

# Data 8 Final Reference Sheet — Fall 2016

## Python Basics

- **Functions:** called with parentheses (i.e., `np.mean(arr)`), defined with `def` statement

```
def spread(values):
    return max(values) - min(values)
```
- **Comparators:** `==`, `!=`, `>`, `<`, `>=`, `<=` compare two values and return `True` or `False`.
- **Conditionals:** A structure to execute different lines of code based on whether certain conditions are true

```
if <if expression>:
    <if body>
elif <elif expression 0>:
    <elif body 0>
else:
    <else body>
```
- **Loops:** a `for` loop iterates through the elements of a sequence

```
two_three_four = make_array()
for x in make_array(1, 2, 3):
    two_three_four = np.append(two_three_four, x + 1)
```

## Distance Between Two Points

- Two attributes  $x$  and  $y$ :  $\sqrt{(x_0 - x_1)^2 + (y_0 - y_1)^2}$
- Three attributes  $x$ ,  $y$ , and  $z$ :  $\sqrt{(x_0 - x_1)^2 + (y_0 - y_1)^2 + (z_0 - z_1)^2}$
- and so on...

## Probability

Probabilities are between 0 and 1.

- If all outcomes are equally likely, then  $P(\text{event happens}) = \text{proportion of outcomes that make the event happen}$
- $P(\text{event happens}) = 1 - P(\text{the event doesn't happen})$
- Chance that two events A and B both happen =  $P(A \text{ happens}) \times P(B \text{ happens given that A has happened})$
- If event A can happen in exactly one of two ways, then  $P(A) = P(\text{first way}) + P(\text{second way})$
- **Bayes' rule:**  $P(\text{first event happens given the second happens}) = \frac{P(\text{first event happens}) \times P(\text{second event happens given the first happens})}{P(\text{second event happens})}$

## Descriptive Statistics

- **Median:** 50th percentile, where  
 $p$ -th percentile = smallest value on list that is at least as large as  $p\%$  of the values
- **Mean** of 5, 7, 8, 8 =  $(5 + 7 + 8 + 8)/4 = 5 \times 0.25 + 7 \times 0.25 + 8 \times 0.5$
- Mean depends on all the values; smoothing operation; center of gravity of histogram; if histogram is skewed, mean is pulled away from median towards the tail
- The mean of a 0/1 population is the proportion of 1s in the population
- **Standard deviation (SD):** The root mean square of deviations from average.
- The SD of a 0/1 population is less than or equal to 0.5
- **Chebychev's Bound:** No matter what the distribution looks like, the proportion in the range average  $\pm z$  SDs is at least  $1 - \frac{1}{z^2}$ .

- If the distribution is normal, about 68% of values are within the range [average  $\pm 1$  SD] and about 95% of values are within the range [average  $\pm 2$  SDs].
- **Total Variation Distance:** A statistic measuring the difference between categorical distributions. The sum of the absolute value of the differences between proportions in each category, divided by two.
- **Standard units (s.u.):** To convert a value to standard units:  $z = \frac{\text{value} - \text{average}}{\text{SD}}$ .  
To convert standard units to original units:  $\text{value} = z \times \text{SD} + \text{average}$ .
- **Correlation ( $r$ ):**  $r = \text{mean}(x \text{ in s.u.} \times y \text{ in s.u.})$
- **Estimate of  $y = r \times x$ ,** when both variables are measured in standard units
- **Slope of the regression line** =  $r \times \frac{\text{SD of } y}{\text{SD of } x}$
- **Intercept of the regression line** = mean of  $y$  - slope  $\times$  mean of  $x$
- **Residual** = observed  $y$  - regression estimate of  $y$
- **Average of residuals** = 0
- **SD of residuals** =  $\sqrt{1 - r^2} \times \text{SD of } y$

## With-Replacement Random Sample Means

The mean of a random sample with replacement is expected to be the population mean.

- **SD of Sample Mean** =  $\frac{\text{Population SD}}{\sqrt{\text{Sample Size}}}$
- **Square Root Law:** If you multiply sample size by a factor, accuracy goes up by the square root of the factor.
- **Central Limit Theorem:** If a sample is large, and drawn at random with replacement, then, regardless of the distribution of the population, the probability distribution of the sample sum (or of the sample mean) is roughly bell-shaped.

Code examples on the other side of this sheet.

Tables

Tables are a data structure used to store tabular (row and column) data. You may assume that these functions exist in Python. `tbl` refers to a table.

