

Gini index

measure of inequality earlier.

0 no inequality
1 inequality 100%.

Lorenz curve — Index = area b/w 45% line ($y=x$) and the curve.

area \propto inequality
area estimated by trapezoidal rule

↓
less area
→ curve aligns with $y=x$
→ less inequality

↓
more area
→ curve no align w/ $y=x$
→ more inequality

Decision Tree :

dataset $\rightarrow S \rightarrow n$ classes.

$$S = n_1 + n_2 \text{ (binary classification)}$$

$$G(S) = \frac{n_1}{S} \cdot G(S_1) + \frac{n_2}{S} \cdot G(S_2)$$

$$G(S) = 1 - \sum p_i^2 \quad \text{relative freq of class } i$$

Ex: ① Unbiased coin:

$$P_{head} = 0.5 \quad P_{tail} = 0.5$$

$$G(S) = 1 - \sum p_i^2 = 1 - [(0.5)^2 + (0.5)^2] = 0.5$$

② Biased coin:

$$P_{head} = 0.3 \quad P_{tail} = 0.7$$

$$G(S) = 1 - \sum p_i^2 = 1 - [(0.3)^2 + (0.7)^2] = 0.42$$

	home owns	Married	Gender	Employed?	Credit rate	Risk class
1	Yes	Yes	M	Yes	A	B ✓
2	No	No	F	Yes	A	A ✓
3	Yes	Yes	F	Yes	B	C
4	Yes	No	M	No	B	B ✓
5	No	Yes	F	Yes	B	C
6	No	No	F	Yes	B	A ✓
7	No	No	M	No	B	B ✓
8	Yes	No	F	Yes	A	A ✓
9	No	Yes	F	Yes	A	C
10	Yes	Yes	F	Yes	A	C

class

③ Attribute: Married

$$(1) Y: A=0 \quad B=1 \quad C=4$$

$$S = 10 \quad | \quad A=3 \quad B=3 \quad C=4 \quad \text{here, final class layer.}$$

$$\textcircled{1} \quad G(S) = 1 - \sum p_i^2 = 1 - \left[\left(\frac{3}{10} \right)^2 + \left(\frac{3}{10} \right)^2 + \left(\frac{4}{10} \right)^2 \right]$$

gini index for overall dist. $G(S) = 0.66$

→ Go attribute by attribute now.

② Attribute: Home

$$\textcircled{1} \quad Y: A=1 \quad B=2 \quad C=2$$

$$\textcircled{15} \quad N: A=2 \quad B=1 \quad C=2$$

$$G(Y_{es}) = 1 - \sum p_i^2 = 1 - \left[\left(\frac{1}{5} \right)^2 + \left(\frac{2}{5} \right)^2 + \left(\frac{2}{5} \right)^2 \right] = 0.64$$

$$G(N_{es}) = 1 - \sum p_i^2 = 1 - \left[\left(\frac{1}{5} \right)^2 + \left(\frac{2}{5} \right)^2 + \left(\frac{2}{5} \right)^2 \right] = 0.64$$

$$G(\text{home}) = \frac{n_1}{S} G(S_1) + \frac{n_2}{S} G(S_2)$$

$$= \frac{5}{10} \cdot 0.64 + \frac{5}{10} \cdot 0.64 = 0.64$$

③ Attribute: Married

$$(1) Y: A=0 \quad B=1 \quad C=4$$

$$(2) N: A=3 \quad B=2 \quad C=0$$

$$G(Y) = 1 - \sum p_i = 1 - \left[\left(\frac{0}{5}\right)^2 + \left(\frac{1}{5}\right)^2 + \left(\frac{4}{5}\right)^2 \right] = 0.32$$

$$G(N) = 1 - \sum p_i = 1 - \left[\left(\frac{3}{5}\right)^2 + \left(\frac{2}{5}\right)^2 + \left(\frac{0}{5}\right)^2 \right] = 0.48$$

$$G(\text{married}) = \frac{n_1}{S} G(S_1) + \frac{n_2}{S} G(S_2) = \frac{5}{10} (0.32) + \frac{5}{10} (0.48) = 0.4$$

