# Shannon Entropy, Information Gain and Tree Construction

(Shannon) Entropy originates from the field of information theory. The basic intuition behind information theory is, that learning that an unlikely event has occurred is more informative than learning that a likely event has occurred.

Let us denote the amount of information <u>of an event</u> with probability $p$ as <u>$i(p)$</u>, and it should fulfill:

(1) $i(p) \gg 0$ decreasing (more info for unlikely events)

(2) $i(pq) = i(p) + i(q)$ (independent events)

It can be shown that [Shannon]:

$$i(p) = \log\left[\frac{1}{p}\right]$$

Note: The base of the logarithm can be chosen freely! If the base 2 is used, then the units of information are **bits**:

$$i(p) = \log_2\left[\frac{1}{p}\right] \quad \text{in bits}$$

binary digit

$[\log_{10} \rightarrow \text{dit}, \log_e \rightarrow \text{nat}]$

Now consider a random variable $X$ with probability distribution $p(x)$. The amount of information of an elementary event $X = x$ is $i(p(x)) = \log\left[\frac{1}{p(x)}\right]$

The <u>average</u> amount of information about $X$ is given by the expectation value:

$$H(X) = \sum_X p(x) \log\left[\frac{1}{p(x)}\right]$$

not a function of $X$ but $p(x)$: $H(p(x))$

Shannon's entropy of the random variable $X$ with distribution $p(x)$.

-2-

We arrive at: $(\log \frac{1}{x} = -\log x)$

$$H(X) = - \sum_x p(x) \log [p(x)]$$

logarithm

$$\log_2[x]$$
$$\ln[x]$$

$$H(X) \sim - \sum \overset{\geqslant 0}{p(x)} \underset{\leqslant 0}{\log[\quad]} \geqslant 0$$

$$\Rightarrow \quad H(X) \geqslant 0$$

Side remark: It can be shown that Boltzmann's entropy $S$ of statistical thermodynamics (= the disorder of a system increases) is formally identical with $H$ (Shannon entropy).

Example: $X$ takes $M$ equally probable values:

$$x = 0, 1, \ldots, M-1 \Rightarrow p(x) = \frac{1}{M}$$

$$H(X) = - \sum_{i=1}^{M} \frac{1}{M} \log\left[\frac{1}{M}\right] = \log[M]$$

$\underbrace{\qquad}$ log of number of possible values

-3-

How does $H(X)$ arise from a given distribution $p(x)$:

independent and identically distr.

i.i.d

$$X = [x_1, x_2, \ldots x_n]$$

$$p(x) = p(x_1)\, p(x_2) \cdots p(x_n)$$

$$n(x) = \begin{cases} 1 & x = x_i \\ 0 \end{cases}$$

$$\boxed{p(x) = \prod p(x_i) = \prod p(x)^{n(x)}}$$

↑ number of $x_i$

$$= \prod \left[ p(x)^{\frac{n(x)}{n}} \right]^n \approx \prod \left[ p(x)^{p(x)} \right]^n$$

**Asymptotic equipartition property**

$$\boxed{= e^{-nH(x)}}$$

for very large $n$ the value of the probability of a sequence $X = (x_1 \ldots x_n)$ is $\approx e^{-nH(x)}$ +

Proof:
$$e^{-nH(x)} = e^{-n\sum_x p(x)\log\left[\frac{1}{p(x)}\right]} = e^{+n\sum p(x)\log[p(x)]}$$

$$= e^{+\sum_x \log\left[\left[\frac{1}{p(x)}\right]^{p(x)n}\right]} =$$

$$= e^{+\sum \log\left[[p(x)]^{p(x)n}\right]}$$

$$= e^{+\left[\log[f(x_1)] + \log[f(x)] + \ldots\right]}$$

$$= e^{+\log\left[\prod[p(x)]^{p(x)\cdot n}\right]}$$

$$= + \prod[p(x)]^{p(x)\cdot n}$$

$n$ large
↓
$$\boxed{p(x) \approx e^{-nH(x)}}$$
for very large $n$
$p$ close to a constant $e^{-nH}$ !

