# Day 11: Naive Bayes and Decision Trees

# Parametric Learning: learn model parameters
(e.g., linear and logistic regression)

# Non-Parametric Learning: models that don't learn parameters
(e.g., KNN and Naïve Bayes)

## Naïve Bayes

Uses Bayes' Rule:

$$P(y \mid x) = \frac{P(x \mid y) P(y)}{P(x)}$$

$\underbrace{\phantom{P(y)}}_{\text{label}}$ $\underbrace{\phantom{P(x)}}_{\text{data}}$

To classify:

$$P(y = \text{label} \mid x = \text{data}) \quad \text{for}$$

- Evaluate ... all possible labels
- Choose label with max probability

In practice, $P(x)$ doesn't make a difference, so we don't bother calculating;

$$\frac{P(x \mid y=dog) P(dog)}{P(x)} \quad vs. \quad \frac{P(x \mid y=cat) P(cat)}{P(x)}$$

Same denomenator, so can ignore

Note on HW3: this is what HW3 means when it says "only calc the numerator"

Using conditional independence assumption;

Using

$$P(X = x_1, x_2, x_3, x_4 \mid y)$$

$$= P(x_1 \mid y) P(x_2 \mid y) P(x_3 \mid y) P(x_4 \mid y)$$

$$\underbrace{P(X \mid y = dog)} P(y = dog)$$

$$= \underbrace{P(x_1 \mid dog) P(x_2 \mid dog) P(x_3 \mid dog) P(x_4 \mid dog)} P(dog)$$

To "train": Precompute probabilities

| $X_1$ | $X_2$ | $X_3$ | $X_4$ | $y$ |
|---|---|---|---|---|
| 0 | 0 | 1 | 1 | car |
| 0 | 1 | 0 | 1 | car |
| 1 | 0 | 1 | 1 | car |
| 1 | 1 | 0 | 0 | bike |
| 1 | 0 | 1 | 0 | bike |

$$P(y = car) = 3/5 \qquad P(y = bike) = 2/5$$

$$P(X_1 \mid car) = 1/3 \qquad P(X_1 \mid bike) = 1$$

$$P(X_2 \mid car) = 1/3 \qquad P(X_2 \mid bike) = 1/2$$

• • •

**Problem:** what is $P(X_i \mid y) = 0$?

$$P(X_1 \mid y) \cdot \ldots \cdot 0 \cdot \ldots \cdot P(X_n \mid y) \cdot P(y) = 0$$

**The Fix:** <u>Laplace Smoothing</u>

* add 1 to the numerator

* add (# classes) to the denomenator

Example: Predict between "spam", "maybe spam", and "definitely not spam":

$$P\left(x_i = \text{"ostrich"} \mid y = \text{"spam"}\right) = \frac{0 + 1}{(\#\text{ spam}) + 3}$$

$$P\left(x_i = \text{"ostrich"} \mid y = \text{"maybe spam"}\right) = \frac{0 + 1}{(\#\text{ maybe spam}) + 3}$$

Common variation:

use a smoothing parameter $\beta$ (hyper)

• adding $\beta$ to the numerator

• adding $\beta(\#\text{ classes})$ to the denominator

hyper

## A practical consideration:

Underflow: $0.0000000000...003 \rightarrow 0$

Standard fix: maximize $\log(P(\cdots))$

$$\log(ab) = \log(a) + \log(b)$$

So, maximizing $P(y_i = c \mid X)$ is the same as

maximizing $\log P(y_i = c \mid X)$

$$\log\left( \prod_{j=1}^{d} \left[ P(X_{ij} \mid y_i = c) \right] P(y_i = c) \right)$$

$$= \sum^{d} \log \left[ P(X_{ij} \mid y_i = c) \right] + \log P(y_i = c)$$

$$= \sum_{j=1}^{} \log \left[ \cdots \right]$$

# Decision Trees for Classification

Animal Decision Tree Classifier:

