

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

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)

unsupervised

learning

Applications

limit 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. are cats?

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

Information gain

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

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

$H(p_i) = -p_i \log_2(p_i) - p_0 \log_2(p_0) = -p_i \log_2(p_i) - (1-p_i) \log_2(1-p_i)$

total - Node added

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 $\rightarrow H_{\text{parent}} - (w_i H(p_i) + w_j H(p_j))$

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

IG is $H(\text{left}) - \text{the calc below}$

Bagged DT \rightarrow given dataset B_1 to B_n $\rightarrow B$ between $64 \rightarrow 128$ when splitting entropy do by feature and weight is num cat in total on

$p_i^L = \frac{1}{10} H(\frac{3}{7})$ $p_i^R = \frac{3}{10} H(\frac{5}{7})$

total cat

Random Forest Algo \Rightarrow at subset node pick random subset of $k \leq n$ (Sampling features available, the algo picks max IG typically) $\text{t} \in \mathcal{T}$ [Boosted Trees] Replacements # Iterations #

Same as bagged DT but most likely to pick examples previously visited

x6 Boost \rightarrow open source of boosted trees, fast efficient, built in stop to overfitting

DT + Tree Ensembles \rightarrow works on structured data, fast, not for unstructured data
works well with multiple models, works slower than DT

when splitting on leaf in regression tree use the min node as the root

$S^2 = \text{sum Variance of parent} \Leftrightarrow S^2 - \text{sum var child} \quad \text{big} = \text{good}$

Bayes Rule $p(y|x) = \frac{p(x|y)p(y)}{p(x)}$ | email representation given bag of words
Is if word shown assigned

Laplace Smoothing $p(F_2=b|y=1) = \frac{1 + (\# \text{examples } F_2=b \& y=1)}{\# \text{possible values } F_2 + (\text{examples } y=1)}$

$p(y=1) = \frac{\# \text{examples } y=1}{m}$

k-fold cross validation split into random folds, validation set to test better train

$p(y=1|x) = \frac{p(x|y=1)p(y=1)}{p(x)} = \frac{p(x|y=1)p(y=1)}{p(x|y=1)p(y=1) + p(x|y=0)p(y=0)}$

$p(A|B) = \frac{P(A \cap B)}{P(B)} = P(B) = P(B|A)P(A) + P(B|\bar{A})P(\bar{A})$ value is greater the

For classification get $P(X, 1+)$ and multiply and negative the one

$F(x_j|y=0) = \frac{1}{\sqrt{2\pi}\sigma_j} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)$ mean $\sigma = \text{std dev} = \frac{1}{n} \sum (x_i - \mu)$

estimate for continuous

law of total probability

$B = B_1 \cup B_2 \cup \dots \cup B_n$

$B_i \cap B_j = \emptyset$ for ever $i \neq j$

parameters estimated by data, define skill, needed for pred

Hyperparameter \rightarrow external can't be estimated
 \rightarrow process tuned, specified by practitioner

Add out cross validation \rightarrow train on S_{train} test on S_{val}

$P(B_i) > 0$ for $i=1 \dots n$ Report error on test only

we can say B_1, \dots, B_n portion \rightarrow then k fold train on all S except S_j

$P(A) = P(A|B_1)p(B_1) + \dots + P(A|B_n)p(B_n)$