

DNA patterns that play the role of Molecular Clocks. (Randomized algs).

- Daily schedules of animals, plants, and even bacteria are controlled by "time keeper" called circadian clock.
- How do individual cells in animals and plants let alone in bacteria know when they should slow down or "increase" the production of certain proteins. Is there a "clock gene"
- genes responsible for "breaking" the circadian clock to cause DIPS.
- Today many circadian genes have been discovered these genes having names like timeless, clock, cycle "orchestrate" the behavior of other genes and display high degree of evolutionary conservation across specie.

- Every plant cell keeps track of day and night separately / independent of other cells, and that just three plant genes, called LHY, CCA1 and TOC1 are clock's master time keepers.

Regulatory genes and the regulatory proteins that they "encode", are often controlled by external factors, in order to allow organisms to adjust their gene expression.

- TOC1 promotes the expression of LHY and CCA1.
- LHY, CCA1 and TOC1 are able to control the "transcription" of other genes because the "regulatory proteins" that they encode are "transcription factors" or "master regulatory proteins" that turn other genes ON/OFF.

o) Master Regulatory proteins that turn other genes ON/OFF. (Also, called Transcription Factors)

o) A transcription factor regulates a gene by binding to a specific short DNA interval called a "regulatory motif" or "transcription

C
G
T

factor binding site" in the gene's upstream region (600 - 1000 bp region preceding the start of gene).

example: CCA1 binds to AATATATCT in upstream regions of many genes regulated by CCA1.

⇒ How to locate regulatory motifs without knowing how they may vary though.
they look:

⇒ Motif finding: discovering the hidden message "shared" by collection of strings. called "Regulatory Motifs" or binding site for transcription factors which are "master regulatory proteins". (can turn other genes ON / OFF).

⇒ Regulatory Motif is a "pattern" that
 "appears at least once" (perhaps with variations or
 mutations) in "each" of the many different
 "regions" that are scattered throughout the genome.

↳ strings

exact d mismatches rather \leq may be

\equiv

d mismatch $\text{Actual kmer looking for}$

(15, d) motif

Generate all k -mers with d mismatch.

Any (k, d) - motif must be at most d mismatches apart from some k -mer appearing in the first string of DNA.

```

def Immediate_Neighbors(pattern):
    neighborhood = set()
    for i in range(0, len(pattern)):
        symbol = pattern[i] #ith Nucleotid.
        for c in "AGTC": # for each nucleotid.
            if c != symbol: x different from symbol
                · add(pattern.replace(symbol, c))

```

return neighborhood

```

def iterative_neighbors(pattern, d):
    neighborhood = set()
    neighborhood.add(pattern)
    for i in range(0, d):
        neighbors_copy = neighborhood.copy()
        for pat in neighbors_copy:
            neighbors = Immediate_Neighbors(pat)

```

Solving neighborhood · update(neighbors)
like union
or

neighborhood |= neighbors

implies union.

return neighborhood

def Neighbors(pattern, d):

if $d = 0$

return {pattern}

if $\text{len}(\text{pattern}) == 1:$

return {A, C, G, T}

A
C
T
G
T

neighborhood = set()

suffix = Neighbors(pattern[1:], d)

for text in suffix:

if hammingdistance(pattern[1:], text) $< d:$

for c in "AGTC":

new = c + text

neighborhood.add(new)

else:

new = pattern[0] + text

neighborhood . add(new)

return neighborhood .

distance b/w all patterns and strings

def Distance Between Pattern AND STRINGS(pattern, dna):

K = len(pattern)

distance = 0

for gene in dna:

hamming_distance = float('inf')

for i in range(0, len(gene) - K + 1):

if hamming_distance > Hamming_Distance

(pattern, gene[i:i+K])

hamming_distance = 2

distance = distance + hamming_distance

return distance

```
def median_string(dna, k):
    distance = float('inf'), median = ""
    all_kmers = generate_kmers(k, ['A', 'T', 'G', 'C'])
    for kmer in all_kmers:
        if distance > d(kmer, dna):
            distance = d(kmer, dna)
            Median = kmer
    return median
```

