

HISTOGRAMS & STEMPLOTS

Describing Distributions
One Quantitative Variable
Unit 1: Exploratory Data Analysis



Although we rarely if ever need to construct graphs by hand, we will discuss the process of creating a histogram and stemplot in order to help you develop a deeper understanding of these visual displays

Histogram Example: Exam Grades

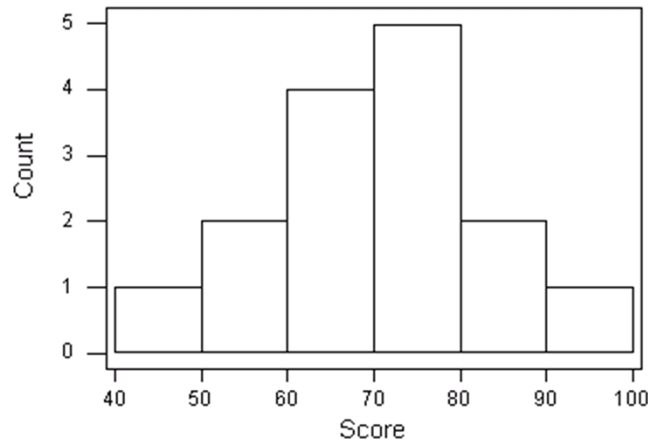
88, 48, 60, 51, 57,	<u>Score</u>	<u>Count</u>
	[40-50)	1
	[50-60)	2
85, 69, 75, 97, 72,	[60-70)	4
	[70-80)	5
	[80-90)	2
71, 79, 65, 63, 73	[90-100)	1

We will use this data on exam grades for 15 students to illustrate creating a histogram.

We break the data into intervals, being sure to define them so that no value is missed or counted in more than one interval. Computers use algorithms to achieve reasonable intervals for a wide range of datasets and many different algorithms are in use. We won't worry about this detail as we will use software to obtain our histograms in practice.

Once we have well-defined intervals, we count how many observations fall into each interval. For this example, we have intervals which begin at 40 and go to 100 by 10. The first interval consists of any value 40 or more but strictly less than 50. Any value of 50 will be placed in the next interval. In this dataset we do have one value, 60, which falls on the "boundary" and you can see that this value is counted in the interval [60, 70). Here we use the algebra notation with a bracket to indicate inclusion of the lower bound (60) and parentheses to indicate the exclusion of the upper bound (70).

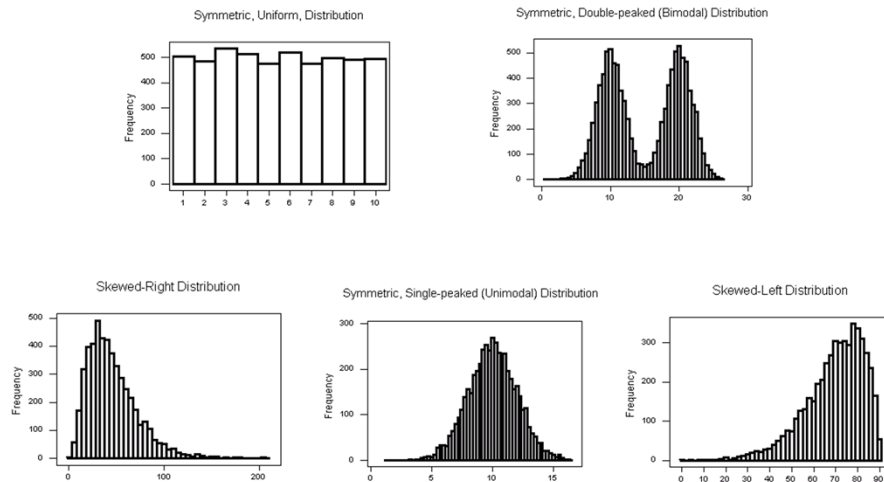
Histogram Example: Exam Grades



Once we have the frequency in each interval, we display this information in a histogram using the height of each bar to represent the frequency or percentage. Unlike bar charts, histograms have no gaps between the bars.

Although some information is lost - in that we can not recover the exact data - we can use these graphs to answer some questions about the data and to help us visualize the distribution.

Histograms



Here are a few other histograms to get you thinking about the possible variety we might see in distributions.

