

Assignment Project Exam Help

Tree Learning

<https://eduassistpro.github.io>

Add WeChat Last revision: 21 Mar 2018 edu_assist_pro

Material derived from slides for the book

“Machine Learning” by T. Mitchell
McGraw-Hill (1997)

<http://www-2.cs.cmu.edu/~tom/mlbook.html>

Material derived from

<http://www-2.cs.cmu.edu/~tom/mlbook.html>

Material derived from slides by Eibe Frank

<http://www.cs.waikato.ac.nz/ml/weka/>

Material derived from slides for the book

“Machine Learning” by P. Flach

Cambridge University Press (2012)

<http://cs.bris.ac.uk/~flach/mlbook>

Aims

This lecture will enable you to describe decision tree learning, the use of entropy and the problem of overfitting. Following it you should be able to:

- define the decision tree representation
- list representation properties of data and models for which decision tree
- repr (T
- define entropy in the context of learning a Boolean examples
- describe the inductive bias of the basic TDIDT algorithm
- define overfitting of a training set by a hypothesis
- describe developments of the basic TDIDT algorithm: pruning, rule generation, numerical attributes, many-valued attributes, costs, missing values
- describe regression and model trees

Brief History of Decision Tree Learning Algorithms

- late 1950's – Bruner et al. in psychology work on modelling concept acquisition
- early 1960's – *Lea* ncept
- late 1960's – *Lea* ncept
- early 1990s – ID3 adds features, develops into C4. “default” machine learning algorithm
- late 1990s – C5.0, commercial version of C4.5 (available at www.rulequest.com)
- current – widely available and applied; influential techniques

<https://eduassistpro.github.io>

Add WeChat [edu_assist_pro](https://eduassistpro.github.io)

Why use decision trees?

Assignment Project Exam Help

- Decision trees are probably the single most popular data mining tool

• <https://eduassistpro.github.io>

- There are some drawbacks, though — e.g., high variance
- They do *classification*, i.e., predict a categorical output from categorical and/or real inputs, or *regression*

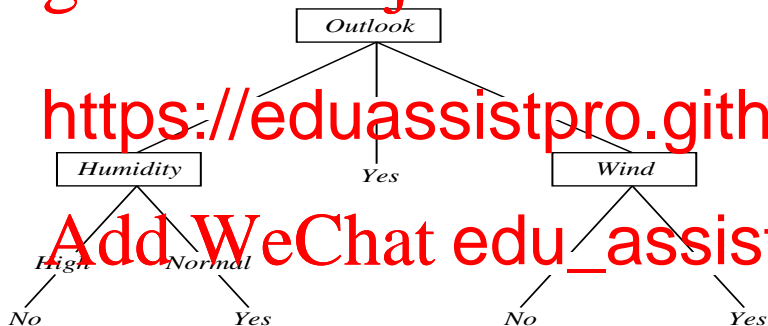
Add WeChat edu_assist_pro

Decision Tree for *PlayTennis*

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr



A Tree to Predict C-Section Risk

Assignment Project Exam Help

Learned from medical records of 1000 women

Negative

```
[833+,167-] .8
Fetal_Presen
| Previous_Ces
| | Primiparous = 0: []
| | Primiparous = 1: [368+,68-] .84+ .16-
| | | Fetal_Distress = 0: [334+,47-] .88+ .12-
| | | | Birth_Weight < 3349: [201+,10.6-] .95+ .05-
| | | | Birth_Weight > 3349: [335+,35.4-] .78+ .22-
| | | Fetal_Distress = 1: [34+,21-] .62+ .38-
| Previous_Cesction = 1: [55+,35-] .61+ .39-
Fetal_Presentation = 2: [3+,29-] .11+ .89-
Fetal_Presentation = 3: [8+,22-] .27+ .73-
```

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Decision Tree for Credit Rating

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Decision Tree for Fisher's Iris data

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Decision Trees

Assignment Project Exam Help

Decision tree representation:

- Each internal node tests an attribute
- Each
- Each

<https://eduassistpro.github.io>

How would we represent the following expressions ?

- \wedge, \vee, XOR
- $(A \wedge B) \vee (C \wedge \neg D \wedge E)$
- M of N

Add WeChat edu_assist_pro

Decision Trees

Assignment Project Exam Help

 $X \wedge Y$

```

X = t:
| Y = t: true
| Y = f: no
X = f: no

```

<https://eduassistpro.github.io>
 $X \vee Y$

```

X = t: true
X = f:
| Y = t: true
| Y = f: no

```

Add WeChat edu_assist_pr

Decision Trees

Assignment Project Exam Help

2 of 5

```

X = t:
| Y = t: true
| Y = f:
| | Z = t: true
| | Z = f: false
X = f:
| Y = t:
| | Z = t: true
| | Z = f: false
| Y = f: false

```

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

So in general decision trees represent a *dis* of constraints on the attributes values of instances.

When are Decision Trees the Right Model?

Assignment Project Exam Help

- With Boolean values for the instances X and class Y , the representation adopted by decision-trees allows us to represent Y as a Boolean function of X .
- Given a Boolean function Y , we can find a decision tree that represents it. The tree assigns $Y = 1$ to some subset of the values for X , and $Y = 0$ to the rest.
- Any Boolean function can be trivially represented by a decision tree. So, for each combination of values with $Y = 1$, have a path from root to a leaf with $Y = 1$. All other leaves have $Y = 0$.

