data →

→ feature 1 ⟶ root node of the tree

decision nodes

[Root] node

No → | feature2 |     Yes → | feature 3 |

feature2: Yes → [✓]   No → [✗]

feature 3: Yes → [✓]   No → [✗]
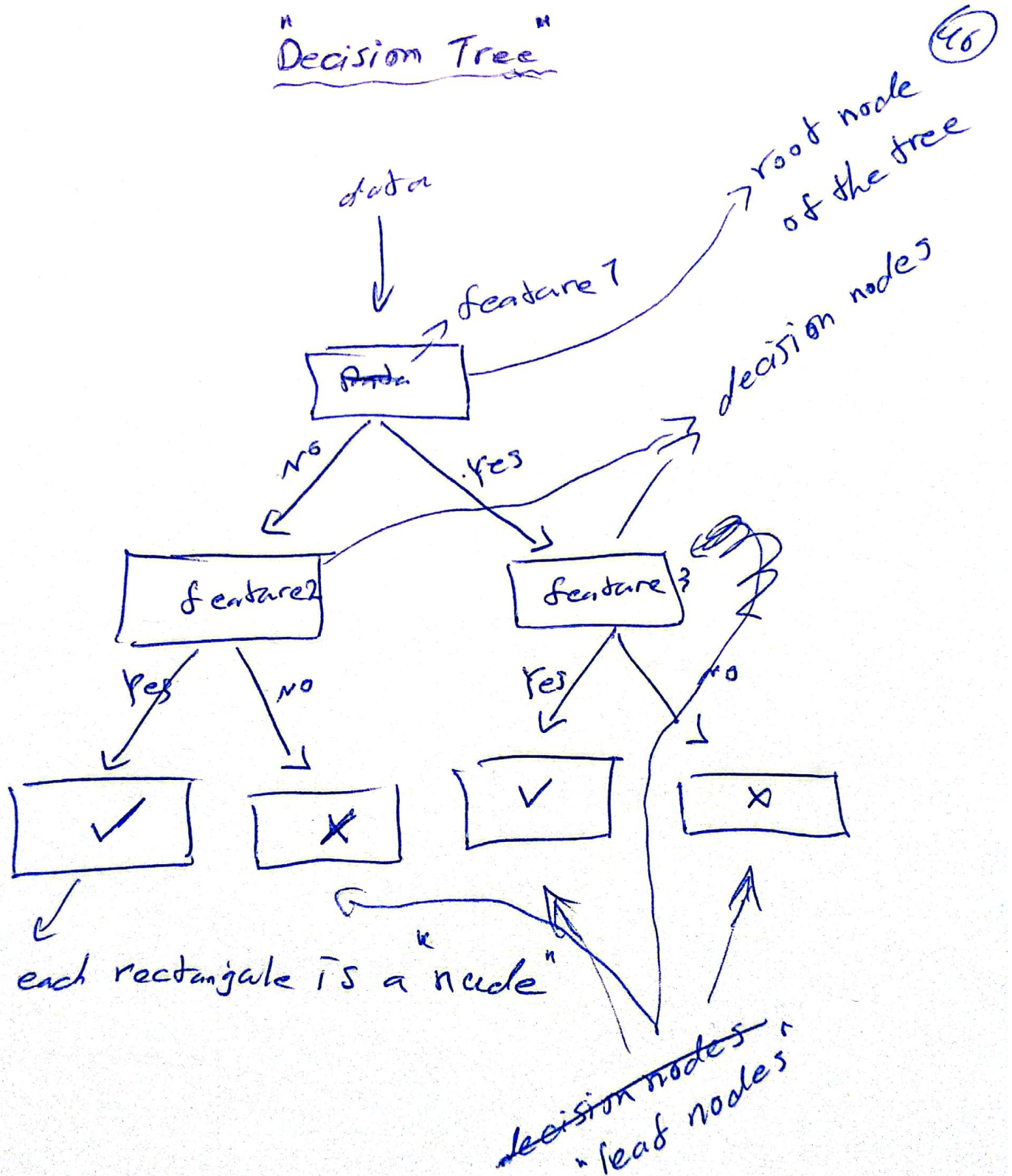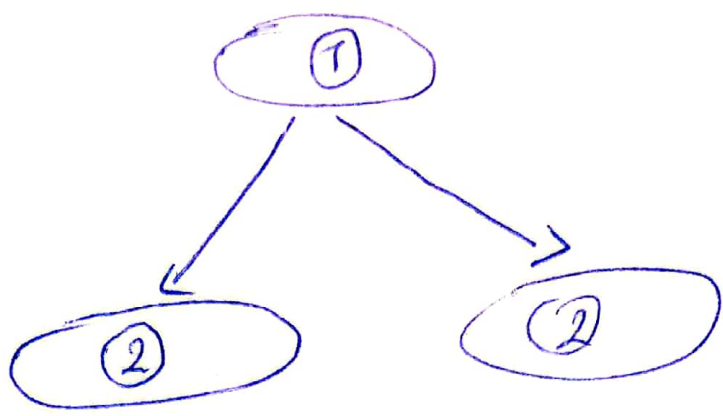
each rectangle is a "node"

