

# Lab 3: Functions and Visualizations

Welcome to Lab 3! In this lab, we'll learn about functions, table methods such as `apply`, and how to generate visualizations!

Recommended Reading:

- [Applying a Function to a Column](#)
- [Visualizations](#)

First, set up the notebook by running the cell below.

```
In [7]: import numpy as np
        from datascience import *

        # These lines set up graphing capabilities.
        import matplotlib
        %matplotlib inline
        import matplotlib.pyplot as plt
        plt.style.use('fivethirtyeight')
        import warnings
        warnings.simplefilter('ignore', FutureWarning)
```

## 1. Defining functions

Let's write a very simple function that converts a proportion to a percentage by multiplying it by 100. For example, the value of `to_percentage(.5)` should be the number 50 (no percent sign).

A function definition has a few parts.

```
def
```

It always starts with `def` (short for **define**):

```
    def
```

**Name**

Next comes the name of the function. Like other names we've defined, it can't start with a number or contain spaces. Let's call our function `to_percentage`:

```
    def to_percentage
```

**Signature**

Next comes something called the *signature* of the function. This tells Python how many arguments your function should have, and what names you'll use to refer to those arguments in the function's code. A function can have any number of arguments (including 0!).

`to_percentage` should take one argument, and we'll call that argument `proportion` since it should be a proportion.

```
def to_percentage(proportion)
```

If we want our function to take more than one argument, we add a comma between each argument name. Note that if we had zero arguments, we'd still place the parentheses () after than name.

We put a colon after the signature to tell Python it's over. If you're getting a syntax error after defining a function, check to make sure you remembered the colon!

```
def to_percentage(proportion):
```

## Documentation

Functions can do complicated things, so you should write an explanation of what your function does. For small functions, this is less important, but it's a good habit to learn from the start. Conventionally, Python functions are documented by writing an **indented** triple-quoted string:

```
def to_percentage(proportion):  
    """Converts a proportion to a percentage."""
```

## Body

Now we start writing code that runs when the function is called. This is called the *body* of the function and every line **must be indented with a tab**. Any lines that are *not* indented and left-aligned with the def statement is considered outside the function.

Some notes about the body of the function:

- We can write code that we would write anywhere else.
- We use the arguments defined in the function signature. We can do this because we assume that when we call the function, values are already assigned to those arguments.
- We generally avoid referencing variables defined *outside* the function. If you would like to reference variables outside of the function, pass them through as arguments!

Now, let's give a name to the number we multiply a proportion by to get a percentage:

```
def to_percentage(proportion):
    """Converts a proportion to a percentage."""
    factor = 100
```

```
    return
```

The special instruction `return` is part of the function's body and tells Python to make the value of the function call equal to whatever comes right after `return`. We want the value of `to_percentage(.5)` to be the proportion .5 times the factor 100, so we write:

```
def to_percentage(proportion):
    """Converts a proportion to a percentage."""
    factor = 100
    return proportion * factor
```

`return` only makes sense in the context of a function, and **can never be used outside of a function**. `return` is always the last line of the function because Python stops executing the body of a function once it hits a `return` statement.

*Note:* `return` inside a function tells Python what value the function evaluates to. However, there are other functions, like `print`, that have no `return` value. For example, `print` simply prints a certain value out to the console.

`return` and `print` are **very** different.

**\*\*Question 1.1.\*\*** Define `to_percentage` in the cell below. Call your function to convert the proportion .2 to a percentage. Name that percentage `twenty_percent`.

```
In [2]: def to_percentage(proportion):
        """Converts a proportion to a percentage."""
        factor = 100
        twenty_percent = proportion * factor
        return twenty_percent

        twenty_percent = to_percentage(.35)
        twenty_percent
```

Out[2]: 35.0

```
In [3]: # TEST
        # Make sure your function has the proper syntax!
        to_percentage(.35) == 35.0
```

Out[3]: True

Like you've done with built-in functions in previous labs (`max`, `abs`, etc.), you can pass in named values as arguments to your function.

**\*\*Question 1.2.\*\*** Use `to_percentage` again to convert the proportion named `a_proportion` (defined below) to a percentage called `a_percentage`.

*Note:* You don't need to define `to_percentage` again! Like other named values, functions stick around after you define them.

```
In [4]: a_proportion = 2**(.5) / 2
a_percentage = a_proportion * 100
a_percentage
```

```
Out[4]: 70.71067811865476
```

```
In [5]: # TEST
a_percentage == 70.71067811865476
```

```
Out[5]: True
```

Here's something important about functions: the names assigned *within* a function body are only accessible within the function body. Once the function has returned, those names are gone. So even if you created a variable called `factor` and defined `factor = 100` inside of the body of the `to_percentage` function and then called `to_percentage`, `factor` would not have a value assigned to it outside of the body of `to_percentage`:

```
In [6]: # You should see an error when you run this. (If you don't, you might
# have defined factor somewhere above.)
factor
```

```
-----
NameError                                Traceback (most recent call last)
Input In [6], in <module>
      1 # You should see an error when you run this. (If you don't, you might
      2 # have defined factor somewhere above.)
----> 3 factor

NameError: name 'factor' is not defined
```

As we've seen with built-in functions, functions can also take strings (or arrays, or tables) as arguments, and they can return those things, too.

