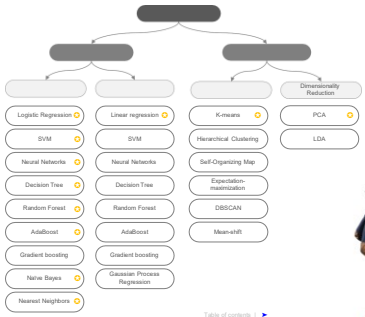


getting started with
Machine Learning

Part 2 | The well-known algorithms

:: Now let's study some popular algorithms



Naïve Bayes

:: Naïve Bayes

Naïve Bayes is a simple but important probabilistic model
Naïve Bayes is a simple multiclass classification algorithm based on applying Bayes' theorem with the "naïve" assumption of independence between the features. It assumes that the conditional probabilities of the independent variables are statistically independent.
It computes the conditional probability distribution of each feature given label, and then it applies Bayes' theorem to compute the conditional probability distribution of label given an observation and use it for prediction.
It classifies new data based on the **highest probability** of its belonging to a particular class.



- Naïve Bayes is a simple classifier model that is :-
- Based on the Bayes theorem
 - Supervised Learning
 - Easy to build
 - Faster to train, compared to other models
 - Often used as a baseline classifier for benchmarking

:: Naïve Bayes

Naïve Assumption
It assume that each input feature x_i is conditionally independent of every other feature x_j given the class C .

$$P(X_1, X_2, \dots, X_n | C) = P(X_1 | C) \cdot P(X_2 | C) \cdot \dots \cdot P(X_n | C)$$



For example, a fruit may be considered to be an apple if it is red, round, and about 10 cm in diameter.

A naïve Bayes classifier considers each of these features to contribute independently to the probability that this fruit is an apple, regardless of any possible correlations between the color, roundness, and diameter features!

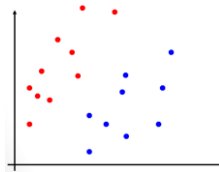
Since 0.9421 is greater than 0.2424 then the answer is 'no', we cannot play a game of tennis today.

- For data that has exactly two classes (you can also use it for multiclass classification with a technique called error-correcting output codes)
- For high-dimensional, nonlinearly separable data
- When you need a classifier that's simple, easy to interpret, and accurate

:: Support Vector Machine (SVM)

How to classify the data set in the following example ?

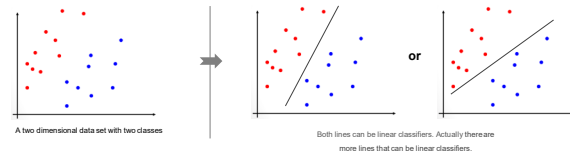
There is a two dimensional data set with two classes. How to classify it into 2 classes correctly?



:: Support Vector Machine (SVM)

How to classify the data set in the following example ?

There is a two dimensional data set with two classes. How to classify it into 2 classes correctly?



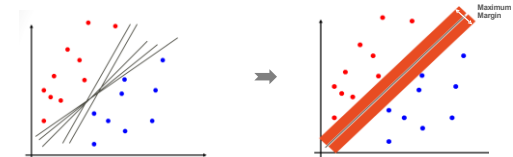
A two dimensional data set with two classes

Both lines can be linear classifiers. Actually there are more lines that can be linear classifiers.

:: Support Vector Machine (SVM)

Which line is the best classifier ?

There is a two dimensional data set with two classes. How to classify it into 2 classes correctly?



There are many lines that can be linear classifiers. Which one is the optimal classifier ?

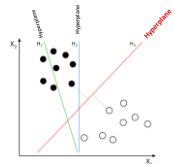
- Define the margin of a linear classifier as the width that the boundary could be increased by before hitting a data point.
- The maximum margin linear classifier is the simplest kind of SVG(called Linear SVM) ¹

<https://www.youtube.com/watch?v=98664gP2a7E>

:: Support Vector Machine (SVM)

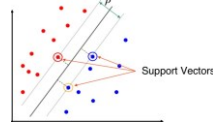
Select the hyper-plane which segregates the two classes better

a good separation is achieved by the hyper-plane that has the largest distance to the nearest training data points of any class (so-called functional margin)



- H₁ does not separate the classes.
- H₂ does, but only with a small margin.
- H₃ separates them with the maximum margin.

Source: https://en.wikipedia.org/wiki/Support_vector_machine

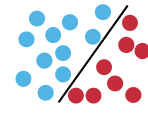


- Examples closest to the hyper-plane are **support vectors**
- Margin** ω of the separator is the distance between support vectors.

:: Support Vector Machine (SVM)

There is a new challenge !

How to classify the nonlinearly separable data by the hyperplane?



Use Linear SVM algorithm
The linear SVM can work well on classifying this dataset.

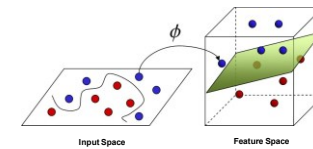


How to classify this data set ?
It seems that it's hard to classify this data set via the linear classifier.

