

Traditional Programming Data, Program \rightarrow output

Machine learning Data, output \rightarrow Program

Supervised learning \rightarrow training data labeled (train on output)

Unsupervised learning \rightarrow training data unlabeled need to find hidden structure

Semisupervised learning \rightarrow only part of input is labeled

Reinforcement learning \rightarrow Learner interacts to find optimal behavior

Deep learning \rightarrow Automatic learning without human interference

unsupervised

learning

Applications

Regression \rightarrow Predict a number

Classification \rightarrow Predict categories

when to split? max purity

- Clustering (group data)

- Dimensionality Reduction (compress data)

- Anomaly detection (find outliers)

init depth of trees to \rightarrow doesn't get too big/widely, prone to overfit

stop splitting when node 100% one class, when purity improvement below threshold, when number of examples below threshold

Entropy (impurity)

when = 1 most impure, = 0 most pure, p_i = fraction of ex. enc. cats?

$H(p_i) = \text{entropy}$, $p_0 = (1-p_i)$ fraction not in ex. not cats

Information gain

Entropy with multiple $\rightarrow H(p_i) = -\sum p_j \log_2(p_j)$

$H(P) = (P_0 H(1) + P_1 H(0))$

$H(p) = -p_1 \log_2(p_1) - p_0 \log_2(p_0) = -p_1 \log_2(p_1) - (1-p_1) \log_2(1-p_1)$

total - Node available

One hot encoding for k categories can have k values 0 or 1

ID3 Algorithm \rightarrow start root calc entropy and if for all start highest IG

use
in
ref
tree
for
node

Regression Tree $\rightarrow H_{\text{parent}} = (w_1 H(p_1) + w_0 H(p_0))$

Variance $s^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}$

Decision tree pros \rightarrow easy understand, little to no data, problems for noise bias, more flexible

Decision tree cons \rightarrow prone to overfit, high variance, costly

DT sensitive to small data changes, building multiple trees solves this issue

Sampling Replacement \rightarrow take random but replace before each grab

Bagged DT \rightarrow given dataset B_1 to B_n $\rightarrow B$ between $64 \rightarrow 128$

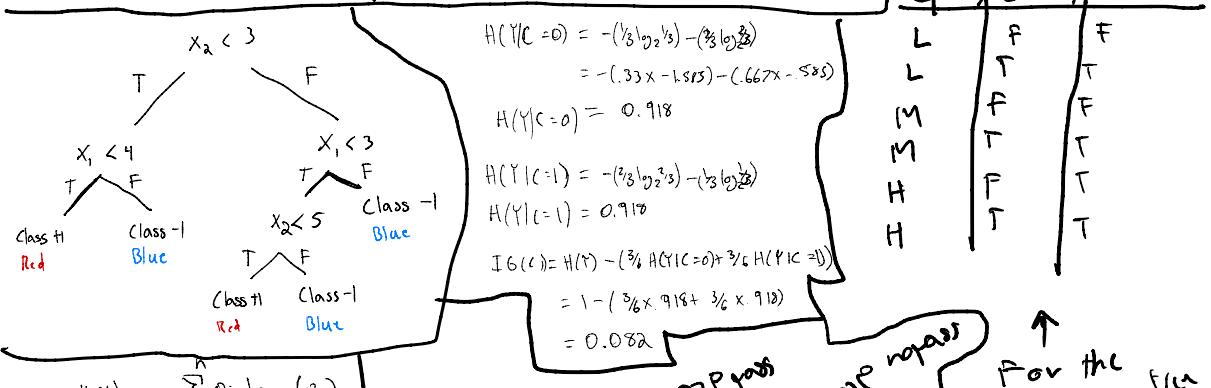
Random Forest Algo → at subset node pick random subset of kN if n features available, the algo picks max IG typically $k \approx n$ {Boosted Trees}

Same as flagged DT but most likely to pick examples previously tested

xgboost → open source of boosted trees, fast efficient, built in stop to overfitting

DT + Tree Ensembles → works on structured data, fast, not for unstructured data
works well with multiple models, works slower than DT } Good, Standard, Pass

Goal	Start	Pass
L	F	F
L	T	T
M	F	F
M	T	T
H	F	T
H	T	T



$$H(Y) = - \sum_{i=1}^n p_i \log_2(p_i)$$

$$H(4) = -(\frac{1}{6} \log_2 \frac{3}{6}) - (\frac{3}{6} \log_2 \frac{3}{6}) \\ = - (.5(-1)) - (.5(1))$$

$$H(Y) = 1$$

$$H(Y|A=0) = -(1/3 \log_2 1/3) - (2/3 \log_2 2/3)$$

$$= -(0.33 \times -1.885) - (0.667 \times -0.885)$$

$$H(Y|A=a) = 0.918$$

$$-1(Y|A=1) = -(z_3 \log z_3) - (z_3 \log \frac{1}{z_3})$$

$$f(Y|A=1) = 0.917$$

$$H(Y) = - \left(\frac{1}{6} \log_2 \frac{1}{6} \right) - \left(\frac{2}{6} \log_2 \frac{2}{6} \right)$$

$$H(x) = -(.0667x - 0.885) - (0.333x - 1.585)$$

$$H(Y) = 0.917$$

$$H(\gamma)_{\text{Studied}} = \text{icos} = 0 \quad \text{All Passed}$$

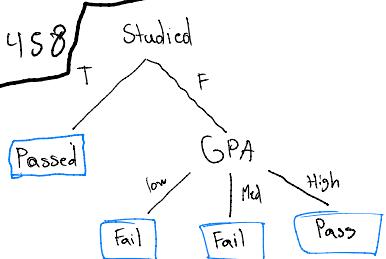
$$H(Y|S_{\text{Studied}} = N_0) = -(1/3 \log_2 1/3) - (2/3 \log_2 2/3)$$

$$= 0.918$$

$$I6(\text{Studied}) = H(Y) - \left(\frac{3}{4}H(Y| \text{Studied}=\text{False}) + \frac{1}{4}H(Y| \text{Studied}=\text{True}) \right)$$

$$= 0.917 \underbrace{(3/6 \times .918 + 3/6 \times 0)}$$

$$= 0.458$$



A	D	C	
0	-1	0	1
-1	0	-1	-1
0	0	0	2
-1	0	-1	2
0	-1	0	2
-1	0	1	2