→ Classification Algorithm
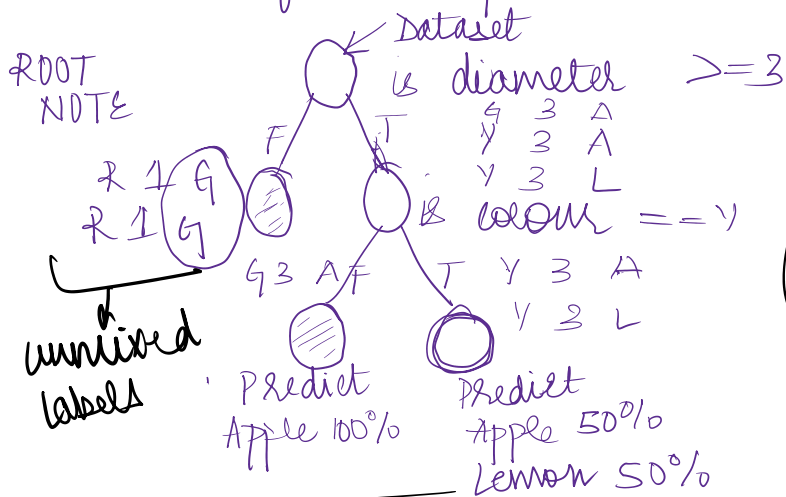
| Color | Diameter | (Fruit) Label |
|-------|----------|---------------|
| Green | 3 | Apple |
| Yellow | 3 | Apple * |
| Red | 1 | Grape |
| Red | 1 | Grape |
| Yellow | 3 | Lemon * |
| C | N | |

→ Goal : Predict this label

categories

No way to separate these 2 data pts. DT handles such case.

decision tree for above problem

Dataset

ROOT NODE

is diameter $\geq = 3$

| | | |
|---|---|---|
| F | G 3 A | T |
| | Y 3 A | |
| | Y 3 L | |

R 1 G
R 1 G
unmixed labels

G 3 A F

is colour $== y$

| | |
|---|---|
| T | Y 3 A |
| | Y 3 L |

Predict Apple 100%

Predict Apple 50% Lemon 50%

→ Gini Impurity.

① We need to quantify how much a past Q helps in unmixing labels.

INFORMATION GAIN

② Quantify how much Q helps Reduce a label uncertainty

As purest tree as poss
↳ largest no of pure Nodes

# Working of a DT

1. Which questions to ask & when?

Data → Tree

CART : Classi & Regr. Trees.

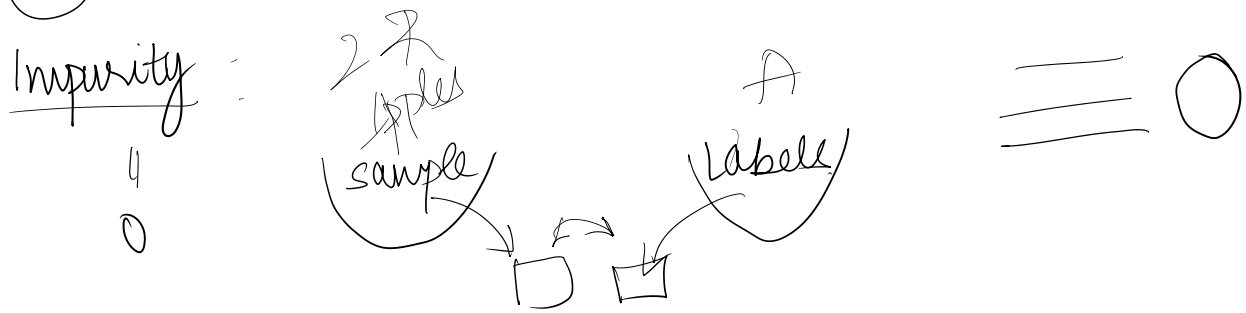gives a procedure to decide when, which Qs to ask.