**\*\*Question 1.3.\*\*** Define a function called `disemvowel`. It should take a single string as its argument. (You can call that argument whatever you want.) It should return a copy of that string, but with all the characters that are vowels removed. (In English, the vowels are the characters "a", "e", "i", "o", and "u".) You can use as many lines inside of the function to do this as you'd like.

*Hint:* To remove all the "a"s from a string, you can use `that_string.replace("a", "")`. The `.replace` method for strings returns a new string, so you can call `replace` multiple times, one after the other.

```
In [8]: def disemvowel(a_string):
a_string = a_string.replace('a', '')
e_string = a_string.replace('e', '')
```

```
i_string = e_string.replace('i', '')
o_string = i_string.replace('o', '')
u_string = o_string.replace('u', '')
```

```
return u_string
```

```
# An example call to your function. (It's often helpful to run
# an example call from time to time while you're writing a function,
# to see how it currently works.)
```

```
disemvowel("Can you read this without vowels?")
```

Out[8]: 'Cn y rd ths wtht vwls?'

In [9]: # TEST

```
disemvowel("Datascience rules!") == "Dtscnc rls!"
```

Out[9]: True

### Calls on calls on calls

Just as you write a series of lines to build up a complex computation, it's useful to define a series of small functions that build on each other. Since you can write any code inside a function's body, you can call other functions you've written.

If a function is like a recipe, defining a function in terms of other functions is like having a recipe for cake telling you to follow another recipe to make the frosting, and another to make the jam filling. This makes the cake recipe shorter and clearer, and it avoids having a bunch of duplicated frosting recipes. It's a foundation of productive programming.

For example, suppose you want to count the number of characters *that aren't* vowels in a piece of text. One way to do that is this to remove all the vowels and count the size of the remaining string.

**Question 1.4.** Write a function called `num_non_vowels`. It should take a string as its argument and return a number. That number should be the number of characters in the argument string that aren't vowels. You should use the `disemvowel` function you wrote above inside of the `num_non_vowels` function.

*Hint:* The function `len` takes a string as its argument and returns the number of characters in it.

In [10]: **def** num\_non\_vowels(a\_string):

```
    """The number of characters in a string, minus the vowels."""
```

```
    return len(disemvowel(a_string))
```

```
# Try calling your function yourself to make sure the output is what
# you expect.
```

In [11]: # TEST

```
num_non_vowels("Go bears!") == 6
```

Out[11]: True

Functions can also encapsulate code that *displays output* instead of computing a value. For example, if you call `print` inside a function, and then call that function, something will get printed.

The `movies_by_year` dataset in the textbook has information about movie sales in recent years. Suppose you'd like to display the year with the 5th-highest total gross movie sales, printed in a human-readable way. You might do this:

```
In [12]: movies_by_year = Table.read_table("movies_by_year.csv")
rank = 5
fifth_from_top_movie_year = movies_by_year.sort("Total Gross", descending=True)
print("Year number", rank, "for total gross movie sales was:", fifth_from_top_n

Year number 5 for total gross movie sales was: 2010
```

After writing this, you realize you also wanted to print out the 2nd and 3rd-highest years. Instead of copying your code, you decide to put it in a function. Since the rank varies, you make that an argument to your function.

**\*\*Question 1.5\*\*** Write a function called `print_kth_top_movie_year`. It should take a single argument, the rank of the year (like 2, 3, or 5 in the above examples). It should print out a message like the one above.

*Note:* Your function shouldn't have a `return` statement.

```
In [13]: def print_kth_top_movie_year(k):
    kth_top_movie_year = movies_by_year.sort('Total Gross', descending=True).co
    print('Year number', k, 'for total gross movie sales was:', kth_top_movie_y

# Example calls to your function:
print_kth_top_movie_year(2)
print_kth_top_movie_year(3)

Year number 2 for total gross movie sales was: 2013
Year number 3 for total gross movie sales was: 2012
```

```
In [14]: # TEST
print_kth_top_movie_year(4) # Year number 4 for total gross movie sales was: 20

Year number 4 for total gross movie sales was: 2009
```

## print is not the same as return

The `print_kth_top_movie_year(k)` function prints the total gross movie sales for the year that was provided! However, since we did not return any value in this function, we can not use it after we call it. Let's look at an example of another function that prints a value but does not return it.

```
In [15]: def print_number_five():
    print(5)
    return 5
```

```
In [16]: print_number_five()
```

```
5
```

```
Out[16]: 5
```

However, if we try to use the output of `print_number_five()`, we see that the value `5` is printed but we get a `TypeError` when we try to add the number 2 to it!

```
In [17]: print_number_five_output = print_number_five()  
print_number_five_output + 2
```

```
5
```

```
Out[17]: 7
```

It may seem that `print_number_five()` is returning a value, 5. In reality, it just displays the number 5 to you without giving you the actual value! If your function prints out a value without returning it and you try to use that value, you will run into errors, so be careful!

