# Gini Index

✓ There are <u>nine</u> tuples belonging to the class buys_computer = yes and the remaining <u>five</u> tuples belong to the class buys_comp = no.

$$Gini(D) = 1 - \left(\frac{9}{14}\right)^2 - \left(\frac{5}{14}\right)^2$$

$$= 0.459$$

✓ To find the splitting criterian for the tuples in D, we need to compute Gini Index for each attribute

✓ let's start with the attribute <u>income</u> and consider each of the possible splitting subset

## Incomes

possible Values are {low, medium, high}

possible subsets are

{low, medium, high}

{low, medium}

{low, high}

{medium, high}

{low},

{medium},

{high}

{ }

⟹ from this we can exclude powerset {low, medium, high} and empty set { } since conceptually they do not represent split.

Consider the subset

{low, medium}

| Income | Yes | No | No. of Instances |
|--------|-----|-----|------------------|
| high | 2 | 2 | 4 |
| medium | 4 | 2 | 6 |
| low | 3 | 1 | 4 |
| | | | 14 |

$$Gini_{income \in \{low, medium\}}^{(D)}$$

$$= \frac{10}{14} Gini(D_1) + \frac{4}{14} Gini(D_2)$$

$$= \frac{10}{14}\left(1 - \left(\frac{7}{10}\right)^2 - \left(\frac{3}{10}\right)^2\right) + \frac{4}{14}\left(1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2\right)$$

$$= 0.443$$

$$= Gini_{income \in \{high\}}^{(D)}$$

$$Gini_{income \in \{low, high\}}^{(D)}$$

$$= \frac{8}{14}\left(1 - \left(\frac{5}{8}\right)^2 - \left(\frac{3}{8}\right)^2\right) + \frac{6}{14}\left(1 - \left(\frac{4}{6}\right)^2 - \left(\frac{2}{6}\right)^2\right)$$

$$= 0.458$$

$$= Gini_{income \in \{medium\}}^{D}$$

Gini

income $\in \{medium, high\}^{(D)}$

$$= \frac{10}{14}\left(1 - \left(\frac{6}{10}\right)^2 - \left(\frac{4}{10}\right)^2\right) + \frac{4}{14}\left(1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2\right)$$

$$= 0.450$$

$$= Gini_{income \in \{low\}}^{(D)}$$

Gini for subsets

$\{low, medium\} = 0.443$ ✓

$\{high\}$ $\neq$ $= 0.443$ ✓

$\{low, high\} = 0.458$

$\{medium\} = 0.458$

$\{medium, high\} = 0.450$

$\{low\} = 0.450$

→ Best binary split for attribute income
is on $\{low, medium\}$ or $\{high\}$ because
it minimizes Gini Index

# Age

## Gini age

**Proper subsets**

$\{$youth, middle-aged, senior$\}$

$\{$youth, middle-aged$\}$

$\{$youth, senior$\}$

$\{$middle-aged, senior$\}$

$\{$youth$\}$, $\{$middle-aged$\}$, $\{$senior$\}$

$\nearrow$  $\{\}$

Exclude powerset and nullset

so we will get

only 6

| age | yes | no | no of Instances |
|---|---|---|---|
| youth | 2 | 3 | 5 |
| middle aged | 4 | 0 | 4 |
| Senior | 3 | 2 | 5 |
| | | | 14    D |

gini

$$\text{age} \in \{\text{youth, middle-aged}\}$$

$$= \frac{9}{14}\left(1 - \left(\frac{6}{9}\right)^2 - \left(\frac{3}{9}\right)^2\right) + \frac{5}{14}\left(1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2\right)$$

$$= 0.4571$$

$$= \text{gini}_{\text{age} \in \{\text{senior}\}}^{(D)}$$

gini

$$\text{age} \in \{\text{youth, Senior}\}^{D}$$

$$= \frac{10}{14}\left(1 - \left(\frac{5}{10}\right)^2 - \left(\frac{5}{10}\right)^2\right) + \frac{4}{14}\left(1 - \left(\frac{4}{4}\right)^2 - \left(\frac{0}{4}\right)^2\right)$$

$$= 0.3571$$

$$= \text{gini}_{\text{age} \in \{\text{middle-aged}\}}^{D}$$

gini                                                    (D)

    age $\in$ {middle-aged, Senior}

$$= \frac{9}{14}\left(1-\left(\frac{7}{9}\right)^2-\left(\frac{2}{9}\right)^2\right) + \frac{5}{14}\left(1-\left(\frac{2}{5}\right)^2-\left(\frac{3}{5}\right)^2\right)$$

$$= 0.3936$$

$$= gini_{age \in \{youth\}}^{(D)}$$

✓ possible  subsets

gini
    {youth, middle-aged} = 0.4571
    {Senior}  = 0.4571
    {youth, Senior}  = 0.3571 ✓
    {middle-aged}  = 0.3571 ✓
    {middle-aged, Senior} = 0.3936
    {youth}  = 0.3936

✓ We  obtain  {youth, Senior} or {middle-aged}
    as  the  best  split  for  age  with  a
    gini  index  of  0.357

## Student

Values are {Yes, no}

| Student | Yes | No | No of Instances |
|---|---|---|---|
| Yes | 6 | 1 | 7 |
| No | 3 | 4 | 7 |
| | | | 14 |

$$gini_{student}(D) = \frac{7}{14}\left(1-\left(\frac{6}{7}\right)^2-\left(\frac{1}{7}\right)^2\right)+$$

$$\left(\frac{7}{14}\right)\left(1-\left(\frac{3}{7}\right)^2-\left(\frac{4}{7}\right)^2\right)$$

$$= 0.367$$

## Credit-rating

Values are {Fair, Excellent}

| Credit-rating | Yes | No | No of Instances |
|---|---|---|---|
| Fair | 6 | 2 | 8 |
| Excellent | 3 | 3 | 6 |

$$gini_{credit.rating}(D) = \frac{8}{14}\left(1-\left(\frac{6}{8}\right)^2-\left(\frac{2}{8}\right)^2\right)+\frac{6}{14}\left(1-\left(\frac{3}{6}\right)^2-\left(\frac{3}{6}\right)^2\right)$$

$$= 0.428$$

| Attribute | split | Gini Index | Reduction in impurity $\Delta gini = gini(0) - gini_A(0)$ |
|---|---|---|---|
| age | {youth, senior} or {middle_aged} | $0.3571$ | $0.459 - 0.357 = 0.102$ |
| Income | {low, medium} or {high} | $0.443$ | $0.459 - 0.443 = 0.016$ |
| student | Binary | $0.367$ | $0.459 - 0.367 = 0.092$ |
| credit_rating | Binary | $0.428$ | $0.459 - 0.428 = 0.031$ |

age ?

{youth, senior}          {middle_aged}

RID: 1, 2, 4, 5, 6, 8, 9, 10, 11, 14          RID: 3, 7, 12, 13

↓ All belongs to Same class