

Problem #1)

a) Ockham's razor: A consistent hypothesis agrees with all the data. How do we choose from the multiple options for consistent data?

\* Choose the Simplest option.

main advantages: Choosing a simple option makes selecting an option much easier. This relieves the computational complexity of learning.

We will want to use the hypothesis ( $h$ ) after we have learned it. "and Computing  $h(x)$  when  $h$  is a linear function is guaranteed to be fast, while Computing an arbitrary Turing machine program is not even guaranteed to terminate. For these reasons, most work on learning has focused on simple representations." - pg. 697 text

To summarize Choosing a more simple representation makes initial computation and future use of the function more easy.

Main disadvantages: Any time we make simplifying assumptions we loose the ability to make an exact fit and model the data perfectly. Simplifying the function will always result in a loss of accuracy. This is like sampling digital music to a CD, or Compressing a photo. There will always be some type of loss when taking an analog signal and making it "digital". There will be loss in the simplification.

b) Main limitations of decision tree learning:

A decision tree is still a model or simplification of the hypothesis space. Because of this simplification, there will be loss of granularity. Some decision trees can not represent the full function without being exponentially large. Some functions can be represented by decision trees while others can not due to the size the hypothesis space is just too large in many cases.

- 1) Choosing the appropriate attribute selection measure
- 2) How deep to grow the decision tree
- 3) Handling continuous attributes
- 4) Handling training data w/ missing attributes
- 5) Handling attributes w/ different costs
- 6) Improving computational efficiency.



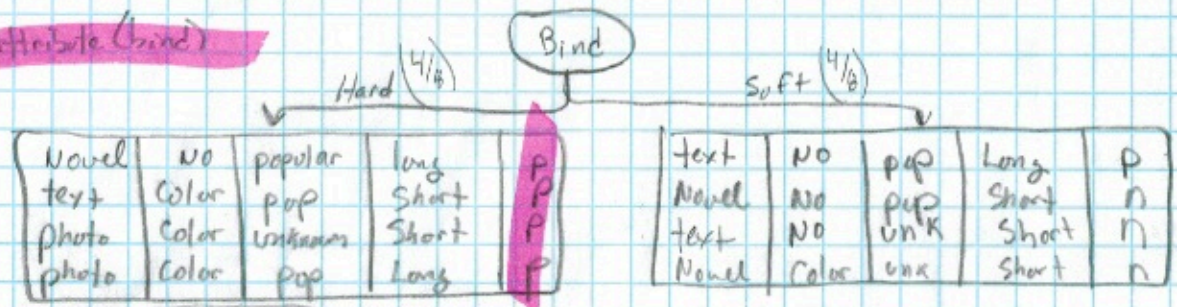
Bind	Style	Color	Known	Length	Classification
Hard	Novel	NO	Popular	Long	P
Soft	Text	NO	Popular	Long	P
Soft	Novel	NO	Popular	Short	n
Hard	Text	Color	Popular	Short	P
Hard	Photo	Color	Unknown	Short	P
Soft	Text	NO	Unknown	Short	n
Hard	Photo	color	Popular	Long	P
Soft	Novel	color	Unknown	Short	n

Step 1: Determine high level entropy

$P = 5$   
 $n = 3$   
 $H(c) = -\left(\frac{5}{8}\right) \log_2 \left(\frac{5}{8}\right) - \left(\frac{3}{8}\right) \log_2 \left(\frac{3}{8}\right)$   
 $H(c) = .95$

Step 2: Calculate info gain for different splits. (5 different splits)

1) Attribute (bind)

