

COMP 204: Sequence alignment examples, more dictionaries

Yue Li

based on material from Mathieu Blanchette, Carlos Oliver,
Christopher J.F. Cameron

Midterm materials coverage and practice midterms

- ▶ Midterm is held on February 22 at 6:30-8:00 pm in LEA 219.
- ▶ Our midterm will cover up to Lecture 17 (Feb 13)
- ▶ Past midterms in COMP 204 Fall 2018 and COMP 364 Fall 2017 are posted on myCourses for practice

A couple more Needleman-Wunsch examples on blackboard

The most important of Assignment 2 is to understand Needleman-Wunsch global sequence alignment algorithm. Let's do a couple of examples together:

Example 1

Sequence 1: G

Sequence 2: GCG

Example 2

Sequence 1: TCGA

Sequence 2: TTCG

A matrix in Python is just a list of lists of the same length

In Assignment #2, we will need to represent two-dimensional tables or matrices, with a fixed number of rows and columns.

Two-dimensional lists can be used to do this in Python.

A 2D list is a list of **lists**, where each of the **lists** is of the same length. Example: A tic-tac-toe grid:

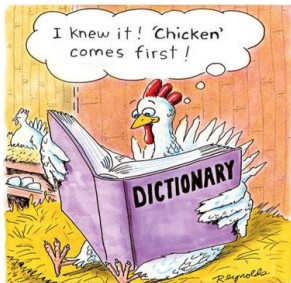
```
1 tictactoe = [ ["X", "", "O"],
2               ["", "X", ""],
3               ["O", "", ""] ]
4
5 print(tictactoe) # [['X', '', 'O'], ['', 'X', ''], ['O', '', '']]
6
7 # to access an element in a 2D list ,
8 # specify the index of the row and column
9 tictactoe[1][2] = "X"
10 print(tictactoe) # [['X', '', 'O'], ['', 'X', 'X'], ['O', '', '']]
```

To create a new matrix with zeros we can use *list comprehension*:

```
1 # Create an alignment scoring grid for two DNA sequences
2 seq1="GATTACA"
3 seq2="GCATGCA"
4 alignmentScoreGrid=[[0 for j in range(len(seq2)+1)]
5                      for i in range(len(seq1)+1)]
6 print(alignmentScoreGrid) # print out the scoring grid
```

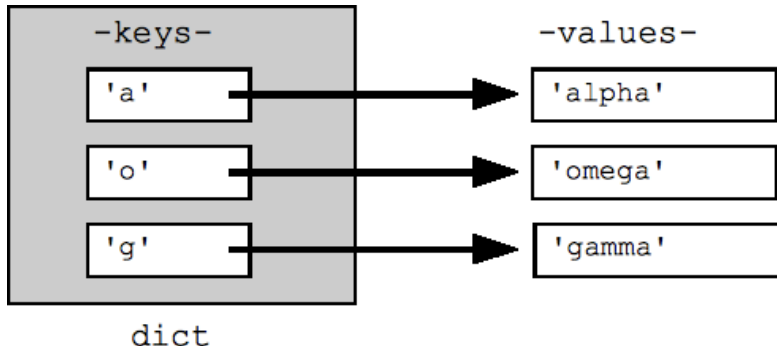
Dictionaries Recap

- ▶ A dictionary is said to be a **mapping** type because it maps *key* objects to *value* objects.
- ▶ Dictionaries are immensely useful and are the magic behind a lot of Python functionality
- ▶ **Syntax:** `my_dict = {[key]: [value], ...,}`
- ▶ The analogy to a real dictionary works. The word you look up is the **key** and the definition is the **value**



Dictionaries: picture

- ▶ Keys **map** to values.
- ▶ We use dictionaries when we want to access data using something other than an index (i.e. lists).



Dictionaries: keys and values

- ▶ A dictionary's keys can be many different types of **immutable** objects (i.e. int, str, tuple)
- ▶ You can access a key's value like a list. **Syntax:**
`my_dict[key]`
- ▶ You can mix and match key types
- ▶ Values can be any object type. You can also mix and match.

```
1 record_sales = {  
2     "Kanye": 2.4,  
3     "Beyonce": 1.5,  
4     "Chance": 1.2,  
5     ("a", 12): "bob"  
6 }  
7 print(record_sales["Beyonce"]) # 1.5  
8 print(record_sales[("a", 12)]) # "bob"
```

Adding keys to a dictionary

- ▶ **Syntax:** `my_dict["key"] = value`
- ▶ If the key does not yet exist, a new key/value pair is created.
- ▶ If the key already exists, its previous value is overwritten

```
1 >>> d = {"bob": 28}
2 >>> print(d)
3 {"bob": 1.2}
4 >>> d["charlie"] = 33
5 >>> print(d)
6 {"bob": 1.2, "charlie": 2.5}
7 >>> d["bob"] = "woooo"
8 {"bob": "woooo", "charlie": 33}
9 >>> del d["bob"] # we can delete keys with the del
  ↪ operator
10 {"charlie": 33}
```

Important properties of dictionaries

- ▶ Dictionaries are **mutable** We can modify the contents of the dictionary as much as we want.

```
1 >>> d = {"bob": 24, "tina": 11}
2 >>> d["tameeka"] = 42
3 >>> d['bob'] = [1, 2, 3, 4]
4 >>> del d["bob"]
5 >>> mystring = 'AAAGGG'
6 >>> mystring[2] = 'T' # this is an error. strings
   ↪ are immutable
```

