# **Sequence Analysis:**

Md. Izhar Ashraf Computational Epigraphy Lab(iCEL), The Institute of Mathematical Sciences(IMSc), Chennai ashraf@imsc.res.in Bits & Scripts II March 13-24, 2025





#### Sequence Analysis in Biology: Examples





Figure: Source: PLoS genetics 2007

Figure: Cell 2019



#### Sequence Analysis in Economics: Examples



Figure: Source: Sitabhra et. al Physica A 2018



#### Figure: Source: Tradingview.com



#### Sequence Analysis in War, Business and Politics



Figure: Source: Wikipedia



#### Figure: Source: Unknown



#### Sequence Analysis in Archaeology



Arabi Chines Dutcl Egyptian (Hi English Greek (Ancient Hausa (Boko) Hebrey Japanese (Kana) Korea LinearB Malay (Burni) Persiar Spanish Sumerian (Cuneiform Tami Turkist Urdu deciphered (Indus) -0.6 -0.2 -0.1 ΔG -0.5 -0.4 0.2

Figure: Source: Wikipedia

#### Figure: Source: In-house



#### Sequence Analysis in Games



#### Figure: Source: nytimes.com



- Sequence analysis is the process of examining and comparing sequences of symbols.
- These sequences can represent:
  - 🔶 Biological data (DNA, RNA, proteins)
  - ✤ Textual data (words, characters)
  - $\diamond$  Any data that can be represented as a linear sequence.
- The goal is to identify patterns, similarities, and differences that reveal underlying relationships.

# Applications:

- $\diamond$  Bioinformatics (evolutionary studies, gene identification)
- $\diamond$  Natural Language Processing (text similarity, error correction)
- $\diamond$  Historical text analysis (Cipher breaking, text comparison)

## Global Alignment: Needleman-Wunsch Objective: Optimal end-to-end sequence alignment. Formula:

Initialization: 
$$F(i,0) = -i \cdot d$$
,  $F(0,j) = -j \cdot d$   
Recurrence:  $F(i,j) = \max \begin{cases} F(i-1,j-1) + s(x_i,y_j) \\ F(i-1,j) - d \\ F(i,j-1) - d \end{cases}$ 

**Example**: Align X = GAT, Y = CAT ( $s_{match} = +1$ ,  $s_{mismatch} = -1$ , d = 1):

Final Alignment:

gene sequences.

Local Alignment: Smith-Waterman **Objective:** Identify high-similarity subregions. **Formula**:

$$S(i,j) = \max egin{cases} 0 \ S(i-1,j-1) + s(a_i,b_j) \ S(i-1,j) - d \ S(i,j-1) - d \end{cases}$$

 Example: Align X = GATT,  $Y = ATT (s_{match} = +2, d = 2)$ :

 -A T T 

 -0 0 0 0 

 G 0 0 0 0 

 G 0 0 2 0 

 T 0 0 2 4 

 T 0 0 0 6 

 T 0 0 0 6 

 Score = 6

names in inscriptions.



Purpose: Minimal edits (insertions, deletions, substitutions).
 Formula:

$$\operatorname{lev}(a,b) = \begin{cases} |a| & \text{if } |b| = 0\\ |b| & \text{if } |a| = 0\\ \operatorname{lev}(\operatorname{tail}(a), \operatorname{tail}(b)) & \text{if } a[0] = b[0]\\ 1 + \min \begin{cases} \operatorname{lev}(\operatorname{tail}(a), b)\\ \operatorname{lev}(a, \operatorname{tail}(b)) & \text{otherwise}\\ \operatorname{lev}(\operatorname{tail}(a), \operatorname{tail}(b)) \end{cases}$$

**Application**: Correcting OCR errors, fuzzy searching.



#### Levenshtein Distance Example

**Example:** Transform "Klng" to "King". Substitution = 1, Insertion/Deletion = 1.

1. Final Matrix:

Γ0	1	2	3	4]	
1	0	1	2	3	
2	1	1	2	3	
3	2	2	1	2	
4	3	3	2	1	
k			~		

Alignment:

K!ng King

Total edits = 1 (substitute "!" "i").