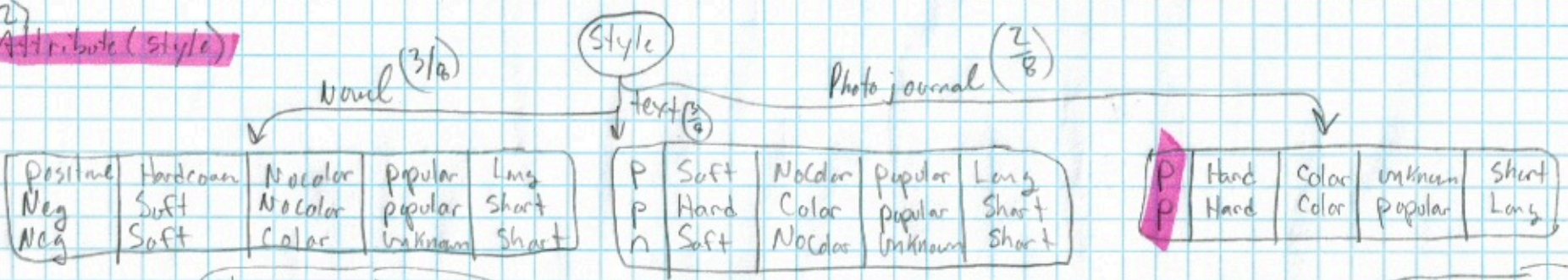


$H_{bind}(hard) = 0$  #pure split

$H_{bind}(soft) = -.25 \log_2(.25) - .75 \log_2(.75)$   
 $= .81$

$Gain_{bind} = .95 - .5(0) - .5(.81)$   
 $= .545$

2) Attribute (style)



$p = .33$   
 $n = .66$   
 $H_{style}(Novel) = .92$

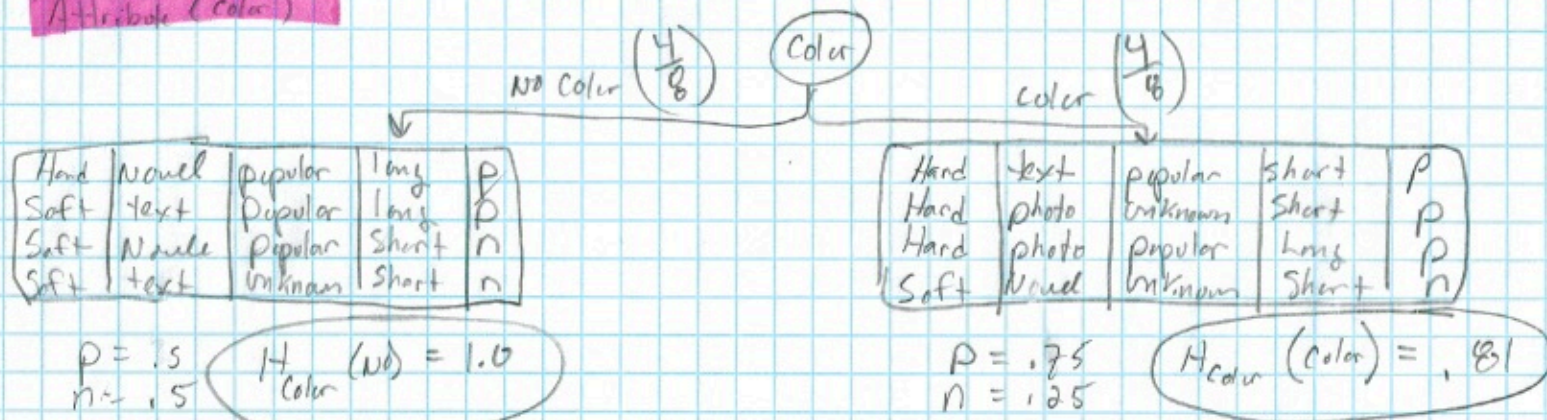
$p = .66$   
 $n = .33$   
 $H_{style}(Text) = .92$

$p = 1.0$   
 $n = 0.0$   
 $H_{style}(photo) = 0$

$Gain_{style} = .95 - \left(\frac{3}{8}\right) \cdot .92 - \left(\frac{3}{8}\right) \cdot .92 - \frac{2}{8}(0)$   
 $= .21$

