

for classification and regression

non parameters (like Knn)

→ the tree will be learned from data

non linear classifier

→ the prediction will be taken in the leaf

→ Rules:

* if features are categorical:

① every root to leafs

1 - no need for one-hot-encoding

② every root to positive leafs

2 - maximum deep = # of features

→ no need for one-hot-encoding

→ no need for feature normalization

→ can overfitting easily (especially with high # of features)

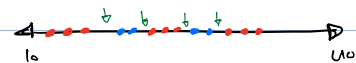
→ for multiclass ⇒ no need for any change

* if features are continuous

→ convert them to binary with thresholds

→ What threshold? the one with less error

between two classes

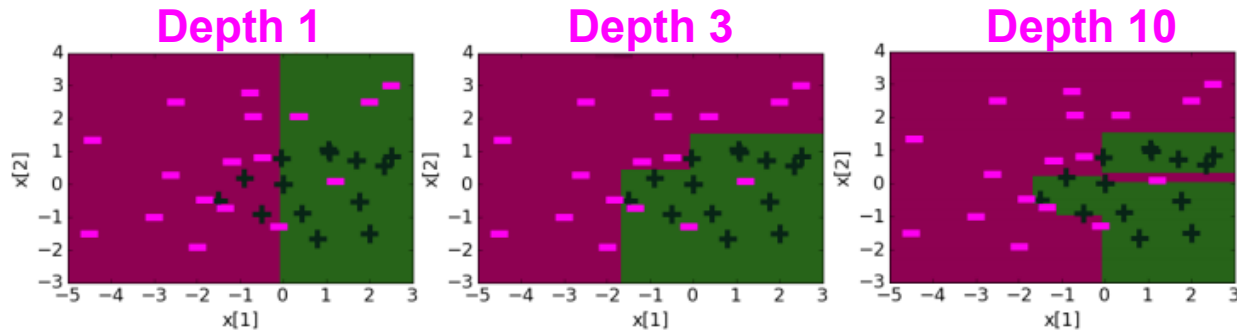


→ you can use the same feature again with another threshold → tend to overfitting

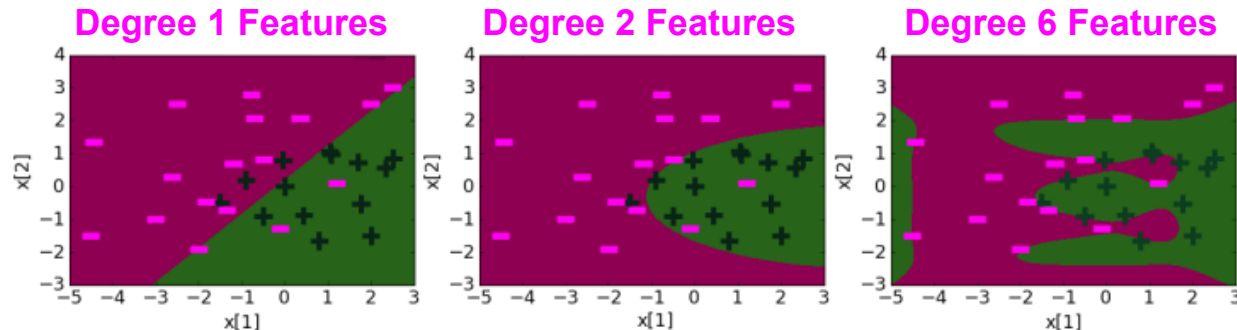
→ no need for normalization

Decision Trees Versus Logistic Regression (6)

Decision Tree *decision boundaries are axis align*



Logistic Regression *not axis align*



✖ How to generate a tree??

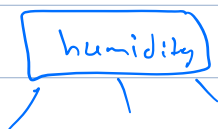
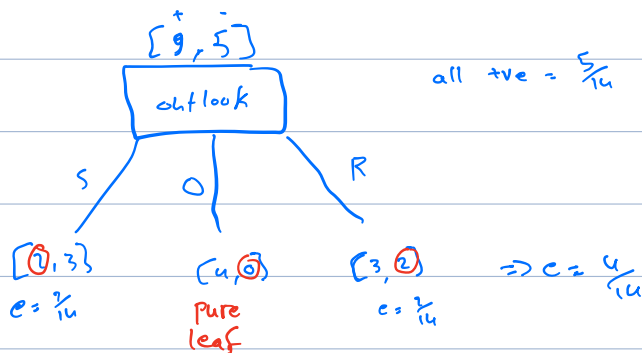
① trial and error \Rightarrow not feasible