<https://eduassistpro.github.io>

Add WeChat [edu_assist_pro](#)

When are Decision Trees the Right Model?

- This is nothing but a re-representation of the truth-table, and will have 2^d leaves. More compact trees may be possible, by taking into account what is common between one or more rows with the same Y value
- But (https://eduassistpro.github.io/ Add WeChat edu_assist_pro are examples)
- In general, although possible in principle to express any function, our search and prior restrictions may not find the correct tree in practice.
- BUT: If you want readable models that combine logical tests with a probability-based decision, then decision trees are a good start

When to Consider Decision Trees?

- Instances described by a mix of numeric features and discrete attribute-value pairs
- Target function is discrete valued (otherwise use regression trees)
- Disj
- Pos
- Inte

<https://eduassistpro.github.io>

Examples are extremely numerous including:

- Equipment or medical diagnosis
- Credit risk analysis
- Modeling calendar scheduling preferences
- etc.

Add WeChat [edu_assist_pro](#)

Top-Down Induction of Decision Trees (TDIDT)

Assignment Project Exam Help

Main loop

- A the “best” decision attribute for next *node*
- Assign
- For each
- Sort training examples to leaf nodes
- If training examples perfectly classified, Then S
new leaf nodes

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Essentially this is the “ID3” algorithm (Quinlan, 198
symbolic Machine Learning algorithm.

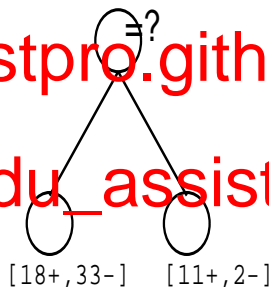
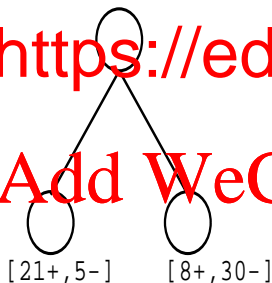
Which attribute is best?

Assignment Project Exam Help

[2

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr



You are watching a set of independent random samples of X

$P(X = 1) = \frac{1}{4}$	https://edupassistpro.github.io
--------------------------	---

You transmit data over a binary serial link. You can assist

0100001001001110110011111100...

Fewer Bits

Assignment Project Exam Help

Someone

$$P(X) = \frac{1}{2} \left(P(X=0) + P(X=1) \right) = \frac{1}{2} (0.5 + 0.5) = 0.5$$

It's possible ...

... to invent a *coding* for your transmission that has a lower average per symbol. How?

Fewer Bits

Someone tells you that the probabilities are not equal

$$P(X = A) = \frac{1}{2} \quad P(X = B) = \frac{1}{4} \quad P(X = C) = \frac{1}{8} \quad P(X = D) = \frac{1}{8}$$

It's possible

... to invent

symbol on average. How?

Add WeChat edu_assist_pr

A	0
B	10
C	110
D	111

(This is just one of several ways)

Fewer Bits

Assignment Project Exam Help

Suppose there are three equally likely values

	–		–		$\frac{1}{3}$
--	---	--	---	--	---------------

Here's a

<https://eduassistpro.github.io>

A	00
B	01
C	10

Add WeChat [edu_assist_pro](#)

Can you think of a coding that would need only 1.6 bits per symbol on average ?

Fewer Bits

Suppose there are three equally likely values

$$P(X = A) = \frac{1}{3} \quad P(X = B) = \frac{1}{3} \quad P(X = C) = \frac{1}{3}$$

Using the s
per symbol

<https://eduassistpro.github.io>

A	0
B	10
C	11

Add WeChat [edu_assist_pro](#)

This gives us, on average $\frac{1}{3} \times 1$ bit for A and $2 \times \frac{1}{3} \times 2$ bits for B and C, which equals $\frac{5}{3} \approx 1.6$ bits.

Is this the best we can do ?

Fewer Bits

Assignment Project Exam Help

Suppose t

<https://eduassistpro.github.io>

From information theory, the optimal number of bits to encode a symbol with probability p is $-\log_2 p \dots$

So the best we can do for this case is $-\log_2 \frac{1}{3} = \log_2 3$, or 1.5849625007211563 bits per symbol

Add WeChat edu_assist_pro

General Case

Assignment Project Exam Help

Suppose X can have one of m values $\dots V_1, V_2, \dots, V_m$

$P(X = V_1) = p_1$	$P(X = V_2) = p_2$	\dots	$P(X = V_m) = p_m$
--------------------	--------------------	---------	--------------------

