

Computing Laboratory

More on Python

Malay Bhattacharyya

Assistant Professor

Machine Intelligence Unit
Indian Statistical Institute, Kolkata

December, 2020



- 1 More on Lists
- 2 Randomization
- 3 Random Sampling
- 4 Data Structures in Python
 - Series
 - DataFrame
- 5 Problems

List comprehension

```
import random
ls = [10, 20, 30, 40, 50]
random.shuffle(ls)
print(ls)
```

Output: [40, 20, 10, 50, 30]

Shuffling of data

```
import random
ls = [10, 20, 30, 40, 50]
random.shuffle(ls)
print(ls)
```

Output: [40, 20, 10, 50, 30]

Randomization of data

```
import random
random.randint(1, 100) # The interval is [1, 100)
```

Output: 15

```
import random
random.random() # The interval is [0.0, 1.0)
```

Output: 0.4289859005273219

Random sampling with numpy

```
import numpy as np
np.random.randint(1, 100) # From uniform distribution
```

Output: 97

The interval is [1, 100)

Random sampling with numpy

```
import numpy as np
np.random.randint(1, 100) # From uniform distribution
```

Output: 97

The interval is [1, 100)

```
import numpy as np
np.random.random(5) # From uniform distribution
```

Output: array([0.37076725, 0.10547203, 0.21298417, 0.96296838,
0.15390583])

The interval is [0.0, 1.0)

Random sampling with numpy

```
import numpy as np
np.random.randint(1, 100) # From uniform distribution
```

Output: 97

The interval is [1, 100)

```
import numpy as np
np.random.random(5) # From uniform distribution
```

Output: array([0.37076725, 0.10547203, 0.21298417, 0.96296838, 0.15390583])

The interval is [0.0, 1.0)

Note: Unlike the `random()` function in `random` module, the one in the `random` submodule of `numpy` module optionally takes size of the array.

Random sampling with numpy

```
import numpy as np
np.random.randn(5) # From normal distribution
```

Output: array([1.62029576, 0.29112406, 1.21198839,
0.26851418, -0.46712281])

```
import numpy as np
np.random.randn(3, 3) # From normal distribution
```

Output: array([[1.13217437, 0.56093212, -2.36792091],
[1.42652606, 0.68953983, 0.521175],
[1.36766496, 0.9488446 , -1.10042531]])

Effect of the seed value on random sampling

Code 1:

```
import numpy as np
print(np.random.randint(1, 100))
print(np.random.randint(1, 100))
print(np.random.randint(1, 100))
```

Effect of the seed value on random sampling

Code 1:

```
import numpy as np
print(np.random.randint(1, 100))
print(np.random.randint(1, 100))
print(np.random.randint(1, 100))
```

Code 2:

```
import numpy as np
np.random.seed(20)
print(np.random.randint(1, 100))
print(np.random.randint(1, 100))
print(np.random.randint(1, 100))
```

Effect of the seed value on random sampling

Code 1:

```
import numpy as np
print(np.random.randint(1, 100))
print(np.random.randint(1, 100))
print(np.random.randint(1, 100))
```

Code 2:

```
import numpy as np
np.random.seed(20)
print(np.random.randint(1, 100))
print(np.random.randint(1, 100))
print(np.random.randint(1, 100))
```

Note: Unlike Code 1, Code 2 will generate the same set of random values whenever executed.

Drawing samples from different distributions using numpy

<code>beta(a, b[, size])</code>	– Samples from Beta distribution
<code>binomial(n, p[, size])</code>	– Samples from binomial distribution
<code>chisquare(df[, size])</code>	– Samples from chi-square distribution
<code>dirichlet(alpha[, size])</code>	– Samples from Dirichlet distribution
<code>exponential([scale, size])</code>	– Samples from exponential distribution
<code>f(dfnum, dfden[, size])</code>	– Samples from F distribution
<code>gamma(shape[, scale, size])</code>	– Samples from Gamma distribution
<code>geometric(p[, size])</code>	– Samples from geometric distribution
<code>gumbel([loc, scale, size])</code>	– Samples from Gumbel distribution

