

# Lecture 6: Data Cleaning and Preparation



# 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 'ffill' or 'bfill'.
<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
```

```
In [19]: data = pd.DataFrame([[1., 6.5, 3.], [1., NA, NA],  
.....:                      [NA, NA, NA], [NA, 6.5, 3.]])
```

```
In [20]: cleaned = data.dropna()
```

```
In [21]: data
```

```
Out[21]:
```

	0	1	2
0	1.0	6.5	3.0
1	1.0	NaN	NaN
2	NaN	NaN	NaN
3	NaN	6.5	3.0

```
In [22]: cleaned
```

```
Out[22]:
```

	0	1	2
0	1.0	6.5	3.0

```
In [23]: data.dropna(how='all')
```

```
Out[23]:
```

	0	1	2
0	1.0	6.5	3.0
1	1.0	NaN	NaN
3	NaN	6.5	3.0

```
In [24]: data[4] = NA
```

```
In [25]: data
```

```
Out[25]:
```

	0	1	2	4
0	1.0	6.5	3.0	NaN
1	1.0	NaN	NaN	NaN
2	NaN	NaN	NaN	NaN
3	NaN	6.5	3.0	NaN

```
In [26]: data.dropna(axis=1, how='all')
```

```
Out[26]:
```

	0	1	2
0	1.0	6.5	3.0
1	1.0	NaN	NaN
2	NaN	NaN	NaN
3	NaN	6.5	3.0

```
In [27]: df = pd.DataFrame(np.random.randn(7, 3))
```

```
In [28]: df.iloc[:4, 1] = NA
```

```
In [29]: df.iloc[:2, 2] = NA
```

```
In [30]: df
```

```
Out[30]:
```

	0	1	2
0	-0.204708	NaN	NaN
1	-0.555730	NaN	NaN
2	0.092908	NaN	0.769023
3	1.246435	NaN	-1.296221
4	0.274992	0.228913	1.352917
5	0.886429	-2.001637	-0.371843
6	1.669025	-0.438570	-0.539741

```
In [31]: df.dropna()
```

```
Out[31]:
```

	0	1	2
4	0.274992	0.228913	1.352917
5	0.886429	-2.001637	-0.371843
6	1.669025	-0.438570	-0.539741

```
In [32]: df.dropna(thresh=2)
```

```
In [31]: df.dropna()
```

```
Out[31]:
```

	0	1	2
4	0.274992	0.228913	1.352917
5	0.886429	-2.001637	-0.371843
6	1.669025	-0.438570	-0.539741

```
In [32]: df.dropna(thresh=2)
```

```
Out[32]:
```

	0	1	2
2	0.092908	NaN	0.769023
3	1.246435	NaN	-1.296221
4	0.274992	0.228913	1.352917
5	0.886429	-2.001637	-0.371843
6	1.669025	-0.438570	-0.539741

# Filling In Missing Data

```
In [33]: df.fillna(0)
```

```
Out[33]:
```

	0	1	2
0	-0.204708	0.000000	0.000000
1	-0.555730	0.000000	0.000000
2	0.092908	0.000000	0.769023
3	1.246435	0.000000	-1.296221
4	0.274992	0.228913	1.352917
5	0.886429	-2.001637	-0.371843
6	1.669025	-0.438570	-0.539741

```
In [34]: df.fillna({1: 0.5, 2: 0})
```

```
Out[34]:
```

	0	1	2
0	-0.204708	0.500000	0.000000
1	-0.555730	0.500000	0.000000
2	0.092908	0.500000	0.769023
3	1.246435	0.500000	-1.296221
4	0.274992	0.228913	1.352917
5	0.886429	-2.001637	-0.371843
6	1.669025	-0.438570	-0.539741

```
In [35]: _ = df.fillna(0, inplace=True)
```

```
In [36]: df
```

```
Out[36]:
```

	0	1	2
0	-0.204708	0.000000	0.000000
1	-0.555730	0.000000	0.000000
2	0.092908	0.000000	0.769023
3	1.246435	0.000000	-1.296221
4	0.274992	0.228913	1.352917
5	0.886429	-2.001637	-0.371843
6	1.669025	-0.438570	-0.539741

