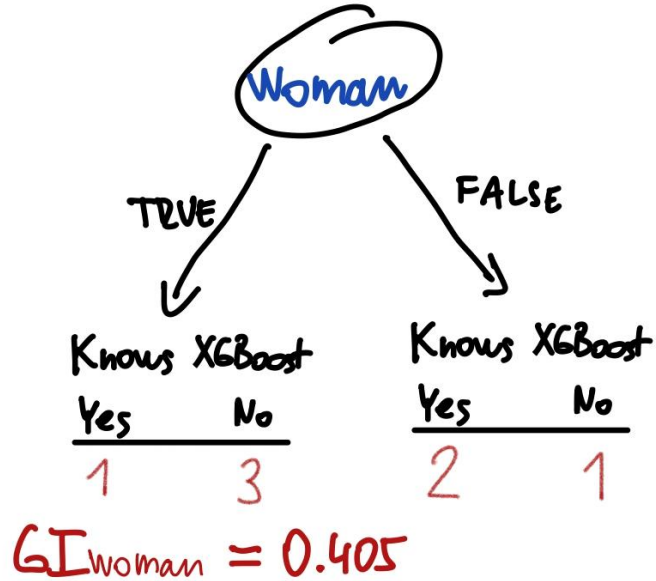
CALCULATING GINI IMPURITY




Total Gini Impurity = weighted average of Gini Impurities for the Leaves

$$= \frac{4}{4+3} \cdot 0.375 + \frac{3}{4+3} \cdot 0.444 = \underline{0.405}$$

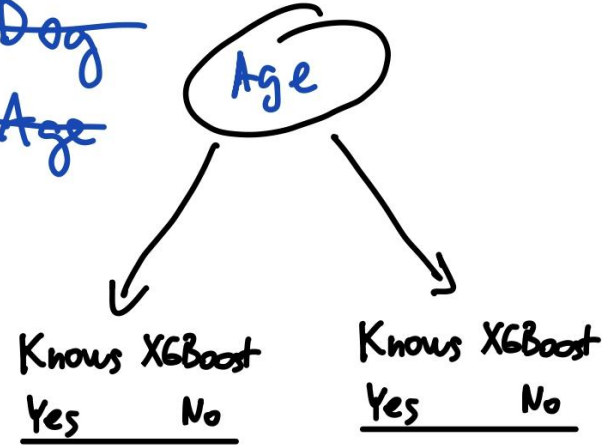
CALCULATING GINI IMPURITY





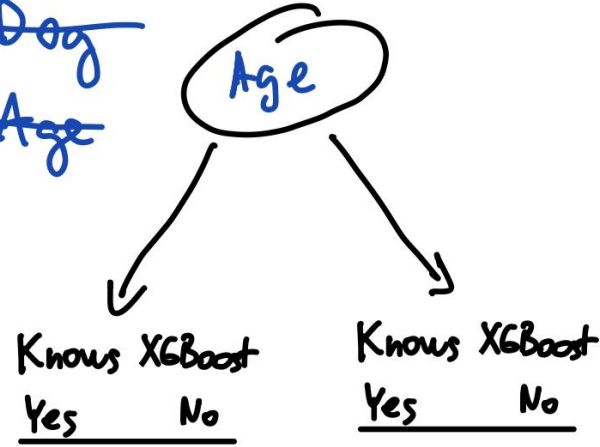
Is a woman	Has a dog	Age	Knows XGBoost
YES	YES	7	NO
YES	NO	12	NO
NO	YES	18	YES
NO	YES	35	YES
YES	YES	38	YES
YES	NO	50	NO
NO	NO	83	NO

~~Woman~~
~~Dog~~
Age



Is a woman	Has a dog	Age	Knows XGBoost
YES	YES	7	NO
YES	NO	12	NO
NO	YES	18	YES
NO	YES	35	YES
YES	YES	38	YES
YES	NO	50	NO
NO	NO	83	NO

~~Woman~~
~~Dog~~
Age



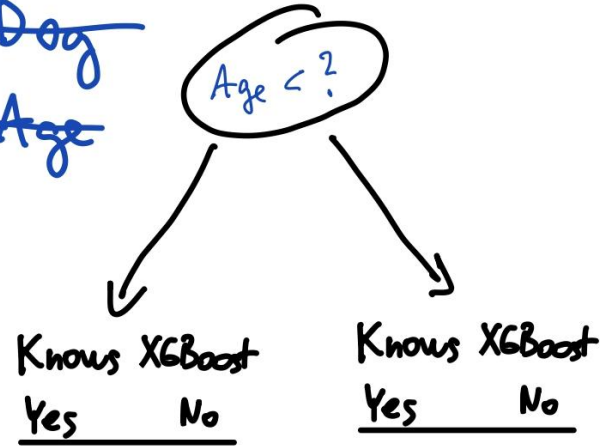
How to calculate Gini
Impurity with numerical
values?

1. SORT AGE IN
ASCENDING ORDER

2. CALCULATE AVG
FOR ADJACENT
PEOPLE

Age	Knows XGBoost
7	NO
12	NO
18	YES
35	YES
38	YES
50	NO
83	NO

~~Woman~~
~~Dog~~
Age

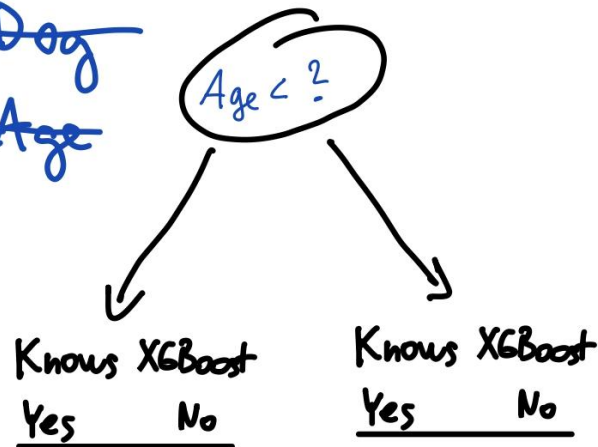


1. SORT AGE IN
ASCENDING ORDER

2. CALCULATE AVG
FOR ADJACENT
PEOPLE

	Age	Knows XGBoost
	7	NO
9.5	12	NO
15	18	YES
26.5	35	YES
36.5	38	YES
44	50	NO
66.5	83	NO

~~Woman~~
~~Dog~~
Age

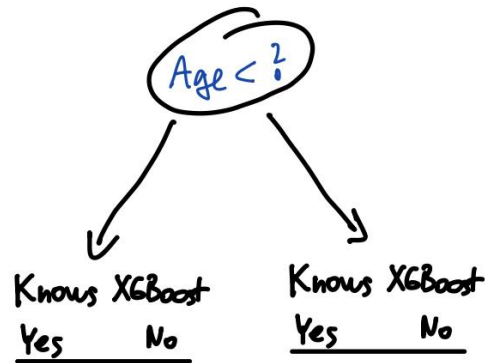


1. SORT AGE IN
ASCENDING ORDER

2. CALCULATE AVG
FOR ADJACENT
PEOPLE

3. CALCULATE GINI
IMPURITY FOR
ALL AVERAGE
AGES.

	Age	Knows XGBoost	
	7	NO	GI
9.5			→ 0.429
	12	NO	
15			0.343
	18	YES	
26.5			0.476
	35	YES	
36.5			0.476
	38	YES	
44			0.343
	50	NO	
66.5			0.429
	83	NO	

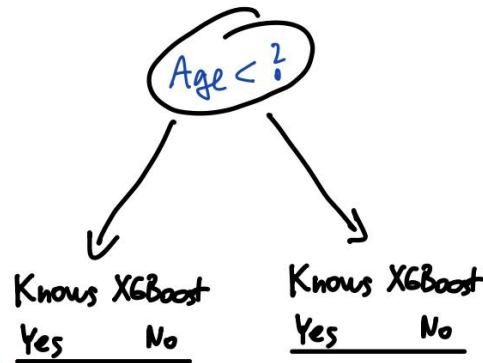


1. SORT AGE IN
ASCENDING ORDER

2. CALCULATE AVG
FOR ADJACENT
PEOPLE

3. CALCULATE GINI
IMPURITY FOR
ALL AVERAGE
AGES.

	Age	Knows XGBoost	
	7	NO	GI 0.429
9.5	12	NO	0.343
15	18	YES	0.476
26.5	35	YES	0.476
36.5	38	YES	0.343
44	50	NO	0.429
66.5	83	NO	



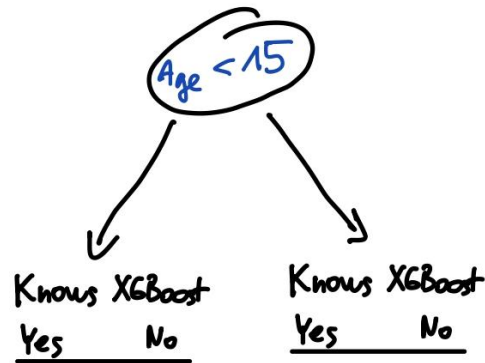
TWO CANDIDATES FOR
THRESHOLD

1. SORT AGE IN
ASCENDING ORDER

2. CALCULATE AVG
FOR ADJACENT
PEOPLE

3. CALCULATE GINI
IMPURITY FOR
ALL AVERAGE
AGES.

	Age	Knows XGBoost	
	7	NO	GI
9.5			→ 0.429
	12	NO	
15			0.343
	18	YES	
26.5			0.476
	35	YES	
36.5			0.476
	38	YES	
44			0.343
	50	NO	
66.5			0.429
	83	NO	

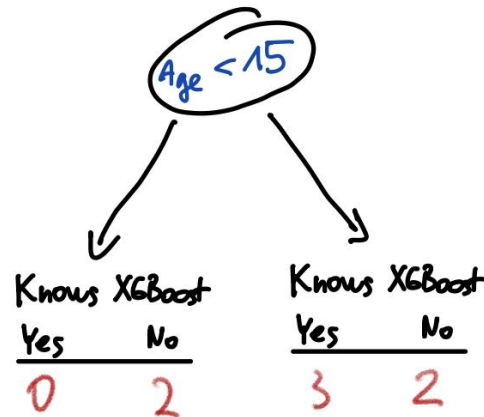


1. SORT AGE IN
ASCENDING ORDER

2. CALCULATE AVG
FOR ADJACENT
PEOPLE

3. CALCULATE GINI
IMPURITY FOR
ALL AVERAGE
AGES.

	Age	Knows XGBoost	
	7	NO	GI
9.5			$\rightarrow 0.429$
	12	NO	
15			0.343
	18	YES	
26.5			0.476
	35	YES	
36.5			0.476
	38	YES	
44			0.343
	50	NO	
66.5			0.429
	83	NO	



$$GI_{Age} = 0.343$$

Is a woman	Has a dog	Age	Knows XGBoost
YES	YES	7	NO
YES	NO	12	NO
NO	YES	18	YES
NO	YES	35	YES
YES	YES	38	YES
YES	NO	50	NO
NO	NO	83	NO

Woman

Dog

Age



WHICH PREDICTOR
SHOULD WE CHOOSE?