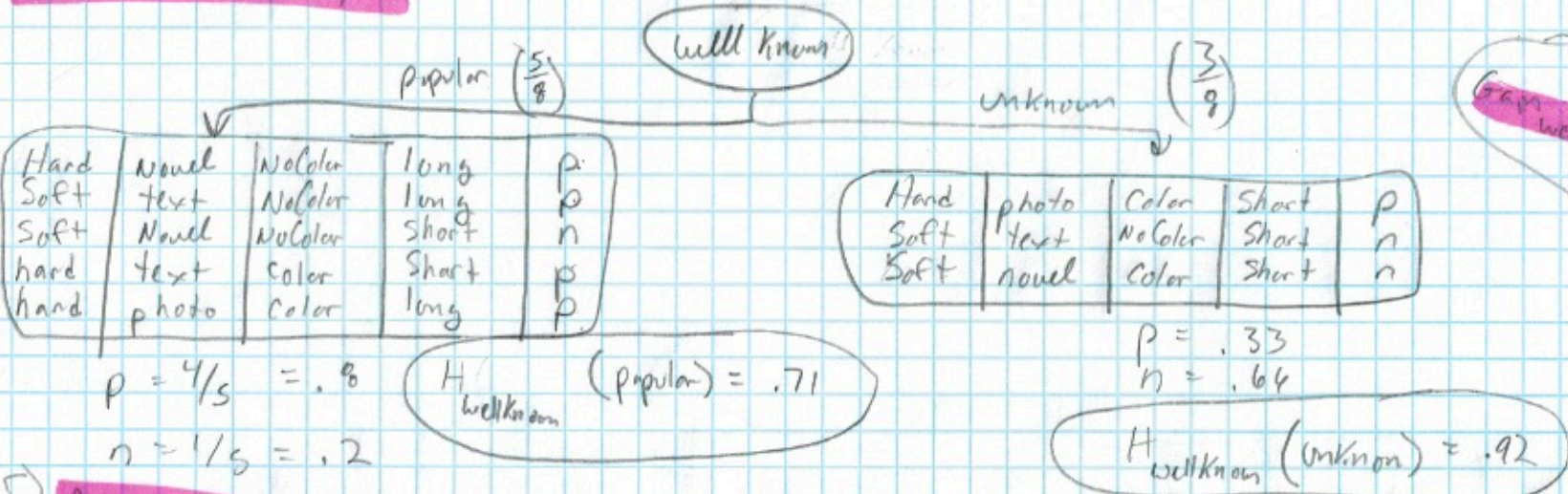


### 3) Attribute (color)



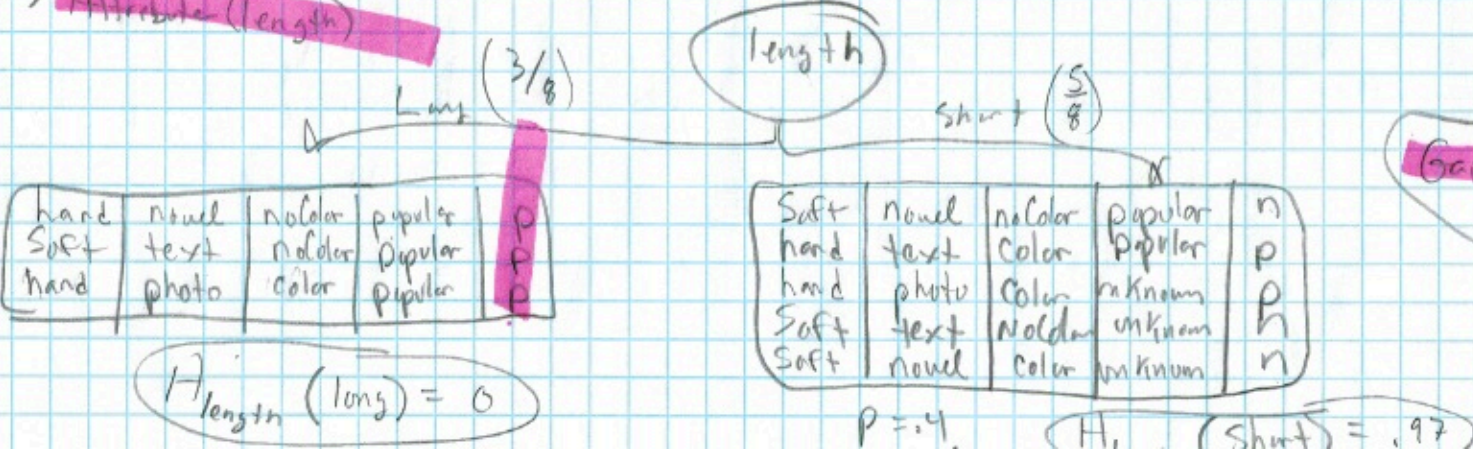
Gain<sub>color</sub> =  $.95 - .5(1) - .5(.81)$   
= .05

### 4) Attribute (well known)



Gain<sub>well known</sub> =  $.95 - (\frac{5}{8}) \cdot .71 - (\frac{3}{8}) \cdot .92$   
= .16