What's the
needed to

<https://eduassistpro.github.io>

tion? It's

$$H(X) = -p_1 \log_2 p_1 - p_2 \log_2 p_2$$

Add WeChat edu_assist_pro

$$= -\sum_{j=1}^m p_j \log_2 p_j$$

$H(X)$ = the *entropy* of X

General Case

Assignment Project Exam Help

“High enr
“Low entr

<https://eduassistpro.github.io>

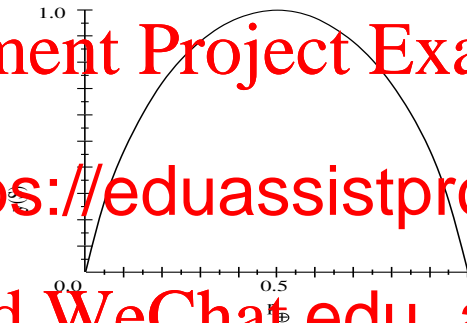
Add WeChat edu_assist_pr

Entropy

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro



Where:

S is a sample of training examples

p_+ is the proportion of positive examples in S

p_- is the proportion of negative examples in S

Entropy

Assignment Project Exam Help

Entropy

<https://eduassistpro.github.io>

A “pure” sample is one in which all examples are of the same

Add WeChat edu_assist_pr

Entropy

Entropy(S) = expected number of bits needed to encode class (\oplus or \ominus) of randomly drawn member of S (under the optimal, shortest-length code)

Why ?

Information
having pr

<https://eduassistpro.github.io>

So, expected number of bits to encode \oplus in S :

$$p_{\oplus}(-\log_2 p_{\oplus}) + p_{\ominus}(-\log_2 p_{\ominus})$$

$$\text{Entropy}(S) \equiv -p_{\oplus} \log_2 p_{\oplus} - p_{\ominus} \log_2 p_{\ominus}$$

Information Gain

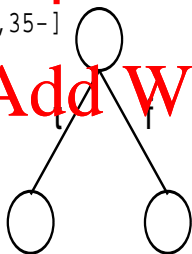
- $\text{Gain}(S, A) = \text{expected reduction in entropy due to sorting on } A$

Assignment Project Exam Help

$$\text{Gain}(S, A) = \text{Entropy}(S) - \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

<https://eduassistpro.github.io>

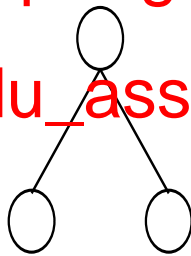
[29+, 35-]



[21+, 5-]

[8+, 30-]

Add WeChat edu_assist_pro



[18+, 33-]

[11+, 2-]

Assignment Project Exam Help

$$Gain(S, \frac{S_t}{S_f} \text{entropy}(S_f))$$

<https://eduassistpro.github.io>

$$\begin{aligned} & \frac{38}{64} \log_2 \left(\frac{38}{64} \right) - \left(\frac{8}{64} \log_2 \left(\frac{8}{64} \right) + \frac{3}{64} \log_2 \left(\frac{3}{64} \right) \right) \\ &= 0.9936 - (0.2869 + 0.1464) \\ &= 0.2658 \end{aligned}$$

Assignment Project Exam Help

<https://eduassistpro.github.io>

= 0.1643

Add WeChat edu_assist_pr

Assignment Project Exam Help

So we choose <https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Training Examples

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

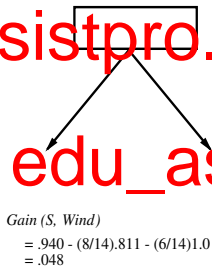
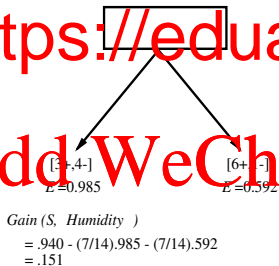
Information gain once more

Assignment Project Exam Help

Which attribute is the best classifier?

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

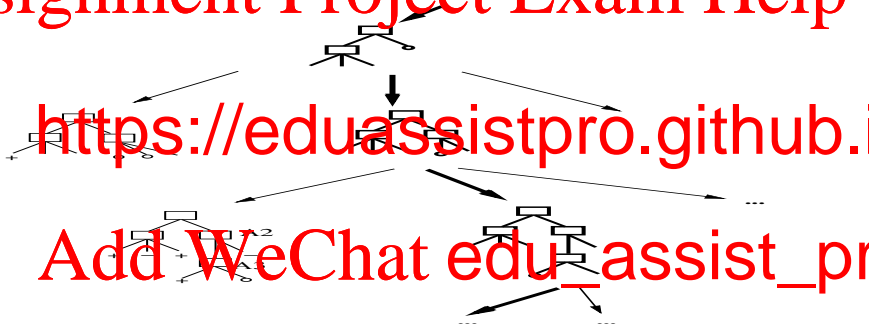


Hypothesis Space Search by ID3

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr



Hypothesis Space Search by ID3

Assignment Project Exam Help

- This can be viewed as a graph-search problem
 - Each vertex in the graph is a decision tree
 - Suppose T is a decision tree, T' is a child of T , and all the features in T' are in T
 - A pair of trees T and T' differ in just the following way: one of the leaf-nodes in T has been replaced by a non-leaf node testing a feature that appeared earlier (and it leaves)
- This is the full space of all decision trees (is it?). We want for a single tree or a small number of trees in this space. How should we do this?

