

Data Sampling and Probability

Sample

Probability/Random Sampling

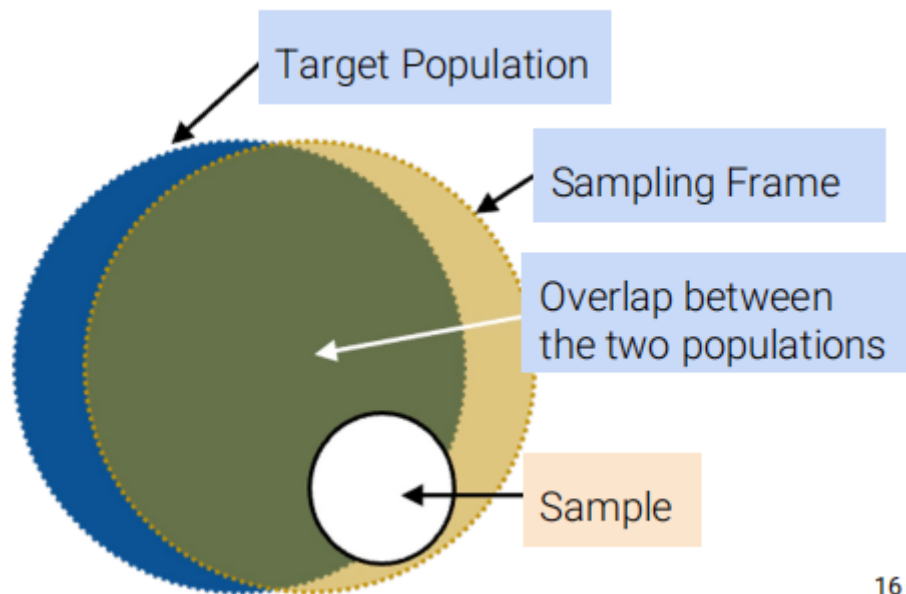
- You **must** be able to provide the **probability** that any specified set of individuals will be in the sample.
- All individuals in the population **do not need to** have the same chance of being selected.
- You will still be able to measure the errors, because you know all the probabilities.

(understand concepts through examples in slides)

probability(random) sample	Non-random sample
simple random sample	convenience sample
systematic sample	voluntary response sample
stratified sample	quota sample
cluster sample	snowball sample

- Convenience samples: **Not random**, select the ease samples. Usually used in pilot testing
- Quota samples: **Not random**. First group by one feature, and then do convenience sampling
- Stratified samples: **Random**. First group by one feature, and then do random sampling in each of the subgroups
- Simple random samples: **Random**. Use random number generator to generate a sample. (without replacement)
 - Every individual (and subset of individuals) has the same chance of being selected.
- Systematic samples: **Random**. Choose one candidates every elements
- Cluster samples: **Random**.
- Voluntary response samples and snowball samples: **not random**

Population, sample, and sampling frame



16

- **Population:** The group that you want to learn something about.
- **Sampling Frame:** The list from which the sample is drawn.
 - If you're sampling people, the sampling frame is the set of all people that could possibly end up in your sample.
- **Sample:** Who you actually end up sampling.
 - A subset of your sampling frame.

Bias vs. Chance Error

- **Bias:** One direction (If your sampling method is biased, those biases will be magnified with a larger sample size)
 - Selection bias: Systematically excluding (or favoring) particular groups.
 - Response bias: People don't always tell the truth
 - Non-response bias: No answers
- **Chance error:** Any direction (get smaller as the sample size get larger, but it is unavoidable)

Sample 1

The City of Feishu wants to hear from its homeowners on issues related to zoning laws.

(For the purposes of this question, homeowners are individuals who own their home, instead of leasing or renting from someone else.)

- (a) (1 pt) One method of surveying would be to have city workers come to Feishu Univ's campus and ask passing by students and faculty members for their thoughts. Suppose for now that the question "Are you a homeowner?" is not asked.

What type of sample is this?

- Convenience sample
- Probability sample, but not simple random sample
- Simple random sample
- Quota sample

- (b) (1 pt) Many students and faculty members aren't homeowners, but will be surveyed anyways.

What form of bias or error is this?

