# Decision Tree tutorial

## Dataset

| age | income | student | credit-rating | target buys_computer |
|---|---|---|---|---|
| ↓ | high | yes | fair | |
| Youth | medium | No | excellent | |
| middleaged | low | | | |
| Senior | | | | |

Total_data = 14

: impurity measure :

— cross entropy | gini-Index

## Multi-way split using cross entropy

Consider Attribute 'Age'

Youth < 3-No / 2-Yes     Senior < 3-Yes / 2-No     middleaged < 4-Yes / 0-No

$$\text{cross-entropy} = -\sum_{k=1}^{k} p_{mk} \log p_{mk}$$

$\underline{K=2}$

initially $\underline{m=1}$

$K=1 \to$ class Yes
$K=2 \to$ class No

$$-(p_{11} \log p_{11} + p_{12} \log p_{12})$$

For attribute age = Youth $p_{11} = \frac{2}{5}$, $p_{12} = \frac{3}{5}$

$$\text{cross} \atop \text{entropy}_{Age} = +(5/14)\left(-\frac{2}{5}\log\frac{2}{5} - \frac{3}{5}\log\frac{3}{5}\right) + \frac{4}{14}\left(-\frac{4}{5}\log\frac{4}{5}\right) + \frac{5}{14}\left(-\frac{3}{5}\log\frac{3}{5}\right.$$

$$\left. -\frac{2}{5}\log\frac{3}{5}\right)$$

$$= \boxed{0.6935}$$

Cross_entropy/credit_range (D) = $\boxed{0.8922}$

Cross_entropy_income : $\boxed{0.9111}$

Cross_entropy_student : $\boxed{0.7885}$

less entropy →

# Decision Tree
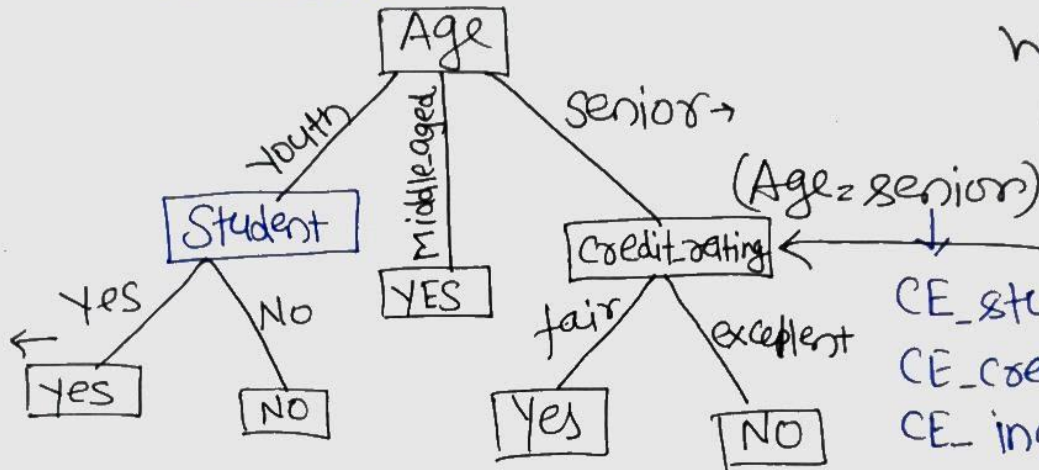
Age
- Youth →
- Middle Aged → **Yes**
- Senior

⌐ You will filter dataset where age = youth

Cross_entropy (Age = youth) = Income

Cross_entropy (Age = youth) = 0 → Minimum value of entropy
student

Cross_entropy (Age = youth) = Credit_score

Age
- youth → **Student**
  - Yes → **Yes** ← Filter data ↓ where (Age=youth) & (Student=yes)
  - No → **NO**
- Middle-aged → **YES**
- senior → **Credit_rating** ← (Age = senior)
  - fair → **Yes**
  - excellent → **NO**

We will stop since entropy = 0

CE_student = 0.95
CE_credit_rating = 0
CE_income = 0.95

# Decision-Tree!

## Binary split using Gini index

Consider 'Age'

Total data points = 14

(+ve) class_proportion: Age = Youth: 2/5

Age = Midd : 1

Age : Senior : 3/5

### Ordering the probability.

Youth, senior, middle

$$\text{Gini-index} = \sum_{R=1}^{K} p_{mk}(1-p_{mk})$$

### Possible splits

{Youth}, {senior, middle}  |  {youth, senior}, {middle}

⇒ ∴ we need to calculate impurity measure for both the above splits.

$$\text{Gini}_{age \in \{youth\}}(D) = 2p(1-p)$$

$$= \frac{5}{14}\left(2 * \frac{2}{5} * \frac{3}{5}\right) + \frac{9}{14}\left(2 * 7/g * 2/g\right)$$

$$= 0.6508$$

For '2' class problem $G = p(1-p) + (1-p)p$

$$= \boxed{2p(1-p)}$$

Gini-age = {youth, senior} =

Gini-student {yes} = $7/14\left(2 * \frac{3}{7} * \frac{4}{7}\right) + 7/14\left(2 \cdot \frac{2}{4} \cdot \frac{1}{2}\right)$

Gini- income → Two splits =

Gini-credit_rating {~~fair~~} = $\frac{8}{14}\left(2 * 6/8 * 2/8\right) + \frac{6}{14}\left(2 * \frac{3}{6} * \frac{3}{6}\right)$

Choose the one having less impure. : lesser value

$(Age)$

{youth, senior} — Age — Middle_age

filter data
↓,
again

Yes

Calculate gini_index for all other attribute other than age.