A
C
G
T

$\Pr(ACTGGGATTAC | \text{Profile})$

↓

probability

from profile matrix

$$0.2 \times 0.6 \times 1 \times \dots = 0.00083988$$

Look for k -mers that have larger values of $\Pr(\text{text} | \text{Profile})$.

def profile_most_probable_kmer(text, k, profile_matrix)

mapping = { 'A': 0, 'C': 1, 'G': 2, 'T': 3 }

$k\text{-mer probability} = \{ \}$

for i in range(len(text) - k + 1):

$k\text{-mer} = \text{text}[i : i+k]$, probability = 1

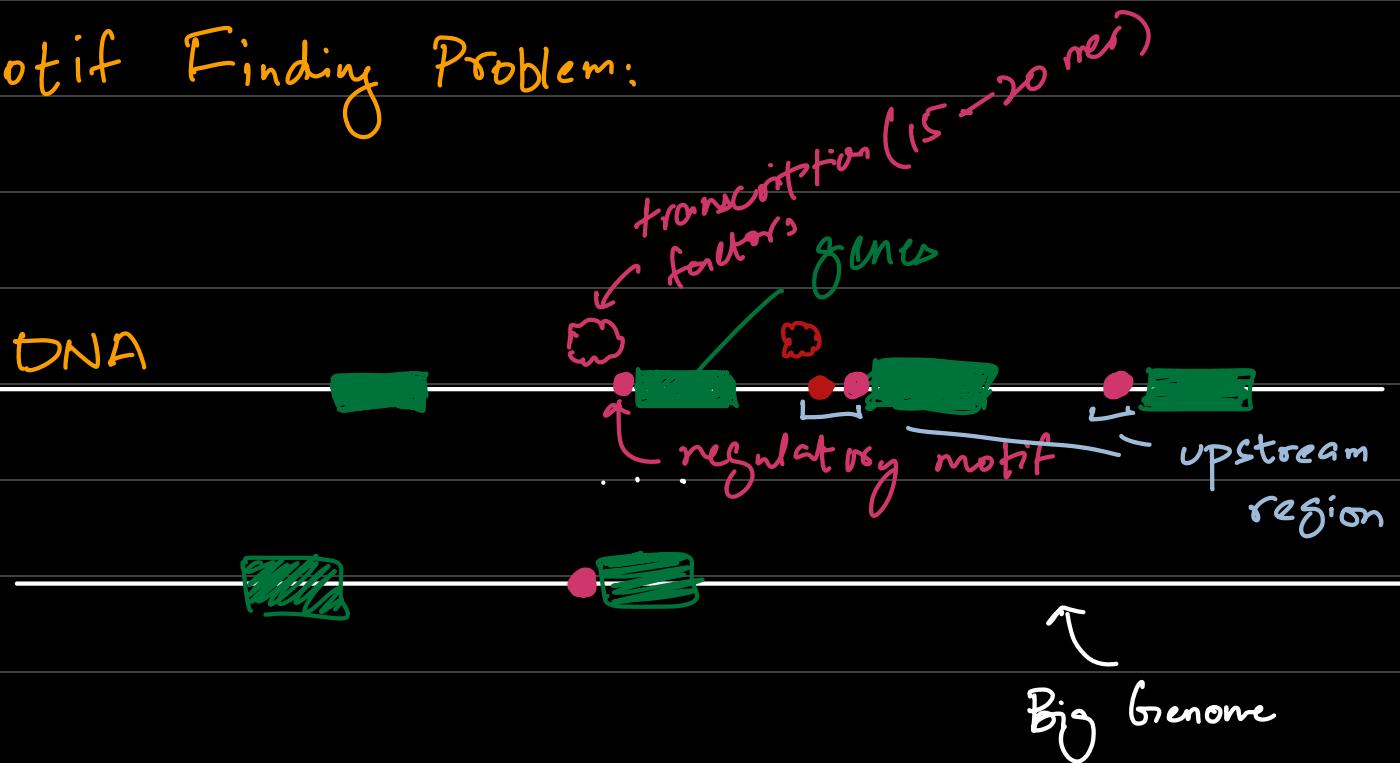
for j in range(len(k-mer)): 0 → 4

probability * = profile_matrix[k-mer[j]][j]

$k\text{-mer_probability}[k\text{-mer}] = \text{probability}$.

maximum = $\max(k\text{-mer_probability} \cdot \text{values}())$
return key with this maximum.