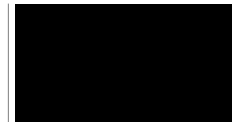
:: Support Vector Machine (SVM)

Create nonlinear classifiers by applying the kernel trick

SVMs sometimes use a kernel transform to transform nonlinearly separable data into higher dimensions where a linear decision boundary can be found.



Kernel trick allows the algorithm to fit the maximum-margin hyperplane in a transformed feature space. Although the classifier is a hyperplane in the transformed feature space, it may be nonlinear in the original input space.



<https://www.youtube.com/watch?v=98664gP2a7E>

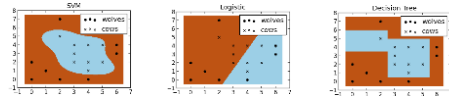


<https://www.youtube.com/watch?v=98664gP2a7E>

:: Support Vector Machine (SVM)

Example: SVM using a non-linear kernel

If you are farmer and you need to set up a fence to protect your cows from packs of wolves, where do you build your fence?



Where do you build your fence?

Well if you're a really data-driven farmer one way you could do it would be to build a classifier based on the position of the cows and wolves in your pasture. Trying a few different types of classifiers, we see that SVM does a great job at separating your cows from the packs of wolves. I thought these plots also do a nice job of illustrating the benefits of using a non-linear classifiers.

You can see the logistic and decision tree models both only make use of straight lines.¹

:: Support Vector Machine (SVM)

Strengths of SVM

Advantages

- The major strengths of SVM are the training is relatively easy. No local optimal, unlike in neural networks.¹
- Non-linear SVMs use a non-linear kernel. Non-linear SVM means that the boundary that the algorithm calculates doesn't have to be a straight line. The benefit is that you can capture much more complex relationships between your data points without having to perform difficult transformations. The downside is that the training time is much longer, because it is much more computationally intensive.²
- SVMs have a regularization parameter, which can help avoid over-fitting.³
- Effective in high dimensional spaces.⁴
- Still effective in cases where number of dimensions is greater than the number of samples.⁵
- Uses a subset of training points in the decision function (called support vectors), so it is also memory efficient.⁶



:: Support Vector Machine (SVM)

Weakness of SVM

Disadvantage

- For classification, the SVM is only directly applicable for two-class tasks. Therefore, algorithms that reduce the multi-class task to several binary problems have to be applied. Well, many SVM packages already have built-in multi-class classification functionality. Otherwise, use one-versus-the-rest approach to build a multiclass classifier.
- Unlike Logistic Regression classifiers, SVMs do not directly provide probability estimates. In many classification problems you actually want the probability of class membership.
- Parameters of a solved model are difficult to interpret.
- Long training time on large data sets
- Choosing a "good" kernel function can be tricky.

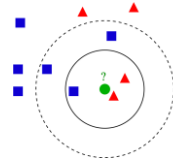


k Nearest Neighbor

:: k Nearest Neighbor (kNN)

kNN categorizes objects based on the classes of their nearest neighbors in the dataset

kNN predictions assume that objects near each other are similar. [Distance metrics](#), such as Euclidean, city block, cosine, and Chebychev, are used to find the nearest neighbor.



Example of k-NN classification

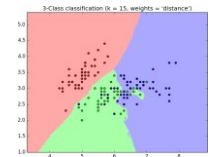
The test sample (green circle) should be classified either to the first class of blue squares or to the second class of red triangles.

- If $k = 3$ (solid line circle) it is assigned to the second class because there are 2 triangles and only 1 square inside the inner circle.
- If $k = 5$ (dashed line circle) it is assigned to the first class (3 squares vs. 2 triangles inside the outer circle).

:: k Nearest Neighbor (kNN)

An object is classified by a majority vote of its neighbors

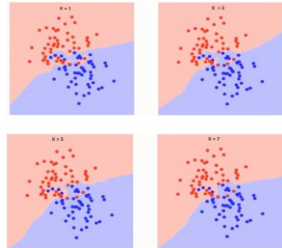
- **kNN algorithm is one of the simplest classification algorithm**
Despite its simplicity, nearest neighbors has been successful in a large number of classification and regression problems, including handwritten digits or satellite image scenes.¹
- **Often used in classification**
Being a non-parametric method, it is often successful in classification situations where the decision boundary is very irregular.²
- **Classification is computed from a simple majority vote of the nearest neighbors of each point**
In k-NN classification, the output is a class membership. An object is classified by a majority vote of its neighbors, with the object being assigned to the class most common among its k nearest neighbors.
- **k is constant specified by user**
In the classification phase, k is a user-defined constant. The best choice of k depends upon the data; generally, larger values of k reduce the effect of noise on the classification, but make boundaries between classes less distinct.
- **kNN is computationally expensive**



:: k Nearest Neighbor (kNN)

How do we choose the factor K?

Let's try to see the effect of value "K" on the class boundaries. Following are the different boundaries separating the two classes with different values of K.



