

*

Filling the Bins

- or -

Turning Numerical Data into Histograms

ISC1057

Janet Peterson and John Burkardt

Computational Thinking

Fall Semester 2016

A histogram is a kind of chart which shows patterns in sets of numeric data.

Strictly speaking, a histogram is different from a bar chart, because a bar chart is based on categories: the number of times each baseball team has won a World Series, for instance. The categories are given, and we just have to assign a number to each one.

For a histogram, we might have measurements of the circumference of 1,000 trees on campus, that is, just series of 1,000 numbers. These numbers might all be different. If we made a bar chart of this data, then we would end up with 1,000 little bars, one for each measurement.

But if the numbers are measurements, then it might make sense to group them together, separating small, medium and large numbers, or, more carefully, dividing the values into ranges, or **bins**.

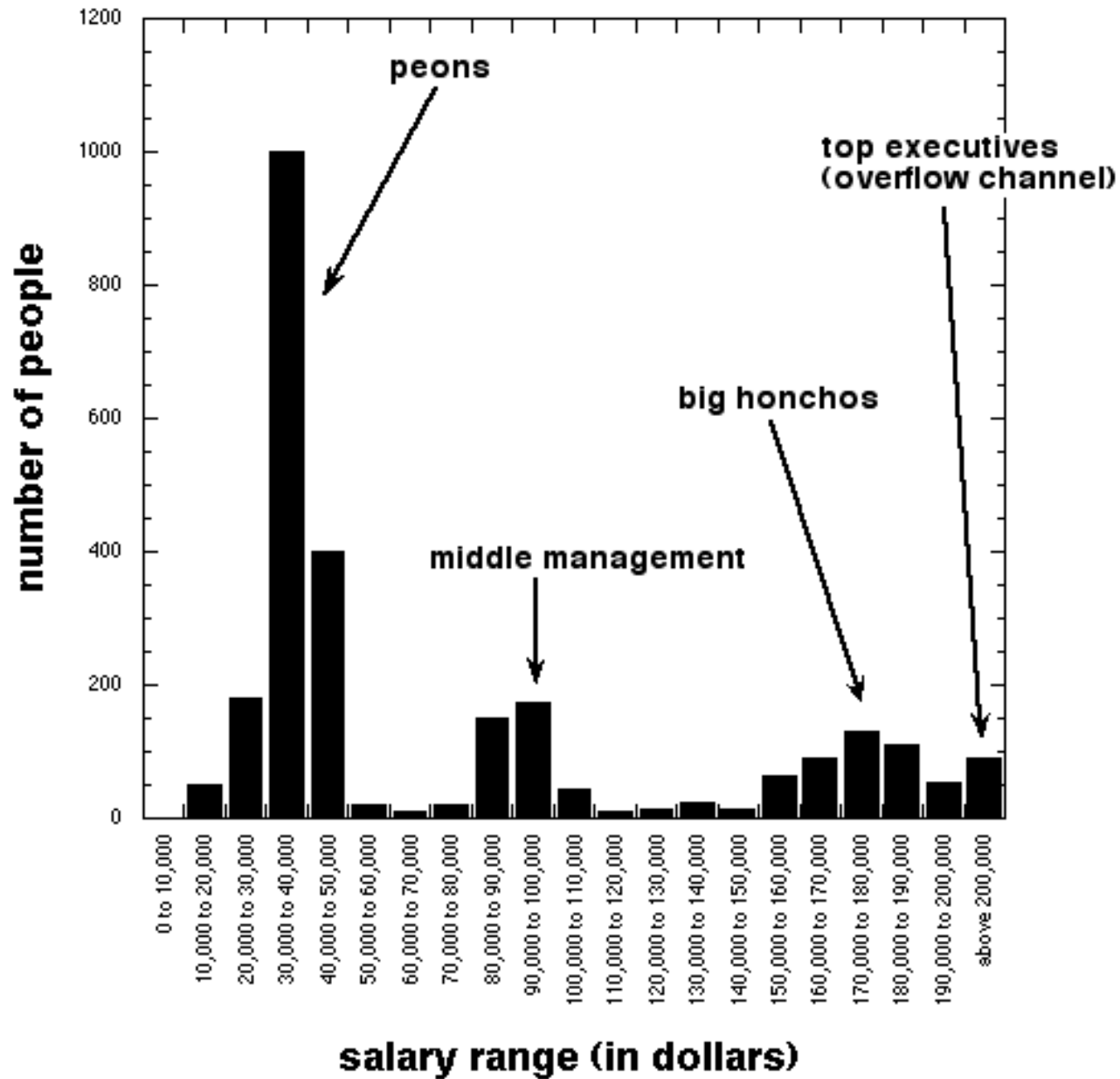
The size and number of bins we use will have a significant effect on whether we can see patterns in the data. We don't want every number in a separate bin, or all the numbers in a single bin. We are looking for natural patterns that the numbers arrange themselves into when we give them a "reasonable" number of bins to choose from.

Sample data:

Employee Number	Position	Annual Salary
1	"Clerk"	32,000
2	"Stocker"	17,125
3	"Senior Manager"	143,000
...	more data	—
2108	"Vice President for Sales"	245,000
2109	"Cleaner"	15,125
2110	"Data Analyst"	74,000

It's possible that no two people in this company earn the same salary. What if we consider salary ranges of \$10,000, and count the number of employees in each range?

MegaMart Salary Histogram



In a histogram, it's usual to make the bin widths or ranges the same. That means we still have to decide how many bins to use, and so sometimes a little experimentation is necessary.

However, just as with bar charts, if the data is properly plotted, we may find that it is:

- pretty much flat or even;
- arranged like a staircase (mostly up, or mostly down);
- arranged like a hill (up, then down);
- arranged like a valley (down, then up).
- having several peak values separated by valleys (*as in our salary data on the previous slide*);
- seemingly random (up and down with no pattern);

A histogram may summarize or illustrate our beliefs, or it may be evidence that asks us to think of an explanation.

Example #1: 200 Height and Weight Measurements

It's a common practice to collect height and weight measurements.

We assume there is an average height, with some range of normal variation. But are there more short people than tall ones?

It's not so easy to guess someone's weight, and the range of variation is probably much greater than for height.

Because both measurements are numbers that are scattered over a numeric range, a histogram can be a good tool to find patterns.

We have 200 height and weight measurements in the file `hw_200.csv`.

"Index",	"Height (inches)",	"Weight (pounds)"
1,	65.78,	112.99
2,	71.52,	136.49
3,	69.40,	153.03
—	<i>more data</i>	—
198,	68.24,	128.30
199,	68.02,	127.47
200,	71.39,	127.88

starter data + NEW GRID IMPORT

ADD DATA SAVE COPY EXPORT UNDO REDO CHOOSE PLOT TYPE DATA TOOLS ANALYSIS

LINE PLOT

Click the column headers to choose x and y columns to graph. Use the different colors to match x columns with y columns.

Line plot

OPTIONS

Error Bars

	Col1	Col2	Col3	Col4
x	choose as x	choose as x	choose as x	choose as x
y	choose as y	choose as y	choose as y	choose as y
1	giraffes	12	20	
2	orangutans	18	14	
3	monkeys	29	23	
4				
5				
6				
7				
8				
9				
10				
11				
12				

starter data

+ NEW GRID ↑ IMPORT

+
ADD DATA

File Upload

Navigation: < jburkardt public_html latex **ct_histogram**

Places: Search, Recently Used, jburkardt, Desktop, File System, BIGBOY

Name	Size	Modified
histogram.tex.backup	10.2 KB	09:26
hw_200.csv	3.6 KB	Yesterday at 22:03
hw_25000.csv	618.5 KB	Yesterday at 22:03
oscar_age_female.csv	4.3 KB	Yesterday at 15:33
oscar_age_female.png	17.9 KB	Sunday
oscar_age_female_bin1...	20.8 KB	Yesterday at 15:14
oscar_age_male.csv	3.8 KB	Sunday
oscar_age_male.png	17.6 KB	Sunday
oscar_female_screen1.p...	201.9 KB	Yesterday at 15:44
oscar_female_screen2.p...	54.0 KB	Yesterday at 15:46

File type filter: All Files

Buttons: Cancel, Open

Click the x and y of different with y co

OPTION

Er

Asymmetric Errors

starter data hw_200.csv + NEW GRID IMPORT

ADD DATA SAVE COPY EXPORT UNDO REDO CHOOSE PLOT TYPE DATA TOOLS ANALYSIS

LINE PLOT

Click the column headers to choose x and y columns to graph. Use the different colors to match x columns with y columns.

Line plot

OPTIONS

- Error Bars
- Asymmetric Errors

	Index	Height(Inches)"	Weight(Pounds)
x	choose as x	choose as x	choose as x
y	choose as y	choose as y	choose as y
1	1	65.78	112.99
2	2	71.52	136.49
3	3	69.4	153.03
4	4	68.22	142.34
5	5	67.79	144.3
6	6	68.7	123.3
7	7	69.8	141.49
8	8	70.01	136.46
9	9	67.9	112.37
10	10	66.78	120.67
11	11	66.49	127.45
12	12	67.62	114.14
13	13	68.3	125.61

LINE PLOT

Click the column headers to choose x and y columns to graph. Use the different colors to match x columns with y columns.

Line plot

OPTI

Scatter plot

Bar chart

Asy Histogram

Area plot

Text

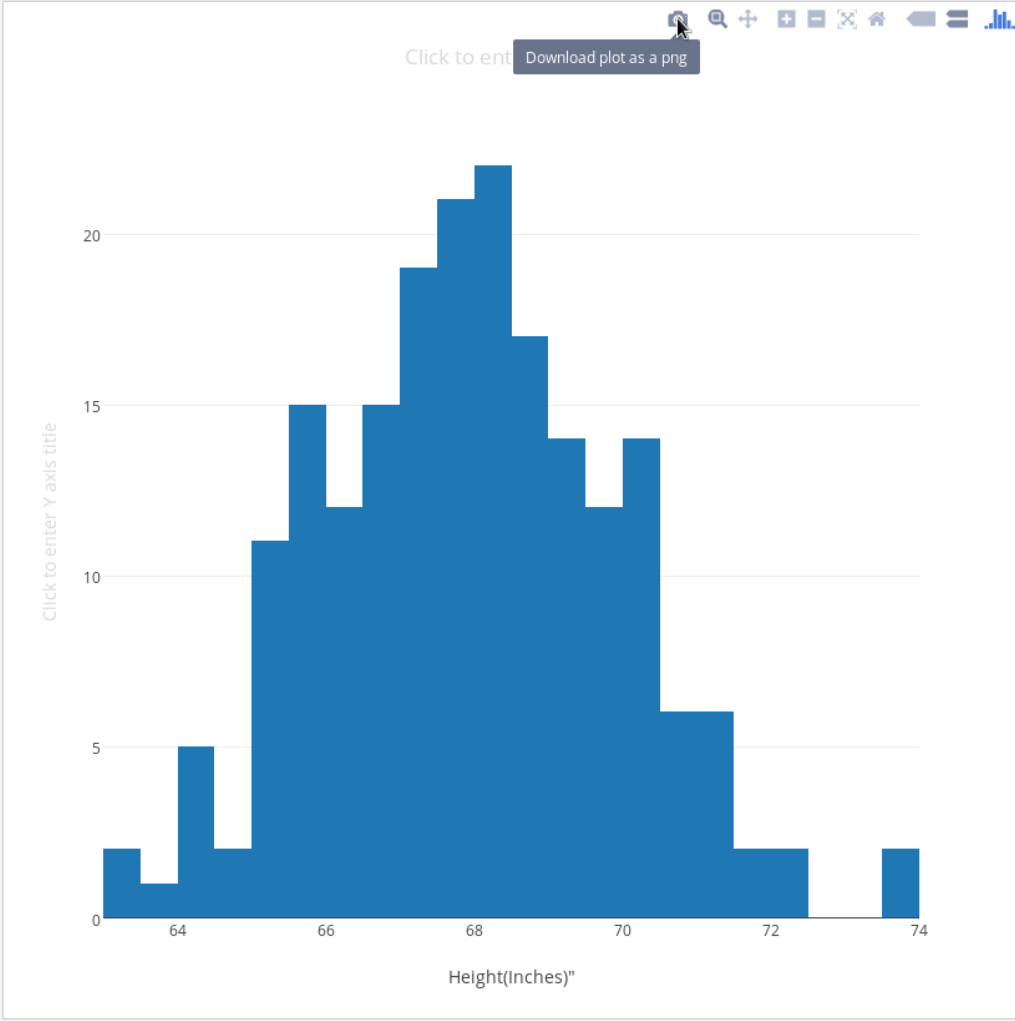
INSERT INTO

Make a new plot

	Index	Height(Inches)"	Weight(Pounds)	Col4	Col5	Col6
x	choose as x	choose as x	choose as x	choose as x	choose as x	choose as x
y	choose as y	choose as y	choose as y	choose as y	choose as y	choose as y
1	1	65.78	112.99			
2	2	71.52	136.49			
3	3	69.4	153.03			
4	4	68.22	142.34			
5	5	67.79	144.3			
6	6	68.7	123.3			
7	7	69.8	141.49			
8	8	70.01	136.46			
9	9	67.9	112.37			
10	10	66.78	120.67			
11	11	66.49	127.45			
12	12	67.62	114.14			
13	13	68.3	125.61			
14	14	67.12	122.46			
15	15	68.28	116.09			
16	16	71.09	140			
17	17	66.46	129.5			
18	18	68.65	142.97			
19	19	71.23	137.9			
20	20	67.13	124.04			
21	21	67.83	141.28			
22	22	68.88	143.54			
23	23	63.48	97.9			
24	24	68.42	129.5			
25	25	67.63	141.85			
26	26	67.21	129.72			
27	27	70.84	142.42			
28	28	67.49	131.55			
29	29	66.53	108.33			
30	30	65.44	113.89			

starter data hw_200.csv Plot Plot + NEW GRID IMPORT

- View data
- View JSON
- Share
- Traces
- Layout
- Axes
- Notes
- Legend
- Fit data
- Themes



starter data hw_200.csv Plot Plot + NEW GRID IMPORT

- View data
- View JSON
- Share
- Traces
- Layout
- Axes
- Notes
- Legend
- Fit data
- Themes

Traces

Height(Inches)"

