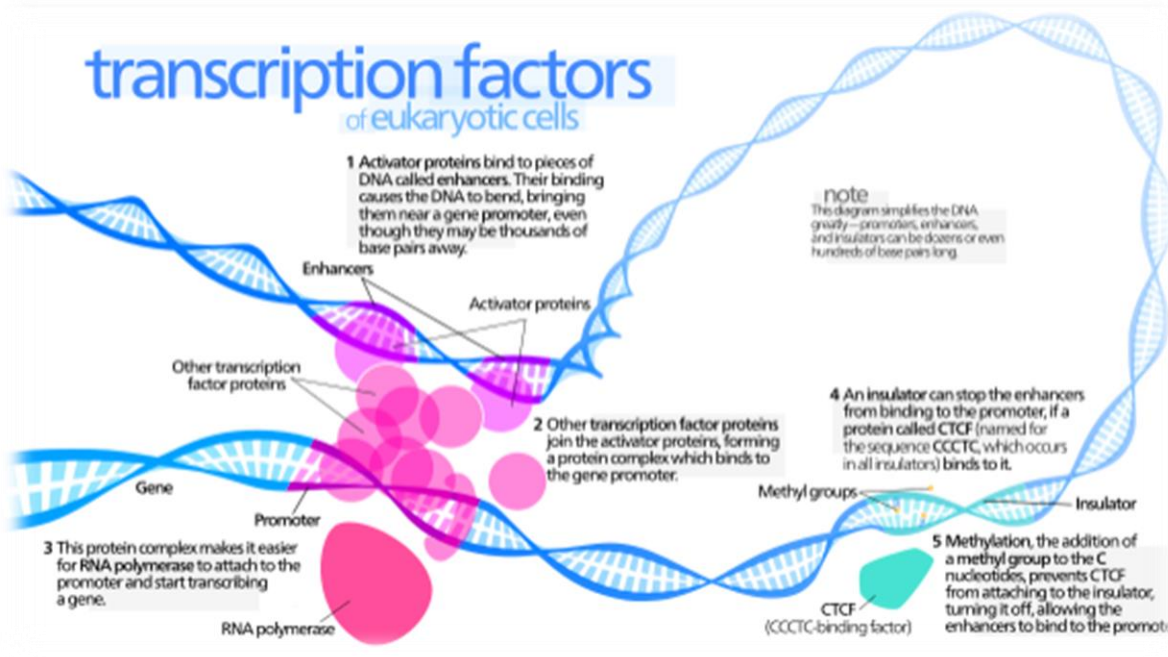


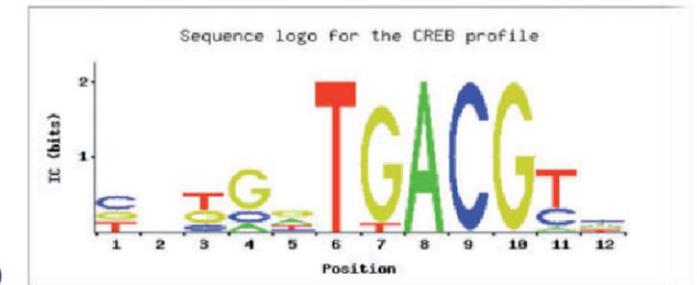
A tool to discover motifs among DNA sequences using Expectation Maximization (EM) algorithm



CTTGGTGACGTG
GTGAGTGACGTC
CGGGTTGACGCA
CCTACTTACGTA
TATGGTGACGTC
TCGGATGACGAT
TAGGATGACGTC
CCTGGTGACGCC
CGCGGTGACGTA
GCCGTTGACGCC
CGCGATGACGCA
CCTGTTGACGTG
TTGCATGACGTC
GTTGGTGACGTC
GAGGATGACGTT
GGTCGTGACGTA

CTGGGTGACGTC (Consensus)