If you watch carefully, you can see that the boundary becomes smoother with increasing value of K. With K increasing to infinity it finally becomes all blue or all red depending on the total majority.

At times, choosing K turns out to be a challenge while performing KNN modeling.

:: k Nearest Neighbor (kNN)

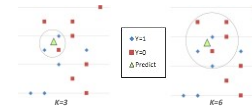
Strengths and Weakness of kNN

Strengths :

- It is beautifully simple and logical

Weaknesses :

- It may be driven by the choice of k, which may be a bad choice.
- Generally, larger values of k reduce the effect of noise on the classification, but make boundaries between classes less distinct.
- The accuracy of the algorithm can be severely degraded by the presence of noisy or irrelevant features, or if the feature scales are not consistent with their importance.
- In binary (two class) classification problems, it is helpful to choose k to be an odd number, as this avoids tied votes.
- It is important to review the sensitivity of the solution to different values of k.



Next Slide...

- When you need a simple algorithm to establish benchmark learning rules
- When memory usage of the trained model is a lesser concern
- When prediction speed of the trained model is a lesser concern

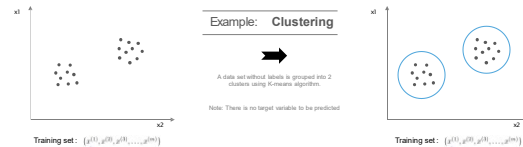
K-mean

Unsupervised Learning

Discover the structure within the unlabeled data

Unsupervised machine learning is the machine learning task of inferring a function to describe hidden structure from "unlabeled" data. Examples of unsupervised learning tasks include:

- clustering** (where we try to discover underlying groupings of examples)
- anomaly detection** (where we try to infer if some examples in a dataset do not conform to some expected pattern).



Example: Clustering

A data set without labels is grouped into 2 clusters using K-means algorithm.

Note: There is no target variable to be predicted

Cluster Analysis

A cluster is a subset of data which are similar

Clustering is a technique for finding similarity groups in a data, called clusters. It attempts to group individuals in a population together by similarity, but not driven by a specific purpose. Clustering is often called an unsupervised learning, as you don't have prescribed labels in the data and no class values denoting a priori grouping of the data instances are given.

A cluster is a collection of data items which are "similar" between them, and "dissimilar" to data items in other clusters.



Source: <http://laplace-machine.com/pooling-the-project-vec-to-evaluate-ml-algorithms>
Source: <http://www.kdd-cup.org/kdd-cup99/kddcup99.html>
Source: <http://www.kdd-cup.org/kdd-cup99/kddcup99.html>

Application of Clustering

Clustering is widely used in many applications

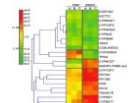
Clustering is hard to evaluate, but very useful in practice. Here are some examples of clustering applications



Market segmentation
Discover customer groups and use them for targeted marketing programs



Social network analysis
Help you find a group of coherent friends on the social network



Genomics
Finding groups of gene with similar expressions

Source: <http://www.kdd-cup.org/kdd-cup99/kddcup99.html>
Source: <http://www.kdd-cup.org/kdd-cup99/kddcup99.html>
Source: <http://www.kdd-cup.org/kdd-cup99/kddcup99.html>

:: k-Means

Weakness of K-means (continued)

- Sensitive to outliers and noise, which results in an inaccurate partition
- When the squared error criterion is used, outliers can unduly influence the clusters that are found¹. Weakness of arithmetic mean is not robust to outliers. Very far data from the centroid may pull the centroid away from the real one¹.
- Because of this, it is often useful to discover outliers and eliminate them beforehand¹.
 - It is important, however, to appreciate that there are certain clustering applications for which outliers should not be eliminated².

:: k-Means

Weakness of K-means (continued)

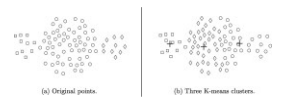
- K-means cannot handle non-globular clusters or clusters of different sizes and densities
- K-means has difficulty detecting the 'natural' clusters, when clusters have non-spherical shapes or widely different sizes or densities.



:: k-Means

Weakness of K-means (continued)

- K-means cannot handle clusters of different sizes

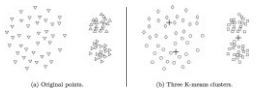


In Figure 1, K-means cannot find the three natural clusters because one of the clusters is much larger than the other two, and hence, the larger cluster is broken, while one of the smaller clusters is combined with a portion of the larger cluster.

:: k-Means

Weakness of K-means (continued)

- K-means cannot handle clusters of densities

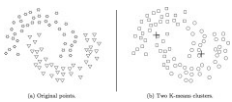


In Figure 2, K-means fails to find the three natural clusters because the two smaller clusters are much denser than the larger cluster.

:: k-Means

Weakness of K-means (continued)

- K-means cannot handle non-globular clusters



In Figure 3, K-means finds two clusters that mix portions of the two natural clusters because the shape of the natural clusters is not globular.

