# ✳ Gini Impurity
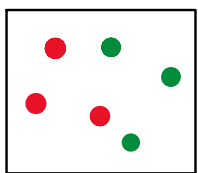
→ Formula: $G(y) = 1 - \sum [P(y_i)]^2$

Background: If our dataset has 'd' features & 'n' datapoints then:

① Calculate Entropy for each feature
(eg, Gender / Education)

② Find out Information Gain for each feature

③ Chose the feature with maximum IG to split our dataset

→ How to compute Gini Impurity in our example?

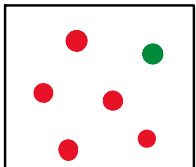$$G(y) = 1 - \sum [P(y_i)]^2 \rightarrow G(y) = 1 - \left[ P(y-g)^2 + P(y-r)^2 \right]$$

$P(-ve) = \frac{1}{2}$        $G(y) = 1 - \left[ (0.5)^2 + (0.5)^2 \right] = 1 - 0.5$
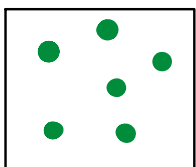
$P(+ve) = \frac{1}{2}$        $G(y) = 0.5$

$P(-ve) = \frac{5}{6}$        $G(y) = 1 - \left[ \left(\frac{5}{6}\right)^2 + \left(\frac{1}{6}\right)^2 \right] = 1 - \left[ \frac{26}{36} \right] = \frac{36-26}{36}$
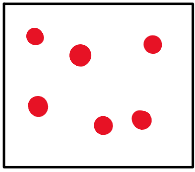
$P(+ve) = \frac{1}{6}$        $G(y) = \frac{10}{36} = 0.28$

$P(-ve) = 0$        $G(y) = 1 - [0^2 + 1^2] = 0$
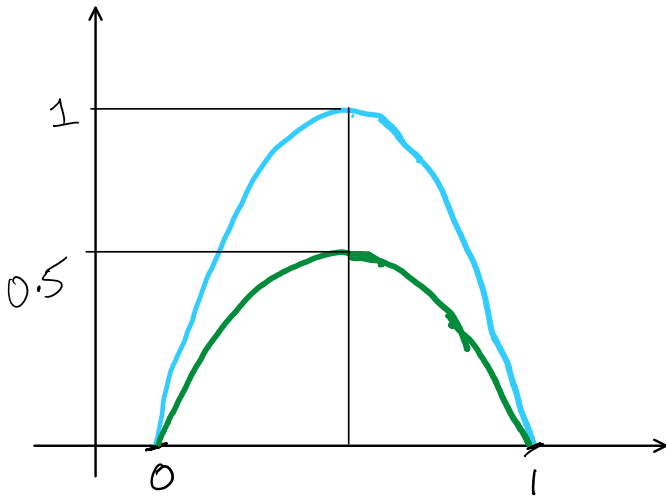
$P(+ve) = 1$

$P(-ve) = 1$
$P(+ve) = 0$

$G(y) = 1 - [1^2 + 0^2] = 0$
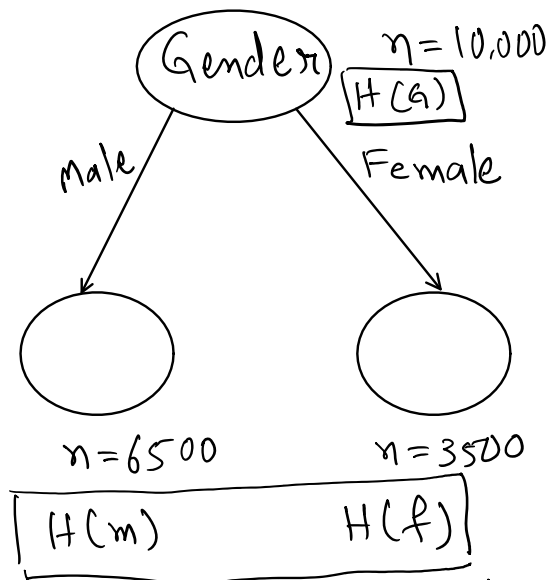


$P(y-g) \rightarrow$

✿ Can we use this same approach to numerical columns? Ans - No. Why?
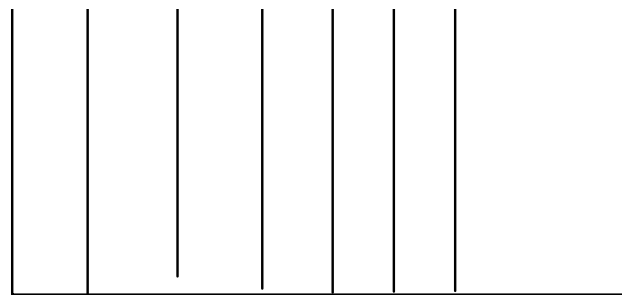
In case of categorical columns:

| Gender | Education |
|--------|-----------|
| male | Non Grad |
| female | Non Grad |
| female | Grad |
| female | Non Grad |
| male | Grad |



For a numeric column:

Price?

price?

