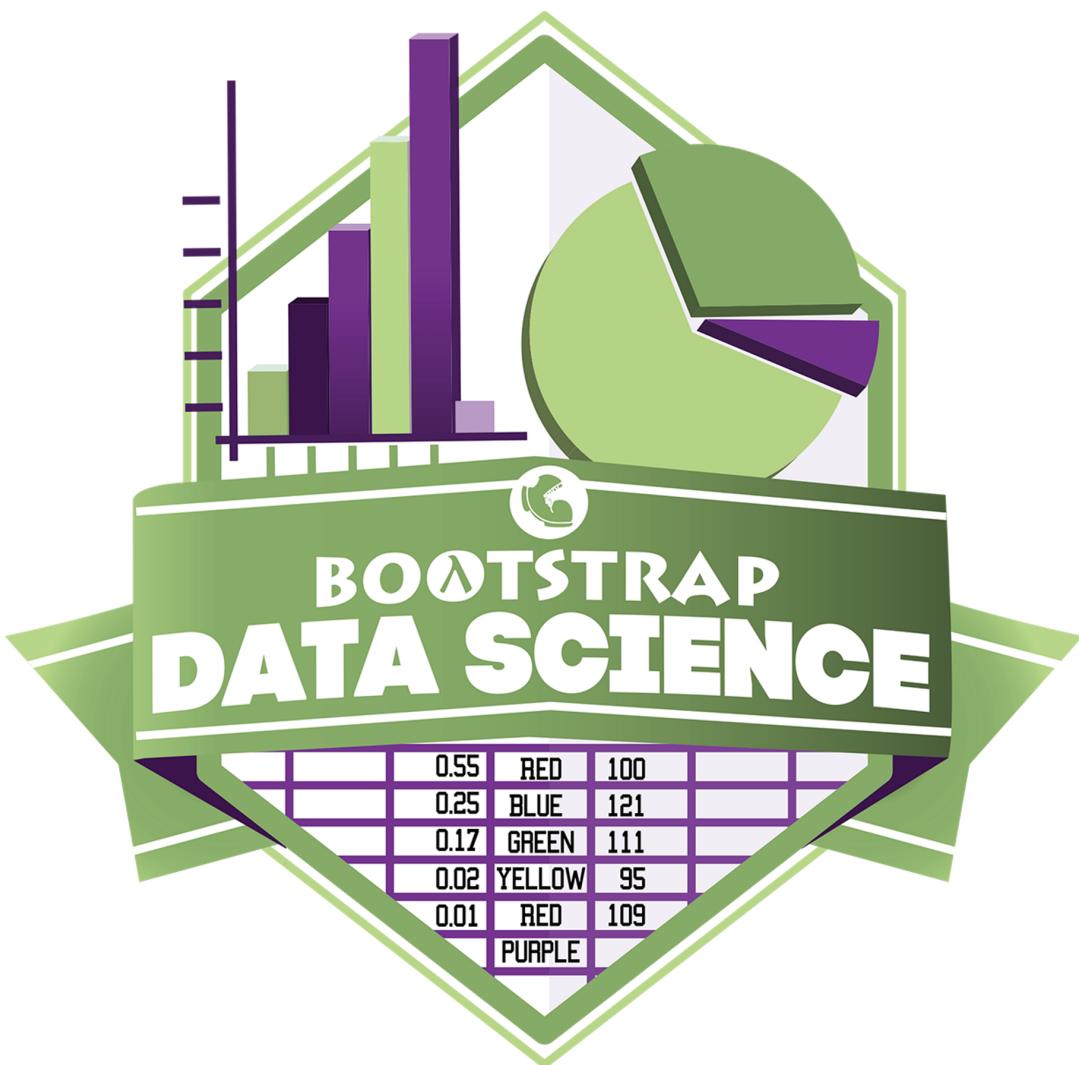


Name: _____



Student Workbook



Workbook v1.5

Brought to you by the Bootstrap team:

- Emmanuel Schanzer
- Kathi Fisler
- Shriram Krishnamurthi
- Dorai Sitaram
- Joe Politz
- Jennifer Poole
- Ed Campos
- Ben Lerner
- Nancy Pfenning
- Flannery Denny

Visual Designer: Colleen Murphy

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 contact@BootstrapWorld.org.

Introduction to Computational Data Science

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 about 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 numeric and measures *an amount*, such as a person's height, a score on a test, distance, etc. A list of quantitative data can be ordered from smallest to largest.
- **Categorical data** is data that specifies *qualities*, such as sex, eye color, country of origin, etc. Categorical data is not subject to the laws of arithmetic—for example, we cannot take the "average" of a list of colors.

Answering questions with 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 of 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 average weight of animals from the shelter?")
- **Relate Questions** require looking for trends across multiple columns. (e.g., "Do cats tend to be adopted sooner than dogs?")

The Animals Dataset

name	species	sex	age	fixed	legs	pounds	weeks
Sasha	cat	female	1	false	4	6.5	3
Snuffles	rabbit	female	3	true	4	3.5	8
Mittens	cat	female	2	true	4	7.4	1
Sunflower	cat	female	5	true	4	8.1	6
Felix	cat	male	16	true	4	9.2	5
Sheba	cat	female	7	true	4	8.4	6
Billie	snail	hermaphrodite	0.5	false	0	0.1	3
Snowcone	cat	female	2	true	4	6.5	5
Wade	cat	male	1	false	4	3.2	1
Hercules	cat	male	3	false	4	13.4	2
Toggle	dog	female	3	true	4	48	1
Boo-boo	dog	male	11	true	4	123	24
Fritz	dog	male	4	true	4	92	3
Midnight	dog	female	5	false	4	112	4
Rex	dog	male	1	false	4	28.9	9
Gir	dog	male	8	false	4	88	5
Max	dog	male	3	false	4	52.8	8
Nori	dog	female	3	true	4	35.3	1
Mr. Peanutbutter	dog	male	10	false	4	161	6
Lucky	dog	male	3	true	3	45.4	9
Kujo	dog	male	8	false	4	172	30
Buddy	lizard	male	2	false	4	0.3	3
Gila	lizard	female	3	true	4	1.2	4
Bo	dog	male	8	true	4	76.1	10
Nibblet	rabbit	male	6	false	4	4.3	2
Snuggles	tarantula	female	2	false	8	0.1	1
Daisy	dog	female	5	true	4	68	8
Ada	dog	female	2	true	4	32	3
Miaulis	cat	male	7	false	4	8.8	4
Heathcliff	cat	male	1	true	4	2.1	2
Tinkles	cat	female	1	true	4	1.7	3
Maple	dog	female	3	true	4	51.6	4

Categorical or Quantitative?

For each piece of data below, circle whether it is **Categorical** or **Quantitative** data.

1 Hair color	categorical	quantitative
2 Age	categorical	quantitative
3 ZIP Code	categorical	quantitative
4 Year	categorical	quantitative
5 Height	categorical	quantitative
6 Sex	categorical	quantitative
7 Street Name	categorical	quantitative

For each question, circle whether it will be answered by **Categorical** or **Quantitative** data.

8 We'd like to find out the average price of cars in a lot.	categorical	quantitative
9 We'd like to find out the most popular color for cars.	categorical	quantitative
10 We'd like to find out which puppy is the youngest.	categorical	quantitative
11 We'd like to find out which cats have been fixed.	categorical	quantitative
12 We want to know which people have a ZIP code of 02907.	categorical	quantitative
13 We'd like to sort a list of phone numbers by area code.	categorical	quantitative

Questions and Column Descriptions

What questions can you ask about the animals dataset? Come up with at least one **Lookup**, **Compute**, **Relate** or **Can't Answer** question, and write them as wonders below. (Note: These question types are defined on Page 1.)

What do you NOTICE about this dataset?	What do you WONDER about this dataset?	Question Type
		<i>Lookup Compute Relate Can't answer</i>

1. This dataset is _____ Animals that came from an animal shelter _____, which contains 31 _____ data rows.

2. Some of the columns are:

a. _____ species _____, which contains _____ categorical _____ data. Some example values are:
_____ "cat", "dog", and "rabbit" _____.

b. _____, which contains _____ data. Some example values are:
_____.

What's on your mind?

Introduction to Programming in Pyret

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

- Numbers are values like `1` , `0.4` , `1/3` , and `-8261.003` .
 - Numbers are *usually* used for quantitative data and other values are *usually* used as categorical data.
 - In Pyret, any decimal *must* start with a `0`. `0.22` is valid, but `.22` is not.
- Strings are values like `"Emma"` , `"Rosanna"` , `"Jen and Ed"` , or even `"08/28/1980"` .
 - In Pyret, all strings *must* be surrounded in quotation marks.
- Booleans are either `true` or `false` .

Operators (like `+` , `-` , `*` , `<` , etc.) work the same way in Pyret that they do in math.

- Operators are written between values, for example: `4 + 2` .
- In Pyret, operators must always have a space around them. `4 + 2` is valid, but `4+2` is not.
- If an expression has different operators, parentheses must be used to show order of operations. `4 + 2 + 6` and `4 + (2 * 6)` are valid, but `4 + 2 * 6` is not.

Applying Functions also works the way it does in math. The function name is first, followed by a list of arguments in parentheses.

- In math this could look like `f(5)` or `f(g(10, 4))` .
- In Pyret this could look like `star(50, "solid", "red")` .
- There are many other Pyret functions, for example `num-sqr` , `num-sqrt` `triangle` , `star` , `string-repeat` , etc.

Functions have contracts, which help explain how a function should be used. Every contract has three parts:

- The Name of the function - literally, what it's called.
- The Domain of the function - what *types of values* the function consumes, and in what order.
- The Range of the function - what *type of value* the function produces.

Value Definitions (like `x = 4` , or `y = 9 + 6`) also work the way they do in math. Every value definition starts with a name, followed by an equals sign, and then an expression. Once a value is defined, it can be referred to by name.

Numbers and Strings

Make sure you've loaded the code.pyret.org editor, 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:

Operators

6. Just like 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:

7. Type in the following expressions, one at a time: `4 + 2 + 6`, `4 + 2 * 6`, `4 + (2 * 6)`. What do you notice?
Write your answer below:

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

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.

1) <code>3 <= 4</code>	<hr/>	7) <code>"a" > "b"</code>	<hr/>
2) <code>3 == 2</code>	<hr/>	8) <code>"a" < "b"</code>	<hr/>
3) <code>2 < 4</code>	<hr/>	9) <code>"a" == "b"</code>	<hr/>
4) <code>3 <> 3</code>	<hr/>	10) <code>"a" <> "b"</code>	<hr/>
5) <code>5 >= 5</code>	<hr/>	11) <code>"a" <> "a"</code>	<hr/>
6) <code>4 >= 6</code>	<hr/>	12) <code>"a" == "a"</code>	<hr/>

13) In your own words, describe what `>` does.

14) In your own words, describe what `<=` does.

15) In your own words, describe what `<>` does.

16) How many **Numbers** are there in the entire universe?

17) How many **Strings** are there in the entire universe?

18) How many **Images** are there in the entire universe?

19) How many **Booleans** are there in the entire universe?

Applying Functions

Type this line of code into the interactions area and hit "Enter": `triangle(50, "solid", "red")`

1)	What is the name of this function?	_____
2)	What did the expression evaluate to?	_____
3)	How many arguments does <code>triangle</code> expect?	_____
4)	What does the <code>triangle</code> function produce? (Numbers? Strings? Booleans?)	_____

Catching Bugs

The following lines of code are all BUGGY! Can you spot the mistake? If you have time, type in the buggy code and see if Pyret agrees with you!

5) `triangle(20, "solid" "red")`

Can you spot the mistake?

What error message does Pyret return?

6) `triangle(20, "solid")`

Can you spot the mistake?

What error message does Pyret return?

7) `triangle(20, 10, "solid", "red")`

Can you spot the mistake?

What error message does Pyret return?

8) `triangle (20, "solid", "red")`

Can you spot the mistake?

What error message does Pyret return?

9) `triangle 20, "solid", "red")`

Can you spot the mistake?

What error message does Pyret return?

Contracts

Consider the following contract:

```
rotate :: (degree :: Number, img :: Image) -> Image
```

What is the **Name** of this function? _____

How many things are in this function's **Domain**? _____

What is the **type** of this function's **first argument**? _____

What is the **name** of this function's **second argument**? _____

What is the **Range** of this function? _____

Circle the expression below that is the correct application of this function, based on its contract.