−4−

# Continuous Entropy:

$$h(X) = E\left(\log\left[\frac{1}{p(X)}\right]\right)$$

$$= \int p(x) \log\left[\frac{1}{p(x)}\right] dx$$

Example: Uniform distribution in interval $(a,b)$

$$p(x) = \frac{1}{b-a} :$$

$$h(U) = \int_a^b \frac{1}{b-a} \log[b-a] \, dx$$

$$= \log[b-a] \frac{1}{b-a} \left[\int_a^b dx\right]$$

$$= \log[b-a] \qquad \left. 1 \right|_a^b = b-a$$

If $b-a < 1 \Rightarrow h(U) < 0$

continuous entropy can become negative!
(one cannot assign an "amount of information" to continuous entropy)

Exercise: $h(x)$ of the normal distribution?

Answer: $h(x) = \log[\sqrt{2\pi e}] + \log a$

## Conditional Entropy: (measure of information provided by Y about X)

Mutual Information / Information Gain

$$I(X;Y) = E\left(\log\left[\frac{p(x|y)}{p(x)}\right]\right)$$

$\underbrace{\phantom{I(X;Y)}}$
Random Variables

with $p(x|y)$ conditional distrib. of $X$ knowing $Y=y$.

$$I(X;Y) = E\left(\log[p(x|y)] - \log[p(x)]\right)$$

$$= E\left(\log\left[\frac{1}{p(x)}\right] - \left(\log\left[\frac{1}{p(x|y)}\right]\right)\right)$$

$$\underbrace{\phantom{= E\left(\log\left[\frac{1}{p(x)}\right]\right)}}_{= H(X)} \quad \underbrace{\phantom{\log\left[\frac{1}{p(x|y)}\right]}}_{H(X|Y)}$$

$$H(X|Y) = \sum\sum_{X} p(x,y) \log\left[\frac{1}{p(x|y)}\right]$$

Let us rewrite $H(X|Y)$:

$$H(X|Y) = \sum_{Y} p(y) \sum_{X} p(x|y) \log\frac{1}{p(x|y)}$$

-6-

$$H(X|Y) = \sum_{Y} p(y)\, H(X|Y=y)$$

$H(X|Y)$ is the average uncertainty about $X$ when $Y$ is known

$$\boxed{H(X|Y) \leq H(X)}$$ knowledge reduces uncertainty!

---

## 2. Entropy Calculation and Decision Trees

|     | height | legs |       |
| id  | h [m]  | $\ell$ | class |
| --- | --- | --- | --- |
| 1   | 0.1 | 0 | Fish  |
| 2   | 0.2 | 2 | Bird  |
| 3   | 1.8 | 2 | Human |
| 4   | 0.2 | 4 | Cat   |
| 5   | 2.1 | 4 | Horse |
| 6   | 1.7 | 2 | Human |
| 7   | 0.1 | 4 | Cat   |
| 8   | 1.6 | 2 | Human |

Entropy of the class information:

$$H(X) = -\sum_X p(x) \log_2 [p(x)]$$

$$= -\sum p(class) \log_2 [p(class)]$$

Class = {F, B, H, C, HO}
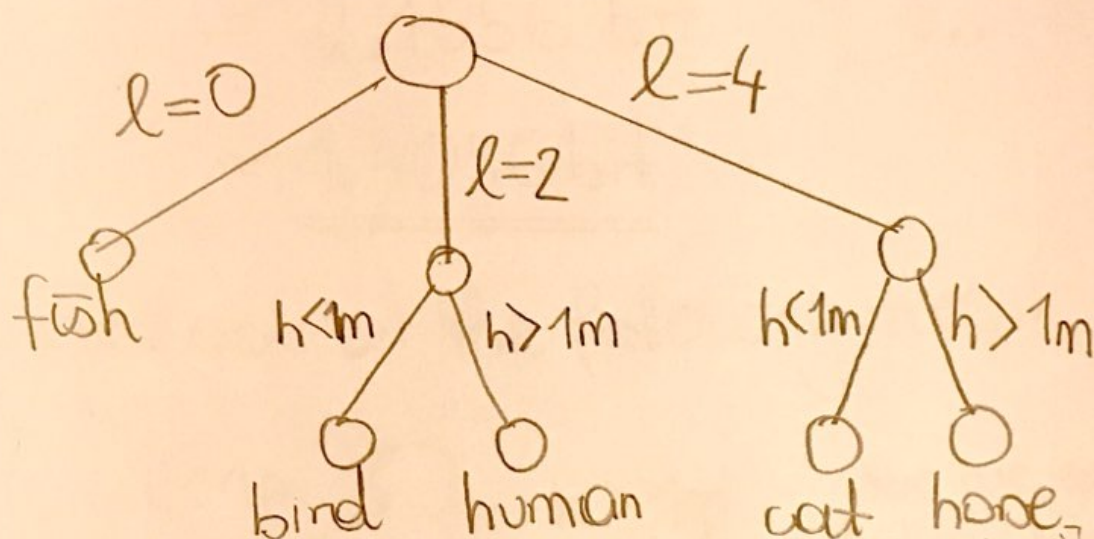
$$= -\left\{ \underbrace{\frac{1}{8} \log_2 \left[\frac{1}{8}\right]}_{F=Fish} + \underbrace{\frac{1}{8} \log_2 \left[\frac{1}{8}\right]}_{B=Bird} \right.$$

$$+ \underbrace{\frac{3}{8} \log_2 \left[\frac{3}{8}\right]}_{H=Human} + \underbrace{\frac{2}{8} \log_2 \left[\frac{2}{8}\right]}_{C=Cat}$$

$$\left. + \underbrace{\frac{1}{8} \log_2 \left[\frac{1}{8}\right]}_{HO=Hose} \right\} \approx 2.1556 \text{ bit}$$

Now let us construct a tree with
root node $\ell$:

$\ell=0$     $\ell=4$

$\ell=2$

fish    h<1m    h>1m    h<1m    h>1m

bird   human    cat   horse

$$-\sum_i \left(Class \mid \ell=0\right) \log\left[class \mid \ell=0\right]$$

$$H(Class \mid \ell=0) = -1 \cdot \log_2[1] = \underline{0}$$

$$H(Class \mid \ell=2) = -\sum_{Class=\{B,H\}} p(Class \mid \ell=2) \log_2(Class \mid \ell=2)$$

$$= -\frac{1}{4}\log_2\frac{1}{4} - \frac{3}{4}\log_2\frac{3}{4} \approx \underline{0.8113\,bit}$$

$$H(Class \mid \ell=4) = -\sum_{Class=\{C,HO\}} p(Class \mid \ell=4) \log(Class \mid \ell=4)$$

$$= -\frac{2}{3}\log_2\frac{2}{3} - \frac{1}{3}\log_2\frac{1}{3} = \underline{0.9183\,bit}$$