Table: Functions in `random` submodule of `numpy`

Drawing samples from different distributions using numpy

<code>hypergeometric(ngood, nbad, nsample[, size])</code>	– Samples from hypergeometric distribution
<code>laplace([loc, scale, size])</code>	– Samples from Laplace distribution
<code>logistic([loc, scale, size])</code>	– Samples from logistic distribution
<code>lognormal([mean, sigma, size])</code>	– Samples from log-normal distribution
<code>logseries(p[, size])</code>	– Samples from logarithmic series distribution

Table: Functions in `random` submodule of `numpy`

Drawing samples from different distributions using numpy

<code>multinomial(n, pvals[, size])</code>	– Samples from multinomial distribution
<code>multivariate_normal(mean, cov[, size, ...])</code>	– Samples from multivariate normal distribution
<code>negative_binomial(n, p[, size])</code>	– Samples from negative binomial distribution
<code>noncentral_chisquare(df, nonc[, size])</code>	– Samples from noncentral chi-square distribution
<code>noncentral_f(dfnum, dfden, nonc[, size])</code>	– Samples from noncentral F distribution
<code>normal([loc, scale, size])</code>	– Samples from normal distribution

Table: Functions in `random` submodule of `numpy`

Drawing samples from different distributions using numpy

<code>pareto(a[, size])</code>	– Samples from pareto distribution
<code>poisson([lam, size])</code>	– Samples from Poisson distribution
<code>power(a[, size])</code>	– Samples from power distribution
<code>rayleigh([scale, size])</code>	– Samples from Rayleigh distribution
<code>standard_cauchy([size])</code>	– Samples from standard Cauchy distribution
<code>standard_exponential([size])</code>	– Samples from standard exponential distribution
<code>standard_gamma(shape[, size])</code>	– Samples from standard Gamma distribution
<code>standard_normal([size])</code>	– Samples from standard normal distribution

Table: Functions in `random` submodule of `numpy`

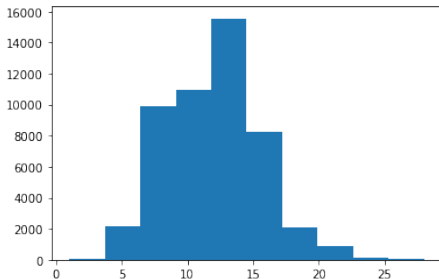
Drawing samples from different distributions using numpy

<code>standard_t(df[, size])</code>	– Samples from standard Student's t distribution
<code>triangular(left, mode, right[, size])</code>	– Samples from triangular distribution
<code>uniform([low, high, size])</code>	– Samples from uniform distribution
<code>vonmises(mu, kappa[, size])</code>	– Samples from von Mises distribution
<code>wald(mean, scale[, size])</code>	– Samples from Wald distribution
<code>weibull(a[, size])</code>	– Samples from Weibull distribution
<code>zipf(a[, size])</code>	– Samples from Zipf distribution

Table: Functions in `random` submodule of `numpy`

Visualizing a distribution

```
import numpy as np
sample = np.random.poisson(10, 50000)
import matplotlib.pyplot as plt
plt.hist(sample, 10)
plt.show()
```



Basics of Series

Series is a one-dimensional and heterogeneous data structure with labeled entries. It is capable of holding any kind of data type (integers, strings, real values, objects, etc.).

Note: The labels are denoted as index.

Using a Series

```
import pandas as pd
pd.Series(Data, index=<Labels>)
```

Using a Series

```
import pandas as pd
pd.Series(Data, index=<Labels>)
```

The Data can be any of the following:

- A scalar value (integer, character, string, etc.)
- An ndarray
- A dictionary

Using a Series

```
import pandas as pd
pd.Series(Data, index=<Labels>)
```

The Data can be any of the following:

- A scalar value (integer, character, string, etc.)
- An ndarray
- A dictionary

Note: If no <Labels> are passed, one will be created having values $[0, \dots, \text{len}(\text{Data}) - 1]$.