### 5) Attribute (length)



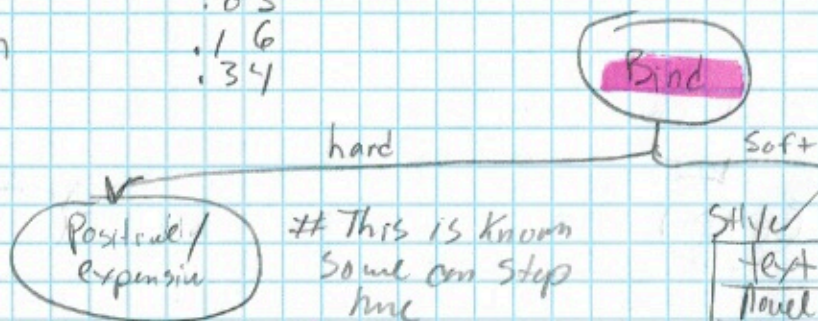
Gain<sub>length</sub> =  $.95 - (\frac{3}{8}) \cdot 0 - (\frac{5}{8}) \cdot .97$   
= .34



Step 3: Break down the next level of splits given the first level split

Split	info Gained
Bind	.55
Style	.26
Color	.05
well Known	.16
length	.34

# This is a recursive step.  
# bind is best split



Novel	noColor	popular	long	P
text	Color	popular	Short	P
photo	Color	unknown	Short	P
photo	Color	popular	long	P

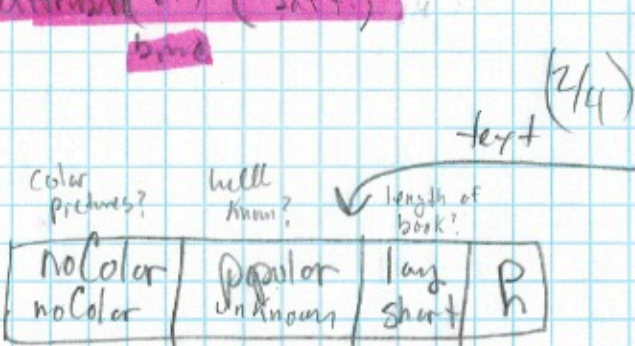
# this is unknown so we need to recurse on the problem

Style	color	well known	length	P
text	noColor	popular	long	P
Novel	noColor	popular	Short	n
text	noColor	unknown	Short	n
Novel	Color	unknown	Short	n

$P = .25$   
 $n = .75$   
 $H_{bind}(soft) = .81$

#  $H_{bind}(soft)$  is the starting entropy for the next level of calculations. Repeat step 2 to calculate gain for different splits.