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Hypothesis Space Search by ID3

- Assignment Project Exam Help
- Usual graph search technique: greedy or beam search, starting with the vertex corresponding to the "empty tree" (single leaf node)
 - Greedy search
 - <https://eduassistpro.github.io/>
 - Most of the calculation will cancel out: so, we will only need local computation at the leaf that was converted
 - RESULT: set of trees with (reasonably) high probability given D : we can now use these to answer questions like $P(y' = \omega_1 | \dots)$? or even make a *decision* or a *classification* that $y' = \omega_1$, given input data x
- Add WeChat edu_assist_pro

Hypothesis Space Search by ID3

Assignment Project Exam Help

- Hypothesis space is complete! (contains all finite discrete-valued functions w.r.t attributes)
- Out
- <https://eduassistpro.github.io>
- No back tracking
 - Local minima
- Statistically based search choices
 - Robust to noisy data...
- Inductive bias: approx “prefer shortest tree”

Add WeChat edu_assist_pro

Inductive Bias in ID3

Note H is the power set of instances X

Assignment Project Exam Help

→ Unbiased?

Not really

- Pre attr
- Bias is a *preference* for some hypotheses *tion* of hypothesis space H
- an incomplete search of a complete hypothesis space *Add WeChat edu_assist_pr*
complete search of an incomplete hypothesis space (as in learning conjunctive concepts)
- Occam's razor: prefer the shortest hypothesis that fits the data

Occam's Razor

Assignment Project Exam Help

William of Occam (c. 1287–1347)

Entities should not be multiplied beyond necessity

Why prefer

<https://eduassistpro.github.io>

Argument in favour:

- Fewer short hypotheses than long hypotheses
- a short hyp that fits data unlikely to be coincidence
- a long hyp that fits data might be coincidence

Add WeChat edu_assist_pro

Occam's Razor

Assignment Project Exam Help

Argument opposed:

- The

- <https://eduassistpro.github.io>

- What's so special about small sets based on *size* of hypothesis??

Look back to linear classification lecture to see how to make
using Minimum Description Length (MDL)

Add WeChat [edu_assist_pro](#)

Why does overfitting occur?

- Greedy search can make mistakes. We know that it can end up in local minima — so a sub-optimal choice earlier might result in a better solution later (i.e. pick a test whose posterior gain (or information gain) is less than the best one

- But the error on this data

- We will see why this is the case later (lectures on Ev
- Suppose we have two models h_1 and h_2 and optimism e_1 and e_2 . Let the true error be $E_1 = e_1 + o_1$ and $E_2 = e_2 + o_2$
- If $e_1 < e_2$ and $E_1 > E_2$, then we will say that h_1 has overfit then training data

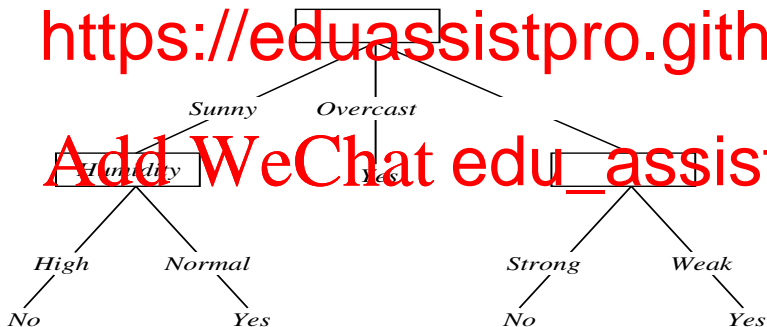
- So, a search method based purely on training data estimates may end overfitting the training data

Overfitting in Decision Tree Learning

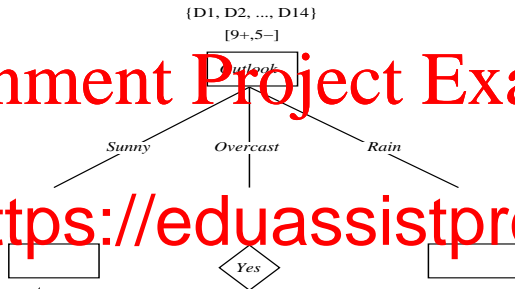
Consider adding noisy training example #15:

Assignment Project Exam Help

What effect



Overfitting in Decision Tree Learning



Which attribute should be tested here?

$$S_{\text{sunny}} = \{D1, D2, D8, D9, D11\}$$

$$\text{Gain}(S_{\text{sunny}}, \text{Humidity}) = .970 - (3/5) 0.0 - (2/5) 0.0 = .970$$

$$\text{Gain}(S_{\text{sunny}}, \text{Temperature}) = .970 - (2/5) 0.0 - (2/5) 1.0 - (1/5) 0.0 = .570$$

$$\text{Gain}(S_{\text{sunny}}, \text{Wind}) = .970 - (2/5) 1.0 - (3/5) .918 = .019$$