④ Attribute: Gender

$$(3) M: A=0 \quad B=3 \quad C=0$$

$$(4) F: A=3 \quad B=0 \quad C=4$$

$$G(\text{male}) = 1 - \sum p_i = 1 - \left[0^2 + \left(\frac{5}{3}\right)^2 + 0 \right] = 0$$

$$G(\text{female}) = 1 - \sum p_i = 1 - \left[\left(\frac{3}{7}\right)^2 + \left(\frac{4}{7}\right)^2 \right] = 0.511$$

$$G(\text{gender}) = \frac{n_1}{S} G(S_1) + \frac{n_2}{S} G(S_2)$$

$$= \frac{3}{7}(0) + \frac{4}{7}(0.511) = 0.358$$

⑤ Attribute: Credit_risk

$$(5) a: A=2 \quad B=2 \quad C=2$$

$$(5) b: A=1 \quad B=1 \quad C=2$$

Same as Home

$$G(\text{credit}) = 0.64$$

⑥ Gini Index Gain:

<u>Attribute</u>	<u>Gini index before split</u>	<u>Gini index after split</u>	<u>Gain</u>
Home	0.66	$G(\text{home})=0.64$	0.02
Married	0.66	$G(\text{married})=0.40$	0.26
Gender	0.66	$G(\text{gender})=0.358$	0.302
Employed	0.66	$G(\text{employed})=0.475$	0.185
Credit	0.66	$G(\text{credit})=0.64$	0.02

⑦ Highest info gain @ gender

↳ split there first.

overall gini index of full dataset S.

$$G(S)$$

Now another example, but w/ info gain + entropy.

NO	Age	Income	Profession	Will See Movie
1	L >30	LOW	Business	Yes
2	H >30	HIGH	Engg	No
3	L <30	MED	Engg	Yes
4	L <30	LOW	Agri	No
5	H >30	HIGH	Business	Yes
6	H >30	MED	Agri	No

① Overall entropy:

$$S = 6 \quad | \quad Y=3 \quad N=3$$

$$H(S) = -\sum p_i \log p_i = -\frac{3}{6} \log_2 \left(\frac{3}{6}\right) - \frac{3}{6} \log_2 \left(\frac{3}{6}\right) = 1$$

② Attribute: Age

$$(3) \text{ Low } : Y=2 \quad N=1$$

$$(3) \text{ High } : Y=1 \quad N=2$$

$$H(\text{low}) = -\frac{2}{3} \log_2 \left(\frac{2}{3}\right) - \frac{1}{3} \log_2 \left(\frac{1}{3}\right) = 0.918$$

5	H>30	High	Business	Yes
6	H>30	MED	Agri	No

Engineer : 23y, low income - movie?

② Attribute: Income

$$(2) \text{ Low : } Y=1 \quad N=1$$

$$(2) \text{ Med : } Y=1 \quad N=1$$

$$(2) \text{ High : } Y=1 \quad N=1$$

$$H(\text{low}) = H(\text{med}) = H(\text{high})$$

$$= -\frac{1}{2} \log\left(\frac{1}{2}\right) - \frac{1}{2} \log\left(\frac{1}{2}\right)$$

$$= 1 \times -\frac{1}{2} \log\left(\frac{1}{2}\right) = -1 \log_2(2)^{-1} = 1$$

Take weighted : (weights are frac. of L, M, H)

$$H(\text{income}) = \frac{2}{6}(1) + \frac{2}{6}(1) + \frac{2}{6}(1) = 1$$

$$\Rightarrow \text{Information gain} = 1 - H(\text{income}) = 0$$

Attribute	Entropy	Information
Age	0.918	0.082

Income 1 0

Profession	0.33	0.67

$$H(\text{low}) = -\frac{2}{3} \log\left(\frac{2}{3}\right) - \frac{1}{3} \log\left(\frac{1}{3}\right) = 0.918$$

$$H(\text{high}) = -\frac{1}{3} \log\left(\frac{1}{3}\right) - \frac{2}{3} \log\left(\frac{2}{3}\right) = 0.918$$

Now take weighted
(weight is fraction of L, H)

$$H(\text{Age}) = \frac{3}{6}(H(\text{low})) + \frac{3}{6}(H(\text{high})) = 0.918$$

Entropy for age

$$\hookrightarrow \text{Information gain} = 1 - 0.918 = 0.082$$

③ Attribute: Profession

$$(2) \text{ Busi : } Y=2 \quad N=0$$

$$(2) \text{ Agri : } Y=0 \quad N=2$$

$$(2) \text{ Eng : } Y=1 \quad N=1$$

$$H(\text{busi}) = H(\text{agri}) = \frac{2}{2} \log\left(\frac{2}{2}\right) = 0$$

$$H(\text{eng}) = -\frac{1}{2} \log\left(\frac{1}{2}\right) - \frac{1}{2} \log\left(\frac{1}{2}\right) = 1$$

Take weighted : weights from B, A, E

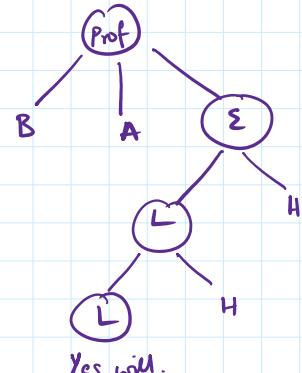
$$H(\text{prof}) = \frac{2}{6}(0) + \frac{2}{6}(0) + \frac{2}{6}(1) = 0.33$$

$$\Rightarrow \text{Info gain} = 1 - 0.33 = 0.67$$

Least entropy
Most info gain] 1st split

Next is age.