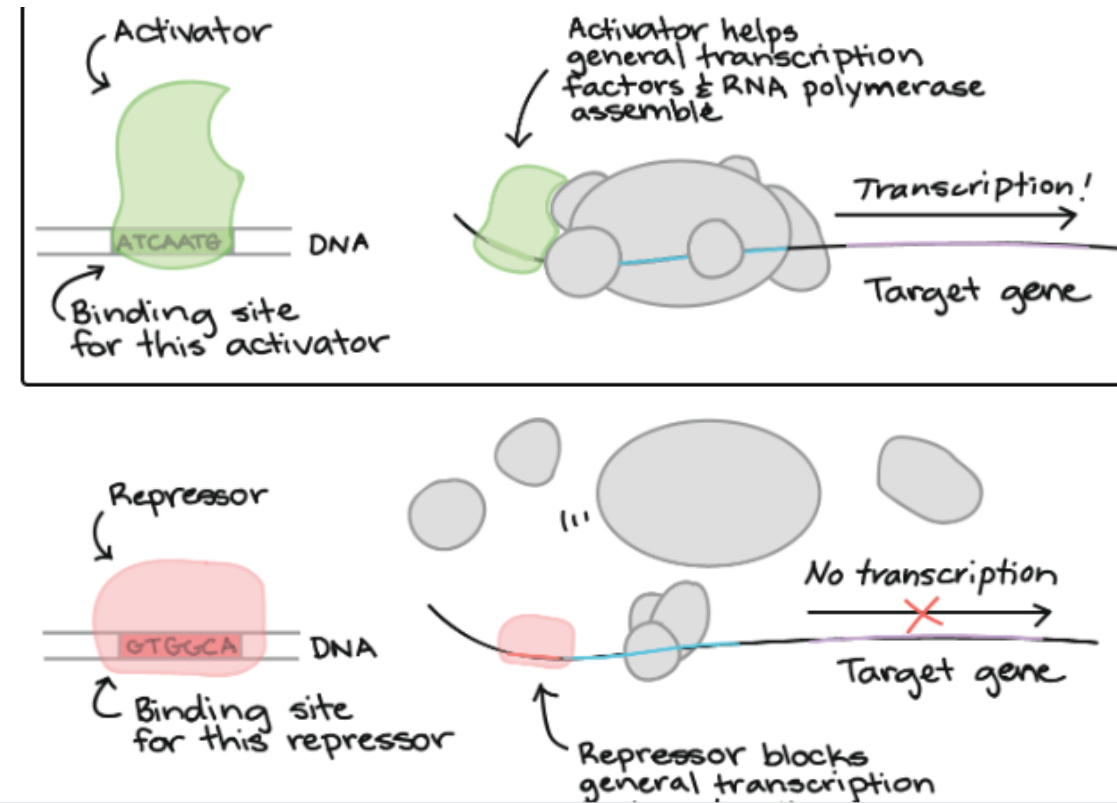
A	[0	3	0	2	5	0	0	16	0	0	1	5]
C	[7	5	3	3	1	0	0	0	16	0	5	6]
G	[5	4	6	11	7	0	15	0	0	16	0	3]
T	[4	4	7	0	3	16	1	0	0	0	10	2]



Huimin Lu
2023-03-14

Background

- In transcription, **transcription factors binding sites** are short sequences that combine transcription factor (TF).
- Transcription factors could activate or deactivate gene transcription
- And transcription factors binding sites are highly conserved.
- So, discovering TF binding sites sequence, motif, could help find **potential drug target** and genetic disease diagnosis.
- And this tool is to find motifs among DNA sequences



Tool presentation

- Three groups of **input dataset**
- One **python script** and import other python software called “Weblogo” .
- Finally getting a **motif logo**,

Package of this project

```
conda    22.9.0
weblogo  3.7
python   3.9
pandas   1.4
numpy    1.21
```

