

STAT 2593

Lecture 002 - Pictorial and Tabular Methods in Statistics

Dylan Spicker

Pictorial and Tabular Methods in Statistics

Learning Objectives

1. Define and characterize the distribution of a dataset.
2. Understand sample notation.
3. Explain the purposes of visualization.
4. Interpret and explain stem-and-leaf plots, dotplots, and histograms.

The Central Problem

Raw data are effectively useless for interpretation.

	id	ht	age	baseht	baseage	logfev1				
1	1,1	20000004768372	9	3415002822876	1	20000004768372	9	3415002822876	0	215110003948212
2	1,1	27999997138977	10	3929004669189	1	20000004768372	9	3415002822876	0	371560007333755
3	1,1	33000004291534	11	4524002075195	1	20000004768372	9	3415002822876	0	488579988479614
4	1,1	41999995708466	12	460000038147	1	20000004768372	9	3415002822876	0	751420021057129
5	1,1	48000001907349	13	4181995391846	1	20000004768372	9	3415002822876	0	832910001277924
6	1,1	5	15	4743003845215	1	20000004768372	9	3415002822876	0	892000019550323
7	1,1	51999998092651	16	3722991943359	1	20000004768372	9	3415002822876	0	871290028095245
8	2,1	12999999523163	6	58729982376099	1	12999999523163	6	58729982376099	0	307480007410049
9	2,1	19000005722046	7	6496000289917	1	12999999523163	6	58729982376099	0	350659996271133
10	2,1	49000000953674	12	7391996383667	1	12999999523163	6	58729982376099	0	756120026111603
11	2,1	52999997138977	13	7741003036499	1	12999999523163	6	58729982376099	0	86710000038147
12	2,1	54999995231628	14	6940002441406	1	12999999523163	6	58729982376099	1	04732000827789
13	2,1	55999994277954	15	8219995498657	1	12999999523163	6	58729982376099	1	15373003482819
14	2,1	57000005245209	16	6679992675781	1	12999999523163	6	58729982376099	0	924260020256042
15	2,1	57000005245209	17	631799697876	1	12999999523163	6	58729982376099	1	13461995124817
16	3,1	17999994754791	6	91309976577759	1	17999994754791	6	91309976577759	0	431780010461807
17	3,1	23000001907349	7	0753999710083	1	17999994754791	6	91309976577759	0	385259985923767
18	3,1	29999995231628	8	9665002822876	1	17999994754791	6	91309976577759	0	598839998245239
19	3,1	35000002384186	9	98770046234131	1	17999994754791	6	91309976577759	0	751420021057129
20	3,1	47000002861023	11	0773000717163	1	17999994754791	6	91309976577759	0	96697998046875
21	3,1	57000005245209	13	0677995681763	1	17999994754791	6	91309976577759	0	89608971065521
22	3,1	5900000333786	14	1027002334595	1	17999994754791	6	91309976577759	1	018884996891022
23	3,1	60000002384186	15	0801000595093	1	17999994754791	6	91309976577759	1	105260014534
24	3,1	60000002384186	16	0163993835449	1	17999994754791	6	91309976577759	1	08519005775452
25	4,1	14999997615814	6	75979995727539	1	14999997615814	6	75979995727539	0	0582699999213219
26	4,1	21000003814697	7	82200002670288	1	14999997615814	6	75979995727539	0	18231999874115
27	4,1	25999999046326	8	81309986114502	1	14999997615814	6	75979995727539	0	277630001306534
28	4,1	30999994277954	9	83440017700195	1	14999997615814	6	75979995727539	0	444689989809966
29	4,1	38999998596489	10	923999786377	1	14999997615814	6	75979995727539	0	576610028743744
30	4,1	46000003814697	11	9315996170044	1	14999997615814	6	75979995727539	0	672940015792847
31	4,1	53999996185303	12	9117002487183	1	14999997615814	6	75979995727539	0	722710013389587
32	4,1	5900000333786	13	9465999603271	1	14999997615814	6	75979995727539	1	02244997024536
33	4,1	60000002384186	14	8664999008179	1	14999997615814	6	75979995727539	1	03673994541168
34	4,1	62999999523163	17	7877998352051	1	14999997615814	6	75979995727539	1	1878399848938
35	5,1	11000001430511	6	50239992141724	1	11000001430511	6	50239992141724	0	029559999704361
36	5,1	14999997615814	7	56470012664795	1	11000001430511	6	50239992141724	0	113329999148846
37	5,1	51999998092651	13	733099937439	1	11000001430511	6	50239992141724	0	896089971065521
38	5,1	53999996185303	14	7049999237061	1	11000001430511	6	50239992141724	0	862890005111694
39	5,1	54999995231628	15	5866003036499	1	11000001430511	6	50239992141724	0	955510020256042
40	5,1	54999995231628	16	6460999095762	1	11000001430511	6	50239992141724	0	924260020256042
41	5,1	55999994277954	17	5221996307373	1	11000001430511	6	50239992141724	0	832910001277924
42	6,1	24000000953674	6	89940023422241	1	24000000953674	6	89940023422241	0	262360006570816
43	6,1	29999995231628	7	96169996261597	1	24000000953674	6	89940023422241	0	47622999548912
44	6,1	36000001430511	8	98560047149658	1	24000000953674	6	89940023422241	0	565310001373291
45	6,1	4099999666214	9	97399997711182	1	24000000953674	6	89940023422241	0	712949991226196
46	6,1	47000002861023	11	06369972229	1	24000000953674	6	89940023422241	0	774730026721954
47	6,1	55999994277954	12	0712003707886	1	24000000953674	6	89940023422241	0	900160014629364
48	6,1	57000005245209	13	0753999710083	1	24000000953674	6	89940023422241	0	980360003836363

