



—

BC COMS 2710: Computational Text Analysis

—

Lecture 2 – Python Overview



- Tutorial 1.1
 - Should be submitted today (Tuesday 05/04)

- Tutorial 1.2
 - Should be submitted tomorrow (Wednesday 05/05)

- Tutorial 1.3
 - Should be submitted Friday (05/07)

- Homework 01:
 - Due Monday 05/10

- Reading Week 1
 - Due Sunday 05/09

Updated Rubric



Participation	5%
4 Homeworks	30%
Reading reflections	15%
Daily Tutorials	20%
Final Project	35%

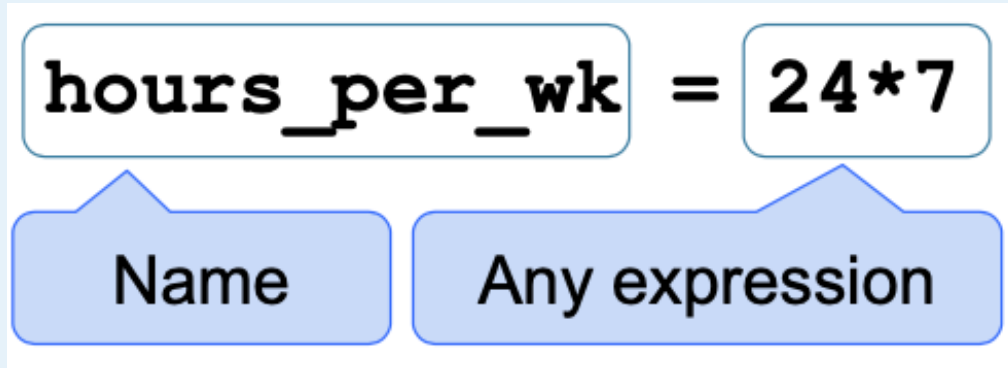


- Popular for data science & software development
- Mature data science and computational text analysis tools
- Learn through practice and doing
- Follow along in the demos



Names & Variables

Assignment Statements



- Statements perform an action
 - don't have a value
- Assignment statement changes the meaning of the name to the left of the = symbol
- The name is bound to a value (not an equation)



Numbers



Two real number types in Python

- **int:** an integer of any size
- **float:** a number with an optional fractional part

An int never has a decimal point; a float does

A float might be printed using scientific notation



- Floats have limited size (the limit is huge)
- Floats have limited precision of 15-16 decimal places
- After arithmetic, the final few decimal places can be wrong



— Strings —



A string value is a snippet of text of any length

- `'a'`
- `'word'`
- `"there can be 2 sentences. Here's the second!"`

Strings consisting of numbers can be converted to numbers

- `int('12')`, `float('1.2')`

Any value can be converted to a string

- `str(5)` becomes `"5"`



Assume you have run the following statements:

$x = 3$

$y = '4'$

$z = '5.6'$

What is the source of the error in each example?

- A. $x + y$
- B. $x + \text{int}(y + z)$
- C. $\text{str}(x) + \text{int}(y)$
- D. $y + \text{float}(z)$

Types – Every value has a type



We've seen 5 types so far:

- int: 2
- float: 2.2
- str: 'Red fish, blue fish'

Types – Every value has a type



The type function tells you the type of a value

- `type(2)`
- `type(2+2)`

An expression's "type" is based on its value

- `x = 2, y="hello"`
- `type(x), type(y) = ???`



Strings that contain numbers can be converted to numbers

- `int("12")`
- `float("1.2")`
- ~~`float("one point two")`~~ # Not a good idea



Any value can be converted to a string

- `str(6)`

Numbers can be converted to other numeric types

- `float(1)`
- `int(2.3)`. # DANGER: why is this a bad idea



Collections



Ordered:

- Lists
- Tuples

Unordered:

- Sets
- Dictionaries



- store multiple items in a single variable
 - `fruit = ["bananas", "apples", "oranges"]`
- Order is preserved
- Access items with brackets
 - `first_fruit = fruit[0]`
 - `second_fruit[= fruit[1]`
 - `last_fruit = fruit[-1]`
 - What are the values assigned to these three names?



```
fruit = ["bananas", "apples", "oranges"]
```

■ Access multiple items:

- `sub_fruit1 = [0:2]`
- `sub_fruit1 = ???`

- `sub_fruit2 = [:2]`
- `sub_fruit2 = ???`

- `sub_fruit3 = [0:]`
- `sub_fruit3 = ???`



```
fruit = ["bananas", "apples", "oranges"]
```