```
In [37]: df = pd.DataFrame(np.random.randn(6, 3))
```

```
In [38]: df.iloc[2:, 1] = NA
```

```
In [39]: df.iloc[4:, 2] = NA
```

```
In [40]: df
```

```
Out[40]:
```

	0	1	2
0	0.476985	3.248944	-1.021228
1	-0.577087	0.124121	0.302614
2	0.523772	NaN	1.343810
3	-0.713544	NaN	-2.370232
4	-1.860761	NaN	NaN
5	-1.265934	NaN	NaN

```
In [41]: df.fillna(method='ffill')
```

```
Out[41]:
```

	0	1	2
0	0.476985	3.248944	-1.021228
1	-0.577087	0.124121	0.302614
2	0.523772	0.124121	1.343810
3	-0.713544	0.124121	-2.370232
4	-1.860761	0.124121	-2.370232
5	-1.265934	0.124121	-2.370232

```
In [42]: df.fillna(method='ffill', limit=2)
```

```
Out[42]:
```

	0	1	2
0	0.476985	3.248944	-1.021228
1	-0.577087	0.124121	0.302614
2	0.523772	0.124121	1.343810
3	-0.713544	0.124121	-2.370232
4	-1.860761	NaN	-2.370232
5	-1.265934	NaN	-2.370232

```
In [43]: data = pd.Series([1., NA, 3.5, NA, 7])
```

```
In [44]: data.fillna(data.mean())
```

```
Out[44]:
```

0	1.000000
1	3.833333
2	3.500000
3	3.833333
4	7.000000

```
dtype: float64
```



## *Table of fillna function arguments*

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

# Data Transformation

## Removing Duplicates

```
In [45]: data = pd.DataFrame({'k1': ['one', 'two'] * 3 + ['two'],  
.....:                      'k2': [1, 1, 2, 3, 3, 4, 4]})
```

```
In [46]: data
```

```
Out[46]:
```

	k1	k2
0	one	1
1	two	1
2	one	2
3	two	3
4	one	3
5	two	4
6	two	4

```
In [47]: data.duplicated()
```

```
Out[47]:
```

0	False
1	False
2	False
3	False
4	False
5	False
6	True

dtype: bool

```
In [48]: data.drop_duplicates()
```

```
Out[48]:
```

	k1	k2
0	one	1
1	two	1
2	one	2
3	two	3
4	one	3
5	two	4

```
In [49]: data['v1'] = range(7)
```

```
In [50]: data.drop_duplicates(['k1'])
```

```
Out[50]:
```

	k1	k2	v1
0	one	1	0
1	two	1	1

```
In [51]: data.drop_duplicates(['k1', 'k2'], keep='last')
```

```
Out[51]:
```

	k1	k2	v1
0	one	1	0
1	two	1	1
2	one	2	2
3	two	3	3
4	one	3	4
6	two	4	6

# Transforming Data Using a Function or Mapping

```
In [52]: 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]})
```

```
In [53]: data  
Out[53]:
```

	food	ounces
0	bacon	4.0
1	pulled pork	3.0
2	bacon	12.0
3	Pastrami	6.0
4	corned beef	7.5
5	Bacon	8.0
6	pastrami	3.0
7	honey ham	5.0
8	nova lox	6.0

```
meat_to_animal = {  
    'bacon': 'pig',  
    'pulled pork': 'pig',  
    'pastrami': 'cow',  
    'corned beef': 'cow',  
    'honey ham': 'pig',  
    'nova lox': 'salmon'  
}
```

```
In [55]: lowercased = data['food'].str.lower()
```

```
In [56]: lowercased  
Out[56]:
```

0	bacon
1	pulled pork
2	bacon
3	pastrami
4	corned beef
5	bacon
6	pastrami
7	honey ham
8	nova lox

