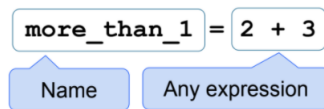


## Statements



- Statements don't have a value; they perform an action.
- An assignment statement changes the meaning of the name to the left of the = symbol.
- The name is bound to the value of the right hand side

## Comparisons

- < and > mean what you expect (less than, greater than)
- <= means "less than or equal; likewise for >=
- == means "equal to"; != means "not equal to"
- Comparing strings compares their alphabetical order

**Arrays:** sequences of the same type that can be manipulated

- Arithmetic and comparisons are applied to each element individually

```
○ make_array(1, 2, 3) > 2 # array([False,
    False, True])
```

- Elementwise operations can be done on arrays of the same size

```
○ make_array(1, 2) * make_array(2, 3) #
    array([2, 6])
```

## Defining a Function

```
def function(arg1, arg2, ...):
    # Body can contain any code
```

Note: Many functions return a value

## for Statements

```
for i in np.arange(12):
    total = total + i
```

- The body is executed **for** every item in a sequence
- The body of the statement can have multiple lines
- The body should do something: assign, sample, etc

## Conditional Statements

```
if <if_expression>:
    <if_body>
elif <elif_expression_0>:
    <elif_body_0>
elif <elif_expression_1>:
    <elif_body_1>
...
else:
    <else_body>
```

## Simulating a Statistic

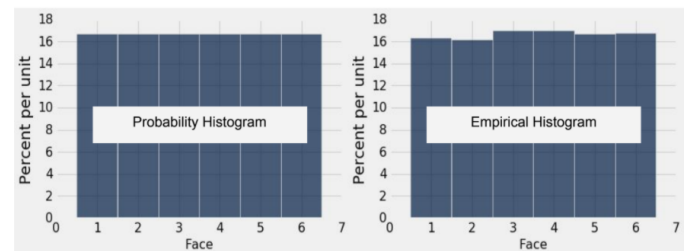
- Define a function to simulate one value of the statistic
- Create an empty collection array
- For each repetition of the process:
  - Call the function to simulate one value
  - Append this value to the collection array
- At the end, all simulated values will be in the collection array

**Total Variation Distance** between two categorical distributions

- For each category, find the difference between the proportions in the two distributions
- Take the absolute value of the differences
- Sum all the absolute values, then divide by 2

A **histogram** has three defining properties:

- Bins are contiguous (though some may be empty) and drawn to scale.
- The area of each bar is the percent of entries in the bin.
- The total area of the histogram is 100%.
- The histogram on the left displays the theoretical probabilities of the number of spots on one roll of a fair die.
- The histogram on the right represents the empirical (or observed) distribution of the numbers of spots on many rolls of a fair die.
- Example of the *Law of Averages*: The more we roll, the more the histogram on the right is likely to resemble the one on the left.



## Finding Probabilities

*Complement Rule:*

$P(\text{event happens}) = 1 - P(\text{event does not happen})$

*Multiplication Rule:*

$P(\text{two events both happen}) = P(\text{first event happens}) * P(\text{second event happens given that the first event happened})$

*Addition Rule:* If an event can happen in only one of two ways, then:

$P(\text{event happens}) = P(\text{happens in the first way}) + P(\text{happens in the second way})$

## p-values

- The *observed significance level* (or *p-value*) of a test is the chance, calculated under the null hypothesis, that the test statistic is equal to the one observed in the sample or more in the direction of the alternative.
- The result of a test of hypotheses is called *statistically significant* if the p-value is less than 5%; *highly statistically significant* if the p-value is less than 1%.

Table Properties and Methods

In the examples in the left column, np refers to the NumPy module, as usual. Everything else is a function, a method, an example of an argument to a function or method, or an example of an object we might call the method on. For example, tbl refers to a table, array refers to an array, and num refers to a number. array.item(0) is an example call for the method item, and in that example, array is the name previously given to some array.