"decision nodes"

"leaf nodes"
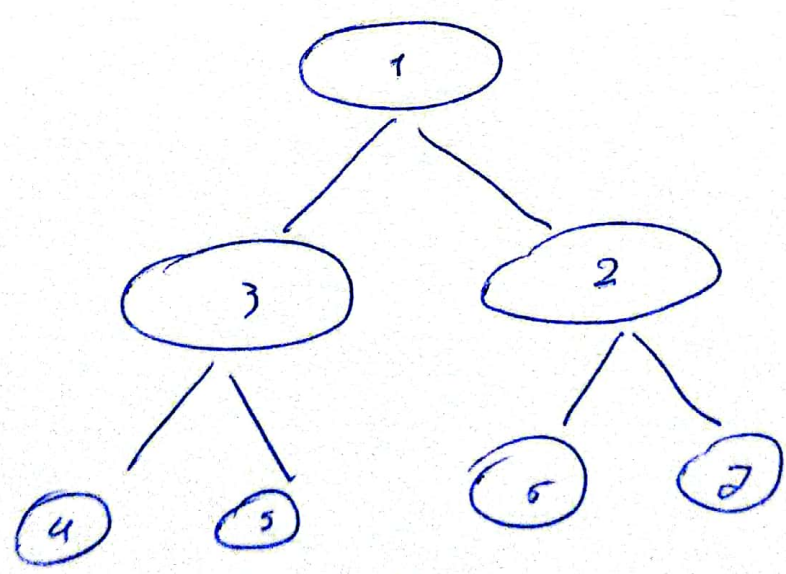
---

## learning Process

1. Pick a feature to use as "root node" and split data based on that feature.

what would be the

2. Second feature to split sub datas based on it.

.
.
.

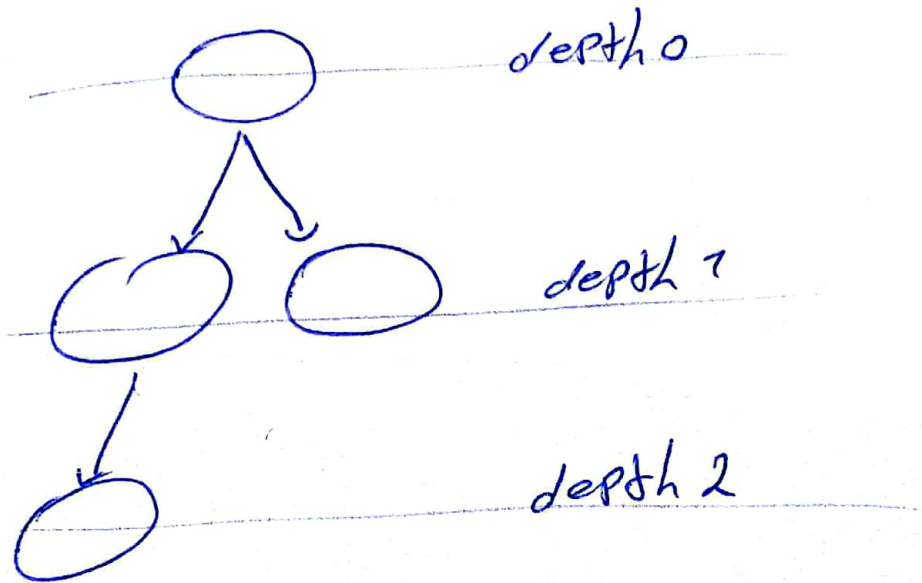3. what would be the nth. feature to split based on it