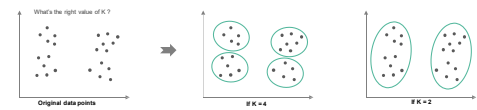
:: k-Means

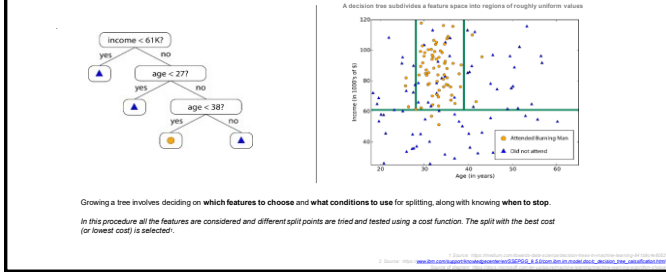
How to choose the value of k ?

In k-means clustering we must choose a value for k . This is still an active area of research and there are no definitive answers¹. By far the most common way of choosing the number of clusters, is still choosing it manually by looking at visualizations or by looking at the output of the clustering algorithm or something else²

Why is it not easy to determine the right value of K ?

Because it is often generally ambiguous about how many clusters there are in the data. For instance, the data points in the following example can be viewed as 4 clusters or 2 clusters. There are no absolutely right answer for that





Decision Tree

The key is to decide which attribute to split on

You need to pick the best attribute for splitting the training examples.

Day	Outlook	Temperature	Humidity	Wind	Play
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Sunny	Hot	High	Weak	No
4	Sunny	Hot	High	Strong	No
5	Sunny	Hot	High	Weak	No
6	Sunny	Hot	High	Strong	No
7	Sunny	Hot	High	Weak	No
8	Sunny	Hot	High	Strong	No
9	Sunny	Hot	High	Weak	No
10	Sunny	Hot	High	Strong	No
11	Sunny	Hot	High	Weak	No
12	Sunny	Hot	High	Strong	No
13	Sunny	Hot	High	Weak	No
14	Sunny	Hot	High	Strong	No
15	Sunny	Hot	High	Weak	No
16	Sunny	Hot	High	Strong	No
17	Sunny	Hot	High	Weak	No
18	Sunny	Hot	High	Strong	No
19	Sunny	Hot	High	Weak	No
20	Sunny	Hot	High	Strong	No
21	Sunny	Hot	High	Weak	No
22	Sunny	Hot	High	Strong	No
23	Sunny	Hot	High	Weak	No
24	Sunny	Hot	High	Strong	No
25	Sunny	Hot	High	Weak	No
26	Sunny	Hot	High	Strong	No
27	Sunny	Hot	High	Weak	No
28	Sunny	Hot	High	Strong	No
29	Sunny	Hot	High	Weak	No
30	Sunny	Hot	High	Strong	No
31	Sunny	Hot	High	Weak	No
32	Sunny	Hot	High	Strong	No
33	Sunny	Hot	High	Weak	No
34	Sunny	Hot	High	Strong	No
35	Sunny	Hot	High	Weak	No
36	Sunny	Hot	High	Strong	No
37	Sunny	Hot	High	Weak	No
38	Sunny	Hot	High	Strong	No
39	Sunny	Hot	High	Weak	No
40	Sunny	Hot	High	Strong	No
41	Sunny	Hot	High	Weak	No
42	Sunny	Hot	High	Strong	No
43	Sunny	Hot	High	Weak	No
44	Sunny	Hot	High	Strong	No
45	Sunny	Hot	High	Weak	No
46	Sunny	Hot	High	Strong	No
47	Sunny	Hot	High	Weak	No
48	Sunny	Hot	High	Strong	No
49	Sunny	Hot	High	Weak	No
50	Sunny	Hot	High	Strong	No
51	Sunny	Hot	High	Weak	No
52	Sunny	Hot	High	Strong	No
53	Sunny	Hot	High	Weak	No
54	Sunny	Hot	High	Strong	No
55	Sunny	Hot	High	Weak	No
56	Sunny	Hot	High	Strong	No
57	Sunny	Hot	High	Weak	No
58	Sunny	Hot	High	Strong	No
59	Sunny	Hot	High	Weak	No
60	Sunny	Hot	High	Strong	No
61	Sunny	Hot	High	Weak	No
62	Sunny	Hot	High	Strong	No
63	Sunny	Hot	High	Weak	No
64	Sunny	Hot	High	Strong	No
65	Sunny	Hot	High	Weak	No
66	Sunny	Hot	High	Strong	No
67	Sunny	Hot	High	Weak	No
68	Sunny	Hot	High	Strong	No
69	Sunny	Hot	High	Weak	No
70	Sunny	Hot	High	Strong	No
71	Sunny	Hot	High	Weak	No
72	Sunny	Hot	High	Strong	No
73	Sunny	Hot	High	Weak	No
74	Sunny	Hot	High	Strong	No
75	Sunny	Hot	High	Weak	No
76	Sunny	Hot	High	Strong	No
77	Sunny	Hot	High	Weak	No
78	Sunny	Hot	High	Strong	No
79	Sunny	Hot	High	Weak	No
80	Sunny	Hot	High	Strong	No
81	Sunny	Hot	High	Weak	No
82	Sunny	Hot	High	Strong	No
83	Sunny	Hot	High	Weak	No
84	Sunny	Hot	High	Strong	No
85	Sunny	Hot	High	Weak	No
86	Sunny	Hot	High	Strong	No
87	Sunny	Hot	High	Weak	No
88	Sunny	Hot	High	Strong	No
89	Sunny	Hot	High	Weak	No
90	Sunny	Hot	High	Strong	No
91	Sunny	Hot	High	Weak	No
92	Sunny	Hot	High	Strong	No
93	Sunny	Hot	High	Weak	No
94	Sunny	Hot	High	Strong	No
95	Sunny	Hot	High	Weak	No
96	Sunny	Hot	High	Strong	No
97	Sunny	Hot	High	Weak	No
98	Sunny	Hot	High	Strong	No
99	Sunny	Hot	High	Weak	No
100	Sunny	Hot	High	Strong	No

