

UNIT-4

* Supervised Learning

- training data is accompanied by labels.
- classification
- data is classified based on training set.

→ Classification

* predicts class labels (discrete/nominal)

→ Numeric Prediction:

* predicts unknown / missing values.

→ Applications:

- Credit/loan approval
- Diagnosis
- fraud detection
- Webpage categorization

Unsupervised learning

- class labels are unknown
- clustering
- with the help of measurements; classes or clusters are formed.

* Classification is a 2-step process.

① Model construction:

↳ a training set formed by assigning each tuple with a label.

② Model usage:

↳ checking for accuracy of the model using test set.

→ with the help of accuracy, we decide whether to use model or not.



* Decision Tree

→ Attributes are selected based on heuristic or statistical measure.

↳ Information gain.

Entropy = measure of uncertainty
$$= -\sum_{i=1}^m P_i \log(P_i)$$

* Attribute Selection measures:

- Information gain: multi-valued attributes
- Gain Ratio: unbalanced splits
- Gini index: multi-valued attributes;
large classes

* Overfitting:

↳ Too many branches, some may reflect anomalies due to noise.

* Approaches to avoid overfitting

① Prepruning: stop tree construction; if splitting a node; if measures falls below threshold.

② Postpruning: remove branches from a fully grown tree.

* Enhancements to Decision Tree Induction:

- continuous valued attributes
- handle missing attribute values
- attribute construction

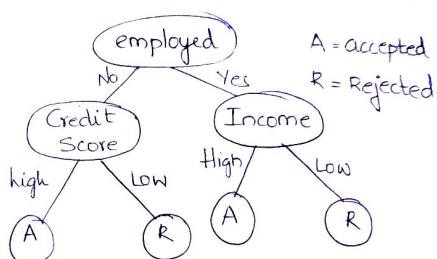
→ AVC list: (Attribute, value, class_label)

* BAT: Bootstrapped Optimistic Algorithm for Tree Construction

→ Statistical technique called bootstrapping to create several ~~smaller~~ smaller samples; each fits in memory.

→ Each subset is used to create a tree.

Ex:



Information Gain:

$$-\frac{P}{P+N} \log_2 \left(\frac{P}{P+N} \right) - \frac{N}{P+N} \log_2 \left(\frac{N}{P+N} \right)$$

Ex:

Age	Completion	Type	Projd
old	Yes	S/W	Down
old	No	S/W	Down
old	No	H/W	Down
mild	Yes	S/W	Down
mild	Yes	H/W	Down
mild	No	H/W	Up
mild	No	S/W	Up

new	Yes	s/w	up
new	No	H/W	up
new	No	s/w	up

→ Information Gain: for Target attribute = profit

$$IG = \frac{-P}{P+N} \log_2 \left(\frac{P}{P+N} \right) - \frac{N}{P+N} \log_2 \left(\frac{N}{P+N} \right)$$

$$P = \text{count(down)} = 5$$

$$N = \text{count(up)} = 5$$

$$= -\frac{5}{10} \log \left(\frac{5}{10} \right) - \frac{5}{10} \log \left(\frac{5}{10} \right)$$

$$= -\frac{1}{2}(-1) - \frac{1}{2}(-1) = 1$$

→ Entropy for remaining attributes:

$$E(A) = \sum \frac{P_i + N_i}{P+N} I(P_i N_i)$$

① Age:

	down	up
old	3	0
mid	2	2
new	0	3

$$P = \text{down count} =$$

$$N = \text{up count} =$$

$$IG(\text{old}) = \frac{-3}{3} \log \left(\frac{3}{3} \right) - \frac{0}{3} \log \left(\frac{0}{3} \right)$$

$$\text{Probability} = \frac{3}{10}$$

$$E(\text{old}) = 0$$

$$IG(\text{mid}) = -\frac{2}{4} \log \left(\frac{2}{4} \right) - \frac{2}{4} \log \left(\frac{2}{4} \right)$$