Last split is income
↳ doesn't really matter.



Next Example.

- Pizza Dataset : using this dataset, draw

Possible Trees (among any 2 variables)

① how many decision trees are at least possible?

② what is the hypothesis space?

N binary attributes $\rightarrow 2^N$ trees

M n-many attributes $\rightarrow n^m$ trees

$n = \text{any of attr.}$
 $m = \text{no. of}$

Entire set of possible trees.

$$\text{Overall entropy} \leq \left\{ \begin{array}{l} P: S \\ N: 3 \end{array} \right\} - \frac{5}{8} \log\left(\frac{5}{8}\right) - \frac{3}{8} \log\left(\frac{3}{8}\right) = H(S)$$

Rest is same as prev. example.

Sample	Crust Size	Shapes	Filling Size	Class
1 e1	Big	Circle	Small	Pos.
2 e2	Small	Circle	Small	Pos.
3 e3	Big	Square	Small	Neg.
4 e4	Big	Triangle	Small	Neg.
5 e5	Big	Square	Big	Pos.
6 e6	Small	Square	Small	Neg.
7 e7	Small	Square	Big	Pos.
8 e8	Big	Circle	Big	Pos.

Q) Using the above dataset, use Entropy Gain and find out which attribute to split first.

One more, no table this time.

(Q2) Consider Training Data Set S with 100 instances and 2 attributes — Height & Weight. S has 70 +ve and 30 -ve examples

$$\text{Values (Height)} = \{\text{Med, Tall}\}$$

$$\text{Values (Weight)} = \{\text{Normal, Fat}\}$$

$$S_{\text{Med}} = [0+, 60-] \quad S_{\text{Tall}} = [20+, 20-]$$

$$S_{\text{Normal}} = [35+, 35-] \quad S_{\text{Fat}} = [30+, 0-]$$

Decide which attribute partitions samples S better.

② Attribute: Weight:

$$(20) \text{ Normal: } P=85 \quad N=35$$

$$(30) \text{ Fat: } P=30 \quad N=0$$

$$H(\text{Norm}) = -\frac{35}{70} \log\left(\frac{35}{70}\right) - \frac{35}{70} \log\left(\frac{35}{70}\right) = 1$$

$$H(\text{Fat}) = -\frac{30}{30} \log\left(\frac{30}{30}\right) - \frac{0}{30} \log\left(\frac{0}{30}\right) = 0$$

Weighted: (weight is fraction of N, F)

$$H(\text{Weight}) = \frac{70}{100} [H(\text{Norm})] + \frac{30}{100} [H(\text{fat})] = 0.7$$

$$\text{Info gain (Weight)} = 0.3$$

100 instances $\Rightarrow S=100$

$$P=70 \quad N=30$$

Height
— Med
— Tall

Weight
— Normal
— Fat

① Attribute: Height

$$(20) \text{ Med: } P=0 \quad N=60$$

$$(40) \text{ Tall: } P=20 \quad N=20$$

$$H(\text{med}) = -\frac{0}{60} \log(0) - \frac{60}{60} \log(1) = 0$$

$$H(\text{tall}) = -\frac{20}{40} \log\left(\frac{20}{40}\right) - \frac{20}{40} \log\left(\frac{20}{40}\right) = 1$$

Weighted: (Weight is fraction of M, T)

$$H(\text{height}) = \frac{60}{100} [H(\text{med})] + \frac{40}{100} [H(\text{tall})] = 0.4$$

$$\text{Info gain (height)} = 0.6$$

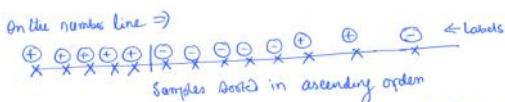
Attribute	Entropy	Info gain
Height	0.4	0.6
Weight	0.7	0.3

Most info gain + least entropy.