Is a woman	Has a dog	Age	Knows XGBoost
YES	YES	7	NO
YES	NO	12	NO
NO	YES	18	YES
NO	YES	35	YES
YES	YES	38	YES
YES	NO	50	NO
NO	NO	83	NO

Woman

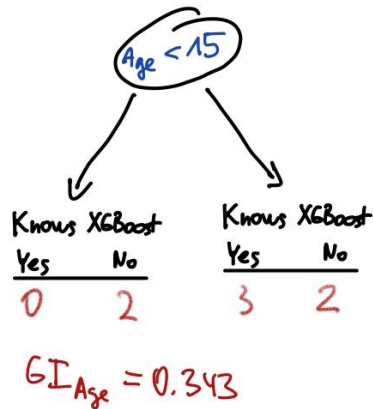
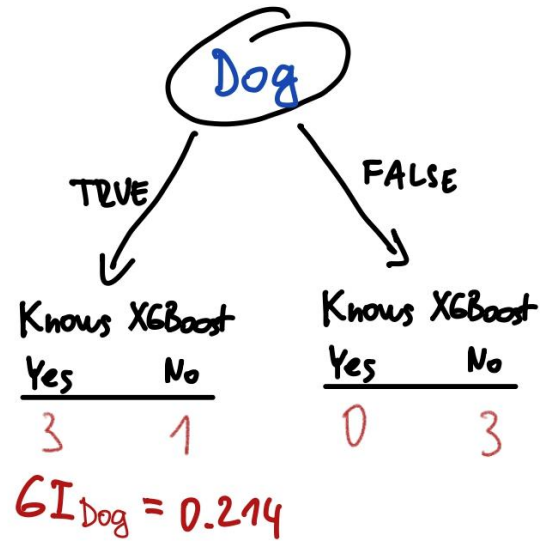
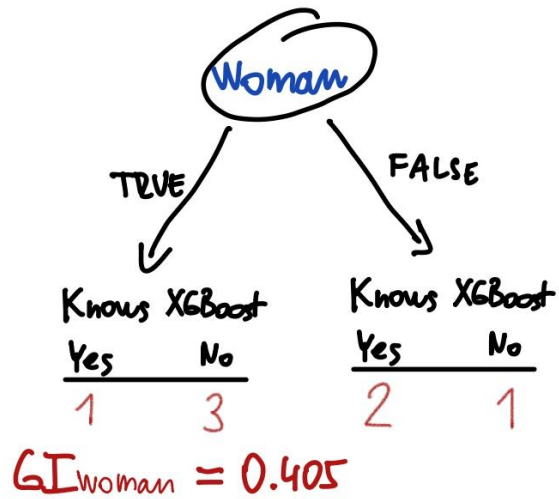
Dog

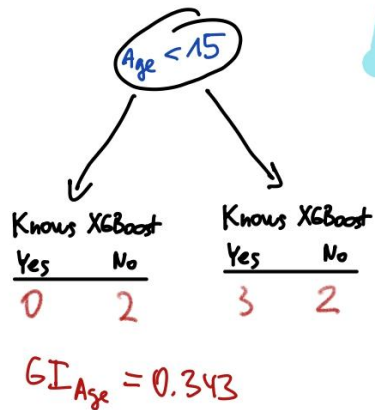
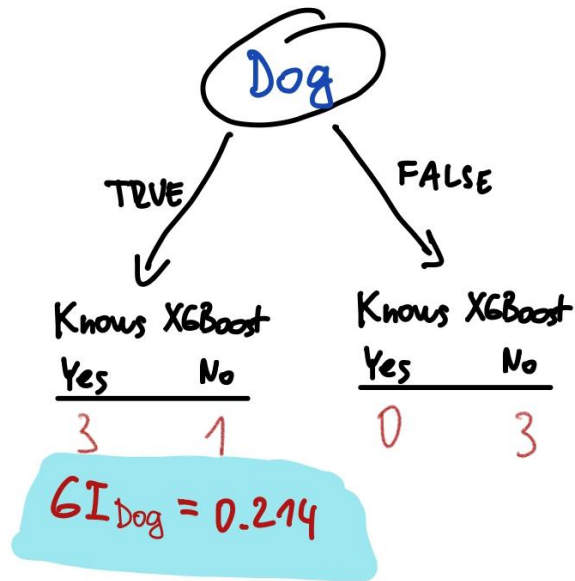
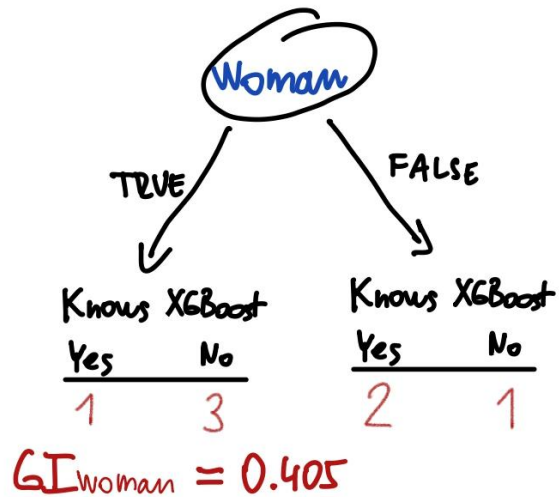
Age



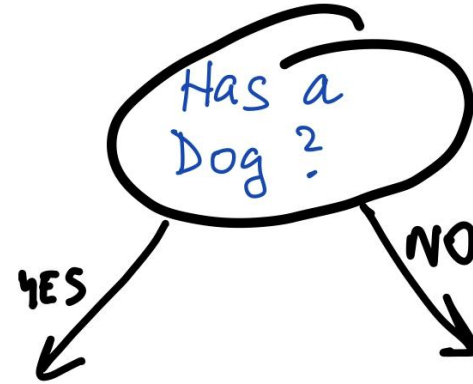
WHICH PREDICTOR
SHOULD WE CHOOSE?

→ THE ONE WITH THE
LOWEST GINI IMPURITY

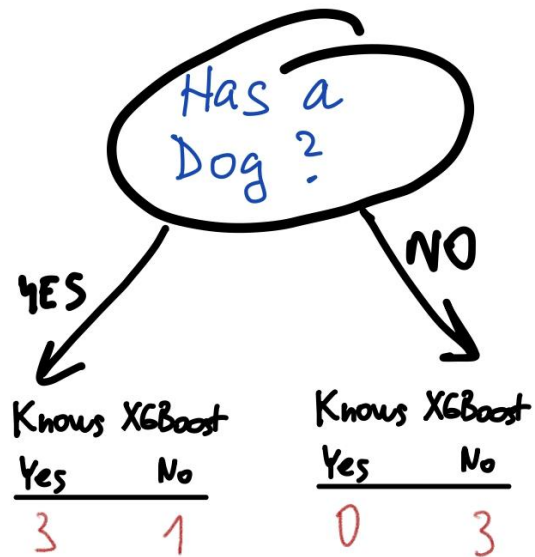




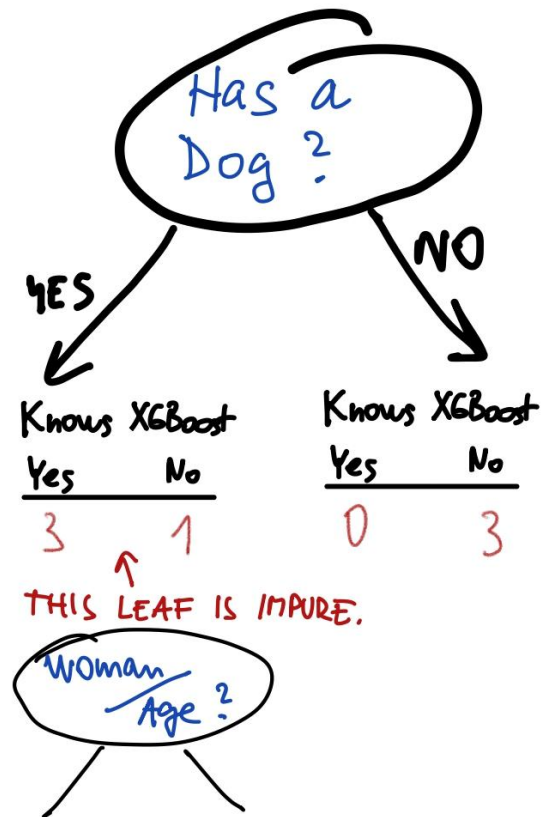
Is a woman	Has a dog	Age	Knows XGBoost
YES	YES	7	NO
YES	NO	12	NO
NO	YES	18	YES
NO	YES	35	YES
YES	YES	38	YES
YES	NO	50	NO
NO	NO	83	NO



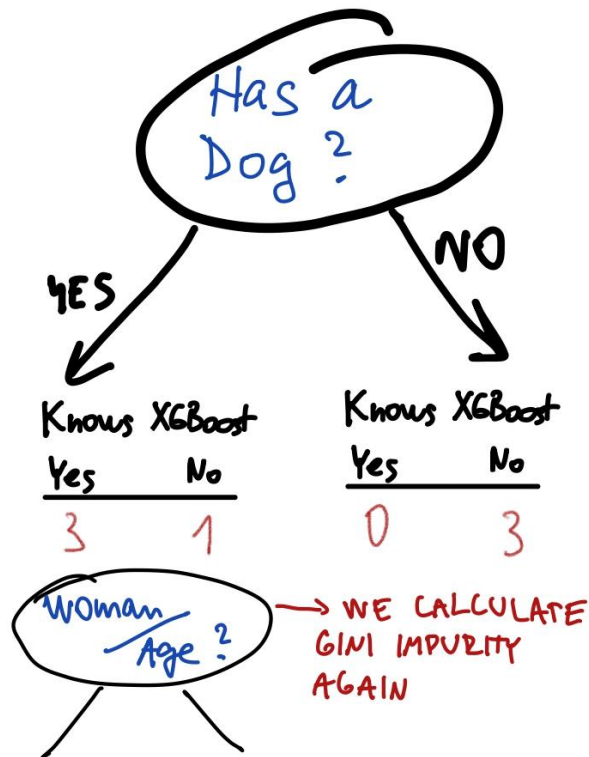
Is a woman	Has a dog	Age	Knows XGBoost
YES	YES	7	NO
YES	NO	12	NO
NO	YES	18	YES
NO	YES	35	YES
YES	YES	38	YES
YES	NO	50	NO
NO	NO	83	NO



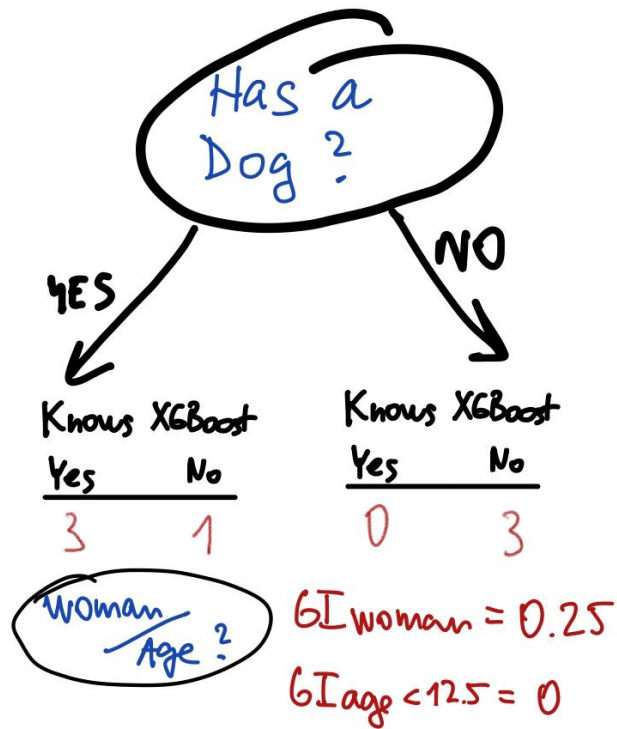
Is a woman	Has a dog	Age	Knows XGBoost
YES	YES	7	NO
YES	NO	12	NO
NO	YES	18	YES
NO	YES	35	YES
YES	YES	38	YES
YES	NO	50	NO
NO	NO	83	NO