Motif Finding Problem:

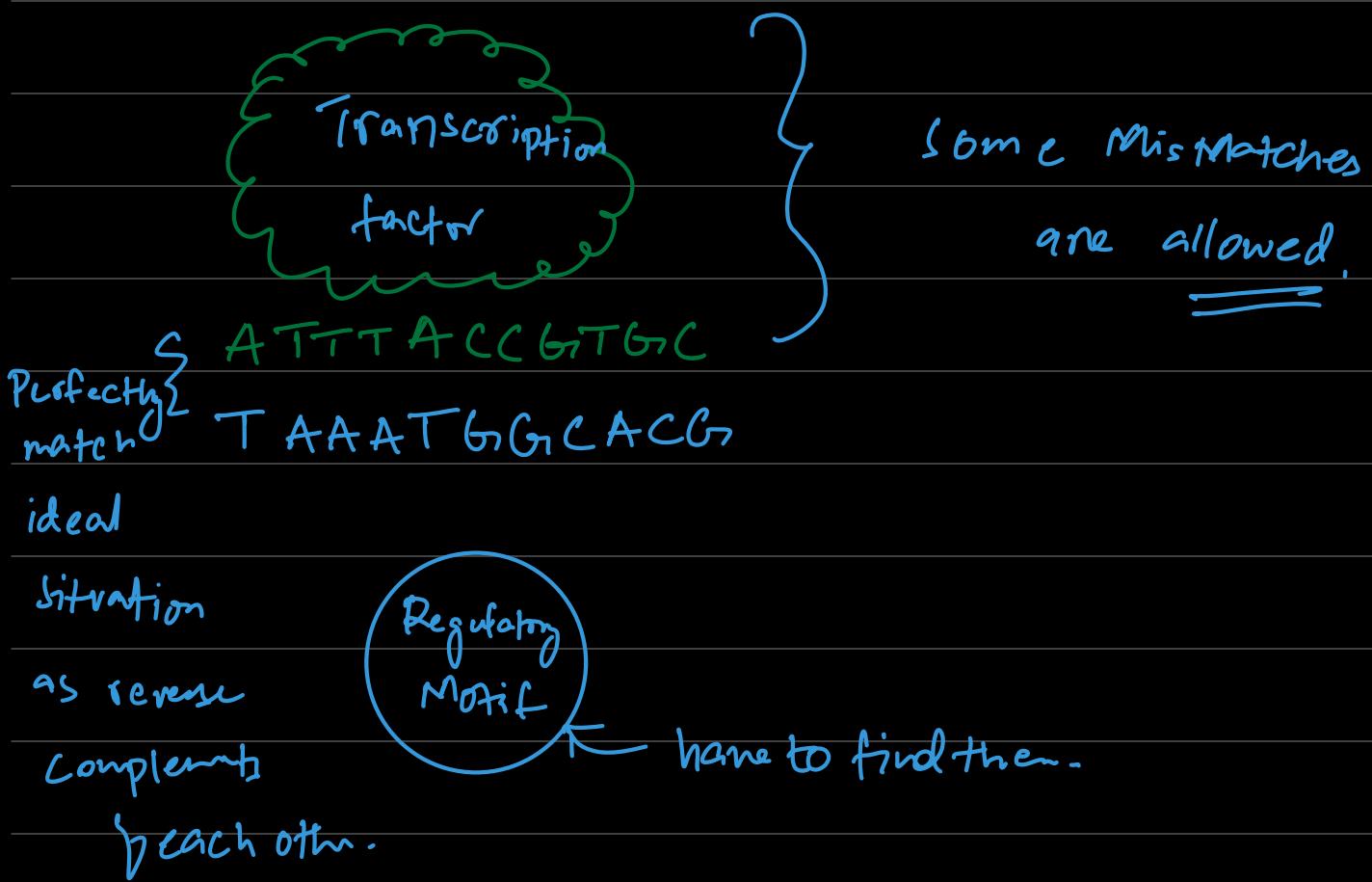


- We might know that some collection of genes are related to function e.g. immunity genes. When organism sick all the immunity genes "turn on" at the same time and "start producing proteins".

