

Changelog: 14 October 2016, 30 October 2016

Decision Trees

Divide and Conquer, Loss Function, Inductive Bias

Lecture 3: Part 1

Marina Santini

Acknowledgements
Slides inspired by Daumé III (2015: 10-18)

2016 Lecture 3: Decision Trees - Part 1 1

Lecture 3: Required Reading

- Handout
- Daumé III (2015: 10-18)
- Witten et al. (2011: 99-108)

- Optional – Self Study: Witten et al. (2011: 192-203)

2016 Lecture 3: Decision Trees - Part 1 2

Outline

- Greedness
- Divide and Conquer
- Inductive Bias of the Decision Tree
- Loss Function
- Expected Loss
- Empirical Error
- Induction

2016 Lecture 3: Decision Trees - Part 1 3

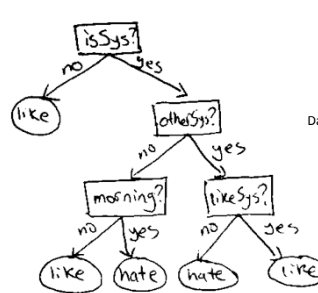
Learning: Generalization-ability

- Predicting the future based on the past

Sepal length	Sepal width	Petal length	Petal width	Species
5.1	3.5	1.4	0.2	<i>I. setosa</i>
4.9	3.0	1.4	0.2	<i>I. setosa</i>
4.7	3.2	1.3	0.2	<i>I. setosa</i>
4.6	3.1	1.5	0.2	<i>I. setosa</i>
5.0	3.6	1.4	0.2	<i>I. setosa</i>
5.4	3.9	1.7	0.4	<i>I. setosa</i>
4.8	3.4	1.4	0.3	<i>I. setosa</i>
5.0	3.4	1.5	0.2	<i>I. setosa</i>
4.4	2.9	1.4	0.2	<i>I. setosa</i>
4.9	3.1	1.5	0.1	<i>I. setosa</i>
5.4	3.7	1.5	0.2	<i>I. setosa</i>
4.8	3.4	1.6	0.2	<i>I. setosa</i>

2016 Lecture 3: Decision Trees - Part 1 4

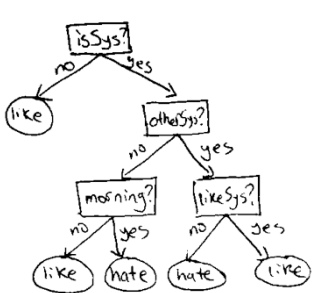
Predict whether a student will like a course



Daumé III (2015)

2016 Lecture 3: Decision Trees - Part 1 5

Training Data



Rating	Easy?	AI?	Sys?	Thy?	Morning?
+2	y	y	n	y	n
+2	y	y	n	y	n
+2	n	y	n	n	n
+2	n	n	n	y	n
+2	n	y	y	n	y
+1	y	y	n	n	n
+1	n	y	n	y	n
0	n	n	n	n	y
0	y	n	n	y	y
0	n	y	n	y	n
0	y	y	y	y	y
-1	y	y	y	n	y
-1	n	n	y	y	n
-1	n	n	y	n	y
-1	y	n	y	n	y
-2	n	n	y	y	n
-2	n	y	y	n	y
-2	y	n	y	n	n
-2	y	n	y	n	y

2016 Lecture 3: Decision Trees - Part 1 6

Problem Representation

- Questions = Features
- Answers = Feature Values
- Ratings = Class Labels

- An example is a set of feature values.
- Training data is a set of examples associated with class labels.

2016 Lecture 3: Decision Trees - Part 1 7

"Greedy" Model: the most useful feature

- Histogram
- Root node

2016 Lecture 3: Decision Trees - Part 1 8

Divide & Conquer

- Divide:
 - Partition the data into 2 parts:
 - Yes part vs No part
- Conquer:
 - Recurse and run the Divide routine

2016 Lecture 3: Decision Trees - Part 1 9

The end of the cycle

... when it becomes useless to query on additional features

2016 Lecture 3: Decision Trees - Part 1 10

Decision Tree: Inductive Bias

The goal of the decision tree learning model is:

- to figure out what questions to ask
- in what order
- what answer to predict once you have asked enough questions

The inductive bias of decision trees:

- The things that we want to learn to predict are more like the root node and less like the other branch nodes.