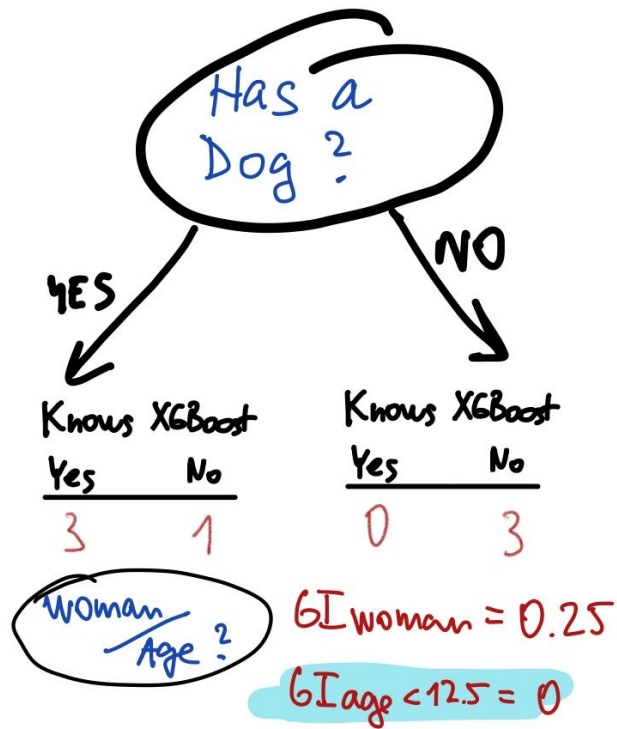
Is a woman	Has a dog	Age	Knows XGBoost
YES	YES	7	NO
YES	NO	12	NO
NO	YES	18	YES
NO	YES	35	YES
YES	YES	38	YES
YES	NO	50	NO
NO	NO	83	NO



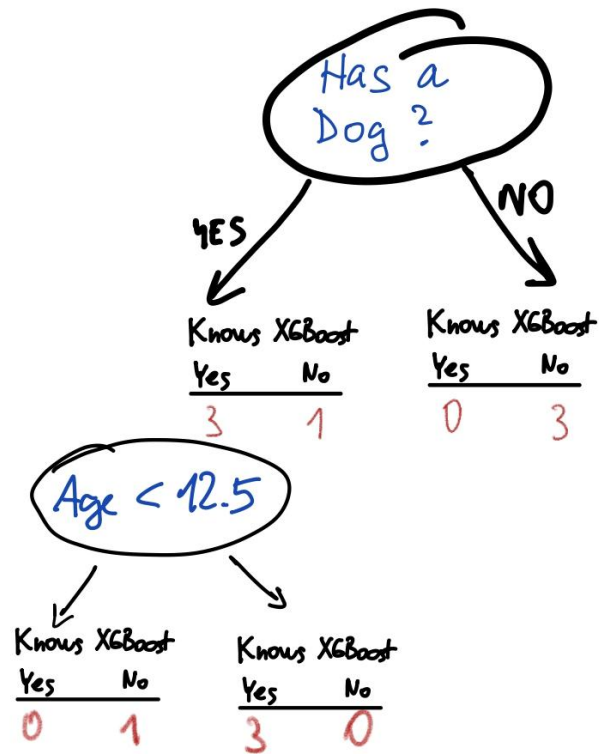
Is a woman	Has a dog	Age	Knows XGBoost
YES	YES	7	NO
YES	NO	12	NO
NO	YES	18	YES
NO	YES	35	YES
YES	YES	38	YES
YES	NO	50	NO
NO	NO	83	NO



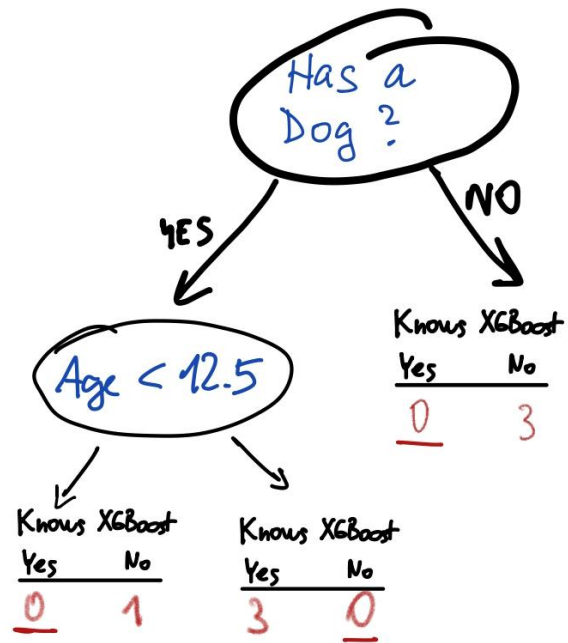
Is a woman	Has a dog	Age	Knows XGBoost
YES	YES	7	NO
YES	NO	12	NO
NO	YES	18	YES
NO	YES	35	YES
YES	YES	38	YES
YES	NO	50	NO
NO	NO	83	NO



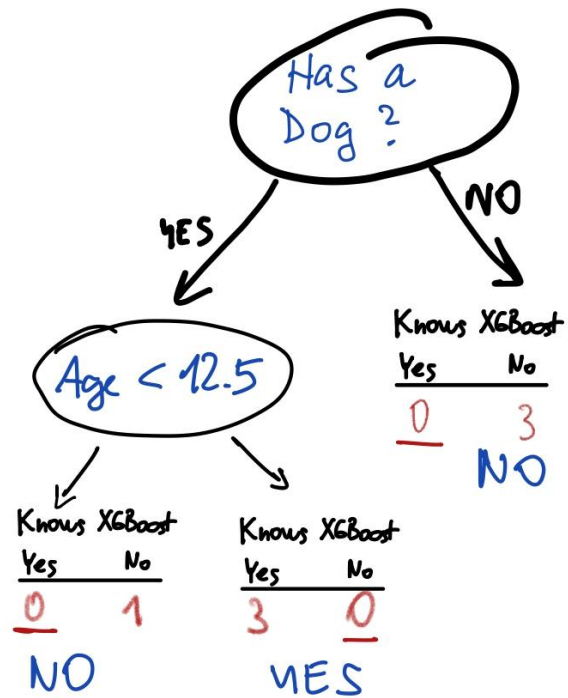
Is a woman	Has a dog	Age	Knows XGBoost
YES	YES	7	NO
YES	NO	12	NO
NO	YES	18	YES
NO	YES	35	YES
YES	YES	38	YES
YES	NO	50	NO
NO	NO	83	NO



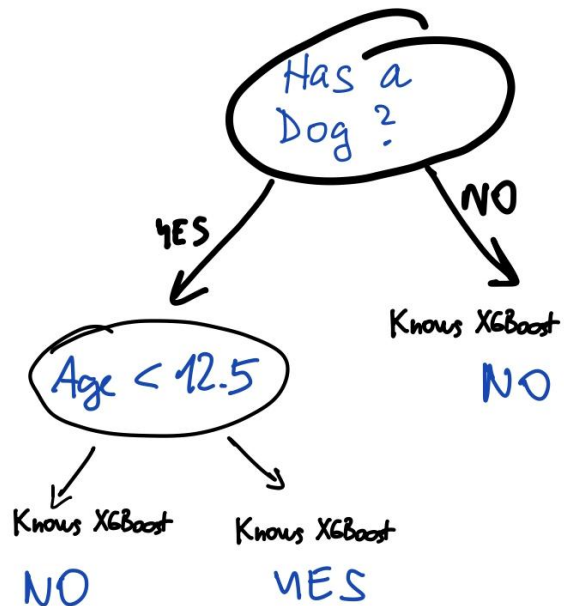
Is a woman	Has a dog	Age	Knows XGBoost
YES	YES	7	NO
YES	NO	12	NO
NO	YES	18	YES
NO	YES	35	YES
YES	YES	38	YES
YES	NO	50	NO
NO	NO	83	NO



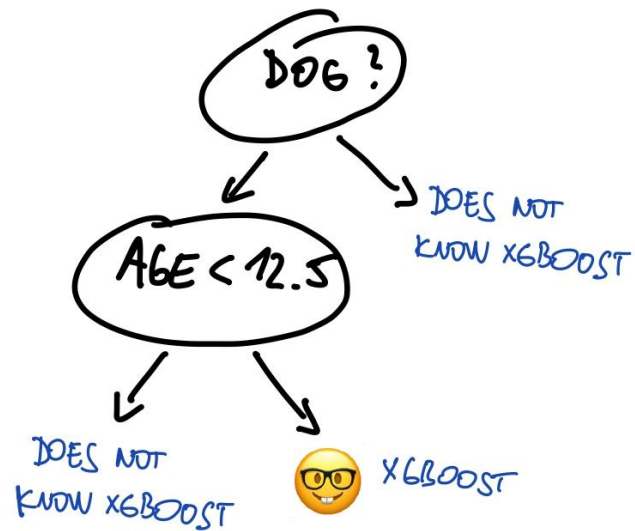
Is a woman	Has a dog	Age	Knows XGBoost
YES	YES	7	NO
YES	NO	12	NO
NO	YES	18	YES
NO	YES	35	YES
YES	YES	38	YES
YES	NO	50	NO
NO	NO	83	NO



Is a woman	Has a dog	Age	Knows XGBoost
YES	YES	7	NO
YES	NO	12	NO
NO	YES	18	YES
NO	YES	35	YES
YES	YES	38	YES
YES	NO	50	NO
NO	NO	83	NO



Is a woman	Has a dog	Age	Knows XGBoost
YES	YES	7	NO
YES	NO	12	NO
NO	YES	18	YES
NO	YES	35	YES
YES	YES	38	YES
YES	NO	50	NO
NO	NO	83	NO



Building trees (summary)

Greedy algorithm:

1. consider all features
2. consider all possible thresholds
3. compute impurity or loss reduction
4. pick the best split
5. recurse on left and right branch
6. stop when: max depth OR not enough samples OR improvement small

Pros and Cons of Decision Trees

- interpretable
- no scaling
- both numerical and categorical
- fast
- nonlinear relationships
- overfitting
- high variance
- instability
- not very accurate



Pros and Cons of Decision Trees

- overfitting
- high variance
- instability
- not very accurate

Solution?

-> Ensemble methods



Ensemble Methods: From Weak Learners to Robust Predictors

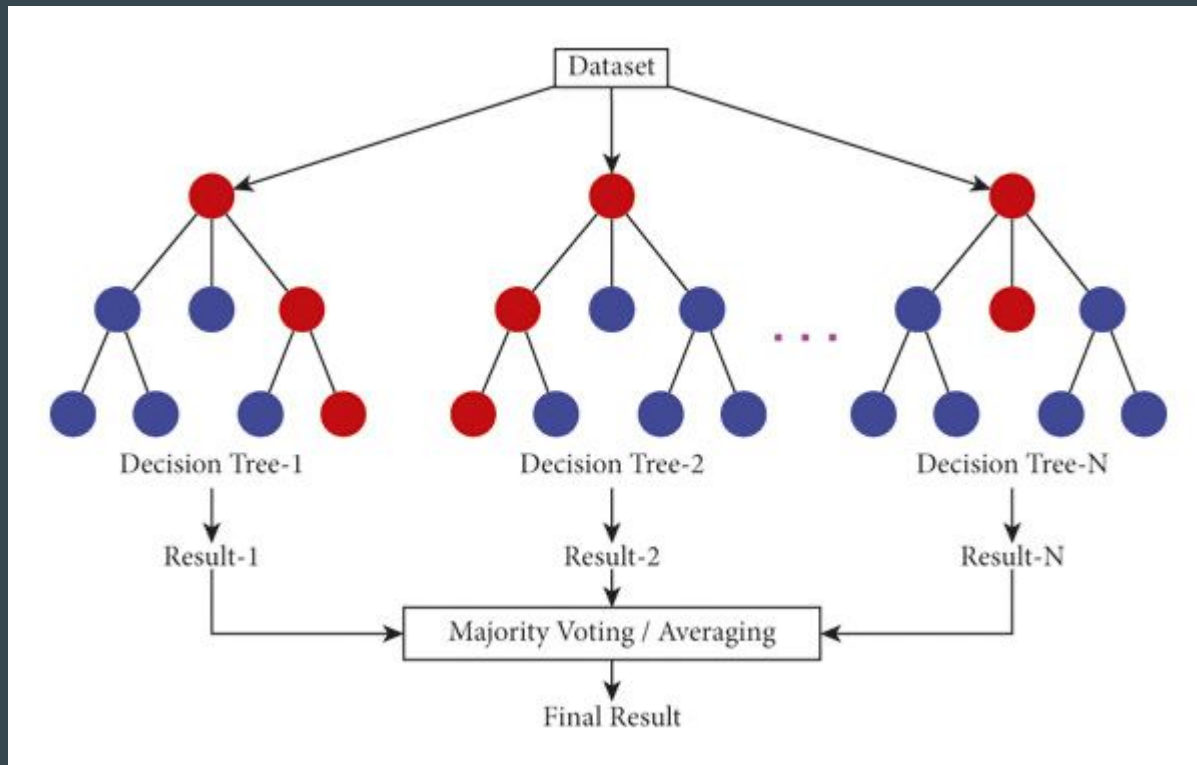
Motivation

- Reduce variance – single tree highly unstable
- Reduce bias –
 - shallow trees → too simple, high bias
 - deep trees → overfit, overly complex
- Bagging (Random Forests)
- Boosting (XGBoost, LightGBM, CatBoost)

