

# STAT4710J SP2024 Midterm RC Part 1

---

Author: Boyuan Zhang

## Data Sampling and Probability (Lec 02)

---

### Sampling

Probability/Random Sampling

- Capable of telling each individual's probability
- No need to have an equal probability for each candidate

Sampling Methods

- Simple random sampling: **Random**
- Random Sampling with replacement: **Random**
- Systematic sampling: **Random**
- Stratified sampling: **Random**
- Cluster sampling: **Random**
- Convenience sampling: **Not random**
- Quota sampling: **Not random**
- Voluntary response sampling and snowball sampling: **Not random**

### Bias vs. Chance Error

- Bias: One direction
  1. Selection bias
  2. Response bias: People don't always tell the truth
  3. Non-response bias: No answers
- Chance error: Any direction

### Population and Sample

- Target Population
- Sample frame
- Sample

## Pandas (Lec 03, 04)

---

### Series

- Creating a Series

```
s = pd.Series([4, -2, 0, 6], index=['a', 'b', 'c', 'd'])
```

- Accessing values

```
# single label
s['a']

# list of labels
s[['a', 'b']]
```

## DataFrame

- Creating a DataFrame
  - Using lists and column names
  - From dictionaries
  - From Series

## Slicing

- `loc`: Selecting by label, inclusive of both sides.

```
# select a single row
df.loc[0]

# select multiple rows
df.loc[0:4]

# select rows and columns
df.loc[0:4, 'Year':'Party']

# select specific rows and columns
df.loc[[87, 25, 179], ['Year', 'Candidate', 'Result']]
```

- `iloc`: Selecting by position, exclusive of the right index.

```
# select a single row
df.iloc[0]

# select multiple rows
df.iloc[0:4]

# select rows and columns
df.iloc[0:4, 0:4]

# select specific rows and columns
df.loc[[1, 2, 3], [0, 1, 2]]
```

## Conditional Selection

- Create boolean arrays and pass them into slicing.

```
df[df['Party'] == 'Republican']
```

- Use `&` and `|` for boolean operations.

```
babynames[(babynames["Sex"] == "F") & (babynames["Year"] < 2000)]
```

## Handy Utility Functions

- `df.size`: Returns the value of rows \* columns.
- `df.shape`: Returns a tuple `(nrow, ncol)` containing the number of rows and columns, where `nrow = df.shape[0]` and `ncol = df.shape[1]`.
- `df.sample(n, replace=..., ...)`: Sampling without replacement by default. Set `replace=True` otherwise.
- `df.sort_values(by=..., key=..., ascending=..., ...)`: Sort values in ascending order by default. Set `ascending=False` otherwise.
- `series.value_counts()`: Counts unique values in descending order.
- `series.unique()`: Returns an array of unique values.
- `series.apply(func)`: Applies the function `func` to all the values in `series`.
- String manipulation
  - `series.str.isin()`
  - `series.str.len()`
  - `series.str.startswith()`

## Column Manipulation

- Addition: `df[newcol] = ...`
- Deletion: `df.drop(colname, axis=1)`
- Modification: `df.rename({col1: newcol1, col2: newcol2, ...})`

## Groupby and Aggregation Functions

`groupby` creates sub-dataframes and `agg` aggregates all the columns with corresponding data types to one row of output to represent the group.

```
# syntax to aggregate all the columns
# may cause errors
df.groupby(colname).agg(aggfunc)

# select specific columns before aggregation
```

```

df.groupby(colname)[[col1, col2, ...]].agg(aggfunc)

# aggregate using lambda functions
df.groupby('Name')[['Count']].agg(lambda x: x.iloc[0])

# aggregate using customized functions
def ratio_to_peak(series):
    return series.iloc[-1]/max(series)

df.groupby('Name')[['Count']].agg(ratio_to_peak)

# can also directly use built-in aggregation functions
# mean(), median(), max(), min(), sum(), count(), size()...
df.groupby(...)[[...]].mean()

# filter subtables that satisfy certain conditions
df.groupby('Name').filter(lambda x: x['num'].sum() > 10)

```

## Groupby and Pivot Table

```

# sometimes we want to group by first feature, and then by second feature
df.groupby(['Year', 'Sex']).agg(sum).head()

# use pivot table instead
df_pivot = df.pivot_table(index='Year',
                           columns='Sex',
                           values=['Count', 'Name'],
                           aggfunc=np.max
                           ).head()

```

## Join Tables

- Inner join: retains only **matched data** from both tables
- Left join: retains all the rows from the **left table** and sets NaN to the unmatched data in the right table
- Right join: retains all the rows from the **right table** and sets NaN to the unmatched data in the left table
- Outer join: retain all the rows from **both tables** and sets NaN to the unmatched data in the corresponding table

```
# basic syntax
pd.merge(
    left=left_df, # the left table to merge
    right=right_df, # the right table to merge
    how=..., # 'inner' by default
    on=None, # if the column names to merge are the same
    left_on=None, # column name to merge in the left table
    right_on=None, # column name to merge in the right table
    left_index=False, # whether to merge on the index of the left table
    right_index=False # whether to merge on the index of the right table
)
```

## Regular Expressions (Lec 05)

```
import re
```

- Syntax

### Basic Regex Syntax

The four basic operations for regular expressions.