Explain to your neighbor how you might add a line of code to the `print_number_five` function (after `print(5)`) so that the code `print_number_five_output + 5` would result in the value `10`, rather than an error.

## 2. Functions and CEO Incomes

In this question, we'll look at the 2015 compensation of CEOs at the 100 largest companies in California. The data was compiled from a [Los Angeles Times analysis](#), and ultimately came from [filings](#) mandated by the SEC from all publicly-traded companies. Two companies have two CEOs, so there are 102 CEOs in the dataset.

We've copied the raw data from the LA Times page into a file called

`raw_compensation.csv`. (The page notes that all dollar amounts are in **millions of dollars**.)

```
In [18]: raw_compensation = Table.read_table('raw_compensation.csv')  
raw_compensation
```

Out [18]:

Rank	Name	Company (Headquarters)	Total Pay	% Change	Cash Pay	Equity Pay	Other Pay	Ratio of CEO pay to average industry worker pay
1	Mark V. Hurd*	Oracle (Redwood City)	\$53.25	(No previous year)	\$0.95	\$52.27	\$0.02	362
2	Safra A. Catz*	Oracle (Redwood City)	\$53.24	(No previous year)	\$0.95	\$52.27	\$0.02	362
3	Robert A. Iger	Walt Disney (Burbank)	\$44.91	-3%	\$24.89	\$17.28	\$2.74	477
4	Marissa A. Mayer	Yahoo! (Sunnyvale)	\$35.98	-15%	\$1.00	\$34.43	\$0.55	342
5	Marc Benioff	salesforce.com (San Francisco)	\$33.36	-16%	\$4.65	\$27.26	\$1.45	338
6	John H. Hammergren	McKesson (San Francisco)	\$24.84	-4%	\$12.10	\$12.37	\$0.37	222
7	John S. Watson	Chevron (San Ramon)	\$22.04	-15%	\$4.31	\$14.68	\$3.05	183
8	Jeffrey Weiner	LinkedIn (Mountain View)	\$19.86	27%	\$2.47	\$17.26	\$0.13	182
9	John T. Chambers**	Cisco Systems (San Jose)	\$19.62	19%	\$5.10	\$14.51	\$0.01	170
10	John G. Stumpf	Wells Fargo (San Francisco)	\$19.32	-10%	\$6.80	\$12.50	\$0.02	256

... (92 rows omitted)

We want to compute the average of the CEOs' pay. Try running the cell below.

```
In [19]: np.average(raw_compensation.column("Total Pay"))
```



```

-----
UFuncTypeError                                Traceback (most recent call last)
Input In [19], in <module>
----> 1 np.average(raw_compensation.column("Total Pay"))

File <__array_function__ internals>:180, in average(*args, **kwargs)

File /Library/Frameworks/Python.framework/Versions/3.10/lib/python3.10/site-packages/numpy/lib/function_base.py:495, in average(a, axis, weights, returned)
    492 a = np.asanyarray(a)
    494 if weights is None:
--> 495     avg = a.mean(axis)
    496     scl = avg.dtype.type(a.size/avg.size)
    497 else:

File /Library/Frameworks/Python.framework/Versions/3.10/lib/python3.10/site-packages/numpy/core/_methods.py:179, in _mean(a, axis, dtype, out, keepdims, where)
    176     dtype = mu.dtype('f4')
    177     is float16 result = True
--> 179 ret = umr_sum(arr, axis, dtype, out, keepdims, where=where)
    180 if isinstance(ret, mu.ndarray):
    181     ret = um.true_divide(
    182         ret, rcount, out=ret, casting='unsafe', subok=False)

UFuncTypeError: ufunc 'add' did not contain a loop with signature matching types (dtype('<U7'), dtype('<U7')) -> None

```

You should see a TypeError. Let's examine why this error occurred by looking at the values in the `Total Pay` column.