Bagging (Random Forest)

= bootstrap aggregating

idea: lower variance by
averaging multiple results



Building Random Forest

1. create a bootstrapped dataset
2. build decision tree using random subset of features at each step

CREATE BOOTSTRAPPED DATASET

Watch TV	Weight	Is a woman	Has a dog	Age	Knows XGBoost
YES	90	YES	YES	7	NO
NO	108	YES	NO	12	NO
NO	172	NO	YES	18	YES
NO	186	NO	YES	35	YES
NO	152	YES	YES	38	YES
YES	149	YES	NO	50	NO
YES	207	NO	NO	83	NO

CREATE BOOTSTRAPPED DATASET

	Watch TV	Weight	Is a woman	Has a dog	Age	Knows XGBoost
1	YES	90	YES	YES	7	NO
2	NO	108	YES	NO	12	NO
3	NO	172	NO	YES	18	YES
4	NO	186	NO	YES	35	YES
5	NO	152	YES	YES	38	YES
6	YES	149	YES	NO	50	NO
7	YES	207	NO	NO	83	NO

= SAMPLING WITH REPLACEMENT N. SAMPLES = #rows

E.G. 7, 2, 1, 4, 7, 3, 1

CREATE BOOTSTRAPPED DATASET

	Watch TV	Weight	Is a woman	Has a dog	Age	Knows XGBoost
1	YES	90	YES	YES	7	NO
2	NO	108	YES	NO	12	NO
3	NO	172	NO	YES	18	YES
4	NO	186	NO	YES	35	YES
5	NO	152	YES	YES	38	YES
6	YES	149	YES	NO	50	NO
7	YES	207	NO	NO	83	NO

= SAMPLING WITH REPLACEMENT N. SAMPLES = #rows

E.G. 7, 2, 1, 4, 7, 3, 1

CREATE BOOTSTRAPPED DATASET

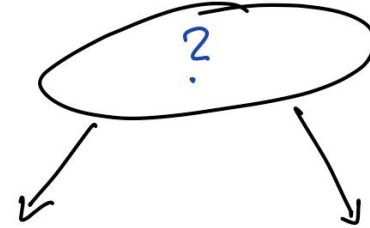
	Watch TV	Weight	Is a woman	Has a dog	Age	Knows XGBoost
1	YES	90	YES	YES	7	NO
2	NO	108	YES	NO	12	NO
3	NO	172	NO	YES	18	YES
4	NO	186	NO	YES	35	YES
7	YES	207	NO	NO	83	NO
1	YES	90	YES	YES	7	NO
7	YES	207	NO	NO	83	NO

= SAMPLING WITH REPLACEMENT N. SAMPLES = #rows

E.G. 7, 2, 1, 4, 7, 3, 1

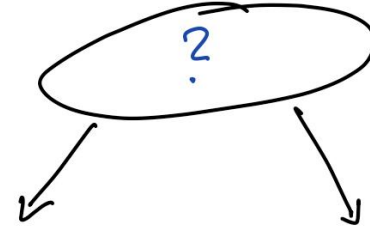
BUILD DECISION TREE

Watch TV	Weight	Is a woman	Has a dog	Age	Knows XGBoost
YES	90	YES	YES	7	NO
NO	108	YES	NO	12	NO
NO	172	NO	YES	18	YES
NO	186	NO	YES	35	YES
YES	207	NO	NO	83	NO
YES	90	YES	YES	7	NO
YES	207	NO	NO	83	NO



BUILD DECISION TREE

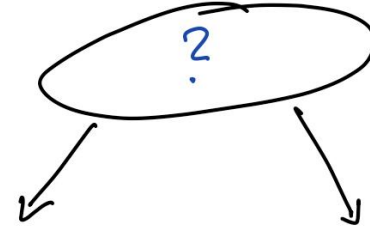
Watch TV	Weight	Is a woman	Has a dog	Age	Knows XGBoost
YES	90	YES	YES	7	NO
NO	108	YES	NO	12	NO
NO	172	NO	YES	18	YES
NO	186	NO	YES	35	YES
YES	207	NO	NO	83	NO
YES	90	YES	YES	7	NO
YES	207	NO	NO	83	NO



→ CHOOSE RANDOM
SUBSET OF VARIABLES
E.G., [Has a Dog, Age]

BUILD DECISION TREE

Watch TV	Weight	Is a woman	Has a dog	Age	Knows XGBoost
YES	90	YES	YES	7	NO
NO	108	YES	NO	12	NO
NO	172	NO	YES	18	YES
NO	186	NO	YES	35	YES
YES	207	NO	NO	83	NO
YES	90	YES	YES	7	NO
YES	207	NO	NO	83	NO



→ CHOOSE RANDOM
SUBSET OF VARIABLES
E.G., [Has a Dog, Age]

→ DETERMINE THE BEST
CANDIDATE BY CALCULATING
GINI IMPURITY

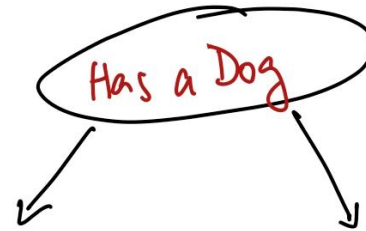
$GI_{\text{Dog}} = \dots$ $GI_{\text{Age}} = \dots$

BUILD DECISION TREE

Watch TV	Weight	Is a woman	Has a dog	Age	Knows XGBoost
YES	90	YES	YES	7	NO
NO	108	YES	NO	12	NO
NO	172	NO	YES	18	YES
NO	186	NO	YES	35	YES
YES	207	NO	NO	83	NO
YES	90	YES	YES	7	NO
YES	207	NO	NO	83	NO

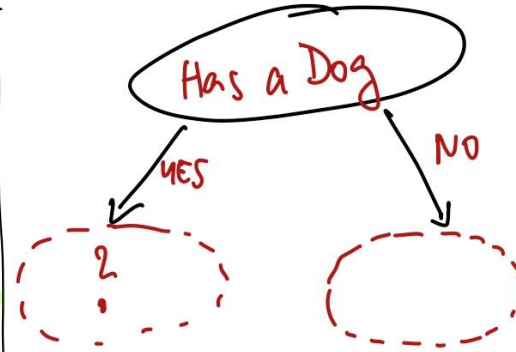
→ DETERMINE THE BEST CANDIDATE BY CALCULATING GINI IMPURITY

Assume $GI_{\text{Dog}} < GI_{\text{Age}}$



BUILD DECISION TREE

Watch TV	Weight	Is a woman	Has a dog	Age	Knows XGBoost
YES	90	YES	YES	7	NO
NO	108	YES	NO	12	NO
NO	172	NO	YES	18	YES
NO	186	NO	YES	35	YES
YES	207	NO	NO	83	NO
YES	90	YES	YES	7	NO
YES	207	NO	NO	83	NO



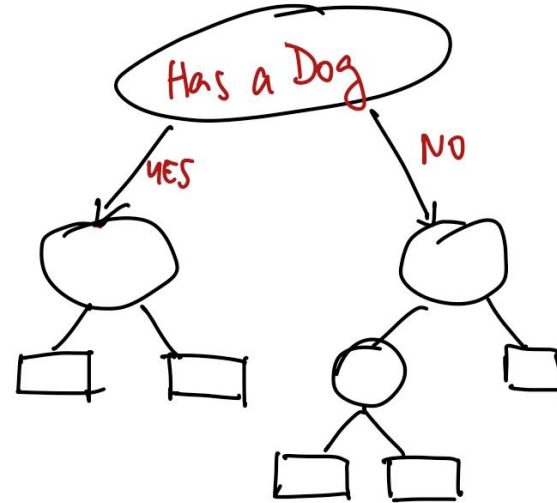
→ CHOOSE RANDOM SUBSET OF VARIABLES AGAIN

→ SELECT THE BEST SPLIT

→ BUILD THE FULL TREE

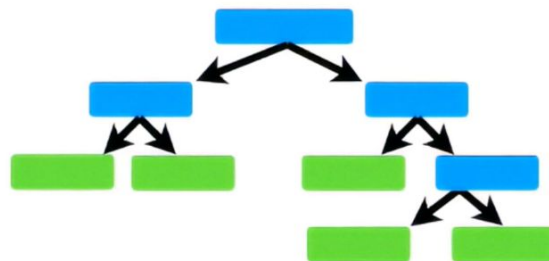
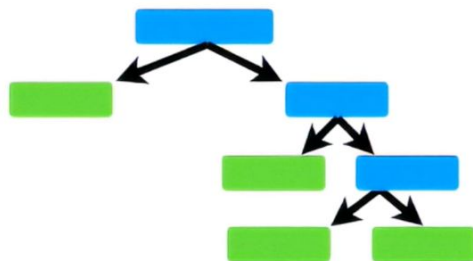
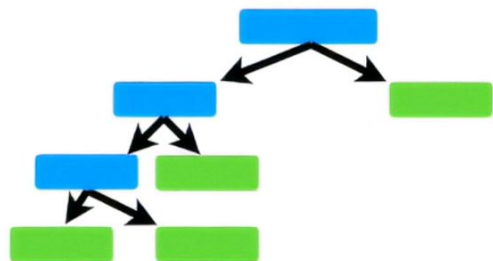
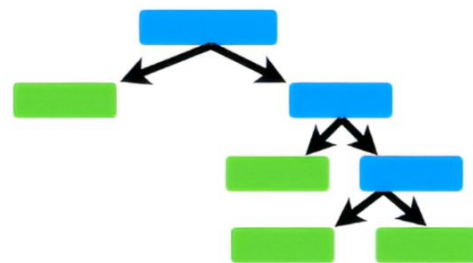
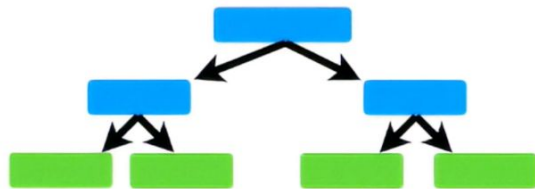
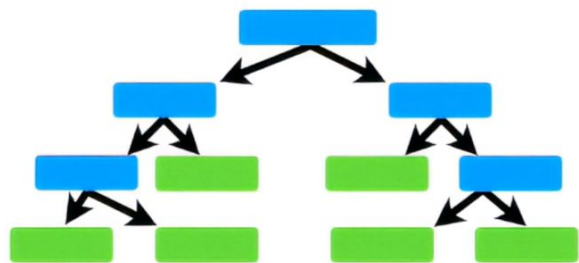
BUILD DECISION TREE

