

Gini Index : $1 - \sum_{i=0}^{C-1} p_i(t)^2$ * $p_i(t)$ is probability of class i at node t , C is total # of classes.

Entropy = $-\sum_{i=0}^{C-1} p_i(t) \log_2(p_i(t))$ + Gain = $P - M$, $P \rightarrow$ impurity before split, $M \rightarrow$ weighted impurity after split

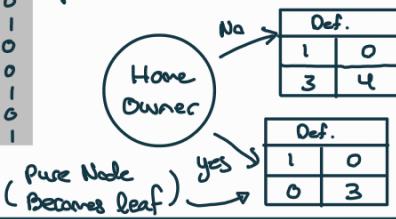
Classification error: $1 - \max[p_i(t)]$ * $\log_2(x) = \ln(x)/\ln(2)$

Hunt's Algorithm with Entropy Impurity:

ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

D Start with single classification node:

② Split on some attribute:



③ Calculate Weighted Impurity for k children

$$M = \sum_{i=0}^{k-1} \frac{n_i}{n} \text{Entropy}(i) = \frac{7}{10} \text{Entropy}(No) + \frac{3}{10} \text{Entropy}(Yes)$$

$$= \frac{7}{10} \left(-\left(\frac{3}{7} \log_2 \frac{3}{7} + \frac{4}{7} \log_2 \frac{4}{7} \right) \right) + \frac{3}{10} \left(-\left(\frac{0}{3} \log_2 \frac{0}{3} + \frac{3}{3} \log_2 \frac{3}{3} \right) \right)$$

$$= \frac{7}{10} (-(-0.5238 - 0.56134)) = 0.3596 = M$$

by convention

④ Next calculate Gain = $P - M = 0.88 - 0.3596 = 0.5202$

↳ Calculate this metric for each attribute and split on feature with highest gain.

↳ Repeat process until desired model fit

Impurity measures prefer many pure child nodes.

• Gain Ratio:

$$\text{Gain Ratio} = \frac{\text{Gain}_{\text{split}}}{\text{Split Info}}$$

$$\text{Split Info} = -\sum_{i=1}^k \frac{n_i}{n} \log_2 \frac{n_i}{n}$$

Parent Node, p is split into k partitions (children)

n_i is number of records in child node i

• Adjusts Information Gain by the entropy of the partitioning (Split Info).

• Higher entropy partitioning (large number of small partitions) is penalized!

CarType			
	Family	Sports	Luxury
C1	1	8	1
C2	3	0	7
Gini			0.163

SplitINFO = 1.52

CarType		
	(Sports, Luxury)	{Family}
C1	9	1
C2	7	3
Gini		
0.468		

SplitINFO = 0.72

CarType		
	{Sports}	{Family, Luxury}
C1	8	2
C2	0	10
Gini		
0.167		

SplitINFO = 0.97

$$-(16/20) \log_2 (16/20) - (4/20) \log_2 (4/20)$$

Sorted Values		Split Positions	Actual Class	No	No	No	Yes	Yes	Yes	No	No	No	No	No	No
0	1	60	No	0	3	0	3	0	3	1	2	2	1	3	0
1	2	65	No	0	3	0	3	0	3	1	2	2	1	3	0
2	3	72	No	0	3	0	3	0	3	1	2	2	1	3	0
3	4	80	No	0	3	0	3	0	3	1	2	2	1	3	0
4	5	87	No	0	3	0	3	0	3	1	2	2	1	3	0
5	6	92	No	0	3	0	3	0	3	1	2	2	1	3	0
6	7	97	No	0	3	0	3	0	3	1	2	2	1	3	0
7	8	102	No	0	3	0	3	0	3	1	2	2	1	3	0
8	9	110	No	0	3	0	3	0	3	1	2	2	1	3	0
9	10	122	No	0	3	0	3	0	3	1	2	2	1	3	0
10	11	132	No	0	3	0	3	0	3	1	2	2	1	3	0
11	12	142	No	0	3	0	3	0	3	1	2	2	1	3	0
12	13	152	No	0	3	0	3	0	3	1	2	2	1	3	0
13	14	162	No	0	3	0	3	0	3	1	2	2	1	3	0
14	15	172	No	0	3	0	3	0	3	1	2	2	1	3	0
15	16	182	No	0	3	0	3	0	3	1	2	2	1	3	0
16	17	192	No	0	3	0	3	0	3	1	2	2	1	3	0
17	18	202	No	0	3	0	3	0	3	1	2	2	1	3	0
18	19	212	No	0	3	0	3	0	3	1	2	2	1	3	0
19	20	222	No	0	3	0	3	0	3	1	2	2	1	3	0

Decision Tree Advantages :- Inexpensive, Fast, interpretable (small trees), robust to noise

Disadvantages : - space of possible trees very large & greedy approach doesn't always work

- inter-attribute interactions ignored, each decision boundary involves only 1 attribute