② heuristic

select the best (forward selection)

what is the best?? error based

Slide 13



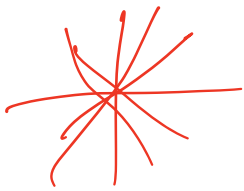
then select the least error rate \Rightarrow then go more deep unless it's pure

✖ when to stop:-

① pure nodes

② we already split on all features (if they are categorical)

③ set the depth of tree hyperparameter



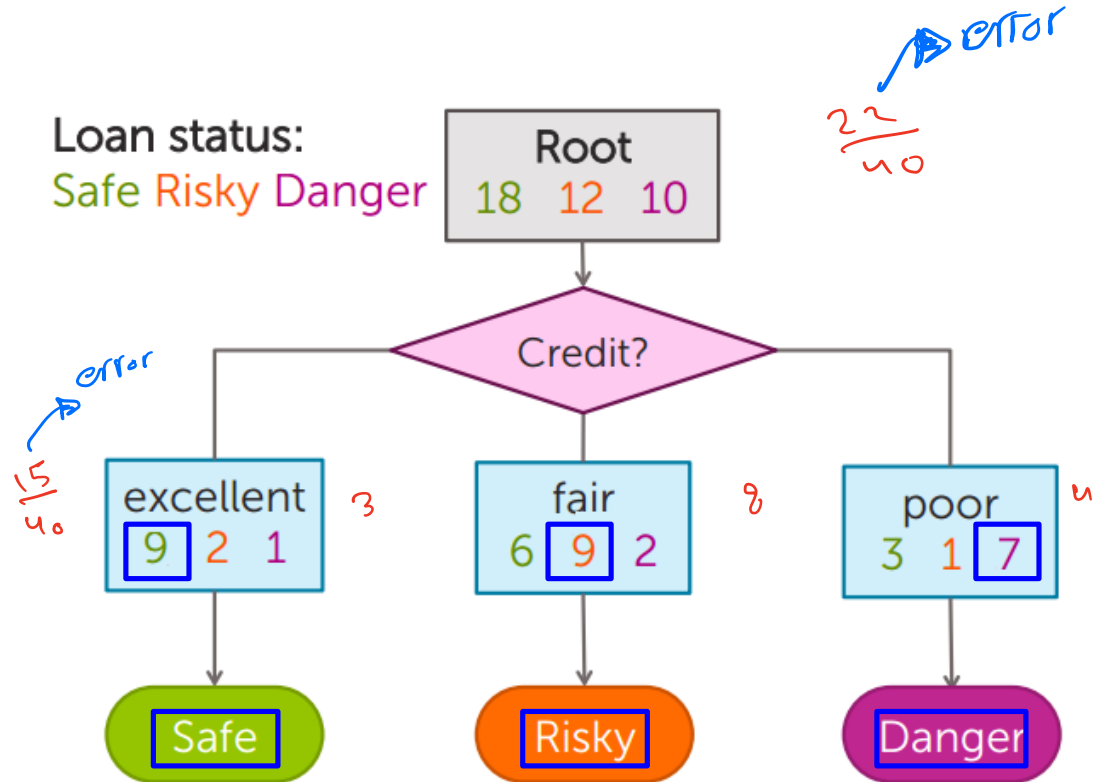
Multiclass Decision Stump

N = 40,
1 feature,
3 classes

Credit	y
excellent	safe
fair	risky
fair	safe
poor	danger
excellent	risky
fair	safe
poor	danger
poor	safe
fair	safe
...	...



Loan status:
 Safe Risky Danger





* of leaf node is indicator of complexity

How to avoid overfitting??

① early stopping

- ① if you split and the error doesn't improve \Rightarrow stop
- ② stop if you have small * of samples that you split in the node \rightarrow hyperparameter
- ③ early stopping at certain depth (longest path) \rightarrow hyperparameter

②

pruning \Rightarrow cutting from bottom until the J starts increasing

$$J = \text{error rate} + \lambda \text{ model complexity} \rightarrow \text{* of leaf nodes}$$

λ hyperparameter

$\lambda = \infty$ decision in root (underfit)

$\lambda = 0$ standard (overfit)

Simpler trees are better!!!

more chance to generalize better



in Regression:-

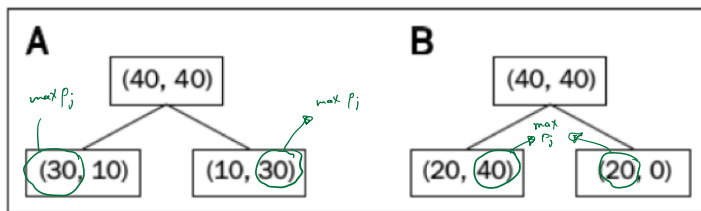
calculate SRR

✖ Different splitting measures

① $\text{classification error} = 1 - \max\{p_j\}$ what we used above probability of class j