outlook

Sunny 2 Yes / 3 No

Overcast 4 Yes / 0 No

Rain 3 Yes / 2 No

Wind

Weak 6 Yes / 2 No

Strong 3 Yes / 3 No

- To compare the different ways to split data in a node
- Want to measure the "purity" of the split:
 - Pure set (5 yes / 0 No) = completely certain (100%)
 - Impure (3 yes / 3 No) = completely uncertain (50%)
- Splits on all attributes are tested
 - Constructing a decision tree is all about finding attribute that returns the **highest information gain** or lowest **Gini Index** (i.e., the most homogeneous branches?)

Decision Tree

Measures that can be used to capture the purity of split

Information Gain, Gain Ratio and Gini Index are the most common methods of attribute selection

Information Gain

Information Gain Ratio


Gini Index

Decision Tree

We need to understand Entropy first before we move on

ID3 uses Entropy and Information Gain to construct a decision tree. It's necessary to understand what Entropy means.

What's Entropy ?



Decision Tree

Generally entropy is a measure of **disorder or uncertainty**.

Entropy is a concept used in Physics, mathematics, computer science (information theory) and other fields of science. The concept of entropy originated in thermodynamics as a measure of molecular disorder: entropy approaches zero when molecules are still and well ordered.

Entropy in Physics:

Low

Medium

High

Entropy, as far, had been a concept in physics. If the particles inside a system have many possible positions to move around, then the system has high entropy, and if they have to stay tight, then the system has low entropy.

For example, water in its three states, solid, liquid, and gas, has different entropies. The molecules in ice have to stay in a lattice, so it is a rigid system, so ice has low entropy. The molecules in water have more freedom to move around, so water in liquid state has medium entropy. The molecules inside water vapor can pretty much go anywhere they want, so water vapor has high entropy.

Decision Tree


Entropy in Information theory

Information entropy is defined as the average amount of information produced by a stochastic source of data... Generally, information entropy is the average amount of information conveyed by an event

- when a low-probability event occurs (i.e. when the data source has a lower-probability value,) the event carries more "information" ("surprise") than when the source data has a higher-probability value. Rarer events provide more information than observed.
- when a high-probability event occurs, the event with high-probability value is less informative than less common event. Entropy is zero when one outcome is certain to occur.

Entropy is a measure of unpredictability of the state, or equivalently, of its average information content. To get an intuitive understanding of these terms, consider the example of a political poll.

- Usually, such polls happen because the outcome of the poll is not already known. In other words, the outcome of the poll is relatively unpredictable, and actually performing the poll and learning the results gives some new information; these are just different ways of saying that the a priori entropy of the poll results is large.
- Now, consider the case that the same poll is performed a second time shortly after the first poll. Since the result of the first poll is already known, the outcome of the second poll can be predicted well and the results obtained will contain much less new information; in this case the a priori entropy of the second poll result is small relative to that of the first.



Decision Tree

More uncertainty, more entropy!

entropy is a measure of disorder or uncertainty

High Entropy

Low Entropy

Increase entropy

- High Disorder
- Messy
- High Randomness
- High uncertainty
- Low certainty

- Well-ordered
- Neat
- Low Randomness
- Low uncertainty
- High certainty

:: Decision Tree

Mathematical definition of entropy

The measure of information entropy associated with each possible data value is the **negative logarithm** of the probability mass function for the value. If we have a set with n different values in it, we can calculate the entropy as follows:

$$\text{Entropy} = - \sum_{i=1}^n p_i \log_2 p_i$$

Where p_i is the **probability** of getting the i^{th} value when randomly selecting one from the set.
In other words, Where there are n classes, and p_i is the probability an object from the i^{th} class appearing.



