## Decision Trees

K-NN, NB, log. reg, lr-reg, SVM . DT ──→ (if...else).

- K-NN, NB → instance based method, probabilistic method
- log.reg, lr.reg, SVM → (geometric, hyperplane).

DT ──→ simply a nested if else condition classifier. — programmatic interpretab.
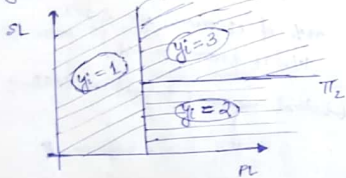
ex → $z_i = <SL, PL, SW, PW>$

model $\begin{bmatrix} \text{if } PL < a: \\ \quad y_i = 1 \\ \text{else} \\ \quad \text{if } SL < b: \\ \quad\quad y_i = 2 \\ \quad \text{else} \\ \quad\quad y_i = 3 \end{bmatrix}$


PL < a → Root-Node.
Y → $y_i = 1$
N → SL < b
Y → $y_i = 2$, N → $y_i = 3$

※ All leaf nodes are decision. (classif$^n$ to a class).

※ At all internal nodes, we make a decision.

### geometric intuition



DT — a set of axis parallel hyperplanes that divide the entire region into hypercubes & hyper-cuboids.

### Building a DT : Entropy.

say Y is a random variable that can take 'k' values — $y_1, y_2 \cdots y_k$.

entropy $\boxed{H(Y) = \sum_{i=1}^{k} P(y_i) \log_b (P(y_i))}$  $b = 2$ or $b = e$.  $\frac{lg}{ln}$

$\boxed{P(y_i) = P(Y = y_i)}$

---

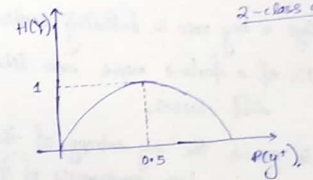say Y: play tennis $\begin{cases} + : 9 \\ - : 5. \end{cases}$  $P(y^+) = 9/14$  $P(y^-) = 5/14$.

$H(Y) = -\left\{ \frac{9}{14} \log\left(\frac{9}{14}\right) + \left(\frac{5}{14}\right) \log\left(\frac{5}{14}\right) \right\}$

### Properties of entropy (for 2-class classificat$^n$).

case 1: $\begin{cases} y_+ = 99\% \\ y_- = 1\% \end{cases}$  $H(Y) \tilde{=} 0.0801$

case 2: $\begin{cases} y_+ = 50\% \\ y_- = 50\% \end{cases}$  $H(Y) \approx 1.0$

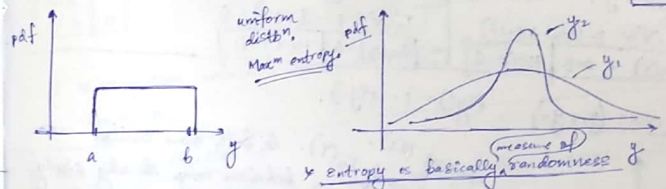case 3: $\begin{cases} y_+ = 0\% \\ y_- = 100\% \end{cases}$  $H(Y) = 0.0$


2-class case

### extending to multi-class →

※ $Y = y_1, y_2 \cdots y_k$. If all are equi-probable, $H(Y) = 1$ (max).

※ If my 1 is most probable & others have prob. $\approx 0$, then $H(Y)$ lower.

$\boxed{H(Y_1) > H(Y_2)}$


uniform distb$^n$. Max$^m$ entropy.
$y_2$, $y_1$
(measure of)
※ entropy is basically randomness

### Information Gain

※ Entropy of a table means basically w.r.t. the class label.

※ we basically break the table on the basis of a feature, then the information gain of that feature is defined as

$\boxed{IG (feature) = H(\text{original table}) - \text{weighted Avg. } H \text{ of all the constituent tables}}$

Let D be the original table.

|D| = no of entries in D.

Let feature $f_o$ breaks D to $D_1, D_2$ & $D_3$.

Then,

$$IG(f_o) = H(D) - \left\{ \frac{|D_1| * H(D_1) + |D_2| * H(D_2) + |D_3| * H(D_3)}{|D|} \right\}$$

IG plays a key role in building decision trees.

* More IG of a feature means more likely is the feature in separating diff. classes.

* More IG means the avg entropy of the constituent tables is less, i.e. randomness of the constituent tables is less.

Gini Impurity ~ similar to entropy.

$$I_G(Y) = 1 - \sum_{i=1}^{k} (P(y_i))^2$$

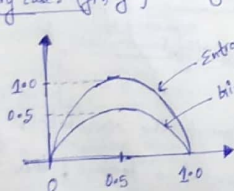Case I: $P(y_+) = P(y_-) = 0.5$.

$I_G(Y) = 1 - (0.25 + 0.25)$

$\boxed{I_G(Y) = 0.5}$  $\boxed{H(Y) = 1}$

Case II: $P(y_+) = 1$, $P(y_-) = 0$.

$I_G(Y) = 1 - (1+0) = 0$.

$\boxed{I_G(Y) = 0}$  $\boxed{H(Y) = 0}$

2-category case :- $(y_+, y_-)$  $P(y^+) = 1 - P(y_-)$.

Entropy H(Y).
Gini Impurity $I_G(Y)$. So, both have similar behaviour only. So, why study $I_G(Y)$?

→ computing "log" takes more time than computing "squares".

* So, $I_G$ is computationally efficient.

---

Constructing a DT

* Recursively at each level, select the feature with maximum IG. to partition the table.

Eventually we may end up in one of the following cases
① pure node → stop growing the node.
② very few pts of the other class. So, further growing tree may lead to overfitting. — So stop growing.
③ we are too deep down the tree. — stop growing.

