

DNA - Double Helix of Nucleotide Strands.

Nucleotide - Half of one "rung" of DNA

Four choices for the nucleic acid of a nucleotide:

1 - Adenine (A)

2 - Cytosine (C)

3 - Guanine (G) - bonds to C

4 - Thymine (T) - bonds to A.

If we know one strand of DNA, we get the other strands for free because of this "complementarity."

- Genome sequencing refers to "Reading" the "genome"
- Genome: The nucleotide sequence read down one side of an organism's chromosomal DNA. A human genome has about 3 billion letters, or nucleotides.
- (An Amoeba) has one of the longest genomes: 670 billion nucleotides.

- For simplicity a Nucleotide is represented by a single character which refers to Nitrogenous base. (A, T, G, or C).

(*) Nucleic Acids (RNA and DNA)

Nucleotide

- * Nitrogenous base
- * Sugar molecule.
- * Phosphate group

are long chains (polymers) made up of many nucleotides linked together. These chains encode genetic information, which is used in cellular processes like protein synthesis, replication, and gene expression.

- Genome: Complete set of genetic information in an organism.

DNA

ATAGCTATTGCCCAATCCCCATACTCTATTGCA

Gene: A DNA sequence with directions for one protein.



An RNA's copy of the gene DNA

- Protein : Molecular Machines of the cell.
- Origin of Replication:
Replication begins in a region called replication origin (denoted Ori)
 - proceeds in both direction.
- Vibro Cholera - bacterium that causes cholera (n 500 Nucleotides).
 - The string of 500 nucleotides must contain a hidden message telling the cell to start replication.
- Hidden Message Problem:
input: String text (representing Ori)
Output: hidden message in text.

Replication Initiation is mediated by a protein called DnaA.

DnaA binds to a short segment in Ori known as a DnaA box, hidden message saying : "bind here"

Multiple DnaA boxes \rightarrow higher chance of binding \rightarrow higher "fitness"

- A K-mer pattern is a most frequent k-mer in a string if no other k-mer is more frequent than this pattern.

K-mer: Substring of length k with in a biological sequence.

2-mers: AA, AT, TA, TT . . .

How many 2-mers are there?

$$4 \times 4 = 16 \text{ 2-mers can exist.}$$

How many K-mers are there?

- x - x - ... K times

$(4)^k$ k -mers are there.

String or pattern : Frequency.

- return one with highest frequency.

- k represent length of substring pattern or k -mer.
 k length pattern.

T - A (mapping)

G - C (mapping)

AGTCGATAGT
TCACTATCA
Reverse

1) First map right
to left.

2) Then reverse from

ACT AT GCG ACT

right to left.

Q: Given Ori, how do we find "the hidden message" (DnaA box) saying that replication should start here?

- Find the most frequent K-mers that appear in this region (Consider their reverse complements as well). (Draft boxes are usually around 9 nucleotides long).

*) Different Genomes \Rightarrow Different hidden messages.

G G T G G T A G G

C C T A C C A C C

G G T G G T A G G

Most Frequent 9-mer

C C T A C C A C C

} are

| | | | | | | | |

} candidate

G G A T G G T G G

} hidden messages.

Check its reverse

complement as well.

G G A T G G T G G

C C T A C C A C C

C C A C C A T C C

CTT GATC AT ←
 ← GAA C TAG TA
 ATG ATC AAG ↘
Reverse Complements
of each other.

A K-mer forms an (L, t) -Clump inside Genome
 if there is short (length L) interval of Genome in
 which it appears many (at least $\geq t$) times.

- Clump Finding Problem:

input: A string Genome and integers K (length of
 a pattern), L (window length) and t (number of
 patterns in a Clump).

Output: All K-mers forming (L, t) -Clumps in
 genome.

```
def construct_lps(pattern):
    lps = [0]
    if len(pattern) == 1: return lps
    i, j = 1, 0 #
    count = 0
    while i <= len(pattern) - 1:
        if pattern[i] == pattern[j]:
            count += 1
            lps.append(count)
            j += 1
        else:
            lps.append(0)
            i += 1
```

$\text{pattern}[i] = \text{Count}$

$j += 1$

else: # not same.

$\text{Count} = 0$

if $j == 0$:

$\text{pattern}[i] = \text{Count}$

else:

$\text{pattern}[i] = \text{Count}$.

$j = 0$ # start from beginning.

Continue.

