

DECISION TREES IN PYTHON

Shashwati Shradha

Contents

1	Introduction	1
1.1	Representation of data	1
1.2	Supervised Learning	1
2	Creating a Decision Tree	1
2.1	Split the Data Set	1
2.2	Recursive Partitioning	1
3	Types of Decision Trees	2
3.1	Classification Trees	2
3.2	Regression Trees	2
4	Classification Tree for Titanic Data set	2
4.1	Decription	2
4.2	Code	3
4.3	Summary	4
5	Splitting criterion for Classification Trees	4
5.1	Entropy	4
5.2	Gini Index	5
6	Regression Tree for ... Data set	5
6.1	Description	5
6.2	Code	5
6.3	Summary	5
7	Splitting criterion for Regression Trees	5
7.1	Mean Square Error	5
8	Accuracy of the Model	5
8.1	K-Folds cross validation	5
8.2	Confusion Matrix	5
8.3	Generalization Error	5
8.4	Bias - Variance trade off	5
8.5	Pruning	5
8.6	Occam Razor	5
9	Conclusion	5
9.1	Advantages	5
9.2	Disadvantages	5

1 Introduction

Decision trees are widely used data structures for data compression and prediction. They are a set of rules used to predict the behavior of a data set, assuming that the similar behavior continues.

The purpose of this document is to summarize decision trees. We will focus on different types of decision trees, ways of measuring the accuracy and improving the model. We will also be using scikit-learn package in Python to code decision trees.

1.1 Representation of data

The goal of classification is to find a better representation of the data. Decision trees, unlike linear models are able to work with non-linear data. They partition the data into rectangular regions in such a way that the concentration of a particular type of class is maximized in that region, along with minimizing impurity.

1.2 Supervised Learning

Decision trees is a form of supervised learning. This means that the model learns through examples. When we are given a data set, we divide it into two groups: Train data and Test data. We train the model using the train data set and hope to get the least predictive error in the test data set.

2 Creating a Decision Tree

2.1 Split the Data Set

The first step in creating a decision tree is to split the given data set into a train data and a test data. We will use the train data to make a model and the accuracy of our model will be tested on the test data.

2.2 Recursive Partitioning

The algorithm used to make the decision trees is called the Top-Down Induction of Decision Trees (TDIDT). The principles guiding this algorithm are:

- The tree is generated using recursive partitioning.
- The same attribute cannot be chosen twice.
- Adequacy Condition. No two instances with the same attributes belong to different classes

According the TDIDT algorithm, for each attribute, we place the instances of the training data into one of the values of the attributes. We repeat this until we have reached a pure class.

3 Types of Decision Trees

3.1 Classification Trees

Classification trees are used to separate data sets into two or more discrete classes. The response variables are categorical in nature. The attributes can be categorical or numeric in nature. To decide the output at a leaf node, we use majority vote.

3.2 Regression Trees

Regression trees are used when the response variables are numeric or continuous. The attributes can be categorical or numeric in nature. To decide the output at a leaf node, we use the average.

4 Classification Tree for Titanic Data set

4.1 Description

- **Pclass**: Passenger Class (1- 1st; 2- 2nd; 3- 3rd)
- **survival**: Survival (0- No; 1- Yes)
- **name**: Name
- **sex**: Sex
- **age**: Age
- **sibsp**: Number of Siblings/Spouses Aboard
- **parch**: Number of Parents/Children Aboard
- **ticket**: Ticket Number
- **fare**: Passenger Fare (British pound)
- **cabin**: Cabin
- **embarked** Port of Embarkation (C- Cherbourg; Q- Queenstown; S- Southampton)
- **boat**: Lifeboat
- **body**: Body Identification Number
- **home.dest**: Home/Destination

4.2 Code

```
from sklearn import preprocessing
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score

import pandas as pd
import numpy as np

dataset = './titanic.csv'
use = ["pclass", "sex", "age", "sibsp", "parch", "embarked", "survived"]

## 1- Getting data
titanic = pd.read_csv(dataset, usecols = use)

## 2- Preprocessing
clean = preprocessing.LabelEncoder()
titanic["sex"] = clean.fit_transform(titanic["sex"])
titanic["embarked"] = clean.fit_transform(titanic["embarked"].astype(str))
titanic[use] = titanic[use].replace(np.NaN, 0)

## 3- Making features and target attributes"
features = titanic[["pclass", "sex", "age", "sibsp", "parch", "embarked"]]
target = titanic["survived"]

## 4- Making the train and test set (80-20)
feat_train, feat_test, tar_train, tar_test =
train_test_split(features, target, test_size=0.20, random_state=2)

## 5- Building the decision tree
clf = DecisionTreeClassifier(max_depth = 6, random_state=1)

clf = clf.fit(feat_train, tar_train)

## 6- Making predictions
y_pred = clf.predict(feat_test)

## 7- Getting accuracy
accuracy = accuracy_score(y_pred, tar_test)

print("Accuracy: ", accuracy)
```

4.3 Summary

The titanic data set contains all the information people who were on board of titanic. The data collected will be used to predict whether a person would have survived on the ship or not. We do this using python's machine learning package scikit learn. We import the necessary packages, import the data and focus on the features variable to construct the model. We train our model using 80% of the data and use the 20% to test our accuracy. Some points to remember:

- **max_depth** can be used to avoid overfitting
- **min_samples_leaf** can be used to avoid overfitting too
- **random_state** is used to make sure our random set is reproducible
- **test_size** should be changed to test the model

5 Splitting criterion for Classification Trees

5.1 Entropy

Entropy is used to get the uncertainty in the dataset.

5.2 Gini Index

6 Regression Tree for ... Data set

6.1 Description

6.2 Code

6.3 Summary

7 Splitting criterion for Regression Trees

7.1 Mean Square Error

8 Accuracy of the Model

8.1 K-Folds cross validation

8.2 Confusion Matrix

8.3 Generalization Error

8.4 Bias - Variance trade off

8.5 Pruning

8.6 Occam Razor

9 Conclusion

9.1 Advantages

- Data compression
- Easy to interpret. It behaves like a white box.
- Does not require attributes to be scaled (e.g. Standardized)
- It identifies non-linear attributes

9.2 Disadvantages