:: Decision Tree

Another example of calculating entropy

Use entropy to measure the "purity" of the split

$$\text{Entropy} = - \sum_{i=1}^n p_i \log_2 p_i$$



Video: Entropy and Information Gain (27min)
<https://www.youtube.com/watch?v=3D5u4wT0EAD0w>

:: Decision Tree

“信息熵”是度量样本集合纯度最常用的一种指标，假定当前样本集合 D 中第 k 类样本所占的比例为 p_k ($k = 1, 2, \dots, |Y|$)，则 D 的信息熵定义为

$$\text{Ent}(D) = - \sum_{k=1}^{|Y|} p_k \log_2 p_k$$

$\text{Ent}(D)$ 的值越小，则 D 的纯度越高

- 计算信息熵时约定：若 $p = 0$ ，则 $p \log_2 p = 0$
- $\text{Ent}(D)$ 的最小值为0，最大值为 $\log_2 |Y|$

:: Decision Tree

Measures that can be used to capture the purity of split

Information Gain, Gain Ratio and Gini Index are the most common methods of attribute selection

Information Gain

Information Gain Ratio

Gini Index

:: Decision Tree

Measures that can be used to capture the purity of split – Information Gain

Information gain increases with the average purity of the subsets. Strategy: choose attribute that gives greatest information gain.

Information Gain

i.e. entropy reduction

离散属性 a 有 v 个可能的取值 $\{a^1, a^2, \dots, a^v\}$ ，用 a 来进行划分，则会产生 v 个分支结点，其中第 a 个分支结点包含 D 中所有在属性 a 上取值为 a^i 的样本，记为 D^i 。则可计算出用属性 a 对样本集 D 进行划分所获得的“信息增益”：

$$\text{Gain}(D, a) = \text{Ent}(D) - \sum_{i=1}^v \frac{|D^i|}{|D|} \text{Ent}(D^i)$$

为分支结点权重，样本数越多的分支结点的影响越大

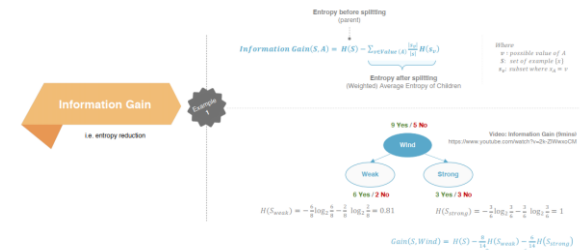
一般而言，信息增益越大，则意味着使用属性 a 来进行划分所获得的“纯度提升”越大

ID3决策树学习算法(Quinlan, 1986)以信息增益为准则来选择划分属性

:: Decision Tree

Measures that can be used to capture the purity of split – Information Gain (continued)

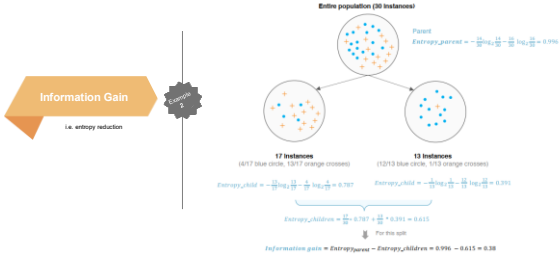
The goal is to decrease in entropy (uncertainty) after splitting. And also to want many items in pure sets



Decision Tree

Measures that can be used to capture the purity of split – Information Gain (continued)

The goal is to decrease in entropy (uncertainty) after splitting. And also to want many items in pure sets



Decision Tree

Measures that can be used to capture the purity of split – Information Gain Ratio

Information gain approach has a problem that it bias towards attributes that have a large number of values over attributes that have a smaller number of values. These "Super Attributes" will easily be selected as the root, resulted in a biased tree that classifies perfectly but performs poorly on unseen instances. We can penalize attributes with large numbers of values by using an alternative method for attribute selection, referred to as Gain Ratio.

Information Gain Ratio

- 增益率定义：
$$\text{Gain_ratio}(D, a) = \frac{\text{Gain}(D, a)}{\text{IV}(a)}$$
其中
$$\text{IV}(a) = - \sum_{v=1}^V \frac{|D^v|}{|D|} \log_2 \frac{|D^v|}{|D|}$$
称为属性a的“固有值” [Quinlan, 1993]，属性a的可能取值数目越多（即V越大），则IV(a)的值通常就越大
- 存在的问题：

增益率准则对可取值数目较少的属性有所偏好
- C4.5 [Quinlan, 1993]使用了一个启发式：先从候选划分属性中找出信息增益高于平均水平的属性，再从中选增益率最高的

Decision Tree

Measures that can be used to capture the purity of split – Gini Index

The impurity measure used in building decision tree in CART algorithm is Gini index. The attribute that maximizes the reduction in impurity is chosen as the splitting attribute