② $\text{Entropy} = -\sum p_j \log_2 p_j$

③ $\text{Gini} = 1 - \sum p_j^2$



in error \Rightarrow they are same $e = \frac{20}{40}$

but B is better since has pure

\Rightarrow So error rate is not good enough

✖ General Formula:-

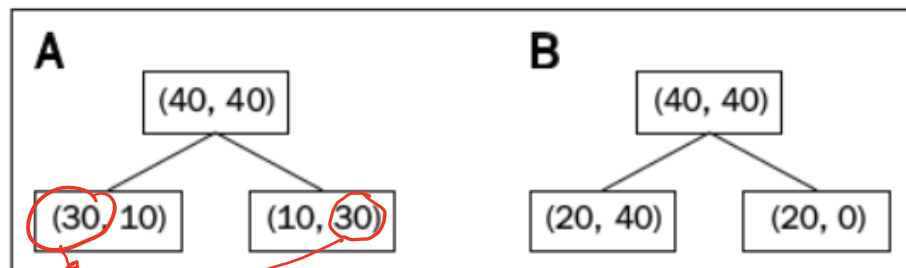
$$IG(D_p) = I(D_p) - \frac{N_{\text{left}}}{N_p} I(D_{\text{left}}) - \frac{N_{\text{right}}}{N_p} I(D_{\text{right}})$$

Information Gain \downarrow Dataset of the parent

any impurity measure (error, entropy, Gini)

\Rightarrow more gain is better

Example-IG with Classification Error



$$\text{error} = 1 - \max p_j$$

$$I_E(D_p) = 1 - \frac{40}{80} = 1 - 0.5 = 0.5$$

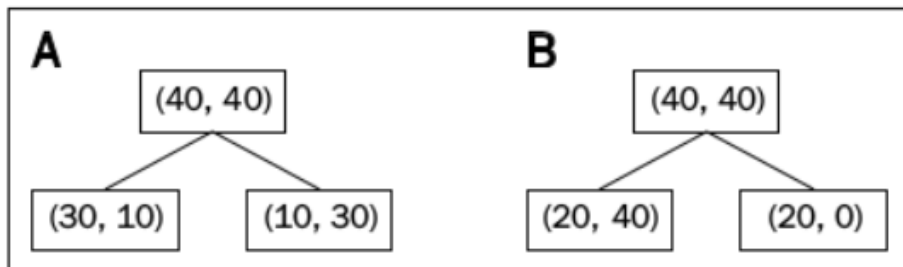
$$A : I_E(D_{\text{left}}) = 1 - \frac{30}{40} = 1 - \frac{3}{4} = 0.25$$

$$A : I_E(D_{\text{right}}) = 1 - \frac{30}{40} = 1 - \frac{3}{4} = 0.25$$

$$IG(D_p) = I(D_p) - \frac{N_{\text{left}}}{N_p} I(D_{\text{left}}) - \frac{N_{\text{right}}}{N_p} I(D_{\text{right}})$$

$$A : IG_E = 0.5 - \frac{40}{80} \times 0.25 - \frac{40}{80} \times 0.25 = 0.5 - 0.125 - 0.125 = 0.25$$

Example-IG with Classification Error



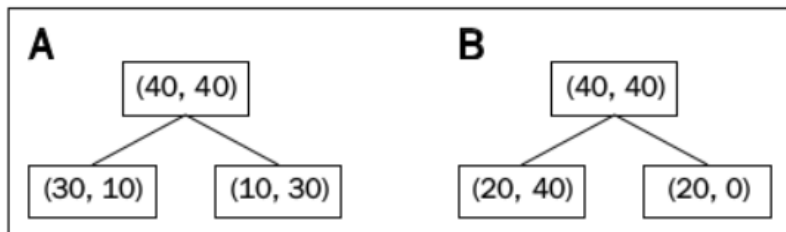
$$error = 1 - \max p_j$$

$$B : I_E(D_{left}) = 1 - \frac{40}{60} = 1 - \frac{2}{3} = \frac{1}{3}$$

$$B : I_E(D_{right}) = 1 - \frac{20}{20} = 1 - 1 = 0$$

$$B : IG_E = 0.5 - \frac{60}{80} \times \frac{1}{3} - \frac{20}{80} \times 0 = 0.5 - 0.25 - 0 = 0.25$$