- Response bias
- Non-response bias
- Chance error
- Selection bias

- (c) (1 pt) The City of Feishu has a list of all the homeowners' email addresses. Instead of the previous surveying technique, now suppose they take the list of all homeowners' email addresses, shuffle it, and send a survey to every other email address. That is, from the shuffled list, they email the first, third, fifth, seventh, and so on.

(You may assume that the shuffling is done uniformly at random, meaning that each email address has the same probability of landing in any particular position. You may also assume that the City of Feishu has the email address for every single homeowner, and that every single homeowner has a unique email address.)

What type of sample is this?

- Quota sample
- Convenience sample
- Probability sample

- (d) (1 pt) Fill in the blank: In this new sampling technique, the sampling frame is _____ the population of interest.

- equal to
- greater than
- smaller than

- (e) (1 pt) In this new sampling technique, some homeowners may see the survey and choose not to respond.

True or False: The only form of bias or error in this new surveying technique is non-response bias.

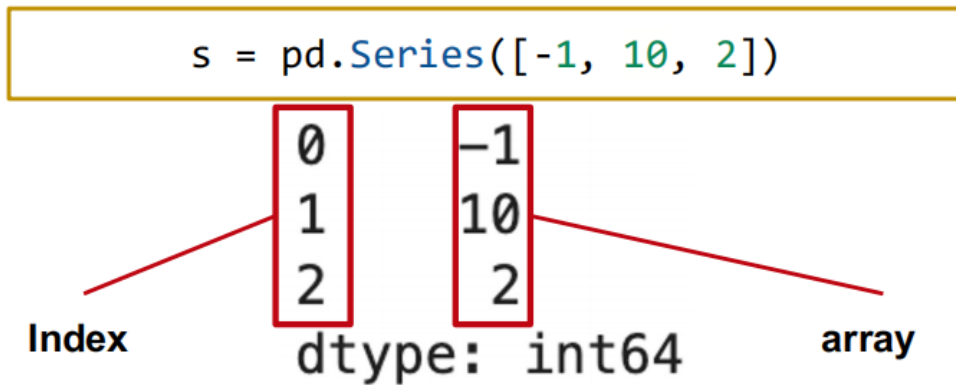
- True
- False

Pandas

Series

- 1D array-like object, value with index

```
1 s = pd.Series([-1,10,2])
2 s.array
3 s.index
```



- custom index (2 ways)

```
1 s = pd.Series([-1,10,2],index = ["a","b","c"])
2 s.index = ["first","second","third"]
```

- selection a single value or a set of values in a Series using:

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

- A single label

```
1 s["a"]
```

- A list of labels

```
1 s[["a","c"]]
```

- A filtering condition

```
1 s[s>0]
```

DataFrame

- Creating a DataFrame

```
1 pandas.DataFrame(data, index, columns)
```

- Using a list and column names

```
1 pd.DataFrame([1,2,3],columns=["Numbers"])
2 pd.DataFrame([[1,"one"],[2,"two"]],columns=["Numbers","Description"])
```

- From a dictionary

```

1 pd.DataFrame({"Fruit":["Strawberry","Orange"],"Price":[5.46,3.99]})
2 pd.DataFrame([{"Fruit":"Strawberry","Price": 5.46},
3               {"Fruit":"Orange","Price":3.99}])

```

- From a Series

```

1 s_a = pd.Series(["a1","b1","c1"],index = ["r1","r2","r3"])
2 s_b = pd.Series(["a2","b2","c2"],index = ["r1","r2","r3"])
3 pd.DataFrame({"A-column":s_a,"B-column":s_b})

```

Indices

- Modifying Indices: select a new column and set it as the index of the DataFrame

```

1 elections.set_index("Candidate", inplace=True)
2 # Resetting the index
3 elections.reset_index(inplace=True)

```

Slicing

	loc	iloc	[]
select by	value	number	one argument (row: number; column: value)
Arguments	a list, a slice (right inclusive), a single value, boolean array	a list, a slice (right exclusive), a single value	a slice of row numbers (right exclusive), a list of column labels, a single column label, boolean array

Conditional Selection

- Create boolean arrays and pass them into slicing.

