# Chapter Three
# Data Handling in Python

Kibret Zewdu (MSc.)

Department of Computer Science
College of Informatics
University of Gondar

January 30, 2023

# Outline

# Data Types

- **Data Type** specifies which type of value a variable can store.
- **Data Types In Python**
  1. Number
  2. String
  3. Boolean
  4. List
  5. Tuple
  6. Set
  7. Dictionary
- ***type()***: function is used to determine a variable's type in Python.

# Number in Python

- used to store numeric values.
- Python has three numeric types:
  - Integers
  - Floating point numbers
  - Complex numbers.

# Number in Python

## Integers

- Integers or int are positive or negative numbers with no decimal point.
- Integers in Python 3 are of unlimited size.

**Example**:

```
1    a= 100
2    b= -100
3    c= 1*20
4    print(a)
5    print(b)
6    print(c)
7
```

# Number in Python

## Type Conversion of Integer

**int()** function converts any data type to integer.

## Example

```
1    a = "101" # string
2    b=int(a) # converts string data type to integer.
3    c=int(122.4) # converts float data type to integer.
4    print(b)
5    print(c)
6
```

## Output
101
122

# Number in Python

## Floating point numbers

- It is a positive or negative real numbers with a decimal point.

## Example

```
1    a = 101.2
2    b = -111.23
3    c = 2.3*3
4    print(a)
5    print(b)
6    print(c)
7
```

## Output

101.2

-111.23

6.8999999999999995

# Number in Python

## Type Conversion of Floating point numbers

**float()** function converts any data type to floating point number.

## Example

```python
a = "120.35" # string
b=float(a) # converts string data type to floating
point number.
c=float(122) # converts integer data type to
floating point number.
print(b)
print(c)
```

## Output

120.35

122.0

# Number in Python

## Complex numbers

- Complex numbers are combination of a real and imaginary part.
- Complex numbers are in the form of X+Yj, where X is a real part and Y is imaginary part.

## Example

```
1   a = complex(5) # convert 5 to a real part  and zero
     imaginary part
2   print(a)
3   b=complex(101,23) # convert 101 with real part and
     23 as imaginary part
4   print(b)
5
```

## Output
(5+0j)
(101+23j)

# String in Python

- A **string** is a sequence of symbols/characters delimited by a single quote ('), double quotes("), triple single quotes ('''), or triple double quotes (""").
- Python 3, all string sequences are Unicode characters by default.

**Example**

```python
str = 'Computer Science'
print('str-', str) # print Computer Science
print('str[0]-', str[0])# print first char 'C'
print('str[1:3]-', str[1:3]) # print string from
position 1 to 3 'om'
print('str[3:]-', str[3:]) # print string staring
from 3rd char 'puter Science'
print('str *2-', str *2 ) # print string two times
print("str +'yes'-", str +'yes') # concatenated
string
```

# String in Python

**Output**

str- Computer Science

str[0]- c

str[1:3]- om

str[3:]- puter Science

str *2- Computer scienceComputer Science

str +'yes'- Computer Scienceyes

# String in Python

- We can insert a single quote in a string delimited for double quote, and vice versa:

```
1    "A single quote (') inside a double quote"
2    'Here "double quotes" inside single quotes'
3
```

- Remember is that if we begin a string with a type of quote, we must finish it with the same type of quote.
- The following string is not valid:

```
1    >>> "Mixing quote types leads to the dark side'
2
```

File ¨<stdin>¨; line 1
"Mixing quote types leads to the dark side'
SyntaxError: EOL while scanning single-quoted string

# String
## String Manipulation

- **Strings** are immutable. Once a string is created, it can't be modified.
- In the following example there is a string that represents an amino-acid sequence and it is called signal_peptide:

```
1   >>> signal_peptide = 'MASKATLLLAFTLLFATCIA'
2   >>> signal_peptide.lower()
3   'maskatlllaftllfatcia'
4   #The original string has not been modifie
5   >>> signal_peptide
6   'MASKATLLLAFTLLFATCIA'
7   #To get new lower case string with the same name
    as the previous one, we need to assign it:
8   >>> signal_peptide = signal_peptide.lower()
9   >>> signal_peptide
10  'maskatlllaftllfatcia'
11
```