Example-IG with Entropy



$$Entropy = - \sum_j p_j \log_2 p_j$$

$$I_H(D_p) = - (0.5 \log_2(0.5) + 0.5 \log_2(0.5)) = 1$$

$$A : I_H(D_{left}) = - \left(\frac{30}{40} \log_2 \left(\frac{30}{40} \right) + \frac{10}{40} \log_2 \left(\frac{10}{40} \right) \right) = 0.81$$

$$A : I_H(D_{right}) = - \left(\frac{10}{40} \log_2 \left(\frac{10}{40} \right) + \frac{30}{40} \log_2 \left(\frac{30}{40} \right) \right) = 0.81$$

$$A : IG_H = 1 - \frac{40}{80} \times 0.81 - \frac{40}{80} \times 0.81 = 0.19$$

Example-IG with Entropy



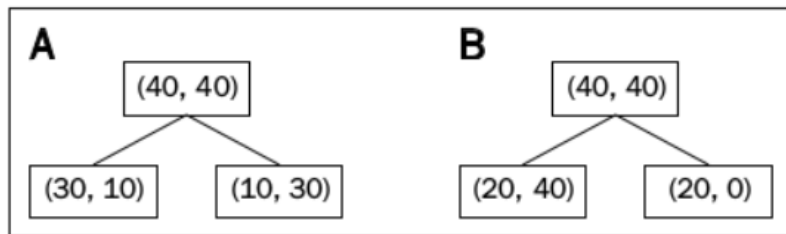
$$Entropy = - \sum_j p_j \log_2 p_j$$

$$B : I_H(D_{left}) = - \left(\frac{20}{60} \log_2 \left(\frac{20}{60} \right) + \frac{40}{60} \log_2 \left(\frac{40}{60} \right) \right) = 0.92$$

$$B : I_H(D_{right}) = - \left(\frac{20}{20} \log_2 \left(\frac{20}{20} \right) + 0 \right) = 0$$

$$B : IG_H = 1 - \frac{60}{80} \times 0.92 - \frac{20}{80} \times 0 = 0.31 \rightarrow \text{more gain} \Rightarrow \text{it's better}$$

Example-IG with Gini index



$$Gini = 1 - \sum_j p_j^2$$

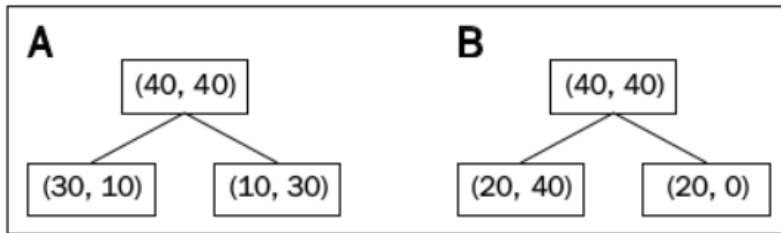
$$I_G(D_p) = 1 - \left(\left(\frac{40}{80} \right)^2 + \left(\frac{40}{80} \right)^2 \right) = 1 - (0.5^2 + 0.5^2) = 0.5$$

$$A : I_G(D_{left}) = 1 - \left(\left(\frac{30}{40} \right)^2 + \left(\frac{10}{40} \right)^2 \right) = 1 - \left(\frac{9}{16} + \frac{1}{16} \right) = \frac{3}{8} = 0.375$$

$$A : I_G(D_{right}) = 1 - \left(\left(\frac{10}{40} \right)^2 + \left(\frac{30}{40} \right)^2 \right) = 1 - \left(\frac{1}{16} + \frac{9}{16} \right) = \frac{3}{8} = 0.375$$

$$A : I_G = 0.5 - \frac{40}{80} \times 0.375 - \frac{40}{80} \times 0.375 = 0.125$$

Example-IG with Gini index



$$Gini = 1 - \sum_j p_j^2$$

$$B : I_G(D_{left}) = 1 - \left(\left(\frac{20}{60} \right)^2 + \left(\frac{40}{60} \right)^2 \right) = 1 - \left(\frac{9}{16} + \frac{1}{16} \right) = 1 - \frac{5}{9} = 0.44$$

$$B : I_G(D_{right}) = 1 - \left(\left(\frac{20}{20} \right)^2 + \left(\frac{0}{20} \right)^2 \right) = 1 - (1 + 0) = 1 - 1 = 0$$

$$B : I_G = 0.5 - \frac{60}{80} \times 0.44 - 0 = 0.5 - 0.33 = 0.17$$

- 13 its better