

Lecture 6: Data Cleaning and Preparation



Alpar Sultan, PhD, Associate professor

Handling Missing Data

```
In [10]: string_data = pd.Series(['aardvark', 'artichoke', np.nan, 'avocado'])
```

```
In [11]: string_data
```

```
Out[11]:
```

```
0    aardvark
1    artichoke
2         NaN
3     avocado
dtype: object
```

```
In [12]: string_data.isnull()
```

```
Out[12]:
```

```
0    False
1    False
2     True
3    False
dtype: bool
```

```
In [13]: string_data[0] = None
```

```
In [14]: string_data.isnull()
```

```
Out[14]:
```

```
0     True
1    False
2     True
3    False
dtype: bool
```

NA handling methods

| Argument | Description |
|----------------------|---|
| <code>dropna</code> | Filter axis labels based on whether values for each label have missing data, with varying thresholds for how much missing data to tolerate. |
| <code>fillna</code> | Fill in missing data with some value or using an interpolation method such as <code>'ffill'</code> or <code>'bfill'</code> . |
| <code>isnull</code> | Return boolean values indicating which values are missing/NA. |
| <code>notnull</code> | Negation of <code>isnull</code> . |

Filtering Out Missing Data

```
In [15]: from numpy import nan as NA
```

```
In [16]: data = pd.Series([1, NA, 3.5, NA, 7])
```

```
In [17]: data.dropna()
```

```
Out[17]:
```

```
0    1.0
```

```
2    3.5
```

```
4    7.0
```

```
dtype: float64
```

```
In [18]: data[data.notnull()]
```

```
Out[18]:
```

```
0    1.0
```

```
2    3.5
```

```
4    7.0
```

```
dtype: float64
```

Filling In Missing Data

- Rather than filtering out missing data (and potentially discarding other data along with it), you may want to fill in the “holes” in any number of ways. For most purposes, the `fillna` method is the workhorse function to use. Calling `fillna` with a constant replaces missing values with that value.

| Argument | Description |
|----------------------|---|
| <code>value</code> | Scalar value or dict-like object to use to fill missing values |
| <code>method</code> | Interpolation; by default <code>'ffill'</code> if function called with no other arguments |
| <code>axis</code> | Axis to fill on; default <code>axis=0</code> |
| <code>inplace</code> | Modify the calling object without producing a copy |
| <code>limit</code> | For forward and backward filling, maximum number of consecutive periods to fill |

Data Transformation

- Transforming Data Using a Function or Mapping
- `data = pd.DataFrame({'food': ['bacon', 'pulled pork', 'bacon', 'Pastrami', 'corned beef', 'Bacon', 'pastrami', 'honey ham', 'nova lox'], 'ounces': [4, 3, 12, 6, 7.5, 8, 3, 5, 6]})`
- `meat_to_animal = { 'bacon': 'pig', 'pulled pork': 'pig', 'pastrami': 'cow', 'corned beef': 'cow', 'honey ham': 'pig', 'nova lox': 'salmon' }`

Replacing Values

- Filling in missing data with the `fillna` method is a special case of more general value replacement. As you've already seen, `map` can be used to modify a subset of values in an object but `replace` provides a simpler and more flexible way to do so

Renaming Axis Indexes

```
In [66]: data = pd.DataFrame(np.arange(12).reshape((3, 4)),  
.....:                      index=['Ohio', 'Colorado', 'New York'],  
.....:                      columns=['one', 'two', 'three', 'four'])
```

```
In [67]: transform = lambda x: x[:4].upper()
```

```
In [68]: data.index.map(transform)
```

```
Out[68]: Index(['OHIO', 'COLO', 'NEW '], dtype='object')
```

```
In [69]: data.index = data.index.map(transform)
```

```
In [70]: data
```

```
Out[70]:
```

| | one | two | three | four |
|------|-----|-----|-------|------|
| OHIO | 0 | 1 | 2 | 3 |
| COLO | 4 | 5 | 6 | 7 |
| NEW | 8 | 9 | 10 | 11 |

Discretization and Binning

- Continuous data is often discretized or otherwise separated into “bins” for analysis. Suppose you have data about a group of people in a study, and you want to group them into discrete age buckets