Watch TV	Weight	Is a woman	Has a dog	Age	Knows XGBoost
YES	90	YES	YES	7	NO
NO	108	YES	NO	12	NO
NO	172	NO	YES	18	YES
NO	186	NO	YES	35	YES
YES	207	NO	NO	83	NO
YES	90	YES	YES	7	NO
YES	207	NO	NO	83	NO



Building Random Forest

1. create a bootstrapped dataset
2. build decision tree using random subset of features at each step
3. repeat for n =number of trees times



Obtaining Prediction

- we evaluate all individual trees in random forest
- aggregate the result
 - majority vote / average value
- bootstrapped dataset + aggregating = bagging

Limitations of Bagging

✓ Pros of Random Forests

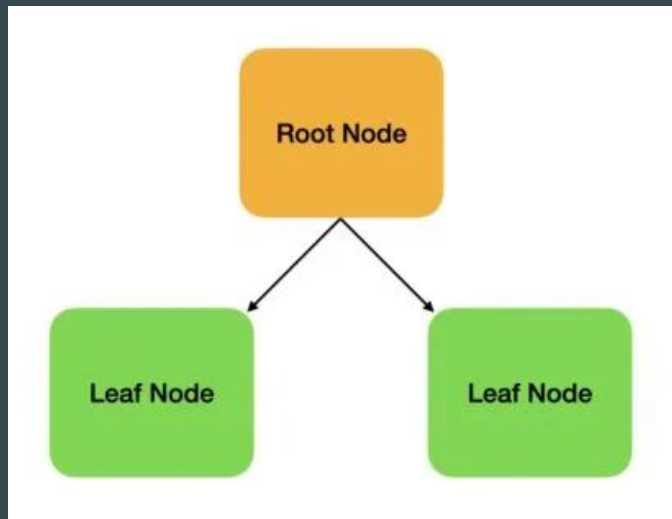
- variance reduction
- + stability

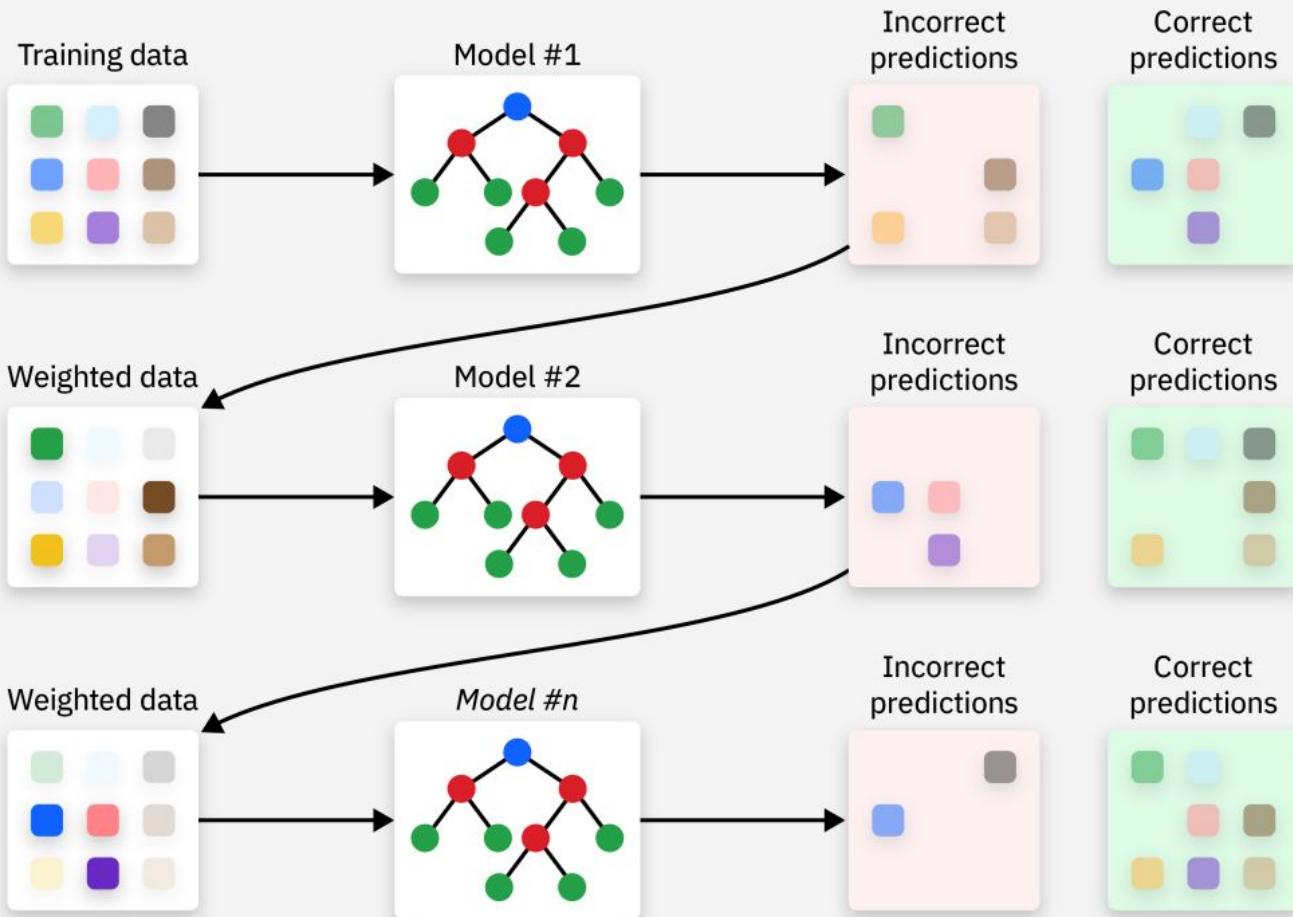
✗ Limitations of Random Forests

- averaging does not reduce **bias**
 - slower inference
 - trees built independently: they don't learn from previous errors
-
- What if we could build trees that *learn from each other*?

Boosting

- idea: sequentially build trees correcting error of previous model
- increase weight of misclassified examples
- usually shallow trees / stumps
- combine predictions
- AdaBoost, Gradient Boosting





Gradient Boosting

- gradient descent to minimize loss function
- -> gradient-guided correction

Step 1: initial prediction (avg)

Step 2: calculate pseudo-residuals

Step 3: build decision tree (8 to 32 leaves) fitting residuals

-> Iterate 2, 3 until stopping criterion

Algorithm 1: Gradient_Boost

```
1  $F_0(\mathbf{x}) = \arg \min_{\rho} \sum_{i=1}^N L(y_i, \rho)$ 
2 For  $m = 1$  to  $M$  do:
3    $\tilde{y}_i = - \left[ \frac{\partial L(y_i, F(\mathbf{x}_i))}{\partial F(\mathbf{x}_i)} \right]_{F(\mathbf{x})=F_{m-1}(\mathbf{x})}, i = 1, N$ 
4    $\mathbf{a}_m = \arg \min_{\mathbf{a}, \beta} \sum_{i=1}^N [\tilde{y}_i - \beta h(\mathbf{x}_i; \mathbf{a})]^2$ 
5    $\rho_m = \arg \min_{\rho} \sum_{i=1}^N L(y_i, F_{m-1}(\mathbf{x}_i) + \rho h(\mathbf{x}_i; \mathbf{a}_m))$ 
6    $F_m(\mathbf{x}) = F_{m-1}(\mathbf{x}) + \rho_m h(\mathbf{x}; \mathbf{a}_m)$ 
7 endFor
end Algorithm
```

Age	Category	Purchase weight (kg)	Amount
25	Electronics	2.5	123.45
34	Clothing	1.3	56.78
42	Electronics	5.0	345.67
19	Hardware	3.2	98.01

$$MSE = \frac{1}{2} (\text{Observed} - \text{Predicted})^2$$

Step 1: Initial prediction

$$\frac{dL}{d\hat{y}} = (\text{predicted} - \text{observed})$$

$$\frac{123.45 + 56.78 + 345.67 + 98.01}{4} = 156$$

Age	Category	Purchase weight (kg)	Amount	PSEUDO RESIDUALS
25	Electronics	2.5	123.45	-32.55
34	Clothing	1.3	56.78	-99.22
42	Electronics	5.0	345.67	189.45
19	Hardware	3.2	98.01	-58.01

$$MSE = \frac{1}{2} (\text{Observed} - \text{Predicted})^2$$

Step 1: Initial prediction

$$\frac{dL}{d\hat{y}} = (\text{predicted} - \text{observed})$$

$$\frac{123.45 + 56.78 + 345.67 + 98.01}{4} = 156$$

Step 2: Calculate pseudo-residuals

Age	Category	Purchase weight (kg)	Amount	PSEUDO RESIDUALS
25	Electronics	2.5	123.45	-32.55
34	Clothing	1.3	56.78	-99.22
42	Electronics	5.0	345.67	189.45
19	Hardware	3.2	98.01	-58.01

$$MSE = \frac{1}{2} (\text{Observed} - \text{Predicted})^2$$

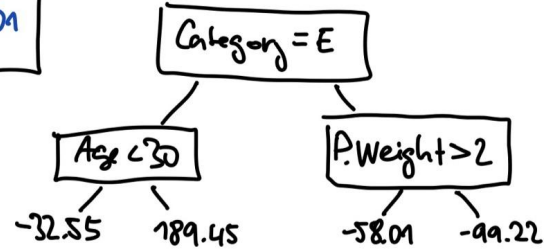
Step 1: Initial prediction

$$\frac{dL}{d\hat{y}} = (\text{predicted} - \text{observed})$$

$$\frac{123.45 + 56.78 + 345.67 + 98.01}{4} = 156$$

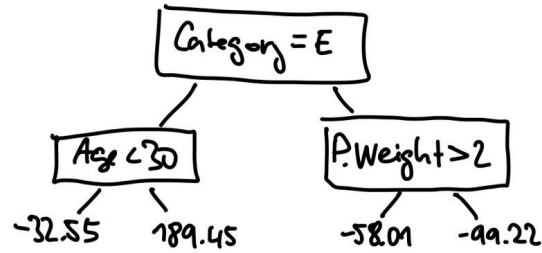
Step 2: Calculate pseudo-residuals

Step 3: Build a weak learner



Step 3: Build a weak learner +

Calculate predictions for each row



Age	Category	Purchase weight (kg)	Amount	PSEUDO RESIDUALS	NEW PREDICTIONS
25	Electronics	2.5	123.45	-32.55	152.745
34	Clothing	1.3	56.78	-99.22	146.078
42	Electronics	5.0	345.67	189.45	174.945
19	Hardware	3.2	98.01	-58.01	150.199

$$\begin{aligned}
 &\text{Initial guess} \downarrow \\
 &= 156 + 0.1 \cdot (-32.55) \\
 &\quad \uparrow \\
 &\quad \text{Learning rate}
 \end{aligned}$$

Leaf value ↓

Step 4: Iterate

→ calculate new pseudo-residuals

→ build new another weak learner

$\hat{y} = \text{Initial guess} + e \cdot (r_1 + r_2 + \dots + r_i)$ after i iterations

Age	Category	Purchase weight (kg)	Amount	PSEUDO RESIDUALS	NEW PREDICTIONS
25	Electronics	2.5	123.45		152.745
34	Clothing	1.3	56.78		146.078
42	Electronics	5.0	345.67		174.945
19	Hardware	3.2	98.01		150.199

XGBoost

- eXtreme Gradient BOOSTting
- GB with several modifications:

XGBoost

- eXtreme Gradient BOOSTting
- GB with several modifications:

Objective function

$$\mathcal{L}(\phi) = \sum_i l(\hat{y}_i, y_i) + \sum_k \Omega(f_k)$$

$$\text{where } \Omega(f) = \gamma T + \frac{1}{2} \lambda \|w\|^2$$

Obj = Loss + Regularization

T = number of leaves in the tree
w_j = weight (output value) of leaf j
γ = penalty for adding a leaf = pruning
λ = L2 regularization on leaf weights

XGBoost

- eXtreme Gradient BOOSTting
- GB with several modifications:

Deciding the split

-> maximizing Gain

$$\text{Gain} = \frac{1}{2} \left(\frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{(G_L + G_R)^2}{H_L + H_R + \lambda} \right) - \gamma$$

G sum of gradients

H sum of Hessians

CART: **unregularized** -> reduce MSE or impurity

XGBoost: regularized -> max Gain

$$\text{Gain} = \frac{1}{2} \left(\underbrace{\frac{G_L^2}{H_L + \lambda}}_{\text{LEFT NODE}} + \underbrace{\frac{G_R^2}{H_R + \lambda}}_{\text{RIGHT NODE}} - \underbrace{\frac{(G_L + G_R)^2}{H_L + H_R + \lambda}}_{\text{ROOT}} \right) - \lambda$$

$$\text{MSE: } L = \frac{1}{2} (y_i - \hat{y}_i)^2$$

$$\text{GRADIENT } g_i = \hat{y}_i - y_i \text{ (residual)}$$

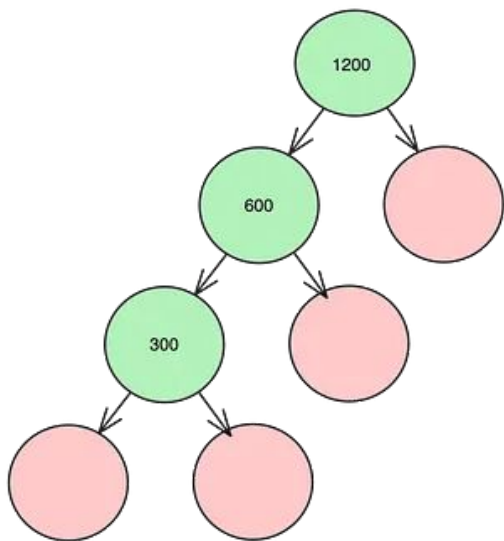
$$\text{HESSIAN } h_i = 1 \text{ (constant)}$$

square of sum
of residuals

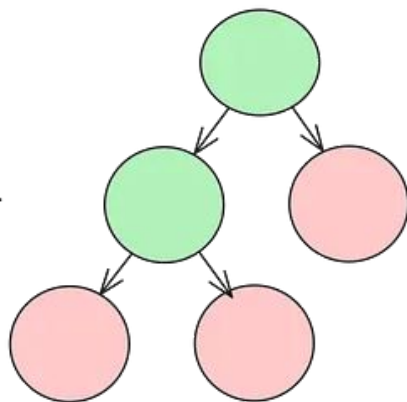
$$\text{NODE SCORE : } \frac{G_L^2}{H_L + \lambda} = \frac{(\sum r_i)^2}{n_L + \lambda} = \frac{(\text{sum of residuals})^2}{\text{number of samples} + \lambda}$$

$$\text{CART: score} = \sum r_i^2 (\text{root}) - (\sum r_i^2 (\text{left}) + \sum r_i^2 (\text{right}))$$

sum of squared residuals



After Pruning
(Gamma = 400)



XGBoost

- eXtreme Gradient BOOSTting
- GB with several modifications:

Output value of a Leaf

$$w_j = -\frac{G_j}{H_j + \lambda}$$

MSE and $\lambda=0$: average of residuals

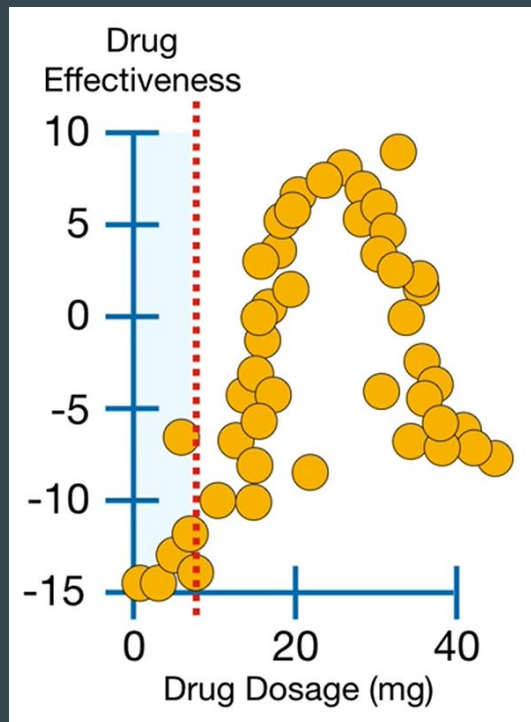
XGBoost

- eXtreme Gradient BOOSTting
- GB with several modifications
- computing predictions: same as GB

XGBoost Optimizations

Approximate Greedy Algorithm

- avoids testing every threshold
- approximation: using quantiles
 - Weighted Quantiles Sketch alg.



XGBoost Optimizations

Approximate Greedy Algorithm

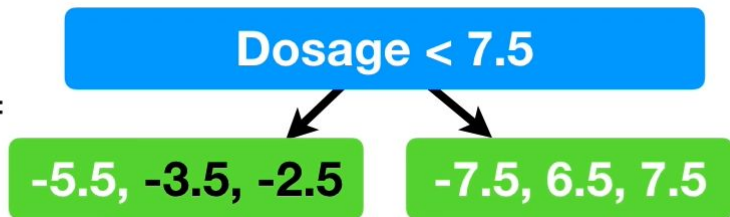
Parallel Learning

Building trees with missing data - Sparsity-Aware Split Finding

Predicted Drug
Effectiveness

0.5

$Gain_{Left} =$



$Gain_{Right} =$



Dosage	Drug Effectiveness	Residuals
10	-7	-7.5
???	-3	-3.5
21	7	6.5
25	8	7.5
5	-5	-5.5
???	-2	-2.5

XGBoost Optimizations

Approximate Greedy Algorithm

Parallel Learning

Sparsity-Aware Split Finding

Cache-Aware access

- storing gradients and hessians in Cache -> speedup

XGBoost Optimizations

Approximate Greedy Algorithm

Parallel Learning

Sparsity-Aware Split Finding

Cache-Aware access

-> other memory optimizations

-> subsampling (rows and columns)

CatBoost

- = CATegorical BOOSTing
- designed for categorical data
- uses Ordered Target Encoding
- symmetric trees - same decision on the same level
 - can be fitted in parallel, very fast with GPU
- -> best for categorical-heavy datasets

LightGBM

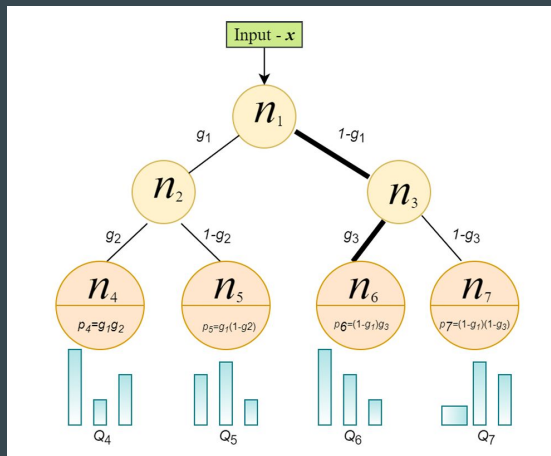
- faster
- leaf-wise tree building - splits the node with the largest split gain first
 - -> unbalanced trees
- low memory usage
- better categorical features handling than XGBoost
- tend to overfit
- -> best for fast training on big datasets

Deep Neural Models for Tabular Data

Two main categories of Deep Models

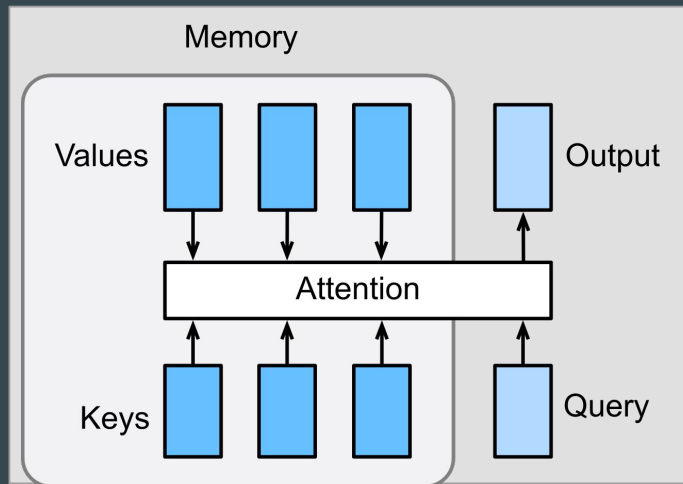
Differentiable Trees

- decision trees are not differentiable
- smoothing the decision function
- softmax, sparsemax