- Adding at the end:
 - `fruit.append("grapefruit")`

- Modifying at a specific location:
 - `fruit[1] = "strawberry"`
 - `fruit ???`



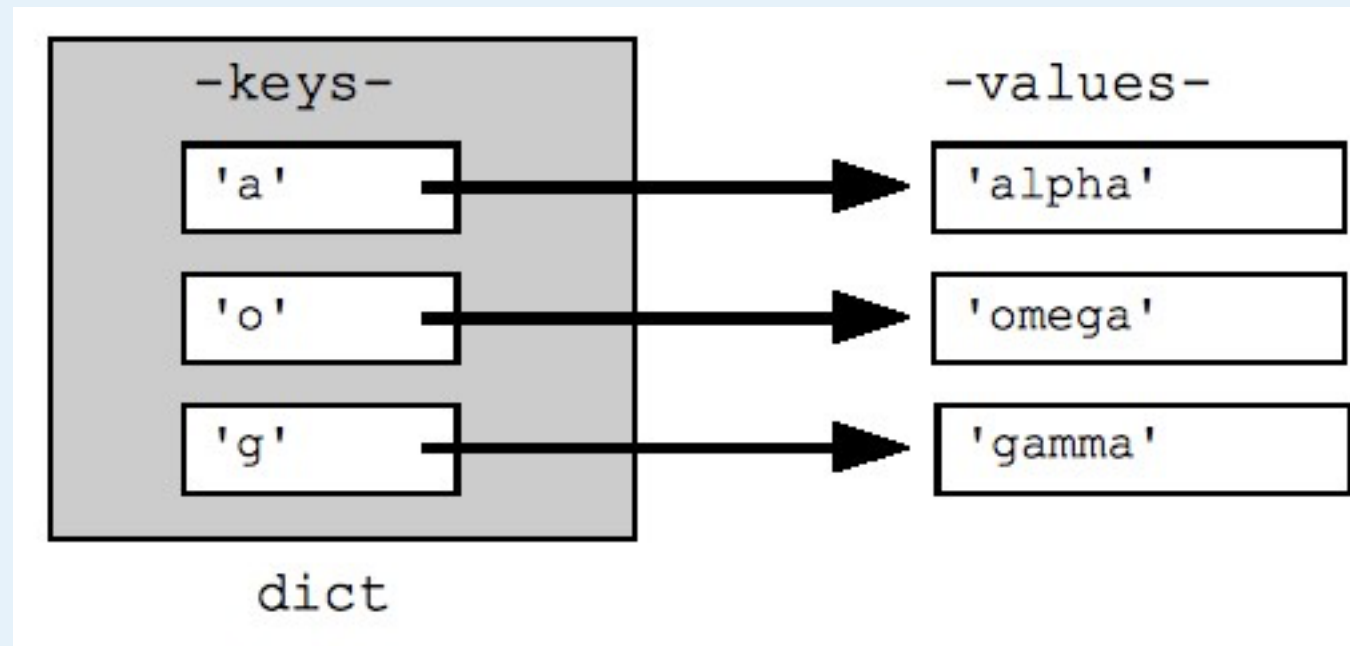
- Immutable lists
- `play = ("Shakespeare", "A Midsummer Night's Dream", 1595)`
- Used to group together related data



- Unordered and unindexed collection
- `authors = set(["Shakespeare", "Austin", "Morrison", "Woolf"])`
- No duplicates

- Store data values in *key:value* pairs.
- Ordered, changeable, no duplicates

- `{"a": 1,
"b": 2,
"c": 3}`





— Functions —

Anatomy of a Call Expression



What
function
to call

Argument to the
function

f (**27**)

"Call f on 27."

Anatomy of a Call Expression



What
function
to call

First argument

Second
argument

max (**15** , **27**)

Python Built-in Functions



		Built-in Functions		
<code>abs()</code>	<code>delattr()</code>	<code>hash()</code>	<code>memoryview()</code>	<code>set()</code>
<code>all()</code>	<code>dict()</code>	<code>help()</code>	<code>min()</code>	<code>setattr()</code>
<code>any()</code>	<code>dir()</code>	<code>hex()</code>	<code>next()</code>	<code>slice()</code>
<code>ascii()</code>	<code>divmod()</code>	<code>id()</code>	<code>object()</code>	<code>sorted()</code>
<code>bin()</code>	<code>enumerate()</code>	<code>input()</code>	<code>oct()</code>	<code>staticmethod()</code>
<code>bool()</code>	<code>eval()</code>	<code>int()</code>	<code>open()</code>	<code>str()</code>
<code>breakpoint()</code>	<code>exec()</code>	<code>isinstance()</code>	<code>ord()</code>	<code>sum()</code>
<code>bytearray()</code>	<code>filter()</code>	<code>issubclass()</code>	<code>pow()</code>	<code>super()</code>
<code>bytes()</code>	<code>float()</code>	<code>iter()</code>	<code>print()</code>	<code>tuple()</code>
<code>callable()</code>	<code>format()</code>	<code>len()</code>	<code>property()</code>	<code>type()</code>
<code>chr()</code>	<code>frozenset()</code>	<code>list()</code>	<code>range()</code>	<code>vars()</code>
<code>classmethod()</code>	<code>getattr()</code>	<code>locals()</code>	<code>repr()</code>	<code>zip()</code>
<code>compile()</code>	<code>globals()</code>	<code>map()</code>	<code>reversed()</code>	<code>__import__()</code>
<code>complex()</code>	<code>hasattr()</code>	<code>max()</code>	<code>round()</code>	



Control Statements



Say we have a list of author names, how can we find the length of each name?

