1) **General supervised learning:**

$$\hat{y} = f(\vec{x}) \approx y \quad \dots \qquad \downarrow \mathcal{L}(y, \hat{y})$$

2) **Success depends on:**

quantity & quality of data & features

3) **Validation is essential:**

CV to tune hyperparameters (IID)
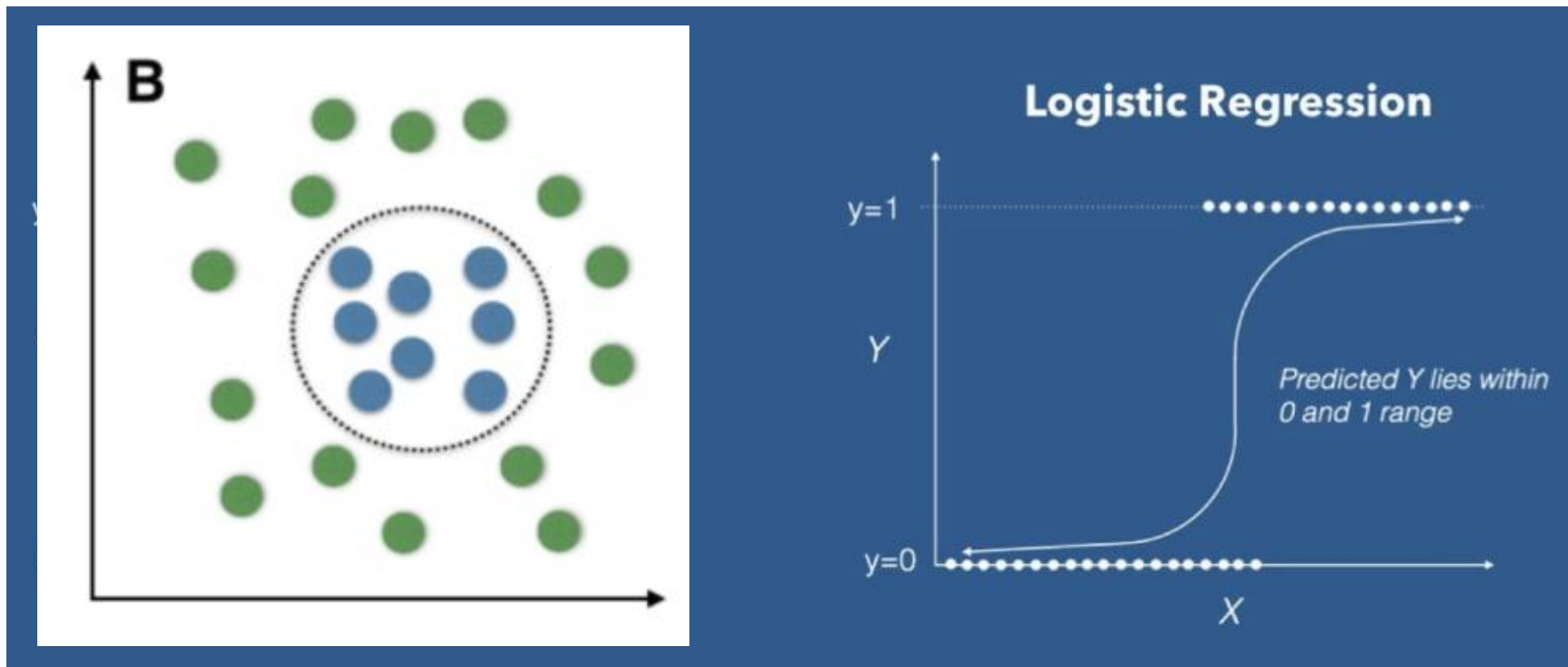↳ test set to further probe generalizability.

