# Algorithms and Data Structures for Data Science

# Lists and Random Data

CS 277

January 29, 2024

Brad Solomon

UNIVERSITY OF ILLINOIS URBANA-CHAMPAIGN

Department of Computer Science

## About ([https://cs.illinois.edu/research/undergraduate-research/srp](https://cs.illinois.edu/research/undergraduate-research/srp))

This hybrid program will take place from June to August (exact dates TBD). Students who participate in the summer program engage with a faculty and their research group as well as attend weekly 'Lunch and Learn' sessions which focus on various professional development topics.

This is not a paid program – it serves as a 'common app' for matching students to faculty and a support and social structure for undergraduate students over the summer. Actual hiring is still handled directly by faculty individually.

## Eligibility [Rolling acceptance – no fixed deadline]

Open to all undergraduates (Illinois CS and external), iCan students, and local high school students.

# CRA-WP Distributed Research Experiences for Undergraduates (DREU)

## About (https://cra.org/cra-wp/dreu/)

A highly selective program that matches students with a faculty mentor for a summer research experience at the faculty mentor's home institution. DREU interns will receive $700 per week for research (up to 10 weeks), and will be directly involved in a research project and interact with graduate students and professors on a daily basis. This experience is invaluable for those who are considering graduate school; DREU will provide a close-up view of what graduate school is really like and increase interns' competitiveness as an applicant for graduate admissions and fellowships.

## Eligibility [Student and mentor application deadline February 15th]

Students who are pursuing an undergraduate degree at an institution in the U.S. or its territories. International Students may apply, however, most of the funds for the DREU program are restricted to US citizens and permanent residents, so the number of non-US student participants will be limited. Priority will be given to persons from populations underrepresented in computing including women, Black/African American, Native American/Alaskan Native/ Pacific Islander, Hispanic/Latinx, LGBTQAI+, Persons with Disabilities, and Veterans.

# DaRin Butz Foundation Research Scholars

## About ([https://isur.engineering.illinois.edu/darin-butz-foundation-research-scholars/](https://isur.engineering.illinois.edu/darin-butz-foundation-research-scholars/))

DaRin Butz Foundation Research Scholars conduct research in the areas of computer science, aerospace, electrical, computer, materials science, nuclear engineering, physics, or astronomy. Scholars will work with faculty mentors who will supervise, guide, and instruct them on their research during the course of the project. They are expected to do research 30–35 hours per week for 10 weeks in summer. Scholars must take ENG 199 UGR in the fall and will present their work in the Fall Engineering Research Fair or the annual ISUR poster expo in the spring semester.

## Eligibility [Student application deadline March 31st]

Must be:
- A University of Illinois woman undergraduate in the Grainger College of Engineering
- U.S. citizen or permanent resident
- Rising sophomore, junior, or senior majoring in CS, aerospace, electrical, computer, materials science, nuclear engineering, physics, or astronomy with a GPA of 3.0 or higher

## About (https://www.siam.org/students-education/programs-initiatives/siam-simons-undergraduate-summer-research-program)

Society for Industrial and Applied Mathematics (SIAM) is pleased to announce we are accepting applications for the SIAM-Simons Undergraduate Summer Research Program, which will provide research, networking, and mentorship opportunities to U.S. students from underrepresented groups. Participating students will receive a stipend of $1,000/week and will have their housing, meals, and travel expenses paid. This is an amazing opportunity for students to immerse themselves in applied math, computational science, and/or data science research while simultaneously participating in professional development and community-building activities designed to foster a strong sense of belonging.

## Eligibility [Student application deadline February 7th]

Must be a U.S. citizen or permanent resident and an undergraduate student enrolled in a US-based college or university in September 2024. Note that while all projects will have an applied math and/or computational science approach, students do not need to have an applied math background to apply. Projects appropriate for students at all undergraduate levels will be available, and prior research experience is not required.

# Reminder: First two labs due today!

If you need an extension, email the request and reason why.

Labs are meant to be completable in the lab section (or shortly after)

Currently labs are due the Monday after they are released

# CS 277 Plagiarism Policy

Github Co-pilot and chatGPT can probably pass this class.