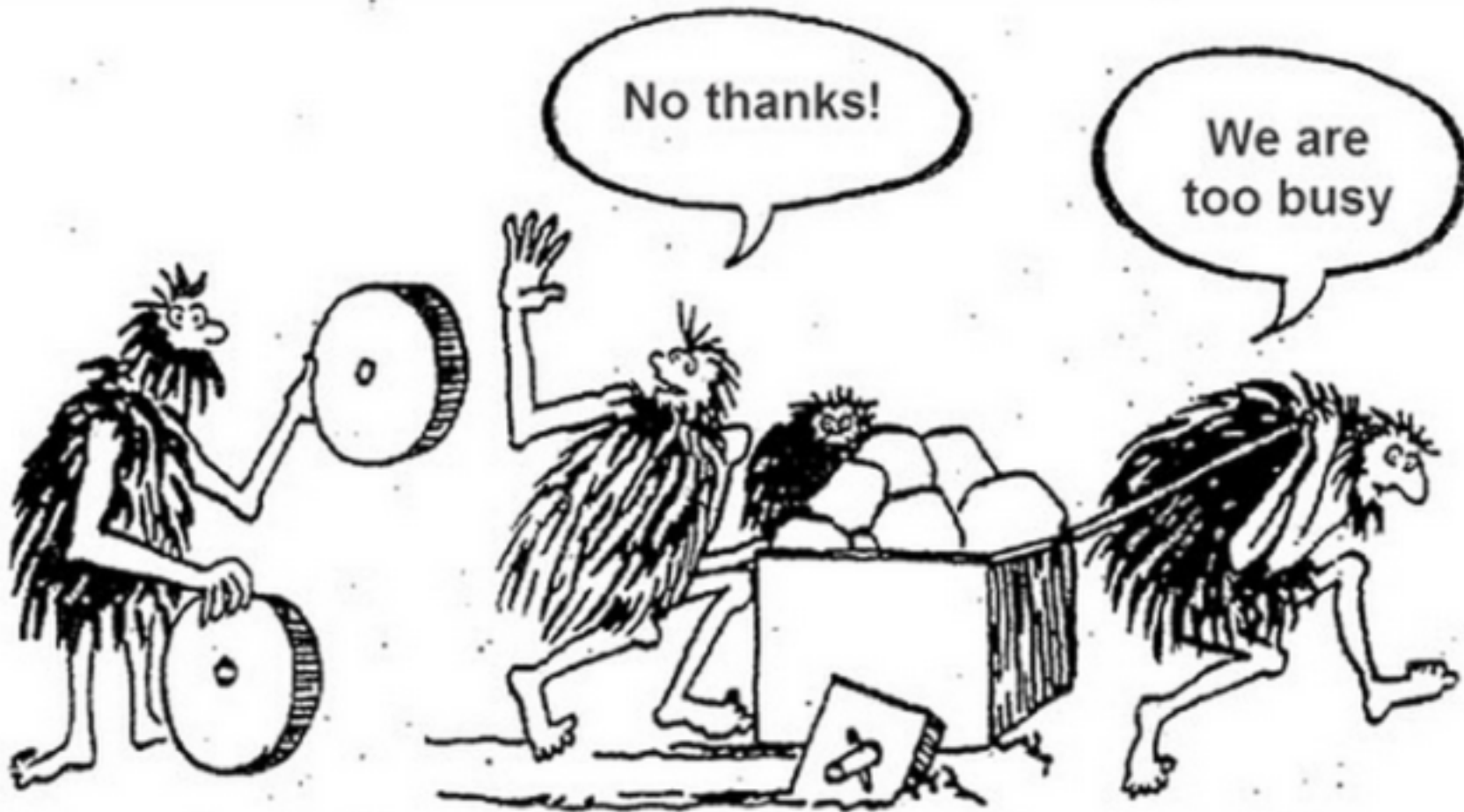


Only apply computation under certain scenario



Libraries

- Don't reinvent the wheel





- Install via command line:
 - `pip install <Library Name>`

- Access library in python:
 - `import <Library Name>`



- BeautifulSoup – webscraping
- Nltk – Processing text
- Spacy - Processing text
- little_mallet – Topic Modeling

- Pandas - Tables
- Matplotlib - Visualization
- Numpy – Vectors

- Sklearn – Machine Learning