Key decisions

1. How to choose what feature to split on at each node?

use the feature which maximizes the purity or minimizes the impurity

2. ~~How~~ when do you stop splitting?

→ when a node is 100% one class.

— when splitting a node will result in the tree exceeding a maximum depth.



depth 0

depth 1

depth 2

— when improvements in purity score are below a threshold.
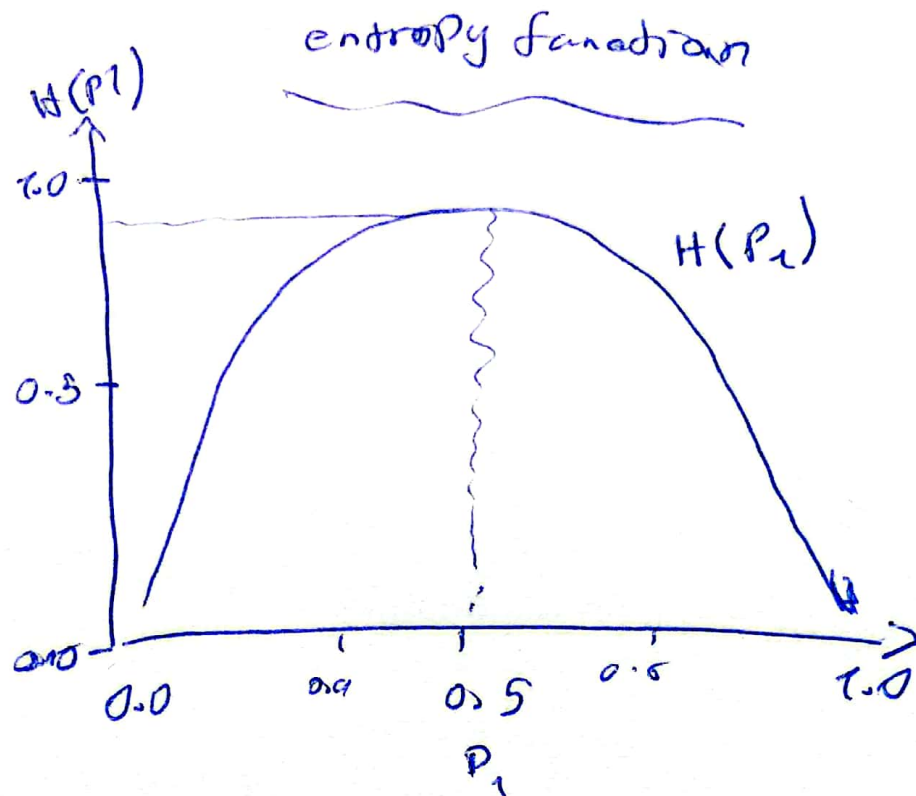
— when number of examples in a node is below a threshold.

## measuring purity

entropy as a measure of impurity.

$P_1$ = fraction of examples that have the same label.

exp:

$P_1 = \frac{3}{6}$



entropy function

$H(P1)$

$H(P_1)$

1.0

0.5

0.0    0.0    0.9    0.5    0.6    1.0

$P_1$

when the data is 50/50 it reaches the highest (1) entropy, means it is in the highest level of entropy (~~~ impurity)

exps:

$P_1 = \frac{3}{6}$      $H(P_1) = 1$

$P_1 = \frac{5}{6}$      $H(P_1) = 0.05$

$P_1 = \frac{6}{6}$      $H(P_1) = 0$

quite impure

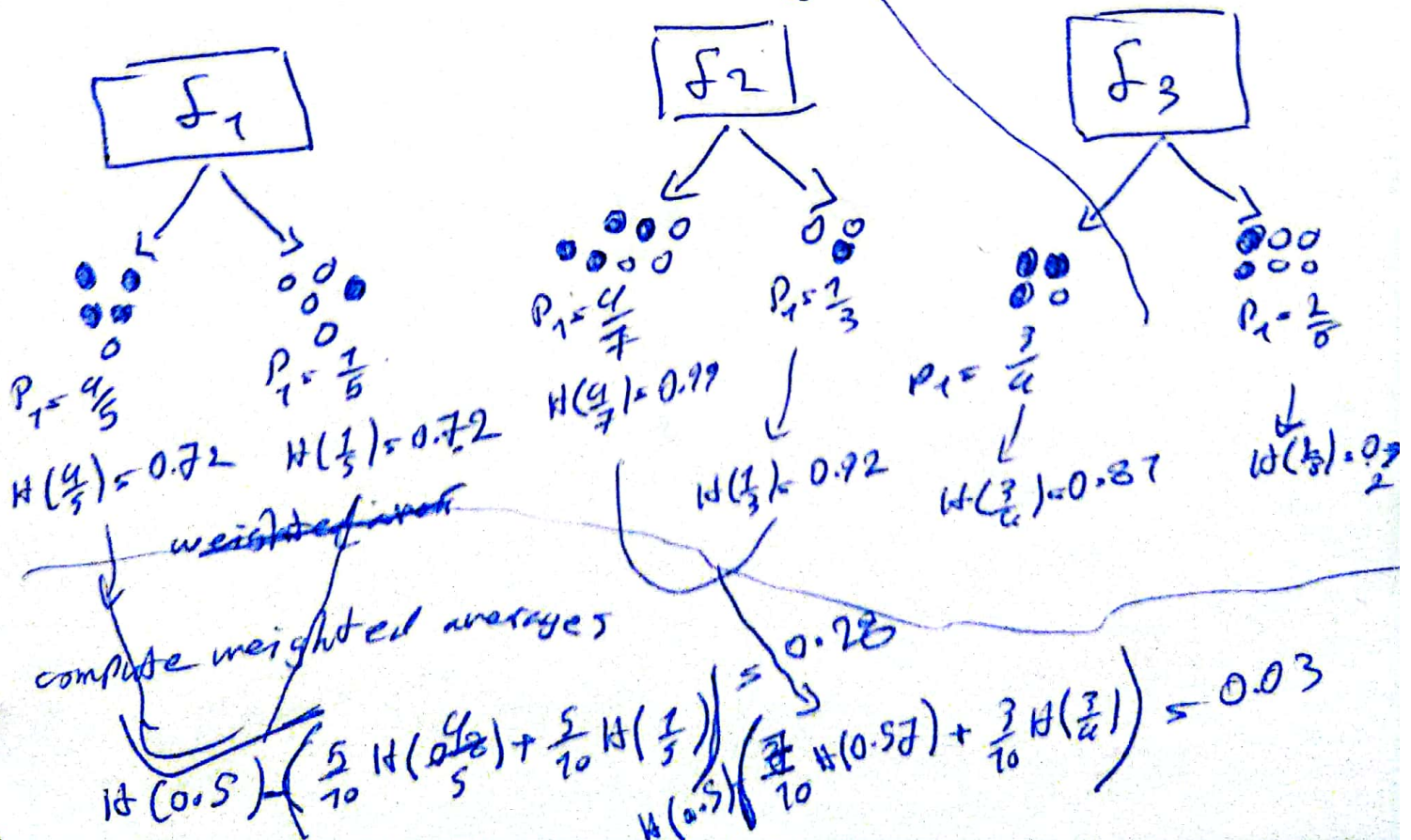$P_1 = \frac{2}{6}$      $H(P_1) = 0.92$

$P_1 = 0$      $H(P_1) = 0$

if $P_0 = 1 - P_1$

$$H(P_1) = -P_1 \log_2(P_1) - P_0 \log_2(P_0)$$

$$= -P_1 \log_2(P_1) - (1-P_1)\log_2(1-P_1)$$

## choosing a split

Information gain,

$$H(0.5) - \left( \frac{4}{10} H(\frac{3}{4}) + \frac{6}{10} H(\frac{2}{6}) \right) = 0.12$$