Metrics for Performance Evaluation

Confusion matrix

		Predicted Class	
		Class=Yes	Class=No
Actual Class	Class=Yes	a (TP)	b (FN)
	Class=No	c (FP)	d (TN)

$$\text{Accuracy (most common)} : (TP+TN) / (TP+TN+FP+FN)$$

$$\text{Precision (accuracy of positive predictions)} = TP / (TP+FP)$$

$$\text{Recall (proportion of true positives model could catch)} = TP / (TP+FN)$$

$$\text{F-measure (?) : } 2TP / (2TP + FN + FP) \rightarrow \text{predicts class 0 ; accuracy = 99.9% but this}$$

Accuracy limitations: Class 0 = 99.9%, class 1 = 0 → If model always is Misleading since model never predicts class 1!

		Predicted Class	
		Class=Yes	Class=No
Actual Class	Class=Yes	a (TP)	b (FN)
	Class=No	c (FP)	d (TN)

Accuracy = 80%

Cost = 3910

Accuracy = 90%

Cost = 4255

Model Evaluation:

↳ Holdout: Reserve k% for training, (100-k)% testing

↳ random subsampling: repeated holdout

↳ Cross-Validation: Train on k-1 partitions, test on remaining partition

Final performance is always Average of all runs

underfitting is when model is too simple, high train & test error

overfitting is when model is too complex, does not generalize well

↳ high test error, low train error

Increasing training set size, reducing model complexity and reducing # of features are all ways to deal with overfitting

Model Selection : Use small validation set to estimate generalization error.

↳ choose simpler of two models with similar generalization

↳ We can use pessimistic error to estimate generalization error of decision tree taking error into account

↳ We can stop splitting early (Pre-pruning)

↳ if impurity measures don't improve

↳ if errgen(T) becomes acceptable

↳ Post-prune → grow tree to max size → trim nodes bottom-up

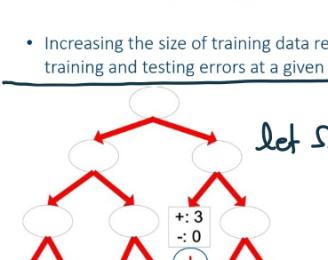
• err(T): error rate on all training records

• Ω : trade-off hyper-parameter (similar to α)

• Relative cost of adding a leaf node

• k : number of leaf nodes

• N_{train} : total number of training records



Nearest Neighbors Classifiers: To classify an unknown record ① compute distance to other records. ② Identify k-nearest neighbors. ③ Use class labels of nearest neighbors to classify unknown record, use a context dependent distance measure and take a weighted ($w = 1/d^2$) majority vote of the neighbors. If k is too small \rightarrow sensitive to noise, too large \rightarrow missclassification. May need feature standardization. Requires presence of all attributes, missing values need to be filled in, irrelevant features add noise, redundant features skew classification to certain class.

INTRO: Why data mining?: increased data collection speed and quantity, cheaper & more powerful compute, strong competitive pressure, helpful for scientists (automated analysis, hypothesis formation). **What is data mining?**: Non-trivial extraction of implicit, previously unknown and potentially useful information from data. **KDD process**: [input data] \rightarrow [preprocessing] \rightarrow [data mining] \rightarrow [postprocessing] \rightarrow [information]. **Why not classical data analysis?**: scalability, high dimensionality, heterogeneous and complex data, data ownership, distributed data, non-traditional analysis. **Origins**:

Draws on machine learning/AI, pattern recognition, statistics and database systems. **Data Mining Tasks**: predictive tasks (use known features to predict unknown features) \rightarrow classification (predict discrete features), regression (predict continuous features), outlier detection (predict anomalous records), descriptive tasks (find human-interpretable patterns to describe data) \rightarrow clustering (find groups of related objects), association rules (discover dependency rules between items in records of itemsets). **DATA: What is Data?**: Collection of objects (records) and their attributes (features).

Attribute Types: Nominal (distinctness $= \neq$, ID, eye color), Ordinal (order has meaning, $<$, $>$, rankings, height), Interval (differences have meaning, $+$, $-$, calendar dates, relative temperature scales), Ratio (ratios are meaningful, \div , \times , absolute scales of measure like kelvin or distance). **Data Characteristics**: Dimensionality (# of features), Sparsity (majority empty features), Resolution (patterns depend on scale), size (dataset size). **Types of Data**: Record (tabular, matrix \rightarrow instances become points in n-D space), Document \rightarrow each document becomes a word count vector, transaction \rightarrow records of ID'd itemsets), Graph (networks), Ordered (sequence \rightarrow braces used to indicate item(s) order, Spatio-Temporal \rightarrow space & time dimensions like in a regional time series temperature heatmap). **Data Quality Problems**: Noise & Outliers (a noise object is an extraneous signal independent of true data, a noisy feature is when the extraneous signal modifies original feature values, an outlier could be noise or could be analysis target \rightarrow ex. instance of fraud), Missing Values (incomplete data or inapplicable features \rightarrow eliminate records, ignore features or estimate values), Duplicate data (ex same person with multiple emails \rightarrow handling is context dependent). **Similarity & Dissimilarity**: similarity $\in [0, 1]$, dissimilarity $\in [0, \infty]$, **Minkowski** (generalized euclidean) $\rightarrow d(x, y) = \left(\sum_{k=1}^n (x_k - y_k)^r \right)^{1/r}$ standardization may be necessary, **Mahalanobis** (takes distribution into account) $\rightarrow M(x, y) = (x - y)^T \Sigma^{-1} (x - y)$ where Σ^{-1} is the feature covariance matrix, **Binary Vectors**: **Simple Matching** (# matches/# attributes) \rightarrow