**Question 2.1.** Use the `type` function and set `total_pay_type` to the type of the first value in the "Total Pay" column.

```
In [20]: total_pay_type = type(raw_compensation.column("Total Pay").item(0))
total_pay_type
```

```
Out[20]: str
```

```
In [21]: # TEST
# Make sure you are examining the values in the column, not the column itself
import numpy
total_pay_type != numpy.ndarray
```

```
Out[21]: True
```

```
In [22]: # TEST
'str' in str(total_pay_type)
```

```
Out[22]: True
```

```
In [23]: # TEST
# Make sure to call the type function on a value in the column
total_pay_type != int
```

```
Out[23]: True
```

**Question 2.2.** You should have found that the values in the `Total Pay` column are strings. It doesn't make sense to take the average of string values, so we need to convert them to numbers if we want to do this. Extract the first value in `Total Pay`. It's Mark Hurd's pay in 2015, in *millions* of dollars. Call it `mark_hurd_pay_string`.

```
In [24]: mark_hurd_pay_string = raw_compensation.column("Total Pay").item(0)
mark_hurd_pay_string
```

```
Out[24]: '$53.25 '
```

```
In [25]: # TEST
str(mark_hurd_pay_string) == '$53.25 '
```

```
Out[25]: True
```

**Question 2.3.** Convert `mark_hurd_pay_string` to a number of *dollars*.

Some hints, as this question requires multiple steps:

- The string method `strip` will be useful for removing the dollar sign; it removes a specified character from the start or end of a string. For example, the value of `"100%".strip("%")` is the string `"100"`.
- You'll also need the function `float`, which converts a string that looks like a number to an actual number.
- Finally, remember that the answer should be in dollars, not millions of dollars.

```
In [26]: mark_hurd_pay = int(float(mark_hurd_pay_string.strip('$')) * 1000000)
mark_hurd_pay
```

```
Out[26]: 53250000
```

```
In [27]: # TEST
# Your answer should be a number
type(mark_hurd_pay) != str
```

```
Out[27]: True
```

```
In [28]: # TEST
# Don't forget to give your answer in dollars, not millions of
# Dollars!
mark_hurd_pay != 5325
```

```
Out[28]: True
```

```
In [29]: # TEST
# Don't forget to give your answer in dollars, not millions of
# Dollars!
mark_hurd_pay == 53250000
```

```
Out[29]: True
```

To compute the average pay, we need to do this for every CEO. But that looks like it would

involve copying this code 102 times.

This is where functions come in. First, we'll define a new function, giving a name to the expression that converts "total pay" strings to numeric values. Later in this lab, we'll see the payoff: we can call that function on every pay string in the dataset at once.

The next section of this lab explains how to define a function. For now, just fill in the ellipses in the cell below.

**\*\*Question 2.4.\*\*** Copy the expression you used to compute `mark_hurd_pay`, and use it as the return expression of the function below. But make sure you replace the specific `mark_hurd_pay_string` with the generic `pay_string` name specified in the first line in the `def` statement.

*Hint:* When dealing with functions, you should generally not be referencing any variable outside of the function. Usually, you want to be working with the arguments that are passed into it, such as `pay_string` for this function. If you're using `mark_hurd_pay_string` within your function, you're referencing an outside variable!

```
In [30]: def convert_pay_string_to_number(pay_string):
          """Converts a pay string like '$100' (in millions) to a number of dollars.
          return int(float(pay_string.strip('$')) * 1000000)
```

```
In [31]: # TEST
          convert_pay_string_to_number("$100 ") == 100000000.0
```

Out[31]: True

```
In [32]: # TEST
          convert_pay_string_to_number("$23 ") == 23000000.0
```

Out[32]: True

Running that cell doesn't convert any particular pay string. Instead, it creates a function called `convert_pay_string_to_number` that can convert *any* string with the right format to a number representing millions of dollars.

We can call our function just like we call the built-in functions we've seen. It takes one argument -- a string -- and it returns a float.

```
In [33]: convert_pay_string_to_number('$42')
```

Out[33]: 42000000

```
In [34]: convert_pay_string_to_number(mark_hurd_pay_string)
```

Out[34]: 53250000

```
In [35]: # We can also compute Safra Catz's pay in the same way:
          convert_pay_string_to_number(raw_compensation.where("Name", are.containing("Safra"))
```

```
Out [35]: 53240000
```

So, what have we gained by defining the `convert_pay_string_to_number` function? Well, without it, we'd have to copy the code `10**6 * float(some_pay_string.strip("$"))` each time we wanted to convert a pay string. Now we just call a function whose name says exactly what it's doing.

### 3. applying functions

Defining a function is a lot like giving a name to a value with `=`. In fact, a function is a value just like the number 1 or the text "data"!