$f_1$

$P_1 = \frac{4}{5}$

$P_1 = \frac{1}{5}$

$H(\frac{4}{5}) = 0.72$    $H(\frac{1}{5}) = 0.72$

~~weighted~~

compute weighted averages

$f_2$

$P_1 = \frac{4}{7}$    $P_1 = \frac{1}{3}$

$H(\frac{4}{7}) = 0.99$

$H(\frac{1}{3}) = 0.92$

$f_3$

$P_1 = \frac{3}{4}$    $P_1 = \frac{1}{6}$

$H(\frac{3}{4}) = 0.87$    $H(\frac{1}{6}) = 0.65$

$$H(0.5) - \left( \frac{5}{10} H(\frac{4}{5}) + \frac{5}{10} H(\frac{1}{5}) \right)$$

$$H(0.5) - \left( \frac{7}{10} H(0.57) + \frac{3}{10} H(\frac{3}{4}) \right) = 0.03$$

$= 0.28$

after computing ~~weighted~~ weighted average

for all features (choose the feature
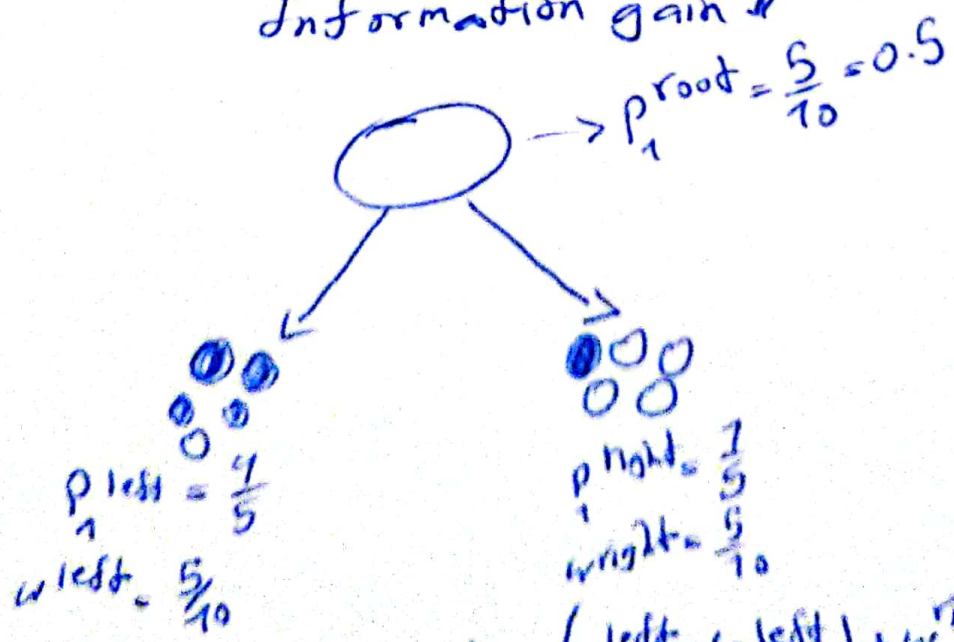
which has the "lowest" weighted average

to split on.

the result of the ~~ave~~ weighted average

* called information gain.

~~Information gain~~

we split the information on a feature which

has the "Bigest" information gain or "lowest"

weighted average.

---

"General formula of computing

Information gain"

if,

$$p_1^{root} = \frac{5}{10} = 0.5$$

$$p_1^{left} = \frac{4}{5}$$

$$w^{left} = \frac{5}{10}$$

$$p_1^{right} = \frac{1}{5}$$

$$w^{right} = \frac{5}{10}$$

$$\text{Information gain} = H(p_1^{root}) - \left( w^{left} \cdot H(p_1^{left}) + w^{right} \cdot H(p_1^{right}) \right)$$

# Decision tree learning (recursive splitting)

**1** start with all examples at the root node

**2** calculate information gain for all possible features, and pick the one with the highest information gain.

**3** split dataset according to selected feature, and create left and right branches of the tree.

**4** keep repeating splitting process until stopping criteria is met:

stopping criteria

- when a node is 100% one class.
- when splitting a node will result in the tree exceeding a maximum depth.
- information gain from additional splits is less than threshold.
- when the number of examples in a node is below a threshold.