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

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

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

Handy Utility Functions

- `np.mean()`, `max`
- shape/size
 - `df.shape`: (row, col)
 - `df.size`: row×col

- `df.describe()` : get statistic information like count, mean, std, ...
`df['column_name'].describe()` : a different set of statistics will be reported if `.describe()` is called on a Series
- `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()` : return an **array** of every unique value in a **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.

```

1 # syntax to aggregate all the columns
2 # may cause errors
3 df.groupby(colname).agg(aggfunc)
4
5 # select specific columns before aggregation
6 df.groupby(colname)[[col1, col2, ...]].agg(aggfunc)
7
8 # aggregate using lambda functions
9 df.groupby('Name')[['Count']].agg(lambda x: x.iloc[0])
10
11 # aggregate using customized functions
12 def ratio_to_peak(series):
13     return series.iloc[-1]/max(series)
14 df.groupby('Name')[['Count']].agg(ratio_to_peak)
15
16 # can also directly use built-in aggregation functions
17 # mean(), median(), max(), min(), sum(), count(), size()...
18 df.groupby(...)[[...]].mean()
19
20 # filter subtables that satisfy certain conditions
21 df.groupby('Name').filter(lambda x: x['num'].sum() > 10)

```

Groupby and Pivot Table

```
1 # sometimes we want to group by first feature, and then by second feature
2 df.groupby(['Year', 'Sex']).agg(sum).head()
3
4 # use pivot table instead
5 df_pivot = df.pivot_table(index='Year',
6                             columns='Sex',
7                             values=['Count', 'Name'],
8                             aggfunc=np.max
9                             ).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

```
1 # basic syntax
2 pd.merge(
3     left=left_df, # the left table to merge
4     right=right_df, # the right table to merge
5     how=..., # 'inner' by default
6     on=None, # if the column names to merge are the same
7     left_on=None, # column name to merge in the left table
8     right_on=None, # column name to merge in the right table
9     left_index=False, # whether to merge on the index of the left table
10    right_index=False # whether to merge on the index of the right table
11 )
```

Regular Expressions

operation	order	example	matches	doesn't match
concatenation	3	AABAAB	AABAAB	every other string
or <code> </code>	4	AA BAAB	AA, BAAB	every other string
closure <code>*</code> (zero or more)	2	b AB*A	AA, ABBBA, ...	AB, ...
group <code>()</code>	1	(AB)*A	A, ABABA, ...	every other string

Common Meta Characters

`^` : Start of string, or start of line in multi-line pattern

`$` : End of string, or end of line in multi-line pattern

`.` : Any character except new line (`\n`)

`*` : Match the previous token 0 or more times. *E.g:* `ah*` matches `a` and `ahhh`

`+` : Match the previous token 1 or more times. *E.g:* `ah+` matches `ah` and `ahhh` , but not `a`

`?` : Match the previous token zero or one time. *E.g:* `bite?` matches `bite` and `bit`

`.*?` : lazy version of zero or more

`{n}` : Match the previous token exactly `n` times.

`{n,}` : Match the previous token `n` or more times.

`{m,n}` : Match the previous token `n` times, $n \in [m,n]$

`|` : Or.

Groups and Ranges

`(...)` : Group

`(?:...)` : Passive (non-capturing) group

`[abc]` : `a` or `b` or `c`

`[^abc]` : Not `a` or `b` or `c`

`[a-k]` : Lower letters from `a` to `k`

`[A-K]` : Upper letter from `A` to `K`

`[a-kA-Q]` : Lower letters from `a` to `k` , upper letter from `A` to `Q`

`[0-8]` : Digits from `0` to `8`

Character Classes

`\s` : White space. Equivalent to `[\r\n\t\f\v]`

`\S` : Not white space. Equivalent to `[^\r\n\t\f\v]`

`\d` : Digit. Equivalent to `[0-9]`

`\D` : Not digit. Equivalent to `[^0-9]`

`\w` : Word character. equivalent to `[a-zA-Z0-9_]`

`\W` : Not word character. equivalent to `[^a-zA-Z0-9_]`

`\c` : Control character

String Manipulation Functions

re	pandas Series										
<p><code>re.sub(pattern, repl, text)</code>: Returns text with all instances of pattern replaced by repl. (patter is a raw string: r"...")</p>	<p><code>ser.str.replace(pattern, repl, regex=True)</code>: Returns Series with all instances of pattern in Series ser replaced by repl.</p>										
<p><code>re.findall(pattern, text)</code>: Return a list of all matches to pattern.</p>	<p><code>ser.str.findall(pattern)</code>: Returns a Series of lists</p> <pre>df["SSN"].str.findall(pattern)</pre> <table border="1"> <thead> <tr> <th></th> <th>SSN</th> </tr> </thead> <tbody> <tr> <td>0</td> <td>987-65-4321</td> </tr> <tr> <td>1</td> <td>forty</td> </tr> <tr> <td>2</td> <td>123-45-6789 bro or 321-45-6789</td> </tr> <tr> <td>3</td> <td>999-99-9999</td> </tr> </tbody> </table> <pre> 0 [987-65-4321] 1 [] 2 [123-45-6789, 321-45-6789] 3 [999-99-9999] Name: SSN, dtype: object </pre>		SSN	0	987-65-4321	1	forty	2	123-45-6789 bro or 321-45-6789	3	999-99-9999
	SSN										
0	987-65-4321										
1	forty										
2	123-45-6789 bro or 321-45-6789										
3	999-99-9999										
	<p><code>ser.str.extract(pattern)</code>: Returns a DataFrame of first match, one group per column</p>										
	<p><code>ser.str.extractall(pattern)</code>: Returns a DataFrame of all matches, one row per match</p>										

Raw String

Regular String	Raw String
'ab*'	r'ab*'
'\\\\section'	r'\\section'
'\\w+\\s+\\1'	r'\\w+\\s+\\1'

Capture Group

- Every set of parentheses specifies a **match/capture group**.