$$= -\frac{1}{2} \log_2(2^{-1}) - \frac{1}{2} \log_2(2^{-1})$$

$$= 1$$

$$\text{Probability} = \frac{4}{10} \Rightarrow E(\text{mid}) = \frac{2}{5}$$

$$IG(\text{new}) = 0$$

$$E(\text{new}) = 0$$

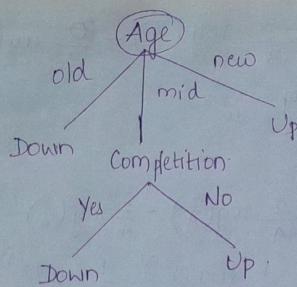
$$E(Age) = E(\text{old}) + E(\text{mid}) + E(\text{new}) = \underline{\underline{0.4}}$$

$$\text{Gain} = IG - E(A) = 0.6$$

$$\text{Gain}(Age) = 0.6 \rightarrow \text{root node} \quad \textcircled{1}$$

$$\text{Gain}(Type) = 0. \quad \textcircled{2}$$

$$\text{Gain}(Completion) = 0.124. \quad \textcircled{2}$$



∴ There is no use of type attribute.

* Bayesian Classification:

→ Bayes Theorem:

Total probability

$$\text{theorem: } P(B) = \sum_{i=1}^M P(B|A_i) P(A_i)$$

Bayes formula: $P(H|x) = \frac{P(x|H) P(H)}{P(x)}$

$$= P(x|H) \times P(H) / P(x)$$

$P(x)$ = probability that sample data is observed.

$P(H)$ = prior probability

$P(x|H)$ = probability of observing sample x , given hypothesis holds.

→ maximum posterior \Rightarrow

$$P(c_i|x)$$

$$P(c_i|x) = \frac{P(x|c_i) P(c_i)}{P(x)}$$

→ $P(x)$ is constant for all classes.

$$\Rightarrow P(c_i|x) = P(x|c_i) P(c_i)$$

EX: $P(\text{King}|\text{Face})$

$$= \frac{P(\text{Face}|\text{King}) \cdot P(\text{King})}{P(\text{Face})}$$

$$= \frac{\frac{1}{12} \times \frac{4}{52}}{\frac{1}{52}} = \frac{\frac{1}{12}}{\frac{3}{52}} = \frac{1}{3}$$

→ Naive Bayes Classifier:

$$P(x|c_i) = \prod_{k=1}^n P(x_k|c_i) = P(x_1|c_i) * P(x_2|c_i) * \dots * P(x_n|c_i)$$

* gaussian distribution:

$$g(x, \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Ex:

Fruit	Yellow	Sweet	Long	Total
Orange	350	450	0	650
Banana	400	300	350	1050
Others	50	100	50	200
Total	800	850	400	1200

$$P(\text{Yellow}|\text{Orange}) = \frac{P(\text{Orange}|\text{Yellow}) P(\text{Yellow})}{P(\text{orange})}$$

$$= \frac{\frac{350}{800} \times \frac{800}{1200}}{\frac{650}{1200}} = 0.53.$$

$$P(\text{Sweet}|\text{Orange}) = \frac{P(\text{Orange}|\text{Sweet}) P(\text{Sweet})}{P(\text{orange})}$$

$$= \frac{\frac{450}{850} \times \frac{850}{1200}}{\frac{650}{1200}} = 0.69.$$

$$P(\text{Long}|\text{Orange}) = 0$$

$$P(\text{Fruit}|\text{orange}) = P(Y|O) * P(S|O) * P(L|O)$$

$$= 0$$

$$P(\text{Fruit}|\text{Banana}) = P(Y|B) * P(S|B) * P(L|B)$$

$$= 0.65$$

$$P(F|\text{others}) = 0.072$$

Given: Fruit = {Yellow, Sweet, Long}

$$\therefore P(F|B) = 0.65$$

$$P(F|O) = 0.072$$

$$\therefore P(F|B) = 0.65$$

* Advantages:

- Easy to implement
- Good results obtained in most of the cases.

* Disadvantages:

- loss of accuracy
- dependencies among variables.