Using these tools is against the course policy and will hurt you on exams

It also hurts your peers by giving a false signal on assessments

The labs in CS 277 are designed to give you programming practice

The MPs in CS 277 are designed to be data science explorations

# Learning Objectives

Motivate the importance of the list data structure

Practice using built-in Python lists

Introduce file I/O in Python

Introduce random() package

# Lists are a great way to store **data**

```
>Sequence_1 assembly1
CCCTAAACCCTAAACCCTAAACCCTAAACCTCTGAATCCTTAATCCCTAAATCCCTAAAT
CTTTAAATCCTACATCCATGAATCCCTAAATACCTAATTCCCTAAACCCGAAACCGGTTT
CTCTGGTTGAAAATCATTGTGTATATAATGATAATTTTATCGTTTTTATGTAATTGCTTA
TTGTTGTGTGTAGATTTTTTAAAAATATCATTTGAGGTCAATACAAATCCTATTTCTTGT
GGTTTTCTTTCCTTCACTTAGCTATGGATGGTTTATCTTCATTTGTTATATTGGATACAA
GCTTTGCTACGATCTACATTTGGGAATGTGAGTCTCTTATTGTAACCTTAGGGTTGGTTT
ATCTCAAGAATCTTATTAATTGTTTGGACTGTTTATGTTTGGACATTTATTGTCATTCTT
>Sequence_2
CCCTAAACCCTAAACCCTAAACCCTAAACCTCTGAATCCTTAATCCCTAAATCCCTAAAT
CTTTAAATCCTACATCCATGAATCCCTAAATACCTAATTCCCTAAACCCGAAACCGGTTT
CTCTGGTTGAAAATCATTGTGTATATAATGATAATTTTATCGTTTTTATGTAATTGCTTA
TTGTTGTGTGTAGATTTTTTAAAAATATCATTTGAGGTCAATACAAATCCTATTTCTTGT
GGTTTTCTTTCCTTCACTTAGCTATGGATGGTTTATCTTCATTTGTTATATTGGATACAA
GCTTTGCTACGATCTACATTTGGGAATGTGAGTCTCTTATTGTAACCTTAGGGTTGGTTT
ATCTCAAGAATCTTATTAATTGTTTGGACTGTTTATGTTTGGACATTTATTGTCATTCTT
```

When Mr. Bilbo Baggins of Bag End announced that he would shortly be celebrating his eleventy-first birthday with a party of special magnificence, there was much talk and excitement in Hobbiton.
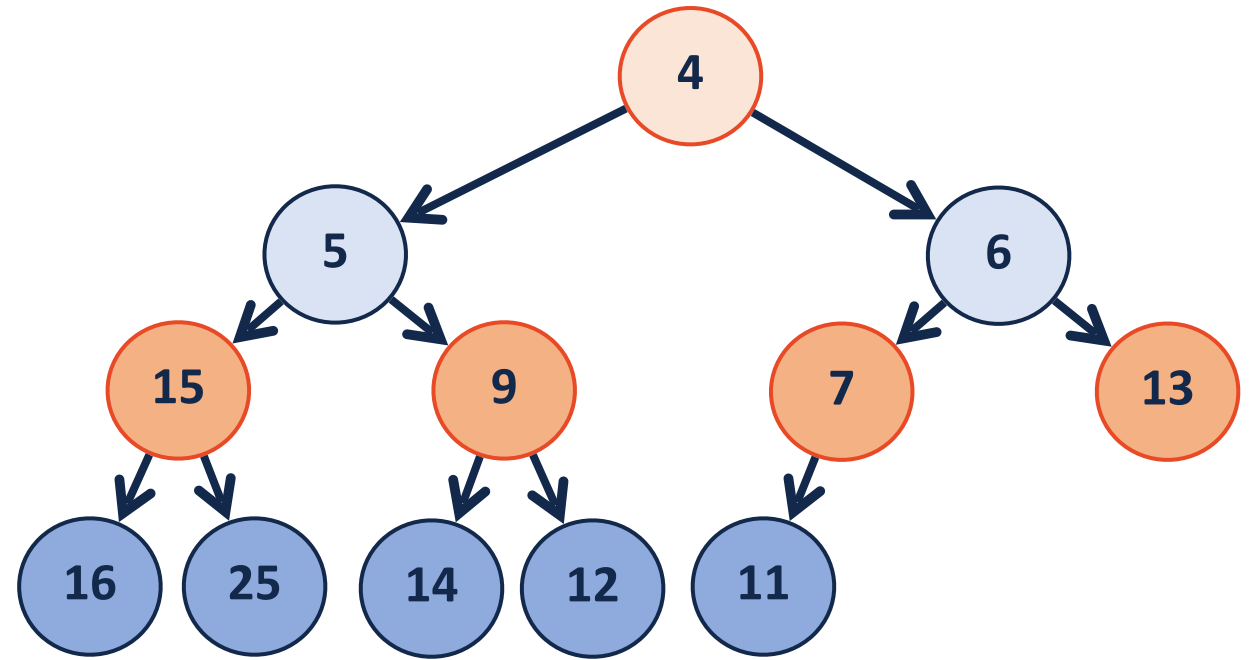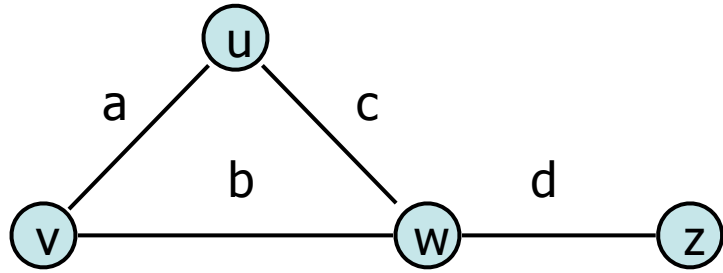
Bilbo was very rich and very peculiar, and had been the wonder of the Shire for sixty years, ever since his remarkable disappearance and unexpected return. The riches he had brought back from his travels had now become a local legend, and it was popularly believed, whatever the old folk might say, that the Hill at Bag End was full of tunnels stuffed with treasure. And if that was not enough for fame, there was also his prolonged vigour to marvel at. Time wore on, but it seemed to have little effect on Mr. Baggins. At ninety he was much the same as at fifty. At ninety-nine they began to call him *well-preserved*; but *unchanged* would have been nearer the mark. There were some that shook their heads and thought this was too much of a good thing; it seemed unfair that anyone should possess (apparently) perpetual youth as well as (reputedly) inexhaustible wealth.