Overfitting in General

Consider error of hypothesis / overfitting

Assignment Project Exam Help

- training data: $error_{train}(h)$
- enti

Definiti

Hypothe

hypothesis $h' \in H$ such that

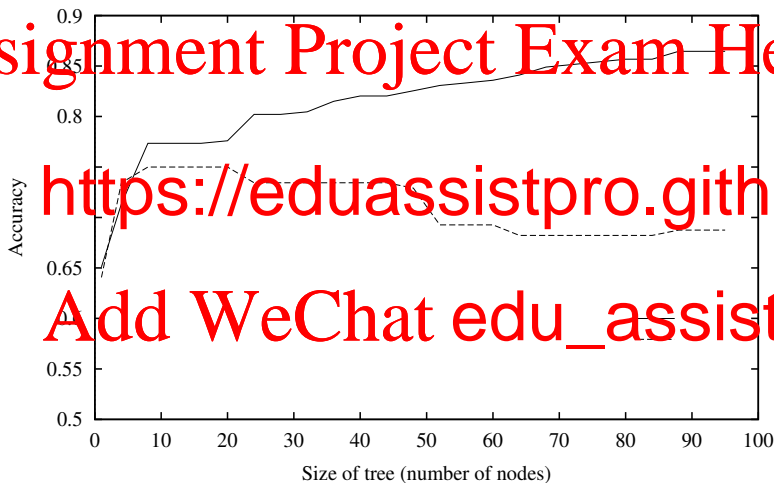
<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

and

$$error_{\mathcal{D}}(h) > error_{\mathcal{D}}(h')$$

Overfitting in Decision Tree Learning



Avoiding Overfitting

Assignment Project Exam Help

- **pre-pruning** stop growing when data split not statistically significant

- **pos**

ove

<https://eduassistpro.github.io>

Post-pru

How to select “best” tree:

- Measure performance over training data?
- Measure performance over separate validation
- MDL: minimize $size(tree) + size(misclassifications(tree))$?

Avoiding Overfitting

Assignment Project Exam Help

Pre-pruning

- Can
- Stop
- For example, in ID3: chi-squared test plus information gain procedure
 - only statistically significant attributes were

<https://eduassistpro.github.io>

Add WeChat [edu_assist_pro](#)

Avoiding Overfitting

Assignment Project Exam Help

Pre-pruning

- Sim low
- In C
- In sklearn, this parameter is `min_samples`
- In sklearn, the parameter `min_impurity` stopping when the this falls below a lower-bound

<https://eduassistpro.github.io>

Add WeChat [edu_assist_pro](#)

Avoiding Overfitting

Assignment Project Exam Help

Early stopping

- Pre-pruning may suffer from early stopping: may stop the growth of tree
- Cla
 - Target structure only visible in fully expanded
 - Prepruning won't expand the root node
- But: XOR-type problems not common in practice
- And: pre-pruning faster than post-pruning

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Avoiding Overfitting

Post-pruning

- Builds full tree first and prunes it afterward
 - Attribute interactions are visible in fully-grown tree
- Pro
- effective
- Two
 - Subtree replacement
 - Subtree raising
- Possible strategies: error estimation, significance principle
- We examine two methods: Reduced-error Pruning and Error-based Pruning

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Reduced-Error Pruning

Assignment Project Exam Help

Split data i

Do until fur

- Eva
(plus those below it)
- Greedily remove the one that most improves accuracy

<https://eduassistpro.github.io>

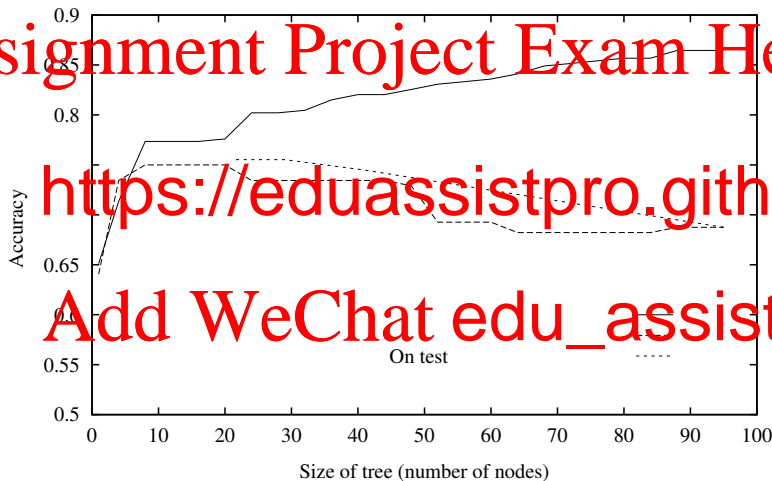
Add WeChat edu_assist_pro

Assignment Project Exam Help

- **Go** <https://eduassistpro.github.io>
- **Not so good** reduces effective size of training set

Add WeChat edu_assist_pr

Effect of Reduced-Error Pruning



Error-based pruning (C4.5 / J48 / C5.0)

Assignment Project Exam Help

Quinlan (1993) describes the successor to ID3 – C4.5

- many extensions – see below
- pos
- incl
- also: pruning by converting tree to rules
- commercial version – C5.0 – is widely used
 - RuleQuest.com
 - now free
- Weka version – J48 – also widely used

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Pruning operator: Sub-tree replacement

Bottom-up:

tree is considered for replacement once all its sub-trees have been considered

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Error-based pruning: error estimate

Assignment Project Exam Help

Goal is to im
from train

But how ca

Make the estimate of error **pessimistic**!

