# Machine Learning: Decision Trees in AIMA and WEKA

# UCI

## Machine Learning Repository
### Center for Machine Learning and Intelligent Systems

About   Citation Policy   Donate a Data Set

Contact

○ Repository ○ Web

Google™

**View ALL Data Sets**

# Zoo Data Set
*Download*: Data Folder, Data Set Description

**Abstract**: Artificial, 7 classes of animals

## http://archive.ics.uci.edu/ml/datasets/Zoo

| Data Set Characteristics: | Multivariate | Number of Instances: | 101 | Area: | Life |
|---|---|---|---|---|---|
| Attribute Characteristics: | Categorical, Integer | Number of Attributes: | 17 | Date Donated | 1990-05-15 |
| Associated Tasks: | Classification | Missing Values? | No | Number of Web Hits: | 18038 |

animal name: string

hair: Boolean

feathers: Boolean

eggs: Boolean

milk: Boolean

airborne: Boolean

aquatic: Boolean

predator: Boolean

toothed: Boolean

backbone: Boolean

breathes: Boolean

venomous: Boolean

fins: Boolean

legs: {0,2,4,5,6,8}

tail: Boolean

domestic: Boolean

catsize: Boolean

type: {mammal, fish, bird, shellfish, insect, reptile, amphibian}

# Zoo data

## 101 examples

aardvark,1,0,0,1,0,0,1,1,1,1,0,0,4,0,0,1,mammal
antelope,1,0,0,1,0,0,0,1,1,1,0,0,4,1,0,1,mammal
bass,0,0,1,0,0,1,1,1,1,0,0,1,0,1,0,0,fish
bear,1,0,0,1,0,0,1,1,1,1,0,0,4,0,0,1,mammal
boar,1,0,0,1,0,0,1,1,1,1,0,0,4,1,0,1,mammal
buffalo,1,0,0,1,0,0,0,1,1,1,0,0,4,1,0,1,mammal
calf,1,0,0,1,0,0,0,1,1,1,0,0,4,1,1,1,mammal
carp,0,0,1,0,0,1,0,1,1,0,0,1,0,1,1,0,fish
catfish,0,0,1,0,0,1,1,1,1,0,0,1,0,1,0,0,fish
cavy,1,0,0,1,0,0,0,1,1,1,0,0,4,0,1,0,mammal
cheetah,1,0,0,1,0,0,1,1,1,1,0,0,4,1,0,1,mammal
chicken,0,1,1,0,1,0,0,0,1,1,0,0,2,1,1,0,bird
chub,0,0,1,0,0,1,1,1,1,0,0,1,0,1,0,0,fish
clam,0,0,1,0,0,0,1,0,0,0,0,0,0,0,0,0,shellfish
crab,0,0,1,0,0,1,1,0,0,0,0,0,4,0,0,0,shellfish
...

# Zoo example

aima-python> python

>>> from learning import *

>>> zoo

<DataSet(zoo): 101 examples, 18 attributes>

>>> dt = DecisionTreeLearner()

>>> dt.train(zoo)

>>> dt.predict(['shark',0,0,1,0,0,1,1,1,1,0,0,1,0,1,0,0]) #eggs=1
'fish'

>>> dt.predict(['shark',0,0,0,0,0,1,1,1,1,0,0,1,0,1,0,0]) #eggs=0
'mammal'

# Zoo example

>> dt.dt

DecisionTree(13, 'legs', {0: DecisionTree(12, 'fins', {0: DecisionTree(8, 'toothed', {0: 'shellfish', 1: 'reptile'}), 1: DecisionTree(3, 'eggs', {0: 'mammal', 1: 'fish'})}), 2: DecisionTree(1, 'hair', {0: 'bird', 1: 'mammal'}), 4: DecisionTree(1, 'hair', {0: DecisionTree(6, 'aquatic', {0: 'reptile', 1: DecisionTree(8, 'toothed', {0: 'shellfish', 1: 'amphibian'})}), 1: 'mammal'}), 5: 'shellfish', 6: DecisionTree(6, 'aquatic', {0: 'insect', 1: 'shellfish'}), 8: 'shellfish'})