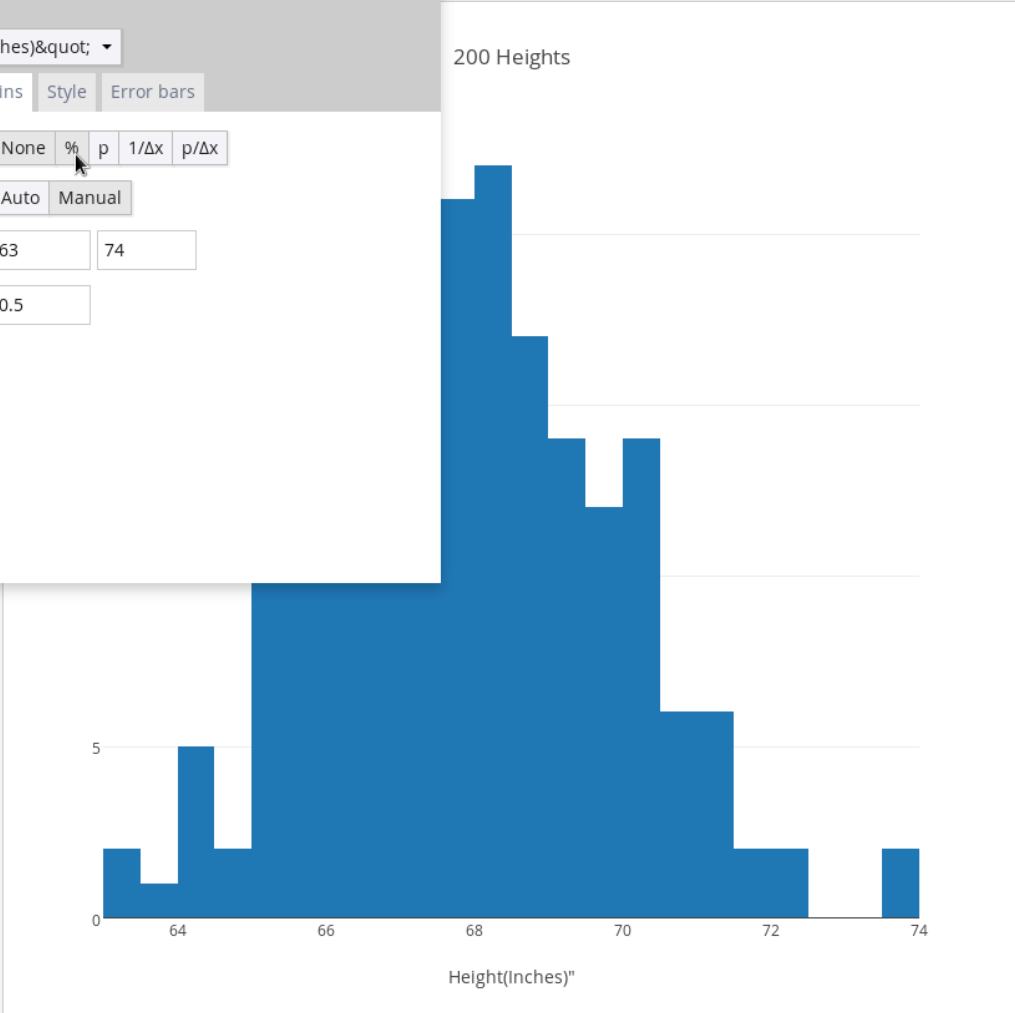
Mode Range/bins Style Error bars

Normalization None % p 1/Δx p/Δx

X bins Auto Manual

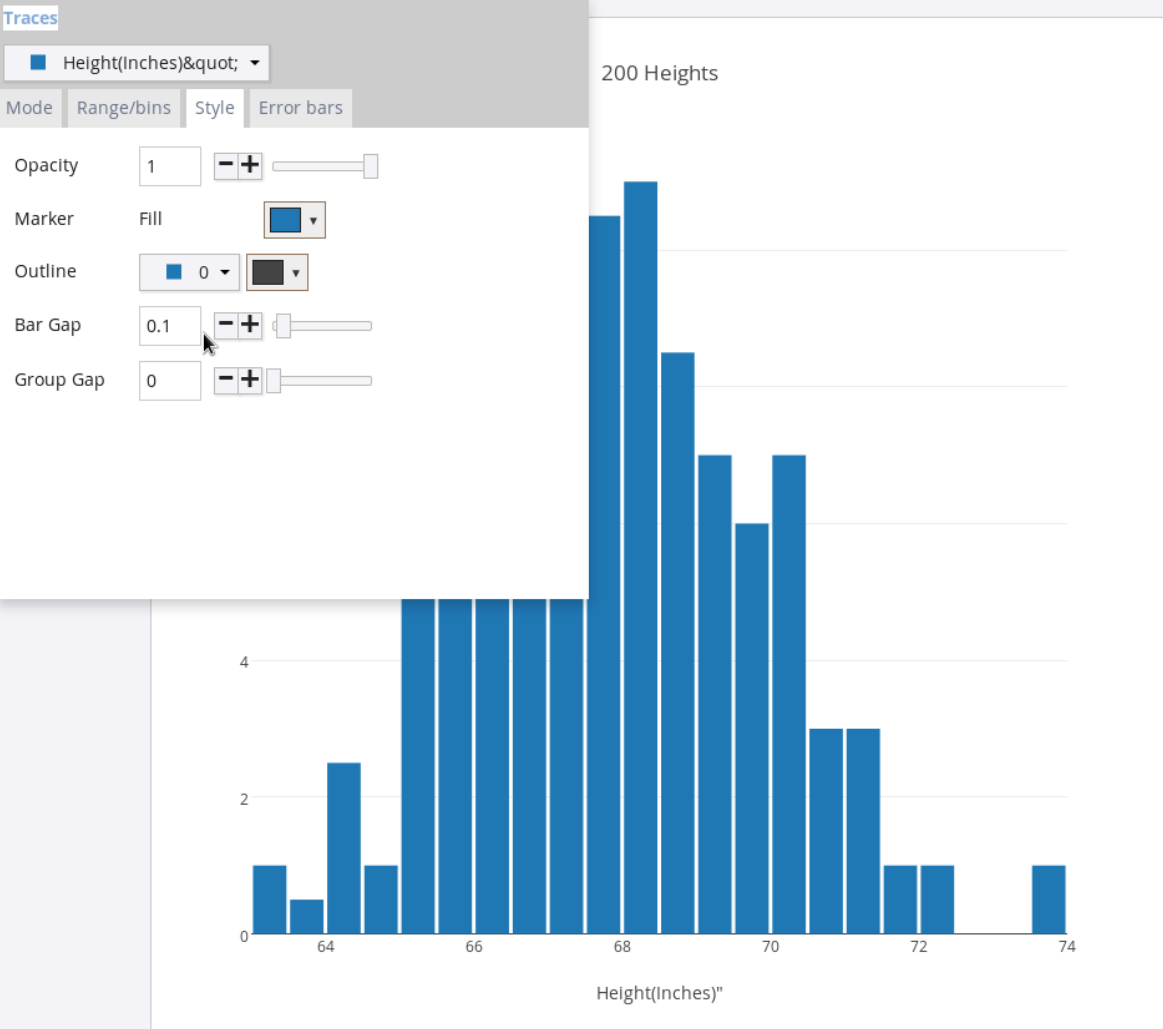
X range 63 74

X bin size 0.5



starter data hw_200.csv Plot Plot + NEW GRID IMPORT

- View data
- View JSON
- Share
- Traces
- Layout
- Axes
- Notes
- Legend
- Fit data
- Themes



starter data hw_200.csv Plot Plot + NEW GRID IMPORT

- View data
- View JSON
- Share
- Traces
- Layout
- Axes
- Notes
- Legend
- Fit data
- Themes

Traces

Height(Inches)"