$$H(Class \mid \ell) = \frac{1}{8} \cdot H(Class \mid \ell=0) + \frac{4}{8} H(Class \mid \ell=2)$$
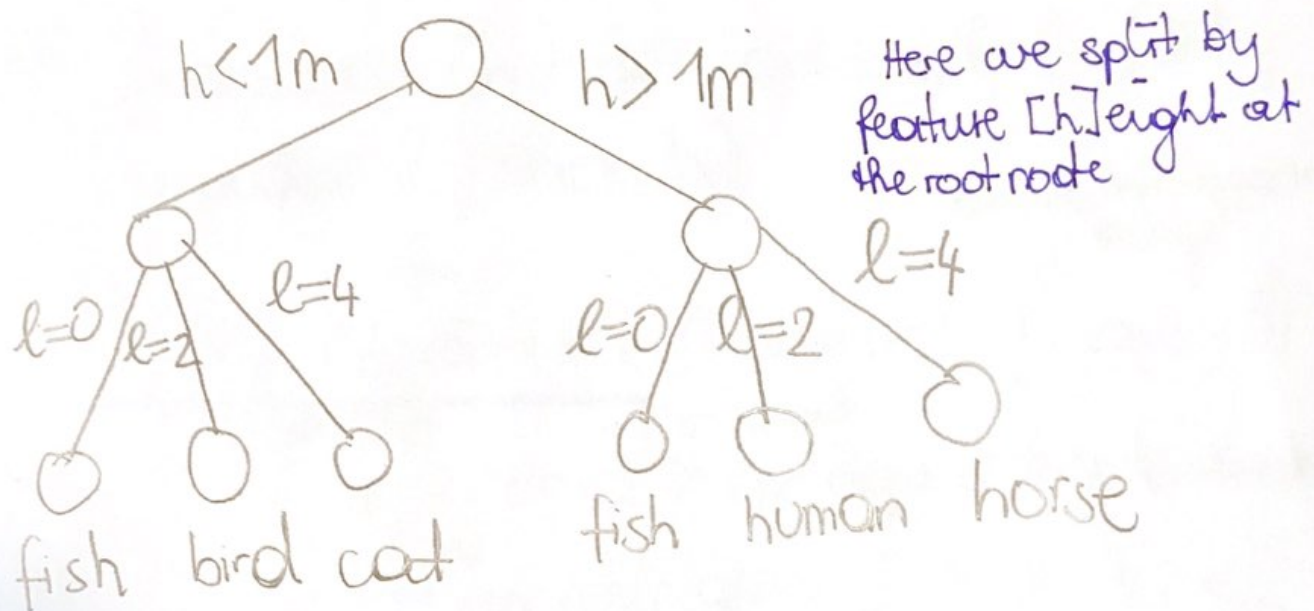$$+ \frac{3}{8} H(Class \mid \ell=4)$$

-9-

Expected Information Gain:

$$I(Class, l) = \underbrace{H(Class)} - H(Class/l)$$

$$= 2.1556 \text{ bit} - \dots =$$

$$\approx 1.4056 \text{ bit}$$

Now look at the following tree:

$h < 1m$ ⃝ $h > 1m$

Here we split by feature [h]eight at the root node

$l=4$

$l=0$ $l=2$ $l=4$

$l=0$ $l=2$

⃝ ⃝ ⃝

⃝ ⃝ ⃝

fish bird cat

fish human horse

$$H(Class | h < 1m) = - \sum p(Class|h<m) \log(Class|h<m)$$

$$= -\frac{1}{4} \log_2 \frac{1}{4} - \frac{1}{4} \log_2 \frac{1}{4} - \frac{2}{4} \log_2 \frac{2}{4}$$

$$\approx 1.5 \text{ bit}$$

$$H(Class | h > 1m) = -\frac{1}{4} \log_2 \frac{1}{4} - \frac{3}{4} \log_2 \frac{3}{4}$$

$$\approx 0.8113 \text{ bit}$$

$$I(class; h) = H(class) - H(class|h)$$

$$= 2.1556 \, bit - \left[\frac{4}{8} H(class|h<m) + \frac{4}{8} H(class|h\geq m)\right]$$

$$\approx 1 \, bit$$

$$\Rightarrow \quad I(class; l) > I(class; h)$$

$\Rightarrow$ first decision tree is optimal (= first consider $l$, then $h$)

