





Workbook v1.4

Brought to you by the Bootstrap team:

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

---

Bootstrap is licensed under a Creative Commons 3.0 Unported License. Based on a work from [www.BootstrapWorld.org](http://www.BootstrapWorld.org). Permissions beyond the scope of this license may be available at [schanzer@BootstrapWorld.org](mailto: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 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 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. Categorical data is not subject to the laws of arithmetic – for example, we cannot take the "average" of a list of colors.

**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



# The Animals Dataset

What do you NOTICE about the animals dataset?	What do you WONDER about the animals dataset?

1. This dataset is Animals from an animal shelter, which contains 31 data rows.
2. Some of the columns are:
  1. name, which contains categorical data, and is of type String. Some example values from this column are: "Toggle", "Fritz", and "Nori".
  2. \_\_\_\_\_, which contains \_\_\_\_\_ data, and is of type \_\_\_\_\_. Some example values from this column are: \_\_\_\_\_.
  3. \_\_\_\_\_, which contains \_\_\_\_\_ data, and is of type \_\_\_\_\_. Some example values from this column are: \_\_\_\_\_.

# 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:
- 

## 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:
- 
7. Typing in the following expressions, one at a time: `4 + 2 + 6`, `4 + 2 * 6`, and `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.

`3 <= 4`

\_\_\_\_\_

`3 == 2`

\_\_\_\_\_

`2 <> 4`

\_\_\_\_\_

`3 <> 3`

\_\_\_\_\_

`"a" > "b"`

\_\_\_\_\_

`"a" <> "b"`

\_\_\_\_\_

`"a" == "b"`

\_\_\_\_\_

`"a" <> "a"`

\_\_\_\_\_

---

## 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!

`(3 <= 4) and (3 == 2)`

\_\_\_\_\_

`("a" == "b") and (3 <> 4)`

\_\_\_\_\_

`(3 <= 4) or (3 == 2)`

\_\_\_\_\_

`("a" == "b") or (3 <> 4)`

\_\_\_\_\_

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





# 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?”)
- **Relate 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.



# Questions about the Animals Dataset

My question is...	This is a... (circle one)
	<ul style="list-style-type: none"><li>• Lookup</li><li>• Compute</li><li>• Relate</li></ul>
	<ul style="list-style-type: none"><li>• Lookup</li><li>• Compute</li><li>• Relate</li></ul>
	<ul style="list-style-type: none"><li>• Lookup</li><li>• Compute</li><li>• Relate</li></ul>
	<ul style="list-style-type: none"><li>• Lookup</li><li>• Compute</li><li>• Relate</li></ul>
	<ul style="list-style-type: none"><li>• Lookup</li><li>• Compute</li><li>• Relate</li></ul>
	<ul style="list-style-type: none"><li>• Lookup</li><li>• Compute</li><li>• Relate</li></ul>

# Lookup Questions

The table below represents four pets at an animal shelter:

## animals-table

name	gender	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?”

```
animals-table.row-n(3)
```

“Which is the last row in the table?”

```
animals-table.row-n(2) ["name"]
```

“What is Fritz’s gender?”

```
animals-table.row-n(1) ["gender"]
```

“What’s the third animal’s name?”

```
animals-table.row-n(3) ["age"]
```

"How much does Nori weigh?"

```
animals-table.row-n(3) [ "pounds" ]
```

“How old is Maple?”

```
animals-table.row-n(0)
```

“What is Toggle’s gender?”

```
animals-table.row-n(2) ["pounds"]
```

“What is the first row in the table?”

```
animals-table.row-n(0) ["gender"]
```

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

```
animals-table.row-n(3)["name"]
```

“Maple”

"male"

4

48

"Nori"

# More Practice with Lookups

Consider the table below, and the four value definitions that follow:

**shapes-table**

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

```
shapeA = shapes-table.row-n(0)
shapeB = shapes-table.row-n(1)
shapeC = shapes-table.row-n(2)
shapeD = shapes-table.row-n(3)
```

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

shapeD	Evaluates to 4
shapeA	Evaluates to the last row in the table
shapeB["corners"]	Evaluates to "square"
shapeC["is-round"]	Evaluates to true
shapeB["name"]	Evaluates to false
shapeA["corners"]	Evaluates to 3
shapeD["name"] == "circle"	Evaluates to the first row in the table

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

a. _____	"rectangle"
b. _____	"triangle"
c. _____	4
d. _____	0
e. _____	true

# 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
    ( ) is
end
fun ( ) :
end
```

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

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

**examples:**

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

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

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

**examples:**

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





## Unit 3

Data Scientists often make **subsets** of data, to group them into logical parts. A dataset of students, for example, might have subsets for each grade, or for each homeroom teacher.

Each subset is a **sample** of the original population. Data Scientists try to make predictions about the whole population based on that sample. However, choosing a good sample instead of a bad one can be tricky!



# Samples from the Animals Dataset

What are some subsets you can create from this dataset? For a given row `r`, what code will identify if that row is in the subset?

Dogs	<code>r["species"] == "dog"</code>
Kittens	<code>(r["age"] &lt; 2) and (r["species"] == "cat")</code>

# My Dataset

What do you NOTICE about your dataset?	What do you WONDER about your dataset?

1. This dataset is \_\_\_\_\_, which contains \_\_\_\_\_ data rows.
2. Some of the columns are:
  1. \_\_\_\_\_, which contains \_\_\_\_\_ data, and is of type \_\_\_\_\_. Some example values from this column are: \_\_\_\_\_.
  2. \_\_\_\_\_, which contains \_\_\_\_\_ data, and is of type \_\_\_\_\_. Some example values from this column are: \_\_\_\_\_.
  3. \_\_\_\_\_, which contains \_\_\_\_\_ data, and is of type \_\_\_\_\_. Some example values from this column are: \_\_\_\_\_.

# Questions about My Dataset

My question is...	This is a...(circle one)
	<ul style="list-style-type: none"> <li>• Lookup</li> <li>• Compute</li> <li>• Relate</li> </ul>
	<ul style="list-style-type: none"> <li>• Lookup</li> <li>• Compute</li> <li>• Relate</li> </ul>
	<ul style="list-style-type: none"> <li>• Lookup</li> <li>• Compute</li> <li>• Relate</li> </ul>
	<ul style="list-style-type: none"> <li>• Lookup</li> <li>• Compute</li> <li>• Relate</li> </ul>
	<ul style="list-style-type: none"> <li>• Lookup</li> <li>• Compute</li> <li>• Relate</li> </ul>
	<ul style="list-style-type: none"> <li>• Lookup</li> <li>• Compute</li> <li>• Relate</li> </ul>

# Samples from My Dataset

What are some subsets you can create from this dataset? For a given row  $x$ , what code will identify if that row is in the subset?


# 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
```

# Design Recipes – Filtering Rows

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 \_\_\_\_\_

$$\# \frac{\quad}{\text{name}} :: \frac{\quad}{\text{domain}} \rightarrow \frac{\quad}{\text{range}}$$

**examples:**

\_\_\_\_\_ ( \_\_\_\_\_ ) **is** \_\_\_\_\_

\_\_\_\_\_ ( \_\_\_\_\_ ) is \_\_\_\_\_

**end**

```
fun _____ ( _____ ) : _____
```

**end**

$$\# \frac{\quad}{\text{name}} :: \frac{\quad}{\text{domain}} \rightarrow \frac{\quad}{\text{range}}$$

**examples:**

\_\_\_\_\_ ( \_\_\_\_\_ ) **is** \_\_\_\_\_

( ) is

**end**



# Unit 4

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

- *Bar charts provide a visual representation of the frequency of values in a **categorical** column.*
- Usually there is no mathematical way to order these bars, but **sometimes there's an order** makes sense. For example, bars for T-Shirt sizes might be presented in order of S, M, L, and XL.

**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 taller the bar.

- *Histograms provide a visual representation of the frequency 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 distribution in the dataset.

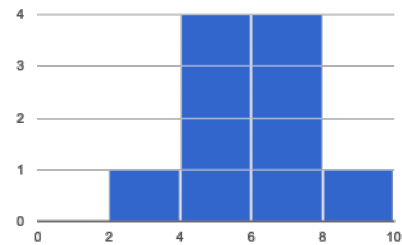


# Reading Histograms

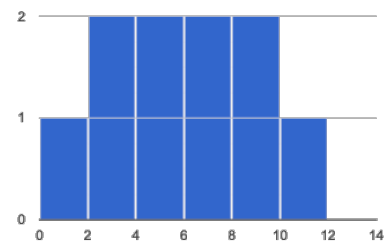
A teacher shows her students five videos, and has them rate how much they liked each one on a scale of 1 to 10. While the **average score** for each video was the same (5.5), the **shapes** of the ratings distributions were very different!

Match the summary description (left) with the histogram of student ratings (right).

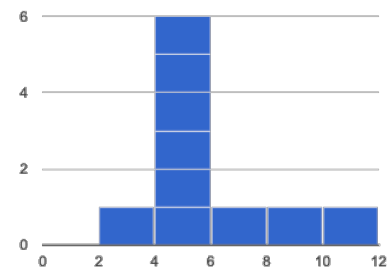
1. Most of the students were fine with the video, but a couple of them gave it an unusually low rating.



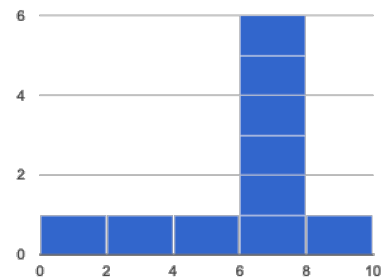
2. Most of the students were OK with the video, but a couple students gave it an unusually high rating.



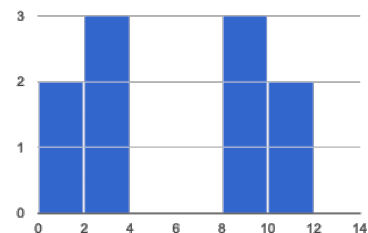
3. Students tended to give the third video an average rating, and they weren't likely to stray far from the average.



4. Students either really liked or really disliked the fourth video.



5. Reactions to the fifth video were all over the place: high ratings and low ratings and in-between ratings were all equally likely.

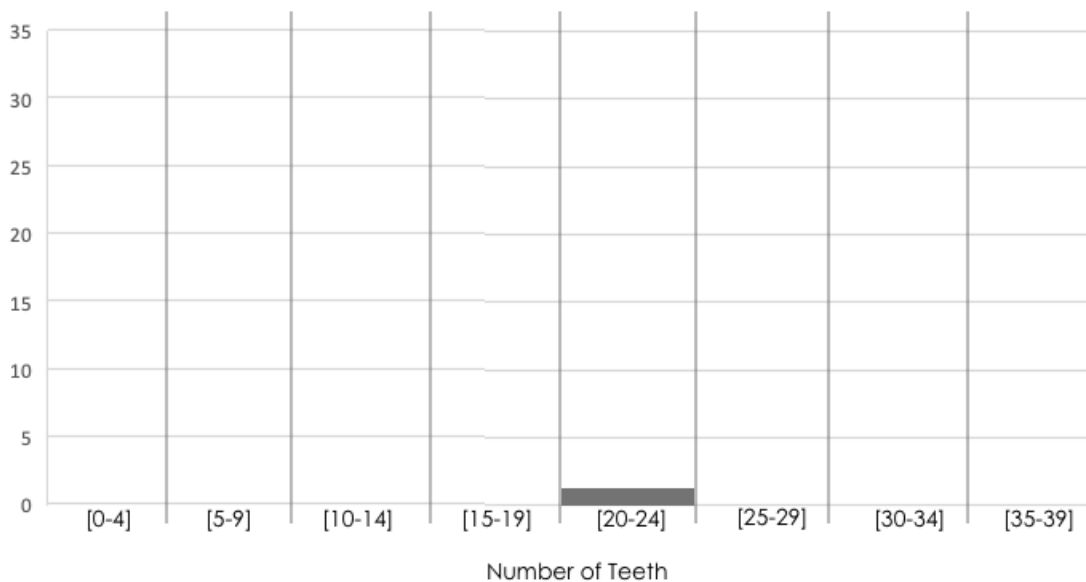


# Constructing Histograms

Suppose we have a data set for number of teeth in a group of 50 adults:

Number of teeth	Count
0	1
22	1
26	2
27	1
28	4
29	3
30	3
31	3
32	33

1. **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.



2. **Circle the statements below that are TRUE**

- The number of teeth in our data set is *skewed left*
- The number of teeth in our data set is *skewed right*
- The number of teeth in our data set has a *low outlier*
- The number of teeth in our data set has a *high outlier*
- The number of teeth in our data set is *symmetric*

# The Shape of the Animals 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 (for example, "fixed dogs at the shelter")

2) I made a histogram, showing the distribution of \_\_\_\_\_ for  
\_\_\_\_\_.

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

# The Shape of My Dataset

Describe two of the histograms you made from your dataset.

3) I made a histogram, showing the distribution of \_\_\_\_\_ for  
column in your dataset

---

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

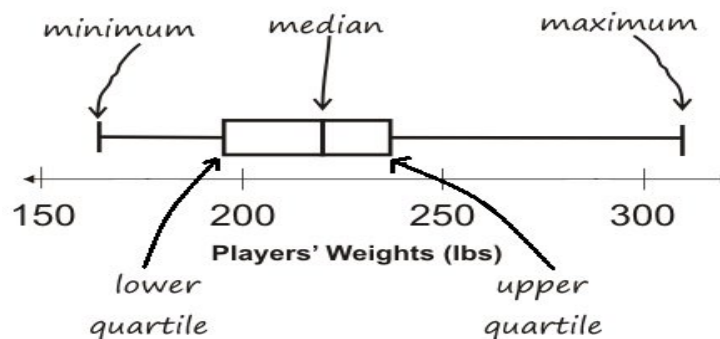
4) I made a histogram, showing the distribution of \_\_\_\_\_ for

---

[illegible]

# Unit 5

- There are three ways to measure the “center” of a dataset, to summarize 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.
- The **shape** of a dataset gives us an idea of which values are more or less common. In a *symmetric* data set, values are just as likely to occur a certain distance below the mean as above it. Outliers or **skew** can shift result in a mean that is higher than the mean (high outliers or right skew) or lower than the mean (low outliers or left skew).
- Data Scientists can also measure the **spread** of a dataset using a **five number summary**:
  - The **minimum** – the smallest value in the dataset
  - The **first, or “lower” quartile (Q1)** – the middle of the smaller half of values, that separates the smallest quarter from the next smallest 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 middle of the larger half of values, that separates the second largest quarter from the largest 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

2) The column I choose to measure is pounds

## 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 or skewness due to values that are unusually \_\_\_\_\_.  
[high / low]

## Measures of Spread

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:

---

---

---

---

# Interpreting Spread

Consider the following 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**.

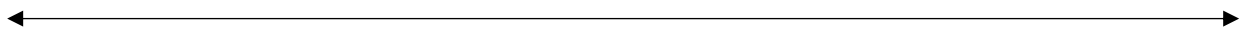
2. In the table below, compute the **measures of center** for this dataset.

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 12k and as much as \$280k!"	

# Summarizing a Column in My Dataset

1) The column I choose to measure is \_\_\_\_\_

## Measures of Center

The three measures for this column are:

Mean (Average)	Median	Mode(s)

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

be outliers or skewness due to values that are unusually \_\_\_\_\_.  
[high / low]

## Measures of Spread

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

**Bar charts** - In bar charts, each bar has a height corresponding to the count or proportion of data values in a given category. Visually, we consider how heights of the bars compare to one another.

**Pie charts** - Pie charts show the relative proportion (or %) of a column's data values that fall into each category. The greater the proportion, the larger the pie slice. Visually, we consider how areas of the slices compare to one another, and to the whole area of 100%.

**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 aren't obviously presented in order



# 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.

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

**examples:**

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

# Design Recipes – Building Columns

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
```

---

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

---

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

**examples:**

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



# Chaining 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).

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

Produces a table containing *only* Toggle and Maple

`t.filter(is-fixed)`

Produces a table, sorted oldest-to-youngest.

`t.build-column("sticker", nametag)`

Produces a table, sorted youngest-to-oldest

`t.filter(is-young)`

Produces a table with an extra column, named "sticker"

`t.order-by("age", false)`

Produces a table containing Maple and Toggle, in that order.

`t.filter(is-young)  
  .order-by("weight", false)`

Produces a table containing the same four animals.

`t.order-by("age", true)  
  .build-column("label", nametag)`

Produces a table with an extra "label" column, sorted youngest-to-oldest



# Unit 7

- **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.
- 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. The points are tightly clustered around a line, it is a **strong** correlation. If they are loosely scattered, it is a **weak** correlation.
- 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 graph 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

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

*Do you agree? If so, why?*

*I hypothesize...*

---

---

---

---

---

---

---

---

---

---

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

---

---

---

---

---

---

---

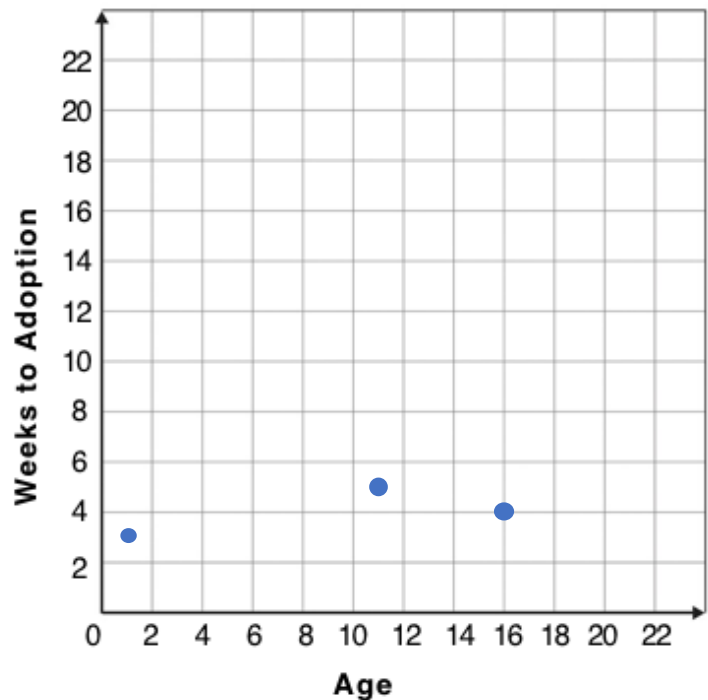
---

---

---

# Creating a Scatter Plot

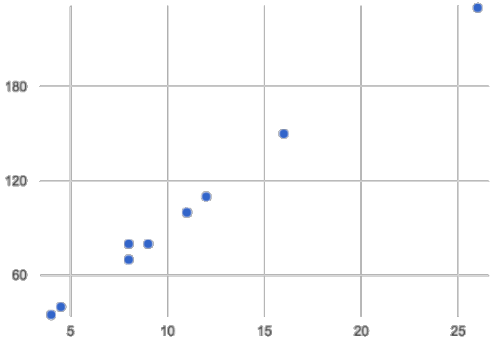
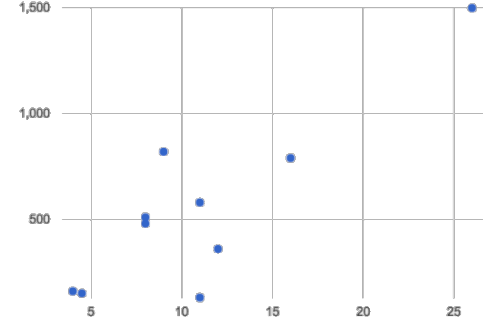
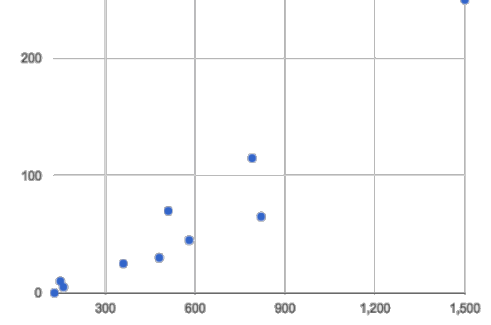
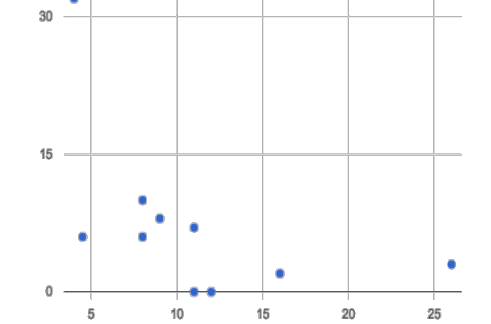
name	species	age	weeks
"Sasha"	"cat"	1	3
"Boo-boo"	"dog"	11	5
"Felix"	"cat"	16	4
"Buddy"	"lizard"	2	24
"Nori"	"dog"	6	9
"Wade"	"cat"	1	2
"Nibblet"	"rabbit"	6	12
"Maple"	"dog"	3	2



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? \_\_\_\_\_
4. Are the points clustered around the line? Loosely scattered? \_\_\_\_\_

# Drawing Predictors

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

<b>A</b>	 <p>fat (g) v. calories-from-fat in common menu items</p>	<p><b>Direction:</b> Positive   Negative   None</p> <p><b>Strength:</b> Strong   Weak</p>
<b>B</b>	 <p>fat (g) v. sodium (g) in common menu items</p>	<p><b>Direction:</b> Positive   Negative   None</p> <p><b>Strength:</b> Strong   Weak</p>
<b>C</b>	 <p>sodium (g) v. cholesterol (mg) in common menu items</p>	<p><b>Direction:</b> Positive   Negative   None</p> <p><b>Strength:</b> Strong   Weak</p>
<b>D</b>	 <p>fat (g) v. sugar (g) in common menu items</p>	<p><b>Direction:</b> Positive   Negative   None</p> <p><b>Strength:</b> Strong   Weak</p>

# 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

- **Linear Regression** is a way of computing the line of best fit, which minimizes the sum of squared 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 to increase, for each unit that the x-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, 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.

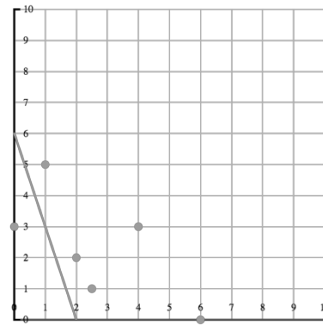
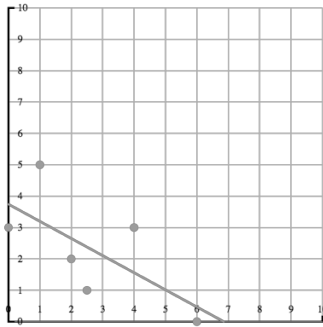


# Grading Predictors

Below are the scatter plots for data sets A-D, with two different 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 **-1** (strong negative correlation) and **+1** (strong positive correlation). A grade of 0 means "there is no correlation".

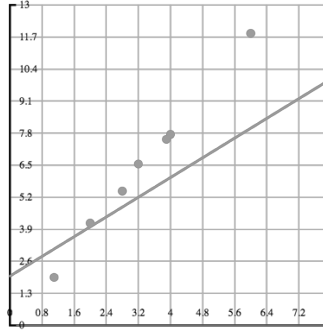
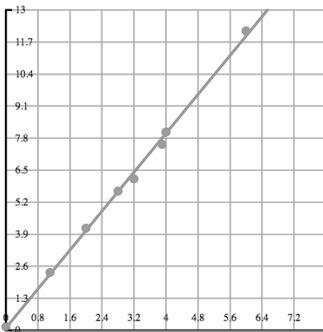
**A**



Correlation:

\_\_\_\_\_

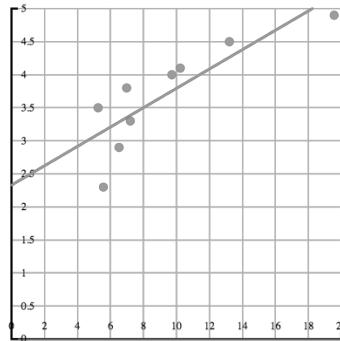
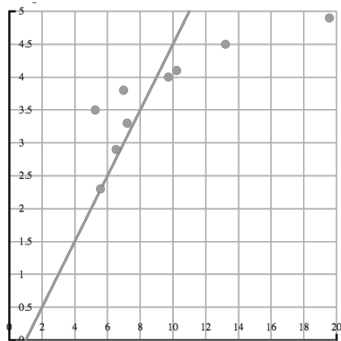
**B**



Correlation:

\_\_\_\_\_

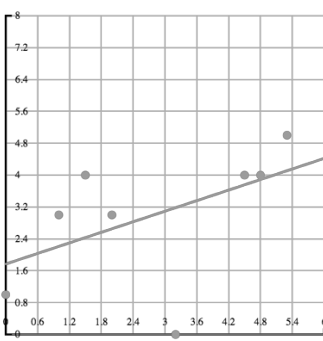
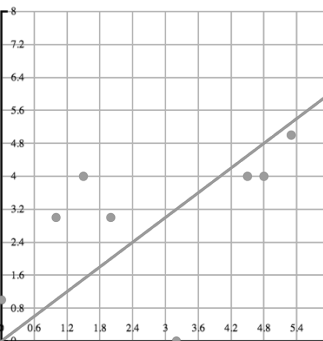
**C**



Correlation:

\_\_\_\_\_

**D**



Correlation:

\_\_\_\_\_

# Regression Analysis in the `animals` Dataset

---

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

---

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

---

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

---

# Regression Analysis in My Dataset

---

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

---

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

---

I performed a linear regression on \_\_\_\_\_, and  
found \_\_\_\_\_ dataset or subset \_\_\_\_\_ correlation between  
\_\_\_\_\_ a weak/strong/moderate, positive/negative, ( $R=$ \_\_\_\_) \_\_\_\_\_ and \_\_\_\_\_. I would predict that  
a 1 \_\_\_\_\_ [x-axis] increase in \_\_\_\_\_ [y-axis] is associated with a \_\_\_\_\_  
\_\_\_\_\_ [x-axis units] 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** – 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 at the shelter longer" means "less likely to be adopted".



# 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?

---

---

---

---

---

---

The animal shelter noticed a large increase in pet adoptions between Christmas and Valentines 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?

---

---

---

---

---

---



# 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 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?

---

---

---

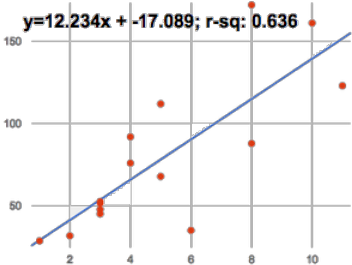
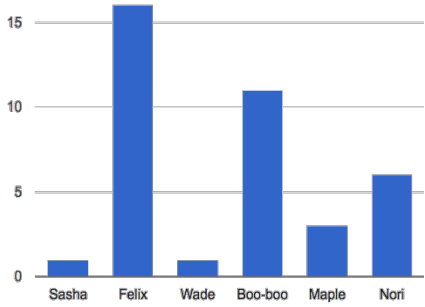
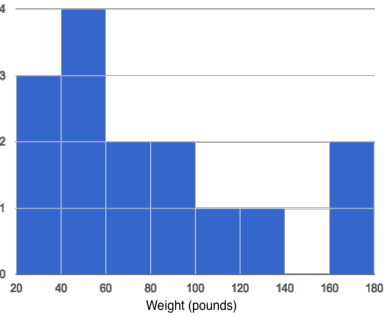
---

---

---

# Fake News!

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

	Data	Claim	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'."	
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."	
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."	
4	 <p>Bar Chart of Pet Ages</p>	"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	 <p>Weight (pounds)</p>	"According to this histogram, most animals weigh between 40 and 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."	

# **Blank Recipes and References**

# 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
```

# Design Recipes

---

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

```
# _____
```

**examples:**

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

---

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

```
# _____
```

**examples:**

```
    _____ ( _____ ) is _____  
    _____ ( _____ ) is _____  
end  
fun _____ ( _____ ) : _____  
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>→</code>	<code>Image</code>
<code>circle</code>	<code>:: (radius :: Number, style :: String, color :: String)</code>	<code>→</code>	<code>Image</code>
<code>star</code>	<code>:: (radius :: Number, style :: String, color :: String)</code>	<code>→</code>	<code>Image</code>
<code>rectangle</code>	<code>:: (width :: Num, height :: Num, style :: Str, color :: Str)</code>	<code>→</code>	<code>Image</code>
<code>ellipse</code>	<code>:: (width :: Num, height :: Num, style :: Str, color :: Str)</code>	<code>→</code>	<code>Image</code>
<code>square</code>	<code>:: (size-length :: Number, style :: String, color :: String)</code>	<code>→</code>	<code>Image</code>
<code>text</code>	<code>:: (str :: String, size :: Number, color :: String)</code>	<code>→</code>	<code>Image</code>
<code>overlay</code>	<code>:: (img1 :: Image, img2 :: Image)</code>	<code>→</code>	<code>Image</code>
<code>rotate</code>	<code>:: (degree :: Number, img :: Image)</code>	<code>→</code>	<code>Image</code>
<code>scale</code>	<code>:: (factor :: Number, img :: Image)</code>	<code>→</code>	<code>Image</code>
<code>string-repeat</code>	<code>:: (text :: String, repeat :: Number)</code>	<code>→</code>	<code>String</code>
<code>string-contains</code>	<code>:: (text :: String, search-for :: String)</code>	<code>→</code>	<code>Boolean</code>
<code>num-sqr</code>	<code>:: (n :: Number)</code>	<code>→</code>	<code>Number</code>
<code>num-sqrt</code>	<code>:: (n :: Number)</code>	<code>→</code>	<code>Number</code>
<code>num-min</code>	<code>:: (a :: Number, b :: Number)</code>	<code>→</code>	<code>Number</code>
<code>num-max</code>	<code>:: (a :: Number, b :: Number)</code>	<code>→</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>count</code>	<code>:: (t :: Table, col :: String)</code>	<code>→ Table</code>
<code>&lt;Table&gt;.row-n</code>	<code>:: (n :: Number)</code>	<code>→ Row</code>
<code>&lt;Table&gt;.order-by</code>	<code>:: (col :: String, increasing :: Boolean)</code>	<code>→ Table</code>
<code>&lt;Table&gt;.filter</code>	<code>:: (test :: (Row → Boolean) )</code>	<code>→ Table</code>
<code>&lt;Table&gt;.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&lt;Number&gt;</code>
<code>bar-chart</code>	<code>:: (t :: Table, col :: String)</code>	<code>→ Image</code>
<code>pie-chart</code>	<code>:: (t :: Table, col :: String)</code>	<code>→ Image</code>
<code>bar-chart-row</code>	<code>:: (t :: Table, labels :: String, values :: String)</code>	<code>→ Image</code>
<code>pie-chart-row</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>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>