

## 3.3.8 Data compression (Huffman coding)

### Lesson plan and printable activities

#### Materials needed

---

1. 3.3.8 (Huffman coding) [Lesson](#).
2. [Quiz 1](#), [Quiz 2](#), [Quiz 3](#).
3. [Tables](#) from worked example.

#### Lesson aims

---

1. To gain an overview of how the Huffman coding method may be used to encode plain ASCII text.
4. To see how the Huffman coding method compares in its application to related compression algorithms, eg Run Length Encoding (RLE)

#### Lesson objectives

---

1. Analyse a piece of plain ASCII text in order to determine the frequency of specific characters occurring.
5. Transform examples of uncompressed ASCII text into a compressed format using the Huffman coding method.
6. Convert a specific example of compressed text into ASCII characters using a given Huffman tree.
7. Describe conditions where the Huffman coding method of compression will provide optimal benefits.

**Please note: This lesson can be split into two parts with slides 1–21 as lesson 1 and slides 22 onwards as lesson 2. To address: ‘Students should be familiar with carrying out calculations to determine the number of bits saved by compressing a piece of data using Huffman coding.’ This could be included as a homework.**

## 3.3 Fundamentals of data representation

### Starter activity (10 minutes)

1. Re-cap what has been discussed in the previous lesson, emphasising that the focus continues to be on the use of data compression techniques. State that we will be looking specifically at Huffman coding in this lesson and that the technique is useful for encoding plain text.

Use [Quiz 1](#) to start the students thinking about the fundamental basis of how Huffman works, ie the fact that characters appear at different frequencies in text, and that there quickly become noticeable differences in the frequency patterns of specific characters as the size of the text increases.

### Main activities (45 minutes)

1. **Slide 4:** Start by explaining that Huffman method is very useful but more complex for students to implement than the Run Length Encoding (RLE) method seen previously. Refer to the teacher notes in regard to the practical complexity of using the method in the classroom.
2. **Slide 5:** Demonstrate to students the frequency of different characters in a piece of text by using the search function on a Word or PDF document. Get students to suggest which the most frequent characters might be. Then get students to attempt [Quiz 1](#).
3. **Slides 7–27:** There is a single worked example that provides sufficient complexity, as the text chosen has a number of characters. It is likely that not all students will follow the diagrams at an equal rate, it is strongly suggested that the accompanying tables are also provided as single-side printed documents in order to help students follow the explanation.
4. It might be useful for more able students to emphasise that the algorithm used is 'lossless' – a term not used on the GCSE specification, the rationale for this is that the encoded text may be reconstituted using the Huffman tree generated.
5. Complete [Quiz 2](#).

### Plenary activity

1. Complete [Quiz 3](#), which asks students to consolidate their thinking after they have seen the algorithm in action.


### Lesson

3.3 Fundamentals of data representation

3.3.8 Data compression (Huffman coding)

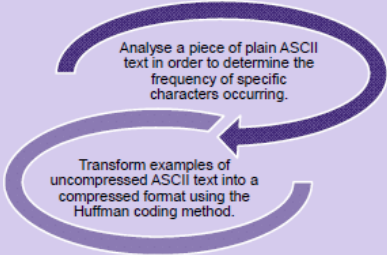
Lesson

© 2016 AQA. Created by Teachit for AQA



**Objectives**

Be able to:



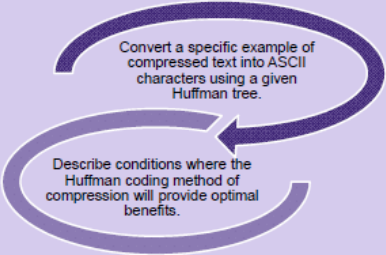
Analyse a piece of plain ASCII text in order to determine the frequency of specific characters occurring.

Transform examples of uncompressed ASCII text into a compressed format using the Huffman coding method.

© 2016 AQA. Created by Teachit for AQA

**Objectives**

Be able to:



Convert a specific example of compressed text into ASCII characters using a given Huffman tree.

Describe conditions where the Huffman coding method of compression will provide optimal benefits.

© 2016 AQA. Created by Teachit for AQA

**Introduction**

Continuing the theme of data compression...

In the previous lesson, we looked at why it is beneficial to be able to compress digital files.

We looked at one specific method of carrying out data compression:

**RLE**

Today, we are going to look at another method:

**Huffman coding**

© 2016 AQA. Created by Teachit for AQA

**How often do different characters occur in text?**

This method works by recording the frequencies by which characters appear in text.

Some characters appear far more often than others.

**Which do you think occur most frequently?**

**How about less frequently?**

**Try Quiz 1**

© 2016 AQA. Created by Teachit for AQA

**Focus on Huffman coding**

There are two steps when using Huffman's technique:

1. Create a Huffman tree – this gives each character used a unique code.
2. Encode the character sequence into a binary stream.

© 2016 AQA. Created by Teachit for AQA

### 3.3 Fundamentals of data representation

**Huffman coding – worked example**

Look at this sentence:

*an easy example*

**What do you notice?**

© 2016 AQA. Created by Teachit for AQA

**Huffman coding – worked example**

Now create a table showing the frequency of occurrence for each character.

*an easy example*

Character	Frequency	Character	Frequency
a	Click to reveal	p	Click to reveal
e			
l			
m			
n			
		space	

© 2016 AQA. Created by Teachit for AQA

**Huffman coding – worked example**

Now we need to create the graphical node representation used by the Huffman method.

To begin this process, we have to start with the lowest frequencies at the top of the list – the actual alphabetical order does not matter.

The **Huffman tree** generated allows you to assign a **unique code** to each character used.

The wider the range of characters used, the bigger the Huffman tree!

© 2016 AQA. Created by Teachit for AQA

**Huffman coding – worked example**

At the **start**, the characters are placed in order of occurrence, from top to bottom. Where the frequency is the same, the order does not matter. ( ) indicates the space character.

**Step 1**  
The top two entries' frequencies are added together and a new node is inserted.

© 2016 AQA. Created by Teachit for AQA

**Huffman coding – worked example**

**Step 2**  
Same as Step 1 with 1 (n) and 1 (p).

© 2016 AQA. Created by Teachit for AQA

**Huffman coding – worked example**

**Step 3**  
Same as Step 2 with 1 (s) and 1 (x).

© 2016 AQA. Created by Teachit for AQA

**Huffman coding – worked example**

**Step 4**  
Look what happens to 1 (y).

© 2016 AQA. Created by Teachit for AQA

**Huffman coding – worked example**

**Step 5**  
The top two nodes combine and move to the end of the list as their combined value exceeds that of the other elements.

© 2016 AQA. Created by Teachit for AQA

### 3.3 Fundamentals of data representation

#### Huffman coding – worked example

**Step 6**  
The pattern continues...

© 2016 AQA. Created by Teachit for AQA

#### Huffman coding – worked example

**Step 7**  
More of the same...

© 2016 AQA. Created by Teachit for AQA

#### Huffman coding – worked example

**Step 8**  
The objective is to end up with one single node.

© 2016 AQA. Created by Teachit for AQA

#### Huffman coding – worked example

**Step 9**  
We made it!

The value of the left-hand node should be equal to the number of characters in the original string.

*an easy example* = 15 characters

Success!

© 2016 AQA. Created by Teachit for AQA

#### Huffman coding – worked example

**Step 10**  
Put '1' on the upward directions and '0' on the downward directions.

How would we denote the 'n'?

Down one, up one, up one, up one = **0111**

© 2016 AQA. Created by Teachit for AQA

#### Huffman coding – worked example

**Step 10 - Add labelling**

Character	Code
t	0001
h	0100
n	0111
e	0110
a	00001
s	00000
y	0001
o	001
x	11
m	10

© 2016 AQA. Created by Teachit for AQA

#### Huffman coding – worked example

a	n	e	a	s	y	e	x	a	m	p	l	e
11			11				11					10
		10				10						
	001				001							
							0100					
								0110				
				0001							0101	
0111												
		00001										
						00000						

Full binary representation =  
11 0111 001 10 11 00001 0001 001 10 00000 11 0100 0110 0101 10

© 2016 AQA. Created by Teachit for AQA

#### Decoding with the Huffman method – worked example

Our encoded text would be received as follows...

110111001101100001000100110000001101000110010110

**You must be in possession of the Huffman tree to decode the text!**

© 2016 AQA. Created by Teachit for AQA

### 3.3 Fundamentals of data representation

**Decoding with the Huffman method – worked example**

**First character**  
Start at beginning of string and follow binary until it terminates:

11011100110110000100  
01001100000011010001  
10010110

Following '1' then '1' takes us to **a**.

© 2016 AQA. Created by Teachit for AQA

**Decoding with the Huffman method – worked example**

**Second character**  
Move to the next set of digits:

44011100110110000100  
01001100000011010001  
10010110

= **an**

© 2016 AQA. Created by Teachit for AQA

**Decoding with the Huffman method – worked example**

**Third character**  
Continue the previous process – this time we find a space character...

44011100110110000100  
01001100000011010001  
10010110

= **an**

© 2016 AQA. Created by Teachit for AQA

**Decoding with the Huffman method – worked example**

**Fourth character**

44011100110110000100  
01001100000011010001  
10010110

= **an e**

Continue until no binary digits remain.

© 2016 AQA. Created by Teachit for AQA

**Huffman coding – conclusion**

Look at the table.

You can see that the characters at the top generally:

- appear most frequently
- have short binary representations.

This is the advantage of Huffman coding.

If conventional ASCII uses eight bits to represent a character.  
Huffman coding can reduce that to as little as two bits.

l	0101
m	0100
n	0111
p	0110
s	00001
x	00000
y	0001
space	001
a	11
e	10

© 2016 AQA. Created by Teachit for AQA

**Huffman coding – conclusion**

Note: Huffman coding method will incur some 'overhead' when used.

What do you think this is?

It is necessary to transmit a representation of the encoding 'tree' along with the encoded data. Each Huffman tree is unique as it is dependent upon the piece of text which it used to encode.

The Huffman method works best when there is a good distribution of frequently occurring 'common' characters.

© 2016 AQA. Created by Teachit for AQA

To round things off...

**Complete Quiz 2 and Quiz 3**

© 2016 AQA. Created by Teachit for AQA

## Quiz 1 – Frequency analysis of characters in text

### Question 1

Analyse the frequency of the occurrence of the characters in text using the example provided below.

Note: You should ignore the case of the characters and include any spaces.

*The quick brown fox jumped over the lazy wolf.*

Character	Frequency	Character	Frequency
<i>a</i>		<i>n</i>	
<i>b</i>		<i>o</i>	
<i>c</i>		<i>p</i>	
<i>d</i>		<i>q</i>	
<i>e</i>		<i>r</i>	
<i>f</i>		<i>s</i>	
<i>g</i>		<i>t</i>	
<i>h</i>		<i>u</i>	
<i>i</i>		<i>v</i>	
<i>j</i>		<i>w</i>	
<i>k</i>		<i>x</i>	
<i>l</i>		<i>y</i>	
<i>m</i>		<i>z</i>	
<i>.</i>		<i>,</i>	
<i>space</i>			

### 3.3 Fundamentals of data representation

#### Question 2

Analyse the frequency of the occurrence of the characters in text using the example provided below.

Note: You should ignore the case of the characters and include any spaces.

*We know the most commonly used passwords are numeric ones, so they are not secure, and people also use the same passwords for multiple sites. If one site gets hacked all the places that you use the same password get compromised, they are a big pain.*

Character	Frequency	Character	Frequency
<i>a</i>		<i>n</i>	
<i>b</i>		<i>o</i>	
<i>c</i>		<i>p</i>	
<i>d</i>		<i>q</i>	
<i>e</i>		<i>r</i>	
<i>f</i>		<i>s</i>	
<i>g</i>		<i>t</i>	
<i>h</i>		<i>u</i>	
<i>i</i>		<i>v</i>	
<i>j</i>		<i>w</i>	
<i>k</i>		<i>x</i>	
<i>l</i>		<i>y</i>	
<i>m</i>		<i>z</i>	
<i>.</i>		<i>,</i>	
<i>space</i>			



**Quiz 1 – Frequency analysis of characters in text – answers**

<b>Question 1</b>			
<b>Character</b>	<b>Frequency</b>	<b>Character</b>	<b>Frequency</b>
<i>a</i>	1	<i>n</i>	1
<i>b</i>	1	<i>o</i>	3
<i>c</i>	1	<i>p</i>	1
<i>d</i>	1	<i>q</i>	1
<i>e</i>	4	<i>r</i>	1
<i>f</i>	2	<i>s</i>	-
<i>g</i>	-	<i>t</i>	-
<i>h</i>	-	<i>u</i>	2
<i>i</i>	1	<i>v</i>	1
<i>j</i>	1	<i>w</i>	2
<i>k</i>	1	<i>x</i>	-
<i>l</i>	1	<i>y</i>	1
<i>m</i>	1	<i>z</i>	1
<i>.</i>	1	<i>,</i>	-
<b>space</b>	7		

### 3.3 Fundamentals of data representation

<b>Question 2</b>			
<b>Character</b>	<b>Frequency</b>	<b>Character</b>	<b>Frequency</b>
<i>a</i>	14	<i>n</i>	6
<i>b</i>	1	<i>o</i>	16
<i>c</i>	5	<i>p</i>	8
<i>d</i>	5	<i>q</i>	-
<i>e</i>	28	<i>r</i>	7
<i>f</i>	2	<i>s</i>	23
<i>g</i>	3	<i>t</i>	11
<i>h</i>	3	<i>u</i>	7
<i>i</i>	7	<i>v</i>	-
<i>j</i>	-	<i>w</i>	5
<i>k</i>	2	<i>x</i>	-
<i>l</i>	6	<i>y</i>	3
<i>m</i>	9	<i>z</i>	-
<i>.</i>	2	<i>,</i>	3
<b>space</b>	45		

## Quiz 2 – Encode a piece of text using Huffman coding

### Question 1

Create a Huffman tree for the following text sample. Your tree should indicate the binary code paths required to be followed in order to obtain each unique character.

*here horse*

## Quiz 2 – Encode a piece of text using Huffman coding – answers

### Question 1

Create a Huffman tree for the following text sample. Your tree should indicate the binary code paths required to be followed in order to obtain each unique character.

*here horse*

Correct analysis as per table

Correct tree L/H total = 10

Correct nodes

Correct binary labelling

### Huffman coding – character frequency record sheet

Character	Frequency	Character	Frequency	Character	Frequency
<b>a</b>	-	<b>n</b>	-	<b>space</b>	1
<b>b</b>	-	<b>o</b>	1	<b>.</b>	-
<b>c</b>	-	<b>p</b>	-	<b>,</b>	-
<b>d</b>	-	<b>q</b>	-		
<b>e</b>	3	<b>r</b>	2		
<b>f</b>	-	<b>s</b>	1		
<b>g</b>	-	<b>t</b>	-		
<b>h</b>	2	<b>u</b>	-		
<b>i</b>	-	<b>v</b>	-		
<b>j</b>	-	<b>w</b>	-		
<b>k</b>	-	<b>x</b>	-		
<b>l</b>	-	<b>y</b>	-		
<b>m</b>	-	<b>z</b>	-		

Diagram

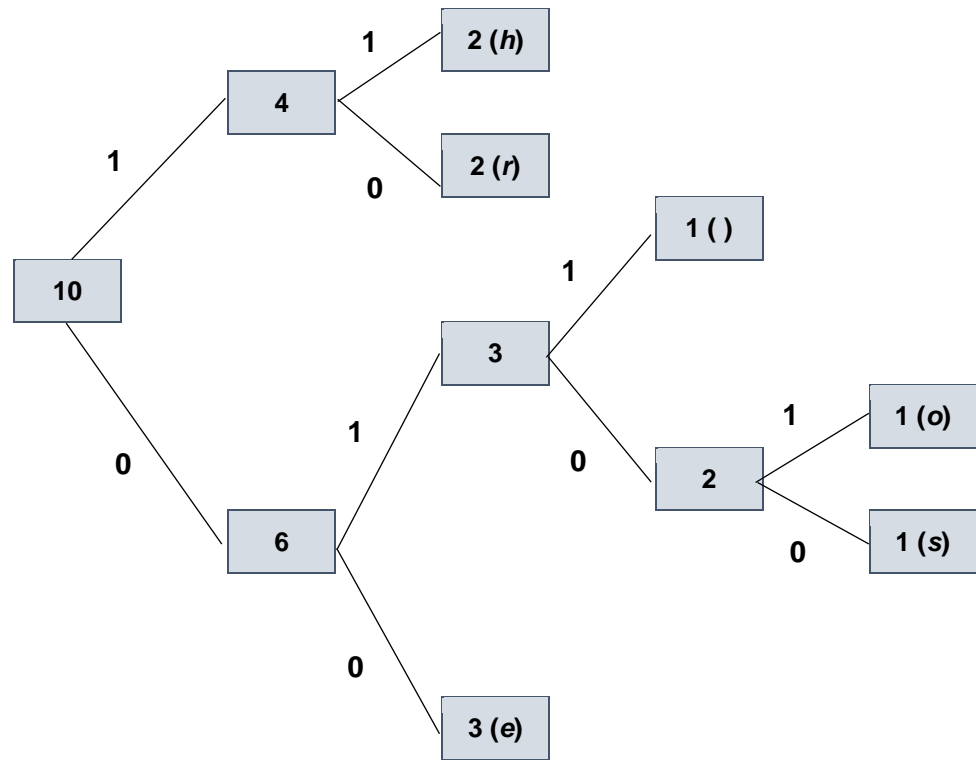


Table – for interest only

<b><i>h</i></b>	11
<b><i>r</i></b>	10
<b><i>e</i></b>	00
<b><i>o</i></b>	0101
<b><i>s</i></b>	0100
<b><i>space</i></b>	011

## Quiz 3 – Some general considerations when using Huffman coding

### Question

Look at the following two extracts of text:

1. *Chemical formula  $C_3H_4Br_2$*
2. *The sun's diameter is 866,000 miles. Rotates every 606 hours. The length of time its current carries the sun over its orbit is unknown. The sun remains a melted mass; its vibration is maintained; has but little vapour and its theme reflected on the surface of its obsequious attendants which gives them heat and light.*

Which of the two pieces of text might benefit from the use of Huffman coding to save space?

Why?

## Quiz 3 – Some general considerations when using Huffman coding – answers

### Question

Look at the following two extracts of text:

1. *Chemical formula  $C_3H_4Br_2$*
2. *The sun's diameter is 866,000 miles. Rotates every 606 hours. The length of time its current carries the sun over its orbit is unknown. The sun remains a melted mass; its vibration is maintained; has but little vapour and its theme reflected on the surface of its obsequious attendants which gives them heat and light.*

Source: [gutenberg.org/cache/epub/1331/pg1331-images.html](http://gutenberg.org/cache/epub/1331/pg1331-images.html)

Which of the two extracts might benefit from the use of Huffman coding to save space?

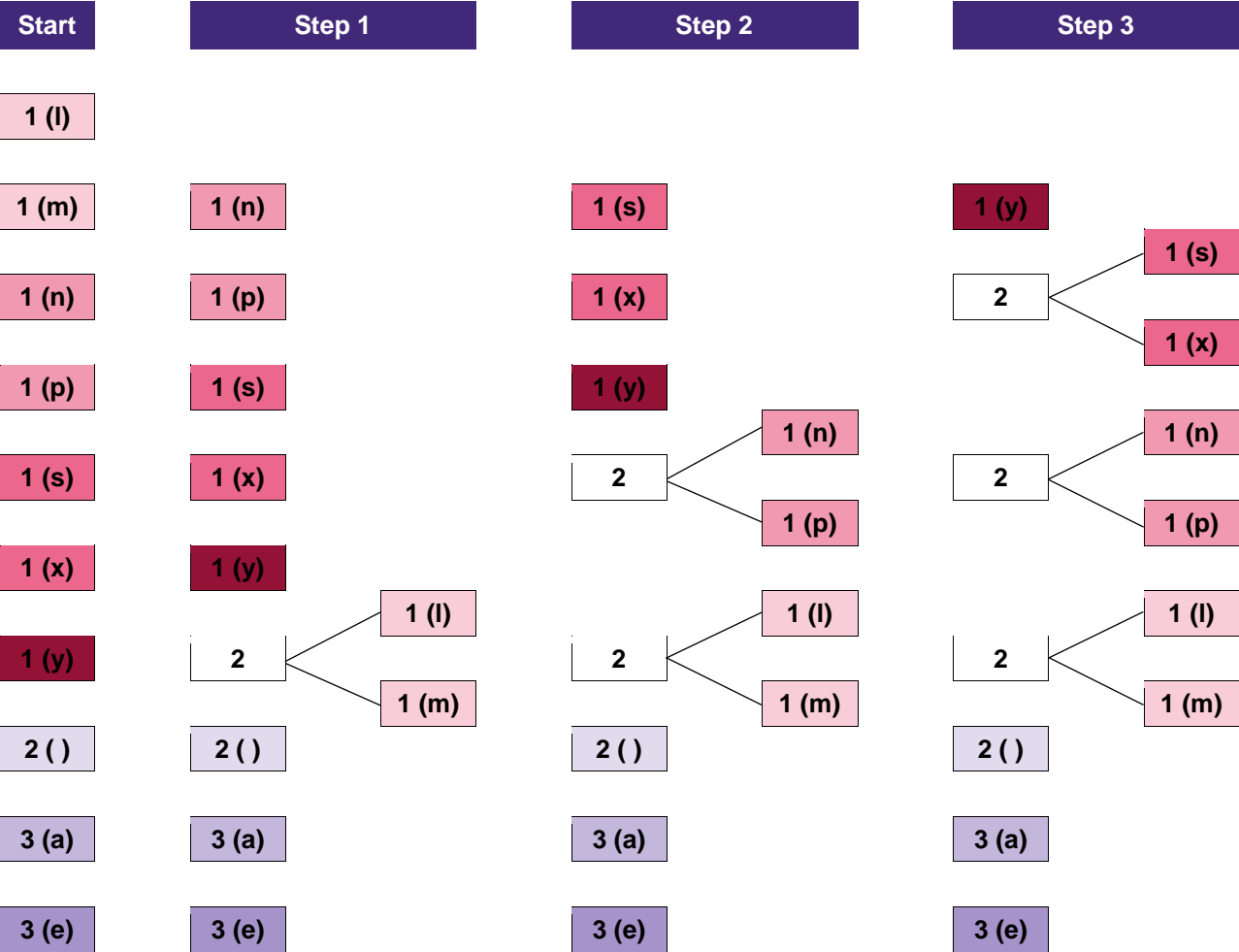
Why?

There is no set answer but you should look for the following thinking:

- Extract 1 is a relatively **short** piece containing a large variety of characters all of which will need their own encoding sequence and lead to a disproportionately large Huffman tree being produced.
- The shortness of the first piece makes it inefficient to encode as also need to include the associated Huffman tree to enable decoding.

Give extra credit if students attempt a (basic) character count to back up their arguments.

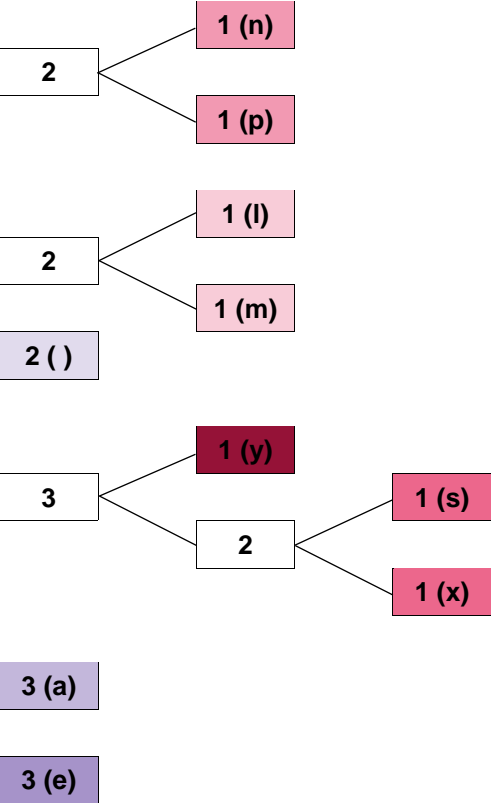
Tables from worked example



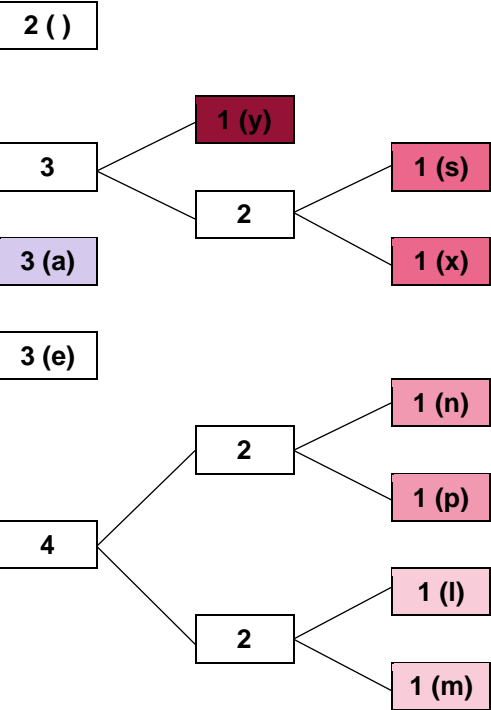


3.3 Fundamentals of data representation

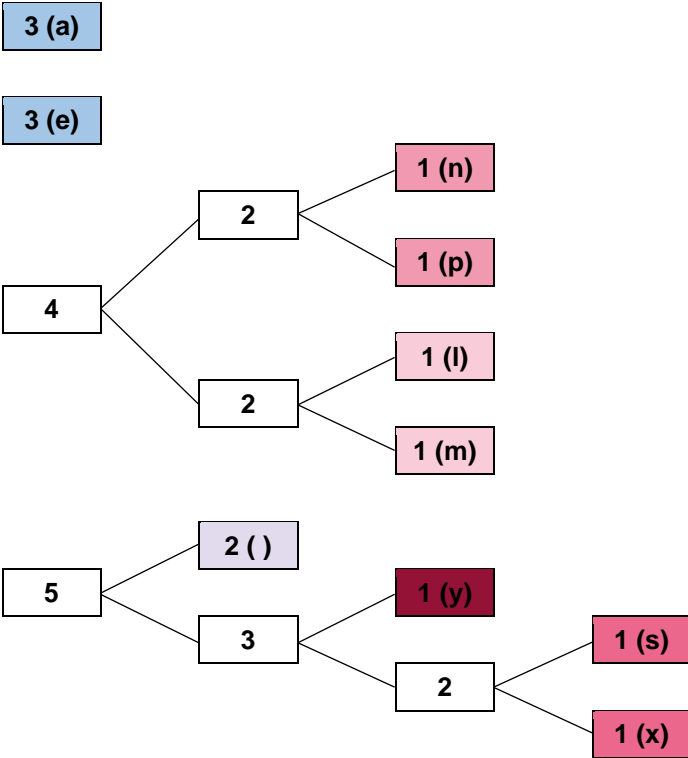
Step 4



Step 5

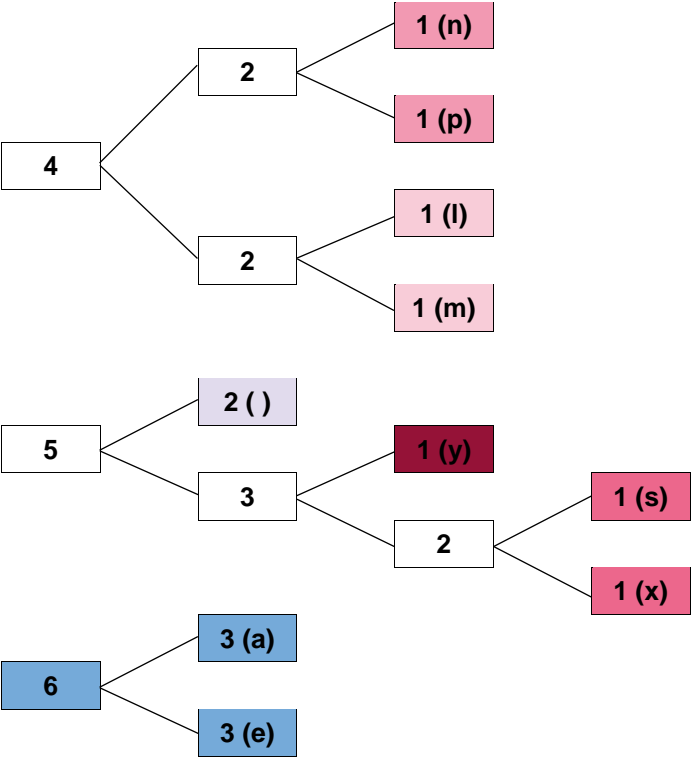


Step 6

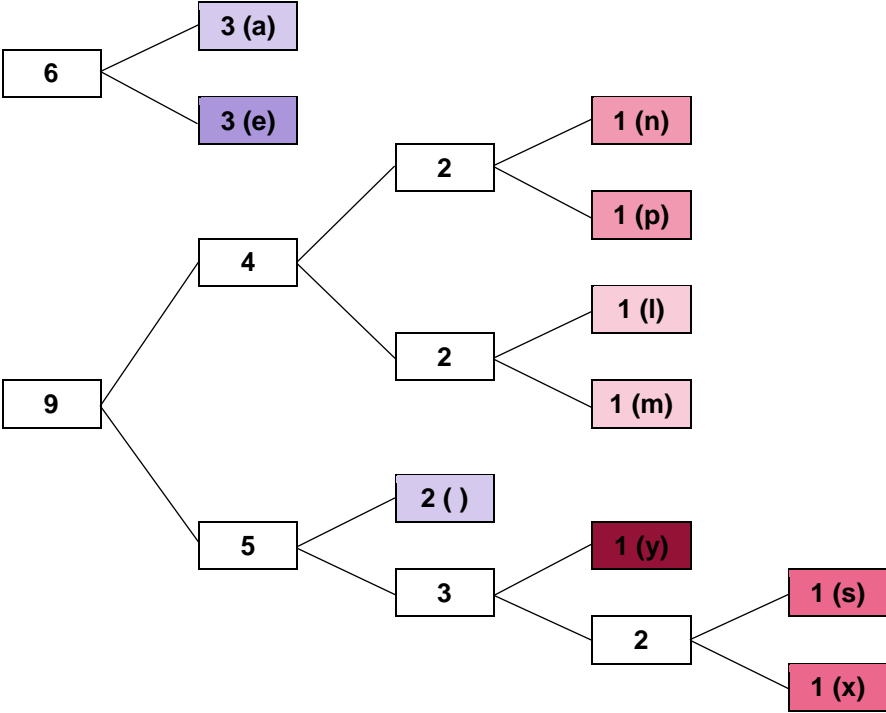


3.3 Fundamentals of data representation

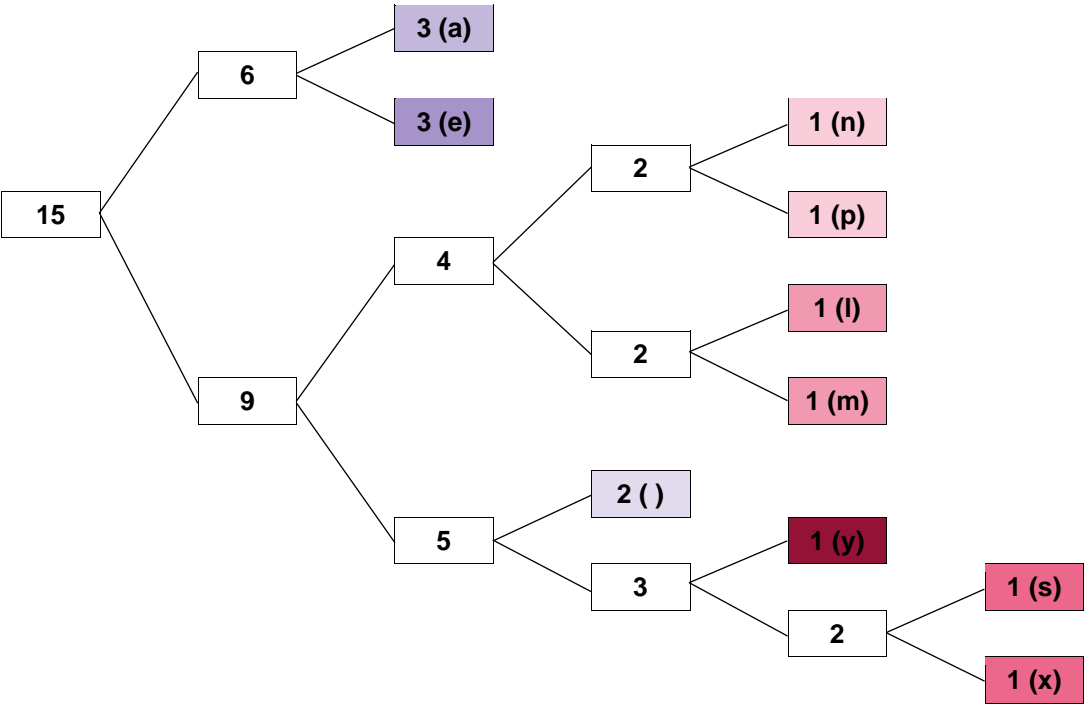
Step 7



Step 8

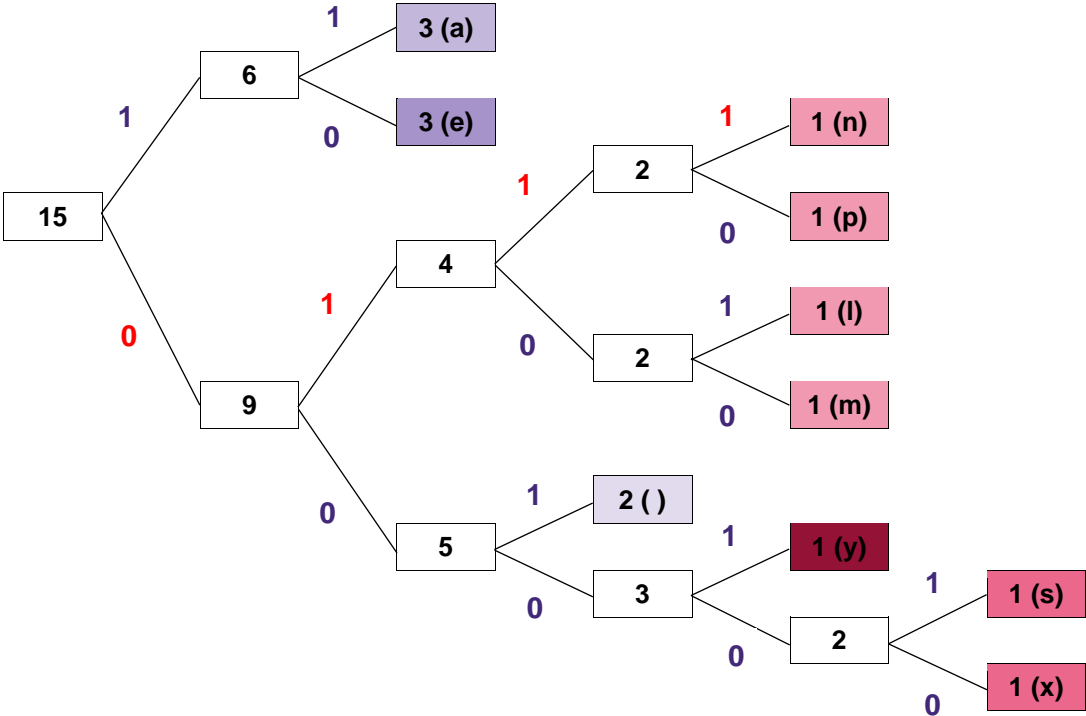


Step 9



### 3.3 Fundamentals of data representation

**Step 10 – Add labelling**



Coding produced	
<i>l</i>	0101
<i>m</i>	0100
<i>n</i>	0111
<i>p</i>	0110
<i>s</i>	00001
<i>x</i>	00000
<i>y</i>	0001
<i>space</i>	001
<i>a</i>	11
<i>e</i>	10