4) **Leverage your domain knowledge:**

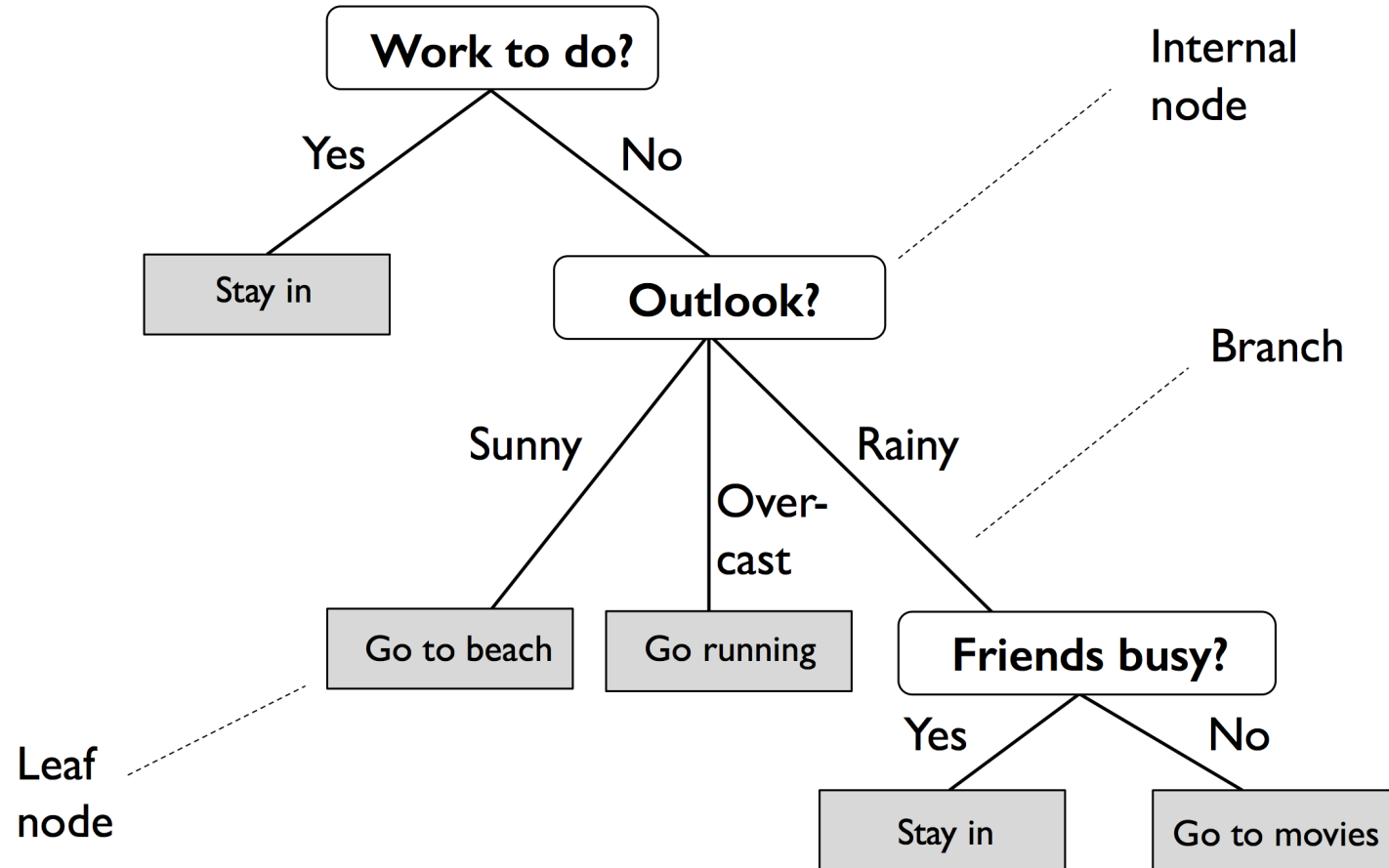better features, more robust data, define success

# What about nonlinear classification?

For binary classification, y is no longer continuous, but binomial:
**y = [1, 1, 1, -1, -1, 1, -1, -1, ...]**

# Brief intro to decision trees

**Work to do?**

Yes     No

Stay in

**Outlook?**

Internal node

Sunny    Over-cast    Rainy

Branch

Go to beach    Go running    **Friends busy?**

Leaf node

Yes    No

Stay in    Go to movies

A

B

if MP > 1000K (40, 40) < 1000K    .f m < 100 (40, 40)   if m > 100

(30, 10)    (10, 30)    (20, 40)    (20, 0)