# String

Methods Associated with Strings

1. **replace**(old,new[,count]): Allows us to replace a portion of a string (old) with another (new). If the optional argument count is used, only the first count occurrences of old will be replaced:

```
1    >>> dna_seq = 'GCTAGTAATGTG'
2    >>> m_rna_seq = dna_seq.replace('T','U')
3    >>> m_rna_seq
4    'GCUAGUAAUGUG'
5
```

2. **count**(sub[, start[, end]]): Counts how many times the substring sub appears, between the start and end positions (if available).

```
1    >>> dna_seq
2    'GCTAGTAATGTG'
3    >>> g = dna_seq.count("G")
4    4
5
```

# String
Methods Associated with Strings

- **find**(sub[,start[,end]]): Returns the position of the substring sub, between the start and end positions (if available). If the substring is not found in the string, it returns the value -**1**:

```
1    >>> m_rna_seq
2    'GCUAGUAAUGUG'
3    >>> m_rna_seq.find('AUG')
4    7
5    >>> m_rna_seq.find('GGG')
6    -1
7
```

- **index**(sub[,start[,end]]): Works like **find**(). **index** will raise a **ValueError** exception when the substring is not found. This method is recommended over **find()** because the value -**1** could be interpreted as a valid value, while a ValueError returned by index() can't be taken as a valid value.

# String
Methods Associated with Strings

- **split**([sep [,maxsplit]]): Separates the "words" of a string and returns them in a list.
  If a separator (sep) is not specified, the default separator will be a white space:

```
1    >>> 'This string has words separated by spaces'
     .split()
2    ['This', 'string', 'has', 'words', 'separated',
     'by', 'spaces']
3    >>> "Tedy Afro, aa-2333, tedyafr@example.com".
     split(",")
4    ['Tedy Afro', 'aa-2333', 'tedyafr@example.com']
5
```

- **Bioinformatics Application**: Parsing *BLAST* Files.
- The BLAST standalone executable can generate output as a "tab separated file".

# String
## Methods Associated with Strings

- **join([seq)**: Joins the sequence using a string as a "glue character":

```
1    ';'.join(['Alex Doe', '5555-2333', '
     nobody@example.com'])
2    'Alex Doe;5555-2333;nobody@example.com'
3
```

- To join a sequence without any glue character, use empty quotes (**""**):

```
1    >>> ''.join(['A','C','A','T'])
2    'ACAT'
3
```

# Boolean In Python

It is used to store two possible values either **true** or **false**

**Example**

```
1   str="comp sc"
2   upper_case=str.isupper() \# test if string contains
      upper case
3   isMale = True \# Boolean true value
4   print(upper_case)
5   print (isMale)
6
7
```

**Output**

```
1   False
2   True
3
```

# **List** In Python

## List

- List are ordered collections of items.
- It is represented by elements separated by commas and enclosed between **square brackets**.
- A list can hold similar or different kinds of data.
- You can also create an empty list and add elements later.

- To define and name a list:

```
1    >>> list1 = [1, 2, 3, 4, 5]
2    >>> list2 = [1, 'two', 3, 4, 'last']
3    >>> nested_list = [1, 'two', list1, 4]
4    >>> empty_list = []
5
```

# **List** In Python
Accessing List Elements In Python

- list elements can be accessed using their index starting at zero.

```
1    >>> first_list = [1, 2, 3, 4, 5]
2    first_list [0]
3    first_list [1]
4
```

- Negative numbers are used to access lists from the right:

```
1    >>> first_list = [1, 2, 3, 4, 5]
2    >>> first_list [-1] \# returns 5
3    >>> first_list [-4] \# returns 2
4
```

