





Workbook v1.2

Brought to you by the Bootstrap team:

- Emmanuel Schanzer
- Kathi Fisler
- Shriram Krishnamurthi
- Ed Campos
- Emma Youndtsmith
- Sam Dooman

Bootstrap is licensed under a Creative Commons 3.0 Unported License. Based on a work from www.BootstrapWorld.org. Permissions beyond the scope of this license may be available at schanzer@BootstrapWorld.org.

Unit 1

Many important questions ("what's the best restaurant in town?", "is this law good for citizens?", etc.) are answered with data. Data Scientists try and answer these questions, by writing *programs that ask questions of data*.

Data of all types can be organized into **Tables**

- Every Table has a **header row**, and some number of **data rows**
- **Quantitative data** is data - usually numeric - that measures *quantity*, such as a person's height, a score on test, a measure of distance, etc. A list of quantitative data can be ordered from smallest to largest.
- **Categorical data** is data that specifies *categories*, such as eye color, country of origin, etc. A list of categorical data has no notion of "smallest" or "largest", and cannot be ordered.

Programming languages involves different *datatypes*, such as Numbers, Strings, Booleans and Images.

- **Operators** (like +, -, *, <, etc.) are written between values. For example: `4 + 2`
- We can use **functions** (like triangle, star, string-repeat, etc.) by writing the function name first, followed by a list of **arguments** in parentheses. For example: `star(50, "solid", "red")`
- **Methods** are special functions that are attached to pieces of data. We use them to manipulate Tables. They are different from functions in several ways:
 - Their names can't be used alone: they can only be used as part of data, separated by a dot. (For example, `shapes.row-n(2)`)
 - Their contracts are different: they include the type of the data as part of their names. (eg, `<table>.row-n :: (index :: Number) → Row`)
 - They have a "secret" argument, which is the data they are attached to
- In this course, we will use three **Table Methods** to manipulate our datasets:
 - `<Table>.order-by` – order the rows of a table based on a column
 - `<Table>.filter` – create a **subset** of the data, with only certain rows
 - `<Table>.build-column` – use the columns of a table to make a new one

Numbers and Strings

Make sure you've loaded the Unit 1 Starter File, and clicked "Run".

1. Try typing `42` into the Interactions Area and hitting "Enter". What happens?
2. Try typing in other Numbers. What happens if you try a decimal like `0.5`? A fraction like `1/3`? Try really big Numbers, and really small ones.
3. String values are always in quotes. Try typing your name (in quotes!). What happens when you hit "Enter"?
4. Try typing your name with the opening quote, but *without* the closing quote. What happens? Now try typing it without *any* quotes.
5. Is `42` the same as `"42"`? Why or why not? Write your answer below:

They are different data types: `42` (without quotes) is a Number, and `"42"` (with quotes) is a string.

Operators

6. Just like in math, Pyret has operators like `+` and `*`. Try typing in `4 + 2`, and then `4+2` (without the spaces). What can you conclude from this? Write your answer below:

Operators (like `+`) need whitespace separating them from their operands.

7. Try typing in `4 + 2 + 6`, `4 + 2 * 6`, and `4 + (2 * 6)`. What can you conclude from this? Write your answer below:

You can use the same operator multiple times without parentheses, but you need parentheses to group order of operations if using different operators (like `+` and `*`) together.

8. Try typing in `4 + "cat"`, and then `"dog" + "cat"`. What can you conclude from this? Write your answer below:

The `+` operator can only be used with Numbers, not Strings.

Booleans

Boolean expressions are yes-or-no questions, and will always evaluate to either `true` ("yes") or `false` ("no"). What will each of the expressions below evaluate to? Write down the result in the blanks provided, and type them into Pyret if you're not sure.

<code>3 <= 4</code>	<u>True</u>	<code>"a" > "b"</code>	<u>False</u>
<code>3 == 2</code>	<u>False</u>	<code>"a" <> "b"</code>	<u>True</u>
<code>2 <> 4</code>	<u>True</u>	<code>"a" == "b"</code>	<u>False</u>
<code>3 <> 3</code>	<u>True</u>	<code>"a" <> "a"</code>	<u>False</u>

Boolean Operators

Pyret also has operators that work on *Booleans*. For each expression below, write down your guess about what it will evaluate to. Then type them in and see if you were right!

<code>(3 <= 4) and (3 == 2)</code>	<u>False</u>
<code>("a" == "b") and (3 <> 4)</code>	<u>False</u>
<code>(3 <= 4) or (3 == 2)</code>	<u>True</u>
<code>("a" == "b") or (3 <> 4)</code>	<u>True</u>

-
1. How many different Number values are there in Pyret? Infinite
 2. How many different String values are there in Pyret? Infinite
 3. How many different Boolean values are there in Pyret? Two

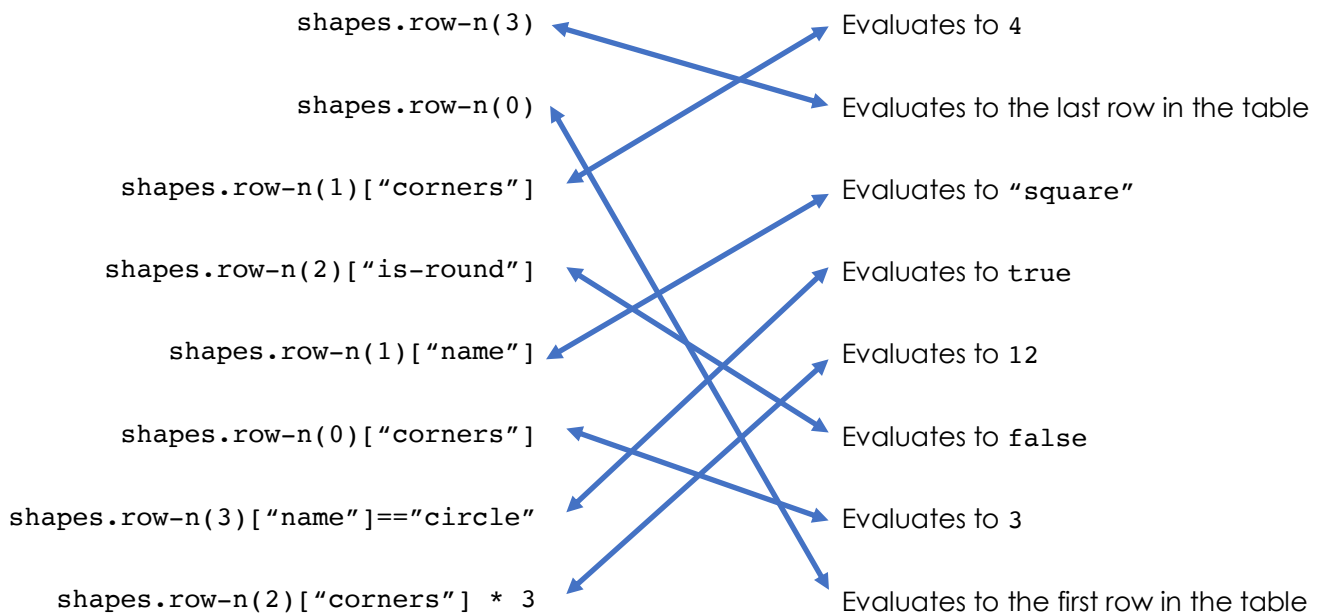
Lookups

The table below represents four shapes in a table:

shapes

name	corners	is-round
"triangle"	3	false
"square"	4	false
"rectangle"	4	false
"circle"	0	true

1. **Match** each Pyret expression (left) to the description of what it looks up (right).



2. Fill in the blanks (left) with the Pyret lookup code that will produce the value (right).

- | | |
|---|-------------|
| a. <u><code>shapes.row-n(2) ["name"]</code></u> | "rectangle" |
| b. <u><code>shapes.row-n(0) ["name"]</code></u> | "triangle" |
| c. <u><code>shapes.row-n(1) ["corners"]</code></u> | 4 |
| d. <u><code>shapes.row-n(3) ["corners"]</code></u> | 0 |
| e. <u><code>shapes.row-n(3) ["is-round"]</code></u> | true |

Unit 2

Answering Questions from Data can take many forms. Here are a few types of questions, each requiring a different kind of analysis:

- **Lookup Questions** can be answered just by finding the right row and column a table. (e.g. – “How old is Toggle?”)
- **Compute Questions** can be answered by computing over a single row or column. (e.g. – “What is the heaviest animal at the shelter?”)
- **Analyze Questions** require looking for trends across multiple rows or columns. (e.g. – “Do cats tend to be adopted sooner than dogs?”)

We can **define our own functions**, using a technique called the **Design Recipe**.

- We use the Design Recipe to help us define functions **without making mistakes**.
- The first step is to write a **Contract** and **Purpose Statement** for the function, which specify the Name, Domain and Range of the function and give a summary of what it does.
- The second step is to **write at least two examples**, which show how the function should work for specific inputs. These examples help us see patterns, and we express those patterns by **circling and labeling** what changes.
- The final step is to **define the function**, which generalizes our examples.

The Animals Dataset

1. This dataset is Animals from an animal shelter, which contains 31 data rows.
2. Some of the columns are:
 - i. name, which contains categorical data, and is of type String. Some example values from this column are: "Toggle", "Fritz", and "Nori".
 - ii. species, which contains categorical data, and is of type String. Some example values from this column are: "cat", "dog".
 - iii. age, which contains quantitative data, and is of type Number. Some example values from this column are: 1, 2, 6.
 - iv. pounds, which contains quantitative data, and is of type Number. Some example values from this column are: 6.5, 35.3, 6.1.