Has wings?
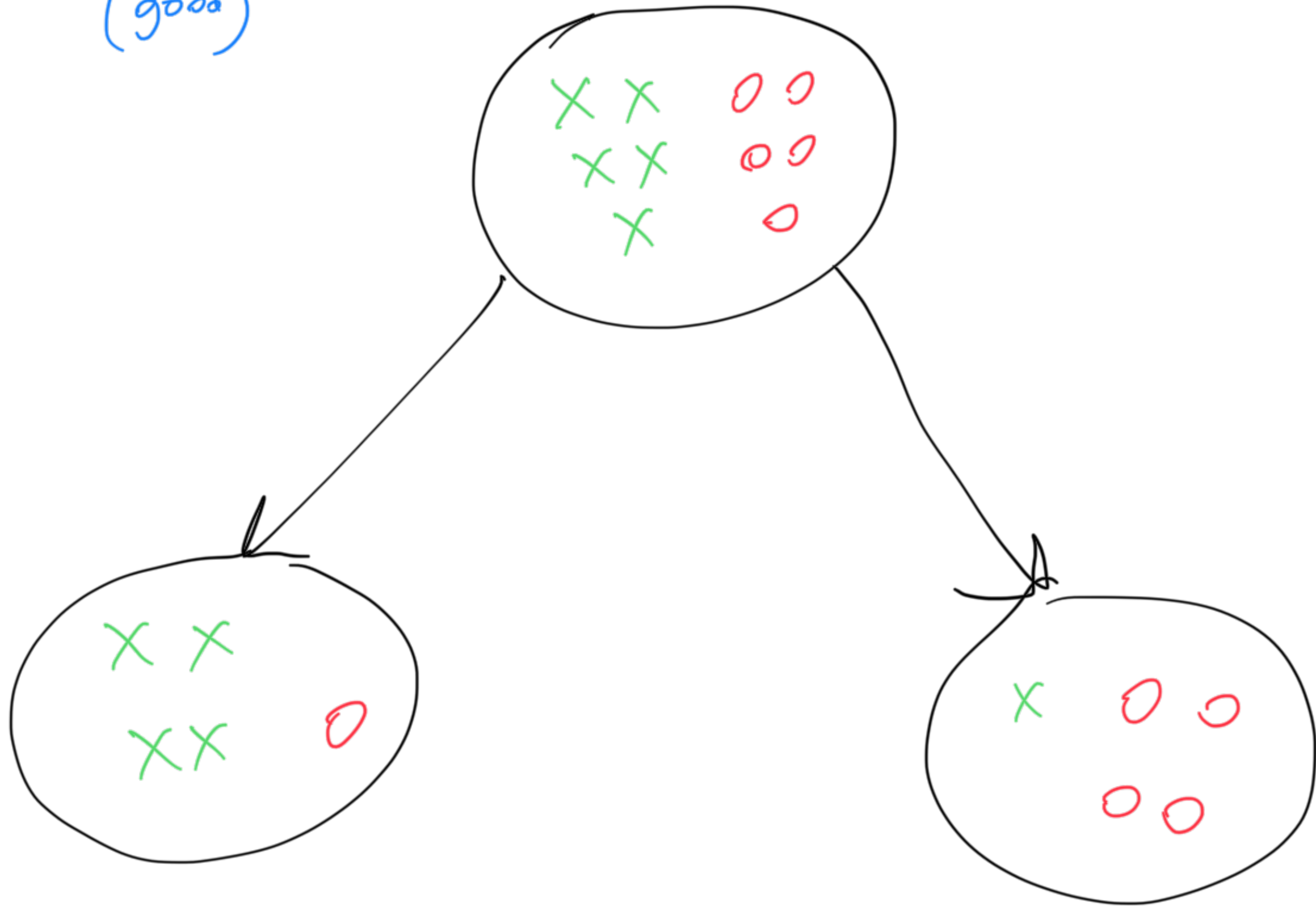