# List
List with Multiple Repeated Items

- You can also turn a non-list object into a list by using the built-in function `list()`:

```
1    >>> aseq = "atggctaggc"
2    >>> list(aseq)
3    ['a', 't', 'g', 'g', 'c', 't', 'a', 'g', 'g', '
     c']
4
```

- To initialize a list with the same item repeated multiple times, you can use the **\*** operator like this:

```
1    >>> samples = ['red'] * 5
2    >>> samples
3    ['red', 'red', 'red', 'red', 'red']
4
```

# List
List Comprehension

- A list can be created from another list.

```
1    >>> [3 * x for x in a]
2    [0, 3, 6, 9, 12, 15]%
3
```

- To find all animals whose name contains letter 'i':

```
1    >>> animals = [' Lion', ' Dog ', 'Tiger ']
2    >>> [x.strip() for x in animals if 'i' in x]
3    ['lion', 'Tiger']
4
```

# List
## Modifying Lists

- Lists can be modified by adding, removing, or changing their elements.
- **Adding**: There are three ways to add elements into a list: **append**, **insert**, and **extend**.

  1. **append(element)**: Adds an element at the end of the list.

     ```
     1    >>> first_list.append(99)
     2
     ```

  2. **insert(position,element)**: Inserts the element element at the position position.

     ```
     1    >>> first_list.insert(2,50)
     2
     ```

  3. **extend(list)**: Extends a list by adding a list to the end of the original list.

     ```
     1    >>> first_list.extend([6,7,8])
     2    >>> [1,2,3] + [4,5] # + is the same
     ```

# List
## Modifying Lists

- **Removing**: There are three ways to remove elements from a list: **pop**, **remove** and **del**
  - **pop(/index/)**: Removes the element in the index position and returns it to the point where it was called.
  - Without parameters, it returns the last element.

```
1    >>> first_list
2    [1, 2, 50, 3, 4, 5, 99, 6, 7, 8]
3    >>> first_list.pop()
4    8
5    >>> first_list.pop(2)
6    50
7    >>> first_list
8    [1, 2, 3, 4, 5, 99, 6, 7]
9
```

# List
## Modifying Lists

- **remove(element)**
    - Removes the element specified in the parameter.
    - In the case where there is more than one copy of the same object in the list, it removes the first one, counting from the left.
    - this function does not return anything.

```
1    >>> first_list.remove(99)
2    >>> first_list
3    [1, 2, 3, 4, 5, 6, 7]
4
```

    - Trying to remove a nonexistent element raises an error
- **del**: it is a command to delete an element from the list.

```
1    del first_list[0]
2
```

# List
## Copying a List

- To copy a list, use the **copy** method from the copy module.

```
1   >>> import copy
2   >>> a = [1, 2, 3]
3   >>> b = copy.copy(a)
4   >>> b.pop()
5   3
6   >>> a
7   [1, 2, 3]
8
```

- Without using the **copy** module

```
1   b = a[:]
2
```

- Note: The assignment operator (=) doesn't copy the values, it copies a reference to the original object.

# Tuples In Python
## Tuples are Immutable Lists

- A **tuple** is a collection of ordered objects with the characteristic that once created, it cannot be modified. They are referred to as "immutable lists."
- **Immutable objects** cannot be modified after they are created.
- **Tuple** elements are enclosed between parentheses:

```
1    point = (23, 56, 11) \# Here elements are
     enclosed by ()
2
```

- When the tuple has only one element, you should use a trailing comma to tell python interpreter that it is a tuple not an expression:

```
1    lone_element_tuple = (5,)
2
```

# Common Properties of of Sequence

Indexing

- You can apply these properties to lists, tuples, and strings.
- **Indexing**: elements in the sequences are ordered, we can gain access to any element through an index that begins at zero.

```
1  >>> point = (23, 56, 11)
2  >>> point[0]
3  23
4  >>> sequence = 'MRVLLVALALLALAASATS'
5  >>> sequence[5]
6  'V'
7  >>> parameters = ['UniGene', 'dna', 'Mm.248907'
   , 5]
8  >>> parameters[2]
9  'Mm.248907'
10
```