'It will have to be paid for,' they said. 'It isn't natural, and trouble will come of it!'
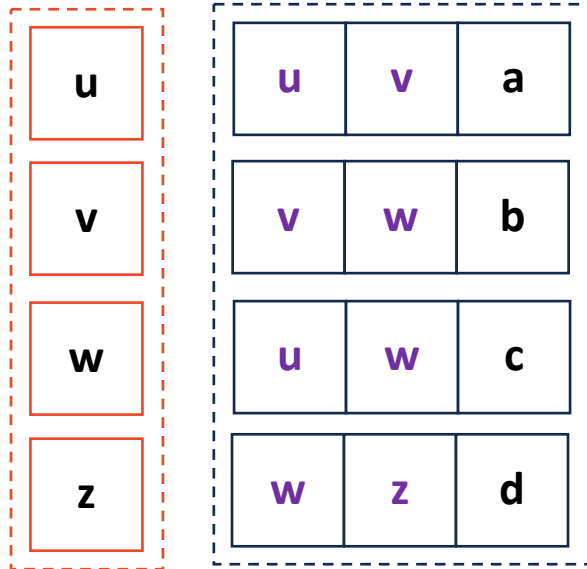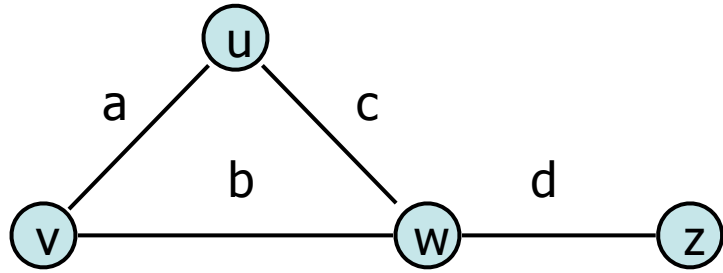
### Chapter I

### A LONG-EXPECTED PARTY

```
"searched_zipcode","searched_lat","searched_lng","searched_address","searched_state","searched_city","searched_metro","is_gh","
latitude","longitude","distance","loc_name","loc_number","url","address","cuisines","delivery_fee_raw","delivery_fee","delivery
_time_raw","delivery_time","service_fee","phone","review_count","review_rating","RunDate","restaurant_tags","delivery_type"
"11216","40.6788319","-73.9506774","300 Madison St, Brooklyn, NY 11216, USA","NY","Brooklyn","NY","$0 with
GH+","40.67882919","-73.94986725","0.04","Chun Vegetarian","1013747","https://www.grubhub.com/chun-vegetarian-582-nostrand-ave-
brooklyn/1013747","582 Nostrand Ave","[""Asian"", ""Vegan"", ""Vegetarian""]","$0
delivery",0.00,"35-45",40.0,0.00,"3476278080",1691,4.57,"2022-04-25 07:01:11","[""PRICING_ELIGIBLE_FOR_STANDARD_ORDER"",
""SUBSCRIPTION_ELIGIBLE_FOR_STANDARD_ORDER""]","SELF"
"11216","40.6788319","-73.9506774","300 Madison St, Brooklyn, NY 11216, USA","NY","Brooklyn","NY","$0 with
GH+","40.67871475","-73.94987488","0.04","India House","1436488","https://www.grubhub.com/india-house-586-nostand-ave-brooklyn/
1436488","586 Nostand Ave","[""Healthy"", ""Indian"", ""Lunch Specials"", ""Vegetarian""]","$0
delivery",0.00,"35-55",45.0,0.00,"7188577011",697,4.77,"2022-04-25 07:01:11","[""PRICING_ELIGIBLE_FOR_STANDARD_ORDER"",
""SUBSCRIPTION_ELIGIBLE_FOR_STANDARD_ORDER""]","SELF"
```