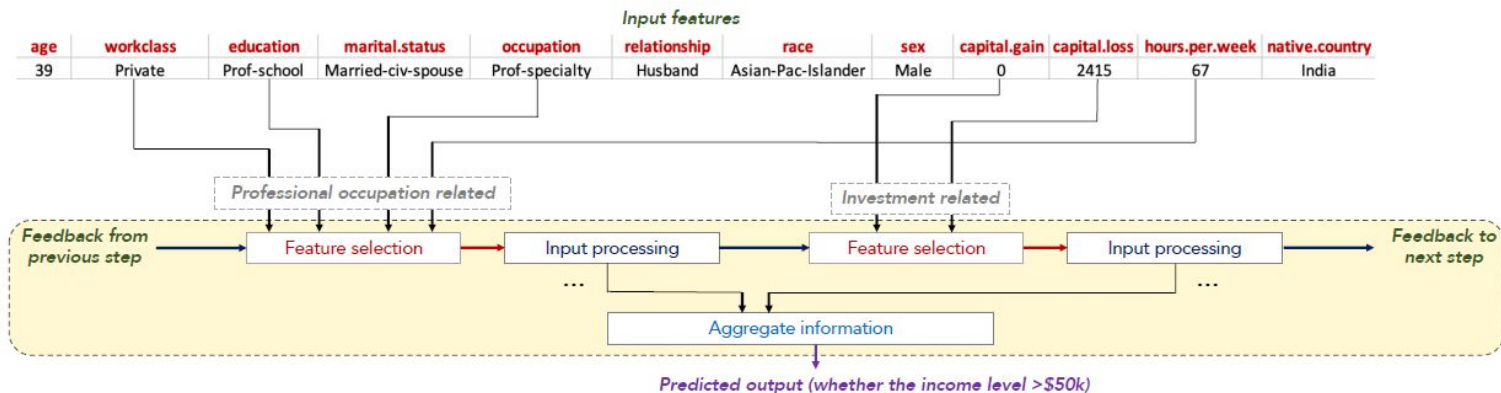
Attention-Based Models

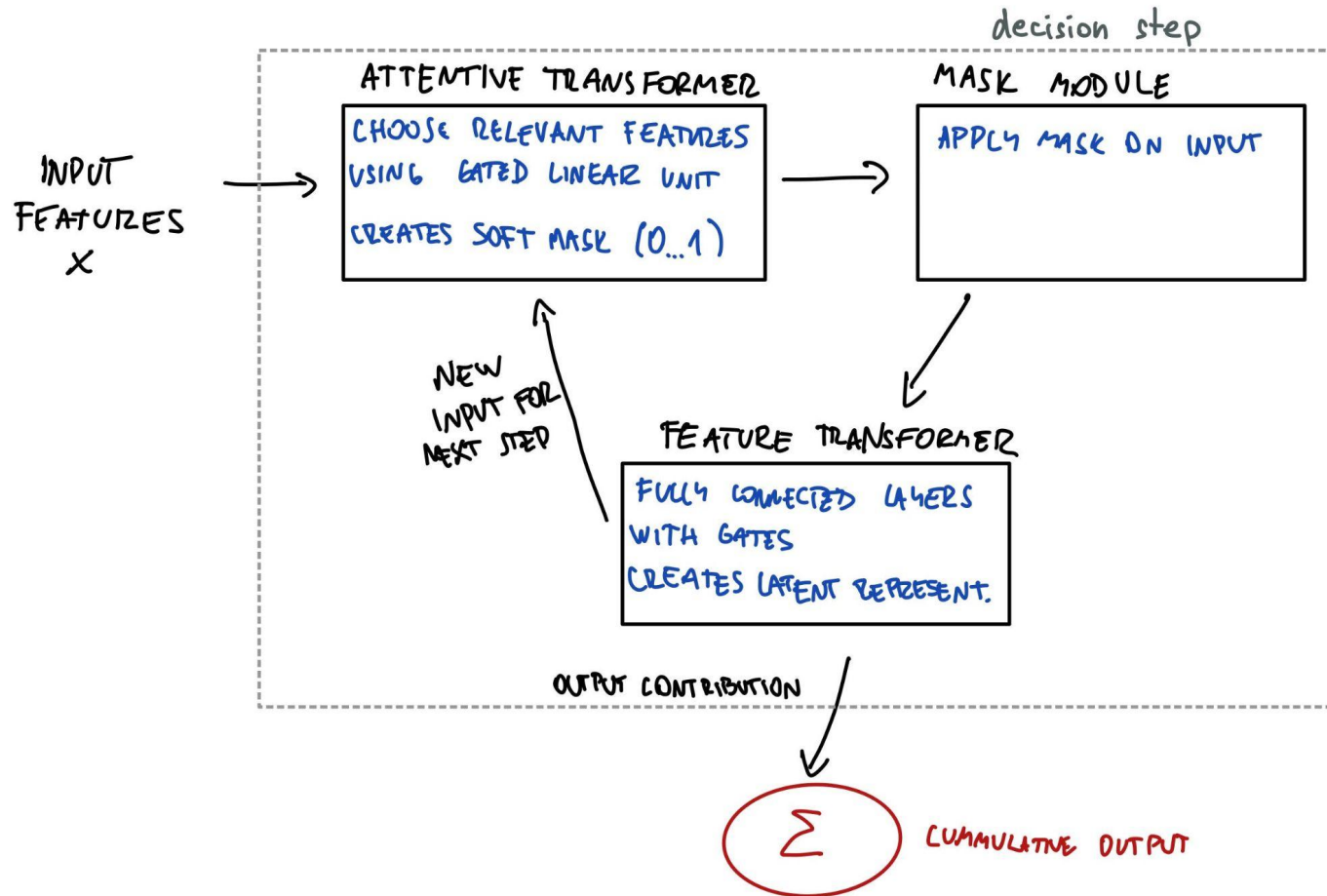
- using attention mechanism detect interactions between features or samples



TabNet

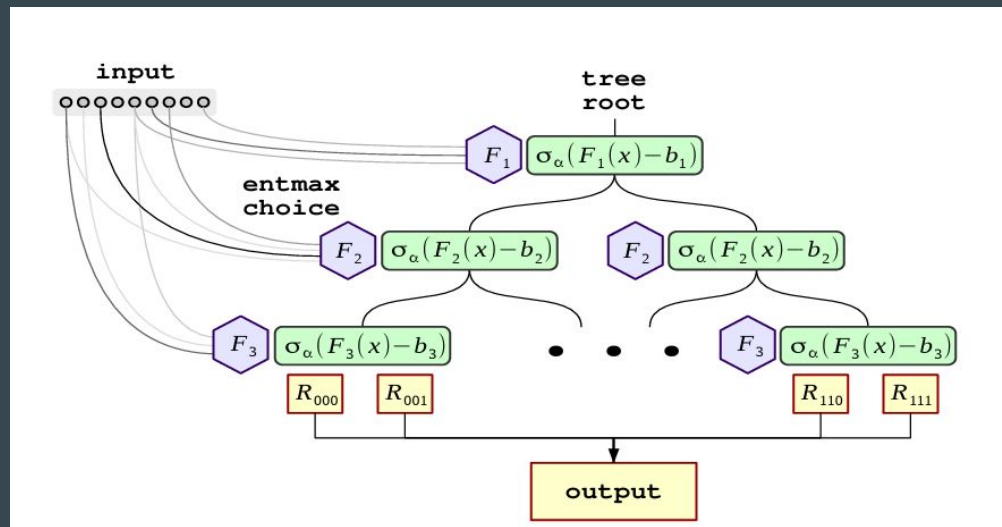
- choosing relevant features for each row using mask and attention
- different relevant features for each data input
- usually 4-10 “decision” steps





Neural Oblivious Decision Ensembles (NODE)

- inspired by CatBoost
- ensemble of **differentiable** “trees”
- not a CART but neural layer
- ODT: all nodes in same depth use the same split rule
- soft routing: softmax instead of 0/1



DNF-Net

- simulating DNFs
- DNF-Net = ensemble of DNNF blocks
- soft AND and soft OR

$$\text{OR}(\mathbf{x}) \triangleq \tanh\left(\sum_{i=1}^d \mathbf{x}_i + d - 1.5\right)$$

$$\text{AND}(\mathbf{x}) \triangleq \tanh\left(\sum_{i=1}^d \mathbf{x}_i - d + 1.5\right)$$

$$(A \wedge \neg B \wedge \neg C) \vee (\neg D \wedge E \wedge F \wedge D \wedge F)$$

INPUT: $x \in \mathbb{R}^d$

LEARNABLE

DNNF Block

LEARN EMBEDDINGS $L(x)$

$$L(x) = \tanh(x^T w + b) \in \mathbb{R}^m$$

LITERALS = m



CONJUNCTIONS OVER L AND

CONJUNCTIONS = k



OUTPUT - NEURAL OR GATE

$$\text{DNNF}(x) = \text{OR}([\text{AND}_{c^1}(L(x)), \text{AND}_{c^2}(L(x)), \dots, \text{AND}_{c^k}(L(x))])$$

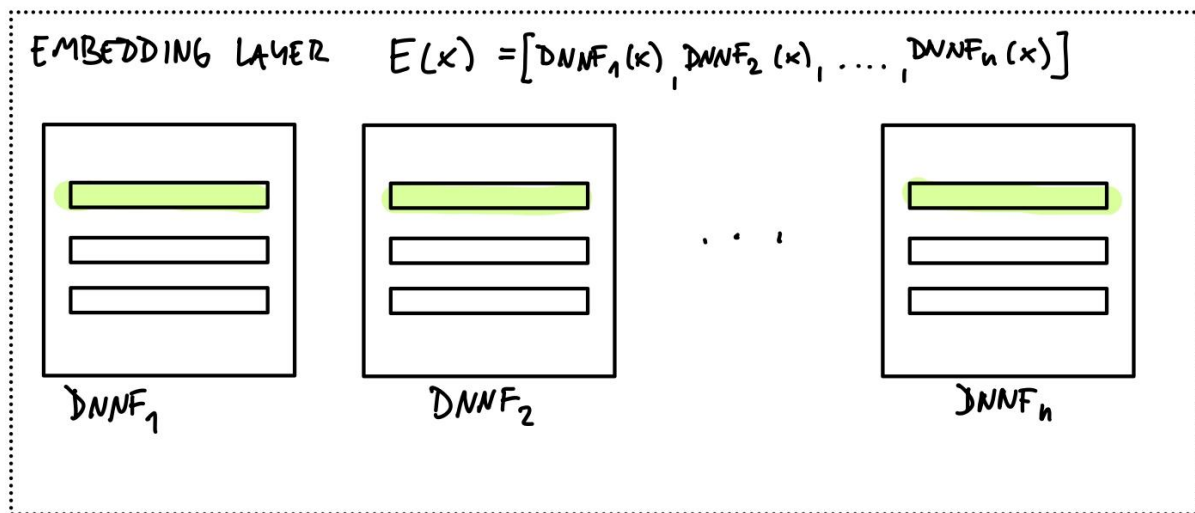
c^1, c^2, \dots, c^k FIXED MASKING VECTORS



OUTPUT: HOW WELL x SATISFIES THE
DNF FORMULA

DNF-Net

LEARNABLE



OUTPUT LAYER OVER $E(x)$

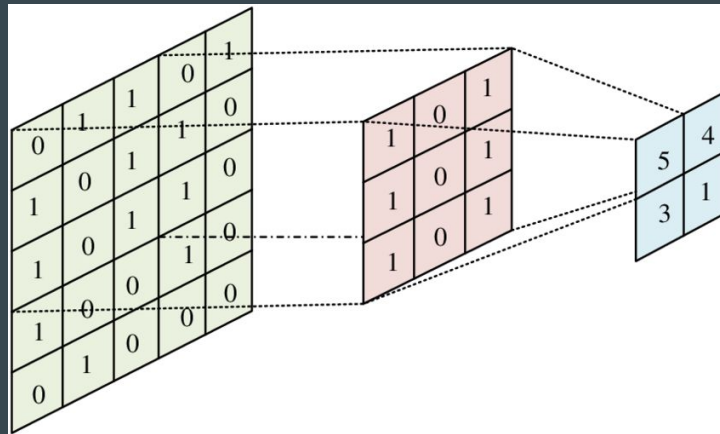
E.G. BINARY CLASSIFICATION $DNF-Net(x) = \sigma \left(\sum_{i=1}^n w_i DNNF_i(x) + b_i \right)$



OUTPUT (E.G. $[0, 1]$)

1D-CNN

- kaggle multi-label classification, probability prediction
 - CNN for feature extraction
 - expect local correlation between features
 - order of columns usually does not matter
-
- idea: re-order features to make them correlated
 - expand dimensionality -> soft-ordering
 - apply CNNs



Performance Comparison: Deep Models and XGBoost

Data Description, Training and Evaluation Metrics

- datasets from DNF-Net, NODE, TabNet papers + 2 Kaggle
- 11 classification and regression tasks
- cross-entropy loss, RMSE
- testing significance of performance differences

Dataset	Features	Classes	Samples	Source	Paper
Gesture Phase	32	5	9.8k	OpenML	DNF-Net
Gas Concentrations	129	6	13.9k	OpenML	DNF-Net
Eye Movements	26	3	10.9k	OpenML	DNF-Net
Epsilon	2000	2	500k	PASCAL Challenge 2008	NODE
YearPrediction	90	1	515k	Million Song Dataset	NODE
Microsoft (MSLR)	136	5	964k	MSLR-WEB10K	NODE
Rossmann Store Sales	10	1	1018K	Kaggle	TabNet
Forest Cover Type	54	7	580k	Kaggle	TabNet
Higgs Boson	30	2	800k	Kaggle	TabNet
Shrutime	11	2	10k	Kaggle	New dataset
Blastchar	20	2	7k	Kaggle	New dataset

Table 1: Description of the tabular datasets

Requirements for Real-World Applications

- high accuracy
- efficient training and inference
- short optimization time

Testing Deep Models on Unseen Datasets

- datasets not included in original paper
- models perform worse than datasets original models
- no deep model consistently outperformed others
- XGBoost generally outperformed deep models (Eye, Gesture, MSLR, Shrutime, Blastchar)
- XGBoost significantly better than not original models in 8 of 11 cases
- ensemble of deep models and XGBoost usually outperformed other models