Function/Method	Description
<code>Table()</code>	Creates an empty table
<code>Table.read_table(filename)</code>	Returns a table read in from a CSV file
<code>tbl.labels</code>	Returns an array of a table's column labels
<code>tbl.num_rows</code> , <code>tbl.num_cols</code>	Returns the number of rows (and columns, respectively)
<code>tbl.column(name)</code>	Returns the values of a column (an array)
<code>tbl.with_column(name, values)</code>	Adds or replaces a column to a table
<code>tbl.with_columns(n1,v1,n2,v2)</code>	Adds or replaces multiple columns
<code>tbl.row(i)</code>	Returns the <i>i</i> -th row of a table
<code>tbl.append(row)</code>	Appends a row to a table
<code>tbl.select(col1,col2,...)</code>	A table with only the selected columns
<code>tbl.drop(col1,col2,...)</code>	A table without the specified set of columns
<code>tbl.take(row_indices)</code>	A table with only the rows at the given indices.
<code>tbl.relabeled(old_label, new_label)</code>	Returns a copy of the table with a column label changed. <code>row_indices</code> is an array of indices.
<code>tbl.apply(function, column)</code>	Returns an array where a function is applied to each item in a column
<code>tbl.sort(column_name)</code> <code>tbl.sort(column_name, descending)</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>	Selects rows from a table based on column values. See "Table.where predicates" below.
<code>tblA.join(colA, tblB, colB)</code>	Generate a table with the columns of self and other, containing rows for all values of a column that appear in both tables. <code>colA</code> is a string specifying a column name, as is <code>colB</code> . Takes the first match found in <code>tblB</code>
<code>tbl.group(column, func)</code>	Group rows by unique values in a column. Other values aggregated by count (default) or optional argument <code>func</code> .
<code>tbl.groups(col_names_array, func)</code>	Group rows by unique combinations of values in some columns. Aggregate/count other values as above.
<code>tbl.pivot(col1, col2)</code> <code>tbl.pivot(col1,col2,vals, collect)</code>	Return 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>vals</code> and <code>collect</code> are optional arguments).
<code>tbl.sample()</code> <code>tbl.sample(n)</code> <code>tbl.sample(n, with_replacement)</code>	Returns a new table with <code>n</code> rows sampled from the original table. Default <code>n</code> is <code>tbl.num_rows</code> . Default sampling is with replacement, otherwise pass in <code>with_replacement = False</code> .
<code>tbl.barh(category_col)</code> <code>tbl.barh(category_col,freq_col)</code>	Displays a bar chart with bars for each category in the column whose name is passed in, with height proportional to the corresponding frequency. <code>freq_col</code> argument unnecessary if table consists just of a column of categories and a column of frequencies.
<code>tbl.hist(columns,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 [a, b).
<code>tbl.scatter(x_col, y_col)</code>	Draws a scatter plot consisting of one point for each row of the table
<code>tbl.plot(x_col, y_col)</code>	Draws a line plot of the data in the columns passed in

Table Predicates

Table predicates are used when a call to `tbl.where(column, predicate)` is made to filter a table by a condition. In the following examples, `x` represents a string or number and `val` represents the value a column takes for a given row. Here is a list of useful predicates:

Predicate	Description
<code>are.equal_to(x)</code>	Selects rows where <code>val == x</code>
<code>are.above(x)</code>	Selects rows where <code>val &gt; x</code>
<code>are.below(x)</code>	Selects rows where <code>val &lt; x</code>
<code>are.between(x, y)</code>	Selects rows where <code>x &lt;= val &lt; y</code>

Arrays

An array is a sequence of elements of the same type. You may assume that these functions exist in Python. In the examples below, `np` refers to the NumPy module, `tbl` refers to a Table object, `arr` refers to a NumPy array, and `num` refers to a number.

Function/Method	Description
<code>make_array()</code> , <code>make_array(val1,val2,...)</code>	Returns an array with the values passed in. If no values are passed in, returns an empty array.
<code>arr.item(i)</code>	Returns the element at the <i>i</i> -th index (the index of the first element is 0)
<code>np.append(arr, item)</code>	Returns a copy of <code>arr</code> with <code>item</code> appended to the end
<code>max(arr)</code> , <code>min(arr)</code>	Returns the maximum or minimum of the sequence
<code>sum(arr)</code>	Returns the sum of all elements in the array
<code>len(arr)</code>	Returns the length (number of elements) in an array
<code>round(num)</code> , <code>np.round(arr)</code>	Returns a number or array rounded to the nearest integer
<code>abs(num)</code> , <code>np.abs(arr)</code>	Returns the absolute value of a number or array
<code>np.mean(arr)</code>	Returns the mean (a.k.a. average) of the values in the array
<code>np.median(arr)</code>	Returns the median of the values in the array
<code>np.std(arr)</code>	Returns the standard deviation of the values in the array
<code>np.arange(start, stop, step)</code> , <code>np.arange(start, stop)</code> , <code>np.arange(stop)</code>	Returns an array of numbers starting from <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, by default, <code>step = 1</code> and <code>start = 0</code>
<code>np.count_nonzero(arr)</code>	Returns the number of nonzero elements in an array (False counts as 0, True as nonzero)
<code>np.random.choice(arr, n)</code> , <code>np.random.choice(arr)</code>	Returns an array of <code>n</code> items sampled with replacement from an array. Default <code>n</code> is 1.

Additional Functions

You may assume that these functions exist in Python.

Function/Method	Description
<code>proportions</code> <code>_from_distribution(tbl,label,n)</code>	Takes in a table, column label corresponding to distribution proportions, and a sample size <i>n</i> . Returns a table augmented with the column "Random Sample" with the sampled proportions.
<code>percentile(n, arr)</code>	Returns the <i>n</i> -th percentile of an array
<code>stats.norm.cdf(x, mean, sd)</code>	Returns the area to the left of <code>x</code> under a normal curve with mean <code>mean</code> and standard deviation <code>sd</code>
<code>minimize(function)</code>	Returns the list of parameters that minimize the function