## 01. Load zoo data. Observe the attributes and their values



## 02. Build C4.5 decision tree

## 03. Visualize the output



Classification accuracy : 92.079%

```
=== Detailed Accuracy By Class ===

                 TP Rate  FP Rate  Precision  Recall  F-Measure  MCC    ROC Area  PRC Area  Class
                 1.000    0.000    1.000      1.000   1.000      1.000  1.000     1.000     mammal
                 1.000    0.011    0.929      1.000   0.963      0.958  0.994     0.929     fish
                 1.000    0.000    1.000      1.000   1.000      1.000  1.000     1.000     bird
                 0.800    0.033    0.727      0.800   0.762      0.735  0.986     0.812     invertebrate
                 0.625    0.032    0.625      0.625   0.625      0.593  0.920     0.677     insect
                 0.750    0.000    1.000      0.750   0.857      0.862  0.872     0.760     amphibian
                 0.600    0.010    0.750      0.600   0.667      0.656  0.793     0.420     reptile
Weighted Avg.    0.921    0.008    0.922      0.921   0.920      0.914  0.976     0.908

=== Confusion Matrix ===

  a  b  c  d  e  f  g   <-- classified as
 41  0  0  0  0  0  0 |  a = mammal
  0 13  0  0  0  0  0 |  b = fish
  0  0 20  0  0  0  0 |  c = bird
  0  0  0  8  2  0  0 |  d = invertebrate
  0  0  0  3  5  0  0 |  e = insect
  0  0  0  0  0  3  1 |  f = amphibian
  0  1  0  0  1  0  3 |  g = reptile
```
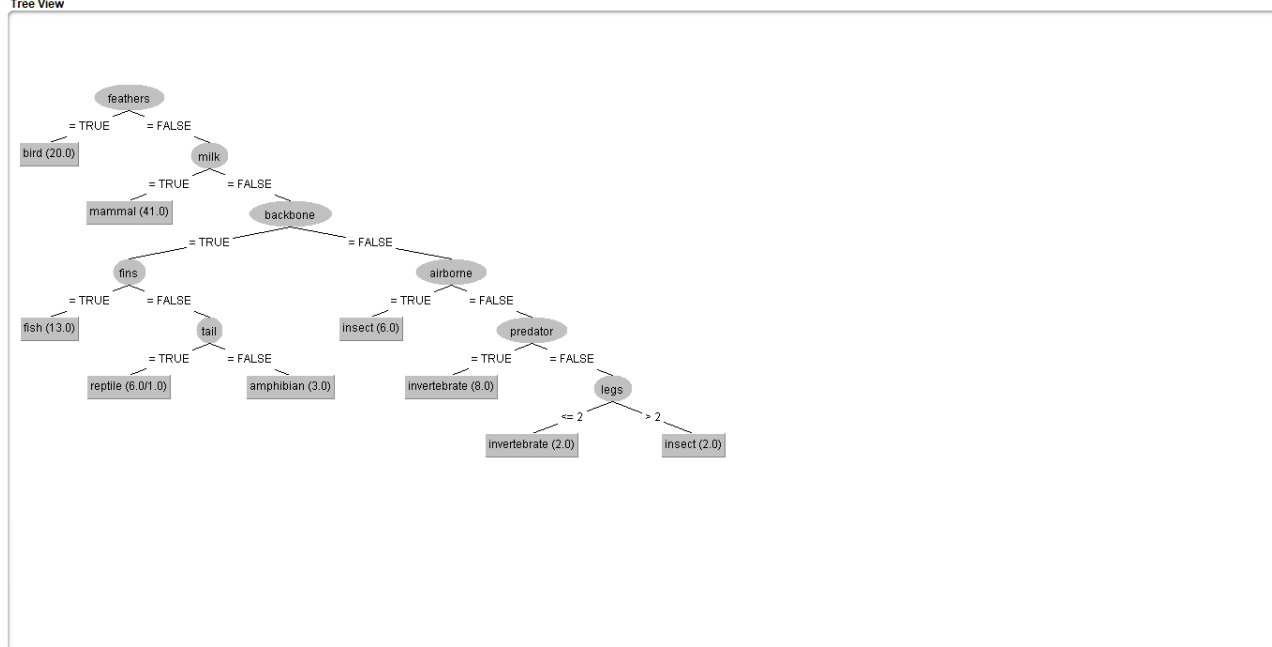
By looking at the confusion matrix we can see that mammals, fish, and birds are perfectly classified. But 2 invertebrates classified as insects, 3 insects classified as invertebrates. It means the identifying insects and invertebrates are not good in this model. And also 2 reptiles classified as 2 different typ so the model is unable to identify the reptiles also.

# 4. Evaluate the C4.5 using

## a. The training set



## b. 10-fold cross validation



**Training set accuracy : 99.008%**
**Cross validation accuracy : 92.07%**

Training set accuracy is higher than the cross validation. That happen because of

## 05. Can we use ID3?

Can't use it because in 2 attributes there are some missing attributes. Therefore we need to remove or fill those missing values.

## 06, Remover 2 instances



## 07. Build ID3

# 08. One R algorithm



In the oneR algorithm accuracy is less which is 60%. Only birds and mammals are correctly classified. And all others were not classified at all. Most of the ones classified as mammals with the whole category. It means that accuracy is 60% because of having 41 mammals. If not the accuracy may fall in to very less.

```
=== Detailed Accuracy By Class ===
```

| | TP Rate | FP Rate | Precision | Recall | F-Measure | MCC | ROC Area | PRC Area | Class |
|---|---|---|---|---|---|---|---|---|---|
| | 1.000 | 0.667 | 0.506 | 1.000 | 0.672 | 0.411 | 0.667 | 0.506 | mammal |
| | 0.000 | 0.000 | ? | 0.000 | ? | ? | 0.500 | 0.129 | fish |
| | 1.000 | 0.000 | 1.000 | 1.000 | 1.000 | 1.000 | 1.000 | 1.000 | bird |
| | 0.000 | 0.000 | ? | 0.000 | ? | ? | 0.500 | 0.099 | invertebrate |
| | 0.000 | 0.000 | ? | 0.000 | ? | ? | 0.500 | 0.079 | insect |
| | 0.000 | 0.000 | ? | 0.000 | ? | ? | 0.500 | 0.040 | amphibian |
| | 0.000 | 0.000 | ? | 0.000 | ? | ? | 0.500 | 0.050 | reptile |
| Weighted Avg. | 0.604 | 0.271 | ? | 0.604 | ? | ? | 0.667 | 0.440 | |

```
=== Confusion Matrix ===

 a  b  c  d  e  f  g   <-- classified as
41  0  0  0  0  0  0 |  a = mammal
13  0  0  0  0  0  0 |  b = fish
 0  0 20  0  0  0  0 |  c = bird
10  0  0  0  0  0  0 |  d = invertebrate
 8  0  0  0  0  0  0 |  e = insect
 4  0  0  0  0  0  0 |  f = amphibian
 5  0  0  0  0  0  0 |  g = reptile
```