# Feature Importance in Decision Trees:

## Impurity-based Importance Calculations & Explainable AI

## Dr. Franziska Boenisch

CISPA
HELMHOLTZ CENTER FOR
INFORMATION SECURITY

SprintML

# You already interact(ed) with decision trees!

What pet should I get?

Cat    Dog

May 11, 2023  2m read

## Random Forests: Netflix Customer Recommendations Improved by 20%

… watching Netflix …

## Spotify — Decision Trees with Music Taste

7 min read · Nov 26, 2020

J  Jinkim  Follow

… or listening to music.

# Outline for Today

- Intuition on classification with decision tree

- Impurity-based feature importance metrics

- Building decision trees based on impurity reduction

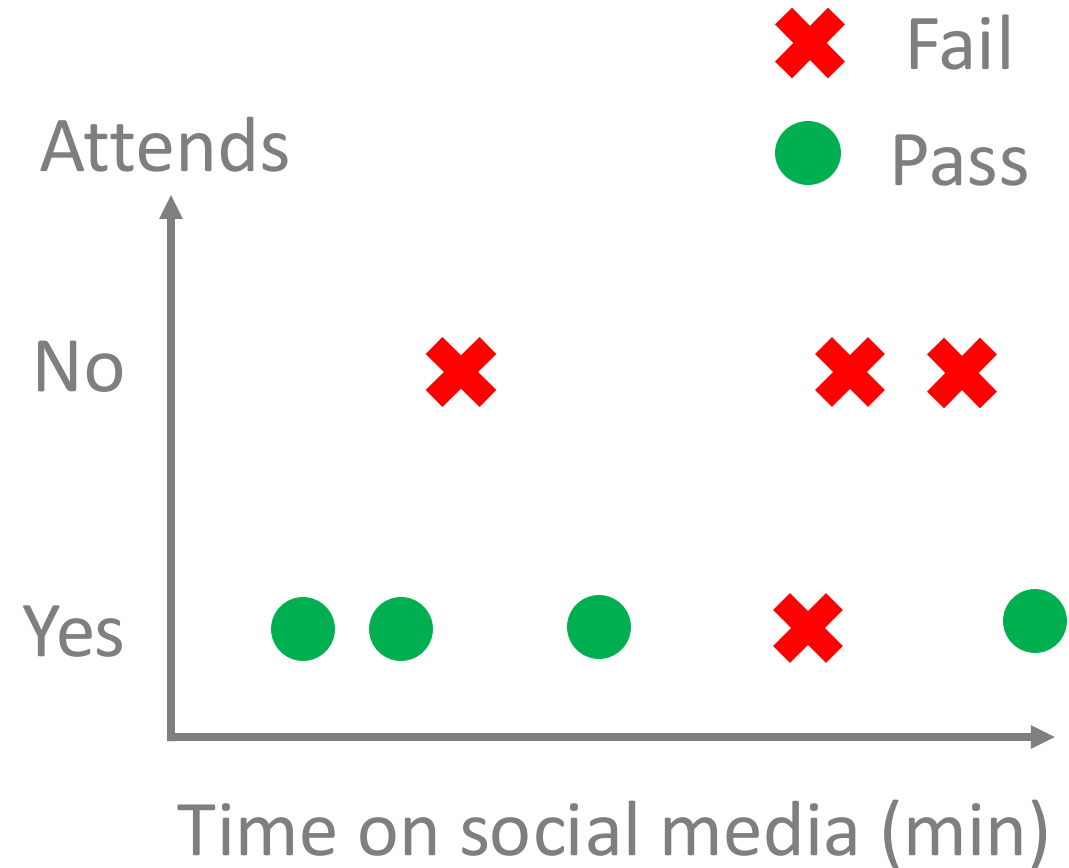- Feature importance and explainable AI

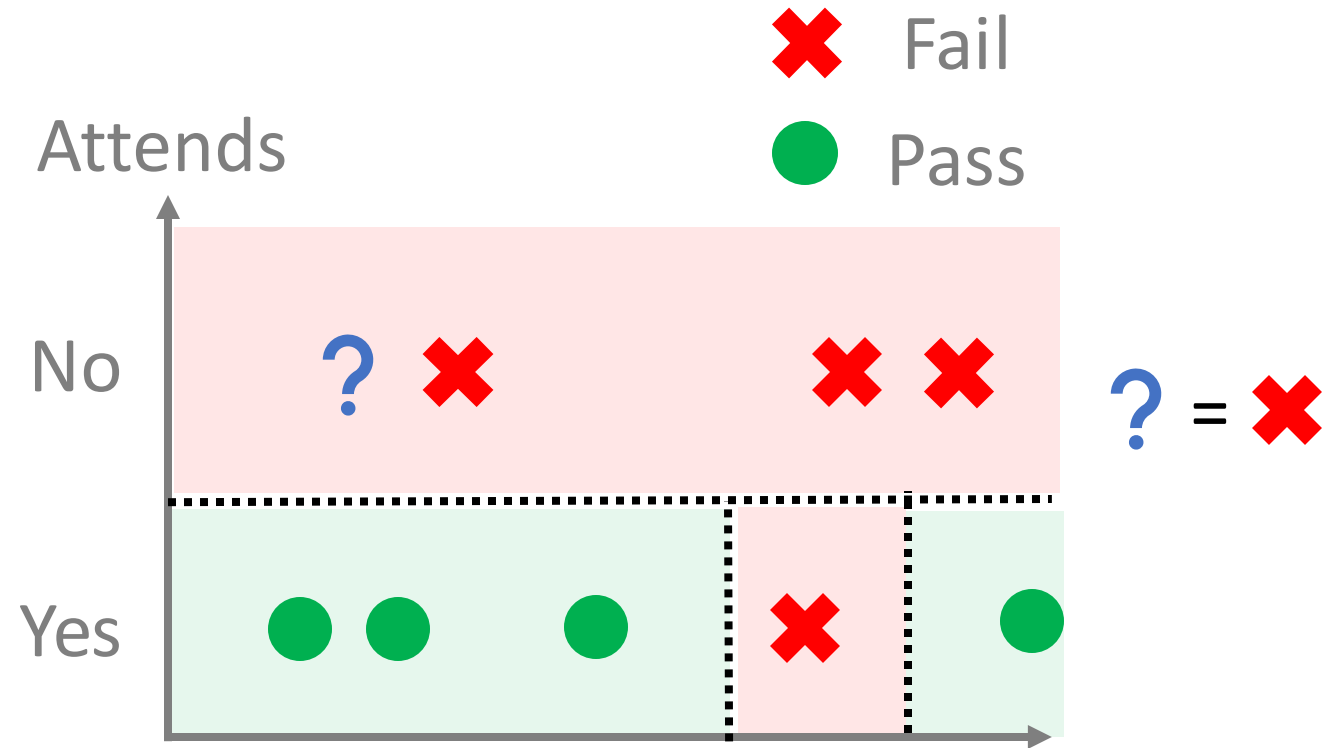# Datasets and Tree-based Classification
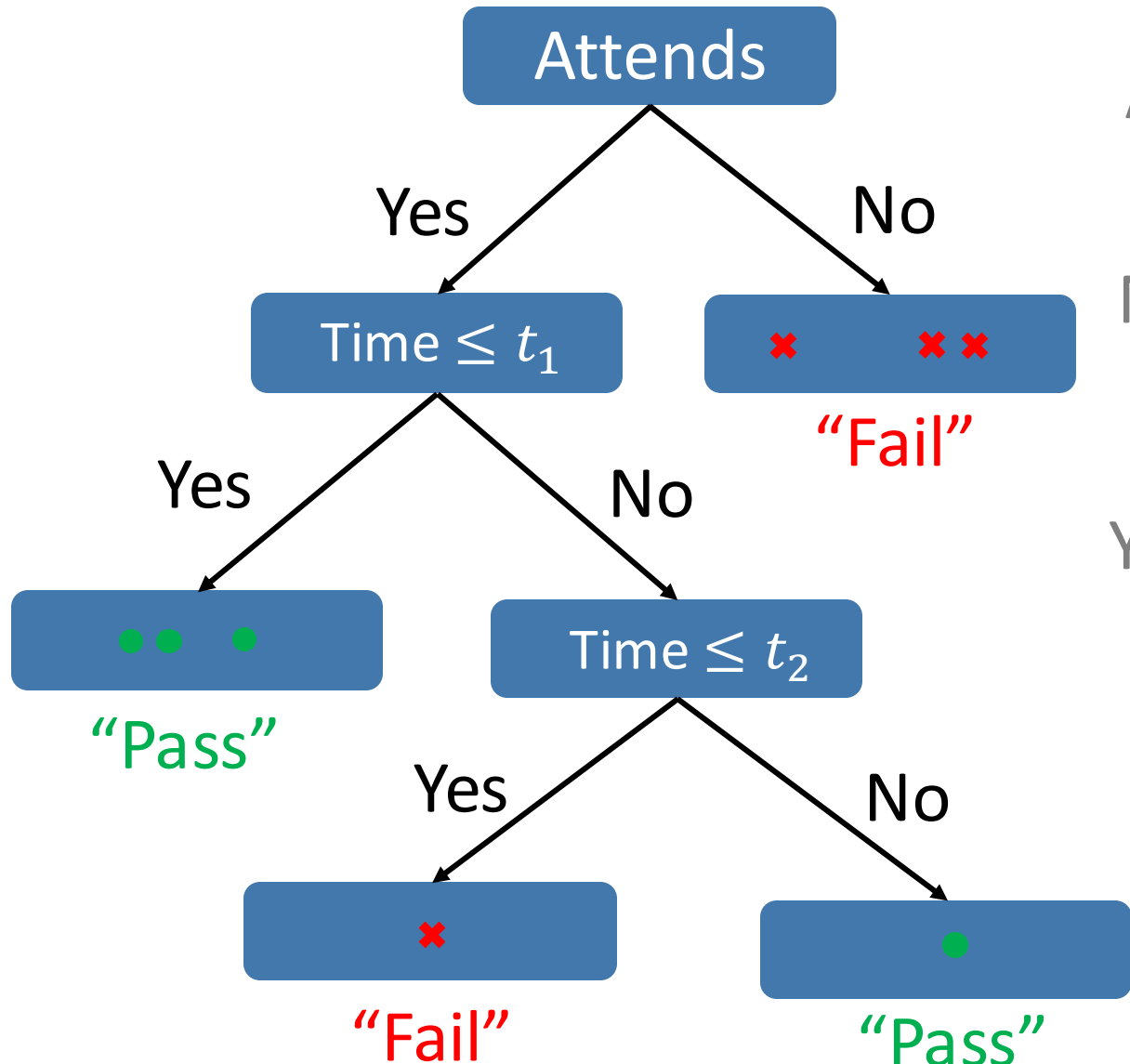
## Dataset from my Lecture