Next example:

(Q3) Decision tree uses Numerical Attribute that is Continuous

In a dataset, if the continuous attribute is sorted in ascending order, it looks like below where the labels are \oplus and \ominus



Q) using this numeric attribute, how do you find the best split point?

Key: * Since continuous number, first sort in ascending order

* Then check where the 'label' flips
for example, between 5th & 6th sample,
10th & 11th sample,
12th & 13th sample

* Here we have 3 split points possible
 $\theta_1, \theta_2, \theta_3$

* Using Information Gain, evaluate these $\theta_1, \theta_2, \theta_3$

③ Numerical + categorical Attribute Hints

Dataset ↳ Find "Best Split" using GINI IMPURITY

Age	Like Fruits	Like KFC	Going to be Healthy
51	1	1	1
42	0	0	0
44	1	1	1
30	1	1	1
29	0	0	0
36	0	1	1
36	0	0	0
47	1	1	1
46	1	0	0
47	0	1	0

Hints

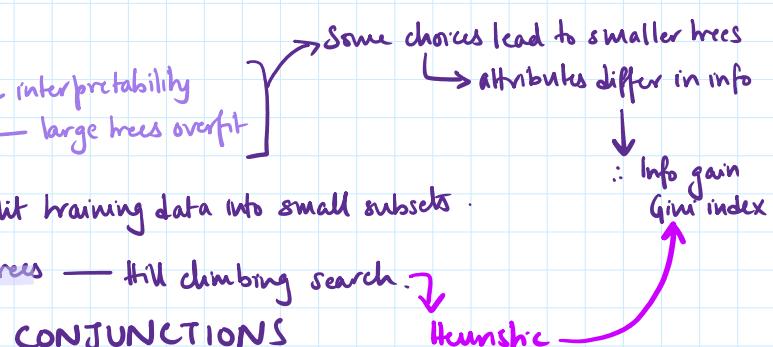
- ① Nominal attribute → Age
So, ascending order arrange
- ② Find the possible split points
of every adjacent pair
↳ for each calculate Gini heuristic
- ③ Do the same for Categorical

Effectively converts Nominal → Categorical or Boolean.

- ④ Choose the Best Split

Some stuff to remember :

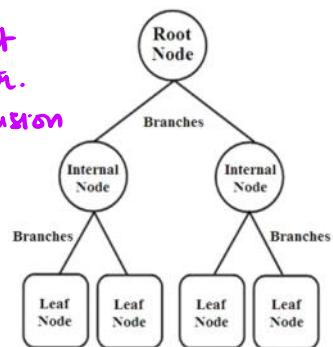
- ① Small tree >> Big tree
- ② Decision tree philosophy: Split training data into small subsets .
- ③ Hypothesis Space = space of trees — hill climbing search .
- ④ DISJUNCTION OF CONJUNCTIONS



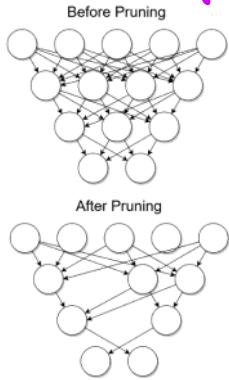
Final takeaways:

- ① DT = flowchart-like struct
- ② Adv: a. highly interpretable
b. little data prep.
c. classification + regression = versatile .
- ③ Disadv: a. more levels dominate — bias.
b. unstable — small variations = diff tree entirely .
c. overfit — hard to generalize

Before Pruning



- b. unstable — small variations = diff tree entirely.
- c. overfit — hard to generalize



Pruning

- reduce no. of connections + nodes.
- a. Prepruning (early stopping)
 - stop tree from growing once criteria met.
- b. Post-pruning
 - remove branches from fully grown tree

④ ID3 algorithm (Iterative dichotomizer 3) — what we did w/ entropy + info gain.

a. Inductive bias

- describe basis by which it chooses one hypothesis

- simple to complex heuristic search (use gain)
- shorter trees \gg longer trees — Occam's razor.
- attr. of higher gain preferred.
- attr. w/ many values preferred

b. Hypothesis space = all possible DT that can be made.

c. Hyperparameter = maximum depth of the tree.

- too low = generalization
- too high = memorization.

d. Advantages =

- simple to understand
- requires little training data.
- discrete + continuous attr.

e. Disadvantages =

class of target predict cont. variables value

⑤ CART (Classification & Regression Trees) → Binary tree — gini index.

a. Advantages =

- simple
- numerical + categorical data
- missing values — impute w/ surrogate splits
- multi-class classification

b. Disadvantages =

- overfitting
- greedy (may not find most optimal)
- biased towards attr. w/ many categories