40 metals
40 insulators

metals                                    insulators

CEMS
Chemical Engineering
& Materials Science

Determine splits by **maximizing information gain (IG)**

**minimizing weighted impurity, I**

$$I(n) = -\sum_{i=0}^{C} p(i|n) \ln\{p(i|n)\}$$

$$c = 0 \quad c = 1$$

Consider $n = 0 \quad (40, 40)$

$$I_0 = -\left\{\left(\frac{1}{2}\ln\frac{1}{2}\right) + \left(\frac{1}{2}\ln\frac{1}{2}\right)\right\}$$

$$= 0.69$$

$$\frac{S}{R_B} \qquad -\frac{S}{R_B} = x\ln x + (1-x)\ln(1-x)$$

$$A_{1-x} B_x$$

5

Determine splits by **maximizing information gain (IG)**

**minimizing weighted impurity, I**

$$IG = \text{how effectively we} \downarrow \text{entropy}$$

$$IG = I_0 - \left( \frac{N_L}{N} I_L + \frac{N_R}{N} I_R \right)$$

node 0

L     R

$$I_0 = 0.69$$

A
(40, 40)
(30, 10)    (10, 30)

B
(40, 40)
(20, 40)    (20, 0)    $[(1, 0)]$

A) $I_L = -\left[\frac{3}{4} \ln \frac{3}{4} + \frac{1}{4} \ln \frac{1}{4}\right] = 0.56$
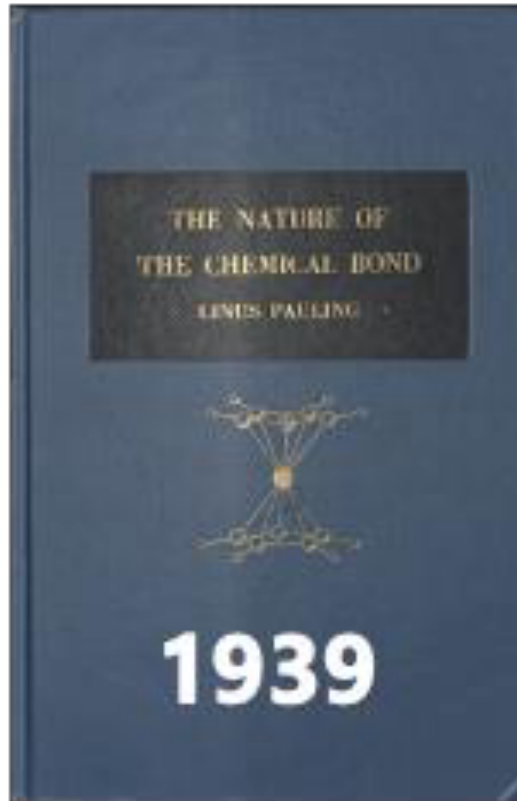
$I_R = I_L = 0.56$