- Measures similarity between two vectors by their orientation or angle in a multidimensional space
- $\square$  Range: [-1,1] (1 = identical, 0 = unrelated, -1 =opposite)
  - Applications: Text/document comparison, recommendation systems

Similarity = 
$$\frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^{n} A_i B_i}{\sqrt{\sum_{i=1}^{n} A_i^2} \sqrt{\sum_{i=1}^{n} B_i^2}}$$

Key Features:Scale-invariant (ignores vector magnitude)

#### Jaro-Winkler Distance: Overview

**Purpose**: Measure string similarity with a bonus for common prefixes. Ideal for short strings (e.g., names, epigraphic terms).

#### **Components**:

- Jaro Similarity: Base score based on matching characters and transpositions.
- Winkler Adjustment: Boosts similarity for shared prefixes.

$$J(s,t) = \begin{cases} 0 & \text{if } m = 0, \\ \frac{1}{3} \left( \frac{m}{|s|} + \frac{m}{|t|} + \frac{m-t}{m} \right) & \text{otherwise.} \end{cases}$$
$$JW(s,t) = J(s,t) + l \cdot p \cdot (1 - J(s,t))$$

#### Variables:

*m*: Matching characters (within window <sup>max(|s|,|t|)</sup>/<sub>2</sub> − 1)

 *t*: Transpositions (half the mismatched order pairs)

 *t*: Length of common prefix (max 4)

 *p*: Scaling factor (default: 0.1)

## Jaro-Winkler Example: MARTHA vs MARHTA

#### Matching Characters:

- 1. Transpositions:
  - ♦ Swapped T/H: 1 transposition ⇒ t = 0.5
- 2. Jaro Calculation:

Matching positions:

М	Α	R	Т	Н	Α
М	Α	R	Н	Т	Α

$$J = \frac{1}{3} \left( \frac{6}{6} + \frac{6}{6} + \frac{6 - 0.5}{6} \right) \approx 0.972^{\blacktriangleright}$$
 Matches: 6/6 characters  
Transpositions: 1 pair (T/H)

3. Winkler Boost (prefix=3):

 $JW = 0.972 + 3.0.1 \cdot (1 - 0.972) \approx 0.980$ 

**Applications**: OCR error correction, deduplication of epigraphic records, name matching in fragmented texts.

# And Hand Sciences

#### Hamming Distance

**Definition:** Hamming Distance measures the number of positions at which two equal-length strings differ.

$$H(s,t) = \sum_{i=1}^{n} \mathbf{1}(s_i \neq t_i)$$

where:

*s<sub>i</sub>* and *t<sub>i</sub>* are symbols at position *i* in sequences *s* and *t*.

 *t<sub>i</sub>* (*s<sub>i</sub>* ≠ *t<sub>i</sub>*) equals 1 if symbols differ, 0 otherwise.

**Example:** Compare binary strings:

$$s = 110101, \quad t = 100111$$

Differences occur at positions: 2,4,5 Hamming Distance: H(s,t) = 3Applications:

Error Detection and Correction (e.g., ECC, QR Codes)

# Caesar Cipher

**Definition:** A substitution cipher that shifts letters by a fixed number k.

#### **Encryption & Decryption:**

$$E(x) = (x + k) \mod 26$$
,  $D(x) = (x - k) \mod 26$ 

#### **Example: Encrypt "HELLO" with** k = 3

- Convert letters to numbers: H = 7, E = 4, L = 11, O = 14.
- Apply E(x):
  E(H) = 10(K), E(E) = 7(H), E(L) = 14(O), E(O) = 17(R)
- Encrypted message: "KHOOR"

#### **Applications:**

- Educational cryptography.
- Historical encryption (e.g., Julius Caesar).

#### Limitations:

• Easily broken using brute force or frequency analysis.



#### Substitution Cipher

**Definition:** A cryptographic technique that replaces each letter of plaintext with another letter based on a fixed mapping.

#### Encryption:

Establish a one-to-one mapping between plaintext and ciphertext alphabets.

- ✤ Example Mapping: A Q, B W, ..., Z M.
- ✤ Plaintext: SIMPLE Ciphertext: ZXQVWU.

#### Decryption:

 $\diamond$  Reverse the mapping to recover plaintext from ciphertext.

#### **Applications:**

Military & Diplomatic Communication (Historical Use)

✤ Teaching Cryptography & Digital Security Concepts

#### Limitations:

- $\diamond$  Vulnerable to Frequency Analysis
- $\diamond$  Security Depends on Key Secrecy

# **Thank You**