For example, we can make a new name for the built-in function `max` if we want:

```
In [36]: our_name_for_max = max
         our_name_for_max(2, 6)
```

```
Out [36]: 6
```

The old name for `max` is still around:

```
In [37]: max(2, 6)
```

```
Out [37]: 6
```

Try just writing `max` or `our_name_for_max` (or the name of any other function) in a cell, and run that cell. Python will print out a (very brief) description of the function.

```
In [38]: max
```

```
Out [38]: <function max>
```

Now try writing `?max` or `?our_name_for_max` (or the name of any other function) in a cell, and run that cell. A information box should show up at the bottom of your screen a longer description of the function

*Note: You can also press Shift+Tab after clicking on a name to see similar information!*

```
In [39]: ?our_name_for_max
```

Let's look at what happens when we set `max` to a non-function value. You'll notice that a `TypeError` will occur when you try calling `max`. Things like integers and strings are not callable. Look out for any functions that might have been renamed when you encounter this type of error

```
In [40]: max = 6
         max(2, 6)
```

```
-----
TypeError                                Traceback (most recent call last)
Input In [40], in <module>
      1 max = 6
----> 2 max(2, 6)

TypeError: 'int' object is not callable
```

```
In [41]: # This cell resets max to the built-in function. Just run this cell, don't cha
import builtins
max = builtins.max
```

Why is this useful? Since functions are just values, it's possible to pass them as arguments to other functions. Here's a simple but not-so-practical example: we can make an array of functions.

```
In [42]: make_array(max, np.average, are.equal_to)
```

```
Out[42]: array([<built-in function max>, <function average at 0x112c4ab00>,
               <function are.equal_to at 0x1255623b0>], dtype=object)
```

**Question 3.1.** Make an array containing any 3 other functions you've seen. Call it `some_functions`.

```
In [43]: some_functions = make_array(max, np.average, are.equal_to)
some_functions
```

```
Out[43]: array([<built-in function max>, <function average at 0x112c4ab00>,
               <function are.equal_to at 0x1255623b0>], dtype=object)
```

```
In [44]: # TEST
len(some_functions)==3
```

```
Out[44]: True
```

```
In [45]: # TEST
# The first thing in your array may not be a function.
callable(some_functions.item(0))
```

```
Out[45]: True
```

```
In [46]: # TEST
# The second thing in your array may not be a function.
callable(some_functions.item(1))
```

```
Out[46]: True
```

```
In [47]: # TEST
# The third thing in your array may not be a function.
callable(some_functions.item(2))
```

```
Out[47]: True
```

Working with functions as values can lead to some funny-looking code. For example, see if you can figure out why the following code works. Check your explanation with a neighbor or

a staff member.

```
In [48]: make_array(max, np.average, are.equal_to).item(0)(4, -2, 7)
```

```
Out[48]: 7
```

A more useful example of passing functions to other functions as arguments is the table method `apply`.

`apply` calls a function many times, once on *each* element in a column of a table. It produces an *array* of the results. Here we use `apply` to convert every CEO's pay to a number, using the function you defined:

```
In [49]: raw_compensation.apply(convert_pay_string_to_number, "Total Pay")
```

```
Out[49]: array([53250000, 53240000, 44910000, 35980000, 33360000, 24840000,
        22040000, 19860000, 19620000, 19320000, 18760000, 18610000,
        18360000, 18090000, 17100000, 16629999, 16329999, 16140000,
        16100000, 16020000, 15100000, 14980000, 14630000, 14510000,
        14440000, 14360000, 14310000, 14090000, 14000000, 13670000,
        12340000, 12200000, 12180000, 12130000, 12050000, 11840000,
        11710000, 11630000, 11160000, 11110000, 11110000, 10730000,
        10500000, 10430000, 10370000, 10280000, 10270000, 10180000,
        10160000, 9970000, 9960000, 9860000, 9740000, 9420000,
        9390000, 9220000, 9060000, 9030000, 8860000, 8760000,
        8570000, 8380000, 8359999, 8350000, 8230000, 7860000,
        7700000, 7580000, 7510000, 7230000, 7210000, 7120000,
        6880000, 6770000, 6640000, 6560000, 6140000, 5920000,
        5900000, 5890000, 5730000, 5420000, 5040000, 4920000,
        4920000, 4470000, 4250000, 4080000, 3930000, 3720000,
        2880000, 2830000, 2820000, 2450000, 1790000, 1680000,
        1530000, 940000, 810000, 70000, 40000, 0])
```

Here's an illustration of what that did:



Note that we didn't write

```
raw_compensation.apply(convert_pay_string_to_number(), "Total Pay") or
raw_compensation.apply(convert_pay_string_to_number("Total Pay"))
```

We just passed the name of the function, with no parentheses, to `apply`, because all we want to do is let `apply` know the name of the function we'd like to use and the name of the column we'd like to use it on. `apply` will then call the function `convert_pay_string_to_number` on each value in the column for us!