The Solution

We use **graphical displays** to summarize the useful information instead.

The **distribution** of a dataset describes the possible values that are in a sample, and the relative frequency of those values.

Notation

- ▶ Our data consist of n observations.
- ▶ Each observation is denoted with a lowercase x .
- ▶ We use subscripts to enumerate sample observations, x_i , for $i = 1, \dots, n$.

Goals for Visualizations

1. Identify what is **typical** in the data.

Goals for Visualizations

1. Identify what is **typical** in the data.
2. Identify the **spread** of the data.

Goals for Visualizations

1. Identify what is **typical** in the data.
2. Identify the **spread** of the data.
3. Determine if there are any **gaps** in the data.

Goals for Visualizations

1. Identify what is **typical** in the data.
2. Identify the **spread** of the data.
3. Determine if there are any **gaps** in the data.
4. Identify the shape of the **distribution** of the data.

Goals for Visualizations

1. Identify what is **typical** in the data.
2. Identify the **spread** of the data.
3. Determine if there are any **gaps** in the data.
4. Identify the shape of the **distribution** of the data.
5. Identify **peaks** in the data.

Goals for Visualizations

1. Identify what is **typical** in the data.
2. Identify the **spread** of the data.
3. Determine if there are any **gaps** in the data.
4. Identify the shape of the **distribution** of the data.
5. Identify **peaks** in the data.
6. Determine whether there are any **outliers**.

Stem-and-Leaf Plots

2, 17, 14, 35, 37, 44, 2, 47, 41, 35, 50,
40, 16, 20, 22, 22, 23, 9, 23, 23, 48, 34,
8, 26, 11, 4, 49, 11, 39, 34, 29, 4, 4, 35,
6, 3, 5, 40, 28, 40, 45, 27, 5, 37, 1, 27,
16, 20, 19, 50, 10, 50, 19, 20, 28, 45,
40, 4, 32, 25

2, 17, 14, 35, 37, 44, 2, 47, 41, 35, 50, 40, 16, 20, 22, 22, 23, 9, 23, 23, 48, 34, 8, 26,
11, 4, 49, 11, 39, 34, 29, 4, 4, 35, 6, 3, 5, 40, 28, 40, 45, 27, 5, 37, 1, 27, 16, 20, 19,
50, 10, 50, 19, 20, 28, 45, 40, 4, 32, 25