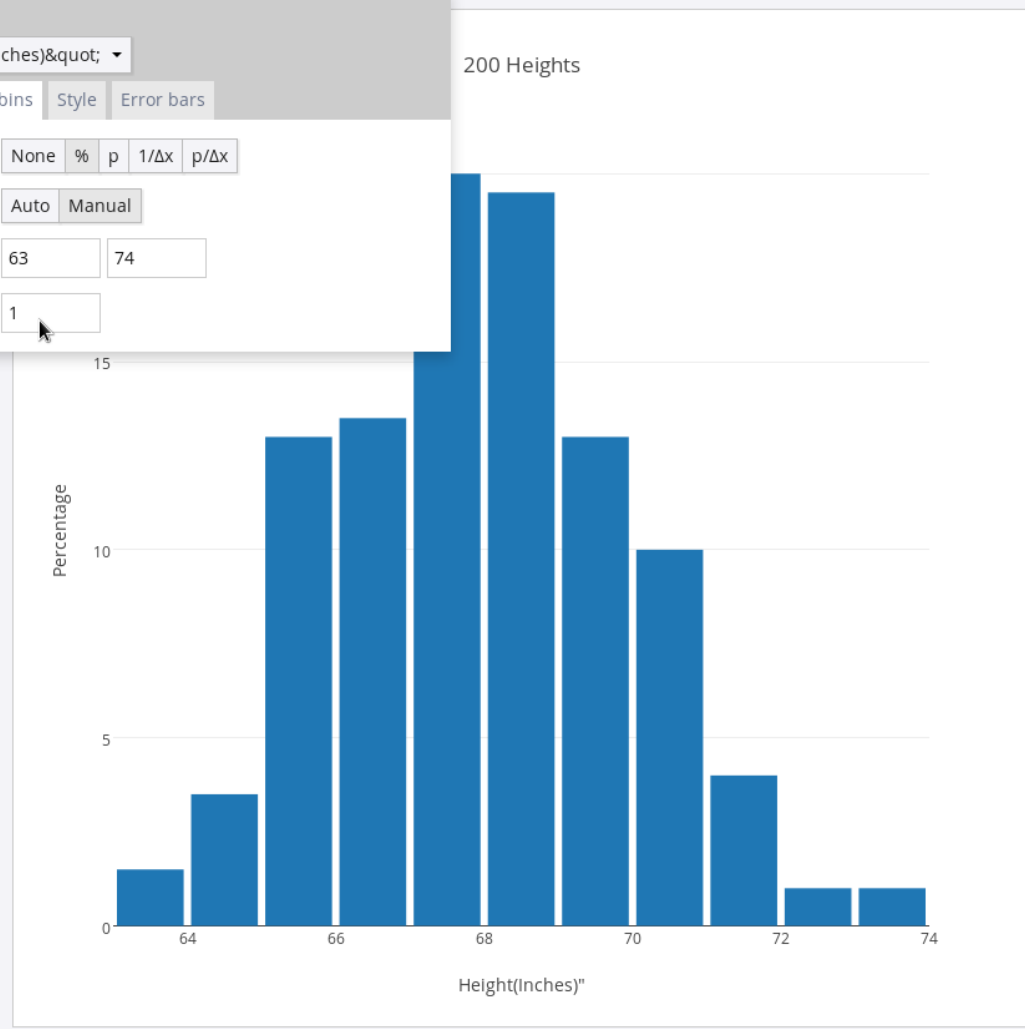
Mode Range/bins Style Error bars

Normalization None % p $1/\Delta x$ $p/\Delta x$

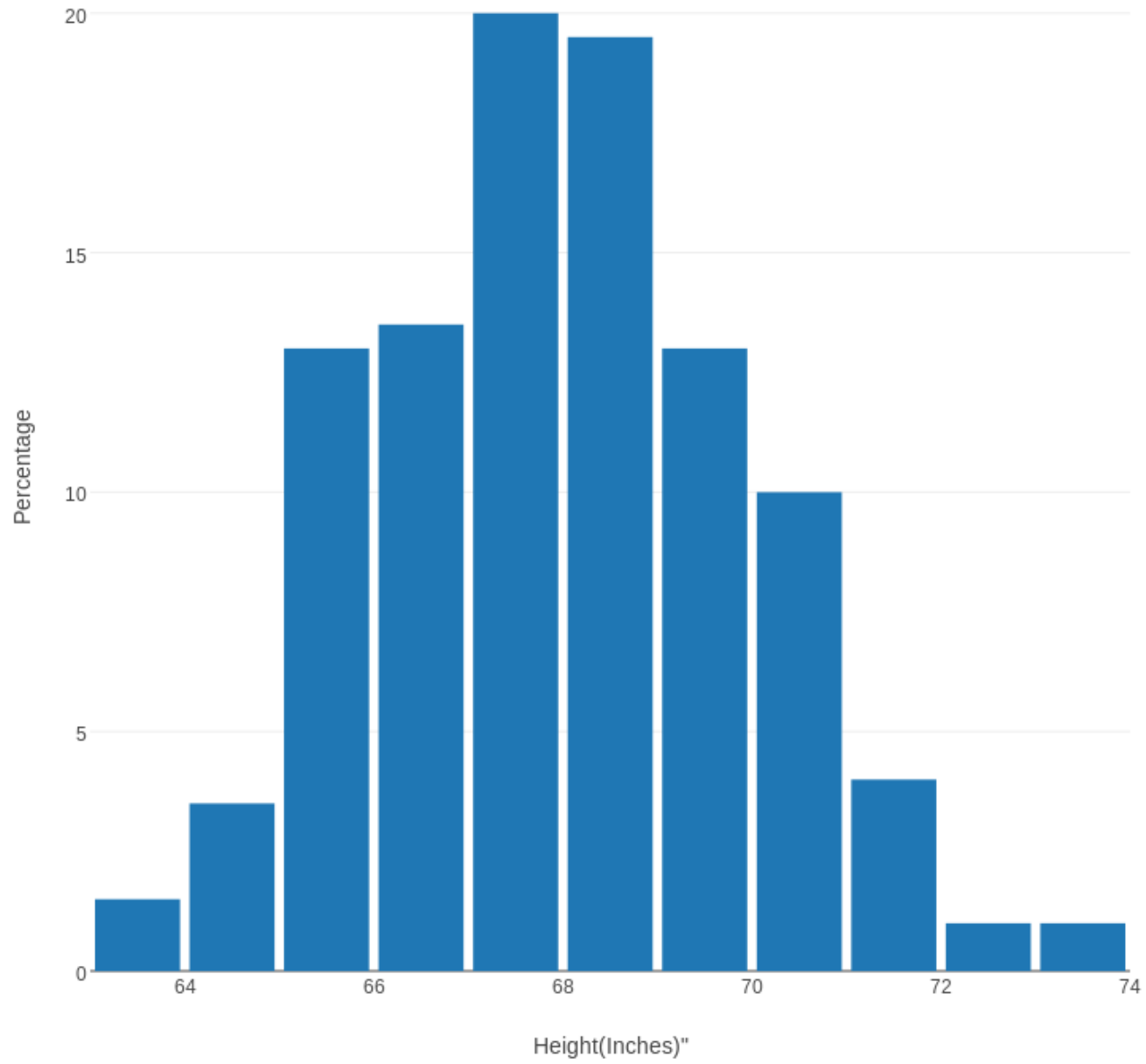
X bins Auto Manual

X range 63 74

X bin size 1



200 Heights (Bin width = 1)



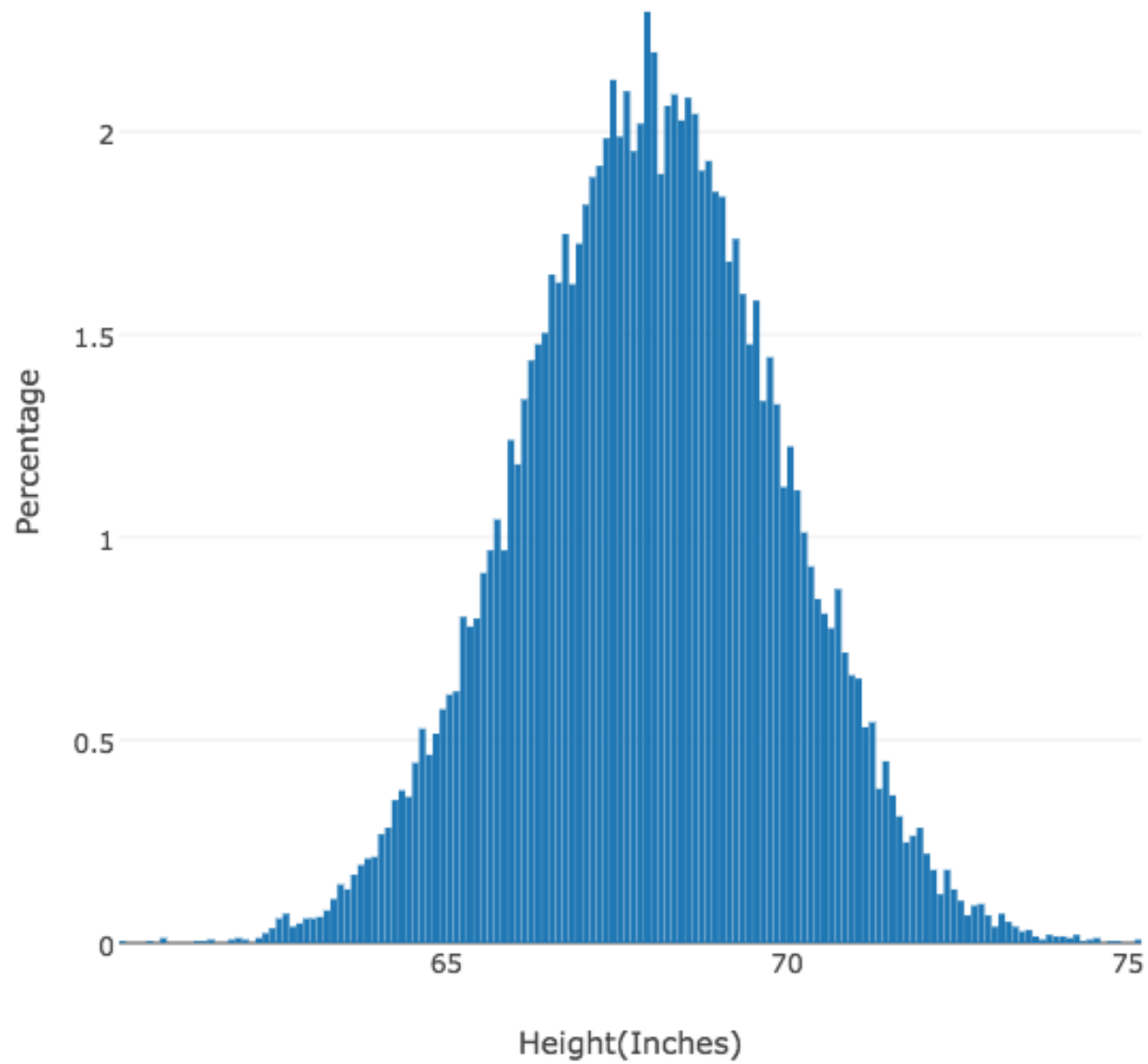
Example #1.5: (DO OVER!) 25,000 Height and Weight Measurements

Our histogram of the height measurements was somewhat disappointing. There was a pattern visible, with the typical height seeming to be 67.5 inches or $5'7\frac{1}{2}"$, and other heights tending to stay close to that value.

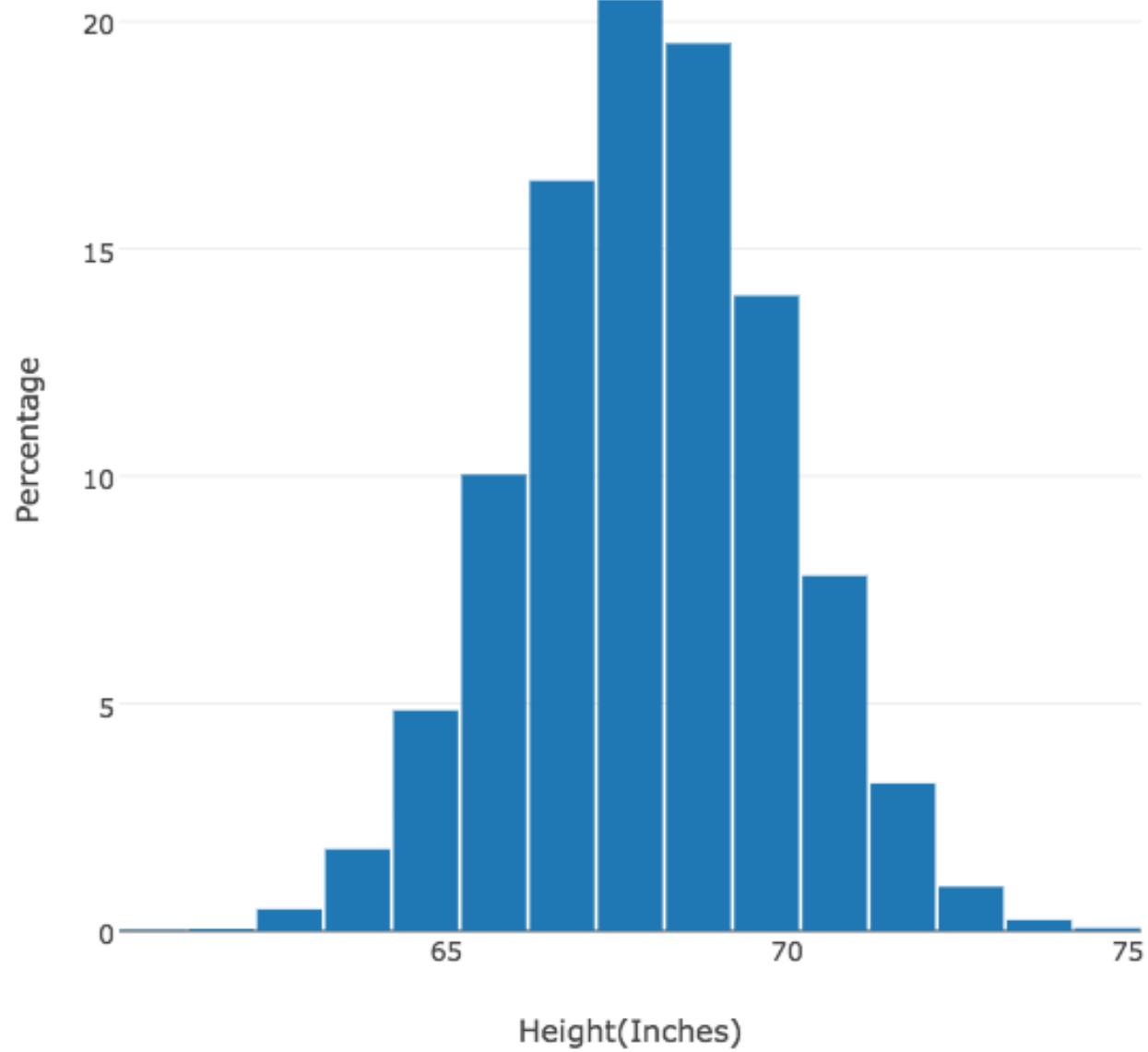
However, the histogram pattern looked very “blocky” and irregular, especially when we had a bin width of $\frac{1}{2}$. We would expect that a law of nature would have a smoother appearance.

However, it's possible that we just haven't looked at enough data for the pattern to become clear. We can test this idea, because we have another dataset of 25,000 height and weight measurements.

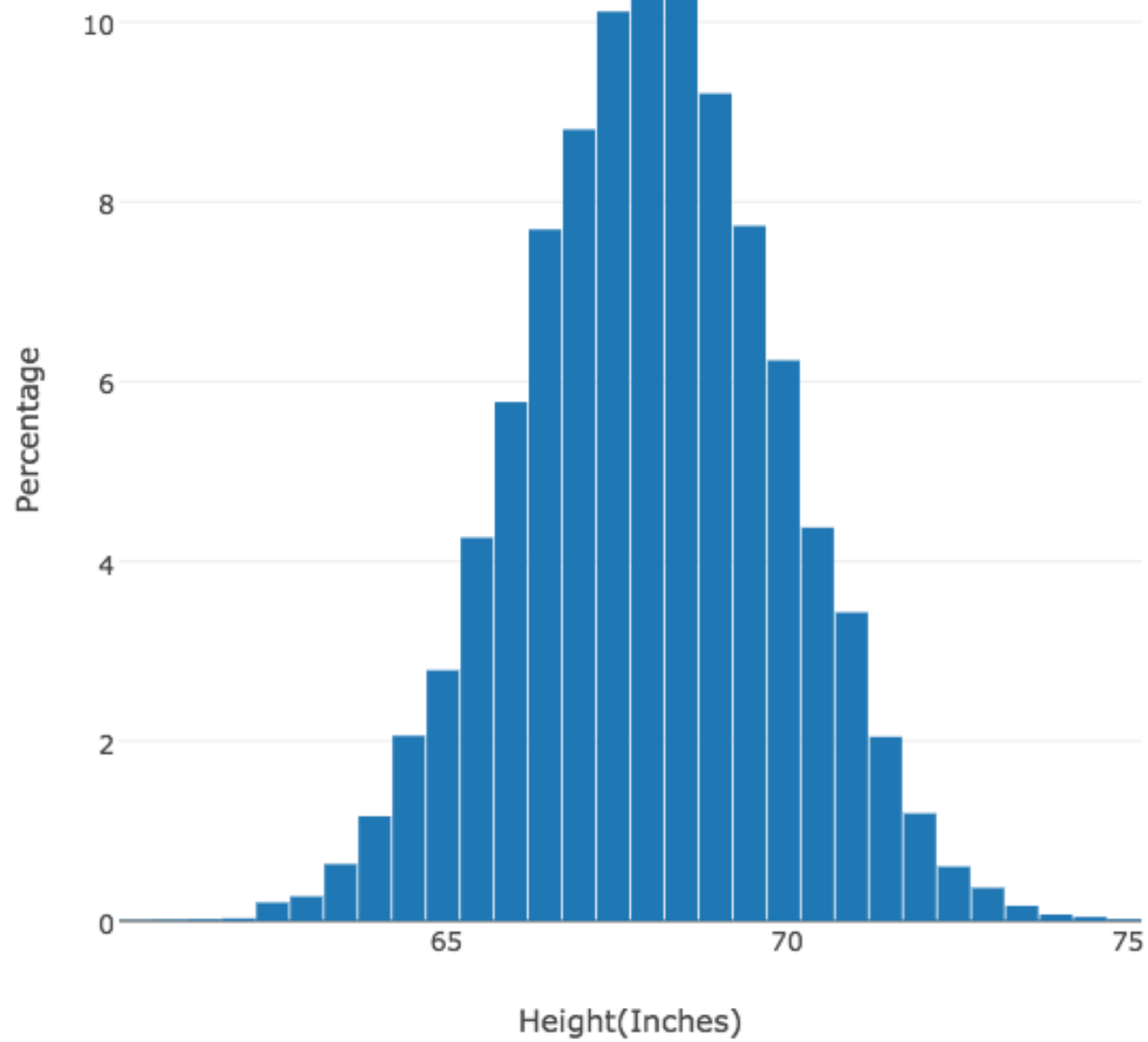
25,000 Heights (bin width = 1/10)



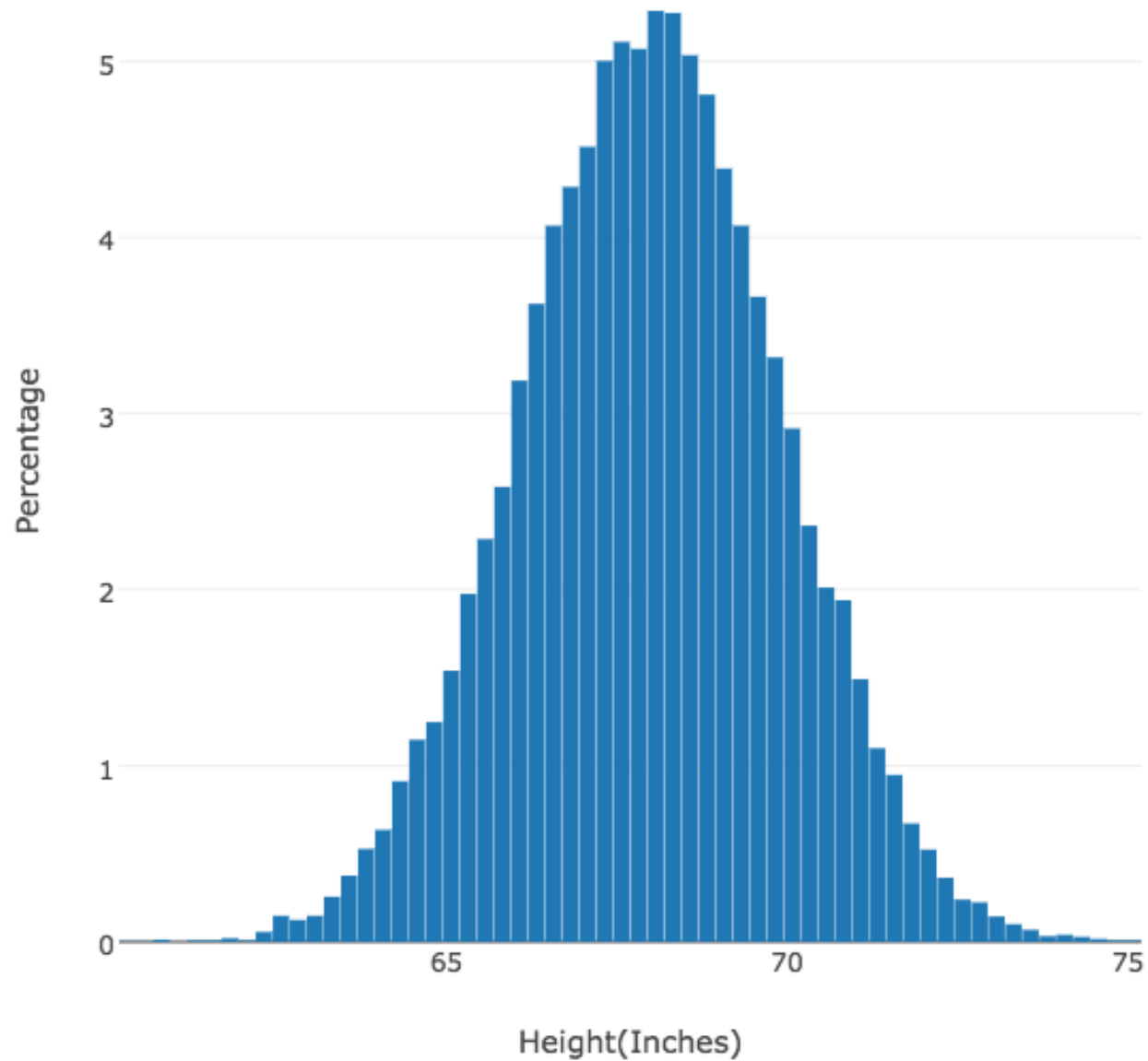
25,000 Heights (bin width = 1)