2016 Lecture 3: Decision Trees - Part 1 11

Decision Tree: Informal Definition

A decision tree is:

- a flow-chart-like structure, where:
 - each internal (non-leaf) node denotes a test on an attribute
 - each branch represents the outcome of a test
 - each leaf (or terminal) node holds a class label

The topmost node in a tree is the root node.

2016 Lecture 3: Decision Trees - Part 1 12

Formalising the learning problem (1)

1. The loss function:

Classification: **zero/one loss** $\ell(y, \hat{y}) = \begin{cases} 0 & \text{if } y = \hat{y} \\ 1 & \text{otherwise} \end{cases}$

2016 Lecture 3: Decision Trees - Part 1 13

Formalising the learning problem (2)

- Data Generating Distribution

$$D(x, y)$$

2016 Lecture 3: Decision Trees - Part 1 14

Expected Loss

- The loss function
- The data generating distribution

$$\epsilon \triangleq \mathbb{E}_{(x,y) \sim D}[\ell(y, f(x))] = \sum_{(x,y)} D(x,y) \ell(y, f(x))$$

2016 Lecture 3: Decision Trees - Part 1 15

How to read formulas: Expected Value

$\epsilon \triangleq \mathbb{E}_{(x,y) \sim D}[\ell(y, f(x))]$

- ϵ = epsilon
- \triangleq = equal by definition to (or: is defined as)
- \mathbb{E} = blackboard-bold E
- (x,y) = sub the pair \mathbf{xy}
- $\sim D$ = overscript D
- $\ell(y, f(x))$ = l of the pair \mathbf{y} f of \mathbf{x}

$= \sum_{(x,y)} D(x,y) \ell(y, f(x))$

= sum over all the pairs \mathbf{xy} in and \mathbf{y} times l of \mathbf{y} and f of \mathbf{x}

2016 Lecture 3: Decision Trees - Part 1 16

Training Error

The training error is the average error over the training data

$$\hat{\epsilon} \triangleq \frac{1}{N} \sum_{n=1}^N \ell(y_n, f(x_n))$$

How to read: the training error epsilon-hat is equal by definition to 1 over N of the Sum from n=1 to capital N of "l" of y and f of x.

2016 Lecture 3: Decision Trees - Part 1 17

Empirical Error (Alpaydin, 2014: 14)

The empirical error is the proportion of training instances where the predictions of h (the hypothesis = the informed guess) do not match the required values given in X (the training set). The error of the the hypothesis h given the training set X is:

$$E(h|X) = \sum_{t=1}^N 1(h(x^t) \neq r^t)$$

2016 Lecture 3: Decision Trees - Part 1 18

Given:

- a loss function l and
- a sample \mathbf{d} from some unknown distribution \mathbf{D}
- compute a function f that has low expected error ϵ over \mathbf{D} with respect to l .

2016

Lecture 3: Decision Trees - Part 1

19

Quiz 1: Training Error

How would you define a training error on a dataset:

1. Training error is the average loss over the training sample
2. Training error is the expected prediction error over an independent test sample
3. None of the above

2016

Lecture 3: Decision Trees - Part 1

20

Quiz 2: Distributions

$$\sum_{(x,y)} \mathcal{D}(x,y) \ell(y, f(x))$$

What kind of distribution is \mathcal{D} in the formula above?

1. Normal
2. Unknown
3. None of the above

2016

Lecture 3: Decision Trees - Part 1

21

Quiz 3: Loss function

How would you define a loss function?

1. The loss function $l(\text{actual value}, \text{predicted value})$ characterizes how bad predictions are.
2. The loss function is an unknown distribution.
3. Both definitions are incorrect.

2016

Lecture 3: Decision Trees - Part 1

22

The end

2016

Lecture 3: Decision Trees - Part 1

23