- How do all those genes get regulated to work together.
- There are short snippets (15 - 20 bp long) of sequence before the gene that can be bound by the same transcription factor.
- When transcription factor binds to a regulatory motif, the "expression's the gene changes (can go up or down).
- A transcription factor is conductor giving signal to bunch of different genes and regulating them.
- A same gene can be regulated by multiple transcription factors (e.g.: may be a gene is related to "recovery from sickness" as well as "Pregnancy").
- If we know that several genes are related to same event. (e.g. their expression

changes after athletic activity) we can search for "common motifs" in "upstream" of these genes.

- If we find correct motif it may enable us to identify additional genes regulated by the "transcription factors".
- Transcription factor don't always bind perfectly to regulatory motifs.



- 10 sequences / strings of DNA
- Where 15 - mers inserted in each sequence at random location with 4 random mismatches.
- The true motif is ["]AAAAAAAAAGGAGGAGGAGG["] Consensus Motif
but exact sequence never appears.

strings or upstream strings \rightarrow generate k -mers of some length ($5, 1$)-motif \rightarrow that k mers nearest

\uparrow \uparrow
 k d

neighbors are generated with d mismatch \Rightarrow one we call such "kmers".

we iterate through every kmer and check if it appears once with d mismatch in every sequence or string of DNA.
Then it is (k, d) motif, but this algo is slow.

- Collection of Motifs \Rightarrow give you a profile matrix
- Using profile matrix we can calculate the probability of every 1Kmer or motif,
- The profile matrix gives the probability of generating the Kmer or motif.
- Kmer has higher probability when it is similar to Consensus String of a profile.

o) Profile - Most Probable 1Kmer.

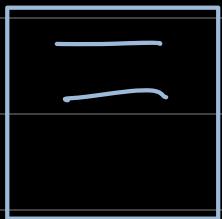
\Rightarrow Greedy Motif Search:

Starts by forming a motif matrix from a arbitrary "Selected Kmers" in each string of DNA. (in our implementation best-motifs are first Kmer in each string). The algorithm then attempts to

improve the initial motif matrix by trying each of the Kmers in Dna_1 as the first motif. For given Kmer Motif₁ in Dna_1 , it

builds a profile vector "Profile" for this

Initial/motifs. lone Kmer, and sets Motif₂ equal to the profile - most probable kmer in Dna_2 . updating profile as the profile matrix found for motif₂



$$motif_3 = \text{Profile most probable} ([motif_1,] [motif_2])$$

After obtaining Kmer from each string to obtain collection motifs. — also tests to see whether Motifs outscore the current best scoring collection of motifs and then moves Motif₁ one symbol over in Dna_1 beginning the entire process of generating motifs again.

Randomise Motif search (Dna, k, t) :

Motifs \leftarrow randomly chosen from each strand
Best motifs \leftarrow motifs

while True:

Profile \leftarrow Profile (Motifs)

New motifs \leftarrow profile most probable (profile, Dna)
motif
append

if score (Motifs) $<$ score (BestMotifs):

BestMotifs \leftarrow motifs

else:

return BestMotifs

Rule of Succession:

Just because we have not observed an event does not mean that we should assign its future probability to zero.

Add pseudocount value to the counts of each type of event before normalizing.

Motifs: T A A C

G T C T

A C T A

A G G T

Profile:

A	2/4	1/4	1/4	1/4
C	0	1/4	1/4	1/4
G	1/4	1/4	1/4	1/4
T	1/4	1/4	1/4	2/4

Adding a Pseudocount of 1 produces following
Count and profile matrix.