| Social Media Time (min) | Attends Class | Passed the Midterm |
|---|---|---|
| 30 | Yes | Pass |
| 80 | Yes | Pass |
| 140 | Yes | Pass |
| 50 | Yes | Pass |
| 110 | No | Fail |
| 60 | No | Fail |
| 100 | Yes | Fail |
| 120 | No | Fail |

Continuous        Categorical

Features        Label

❌ Fail

🟢 Pass

Attends
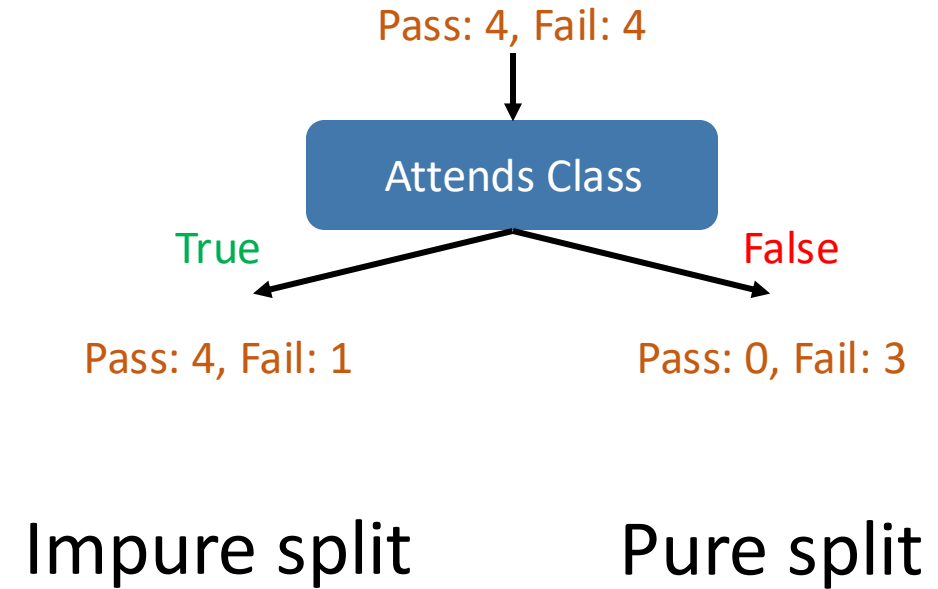
No

Yes

Time on social media (min)

# Datasets and Tree-based Classification



Successively find splits of the data to create decision regions

Greedy, recursive partitioning

# Finding the Best Split Criterion

| Social Media Time (min) | Attends Class | Passed the Midterm |
|:---:|:---:|:---:|
| 30 | Yes | Pass |
| 80 | Yes | Pass |
| 140 | Yes | Pass |
| 50 | Yes | Pass |
| 110 | No | Fail |
| 60 | No | Fail |
| 100 | Yes | Fail |
| 120 | No | Fail |

Is attendance our best spilt?