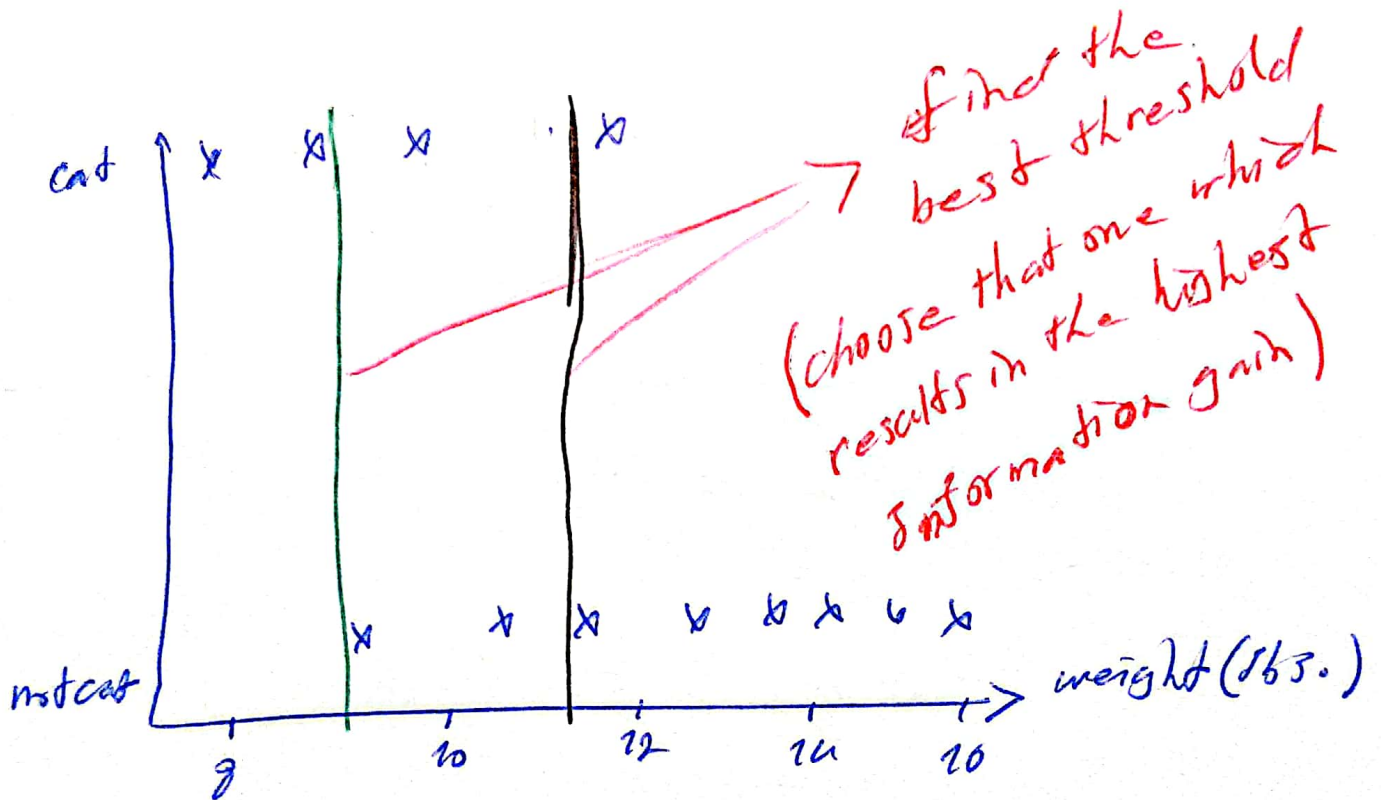using one-hot encoding of categorical features

exps of one-hot encoding

| Ear shape | pointy ear | floppy ear | oral ear |
|-----------|------------|------------|----------|
| pointy | 1 | 0 | 0 |
| oral | 0 | 0 | 1 |
| oral | 0 | 0 | 1 |
| floppy | 0 | 1 | 0 |
| pointy | 1 | 0 | 0 |
| oral | 0 | 0 | 1 |
| floppy | 0 | 1 | 0 |
| floppy | 0 | 1 | 0 |

one-Hot encoding : If a categorical feature can take on k values, create k binary features (0 or 1 variables).