Count (matrix) original

A: 2+1 1+1 1+1 1+1

C: 0+1 1+1

G: ... - - - -

T: . - - - -

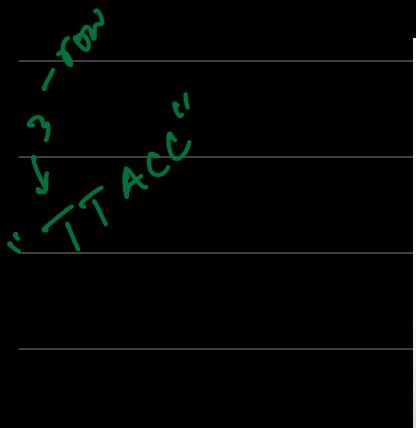
Profile Matrix

$\left\{ \begin{array}{cccc} \frac{3}{8} & \frac{2}{8} & \frac{2}{8} & \frac{2}{8} \\ \cdot & \cdot & - & - \\ \cdot & \cdot & - & - \\ \cdot & - & - & - \end{array} \right.$

Modify Create Profile Matrix to account for pseudocounts.

→ add one to each entry of counts.

$$\text{profile} = \text{count}/(t+4)$$



* After we have found collection of motifs \Rightarrow we can find consensus of those motifs for true motif.

*> Looking for hidden messages:

If collection of genes are implicated in the same function (e.g. circadian clock or responsible), then a single transcription factor may "bind" to the same "Keyword" in many of the genes' upstream regions, perhaps with "minor variations".

Regulatory motifs

↳ they will have slight variations.

⇒ Randomized Algorithms:

o) Randomized Algorithms use randomness to achieve their desired outcomes.

Las Vegas Algorithms: Results are guaranteed to be correct, but the amount of time the algo will take is random

Ex: Quick Sort Algos

Monte Carlo Algorithms rely on repeated sampling of results. Not guaranteed to be correct. (fast).

Ex: Finding the probability of biased die landing on "each side" by rolling it many times.

↳ estimating probability of landing on each side.

def Randomized Motif Search (Dna, k, t):

Motifs → randomly select kmer motifs from each sequence
best motifs = ↳

while (forever):

Profile ← Profile (Motifs) → iterate through each sequence one by one.

Motifs ← Motifs (profile, Dna)

↳ Profile now contains kmer

new set of

Motif or profile motif if score(Motif) < score(Best Motif)
probable
kmers.
Best Motifs \leftarrow motif

else:

return Best Motifs

(*) If initial random set of motifs

- Run 1000 includes at least one of the "correct" motifs, the profile matrix (slightly) biased towards the correct motif
- This bias points the search in

\Rightarrow Gibbs Sampling: the right direction.

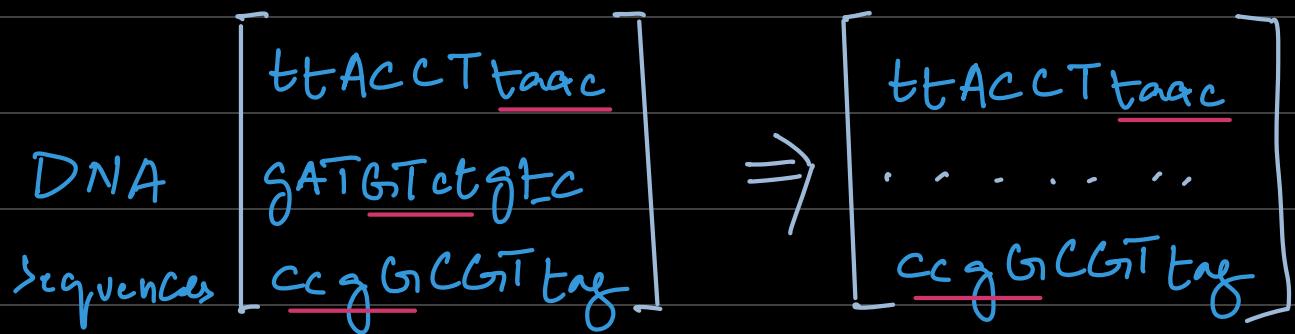
- This is a Markov - Chain Monte Carlo (MCMC) method.
- is method to "draw samples" from a "Complex probability distribution" by breaking it down into "Simpler pieces"
 -) Have Unknown things to estimate (variables $x_1, x_2, x_3 \dots$)
 -) Instead of trying to figure out all at once, fix all but one variable, and draw sample for that single variable based on the others.