Pass: 4, Fail: 4

Attends Class

True                                    False

Pass: 4, Fail: 1          Pass: 0, Fail: 3

Impure split              Pure split
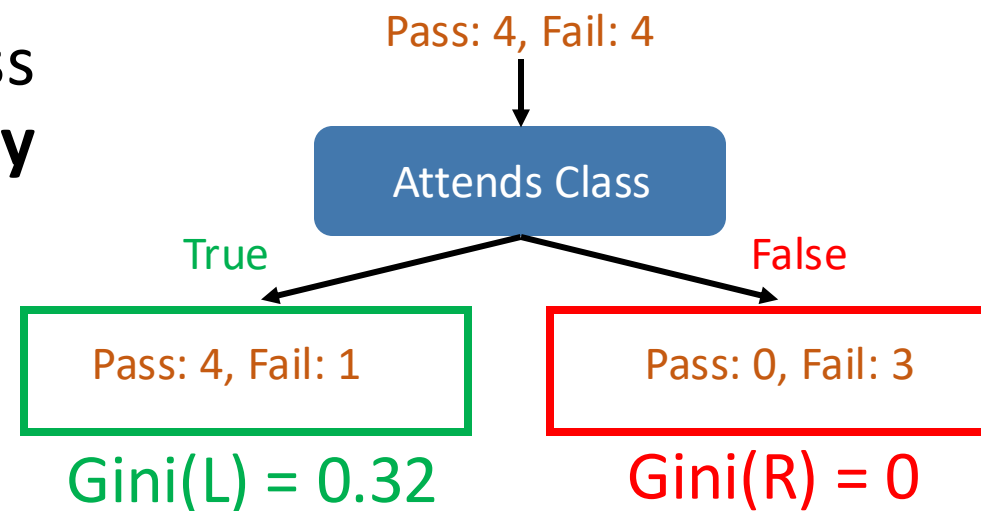
# Gini Impurity: Definition

Given a node with $K$ classes and class probabilities $p_1, \ldots, p_k$. The **Gini Impurity** is defined as

$$Gini(Node) = 1 - \sum_{k=1}^{K} p_k^2.$$

Here: $1 - \left( p_{pass}^2 + p_{\text{fail}}^2 \right)$

Attends Class

True           False

| Pass: 4, Fail: 1 | Pass: 0, Fail: 3 |
|---|---|
| Gini(L) = 0.32 | Gini(R) = 0 |

$p_{pass} = \frac{4}{5}, p_{fail} = \frac{1}{5}$      $p_{pass} = \frac{0}{3}, p_{fail} = \frac{3}{3}$
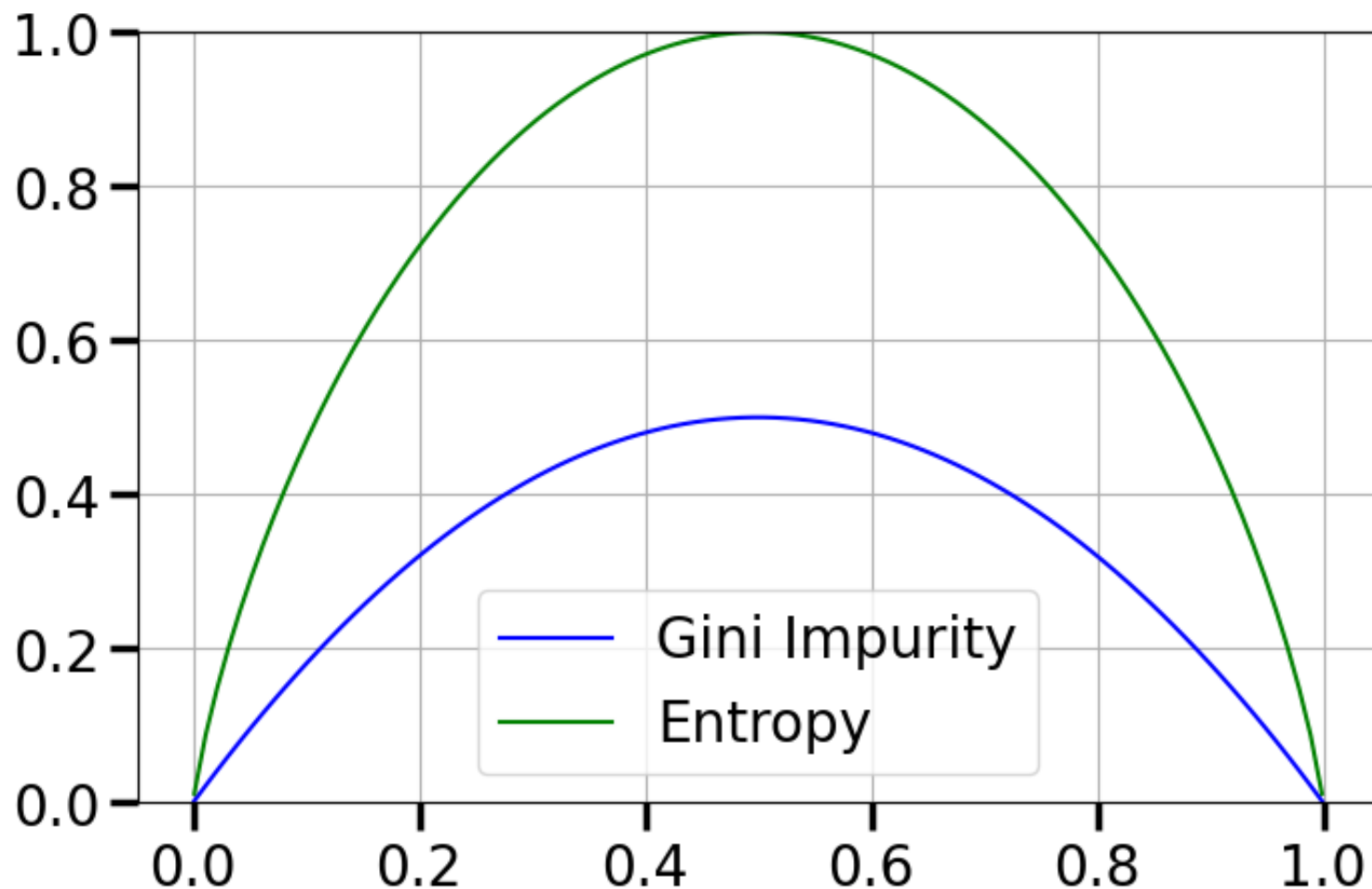
In our example:

- True branch: Gini(L) $= 1 - \left(\frac{4}{5}\right)^2 - \left(\frac{1}{5}\right)^2 = 0.32$

- False branch: Gini(R) $= 1 - (0)^2 - (1)^2 = 0$

*Gini of a pure split is zero*
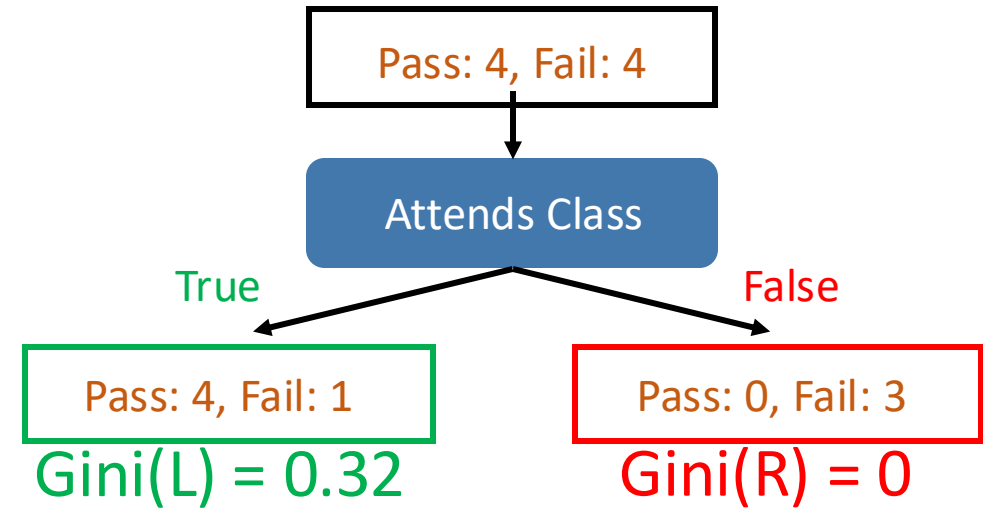
7

# Gini Impurity and Entropy

Impurity



Proportion of Positive Class

Alternative to Gini:

**Entropy** $= -\sum_k p_k \log_2 p_k$,

with $p_k$ : proportion of data from class $k$ in the node.

# Gini Impurity of the Entire Split

When evaluating a split, we compute the **weighted Gini of the children**:

$$Gini_{split} = \frac{n_L}{n} Gini(L) + \frac{n_R}{n} Gini(R)$$

Pass: 4, Fail: 4

Attends Class

True          False

Pass: 4, Fail: 1          Pass: 0, Fail: 3

Gini(L) = 0.32          Gini(R) = 0

$n$ : total number of data points in split → 8          $Gini_{split} = 0.2$

$n_L$: number of points in L(eft) node → 5

$n_R$: number of points in R(right) node → 3

$$Gini_{split} = \frac{5}{8} * 0.32 + \frac{3}{8} * 0 = 0.2$$
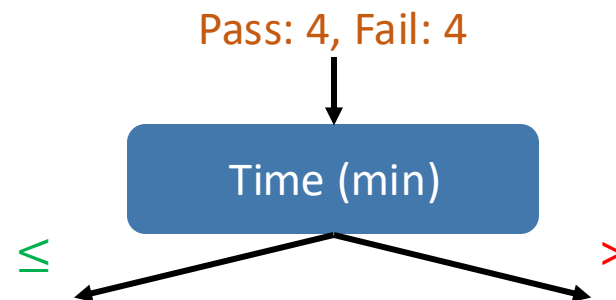
# Gini Impurity on Continuous Values

Goal: Identify the best splitting threshold

| Social Media Time (min) | Passed the Midterm |
|---|---|
| 30 | Pass |
| 50 | Pass |
| 60 | Fail |
| 80 | Pass |
| 100 | Fail |
| 110 | Fail |
| 120 | Fail |
| 140 | Pass |