Important properties of dictionaries

- ▶ Key-value pairs are **NOT** always stored in order. (for the current Python 3.7 they are, but assume it won't be like this forever)
- ▶ If you want to iterate over the keys in a dictionary use the `dict.keys()` function.

```
1 >>> d = {"bob": 24, "tina": 11}
2 >>> for k in d.keys():
3 >>> ... print(k)
4 "tina"
5 "bob"
```

Useful dictionary methods and operators

- ▶ `d.items()` produces an iterator which yields tuples of the form `(key, value)`

```
1 >>> for k,v in d.items():
2 >>> ... print("key:", k, "value:", v)
3 >>>
4     key: bob, value 24
5     key: tina, value 11
```

- ▶ `k in d` evaluates to `True` if the key exists in the dictionary and `False` otherwise.
- ▶ `d.update(d2)` “merges” two dictionaries into one.

```
1 >>> d = {"a": 3, "b": 4}
2 >>> d.update({"c": 5})
3 {"a": 3, "c": 5, "b": 4}
```

Quick dictionary example: mini BLAST

- ▶ BLAST is a very popular bioinformatics tool used to compare DNA sequences. One of the main innovations is to index a genome by 'words'.
- ▶ *words* are short sequences. AT, CG, CC, GG, AA
- ▶ **Goal:** Given a genome and a list of words return a dictionary with a list of positions where each given word occurs.
- ▶ **Example:** for words AAG, AAT in genome GAAGAAGGGAATGGAAGAAT we should return 'AAG': [1,4,14], 'AAT': [9,17].

Note BLAST is a *heuristic approach* to do fast sequence search but Needleman-Wunsch global alignment algorithm (or Smith-Waterman local alignment) is a more principled way to find optimal match(es) at the cost of speed.

Building genomic dictionary

```
1 #Args: genome_seq: a DNA sequence as a string
2 #      words: an iterable of sequences
3 #Returns:
4 # genomeDict: a dict with a key for each word mapping to
5 # list of indices.
6 def buildGenomeDict(genome_seq, words):
7     genomeDict = {}
8     for w in words:
9         for i in range(len(genome_seq)-len(w)+1):
10             if genome_seq[i:i+len(w)] == w:
11                 if w not in genomeDict:
12                     genomeDict[w] = []
13                 genomeDict[w].append(i)
14     return genomeDict
15 genome_seq = "AGCGACGTATAATCGACTA"
16 words=["CG", "TATA"]
17 genomeDict = buildGenomeDict(genome_seq, words)
18 print(genomeDict)
19 # {'CG': [2, 13], 'TATA': [7]}
```

Searching genomic dictionary

```
1 #Args: genomeDict: build from genome_index
2 #     genome_seq: DNA sequence corresponding to genomeDict
3 #     queries: a list of query sequences
4 #Returns: blasthits: a dict with a key for each query and
5 #         their genomic location(s)
6 def searchGenomeDict(genomeDict, genome_seq, queries):
7     blasthits={}
8     for q in queries:
9         blasthits[q]=[] # initialize query hit list
10        for genomeDictKey in genomeDict.keys():
11            for i in range(len(q)-len(genomeDictKey)+1):
12                wordlen = len(genomeDictKey)
13                querySubstr = q[i:i+wordlen]
14                if querySubstr == genomeDictKey:
15                    genomePosList = genomeDict[querySubstr]
16                    for pos in genomePosList:
17                        if genome_seq[pos-i:pos+len(q)-i] ==
18                            q: # mistake: genome_seq[pos:pos+len(q)] == q:
19                            blasthits[q].append(pos-i)
20                if len(blasthits[q])>0:
21                    blasthits[q] = set(blasthits[q]) # set returns
22                    unique values (more info on set in the next lecture)
23                return blasthits
24 queries = ["ACGT", "CGACGT", "TATAAT", "CGACT", "XYZ"]
25 myhits = searchGenomeDict(genomeDict, genome_seq, queries)
26 print(myhits) # {'ACGT': {4}, 'CGACGT': {2}, 'TATAAT': {7}}
```

A convenient method:.setdefault

- ▶ Let's look at line 15 in the previous example:

```
1 if w not in word_index:  
2     word_index[w] = []  
3 word_index[w].append(i)
```

- ▶ You will find yourself writing this statement many times.
- ▶ `mydict.setdefault(key, [default])` If key is in the dictionary, return its value. If not, insert key with a value of `default` and return `default`.
- ▶ We can replace it with one line using `setdeafult`

```
1 word_index.setdefault(w, []).append(i)
```

Dictionaries Pop Quiz

- ▶ True or False: dictionaries are immutable.
- ▶ Error? `myd = {[1,2]: "hello"}`
- ▶ True or False: dictionary keys must be unique.
- ▶ Error? `d2 = {'bob': 2, 'susan': 3, 'bob':4}`
- ▶ Error? `d = {}; d['bob'].append(3)`
- ▶ True/False: once a key-value is stored we can't update it.

Dictionaries Pop Quiz Answers

- ▶ True or **False**: dictionaries are immutable.
- ▶ Error? `myd = {[1,2]}: "hello"` – **Yes. Keys must be immutable.**
- ▶ **True** or False: dictionary keys must be unique.
- ▶ Error? `d2 = {'bob': 2, 'susan': 3, 'bob':4}` – **No. Duplicate keys are overwritten**
- ▶ Error? `d = {}; d['bob'].append(3)` – **Yes. Key 'bob' has not been created.**
- ▶ True/**False**: once a key-value is stored we can't update it.
`d['bob'] = 3; d['bob']='hi'` is valid.