Python program

```
usage: main.py [-h] [-kmer KMER] [-A A] [-T T] [-C C] [-G G] [-X X] input_file output_file

Replication EM alorgrithm in DNA motif discovery

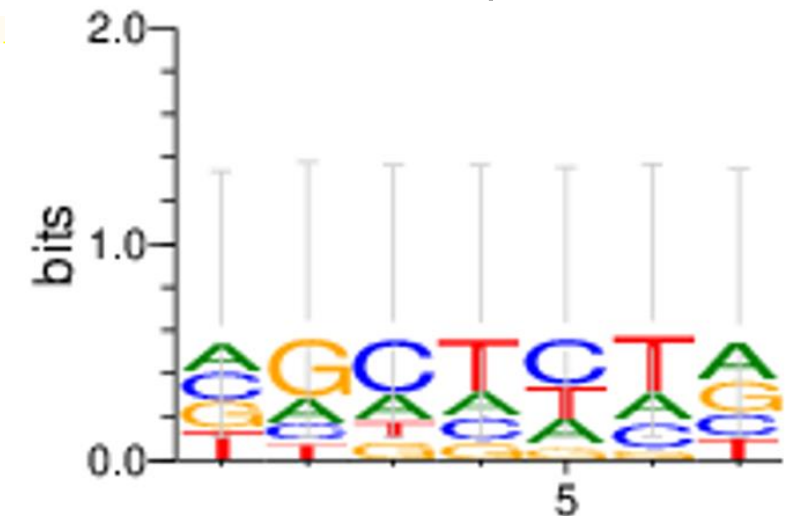
positional arguments:
  input_file  input fasta file, it must be standard fasta file
  output_file output txt file, this file will be used to make motif logo

options:
  -h, --help  show this help message and exit
  -kmer KMER  the length of motif
  -A A        background(non motif probability)
  -T T        background(non motif probability)
  -C C        background(non motif probability)
  -G G        background(non motif)probability
  -X X        probability of motif among dataset, used in calculate initial p matrix
```

Input file

```
>seq1
CTATAAACGTTACA
>seq2
ATAGCGATTGACT
>seq3
CAGCCCAGAACCT
>seq4
CGGTGAACCTTACA
>seq5
TGCATTCAATAGCT
>seq6
TGCTCTGTCCACTC
>seq7
TGCTCTGTCCACTC
>seq8
GGTCTACCTTTATC
```

output



Significance and novelty

Complexity:

- This tool does not use regular expression, instead using Expectation Maximization Algorithm, **iterating update**.
- This helps to find more potent motif whose patterns are difficult to predict using Regtex.
- The most difficult part is **to understand this algorithm and compute by my own script**

Feature:

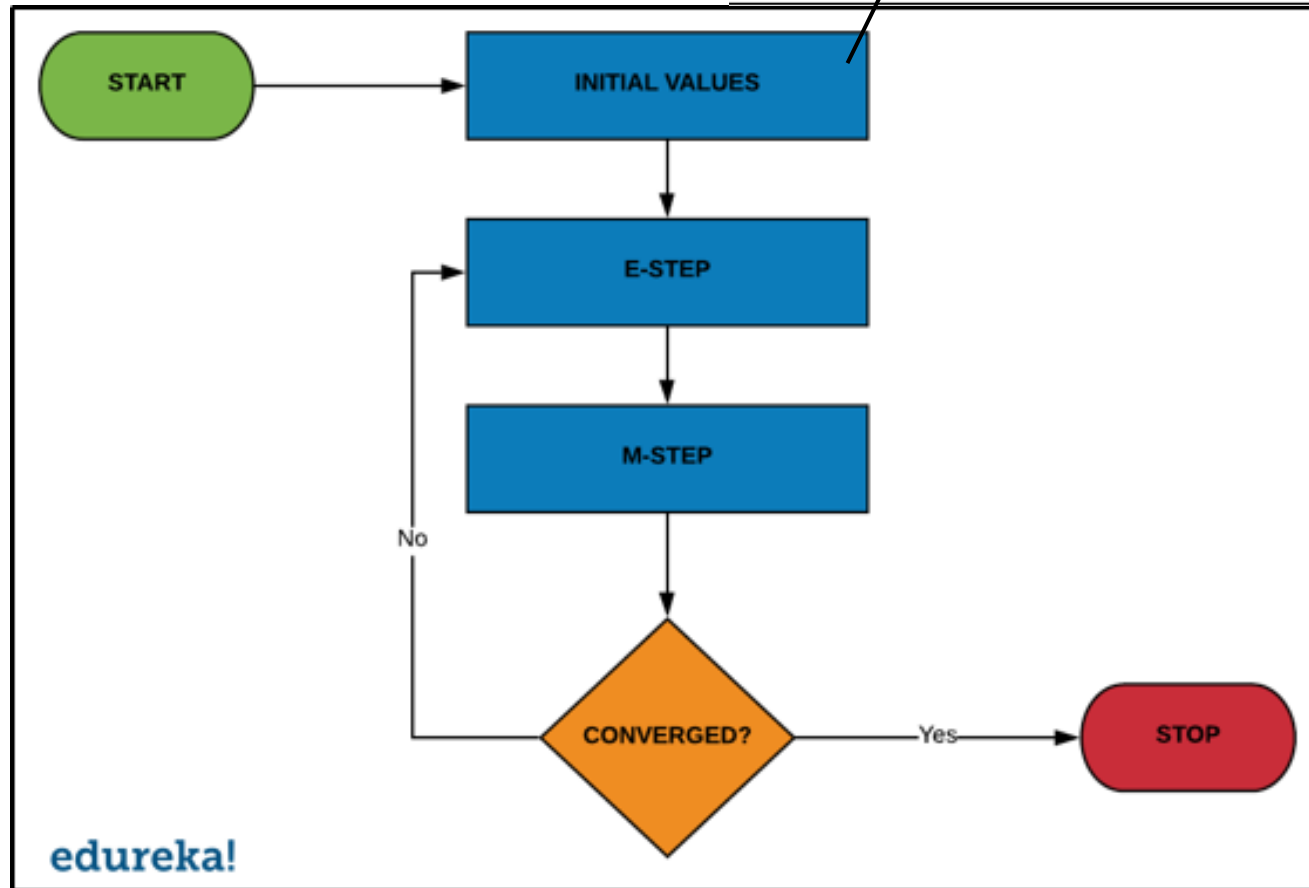
- This tool can calculate motif profile matrix and make motif logo