The decimal point is 1 digit(s) to the right of the |

0 | 12234444

0 | 55689

1 | 0114

1 | 66799

2 | 00022333

2 | 5677889

3 | 244

3 | 555779

4 | 000014

4 | 55789

5 | 000

Constructing Stem-and-Leaf Plots

1. Identify the stem (leading digit) and leaves (remaining digits)

Constructing Stem-and-Leaf Plots

1. Identify the stem (leading digit) and leaves (remaining digits)
2. Place the stems from smallest to largest in a vertical column

Constructing Stem-and-Leaf Plots

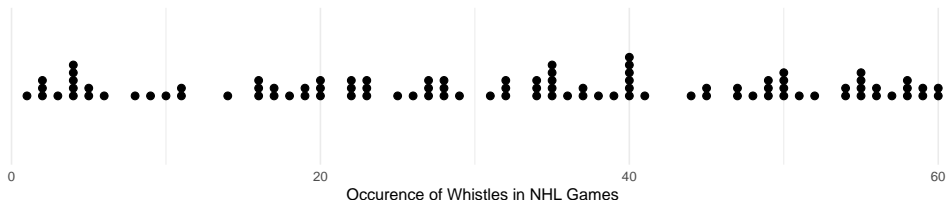
1. Identify the stem (leading digit) and leaves (remaining digits)
2. Place the stems from smallest to largest in a vertical column
3. Place each leaf with the corresponding stem in ascending order

Constructing Stem-and-Leaf Plots

1. Identify the stem (leading digit) and leaves (remaining digits)
2. Place the stems from smallest to largest in a vertical column
3. Place each leaf with the corresponding stem in ascending order
4. Indicate the units (where is the decimal point?)

Dot Plots

2, 51, 17, 56, 14, 35, 37, 44, 2, 47, 41, 35, 50, 40, 16, 20, 22, 22, 23, 9, 23, 23, 48,
34, 8, 26, 56, 11, 4, 58, 58, 49, 11, 60, 39, 34, 29, 4, 4, 55, 35, 52, 6, 3, 5, 40, 28, 40,
45, 59, 58, 27, 5, 37, 1, 27, 55, 16, 20, 19, 50, 10, 50, 19, 20, 28, 55, 54, 45, 40, 57,
54, 4, 32, 59, 25, 35, 60, 27, 28, 4, 32, 50, 55, 47, 31, 49, 2, 22, 32, 38, 18, 17, 40,
35, 34, 16, 49, 40, 36



Constructing Dot Plots

1. Plot all observations on the vertical axis at their value.

Constructing Dot Plots

1. Plot all observations on the vertical axis at their value.
2. Stack dots for repeated observations.

Histograms

- ▶ Histograms represent the **frequency distribution** using bar charts.

Histograms

- ▶ Histograms represent the **frequency distribution** using bar charts.
- ▶ For discrete variables, we use each category, and count the observations in each category.

Histograms

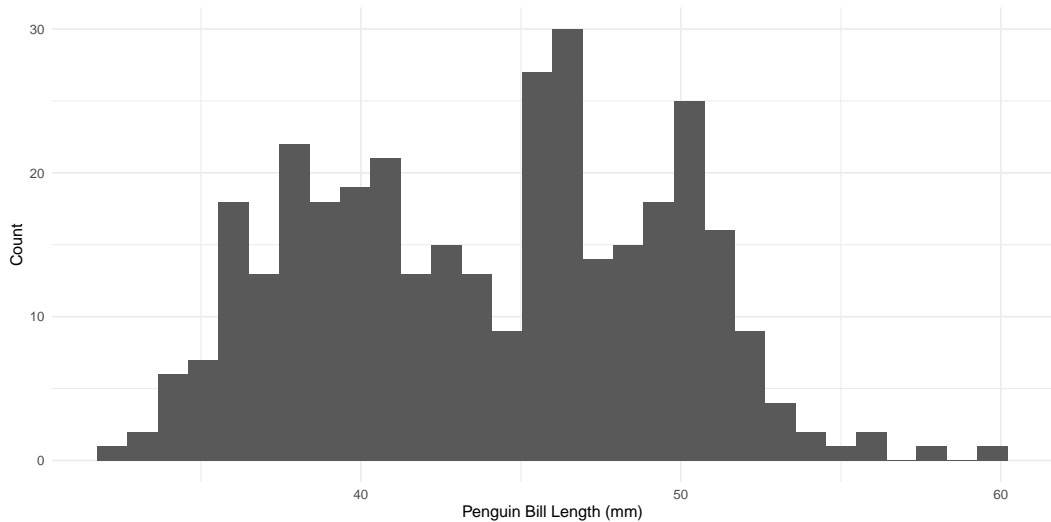
- ▶ Histograms represent the **frequency distribution** using bar charts.
- ▶ For discrete variables, we use each category, and count the observations in each category.
- ▶ For continuous variables, we **bin** the observations, and count the observations in each bin.