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Error-based pruning: error estimate

Assignment Project Exam Help

- App

- C4.5

deri

<https://eduassistpro.github.io>

- Standard Bernoulli-process-based method
- **Note:** statistically motivated, but not statis
- **However** works well in practice

Add WeChat edu_assist_pr

Error-based pruning: error estimate

Assignment Project Exam Help

- The error estimate for a tree node is the weighted sum of error estimates for all its subtrees (possibly leaves).
- Upp

<https://eduassistpro.github.io>

- f is actual (empirical) error of tree on examples a
- N is the number of examples at the tree node
- Z_c is a constant whose value depends on c
- C4.5's default value for confidence $c = 0.25$
- If $c = 0.25$ then $Z_c = 0.69$ (from standardized normal distribution)

Add WeChat edu_assist_pro

Error-based pruning: error estimate

Assignment Project Exam Help

- Ho
- Wh
- As
- See example on next slide (note: values not calculated using the above formula)

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Error-based pruning: error estimate

Assignment Project Exam Help

- health plan contribution: node 46

<https://eduassistpro.github.io/>

- f

Add WeChat edu_assist_pro

- sub-tr 0.51
- sub-trees estimated to give *greater* error so prune away

Rule Post-Pruning

Assignment Project Exam Help

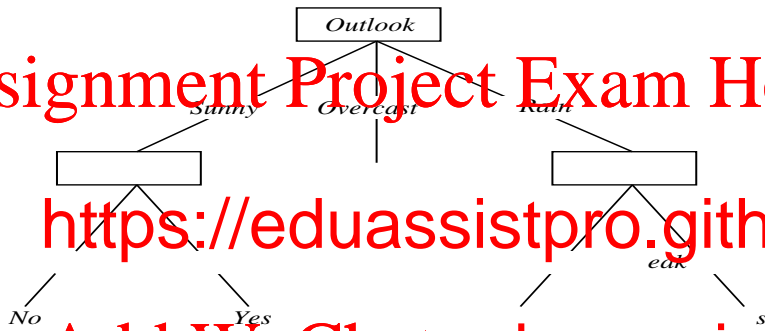
This method was introduced in Quinlan's C4.5

- Con
- Pru
- Sort final rules into desired sequence for use

For: simpler classifiers, people prefer rules to trees

Against: does not scale well, slow for large trees & data

Converting A Tree to Rules



IF $(\text{Outlook} = \text{Sunny}) \wedge (\text{Humidity} = \text{High})$

THEN $\text{PlayTennis} = \text{No}$

IF $(\text{Outlook} = \text{Sunny}) \wedge (\text{Humidity} = \text{Normal})$

THEN $\text{PlayTennis} = \text{Yes}$

Rules from Trees (Rule Post-Pruning)

Assignment Project Exam Help

Rules can be simpler than trees but just as accurate, e.g., in C4.5Rules:

- pat
- can s
 - i.e., rules can be generalized while maintaining accuracy
- greedy rule simplification algorithm
 - drop the condition giving lowest estimated error
 - continue while estimated error does not increase

<https://eduassistpro.github.io>
Add WeChat: edu_assist_pro

Rules from Trees

Assignment Project Exam Help

Select a 'good' subset of rules within a class (C4.5Rules):

- goal: remove rules not useful in terms of accuracy
- find a
- trad
- stoc

<https://eduassistpro.github.io>

Sets of rules can be ordered by class (C4.5Rules):

- order classes by increasing chance of making errors
- set as a default the class with the most training instances by any rule

Add WeChat edu_assist_pro

Continuous Valued Attributes

Assignment Project Exam Help

Decision trees originated for discrete attributes only. Now: continuous attributes.

Can creat

- T_e
- $(T_e$

<https://eduassistpro.github.io>

- Usual method: continuous attributes have a bin
- Note:
 - discrete attributes – one split exhausts all values
 - continuous attributes – can have many splits in a tree

Add WeChat edu_assist_pro

Continuous Valued Attributes

- Splits evaluated on all possible split points
- More computation: $n - 1$ possible splits for n values of an attribute in training set

- Fay

<https://eduassistpro.github.io>

$$\frac{(48+60)}{2} \text{ and } \frac{(80+90)}{2}$$

- Choose best split point by info gain (or evaluation o
- Note: C4.5 uses actual values in data

<i>Temperature:</i>	40	48	60	72	80	90
<i>PlayTennis:</i>	No	No	Yes	Yes	Yes	No

Axis-parallel Splitting

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Fitting data that is not a good “match” to the possible splits in a tree.

“Pattern Classification” Duda, Hart, and Stork, (2001)

Splitting on Linear Combinations of Features

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Reduced tree size by allowing splits that are a better “match” to the data.

“Pattern Classification” Duda, Hart, and Stork, (2001)

Attributes with Many Values

Assignment Project Exam Help

Problem:

- If att
- Wh
-
- Imagine using *Date* = March 21, 2018 as a
- High gain on training set, useless for prediction

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Attributes with Many Values

Assignment Project Exam Help

One approach: use *GainRatio* instead

<https://eduassistpro.github.io>

$$SplitInformation(S, A) \equiv - \sum_{i=1}^c \frac{|S_i|}{|S|} \log_2 \frac{|S_i|}{|S|}$$

Add WeChat edu_assist_pro

where S_i is subset of S for which A has val i

Attributes with Many Values

Assignment Project Exam Help

Why does this help ?