1. `rotate(45, 90)`
2. `rotate(circle(99, "solid", "green"))`
3. `rotate(25, rectangle(7, 10, "outline", "black"))`
4. `rotate(rectangle(7, 10, "outline", "black"), 25)`

Matching Expressions and Contracts

Match the contract (left) with the expression described by the function being used (right).

Contract	Expression
make-id :: (name :: String, age :: Number) -> Image 1	A make-id("Hannah", "Smith")
phone-bill :: (minutes :: Number, texts :: Number) -> Number 2	B make-id("George", 17)
phone-bill :: (minutes :: Number) -> Number 3	C phone-bill(31, 287)
make-id :: (first :: String, last :: String) -> Image 4	D make-id("Jessica", "Jones", 32)
make-id :: (first :: String, last :: String, age :: Number) -> Image 5	E phone-bill(55)

What's on your mind?

Plotting and Displaying Data

Data Scientists use **displays** to visualize data. You've probably seen some of these charts, graphs and plots yourselves! When it comes to displaying **Categorical Data**, there are two displays that are especially useful.

1. **Bar charts** show the *count or percentage* of rows in each category.

- Bar charts provide a visual representation of the frequency of values in a categorical column.
- Bar charts have a bar for every category in a column.
- The more rows in a category, the taller the bar.
- Bars in a bar chart can be shown in *any order*, without changing the meaning of the chart. However, bars are usually shown in some sensible order (bars for the number of orders for different t-shirt sizes might be presented in order of smallest to largest shirt).

2. **Pie charts** show the *percentage* of rows in each category.

- Pie charts provide a visual representation of the relative frequency of values in a categorical column.
- Pie charts have a slice for every category in a column.
- The more rows in a category, the larger the slice.
- Slices in a pie chart can be shown in *any order*, without changing the meaning of the chart. However, slices are usually shown in some sensible order (e.g. slices might be shown in alphabetical order or from the smallest to largest slice).

Exploring Displays

Using your Contracts page and the Animals Starter File, make each type of display below in pyret. Then sketch the displays and answer the questions. Be sure to add examples of the code you use to your contracts page!

Pie Charts	Bar Charts
Sketch a pie chart here.	Sketch a bar chart here.
Pie charts are constructed from _____ 1 column(s).	Bar charts are constructed from _____ column(s).
They show _____ categorical data.	They show _____ data.
What does this display tell us? _____ _____	What does this display tell us? _____ _____
Box Plots	Histograms
Sketch a box plot here.	Sketch a histogram here.
Box plots are constructed from _____ column(s).	Histograms are constructed from _____ column(s).
They show _____ data.	They show _____ data.
What does this display tell us? _____ _____	What does this display tell us? _____ _____

(More) Exploring Displays

For each type of display, fill in the information below.

Scatter Plots	Linear Regression Plots
Sketch a box plot here.	Sketch a histogram here.
Box plots are constructed from _____ column(s).	Histograms are constructed from _____ column(s).
They show _____ data.	They show _____ data.
What does this display tell us? _____ _____	What does this display tell us? _____ _____

What's on your mind?

Data Displays and Lookups

Data scientists use data visualizations to gain better insights into their data, and to communicate their findings with others.

Making a display requires answering three questions:

1. **What data** is being displayed? This could be "a random sample of 2000 people", "every animal from the shelter", or "students' aged 14-17".
2. **What variables** are being explored? Are we looking at the `species` column? The number of kilograms that an animal weighs? Searching for a relationship between a person's income and their height ?
3. **What display** is being used, given the variables being explored? If it's a quantitative variable, we might use a histogram or box plot. If it's categorical, we could use a pie or bar chart. If it's two quantitative variables, we probably want a scatter plot.

When **looking up a data Row** from a Table, programmers use the `row-n` method. This method takes a single number as its input, which tells the computer which Row we want. *Note: Rows are numbered starting at zero!*

For example:

```
animals-table.row-n(0) # access the 1st data row  
animals-table.row-n(16) # access the 17th data row
```

When **looking up a column** from a Row, programmers use square brackets and the name of the column they want.

For example:

```
animals-table.row-n(11) ["age"]      # look up the age of the animal in the 12st data row  
animals-table.row-n(14) ["species"]   # look up the species of the animal in the 15th data row
```

Throughout the rest of the workbook, we will sometimes refer to `animalA` and `animalB`.

```
animalA = animals-table.row-n(4)  
animalB = animals-table.row-n(13)
```

What Display Goes with Which Data?

Match the Display with the description of the data being plotted. Some descriptions may go with more than one display!

Pie Charts 1

A 1 column of Quantitative Data

Bar Charts 2

Histograms 3

B 2 columns of Quantitative Data

Box Plots 4

Scatter Plots 5

C 1 column of Categorical Data

Data Displays

Fill in the tables below, then write the Pyret code that will make that display. The first column has been filled in for you.

- 1) A pie-chart showing the species of animals from the shelter.

Which Rows?	Which Column(s)?	What Display?
All the animals		

code: _____

- 2) A bar-chart showing the sex of animals from the shelter.

Which Rows?	Which Column(s)?	What Display?
All the animals		

code: _____

- 3) A histogram of the number of pounds that animals weigh.

Which Rows?	Which Column(s)?	What Display?
All the animals		

code: _____

- 4) A box-plot of the number of pounds that animals weigh.

Which Rows?	Which Column(s)?	What Display?
All the animals		

code: _____

- 5) A scatter-plot , using the animals' species as the labels, age as the x-axis, and pounds as the y-axis.

Which Rows?	Which Column(s)?	What Display?
All the animals		

code: _____

- 6) A scatterplot , using the animals' name as the labels, pounds as the x-axis, and weeks as the y-axis.

Which Rows?	Which Column(s)?	What Display?
All the animals		

code: _____

Lookup Questions

The table below represents four pets:

pets-table

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

1) Match each Lookup Question (left) to the code that will give the answer (right).

- | | | | |
|---------------------------------------|---|---|--------------------------------|
| "How much does Maple weigh?" | 1 | A | pets-table.row-n(3) |
| "Which is the last row in the table?" | 2 | B | pets-table.row-n(2) ["name"] |
| "What is Fritz's sex?" | 3 | C | pets-table.row-n(1) ["sex"] |
| "What's the third animal's name?" | 4 | D | pets-table.row-n(3) ["age"] |
| "How much does Nori weigh?" | 5 | E | pets-table.row-n(3) ["pounds"] |
| "How old is Maple?" | 6 | F | pets-table.row-n(0) |
| "What is Toggle's sex?" | 7 | G | pets-table.row-n(2) ["pounds"] |
| "What is the first row in the table?" | 8 | H | pets-table.row-n(0) ["sex"] |

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

a.	<code>pets-table.row-n(3)[\"name\"]</code>	"Maple"
b.	<code></code>	"male"
c.	<code></code>	4
d.	<code></code>	48
e.	<code></code>	"Nori"

What's on your mind?

Defining Functions

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

- We use the Design Recipe to help us define functions **and think through problems clearly**.
- 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 Design Recipe

Directions: Define a function called `gt`, which makes solid green triangles of whatever size we want.

Contract and Purpose Statement

Every contract has three parts...

#	gt::	(size :: Number)	->	Image
	function name	domain		range

Consumes a size, and produces a solid green triangle of that size.

what does the function do?

Examples

Write some examples, then circle and label what changes...

examples:

_____	(_____)	is _____	what the function produces
function name	input(s)		

_____	(_____)	is _____	what the function produces
function name	input(s)		

end

Definition

Write the definition, giving variable names to all your input values...

```
fun      gt(   size   ):  
        _____  
        variable(s)  
        triangle(size, "solid", "green")  
        _____  
        what the function does with those variable(s)  
end
```

Directions: Define a function called `bc`, which makes solid blue circles of whatever radius we want.

Contract and Purpose Statement

Every contract has three parts...

#	::	->	
	function name	domain	range
#			<i>what does the function do?</i>

Examples

Write some examples, then circle and label what changes...

examples:

_____	(_____)	is _____	what the function produces
function name	input(s)		

_____	(_____)	is _____	what the function produces
function name	input(s)		

end

Definition

Write the definition, giving variable names to all your input values...

```
fun      ( _____ ):  
        _____  
        variable(s)  
        _____  
        what the function does with those variable(s)  
end
```

The Design Recipe

Directions: Define a function called `sticker`, which draws 50px stars in whatever color is input.

Contract and Purpose Statement

Every contract has three parts...

_____ :: _____ -> _____

_____ what does the function do?

Examples

Write some examples, then circle and label what changes...

examples:

_____ (_____) is _____
function name input(s) what the function produces
_____ (_____) is _____
function name input(s) what the function produces

end

Definition

Write the definition, giving variable names to all your input values...

fun _____ (_____):
function name variable(s)

_____ what the function does with those variable(s)

end

Directions: Define a function called `nametag`, which consumes a `Row` of the animals table and draws their name in purple, 10px letters. (Assume you have rows `animalA` and `animalB` defined.)

Contract and Purpose Statement

Every contract has three parts...

nametag:: (r :: Row) -> Image
function name domain range

Consumes an animal, and produces that animal's name in purple, 10px letters.

what does the function do?

Examples

Write some examples, then circle and label what changes...

examples:

nametag ("animalA") is _____
function name input(s) what the function produces
(_____) is _____
function name input(s) what the function produces

end

Definition

Write the definition, giving variable names to all your input values...

fun nametag(r):
function name variable(s)

text(r["name"], 10, "purple")

what the function does with those variable(s)

end

What's on your mind?

Defining Row Functions & Using Table Methods

Methods are special functions that are attached to pieces of data. We use them to manipulate Tables.

- In this course, the methods we'll be using are
 - `row-n` - consumes an index (starting with zero!) and produces a row from a table
 - `order-by` - consumes the name of a column and a Boolean value to determine if that table should be sorted by that column in ascending order
 - `filter` - consumes a *Boolean-producing function*, and produces a table containing only rows for which the function returns `true`
 - `build-column` - consumes the name of a new column, and a function that produces the values in that column for each Row
- Unlike functions, methods can't be used alone. They have a "secret" argument, which is the data they are attached to. They are written as part of that data, separated by a dot. For example:

```
shapes.row-n(2)
```

- Contracts for methods are different from other functions. They include the type of the data as part of their names. For example:

```
<table>.row-n :: (index :: Number) -> Row
```

Reading Function Definitions

Make sure you have the "Table Methods Starter File" open on your computer, and click "Run".

1	How many functions are defined here?	
2	What are their names?	
3	What is the domain of <code>is-dog</code> ?	
4	What is the range of <code>is-old</code> ?	
5	What is the range of <code>lookup-name</code> ?	
6	What does <code>is-fixed(animalA)</code> evaluate to?	
7	What does <code>lookup-name(animalB)</code> evaluate to?	
8	What does <code>is-old(animalA)</code> evaluate to?	
9	What does <code>is-dog(animalA)</code> evaluate to?	
10	What does <code>is-fixed</code> do?	
11	What does <code>lookup-name</code> do?	
12	What does <code>is-old</code> do?	

The Design Recipe

For the word problems below, assume `animalA` and `animalB` are defined as the data rows for Felix and Midnight, respectively.

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

Contract and Purpose Statement

Every contract has three parts...

`lookup-fixed::` _____ (`r :: Row`) _____ -> Boolean
function name domain range

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

what does the function do?

Examples

Write some examples, then circle and label what changes...

examples:

_____ (_____) is _____
function name input(s) what the function produces
_____ (_____) is _____
function name input(s) what the function produces

end

Definition

Write the definition, giving variable names to all your input values...

fun `lookup-fixed(` _____ `r` _____ `):`

function name

variable(s)

what the function does with those variable(s)

end

Directions: Define a function called `lookup-sex`, which consumes a Row of the animals table and looks up the sex of that animal.

Contract and Purpose Statement

Every contract has three parts...

_____ :: _____ -> _____
function name domain range

what does the function do?

Examples

Write some examples, then circle and label what changes...

examples:

_____ (_____) is _____
function name input(s) what the function produces
_____ (_____) is _____
function name input(s) what the function produces

end

Definition

Write the definition, giving variable names to all your input values...

fun _____ (_____) `:`

function name

variable(s)

what the function does with those variable(s)

end

The Design Recipe

For the word problems below, assume `animalA` and `animalB` are defined as the data rows for Felix and Midnight, respectively.

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

Contract and Purpose Statement

Every contract has three parts...

is-cat:: (r :: Row) -> Boolean
function name domain range

Consumes an animal, and computes whether the species == "cat"

what does the function do?

Examples

Write some examples, then circle and label what changes...

examples:

is-cat ("animalA") is _____
function name input(s) what the function produces
(_____) is _____
function name input(s) what the function produces

end

Definition

Write the definition, giving variable names to all your input values...

fun is-cat(r):
function name variable(s)
r["species"] == "cat"
what the function does with those variable(s)
end

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

Contract and Purpose Statement

Every contract has three parts...

:: -> range
function name domain range

what does the function do?

Examples

Write some examples, then circle and label what changes...

examples:

(_____) is _____
function name input(s) what the function produces
(_____) is _____
function name input(s) what the function produces

end

Definition

Write the definition, giving variable names to all your input values...

fun ():
function name variable(s)
what the function does with those variable(s)
end

What's on your mind?

Method Chaining

Method chaining allows us to apply multiple methods with less code.

For example, instead of using multiple definitions, like this:

```
with-labels = animals-table.build-column("labels", nametag)
cats = with-labels.filter(is-cat)
cats.order-by("age", true)
```

We can use method-chaining to write it all on one line, like this:

```
animals-table.build-column("labels", nametag).filter(is-cat).order-by("age", true)
```

Order Matters! The methods are applied in the order they appear. For example, trying to order a table by a column that hasn't been built will result in an error.

The Design Recipe

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

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

Contract and Purpose Statement

Every contract has three parts...

is-dog:: _____ (r :: Row) -> Boolean
function name domain range
Consumes an animal, and computes whether the species == "dog"
what does the function do?

Examples

Write some examples, then circle and label what changes...

examples:

is-dog ("animalA") is animalA["species"] == "dog"
function name input(s) what the function produces
is-dog ("animalB") is _____
function name input(s) what the function produces

end

Definition

Write the definition, giving variable names to all your input values...

fun is-dog(r):
function name variable(s)
r["species"] == "dog"
what the function does with those variable(s)

end

Directions: Define a function called `is-female`, which consumes a Row of the animals table and returns true if the animal is female.

Contract and Purpose Statement

Every contract has three parts...

:: _____ -> _____
function name domain range

what does the function do?

Examples

Write some examples, then circle and label what changes...

examples:

(_____) is _____ what the function produces
function name input(s)
(_____) is _____ what the function produces
function name input(s)

end

Definition

Write the definition, giving variable names to all your input values...

fun ():
function name variable(s)
what the function does with those variable(s)

end

The Design Recipe

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

Directions: Define a function called `is-old`, which consumes a Row of the animals table and computes whether it is more than 12 years old.

Contract and Purpose Statement

Every contract has three parts...

_____ :: _____ -> _____
function name domain range

what does the function do?

Examples

Write some examples, then circle and label what changes...

examples:

_____ (_____) is _____
function name input(s) what the function produces
_____ (_____) is _____
function name input(s) what the function produces

end

Definition

Write the definition, giving variable names to all your input values...

fun _____ (_____):
function name variable(s)

what the function does with those variable(s)

end

Directions: Define a function called `name-has-s`, which returns true if an animal's name contains the letter "s"

Contract and Purpose Statement

Every contract has three parts...

name-has-s:: _____ -> _____
function name domain range

what does the function do?

Examples

Write some examples, then circle and label what changes...

examples:

_____ (_____) is _____
function name input(s) what the function produces
_____ (_____) is _____
function name input(s) what the function produces

end

Definition

Write the definition, giving variable names to all your input values...

fun name-has-s(_____ r):
function name variable(s)

string-contains(r["name"], "s")

what the function does with those variable(s)

end

Chaining Methods

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

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

The table `t` below represents four animals from the shelter:

name	sex	age	fixed	pounds
"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).

- | | | | |
|---|---|---|--|
| <code>t.order-by("age", true)</code> | 1 | A | Produces a table containing only Toggle and Maple |
| <code>t.filter(is-fixed)</code> | 2 | B | Produces a table of only young, fixed animals |
| <code>t.build-column("sticker", nametag)</code> | 3 | C | Produces a table, sorted youngest-to-oldest |
| <code>t.filter(is-young)</code> | 4 | D | Produces a table with an extra column, named "sticker" |
| <code>t.filter(is-young).filter(is-fixed)</code> | 5 | E | Produces a table containing Maple and Toggle, in that order |
| <code>t.filter(is-young).order-by("pounds", false)</code> | 6 | F | Produces a table containing the same four animals |
| <code>t.build-column("label", nametag).order-by("age", true)</code> | 7 | G | Won't run: will produce an error |
| <code>t.order-by("agee", false)</code> | 8 | H | Produces a table with an extra "label" column, sorted youngest-to-oldest |

Chaining Methods 2: Order Matters!

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

```
fun is-female(r): r["sex"] == "female"    end
fun kilograms(r): r["pounds"] / 2.2        end
fun is-heavy(r):  r["kilos"] > 25          end
```

The table `t` below represents four animals from the shelter:

name	sex	age	fixed	pounds
"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). Note: one description might match multiple expressions!

`t.order-by("kilos", true)`

1

A Produces a table containing Toggle, Nori and Maple, with an extra column showing their weight in kilograms

`t.filter(is-female)
.build-column("kilos", kilograms)`

2

B Produces a table containing Maple, Nori and Toggle (in that order)

`t.build-column("kilos", kilograms)
.filter(is-heavy)`

3

C Produces a table containing only Fritz, with a single extra column called kilos

`t.filter(is-heavy)
.build-column("kilos", kilograms)`

4

D Won't run: will produce an error

`t.build-column("kilos", kilograms)
.filter(is-heavy)
.order-by("sex", true)`

5

E Produces a table containing only Fritz, with two extra columns

`t.build-column("female", is-female)
.build-column("kilos", kilograms)
.filter(is-heavy)`

6

F Produces a table containing Maple and Fritz

What's on your mind?

Mood Generator

1) Open the Mood Generator starter file, and read through the code you find there. This code contains new programming that you haven't seen yet! Take a moment to list everything you Notice, and then everything you Wonder...

Notice	Wonder

2) Add another line of code to the definition, so that `mood("mad")` produces the *same emoji* as `mood("angry")`.

3) Add **another example** to the `examples:` section for "laughing", using the appropriate emoji. (To bring up the emojis on your computer, type `Cmd-Ctrl-Space` on a Mac, or `Windows-Period` on Windows 10)

4) Come up with some new moods, and add them to the code. Make sure you include `examples: !`

5) In your own words, how do if-expressions work in Pyret? Write your answer below.

6) Write down at least 2 ways you could use if-expressions when analyzing the Animals Dataset.

Word Problem: species-color

Directions: We want to generate a custom dot for our `image-scatter-plot`, such that every species gets a unique color. Write a function called `species-color`, which takes in a Row from the animals table and returns a solid, 5px circle using a color you've chosen.

Contract and Purpose Statement

Every contract has three parts...

_____ :: _____ -> _____
_____ _____ _____
_____ _____ _____
_____ _____ _____

function name *domain* *range*
what does the function do?

Examples

Write some examples, then circle and label what changes...

examples:

_____ (_____) is _____
function name input(s) what the function produces
_____ (_____) is _____
function name input(s) what the function produces
_____ (_____) is _____
function name input(s) what the function produces
_____ (_____) is _____
function name input(s) what the function produces
_____ (_____) is _____
function name input(s) what the function produces
end

Definition

Write the definition, giving variable names to all your input values...

fun _____ (_____):
function name variable(s)

what the function does with those variable(s)

end

end

Randomness and Sample Size

Computer Scientists may take **samples** that are subsets of a data set. If their sample is well chosen, they can use it to test if their code does what it's supposed to do. However, choosing a good sample can be tricky!

Random Samples are a subset of a population in which each member of the subset has an equal chance of being chosen. A random sample is intended to be a representative subset of the population. The larger the random sample, the more closely it will represent the population and the better our inferences about the population will tend to be.

Grouped Samples are a subset of a population in which each member of the subset was chosen for a specific reason. For example, we might want to look at the difference in trends between two groups ("Is the age of a dog a bigger factor in adoption time v. the age of a cat?"). This would require making grouped samples of *just the dogs* and *just the cats*.

Sampling and Inference

1) Evaluate the `big-animals-table` in the Interactions Area. This is the *complete* population of animals from the shelter! Below is a true statement about that population:

The population is 47.7% fixed and 52.3% unfixed.

2) How close to these percentages do we get with random samples?

Type each of the following lines into the Interactions Area and hit "Enter".

```
random-rows(big-animals-table, 10)  
random-rows(big-animals-table, 40)
```

3) What do you get?

4) What is the contract for `random-rows` ? _____

5) What does the `random-rows` function do?

6) In the Definitions Area, define `tiny-sample` and `small-sample` to be these two random samples.

7) Make a `pie-chart` for the animals in each sample, showing percentages of fixed and unfixed.

- The percentage of fixed animals in the entire populations is 47.7%.
- The percentage of fixed animals in `tiny-sample` is _____.
- The percentage of fixed animals in `small-sample` is _____.

8) Make a `pie-chart` for the animals in each sample, showing percentages for each species.

- The percentage of tarantulas in the entire population is roughly 5%.
- The percentage of tarantulas in `tiny-sample` is _____.
- The percentage of tarantulas in `small-sample` is _____.

9) Click "Run" to direct the computer to generate a different set of random samples of these sizes. Make a new `pie-chart` for each sample, showing percentages for each species.

- The percentage of tarantulas in the entire population is roughly 5%.
- The percentage of tarantulas in `tiny-sample` is _____.
- The percentage of tarantulas in `small-sample` is _____.

10) Which repeated sample gave us a more accurate inference about the whole population? Why?

Grouped Samples from the Animals Dataset

Use method chaining to define the **grouped samples** below, using the helper functions that you've already defined: `is-old`, `is-young`, `is-cat`, `is-dog`, `is-female`, `lookup-fixed`, and `has-s-name`. We've given you the solution for the first sample, to get you started.

Subset	The code to define that subset
Kittens	<code>kittens = animals-table.filter(is-cat) .filter(is-young)</code>
Puppies	<code>young-dogs = animals-table.</code> _____
Fixed Cats	<code>fixed-cats = animals-table.</code> _____
Cats with "s" in their name	<code>s-cats = animals-table.</code> _____
Old Dogs	<code>old = animals-table.</code> _____
Fixed Animals	<code>fixed = animals-table.</code> _____
Old Female Cats	<code>old-cats = animals-table.</code> _____
Fixed Kittens	<code>young-fixed-cats = animals-table.</code> _____
Fixed Female Dogs	<code>fixed-female-dogs = animals-table.</code> _____
Old Fixed Female Cats	<code>old-fixed-female-cats = animals-table.</code> _____

Displaying Data

Fill in the tables below, then use Pyret to make the following displays. Record the code you used.

The first table has been filled in for you.

1) A bar-chart showing how many puppies are fixed or not.

What Rows?	Which Column(s)?	What Display?
puppies	fixed	bar-chart

code: _____

2) A pie-chart showing how many heavy dogs are fixed or not.

What Rows?	Which Column(s)?	What Display?

code: _____

3) A histogram of the number of weeks it takes for a random sample of animals to be adopted.

What Rows?	Which Column(s)?	What Display?

code: _____

4) A box-plot of the number of pounds that kittens weigh.

What Rows?	Which Column(s)?	What Display?

code: _____

5) A scatter-plot of a random sample using name as the labels, age as the x-axis, and weeks as the y-axis.

What Rows?	Which Column(s)?	What Display?

code: _____

6) A scatter-plot of fixed cats, using species as the labels, pounds as the x-axis, and weeks as the y-axis.

What Rows?	Which Column(s)?	What Display?

code: _____

What's on your mind?

Choosing Your Dataset

When selecting a dataset to explore, *pick something that matters to you!* You'll be working with this data for a while, so you don't want to pick something at random just to get it done.

When choosing a dataset, it's a good idea to consider a few factors:

1. Is it **interesting**? This should be data you are curious about, that answers questions you'd want to ask. Pick a dataset you're genuinely interested in, so that you can explore questions that matter to you!
2. Is it **relevant**? Does this data impact you in any way? Are there questions you have about the dataset that mean something to you or someone you know? Pick a dataset that deals with something personally relevant to you!
3. Is it **familiar**? You wouldn't be able to make samples of the Animals Dataset properly if you didn't know that some animals are much bigger or longer-lived than others. Pick a dataset you know about, so you can use your expertise to deepen your analysis!

My Dataset

I chose to work with the _____ dataset, which contains _____ data rows.

What do you NOTICE?	What do you WONDER?	Question Type
		<i>Lookup</i> <i>Compute</i> <i>Relate</i> <i>Can't answer</i>

Some of the columns are:

1) _____, which contains _____ data. Some example values from this column are:

_____.

2) _____, which contains _____ data. Some example values from this column are:

_____.

Samples from My Dataset

How can we define grouped samples? For a given row `r`, what function will identify if that row is in the sample?

Subset	A function that returns true if a row <code>r</code> is in the subset
	fun _____ (<code>r</code>) : end
	fun _____ (<code>r</code>) : end
	fun _____ (<code>r</code>) : end
	fun _____ (<code>r</code>) : end
	fun _____ (<code>r</code>) : end

The Design Recipe

Write helper functions for **your** dataset, which you can use to define subsets.

Directions : Define a function called _____, which consumes a Row of the _____ table and produces _____.

Contract and Purpose Statement

Every contract has three parts...

_____ :: _____ Row _____ -> Boolean
function name domain range

_____ what does the function do?

Examples

Write some examples, then circle and label what changes...

examples :

_____ (_____) is _____ what the function produces
function name input(s)
_____ (_____) is _____ what the function produces
function name input(s)

end

Definition

Write the definition, giving variable names to all your input values...

fun _____ (_____):
function name variable(s)

what the function does with those variable(s)

end

Directions : Define a function called _____, which consumes a Row of the _____ table and produces _____.

Contract and Purpose Statement

Every contract has three parts...

_____ :: _____ Row _____ -> Boolean
function name domain range

_____ what does the function do?

Examples

Write some examples, then circle and label what changes...

examples :

_____ (_____) is _____ what the function produces
function name input(s)
_____ (_____) is _____ what the function produces
function name input(s)

end

Definition

Write the definition, giving variable names to all your input values...

fun _____ (_____):
function name variable(s)

what the function does with those variable(s)

end

The Design Recipe

Write your own word problems below, and solve them using the Design Recipe.

Directions : Define a function called _____, which consumes a Row of the table and produces _____.

Contract and Purpose Statement

Every contract has three parts...

_____ :: _____ Row -> Boolean
function name domain range

what does the function do?

Examples

Write some examples, then circle and label what changes...

examples :

_____ (_____) is _____
function name input(s) what the function produces
_____ (_____) is _____
function name input(s) what the function produces

end

Definition

Write the definition, giving variable names to all your input values...

fun _____ (_____):
function name variable(s)

what the function does with those variable(s)

end

Directions : Define a function called _____, which consumes a Row of the table and produces _____.

Contract and Purpose Statement

Every contract has three parts...

_____ :: _____ Row -> Boolean
function name domain range

what does the function do?

Examples

Write some examples, then circle and label what changes...

examples :

_____ (_____) is _____
function name input(s) what the function produces
_____ (_____) is _____
function name input(s) what the function produces

end

Definition

Write the definition, giving variable names to all your input values...

fun _____ (_____):
function name variable(s)

what the function does with those variable(s)

end

What's on your mind?

Histograms

To best understand histograms, it's helpful to contrast them first with bar charts.

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

- *Bar charts provide a visual representation of the frequency of values in a categorical column.*
- There's no strict numerical way to order these bars, but **sometimes there's an order** that makes sense. For example, bars for the sales of different t-shirt sizes might be presented in order of smallest to largest shirt.

Histograms show the number of rows that fall within certain intervals, or "bins", on a horizontal axis. The more rows that fall within a particular "bin", the taller the bar.

- *Histograms provide a visual representation of the frequencies (or relative frequencies) of values in a quantitative column.*
- Quantitative data **can always be ordered**, so the bars of a histogram always progress from smallest (on the left) to largest (on the right).
- 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 shape of the dataset. Choosing a good bin size can take some trial and error!

The **shape** of a data set tells us which values are more or less common.

- In a **symmetric** data set, values are just as likely to occur a certain distance above the mean as below the mean.
- A data set that is **skewed left** and/or has low outliers has a few values that are unusually low. The histogram for a skewed left dataset has a few data points that are stretched out to the left (lower) end of the x-axis.
- A data set that is **skewed right** and/or high outliers means there are a few values that are unusually high. The histogram for a skewed right dataset has a few data points that are stretched out to the right (higher) end of the x-axis.
- One way to visualize the difference between a histogram of data that is **skewed left** or **skewed right** is to think about the lengths of our toes on our left and right feet. Much like a histogram that is "skewed left", our left feet have smaller toes on the left and a bigger toe on the right. Our right feet have the big toe on the left and smaller toes on the right, more closely resembling the shape of a histogram of "skewed right" data.

The Design Recipe

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

Directions: Define a function called `kilos`, which consumes a Row of the animals table and divides the pounds column by 2.2 to compute the animal's weight in kilograms.

Contract and Purpose Statement

Every contract has three parts...

_____ :: _____ (r :: Row) -> _____
function name domain range

what does the function do?

Examples

Write some examples, then circle and label what changes...

examples:

(_____) is _____
function name input(s) what the function produces
(_____) is _____
function name input(s) what the function produces

end

Definition

Write the definition, giving variable names to all your input values...

fun _____ (_____):
function name variable(s)

what the function does with those variable(s)

end

Directions: Define a function called `smart-dot`, which consumes a Row of the animals table and computes the image of a solid red circle using the animal's `pounds` as the radius.

Contract and Purpose Statement

Every contract has three parts...

_____ :: _____ -> Image
function name domain range
Consumes an animal, and computes a solid red circle using the weight in pounds as the radius
what does the function do?

Examples

Write some examples, then circle and label what changes...

examples:

smart-dot ("animalA") is _____
function name input(s) what the function produces
(_____) is _____
function name input(s) what the function produces

end

Definition

Write the definition, giving variable names to all your input values...

fun _____ (_____):
function name variable(s)

what the function does with those variable(s)

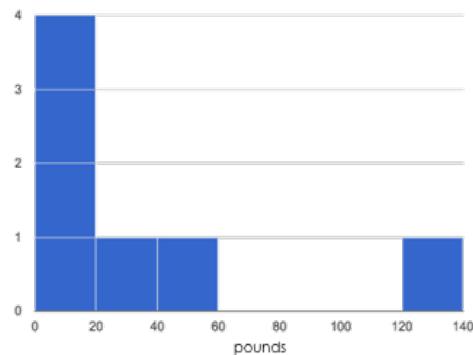
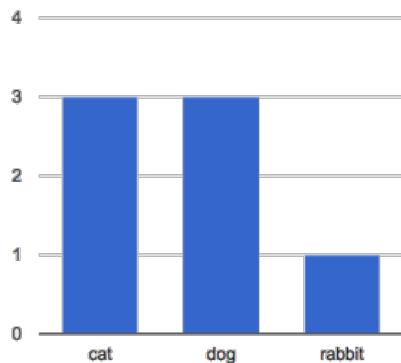
end

Summarizing Columns

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

1	How many cats are there in the table above?	
2	How many dogs are there?	
3	How many animals weigh between 0-20 pounds?	
4	How many animals weigh between 20-40 pounds?	
5	Are there more animals weighing 40-60 than 60-140 pounds?	

The charts below are both based on this table. What is similar about them? What is different?



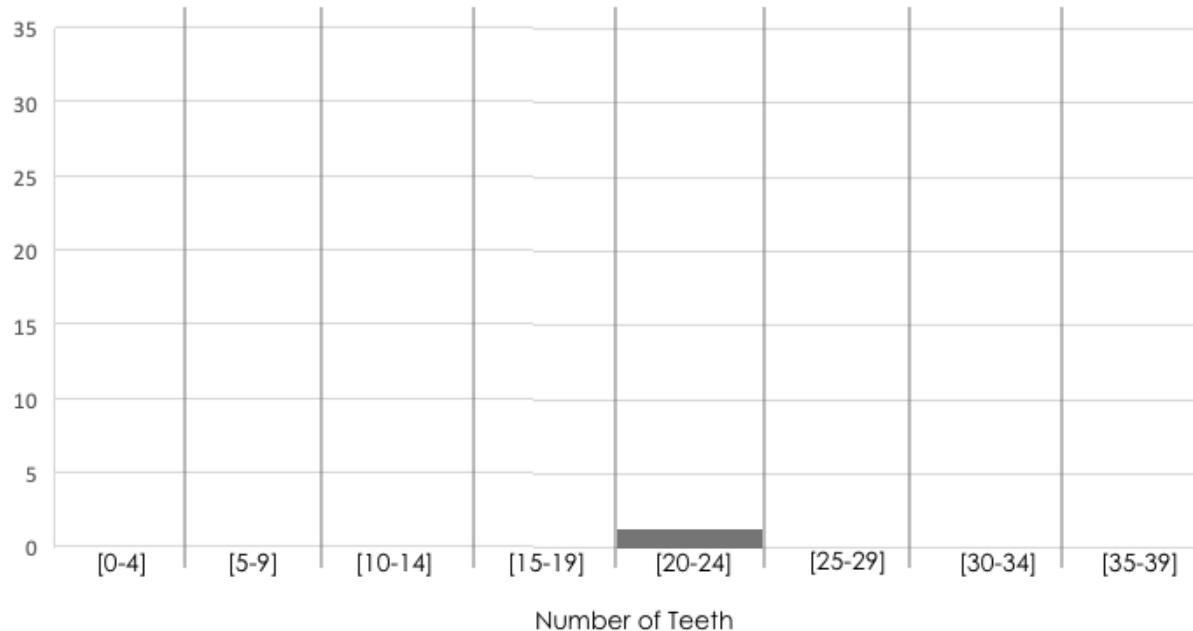
Similarities	Differences

Making Histograms

Suppose we have a data set for a group of 50 adults, showing the number of teeth each person has:

Number of teeth	Count
0	5
22	1
26	1
27	1
28	4
29	3
30	5
31	3
32	27

Draw a histogram for the table in the space below. For each row, find which interval (or "bin") on the x-axis represents the right number of teeth. Then fill in the box so that the height of the box is equal to the *sum of the counts* that fit into that interval. One of the intervals has been completed for you.



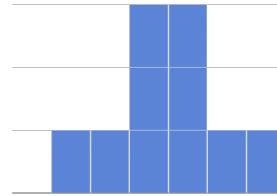
Reading Histograms

Students watched 5 videos, and rated them on a scale of 1 to 10. While the **average score** for every video is the same (5.5), the **shapes** of the ratings distributions were very different! Match the summary description (left) with the **shape** of the histogram of student ratings (right). For each histogram, the **x-axis is the score**, and the **y-axis is the number of students who gave it that score**. These axes are intentionally unlabeled - focusing on the **shape** is what matters here!

Most of the students were fine with the video, but a couple of them gave it an unusually low rating.

1

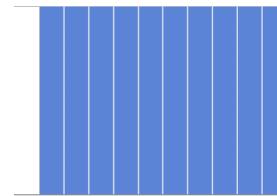
A



Most of the students were okay with the video, but a couple students gave it an unusually high rating.

2

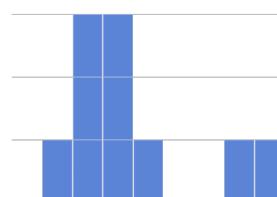
B



Students tended to give the video an average rating, and they weren't likely to stray far from the average.

3

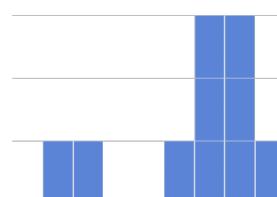
C



Students either really liked or really disliked the video.

4

D



Reactions to the video were all over the place: high ratings and low ratings and inbetween ratings were all equally likely.

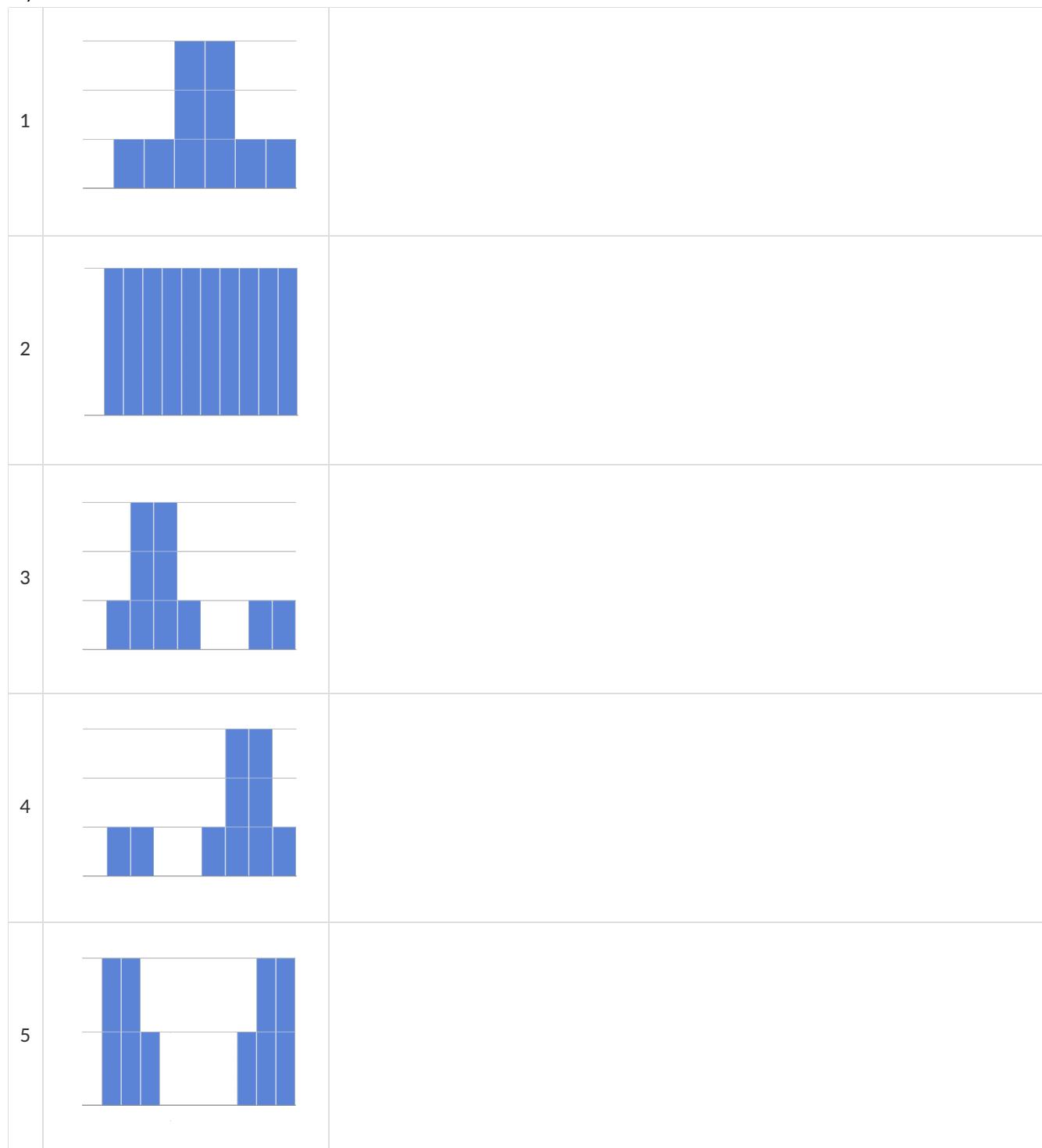
5

E



Identifying Shape

Describe the shape of histograms on the left in complete sentences, using vocabulary like "Skewed Left", "Skewed Right", or "Symmetric".



The Shape of the Animals Dataset

Describe two histograms made from columns of the animals dataset.

1) Make a histogram, showing the distribution of _____ pounds for _____ column in your dataset

animals from the shelter .

your subset, e.g., "fixed dogs from the shelter"

2) Make another histogram, showing the distribution of _____ for _____ column in your dataset

your subset, e.g., "fixed dogs from the shelter"

3) What do you Notice and Wonder about these two histograms? What shape do they have?

What do you NOTICE?	What do you WONDER?

The Shape of My Dataset

Describe two of the histograms you made from your dataset.

1) I made a histogram, showing the distribution of _____ for _____ column in your dataset

your subset, e.g., "fixed dogs from the shelter"

2) I made a histogram, showing the distribution of _____ for _____ column in your dataset

your subset, e.g., "fixed dogs from the shelter"

3) In the table below, describe the histograms. **Are they symmetric?** Do they show left skewness and/or low outliers? ** Do they show Right skewness and/or high outliers?

What do you NOTICE about these displays?	What do you WONDER about these displays?

What's on your mind?

Measures of Center and Spread

There are three ways to measure the **center** of a dataset, to summarize a whole column of quantitative 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. In an ordered list the median will either be the middle number or the average of the two middle numbers.
- The **mode(s)** of a data set is the value (or values) occurring most often. When all of the values occur equally often, a dataset has no mode.

In a **symmetric** dataset, values are just as likely to occur a certain distance above the mean as below the mean, and the median and mean are usually close together.

When a dataset is asymmetric, the median is a more descriptive measure of center than the mean.

- A dataset with **left skew**, and/or low outliers, has a few values that are unusually low, pulling the mean *below* the median.
- A dataset with **right skew**, and/or high outliers, means there are a few values that are unusually high, pulling the mean *above* the median.

When a dataset contains a small number of values, the mode may be the most descriptive measure of center.

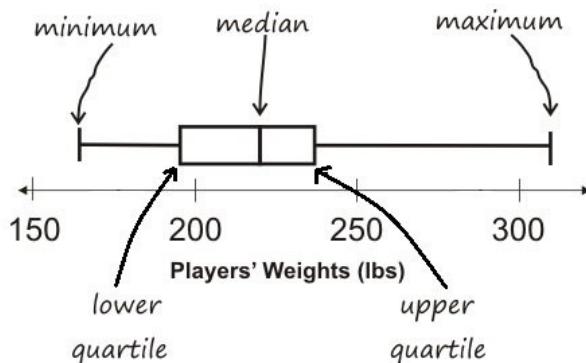
Data Scientists can also measure the **spread** of a dataset using a **five-number summary**:

- The **minimum** – the lowest value in the dataset
- The **first, or “lower” quartile (Q1)** – the middle of the lower half of values, which separates the lowest quarter from the next smallest quarter
- The **second quartile (Q2)** – the middle value, which separates the entire dataset into “top” and “bottom” halves
- The **third, or “upper” quartile (Q3)** – the middle of the higher half of values which separates the second highest quarter from the highest quarter
- The **maximum** – the largest value in the dataset

Measures of Center and Spread (continued)

The **five-number summary** can be used to draw a **box plot**.

- Each of the four sections of the box plot contains 25% of the data. *If the values are distributed evenly across the range, the four sections of the box plot will be equal in width.* Uneven distributions will show up as differently-sized sections of a box plot.
- The left **whisker** extends from the minimum to Q1.
- The **box**, or **interquartile range**, extends from Q1 to Q3. It is divided into 2 parts by the **median**. Each of those parts contains 25% of the data, so the whole box contains the central 50% of the data.
- The right **whisker** extends from Q3 to the maximum.



The box plot above, for example, tells us that:

- The minimum weight is about 165 pounds. The median weight is about 220 pounds. The maximum weight is about 310 pounds.
 - 1/4 of the players weigh roughly between 165 and 195 pounds
 - 1/4 of the players weigh roughly between 195 and 220 pounds
 - 1/4 of the players weigh roughly between 220 and 235 pounds
 - 1/4 of the players weigh roughly between 235 and 310 pounds
 - 50% of the players weigh roughly between 165 and 220 pounds
 - 50% of the players weigh roughly between 195 and 235 pounds
 - 50% of the players weigh roughly between 220 and 310 pounds
- The densest concentration of players' weights is between 220 and 235 pounds.
- Because the widest section of the box plot is between 235 and 310 pounds, we understand that the weights of the heaviest 25% fall across a wider span than the others. 310 may be an outlier, the weights of the players weighing between 235 pounds and 310 pound could be evenly distributed across the range, or all of the players weighing over 235 pounds may weigh around 310 pounds.

Summarizing Columns in the Animals Dataset

Find the measures of center and spread to summarize the _____ pounds column of the Animals Table.

Be sure to add examples to your Contracts page as you work.

Measures of Center

The three measures of center for this column are:

Mean (Average)	Median	Mode(s)

Since the mean is _____ compared to the median, this suggests the shape is
[higher/lower/about equal]

[skewed right (or high outliers) / skewed left (or low outliers) / symmetric]

Measures of Spread

My five-number summary is:

Minimum	Q1	Median	Q3	Maximum

Displaying Center and Spread with a Box Plot

Draw a box plot from this summary on the number line below.

Be sure to label the number line with consistent intervals.



From this summary and box plot, I conclude:

Interpreting Spread

Consider the following dataset, representing the annual income of ten people.

All numbers represent *thousands of dollars* (so 14 means "\$14,000"):

60, 10, 21, 180, 14, 20, 45, 35, 45, 170

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

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

Mean (Average)	Median	Mode(s)

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

Minimum	Q1	Q2 (Median)	Q3	Maximum

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

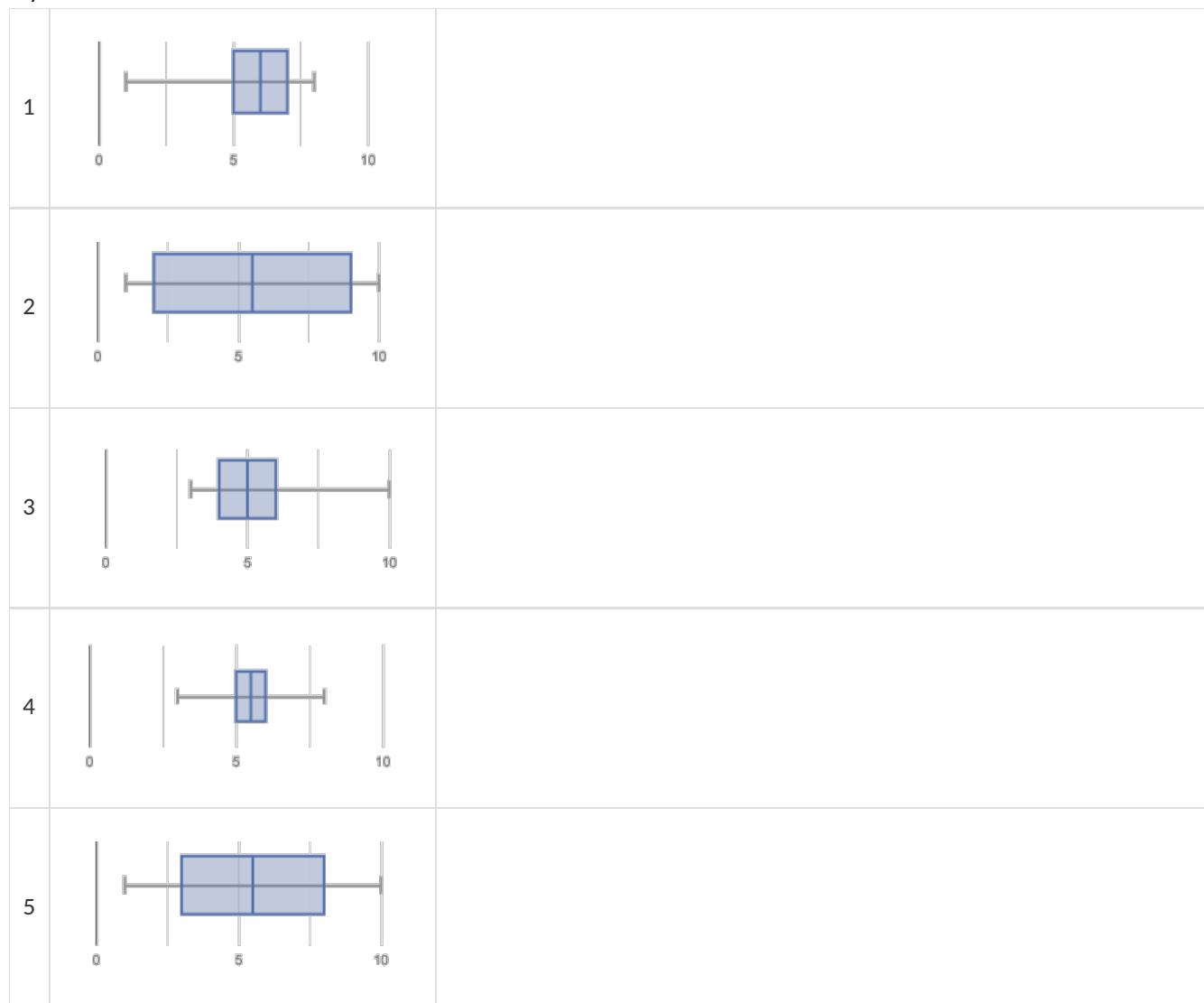


5) The following statements are *correct ... but misleading*. Write down the reason why.

Statement	Why it's misleading
"They're rich! The average person makes more than \$70k dollars!"	
"It's a middle-income list: the most common salary is \$45k/yr!"	
"This group is really diverse, with people making as little as \$14,000 and as much as \$280,000!"	

Identifying Shape

Describe the shape of the box plots below in complete sentences, using vocabulary like "Skewed Left", "Skewed Right", or "Symmetric".



Shape of My Dataset

Find the measures of center and spread to summarize a column of your dataset.

The column I chose to summarize is _____.

Measures of Center

The three measures of center for this column are:

Mean (Average)	Median	Mode(s)

Since the mean is _____ compared to the median, this suggests the shape is
[higher/lower/about equal]

[skewed right (or high outliers) / skewed left (or low outliers) / symmetric]

Measures of Spread

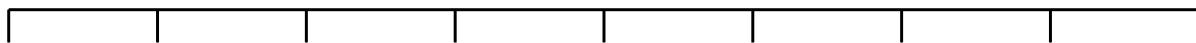
My five-number summary is:

Minimum	Q1	Q2 (Median)	Q3	Maximum

Displaying Center and Spread with a Box Plot

Draw a box plot from this summary on the number line below.

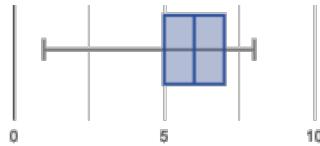
Be sure to label the number line with consistent intervals.



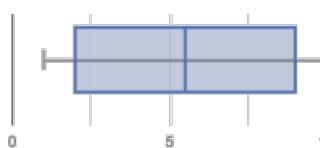
From this summary and box plot, I conclude:

Matching Box-Plots to Histograms

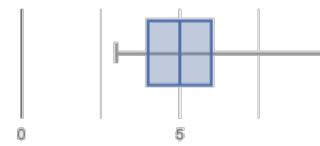
Students watched 5 videos, and rated them on a scale of 1 to 10. For each video, their ratings were used to generate box-plots and histograms. Match the box-plot to the histogram that displays the same data.



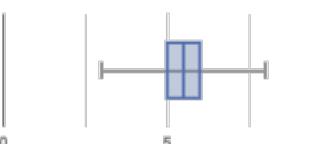
1



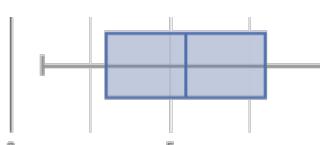
2



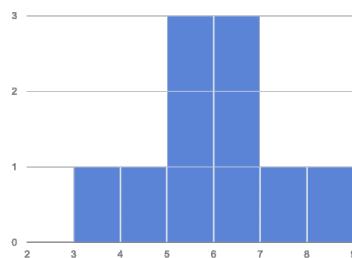
3



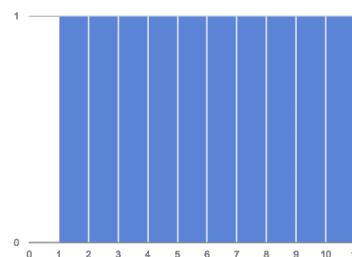
4



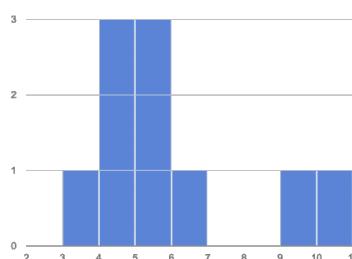
5



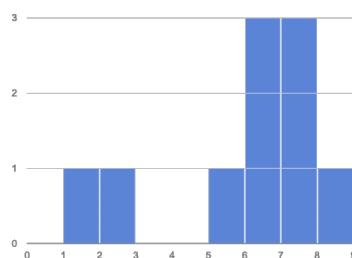
A



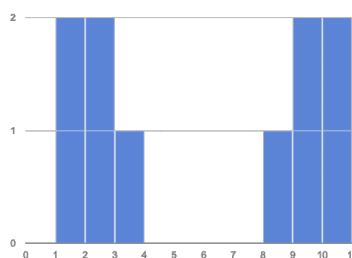
B



C



D



E

What's on your mind?

“Trust, but verify ...”

A “helpful” Data Scientist gives you access to the following functions:

```
# fixed-cats :: (animals :: Table) -> Table
# consumes a table of animals, and produces a table containing only
# cats that have been fixed, sorted from youngest-to-oldest
```

You can use the function, *but you can't see the code for it!* **How do you know if you can trust their code?**

HINT:

- You could make a *verification subset* that contains one of every species, and make sure that the function filters out everything but cats.
- You could make sure this subset has multiple cats not already ordered of youngest-to-oldest, and make sure the function puts them in the right order.

1) What other qualities would this subset need to have?

2) Create your verification subset! In the space below, list the name of each animal in your subset.

Name

“Trust, but verify...”

A “helpful” Data Scientist gives you access to the following functions:

```
# old-dogs-nametags:: (animals :: Table) -> Table
# consumes a table of animals, and produces a table containing only
# dogs 5 years or older, with an extra column showing their name in red
```

You can use the function, *but you can't see the code for it!* **How do you know if you can trust their code?**

- 1) What qualities would a verification subset need to have?

- 2) Create your verification subset! In the space below, list the name and index of each animal in your subset.

Name

What's on your mind?

Scatter Plots

Scatter Plots can be used to show a 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 straight line, it is possible that a linear relationship exists between those two columns. A number called a **correlation** can be used to summarize this relationship.
- r is the name of the **correlation statistic**. The r -value will always fall between -1 and $+1$. The sign tells us whether the correlation is positive or negative. Distance from 0 tells us the strength of the correlation.
 - -1 or $+1$ is really strong.
 - 0 means no correlation.
- The correlation is **positive** if the point cloud slopes up as it goes farther to the right. It is **negative** if it slopes down as it goes farther to the right. If the points are tightly clustered around a line, it is a **strong** correlation. If they are loosely scattered, it is a **weak** correlation.
- Points that are far above or below the cloud of points in a scatter plot are called **outliers**.
- We graphically summarize this relationship by drawing a straight line through the data cloud, so that the vertical distance between the line and each of the points is as small as possible. This line is called the **line of best fit** and allows us to predict y-values based on x-values.

(Dis)Proving a Claim

Smaller animals get adopted faster because they're easier to care for."

Do you agree? If so, why?

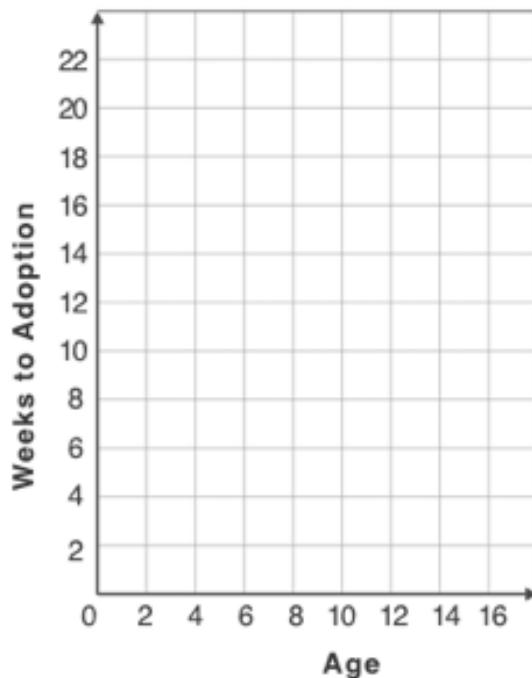
I hypothesize ...

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

Creating a Scatter Plot

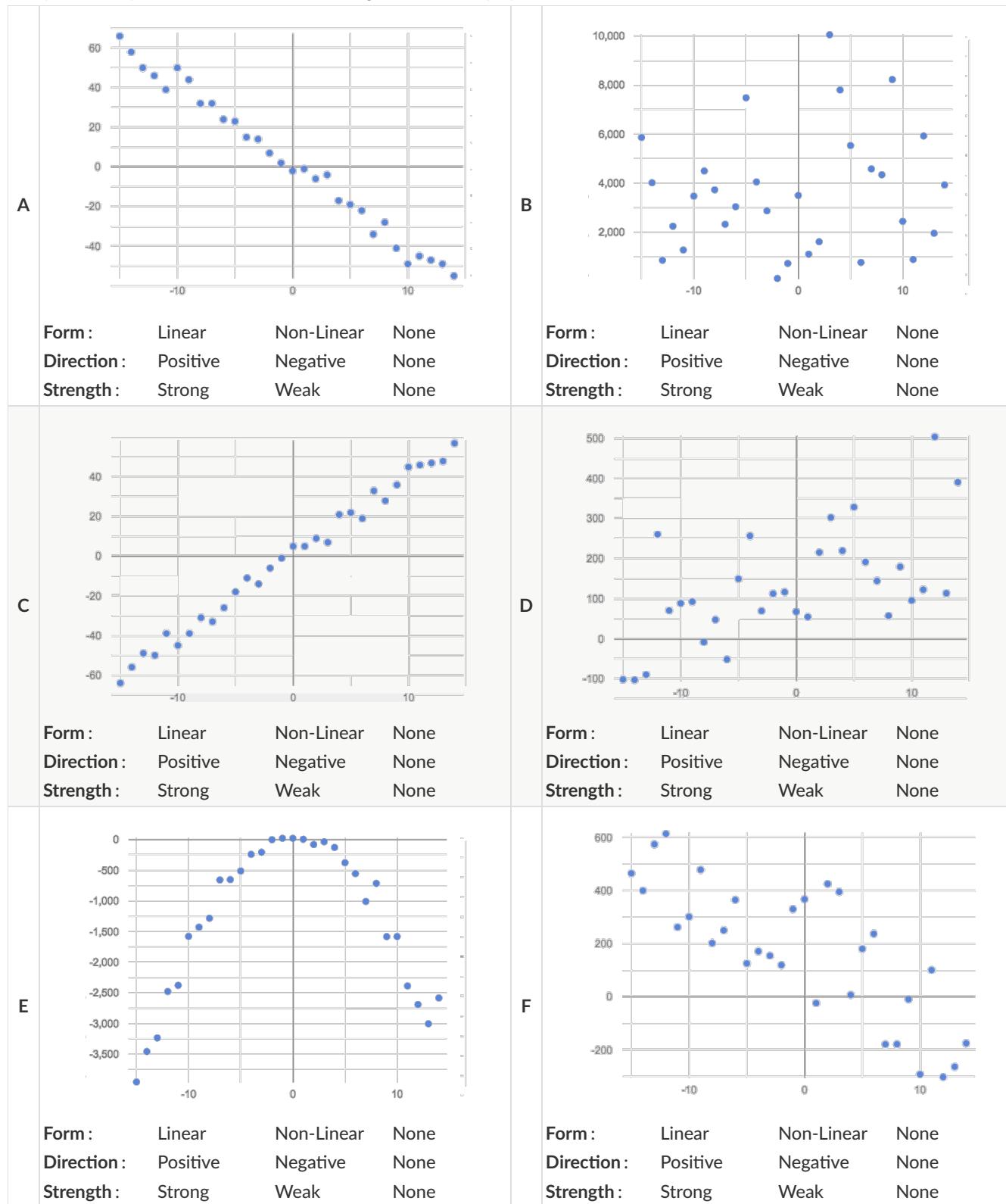
- For each row in the Sample Table on the left, add a point to the scatter plot on the right. Use the values from the age column for the x-axis, and values from the weeks column for the y-axis.
 - Do you see a pattern? Do the points seem to go up or down as age increases to the right?
 - Draw a line on the scatter plot to show this pattern.
 - Does the line slope upwards or downwards?
-
- Are the points tightly clustered around the line or loosely scattered?
-

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



Identifying Form, Direction and Strength

Can you identify the Form, Direction, & Strength of these displays? Note: If the form is non-linear, there is no direction!

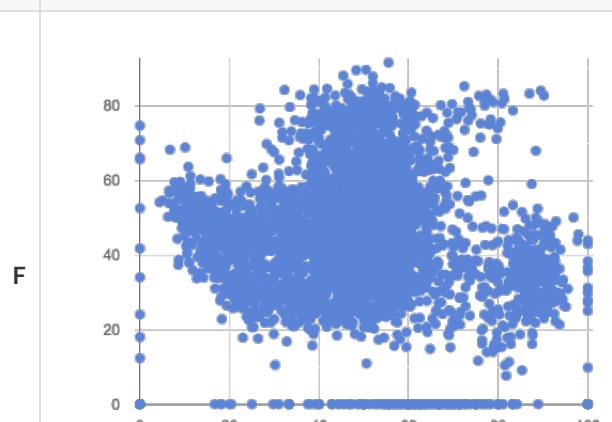
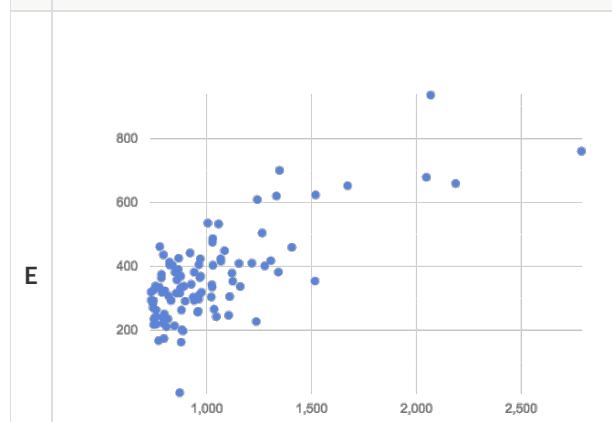
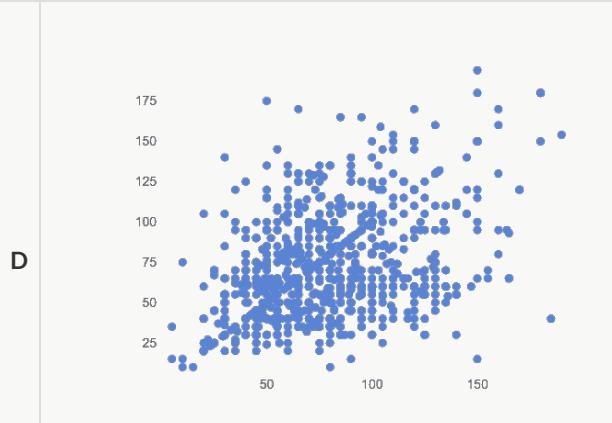
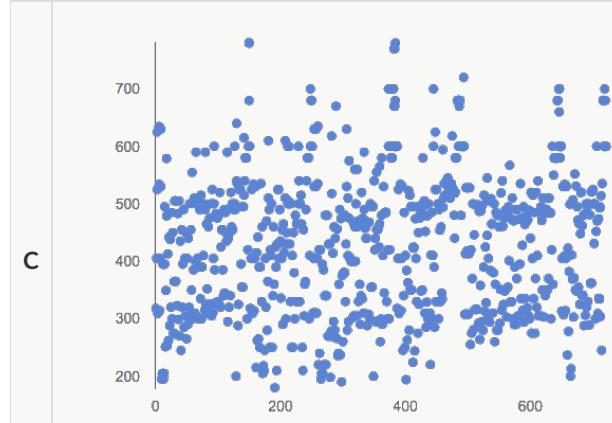
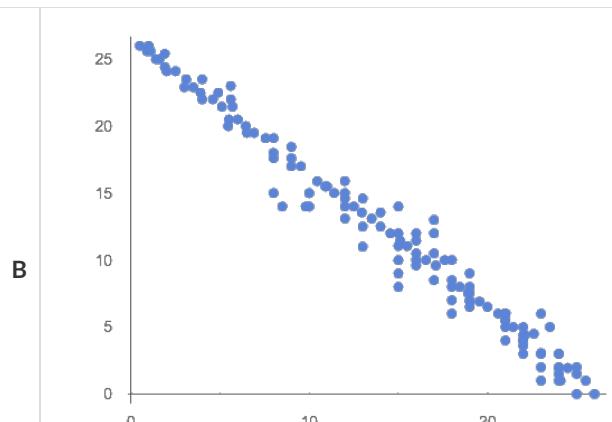
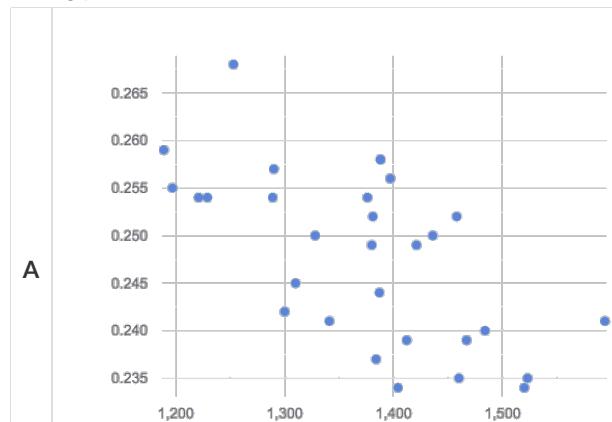


Identifying Form and r-Values

Can you identify the Form, Direction, and Strength of these displays?

If the form is linear, approximate the r -value to express Direction and Strength.

Reminder: An r -value close to -1 is a strong negative relationship, an r -value close to 0 is weak, and an r -value close to +1 is a strong positive!



Correlations in My Dataset

1) There may be a correlation between _____ and _____.

column

column

I think it is a _____, _____ correlation,
strong/weak positive/negative

because _____.

It might be stronger if I looked at _____.
a sample or extension of my data

2) There may be a correlation between _____ and _____.

column

column

I think it is a _____, _____ correlation,
strong/weak positive/negative

because _____.

It might be stronger if I looked at _____.
a sample or extension of my data

3) There may be a correlation between _____ and _____.

column

column

I think it is a _____, _____ correlation,
strong/weak positive/negative

because _____.

It might be stronger if I looked at _____.
a sample or extension of my data

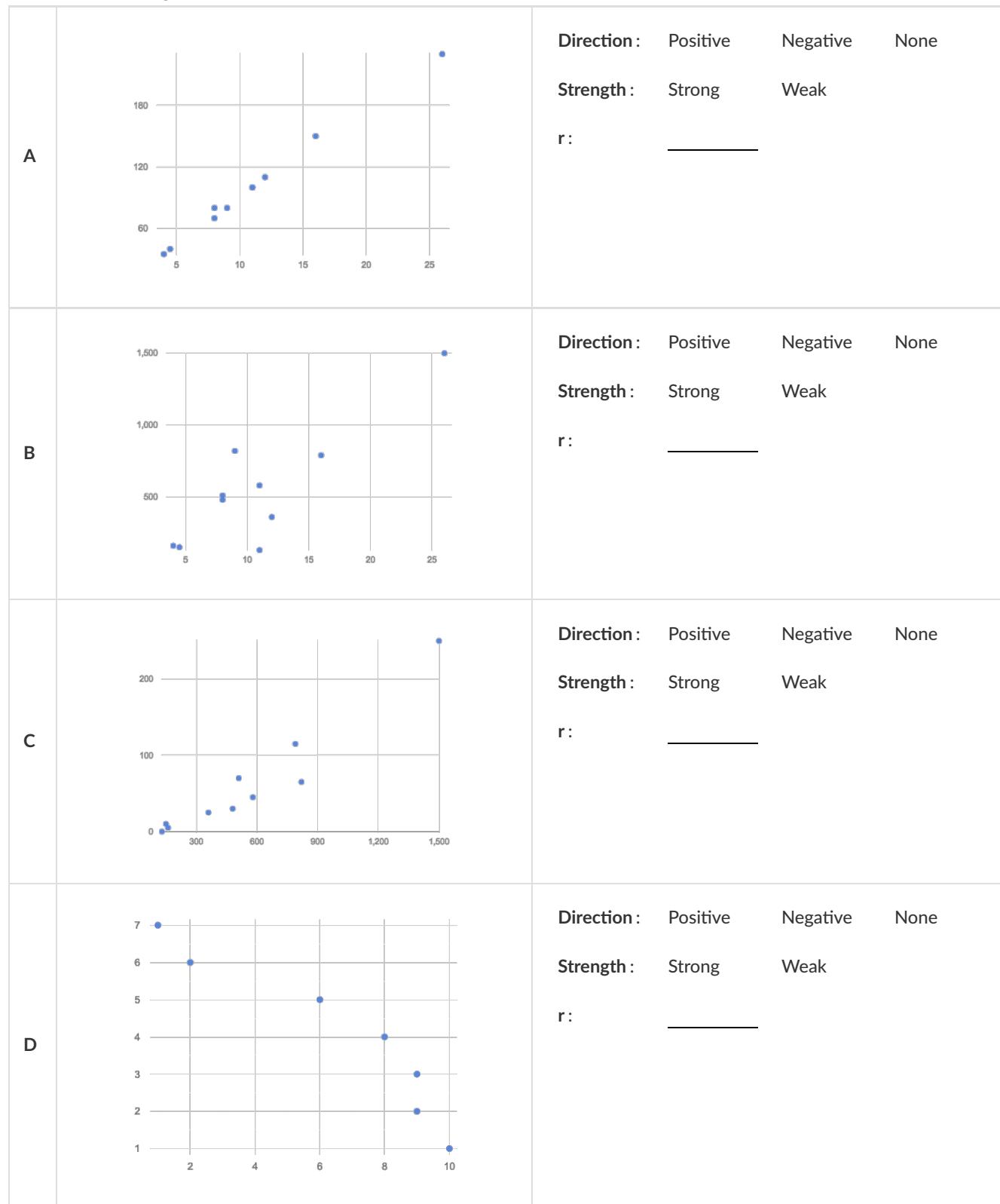
Computing Relationships

Linear Regression is a way of computing the **line of best fit**, which minimizes the sum of vertical distances of all scatter plot points from the line. Calculating the slope and intercept of this line is a task best left to computing or statistical software.

- **Slope** provides us with the easiest summary to grasp: it's how much we predict the y-variable (response variable) will increase or decrease for each unit that the x-variable (explanatory variable) increases.
- **Correlation is not causation!** Correlation only suggests that two column variables are related, but does not tell us if one causes the other. For example, hot days are correlated with people running their air conditioners, but air conditioners do not cause hot days!
- **Sample size matters!** The number of data values is also relevant. We'd be more convinced of a positive relationship in general between cat age and time to adoption if a correlation of +0.57 were based on 50 cats instead of 5.

Drawing Predictors

For each of the scatter plots below, draw a **predictor line** that seems like the best fit. Describe the correlation in terms of Direction and Strength, then estimate the r -value.



Interpreting Regression Lines & r-Values

Each description on the left is written about the linear regression findings on the right. Fill in the blanks using the information in the line of best fit and the r-value.

<p>1</p> <p>For every additional Marvel Universe movie released each year, the average person is predicted to consume _____ pounds of sugar! <small>[amount] [more / fewer]</small></p> <p>This correlation is _____. <small>[strong, moderate, weak, non-existent]</small></p>	$y = -3.19x + 12$ $r = -0.05$
<p>2</p> <p>Shoe size and height are _____, <small>[strongly, moderately, weakly, not]</small> correlated. If person A is one size bigger than person B, we predict that they will be roughly _____ inches taller than person B as well. <small>[positively / negatively]</small> <small>[amount]</small></p>	$y = 1.65x + 52$ $r = 0.89$
<p>3</p> <p>There is _____ relationship found between the number of Uber drivers in a city and the number of babies born each year. <small>[a strong, a moderate, a weak, no]</small></p>	$y = -15.3x + 1150$ $r = 0.01$
<p>4</p> <p>The correlation between weeks-of-school-missed and SAT score is _____ and _____. For every week a student misses, we predict a more than a _____ point _____ in their SAT score. <small>[strong, moderate, weak, non-existent]</small> <small>[positive / negative]</small> <small>[amount]</small> <small>[gain / drop]</small></p>	$y = -5.35x - 16$ $r = -0.65$
<p>5</p> <p>There is a _____, _____ correlation between the number of streaming video services someone has, and how much they weigh. For each service, we expect them to be roughly _____ pounds heavier. <small>[amount]</small></p>	$y = 1.6x + 160$ $r = 0.12$

Regression Analysis in the Animals Dataset

1) I performed a linear regression on a sample of _____ cats from the shelter
dataset or subset
and found _____ a moderate ($r=0.566$), positive correlation
between _____ age of the cats (in years) and _____ number of weeks to adoption.
I would predict that a 1 _____ year increase in _____ age is associated with a
0.23 week increase in _____ adoption time.
[slope, y-units] [x-axis units] [increase/decrease] [x-axis] [y-axis]

2) I performed a linear regression on a sample of _____ and
dataset or subset
found _____ a weak/strong/moderate ($R=...$), positive/negative correlation between
and _____.
[x-axis] [y-axis]
I would predict that a 1 _____ increase in _____ is associated with a
[x-axis units] [x-axis]
in _____.
[slope, y-units] [increase/decrease] [y-axis]

3) I performed a linear regression on a sample of _____ dataset or subset
and found _____ correlation
between _____ a weak/strong/moderate ($R=...$), positive/negative
and _____.
[x-axis] [y-axis]
I would predict that a 1 _____ increase in _____ is associated with a
[x-axis units] [x-axis]
in _____.
[slope, y-units] [increase/decrease] [y-axis]

Regression Analysis in Your Dataset

My Dataset is _____.

1) I performed a linear regression on _____ and found
dataset or subset

correlation between

a weak/strong/moderate ($R=...$), positive/negative

and _____.
[x-axis] [y-axis]

I would predict that a 1 _____ increase in _____ is associated with a
[x-axis units] [x-axis]

in _____.
[slope, y-units] [increase/decrease] [y-axis]

2) I performed a linear regression on _____ and found
dataset or subset

correlation between

a weak/strong/moderate ($R=...$), positive/negative

and _____.
[x-axis] [y-axis]

I would predict that a 1 _____ increase in _____ is associated with a
[x-axis units] [x-axis]

in _____.
[slope, y-units] [increase/decrease] [y-axis]

3) I performed a linear regression on _____ and found
dataset or subset

correlation between

a weak/strong/moderate ($R=...$), positive/negative

and _____.
[x-axis] [y-axis]

I would predict that a 1 _____ increase in _____ is associated with a
[x-axis units] [x-axis]

in _____.
[slope, y-units] [increase/decrease] [y-axis]

What's on your mind?

Case Study: Ethics, Privacy, and Bias

My Case Study is _____

- 1) Read the case study you or your group was assigned, and write your summary here.

- 2) Is this a good thing or a bad thing? Why?

- 3) What are the arguments on *each* side?

Data Science used for this purpose is good because...

Data Science used for this purpose is bad because...

Threats to Validity

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.
- **Study bias** - If someone is supposed to assess how much cat food is eaten each day on average, but they only measure how much cat food is put in the bowls (instead of how much is actually consumed), they'll end up with an over-estimate.
- **Poor choice of summary** - Suppose a different shelter that had 10 animals recorded adoption times (in weeks) as 1, 1, 1, 7, 7, 8, 8, 9, 9, 10. Using the mode (1) to report what's typical would make it seem like the animals were adopted much quicker than they really were, since 7 out of 10 animals took at least 7 weeks to be adopted.
- **Confounding variables** - Shelter workers might steer people towards newer animals, because they've become attached to the animals that have been there for a while, making it appear that "staying in the shelter longer" means "less likely to be adopted".

Identifying 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. From this survey, they concluded that dogs are the most popular pet in the state.

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

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

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

Identifying 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 they found that spider and rabbit food was by far the most popular cuisine!

Explain why sampling just two animals can result in unreliable conclusions about what kind of food is needed.

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?

Fake News!

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

	Data	Claim	What's Wrong
1	The average player on a basketball team is 6'1".	"Most of the players are taller than 6 feet."	
2	Linear regression found a positive correlation ($r=0.18$) between people's height and salary.	"Higher salaries can make people taller!."	
3	$y=12.234x + -17.089; r-sq: 0.636$	"According to the predictor function indicated here, the value on the x-axis will predict the value on the y-axis 63.6% of the time."	
4		"According to this bar chart, Felix makes up a little more than 15% of the total ages of all the animals in the dataset."	
5		"According to this histogram, most animals weigh between 40 and 60 pounds."	
6	Linear regression found a negative correlation ($r= -0.91$) between the number of hairs on a person's head and their likelihood of owning a wig.	"Owning wigs causes people to go bald."	

Lies, Darned Lies, and Statistics

1) Using real data and displays from your dataset, come up with a misleading claim.

2) Trade papers with someone and figure out why their claims are wrong!

Data	Claim	Why it's wrong
1		
2		
3		
4		

What's on your mind?

Design Recipe

Directions :

Contract and Purpose Statement

Every contract has three parts...

_____ :: _____ -> _____
function name domain range

what does the function do?

Examples

Write some examples, then circle and label what changes...

examples:

_____ (_____) is _____
function name input(s) what the function produces

_____ (_____) is _____
function name input(s) what the function produces

end

Definition

Write the definition, giving variable names to all your input values...

fun _____ (_____):
function name variable(s)

what the function does with those variable(s)

end

Directions :

Contract and Purpose Statement

Every contract has three parts...

_____ :: _____ -> _____
function name domain range

what does the function do?

Examples

Write some examples, then circle and label what changes...

examples:

_____ (_____) is _____
function name input(s) what the function produces

_____ (_____) is _____
function name input(s) what the function produces

end

Definition

Write the definition, giving variable names to all your input values...

fun _____ (_____):
function name variable(s)

what the function does with those variable(s)

end

Design Recipe

Directions :

Contract and Purpose Statement

Every contract has three parts...

_____ :: _____ -> _____
function name domain range

what does the function do?

Examples

Write some examples, then circle and label what changes...

examples:

_____ (_____) is _____
function name input(s) what the function produces

_____ (_____) is _____
function name input(s) what the function produces

end

Definition

Write the definition, giving variable names to all your input values...

fun _____ (_____):
function name variable(s)

what the function does with those variable(s)

end

Directions :

Contract and Purpose Statement

Every contract has three parts...

_____ :: _____ -> _____
function name domain range

what does the function do?

Examples

Write some examples, then circle and label what changes...

examples:

_____ (_____) is _____
function name input(s) what the function produces

_____ (_____) is _____
function name input(s) what the function produces

end

Definition

Write the definition, giving variable names to all your input values...

fun _____ (_____):
function name variable(s)

what the function does with those variable(s)

end

Design Recipe

Directions :

Contract and Purpose Statement

Every contract has three parts...

_____ :: _____ -> _____
function name domain range

what does the function do?

Examples

Write some examples, then circle and label what changes...

examples:

_____ (_____) is _____
function name input(s) what the function produces

_____ (_____) is _____
function name input(s) what the function produces

end

Definition

Write the definition, giving variable names to all your input values...

fun _____ (_____):
function name variable(s)

what the function does with those variable(s)

end

Directions :

Contract and Purpose Statement

Every contract has three parts...

_____ :: _____ -> _____
function name domain range

what does the function do?

Examples

Write some examples, then circle and label what changes...

examples:

_____ (_____) is _____
function name input(s) what the function produces

_____ (_____) is _____
function name input(s) what the function produces

end

Definition

Write the definition, giving variable names to all your input values...

fun _____ (_____):
function name variable(s)

what the function does with those variable(s)

end

Contracts

Contracts tell us how to use a function. For example: num-min :: (a :: Number, b :: Number) -> Number tells us that the name of the function is num-min , it takes two inputs (both Numbers), and it evaluates to a Number . From the contract, we know num-min(4, 6) will evaluate to a Number . Use the blank line under each contract for notes or sample code for that function!

Name	Domain	Range
triangle	:: (side-length :: Number, style :: String, color :: String)	-> Image
circle	:: (radius :: Number, style :: String, color :: String)	-> Image
star	:: (radius :: Number, style :: String, color :: String)	-> Image
rectangle	:: (width :: Num, height :: Num, style :: Str, color :: Str)	-> Image
ellipse	:: (width :: Num, height :: Num, style :: Str, color :: Str)	-> Image
square	:: (size-length :: Number, style :: String, color :: String)	-> Image
text	:: (str :: String, size :: Number, color :: String)	-> Image
overlay	:: (img1 :: Image, img2 :: Image)	-> Image
beside	:: (img1 :: Image, img2 :: Image)	-> Image
above	:: (img1 :: Image, img2 :: Image)	-> Image
put-image	:: (img1 :: Image, x :: Number, y :: Number, img2 :: Image)	-> Image
rotate	:: (degree :: Number, img :: Image)	-> Image
scale	:: (factor :: Number, img :: Image)	-> Image

Contracts

Contracts tell us how to use a function. For example: num-min :: (a :: Number, b :: Number) -> Number tells us that the name of the function is num-min , it takes two inputs (both Numbers), and it evaluates to a Number . From the contract, we know num-min(4, 6) will evaluate to a Number . Use the blank line under each contract for notes or sample code for that function!

Name	Domain	Range
string-repeat	:: (text :: String, repeat :: Number)	-> String
string-contains	:: (text :: String, search-for :: String)	-> Boolean
num-sqr	:: (n :: Number)	-> Number
num-sqrt	:: (n :: Number)	-> Number
num-min	:: (a :: Number, b:: Number)	-> Number
num-max	:: (a :: Number, b:: Number)	-> Number
count	:: (t :: Table, col :: String)	-> Table
mean	:: (t :: Table, col :: String)	-> Number
median	:: (t :: Table, col :: String)	-> Number
modes	:: (t :: Table, col :: String)	-> List<Number>
bar-chart	:: (t :: Table, col :: String)	-> Image
pie-chart	:: (t :: Table, col :: String)	-> Image
histogram	:: (t :: Table, values :: String, bin-width :: Number)	-> Image

Contracts

Contracts tell us how to use a function. For example: num-min :: (a :: Number, b :: Number) -> Number tells us that the name of the function is num-min , it takes two inputs (both Numbers), and it evaluates to a Number . From the contract, we know num-min(4, 6) will evaluate to a Number . Use the blank line under each contract for notes or sample code for that function!

Name	Domain	Range
box-plot	:: (t :: Table, col :: String)	-> Image
modified-box-plot	:: (t :: Table, col :: String)	-> Image
scatter-plot	:: (t :: Table, labels :: String, xs :: String, ys :: String)	-> Image
image-scatter-plot	:: (t :: Table, xs :: String, ys :: String, f :: (Row -> Image))	-> Image
r-value	:: (t :: Table, xs :: String, ys :: String)	-> Number
lr-plot	:: (t :: Table, labels :: String, xs :: String, ys :: String)	-> Image
random-rows	:: (t :: Table, num-rows :: Number)	-> Table
<Table>.row-n	:: (n :: Number)	-> Row
<Table>.order-by	:: (col :: String, increasing :: Boolean)	-> Table
<Table>.filter	:: (test :: (Row -> Boolean))	-> Table
<Table>.build-column	:: (col :: String, builder :: (Row -> Any))	-> Table
bar-chart-summarized	:: (t :: Table, labels :: String, values :: String)	-> Image
pie-chart-summarized	:: (t :: Table, labels :: String, values :: String)	-> Image