```


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

```

```

[('03', '04', '53', 'Horse awakens.'),
 ('03', '05', '14', 'Horse goes back to sleep.')]

```



Sample3

For this question, you're given the following code:

```
re.findall(pattern, "godoggogo100")
```

For each possible pattern, list the number of times that the string "go" appears as an item in the list returned by the above code. The first two have been completed for you: Pattern 1 returns ["go", "go", "go"], so we wrote 3; pattern 2 returns ["godo"] and does not contain the string "go" as an item, so we wrote 0.

Each response is worth 1 point.

- | | |
|---------------------------|-------|
| 1. pattern = r'go' | 3 |
| 2. pattern = r'godo' | 0 |
| 3. pattern = r'go.*' | _____ |
| 4. pattern = r'.*go.*' | _____ |
| 5. pattern = r'go{2}' | _____ |
| 6. pattern = r'(go){1}' | _____ |
| 7. pattern = r'(go)[dg1]' | _____ |
| 8. pattern = r'go' | _____ |
| 9. pattern = r'[go]*(go)' | _____ |

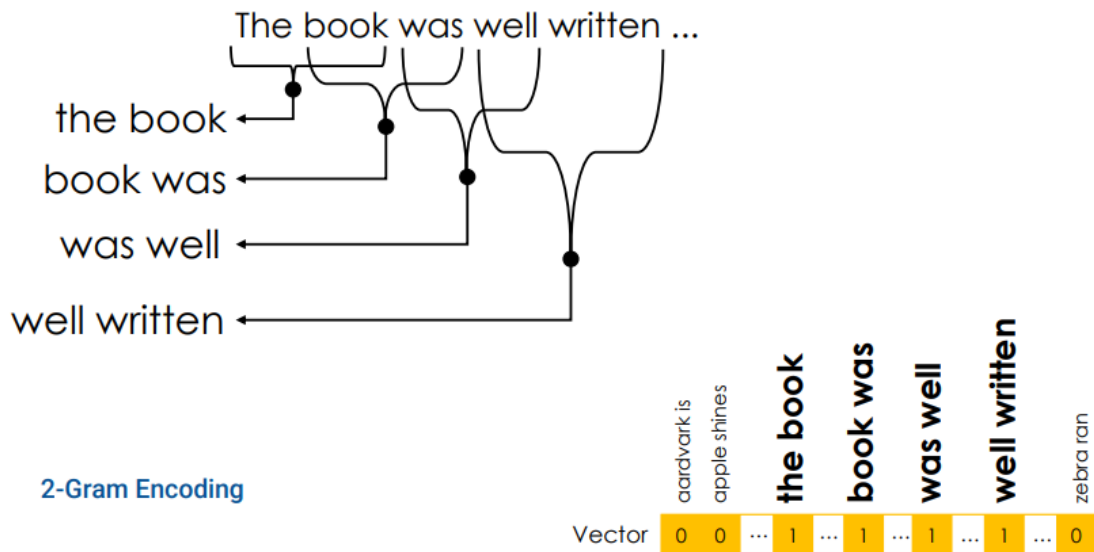
Word Embedding

Bag of words

- A **bag** is another term for a multiset: an unordered collection which may contain multiple instances of each element
- Represents the text as the frequency with which each term appears in the vocabulary
- **Key idea:** The **more similar** two unit vectors are, the **larger their dot product** is

N-Gram Encoding

- N-Gram: "Bag-of- sequences-of-words"



TF-IDF

Motivation

- The bag of words model doesn't know which words are "important" in a document.
- How do we determine which words are important?
 - The **most common words** ("the", "has") often don't have much meaning!
 - The **very rare words** are also less important!

TF (Term frequency)

- Importance in a single document

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

IDF (Inverse document frequency)

- the "rarity factor" of t across documents – **the larger $\text{idf}(t)$ is, the more rare t is.**

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

TF-IDF

- Goal: Quantify how well word t summarizes document d

$$\begin{aligned} \text{tfidf}(t, d) &= \text{tf}(t, d) \cdot \text{idf}(t) \\ &= \frac{\text{\# of occurrences of } t \text{ in } d}{\text{total \# of words in } d} \cdot \log\left(\frac{\text{total \# of documents}}{\text{\# of documents in which } t \text{ appears}}\right) \end{aligned}$$

Data on the Internet

HTTP

- Hypertext Transfer Protocol
- follows the **request-response** model

Request methods

two ways to make HTTP requests outside of a browser:

- From the command line, with `curl`.
- **From Python, with the `requests` package.**
 - `get`: used to request data **from** a specified resource.

