

Discussion #4

Name:

Pandas: Grouping Multiple Columns

Throughout this section you'll be working with the babynames (left) and elections (right) datasets as shown below:

	State	Sex	Year	Name	Count
0	CA	F	1910	Mary	295
1	CA	F	1910	Helen	239
2	CA	F	1910	Dorothy	220
3	CA	F	1910	Margaret	163
4	CA	F	1910	Frances	134

	Year	Candidate	Party	Popular vote	Result	%
0	1824	Andrew Jackson	Democratic-Republican	151271	loss	57.210122
1	1824	John Quincy Adams	Democratic-Republican	113142	win	42.789878
2	1828	Andrew Jackson	Democratic	642806	win	56.203927
3	1828	John Quincy Adams	National Republican	500897	loss	43.796073
4	1832	Andrew Jackson	Democratic	702735	win	54.574789

- (a) Which of the following lines of code will output the following dataframe? Recall that the arguments to `pd.pivot_table` are as follows: `data` is the input dataframe, `index` includes the values we use as rows, `columns` are the columns of the pivot table, `values` are the values in the pivot table, and `aggfunc` is the aggregation function that we use to aggregate values.

	Result	loss	win
Party			
Democratic	43.697060	51.441864	
Democratic-Republican	57.210122	42.789878	
National Union	NaN	54.951512	
Republican	42.047791	52.366967	
Whig	35.258650	50.180255	

- A. `pd.pivot_table(data=winners_only, index='Party', columns='Result', values='%', aggfunc=np.mean)`
- B. `winners_only.groupby(['Party', 'Result'])['%'].mean()`
- C. `pd.pivot_table(data=winners_only, index='Result', columns='Party', values='%', aggfunc=np.mean)`
- D. `winners_only.groupby('%')[['Party', 'Result']].mean()`

- (b) `name_counts_since_1940 = babynames[babynames["Year"] >= 1940].groupby(["Name", "Year"]).sum()` generates the multi-indexed DataFrame below.

Name	Year	
Aadan	2008	7
	2009	6
	2014	5
Aaden	2007	20
	2008	135
	2009	158
	2010	62
	2011	39
	2012	38

We can index into multi-indexed DataFrames using `loc` with slightly different syntax. For example `name_counts_since_1940.loc[("Aahna", 2008):("Aaiden", 2014)]` yields the DataFrame below.

		Count
Name	Year	
Aahna	2014	7
Aaiden	2009	11
	2010	11
	2011	8
	2013	13
	2014	12

Using `name_counts_since_1940`, set `imani_2013_count` equal to the number of babies born with the name 'Imani' in the year 2013. You may use either `.loc`. Make sure you're returning a value and not a Series or DataFrame.

```
imani_2013_count =
```

Pandas: String Operations and Table Joining

2. (a) Create a new DataFrame called `elections_with_first_name` with a new column 'First Name' that is equal to the Candidate's first name. Hint: Use `.str.split`.
- ```
elections_with_first_name =
```

- (b) Now create `elections_and_names` by joining `elections_with_first_name` with `name_counts_since_1940_numerical_index` (the modified version of `name_counts_since_1940` with the index reset) on both the first names of each person along and the year.

```
elections_and_names =
```

## Regular Expressions

Here's a complete list of metacharacters:

. ^ \$ \* + ? { } [ ] \ | ( )

Some reminders on what each can do (this is not exhaustive):

|                                                                                      |                                                                                               |
|--------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------|
| "^" matches the position at the beginning of string (unless used for negation "[^"]) | "\d" match any <i>digit</i> character. "\D" is the complement.                                |
| "\$" matches the position at the end of string character.                            | "\w" match any <i>word</i> character (letters, digits, underscore). "\W" is the complement.   |
| "?" match preceding literal or sub-expression 0 or 1 times.                          | "\s" match any <i>whitespace</i> character including tabs and newlines. \S is the complement. |
| "+" match preceding literal or sub-expression <i>one</i> or more times.              | "*?" Non-greedy version of *. Not fully discussed in class.                                   |
| "*" match preceding literal or sub-expression <i>zero</i> or more times              | "\b" match boundary between words. Not discussed in class.                                    |
| "." match any character except new line.                                             | "+?" Non-greedy version of +. Not discussed in class.                                         |
| "[ ]" match any one of the characters inside, accepts a range, e.g., "[a-c]"         |                                                                                               |
| "( )" used to create a sub-expression                                                |                                                                                               |

Some useful re package functions:

|                                                                                                  |                                                                                                |
|--------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------|
| re.split(pattern, string) split the string at substrings that match the pattern. Returns a list. | placing matching substrings with replace. Returns a string.                                    |
| re.sub(pattern, replace, string) apply the pattern to string re-                                 | re.findall(pattern, string) Returns a list of all matches for the given pattern in the string. |

## Regular Expressions

3. For each pattern specify the starting and ending position of the first match in the string. The index starts at zero and we are using closed intervals (both endpoints are included).

|           | abcdefg | abcs! | ab_abc | abc, _123 |
|-----------|---------|-------|--------|-----------|
| abc*      | [0, 2]  |       |        |           |
| [^\s]+    |         |       |        |           |
| ab.*c     |         |       |        |           |
| [a-z1,9]+ |         |       |        |           |

4. Given the text:

```
"<record>_Joey_Gonzalez_<jegonzal@cs.berkeley.edu>_Faculty_</record>"
"<record>_Manana_Hakobyan_<manana.hakobyan@berkeley.edu>_TA_</record>"
```

Which of the following matches exactly to the email addresses (including angle brackets)?

- A. <.\*@.\*>     B. <[^>]\*@[^>]\*>     C. <.\*@\w+\..\*>

5. Write a regular expression that matches strings that contain exactly 5 vowels.

6. Given that `sometext` is a string, use `re.sub` to replace all clusters of non-vowel characters with a single period. For example "a\_big\_moon,\_between\_us..." would be changed to "a.i.oo.e.ee.u.".

7. Given the following text in a variable `log`:

```
169.237.46.168 - - [26/Jan/2014:10:47:58 -0800]
"GET_/stat141/Winter04/_HTTP/1.1" 200 2585
"http://anson.ucdavis.edu/courses/"
```



11. Given `sometext = "I've got 10 eggs, 20 geese, and 30 giants."`, use `re.findall` to extract all the items and quantities from the string. The result should look like `['10 eggs', '20 geese', '30 giants']`. You may assume that a space separates quantity and type, and that each item ends in `s`.