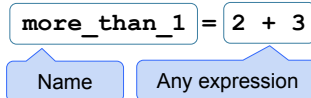


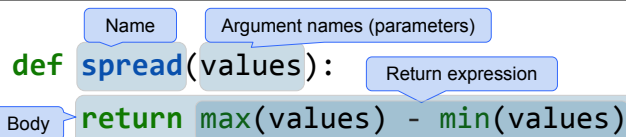
**Statements**

- 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 a value (not an equation)

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

**Arrays** - sequences that can be manipulated easily

- All elements of an array should have the same type
- Arithmetic is applied to each element of an array individually
- Elementwise operations can be done on arrays of the same size



```
for i in np.arange(12):
    print(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: print, assign, hist, 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>
```

**Growth Rate:** the rate of increase per unit time

- After one time unit, a quantity  $x$  growing at rate  $g$  will be  $x * (1 + g)$
- After  $t$  time units, a quantity  $x$  growing at rate  $g$  will be  $x * (1 + g) ** t$
- If **after** and **before** are measurements of the same quantity taken  $t$  time units apart, then the growth rate is  $(\text{after/before}) ** (1/t) - 1$

**Total Variation Distance:** Measures the difference between two categorical distributions

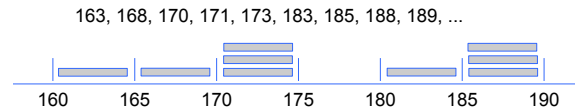
- For each category, compute the difference in proportions between two distributions
- Take the absolute value of each difference
- Sum and divide by 2

**Values in Tables:** Every column of a table is an array.

- **Categorical**
  - May or may not have an ordering
  - Categories are the same or different
  - Allows grouping by value (**group**, **pivot**, **join**)
- **Numerical**
  - Ordered
  - Allows binning by value (**bin**, **hist**)

**Binning** is counting the number of numerical values that lie within ranges, called bins.

- Bins include the lower bound and exclude the upper bound
- Values equal to the upper bound of a bin go into the next bin
- The upper bound of a bin is the lower bound of the next bin



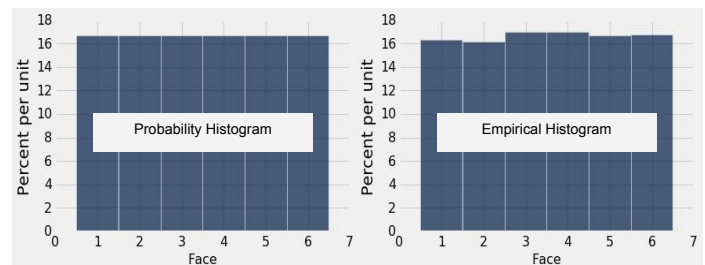
A **histogram** has two defining properties:

- The bins are contiguous (though some might be empty) and are drawn to scale
- The **area** of each bar is equal to the proportion of entries in the bin

Has total area 1 (or 100%)

Vertical axis units: Proportion / Unit on the horizontal axis

- A histogram of proportions of all possible outcomes of a *known* random process is called a *probability histogram*
- A histogram is a summary visualization of a *distribution*
- A histogram of proportions of actual outcomes generated by sampling or actual data is called an *empirical histogram*

**Calculating Probabilities**

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

**Multiplication Rule:**  $P(\text{two events happen}) = P(\text{one happens}) * P(\text{other happens, given the first happened})$

**Addition Rule:**  $P(\text{an event happens}) = P(\text{first way it can happen}) + P(\text{second way it can happen})$  IF it can happen in **ONLY** one of two ways

**Testing a Hypothesis****Step 1: The Hypotheses**

- A test chooses between two views of how data were generated
- *Null hypothesis* proposes that data were generated at random
- *Alternative hypothesis* proposes some effect other than chance

**Step 2: The Test Statistic**

- A value that can be computed for the data and for samples

**Step 3: The Sampling Distribution of the Test Statistic**

- What the test statistic might be if the null hypothesis were true
- Approximate the sampling distribution by an empirical distribution

## Data 8 Midterm Reference Sheet — Page 2

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.

Example function call	Value of a call to the function
<code>max(array); min(array)</code>	Maximum or minimum of a sequence
<code>sum(array)</code>	Sum of all elements in an array
<code>len(array)</code>	Length (num elements) in an array
<code>round(num); np.round(array)</code>	Round number or array of numbers to the nearest integer
<code>abs(num); np.abs(array)</code>	Take the absolute value of number or each number in an array
<code>np.average(array), np.mean(array)</code>	The average of the values in an array
<code>np.arange(start, stop, step)</code> <code>np.arange(start, stop)</code> <code>np.arange(stop)</code>	An array of numbers starting with <code>start</code> , going up in increments of <code>step</code> , and going up to but excluding <code>stop</code> . When <code>start</code> and/or <code>step</code> are left out, default values are used in their place. Default <code>step</code> is 1; default <code>start</code> is 0.
<code>np.count_nonzero(array)</code>	Count the number of non-zero elements in an array ( <code>False</code> counts as zero, <code>True</code> as non-zero)
<code>array.item(index)</code>	The item in the array at some index. <code>array.item(0)</code> is the first item of array.
<code>np.append(array, item)</code>	A copy of the array with <code>item</code> appended to the end.
<code>np.random.choice(array, n)</code> <code>np.random.choice(array)</code>	An array of items selected at random with replacement from an array. Default number of items is 1 if <code>n</code> not specified.
<code>Table()</code>	An empty table.
<code>Table.read_table(filename)</code>	A table with data from a file.
<code>tbl.num_rows</code>	The number of rows in a table.
<code>tbl.num_columns</code>	The number of columns in a table.
<code>tbl.labels</code>	A list of the column labels of a table.
<code>tbl.with_column(name, values)</code> <code>tbl.with_columns(n1, v1, n2, v2...)</code>	A table with an additional or replaced column or columns. <code>name</code> is a string for the name of a column, <code>values</code> is an array.
<code>tbl.column(column_name_or_index)</code>	The values of a column (an array).
<code>tbl.select(col1, col2, ...)</code>	A table with only the selected columns. (Each argument is the name of a column, or a column index.)
<code>tbl.drop(col1, col2, ...)</code>	A table without the selected columns. (Each argument is the name of a column, or a column index.)
<code>tbl.relabelled(old_label, new_label)</code>	A new table with a label changed.
<code>tbl.take(row_indices)</code>	A table with only the rows at the given indices. <code>row_indices</code> is an array of indices.
<code>tbl.sort(column_name_or_index)</code>	A table of rows sorted according to the values in a column (specified by name/index). Default order is ascending. For descending order, use argument <code>descending=True</code> .
<code>tbl.where(column, predicate)</code>	A table of the rows for which the column satisfies some predicate. See “Table.where predicates” below.
<code>tbl.apply(function, column)</code>	An array where a function is applied to each item in a column.
<code>tbl.group(column_or_columns, func)</code>	Group rows by unique values or combinations of values in a column. Other values aggregated by count (default) or optional argument <code>func</code> .
<code>tblA.join(colA, tblB, colB)</code> <code>tblA.join(colA, tblB)</code>	Generate a table with the columns of self and other, containing rows for all values of a column that appear in both tables. Default <code>colB</code> is <code>colA</code> . <code>colA</code> is a string specifying a column name, as is <code>colB</code> .
<code>tbl.pivot(col1, col2, vals, 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>vals</code> and <code>collect</code> return counts in cells.
<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. Default is with replacement. For sampling without replacement, use argument <code>with_replacement=False</code> . For non-uniform sample, provide <code>weights=distribution</code> where <code>distribution</code> is an array containing the probability of each row.
<code>tbl.scatter(x_column, y_column)</code>	Draws a scatter plot consisting of one point for each row of the table.
<code>tbl.barh(categories)</code> <code>tbl.barh(categories, values)</code>	Displays a bar chart with bars for each category in a column, with height proportional to the corresponding frequency. <code>values</code> argument unnecessary if table has only a column of categories and a column of values.
<code>tbl.hist(column, units, bins)</code>	Generates a histogram of the numerical values in a column. <code>units</code> and <code>bins</code> are optional arguments, used to label the axes and group the values into intervals (bins), respectively. Bins have the form <code>[a, b)</code> .

**Operations:** addition `2+3=5`; subtraction `4-2=2`; division `9/2=4.5`  
multiplication `2*3=6`; division remainder `11%3=2`; exponent `2**3=8`

**Data Types:** string `'hello'`; boolean `True, False`;  
int `1, -5`; float `- 2.3, -52.52, 7.9`

Arithmetic with arrays is elementwise:  
`make_array(1,2,3) ** 2 # [1, 4, 9]`

**Table.where predicates** (`x` is a string or number)  
`are.equal_to(x) # [2, 3, 4]`  
`are.above(x) # val > x`  
`are.below(x) # val < x`  
`are.between(x, y) # x <= val < y`