FOR a node to be pure:

AT EVERY NODE  GI = 0
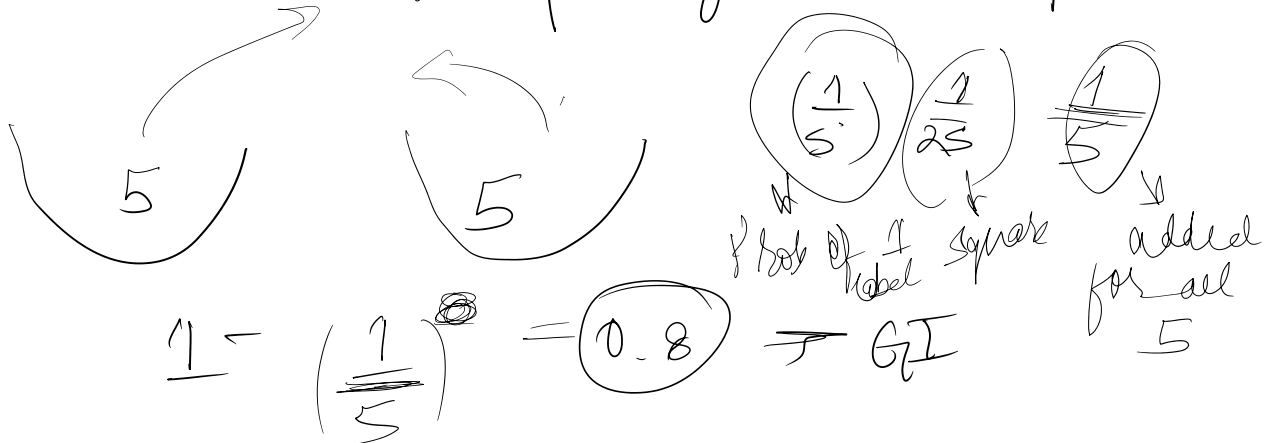Creating a list of all possible Qs.
& iterating over them

QUANTIFY UNCERTAINITY

The BEST Q: it is the one that reduces the uncertainty the most

(GI): how much uncertainity is there at every Node

(IG): tells us how much a Q reduces the uncertainity

Impurity : 2

apples

sample          A

          Labels

4

0

GI → incorrect label prediction at a node
          if at random a sample (eg) &
          a corresponding label is picked.

5          5

$\left(\dfrac{1}{5}\right)$  $\dfrac{1}{25}$  $\dfrac{1}{5}$

prob of
label  square    added
                    for all
                    5

$1 - \left(\dfrac{1}{5}\right)^2 = \boxed{0.8}$ = GI

$$GINI(t) = 1 - \sum_{K} \left[ P(c_K \mid t) \right]$$
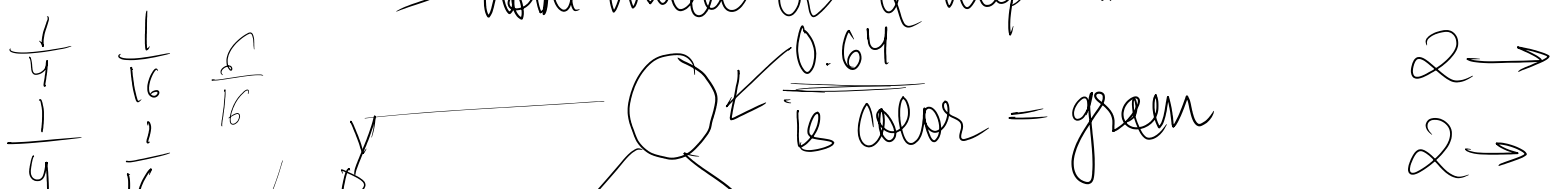
$P(c_K \mid t) = $ prob of a node t being categorized
                    as $c_K$.