```
Name: food, dtype: object
```

```
In [57]: data['animal'] = lowercased.map(meat_to_animal)
```

```
In [58]: data
```

```
Out[58]:
```

	food	ounces	animal
0	bacon	4.0	pig
1	pulled pork	3.0	pig
2	bacon	12.0	pig
3	Pastrami	6.0	cow
4	corned beef	7.5	cow
5	Bacon	8.0	pig
6	pastrami	3.0	cow
7	honey ham	5.0	pig
8	nova lox	6.0	salmon

```
In [59]: data['food'].map(lambda x: meat_to_animal[x.lower()])
```

```
Out[59]:
```

0	pig
1	pig
2	pig
3	cow
4	cow
5	pig
6	cow
7	pig
8	salmon

```
Name: food, dtype: object
```

# Replacing Values

```
In [60]: data = pd.Series([1., -999., 2., -999., -1000., 3.])
```

```
In [61]: data
```

```
Out[61]:
```

```
0      1.0
```

```
1    -999.0
```

```
2      2.0
```

```
3    -999.0
```

```
4   -1000.0
```

```
5      3.0
```

```
dtype: float64
```

```
In [62]: data.replace(-999, np.nan)
```

```
Out[62]:
```

```
0      1.0
```

```
1      NaN
```

```
2      2.0
```

```
3      NaN
```

```
4   -1000.0
```

```
5      3.0
```

```
dtype: float64
```

```
In [63]: data.replace([-999, -1000], np.nan)
```

```
Out[63]:
```

```
0      1.0
```

```
1      NaN
```

```
2      2.0
```

```
3      NaN
```

```
4      NaN
```

```
5      3.0
```

```
dtype: float64
```

```
In [64]: data.replace([-999, -1000], [np.nan, 0])
```

```
Out[64]:
```

```
0      1.0
```

```
1      NaN
```

```
2      2.0
```

```
3      NaN
```

```
4      0.0
```

```
5      3.0
```

```
dtype: float64
```

```
In [65]: data.replace({-999: np.nan, -1000: 0})
```

```
Out[65]:
```

```
0      1.0
```

```
1      NaN
```

```
2      2.0
```

```
3      NaN
```

```
4      0.0
```

```
5      3.0
```

```
dtype: float64
```

# 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

```
In [71]: data.rename(index=str.title, columns=str.upper)
```

```
Out[71]:
```

	ONE	TWO	THREE	FOUR
Ohio	0	1	2	3
Colo	4	5	6	7
New	8	9	10	11

# Detecting and Filtering Outliers

```
In [92]: data = pd.DataFrame(np.random.randn(1000, 4))
```

```
In [93]: data.describe()
```

```
Out[93]:
```

	0	1	2	3
count	1000.000000	1000.000000	1000.000000	1000.000000
mean	0.049091	0.026112	-0.002544	-0.051827
std	0.996947	1.007458	0.995232	0.998311
min	-3.645860	-3.184377	-3.745356	-3.428254
25%	-0.599807	-0.612162	-0.687373	-0.747478
50%	0.047101	-0.013609	-0.022158	-0.088274
75%	0.756646	0.695298	0.699046	0.623331
max	2.653656	3.525865	2.735527	3.366626

```
In [94]: col = data[2]
```

```
In [95]: col[np.abs(col) > 3]
```

```
Out[95]:
```

```
41      -3.399312
```

```
136     -3.745356
```

```
Name: 2, dtype: float64
```

```
In [96]: data[(np.abs(data) > 3).any(1)]
```

```
Out[96]:
```

	0	1	2	3
41	0.457246	-0.025907	-3.399312	-0.974657
60	1.951312	3.260383	0.963301	1.201206
136	0.508391	-0.196713	-3.745356	-1.520113
235	-0.242459	-3.056990	1.918403	-0.578828
258	0.682841	0.326045	0.425384	-3.428254
322	1.179227	-3.184377	1.369891	-1.074833
544	-3.548824	1.553205	-2.186301	1.277104
635	-0.578093	0.193299	1.397822	3.366626
782	-0.207434	3.525865	0.283070	0.544635
803	-3.645860	0.255475	-0.549574	-1.907459