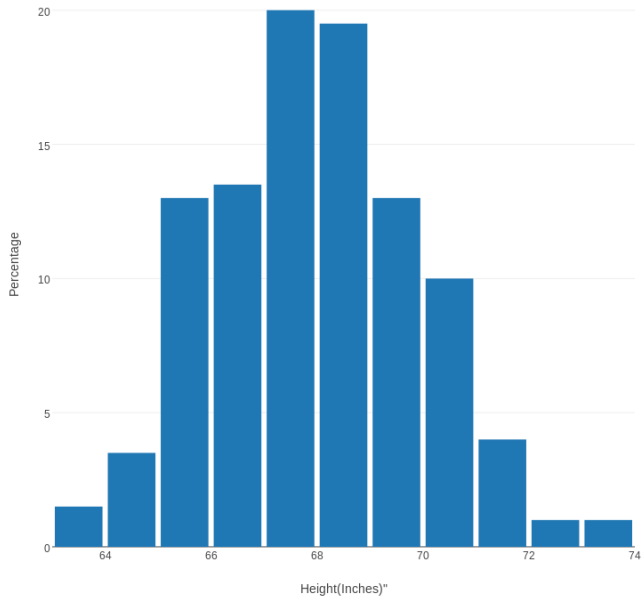
25,000 Heights (bin width = 1/2)



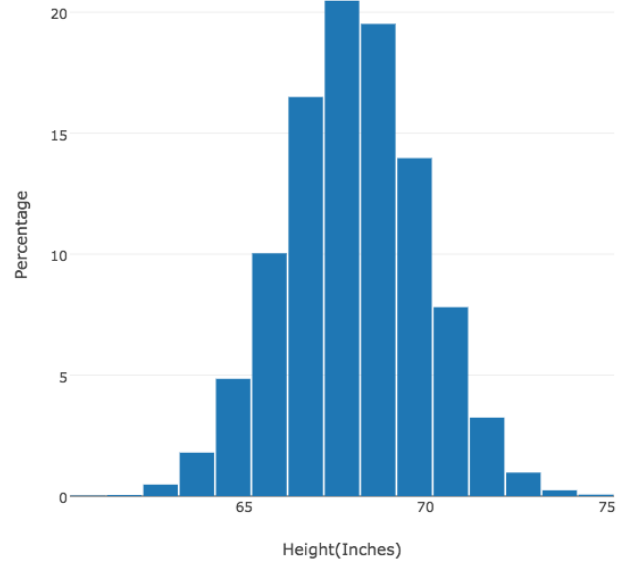
25,000 Heights (bin width = 1/4)



200 Heights (Bin width = 1)



25,000 Heights (bin width = 1)



Example #2: Female Oscar Winners

The Academy awards or “Oscars” for best Actor and Actress have been given every year since 1928, a total of 89 years.

It has often been said that Hollywood treats men and women stars very differently: male actors have long careers, and can have leading roles into their 60’s, while actresses are often allowed only a few hit roles before being dropped in favor of younger stars.

As a small test of this statement, we can plot the ages of female Oscar winners, and see if the plot has any story to tell.

Our data is stored in the file `oscar_female.csv`.

"Index",	"Year",	"Age",	"Name",	"Movie"
1,	1928,	22,	"Janet Gaynor",	"Seventh Heaven"
2,	1929,	37,	"Mary Pickford",	"Coquette"
3,	1930,	28,	"Norma Shearer",	"The Divorcee"
—	—	<i>more data</i>	—	—
87,	2014,	44,	"Cate Blanchett",	"Blue Jasmine"
88,	2015,	54,	"Julianne Moore",	"Still Alice"
89,	2016,	26,	"Brie Larson",	"Room"

We begin our histogram creation by logging into Plotly. There is already some sample data entered in the Plotly grid, but we want to replace this with the data from our file. We do this by going to the **IMPORT** menu item, choosing **Upload file**, then using the browse menu to find our file and clicking on it.

LINE PLOT

Click the column headers to choose x and y columns to graph. Use the different colors to match x columns with y columns.

Line plot

OPTIONS

- Error Bars
- Asymmetric Errors
- Group By
- Text

INSERT INTO

Make a new plot

File Upload

Places: jburkardt, public_html, latex, ct_histogram

Name	Size	Modified
histogram.aux	8 bytes	15:36
histogram.log	10.5 KB	15:36
histogram.pdf	39.5 KB	15:36
histogram.tex	6.8 KB	15:36
histogram.tex.backup	5.6 KB	Yesterday at 13:47
oscar_age_female.csv	4.3 KB	15:33
oscar_age_female.png	17.9 KB	Yesterday at 11:17
oscar_age_female_bin1.png	20.8 KB	15:14
oscar_age_male.csv	3.8 KB	Yesterday at 11:36
oscar_age_male.png	17.6 KB	Yesterday at 11:20
salaries.png	5.6 KB	Yesterday at 13:38
trees.csv	860 bytes	Yesterday at 13:28

Cancel | Open

Col10	Col11	Col12	Col13	Col14	Col15
choose as x	choose as x	choose as x	choose as x	choose as x	choose as x
choose as y	choose as y	choose as y	choose as y	choose as y	choose as y

LINE PLOT

Click the column headers to choose x and y columns to graph. Use the different colors to match x columns with y columns.

Line plot

OPTIONS

Error Bars

	Index	Year	Age	Name	Movie
x	choose as x	choose as x	choose as x	choose as x	choose as x
y	choose as y	choose as y	choose as y	choose as y	choose as y
1	1	1928	22	Janet Gaynor	Seventh Heaven, Street Angel and Sunrise: A Song of Two Humans
2	2	1929	37	Mary Pickford	Coquette
3	3	1930	28	Norma Shearer	The Divorcee
4	4	1931	63	Marie Dressler	Min and Bill
5	5	1932	32	Helen Hayes	The Sin of Madelon Claudet
6	6	1933	26	Katharine Hepburn	Morning Glory

LINE PLOT

Click the column headers to choose x and y columns to graph. Use the different colors to match x columns with y columns.

Line plot

OPTIONS

Error Bars

	Index	Year	Age	Name	Movie	C
x	choose as x	choose as x	choose as x	choose as x	choose as x	choo
y	choose as y	choose as y	choose as y	choose as y	choose as y	choo
1	1	1928	22	Janet Gaynor	Seventh Heaven, Street Angel and Sunrise: A Song of Two Humans	
2	2	1929	37	Mary Pickford	Coquette	
3	3	1930	28	Norma Shearer	The Divorcee	
4	4	1931	63	Marie Dressler	Min and Bill	
5	5	1932	32	Helen Hayes	The Sin of Madelon Claudet	
6	6	1933	26	Katharine Hepburn	Morning Glory	

Plotly | make charts and dashboards online

starter data | oscar_age_fem... | + NEW GRID | IMPORT

ADD DATA | SAVE | COPY | EXPORT | UNDO | REDO | CHOOSE PLOT TYPE | DATA TOOLS | ANALYSIS | Share

LINE PLOT

Click the column headers to choose x and y columns to graph. Use the different colors to match x columns with y columns.

Line plot

OPTI

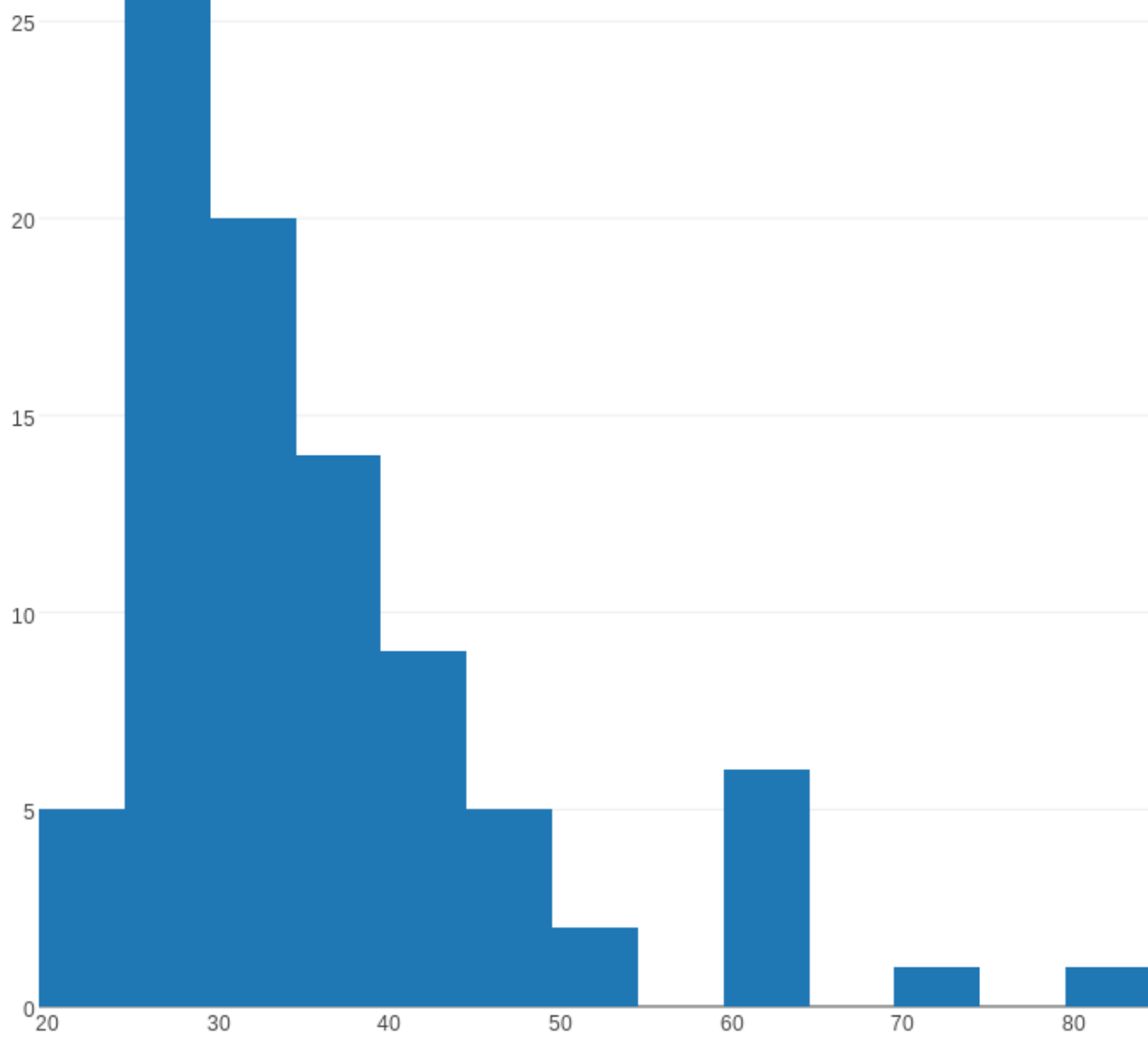
- Scatter plot
- Bar chart
- Histogram
- Area plot

Text

INSERT INTO

Make a new plot

	Index	Year	Age	Name	Movie	Col6	Col7	Col8	Col9	Col10	Col11	Col12	Col13	Col14	Col15
x	choose as x	choose as x	choose as x	choose as x	choose as x	choose as x	choose as x	choose as x	choose as x	choose as x	choose as x	choose as x	choose as x	choose as x	choose as x
y	choose as y	choose as y	choose as y	choose as y	choose as y	choose as y	choose as y	choose as y	choose as y	choose as y	choose as y	choose as y	choose as y	choose as y	choose as y
1					Seventh Heaven, Street Angel and Sunrise: A Song of Two Humans										
2	1	1928	22	Janet Gaynor											
3	2	1929	37	Mary Pickford	Coquette										
4	3	1930	28	Norma Shearer	The Divorcee										
5	4	1931	63	Marie Dressler	Min and Bill										
6	5	1932	32	Helen Hayes	The Sin of Madelon Claudet										
7	6	1933	26	Katharine Hepburn	Morning Glory										
8	7	1934	31	Claudette Colbert	It Happened One Night										
9	8	1935	27	Bette Davis	Dangerous										
10	9	1936	27	Luise Rainer	The Great Ziegfeld										
11	10	1937	28	Luise Rainer	The Good Earth										
12	11	1938	30	Bette Davis	Jezebel										
13	12	1939	26	Vivien Leigh	Gone with the Wind										
14	13	1940	29	Ginger Rogers	Kitty Foyle										
15	14	1941	24	Joan Fontaine	Suspicion										
16	15	1942	38	Greer Garson	Mrs. Miniver										
17	16	1943	25	Jennifer Jones	The Song of Bernadette										
18	17	1944	29	Ingrid Bergman	Gaslight										
19	18	1945	40	Joan Crawford	Mildred Pierce										
20	19	1946	30	Olivia de Havilland	To Each His Own										



- View data
- View JSON
- Share
- Traces
- Layout
- Axes
- Notes
- Legend
- Fit data
- Themes