attribute (color) (soft)  
bind



noColor	popular	long	P
noColor	unknown	short	n

$E_{style}(text) = 1.0$

noColor	popular	Short	n
Color	unknown	Short	n

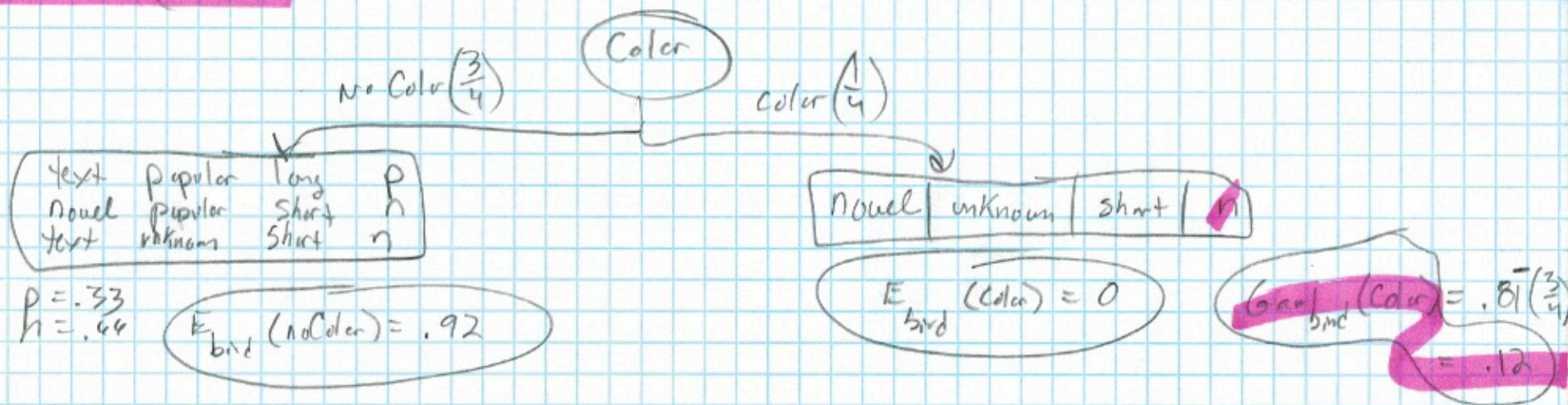
$E_{style}(novel) = 0$

$gain = .81 - .5(1) = .31$

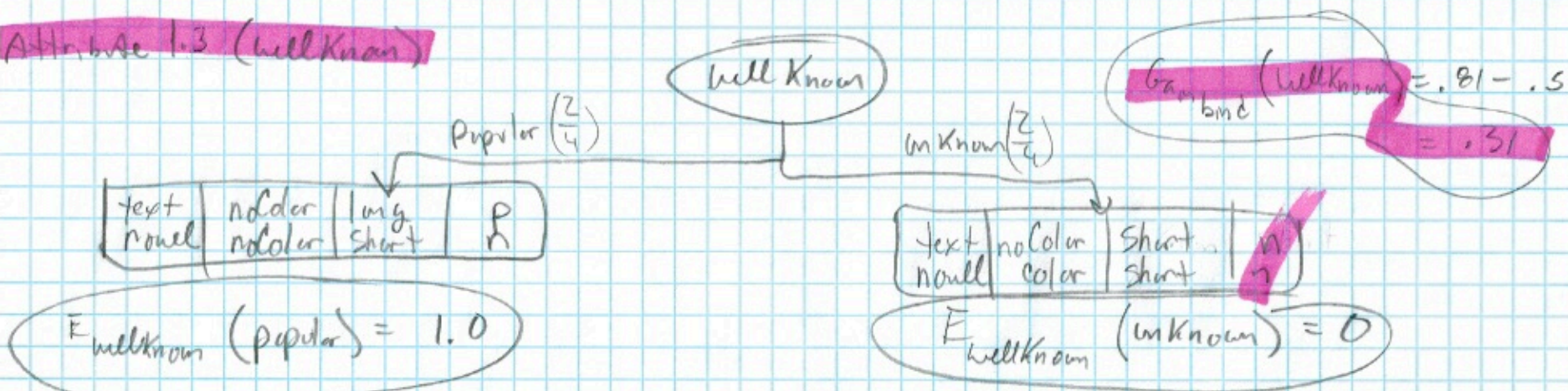
$gain_{bind}(style) = .31$



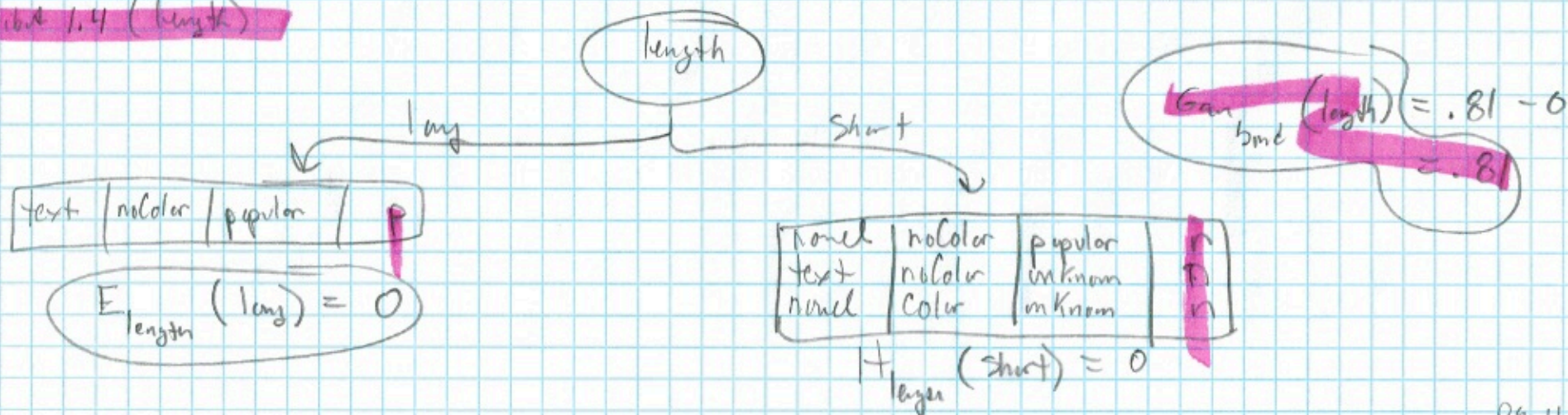
### Attribute 1.2 (Color)