Using one-hot encoding, categorical features can be used to train other models like neural networks or logistic regression.

# "Decision tree for continuous valued features"



find the best threshold
(choose that one which results in the highest information gain)

cat — weight (lbs.)

$$H(0.5) - \left(\frac{2}{12} H\left(\frac{2}{2}\right) + \frac{10}{12} H\left(\frac{2}{10}\right)\right) \approx 0.3$$

$$H(0.5) - \left(\frac{5}{12} H\left(\frac{3}{5}\right) + \frac{7}{12} H\left(\frac{3}{7}\right)\right) \approx 0.4$$

---

## regression tree

### choosing a split
based on variance of targets in each node.

eg, L variance = 1.47   R variance = 21.87 → information gain

$$\left(\text{overall var}\right) - \left(\frac{5}{10} \times 1.47 + \frac{5}{10} \times 21.87\right)$$

w left $\frac{5}{10}$   w right $\frac{5}{10}$

(Tree ensemble)

"Sampling with replacement"

choose a subset of original dataset with replace-
ment (you may choose ~~each~~ examples one by one ~~and~~
with replacement) and train a tree with that
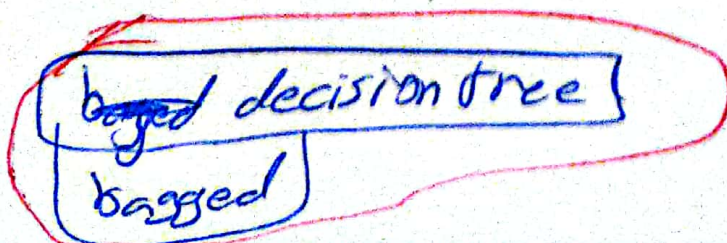subset.

"random forest algorithm"

(bagged decision Tree)

Given training set of size $m$:

For $b = 1$ to $B$:

T=use sampling with replacement to create a
new training set of size $m$.

$(T.size = m)$

Train a tree with T ~~size(m)~~ -

trained
~~x~~
at the end there will be $b$ different tree.

bagged decision tree
bagged

# random forest

At each node, when choosing a feature to use to split, if n features are available, Pick a random subset of $K < n$ features and allow the algorithm to only choose from that subset of features.

$$K = \sqrt{n}$$

← a typical choice for $K$ when $n$ is large.

---

# "XG Boost"

(implementation)

(extreme Gradient Boosting)

use decision tree

- The most common way to algorithm.

Given training set of size m:

- for b = 1 to B:
  - use sampling with replacement to create a new training set of size m (But instead of picking from all examples with equal ($1/m$) probability, make it more likely to pick misclassified examples from previously trained trees. Train a tree on the new dataset.

## Classification

```
from xgboost import XGBClassifier, XGBRegressor

model = XGBClassifier()
model.fit( X_train, Y_train)
y_pred = model.predict(X_test)
```

## regression

```
model = XGBRegressor()
model.fit(X_train, Y_Train)
y_pred = model.predict(X_test)
```

---

## "when to use decision Trees"

### "Decision Trees" and tree ensembles

- works well on tabular (structured) data.
- not reconended for unstructured data (images, audio, text).
- decision trees are fast.
- small decision trees may be human interpretable.

### "neural networks"

- works well on all types of data, including tabular (structured) and unstructured data.

- maybe slower than a decision tree.

- NN works with transfer learning.

- when building a system of multiple models working together, it might be easier to string together multiple neural networks.

---

2024 . 02 .15

L. Samani