Histograms

- ▶ Histograms represent the **frequency distribution** using bar charts.
- ▶ For discrete variables, we use each category, and count the observations in each category.
- ▶ For continuous variables, we **bin** the observations, and count the observations in each bin.
- ▶ Can also use the **relative frequency**, which is given by

$$\text{relative frequency} = \frac{\text{frequency}}{\text{total number of observations}}.$$

Histograms



Considerations for Histograms

- ▶ Classes should be the same width.

Considerations for Histograms

- ▶ Classes should be the same width.
- ▶ Classes should not overlap.

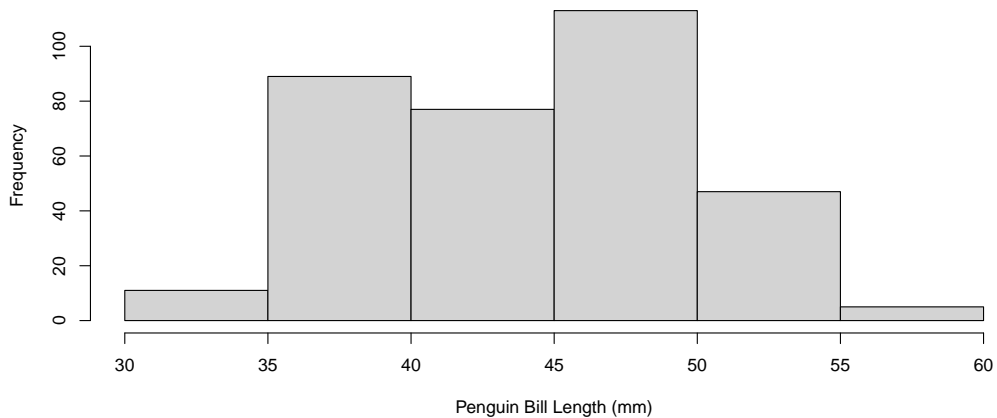
Considerations for Histograms

- ▶ Classes should be the same width.
- ▶ Classes should not overlap.
- ▶ Classes should include all possible values.

Histograms

	Bins	Counts	Relative Frequency
1	30 - 35	11	0.03216374
2	35 - 40	89	0.26023392
3	40 - 45	77	0.22514620
4	45 - 50	113	0.33040936
5	50 - 55	47	0.13742690
6	55 - 60	5	0.01461988

Histograms



What to Look for in a Histogram

- ▶ Shape, including **modality**, **symmetry**, and **skewness**

What to Look for in a Histogram

- ▶ Shape, including **modality**, **symmetry**, and **skewness**
 - ▶ How many peaks (modes) are there: one is **unimodal**, two is **bimodal**, three or more is **multimodal**

What to Look for in a Histogram

- ▶ Shape, including **modality**, **symmetry**, and **skewness**
 - ▶ How many peaks (modes) are there: one is **unimodal**, two is **bimodal**, three or more is **multimodal**
 - ▶ Does the distribution look like a mirror image? If not, it is skewed.