2. Identify where class changes

3. Take average as threshold

1. Sort ascending

Pass: 4, Fail: 4

Time (min)

≤      >

≤ 55

≤ 70

≤ 90

≤ 130

# Gini Impurity on Continuous Values

| Social Media Time (min) | Passed the Midterm |
|---|---|
| 30 | Pass |
| 50 | Pass |
| 60 | Fail |
| 80 | Pass |
| 100 | Fail |
| 110 | Fail |
| 120 | Fail |
| 140 | Pass |

0.33
0.48
0.38

0.44

Pass: 4, Fail: 4

Time (min)    $\leq 55$

$\leq$                    $>$

Pass: 2, Fail: 0          Pass: 2, Fail: 4

$Gini(L) = 0$    $Gini(R) \approx 0.44$

$Gini_{split} \approx 0.33$

# Impurity Reduction to Choose the Best Split

Choose the split that causes the maximum **impurity reduction** $\Delta i$(split):

$$\Delta i \, (split) = \max(Gini_{parent} - Gini_{split})$$

Attends

$$Gini_{parent} = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = 0.5$$

*is fixed* → *find lowest* $Gini_{split}$

$$\Delta i \, (Attends) = 0.5 - 0.2 = 0.3$$

Gini impurity over all possible splits:

$$Gini_{Attends} = 0.2$$
$$Gini_{Time \leq 55} = 0.33$$
$$Gini_{Time \leq 70} = 0.48$$
$$Gini_{Time \leq 90} = 0.38$$
$$Gini_{Time \leq 130} = 0.44$$
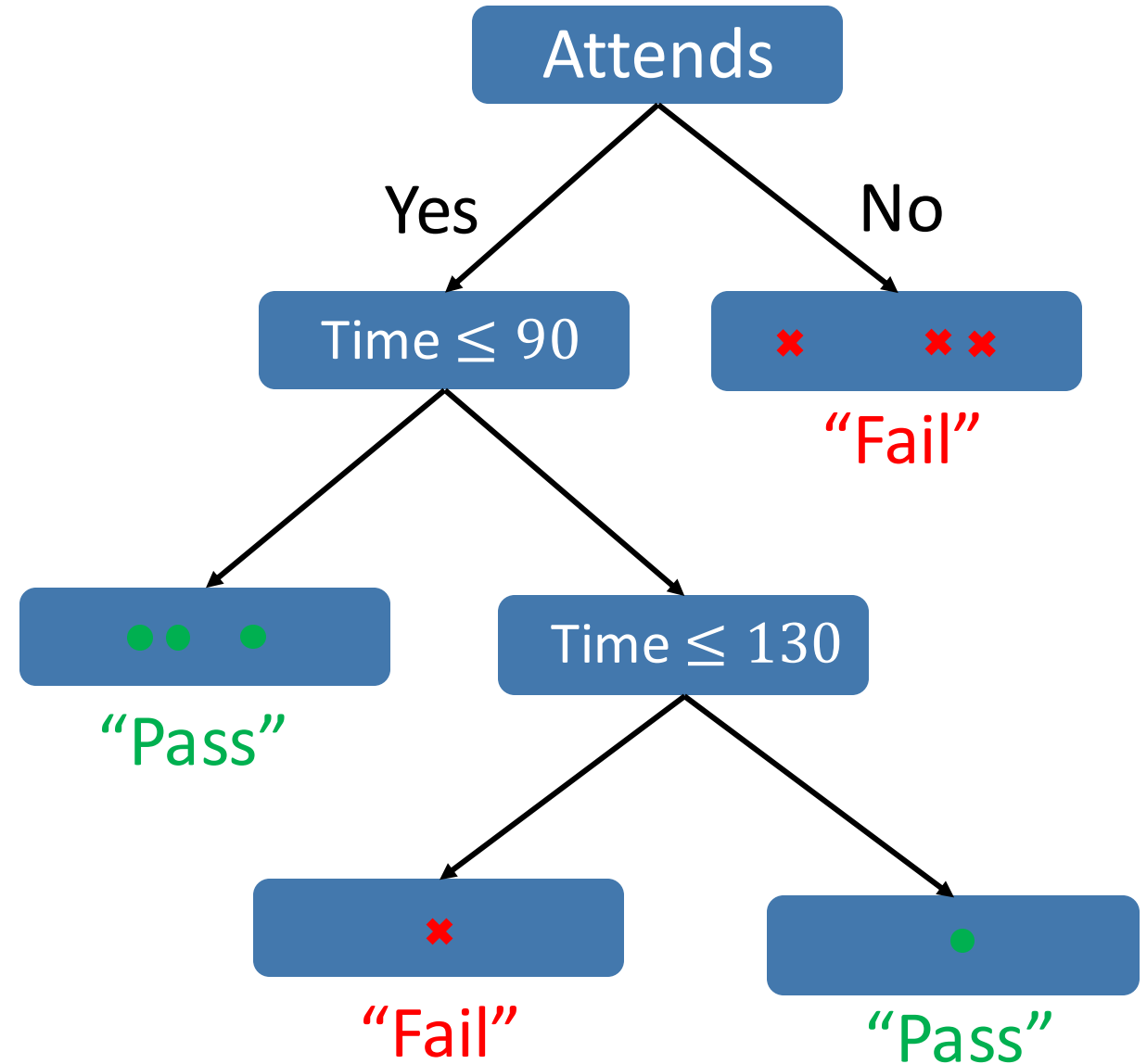
# From Trees to Explainable AI