Defining Series with a scalar value

```
import pandas as pd
s = pd.Series(10., index=['r1', 'r2', 'r3', 'r4', 'r5'])
print(s)
```

Defining Series with a scalar value

```
import pandas as pd
s = pd.Series(10., index=['r1', 'r2', 'r3', 'r4', 'r5'])
print(s)
```

Output:

```
r1    10.0
r2    10.0
r3    10.0
r4    10.0
r5    10.0
dtype: float64
```


Defining Series with an ndarray

```
import pandas as pd
s = pd.Series([1, 2, 3, 4], index=['a', 'b', 'c', 'd'])
print(s)
```

Defining Series with an ndarray

```
import pandas as pd
s = pd.Series([1, 2, 3, 4], index=['a', 'b', 'c', 'd'])
print(s)
```

Output:

```
a      1
b      2
c      3
d      4
dtype: int64
```

Defining Series with an ndarray

```
import pandas as pd
s = pd.Series(['malay', 'amitava', 'utpal'])
print(s)
```

Defining Series with an ndarray

```
import pandas as pd
s = pd.Series(['malay', 'amitava', 'utpal'])
print(s)
```

Output:

```
0      malay
1  amitava
2      utpal
dtype: object
```

Defining Series with an ndarray

```
import numpy as np
import pandas as pd
s = pd.Series(np.random.randn(5))
print(s)
```

Defining Series with an ndarray

```
import numpy as np
import pandas as pd
s = pd.Series(np.random.randn(5))
print(s)
```

Output:

```
0    -0.594763
1    -0.410549
2     1.153695
3     0.454634
4    -0.665784
dtype: float64
```

Defining Series with a dictionary

```
import pandas as pd
s = pd.Series({'name': 'MB', 'age': 37, 'gender': 'M'})
print(s)
```

Defining Series with a dictionary

```
import pandas as pd
s = pd.Series({'name': 'MB', 'age': 37, 'gender': 'M'})
print(s)
```

Output:

```
name      MB
age       37
gender    M
dtype: object
```


Pulling data from a Series

```
import pandas as pd
s = pd.Series({'name': 'MB', 'age': 37, 'gender': 'M'})
print(s)
s0 = pd.Series(s, index=['name', 'gender', 'height'])
print(s0)
```

Pulling data from a Series

```
import pandas as pd
s = pd.Series({'name': 'MB', 'age': 37, 'gender': 'M'})
print(s)
s0 = pd.Series(s, index=['name', 'gender', 'height'])
print(s0)
```

Output:

```
name      MB
age       37
gender    M
dtype: object
name      MB
gender    M
height   NaN
dtype: object
```

Pulling data from a Series

```
import pandas as pd
s = pd.Series({'name': 'MB', 'age': 37, 'gender': 'M'})
print(s[1])
print(s[1:2])
print(s[[2, 1]])
```

Pulling data from a Series

```
import pandas as pd
s = pd.Series({'name': 'MB', 'age': 37, 'gender': 'M'})
print(s[1])
print(s[1:2])
print(s[[2, 1]])
```

Output:

```
37
age      37
dtype: object
gender    M
age      37
dtype: object
```

Searching in a Series

```
import pandas as pd
s = pd.Series({'name': 'MB', 'age': 37, 'gender': 'M'})
print(s['name'])
print('age' in s)
print('height' in s)
print('MB' in s) # Searching for the index 'MB'
```

Searching in a Series

```
import pandas as pd
s = pd.Series({'name': 'MB', 'age': 37, 'gender': 'M'})
print(s['name'])
print('age' in s)
print('height' in s)
print('MB' in s) # Searching for the index 'MB'
```

Output:

MB

True

False

False

Basics of DataFrame

DataFrame is a two-dimensional, heterogeneous and size-mutable data structure with entries labeled by rows and columns.

Note: The labels are denoted as index.

Using a DataFrame

```
import pandas as pd
pd.DataFrame(Data, index=<Labels>)
```


Using a DataFrame

```
import pandas as pd
pd.DataFrame(Data, index=<Labels>)
```

A DataFrame can have any of the following inputs:

- A dictionary of 1D ndarrays, lists, dicts, or Series
- 2-D ndarray
- Structured or record ndarray
- A Series
- Another DataFrame

Using a DataFrame

```
import pandas as pd
pd.DataFrame(Data, index=<Labels>)
```

A DataFrame can have any of the following inputs:

- A dictionary of 1D ndarrays, lists, dicts, or Series
- 2-D ndarray
- Structured or record ndarray
- A Series
- Another DataFrame

Note: If no <Labels> are passed, one will be created having values $[0, \dots, \text{len}(\text{Data}) - 1]$.