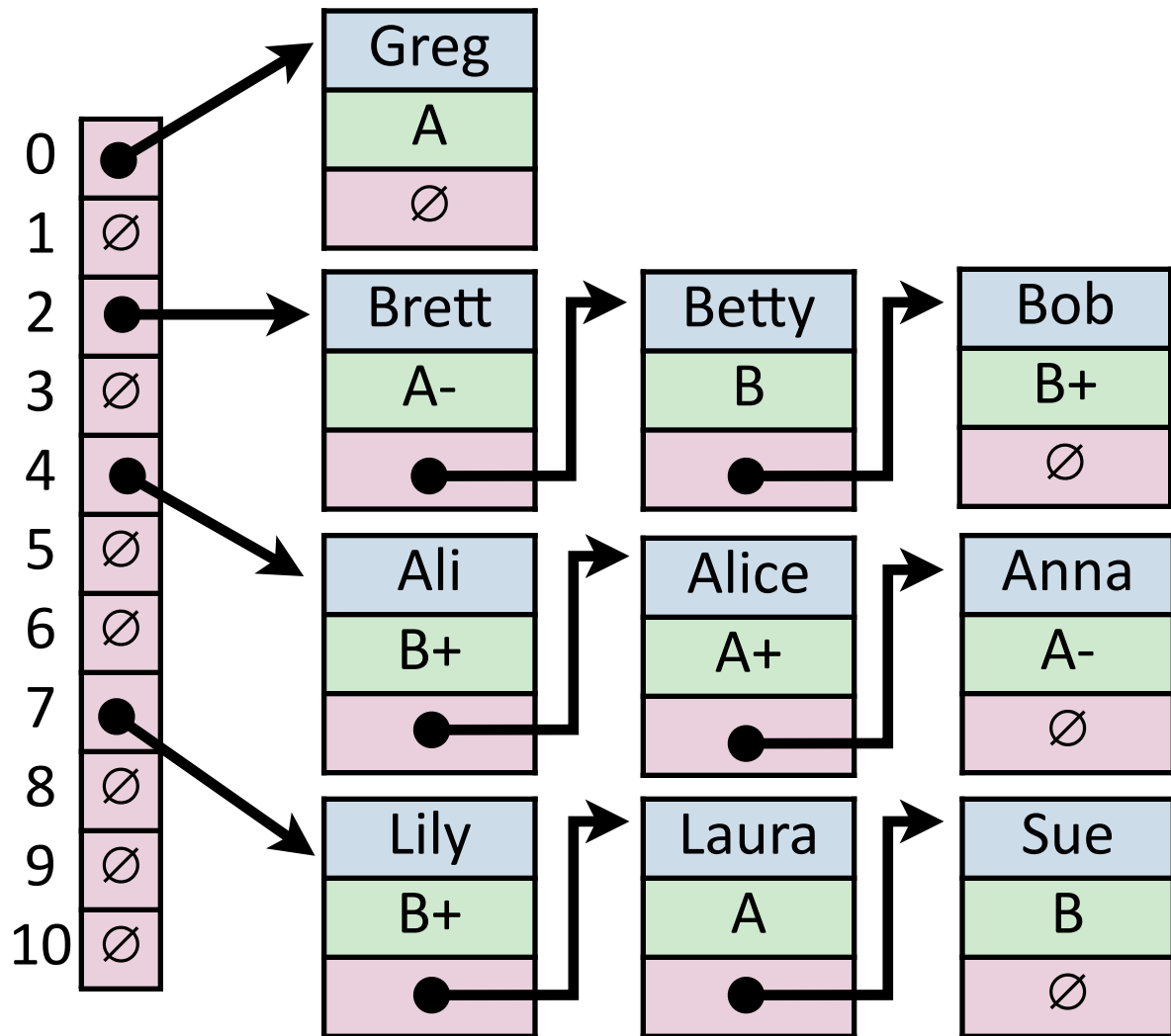# Lists are a great way to store **data structures**