- sen
- act
-
- therefore higher for many-valued attributes, e
uniformly distributed across possible values

<https://eduassistpro.github.io>
Add WeChat edu_assist_pr

Attributes with Costs

Assignment Project Exam Help

Consider

- me
- rob

<https://eduassistpro.github.io>

How to learn a consistent tree with low expected cost?

Add WeChat edu_assist_pro

Attributes with Costs

Assignment Project Exam Help

One appr

- Exa

<https://eduassistpro.github.io>

$\overline{Cost($

Preference for decision trees using lower cost attrib

Add WeChat edu_assist_pr

Attributes with Costs

Assignment Project Exam Help

Also: class (misclassification) costs, instance costs, ...

SEE5 /

<https://eduassistpro.github.io>

Can give *false positives* a different cost to *false negatives*

Forces a different tree structure to be learned to minimize misclassification costs – can help if class distribution is skewed

Add WeChat [edu_assist_pro](#)

Unknown Attribute Values

Assignment Project Exam Help

What if some examples missing values of A ?

Use training example anyway, sort through tree. Here are 3 possible approaches

- If no other example
- assign most common value of A among target value
- assign probability p_i to each possible value
 - assign fraction p_i of example to each desc

Note: need to classify new (unseen) examples in same fashion

Windowing

Early implementations – training sets too large for memory

Assignment Project Exam Help

As a solution ID3 implemented *windowing*:

1. select subs
2. construct
3. use tree to cla
4. if all instances correctly classified then halt, else
5. add selected misclassified instances to the window
6. go to step 2

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Windowing retained in C4.5 because it can lead to *more accurate* trees.
Related to *ensemble learning*.

Non-linear Regression with Trees

Assignment Project Exam Help

Despite some nice properties of Neural Networks, such as generalization to deal sensibly with unseen input patterns and robustness to losing neurons (predicti problem

- Back computing time; may have to be partitioned into separate modules that can be trained independently, e.g. NetTall
- Neural Networks are not very *transparent* representation of what has been learned

Possible solution: exploit success of tree-structured approaches in ML

Regression trees

Assignment Project Exam Help

- Differences to decision trees:

-

-

- <https://eduassistpro.github.io>

- Can approximate piecewise constant function
- Easy to interpret
- More sophisticated version: model trees

Add WeChat edu_assist_pro

A Regression Tree and its Prediction Surface

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

“Elements of Statistical Learning” Hastie, Tibshirani & Friedman (2001)

Regression Tree on sine dataset

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Regression Tree on CPU dataset

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Tree learning as variance reduction

- The variance of a Boolean (i.e., Bernoulli) variable with success probability \hat{p} is $\hat{p}(1 - \hat{p})$, which is half the Gini index. So we could interpret the goal of tree learning as minimising the class variance (or standard deviation, in case of $\sqrt{\text{Gini}}$) in the leaves.
- In re

<https://eduassistpro.github.io>

If a split partitions the set of target values into sets $\{Y_1, \dots, Y_l\}$, the weighted average variance

$$\text{Var}(\{Y_1, \dots, Y_l\}) = \sum_{j=1}^l \frac{|Y_j|}{|Y|} \text{Var}(Y_j) = \dots = \frac{1}{|Y|} \sum_{y \in Y} y^2 - \sum_{j=1}^l \frac{|Y_j|}{|Y|} \bar{y}_j^2$$

The first term is constant for a given set Y and so we want to maximise the weighted average of squared means in the children.

Learning a regression tree

Imagine you are a collector of vintage Hammond tonewheel organs. You have been monitoring an online auction site from which you collected some data about interesting transactions:

3.	A100	good	no	1051
4.	T202	good	no	270
5.	M102	good	yes	870
6.	A100	excellent	no	1770
7.	T202	fair	no	99
8.	A100	good	yes	1900
9.	E112	fair	no	77

Learning a regression tree

From this data, you want to construct a regression tree that will help you determine a reasonable price for your next purchase.

There are three features, hence three possible splits:

Model = [A100, B3, E112, M102, T202]

Condition

Leslie = [yes, no] [625, 870, 1900][77, 99, 2

The means of the first split are 1574, 4513, 77, 870 and 331, a weighted average of squared means is 321

split are 3142, 1023 and 267, with weighted average of sq

$2.68 \cdot 10^6$; for the third split the means are 1132 and 1297, with weighted average of squared means $1.55 \cdot 10^6$. We therefore branch on Model at the top level. This gives us three single-instance leaves, as well as three A100s and three T202s.