3. Some questions I have about this dataset:

My question is...	Lookup, Compute or Analyze?

Practicing Lookups

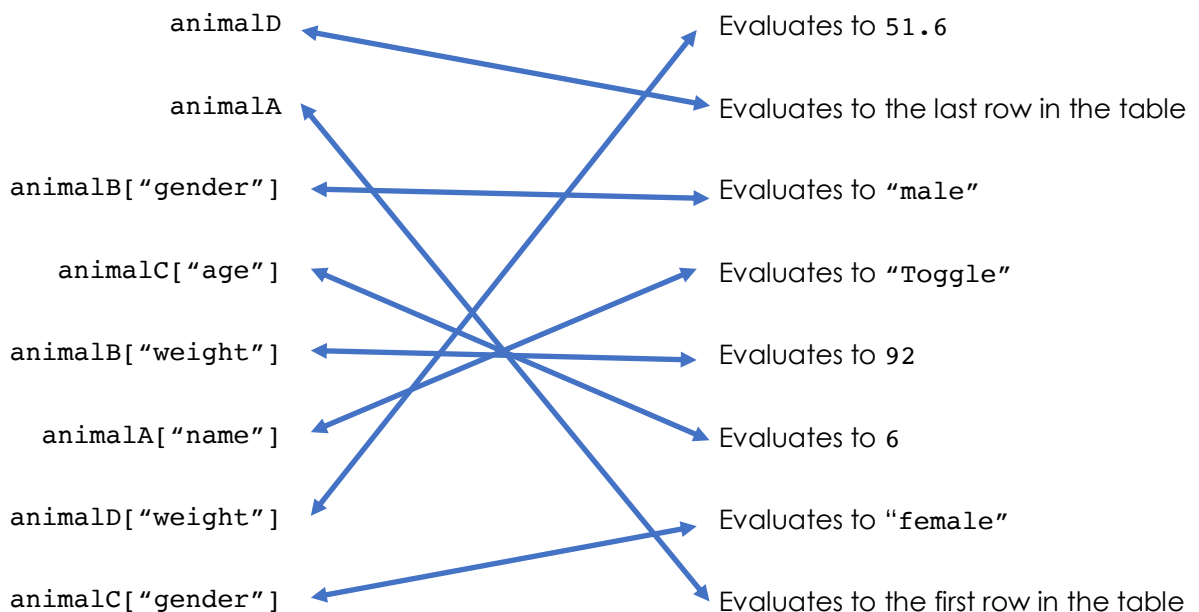
The table below represents four pets at an animal shelter, and four value definitions for rows in that table:

animals-table

name	gender	age	Weight
"Toggle"	"female"	3	48
"Fritz"	"male"	4	92
"Nori"	"female"	6	35.3
"Maple"	"female"	3	51.6

```
animalA = animals-table.row-n(0)
animalB = animals-table.row-n(1)
animalC = animals-table.row-n(2)
animalD = animals-table.row-n(3)
```

v. Match each Pyret expression (left) to the description of what it looks up(right).



vi. Fill in the blanks (left) with the Pyret lookup code that will produce the value (right).

<u><code>animalD["name"]</code></u>	<code>"Maple"</code>
<u><code>animalB["gender"]</code></u>	<code>"male"</code>
<u><code>animalB["age"]</code></u>	<code>4</code>
<u><code>animalA["weight"]</code></u>	<code>48</code>
<u><code>animalC["name"]</code></u>	<code>"Nori"</code>

The Design Recipe

For the word problems below, assume you have `animalA` and `animalB` defined in your code.

Define a function called `is-fixed`, which looks up whether or not an animal is fixed

```
# is-fixed :: (animal :: Row) → Boolean
   name      domain      range
```

```
# Consumes an animal, and looks up the value in the fixed column
```

examples:

```
    is-fixed ( animalA ) is animalA["fixed"]
    is-fixed ( animalB ) is animalB["fixed"]
end
fun is-fixed ( animal ) : animal["fixed"]
end
```

Define a function called `gender`, which consumes a Row of the animals table and looks up the gender of that animal

```
# gender :: (animal :: Row) → String
   name      domain      range
```

```
# Consumes an animal, and produces the value in the gender column
```

examples:

```
    gender ( animalA ) is animalA["gender"]
    gender ( animalB ) is animalB["gender"]
end
fun gender ( animal ) : animal["gender"]
end
```

The Design Recipe

For the word problems below, assume you have `animalA` and `animalB` defined in your code.

Define a function called `is-cat`, which consumes a Row of the `animals` table and computes whether the animal is a cat.

```
# is-cat :: (animal :: Row) → Boolean
   name           domain           range
# Consumes an animal, look up the species column, and computer if species = "cat"
```

examples:

```
    is-cat ( animalA ) is animalA["species"] == "cat"
    is-cat ( animalB ) is animalB["species"] == "cat"
end
fun is-cat ( animal ) : animal["species"] == "cat"
end
```

Define a function called `is-young`, which consumes a Row of the `animals` table and computers whether it is less than four years old.

```
# is-young :: (animal :: Row) → Boolean
   name           domain           range
# Consumes an animal, returns true if the animal is less than 4 years old
```

examples:

```
    is-young ( animalA ) is animalA["age"] < 4
    is-young ( animalB ) is animalB["age"] < 4
end
fun is-young ( animal ) : animal["age"] < 4
end
```

Unit 3

Functions can contain value definitions

We use **Table Plans** to help us use table methods correctly, without making mistakes:

- Like functions, we start with a Contract and Purpose Statement
- But instead of writing *programmed examples*, we sketch out **Sample Tables** and **Results**, based on the Contract and Purpose.
- Then we define the function based on our Sample Table and Result. Every function includes both the table definition (using methods) and a table expression.

Design Recipe

For the word problems below, assume you have `animalA` and `animalB` defined in your code.

Define a function called `birth-year`, which consumes a Row of the `animals` table and produces the year that animal was born.

```
# birth-year :: (animal :: Row) → Number
   name           domain           range
```

```
# Consumes an animal, and produces the year that they were born, subtracting age from
   the current year
```

examples:

```
birth-year ( animalA ) is 2019 - animalA["age"]
birth-year ( animalB ) is 2019 - animalB["age"]
end
fun birth-year ( animal ) : 2019 - animal["age"]
end
```

Define a function called `nametag`, prints out each animal's name in big red letters.

```
# nametag :: (animal :: Row) → Image
   name           domain           range
```

```
# Consumes an animal, and produces an image of their name in big, red letters
```

examples:

```
nametag ( animalA ) is text(animalA["name"], 50, "red")
nametag ( animalB ) is text(animalB["name"], 50, "red")
end
fun nametag ( animal ) : text(animal["name"], 50, "red")
end
```


Playing with Methods

You have the following functions defined below (read them *carefully!*):

```
fun is-fixed(animal): animal["fixed"] end  
fun is-young(animal): animal["age"] < 4 end  
fun nametag(animal): text(animal["name"], 20, "red") end
```

The table **t** below represents four animals at the shelter:

name	gender	age	fixed	weight
"Toggle"	"female"	3	true	48
"Fritz"	"male"	4	true	92
"Nori"	"female"	6	true	35.3
"Maple"	"female"	3	true	51.6

Match each Pyret expression (left) to the description of what it does (right).