# Common Properties of of Sequence
Indexing

- To access the elements of a sequence from the right by using **negative numbers**

```
1  >>> parameters [-3]
2  'dna'
3  >>> sequence [-1]
4  'S'
5
```

- To access an element that is inside a sequence, which is itself inside another sequence, you need to use another index:

```
1  >>> seqdata = ('MRVLLVALALLA', 12, '5
    FE9EEE8EE2DC2C7')
2  >>> seqdata[0][5] # The 6th element of the first
    element in the sequence
3  'V'
4
```

# Common Properties of of Sequence
**Slicing**

- used to select a portion of a sequence.
- consists of using two indexes separated by a colon (**:**).

```
1   >>> string ="Python"
2   >>> string[0:2]
3   'Py'
4
```

- If the first sub-index is omitted, the index value defaults to the first position (**0**):

```
1   >>> string[:2]
2   'Py'
3
```

- If the second sub-index is omitted, the index value defaults to the last position (-1)

```
1   >>> string[4:] # from index 4 to the end
2   'on'
```

# Common Properties of of Sequence
**Slicing**

- There is a third, optional index to skip positions (step argument):

```
1    >>> string[1:5]
2    'ytho'
3    >>> string[1:5:2]
4    'yh'
5
```

- A step with a negative number is used to count backwards. So -1 (in the third position) can be used to invert a sequence:

```
1    >>> string[::-1]
2    'nohtyP'
3
```

- Note that slicing always returns another sequence

# Common Properties of of Sequence
**Membership Test**

- To verify whether an element belongs to a sequence, using the
  **in** keyword:

```
1   >>> point = (23, 56, 11)
2   >>> 11 in point
3   True
4   >>> my_sequence = 'MRVLLVALALLALAASATS'
5   >>> 'X' in my_sequence
6   False
7
```

# Common Properties of of Sequence
**Concatenation**

- concatenate two or more sequences of the same class using the "$+$" sign:

```
1    >>> point = (23, 56, 11)
2    >>> point2 = (2, 6, 7)
3    >>> point + point2
4    (23, 56, 11, 2, 6, 7)
5    >>> dna_seq = 'ATGCTAGACGTCCTCAGATAGCCG'
6    >>> tata_box = 'TATAAA'
7    >>> tata_box + dna_seq
8    'TATAAAATGCTAGACGTCCTCAGATAGCCG'
9
```

- Sequences of different types can't be concatenated.

# Common Properties of of Sequence

**len, max, and min**

- **len()** returns the length (the number of items) of a sequence:

```
1   >>> my_sequence = 'MRVLLVALALLALAASATS'
2   >>> len(my_sequence)
3   19
4
```

- **max()** and **min()** applied over a sequence of numbers return, as expected, the maximum and the minimum value:

```
1   >>> min(point)
2   11
3
```

- **max()** and **min()** applied to strings return a character according to the maximum or minimum value of its ASCII code:

```
1   >>> my_sequence = 'MRVLLVALALLALAASATS'
2   >>> max(my_sequence)
3   'V'
```

# Common Properties of of Sequence
**Turn a Sequence into a List**

- To convert a sequence (like a tuple or a string) into a list, use the list() method:

```
1  >>> tata_box = 'TATAAA'
2
3  >>> list(tata_box)
4
5  ['T', 'A', 'T', 'A', 'A', 'A']
```

- Using a list provides us with methods to indirectly modify a string

# Dictionary In Python

## Dictionary

- are **unordered collection** of items and each item consist of a key and a value.
- They are defined by enclosing is *key:value* pairs between curly brackets **{}**
- The key is the index used to retrieve the value.