Defining DataFrame with a dictionary of 1D ndarrays/lists

```
import pandas as pd
s = pd.DataFrame({'col1': [1, 2, 3, 4], 'col2': [5, 6,
7, 8]}, index=['a', 'b', 'c', 'd'])
print(s)
```

Defining DataFrame with a dictionary of 1D ndarrays/lists

```
import pandas as pd
s = pd.DataFrame({'col1': [1, 2, 3, 4], 'col2': [5, 6,
7, 8]}, index=['a', 'b', 'c', 'd'])
print(s)
```

Output:

```
      col1  col2
a         1     5
b         2     6
c         3     7
d         4     8
dtype: int64
```

Defining DataFrame with a dictionary of dicts/series

```
import pandas as pd
s = pd.DataFrame({'col1': pd.Series([.1, .2, .3, .4,
    .5], index=['a', 'b', 'c', 'd', 'e']), 'col2':
    pd.Series([1., 2., 3.], index=['a', 'b', 'd'])})
print(s)
```

Defining DataFrame with a dictionary of dicts/series

```
import pandas as pd
s = pd.DataFrame({'col1': pd.Series([.1, .2, .3, .4,
    .5], index=['a', 'b', 'c', 'd', 'e']), 'col2':
    pd.Series([1., 2., 3.], index=['a', 'b', 'd'])})
print(s)
```

Output:

	col1	col2
a	0.1	1.0
b	0.2	2.0
c	0.3	NaN
d	0.4	3.0
e	0.5	NaN

Defining DataFrame with a structured or record ndarray

```
import numpy as np
import pandas as pd
data = np.zeros((2, ), dtype=[('Year', 'i4'),
                              ('Company', 'a12'), ('Share/Mutual Fund', 'f4')])
print(data)
data[:] = [(1907, 'TATA STEEL', 344.00), (2019, "TATA
DIGITAL INDIA FUND", '19.14')]
s = pd.DataFrame(data)
print(s)
```

Defining DataFrame with a structured or record ndarray

```
import numpy as np
import pandas as pd
data = np.zeros((2, ), dtype=[('Year', 'i4'),
                              ('Company', 'a12'), ('Share/Mutual Fund', 'f4')])
print(data)
data[:] = [(1907, 'TATA STEEL', 344.00), (2019, "TATA
DIGITAL INDIA FUND", '19.14')]
s = pd.DataFrame(data)
print(s)
```

Output:

```
[(0, b'', 0.) (0, b'', 0.)]
```

	Year	Company	Share/Mutual Fund
0	1907	b'TATA STEEL'	344.000000
1	2019	b'TATA DIGITAL'	19.139999

Problems – Day 4

- 1 Write a function `uniquify` that
 - takes two arguments: an array `A` containing at most 100 non-negative integers, and `n`, the number of integers contained in the array, and
 - stores in `A`, a list of the distinct integers contained in `A` if the integers in `A` are sorted in increasing order, and returns the number of distinct integers in `A`;
 - returns -1 otherwise.

- 2 Suppose that we choose to represent sets by arrays of non-negative integers, containing at most 100 elements each. Let `A` and `B` be 2 such sets of integers. Implement the following set operations:
 - $A \cup B$;
 - $A \cap B$;
 - $A \subseteq B$ (returns 1 if `A` is a subset of `B`, 0 otherwise);
 - prints the power set of `A` on the screen.

Problems – Day 4

- 3** Write a program that will print '0' with a probability of 0.5 and '1' in rest of the cases. Now write a function RUN to count the number of runs in the binary stream generated with your program. The number of runs will simply explore and count dichotomous (reflecting contrast between two things) sequence of values.

Note: The following sequence of 0's and 1's has a run count of nine.

0 0 1 1 0 1 0 0 0 0 1 0 0 1 1 1 1 1 0 0

- 4** Write a program to find out the summation of the following series given n as the user input.

$$\frac{1}{1} + \frac{11}{12} + \frac{111}{123} + \dots + \frac{111 \dots n \text{ times}}{123 \dots n}$$