Limitation:

- Only discover DNA sequences not protein sequences.

Methodology

$$p = \begin{array}{c} \begin{array}{rcc} & 1 & 2 & 3 \\ \text{A} & 0.17 & 0.5 & 0.17 \\ \text{C} & 0.17 & 0.17 & 0.17 \\ \text{G} & 0.17 & 0.17 & 0.17 \\ \text{T} & 0.5 & 0.17 & 0.5 \end{array} \end{array}$$



edureka!

● E-step

Example: Estimating Z

$X_i = \text{G C T G T A G}$

$$p = \begin{array}{c} \begin{array}{rcccc} & 0 & 1 & 2 & 3 \\ \text{A} & 0.25 & 0.1 & 0.5 & 0.2 \\ \text{C} & 0.25 & 0.4 & 0.2 & 0.1 \\ \text{G} & 0.25 & 0.3 & 0.1 & 0.6 \\ \text{T} & 0.25 & 0.2 & 0.2 & 0.1 \end{array} \end{array}$$

$$Z_{i1} = 0.3 \times 0.2 \times 0.1 \times 0.25 \times 0.25 \times 0.25 \times 0.25$$

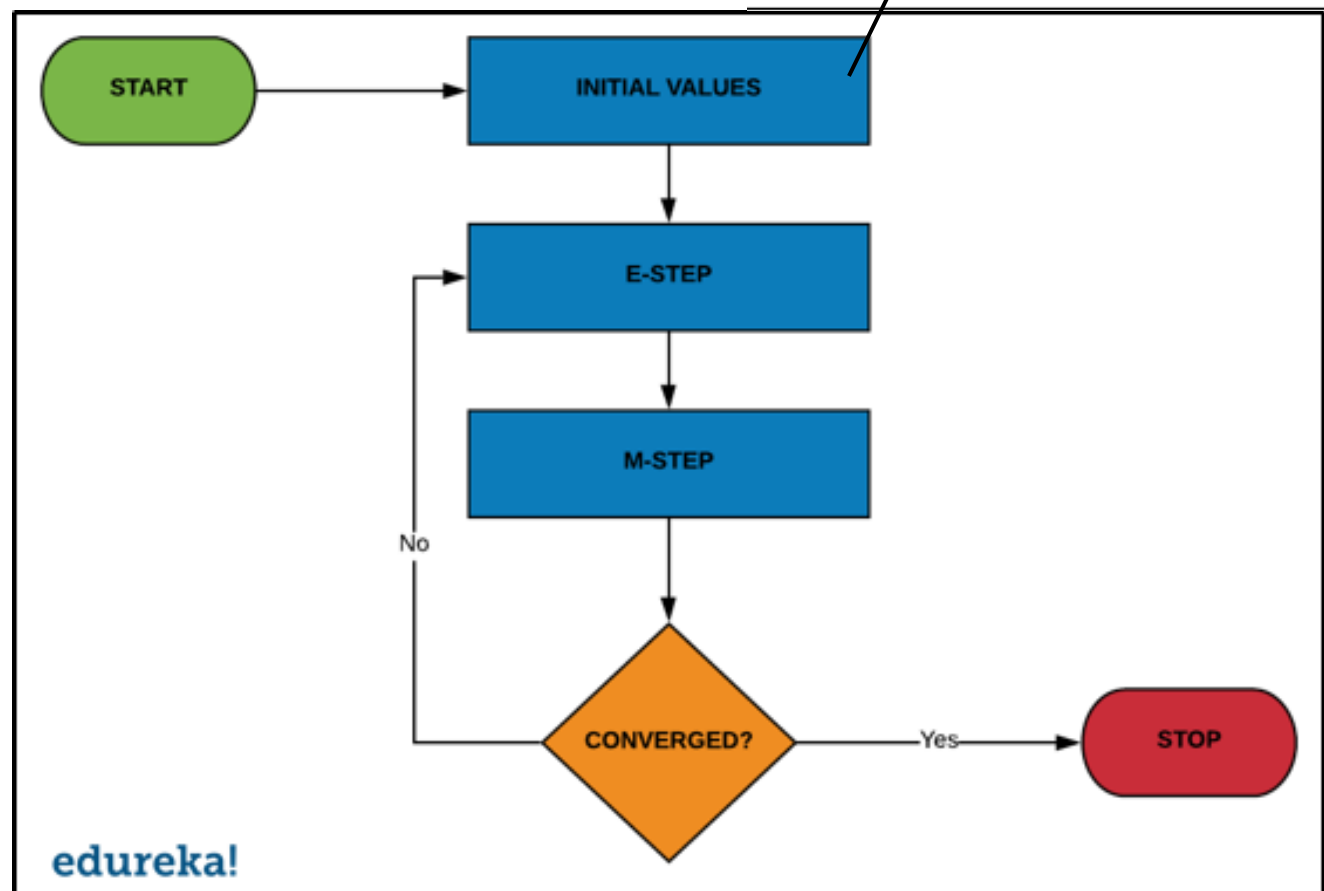
$$Z_{i2} = 0.25 \times 0.4 \times 0.2 \times 0.6 \times 0.25 \times 0.25 \times 0.25$$

⋮

- then normalize so that $\sum_{j=1}^{L-W+1} Z_{ij} = 1$

$$Z = \begin{array}{c} \begin{array}{rcccc} & 1 & 2 & 3 & 4 \\ \text{seq1} & 0.1 & 0.1 & 0.2 & 0.6 \\ \text{seq2} & 0.4 & 0.2 & 0.1 & 0.3 \\ \text{seq3} & 0.3 & 0.1 & 0.5 & 0.1 \\ \text{seq4} & 0.1 & 0.5 & 0.1 & 0.3 \end{array} \end{array}$$