Traces

Age

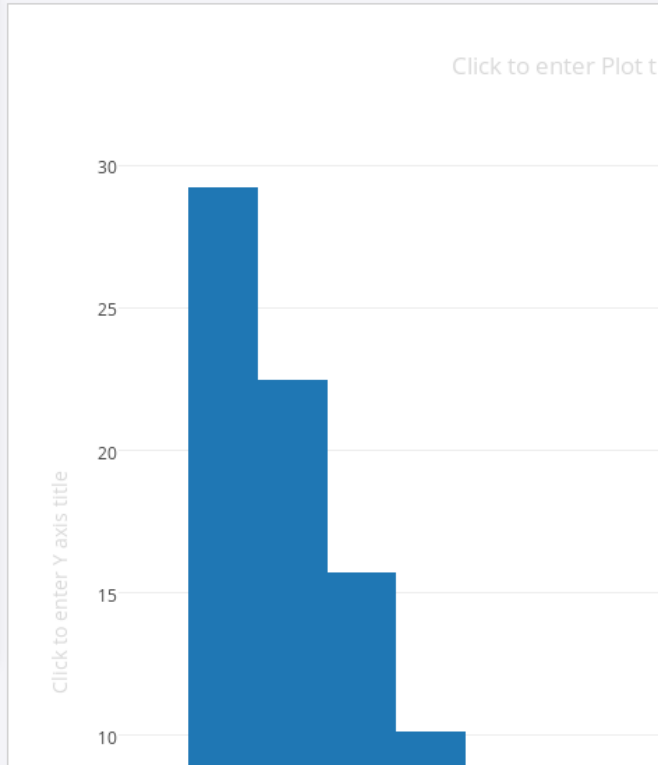
Mode Range/bins Style Error bars

Normalization None % p $1/\Delta x$ $p/\Delta x$

X bins Auto Manual

X range 19.5 84.5

X bin size 5



Traces

Age

Mode Range/bins Style Error bars

Opacity 1

Marker Fill

Outline 0

Bar Gap 0.1

Group Gap 0

View data

View JSON

Share

Traces

Layout

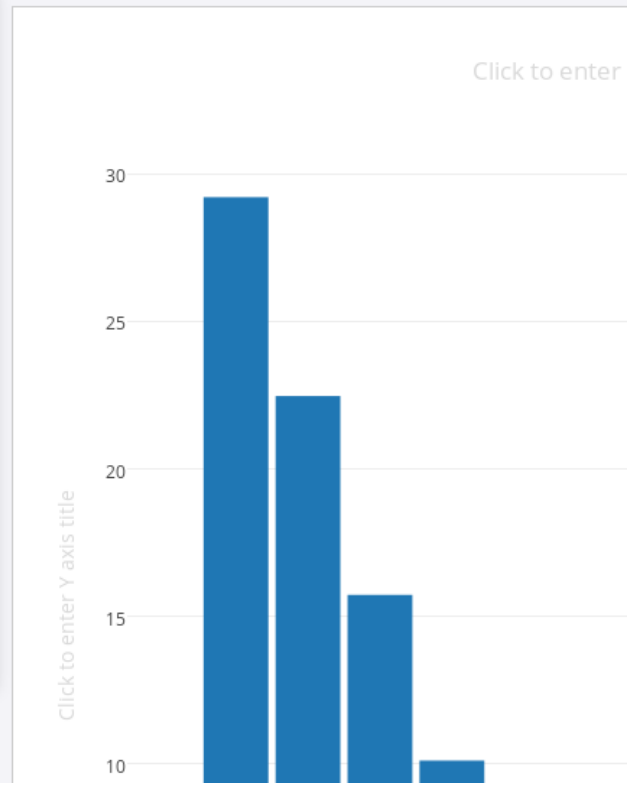
Axes

Notes

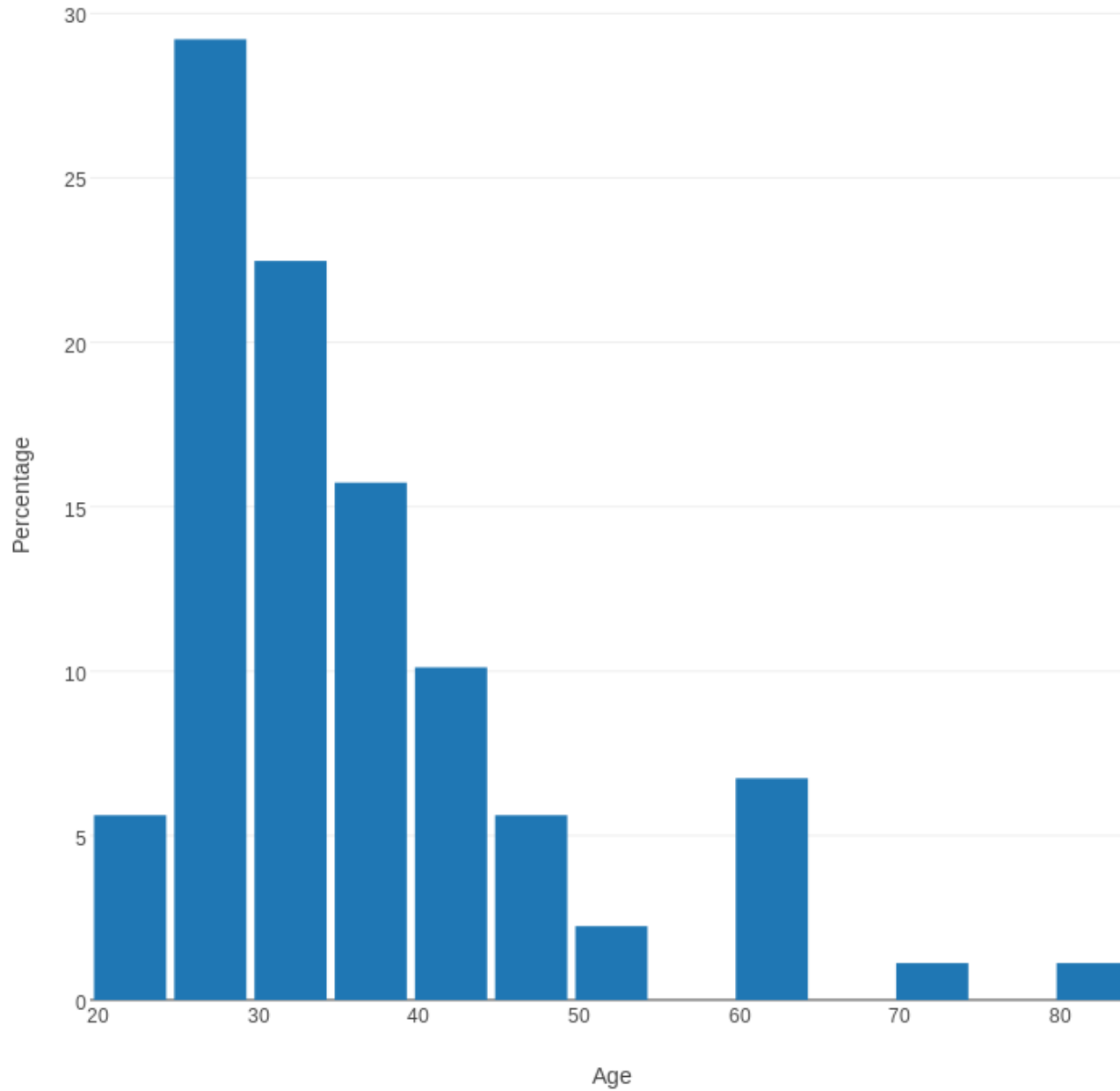
Legend

Fit data

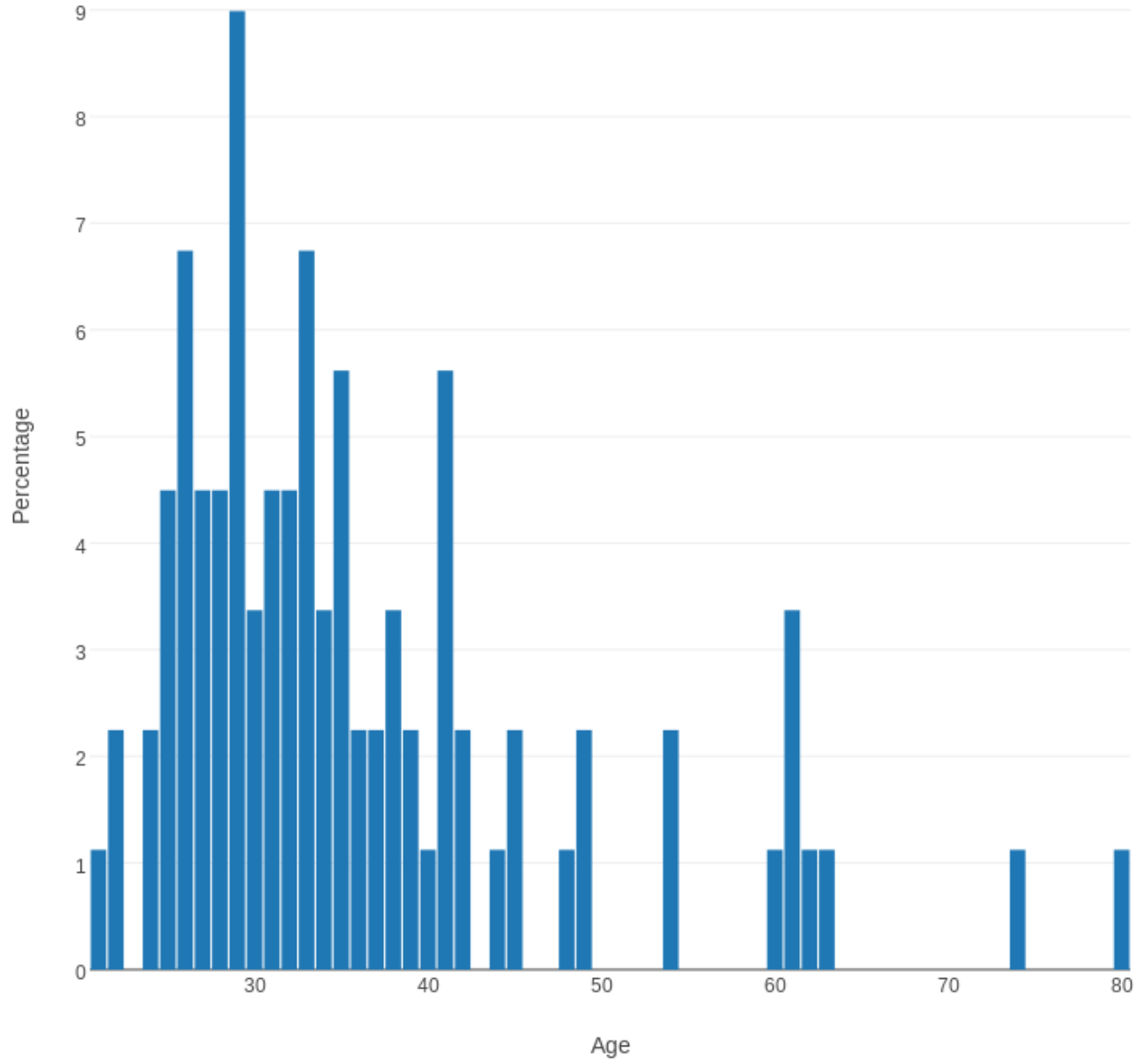
Themes



Female Oscar Winners by Age



Female Oscar Winners by Age



Example #3: Male Oscar Winners

The female Academy Award winners represent a small and untypical sample of the acting careers of women in Hollywood, but our histogram does suggest that women have the greatest success at a very early age, and seem to decrease in visibility each time we move 5 years further in age range, although there are some “blips” popping up in the later end of the plot.

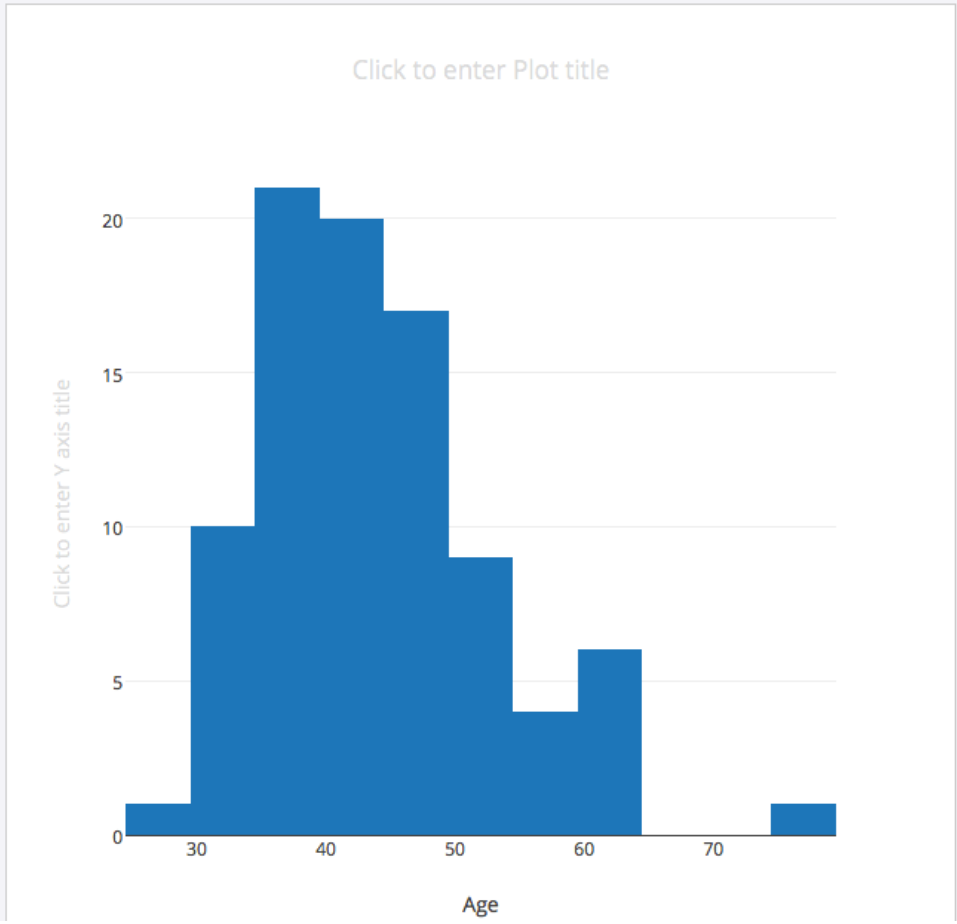
If our data is correctly showing us the “story”, then it’s natural to wonder about what we would see if we examined the same data for male Oscar winners.

First we need to grab the corresponding data file!

Our data is stored in the file `oscar_male.csv`.