What to Look for in a Histogram

- ▶ Shape, including **modality**, **symmetry**, and **skewness**
 - ▶ How many peaks (modes) are there: one is **unimodal**, two is **bimodal**, three or more is **multimodal**
 - ▶ Does the distribution look like a mirror image? If not, it is skewed.
 - ▶ Left (or negative) skew has a tail pointing left; right (or positive) skew has a tail pointing right.

What to Look for in a Histogram

- ▶ Shape, including **modality**, **symmetry**, and **skewness**
 - ▶ How many peaks (modes) are there: one is **unimodal**, two is **bimodal**, three or more is **multimodal**
 - ▶ Does the distribution look like a mirror image? If not, it is skewed.
 - ▶ Left (or negative) skew has a tail pointing left; right (or positive) skew has a tail pointing right.
- ▶ The **centre** of the distribution

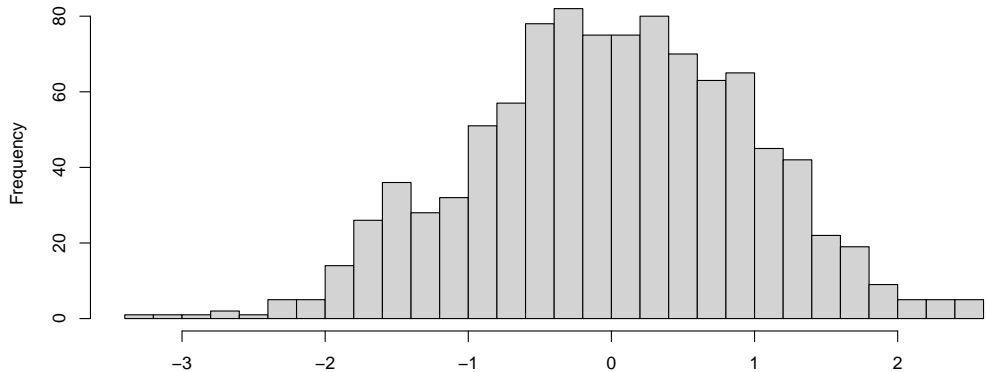
What to Look for in a Histogram

- ▶ Shape, including **modality**, **symmetry**, and **skewness**
 - ▶ How many peaks (modes) are there: one is **unimodal**, two is **bimodal**, three or more is **multimodal**
 - ▶ Does the distribution look like a mirror image? If not, it is skewed.
 - ▶ Left (or negative) skew has a tail pointing left; right (or positive) skew has a tail pointing right.
- ▶ The **centre** of the distribution
- ▶ The **spread** of the distribution

What to Look for in a Histogram

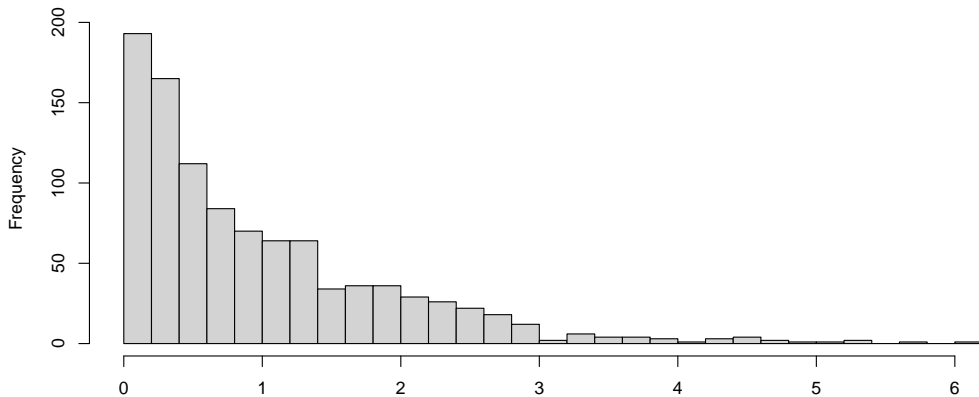
- ▶ Shape, including **modality**, **symmetry**, and **skewness**
 - ▶ How many peaks (modes) are there: one is **unimodal**, two is **bimodal**, three or more is **multimodal**
 - ▶ Does the distribution look like a mirror image? If not, it is skewed.
 - ▶ Left (or negative) skew has a tail pointing left; right (or positive) skew has a tail pointing right.
- ▶ The **centre** of the distribution
- ▶ The **spread** of the distribution
- ▶ **Outliers** or deviations from the general pattern