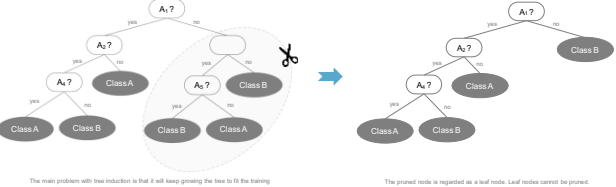
Gini Index

- 数据集D的纯度可用“基尼值”来度量
$$\text{Gini}(D) = \sum_{k=1}^K \sum_{j \neq k} p_k p_j = 1 - \sum_{k=1}^K p_k^2$$
反映了从D中随机抽取两个样本，其类别标记不一致的概率
- Gini(D)越小，数据集D的纯度越高
- 属性a的基尼指数定义为：
$$\text{Gini_index}(D, a) = \sum_{v=1}^V \frac{|D^v|}{|D|} \text{Gini}(D^v)$$
- 应选择那个使划分后基尼指数最小的属性作为最优划分属性，即
$$a_* = \underset{a \in A}{\text{argmin}} \text{Gini_index}(D, a)$$
- CART [Breiman et al., 1984]采用“基尼指数”来选择划分属性

Decision Tree

Address overfitting in Decision Trees by Pruning

Pruning means to change the model by deleting the child nodes of a branch node. Pruning usually results in reducing size of tree, avoids unnecessary complexity, and to avoid overfitting.



The main problem with tree induction is that it will keep growing the tree to fit the training data until it creates pure leaf nodes. This will result in large, overly complex trees that overfit the data.

Decision Tree

Approaches for Pruning

pre-pruning

One strategy to prevent overfitting is to prevent the tree from becoming really detailed and complex by stopping its growth early.

- Halt tree construction early—do not split a node if this would result in the goodness measure falling below a threshold
- Upon halting the node becomes a leaf! Upon halting, the node becomes a leaf
- The leaf may hold the most frequent class among the subset tuples

Problem: difficult to choose an appropriate threshold!

post-pruning

Another strategy is to build a complete tree with pure leaves but then to prune back the tree into a simpler form

- Remove non-significant branches from a "fully grown" tree
- After trimming, replace subtree by a leaf node
- The leaf is labeled with the most frequent class

Practically, post-pruning strategy is more successful because it is not easy to precisely estimate when to stop growing the tree!

Decision Tree

Advantages of Decision Tree

- Simple to understand and to interpret
Decision tree is a "white box" model. Decision Trees are able to produce "understandable" rules. The trees can be also visualized. In contrast, in a black box model (e.g. in neural network, random forest), it is usually hard to explain in simple terms why the predictions were made.
- To build decision tree requires little data preparation
Decision trees do not need feature scaling. Decision trees can deal with a reasonable amount of missing values. Decision trees are also not sensitive to outliers. In contrast, other algorithms often require data normalization, dummy variables need to be created and blank values to be removed.
- Decision trees are able to handle both continuous and categorical variables.
Other techniques are usually specialized in analyzing datasets that have only one type of variables. Able to handle multiclass problems.
- Implicitly perform feature selection
Decision trees such as CART, for instance, have a built-in mechanism to perform variable selection. Decision trees provide a clear indication of which fields are most important for prediction or classification



:: Decision Tree

Disadvantages of Decision Tree

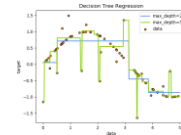
- **They are prone to over-fitting.**
Decision-tree learners can create over-complex trees that do not generalize to the data well.
- **Decision Tree learner can create biased trees in case of unbalanced data**
Decision-tree learners create biased trees if some classes dominate. It is therefore recommended to balance the dataset prior to fitting with the decision tree.
- **Instability**
Decision trees can be unstable because small variations in the data might result in a completely different tree being generated. This problem is mitigated by using decision trees within an ensemble. Random Forests can limit this instability by averaging predictions over many trees.
- **Greedy approach used by Decision tree doesn't guarantee best solution**
Greedy algorithm where locally optimal decisions are made at each node cannot guarantee to return the globally optimal decision tree. This can be mitigated by training multiple trees, where the features and samples are randomly sampled with replacement.



:: Decision Tree

Disadvantages of Decision Tree - continued

- **Standard decision trees are restricted by hard, axis-aligned splits of the input space¹**
The splits are aligned with the axes of the feature space and this leads to rectilinear decision boundaries.



The deeper the tree, the more complex the decision rules and the better the model.

We can see that if the maximum depth of the tree is set too high, the decision trees learn too few details of the training data and learn from the noise, i.e. they overfit.