Decisions in the tree are:

- Human-interpretable
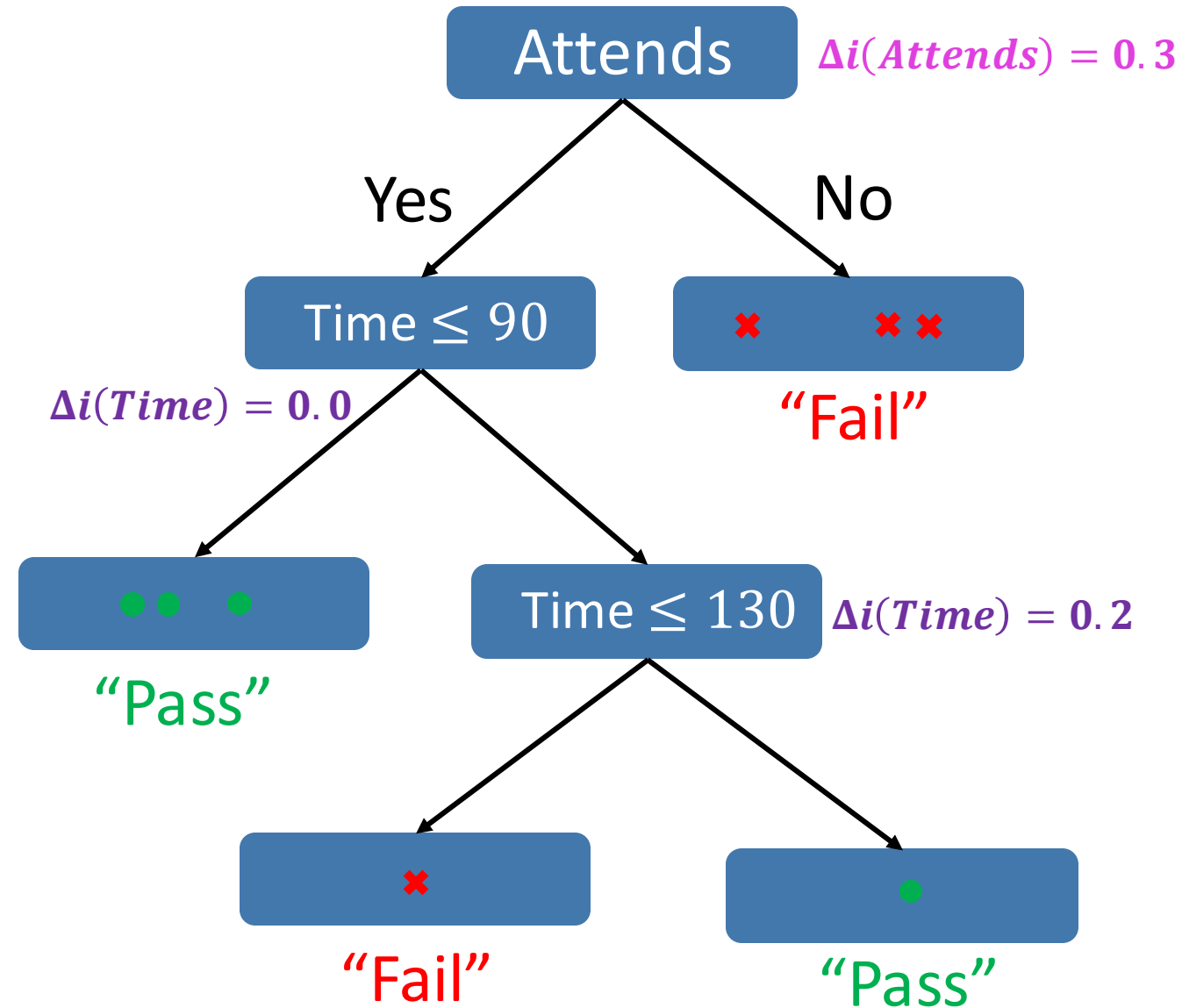- Verifiable
- We can ask "What if?" (Counterfactuals)