```
1 get_res = requests.get('https://bing.com')
```

- `post`: is used to **send** data to the server.

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

HTTP status codes

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

Data formats

Responses typically come in one of two formats: **HTML or JSON**.

- The response body of a `GET` request is usually either JSON (when using an API) or HTML (when accessing a webpage).
- The response body of a `POST` request is usually JSON.

JSON

- **JavaScript Object Notation**
- JSON objects *resemble* Python dictionaries (but are not the same!).

```

1 family_tree = {'name': 'Grandma',
2   'age': 94,
3   'children': [{'name': 'Dad',
4     'age': 60,
5     'children': [{'name': 'Me', 'age': 24}, {'name': 'Brother', 'age':
22}]},
6   {'name': 'My Aunt',
7     'children': [{'name': 'Cousin 1', 'age': 34},
8     {'name': 'Cousin 2',
9       'age': 36,
10      'children': [{'name': 'Cousin 2 Jr.', 'age': 2}]}]}]}

```

- o Accessing values

```
1 family_tree['children'][0]['children'][1]['age']
```

- `json.load(f)` loads a JSON file from a file object.
- `json.loads(f)` loads a JSON file from a string.

APIs and scraping

- There are two ways of collecting data through a request:
 - o By using a published **API** (application programming interface).
 - o By **scraping** a webpage to collect its HTML source code.

HTML

- HTML (HyperText Markup Language) is **the** basic building block of the internet.

Useful tags:

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

- tags can have **attributes**, which further specify how to display information on a webpage.

```
1 
2 <!-- `src`: source of image-->
3 <!-- `alt`: text to display when the image fails to display-->
```

```
1 Click <a href="https://bing.com">this link</a> to search your keywords.
2 <!--destination of the hyperlink-->
```

```
1 <input id="username"> <!--usually unique-->
2 <div class="stat4710j">some text</div> <!--not necessarily unique-->
```

Parsing HTML using BeautifulSoup

- BeautifulSoup 4 is a Python HTML parser.
 - To "parse" means to "extract meaning from a sequence of symbols".