**\*\*Question 3.2.\*\*** Using `apply`, make a table that's a copy of `raw_compensation` with one additional column called `Total Pay ($)`. That column should contain the result of applying `convert_pay_string_to_number` to the `Total Pay` column (as we did above). Call the new table `compensation`.

```
In [50]: compensation = raw_compensation.with_column(
        "Total Pay ($)",
        raw_compensation.apply(convert_pay_string_to_number, "Total Pay")
```

)  
compensation

Out[50]:

Rank	Name	Company (Headquarters)	Total Pay	% Change	Cash Pay	Equity Pay	Other Pay	Ratio of CEO pay to average industry worker pay	Total P (
1	Mark V. Hurd*	Oracle (Redwood City)	\$53.25	(No previous year)	\$0.95	\$52.27	\$0.02	362	5325000
2	Safra A. Catz*	Oracle (Redwood City)	\$53.24	(No previous year)	\$0.95	\$52.27	\$0.02	362	5324000
3	Robert A. Iger	Walt Disney (Burbank)	\$44.91	-3%	\$24.89	\$17.28	\$2.74	477	4491000
4	Marissa A. Mayer	Yahoo! (Sunnyvale)	\$35.98	-15%	\$1.00	\$34.43	\$0.55	342	3598000
5	Marc Benioff	salesforce.com (San Francisco)	\$33.36	-16%	\$4.65	\$27.26	\$1.45	338	3336000
6	John H. Hammergren	McKesson (San Francisco)	\$24.84	-4%	\$12.10	\$12.37	\$0.37	222	2484000
7	John S. Watson	Chevron (San Ramon)	\$22.04	-15%	\$4.31	\$14.68	\$3.05	183	2204000
8	Jeffrey Weiner	LinkedIn (Mountain View)	\$19.86	27%	\$2.47	\$17.26	\$0.13	182	1986000
9	John T. Chambers**	Cisco Systems (San Jose)	\$19.62	19%	\$5.10	\$14.51	\$0.01	170	1962000
10	John G. Stumpf	Wells Fargo (San Francisco)	\$19.32	-10%	\$6.80	\$12.50	\$0.02	256	1932000

... (92 rows omitted)

```
In [51]: # TEST
# You either didn't add the 'Total Pay ($)' column, or you mislabeled it
'Total Pay ($)' in compensation.column_labels
```

Out[51]: True

```
In [52]: # TEST
# You have the column in your table,
# but the values may be wrong
t = compensation.sort('Total Pay ($)', descending = True)
t.column('Total Pay ($').item(0) == 53250000.0
```

Out[52]: True

Now that we have all the pays as numbers, we can learn more about them through computation.

**Question 3.3.** Compute the average total pay of the CEOs in the dataset.

```
In [53]: average_total_pay = np.average(compensation.column('Total Pay ($)'))  
average_total_pay
```

```
Out[53]: 11445294.088235294
```

```
In [54]: # TEST  
import math  
math.isclose(average_total_pay, 11445294.11764706, rel_tol = 0.1)
```

```
Out[54]: True
```

**Question 3.4.** Companies pay executives in a variety of ways: in cash, by granting stock or other equity in the company, or with ancillary benefits (like private jets). Compute the proportion of each CEO's pay that was cash. (Your answer should be an array of numbers, one for each CEO in the dataset.)

*Note:* When you answer this question, you'll encounter a red box appearing below your code cell that says something like `RuntimeWarning: invalid value encountered in true_divide`. Don't worry too much about the message. Warnings are raised by Python when it encounters an unusual condition in your code, but the condition is not severe enough to warrant throwing an error.

The warning below is Python's cryptic way of telling you that you're dividing a number by zero. If you extract the values in `Total Pay ($)` as an array, you'll see that the last element is 0.

```
In [55]: cash_proportion = compensation.apply(convert_pay_string_to_number, 'Cash Pay')  
cash_proportion
```

```
/var/folders/hc/231thybn1pd8ng7b0dfdzgr40000gn/T/ipykernel_1976/1002381804.py:  
1: RuntimeWarning: invalid value encountered in true_divide  
    cash_proportion = compensation.apply(convert_pay_string_to_number, 'Cash Pa  
y') / compensation.apply(convert_pay_string_to_number, 'Total Pay')
```



```
Out[55]: array([0.01784038, 0.01784373, 0.55421955, 0.02779322, 0.13938849,
0.48711755, 0.19555354, 0.12437059, 0.25993884, 0.35196687,
0.3075693 , 0.22138635, 0.13126362, 0.1708126 , 0.23099415,
0.06734817, 0.13043479, 0.28004957, 0.33229814, 0.15355805,
0.29337748, 0.21829105, 0.31100478, 0.25086147, 0.2299169 ,
0.16991643, 0.31795947, 0.26188786, 0.28357143, 0.15654718,
0.38168558, 0.28934426, 0.20361248, 0.47650453, 0.45643154,
0.36402027, 0.2177626 , 0.24763543, 0.42562724, 0.2610261 ,
0.18361836, 0.1444548 , 0.33333333, 0.10834132, 0.20925747,
0.97276265, 0.22979552, 0.22789784, 0.37893701, 0.25175527,
0.73895582, 0.37018256, 0.2412731 , 0.21337569, 0.20553781,
0.23318872, 0.33664459, 0.3875969 , 0.56094808, 0.11757991,
0.35239207, 0.24462995, 0.24999991, 0.23712575, 0.43377886,
0.31424936, 0.46363636, 0.32585752, 0.24766977, 0.98755187,
0.27184466, 0.96207865, 0.31831395, 0.81979321, 0.23795181,
0.17530488, 0.21172638, 0.37162162, 0.27288136, 0.26994907,
0.55148342, 0.3597786 , 0. , 0.47154472, 0.47154472,
0.29753915, 0.16235294, 0.48529412, 0.46819338, 0.32526882,
0.98958333, 0.61130742, 0.67021277, 0.75510204, 0.50837989,
0.98809524, 0.98039216, 0.9893617 , 0.87654321, 0. ,
1. , nan])
```

```
In [56]: # TEST
# Checking length of cash_proportion
len(cash_proportion) == 102
```

```
Out[56]: True
```

```
In [57]: # TEST
import math
math.isclose(cash_proportion.item(0), 0.01784038, rel_tol = .001)
```

```
Out[57]: True
```