Python built-in string methods

| Argument | Description |
|--|---|
| <code>count</code> | Return the number of non-overlapping occurrences of substring in the string. |
| <code>endswith</code> | Returns <code>True</code> if string ends with suffix. |
| <code>startswith</code> | Returns <code>True</code> if string starts with prefix. |
| <code>join</code> | Use string as delimiter for concatenating a sequence of other strings. |
| <code>index</code> | Return position of first character in substring if found in the string; raises <code>ValueError</code> if not found. |
| <code>find</code> | Return position of first character of <i>first</i> occurrence of substring in the string; like <code>index</code> , but returns <code>-1</code> if not found. |
| <code>rfind</code> | Return position of first character of <i>last</i> occurrence of substring in the string; returns <code>-1</code> if not found. |
| <code>replace</code> | Replace occurrences of string with another string. |
| <code>strip</code> , <code>rstrip</code> , <code>rstrip</code> | Trim whitespace, including newlines; equivalent to <code>x.strip()</code> (and <code>rstrip</code> , <code>rstrip</code> , respectively) for each element. |
| <code>split</code> | Break string into list of substrings using passed delimiter. |
| <code>lower</code> | Convert alphabet characters to lowercase. |
| <code>upper</code> | Convert alphabet characters to uppercase. |
| <code>casefold</code> | Convert characters to lowercase, and convert any region-specific variable character combinations to a common comparable form. |
| <code>ljust</code> , <code>rjust</code> | Left justify or right justify, respectively; pad opposite side of string with spaces (or some other fill character) to return a string with a minimum width. |

Regular expression methods

| Argument | Description |
|--------------------------------------|--|
| <code>findall</code> | Return all non-overlapping matching patterns in a string as a list |
| <code>finditer</code> | Like <code>findall</code> , but returns an iterator |
| <code>match</code> | Match pattern at start of string and optionally segment pattern components into groups; if the pattern matches, returns a match object, and otherwise <code>None</code> |
| <code>search</code> | Scan string for match to pattern; returning a match object if so; unlike <code>match</code> , the match can be anywhere in the string as opposed to only at the beginning |
| <code>split</code> | Break string into pieces at each occurrence of pattern |
| <code>sub</code> , <code>subn</code> | Replace all (<code>sub</code>) or first <code>n</code> occurrences (<code>subn</code>) of pattern in string with replacement expression; use symbols <code>\1</code> , <code>\2</code> , <code>...</code> to refer to match group elements in the replacement string |

| Method | Description |
|---|--|
| <code>cat</code> | Concatenate strings element-wise with optional delimiter |
| <code>contains</code> | Return boolean array if each string contains pattern/regex |
| <code>count</code> | Count occurrences of pattern |
| <code>extract</code> | Use a regular expression with groups to extract one or more strings from a Series of strings; the result will be a DataFrame with one column per group |
| <code>endswith</code> | Equivalent to <code>x.endswith(pattern)</code> for each element |
| <code>startswith</code> | Equivalent to <code>x.startswith(pattern)</code> for each element |
| <code>findall</code> | Compute list of all occurrences of pattern/regex for each string |
| <code>get</code> | Index into each element (retrieve <i>i</i> -th element) |
| <code>isalnum</code> | Equivalent to built-in <code>str.isalnum</code> |
| <code>isalpha</code> | Equivalent to built-in <code>str.isalpha</code> |
| <code>isdecimal</code> | Equivalent to built-in <code>str.isdecimal</code> |
| <code>isdigit</code> | Equivalent to built-in <code>str.isdigit</code> |
| <code>islower</code> | Equivalent to built-in <code>str.islower</code> |
| <code>isnumeric</code> | Equivalent to built-in <code>str.isnumeric</code> |
| <code>isupper</code> | Equivalent to built-in <code>str.isupper</code> |
| <code>join</code> | Join strings in each element of the Series with passed separator |
| <code>len</code> | Compute length of each string |
| <code>lower</code> , <code>upper</code> | Convert cases; equivalent to <code>x.lower()</code> or <code>x.upper()</code> for each element |

| Method | Description |
|----------------------|--|
| <code>match</code> | Use <code>re.match</code> with the passed regular expression on each element, returning matched groups as list |
| <code>pad</code> | Add whitespace to left, right, or both sides of strings |
| <code>center</code> | Equivalent to <code>pad(side='both')</code> |
| <code>repeat</code> | Duplicate values (e.g., <code>s.str.repeat(3)</code> is equivalent to <code>x * 3</code> for each string) |
| <code>replace</code> | Replace occurrences of pattern/regex with some other string |
| <code>slice</code> | Slice each string in the Series |
| <code>split</code> | Split strings on delimiter or regular expression |
| <code>strip</code> | Trim whitespace from both sides, including newlines |
| <code>rstrip</code> | Trim whitespace on right side |
| <code>lstrip</code> | Trim whitespace on left side |