1.5    (2.3)    2.7    105    75.7

H(1.5)   H(2.3)   H(2.7)   H(105)   H(75.7)

How many Entropy calculations are needed?

$F_1$   $F_2$   $F_3$   $F_4$   $F_5$   $F_6$

2  +  2  +  2  +  1000  +  1000  +  1000  = 3006

**Steps** — ① Sort the data in ascending order of numerical column

② For each unique value, calculate Entropy

③ Compute IG for all the thresholds

④ Find the question with maximum IG.

**Disadvantage** —

A lot expensive computationally.

Then how can we find Entropy of numerical columns?

Ans: Creating bins on the numerical column and calculating Entropy of each bin rather than calculating for each unique value.

# ☆ An entire view:



In context of Underfitting & Overfitting:

**Underfitting:**
→ The model doesn't learn enough.

**Overfitting:**
→ Our model tries to learn each datapoint



← Each question creates an axis-parallel boundary.

Very Shallow Tree leads to Underfitting

Extremely Deep trees lead to overfitting.

→ ∴ Won't it be a good idea to control 'd'?

Hence, d is one of the Hyperparameter while implementing Decision Trees

implementing Decision Trees.

$$d \uparrow \uparrow \rightarrow \text{overfitting} \qquad d \downarrow \downarrow \rightarrow \text{Underfitting.}$$

| Depth ('d') | Training Accuracy | Validation Accuracy | Comments |
|---|---|---|---|
| 1 | ↓↓ | ↓↓ | } underfit |
| 2 | ↓↓ | ↓↓ | |
| 3 | ↓ | ↓↓ | |
| 5 | ↑ | ↑ | } Best accuracy at validation (Good choices of d) |
| 7 | ↑ | ↑↑ | |
| 10 | ↑ | ↑↑ | |
| 50 | ↑↑ | ↓ | } Overfit |
| 100 | ↑↑ | ↓ | |

## ✡ Do outliers impact Decision Tree?



$d = 1$

→ Outlier

→ for shallow trees, impact of outliers is negligible

→ For extreemly deep trees, impact of outliers will be significantly high.

→ For 'best' values of d impact of outliers will be there but

|

30          45
→ Outlier

→ For 'best' values of d impact of outliers will be there but not very significant.

## ✳ Do we need feature scaling in decision trees? (Normalisation | Standardisation)

Ans – No. Why?

① We divide the data by asking a question that splits the data into two parts instead of finding distances of each point and hence, we don't need to shrink the numbers.

② While calculating Entropy | Gini Impurity, we compute probabilities

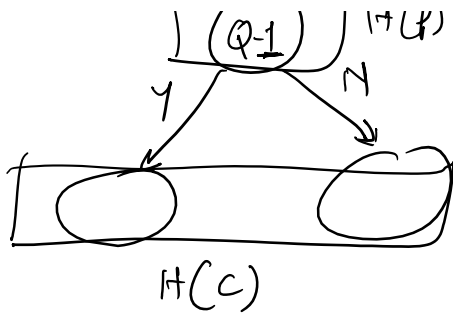$$H(Y) = -\sum_{i=1}^{n} P(Y) \cdot log(P(Y)) \qquad G(Y) = 1 - \sum P(Y)^2$$

And probability considers 'frequency' of the point not the value of that point.

But, should we Normalise | Standardise? → Yes!

## ✳ Should we use Decision Trees for high dimension data? eg, d = 1,00,000 ⇒ Ans: NO!
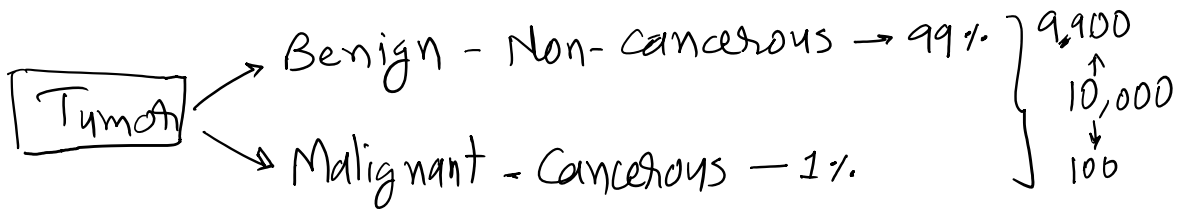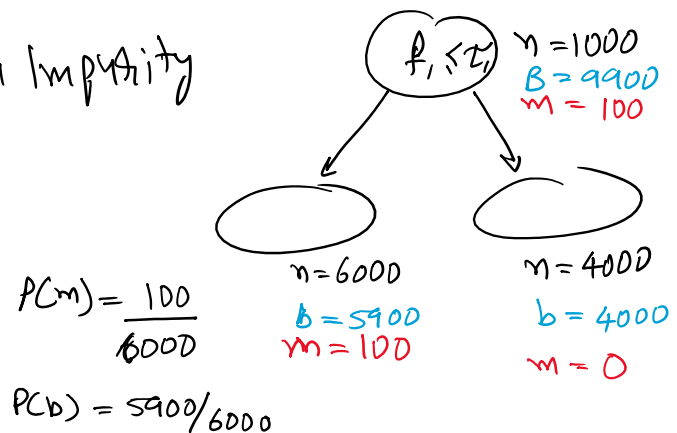
→ Because it will be very slow