```
1 import requests
2 import bs4
3
4 html_string = requests.get(url)
5 soup = bs4.BeautifulSoup(html_string)
6
7 # Finding elements in a tree
8 soup.find(tag, attrs={...}) # finds the first instance of a tag
9 soup.find_all(tag, attrs={...}) # finds all instances of a tag
10
11 # Node Attributes
12 soup.find('p').text # gets the text between the opening and closing tags
13 soup.find('div').attrs # lists all attributes of a tag
14 soup.find('div').get('id') # gets the value of a tag attribute (must be
    called directly on the node that contains the attribute you're looking for)
```

Data Erangling and EDA

Key Data Properties to Consider in EDA

Structure -- the "shape" of a data file

Granularity -- how fine/coarse is each datum

Scope -- how (in)complete is the data

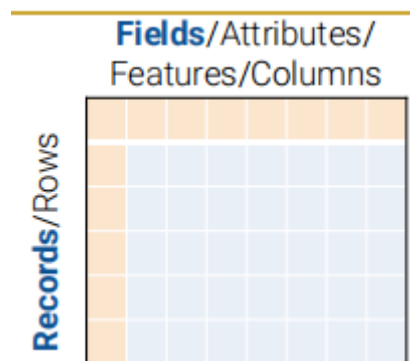
Temporality -- how is the data situated in time

Faithfulness -- how well does the data capture "reality"

Structure

- the “shape” of a data file

Rectangular Data

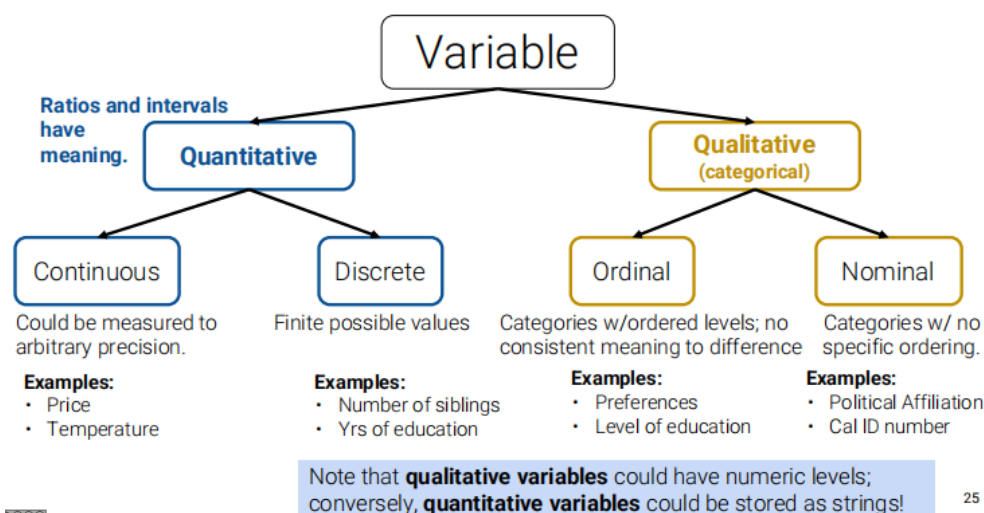


- Two kinds of rectangular data: **Tables** and **Matrices**.
 - Table: dataframes in R/Python and relations in SQL
 - Named columns with **different types**
 - Manipulated using data transformation languages (map, filter, group by, join, ...)
 - Matrices:
 - Numeric data of the **same type** (float, int, etc.)
 - Manipulated using linear algebra

File Format

- TSV
- CSV
- JSON (Not rectangular)

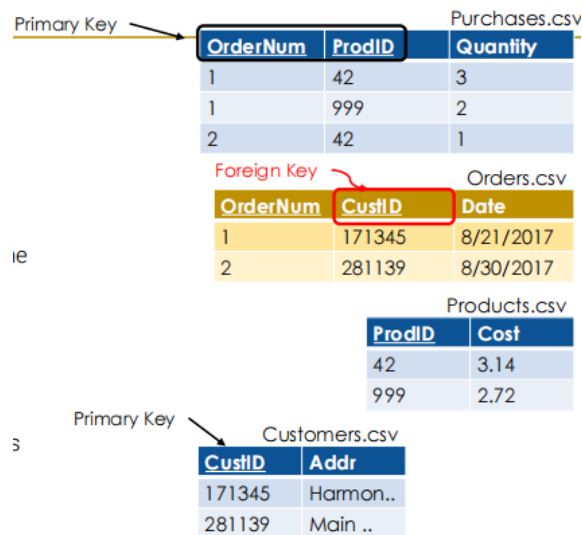
Variable Type



25

Join Multiple files: Keys

- **Primary key:** the column or set of columns in a table that determine the values of the remaining columns
 - Primary keys are unique
- **Foreign keys:** the column or sets of columns that reference primary keys in other tables.



Granularity, Scope, Temporality

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

Boyuan Zhang, midtermRC_Part1, SP2024

Jingye Lin, midtermRC_Part1, SP2023

Zhitong Tang, Mid_RC_Part1, SU2023

STAT4710J Slides, SP2024

It's the annual Monopoly World Championship! The finalists: Shawn, Amanda, Neil, and Annie are playing Monopoly, a board game where players pay a price to buy properties, which can then generate income for them. Each property can be owned by only one player at a time. At the end of the game, the player with the most money wins.

Shawn wants to figure out which properties are most worth buying. He creates a DataFrame `income` with data on the current game state, shown on the **left**. He also finds a DataFrame `properties` with data on Monopoly properties, shown on the **right**.

Both tables have 28 rows. For brevity, only the first few rows of each DataFrame are shown.

	Player	Property	Income Generated
0	Shawn	Boardwalk	\$425
1	Amanda	Park Place	\$375
2	Neil	Marvin Gardens	\$200
3	NaN	Kentucky Ave	NaN
4	Shawn	Pennsylvania Ave	\$150
5	Annie	Oriental Ave	\$50
6	Amanda	Baltic Ave	\$60

`income`