$IG = 0.69 - \frac{1}{2} 0.56 - \frac{1}{2} 0.56 = 0.13$

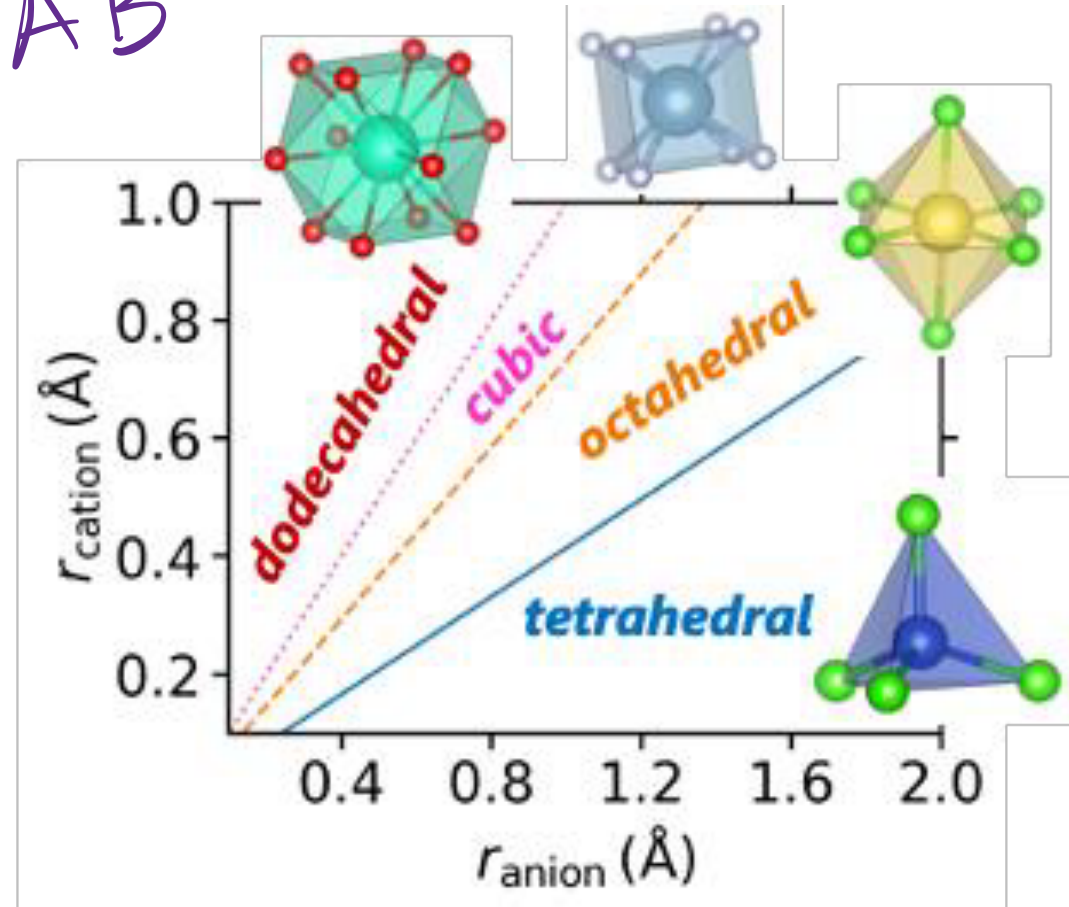B) $I_L = -\left[\frac{1}{3} \ln \frac{1}{3} + \frac{2}{3} \ln \frac{2}{3}\right] = 0.64$