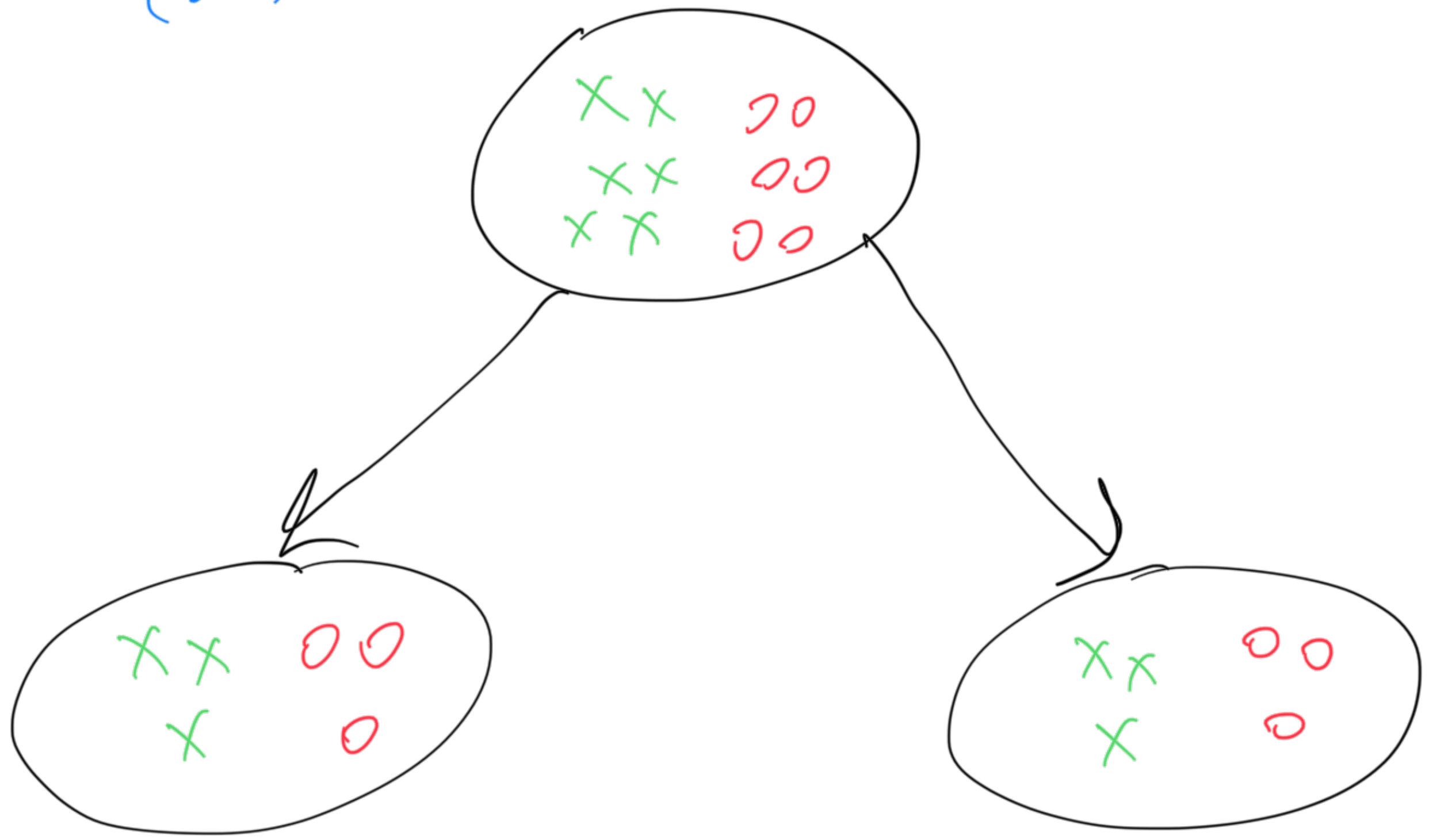
Yes → Bird