## Example

```python
my_dict = {'subject': 'Python', 'class': 'Pg
 Building'}
print(my_dict)
print ("Subject : ", my_dict['subject'])

```

**Output**: {'subject': 'Python ', 'class': 'Pg Building'}
Subject : 'Python'

# Dictionary In Python

- This dictionary works as a translation table that allows us to translate between the one-letter amino acid code to a three-letter code.

```
1   >>> iupac = {'A':'Ala','C':'Cys','E':'Glu'}
2   >>> iupac['E']
3   'Glu'
4
```

- Not every object can be used as a dictionary key.
- Only immutable objects like strings, tuples and numbers can be used as keys.

# Dictionary In Python

- A dictionary can also be created from a sequence with **dict()**:

```
1    >>> rgb = [('red','ff0000'), ('green','00ff00')
     , ('blue','0000ff')]
2    >>> colors = dict(rgb)
3    >>> colors
4    {'red': 'ff0000', 'green': '00ff00', 'blue': '
     0000ff'}
5
```

- **dict** also accepts `name=value` pairs in the argument list:

```
1    >>> rgb = dict('red'='ff0000', 'green'='00ff00'
     , 'blue'='0000ff')
2
```

- To create an empty dictionary and add elements later:

```
1    >>> rgb = {}
2    >>> rgb['red'] = 'ff0000'
3    >>> rgb['green'] = '00ff00'
4
```

# Dictionary In Python

- **len()**, returns the number of elements in the dictionary:

```
1   >>> len(iupac)
2   3
3
```

- To add values to a dictionary,

```
1   >>> iupac['S'] = 'Ser'
2
```

- Dictionaries are unordered because they don't keep track of the order of their elements.

- When you request to see the contents of the dictionary, you may or may not get the elements in the same order as they were entered

# Dictionary In Python
## Dictionary Methods

- To get the keys or values of a dictionary, there are methods like **keys()** and **values()**:

```
1   >>> iupac = {'E': 'Glu', 'X': 'Xaa', 'C': 'Cys'
    , 'A': 'Ala'}
2   >>> iupac.keys()
3   dict_keys(['E', 'X', 'C', 'A'])
4   >>> iupac.values()
5   dict_values(['Glu', 'Xaa', 'Cys', 'Ala'])
6
```

- **items()** is another method to access elements of dictionaries which returns a **dictionary view** a tuple of every key-value pair.

```
1   >>> iupac.items()
2   dict_items([('E', 'Glu'), ('A', 'Ala'), ('C', '
    Cys'), ('X', 'Xaa')])
3
```

# Dictionary In Python
## Dictionary Methods

| Properties | Description |
|---|---|
| len(d) | Number of elements of $d$ |
| d[k] | The element from $d$ that has a $k$ key |
| d[k] = v | Set d[k] to $v$ |
| del d[k] | Remove d[k] from $d$ |
| d.clear() | Remove all items from $d$ |
| d.copy() | Copy $d$ |
| k in d | True if d has a key $k$, else False |
| k not in d | Equivalent to not $k$ in $d$ |
| d.has_key(k) | Equivalent to $k$ in $d$, use that form in new code |
| d.items() | A copy of $d$'s list of (key, value) pairs |
| d.keys() | A copy of $d$'s list of keys |
| d.update([b]) | Updates (and overwrites) key/value pairs from $b$ |
| d.fromkeys(seq[, value]) | Creates a new dictionary with keys from $seq$ and values set to $value$ |
| d.values() | A copy of $d$'s list of values |
| d.get(k[, x]) | a[k] if $k$ in $d$, else $x$ |
| d.setdefault(k[, x]) | a[k] if $k$ in $d$, else $x$ (also setting it) |
| d.pop(k[, x]) | d[k] if $k$ in $d$, else $x$ (and remove $k$) |
| d.popitem() | Remove and return an arbitrary (key, value) pair |

# Dictionary In Python
## Query Dictionary Values

- To query a value from a dictionary without the risk of invoking an exception, use **get(k,x)**, where k is the key of the element to extract, while x is the element that will be returned in case k is not found as a key of the dictionary.

```
1  >>> iupac = {'E': 'Glu', 'X': 'Xaa', 'C': 'Cys', 'A'
   : 'Ala'}
2  >>> iupac.get('A','No translation available')
3  'Ala'
4  >>> iupac.get('Z','No translation available')
5  'No translation available'
6
```