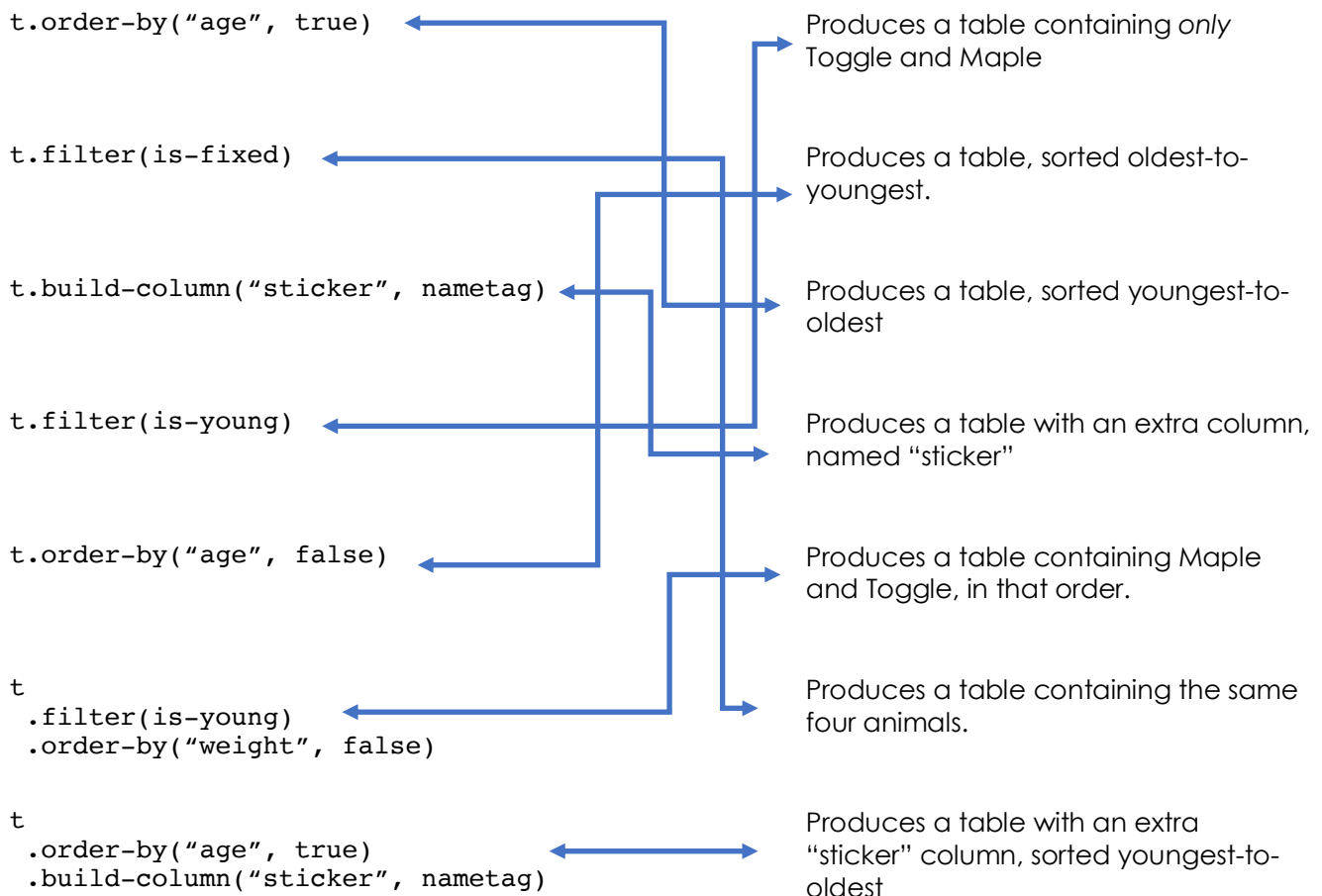


Table Plan

The shelter wants to print up bar charts showing young animal's ages, in alphabetical order. Sometimes they want to do this for every animal, but sometimes they just need it for the cats, or for animals that are fixed.

Define a function `sorted-age-bar`, which takes in a table of animals and computes a bar-chart showing their ages (in alphabetical order), for only the young animals.

Contract and Purpose

`sorted-age-bar` :: `(animals :: Table)` → `Image`

Consume a table of animals, and compute a bar chart showing their ages, in alphabetical order

Where I start, what I type, and what I get back

An example table to start with:

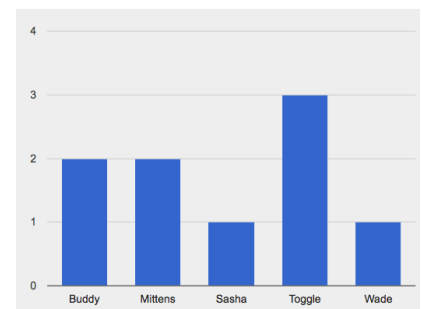
`example-table`

name	...	age
Sasha		1
Toggle		3
Buddy		2
Wade		1
Mittens		2



To use the function, I would type:

`sorted-age-bar(example-table)`



Define the function

Use the relevant methods (circle your helper functions!), then produce a result with the new table.

fun `sorted-age-bar` (`animals`) :

`t = animals`

`.build-column()`

`.filter()`

`.order-by("age", true)`

`bar-chart(t, "name", "age")`

end

Define the table

Are there more columns?

Are there fewer rows?

Are the rows ordered?

Produce the result

Table Plan

The shelter wants to see if there's a relationship between how old an animal is, and how long it takes them to be adopted. Sometimes they want to do this for every animal, but sometimes they just need it for the cats, or for animals that are young. Define a function `age-adopted-scatter`, which takes in a table of animals and computes a scatter-plot showing only the fixed animals, with their ages on the x-axis and weeks to be adopted on the y-axis.

Contract and Purpose

`age-adopted-scatter` :: `(animals :: Table)` → `Image`

Consume a table of animals, and compute a scatterplot showing their ages on the x-axis, and weeks be adopted on the y-axis

Where I start, what I type, and what I get back

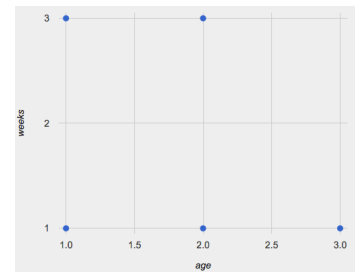
A sample table to start with:

name	...	age	weeks
Sasha		1	3
Toggle		3	1
Buddy		2	3
Wade		1	1
Mittens		2	1



To use the function, I would type:

`age-adopted-scatter(sample)`



Define the function

Use the relevant methods (circle your helper functions!), then produce a result with the new table.

fun `age-adopted-scatter` (`animals`) :

`t = animals`

`.build-column(`

`.filter(`

`.order-by(`

`scatter-plot(t, "name", "age", "weeks")`

end

Define the table

Are there more columns?

Are there fewer rows?

Are the rows ordered?

Produce the result

Unit 4

Bar charts show the *absolute* quantity of each row in a dataset. The larger the quantity, the longer the bar. Bar charts provide a visual representation of values in a dataset.

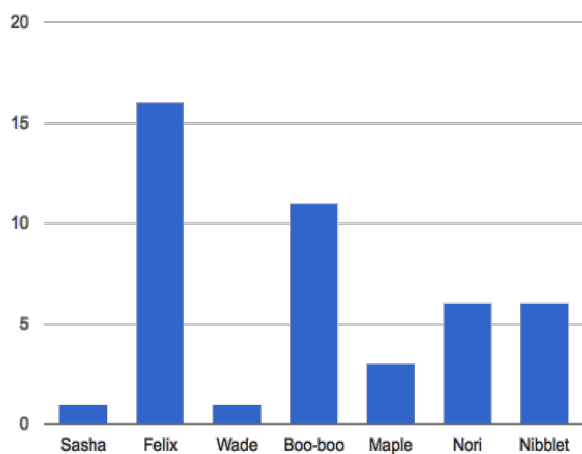
Pie charts show the *relative* quantity of each row in a dataset. The greater the percentage, the larger the pie slice. Pie charts provide a visual representation of proportions in a dataset.

Choosing a Sample Table is important when coming up with small examples for Table Plans. A good sample table has:

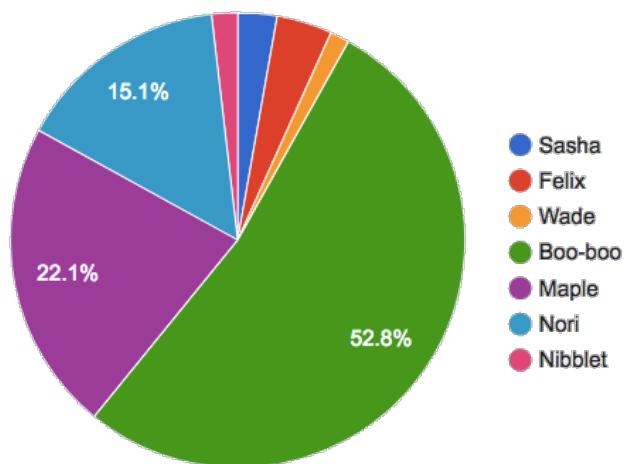
- At least all the relevant columns
- Enough rows to accurately represent the dataset
- Rows that are randomly-ordered

Quantity Charts in the `Animals` Dataset

Below are two **quantity charts** made from subsets of the animals table



Animals Ages (yrs)



Animals Weights (lbs)

[illegible]

Why are some questions easier to answer with one kind of chart or another?

Bad Sample Tables!

For each word problem, a Sample Table must have (1) all the columns that matter, (2) a representative sample of the rows, and be in (3) random order. For each problem below, check the boxes if the Sample Table meets those criteria.

1. The shelter wants to a scatter plot showing the age of the cats v. their weight

name	species	age	fixed	legs	pounds	weeks
Sasha	cat	1	FALSE	4	6.5	3
Mittens	cat	2	TRUE	4	7.4	5
Sunflower	cat	5	TRUE	4	8.1	10

- ✓ Relevant columns
- ✓ Representative sample of rows
- ✓ Random order

2. The shelter wants a pie chart showing all the dogs' weight

name	species	age
Fritz	dog	4
Wade	cat	2
Nibblet	rabbit	6
Daisy	dog	5

- ☐ Relevant columns
- ☐ Representative sample of rows
- ✓ Random order

3. Sort all the animals alphabetically by name

name	species	age	fixed	legs	pounds	weeks
Ada	dog	2	TRUE	4	32	3
Bo	dog	4	TRUE	4	76.1	10
Boo-boo	dog	11	TRUE	4	123	10

- ✓ Relevant columns
- ☐ Representative sample of rows
- ☐ Random order

4. Make a bar chart for all the fixed animals

name	species	age	fixed	legs	pounds	weeks
Sasha	cat	1	FALSE	4	6.5	3

- ✓ Relevant columns
- ☐ Representative sample of rows
- ☐ Random order

Table Plan

Define a function `pie-pounds-young`, which takes in a Table of animals and creates a pie chart of the animals' weight, but only for animals that are young.