$SMC(x, y) = (f_{11} + f_{00}) / (f_{11} + f_{00} + f_{01} + f_{10})$, **Jaccard** (# 11 matches/# nonzero attributes) $\rightarrow J(x, y) = (f_{11}) / (f_{11} + f_{01} + f_{10})$, **Cosine** $\rightarrow \cos(d_1, d_2) = \langle d_1, d_2 \rangle / \|d_1\| \|d_2\|$, **Extended Jaccard** (for continuous or count variables) $\rightarrow EJ(x, y) = (x \cdot y) / (\|x\|^2 + \|y\|^2 - x \cdot y)$. **Preprocessing: Aggregation** (combining multiple attributes/objects into a single attribute/object for data reduction, change of scale or stabilize (smooth) data), Sampling (randomly selected representative subset of the data considering sample size, with/without replacement or stratified sampling (sampling independently from partitions)), **Dimensionality Reduction** (curse of dimensionality: high dimensionality results in extremely sparse data making density/distance measures less meaningful) \rightarrow **PCA** (principal component analysis involves finding eigenvalues & eigenvectors of covariance matrix and then projecting data onto eigenvector with largest eigenvalue), **Feature Reduction** (remove redundant/irrelevant features), **Discretization** (continuous variable mapped to ordinal). **EXPLORATION**: Preliminary exploration to better understand data, select best future preprocessing/analysis tools. **EDA**: focus was mainly visualization using clustering and anomaly detection, today these are major areas of interest. **Summary Statistics: Spread** \rightarrow (mean, std, variance, range, iqr) $\rightarrow AAD(x) = \frac{1}{m} \sum_{i=1}^m |x_i - \bar{x}|$, $MAD(x) = \text{median}(\{|x_1 - \bar{x}|, \dots, |x_m - \bar{x}|\})$, $\text{interquartile range}(x) = x_{75\%} - x_{25\%}$, **frequency** (percentage occurrence), **mode** (most common occurrence), **Percentiles** (p^{th} percentile of an ordinal or continuous feature is the value strictly greater than p percent of all its other values), **Location** (mean is very sensitive to outliers, median or trimmed mean is better location measure). **Visualization**: converting data to visual or tabular formats for summarization, leveraging humans' well developed visual pattern recognition abilities, **Representation** (mapping of information to visual format), **Arrangement** (placement of visual elements within a display), **Selection for visualization** (elimination or de-emphasis of certain objects/features), **Histograms** (visualize distribution of continuous variables using "bins", 2-d for single feature distribution, 3-d for two feature joint distribution), **Box plots** (good for iqr, outliers, percentiles and median), **Scatter Plots** (features determine position, others can be used to change point shape, color, size, etc...), **Contour Plots** (partition spatial grid into regions where continuous feature values are similar, contour lines form region boundaries), **Data Matrix** (sort instances by class along one axis, arrange features along the other axis, calculate some metric for each instance-feature intersection), **Parallel Coordinates** (visualize each instance as a line colored by class, visualize movement of lines through sequential parallel axis observing patterns between classes), **Star Plot** (multiple features radiate from central point as individual axes, a line intersects the axes based on an objects feature values), **Chernoff Faces** (data features are mapped to face features). **On-Line Analytical Processing**: multidimensional array data representation (**Data Cube**) effective for aggregation operations, **1.** Identify and discretize selected features (features \rightarrow cube dimensions, unique feature values \rightarrow dimension indices), **2.** group-by unique tuples and sum target variable values or calculate group count, **3.** place target values in cube using the unique feature tuple values as cell indices, **example**: for sales, products and locations cube dimensions, select Product ID and Location, aggregating along sales to get total sales for each product at each location, many more combinations. **Slicing**: subset of cube specifying one value for one/more features, **Dicing**: subset cube specifying range of attribute values across dimensions, **Roll-up**: (zoom out), move up attribute hierarchy, a kind of aggregation (cities \rightarrow countries), **Drill-Down**: (zoom in), move down attribute hierarchy increasing granularity/resolution, a kind of unravel (quarter \rightarrow months). |