Check out the `% Change` column in `compensation`. It shows the percentage increase in the CEO's pay from the previous year. For CEOs with no previous year on record, it instead says "(No previous year)". The values in this column are *strings*, not numbers, so like the `Total Pay` column, it's not usable without a bit of extra work.

Given your current pay and the percentage increase from the previous year, you can compute your previous year's pay. For example, if your pay is \$120 this year, and that's an increase of 50% from the previous year, then your previous year's pay was  $\frac{\$120}{1 + \frac{50}{100}}$ , or \$80.

**\*\*Question 3.5.\*\*** Create a new table called `with_previous_compensation`. It should be a copy of `compensation`, but with the "(No previous year)" CEOs filtered out, and with an extra column called `2014 Total Pay ($)`. That column should have each CEO's pay in 2014.

*Hint 1:* You can print out your results after each step to make sure you're on the right track.

*Hint 2:* We've provided a structure that you can use to get to the answer. However, if it's confusing, feel free to delete the current structure and approach the problem your own way!

```
In [61]: # Definition to turn percent to number
def percent_string_to_num(percent_string):
    """Converts a percentage string to a number."""
    return int(float(percent_string.strip('%')))

def prev_compensation(total_pay, change):
    if total_pay == 0:
        return 0
    else:
        return total_pay / (1+percent_string_to_num(change) / total_pay)

# Compensation table where there is a previous year
having_previous_year = compensation.where('% Change', are.not_equal_to('(No pre

# Get the percent changes as numbers instead of strings
# We're still working off the table having_previous_year
percent_changes = having_previous_year.apply(percent_string_to_num, '% Change')

# Calculate the previous year's pay
# We're still working off the table having_previous_year
previous_pay = having_previous_year.apply(prev_compensation, 'Total Pay ($)', '

# Put the previous pay column into the having_previous_year table
with_previous_compensation = having_previous_year.with_column('2014 Total Pay (

with_previous_compensation
```

Out [61]:

Rank	Name	Company (Headquarters)	Total Pay	% Change	Cash Pay	Equity Pay	Other Pay	Ratio of CEO pay to average industry worker pay	Total Pa (
3	Robert A. Iger	Walt Disney (Burbank)	\$44.91	-3%	\$24.89	\$17.28	\$2.74	477	4491000
4	Marissa A. Mayer	Yahoo! (Sunnyvale)	\$35.98	-15%	\$1.00	\$34.43	\$0.55	342	3598000
5	Marc Benioff	salesforce.com (San Francisco)	\$33.36	-16%	\$4.65	\$27.26	\$1.45	338	3336000
6	John H. Hammergren	McKesson (San Francisco)	\$24.84	-4%	\$12.10	\$12.37	\$0.37	222	2484000
7	John S. Watson	Chevron (San Ramon)	\$22.04	-15%	\$4.31	\$14.68	\$3.05	183	2204000
8	Jeffrey Weiner	LinkedIn (Mountain View)	\$19.86	27%	\$2.47	\$17.26	\$0.13	182	1986000
9	John T. Chambers**	Cisco Systems (San Jose)	\$19.62	19%	\$5.10	\$14.51	\$0.01	170	1962000
10	John G. Stumpf	Wells Fargo (San Francisco)	\$19.32	-10%	\$6.80	\$12.50	\$0.02	256	1932000
11	John C. Martin**	Gilead Sciences (Foster City)	\$18.76	-1%	\$5.77	\$12.98	\$0.01	117	1876000
13	Shantanu Narayen	Adobe Systems (San Jose)	\$18.36	3%	\$2.41	\$15.85	\$0.09	125	1836000

... (71 rows omitted)

```
In [62]: # TEST
# Make sure to remove the "(No previous year)" CEOs
"(No previous year)" not in with_previous_compensation.column("% Change")
```