$\Downarrow$

divide into opposite
$\downarrow$ things

# The ID3 algorithm (Iterative Dichotomiser 3)
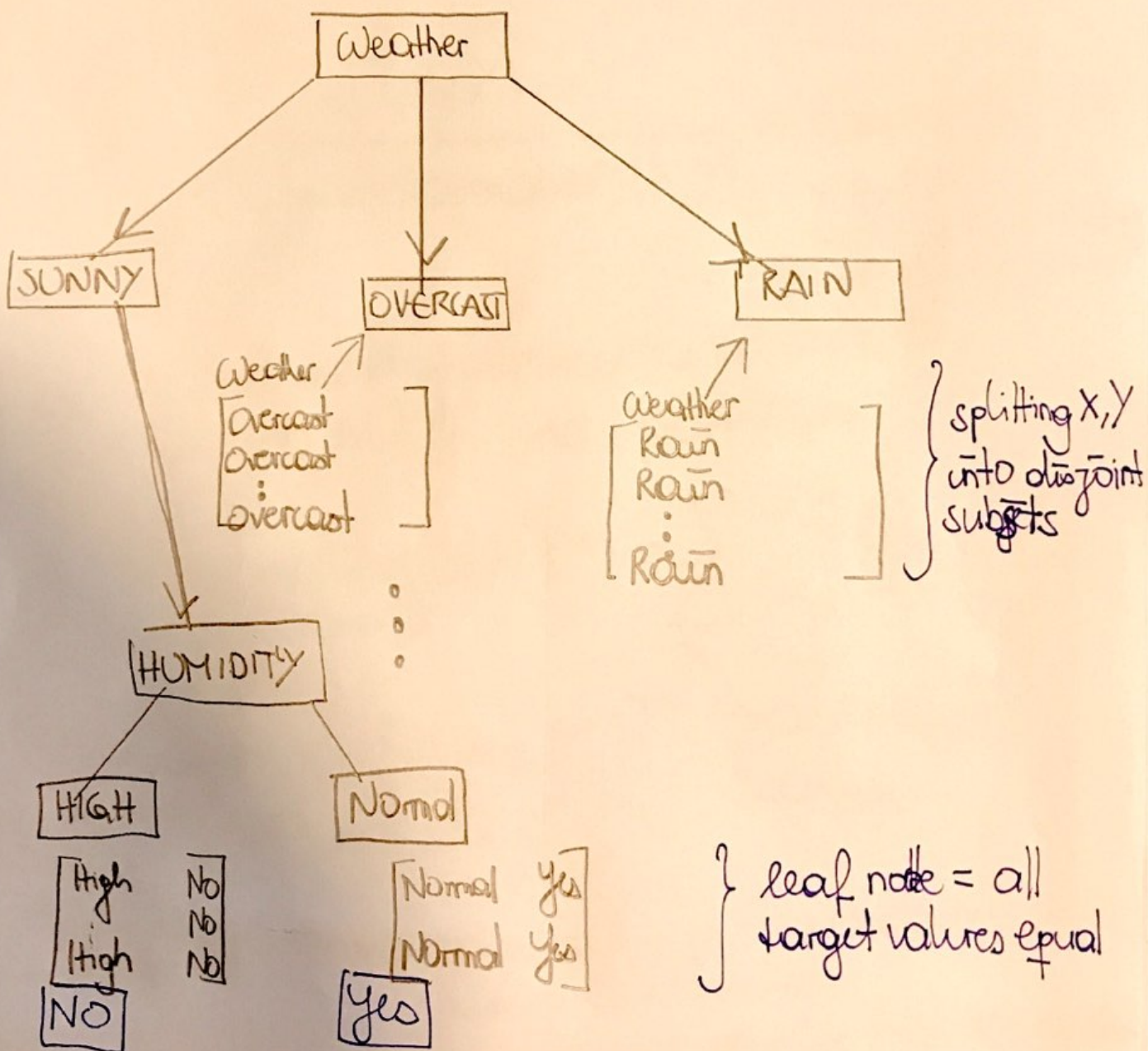
$\downarrow$ feature set    $\downarrow$ target

1. If $I$ is empty or all $Y$ are equal then terminate

2. Compute Information gain $I_i(X;Y) \; \forall i \in I$

3. determine winner feature $j = \arg\max\{I_i(X,Y)\}$

4. partition $X,Y$ into disjoint subsets for each root node
$$X_i = \{x_k \in X | x_k^{(j)} = i\} \qquad i = 1, \dots v_j$$