(Q-1)   |+(P)
Y /    \ N
TC  [a]

$IG = \boxed{9}$

slow
$\rightarrow$ Lot's of computation are needed.

⚝ **Will imballanced data affect decision trees?**
**(Is it needed to do data-rebalancing while using decision trees?) ⟹ Yes!**

Tumor → Benign - Non-Cancerous → 99% ⎤ 9,900
                                        ⎬ 10,000 ↑
        → Malignant - Cancerous — 1%   ⎦ 100

Question → Entropy | Gini Impurity
           ↳ P(y) ↵



$f_1 < \tau$  n = 1000
              B = 9900
              M = 100

n = 6000          n = 4000
b = 5900          b = 4000
m = 100           m = 0
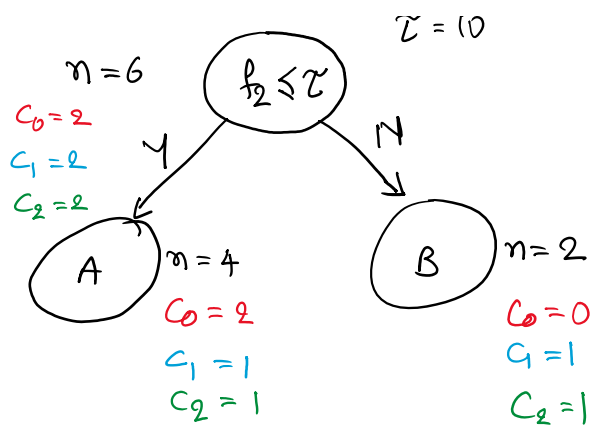
$P(m) = \dfrac{100}{6000}$

$P(b) = 5900/6000$

Data Rebalancing — Under sampling, Over Sampling, Class Weights, SMOTE

⚝ **Can we use Decision Trees in multiclass-classification?**

n = 6    $f_2 < \tau$     $\tau = 10$

| $f_1$ | $f_2$ | $f_3$ | $f_4$ | Y |
|---|---|---|---|---|
| – | 2 | – | – | 0 |
| – | 5 | – | – | 1 |

Tree (top left):

$n=6$
$C_0 = 2$
$C_1 = 2$
$C_2 = 2$

$f_2 \leq \tau$   $\tau = 10$

Y (left) → A: $n = 4$, $C_0 = 2$, $C_1 = 1$, $C_2 = 1$

N (right) → B: $n = 2$, $C_0 = 0$, $C_1 = 1$, $C_2 = 1$

Table (top right):

| | | | | |
|---|---|---|---|---|
| – 2 | – | – | . | |
| – 5 | – | – | 1 | |
| – 11 | – | – | 2 | |
| – 15 | – | – | 1 | |
| – 10 | – | – | 2 | |
| – 8 | – | – | 0 | |

$$G(A) = 1 - \left[ P(C_0)^2 + P(C_1)^2 + P(C_2)^2 \right]$$

Yes! We can use Decision Trees for multiclass classification.

☆ How will we calculate feature importance through Decision Trees?

Importance:

$$f_1 > f_2 > f_3 > f_9$$



$f_1 \leq \tau_1$   $n = 10,000$

left → $f_3 \leq \tau_3$   $n_1 = 6000$

right → $f_2 \leq \tau_2$   $n_2 = 4000$

$f_3$ children: $f_1 \leq \tau_5$   $n_3 = 3000$ ; $f_2 \leq \tau_6$   $n_4 = 3000$

$f_2$ children: $f_3 \leq \tau_9$   $n_5 = 1000$ ; $f_1 \leq \tau_2$   $n_6 = 3000$

We compute **Normalised Information Gain** of each feature and then the feature with highest NIG is the most important feature.

$$\text{NIG of } f_2 = \frac{n_2}{n} \cdot \text{IG of } f_2 + \frac{n_4}{n} \cdot \text{IG of } f_2$$

(at d=3)

$$\text{NIG of } f_2 = \frac{n_2}{n} \cdot \text{IG of } f_2 + \frac{n_{14}}{n} \cdot \text{IG of } f_2$$
$$(\text{at } d = 2) \qquad\qquad (\text{at } d = 3)$$