$I_R = 0$

$IG = 0.69 - \frac{3}{4} 0.64 - \frac{1}{4} 0 = 0.21$ ✓

A B

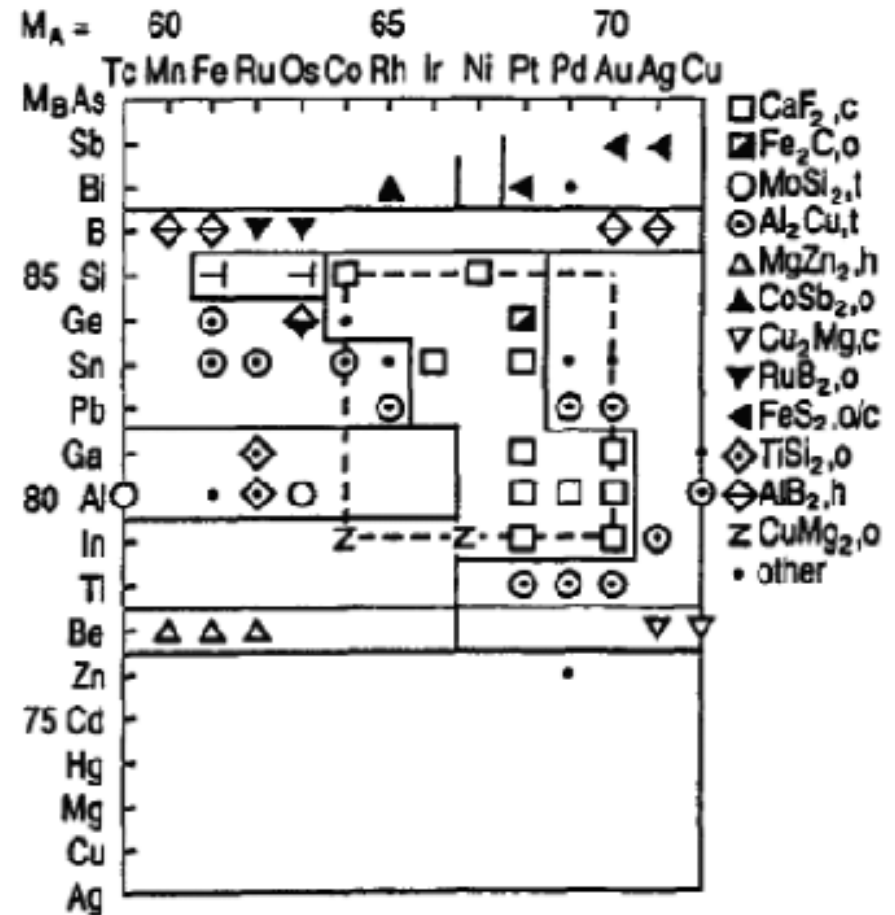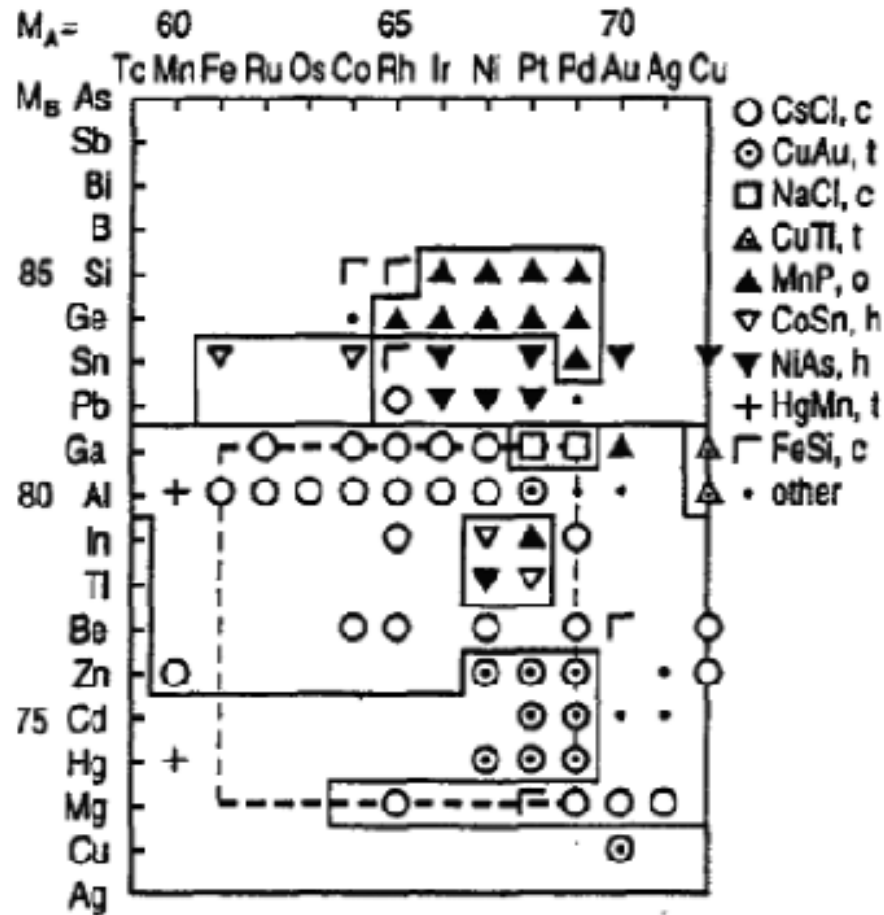**y –** *target property (observable)*          ***y*** *– data you find or generate*

CEMS
Chemical Engineering
& Materials Science

**y –** *target property (observable)*
**X –** *feature space (representation)*

**y** – *data you find or generate*
**X** – *stuff you hope relates to* **y**



$$
\mathbf{y} \qquad \mathbf{X}
$$

| | $r_{cation}$ | $r_{anion}$ |
|---|---|---|
| | 0.5 | 1.7 |
| | 0.7 | 1.1 |
| | 0.3 | 1.2 |
| | ... | ... |

**y** – *target property (observable)*
**X** – *feature space (representation)*
**f(X)** – *model (descriptor)*
**ŷ** – *prediction (model output)*

**y** – *data you find or generate*
**X** – *stuff you hope relates to* **y**
**f** – *the learned mapping of* **X** *to* **y**

$f(X)$: min(y-ŷ)