# Impurity-based Feature importance

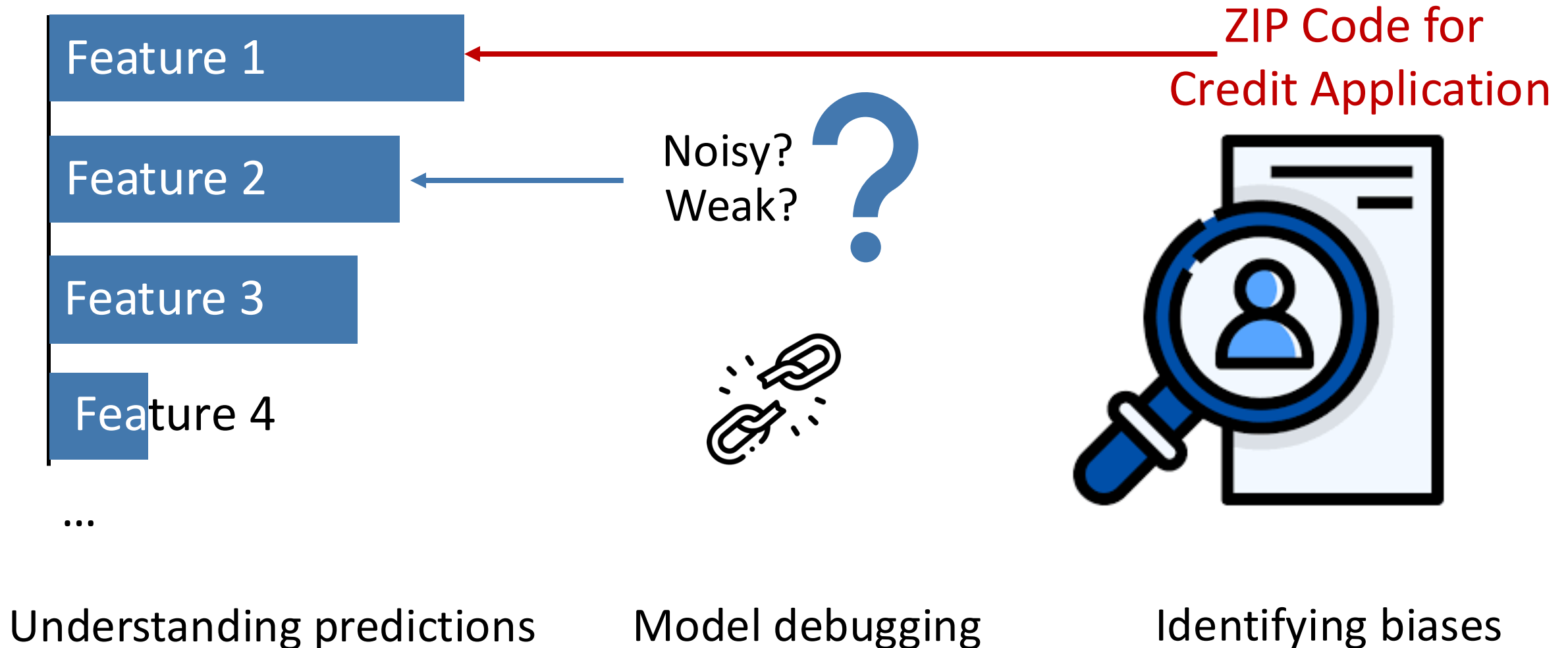We calculate the **importance** of a feature $f$ in a tree as:

$$Importance(f) = \frac{\sum_{t \in Splits\ on\ f} \Delta i(t)}{\sum_{s \in All\ splits} \Delta i(s)}.$$

$Importance(Attends)$

$$= \frac{0.3}{0.3 + 0.0 + 0.2} = 0.6 = 60\%$$

$Importance(Time)$

$$= \frac{0.0 + 0.2}{0.3 + 0.0 + 0.2} = 0.4 = 40\%$$



$\Delta i(Attends) = 0.3$

Yes      No

Attends

Time $\leq 90$

$\Delta i(Time) = 0.0$

✖   ✖ ✖

"Fail"

● ●   ●

"Pass"

Time $\leq 130$   $\Delta i(Time) = 0.2$

✖

"Fail"

●

"Pass"

# Feature Importance for Explainability



Understanding predictions        Model debugging        Identifying biases

# Summary & Lecture Materials



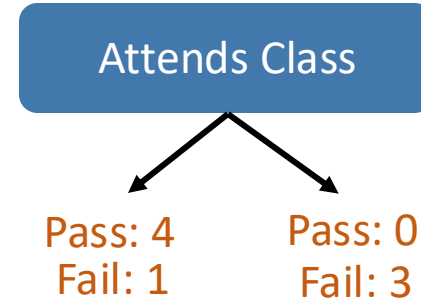**Decision Trees: Omnipresent**

**Divide Data in Regions**

**Impurity-based Feature Splits**

**Serve Explainable AI**

Lecture Materials: