

Decision trees and ensemble methods

Mauricio A. Álvarez

Machine Learning and Adaptive Intelligence
The University of Sheffield



The
University
Of
Sheffield.

Contents

Decision trees

- Definitions

- Shannon's Entropy

- Information Gain

Handling continuous features

Growing a tree

Regression trees

Beyond basic decision trees

- Alternative impurity measures

- Overfitting and Tree Pruning

- Model Ensembles

Contents

Decision trees

Definitions

Shannon's Entropy

Information Gain

Handling continuous features

Growing a tree

Regression trees

Beyond basic decision trees

Alternative impurity measures

Overfitting and Tree Pruning

Model Ensembles

Nodes in a decision tree

- ❑ A decision tree consists of:
 1. a **root node** (or starting node),
 2. **interior nodes**
 3. and **leaf nodes** (or terminating nodes).

- ❑ Each of the non-leaf nodes (root and interior) in the tree specifies a test to be carried out on one of the query's descriptive features.

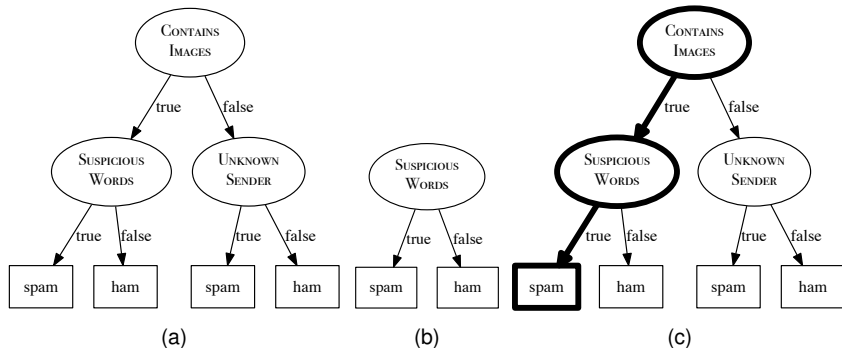
- ❑ Each of the leaf nodes specifies a predicted classification for the query.

An email spam prediction dataset

An email spam prediction dataset.

ID	SUSPICIOUS WORDS	UNKNOWN SENDER	CONTAINS IMAGES	CLASS
376	true	false	true	spam
489	true	true	false	spam
541	true	true	false	spam
693	false	true	true	ham
782	false	false	false	ham
976	false	false	false	ham

Two decision trees and a query instance



(a) and (b) show two decision trees that are consistent with the instances in the spam dataset. (c) shows the path taken through the tree shown in (a) to make a prediction for the query instance: SUSPICIOUS WORDS = 'true', UNKNOWN SENDER = 'true', CONTAINS IMAGES = 'true'.

Which tree to use?

- Both of these trees will return identical predictions for all the examples in the dataset.
- So, which tree should we use?
- Occam's Razor: *"Among competing hypotheses, the one with the fewest assumptions should be selected"*.

Informative features

- ❑ A decision tree is a machine learning algorithm that tries to build predictive models **using the most informative features**.
- ❑ An informative feature is a feature whose values split the instances in the dataset into **homogeneous sets** with respect to the target value.
- ❑ So the problem is to figure out which features are the most informative ones to ask questions about.

How do we create shallow trees?

- ❑ The tree that tests **SUSPICIOUS WORDS** at the root is very shallow because the **SUSPICIOUS WORDS** feature perfectly splits the data into pure groups of '*spam*' and '*ham*'.
- ❑ Descriptive features that split the dataset into pure sets with respect to the target feature provide information about the target feature.
- ❑ So we can make shallow trees by testing the informative features early on in the tree.
- ❑ All we need to do that is a computational metric of the purity of a set: **entropy**

Contents

Decision trees

Definitions

Shannon's Entropy

Information Gain

Handling continuous features

Growing a tree

Regression trees

Beyond basic decision trees

Alternative impurity measures

Overfitting and Tree Pruning

Model Ensembles

Entropy

- Claude Shannon's entropy model defines a computational measure of the impurity of the elements of a set.
- An easy way to understand the entropy of a set is to think in terms of the uncertainty associated with guessing the result if you were to make a random selection from the set.

Probability and entropy

- Entropy is related to the probability of an outcome.
 - High probability \rightarrow Low entropy
 - Low probability \rightarrow High entropy
- If we take the **log** of a probability and multiply it by -1 we get this mapping!

Mathematical definition of entropy

- Shannon's model of entropy is a weighted sum of the logs of the probabilities of each of the possible outcomes when we make a random selection from a set.

$$H(t) = - \sum_{i=1}^I (P(t = i) \times \log_s(P(t = i)))$$

Entropy in poker cards

- What is the entropy of a set of 52 different playing cards?

$$\begin{aligned} H(card) &= - \sum_{i=1}^{52} P(card = i) \times \log_2(P(card = i)) \\ &= - \sum_{i=1}^{52} 0.019 \times \log_2(0.019) = - \sum_{i=1}^{52} -0.1096 \\ &= 5.700 \text{ bits} \end{aligned}$$

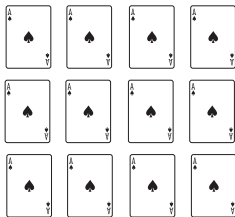
Entropy in poker cards

- What is the entropy of a set of 52 playing cards if we only distinguish between the cards based on their suit $\{\heartsuit, \clubsuit, \diamondsuit, \spadesuit\}$?

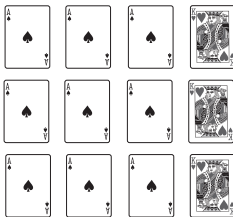
Entropy in poker cards

$$\begin{aligned} H(\text{suit}) &= - \sum_{l \in \{\heartsuit, \clubsuit, \diamondsuit, \spadesuit\}} P(\text{suit} = l) \times \log_2(P(\text{suit} = l)) \\ &= - \left((P(\heartsuit) \times \log_2(P(\heartsuit))) + (P(\clubsuit) \times \log_2(P(\clubsuit))) \right. \\ &\quad \left. + (P(\diamondsuit) \times \log_2(P(\diamondsuit))) + (P(\spadesuit) \times \log_2(P(\spadesuit))) \right) \\ &= - \left(\left(\frac{13}{52} \times \log_2\left(\frac{13}{52}\right) \right) + \left(\frac{13}{52} \times \log_2\left(\frac{13}{52}\right) \right) \right. \\ &\quad \left. + \left(\frac{13}{52} \times \log_2\left(\frac{13}{52}\right) \right) + \left(\frac{13}{52} \times \log_2\left(\frac{13}{52}\right) \right) \right) \\ &= - \left((0.25 \times -2) + (0.25 \times -2) \right. \\ &\quad \left. + (0.25 \times -2) + (0.25 \times -2) \right) \\ &= 2 \text{ bits} \end{aligned}$$

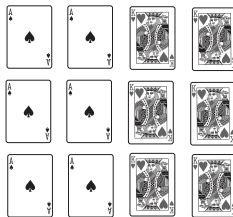
Different entropies



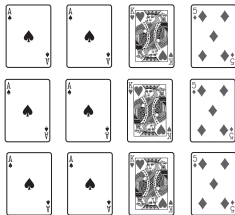
(a) $H(card) = 0.00$



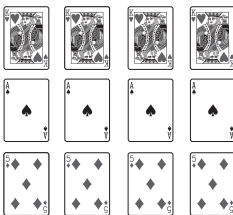
(b) $H(card) = 0.81$



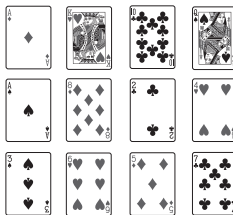
(c) $H(card) = 1.00$



(d) $H(card) = 1.50$



(e) $H(card) = 1.58$



(f) $H(card) = 3.58$

The entropy of different sets of playing cards measured in bits.

Entropy of a message

The relationship between the entropy of a message and the set it was selected from.

Entropy of a Message	Properties of the Message Set
High	A large set of equally likely messages.
Medium	A large set of messages, some more likely than others.
Medium	A small set of equally likely messages.
Low	A small set of messages with one very likely message.

Contents

Decision trees

- Definitions

- Shannon's Entropy

- Information Gain**

Handling continuous features

Growing a tree

Regression trees

Beyond basic decision trees

- Alternative impurity measures

- Overfitting and Tree Pruning

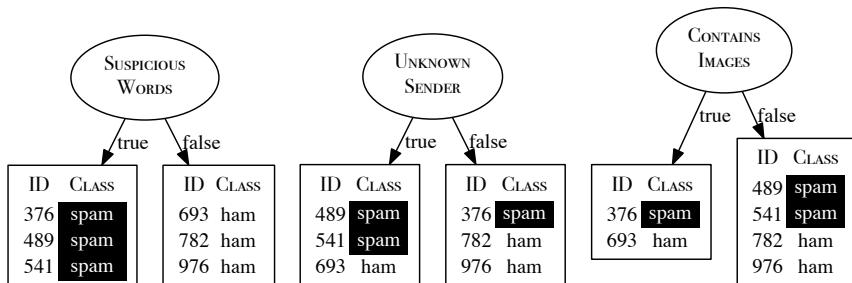
- Model Ensembles

The spam dataset again

An email spam prediction dataset.

ID	SUSPICIOUS WORDS	UNKNOWN SENDER	CONTAINS IMAGES	CLASS
376	true	false	true	spam
489	true	true	false	spam
541	true	true	false	spam
693	false	true	true	ham
782	false	false	false	ham
976	false	false	false	ham

Partitions in the spam dataset



How the instances in the spam dataset split when we partition using each of the different descriptive features from the spam dataset table.

Pure subsets

- Our intuition is that the ideal discriminatory feature will partition the data into **pure** subsets where all the instances in each subset have the same classification.
 - SUSPICIOUS WORDS perfect split.
 - UNKNOWN SENDER mixture but some information (when *'true'* most instances are *'spam'*).
 - CONTAINS IMAGES no information.
- One way to implement this idea is to use a metric called **information gain**.

Information Gain

- The information gain of a descriptive feature can be understood as a measure of the reduction in the overall entropy of a prediction task by testing on that feature.

Computing the information gain

1. Compute the entropy of the original dataset with respect to the target feature. This gives us a measure of how much information is required in order to organize the dataset into pure sets.
2. For each descriptive feature, create the sets that result by partitioning the instances in the dataset using their feature values, and then sum the entropy scores of each of these sets. This gives a measure of the information that remains required to organize the instances into pure sets after we have split them using the descriptive feature.
3. Subtract the remaining entropy value (computed in step 2) from the original entropy value (computed in step 1) to give the information gain.

Computing the information gain

Computing information gain involves the following three equations:

$$H(t, \mathcal{D}) = - \sum_{l \in \text{levels}(t)} (P(t = l) \times \log_2(P(t = l)))$$

$$\text{rem}(d, \mathcal{D}) = \sum_{l \in \text{levels}(d)} \underbrace{\frac{|\mathcal{D}_{d=l}|}{|\mathcal{D}|}}_{\text{weighting}} \times \underbrace{H(t, \mathcal{D}_{d=l})}_{\text{entropy of partition } \mathcal{D}_{d=l}}$$

$$IG(d, \mathcal{D}) = H(t, \mathcal{D}) - \text{rem}(d, \mathcal{D})$$

Example with the spam dataset

- As an example we will calculate the information gain for each of the descriptive features in the spam email dataset.

Step 1: Entropy for the target feature

- Calculate the **entropy** for the target feature in the dataset.

$$H(t, \mathcal{D}) = - \sum_{l \in \text{levels}(t)} (P(t = l) \times \log_2(P(t = l)))$$

ID	SUSPICIOUS WORDS	UNKNOWN SENDER	CONTAINS IMAGES	CLASS
376	true	false	true	spam
489	true	true	false	spam
541	true	true	false	spam
693	false	true	true	ham
782	false	false	false	ham
976	false	false	false	ham

Step 1: Entropy for the target feature

$$\begin{aligned} H(t, \mathcal{D}) &= - \sum_{l \in \{\text{'spam'}, \text{'ham'}\}} (P(t = l) \times \log_2(P(t = l))) \\ &= - ((P(t = \text{'spam'}) \times \log_2(P(t = \text{'spam'})) \\ &\quad + (P(t = \text{'ham'}) \times \log_2(P(t = \text{'ham'}))) \\ &= - \left(\left(\frac{3}{6} \times \log_2\left(\frac{3}{6}\right) \right) + \left(\frac{3}{6} \times \log_2\left(\frac{3}{6}\right) \right) \right) \\ &= 1 \text{ bit} \end{aligned}$$

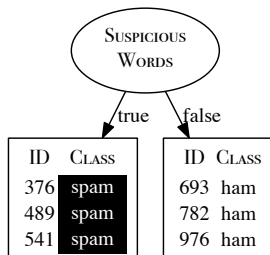
Step 2: Remainder for the SUSPICIOUS W. feature

- Calculate the **remainder** for the SUSPICIOUS WORDS feature in the dataset.

$$rem(d, \mathcal{D}) = \sum_{l \in levels(d)} \underbrace{\frac{|\mathcal{D}_{d=l}|}{|\mathcal{D}|}}_{\text{weighting}} \times \underbrace{H(t, \mathcal{D}_{d=l})}_{\text{entropy of partition } \mathcal{D}_{d=l}}$$

ID	SUSPICIOUS WORDS	UNKNOWN SENDER	CONTAINS IMAGES	CLASS
376	true	false	true	spam
489	true	true	false	spam
541	true	true	false	spam
693	false	true	true	ham
782	false	false	false	ham
976	false	false	false	ham

Step 2: Partitions for the SUSPICIOUS W. feature



Partitions for the SUSPICIOUS W. feature

Step 2: Remainder for the SUSPICIOUS W. feature

$$\begin{aligned} & \text{rem}(\text{WORDS}, \mathcal{D}) \\ &= \left(\frac{|\mathcal{D}_{\text{WORDS}=T}|}{|\mathcal{D}|} \times H(t, \mathcal{D}_{\text{WORDS}=T}) \right) + \left(\frac{|\mathcal{D}_{\text{WORDS}=F}|}{|\mathcal{D}|} \times H(t, \mathcal{D}_{\text{WORDS}=F}) \right) \\ &= \left(\frac{3}{6} \times \left(- \sum_{l \in \{\text{'spam'}, \text{'ham'}\}} P(t=l) \times \log_2(P(t=l)) \right) \right) \\ &+ \left(\frac{3}{6} \times \left(- \sum_{l \in \{\text{'spam'}, \text{'ham'}\}} P(t=l) \times \log_2(P(t=l)) \right) \right) \\ &= \left(\frac{3}{6} \times \left(- \left(\left(\frac{3}{3} \times \log_2\left(\frac{3}{3}\right) \right) + \left(\frac{0}{3} \times \log_2\left(\frac{0}{3}\right) \right) \right) \right) \right) \\ &+ \left(\frac{3}{6} \times \left(- \left(\left(\frac{0}{3} \times \log_2\left(\frac{0}{3}\right) \right) + \left(\frac{3}{3} \times \log_2\left(\frac{3}{3}\right) \right) \right) \right) \right) = 0 \text{ bits} \end{aligned}$$

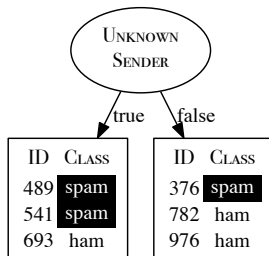
Step 2: Remainder for the UNKNOWN SENDER feature

- Calculate the **remainder** for the UNKNOWN SENDER feature in the dataset.

$$rem(d, \mathcal{D}) = \sum_{l \in levels(d)} \underbrace{\frac{|\mathcal{D}_{d=l}|}{|\mathcal{D}|}}_{\text{weighting}} \times \underbrace{H(t, \mathcal{D}_{d=l})}_{\text{entropy of partition } \mathcal{D}_{d=l}}$$

ID	SUSPICIOUS WORDS	UNKNOWN SENDER	CONTAINS IMAGES	CLASS
376	true	false	true	spam
489	true	true	false	spam
541	true	true	false	spam
693	false	true	true	ham
782	false	false	false	ham
976	false	false	false	ham

Step 2: Partitions for the UNKNOWN SENDER feature



Partitions for the UNKNOWN SENDER feature

Step 2: Remainder for the UNKNOWN SENDER feature

$$\begin{aligned} \text{rem}(\text{SENDER}, \mathcal{D}) &= \left(\frac{|\mathcal{D}_{\text{SENDER}=T}|}{|\mathcal{D}|} \times H(t, \mathcal{D}_{\text{SENDER}=T}) \right) + \left(\frac{|\mathcal{D}_{\text{SENDER}=F}|}{|\mathcal{D}|} \times H(t, \mathcal{D}_{\text{SENDER}=F}) \right) \\ &= \left(\frac{3}{6} \times \left(- \sum_{l \in \{\text{'spam'}, \text{'ham'}\}} P(t=l) \times \log_2(P(t=l)) \right) \right) \\ &\quad + \left(\frac{3}{6} \times \left(- \sum_{l \in \{\text{'spam'}, \text{'ham'}\}} P(t=l) \times \log_2(P(t=l)) \right) \right) \\ &= \left(\frac{3}{6} \times \left(- \left(\left(\frac{2}{3} \times \log_2\left(\frac{2}{3}\right) \right) + \left(\frac{1}{3} \times \log_2\left(\frac{1}{3}\right) \right) \right) \right) \right) \\ &\quad + \left(\frac{3}{6} \times \left(- \left(\left(\frac{1}{3} \times \log_2\left(\frac{1}{3}\right) \right) + \left(\frac{2}{3} \times \log_2\left(\frac{2}{3}\right) \right) \right) \right) \right) = 0.9183 \text{ bits} \end{aligned}$$

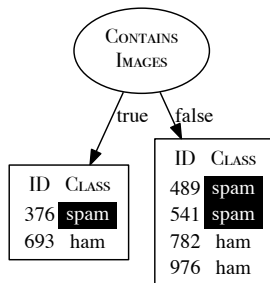
Step 2: Remainder for the CONTAINS IMAGES feature

- Calculate the **remainder** for the CONTAINS IMAGES feature in the dataset.

$$rem(d, \mathcal{D}) = \sum_{l \in levels(d)} \underbrace{\frac{|\mathcal{D}_{d=l}|}{|\mathcal{D}|}}_{\text{weighting}} \times \underbrace{H(t, \mathcal{D}_{d=l})}_{\text{entropy of partition } \mathcal{D}_{d=l}}$$

ID	SUSPICIOUS WORDS	UNKNOWN SENDER	CONTAINS IMAGES	CLASS
376	true	false	true	spam
489	true	true	false	spam
541	true	true	false	spam
693	false	true	true	ham
782	false	false	false	ham
976	false	false	false	ham

Step 2: Partitions for the CONTAINS IMAGES feature



Partitions for the CONTAINS IMAGES feature

Step 2: Remainder for the CONTAINS IMAGES feature

$$\begin{aligned} & \text{rem}(\text{IMAGES}, \mathcal{D}) \\ &= \left(\frac{|\mathcal{D}_{\text{IMAGES}=T}|}{|\mathcal{D}|} \times H(t, \mathcal{D}_{\text{IMAGES}=T}) \right) + \left(\frac{|\mathcal{D}_{\text{IMAGES}=F}|}{|\mathcal{D}|} \times H(t, \mathcal{D}_{\text{IMAGES}=F}) \right) \\ &= \left(\frac{2}{6} \times \left(- \sum_{l \in \{\text{'spam'}, \text{'ham'}\}} P(t=l) \times \log_2(P(t=l)) \right) \right) \\ &+ \left(\frac{4}{6} \times \left(- \sum_{l \in \{\text{'spam'}, \text{'ham'}\}} P(t=l) \times \log_2(P(t=l)) \right) \right) \\ &= \left(\frac{2}{6} \times \left(- \left(\left(\frac{1}{2} \times \log_2\left(\frac{1}{2}\right) \right) + \left(\frac{1}{2} \times \log_2\left(\frac{1}{2}\right) \right) \right) \right) \right) \\ &+ \left(\frac{4}{6} \times \left(- \left(\left(\frac{2}{4} \times \log_2\left(\frac{2}{4}\right) \right) + \left(\frac{2}{4} \times \log_2\left(\frac{2}{4}\right) \right) \right) \right) \right) = 1 \text{ bit} \end{aligned}$$

Step 3: Compute the Information Gain

- Calculate the **information gain** for the three descriptive feature in the dataset.

$$IG(d, \mathcal{D}) = H(t, \mathcal{D}) - \text{rem}(d, \mathcal{D})$$

Step 3: Compute the Information Gain

$$\begin{aligned}IG(\text{SUSPICIOUS WORDS}, \mathcal{D}) &= H(\text{CLASS}, \mathcal{D}) - \text{rem}(\text{SUSPICIOUS WORDS}, \mathcal{D}) \\ &= 1 - 0 = 1 \text{ bit}\end{aligned}$$

$$\begin{aligned}IG(\text{UNKNOWN SENDER}, \mathcal{D}) &= H(\text{CLASS}, \mathcal{D}) - \text{rem}(\text{UNKNOWN SENDER}, \mathcal{D}) \\ &= 1 - 0.9183 = 0.0817 \text{ bits}\end{aligned}$$

$$\begin{aligned}IG(\text{CONTAINS IMAGES}, \mathcal{D}) &= H(\text{CLASS}, \mathcal{D}) - \text{rem}(\text{CONTAINS IMAGES}, \mathcal{D}) \\ &= 1 - 1 = 0 \text{ bits}\end{aligned}$$

- The results of these calculations match our intuitions.

Contents

Decision trees

- Definitions

- Shannon's Entropy

- Information Gain

Handling continuous features

Growing a tree

Regression trees

Beyond basic decision trees

- Alternative impurity measures

- Overfitting and Tree Pruning

- Model Ensembles

Turn continuous features into boolean features

- ❑ The easiest way to handle continuous valued descriptive features is to turn them into boolean features by defining a threshold.
- ❑ The threshold is used to partition the instances based on their value of the continuous descriptive feature.
- ❑ How do we set the threshold?

Sorting the instances

1. The instances in the dataset are sorted according to the continuous feature values.
2. The adjacent instances in the ordering that have different classifications are then selected as possible threshold points.
3. The optimal threshold is found by computing the information gain for each of these classification transition boundaries and selecting the boundary with the highest information gain as the threshold.

Treat the feature as a categorical feature

- Once a threshold has been set the dynamically created new boolean feature can compete with the other categorical features for selection as the splitting feature at that node.
- This process can be repeated at each node as the tree grows.

Vegetation classification dataset



(a) chaparral veg.



(b) riparian veg.



(c) conifer veg.

Dataset for predicting the vegetation in an area with a continuous ELEVATION feature (measured in feet).

ID	STREAM	SLOPE	ELEVATION	VEGETATION
1	false	steep	3 900	chaparral
2	true	moderate	300	riparian
3	true	steep	1 500	riparian
4	false	steep	1 200	chaparral
5	false	flat	4 450	conifer
6	true	steep	5 000	conifer
7	true	steep	3 000	chaparral

Sorted instances

Dataset for predicting the vegetation in an area sorted by the continuous ELEVATION feature.

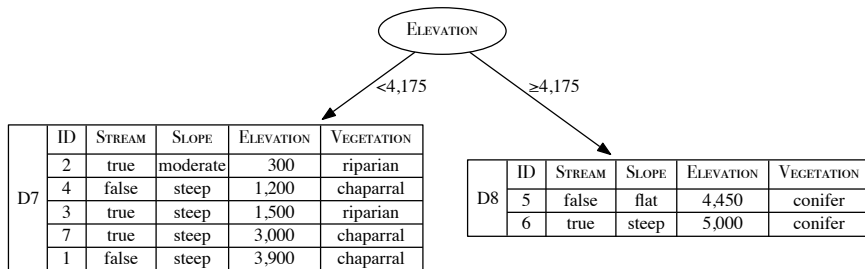
ID	STREAM	SLOPE	ELEVATION	VEGETATION
2	true	moderate	300	riparian
4	false	steep	1 200	chapparal
3	true	steep	1 500	riparian
7	true	steep	3 000	chapparal
1	false	steep	3 900	chapparal
5	false	flat	4 450	conifer
6	true	steep	5 000	conifer

Thresholds and partitions

Partition sets (Part.), entropy, remainder (Rem.), and information gain (Info. Gain) for the candidate ELEVATION thresholds: ≥ 750 , $\geq 1\ 350$, $\geq 2\ 250$ and $\geq 4\ 175$.

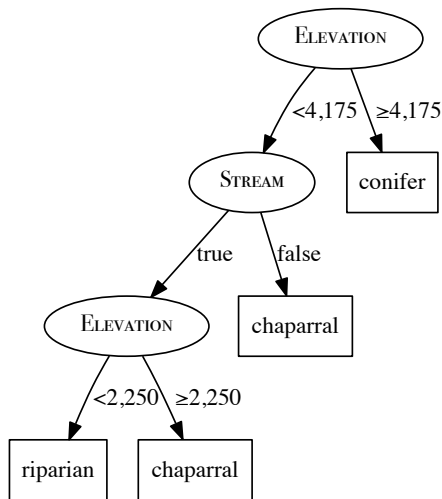
Split by Threshold	Part.	Instances	Partition Entropy	Rem.	Info. Gain
≥ 750	\mathcal{D}_1	d_2	0.0	1.2507	0.3060
	\mathcal{D}_2	$d_4, d_3, d_7, d_1, d_5, d_6$	1.4591		
$\geq 1\ 350$	\mathcal{D}_3	d_2, d_4	1.0	1.3728	0.1839
	\mathcal{D}_4	d_3, d_7, d_1, d_5, d_6	1.5219		
$\geq 2\ 250$	\mathcal{D}_5	d_2, d_4, d_3	0.9183	0.9650	0.5917
	\mathcal{D}_6	d_7, d_1, d_5, d_6	1.0		
$\geq 4\ 175$	\mathcal{D}_7	d_2, d_4, d_3, d_7, d_1	0.9710	0.6935	0.8631
	\mathcal{D}_8	d_5, d_6	0.0		

First split



The vegetation classification decision tree after the dataset has been split using $ELEVATION \geq 4\,175$.

Second split



The decision tree that would be generated for the vegetation classification dataset using information gain.

Contents

Decision trees

- Definitions

- Shannon's Entropy

- Information Gain

Handling continuous features

Growing a tree

Regression trees

Beyond basic decision trees

- Alternative impurity measures

- Overfitting and Tree Pruning

- Model Ensembles

The ID3 Algorithm

- ❑ ID3 Algorithm (Iterative Dichotomizer 3)
- ❑ Attempts to create the shallowest tree that is consistent with the data that it is given.
- ❑ The ID3 algorithm builds the tree in a recursive, depth-first manner, beginning at the root node and working down to the leaf nodes.

How the different nodes are treated?

1. The algorithm begins by choosing the best descriptive feature to test (i.e., the best question to ask first) using **information gain**.
2. A root node is then added to the tree and labelled with the selected test feature.
3. The training dataset is then partitioned using the test.
4. For each partition a branch is grown from the node.
5. The process is then repeated for each of these branches using the relevant partition of the training set in place of the full training set and with the selected test feature excluded from further testing.

The CART Algorithm

- ❑ CART (Classification and Regression Trees).
- ❑ Constructs binary trees.
- ❑ At each node it looks for the feature and threshold in that feature that provides the largest information gain at that node.
- ❑ This is the one implemented in scikit-learn.

Greedy algorithms

- Both ID3 and CART are examples of *greedy* algorithms.
- They find the “best tree” by greedily finding the best partitions at each level of the tree.
- They don't check whether the best splits done at the higher levels of the tree will lead to the lowest possible impurity several levels down.

Contents

Decision trees

- Definitions

- Shannon's Entropy

- Information Gain

Handling continuous features

Growing a tree

Regression trees

Beyond basic decision trees

- Alternative impurity measures

- Overfitting and Tree Pruning

- Model Ensembles

Regression trees

- ❑ Regression trees are constructed so as to reduce the **variance** in the set of training examples at each of the leaf nodes in the tree
- ❑ We can do this by adapting the ID3 algorithm to use a measure of variance rather than a measure of classification impurity (entropy) when selecting the best attribute

Impurity for regression

- The impurity (variance) at a node can be calculated using the following equation:

$$\text{var}(t, \mathcal{D}) = \frac{\sum_{i=1}^n (t_i - \bar{t})^2}{n - 1}$$

- We select the feature that minimizes the weighted variance across the resulting partitions:

$$\mathbf{d}[\text{best}] = \underset{d \in \mathbf{d}}{\operatorname{argmin}} \sum_{l \in \text{levels}(d)} \frac{|\mathcal{D}_{d=l}|}{|\mathcal{D}|} \times \text{var}(t, \mathcal{D}_{d=l})$$

A dataset listing the number of bike rentals per day

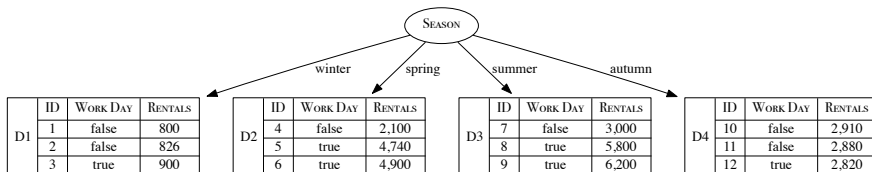
ID	SEASON	WORK DAY	RENTALS
1	winter	false	800
2	winter	false	826
3	winter	true	900
4	spring	false	2 100
5	spring	true	4 740
6	spring	true	4 900
7	summer	false	3 000
8	summer	true	5 800
9	summer	true	6 200
10	autumn	false	2 910
11	autumn	false	2 880
12	autumn	true	2 820

Weighted variance according to the partitions

Partitioning of the dataset for bike rentals based on SEASON and WORK DAY features and the computation of the weighted variance for each partitioning.

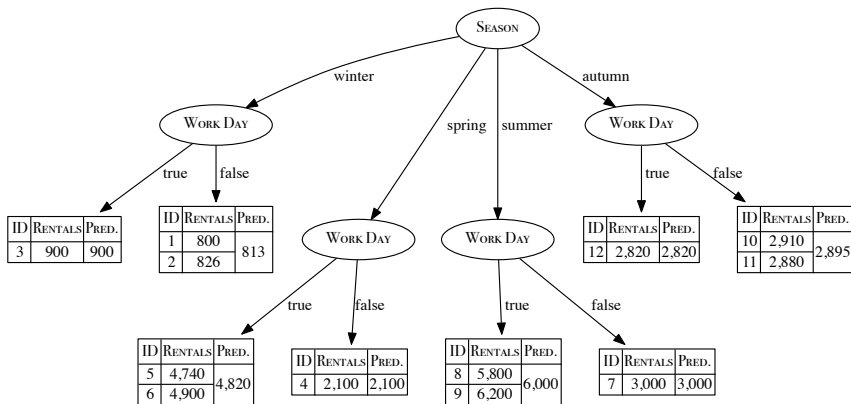
Split by Feature	Level	Part.	Instances	$\frac{ \mathcal{D}_{d=l} }{ \mathcal{D} }$	$var(t, \mathcal{D})$	Weighted Variance
SEASON	'winter'	\mathcal{D}_1	$\mathbf{d}_1, \mathbf{d}_2, \mathbf{d}_3$	0.25	2 692	$1\,379\,331\frac{1}{3}$
	'spring'	\mathcal{D}_2	$\mathbf{d}_4, \mathbf{d}_5, \mathbf{d}_6$	0.25	$2\,472\,533\frac{1}{3}$	
	'summer'	\mathcal{D}_3	$\mathbf{d}_7, \mathbf{d}_8, \mathbf{d}_9$	0.25	3 040 000	
	'autumn'	\mathcal{D}_4	$\mathbf{d}_{10}, \mathbf{d}_{11}, \mathbf{d}_{12}$	0.25	2 100	
WORK DAY	'true'	\mathcal{D}_5	$\mathbf{d}_3, \mathbf{d}_5, \mathbf{d}_6, \mathbf{d}_8, \mathbf{d}_9, \mathbf{d}_{12}$	0.50	$4\,026\,346\frac{1}{3}$	$2\,551\,813\frac{1}{3}$
	'false'	\mathcal{D}_6	$\mathbf{d}_1, \mathbf{d}_2, \mathbf{d}_4, \mathbf{d}_7, \mathbf{d}_{10}, \mathbf{d}_{11}$	0.50	1 077 280	

Resulting decision tree



The decision tree resulting from splitting the data using the feature SEASON.

Final decision tree



The final decision tree induced from the dataset. To illustrate how the tree generates predictions, this tree lists the instances that ended up at each leaf node and the prediction (PRED.) made by each leaf node.

Contents

Decision trees

- Definitions

- Shannon's Entropy

- Information Gain

Handling continuous features

Growing a tree

Regression trees

Beyond basic decision trees

- Alternative impurity measures

- Overfitting and Tree Pruning

- Model Ensembles

Contents

Decision trees

- Definitions

- Shannon's Entropy

- Information Gain

Handling continuous features

Growing a tree

Regression trees

Beyond basic decision trees

- Alternative impurity measures**

- Overfitting and Tree Pruning

- Model Ensembles

Gini index

- Another commonly used measure of impurity is the **Gini index**:

$$Gini(t, \mathcal{D}) = 1 - \sum_{l \in levels(t)} P(t = l)^2$$

- The Gini index can be thought of as calculating how often you would misclassify an instance in the dataset if you classified it based on the probability distribution of the target feature in the dataset.
- Information gain can be calculated using the Gini index by replacing the entropy measure with the Gini index.

