Random forest $\longrightarrow$ Ensemble Technique

$\quad\quad\quad\quad \hookrightarrow$ Bagging (Bootstrapped

$\quad\quad\quad\quad\quad\quad\quad$ Aggregation)



Predictors (e.g., classifiers)

Training

$M_1$ $\quad$ $M_2$ $\quad$ $M_3$ $\quad$ $M_4$

Random sampling (with replacement = bootstrap)

Training set

$\longrightarrow$ Decision Tree

$\quad\quad \hookrightarrow$ Random forest

**Bagging Classifier ( algorithm)**

$\quad\quad\quad \hookrightarrow$ Real

Advantage (RF) $\quad$ time $\quad$ $M_1$ $\quad$ Decision Tree

$\quad\quad\quad\quad\quad\quad$ industry project $\quad\quad\quad\downarrow$

$\quad\quad \hookrightarrow$ avoid overfitting ① $\quad\quad$ Overfitting

$\quad\quad\quad\quad\quad$ (DT) $\quad\quad\quad\quad\quad \hookrightarrow$ Pruning

$\quad\quad \hookrightarrow$ increase accuracy ②

classification

Regression

n_estimators = 5 , subset of data

$\quad\quad\quad$ (randomly) $\rightarrow$ by replacement

$M_1$ — $DT_1$

$M_2$ — $DT_2$

$M_3$ — $DT_3$

$M_4$ — $DT_4$

$M_5$ — $DT_5$

$0 \quad$ 2cr

$0 \quad$ 1cr $\quad$ Majority $\Rightarrow 0$

$0 \quad$ 3cr $\quad$ voting

$1 \quad$ 2.5cr

$0 \quad$ 4cr

Regression

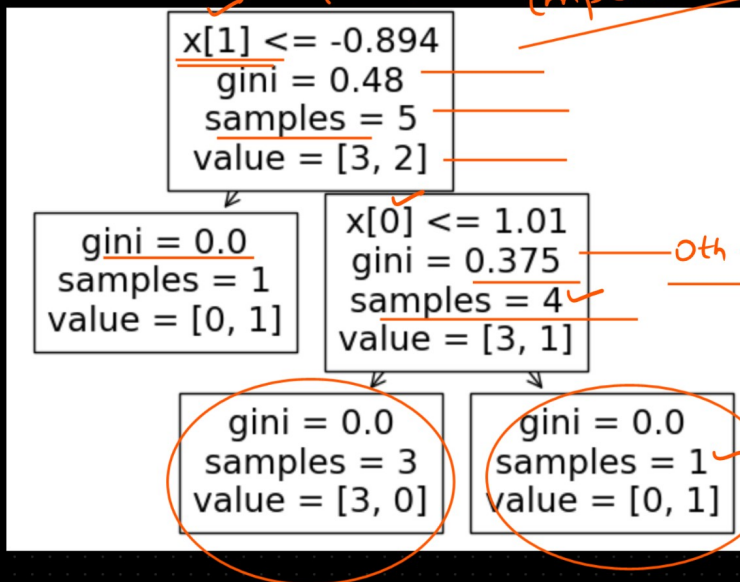$\quad \hookrightarrow$ average of all the given prediction

$$(2cr + 1cr + 3cr + 2.5cr + 4cr)/5$$

$$n_i \implies \frac{N\_t}{N}\left(\text{impurity} - \left(\frac{N\_t\_r}{N\_t} * \text{right impurity}\right) - \right.$$

right subtree

$$\left(\frac{N\_t\_l}{N\_t} * \text{left impurity}\right)$$

left subtree

1st feature importance



x[1] <= -0.894
gini = 0.48
samples = 5
value = [3, 2]

gini = 0.0
samples = 1
value = [0, 1]

x[0] <= 1.01
gini = 0.375
samples = 4
value = [3, 1]

gini = 0.0
samples = 3
value = [3, 0]

gini = 0.0
samples = 1
value = [0, 1]