5. create new node $N_i$

6. back to 1.    $-11-$

| DAY | Weather | Humidity | Wind | Play |
|-----|---------|----------|------|------|
| D1 | Sunny | high | weak | No |
| | ⋮ | ⋮ | ⋮ | Yes |
| | | | | ⋮ |
| | | | | Yes |

```
                    ┌──────────┐
                    │ Weather  │
                    └──────────┘
          ╱              │              ╲
        ╱                │                ╲
      ╱                  ▼                  ╲
   ▼                ┌──────────┐              ▼
┌────────┐         │ OVERCAST │         ┌────────┐
│ SUNNY  │         └──────────┘         │ RAIN   │
└────────┘                              └────────┘
     │       Weather ↗                   Weather ↗
     │       ┌                 ┐         ┌              ┐
     │       │ Overcast        │         │ Rain         │
     │       │ Overcast        │         │ Rain         │
     │       │   ⋮             │         │   ⋮          │
     ▼       │ Overcast        │         │ Rain         │
┌──────────┐ └                 ┘         └              ┘
│ HUMIDITY │
└──────────┘              ⋮
   │      ╲
```

{ splitting X,Y into disjoint subsets

```
┌────────┐          ┌────────┐
│ HIGH   │          │ Normal │
└────────┘          └────────┘
```

| High | No |
|------|----|
| High | No |
| | No |

**NO**

| Normal | Yes |
|--------|-----|
| Normal | Yes |

**Yes**

{ leaf node = all target values equal

ID3 has bottlenecks, like requiring discrete features and a tendency to overfit. ← $I(X;Y)$ maximal if feature has all attributes different

An improvement is C4.5. Here instead of $I(X;Y)$ a normalized version is used:

$$\frac{I(X;Y)}{\text{SplitInformation}(X;Y)} = \text{GainRatio}(X;Y)$$

~ expected reduction in entropy by knowing attribute Y

information generated by splitting the training set X into i partitions

$$\text{SplitInformation}(X;Y) = -\sum_{X} \frac{|X_i|}{|X|} \log_2 \left[\frac{|X_i|}{|X|}\right]$$
(= Intrinsic Value)

fraction of examples

for example: $14 \cdot \left[-\frac{1}{14} \log_2 \left[\frac{1}{14}\right]\right] \approx 3.807$

| day | class |
|-----|-------|
| 1.6 | C0 |
| 2.6 | C0 |
| ⋮ | i |
| 14.6 | C1 |

⇒ reduce bias towards multi-valued attributes
⇒ more reliable than Information Gain

-13-

| Day | Wind | Class |
| --- | --- | --- |
| 1 | Weak | No |
| 2 | Strong | No |
| 3 | Weak | yes |
| 4 | weak | yes |
| 5 | Weak | yes |
| 6 | Strong | No |
| 7 | Strong | yes |
| 8 | Weak | No |
| 9 | Weak | yes |
| 10 | Weak | yes |
| 11 | Strong | yes |
| 12 | Strong | yes |
| 13 | Weak | yes |
| 14 | Strong | No |

## Example:

$$H(X) = -\frac{5}{14}\log_2\frac{5}{14} - \frac{9}{14}\log_2\frac{9}{14}$$

$$= 0.940 \text{ bits}$$

$$H(\text{Class} \mid \text{Wind} = \text{weak}) = -\frac{2}{8}\log_2\frac{2}{8} - \frac{6}{8}\log_2\frac{6}{8}$$

$$= 0.811 \text{ bits}$$

$$H(\text{Class} \mid \text{Wind} = \text{strong}) = -\frac{3}{6}\log_2\frac{3}{6} - \frac{3}{6}\log_2\frac{3}{6}$$

$$= 1 \text{ bits}$$

$$H(\text{Class} \mid \text{Wind}) = \frac{8}{14} H(\text{Class} \mid \text{Wind} = \text{weak}) + \frac{6}{14} H(\text{1st})$$

$$\Rightarrow I(X;Y) = 0.940 \text{ bits} - \ldots = \underline{0.049 \text{ bit}}$$

$$\text{SplitInfo}(\text{Class};\text{Wind}) = -\underbrace{\frac{8}{14}}_{\text{8 weak wind}}\log_2\frac{8}{14} - \underbrace{\frac{6}{14}}_{\text{6 strong wind}}\log_2\frac{6}{14}$$

$$\Rightarrow \text{GainRatio}(X;Y) = \frac{0.049}{0.985} \approx 0.049$$

CART (Classification and regression tree) uses <u>Gini index</u> as an attribute selection measure to build a decision tree:

Idea: $p_i = 1 \Rightarrow 0$; $p_i = 0 \Rightarrow 0 \Rightarrow 0$ for pure splits

$$\sum_{i=1}^{m} \left[ p_i (1 - p_i) \right] = \sum (p_i - p_i^2) =$$

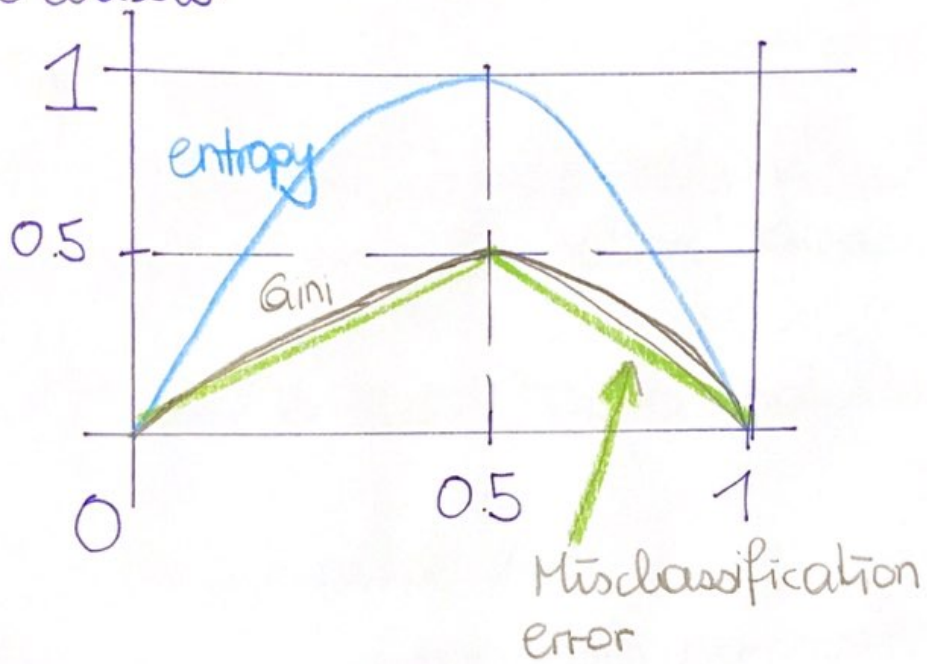$$= \sum_{i=1}^{m} p_i - \sum_{i=1}^{m} p_i^2 = 1 - \sum_{i=1}^{m} p_i^2 =$$

$$= \text{Gini Impurity}(X;Y)$$

$$\text{Gini Index} = \sum \frac{|X_i|}{|X|} \text{Gini Impurity}(X;Y)$$

| Attribute | Class |
|-----------|-------|
| M | C0 |
| M | C0 |
| F | C1 |
| F | C0 |
| M | C1 |

for M:
$$1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 = A \text{ Impurity}$$

for F:
$$1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = B$$

Gini Index: $\boxed{\frac{3}{5}}$ A $+ \boxed{\frac{2}{5}}$ B

3×C0    2×C1

two classes:

# Sources:

[1] O. Rioul, "This is IT: A primer on Shannon's entropy and information." Seminaire Poincaré XXIII (2018) 43-77.

[2] Runkler, T.A. (2020) "Data Analytics"

[3] aitimejournal.com

[4] sefiks.com/2018/05/13/a-step-by-step-c4-5-decision-tree-example/