Another issue is particularly problematic in regression, where we are typically aiming to model smooth functions, and yet the tree model produces piecewise-constant predictions with discontinuities at the split boundaries.²



:: Decision Tree Example

编号	色泽	根蒂	敲声	纹理	脐部	触感	好瓜
1	青绿	蜷缩	浊响	清晰	凹脐	硬滑	是
2	乌黑	蜷缩	沉闷	清晰	凹脐	硬滑	是
3	乌黑	蜷缩	浊响	清晰	凹脐	硬滑	是
4	青绿	蜷缩	沉闷	清晰	凹脐	硬滑	是
5	浅白	蜷缩	浊响	清晰	凹脐	硬滑	是
6	青绿	蜷缩	浊响	清晰	脐凹	软粘	是
7	乌黑	蜷缩	浊响	清晰	脐凹	硬滑	是
8	乌黑	蜷缩	浊响	清晰	脐凹	硬滑	是
9	乌黑	蜷缩	沉闷	清晰	脐凹	硬滑	否
10	青绿	硬挺	清脆	清晰	平脐	软粘	否
11	浅白	硬挺	清脆	模糊	平脐	硬滑	否
12	浅白	蜷缩	浊响	模糊	平脐	软粘	否
13	青绿	蜷缩	浊响	清晰	凹脐	硬滑	否
14	浅白	蜷缩	沉闷	清晰	凹脐	硬滑	否
15	乌黑	蜷缩	浊响	清晰	脐凹	软粘	否
16	浅白	蜷缩	浊响	模糊	平脐	硬滑	否
17	青绿	蜷缩	沉闷	清晰	脐凹	硬滑	否

该数据集包含17个训练样本, $|D|=17$, 其中正例占 $p_1 = \frac{8}{17}$, 反例占 $p_2 = \frac{9}{17}$, 计算得到根结点的信息熵为

$$\text{Ent}(D) = -\sum_{k=1}^2 p_k \log_2 p_k = -\left(\frac{8}{17} \log_2 \frac{8}{17} + \frac{9}{17} \log_2 \frac{9}{17}\right) = 0.998$$

:: Decision Tree Example

- 以属性“色泽”为例, 其对应的3个数据子集分别为 D^1 (色泽=青绿), D^2 (色泽=乌黑), D^3 (色泽=浅白)
- 子集 D^1 包含编号为{1, 4, 6, 10, 13, 17}的6个样例, 其中正例占 $p_1 = \frac{3}{6}$, 反例占 $p_2 = \frac{3}{6}$, D^2 , D^3 同理, 3个结点的信息熵为:

$$\text{Ent}(D^1) = -\left(\frac{3}{6} \log_2 \frac{3}{6} + \frac{3}{6} \log_2 \frac{3}{6}\right) = 1.000$$

$$\text{Ent}(D^2) = -\left(\frac{4}{6} \log_2 \frac{4}{6} + \frac{2}{6} \log_2 \frac{2}{6}\right) = 0.918$$

$$\text{Ent}(D^3) = -\left(\frac{1}{3} \log_2 \frac{1}{3} + \frac{2}{3} \log_2 \frac{2}{3}\right) = 0.722$$

- 属性“色泽”的信息增益为

$$\begin{aligned} \text{Gain}(D, \text{色泽}) &= \text{Ent}(D) - \sum_{i=1}^3 \frac{|D^i|}{|D|} \text{Ent}(D^i) \\ &= 0.998 - \left(\frac{6}{17} \times 1.000 + \frac{6}{17} \times 0.918 + \frac{5}{17} \times 0.722\right) \\ &= 0.109 \end{aligned}$$

:: Decision Tree Example

- 类似的, 其他属性的信息增益为

$$\text{Gain}(D, \text{根蒂}) = 0.143$$

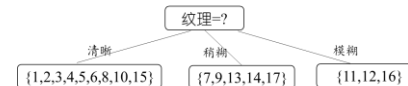
$$\text{Gain}(D, \text{敲声}) = 0.141$$

$$\text{Gain}(D, \text{纹理}) = 0.381$$

$$\text{Gain}(D, \text{脐部}) = 0.289$$

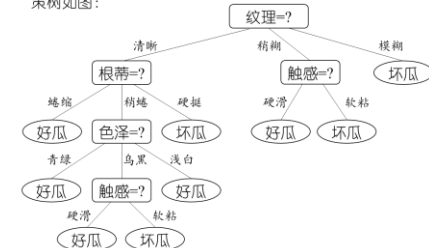
$$\text{Gain}(D, \text{触感}) = 0.006$$

- 显然, 属性“纹理”的信息增益最大, 其被选为划分属性



:: Decision Tree Example

- 决策树学习算法将对每个分支结点做进一步划分, 最终得到的决策树如图:



:: Decision Tree 练习

ID	年龄	有工作	有自己的房子	信贷情况	类别
1	青年	否	否	一般	否
2	青年	否	否	好	否
3	青年	是	否	好	是
4	青年	是	是	一般	是
5	青年	否	否	一般	否
6	中年	否	否	一般	否
7	中年	否	否	好	否
8	中年	是	是	好	是
9	中年	否	是	非常好	是
10	中年	否	是	非常好	是
11	老年	否	是	非常好	是
12	老年	否	是	好	是
13	老年	是	否	好	是
14	老年	是	否	非常好	是
15	老年	否	否	一般	否

通过所给的训练数据学习一个贷款申请的决策树，用以对未来的贷款申请进行分类，即当新的客户提出贷款申请时，根据申请人的特征利用决策树决定是否批准贷款申请。（自由选择方法）