# Goldschmidt's tolerance factor for perovskite stability

## For 576 experimentally characterized $ABX_3$ compounds

$$t = \frac{r_A + r_X}{\sqrt{2}(r_B + r_X)}$$



**74% acc (1921)**

# New tolerance factor!

## For 576 experimentally characterized $ABX_3$ compounds

$$t = \frac{r_A + r_X}{\sqrt{2}(r_B + r_X)}$$

**SISSO**

$$\tau = \frac{r_X}{r_B} - n_A \left( n_A - \frac{r_A/r_B}{\ln r_A/r_B} \right)$$



Left plot: Counts vs $t$, with legend: perovskite, nonperovskite, predicted perovskite, predicted nonperovskite. **74% acc (1921)**

Right plot: Counts vs $\tau$, with legend: perovskite, nonperovskite, predicted perovskite, predicted nonperovskite. **92% acc (2019)**

# Decision trees w/ Goldschmidt's t

**576 ABX$_3$**
(313, 263)

t < 0.825          t > 0.825

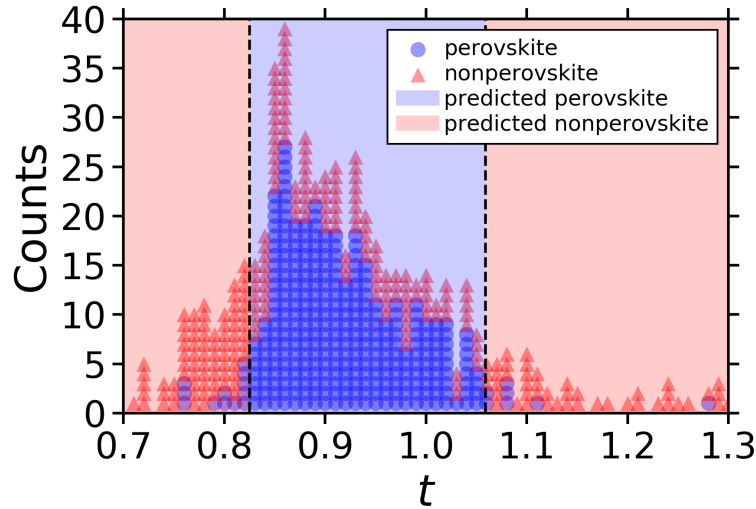**89 ABX$_3$**          **487 ABX$_3$**
(12, 77)          (301, 186)

t < 1.059          t > 1.059

**427 ABX$_3$**          **60 ABX$_3$**
(294, 133)          (7, 53)

$$y \begin{bmatrix} 1 \\ -1 \\ 1 \\ \ldots \end{bmatrix} \quad X \begin{bmatrix} t = \dfrac{r_A + r_X}{\sqrt{2}(r_B + r_X)} \end{bmatrix}$$

$$Acc = \frac{77 + 53 + 294}{576} = 0.74$$

$$FPR = \frac{133}{263} \approx 50\%$$

# Decision trees w/ $\tau$

**576 ABX$_3$**
(313, 263)

$\tau < 4.18$        $\tau > 4.18$

**321 ABX$_3$**          **255 ABX$_3$**
(293, 28)              (20, 235)
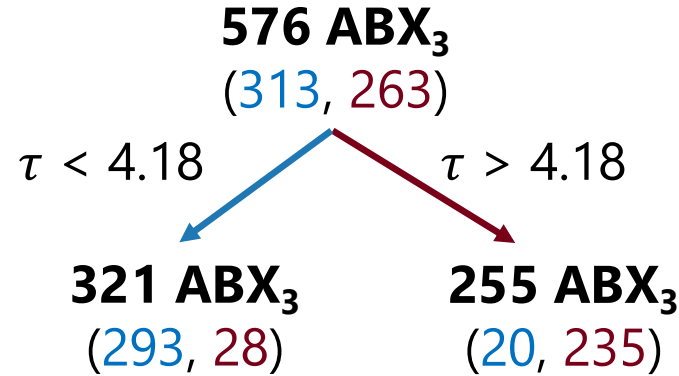
**y**                    **X**

$$\begin{bmatrix} 1 \\ -1 \\ 1 \\ \dots \end{bmatrix} \quad \begin{bmatrix} \tau = \dfrac{r_X}{r_B} - n_A \left( n_A - \dfrac{r_A/r_B}{\ln r_A/r_B} \right) \end{bmatrix}$$

$$Acc = \frac{293 + 235}{576} \approx 0.92$$

$$FPR = \frac{28}{263} \approx 10\%$$