"Index",	"Year",	"Age",	"Name",	"Movie"
1,	1928,	44,	"Emil Jannings",	"The Last Command"
2,	1929,	41,	"Warner Baxter",	"In Old Arizona"
3,	1930,	62,	"George Arliss",	"Disraeli"
—	—	<i>more data</i>	—	—
87,	2014,	44,	"Matthew McConaughey",	"Dallas Buyers Club"
88,	2015,	33,	"Eddie Redmayne",	"The Theory of Everything"
89,	2016,	41,	"Leonardo DiCaprio",	"The Revenant"

- 📄
- 📄
- ⬇️
- ↶
- Ⓜ️
- 📄 View data
- 📄 View JSON
- 🔗 Share
- 📄 Traces
- 📄 Layout
- 📄 Axes
- 📄 Notes
- ☰ Legend
- 📄 Fit data
- 📄 Themes



Our first plot has an interesting shape, but if we are going to compare it to the same plot for the female Oscar winners, we need to make some modifications. In particular, note that the horizontal data (the "Age" values) doesn't seem to have the same extent as in the plot of the female data.

Unless we say so, Plotly will make the graph just wide enough for the data it sees. But we can go into the **Axes** menu item, and pick the **Range** item and then change the range from the values

```
Range 24.5 79.5
```

to match the range on the female plot:

```
Range 19.5 84.5
```

starter data oscar_age_male... Plot Plot + NEW GRID IMPORT

Axes

X Axis + -

View data Range Lines Ticks Labels Layout

View JSON

Share

Traces

Layout

Axes

Notes

Legend

Fit data

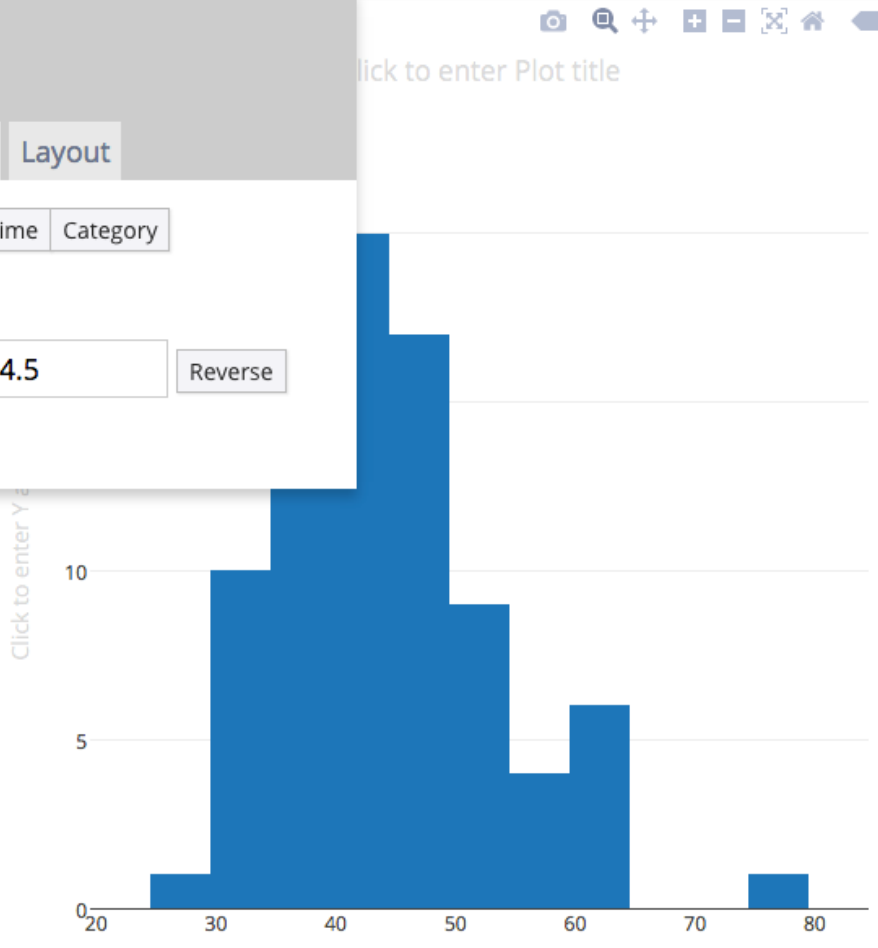
Themes

Type Linear Log DateTime Category

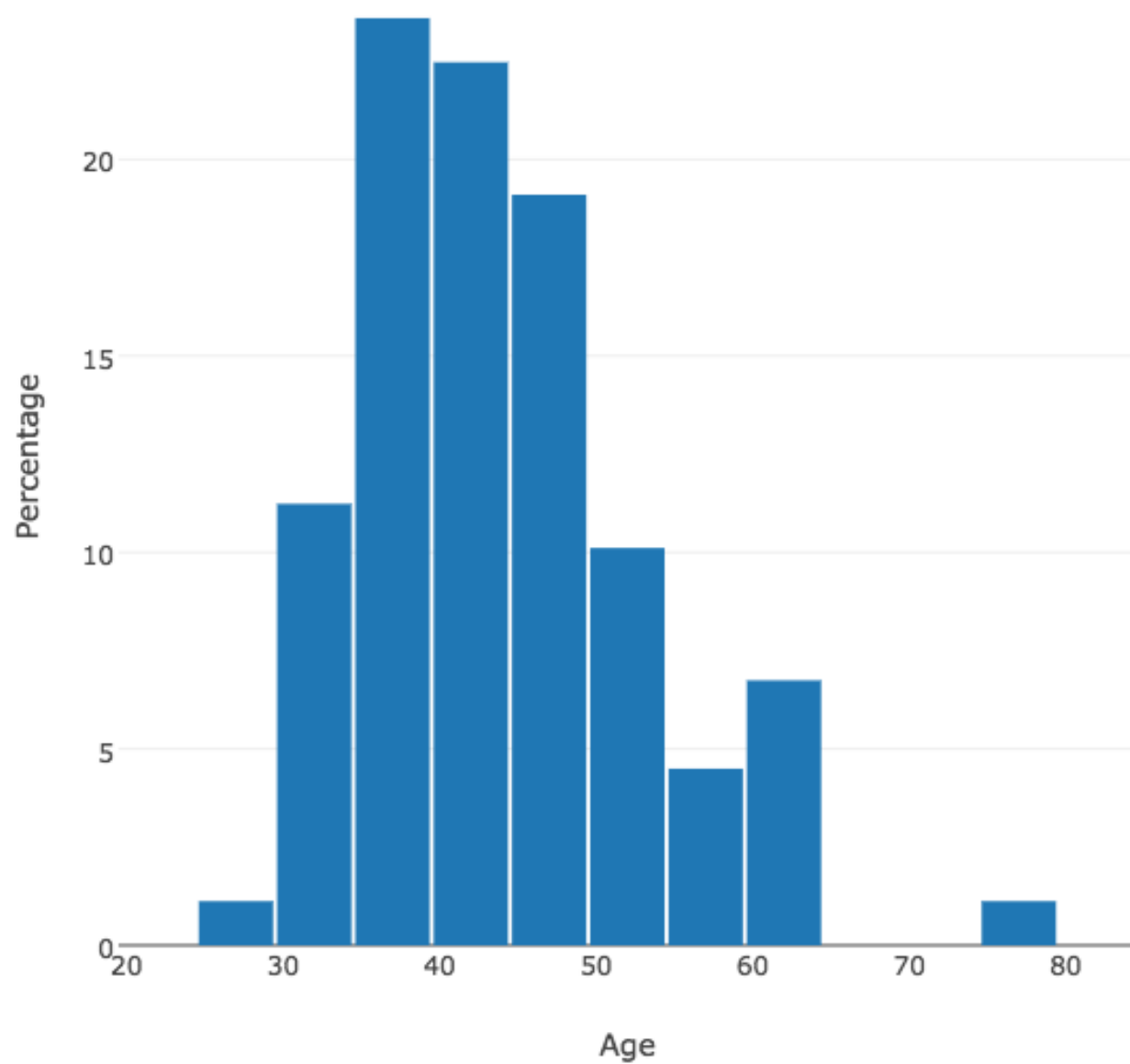
Autorange On Off

Range 19.5 84.5 Reverse

Zoom / pan Interactive Fixed



Male Oscar Winners by Age

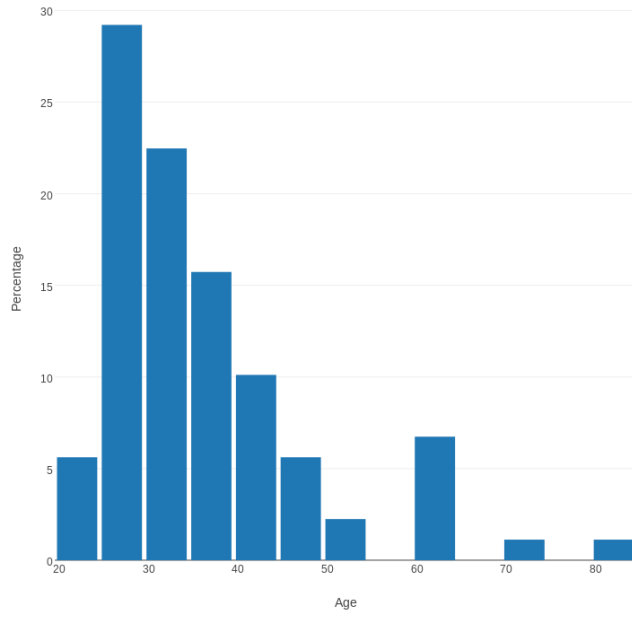


Now we have both sets of data plotted with the same age range matched up.

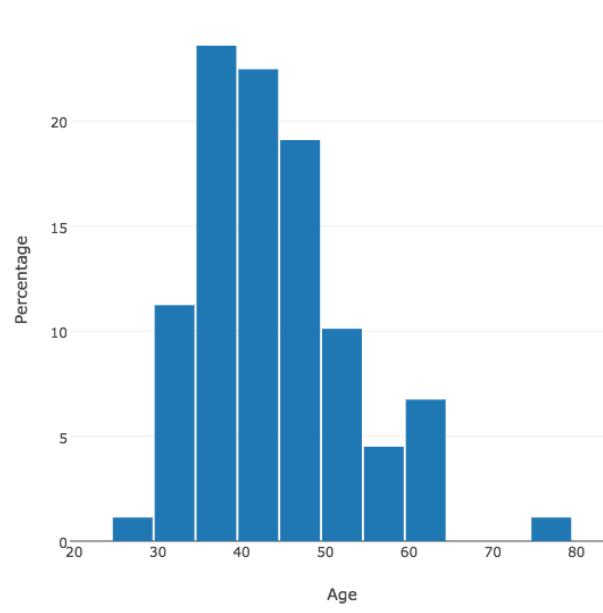
We can concentrate on the “story” that the two sets of histogram bars seem to be telling. One reading of this story is:

- Male actors have their peak popularity in their late 30’s, and their careers don’t start to drop off until around 50, continuing til 65.
- Female actores have their peak popularity in their late 20’s, and their careers immediately begin to drop off substantially every 5 years thereafter til about 55.

Female Oscar Winners by Age



Male Oscar Winners by Age



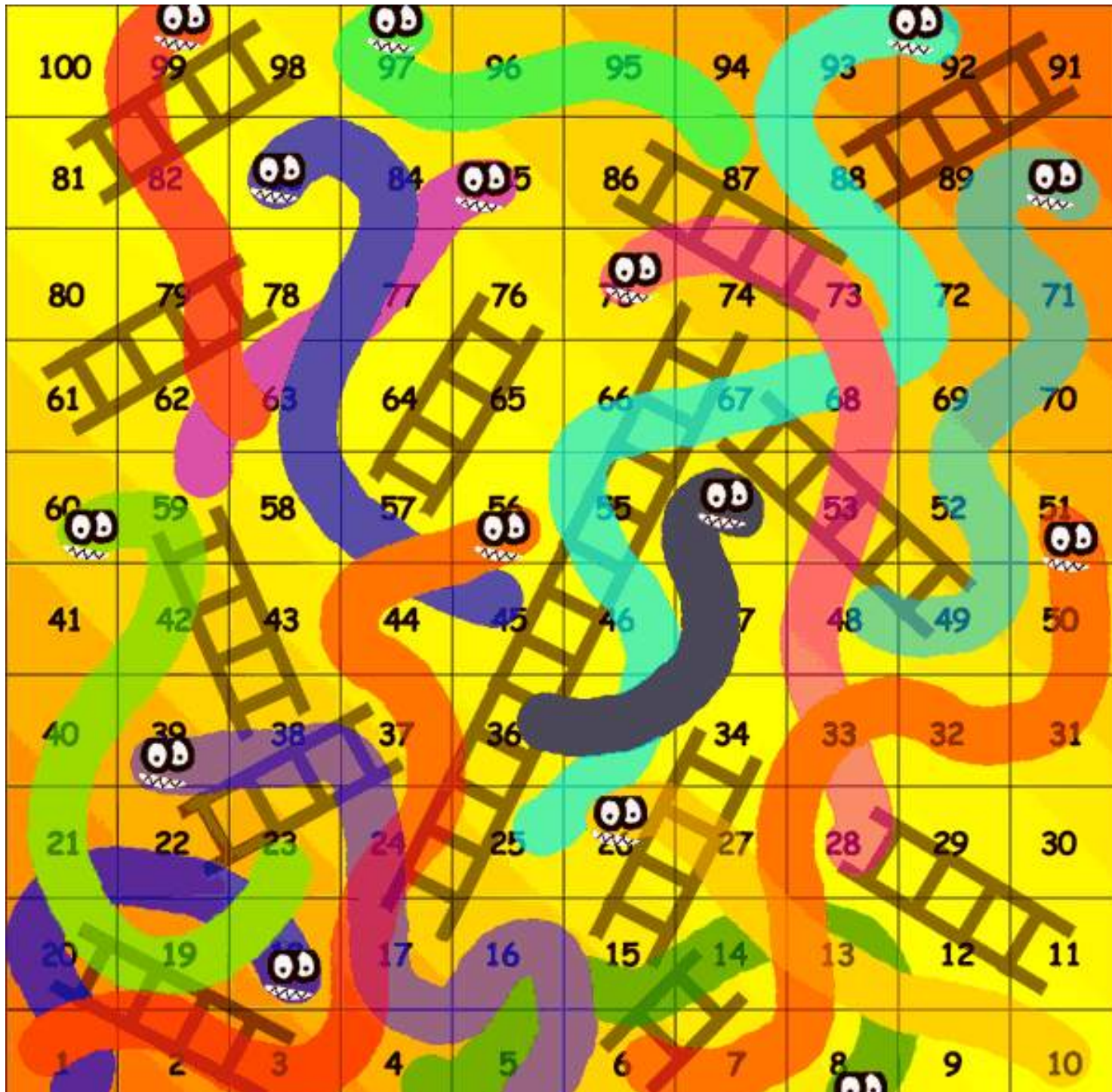
Example #4: Probabilities (Number of Turns in Snakes and Ladders)

Snakes and Ladders is a game played on a board of 100 numbered squares.