No → Length > 3ft?

No → Mouse

Yes → Snake

Prediction: go down the tree until you reach a leaf node

"Training": many variants, but general idea: recursively split tree on input features with the largest "information gain"

High Information Gain:
(good)



Low Information Gain:
(Bad)

# Basic Algorithm: ID3

"Iterative Dichotomiser 3"

ID3 (data D, features F):

(1) Calculate splitting metric for every feature $f \in F$ of dataset $D$

(2) split $D$ into subsets based on the splitting metric

(3) make a node for the selected feature $f$ which minimizes/maximizes the metric

(4) recurse on the subsets using the non-selected features

Many possible splitting metrics;

one popular one:

Gini Impurity

$$G = 1 - \sum_{i=1}^{K} p_i^2, \quad K = \text{\# of classes}$$

$p_i = $ probability of being in class $i$

(split on <u>lowest</u> impurity)

High impurity; same probability per class

$$1 - \underbrace{(0.5)^2}_{0.25} - \underbrace{(0.5)^2}_{0.25} = 0.5$$

Low impurity; high probability for one class, low probability for the other

$$1 - \underbrace{(0.99)^2}_{0.9801} - \underbrace{(0.01)^2}_{0.0001} = 0.0198$$
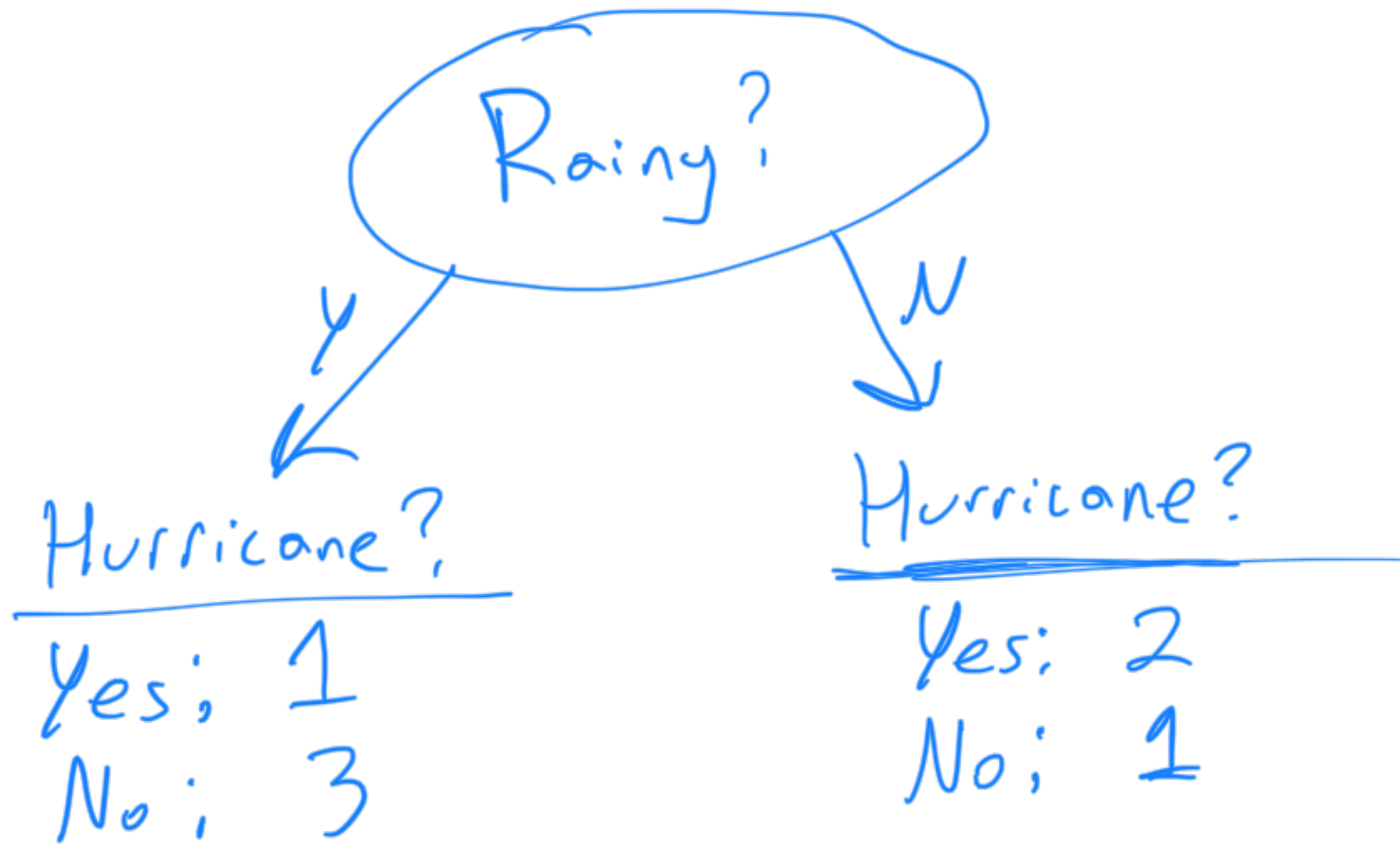
# Example (using Gini impurity as our splitting metric):

| Raing? | Humid? | Temp. | Hurricane? |
|--------|--------|-------|------------|
| Y | Y | 7 | No |
| Y | N | 12 | No |
| N | Y | 18 | Yes |
| N | Y | 35 | Yes |
| Y | Y | 38 | Yes |
| Y | N | 50 | No |
| N | N | 83 | No |

First, see which input has the ↓ impurity.

Rainy

Rainy?

Y → Hurricane?
N → Hurricane?

Hurricane?

Yes: 1
No: 3

Hurricane?

Yes: 2
No: 1

$$G\left(Hurricane = yes\right)$$
$$= 1 - \left(\frac{1}{1+3}\right)^2 - \left(\frac{3}{1+3}\right)^2 = 0.375$$

$$G\left(Hurricane = no\right)$$
$$= 1 - \left(\frac{2}{2+1}\right)^2 - \left(\frac{1}{2+1}\right)^2 = 0.444$$

$$\left( \qquad - \left(\frac{4}{}\right) 0.375 + \left(\frac{3}{4+3}\right) 0.444 \right)$$

$G_{total} = \left(\frac{}{4+3}\right)$ ... $\left( \quad \right)$

$= 0.405$

Humid

Humid?