Methodology



● M-step

The M-step: Estimating p

- recall $p_{c,k}$ represents the probability of character c in position k ; values for position 0 represent the background

$$p_{c,k}^{(t+1)} = \frac{n_{c,k} + d_{c,k}}{\sum_b (n_{b,k} + d_{b,k})}$$

pseudo-counts

$$n_{c,k} = \begin{cases} \sum_i \mathbb{1}_{\{j | X_{i,j+k-1} = c\}} Z_{ij} & k > 0 \\ n_c - \sum_{j=1}^W n_{c,j} & k = 0 \end{cases}$$

total # of c's in data set

Example: Estimating p

A C A G C A

$$Z_{1,1} = 0.1, Z_{1,2} = 0.7, Z_{1,3} = 0.1, Z_{1,4} = 0.1$$

A G G C A G

$$Z_{2,1} = 0.4, Z_{2,2} = 0.1, Z_{2,3} = 0.1, Z_{2,4} = 0.4$$

T C A G T C

$$Z_{3,1} = 0.2, Z_{3,2} = 0.6, Z_{3,3} = 0.1, Z_{3,4} = 0.1$$

$$p_{A,1} = \frac{Z_{1,1} + Z_{1,3} + Z_{2,1} + Z_{3,3} + 1}{Z_{1,1} + Z_{1,2} + \dots + Z_{3,3} + Z_{3,4} + 4}$$

Discuss

➤ **Pro:**

First time replicating algorithm

Learn skills in python programming

➤ **Cons:**

Did not find good method to calculating convergence

The result is not that precise