# Zoo example

```
>>> dt.dt.display()
Test legs
 legs = 0 ==> Test fins
    fins = 0 ==> Test toothed
       toothed = 0 ==> RESULT =  shellfish
       toothed = 1 ==> RESULT =  reptile
    fins = 1 ==> Test eggs
       eggs = 0 ==> RESULT =  mammal
       eggs = 1 ==> RESULT =  fish
 legs = 2 ==> Test hair
    hair = 0 ==> RESULT =  bird
    hair = 1 ==> RESULT =  mammal
 legs = 4 ==> Test hair
    hair = 0 ==> Test aquatic
       aquatic = 0 ==> RESULT =  reptile
       aquatic = 1 ==> Test toothed
          toothed = 0 ==> RESULT =  shellfish
          toothed = 1 ==> RESULT =  amphibian
    hair = 1 ==> RESULT =  mammal
 legs = 5 ==> RESULT =  shellfish
 legs = 6 ==> Test aquatic
    aquatic = 0 ==> RESULT =  insect
    aquatic = 1 ==> RESULT =  shellfish
 legs = 8 ==> RESULT =  shellfish
```
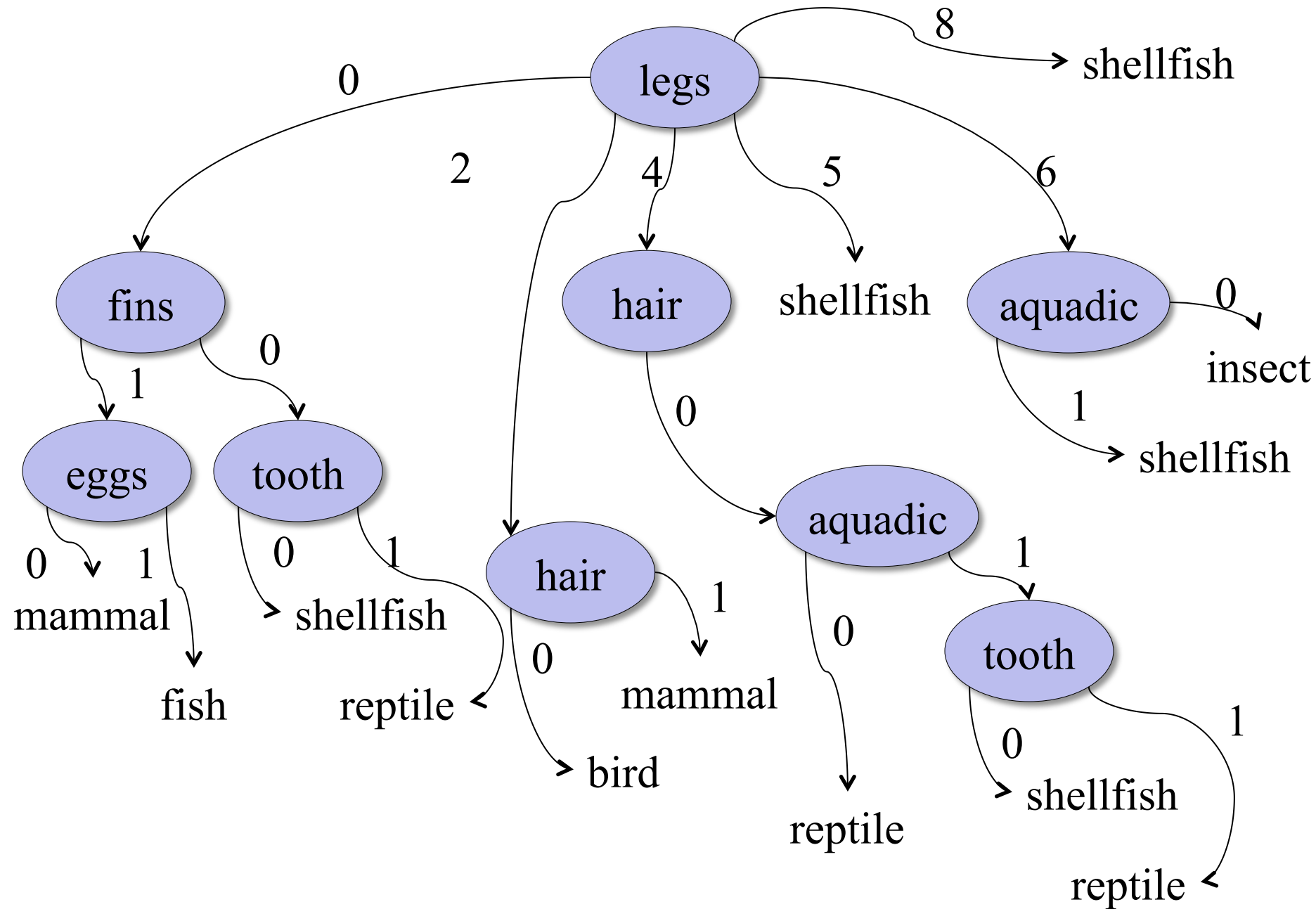
# Zoo example

```
>>> dt.dt.display()
Test legs
 legs = 0 ==> Test fins
    fins = 0 ==> Test toothed
       toothed = 0 ==> RESULT =  shellfish
       toothed = 1 ==> RESULT =  reptile
    fins = 1 ==> Test milk
      milk = 0 ==> RESULT =  fish
      milk = 1 ==> RESULT =  mammal
 legs = 2 ==> Test hair
    hair = 0 ==> RESULT =  bird
    hair = 1 ==> RESULT =  mammal
 legs = 4 ==> Test hair
    hair = 0 ==> Test aquatic
      aquatic = 0 ==> RESULT =  reptile
      aquatic = 1 ==> Test toothed
        toothed = 0 ==> RESULT =  shellfish
        toothed = 1 ==> RESULT =  amphibian
    hair = 1 ==> RESULT =  mammal
 legs = 5 ==> RESULT =  shellfish
 legs = 6 ==> Test aquatic
    aquatic = 0 ==> RESULT =  insect
    aquatic = 1 ==> RESULT =  shellfish
 legs = 8 ==> RESULT =  shellfish
```
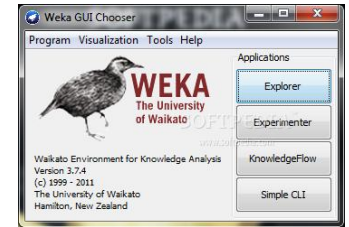
Add the shark example to the training set and retrain