Name	Description	Input	Output
Table()	Create an empty table, usually to extend with data	None	An empty Table
tbl.with_columns(name, values) tbl.with_columns(n1, v1, n2, v2,...)	A table with an additional or replaced column or columns. name is a string for the name of a column, values is an array	1.string: the name of the new column; 2.array: the values in that column	Table: a copy of the original Table with the new columns added
tbl.column(column_name or index)	The values of a column (an array)	string or int: the column_name or index	array: the values in that column
tbl.num_rows	Compute the number of rows in a table	None	int: the number of rows in the table
tbl.num_columns	Compute the number of columns in a table	None	int: the number of columns in the table
tbl.labels	Lists the column labels in a table	None	array: the names of each column (as strings) in the table
tbl.select(col1, col2, ...)	Create a copy of a table with only some of the columns.	string or int: column name(s) or index(es)	Table with the selected columns
tbl.drop(col1, col2, ...)	Create a copy of a table without some of the columns.	string or int: column name(s) or index(es)	Table without the selected columns
tbl.relabeled(old+label, new_label)	Creates a new table, changing the column name specified by the old label to the new label, and leaves the original table unchanged.	1.string: the old column name 2.string: the new column name	Table: a new table
tbl.show(n)	Display n rows of a table. If no argument is specified, defaults to displaying the entire table.	(Optional) int: number of rows you want to display	None: Displays a table with n rows
tbl.sort(column_name or index)	Create a copy of a table sorted by the values in a column. Defaults to ascending order unless descending=True is included.	1.string or int: column index or name 2. (Optional) boolean: descending=True	Table: a copy of the original table with the column sorted
tbl.where(column, predicate)	Create a copy of a table with only the rows that match some predicate See Table.where predicates table.	1.string or int: column index or name 2.are(...) predicate	Table: a copy of the original table with only the rows that match the predicate
tbl.take(row_indices)	A table with only the rows at the given indices. row_indices is either an array of indices or an integer corresponding to one index.	array of ints: the indices of the rows to be included in the Table OR int: the index of the row to be included	Table: a copy of the original table with only the rows at the given indices
tbl.scatter(x_column, y_column)	Draws a scatter plot consisting of one point for each row of the table. Note that x_column and y_column must be strings specifying column names.	1.string or int: name or index of the column on x-axis 2.string or int: name or index of the column on y-axis 3.(Optional) fit line=True	None: Draws a scatter plot
tbl.plot(x_column, y_column) tbl.plot(x_column)	Draws a line graph consisting of one point for each row of the table. If you only specify one column, it will plot the rest of the columns on the y-axis as different colored lines.	1.string or int: name or index of the column on x-axis 2.string or int: name or index of the column on y-axis	None: Draws a line graph
tbl.barh(categories) tbl.barh(categories, values)	Displays a bar chart with bars for each category in a column, with height proportional to the corresponding frequency. values argument unnecessary if table has only a column of categories and a column of values.	1.string or int: name or index of the column with categories 2. (Optional) string or int: name or index of the column with values for corresponding categories	None: Draws a bar chart

<code>tbl.hist(column, unit, bins, group)</code>	Generates a histogram of the numerical values in a column. <code>unit</code> , <code>bins</code> and <code>group</code> are optional arguments, used to label the axes, specify the intervals (bins) and plot separate histograms per group, respectively	1. <code>string</code> or <code>int</code> : name or index of the column with categories 2. (Optional) <code>string</code> : units of x-axis 3. (Optional) <code>array</code> : of ints/floats denoting bin boundaries 4. (Optional) <code>str</code> : name of column to group by	None: Draws a histogram
<code>tbl.bin(column_name or index)</code> <code>tbl.bin(column_name or index, bins)</code>	Groups values into intervals, known as bins. Results in a two-column table that contains the number of rows in each bin. The first column lists the left endpoints of the bins, except in the last row.	1. <code>string</code> or <code>int</code> : column name(s) or index(es) 2. (Optional) <code>array</code> of ints/floats denoting bin boundaries or an <code>int</code> of the number of bins you want	Table: A new table
<code>tbl.apply(function)</code> <code>tbl.apply(function, col1, col2, ...)</code>	Returns an array of values resulting from applying a function to each item in a column.	1. <code>function</code> : function to apply to column 2. (Optional) <code>string</code> or <code>int</code> : name or index of the column to apply function to (if you have multiple columns, the respective column's values will be passed as the corresponding argument to the function), and if there is no argument, your function will be applied to every row (Row object) in <code>tbl</code>	<code>array</code> : contains an element for each value in the original column after applying the function to it.
<code>tbl.group(column_or_columns, collect)</code>	Group rows by unique values or combinations of values in a column(s). Multiple columns must be entered in array or list form. Other values aggregated by count (default) or optional argument <code>collect</code> .	1. <code>string/int</code> or <code>array</code> of strings/ints: column(s) on which to group 2. (Optional) <code>collect</code> : function to aggregate values in cells (defaults to count)	Table: A new table
<code>tbl.pivot(col1, col2, values, collect)</code> <code>tbl.pivot(col1, col2)</code>	A pivot table where each unique value in <code>col1</code> has its own column and each unique value in <code>col2</code> has its own row. Count or aggregate values from a third column, <code>collect</code> with some function. Default <code>values</code> and <code>collect</code> return counts in cells.	1. <code>string</code> or <code>int</code> : name or index of column whose unique values will make up columns of the pivot table 2. <code>string</code> or <code>int</code> : name or index of column whose unique values will make up rows of the pivot table 3. (Optional) <code>string</code> or <code>int</code> : name or index of column containing the values of cell 4. (Optional) <code>function</code> : how the values are collected; e.g. <code>np.mean</code>	Table: A new table
<code>tblA.join(colA, tblB, colB)</code> <code>tblA.join(colA, tblB)</code>	Generate a table with the columns of <code>tblA</code> and <code>tblB</code> , containing rows for all values of a column that appear in both tables. Default <code>colB</code> is <code>colA</code> . <code>colA</code> and <code>colB</code> must be strings specifying column names.	1. <code>string</code> : name of column in <code>tblA</code> with values to join on 2. Table: other Table 3. (Optional) <code>string</code> : if column names are different between Tables, the name of the shared column in <code>tblB</code>	Table: A new table
<code>tbl.sample(n)</code> <code>tbl.sample(n, with_replacement)</code>	A new table where <code>n</code> rows are randomly sampled from the original table; by default, <code>n=tbl.num_rows</code> . Default is with replacement. For sampling without replacement, use argument <code>with_replacement=False</code> . For a non-uniform sample, provide a third argument <code>weights=distribution</code> where <code>distribution</code> is an array or list containing the probability of each row.	1. <code>int</code> : sample size 2. (Optional) <code>with_replacement=True</code>	Table: A new table with <code>n</code> rows
<code>tbl.row(row_index)</code>	Accesses the row of a table by taking the index of the row as its argument. Note that rows are in general not arrays, as their elements can be of different types. However, you can use <code>.item</code> to access a particular element of a row using <code>row.item(label)</code> .	<code>int</code> : row index	Row object with the values of the row and labels of the corresponding columns
<code>tbl.rows</code>	Can use to access all of the rows of a table.	None	Row object made up of all rows as individual row objects

Miscellaneous Functions

These are functions in the datascience library that are used in the course that don't fall into any of the categories above.

Name	Description	Input	Output
<code>sample_proportions(sample_size, model_proportions)</code>	sample size should be an integer, model proportions an array of probabilities that sum up to 1. The function samples sample size objects from the distribution specified by model proportions. It returns an array with the same size as model proportions. Each item in the array corresponds to the proportion of times it was sampled out of the sample size times.	1.int: sample size 2. array: an array of proportions that should sum to 1	array: each item corresponds to the proportion of times that corresponding item was sampled from model proportions in sample size draws. Should sum to 1.
<code>minimize(function)</code>	Returns an array of values such that if each value in the array was passed into function as arguments, it would minimize the output value of function.	function: name of a function that will be minimized	array: An array in which each element corresponds to an argument that minimizes the output of the function. Values in the array are listed based on the order they are passed into the function; the first element in the array is also going to be the first value passed into the function.

Array Functions and Methods

Table.where Predicates

Any of these predicates can be negated by adding not in front of them, e.g. `are.not_equal_to(Z)` or `are.not_containing(S)`.

Function	Description	Function	Description
<code>max(array)</code>	Returns the maximum value of an array	<code>are.equal_to(Z)</code>	Equal to Z
<code>min(array)</code>	Returns the minimum value of an array	<code>are.above(x)</code>	Greater than x
<code>sum(array), np.sum(array)</code>	Returns the sum of the values in an array	<code>are.above_or_equal_to(x)</code>	Greater than or equal to x
<code>abs(num), np.abs(array)</code>	Take the absolute value of a number or each number in an array	<code>are.below(x)</code>	Less than x
<code>np.round(num), np.round(array)</code>	Round number or array of numbers to the nearest integer	<code>are.below_or_equal_to(x)</code>	Less than or equal to x
<code>len(array)</code>	Returns the length (number of elements) of an array	<code>are.between(x,y)</code>	Greater than or equal to x and less than y
<code>make_array(val1, val2, ...)</code>	Makes a numpy array with the values passed in	<code>are.between_or_equal_to(x,y)</code>	Greater than or equal to x and less than or equal to y
<code>np.average(array), np.mean(array)</code>	Returns the mean value of an array	<code>are.contained_in(A)</code>	Is a substring of A (if A is a string) or an element of A (if A is a list/array)
<code>np.std(array)</code>	Returns the standard deviation of an array	<code>are.containing(S)</code>	Contains the string S
<code>np.diff(array)</code>	Returns a new array of size len(arr)-1 with elements equal to the difference between adjacent elements; val 2 - val 1, val 3 - val 2, etc.	<code>are.strictly_between(x,y)</code>	Greater than x and less than y
<code>np.sqrt(array)</code>	Returns an array with the square root of each element		
<code>np.arange(start, stop, step)</code> <code>np.arange(start, stop)</code> <code>np.arange(stop)</code>	An array of numbers starting with start, going up in increments of step, and going up to but excluding stop. When start and/or step are left out, default values are used in their place. Default step is 1; default start is 0		
<code>array.item(index)</code>	Returns the (index-1)-th item in an array (remember Python indices start at 0!)		
<code>np.random.choice(array, n)</code> <code>np.random.choice(array)</code> <code>np.random.choice(array, n, replace)</code>	Picks one (by default) or some number 'n' of items from an array at random. Default is with replacement. For sampling without replacement, use argument <code>replace=False</code> .		
<code>np.count_nonzero(array)</code>	Returns the number of non-zero (or True) elements in an array		
<code>np.append(array, item)</code>	Returns a copy of the input array with item (must be the same type as the other entries in in the array) appended to the end		
<code>percentile(p, array)</code>	Returns the pth percentile of an array		