* Rule-based classification:

→ using IF-THEN rules:

R: if age = youth and student = yes

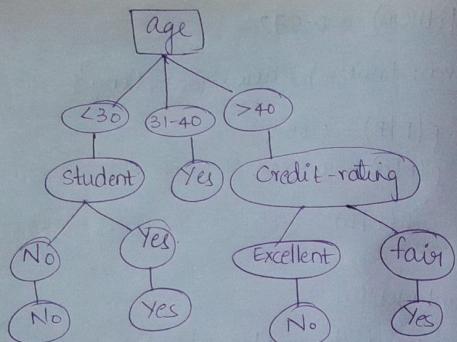
then buys_computer = yes

* Rule extraction from a decision tree:

→ One rule is created for each path from the root to a leaf.

→ Leaf holds the class label.

Ex:



if age = young and student = no
then buys_computer = no

if age = old and credit_rating = fair
then buys_computer = yes.

Rule Induction:

Sequential covering method:

Steps:

- Rules are learned one at a time.
- Each time a rule is learned, tuples covered are removed.
- Repeat until quality of rule is below threshold

* Model Evaluation and Selection:

→ Accuracy is assessed using validation test set.

→ Methods to estimate classifier's accuracy:

- Holdout method, random subsampling
- Cross-validation
- Bootstrap.

* Confusion Matrix:

Actual / Predicted	C_1	$\neg C_1$	
C_1	True Positives	False Neg.	P
$\neg C_1$	False Positives	True Neg.	N
	P'	N'	All

Metrics:

$$\text{Accuracy} = \frac{TP + TN}{\text{All}}$$

$$\text{Error Rate} = \frac{FP + FN}{\text{All}}$$

$$\text{Sensitivity} = \frac{TP}{P}$$

$$\text{Specificity} = \frac{TN}{N}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

$$F\text{-score} = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

* Evaluating Classifier Accuracy:

→ Holdout method:

Given data is partitioned into 2 sets.
 → Training set $\Rightarrow 2/3$ for construction
 → Test set $\Rightarrow 1/3$ for accuracy estimation

Random sampling:

Repeat holdout k times
 accuracy = avg of accuracies obtained

→ Cross-validation:

→ k-fold where $k=10$.

→ Randomly partition the data into k mutually exclusive subsets
 each of equal size.

* t-test:

$$t = \frac{\bar{e\pi}(M_1) - \bar{e\pi}(M_2)}{\sqrt{\text{Var}(M_1 - M_2) / k}}$$

$$\text{Var}(M_1 - M_2) = \frac{1}{2} \sum_{i=1}^k [\bar{e\pi}(M_1) - \bar{e\pi}(M_2)]^2 - (\bar{\bar{e\pi}}(M_1) - \bar{\bar{e\pi}}(M_2))^2$$

→ non-paired t-test

$$\text{Var}(M_1 - M_2) = \sqrt{\frac{\text{Var}(M_1)}{k_1} + \frac{\text{Var}(M_2)}{k_2}}$$

$M_1, M_2 \rightarrow$ Null hypothesis

→ select a significance level (sig)

$$\text{confidence limit} = \frac{\text{sig}}{2} = z$$

→ $t > z$ or $t < -z \Rightarrow$ Reject null hypothesis
 ⇒ statistically significant diff.

else: any diff. is chance.

→ Receiver Operating Characteristics (ROC) curve
 used for visual comparison of classification models.

* Issues affecting model selection:

- ① Accuracy
- ② Speed
- ③ Robustness
- ④ Scalability
- ⑤ Interpretability

* Techniques to improve classification accuracy.

→ Ensemble Methods:

- * Use a combination of models.
- * combine a series of k learned models to create a improved M^* model.

→ Popular methods:

- Bagging
- Boosting
- Ensemble

* Bagging: averaging the prediction over a collection of classifiers

→ Analogy: diagnosis based on majority vote.

- Training
- Classification
- Prediction

→ Accuracy:

- Often significantly better than single classifier.
- for noise data.
- improved accuracy in prediction

* Boosting: weighted vote with a collection of classifiers

Steps:

- weights are assigned to each training tuple.
- A series of k classifiers is iteratively learned.
- After a classifier M_i is learned, weights are updated.
- final M^* combines the votes of each classifier.
- * greater accuracy compared to bagging.
- * But risks overfitting the model to misclassify data.

* Classification of class - Imbalanced data

Methods:

- Oversampling: resampling of data from +ve class.
- Undersampling: eliminate tuples from -ve class.
- Threshold moving: moves decision threshold; so that rare class tuples are easier to classify.