If depth is small ⟹ underfit.

DT :- hyperparameter ⟹ depth.

splitting numerical features

| $f_1$ | $y$ |
|-------|-----|
| 2.2 | 1 |
| 2.6 | 1 |
| 3.5 | 0 |
| 3.8 | 0 |
| 4.6 | 1 |
| 5.3 | 0 |

. for numerical features
① sort the entries in ascending order of that numerical feature.
② serially take each value in the feature as threshold τ & splits on the decision

$$f_1 < \tau$$

③ select the τ with maximum $I_G$.

It is very time consuming as we have to go through all n entries as threshold.
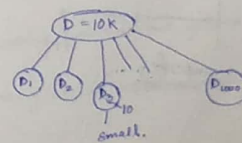
Feature Standardization

In DT, we don't depend on actual values of the features, but rather depend on their relative order.

So, we don't need to perform any kind of feature standardization.

categorical features with many categories
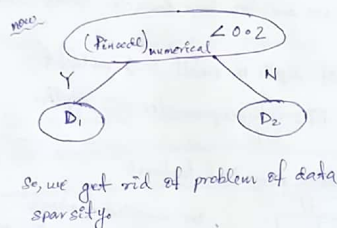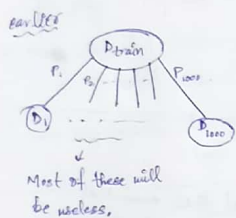
Pincode/ zipcode — 1000's of classes.

So splitting to these many classes may cause some classes to have very few data pts

So, instead of using pincode as a categorical feature, we make use of a feature engineering hack to use it as a numerical feature.

We replace every $P_j$ (pincode j) with $P(y_i=1|P_j)$. i.e,

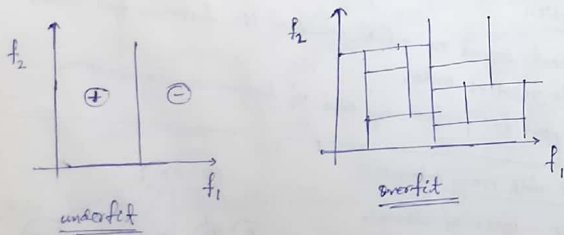probability of $y_i = 1$ given pincode j. $= \dfrac{\# y_i=1 \text{ and } P_j}{\# P_j}$

So, now we have a improved situation.

earlier

```
        D-train
   P₁  P₂  ⋯   P₁₀₀₀
  (D₁) ⋯⋯ (D₁₀₀₀)
```

Most of these will be useless.

now

```
     (Pincode)numerical  ∠0·02
       Y              N
     (D₁)            (D₂)
```

So, we get rid of problem of data sparsity.

---

## Overfitting and Underfitting

As depth↑ → possibility of having very few pts @ a leaf node ↑ (sparsity).
Also interpretability ↓ due to lots of constraints.
↳ we maybe overfitting to noise.

depth determined via CV (cross validation)



underfit

overfit

---

## Train & Run-time complexity

train ~ $O(n\log n \; d)$  
  ↳ sorting  ↳ evaluating  
$n = \# \text{pts in } D_{train}$  
$d = \dim.$

Since $d$ is a direct factor in time complexity, for large dimensions, DT mayn't be good.

After training :— storing $D_{tree}$ → nested if-else conditions.
  space $= O(\text{internal nodes}) + O(\text{leaves})$.

Runtime complexity $= O(k)$  $k = $ max depth of any leaf node.

DT → large data, dimension small  ✓
  ↳ (RF, GBDT).  low latency. used by Amazon, Google and many diff internet applications.

---

## Regression using decision tree

Instead of entropy/Gini Impurity, we use MSE (mean squared error) or MAE (median absolute error).

Just like we find IG by the difference of entropy and weighted avg entropy, here also everything remains same.

* Prediction at each level $\hat{y}_i = \text{mean}(y_i)$ in that table.

---

## Real world cases

① Imbalanced data → balance it, (upsampling/weights etc..).
  impacts entropy/MSE calcn.

② large 'd' :— @ each node, split & check for max. IG feature.
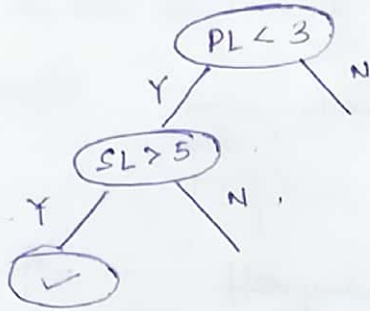  So, time complexity to train ↑.

✗ For categorical feature, avoid one-hot encoding, as 1 feature gets converted to many features.
  Instead convert to numerical features just as explained earlier.

③ Similarity Matrix — DT can't work.
Need to give features explicitly.

④ Multi-class classif^n — Naturally can be extended, don't require <u>One Vs Rest</u>

⑤ Decision Surfaces — Non-linear axis parallel hyperplanes.

⑥ Feature interactions — logical feature interactions are inbuilt in DT.

```
        ( PL < 3 )
       Y╱      ╲N
  ( SL > 5 )          (PL < 3) AND (SL > 5).
 Y╱    ╲N
( ✓ )
```

⑦ Outliers — depth ↑, outliers will impact.
Tree unstable.

⑧ Interpretability. — very much interpretable. unless depth not too much.

⑨ Feature Importance:—

★ normalize sum up reductions in $H / I_u$ due to $f_i$ — Measure of feature importance

★ Basically how much does $f_i$ helps in reducing the entropy.
or how much does $f_i$ helps in building D.T.

So, computing feature importance is straight forward.