# Lists are a great way to store **data structures**

# Lists are a great way to store **data structures**

$$H = \{h_1, h_2, \ldots, h_k\}$$

# List Abstract Data Type

A list is an **ordered** collection of items

Items can be either **heterogeneous** or **homogenous**

The list can be of a **fixed size** or is **resizable**

Can be **localized** or **distributed** in memory

# List Abstract Data Type

A minimally functional list must have the following functions:

**Constructor:** Create a new empty list

**Insert:** Add an item to a list

**Delete:** Remove an item from the list

**Index:** Look up the value at a specific index

**Size()**:** Get the size of the list

# Lists in Python: Constructor

There are many ways to construct (or initialize) a Python list.

```
 1  l1 = [1, 2, 3]
 2
 3
 4  l2 = list( (3, 4, 5) )
 5
 6
 7  l3 = l1+l2
 8
 9
10  l4 = l1.copy()
11
12
13  l5 = l1
14
15
16
```

# Lists in Python: Insert

**Insert()** tries to add the object at a particular index in the list

**Append()** adds the item to the end of the list

```
 1  l1 = []
 2
 3
 4  l1.insert(2, "A")
 5
 6
 7  l1.append("C")
 8
 9
10  l1.insert(1, "B")
11
12
13  l1.insert(5, "D")
14
15
16
```

# Lists in Python: Remove

**Remove()** removes the first instance of the object

**Remove()** crashes if the object doesn't exist!

**Pop()** removes and returns the object at a specific index (default is last)

```
1  l1 = [1,2,1,3,1,4]
2
3
4  l1.remove(7)
5
6
7  l1.remove(1)
8
9
10 x = l1.pop()
11
12
13 y = l1.pop(2)
14
15
16
```

# Lists in Python: Index

**Index()** returns the index of the first matching item in the list

**__getItem__(), []** returns the index of the first matching item in the list

```
1  l1 = [1,2,1,3,1,4]
2
3  print(l1.index(1))
4
5
6  print(l1.index(5))
7
8
9  print(l1.index(4))
10
11
12  print(l1[0])
13
14
15  print(l1[10])
16
17
18  print(l1[-1])
```

# Lists in Python: Size

**len()** is not a class method but works for most forms of a list in Python

**Remember:** Python indexing starts at 0

```
1  l1 = [1,2,1,3,1,4]
2
3  print(len(l1))
4
5
6
7
8
```

# List Abstract Data Type

A minimally functional list must have the following functions:

**Constructor:** `__init__()`

**Insert:** `append(x)`      `insert(i, x)`

**Delete:** `remove(x)`      `pop()`

**Index**     `__getitem__() []`      `index(x)`

**Size()**\*\*    `len(list)`

# Programming Practice: Lists

What is the output of the print statement after this function is run?

```
 1  def inclass():
 2      l = [1,2,3,4,5,6]
 3      x = l.pop()
 4      l.pop()
 5      l.pop()
 6      l.insert(0, x)
 7      l.append(8)
 8      l += [1, 2]
 9      return l
10
11  print(inclass())
12
13
14
15
16
17
18
```

# Programming Practice: Index of Every X

Given a list and a character as input, return a list containing the index of every instance of that character.

# Programming Practice: Index of Every X

Given a list and a character as input, return a list containing the index of every instance of that character.

Now lets remove every third value from the output list.

# Python List

There are many implementations of lists in Python. Here are three*:

```
1 myList = [1, 2, 3, 4, 5]
2
3 print(myList)
4
5 print(len(myList))
6
7 print(myList[2])
8
9
```

```
1 myTuple = (1, 2, 3, 4,
2 5)
3
4 print(myTuple)
5
6 print(len(myTuple))
7
8 print(myTuple[2])
9
```

```
1 import numpy as np
2 myNP =
3 np.array([1,2,3,4,5])
4
5 print(myNP)
6
7 print(len(myNP))
8
9 print(myNP[2])
```

# Why are there so many different implementations?

Lets find out together in our first class mini-project!

**P1)** Generate random datasets

**P2)** Code various analysis functions on datasets using lists

**P3)** Measure efficiency of **P2** using varying size datasets from **P1**

# Mini-Project 0: Random Data Generation

**Learning Objectives:**

Practice string and list manipulations in the context of file I/O

Introduce seeded random data generation

Introduce how to measure runtime efficiency in Python

Investigate the efficiency of different Python implementations of lists

# Programming Toolbox: File I/O

```
 1  readableFile = open('inputFile.txt', 'r')
 2
 3
 4
 5
 6  writableFile = open('outputFile.txt', 'w')
 7
 8
 9
10
11  carefulWriteFile = open('outputFile.txt', 'x')
12
13
14
15
16  appendableFile = open('outputFile.txt', 'a')
17
18
19
20
21
22
23
```

# Programming Toolbox: File I/O

**Which approach do you prefer?**

```python
1   # Approach 1
2
3   readableFile = open('inputFile.txt', 'r')
4
5
6   fileData = readableFile.read()
7
8
9   readableFile.close()
10
11
12  # Approach 2
13
14  with open('inputFile.txt', 'r') as myFile:
15      fileData = myFile.read()
16
17
18
```

# Programming Toolbox: File I/O

```python
1  with open('data/temp1.txt', 'x') as myFile:
2      for i in range(10):
3          myFile.write(str(i))
4      myFile.write("\n")
5      myFile.write("Line 2")
6
7
8
9  myFile = open('data/temp2.txt', 'w')
10 for i in range(5):
11     myFile.write(str(i) + "\n")
12 myFile.close()
13
14
15
16 with open('data/temp2.txt', 'a') as myFile:
17     myFile.write("Hello World!\n")
18
19 with open('data/temp2.txt', 'a') as myFile:
20     myFile.write("Hello World!\n")
21
22
23
```

# Programming Toolbox: File I/O

```
1  with open('data/temp1.txt', 'x') as myFile:
2      for i in range(10):
3          myFile.write(str(i))
4      myFile.write("\n")
5      myFile.write("Line 2")
6
7
8
9  myFile = open('data/temp2.txt', 'w')
10 for i in range(5):
11     myFile.write(str(i) + "\n")
12 myFile.close()
13
14
15
16 with open('data/temp2.txt', 'a') as myFile:
17     myFile.write("Hello World!\n")
18
19 with open('data/temp2.txt', 'a') as myFile:
20     myFile.write("Hello World!\n")
21
22
23
```

```
1  0123456789
2  Line 2
```

temp1.txt

```
1  0
2  1
3  2
4  3
5  4
6
```

temp2.txt

# Programming Toolbox: File I/O

```
1  with open('data/temp1.txt', 'x') as myFile:
2      for i in range(10):
3          myFile.write(str(i))
4      myFile.write("\n")
5      myFile.write("Line 2")
6
7
8
9  myFile = open('data/temp2.txt', 'w')
10 for i in range(5):
11     myFile.write(str(i) + "\n")
12 myFile.close()
13
14
15
16 with open('data/temp2.txt', 'a') as myFile:
17     myFile.write("Hello World!\n")
18
19 with open('data/temp2.txt', 'a') as myFile:
20     myFile.write("Hello World!\n")
21
22
23
```

```
1 0123456789
2 Line 2
```
temp1.txt

```
1 0
2 1
3 2
4 3
5 4
6 Hello World!
7
```
temp2.txt

# Programming Toolbox: File I/O

```
1  with open('data/temp1.txt', 'x') as myFile:
2      for i in range(10):
3          myFile.write(str(i))
4      myFile.write("\n")
5      myFile.write("Line 2")
6
7
8
9  myFile = open('data/temp2.txt', 'w')
10 for i in range(5):
11     myFile.write(str(i) + "\n")
12 myFile.close()
13
14
15
16 with open('data/temp2.txt', 'a') as myFile:
17     myFile.write("Hello World!\n")
18
19 with open('data/temp2.txt', 'a') as myFile:
20     myFile.write("Hello World!\n")
21
22
23
```

```
1 0123456789
2 Line 2
```

temp1.txt

```
1 0
2 1
3 2
4 3
5 4
6 Hello World!
7 Hello World!
8
```

temp2.txt

# Programming Toolbox: File I/O

```
1   with open('data/temp1.txt') as myFile:
2       inList = myFile.readlines()
3   print(inList)
4
5
6
7
8
9   myFile = open('data/temp2.txt')
10  for i in range(10):
11      print("Line Content: {}".format(myFile.readline()))
12  myFile.close()
13
14
15
16
17
18  with open('data/temp1.txt') as myFile:
19      print(myFile.read())
20
21
22
23
```

```
1  0123456789
2  Line 2
3
4
5
6
```

temp1.txt

```
1  0
2  1
3  2
4  3
5  4
6
```

temp2.txt

# Programming Toolbox: File I/O

```
1  with open('data/temp1.txt') as myFile:
2      inList = myFile.readlines()
3  print(inList)
4
5
6
7
8
9  myFile = open('data/temp2.txt')
10 for i in range(10):
11     print("Line Content: {}".format(myFile.readline()))
12 myFile.close()
13
14
15
16
17
18 with open('data/temp1.txt') as myFile:
19     print(myFile.read())
20
21
22
23
```

temp1.txt

| 1 | 0123456789 |
|---|---|
| 2 | Line 2 |
| 3 | |
| 4 | |
| 5 | |
| 6 | |

temp2.txt

| 1 | 0 |
|---|---|
| 2 | 1 |
| 3 | 2 |
| 4 | 3 |
| 5 | 4 |
| 6 | |

['0123456789\n', 'Line 2']
Line Content: 0

Line Content: 1

Line Content: 2

Line Content: 3

Line Content: 4

Line Content:
Line Content:
Line Content:
Line Content:
Line Content:
0123456789
Line 2

# Programming Toolbox: File I/O

`str.strip()` will remove whitespace from the string.

| 1 | 0 |
|---|---|
| 2 | 1 |
| 3 | 2 |
| 4 | 3 |
| 5 | 4 |
| 6 | |

temp2.txt

```
1  myFile = open('data/temp2.txt')
2  for i in range(6):
3      print("Line Content: {}".format(myFile.readline().strip()))
4  myFile.close()
5
6  with open('data/temp2.txt') as myFile:
7      for line in myFile:
8          print(line.strip())
```

Line Content: 0
Line Content: 1
Line Content: 2
Line Content: 3
Line Content: 4
Line Content:
0
1
2
3
4

# Programming Toolbox: File I/O

`str.strip()` will remove whitespace from the string.

```
 1  tmp = "1, 2, 3"
 2  tmp2 = "1,2,3"
 3
 4  x = tmp.split(",")
 5  y = tmp2.split(",")
 6  for i in range(len(x)):
 7      if x[i]!=y[i]:
 8          print("No match!", x[i], y[i])
 9      else:
10          print("Match!")
11
12
```

# Programming Toolbox: File I/O

**Four ways to open a file:** read, write, carefulWrite, append

**Three ways to read a file:** all lines, line-by-line, text as string

**Writing to file:** String formatting is key, file writes are literal copies

# Python File I/O

**Whats wrong here**:

```
1
2  with open('dataFile2.txt','w') as myFile:
3      data = myFile.readlines()
4
5  for line in data:
6      print(line)
7
8
9
10
11
12
13
14
15
16
17
18
```

# Programming Practice: File I/O

Given an input filename, write a new file that is the same file with reversed lines.

# Mini-Project 0: Random Data Generation

**Learning Objectives:**

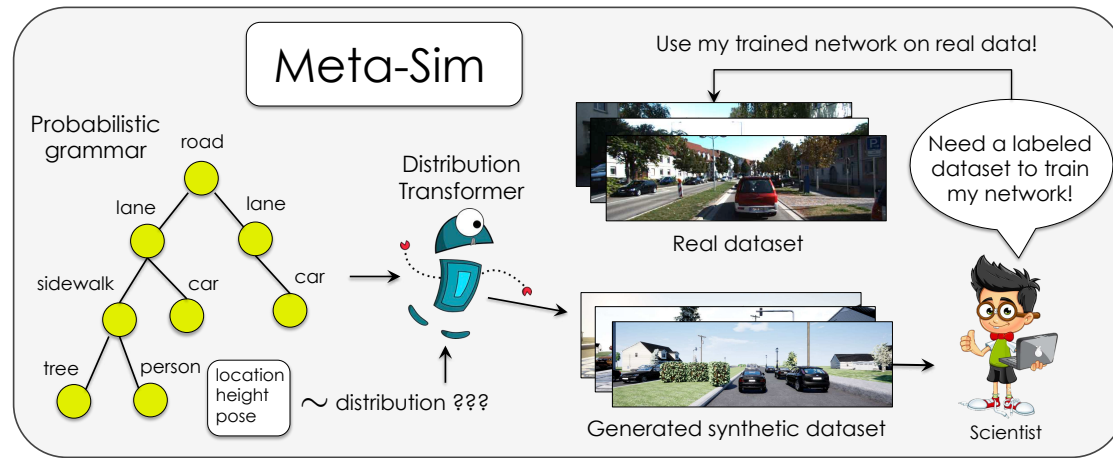Practice string and list manipulations in the context of file I/O

Introduce **seeded random data generation**

Introduce how to measure runtime efficiency in Python

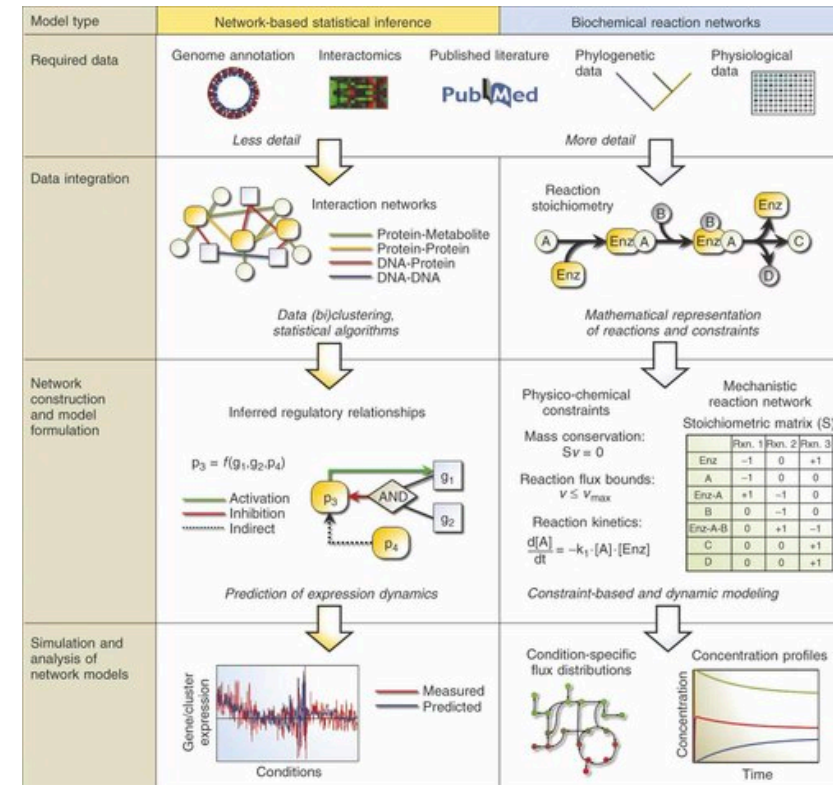Investigate the efficiency of different Python implementations of lists

# Programming Toolbox: Synthetic Datasets

**Performance Analysis of K-Means and K-Medoids Clustering Algorithms For A Randomly Generated Data Set.** T Velmurugan et al 2008



**Meta-Sim: Learning to Generate Synthetic Datasets.** A Kar et al. 2019



**In silicon models of cancer.** Edelman et al. 2010

"In many applications of computers and other electronic devices, there is a need for a physical source of true random numbers, for example, **in computer simulations of various probabilistic algorithms,** computer games, and, most notably, in possibly randomized cryptographic algorithms and protocols whose security relies on the ability to generate unpredictable secret keys and random numbers."

**New Methods for Digital Generation and Postprocessing of Random Data.** J.D.J Golic 2006

# Programming Toolbox: Random.Random()

Python's random module implements **pseudo-random** number generators

What does pseudo-random mean? Well try running the following:

```
Random.random()
```

```
Random.seed(1)
```

```
Random.random()
```

# Programming Toolbox: Random.Seed()

Python's random module implements **pseudo-random** number generators

The values produced are not actually random!

Given a start value, the output is entirely deterministic!

```
Random.random()
```

```
Random.seed(1)
```

```
Random.random()            0.134642441240122
```

# Programming Toolbox: Random.Seed()

Python's random module implements **pseudo-random** number generators

The values produced are not actually random!

Given a start value, the output is entirely deterministic!

This is actually quite useful (for us):

1. Total reproducibility

2. Autogradable randomness

3. It's built-in

# Programming Toolbox: Random.RandInt()

While most uniform random data generation can be done using random(), some other functions are useful for specific tasks.

```
Random.randint(min, max)
```

# Programming Toolbox: Random.Choice()

While most uniform random data generation can be done using random(), some other functions are useful for specific tasks.

```
Random.choice(list)
```

# Programming Toolbox: Random.Sample()

While most uniform random data generation can be done using random(), some other functions are useful for specific tasks.

```
Random.sample(list, num)
```

# Programming Toolbox: Random.shuffle()

While most uniform random data generation can be done using random(), some other functions are useful for specific tasks.

```
Random.shuffle(list)
```

# Programming Toolbox: Lots more to random

https://docs.python.org/3/library/random.html

# Programming Practice: Uniform Random

Lets visualize uniform randomness using lists!

Write a function that takes in two numbers, a length and count, and generates count random variables across length. Instead of storing the sequence, we will only store the counts assigned to each value.