Learning a regression tree

For the A100s we obtain the following splits:

Condition = [excellent, good, fair] [1770][1051, 1900][]

Leslie = [

Without g

results in l

variance e

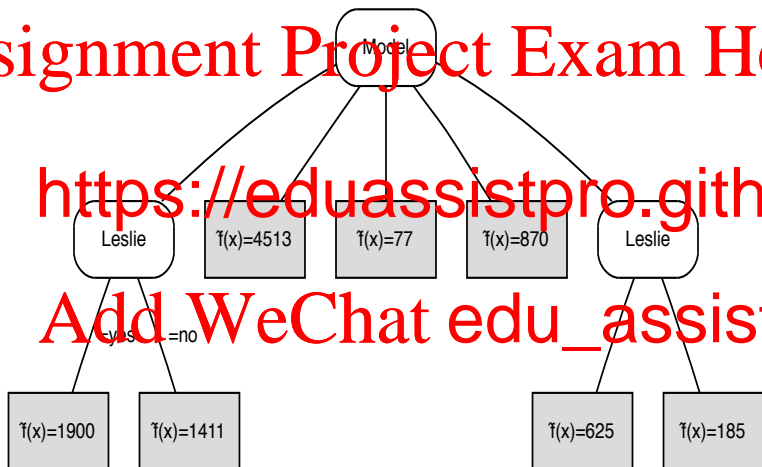
follows:

Condition = [excellent, good, fair] [] [270][99]

Leslie = [yes, no] [625][99, 270]

Again we see that splitting on Leslie gives tighter clusters of values. The learned regression tree is depicted on the next slide.

A regression tree



A regression tree learned from the Hammond organ dataset.

Model trees

Assignment Project Exam Help

- Like regression trees but with linear regression functions at each node
- Linear
- Only has to be linear in a small subset of attributes
 - Attributes occurring in subtree (+maybe at the root)
- Fast: overhead for Linear Regression (LR) not large as only a small subset of attributes is used in tree

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

Two uses of features

Assignment Project Exam Help

Suppose

$$-1 \leq x \leq 1.$$

A linear ap

$y = 0$. Ho

$0 \leq x \leq$

interval. We can achieve this by using x b

a regression variable (next slide).

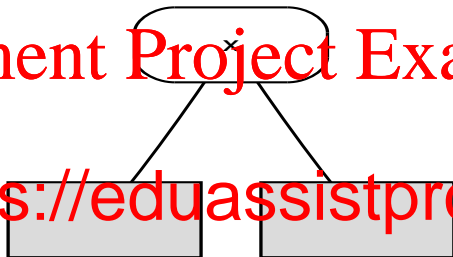
<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

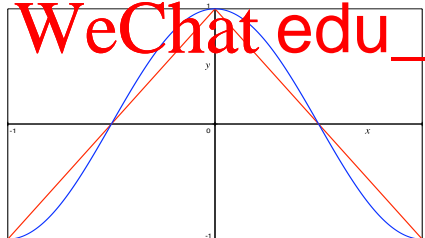
A small model tree

Assignment Project Exam Help

<https://eduassistpro.github.io>



Add WeChat edu_assist_pro



Model Tree on CPU dataset

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Smoothing

- Naïve prediction method – output value of LR model at corresponding leaf node
- Improve performance by *smoothing* predictions with *internal* LR model

<https://eduassistpro.github.io>

- Sm

$$\frac{p' + q}{n + k}$$

- p' prediction passed up to next higher node
- q prediction passed to this node from below
- q value predicted by model at this node
- n number of instances that reach node below
- k smoothing constant

- Same effect can be achieved by incorporating the internal models into the leaf nodes

Building the tree

Assignment Project Exam Help

- Splitting criterion: *standard deviation reduction*

<https://eduassistpro.github.io>

where T_1, T_2, \dots are the sets from splits of data

- Termination criteria (important when building prediction)
 - Standard deviation becomes smaller than ce training set (e.g. 5%)
 - Too few instances remain (e.g. less than four)

Pruning the tree

- Pruning is based on estimated absolute error of LR models
- Heuristic estimate:

<https://eduassistpro.github.io>

where

the number of parameters in the linear model

- LR models are pruned by greedily removing term estimated error
- Model trees allow for heavy pruning: often a single LR model can replace a whole subtree
- Pruning proceeds bottom up: error for LR model at internal node is compared to error for subtree

Discrete (nominal) attributes

Assignment Project Exam Help

- Nominal attributes converted to binary attributes and treated as numerical

• <https://eduassistpro.github.io>

- the i th binary attribute is 0 if an instance's value is one of the first i in the ordering, 1 otherwise
- Best binary split with original attribute provably on one of the new attributes

Add WeChat edu_assist_pro

Summary – decision trees

- Decision tree learning is a practical method for many classifier learning tasks – still a “Top 10” data mining algorithm – see `sklearn.tree.DecisionTreeClassifier`
- TDIDT family descended from ID3 searches complete hypothesis space
- Use not t
- Overfitting is inevitable with an expressive hypothesis data, so pruning is important
- Decades of research into extensions and refined approach, e.g., for numerical prediction, logical trees
- Often the “try-first” machine learning method in applications, illustrates many general issues
- Performance can be improved with use of “ensemble” methods

Summary – regression and model trees

- Regression trees were introduced in CART — R's implementation is close to CART, but see `sklearn.tree.DecisionTreeRegressor` for a basic version

- Qui

- M5' is bas

- Quinlan also investigated combining instanc

- CUBIST: Quinlan's rule learner for numeric pr
www.rulequest.com

- Interesting comparison: Neural nets vs. model trees — both do *non-linear regression*
- other methods also can learn non-linear models