Information gain ratio

- Entropy based information gain has a preference for features with many values.
- One way of addressing this issue is to use **information gain ratio** which is computed by dividing the information gain of a feature by the amount of information used to determine the value of the feature:

$$GR(d, \mathcal{D}) = \frac{IG(d, \mathcal{D})}{-\sum_{l \in levels(d)} (P(d = l) \times \log_2(P(d = l)))}$$

Contents

Decision trees

- Definitions

- Shannon's Entropy

- Information Gain

Handling continuous features

Growing a tree

Regression trees

Beyond basic decision trees

- Alternative impurity measures

- Overfitting and Tree Pruning

- Model Ensembles

Why overfitting?

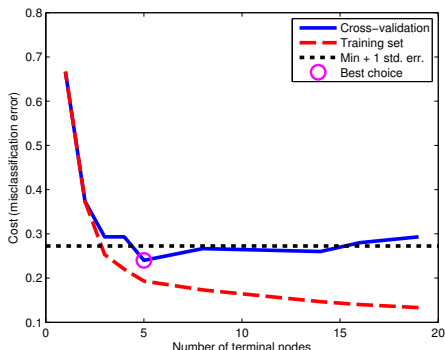
The likelihood of over-fitting occurring increases as a tree gets deeper because the resulting classifications are based on smaller and smaller subsets as the dataset is partitioned after each feature test in the path.

Pruning strategies

- **Pre-pruning:** stop the recursive partitioning early. Pre-pruning is also known as **forward pruning**.
- We can stop growing the tree if the decrease in the error is not sufficient to justify the extra complexity of adding an extra subtree.
- **Post-pruning:** allow the algorithm to grow the tree as much as it likes and then prune the tree of the branches that cause over-fitting.

Common **Post**-pruning Approach

- Using the validation set evaluate the prediction accuracy achieved by both the fully grown tree and the pruned copy of the tree.
- If the pruned copy of the tree performs no worse than the fully grown tree the node is a candidate for pruning.



Contents

Decision trees

- Definitions

- Shannon's Entropy

- Information Gain

Handling continuous features

Growing a tree

Regression trees

Beyond basic decision trees

- Alternative impurity measures

- Overfitting and Tree Pruning

- Model Ensembles

Definition

- Rather than creating a single model, we can generate a set of models and then make predictions by aggregating the outputs of these models.
- A prediction model that is composed of a set of models is called a **model ensemble**.
- In order for this approach to work the models that are in the ensemble must be different from each other.

Approaches

There are two standard approaches to creating ensembles:

1. **boosting**
2. **bagging.**

Boosting: How does it work?

- ❑ Boosting works by iteratively creating models and adding them to the ensemble.
- ❑ The iteration stops when a predefined number of models have been added.
- ❑ When we use **boosting** each new model added to the ensemble is biased to pay more attention to instances that previous models miss-classified.
- ❑ This is done by incrementally adapting the dataset used to train the models. To do this we use a **weighted dataset**

Boosting: Weighted Dataset

- Each instance has an associated weight $\mathbf{w}_i \geq 0$,
- Initially set to $\frac{1}{n}$ where n is the number of instances in the dataset.
- After each model is added to the ensemble it is tested on the **training data** and the weights of the instances the model gets correct are decreased and the weights of the instances the model gets incorrect are increased.
- These weights are used as a distribution over which the dataset is sampled to create a **replicated training set**, where the replication of an instance is proportional to its weight.

Predictions

- Once the set of models have been created the ensemble makes **predictions** using a weighted aggregate of the predictions made by the individual models.
- The weights used in this aggregation are simply the confidence factors associated with each model.

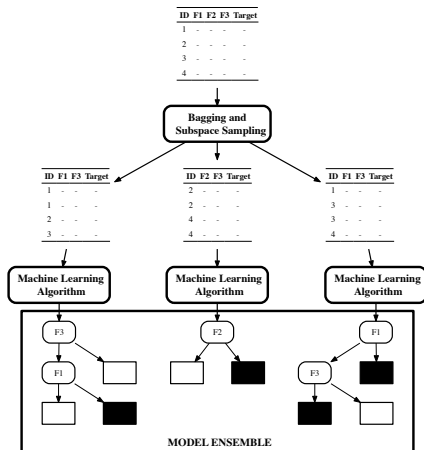
Bagging: Definition

- When we use **bagging** (or **bootstrap aggregating**) each model in the ensemble is trained on a random sample of the dataset known as **bootstrap samples**.
- Each random sample is the same size as the dataset and **sampling with replacement** is used.
- Consequently, every bootstrap sample will be missing some of the instances from the dataset so each bootstrap sample will be different and this means that models trained on different bootstrap samples will also be different

Bagging: Random forest

- When bagging is used with decision trees each bootstrap sample only uses a randomly selected subset of the descriptive features in the dataset. This is known as **subspace sampling**.
- The combination of bagging, subspace sampling, and decision trees is known as a **random forest** model.

Example of using bagging



The process of creating a model ensemble using bagging and subspace sampling.