Y         N

Hurricane?
_____
Y:   3
N:   1

Hurricane?
_____
Y:   0
N:   3

• • •

$G_{total} = 0.214$

Temp.

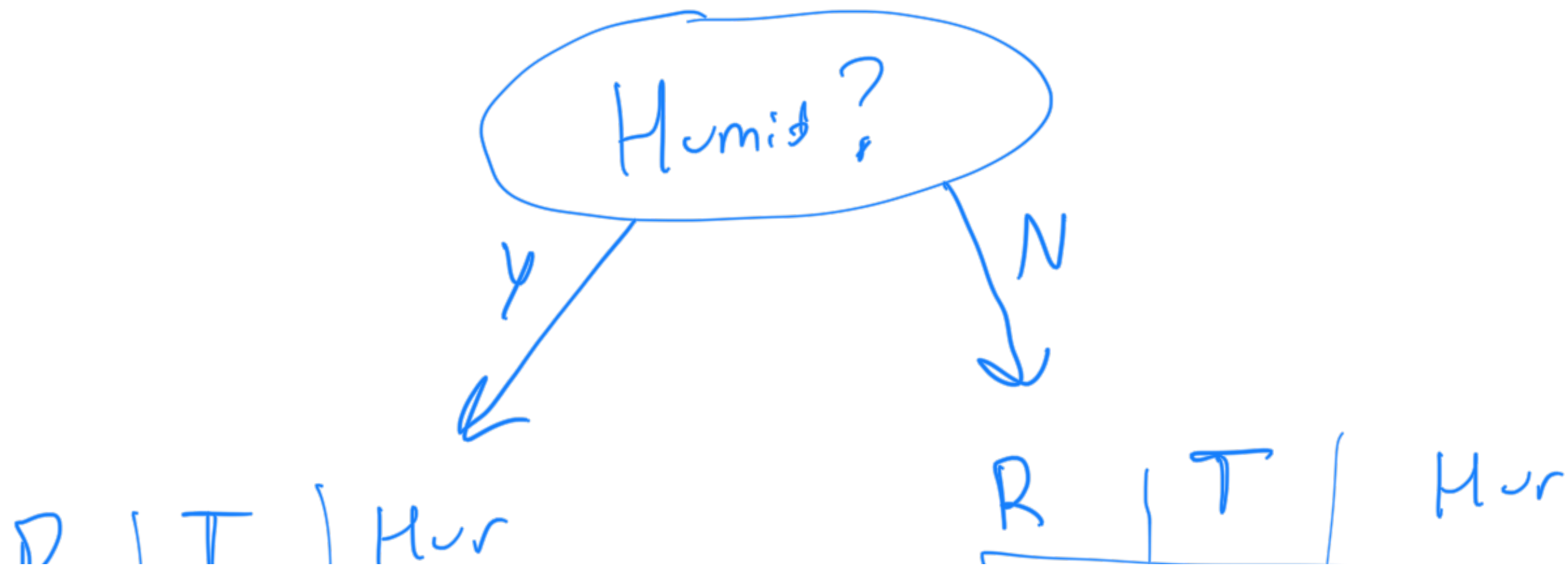| Temp | Hurricane: |
|---|---|
| | N   G = 0.429 |
| 7 | N   G = 0.343 |
| 9.5 | |
| 12 | Y   G = 0.476 |
| 15 | |
| 18 | Y   G = 0.476 |
| 26.5 | |
| 35 | Y |
| 36.5 | |
| 38 | G = 0.343 |
| 44 | |
| 50 | N   G = 0.429 |
| 66.5 | |
| 83 | N |

Temp. < 9.5 ?

Y → Hurricane?
Y: 0

N → Hurricane?
Y: 3
N: 3

$$G_{total} = \left(\frac{1}{1+6}\right) G_{temp < 9.5} + \left(\frac{6}{1+6}\right) G_{temp \geq 9.5}$$

$$= 0.429$$

Deciding on the root node:

Since "Humid?" has lowest $G$,

"Humid?" is our root node:



Humid?

Y            N

R | T | Hur

R | T | Hur

| R | | N |
|---|---|---|
| Y | 7 | N |
| N | 18 | Y |
| N | 35 | Y |
| Y | 38 | Y |

| Y | 12 | N |
|---|---|---|
| Y | 50 | N |
| N | 83 | N |

(Recursively)

Run through the same process but on this smaller sub-dataset

100% of the data points have "Hurricane = no", so stop splitting (make this a leaf node)