### Attribute 1.3 (Well Known)



### Attribute 1.4 (Length)

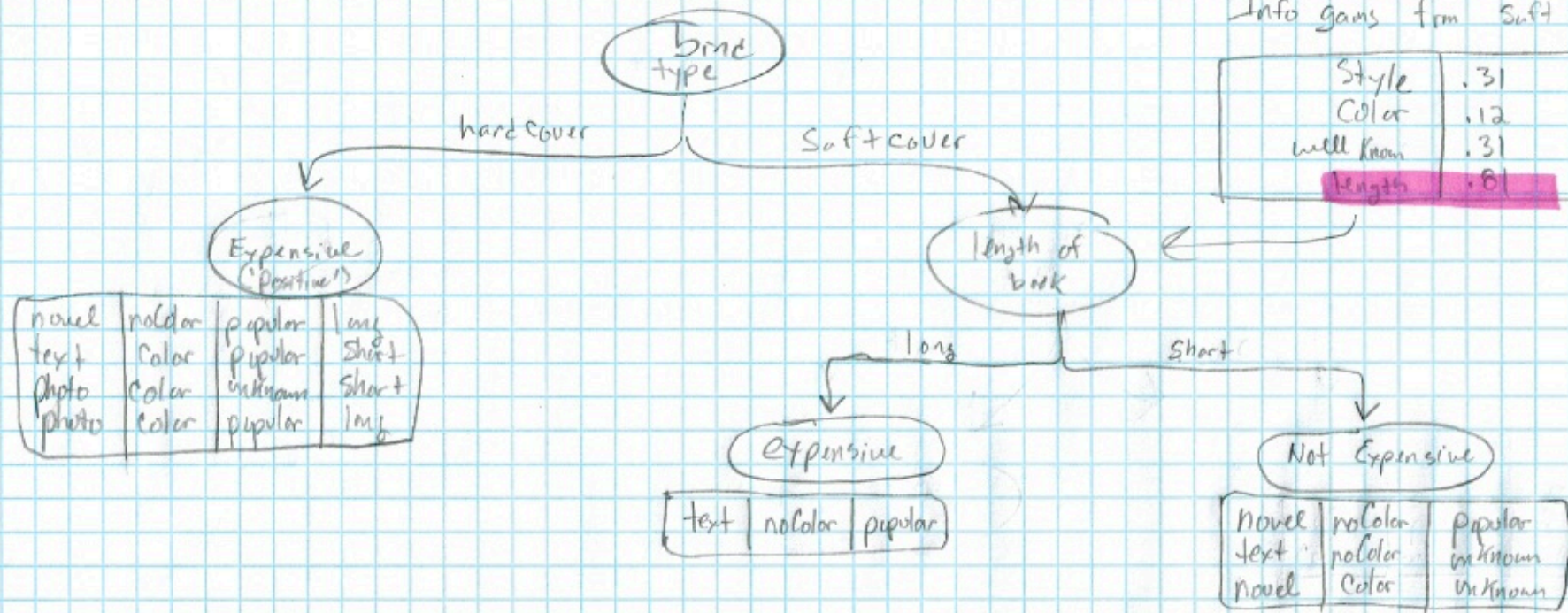




# length is the next most important attribute and produces a pure set.

Info gains from Soft cover

Style	.31
Color	.12
well known	.31
length	.81





# Artificial Intelligence

Name and ID \_\_\_\_\_

## Assignment 5

**Due date: November 30 at 11:59pm**

### Problem 1 (4 points)

Read Chapter 18 and answer the following questions:

- What are the main advantages and drawbacks of using Ockham's razor in learning
- What are the main limitation of the decision-tree learning?

### Problem 2 (6 points)

This file contains eight training data and three test data, and the attributes refer to whether a book will be expensive at the local bookstore.

Classification Bind type Style of book Color pictures? Is the book well known? Length of book

Positive -	Hardcover Novel	Nocolor	Popular	Long
Positive -	Softcover Textbook	Nocolor	Popular	Long
- Negative	Softcover Novel	Nocolor	Popular	Short
Positive -	Hardcover Textbook	Color	Popular	Short
Positive -	Hardcover Photojournal	Color	Unknown	Short
- Negative	Softcover Textbook	Nocolor	Unknown	Short
Positive -	Hardcover Photojournal	Color	Popular	Long
- Negative	Softcover Novel	Color	Unknown	Short

Use these data to construct a decision tree; you should compute the information gains to decide which attributes are more important. For each node of the tree, indicate the corresponding information gain.

### Problem 3 (10 points)

Implement a program for building decision trees. It should read a file with training and test examples, use the training examples to build a tree, and then classify the test examples. The only required output is the classification of the test examples; it does not have to include the tree itself (if you output also the tree, you will have 5 points bonus). The input format is as follows: