Lecturer: Yu LI

Scriber: LEE, Michelle (1155179462)

# 1. Different data types

### a. Sequential data

- Displays the sequence of nucleic acid or amino acids (DNA, RNA, and proteins).
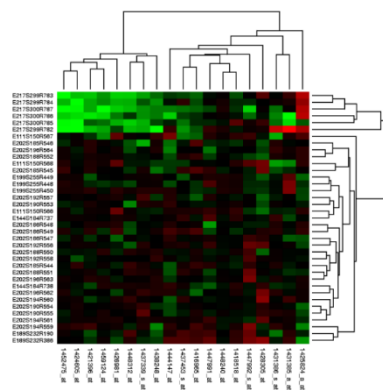- Example: plasmid sequence (pSB1C3 from https://parts.igem.org/Part:pSB1C3)

```
tactagtagcggccgctgcagtccggcaaaaaagggcaaggtgtcaccaccctgcccttttttctttaaaaccgaaaagattacttcgcgttatgcaggcttcctcgctcactgactcgctgcgctcggtc
gttcggctgcggcgagcggtatcagctcactcaaaggcggtaatacggttatccacagaatcaggggataacgcaggaaagaacatgtgagcaaaaggccagcaaaaggccaggaaccgtaaaaaggccg
cgttgctggcgtttttccacaggctccgcccccctgacgagcatcacaaaaatcgacgctcaagtcagaggtggcgaaacccgacaggactataaagataccaggcgtttccccctggaagctccctcgt
gcgctctcctgttccgaccctgccgcttaccggatacctgtccgcctttctcccttcgggaagcgtggcgctttctcatagctcacgctgtaggtatctcagttcggtgtaggtcgttcgctccaagctg
ggctgtgtgcacgaaccccccgttcagcccgaccgctgcgccttatccggtaactatcgtcttgagtccaacccggtaagacacgacttatcgccactggcagcagccactggtaacaggattagcagag
cgaggtatgtaggcggtgctacagagttcttgaagtggtggcctaactacggctacactagaagaacagtatttggtatctgcgctctgctgaagccagttaccttcggaaaaagagttggtagctcttg
atccggcaaacaaaccaccgctggtagcggtggtttttttgtttgcaagcagcagattacgcgcagaaaaaaaggatctcaagaagatcctttgatcttttctacggggtctgacgctcagtggaacgaa
aactcacgttaagggattttggtcatgagattatcaaaaaggatcttcacctagatcctttaaattaaaaatgaagttttaaatcaatctaaagtatatatgagtaaacttggtctgacagctcgaggc
ttggattctcaccaataaaaaacgcccggcggcaaccgagcgttctgaacaaatccagatggagttctgaggtcattactggatctatcaacaggagtccaagcgagctcgatatcaaattacgccccgc
cctgccactcatcgcagtactgttgtaattcattaagcattctgccgacatggaagccatcacaaacggcatgatgaacctgaatcgccagcggcatcagcaccttgtcgccttgcgtataatatttgcc
catggtgaaaacgggggcgaagaagttgtccatattggccacgtttaaatcaaaactggtgaaactcacccagggattggctgagacgaaaaacatattctcaataaaccctttaggggaaataggccagg
ttttcaccgtaacacgccacatcttgcgaatatatgtgtagaaactgccggaaatcgtcgtggtattcactccagcgatgaaaacgtttcagtttgctcatggaaaacggtgtaacaagggtgaacac
tatcccatatcaccagctcaccgtctttcattgccatacgaaattccggatgagcattcatcaggcgggcaagaatgtgaataaaggccggataaaacttgtgcttatttttctttacggtcttaaaaa
ggccgtaatatccagctgaacggtctggttataggtacattgagcaactgactgaaatgcctcaaaatgttctttacgatgccattgggatatatcaacggtggtatatccagtgatttttttctccatt
ttagcttccttagctcctgaaaatctcgataactcaaaaaatacgcccggtagtgatcttatttcattatgtgaaagttggaacctcttacgtgcccgatcaactcgagtgccacctgacgtctaagaa
accattattatcatgacattaacctataaaaataggcgtatcacgaggcagaatttcagataaaaaaaatccttagctttcgctaaggatgatttctggaattcgcggccgcttctagag
```

### b. Data matrix

- Collection of records; each consists of a fixed set of attributes.
- Can be represented by row(n) x columns(m) matrix, each row for each object and each column for each attribute.
- Swapping the entire column or the entire row at one time will not change the data.
- Commonly found in spreadsheets forms.
- Examples:

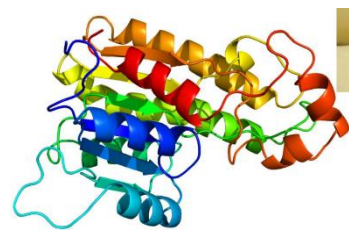| | | Attributes | |
|---|---|---|---|
| | Person | Height (m) | Weight (kg) |
| Objects | P1 | 1.79 | 75 |
| | P2 | 1.64 | 54 |
| | P3 | 1.70 | 63 |
| | P4 | 1.88 | 78 |

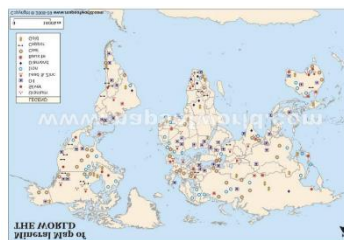(4x2) matrix with 4 people and 2 attributes
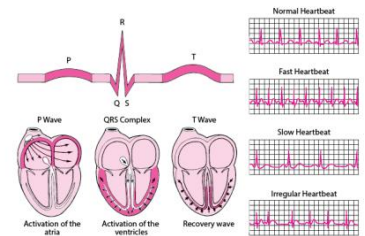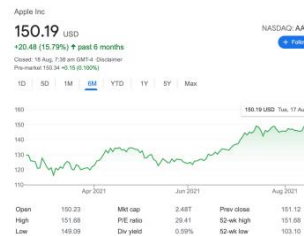


Gene expression in cells

### c. Spatial data

- Displays geographic locations or spatial information.
- Changing the rows or columns will change the data.
- Examples:
  - World mineral map
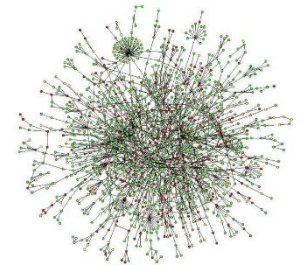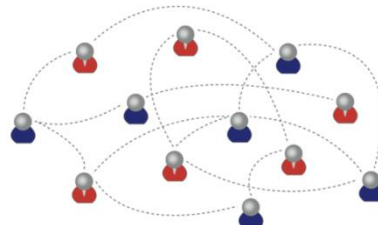  - Coordinates of atoms in a protein structure (3D data)

d. Temporal data
- Build-in support for data that involves time.
- Change of data over time.
- Examples:
  o Stock market graphs
  o ECG signals

e. Graph or networks
- Displays object and its connections
- Examples:
  o Social network
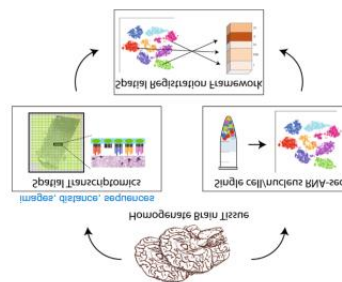  o Protein-protein interaction (PPI) networks

f. Text
- Data in text format; usually for describing something.
- Such as short and long sentences, documents.
- Examples:
  o Wikipedia articles
  o PubMed articles
  o Tweets in twitter

g. Multi-modality data
- Combinations between multiple data types.
- Examples:
  o Video; a combination of temporal images, audio, and transcript
  o Electronic health records; a combination of data matrix, images, and text
  o Spatial transcriptomics; a combination of spatial data and data matrix

h. Unknown data type
- The data has not been shown; only when data is shown, we will know the data type.
- Examples: diet and exercise.

# 2. Introduction to Python programming
- Programming is a **way of communicating with a computer** so that it can do something for us, and this is achieved through codes that we feed into the computer.
- Python is the **software** to send codes to the computer as well as **one of the languages** used to communicate with the computer.
- Plug-ins are available in Python to make Python more powerful; need to be loaded before its usage in the codes.

- o NumPy
- o SciPy
- o Pandas
- Code examples:
  - o Calculating the mean of some values:

| Codes | Returns |
|---|---|
| `import numpy`<br>`numpy.mean([1,2,3])` | 2.0 |

  - o Storing values in array and calculate the mean, variance, median, max, mean, etc.:

| Codes | Returns |
|---|---|
| `import numpy`<br>`a = [1,2,3,4,5,6,7,8,9]`<br>`numpy.mean(a)` | 4.916666666666667 |
| `numpy.std(a)` | 2.253084305765962 |
| `numpy.median(a)` | 5.0 |
| `numpy.max(a)` | 9 |
| `print(a)` | [1, 2, 3, 4, 5, 6, 7, 8, 9] |

  - o Printing strings and variables

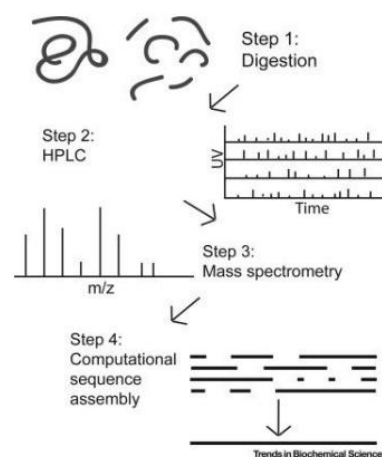| Codes | Returns |
|---|---|
| `print("The a array is ", a)` | The a array is [1, 2, 3, 4, 5, 6, 7, 8, 9] |
| Storing values in variables:<br>`import numpy`<br>`a = [1,2,3,4,4,5,5,5,6,7,8,9]`<br>`a_mean = numpy.mean(a)`<br>`a_std = numpy.std(a)`<br>`a_med = numpy.median(a)`<br>`a_max = numpy.max(a)` | |
| Printing the stored values:<br>`print("The a array is ", a, "Its mean is ", a_mean)` | The a array is [1, 2, 3, 4, 4, 5, 5, 5, 6, 7, 8, 9] Its mean is 4.916666666666667 |
| Printing without storing the value into a variable:<br>`print("The a array is ", a, "Its mean is ", numpy.mean(a))` | The a array is [1, 2, 3, 4, 4, 5, 5, 5, 6, 7, 8, 9] Its mean is 4.916666666666667 |

## 3. Sequence data

- DNA, RNA, and protein are part of the central dogma of biology.
- Genetic information is hidden in DNA sequences.
- Phenotype results from the combination of genotype (believed to be the sequences by biologists) and environmental factors.
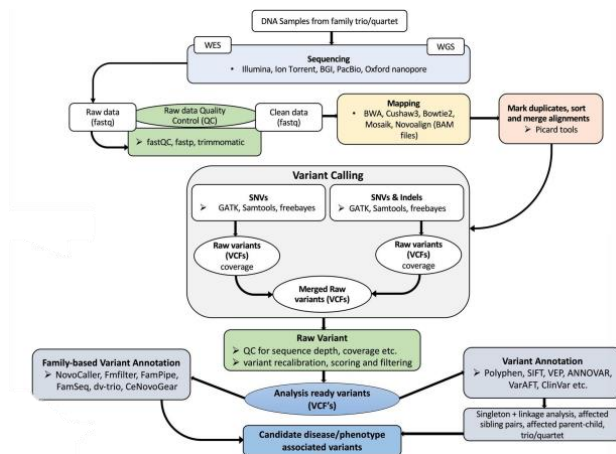
| DNA sequence | o Composed of A, T, C, G.<br>o Normally exist in double-stranded form.<br>o Approximately 3 billion base pairs in the human genome. |
|---|---|
| RNA sequence | o Composed of A, U, C, G.<br>o Single-stranded. |
| Protein sequence | o Usually composed of 20 amino acids. |

| | o   Multiple sequence alignment: technique to compare and analyze multiple protein sequences to study its evolution and function. |
|---|---|

- DNA/RNA sequencing can be used to obtain short to long reads of its sequences; the technology is still under active development. Milestones in this field are the following:
  - o  Sanger method (1977)
  - o  Human mitochondrial genome sequence (1981)
  - o  Human genome project (1990)
  - o  Complete cell genome (1995)
  - o  Complete eukaryotic genome (1996)
  - o  Complete the human genome project (2001)
  - o  Second generation sequencer (2005–2007)
  - o  Research human microbiome project (2008)
  - o  Third generation sequencer (2011)
  - o  Nanospace sequencing (2014)
  - o  Third stage human microbiome project (2019)
- Nanopore sequencing
  - o  Exploits different electrical current change generated by different bases of DNA when it goes through a chemical pore.
  - o  Using MinION, an enzyme that unwinds DNA and feeds one strand of DNA through the protein pore.
  - o  Able to sequence long sequences up to 3Mb, compared to older generation technologies that can only sequence up to 1000bp.
  - o  The error rate is relatively high compared to the older technology (5% and 0.001% respectively).
  - o  Still under active development.
- Protein sequencing
  - o  Through the process of digestion into shorter pieces, HPLC, mass spectrometry to determine the weight, and computational sequence assembly.



Trends in Biochemical Sciences

- Raw data, such as DNA sequences, will go through quality control, followed by mapping the reads to reference genome, variant calling, and finally processed to phenotype associated variants.

- Protein sequences are to be compared, for example, through multiple sequence alignment. Similar sequences correspond to similar structure and function. Similar sequences also indicate common ancestor.

# 4. Sequence comparison and alignment score

a. Sequence alignment and similarity
- Sequence alignment is to determine the similarity between sequences and identify the regions of similarity.
- Pairwise sequence alignment: arranging two sequences to maximize their similarity; inserting gaps is allowed.

b. Sequence alignment score
- To define sequence similarity by quantity depending on match (e.g., A with A), mismatch (e.g., G with T), or gap (e.g., C with gap).
- Using scoring matrix to calculate the score:

|   | A | C | G | T |
|---|---|---|---|---|
| A | 2 | −7 | −5 | −7 |
| C | −7 | 2 | −7 | −5 |
| G | −5 | −7 | 2 | −7 |
| T | −7 | −5 | −7 | 2 |

Gap penalty $= -10$

- Examples, by enumeration:

| Alignment | Score |
|---|---|
| AGGCCG<br>ATGC_G | 2 – 7 + 2 + 2 – 10 + 2<br>$= -9$ |
| AGGCCG<br>ATGCG_ | 2 – 7 + 2 + 2 – 7 – 10<br>$= -18$ |

- To find the best pairwise alignment, the straightforward solution will be by enumeration, calculating scores for all possible alignments and selecting the one with the highest score.
  - Problem: too many possible alignments.
  - No of possible alignments $= \binom{2n}{n} = \frac{(2n)!}{(n!)^2}$
  - For n = 300, there are 7 x $10^{88}$ possible alignments.
  - Solution: dynamic programming

*Disclaimer: all figures are adapted from Prof. Li BMEG3105 Lecture Notes*