```
In [97]: data[np.abs(data) > 3] = np.sign(data) * 3
```

```
In [98]: data.describe()
```

```
Out[98]:
```

	0	1	2	3
count	1000.000000	1000.000000	1000.000000	1000.000000
mean	0.050286	0.025567	-0.001399	-0.051765
std	0.992920	1.004214	0.991414	0.995761
min	-3.000000	-3.000000	-3.000000	-3.000000
25%	-0.599807	-0.612162	-0.687373	-0.747478
50%	0.047101	-0.013609	-0.022158	-0.088274
75%	0.756646	0.695298	0.699046	0.623331
max	2.653656	3.000000	2.735527	3.000000

```
In [99]: np.sign(data).head()
```

```
Out[99]:
```

	0	1	2	3
0	-1.0	1.0	-1.0	1.0
1	1.0	-1.0	1.0	-1.0
2	1.0	1.0	1.0	-1.0
3	-1.0	-1.0	1.0	-1.0
4	-1.0	1.0	-1.0	-1.0

# Computing Indicator/Dummy Variables

```
In [109]: df = pd.DataFrame({'key': ['b', 'b', 'a', 'c', 'a', 'b'],  
.....:                      'data1': range(6)})
```

```
In [110]: pd.get_dummies(df['key'])
```

```
Out[110]:
```

	a	b	c
0	0	1	0
1	0	1	0
2	1	0	0
3	0	0	1
4	1	0	0
5	0	1	0

```
In [111]: dummies = pd.get_dummies(df['key'], prefix='key')
```

```
In [112]: df_with_dummy = df[['data1']].join(dummies)
```

```
In [113]: df_with_dummy
```

```
Out[113]:
```

	data1	key_a	key_b	key_c
0	0	0	1	0
1	1	0	1	0
2	2	1	0	0
3	3	0	0	1
4	4	1	0	0
5	5	0	1	0

```
In [114]: mnames = ['movie_id', 'title', 'genres']
```

```
In [115]: movies = pd.read_table('datasets/movielens/movies.dat', sep='::',  
.....:                          header=None, names=mnames)
```

```
In [116]: movies[:10]
```

```
Out[116]:
```

	movie_id	title	genres
0	1	Toy Story (1995)	Animation Children's Comedy
1	2	Jumanji (1995)	Adventure Children's Fantasy
2	3	Grumpier Old Men (1995)	Comedy Romance
3	4	Waiting to Exhale (1995)	Comedy Drama
4	5	Father of the Bride Part II (1995)	Comedy
5	6	Heat (1995)	Action Crime Thriller
6	7	Sabrina (1995)	Comedy Romance
7	8	Tom and Huck (1995)	Adventure Children's
8	9	Sudden Death (1995)	Action
9	10	GoldenEye (1995)	Action Adventure Thriller

```
In [117]: all_genres = []
```

```
In [118]: for x in movies.genres:  
.....:     all_genres.extend(x.split('|'))
```

```
In [119]: genres = pd.unique(all_genres)
```

```
In [120]: genres
```

```
Out[120]:
```

```
array(['Animation', 'Children's', 'Comedy', 'Adventure', 'Fantasy',  
      'Romance', 'Drama', 'Action', 'Crime', 'Thriller', 'Horror',  
      'Sci-Fi', 'Documentary', 'War', 'Musical', 'Mystery', 'Film-Noir',  
      'Western'], dtype=object)
```

```
In [121]: zero_matrix = np.zeros((len(movies), len(genres)))
```

```
In [122]: dummies = pd.DataFrame(zero_matrix, columns=genres)
```

```
In [123]: gen = movies.genres[0]
```

```
In [124]: gen.split('|')
```

```
Out[124]: ['Animation', 'Children's', 'Comedy']
```

```
In [125]: dummies.columns.get_indexer(gen.split('|'))
```

```
Out[125]: array([0, 1, 2])
```

```
In [126]: for i, gen in enumerate(movies.genres):
.....:     indices = dummies.columns.get_indexer(gen.split('|'))
.....:     dummies.iloc[i, indices] = 1
.....:
```