Model Name	Rossmann	CoverType	Higgs	Gas	Eye	Gesture
XGBoost	490.18 \pm 1.19	3.13 \pm 0.09	21.62 \pm 0.33	2.18 \pm 0.20	56.07 \pm 0.65	80.64 \pm 0.80
NODE	488.59 \pm 1.24	4.15 \pm 0.13	21.19 \pm 0.69	2.17 \pm 0.18	68.35 \pm 0.66	92.12 \pm 0.82
DNF-Net	503.83 \pm 1.41	3.96 \pm 0.11	23.68 \pm 0.83	1.44 \pm 0.09	68.38 \pm 0.65	86.98 \pm 0.74
TabNet	485.12 \pm 1.93	3.01 \pm 0.08	21.14 \pm 0.20	1.92 \pm 0.14	67.13 \pm 0.69	96.42 \pm 0.87
1D-CNN	493.81 \pm 2.23	3.51 \pm 0.13	22.33 \pm 0.73	1.79 \pm 0.19	67.9 \pm 0.64	97.89 \pm 0.82
Simple Ensemble	488.57 \pm 2.14	3.19 \pm 0.18	22.46 \pm 0.38	2.36 \pm 0.13	58.72 \pm 0.67	89.45 \pm 0.89
Deep Ensemble w/o XGBoost	489.94 \pm 2.09	3.52 \pm 0.10	22.41 \pm 0.54	1.98 \pm 0.13	69.28 \pm 0.62	93.50 \pm 0.75
Deep Ensemble w XGBoost	485.33 \pm 1.29	2.99 \pm 0.08	22.34 \pm 0.81	1.69 \pm 0.10	59.43 \pm 0.60	78.93 \pm 0.73

TabNet

DNF-Net

Model Name	YearPrediction	MSLR	Epsilon	Shrutime	Blastchar
XGBoost	77.98 \pm 0.11	55.43 \pm 2e-2	11.12 \pm 3e-2	13.82 \pm 0.19	20.39 \pm 0.21
NODE	76.39 \pm 0.13	55.72 \pm 3e-2	10.39 \pm 1e-2	14.61 \pm 0.10	21.40 \pm 0.25
DNF-Net	81.21 \pm 0.18	56.83 \pm 3e-2	12.23 \pm 4e-2	16.8 \pm 0.09	27.91 \pm 0.17
TabNet	83.19 \pm 0.19	56.04 \pm 1e-2	11.92 \pm 3e-2	14.94 \pm , 0.13	23.72 \pm 0.19
1D-CNN	78.94 \pm 0.14	55.97 \pm 4e-2	11.08 \pm 6e-2	15.31 \pm 0.16	24.68 \pm 0.22
Simple Ensemble	78.01 \pm 0.17	55.46 \pm 4e-2	11.07 \pm 4e-2	13.61 \pm , 0.14	21.18 \pm 0.17
Deep Ensemble w/o XGBoost	78.99 \pm 0.11	55.59 \pm 3e-2	10.95 \pm 1e-2	14.69 \pm 0.11	24.25 \pm 0.22
Deep Ensemble w XGBoost	76.19 \pm 0.21	55.38 \pm 1e-2	11.18 \pm 1e-2	13.10 \pm 0.15	20.18 \pm 0.16

NODE

New datasets

Direct Comparison

- calculating relative performance to the best model
- average relative performance per model on unseen datasets
- best: Ensemble w XGBoost
- best single model: XGBoost
- deep models: last 4 places

Name	Average Relative Performance (%)
XGBoost	3.34
NODE	14.21
DNF-Net	11.96
TabNet	10.51
1D-CNN	7.56
Simple Ensemble	3.15
Deep Ensemble w/o XGBoost	6.91
Deep Ensemble w XGBoost	2.32

Why are the Deep Models Worse?

- selection bias
- hyperparameters optimization (more extensive search could lead to better performance)

Do we need both?

- ensemble of deep models with XGBoost were the best
- which component of ensemble is mandatory?
- do we need deep models?
- (XGBoost, SVM, CatBoost)
- (deep model ensemble)
- (deep models + XGBoost)

Comparison of Ensemble Models

- deep models ensemble did not perform well
- classical ensemble better than deep models
- however, deep models combined with XGBoost are significantly better than classical ensemble or XGBoost itself

Name	Average Relative Performance (%)
XGBoost	3.34
NODE	14.21
DNF-Net	11.96
TabNet	10.51
1D-CNN	7.56
Simple Ensemble	3.15
Deep Ensemble w/o XGBoost	6.91
Deep Ensemble w XGBoost	2.32

Performance & Computation Demand Trade-off

- + performance
 - + expensive computation
 - Do we need all models in the ensemble?
 - How to choose the subset?
-
- 3 models \approx full ensemble
 - significant differences between selection approaches

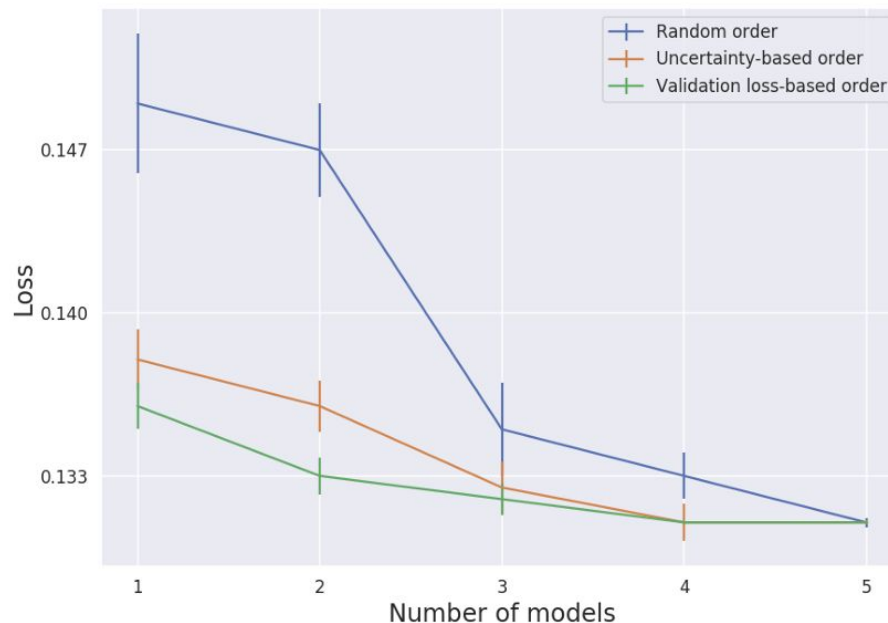
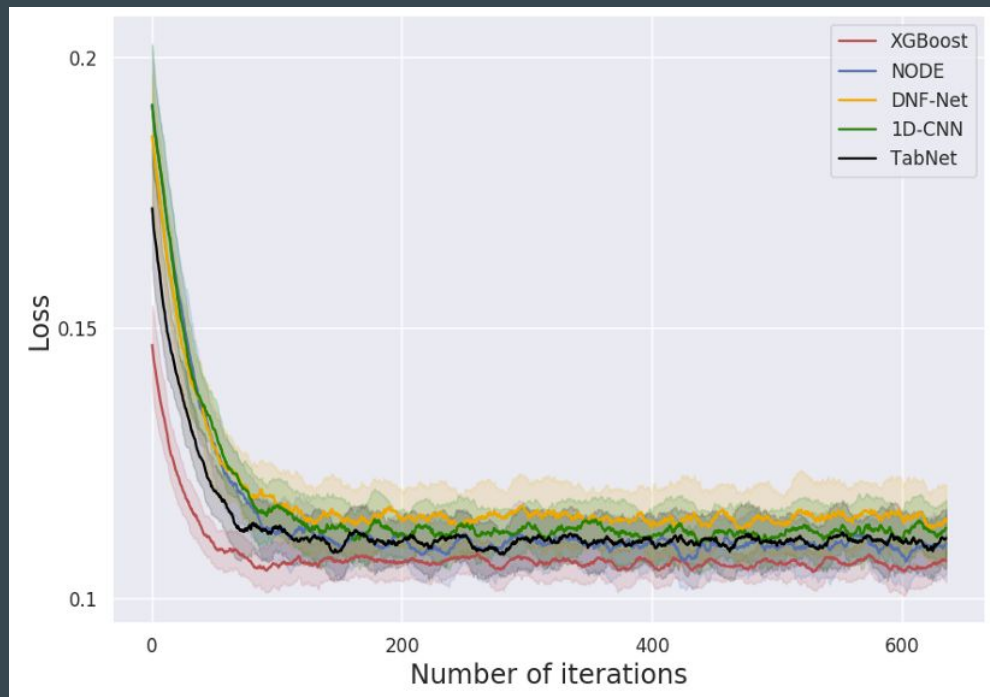


Figure 1: The impact of selecting a subset of models in the ensemble.

How difficult is the hyperparameter optimization?

- limited time for training in real applications
- FLOPS not suitable - different architectures
- runtime affected by software optimizations (XGBoost)
- proxy measure: number of optimization iterations
- XGBoost much faster than deep models
-



Summary

- Deep models often underperform on unseen datasets.
- XGBoost remains a strong baseline.

Summary

- Deep models often underperform on unseen datasets.
- XGBoost remains a strong baseline.
- Ensembles combining DL with XGBoost outperform individual models.

Summary

- Deep models often underperform on unseen datasets.
- XGBoost remains a strong baseline.
- Ensembles combining DL with XGBoost outperform individual models.
- DL models harder to optimize

Summary

- Deep models often underperform on unseen datasets.
- XGBoost remains a strong baseline.
- Ensembles combining DL with XGBoost outperform individual models.
- DL models harder to optimize

Limited time? -> XGBoost

Maximal performance? -> Add deep models to an ensemble

Grinsztajn et al.: Why do tree-based models still outperform deep learning on typical tabular data?

Random Forest

XGBoost

GBTree

SAINT

Resnet

FT_Transformer

MLP

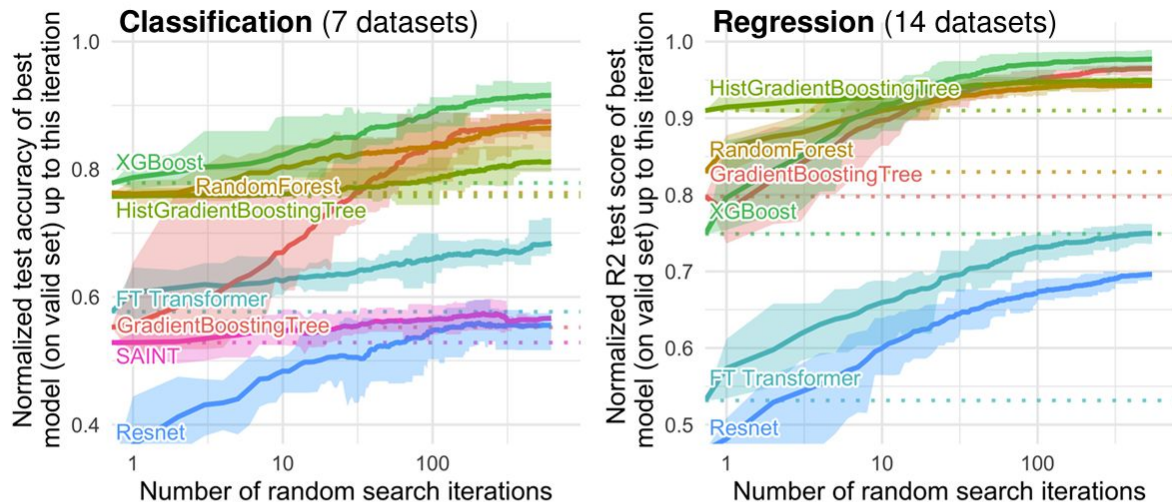


Figure 2: **Benchmark on medium-sized datasets, with both numerical and categorical features.** Dotted lines correspond to the score of the default hyperparameters, which is also the first random search iteration. Each value corresponds to the test score of the best model (on the validation set) after a specific number of random search iterations, averaged on 15 shuffles of the random search order. The ribbon corresponds to the minimum and maximum scores on these 15 shuffles.