Contract and Purpose																				
#	<code>pie-pounds-young</code>	:: (animals :: Table) → Image																		
#	<i>Consumes a table of animals, filters to show only young animals, and produces a pie chart of their weight</i>																			
Where I start, what I type, and what I get back																				
A sample table to start with:		To use the function, I would type:																		
<code>sample-table</code>	→	<code>pie-pounds-young(sample-table)</code>																		
<table border="1"> <thead> <tr> <th>name</th> <th>age</th> <th>pounds</th> </tr> </thead> <tbody> <tr> <td>Snowcone</td> <td>...</td> <td>6.1</td> </tr> <tr> <td>Lucky</td> <td>...</td> <td>45.4</td> </tr> <tr> <td>Hercules</td> <td>...</td> <td>13.4</td> </tr> <tr> <td>Toggle</td> <td>...</td> <td>48</td> </tr> <tr> <td>Snuggles</td> <td>...</td> <td>0.1</td> </tr> </tbody> </table>	name	age	pounds	Snowcone	...	6.1	Lucky	...	45.4	Hercules	...	13.4	Toggle	...	48	Snuggles	...	0.1		
name	age	pounds																		
Snowcone	...	6.1																		
Lucky	...	45.4																		
Hercules	...	13.4																		
Toggle	...	48																		
Snuggles	...	0.1																		
Define the function																				
Use the relevant methods (circle your helper functions!), then produce a result with the new table.																				
fun	<code>pie-pounds-young</code>	(<code>animals</code>) :																		
	<code>t = animals</code>	<u>Define the table</u>																		
	<code>.filter(is-young)</code>	Are there more columns? Are there fewer rows? Are the rows ordered?																		
	<code>pie-chart(t, "name", "pounds")</code>	<u>Produce the result</u>																		
end																				

My Dataset

1. This dataset is _____, which contains _____ data rows.
2. Some of the columns are:
 - i. _____, which contains _____ data, and is of type _____. Some example values from this column are: _____.
 - ii. _____, which contains _____ data, and is of type _____. Some example values from this column are: _____.
 - iii. _____, which contains _____ data, and is of type _____. Some example values from this column are: _____.
 - iv. _____, which contains _____ data, and is of type _____. Some example values from this column are: _____.

3. Some questions I have about this dataset:

My question is...	Lookup, Compute or Analyze?

My Dataset

What are two ways you might want to *order* this dataset?

1) _____

2) _____

What are two subsets into which you might *filter* this dataset?

1) _____

2) _____

What are two new columns you might want to *build* from this dataset?

1) _____

2) _____

Design Recipes – Filtering Rows

What are two criteria you might want to *filter* by? Write your own word problems below, and solve them using the Design Recipe.

Define a function called _____ **, which consumes a Row of the**
_____ **table and** _____

```
# _____ :: _____ → _____  
      name          domain          range  
# _____
```

examples:

```
    _____ ( _____ ) is _____  
    _____ ( _____ ) is _____  
end  
fun _____ ( _____ ) : _____  
end
```

```
# _____ :: _____ → _____  
      name          domain          range  
# _____
```

examples:

```
    _____ ( _____ ) is _____  
    _____ ( _____ ) is _____  
end  
fun _____ ( _____ ) : _____  
end
```

Design Recipes – Building Columns

What are two columns you might want to *build* for your dataset? Write your own word problems below, and solve them using the Design Recipe.

```
# _____ :: _____ → _____  
      name                domain                range
```

```
# _____
```

examples:

```
    _____ ( _____ ) is _____  
    _____ ( _____ ) is _____  
end  
fun _____ ( _____ ) : _____  
end
```

```
# _____ :: _____ → _____  
      name                domain                range
```

```
# _____
```

examples:

```
    _____ ( _____ ) is _____  
    _____ ( _____ ) is _____  
end  
fun _____ ( _____ ) : _____  
end
```

Quantity Charts in My Dataset

Describe two of the pie or bar charts you made from your dataset.

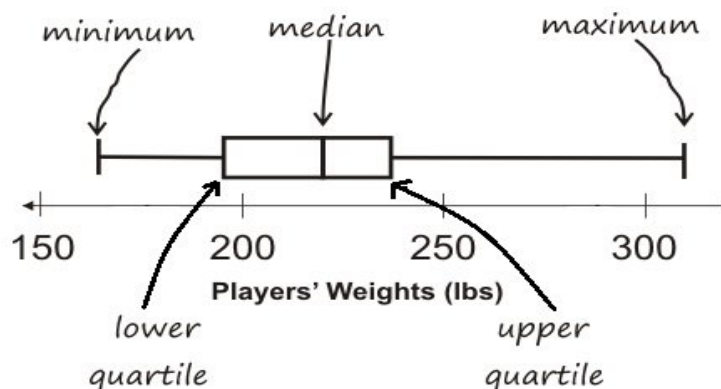
1) I made a pie / bar chart, showing the column in your dataset for your subset (for example, "fixed dogs at the shelter").

2) I made a _____ chart, showing the _____ for _____.

[illegible]

Unit 5

- There are three ways to measure the “center” of a dataset, to talk about a whole column of data using just one number:
 - The **mean** of a dataset is the average of all the numbers
 - The **median** of a dataset is a value that is smaller than half the dataset, and larger than the other half
 - The **modes** of a dataset are the numbers that appear the most often.
- Data Scientists can also measure the “variation” of a dataset using a **five number summary**:
 - The **minimum** – the smallest value in the dataset
 - The **first, or “lower” quartile (Q1)** – the median value that separates the first quarter of the values in the dataset from the second quarter
 - The **second quartile (Q2)** – the median value which separates the entire dataset into “top” and “bottom” halves.
 - The **third, or “upper” quartile (Q3)** – the median value that separates the third quarter of the values in the dataset from the fourth quarter
 - The **maximum** – the largest value in the dataset
- The **five number summary** can be used to draw a **box-and-whisker plot**.



Summarizing Columns in Animals

1) The column I choose to measure is weeks

Measures of Center

The three measures for this column are:

Mean (Average)	Median	Mode(s)
6.0689	4	1

2) Since the mean is higher than the median, this suggests that there may
[higher/lower]

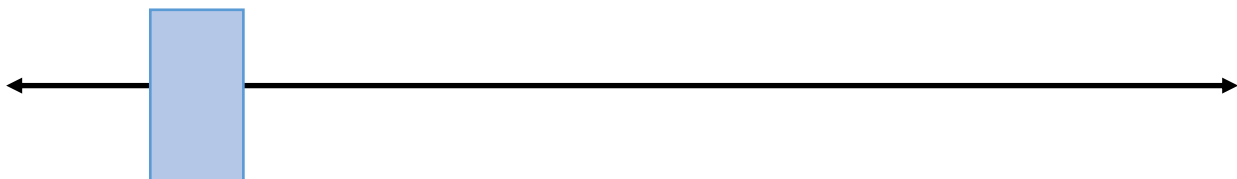
be outliers representing a few animals who took a long time to be adopted.
[explain your outliers!]

Measures of Variation

My five-number summary is:

Minimum	Q1	Q2 (Median)	Q3	Maximum
1	2.5	4	8	30

A box plot can be drawn from this summary on the number line below:



From this summary and box-plot, I conclude:

The vast majority of animals are adopted before 8 weeks in the shelter, but there are a number of outliers (such as the maximum of 30).

Interpreting Variation

Consider the following list dataset, representing the annual income of ten people:

\$65k, \$12k, \$14k, \$280k, \$15k, \$22k, \$45k, \$34k, \$45k, \$175k

1. In the space below, rewrite this dataset in **sorted order**.

\$12k, \$14k, \$15k, \$22k, \$34k, \$45k, \$45k, \$65k, \$175k, \$280k

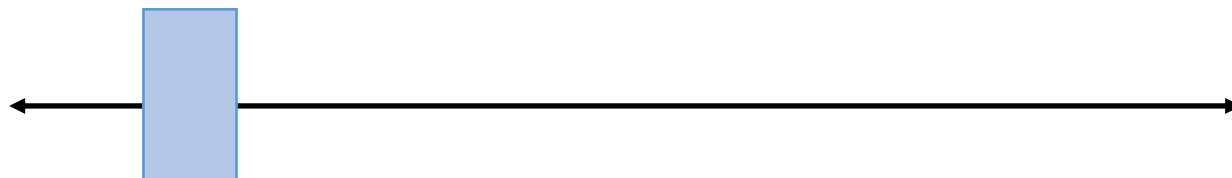
2. In the table below, compute the **measures of center** for this dataset.

Mean (Average)	Median	Mode(s)
70,700	39,500	45,000

3. In the table below, compute the **five number summary** of this dataset.

Minimum	Q1	Q2 (Median)	Q3	Maximum
12,000	15,000	39,500	65,000	280,000

4. On the number line below, draw a **box plot** for this dataset.



5. The following statements are *correct*...but misleading. Write down the reason why.

Statement	Why it's misleading
"They're rich! The average person makes more than \$70k dollars!"	While the mean is close to \$70k, there are some very high earning outliers pushing the average up.
"It's a middle-income list: the most common salary is \$45k/yr!"	In the full dataset, more than half of the entries are people making less than \$45k, making the mode misleading.
"This group is really diverse, with people making as little as 12k and as much as \$280k!"	While the spread of incomes is large, the vast majority are still making less than \$65k, with very high earning outliers.

Table Plan

The Animal Shelter Bureau would like to study the distribution of weeks-until-adoption for fixed animals housed at shelters around the country. They need a function that consumes an Animals table, filters to show only the fixed animals, and produces a box-plot for the weeks column. Define a function called `fixed-weeks-box` below.

Contract and Purpose

`fixed-weeks-box` :: `(animals :: Table)` → `Image`

Consumes a table of animals, filters only the fixed animals, and produces a box plot of their weeks until adoption

Where I start, what I type, and what I get back

A sample table to start with:

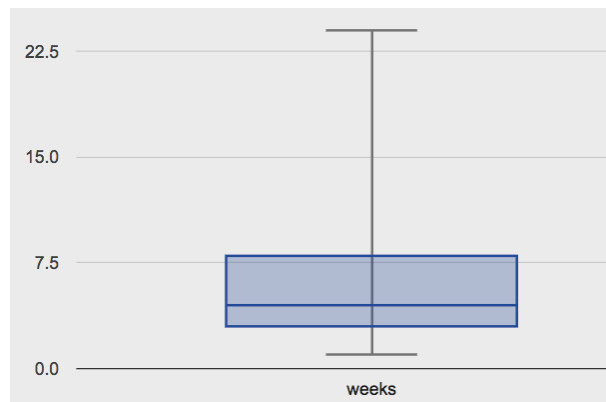
To use the function, I would type:

sample table



fixed-weeks-box(sample table)

name	species	age	fixed	legs	weight	weeks
Snowcone	cat	2	TRUE	4	6.1	5
Lucky	dog	3	TRUE	3	45.4	9
Hercules	cat	3	FALSE	4	13.4	7
Toggle	dog	3	TRUE	4	48	3
Snuggles	tarantula	2	FALSE	8	0.1	1



Define the function

Use the relevant methods (circle your helper functions!), then produce a result with the new table.

fun `fixed-weeks-box` (`animals`) :

`t = animals-table`

`.filter(is-fixed)`

`box-plot(t, "weeks")`

end

Define the table

Are there more columns?

Are there fewer rows?

Are the rows ordered?

Produce the result

Summarizing a Column in My Dataset

The column I choose to measure is _____

Measures of Center

The three measures for this column are:

Mean (Average)	Median	Mode(s)

3) Since the mean is _____ than the median, this suggests that there may
[higher/lower]

be outliers representing_____.
[explain your outliers!]

Measures of Variation

My five-number summary is:

Minimum	Q1	Q2 (Median)	Q3	Maximum

A box plot can be drawn from this summary on the number line below:



From this summary and box-plot, I conclude:

Unit 6

Frequency Bar charts show the number of rows belonging to a given category. The more rows in each category, the longer the bar.

- *Frequency bar charts provide a visual representation of the frequency of values in a **categorical** column.*
- Since categorical data cannot be ordered, there is no strict ordering of bars in a frequency bar chart.

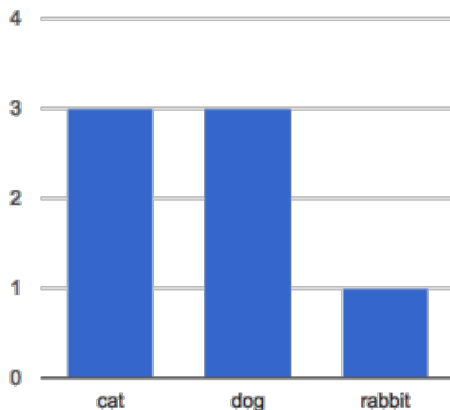
Histograms show the number of rows that fall within certain ranges, or "bins" of a dataset. The more rows that fall within a particular "bin", the longer the bar.

- *Histograms provide a visual representation of the frequency of values in a **quantitative** column.*
- Quantitative data can be ordered, so the bars of a histogram are always sorted.
- When dealing with histograms, it's important to select a good **bin size**. If the bins are too small or too large, it is difficult to see the distribution in the dataset.

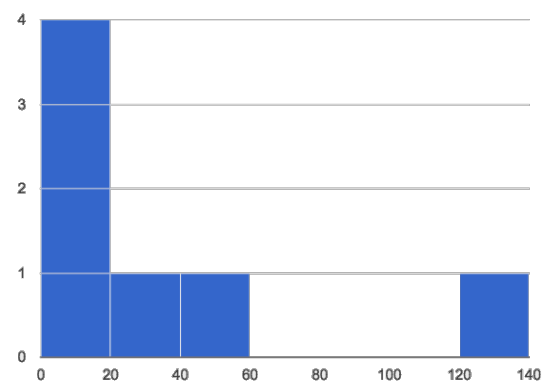
Frequency Charts in the Animals Dataset

name	species	age	pounds
"Sasha"	"cat"	1	6.5
"Boo-boo"	"dog"	11	123
"Felix"	"cat"	16	9.2
"Nori"	"dog"	6	35.3
"Wade"	"cat"	1	3.2
"Nibblet"	"rabbit"	6	4.3
"Maple"	"dog"	3	51.6

- How many cats are there? 3
- How many dogs are there? 3
- How many animals are between 3-6 years old? 3
- How many weigh between 0-5 pounds? 2
- Are there more animals weighing 0-5 than 6-10 pounds? Yes
- The charts below are based on the Sample Table above. What is each one measuring? Write down your guess underneath each one.



Amount of each species



Frequency of animal weights

Table Plan

Define a function `freq-bar-gender`, which takes in a Table of animals and creates a frequency bar chart showing how many animals are male v. female.

Contract and Purpose

`freq-bar-gender` :: (animals :: Table) → Image

Consumes a table of animals and produces a frequency bar chart of their genders, for *fixed* animals

Examples

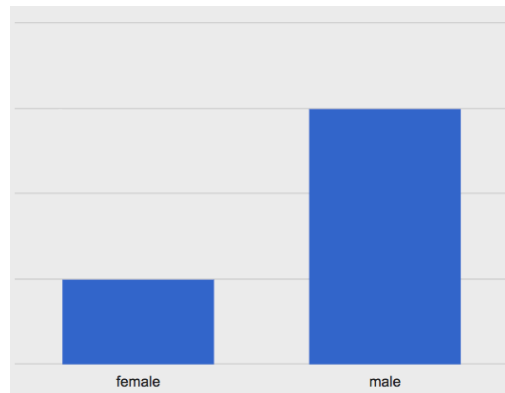
Make a Start Table and a result based on that table.

animals-table



freq-bar-gender(animals-table)

name	species	age	gender
Fritz	dog	4	male
Wade	cat	2	male
Nibblet	rabbit	6	male
Daisy	dog	5	female



Define the function

Use the relevant methods (circle your helper functions!), then produce a result with the new table.

fun `freq-bar-gender` (`animal`) :

`t = animals`

`freq-bar-chart(t, "gender")`

end

Define the table

Are there more columns?

Are there fewer rows?

Are the rows ordered?

Produce the result

Table Plan

Define a function `histogram-adoption`, which takes in a Table of animals and creates a histogram showing how long it took for animals to get adopted

Contract and Purpose

`histogram-adoption` :: (animals :: Table) → Image

Consumes a table of animals and produces a histogram showing how long it took for the animals to get adopted

Examples

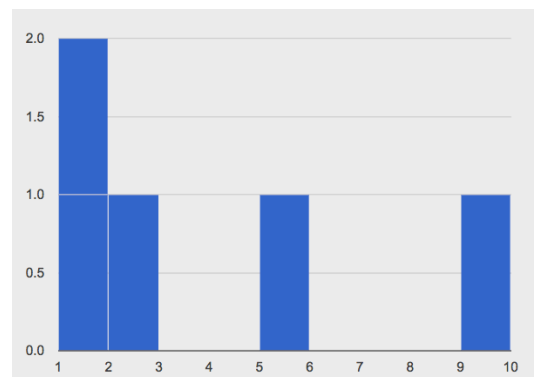
Make a Start Table and a result based on that table.

animals-table



histogram-adoption(animals-table)

name	species	age	fixed	legs	weight	weeks
Snowcone	cat	2	TRUE	4	6.1	5
Lucky	dog	3	TRUE	3	45.4	9
Hercules	cat	3	FALSE	4	13.4	7
Toggle	dog	3	TRUE	4	48	3
Snuggles	tarantula	2	FALSE	8	0.1	1



Define the function

Use the relevant methods (circle your helper functions!), then produce a result with the new table.

fun `histogram-adoption` (`animals`) :

`t = animals`

`histogram(t, "weeks", 1)`

end

Define the table

Are there more columns?

Are there fewer rows?

Are the rows ordered?

Produce the result

Visualizing My Dataset

Describe two of the histograms or frequency bar charts you made from your dataset.

1) I made a _____, showing the _____ for
 histogram / frequency bar chart column in your dataset

your subset (for example, "fixed dogs at the shelter")

2) I made a _____, showing the _____ for _____.

[illegible]

Matching Charts to Questions

For each of the questions below, draw a line to the chart that will best answer it. (You may find that more than one question is best answered by the same chart!)

1. Are there more of the animals at the shelter fixed or unfixed?	←
2. How many weeks did each cat wait to be adopted?	←
3. How many male v. female dogs are there?	←
4. How many animals have 4 legs? 8? 3?	←
5. What percent of the total weight at the shelter is made up by Boo-boo?	←
6. What is the distribution of weights across all the animals older than 3?	←
7. How many animals are there of each species?	←
8. Who waited the longest to be adopted?	←

→	Pie Chart
→	Bar Chart
→	Frequency Bar Chart
→	Histogram

```
graph LR; Q1[1. Are there more of the animals at the shelter fixed or unfixed?] --> PC[Pie Chart]; Q2[2. How many weeks did each cat wait to be adopted?] --> BC[Bar Chart]; Q3[3. How many male v. female dogs are there?] --> BC; Q4[4. How many animals have 4 legs? 8? 3?] --> BC; Q5[5. What percent of the total weight at the shelter is made up by Boo-boo?] --> FBC[Frequency Bar Chart]; Q6[6. What is the distribution of weights across all the animals older than 3?] --> H[Histogram]; Q7[7. How many animals are there of each species?] --> H; Q8[8. Who waited the longest to be adopted?] --> H;
```


Unit 7

- **Scatter Plots** show the relationship between two quantitative columns. Each row in the dataset is represented by a point, with one column providing the x-value and the other providing the y-value. The resulting “point cloud” makes it possible to look for a relationship between those two columns.
- If the points in a scatter plot appear to follow a pattern, it is possible that a relationship – or **correlation** – exists between those two columns.
- If there is a pattern to the points in a scatter plot, points that are far away from the pattern are called **outliers**.
- We can express this correlation by drawing line through the data cloud, so that the distance between the line and each of the points is as small as possible. This line is called the **line of best fit** – or **predictor function** - and allows us to make predictions based on the dataset.

(Dis)Proving a Claim

“Younger animals are cuter, so they get adopted faster.”

Do you agree? If so, why?

I hypothesize...

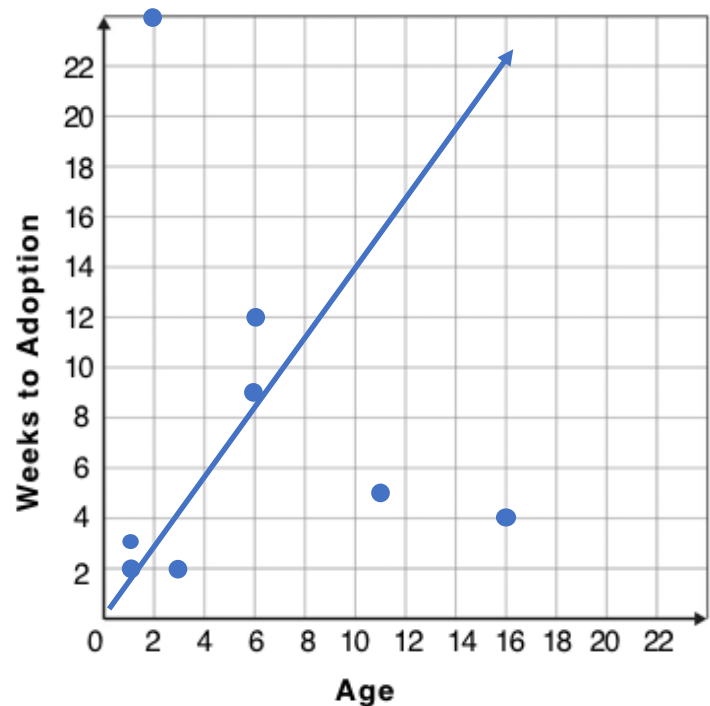
that younger animals *will* get adopted faster, possibly because
they are considered cuter, but there may be other factors
causing them to get adopted faster.

What would you look for in the dataset to see if you are right?

I would look at both the ages and number of weeks until adoption
for each animal to see if there was a correlation. I would also
want to collect more data, such as conduct a survey of adopters.

Creating a Scatter Plot

name	species	age	weeks
"Sasha"	"cat"	1	3
"Boo-boo"	"dog"	11	5
"Felix"	"cat"	16	4
"Buddy"	"lizard"	2	24
"Nori"	"dog"	6	9
"Wade"	"cat"	1	2
"Nibblet"	"rabbit"	6	12
"Maple"	"dog"	3	2



1. **For each row in the Sample Table on the left, add a point to the scatter plot on the right.** The first 3 rows have been completed for you. Use the values from the age column for the x-axis, and values from the weeks column for the y-axis.
2. Do you see a pattern? Do the points seem to shift up or down as age increases?
Draw a line on the scatter plot to show this pattern.

3. Does the line slope upwards or downwards?

Slightly upwards

4. Are the points mostly close to the line?

A few points are close to the line, but as ages increase the points get much farther apart.

Table Plan

Define a function `cats-age-weeks`, which takes in a Table of animals and creates a scatter plot of all the cats, tracking their `age` on the x-axis and the number of `weeks` it took for them to be adopted on the y-axis.

Contract and Purpose

`cats-age-weeks` :: `(animals :: Table)` → `Image`

Consumes a table of animals, creates a scatter plot of only the cat's ages and their weeks to adoption

Examples Where I start, what I type, and what I get back

A sample table to start with:

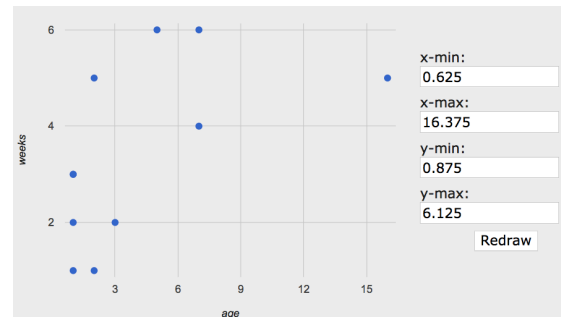
To use the function, I would type:

`animals-table`



`cats-age-weeks(animals-table)`

name	species	age	fixed	legs	weight	weeks
Snowcone	cat	2	TRUE	4	6.1	5
Lucky	dog	3	TRUE	3	45.4	9
Hercules	cat	3	FALSE	4	13.4	7
Toggle	dog	3	TRUE	4	48	3
Snuggles	tarantula	2	FALSE	8	0.1	1



Define the function

Use the relevant methods (circle your helper functions!), then produce a result with the new table.

`fun` `cats-age-weeks` (`animals`) :

`t = animals-table`

`.filter(is-cat)`

`scatter-plot(t, "name", "age", "weeks")`

`end`

Define the table

Are there more columns?

Are there fewer rows?

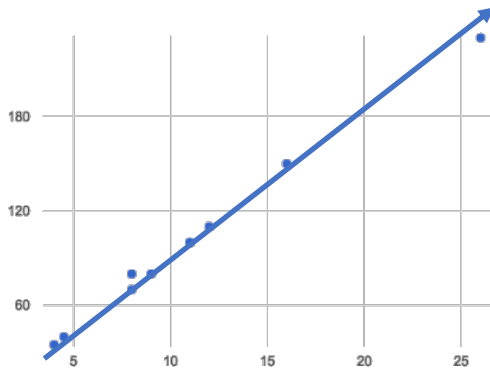
Are the rows ordered?

Produce the result

Drawing Predictors

For each of the scatter plots below, draw a **predictor line** that fits best.

A

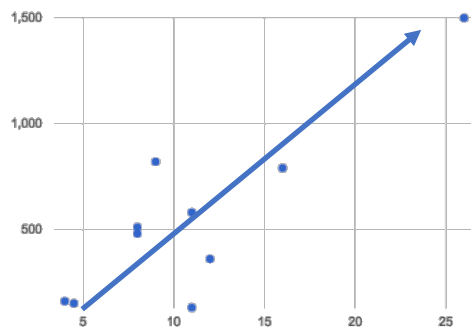


fat (g) v. calories-from-fat in common menu items

Direction: Positive Negative None

Strength: Strong Weak

B

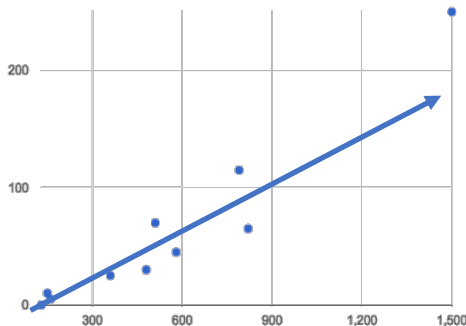


fat (g) v. sodium (g) in common menu items

Direction: Positive Negative None

Strength: Strong Weak

C

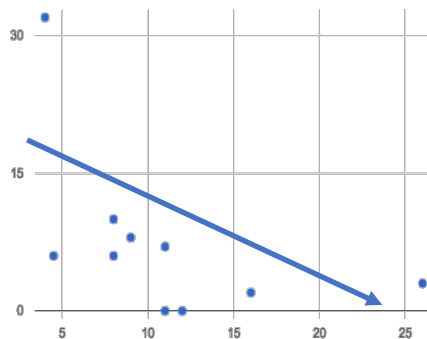


sodium (g) v. cholesterol (mg) in common menu items

Direction: Positive Negative None

Strength: Strong Weak

D



fat (g) v. sugar (g) in common menu items

Direction: Positive Negative None

Strength: Strong Weak

Correlations in My Dataset

1) There may be a correlation between _____ and
column
_____. I think it is a _____,
column strong / weak positive / negative
correlation, because _____
_____. It might be stronger if I looked
at _____.
a subset or extension of my data

2) There may be a correlation between _____ and
column
_____. I think it is a _____,
column strong / weak positive / negative
correlation, because _____
_____. It might be stronger if I looked
at _____.
a subset or extension of my data

3) There may be a correlation between _____ and
column
_____. I think it is a _____,
column strong / weak positive / negative
correlation, because _____
_____. It might be stronger if I looked
at _____.
a subset or extension of my data

Unit 8

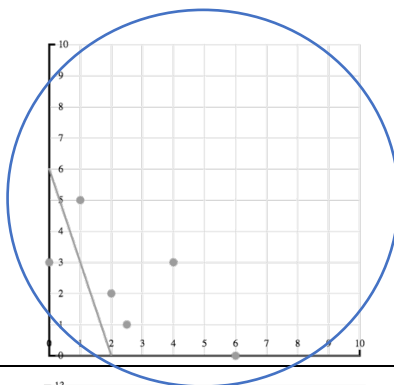
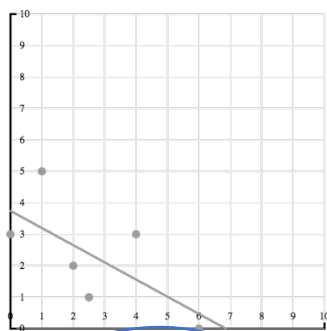
- Given a **predictor function** and a scatter plot, we can compute the error by adding the squares of all the distances between the function and each point in the plot. The error is called the **r^2 statistic**, which tells us *how much of the variation in the y-axis can be explained by the x-axis*.
- A **strong correlation** will have a large r^2 . A **weak correlation** will have a small r^2 .
- A **positive correlation** means the slope of the line of best fit is positive. A **negative correlation** means the slope is negative.
- **Linear Regression** is a way of computing the **line of best fit**, by taking a scatter plot and deriving the slope and y-intercept for a line that has the smallest possible r^2 .
- **Correlation is not causation!** Correlation only suggests that two measures are *related*, but does not tell us if one *causes* the other. For example, hot days are *correlated* with people running their air conditioners, air conditioners do not *cause* hot days!

Grading Predictors

Below are the scatter plots for data sets A-D, with two different lines predictor lines drawn on top. For plots A-D:

1. Circle the plot with the line that fits better
2. Give the plot you circled a grade between 0 (no correlation) and 1 (perfect correlation)

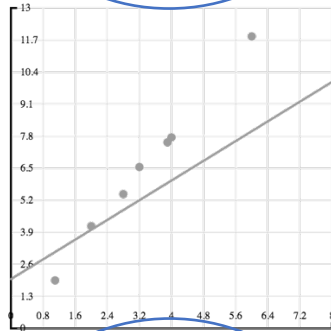
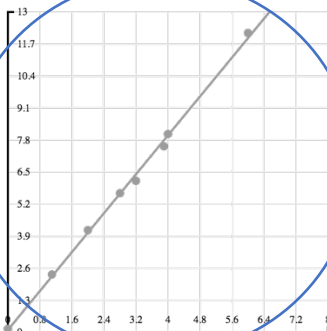
A



Strength of
Correlation:

0.2

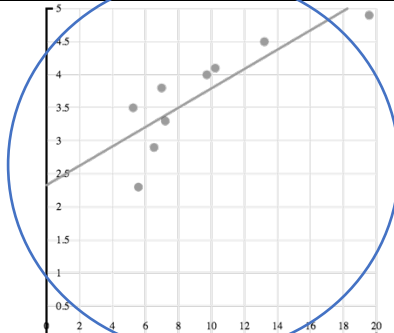
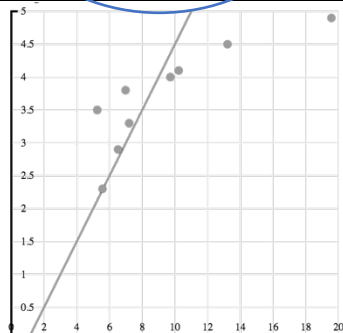
B



Strength of
Correlation:

0.95

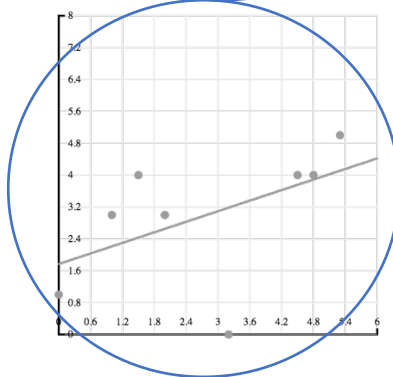
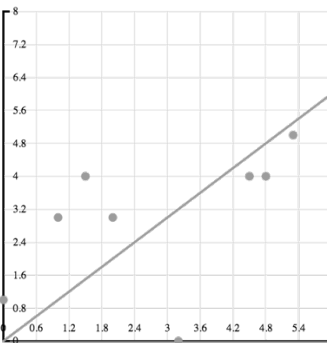
C



Strength of
Correlation:

0.65

D



Strength of
Correlation:

0.4

Regression Analysis in the animals Dataset

I performed a linear regression on cats at the shelter, and
dataset or subset

found a weak ($r^2=0.321$), positive correlation between
a strong/weak ($r^2=$), positive/negative

age of the cats (in weeks) and number of weeks to adoption. From this, I
[x-axis] [y-axis]

conclude that 32.1% of the variability in adoption time is explained by the
 r^2 % of the variation in [y-axis] is explained by [x-axis]

age of the cat. I would predict that a 1 year increase in
[x-axis units]

age is associated with a 0.23 week increase in adoption time.
[x-axis] [slope, y-units] [increase/decrease] [y-axis]

I performed a linear regression on _____, and
dataset or subset

found _____ correlation between
a strong/weak ($r^2=$), positive/negative

_____ and _____. From this, I
[x-axis] [y-axis]

conclude that _____
 r^2 % of the variation in [y-axis] is explained by [x-axis]

_____. I would predict that a 1 _____ increase in
[x-axis units]

_____ is associated with a _____ in _____.
[x-axis] [slope, y-units] [increase/decrease] [y-axis]

I performed a linear regression on _____, and
dataset or subset

found _____ correlation between
a strong/weak ($r^2=$), positive/negative

_____ and _____. From this, I
[x-axis] [y-axis]

conclude that _____
 r^2 % of the variation in [y-axis] is explained by [x-axis]

_____. I would predict that a 1 _____ increase in
[x-axis units]

_____ is associated with a _____ in _____.
[x-axis] [slope, y-units] [increase/decrease] [y-axis]

Regression Analysis in My Dataset

I performed a linear regression on _____, and
dataset or subset

found _____ correlation between
a strong/weak ($r^2=$ ____), positive/negative

_____ and _____. From this, I
[x-axis] [y-axis]

conclude that _____
 r^2 % of the variation in [y-axis] is explained by [x-axis]

_____. I would predict that a 1 _____ increase in
[x-axis units]

_____ is associated with a _____ in _____.
[x-axis] [slope, y-units] [increase/decrease] [y-axis]

I performed a linear regression on _____, and
dataset or subset

found _____ correlation between
a strong/weak ($r^2=$ ____), positive/negative

_____ and _____. From this, I
[x-axis] [y-axis]

conclude that _____
 r^2 % of the variation in [y-axis] is explained by [x-axis]

_____. I would predict that a 1 _____ increase in
[x-axis units]

_____ is associated with a _____ in _____.
[x-axis] [slope, y-units] [increase/decrease] [y-axis]

I performed a linear regression on _____, and
dataset or subset

found _____ correlation between
a strong/weak ($r^2=$ ____), positive/negative

_____ and _____. From this, I
[x-axis] [y-axis]

conclude that _____
 r^2 % of the variation in [y-axis] is explained by [x-axis]

_____. I would predict that a 1 _____ increase in
[x-axis units]

_____ is associated with a _____ in _____.
[x-axis] [slope, y-units] [increase/decrease] [y-axis]

Unit 9

Threats to Validity can undermine a conclusion, even if the analysis was done correctly. Some examples of threats are:

- **Selection bias** – identifying the favorite food of the rabbits won't tell us anything reliable about what all the animals eat.
- **Sample size** – averaging the age of only three animals won't tell us anything reliable about the age of animals at the shelter!
- **Sample error** – surveying dogs when they are puppies won't tell us anything reliable about overall dog behavior, since their behavior changes as they age.
- **Confounding variables** – if the person surveying the animals has a piece of bacon in their pocket, they will incorrectly find that all dogs are friendly!

Threats to Validity

Some volunteers from the animal shelter surveyed a group of pet owners at a local dog park. They found that almost all of the owners were there with their dogs, and from this survey they concluded that dogs are the most popular pet in the region.

What are some possible threats to the validity of this conclusion?

Not many people are likely to walk their cats at the park, so if the volunteers only surveyed pet owners at the park, dogs are likely to be more highly represented in their sampling.

The animal shelter noticed a large increase in pet adoptions between Thanksgiving and Valentine's Day. They conclude that at this current rate, there will be a huge demand for pets this Spring.

What are some possible threats to the validity of this conclusion?

Lots of people may be adopting animals during the holiday season, so these past patterns are unlikely to predict future patterns in adoption rates.

Threats to Validity

The animal shelter wanted to find out what kind of food to buy for their animals. They took a random sample of two animals and the food they eat, and found that spider and rabbit food was by far the most popular cuisine!

What are some possible threats to the validity of this conclusion?

A random sample may not be representative of the whole group of pets. In this case, there are many more dogs and cats than spiders and rabbits at the shelter, so using this random sample to draw conclusions about the whole group is wrong!

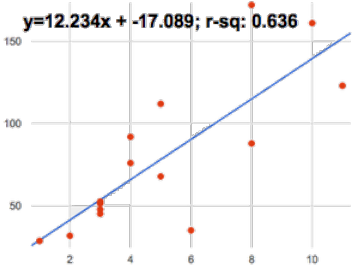
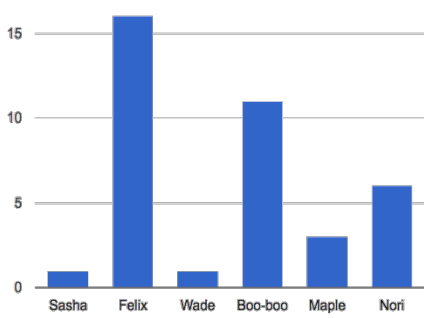
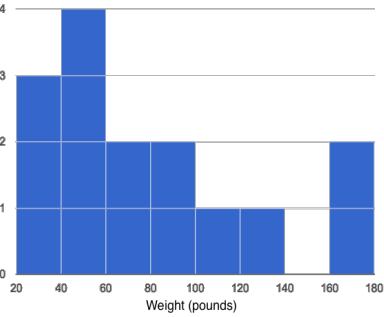
A volunteer opens the shelter in the morning and walks all the dogs. At mid-day, another volunteer feeds all the dogs and walks them again. In the evening, a third volunteer walks the dogs a final time, and closes the shelter. The volunteers report that the dogs are much friendlier and more active at mid-day, so the shelter staff assume the second volunteer must be better with animals than the others.

What are some possible threats to the validity of this conclusion?

There may be other reasons the dogs are happier at mid-day than morning and evening- for instance, mid-day is when they eat lunch, which is likely to make the dogs very excited!

Fake News!

Every claim below is *wrong*! Your job is to figure out why, by looking at the data.

	Data	Claim	Why it's wrong
1	The average player on a basketball team is 6'1".	"Most of the players on the team are taller than 6'."	The average is based on all the players, and there may be outliers pushing the average height up-average tells you nothing about the majority of the players.
2	After performing linear regression on census data, a positive correlation ($r^2=0.18$) was found between people's height and salary.	"Taller people get paid more."	Only 18% of the variation in salary is based on height, which is not a large enough r-squared value to say that taller people get paid more.
3		"According to the predictor function indicated here, the value on the x-axis is will predict the value on the y-axis 63.6% of the time."	The r-squared value of 0.636 does not mean how often the y-value will be predicted, rather what percent of variation in the y-value is based on the x-value.
4	 Bar Chart of Pet Ages	"According to this bar chart, Felix makes up a little more than 15% of the total ages of all the animals in the dataset."	Bar charts are not the most appropriate image for showing the percentage of each measurement based on the total- pie charts should be used for that info. This bar chart shows that Felix is a little more than 15 years old.
5		"According to this histogram, most animals weigh between 40 and 60 pounds."	More animals fit into the histogram bin between 40-60 pounds than any other bin, but that doesn't mean that most animals weigh between 40-60 pounds.
6	After performing linear regression, a negative correlation ($r^2=0.91$) was found between the number of hairs on a person's head and their likelihood of owning a wig.	"Owning wigs causes people to go bald."	Though there is a strong correlation between hair and owning a wig, correlation does NOT equal causation.

Blank Recipes, Table Plans, and References

Design Recipes

```
# _____ :: _____ → _____  
      name                domain                range
```

```
# _____
```

examples:

```
    _____ ( _____ ) is _____  
end _____ ( _____ ) is _____  
fun _____ ( _____ ) : _____  
end
```

```
# _____ :: _____ → _____  
      name                domain                range
```

```
# _____
```

examples:

```
    _____ ( _____ ) is _____  
end _____ ( _____ ) is _____  
fun _____ ( _____ ) : _____  
end
```

Design Recipes

```
# _____ :: _____ → _____  
      name                domain                range
```

```
# _____
```

examples:

```
      _____ ( _____ ) is _____  
      _____ ( _____ ) is _____  
end  
fun _____ ( _____ ) : _____  
end
```

```
# _____ :: _____ → _____  
      name                domain                range
```

```
# _____
```

examples:

```
      _____ ( _____ ) is _____  
      _____ ( _____ ) is _____  
end  
fun _____ ( _____ ) : _____  
end
```

Design Recipes

```
# _____ :: _____ → _____  
      name                domain                range
```

```
# _____
```

examples:

```
      _____ ( _____ ) is _____  
      _____ ( _____ ) is _____  
end  
fun _____ ( _____ ) : _____  
end
```

```
# _____ :: _____ → _____  
      name                domain                range
```

```
# _____
```

examples:

```
      _____ ( _____ ) is _____  
      _____ ( _____ ) is _____  
end  
fun _____ ( _____ ) : _____  
end
```

Table Plan

Contract and Purpose

_____ :: _____ → _____

Examples

Make a Start Table and a result based on that table.

_____ → _____

Define the function

Use the relevant methods (circle your helper functions!), then produce a result with the new table.

fun _____ (_____) :

t = _____

Define the table

Are there more columns?

Are there fewer rows?

Are the rows ordered?

Produce the result

end

Table Plan

Contract and Purpose

_____ :: _____ → _____

Examples

Make a Start Table and a result based on that table.

_____ → _____

Define the function

Use the relevant methods (circle your helper functions!), then produce a result with the new table.

fun _____ (_____) :

t = _____

Define the table

Are there more columns?

Are there fewer rows?

Are the rows ordered?

Produce the result

end

Table Plan

Contract and Purpose

_____ :: _____ → _____

Examples

Make a Start Table and a result based on that table.

	→	

Define the function

Use the relevant methods (circle your helper functions!), then produce a result with the new table.

fun _____ (_____) :

t =

Define the table

Are there more columns?

Are there fewer rows?

Are the rows ordered?

Produce the result

end

Contracts

Contracts tell us how to use a function. For example: `num-sqr :: (n :: Number) → Number` tells us that the name of the function is `num-sqr`, that it takes one input (a `Number`), and that it evaluates to a number. From the contract, we know `num-sqr(4)` will evaluate to a `Number`.

Name	Domain	Range
<code>triangle</code>	<code>:: (side-length :: Number, style :: String, color :: String) →</code>	<code>Image</code>
<code>circle</code>	<code>:: (radius :: Number, style :: String, color :: String) →</code>	<code>Image</code>
<code>star</code>	<code>:: (radius :: Number, style :: String, color :: String) →</code>	<code>Image</code>
<code>rectangle</code>	<code>:: (width :: Num, height :: Num, style :: Str, color :: Str) →</code>	<code>Image</code>
<code>ellipse</code>	<code>:: (width :: Num, height :: Num, style :: Str, color :: Str) →</code>	<code>Image</code>
<code>square</code>	<code>:: (size-length :: Number, style :: String, color :: String) →</code>	<code>Image</code>
<code>text</code>	<code>:: (str :: String, size :: Number, color :: String) →</code>	<code>Image</code>
<code>overlay</code>	<code>:: (img1 :: Image, img2 :: Image) →</code>	<code>Image</code>
<code>rotate</code>	<code>:: (degree :: Number, img :: Image) →</code>	<code>Image</code>
<code>scale</code>	<code>:: (factor :: Number, img :: Image) →</code>	<code>Image</code>
<code>string-repeat</code>	<code>:: (text :: String, repeat :: Number) →</code>	<code>String</code>
<code>string-contains</code>	<code>:: (text :: String, search-for :: String) →</code>	<code>Boolean</code>
<code>num-sqr</code>	<code>:: (n :: Number) →</code>	<code>Number</code>
<code>num-sqrt</code>	<code>:: (n :: Number) →</code>	<code>Number</code>
<code>num-min</code>	<code>:: (a :: Number, b :: Number) →</code>	<code>Number</code>
<code>num-max</code>	<code>:: (a :: Number, b :: Number) →</code>	<code>Number</code>

Contracts

Contracts tell us how to use a function. For example: `<Table>.filter :: (test :: (Row → Boolean) → Row` tells us that the name of the function is `.filter` and that it is a `Table` method. The domain says it one input (a function that consumes `Rows` and produces `Booleans`), and that the method evaluates to a `Table`. From the contract, we know `animals-table.filter(is-cat)` will evaluate to a `Table`.

Name	Domain		Range
<code><Table>.row-n</code>	<code>:: (n :: Number)</code>	→	<code>Row</code>
<code><Table>.order-by</code>	<code>:: (col :: String, increasing :: Boolean)</code>	→	<code>Table</code>
<code><Table>.filter</code>	<code>:: (test :: (Row → Boolean))</code>	→	<code>Table</code>
<code><Table>.build-column</code>	<code>:: (col :: String, builder :: (Row → Value))</code>	→	<code>Table</code>
<code>mean</code>	<code>:: (t :: Table, col :: String)</code>	→	<code>Number</code>
<code>median</code>	<code>:: (t :: Table, col :: String)</code>	→	<code>Number</code>
<code>modes</code>	<code>:: (t :: Table, col :: String)</code>	→	<code>List<Number></code>
<code>bar-chart</code>	<code>:: (t :: Table, labels :: String, values :: String)</code>	→	<code>Image</code>
<code>pie-chart</code>	<code>:: (t :: Table, labels :: String, values :: String)</code>	→	<code>Image</code>
<code>box-plot</code>	<code>:: (t :: Table, col :: String)</code>	→	<code>Image</code>
<code>freq-bar-chart</code>	<code>:: (t :: Table, values :: String)</code>	→	<code>Image</code>
<code>histogram</code>	<code>:: (t :: Table, values :: String, bin-width :: Number)</code>	→	<code>Image</code>
<code>scatter-plot</code>	<code>:: (t :: Table, labels :: String, xs :: String, ys :: String)</code>	→	<code>Image</code>
<code>lr-plot</code>	<code>:: (t :: Table, labels :: String, xs :: String, ys :: String)</code>	→	<code>Image</code>