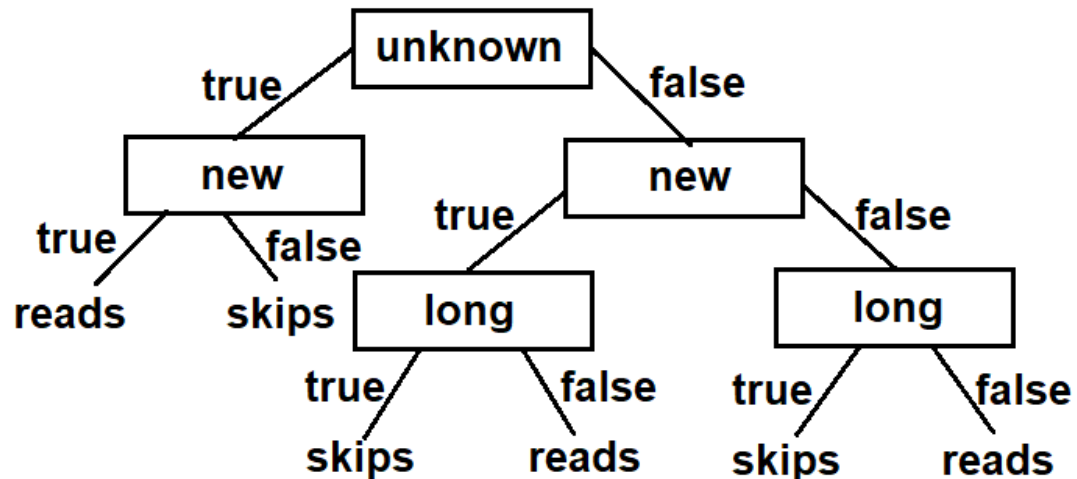


Assignment 2

Zhichao He z5282955

Question 1:

a)



They are different. The function with the maximum information gain split looks like this.

```
if long(e): return skips
else if new(e): return reads
else if unknown(e): return skips
else: return reads
```

(function 1)

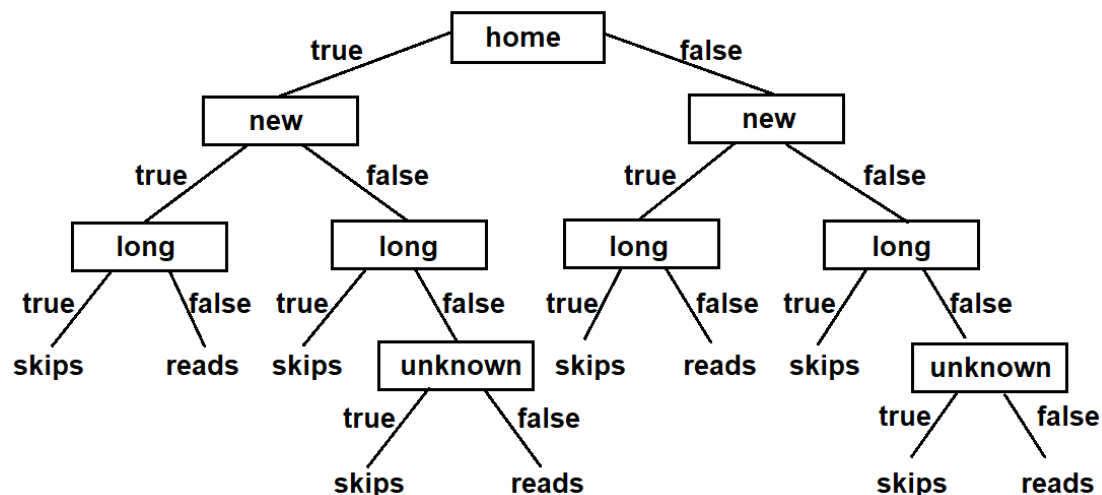
The function from the new decision tree is like

```
if unknown(e): if new(e): return reads
                else: return skips
else: if new(e): if long(e): return skips
                else: return reads
                else: if long(e): return skips
                     else: return reads
```

(function 2)

Structurally, they are different functions. And one function cannot be transformed to the other. Additionally, there's a counter example e_{19} from Figure 7.1 which is $\langle unknown, new, long, work, User_action? \rangle$. If we use this to test both functions, the function 1 gives the output $User_action = \mathbf{skips}$ while function 2 gives \mathbf{reads} . Thus, they are completely different.

b)



This tree represents the following function 3 which is substantially the same as function 1.

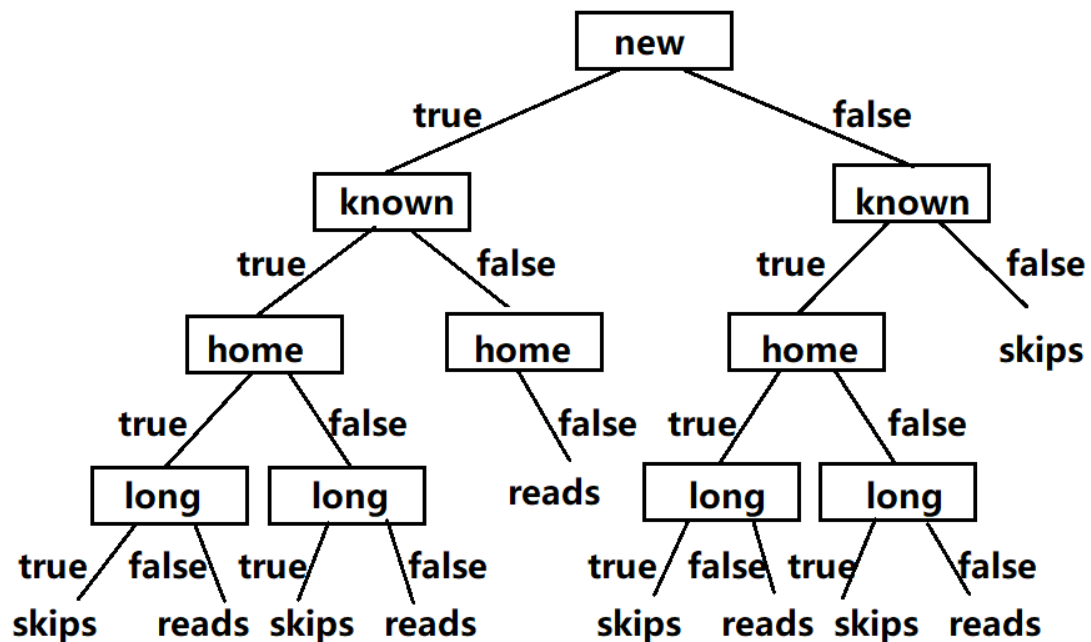
```

if home(e): if new(e): if long(e): return skips
                    else: return reads
            else: if long(e): return skips
                  else: if unknown(e): return skips
                        else: return reads
else: if new(e): if long(e): return skips
                else: return reads
      else: if long(e): return skips
            else: if unknown(e): return skips
                  else: return false  (function 3)
  
```

Though the structure of the 2 functions look different, they are substantially identical. As we can transform one function to the other. In the new decision tree, we can see that whatever values *home(e)* and *new(e)* are, if *long(e)* is true, the result is **skips**. This corresponds to 1st line of function 1. Similarly, whatever value *home(e)* is, if *long(e)* is false and *new(e)* is true, the result is **reads**. Corresponding to 2nd line of function 1. Whatever value *home(e)* is, if *long(e)* is false and *new(e)* is false and *unknown(e)* is true, the result is **skips**. Corresponding to 3rd line of function 1. Finally, whatever value *home(e)* is, if *long(e)*, *new(e)* and *unknown(e)* are all false, the result is **reads**. Corresponding to 4th line of function 1. Thus, they are the same function.

c)

Yes, there is such a tree.



This tree is generated if the order of features is [*Thread*, *Author*, *Where*, *Length*]. This tree correctly classifies the training examples as each example can obtain the correct *User_action* using this tree. And the tree represents different functions from those 2 in a) and b) Considering a test case *<unknown, new, home, long, User_action?>*, the tree in a) gives **reads**, the tree in b) gives **skips**, but the tree in c) doesn't define this scenario because there is no such kind of data to generate the that branch of the tree. Thus, although the new tree correctly classifies the training data, the function represented is different from the preceding ones.

Question 2:

I used Weka to create the decision tree. As Weka only accepts .arff format files, I combined adult.data and adult.test to make a complete data file and add labels @relation, @attribute and @data to make it a correct .arff file. Before training and testing, I split the data into first 2/3 to be training data and the rest 1/3 to be testing data.

Then I used J48 algorithm to build a few decision trees. I changed some parameters and compare the results in the table.

Algorithm	J48						
ConfFactor	-	0.75	0.50	0.25	0.1	0.25	0.25
MinNumObj	2	2	2	2	2	3	4
Model time	2.68s	73.68s	3.16s	5.1s	3.3s	2.72s	2.35s
Test time	0.07s	0.03s	0.02s	0.1s	0.01s	0.02s	0.02s
# leaves	9129	8025	2386	696	205	503	401
Tree size	10657	9365	2946	911	275	651	530
Accuracy	84.1%	84.1%	85.3%	85.9%	85.5%	85.8%	85.7%

ConfidenceFactor is the confidence factor for pruning (smaller values incur more pruning). The tree become smaller as this parameter decreases, but the accuracy isn't affected much.

MinNumObj is the minimum number of instances per leaf. The tree become smaller as this parameter increases but doesn't affect the accuracy much.

Considering the accuracy and efficiency, the optimal choice of parameters is the shaded one which has relatively high accuracy and the shortest test time. The decision tree of it is below.

```
capital-gain <= 6849
|   matital-status = Married-civ-spouse
|   |   capital-loss <= 1844
|   |   |   education-num <= 11
|   |   |   |   capital-gain <= 5060
|   |   |   |   |   age <= 29: <=50K (1999.0/241.0)
|   |   |   |   |   age > 29
|   |   |   |   |   |   hours-per-week <= 34: <=50K (1243.0/155.0)
|   |   |   |   |   |   hours-per-week > 34
|   |   |   |   |   |   |   education-num <= 9: <=50K (6969.0/1820.0)
|   |   |   |   |   |   |   education-num > 9
|   |   |   |   |   |   |   |   capital-loss <= 1510
|   |   |   |   |   |   |   |   |   occupation = Tech-support
|   |   |   |   |   |   |   |   |   capital-gain <= 3103: >50K (169.69/71.36)
```

[illegible]

[illegible]

[illegible]

[illegible]


```

| | | capital-loss <= 1980: >50K (857.0/18.0)
| | | capital-loss > 1980
| | | | capital-loss <= 2163: <=50K (104.0)
| | | | capital-loss > 2163
| | | | | capital-loss <= 2415
| | | | | capital-loss <= 2377
| | | | | | age <= 64: <=50K (38.0/4.0)
| | | | | | age > 64: >50K (30.0/3.0)
| | | | | capital-loss > 2377: >50K (82.0)
| | | | | capital-loss > 2415: <=50K (14.0)
| matital-status = Divorced: <=50K (6454.0/498.0)
| matital-status = Never-married
| | capital-loss <= 2206: <=50K (15843.0/495.0)
| | capital-loss > 2206
| | | capital-loss <= 2377: <=50K (38.0/9.0)
| | | capital-loss > 2377: >50K (27.0/1.0)
| matital-status = Separated: <=50K (1505.0/76.0)
| matital-status = Widowed
| | capital-loss <= 2205: <=50K (1460.0/82.0)
| | capital-loss > 2205
| | | race = White: >50K (23.0/9.0)
| | | race = Asian-Pac-Islander: >50K (0.0)
| | | race = Amer-Indian-Eskimo: >50K (0.0)
| | | race = Other: >50K (0.0)
| | | race = Black: <=50K (2.0)
| matital-status = Married-spouse-absent: <=50K (613.0/44.0)
| matital-status = Married-AF-spouse: <=50K (35.0/12.0)
capital-gain > 6849: >50K (2055.0/28.0)

```

Then I used REPTree algorithm to build a few other decision trees.
I changed some parameters and compare the results in the table.

Algorithm	REPTree						
MaxDepth	noLimit	7	4	2	4	4	4
MinNum	2	2	2	2	3	2	2
numFolds	3	3	3	3	3	4	5
Model time	0.83s	0.7s	0.37s	0.35s	0.77s	0.46s	0.64s
Test time	0.02s	0.02s	0.01s	0.02s	0.01s	0.01s	0.01s
Tree size	2041	1094	238	45	200	217	209
Accuracy	84.9%	85.1%	84.3%	82.7%	84.3%	84.5%	84.4%

MaxDepth is the max depth of the decision tree, as we decrease the MaxDepth, the tree size will decrease, and accuracy will also decrease slightly.

MinNum is the minimum total weight of the instances in a leaf, as it increases, the tree size decreases. But the accuracy doesn't change much.

numFolds determines the amount used for pruning. As it increases,

the tree size decreases slightly. But the accuracy doesn't change much.

Following is a decision tree with the minimum tree size in the table (not the most accurate).

```
relationship = Wife
| education = Bachelors : >50K (298/90) [166/64]
| education = Some-college : <=50K (306/132) [162/69]
| education = 11th : <=50K (36/3) [8/1]
| education = HS-grad : <=50K (485/153) [239/84]
| education = Prof-school : >50K (25/4) [13/1]
| education = Assoc-acdm : >50K (78/36) [30/14]
| education = Assoc-voc : >50K (76/33) [42/20]
| education = 9th : <=50K (25/2) [10/1]
| education = 7th-8th : <=50K (22/1) [7/0]
| education = 12th : <=50K (14/3) [7/0]
| education = Masters : >50K (117/20) [63/12]
| education = 1st-4th : <=50K (1/0) [5/0]
| education = 10th : <=50K (27/1) [19/2]
| education = Doctorate : >50K (18/1) [8/2]
| education = 5th-6th : <=50K (17/4) [4/0]
| education = Preschool : <=50K (1/0) [2/0]
relationship = Own-child
| capital-gain < 4718.5 : <=50K (4991/46) [2541/29]
| capital-gain >= 4718.5 : >50K (24/5) [25/8]
relationship = Husband
| education = Bachelors : >50K (2375/758) [1261/416]
| education = Some-college : <=50K (2486/1128) [1172/480]
| education = 11th : <=50K (322/44) [157/24]
| education = HS-grad : <=50K (4297/1328) [2091/677]
| education = Prof-school : >50K (362/53) [196/33]
| education = Assoc-acdm : >50K (380/194) [201/93]
| education = Assoc-voc : <=50K (576/250) [307/137]
| education = 9th : <=50K (186/17) [119/11]
| education = 7th-8th : <=50K (337/39) [166/17]
| education = 12th : <=50K (114/26) [52/11]
| education = Masters : >50K (876/212) [463/113]
| education = 1st-4th : <=50K (74/6) [37/2]
| education = 10th : <=50K (303/48) [168/27]
| education = Doctorate : >50K (256/44) [121/23]
| education = 5th-6th : <=50K (164/17) [74/5]
| education = Preschool : <=50K (10/0) [13/0]
relationship = Not-in-family
| capital-gain < 8296 : <=50K (8237/643) [4007/295]
| capital-gain >= 8296 : >50K (236/0) [103/1]
relationship = Other-relative : <=50K (1006/33) [500/19]
relationship = Unmarried
| capital-gain < 7139.5 : <=50K (3352/148) [1695/89]
| capital-gain >= 7139.5 : >50K (51/3) [27/3]
```

In conclusion, the best accuracy of predicting this dataset is 85.5% using J48 algorithm. Appropriate pruning and limitation of max tree depth will adjust the decision tree to a smaller size and the accuracy

will remain much the same or improve. Increasing MinNumObj in J48 may result in a smaller tree. Increasing MinNum or numFold in REPTree may also result a smaller tree. The best Max tree depth of REPTree is 7.