Stemplot

34 34 26 37 42 41 35 31

41 33 30 74 33 49 38 61

21 41 26 80 43 29 33 35

45 49 39 34 26 25 35 33

To create a stemplot, we first need to have data in which all values are rounded the same, we have that here with the data containing the ages of best actress Oscar winners.

Stemplot

34 34 26 37 42 41 35 31

41 33 30 74 33 49 38 61

21 41 26 80 43 29 33 35

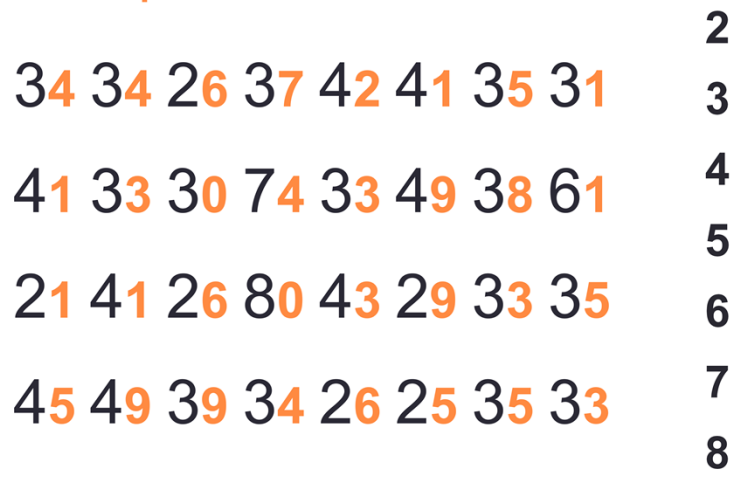
45 49 39 34 26 25 35 33

The right most digit becomes the leaf – presented here in bold and orange

All remaining digits become the stems – seen here in black with a larger font.

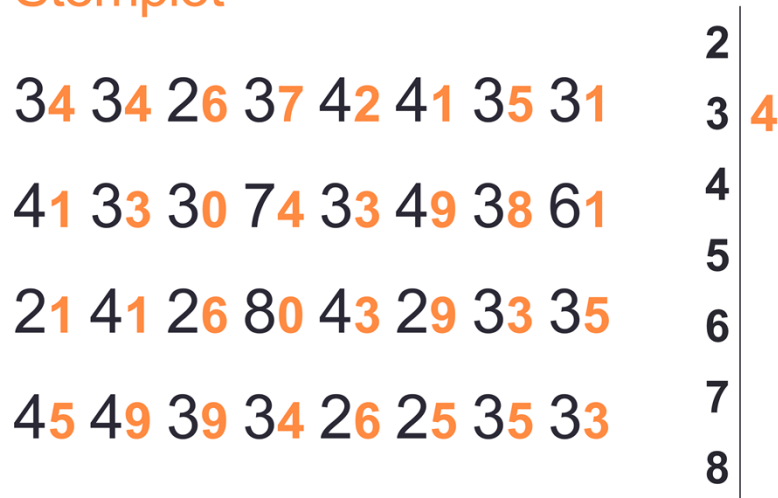
Our stems range from 2 through 8

Stemplot



We write the stems to the left of a vertical line, starting at the top with the smallest stem, we put the stems in order. Don't skip any stems in the range of your data – even if there is no data for that stem.

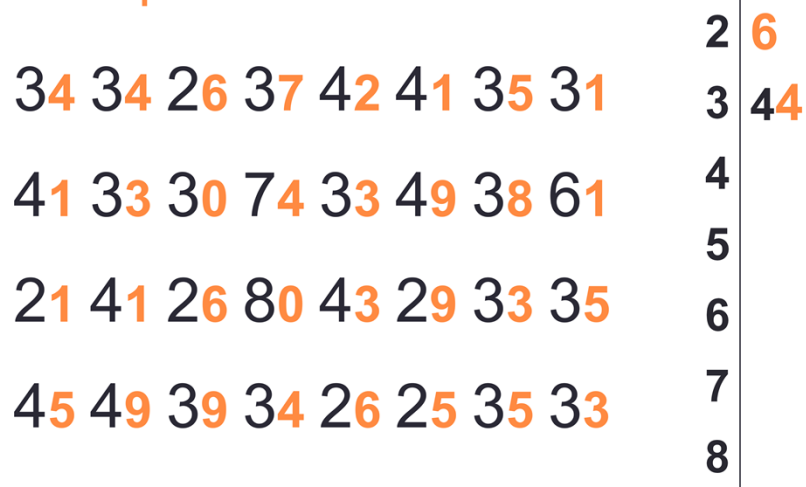
Stemplot



Then we add the leaves for each number to the right of the stem.

For 34, we place a leaf of 4 next to the stem for 3.

Stemplot



There was another 34 followed by a 26.

Stemplot

34 34 26 37 42 41 35 31

41 33 30 74 33 49 38 61

21 41 26 80 43 29 33 35

45 49 39 34 26 25 35 33

2	6
3	44751
4	21
5	
6	
7	
8	

Continuing in this way through the first row we would get what you see here on the right

You can pause the video and finish the rest for yourself. Next we will show the final steps.

Stemplot

- 34 34 26 37 42 41 35 31 41 33 30 74 33 49 38 61 21 41
26 80 43 29 33 35 45 49 39 34 26 25 35 33

steps 1, 2 and 3

```
2|616965
3|447513038359453
4|21191359
5|
6|1
7|4
8|0
```

==>

step 4

```
2|156669
3|013333444555789
4|11123599
5|
6|1
7|4
8|0
```

When we put all leaves in the display for our data we will have the stemplot on the lower left. The leaves are not in order (unless our original data was in order).

Next we sort the leaves to provide the final stemplot.

Note that it preserves the original data AND sorts the data.

These are two nice features of using a stemplot

(definitely helpful if you are ever stuck on a deserted island needing to quickly visualize the distribution of one quantitative variable!)

Stemplot – Split Stems

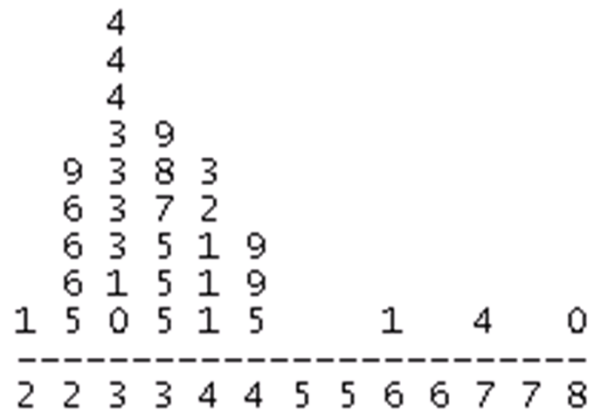
2 156669		2 1
3 013333444555789		2 56669
4 11123599		3 013333444
5	==>	3 555789
6 1	*	4 11123
7 4		4 599
8 0		5
		5
		6 1
		6
		7 4
		7
		8 0

One additional option is to split the stems. This is useful whenever there are a large number of values in only a few stems.

There are a few easy ways to split stems. The easiest is to split the leaves for each stem equally into two groups of 5 digits; 0-4, and 5-9 which we have done here in the stemplot on the right. This gives two sets of leaves for each stem. You can see that we repeat the stem on the left. If a lower or upper stem is not needed for the data it does not need to be displayed, as for the stem of 8 in our display – we only needed the first 8 as the largest observation was 80.

Leaves for each stem can also be split into 5 groups of 2 digits if there is a large amount of data (0-1, 2-3, 4-5, 6-7, 8-9) and some software packages split stems in almost any way they choose!

Stemplots



When we rotate the stemplot 90 degrees counterclockwise, the plot resembles a histogram.

Interactive Applet

- We can Analyze One Quantitative Variable with this

[One-Variable Statistical Calculator](#)

- From online content for [Introduction to the Practice of Statistics](#), Seventh Edition (Moore and McCabe)



Remember the interactive applet for analyzing one variable, it might help you pull together the ideas in this section on describing the distribution of one quantitative variable.

{web address of One-Variable Statistical Calculator applet:

http://content.bfwpub.com/webroot_pubcontent/Content/BCS_4/IPS7e/Student/Statistical%20Applets/onevar.html

}

{web address of textbook: <http://bcs.whfreeman.com/ips7e> }



HISTOGRAMS & STEMPLOTS

**Describing Distributions
One Quantitative Variable**
Unit 1: Exploratory Data Analysis

Histograms and to a lesser degree stemplots are helpful displays for visualizing the distribution of one quantitative variable.