	Property	Property Color	Purchase Price
0	Park Place	Dark Blue	350.0
1	Oriental Ave	Light Blue	100.0
2	Vermont Ave	Light Blue	100.0
3	Pacific Ave	Green	300.0
4	Boardwalk	Dark Blue	400.0
5	Illinois Ave	Red	240.0
6	Atlantic Ave	Yellow	260.0

`properties`

Left table:

- Player is the name of the player, as a str.
- Property is a property currently owned by the player, as a str.
- Income Generated is the amount of income a player has earned from that property so far, as a str.

Right table:

- Property is the name of the property, as a str. There are 28 unique properties.
- Property Color is a color group that the property belongs to, as a str. There are 10 unique color groups, and each property belongs to a single group.
- Purchase Price is the price to buy the property, as a float.

(a) What is the granularity of the income table?

(b) Which of the following line(s) of code successfully returns a Series with the number of properties each player owns? **Select all that apply.**

- `income.groupby('Player').agg(pd.value_counts)`
- `income['Player'].value_counts()`
- `income['Player', 'Property'].groupby('Player').size()`
- `income.groupby('Player').size()`

(c) He now decides to calculate the amount of profit from each property. He wants to store this in a column called Profit in the income DataFrame.

i. To do this, he first has to transform the Income Generated column to be of a float datatype. Write one line of code to replace the old column with a new column, also called Income Generated, with the datatype modification described above. You may assume that each entry in Income Generated consists of a dollar sign (\$) followed by a number, except for the NaN values.

ii. [3 Pts] Assuming that the answer to the last sub-part is correct, let's add a Profit column to the `income` DataFrame. **Fill in the following blanks to do this**, and please add arguments to function calls as you see appropriate.

Note: Profit is calculated by subtracting the purchase price from generated income.

```
combined_df = income.____A____ (_____B_____)  
income["Profit"] = _____C_____
```

(d) [2 Pts] Regardless of your answer to the previous sub-part, assume we've successfully created the Profit column. Let's help Shawn see how well he's doing by finding the average profit he's made on the properties he owns. Which of the following line(s) of code does this? Select all that apply.

- `income.groupby("Player")["Shawn"].agg(np.average)`
- `income[income["Player"]=="Shawn"]["Profit"].mean()`
- `income.loc[income["Player"]=="Shawn", "Profit"].mean()`
- `income.iloc[income["Player"]=="Shawn", "Profit"].mean()`

1

(a) Property

(b) BD

(c) i. `income['Income Generated'] = income['Income
Generated'].str[1:].astype(float)`

ii. `combined_df = income.merge(properties, on = "Property")`

`income["Profit"] = combined_df["Income Generated"] -
combined_df["Purchase Price"]`

(d) BC