Examples



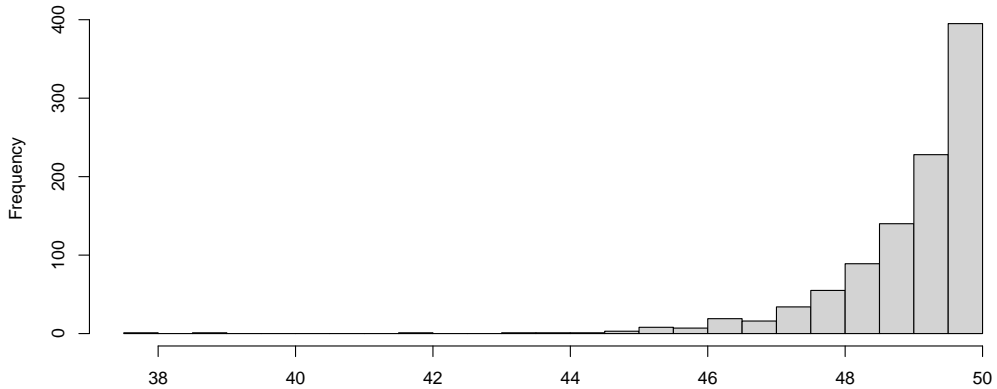
Histogram Example 1

Examples



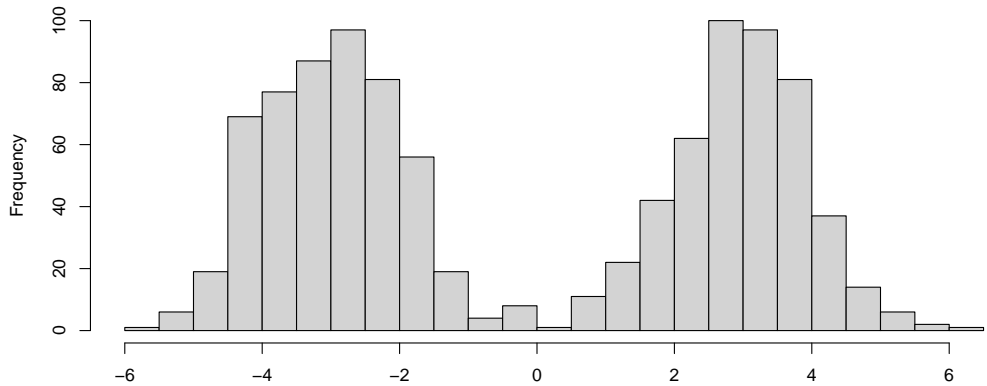
Histogram Example 2

Examples



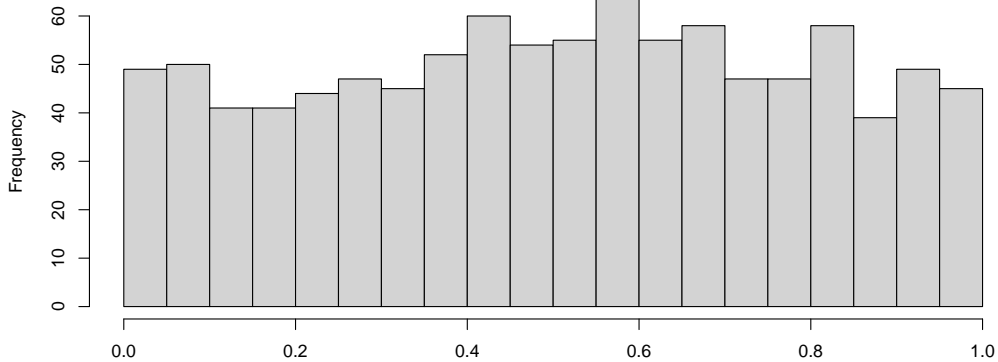
Histogram Example 3

Examples



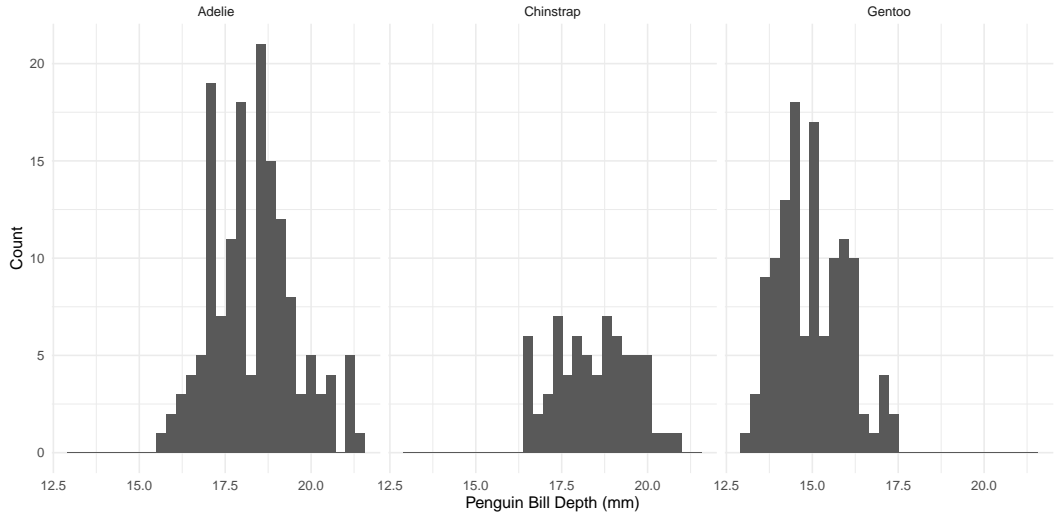
Histogram Example 4

Examples

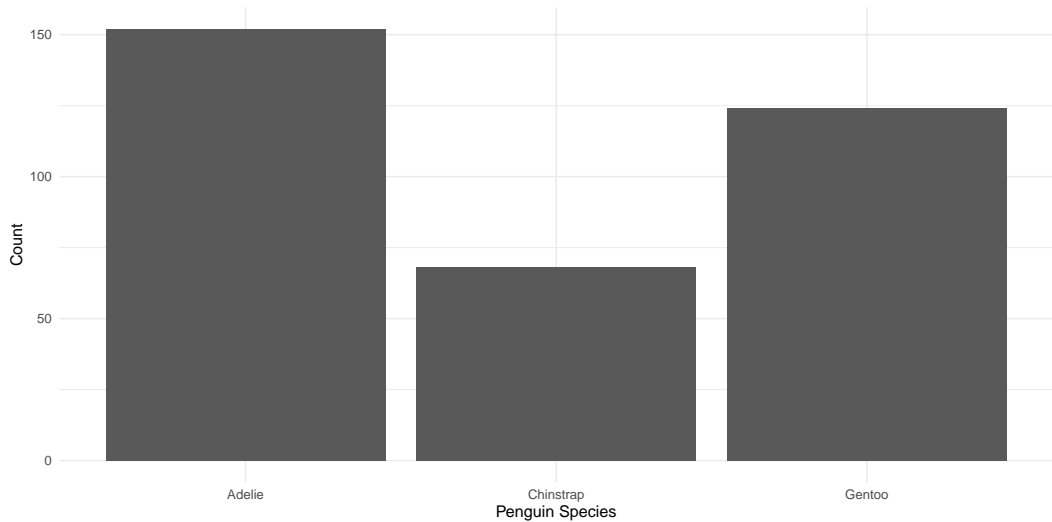


Histogram Example 5

Histograms: Comparing Distributions



Bar Charts: Histograms for Categorical Data



Summary

- ▶ Raw data are difficult to interpret on their own.
- ▶ Visualizations can help display the character of a distribution.
- ▶ Stem-and-leaf plots, dot plots, and histograms are all useful for quantitative variables.
- ▶ Histograms (as bar plots) can be used for categorical variables.