You can technically do anything with just these basic four (albeit tediously).

|, \*, () are **metacharacters**. They manipulate adjacent characters.

operation	order	example	matches	doesn't match
<b>concatenation</b>	3	AABAAB	AABAAB	every other string
<b>or</b>	4	AA BAAB	AA BAAB	every other string
<b>closure</b> (zero or more)	2	AB*A	AA ABBBBBBA	AB ABABA
<b>group</b> (parenthesis)	1	A(A B)AAB	AAAAB ABAAB	every other string
		(AB)*A	A ABABABABA	AA ABBA

AB\*: A then zero or more copies of B:

(AB)\*: Zero or more copies of AB:

A, AB, ABB, AB BB

ABABABAB, ABAB, ,AB,

matches the empty string!



15



## Expanded Regex Syntax

**wildcard** `.`  
Consider: `.*`

**character class:**  
Match one character in `[]`

Repeat preceding item `{...}` times

Compare/contrast:  
`o*`, `o+`, `o?`

operation	example	matches	doesn't match
<b>any character</b> (except newline)	<code>.U.U.U.</code>	CUMULUS JUGULUM	SUCCUBUS TUMULTUOUS
<b>character class</b>	<code>[A-Za-z][a-z]*</code>	word Capitalized	camelCase 4illegal
<b>repeated exactly a times:</b> <code>{a}</code>	<code>j[aeiou]{3}hn</code>	jaoehn jooohn	jhn jaeiouhn
<b>repeated from a to b times:</b> <code>{a,b}</code>	<code>j[ou]{1,2}hn</code>	john juohn	jhn jooohn
<b>at least one</b>	<code>jo+hn</code>	john joooooooohn	jhn jjohn
<b>zero or one</b>	<code>joh?n</code>	jon john	any other string



23



## Convenient Regex Syntax

`\w` `[A-Za-z0-9_]`  
`\d` `[0-9]`  
`\s` whitespace  
`+` at least one

`[^...]` negates entire character class

"take this next character literally"

operation	example	matches	doesn't match
<b>built-in character classes</b>	<code>\w+</code> <code>\d+</code> <code>\s+</code>	Fawef_03 231231 whitespace	this person 423 people non-whitespace
<b>character class negation</b>	<code>[^a-z]+</code>	PEPPERS3982 17211!↑å	porch CLAmS
<b>escape character</b>	<code>cow\.com</code>	cow.com	cowscom



33



## Even More Regular Expression Features

A few additional common regex features are listed above.

- Won't discuss these in lecture, but **might come up** in discussion or hw.
- There are even more features out there!

operation	example	matches	doesn't match
<b>beginning of line</b>	<code>^ark</code>	ark two ark o ark	dark
<b>end of line</b>	<code>ark\$</code>	dark ark o ark	ark two
<b>lazy version of zero or more</b> <code>*?</code>	<code>5.*?5</code>	5005 55	5005005

hell

Greedy: `h.+l` matches hell

Lazy: `h.+?l` matches hel.

`<em>Hello World</em>`

Greedy: `<.+>` will match `<em>Hello World</em>`

Lazy: `<.+?>` will match `<em>` and `</em>`

Again—The official guide is good!

<https://docs.python.org/3/howto/regex.html>



36



- String Manipulation Functions

Python String	re	pandas Series
<code>s.lower()</code> <code>s.upper()</code>		<code>ser.str.lower()</code> <code>ser.str.upper()</code>
<code>s.replace(...)</code>	<code>re.sub(pattern, repl, text)</code>	<code>ser.str.replace(pattern, repl, regex=True)</code>
<code>s.split(...)</code>	<code>re.split(...)</code>	<code>ser.str.split(...)</code>
<code>s[1:4]</code>		<code>ser.str[1:4]</code>
	<code>re.findall(pattern, text)</code>	<code>ser.str.findall(pattern)</code> <code>ser.str.extract(pattern)</code> <code>ser.str.extractall(pattern)</code>
<code>'ab' in s</code>	<code>re.search(...)</code>	<code>ser.str.contains(...)</code>
<code>len(s)</code>		<code>ser.str.len()</code>
<code>s.strip()</code>		<code>ser.str.strip()</code>