- If you omit x, and there is no k key in the dictionary, it returns **None**.

```
1  >>> iupac.get('Z')
2  None
```

## Erasing Elements

- To erase elements from a dictionary, use the **del** instruction:

```
1   >>> iupac = {'E': 'Glu', 'X': 'Xaa', 'C': 'Cys', 'A'
     : 'Ala'}
2   >>> del iupac['A']
3   >>> iupac
4   {'C': 'Cys', 'X': 'Xaa', 'E': 'Glu'}
5
```

# Sets

## Set

- It is an unordered collection of unique and immutable (which cannot be modified)items.
- Every element is unique and is defined by **{}**
- The most common uses of sets are **membership testing**, **duplicate removal**, and the application of **mathematical operations**:
    - **intersections**,
    - **unions**,
    - **differences**, and
    - **symmetrical differences**.

# Sets
## Set In Python

### Example

```
1    set1={11,22,33,22}
2    print(set1)
3    set2={11,22,33,22, 44, 55}
4    print(set2)
5
```

### Output

```
1    {33, 11, 22}
2    {33, 22, 55, 11, 44}
3
```

# Set In Python
## Creating a Set

- Sets are created with the instruction **set()**:

```
1    >>> first_set = {'CP0140.1','XJ8113.5','EF3616
     .3'}
2
```

- To create an empty set and then add the elements as needed:

```
1    >>> first_set = set()
2    >>> first_set.add('CP0140.1')
3    >>> first_set.add('XJ8113.5')
4    >>> first_set
5    {'CP0140.1','XJ8113.5'}
6
```

# Set In Python
## Creating a Set

- set by comprehension:

```
1  >>> {2 * x for x in [1,2,3]}
2  {2, 4, 6}
3
```

- A **set** does not accept repeated elements

```
1  >>> {2*x for x in [1,1,2,2,3,3]}
2  {2, 4, 6}
3  >>> uniques = {2,2,3,4,5,3} # Removes duplicate
    elements
4
```

# Set Operations

**Intersection**

- To get the common elements in two sets, use the operator **intersection()**:

```
1   >>> set1 = {1, 3, 5, 7}
2   >>> set2 = {4, 5, 2, 7, 8}
3   >>> common = set1.intersect(set2)
4   {5, 7}
5
```

- It is equivalent to &:

```
1   >>> common = set1 & set2
2   {5, 7}
3
```

# Set Operations

**Union**

- The union of two (or more) sets is the operator **union()**

- its abbreviated form is |

```
1  >>> set1.union(set2)
2  {1, 2, 3, 4, 5, 7, 8}
3
```

- It is equivalent to |:

```
1  >>> set1 | set2
2  {1, 2, 3, 4, 5, 7, 8}
3
```

# Set Operations

**Difference**

- A **difference** is the resulting set of elements that belongs to one set but not to the other and it achieved in python using the method name **difference()**.
- Its shorthand is -

```
1   >>> set1.difference(set2)
2   {1, 3}
3
```

- It is equivalent to -:

```
1   >>> set1 - set2
2    {1, 3}
3
```

# Set Operations

## Symmetric Difference

- A symmetric difference refers to those elements that are not a part of the intersection and its operator is **symmetric_difference()**

- and it is shortened as ^

```
1  >>> set1.symmetric_difference(set2)
2  {1, 2, 3, 4, 8}
3
```

- It is equivalent to ^:

```
1  >>> set1 ^ set2
2  {1, 3}
3
```

# Set Operations

Shared Operations with Other Data Types

**Maximum, Minimum, and Length**

- Sets share some properties with sequences, such as :
  - **max**,
  - **min**,
  - **len**,
  - **in**, etc.

- As we can expect, these properties work in the same way.

- Type ***help(set())*** in the console to see all methods associated with sets.

# End of The Chapter



Figure