0th feature importance

$$n_1 \implies$$

$$\frac{5}{5}\left(0.48 - \frac{4}{5} * 0.375 - \frac{1}{5} * 0\right)$$

$$= 0.18$$

$$n_0 \rightarrow \frac{4}{5}\left(0.375 - \frac{1}{4} * 0 - \frac{3}{4} * 0\right) = 0.30$$

$$\frac{4}{5} * 0.375$$
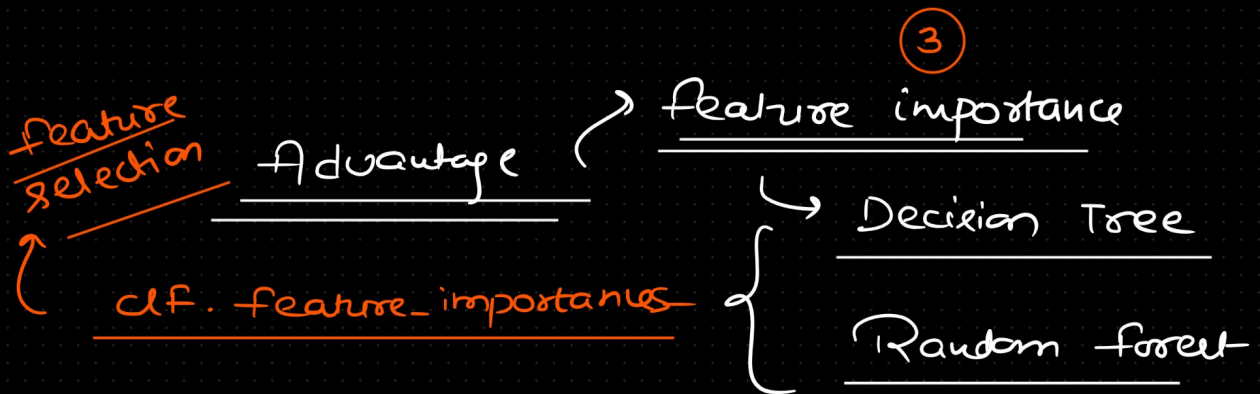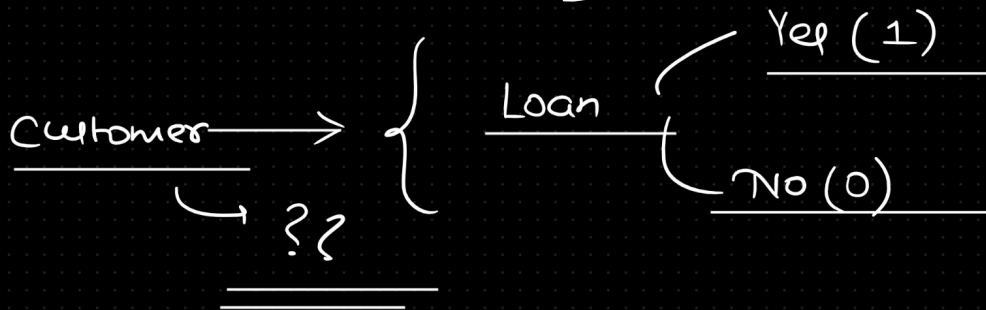
Normalization

$$0.48 +$$

$$= 1 \quad \longleftarrow$$

$$f_0 \rightarrow \frac{0.30}{0.30 + 0.18} = 0.625 ✓$$

$$f_1 \rightarrow \frac{0.18}{0.18 + 0.30} = 0.375$$

## Disadvantage

① More training time

② Interpretability is bit difficult

Customer → { Loan { Yes (1)
                      No (0)

↳ ??

③

feature
selection ___ Advantage ⟨ → feature importance

↳ Decision Tree

clf. feature_importances { Random forest

| $f_1$ | ~~$f_2$~~ | ~~$f_3$~~ | $f_4$ | $f_5$ | target (Loan) |
|---|---|---|---|---|---|
| | | | | | 0 |
| | | | | | 1 |
| | | | | | 0 |
| | | | | | 0 |
| | | | | | 1 |