t → no of data points

# INFORMATION GAIN

→ helps us find Q that reduces uncertainity
          the most
                    ↳ how much a Q helps unmix the labels
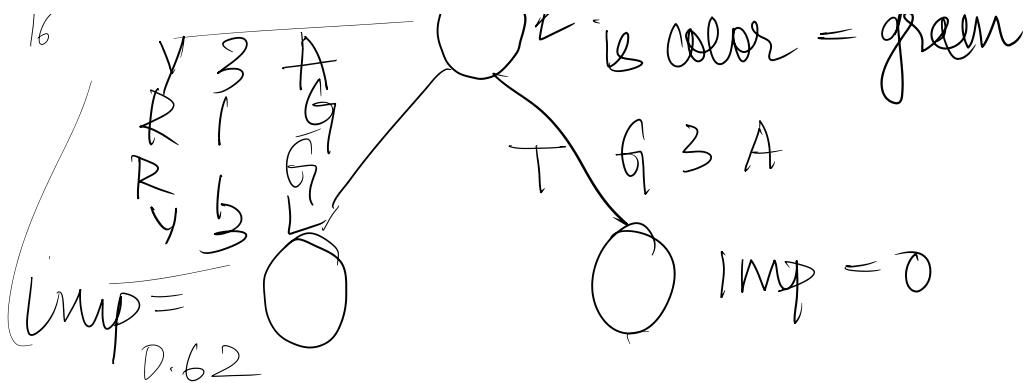
$\dfrac{1}{4}$  $\dfrac{1}{16}$  $\dfrac{6}{16}$

$\dfrac{1}{4}$  $\dfrac{1}{16}$

$\dfrac{0.64}{16}$ color = green          2→

2→

A  $2/5$    $4/25$

G  $2/5$    $4/25$

$$\frac{1}{4} \quad \frac{i}{16} \quad 16$$

$$\frac{2}{4} \quad \frac{4}{16}$$

$$\frac{4}{4} \quad \frac{4}{16}$$

$$1 - \frac{6}{16}$$

$$\frac{10}{16}$$

Y 3 A
R 1 G
R y G
R y L

Imp = 0.62

Is color = green

T   G 3 A

Imp = 0

$2 \rightarrow$

$1 \rightarrow$

$5$

$\frac{06}{25}$

$$\text{Avg Imp} = \frac{4}{5} \times 0.62 + \frac{1}{5} \times 0$$

$$= 0.5$$

$$\text{INFORMATION gain} = \text{Starting GI} - A$$

$$= 0.64 - 0.5 = \boxed{0.}$$

## Dry Run

```python
def build_tree(rows):
    """Builds the tree.

    Rules of recursion: 1) Believe that it works. 2) Start by checking
    for the base case (no further information gain). 3) Prepare for
    giant stack traces.
    """

    # Try partitioing the dataset on each of the unique attribute,
    # calculate the information gain,
    # and return the question that produces the highest gain.
    gain, question = find_best_split(rows)

    # Base case: no further info gain
    # Since we can ask no further questions,
    # we'll return a leaf.
    if gain == 0:
        return Leaf(rows)

    # If we reach here, we have found a useful feature / value
    # to partition on.
    true_rows, false_rows = partition(rows, question)
```

$G \cdot 2/5 \quad 4/25$

$L \cdot 1/5 \quad 4/25$

$1 - \dfrac{9}{25}$

very 41 of
next step

$14 \rightarrow$ max info gain

$\dfrac{2}{3} \quad \dfrac{1}{3}$

$\dfrac{4}{9} \quad 1 \quad \dfrac{5}{9}$

$\dfrac{5}{9}$

```
    # to partition on.
    true_rows, false_rows = partition(rows, question)

    # Recursively build the true branch.
    true_branch = build_tree(true_rows)

    # Recursively build the false branch.
    false_branch = build_tree(false_rows)

    # Return a Question node.
    # This records the best feature / value to ask at this point,
    # as well as the branches to follow
    # dependingo on the answer.
    return Decision_Node(question, true_branch, false_branch)
```

```
training_data = [
    ['Green', 3, 'Apple'],
    ['Yellow', 3, 'Apple'],
    ['Red', 1, 'Grape'],
    ['Red', 1, 'Grape'],
    ['Yellow', 3, 'Lemon'],
]
```

> 0.64

G 3 A
3A=3 Y 3 A
F Y 3 L

F RIG T=01

S
G 3
T F Y 3 A
Y 3 C
0.5