```
In [127]: movies_windic = movies.join(dummies.add_prefix('Genre_'))
```

```
In [128]: movies_windic.iloc[0]
```

```
Out[128]:
```

movie_id		1
title	Toy Story (1995)	
genres	Animation Children's Comedy	
Genre_Animation		1
Genre_Children's		1
Genre_Comedy		1
Genre_Adventure		0
Genre_Fantasy		0
Genre_Romance		0
Genre_Drama		0
	...	
Genre_Crime		0
Genre_Thriller		0
Genre_Horror		0
Genre_Sci-Fi		0
Genre_Documentary		0
Genre_War		0
Genre_Musical		0
Genre_Mystery		0
Genre_Film-Noir		0
Genre_Western		0
Name: 0, Length: 21, dtype: object		

# String Manipulation

## String Object Methods

```
In [134]: val = 'a,b, guido'
```

```
In [135]: val.split(',')  
Out[135]: ['a', 'b', ' guido']
```

```
In [136]: pieces = [x.strip() for x in val.split(',')] 
```

```
In [137]: pieces  
Out[137]: ['a', 'b', 'guido']
```

```
In [138]: first, second, third = pieces
```

```
In [139]: first + '::' + second + '::' + third  
Out[139]: 'a::b::guido'
```

```
In [140]: '::'.join(pieces)  
Out[140]: 'a::b::guido'
```

```
In [141]: 'guido' in val  
Out[141]: True
```

```
In [142]: val.index(',')  
Out[142]: 1
```

```
In [143]: val.find(':')  
Out[143]: -1
```

```
In [145]: val.count(',')  
Out[145]: 2
```

```
In [146]: val.replace(',', '::')  
Out[146]: 'a::b:: guido'
```

```
In [147]: val.replace(',', ' ')  
Out[147]: 'ab guido'
```

# 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>lstrip</code>	Trim whitespace, including newlines; equivalent to <code>x.strip()</code> (and <code>rstrip</code> , <code>lstrip</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> , ... to refer to match group elements in the replacement string

---

# Vectorized String Functions in pandas

```
In [167]: data = {'Dave': 'dave@google.com', 'Steve': 'steve@gmail.com',  
.....:          'Rob': 'rob@gmail.com', 'Wes': np.nan}
```

```
In [168]: data = pd.Series(data)
```

```
In [169]: data
```

```
Out[169]:
```

```
Dave    dave@google.com  
Rob      rob@gmail.com  
Steve   steve@gmail.com  
Wes                      NaN  
dtype: object
```

```
In [170]: data.isnull()
```

```
Out[170]:
```

```
Dave    False  
Rob      False  
Steve   False  
Wes      True  
dtype: bool
```

```
In [171]: data.str.contains('gmail')
```

```
Out[171]:
```

```
Dave    False  
Rob      True  
Steve   True  
Wes      NaN  
dtype: object
```

```
In [172]: pattern
```

```
Out[172]: '([A-Z0-9._%+-]+)@([A-Z0-9.-]+\.[A-Z]{2,4})'
```

```
In [173]: data.str.findall(pattern, flags=re.IGNORECASE)
```

```
Out[173]:
```

```
Dave    [(dave, google, com)]  
Rob      [(rob, gmail, com)]  
Steve   [(steve, gmail, com)]  
Wes                      NaN  
dtype: object
```



```
In [174]: matches = data.str.match(pattern, flags=re.IGNORECASE)
```

```
In [175]: matches
```

```
Out[175]:
```

```
Dave      True
```

```
Rob       True
```

```
Steve     True
```

```
Wes       NaN
```

```
dtype: object
```

```
In [176]: matches.str.get(1)
```

```
Out[176]:
```

```
Dave      NaN
```

```
Rob       NaN
```

```
Steve     NaN
```

```
Wes       NaN
```

```
dtype: float64
```

```
In [177]: matches.str[0]
```

```
Out[177]:
```

```
Dave      NaN
```

```
Rob       NaN
```

```
Steve     NaN
```

```
Wes       NaN
```

```
dtype: float64
```

```
In [178]: data.str[:5]
```

```
Out[178]:
```

```
Dave      dave@
```

```
Rob       rob@g
```

```
Steve     steve
```

```
Wes       NaN
```

```
dtype: object
```

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, 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