- Raw String

Regular String	Raw String
<code>'ab*'</code>	<code>r'ab*'</code>
<code>'\\\\section'</code>	<code>r'\\section'</code>
<code>'\\w+\\s+\\1'</code>	<code>r'\\w+\\s+\\1'</code>

- Capture Group

```
text = '''Observations: 03:04:53 - Horse awakens.
03:05:14 - Horse goes back to sleep.'''
pattern = r'(\d\d):(\d\d):(\d\d) - (.*)'
matches = re.findall(pattern, text)
```

## Word Embedding (Lec 06)

- Bag-of-Words Encoding
- N-Gram Encoding
- TF-IDF

- Term Frequency (Importance in a single document)

$$tf(t, d) = \frac{\# \text{ of occurrences of } t \text{ in } d}{\text{total } \# \text{ of words in } d}$$

- Inverse Document Frequency (Rarity factor across documents)

$$idf(t) = \log \left( \frac{\text{total } \# \text{ of documents}}{\# \text{ of documents in which } t \text{ appears}} \right)$$

- TF-IDF

$$tfidf(t, d) = tf(t, d) \cdot idf(t)$$

## Data on the Internet (Lec 07, 08)

---

### Introduction to HTTP

- HTTP: Hypertext Transfer Protocol
- Request Methods
  - `GET`: Request data from a specified source
  - `POST`: Send data to the server

```
import requests

# get data
get_res = requests.get('https://bing.com')

# post data
post_res = requests.post('https://httpbin.org/post', data={'name': 'King Triton'})
```

- HTTP Status Code
  - `200`: Successful request
  - `400`: Bad request
  - `404`: Page not found
  - `500`: Internal server error
  - ...

## JSON

- JSON: JavaScript Object Notation
- Structure: Resemble Python dictionaries



```
family_tree =
{'name': 'Grandma',
 'age': 94,
 'children': [{'name': 'Dad',
                 'age': 60,
                 'children': [{'name': 'Me', 'age': 24}, {'name': 'Brother', 'age': 22}]},
               {'name': 'My Aunt',
                 'children': [{'name': 'Cousin 1', 'age': 34},
                              {'name': 'Cousin 2',
                               'age': 36,
                               'children': [{'name': 'Cousin 2 Jr.', 'age': 2}]}]}}]}
```

- Accessing values

```
family_tree['children'][0]['children'][0]['age']
```

## HTML

- HTML: HyperText Markup Language
- Useful Tags to Know

Element	Description
<code>&lt;html&gt;</code>	the document
<code>&lt;head&gt;</code>	the header
<code>&lt;body&gt;</code>	the body
<code>&lt;div&gt;</code>	a logical division of the document
<code>&lt;span&gt;</code>	an <i>inline</i> logical division
<code>&lt;p&gt;</code>	a paragraph
<code>&lt;a&gt;</code>	an anchor (hyperlink)
<code>&lt;h1&gt;, &lt;h2&gt;, ...</code>	header(s)
<code>&lt;img&gt;</code>	an image

- Attributes

```

<!--`src`: source of image-->
<!--`alt`: text to display when the image fails to display-->

<a href="https://bing.com">this link</a> <!--destination of the hyperlink-->

<input id="username"> <!--usually unique-->

<div class="stat4710j">some text</div> <!--not necessarily unique-->
```

# Parsing HTML Using BeautifulSoup

```
import requests
import bs4

html_string = requests.get(url)
soup = bs4.BeautifulSoup(html_string)

soup.find(tag, attrs={...}) # finds the first instance of a tag
soup.find_all(tag, attrs={...}) # finds all instances of a tag

# Attributes
soup.find('p').text
soup.find('div').attrs
soup.find('div').get('id')
```

## Data Wrangling and EDA (Lec 09)

---

### Structure

- The "shape" of a data file
- Variable Feature Types
  - Quantitative
    - Continuous
    - Discrete
  - Qualitative (Categorical)
    - Ordinal
    - Nominal
- Keys
  - Primary key: Unique, determines the values of the remaining columns
  - Foreign key: references primary keys in other tables

### Granularity

- How fine/coarse is each datum?
- What does each record represent?

### Scope

- How (in)complete is the data?
  - Does my data cover my area of interest?
  - Are my data too expensive?
  - Does my data cover the right time frame?

## Temporality

- How is the data situated in time?
  - Data changes
  - Periodicity
  - Time zone
  - Null values: January 1st 1970, January 1st 1900, ...

## Faithfulness (Missing Values)

- How well does the data capture "reality"?
- Deal with missing data
  - Drop records with missing values
  - Keep as `NaN`
  - Imputation/Interpolation
    - Average value
    - Random value
    - Predicted value

## Reference

---

Jingye Lin, midtermRC\_Part1, SP2023

Zhitong Tang, Mid\_RC\_Part1, SU2023

STAT4710J Slides, SP2024