# Weka



- Open-source Java machine learning tool
- [http://www.cs.waikato.ac.nz/ml/weka/](http://www.cs.waikato.ac.nz/ml/weka/)
- Implements many classifiers & ML algorithms
- Uses common data representation format, making comparisons easy
- Comprehensive set of data pre-processing tools and evaluation methods
- Three modes of operation: GUI, command line, Java API

**Weka GUI Chooser**

Program  Visualization  Tools  Help

WEKA
The University
of Waikato

**Applications**

Explorer

Experimenter

KnowledgeFlow

Workbench

Simple CLI

Waikato Environment for Knowledge Analysis
Version 3.8.0
(c) 1999 – 2016
The University of Waikato
Hamilton, New Zealand

---

**Weka Explorer**

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

**Classifier**

Choose    J48 -C 1.0 -M 0

**Test options**

- ● Use training set
- ○ Supplied test set    Set...
- ○ Cross-validation    Folds  10
- ○ Percentage split    %  66

More options...

(Nom) WillWait

Start    Stop

**Result list (right-click for options)**

20:32:20 – trees.J48
20:32:38 – trees.J48
20:32:40 – trees.J48
20:33:06 – trees.J48
20:44:28 – trees.J48

**Classifier output**

```
J48 pruned tree
------------------

HowCrowded = None: No (2.0)
HowCrowded = Some: Yes (4.0)
HowCrowded = Full
|   Hungry = Yes
|   |   IsFridayOrSaturday = Yes
|   |   |   Price = $: Yes (2.0)
|   |   |   Price = $$: Yes (0.0)
|   |   |   Price = $$$: No (1.0)
|   |   IsFridayOrSaturday = No: No (1.0)
|   Hungry = No: No (2.0)

Number of Leaves  :      7

Size of the tree :      11


Time taken to build model: 0.11 seconds

=== Evaluation on training set ===
```

**Status**

OK

Log    x 0

# Common .arff* data format

@relation heart-disease-simplified

@attribute age numeric &larr; *Numeric attribute*
@attribute sex { female, male }
@attribute chest_pain_type { typ_angina, asympt, non_anginal, atyp_angina}
@attribute cholesterol numeric
@attribute exercise_induced_angina {no, yes}
@attribute class {present, not_present} &larr; *Nominal attribute*

@data &larr; *Training data*
63,male,typ_angina,233,no,not_present
67,male,asympt,286,yes,present
67,male,asympt,229,yes,present
38,female,non_anginal,?,no,not_present
...

*ARFF = Attribute-Relation File Format

# Weka demo

# Open the Weka GUI

# Load the restaurant .arff data

# Select J48 tree classifier

# Click Start to train
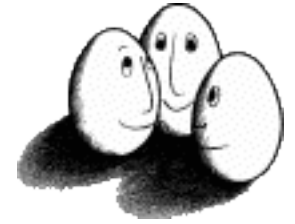
# See the training results

# Compare results

HowCrowded = None: No (2.0)

HowCrowded = Some: Yes (4.0)

HowCrowded = Full

|   Hungry = Yes

|   |   IsFridayOrSaturday = Yes

|   |   |   Price = $: Yes (2.0)

|   |   |   Price = $$: Yes (0.0)

|   |   |   Price = $$$: No (1.0)

|   |   IsFridayOrSaturday = No: No (1.0)

|   Hungry = No: No (2.0)



**J48 pruned tree: nodes:11; leaves:7, max depth:4**

**ID3 tree: nodes:12; leaves:8, max depth:4**

# Weka vs. svm_light vs. …

- Weka: good for experimenting with different ML algorithms
- Other tools are much more efficient &scalable
- Scikit-learn is a popular suite of open-source machine-learning tools in Python
  - Built on NumPy, SciPy, and matplotlib for efficiency
  - Use anaconda or do pip install scikit-learn
- For SVMs many use svm_light