From the starting square, players repeatedly roll a die, trying to reach the final square first.

Certain squares have a ladder that moves you ahead; others are snakes that drag you backwards.

A simple question to investigate: *How many turns does it take for one player to reach the final square?*



Since dice are involved, there's an element of chance, so we have to revise our question and ask, what is the **typical** or **average** number of turns required for one player?

It's not easy to see a way to answer this question, except to play a lot of games and see what happens.

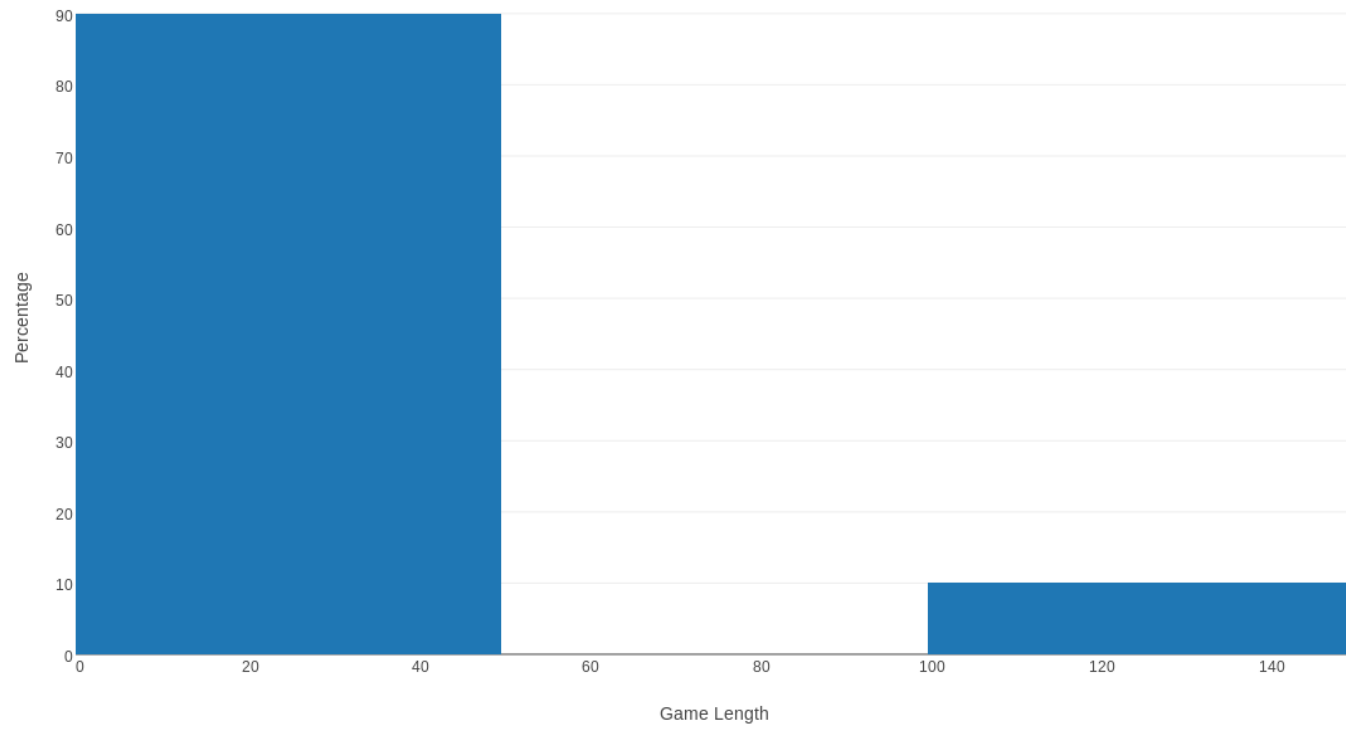
To keep things simple, we imagine a single player, who rolls the die over and over, following the rules, until reaching the end, and writing down the number of turns required.

We could store our data in a file called `snakes_10.csv`.

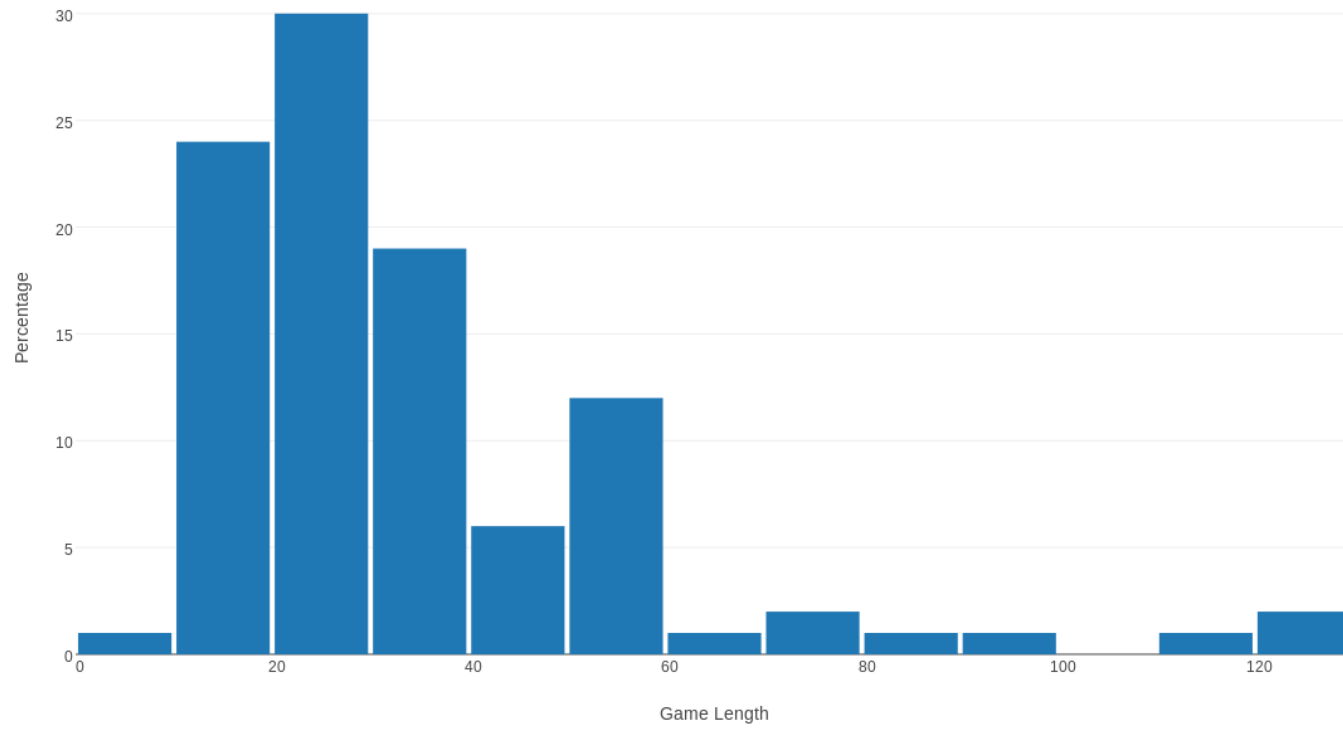
"Game Number"	"Game Length"
1,	30
2,	29
3,	31
4,	16
5,	24
6,	29
7,	28
8,	117
9,	42
10,	23

but that one long 117 turn game suggests we need more data!

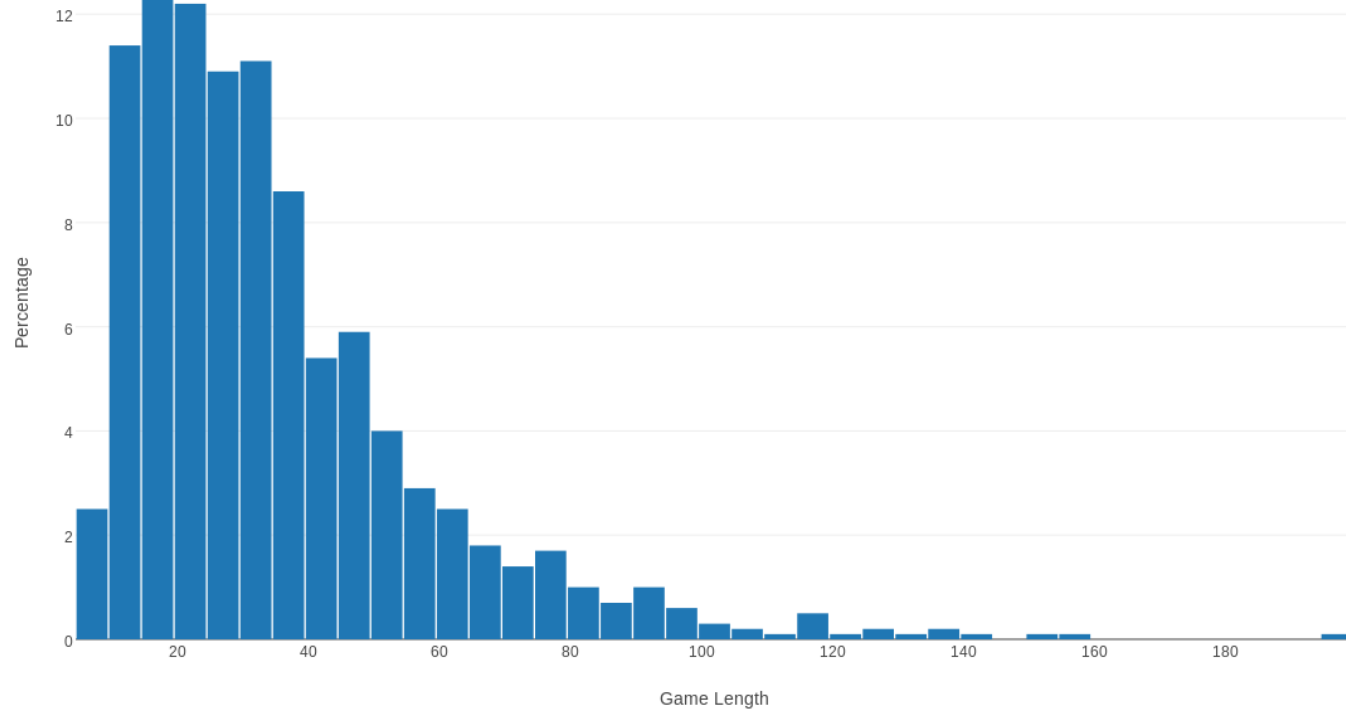
10 Games of Snakes and Ladders



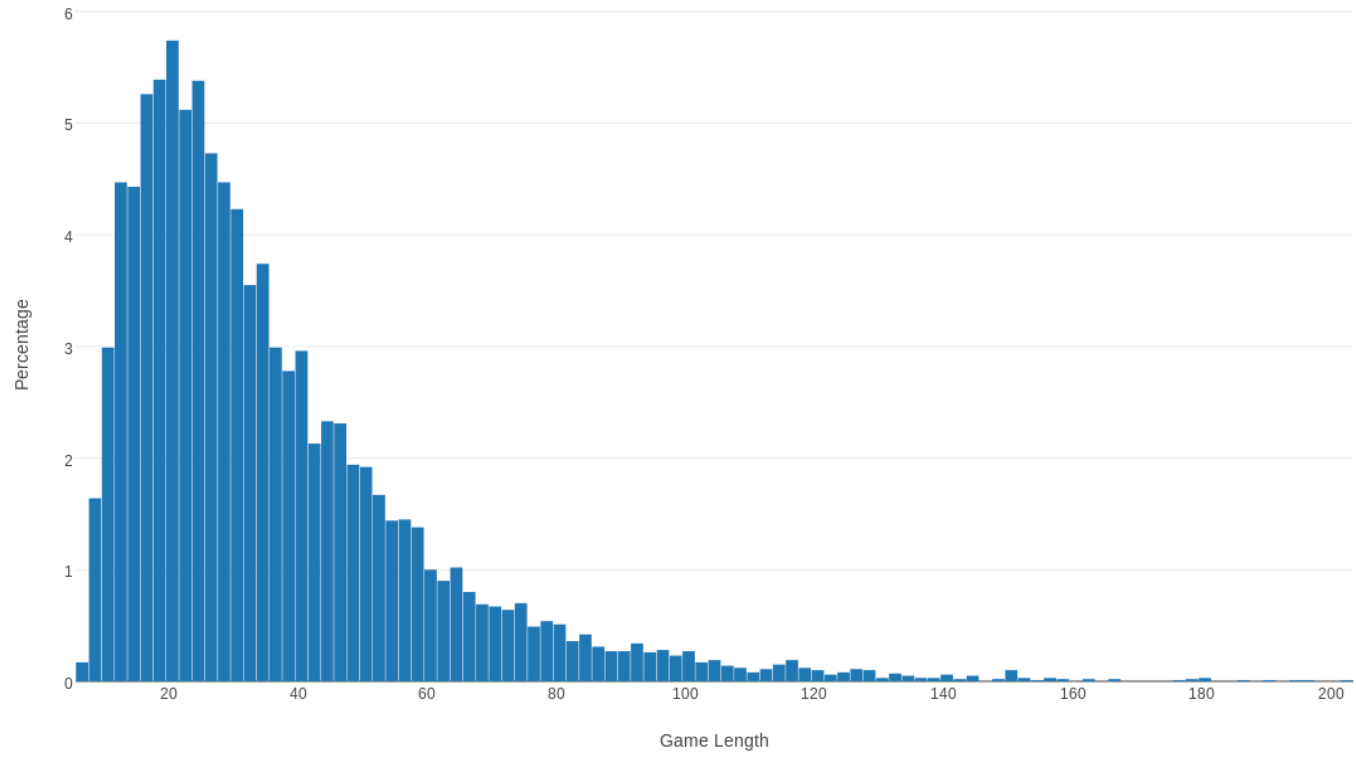
100 Games of Snakes and Ladders



1,000 Games of Snakes and Ladders



10,000 Games of Snakes and Ladders

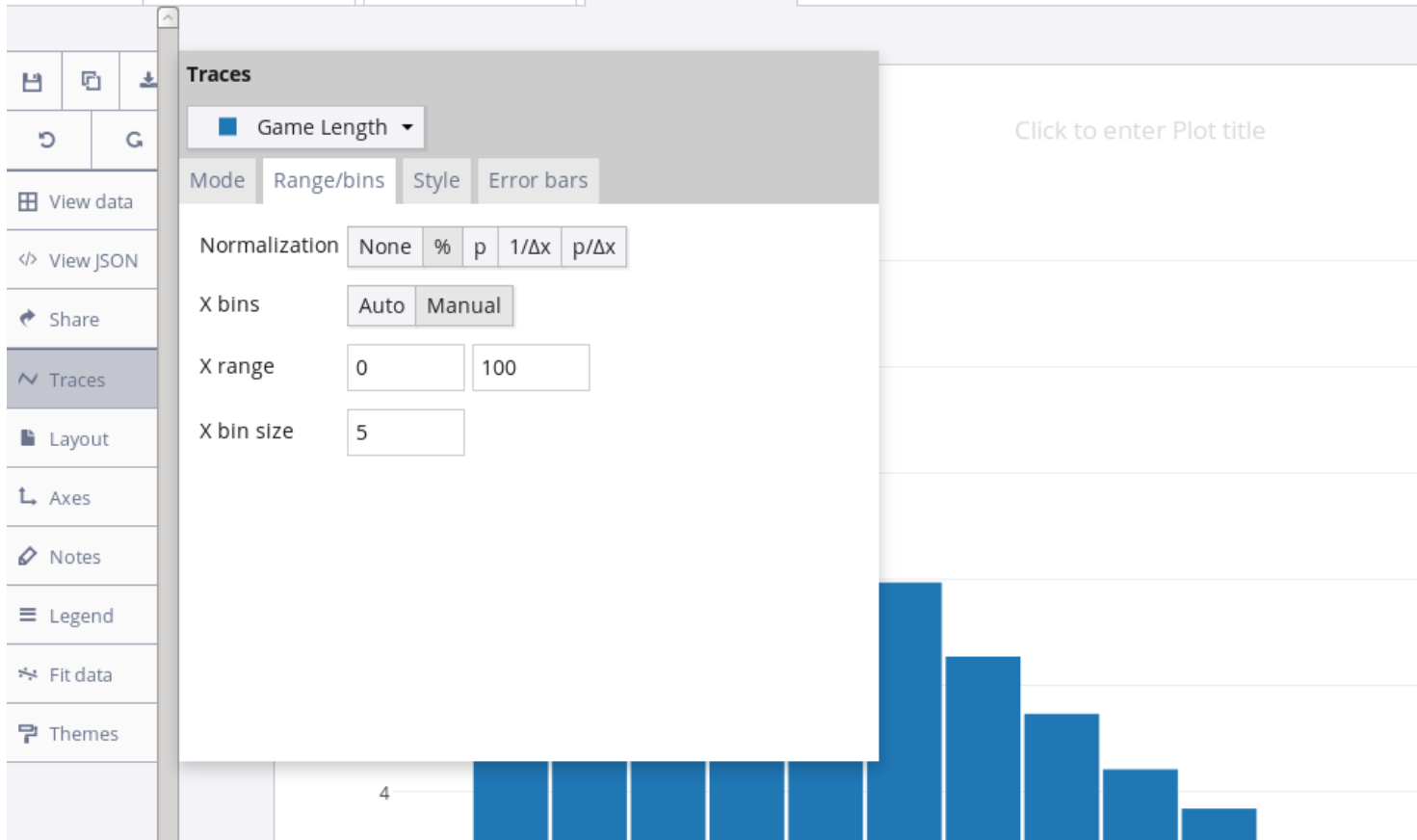


Using data from 10,000 games, our histogram seems to have a fairly smooth shape, which suggests that the data has “settled down”, and that the effects of chance have been averaged out.

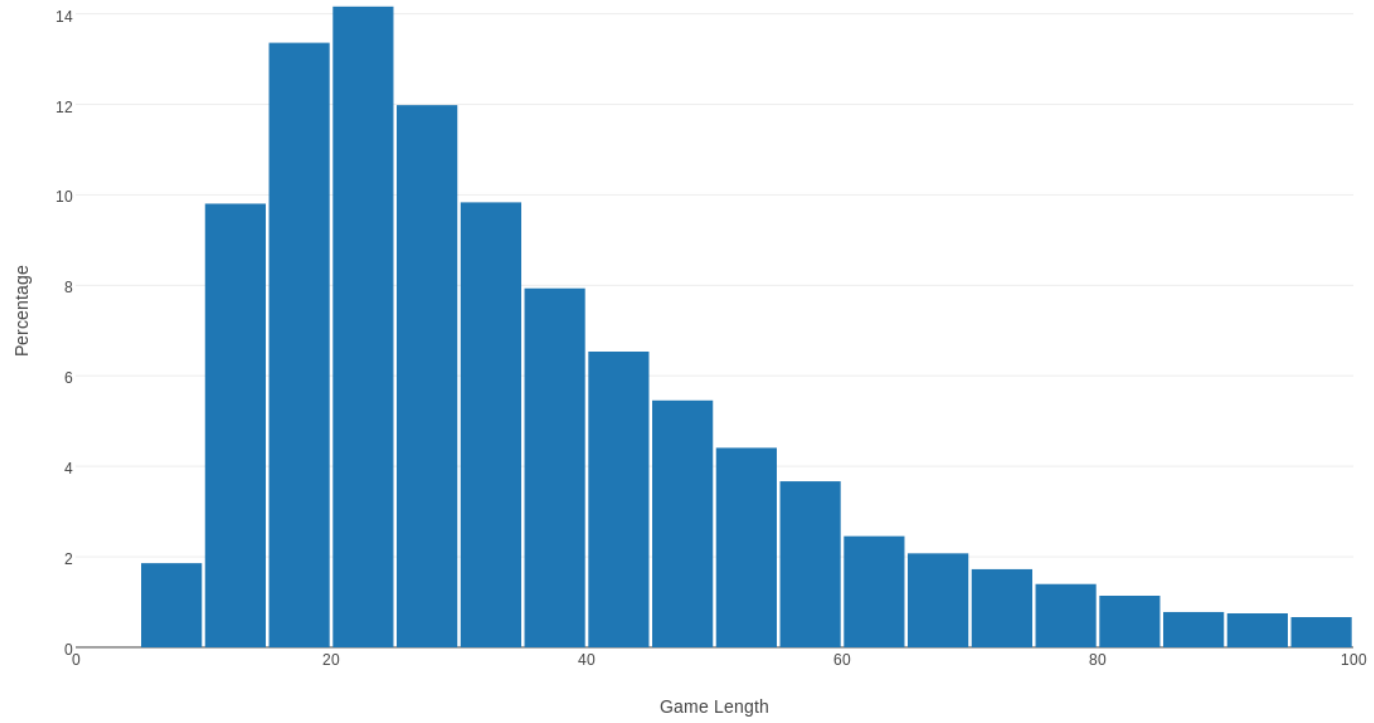
However, we still can see a certain amount of saw-tooth or up-and-down variations in the plot. If we look in the **Traces** menu, we can see that our bin width is 2. If we use wider bins, we may be able to smooth out our plot. At the same time, we can reduce the maximum interesting value of Game Length to 100, which will allow us to focus on the most commonly occurring values.

The result is a simplified plot that is smoother and easier to read.

starter data snakes_10000.c... Plot + NEW GRID ↑ IMPORT



10,000 Games of Snakes and Ladders



Example #5: A Camel with Two Humps (Bimodal data)

Sometimes a histogram can show an interesting pattern in the data. It may look like a double mountain range, having two peaks.

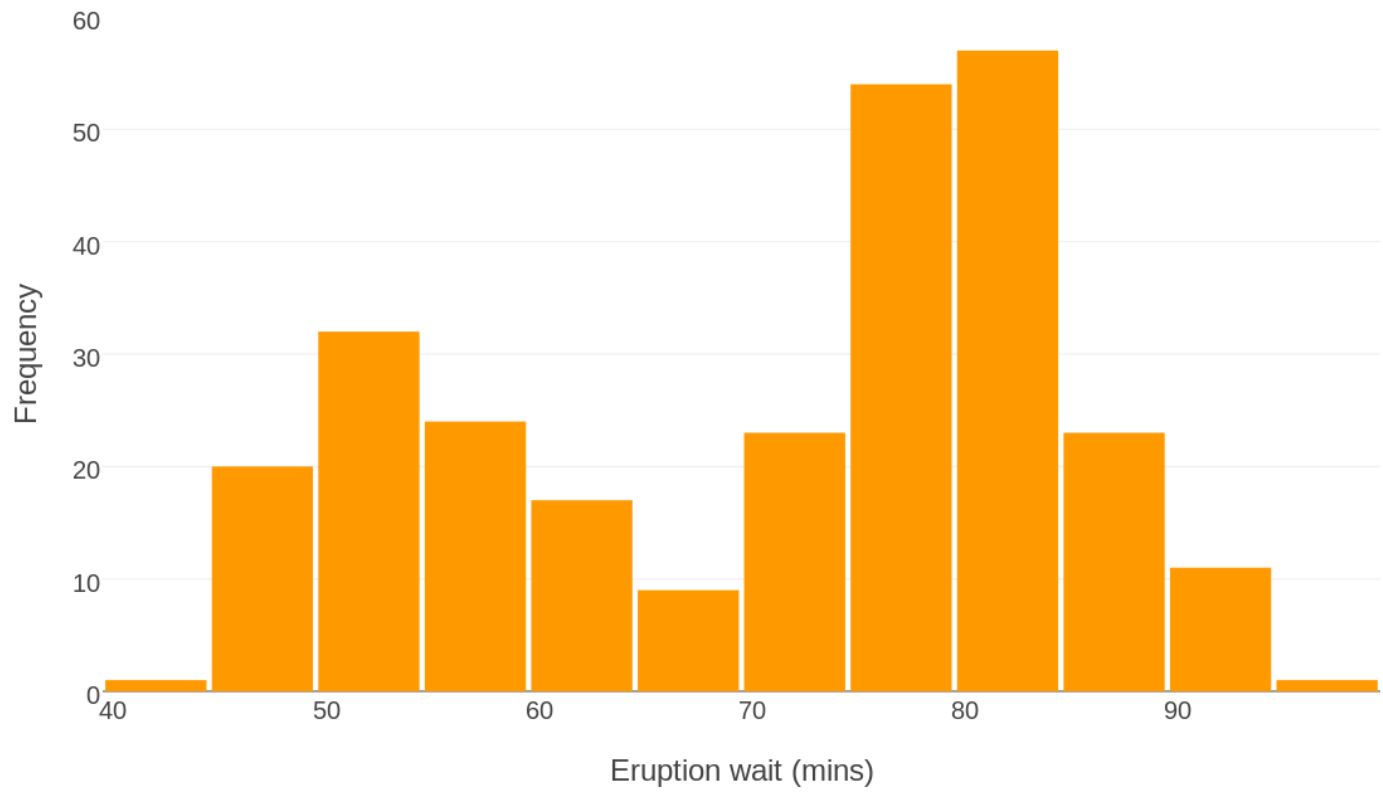
Often, a peak represents a typical behavior, and if the data has two peaks, then it suggests that sometimes the system wants to do one thing, and sometimes the other.

A peak in a histogram is sometimes called a **mode** and a histogram that has two peaks is called a **bimodal** histogram.

Our data is stored in the file `faithful.csv`.

"Index"	"Eruption length (mins)"	"Eruption wait (mins)"
1,	3.600,	79
2,	1.800,	54
3,	3.333,	74
—	<i>more data</i>	—
270,	4.417,	90
271,	1.817,	46
272,	4.467,	74

Interval Between Eruptions of Old Faithful Geyser



Now we will do a histogram of the duration of an eruption, that is, after however long we had to wait, how long the geyser actually sprays water upwards in a spectacular show.

Since the geyser has a short and a long waiting period, we might expect that sometimes the water doesn't have so much time to heat up, and other times it does, so that we have a corresponding short and long eruption pattern as well.

By default, the horizontal plot range starts with the smallest observed time. In order to better suggest the time scale, I went to the `coloredAxes` left hand menu item and reset the plot range to start at 0.

- View data
- View JSON
- Share
- Traces
- Layout
- Axes**
- Notes
- Legend
- Fit data
- Themes

Axes

X Axis

Range Lines Ticks Labels Layout

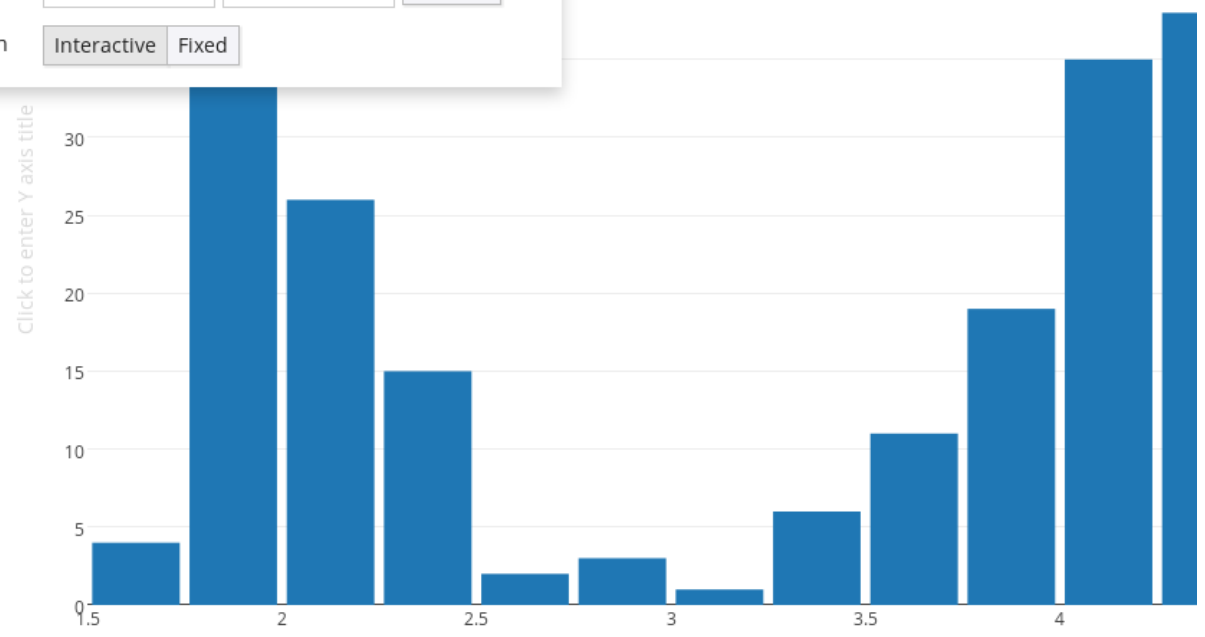
Type

Autorange

Range

Zoom / pan

Click to enter Plot title



Length of Old Faithful Geyser Eruptions

