

Reference Sheet

Pandas

In the entries below suppose `df` is a `DataFrame`, `s` is a `Series` and `pd` is the `Pandas` package

Attribute	Description
<code>df.index</code>	Returns the index labels of the <code>DataFrame</code> .
<code>df.columns</code>	Returns the column labels of the <code>DataFrame</code> .
<code>df.shape</code> or <code>s.shape</code>	Returns the shape of a <code>DataFrame</code> or <code>Series</code> in the form number of rows, number of columns
<code>df.size</code> or <code>s.size</code>	Returns the total number of entries in a <code>DataFrame</code> or <code>Series</code> (number of rows times number of columns)

Function	Description
<code>df[col]</code>	Returns the column labeled <code>col</code> from <code>df</code> as a <code>Series</code> .
<code>df[[col1, col2]]</code>	Returns a <code>DataFrame</code> containing the columns labeled <code>col1</code> and <code>col2</code> .
<code>s.loc[rows]</code> or <code>df.loc[rows, cols]</code>	Returns a <code>Series/DataFrame</code> with rows (and columns) selected by their index values.
<code>s.iloc[rows]</code> or <code>df.iloc[rows, cols]</code>	Returns a <code>Series/DataFrame</code> with rows (and columns) selected by their positions.
<code>s.isnull()</code> or <code>df.isnull()</code>	Returns boolean <code>Series/DataFrame</code> identifying missing values
<code>df.info()</code>	displays name and type of each column, the number of non-null entries, and size of dataframe
<code>s.fillna(value)</code> or <code>df.fillna(value)</code>	Returns a <code>Series/DataFrame</code> where missing values are replaced by <code>value</code>
<code>df.drop(labels, axis)</code>	Returns a <code>DataFrame</code> without the rows or columns named <code>labels</code> along <code>axis</code> (either 0 or 1)
<code>df.rename(index=None, columns=None)</code>	Returns a <code>DataFrame</code> with renamed columns from a dictionary <code>index</code> and/or <code>columns</code>
<code>df.sort_values(by, ascending=True)</code>	Returns a <code>DataFrame</code> where rows are sorted by the values in columns <code>by</code>
<code>s.sort_values(ascending=True)</code>	Returns a sorted <code>Series</code> .
<code>s.unique()</code>	Returns a <code>NumPy</code> array of the unique values
<code>s.value_counts()</code>	Returns the number of times each unique value appears in a <code>Series</code>
<code>pd.merge(left, right, how='inner', left_on='a', right_on='b')</code>	Returns a <code>DataFrame</code> joining <code>DataFrames</code> <code>left</code> and <code>right</code> on the columns labeled <code>a</code> in the left database and <code>'b'</code> in the right database; the join is of type <code>inner</code>
<code>df.pivot_table(index, columns, values=None, aggfunc='mean')</code>	Returns a <code>DataFrame</code> pivot table where columns are unique values from <code>columns</code> (column name or list), and rows are unique values from <code>index</code> (column name or list); cells are collected values using <code>aggfunc</code> . If <code>values</code> is not provided, cells are collected for each remaining column with multi-level column indexing.
<code>df.set_index(col)</code>	Returns a <code>DataFrame</code> that uses the values in the column labeled <code>col</code> as the row index.
<code>df.reset_index()</code>	Returns a <code>DataFrame</code> that has row index 0, 1, etc., and adds the current index as a column.
<code>s.str.len()</code>	Returns a <code>Series</code> containing length of each string
<code>s.str.lower()</code> or <code>s.str.upper()</code>	Returns a <code>Series</code> containing lowercase/uppercase version of each string
<code>s.str.split(pat)</code>	Split strings around given separator/delimiter <code>pat</code> . If not specified, split on whitespace

Groupby

In the groupby entries below, `col` can be a column label or a list of column labels:

Function	Description
<code>df.groupby(col).count()</code>	Returns a <code>Series/DataFrame</code> with the counts of non-missing values in each column
<code>df.groupby(col).size()</code>	Returns a <code>Series</code> counting the number of rows in each group, including missing values
<code>df.groupby(col).mean()</code> or <code>df.groupby(col).min()</code> or <code>df.groupby(col).max()</code>	Returns a <code>Series/DataFrame</code> containing mean/min/max of each group for each column, excluding missing values
<code>df.groupby(col).first()</code> or <code>df.groupby(col).last()</code>	Returns a <code>Series/DataFrame</code> containing the first/last non-null entry of each group for each column
<code>df.groupby(col).agg(f)</code>	Returns a <code>DataFrame</code> with index <code>col</code> . Aggregates other columns using the given function <code>f</code> .

You are given a Pandas DataFrame `cereal` with **information per serving** about 80 different breakfast cereals. Here are the first 5 rows of the DataFrame:

	name	manufacturer	type	calories	protein	fat	fiber	carbo	sugars
0	100% Bran	Nabisco	cold	70	4	1	10.03	5.0	6
1	100% Natural Bran	Quaker Oats	cold	120	3	5	1.93	8.0	8
2	All-Bran	Kelloggs	cold	70	4	1	8.80	7.0	5
3	All-Bran with Extra Fiber	Kelloggs	cold	50	4	0	14.04	8.0	0
4	Almond Delight	Ralston Purina	cold	110	2	2	1.00	14.0	8

Fill in blank A:

1. (5 pts) Fill in the blanks labeled below to add a new column to `cereal` named `low_calorie` which has the boolean value `True` if the cereal is low-calorie and `False` otherwise. Assume a cereal is low-calorie if the value of `calories` in the `cereal` DataFrame is less than or equal to 100.

`cereal[___A___] = ___B___`

`'low_calorie'`

Fill in blank B:

`cereal['calories'] <= 100`

2. (12 pts) Write code below to construct a DataFrame called `max_sugar` indexed by manufacturer, that has one column whose value is equal to the the maximum `sugars` value of all cereals by that manufacturer. Your DataFrame should be sorted by the max sugar value in **decreasing** order.

For example, the first few entries of the DataFrame would be:

	sugars
Post	16
Kelloggs	15
Quaker Oats	14

Write your code directly in the blanks provided below. **You can leave lines inside parentheses blank to represent a function call with no arguments.** You may not need all lines.

`max_sugar = (cereal[_____ 'sugars' _____] . groupby ('manufacturer' _____)`
`_____ . max (_____)`
`_____ . sort_values (Ascending = False _____)`
`)`

3. (3 pts)

Suppose you are given a function named `test` that takes a DataFrame as its argument and returns a Boolean.

What will be the output of `cereal.groupby("type").filter(test)` if the `test` function returns `True` for a group?

- ☒ All rows in the group will be included in the result
- ☐ All rows in the group will be excluded from the result
- ☐ Only the first row in the group will be included in the result
- ☐ Only the last row in the group will be included in the result
- ☐ None of the above.

END OF QUIZ