Out[62]: True

```
In [63]: # TEST
import math
# You have the column, but some of
# your values may be wrong.
t = with_previous_compensation.sort("2014 Total Pay ($)", descending = True)
value = t.column("2014 Total Pay ($)").item(0)
math.isclose(value, 67700000.0, rel_tol = 1000)
```

Out[63]: True

```
In [64]: # TEST
# You have the column, but your number of rows is off
with_previous_compensation.num_rows==81
```

Out[64]: True

**\*\*Question 3.6.\*\*** What was the average pay of these CEOs in 2014?

```
In [65]: average_pay_2014 = np.average(with_previous_compensation.column('2014 Total Pay'))
         average_pay_2014
```

```
Out[65]: 10787401.703941908
```

```
In [66]: # TEST
         import math
         math.isclose(average_pay_2014, 11649176.115603436, rel_tol = 10)
```

```
Out[66]: True
```

### Why is `apply` useful?

For operations like arithmetic, or the functions in the NumPy library, you don't need to use `apply`, because they automatically work on each element of an array. But there are many things that don't. The string manipulation we did in today's lab is one example. Since you can write any code you want in a function, `apply` gives you total control over how you operate on data.

## 4. Histograms

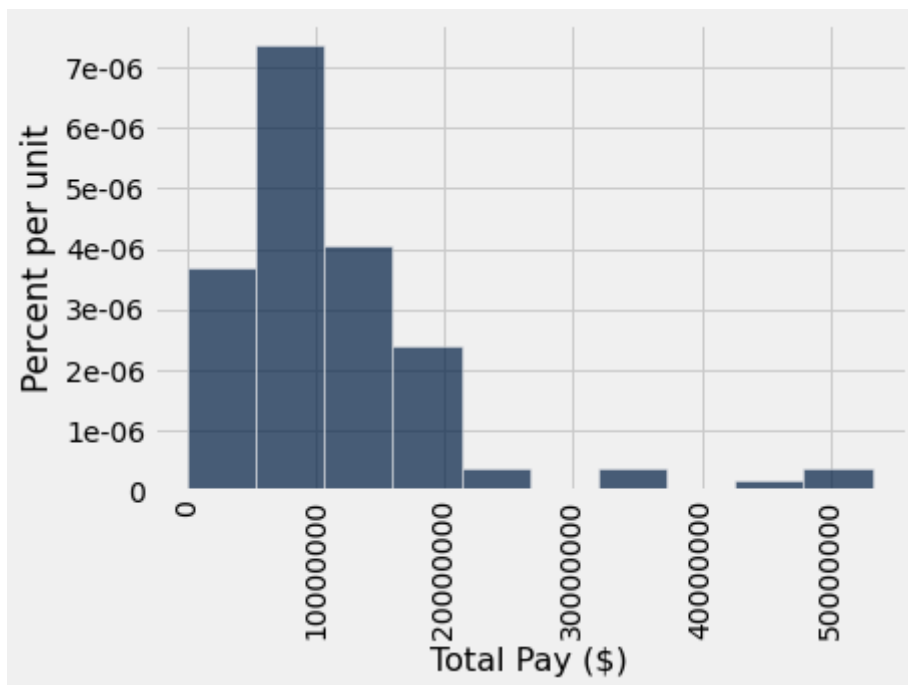
Earlier, we computed the average pay among the CEOs in our 102-CEO dataset. The average doesn't tell us everything about the amounts CEOs are paid, though. Maybe just a few CEOs make the bulk of the money, even among these 102.

We can use a *histogram* method to display the *distribution* of a set of numbers. The table method `hist` takes a single argument, the name of a column of numbers. It produces a histogram of the numbers in that column.

**\*\*Question 4.1.\*\*** Make a histogram of the total pay of the CEOs in `compensation`. Check with your neighbor or a staff member to make sure you have the right plot.

```
In [67]: compensation.hist('Total Pay ($)')
```

```
/Library/Frameworks/Python.framework/Versions/3.10/lib/python3.10/site-packages/datascience/tables.py:5800: UserWarning: FixedFormatter should only be used together with FixedLocator
  axis.set_xticklabels(ticks, rotation='vertical')
```



**\*\*Question 4.2.\*\*** How many CEOs made more than \$30 million in total pay? Find the value using code, then check that the value you found is consistent with what you see in the histogram.

*Hint:* Use the table method `where` and the property `num_rows`.

```
In [68]: num_ceos_more_than_30_million_2 = compensation.where('Total Pay ($)', are.above
num_ceos_more_than_30_million_2
```

Out[68]: 5

```
In [69]: # TEST
num_ceos_more_than_30_million_2 == 5
```

Out[69]: True

Great job! You're finished with lab 3! Be sure to...

- **run all the tests,**
- **print the notebook as a PDF,**
- **and submit both the notebook and the PDF to Canvas.**