$i += 1$

j

i

ababa

Count

return LPS

\downarrow

\downarrow

a	b	a	b	a	a	b
0	0	1	2	3	1	2

$\text{Count} += 1$

- Move both pointers forward

} if same.

if not same

Set Count = 0

and set that

index element to count.

if $j == 0$: # already left most then do nothing

Otherwise move it to the left by 1.

a b a b a a b
0 0 1 2 3 1 2

```
def kmp_search(text, pattern):  
    n = len(text)  
    m = len(pattern)
```

lps = construct_lps(pattern)

res = []

i, j = 0, 0

while i < n

if text[i] == pattern[j]:

i += 1

j += 1

if j == m:

res.append(i - j)

j = lps[j - 1]

else:

if j != 0:

$$j = \text{lps}[j - 1]$$

else:

$$i + 1$$

return res.

a b c d e f g h i j
i .
0 + 4
3 1 + 4 1 : 5

BA2B

def Kmp_Search(text, pattern): # How many times
return count. string appeared in
the text.

def frequent_words(text, k):
freq_table = {} words = []
for i in range(0, len(text)):
K_mer = text[i:i+k] 0: 4
1: 5

freq = Kmp_Search(text, K_mer)
freq_table[K_mer] = freq

maximum = max(freq_table.values())

for key in freq_table:

if $\text{freq_table}[\text{key}] == \text{maximum}$:
 words.append(key)

return words.

A A A A C C C G T

T T T T G G C C A

2

A C C G G G T T T T

G T C A

C A G T

B A I C

T G A C

def reverse_complement(pattern):

mapping = { 'G': 'C', 'C': 'G', 'A': 'T',
 'T': 'A' }

new_string = ""

for i in range(0, len(pattern)):

new_string += mapping[pattern[i]]

```
new_string = new_string[::-1] # reverse  
return new_string.
```

```
def clump_finding (Genome, K, L, t):  
    k_mers = [] set  
    k-mer      interval  
    at least this  
    many times  
    appeared in  
    interval.
```

- Start with first window of size L, create all kmers in this window and store frequency in a dictionary.
- Simultaneously for that particular window check if any kmer appeared at least t times and append it to list of kmers.
- return kmers.

$A \overset{1+1}{C} C C$

```
def skew (Genome):  
    arr = [0] * (len(Genome) + 1)  
    J_S, C_S = 0, 0, indices = []  
    for i in range(0, len(Genome)):  
        if Genome[i] == 'A' or genome[i]  
            = = 'T':  
                arr[i+1] = arr[i-1]  
            else:  
                if Genome[i] == 'G':  
                    J_S += 1  
                elif Genome[i] == 'C':  
                    C_S += 1
```

$$arr[i+1] = J_S - C_S$$

$$\text{minimum} = \min(arr)$$

```
indices = [i for i in range(0, len(arr))  
          if arr[i] == minimum]
```

return indices.

$[0]$

$\downarrow \downarrow \uparrow$
 $C A T G C <$

$[0, -1]$

$(0, -1, -1)$

$0, -1 -1 -1$
 $0, -1 -1 -1 0 \}$

$\{0, -1, -\epsilon, -t, 0, \gamma, -2\}$

```
def Clump_finding( genome, k, l, t):
    k_mers = set()
```

```
for i in range( len(genome) ):  
    interval = genome[ i : i + L ]  
    K_mer_count = defaultdict( int )
```

```

for j in range(len(interval)):
    k_mer = window[j:j+k]
    k_mer_count[k_mer] += 1

```

for K_{-ver} , cont in K_{-ver_cont} .
itens

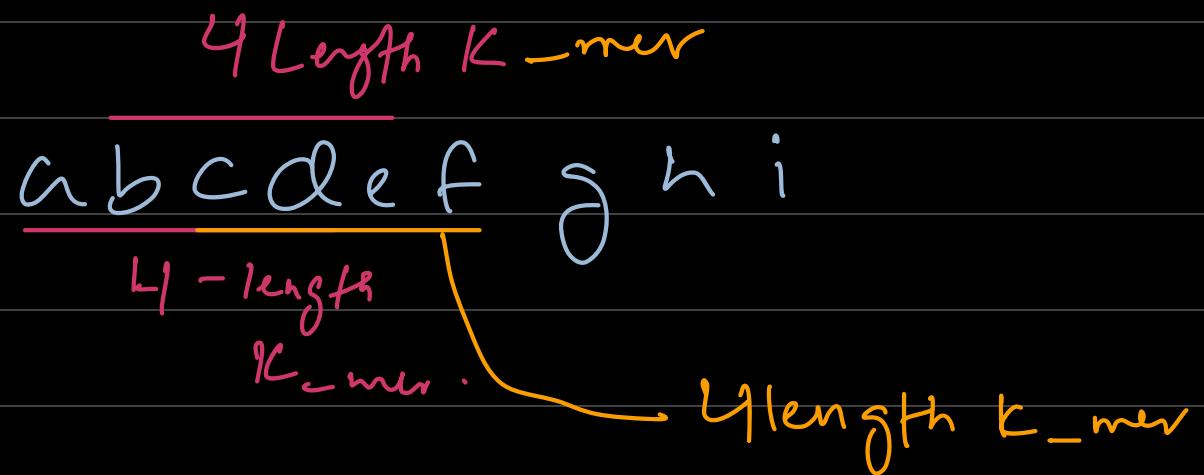
if count >= t :

K_mers.add(K_mer)

return k-mer

* Find most frequent words in a string - basically just use a dictionary.

$$K_mer = [i : i + k]$$



K-mer - key
Count - value

and then return the one with largest Count.