•) rotate through each variable, each time sampling the new value from its "Conditional distribution" given the latest values of all the other variables.

•) By cycling the variables many times, the collected samples start to reflect the true underlying (joint) distribution you're interested in.



*) Gibbs Sampling will "change" only a "single kmer" in each step, as well as changing this kmer more liberally.



- pick random motifs of some length k from each sequence. Gibbs Sampling randomly selects one of the strings to be replaced.
- Generate Profile Matrix for the remaining motifs. Use that profile matrix and the sequence to get probability distribution for each kmer for that sequence and select one randomly weighted by the probabilities after "normalizing" them so that they sum to 1.

$$\left(\frac{1}{8}\right) \cdot \left(\frac{2}{8}\right) \cdot \left(\frac{2}{8}\right) \cdot \left(\frac{1}{6}\right) = \frac{4}{8^4}$$

o) Set of probabilities for each kmer.

o) Randomly sample or choose kmer randomly
but it is proportional to probability
given here sort of like roulette wheel.
high probability more likely to get
selected.

o) N times.

GibbsSampler(Dna, k, t, N)

randomly select k -mers $Motifs = (Motif_1, \dots, Motif_t)$ from Dna
 $BestMotifs \leftarrow Motifs$

for $j \leftarrow 1$ to N

$i \leftarrow$ randomly generated integer between 1 and t

$Profile \leftarrow$ profile formed from all $Motifs$ other than $Motif_i$

$Motif_i \leftarrow Profile$ -randomly generated k -mer in Dna ;

if $Score(Motifs) < Score(BestMotifs)$

$BestMotifs \leftarrow Motifs$

return $BestMotifs$

$n=1000$

run this whole thing
20 times

HW 2:

o) Increase g. length 15 - 20

o) More iterations.

* Gibbs Sampling:

— Markov Chain
Monte Carlo

Sample from a joint distribution that is hard to directly sample from.

- Get values for bunch of variables that are related to each other in some way.
- Health Metrics: height, age, weight, blood pressure, etc
 - ↳ generate a collection of these variable values that would be realistic for real human, without being same as data we have access to.

Gibbs Sampling: want to generate values for x_1, \dots, x_n

① Pick initial guesses for each variable.

② Update the value of x_i by sampling metric

Some sample

Set v random values for all

the probability $P(x_1 | x_2, x_3 \dots x_n)$

Find the height that's most probable given the other measurement values.

- ③ Update the value of x_2 by sampling $P(x_2 | x_1, x_3 \dots x_n)$

Find the most likely age for this person.

- ④ Do this for all variables.

- ⑤ Iterate many times ↗

* Wait for samples to "converge"
for human dataset.
 .

Problems: skewed may be our sequence is all As or background all A's.

AAAA GT CG AAA
 AA GCTGA AAA
AAAAA GCCTAA
ACCCGA AAA
AA AAA GCG

⇒ Scoring function: "relative entropy"

⇒ Let $w = A \text{ or } T$ WC6w

⇒ Let's observe protein binding to DNA sequence.

- o Try identifying a set of DNA sequences with similar function, and then find motifs.
- o Try analyzing chip-seq data.

• Try different scoring functions (entropy, relative entropy).

Sequences - strings of DNA

best motifs = all first kmers of all sequences of DNA

Start with first kmer of seq₁ (we iterate through all)

motifs = [kmer]

profile = profile(kmer) # create profile

iterate through every sequence seq₂ ... seq_t:

motif = profile-work-probable(profile,

add this new motif to motifs^{seq_t}

profile = profile(motifs) # create

if score(Motifs) < score(BestMotifs):

BestMotifs = Motifs

return best motifs.

LCP choosing local optimum
at each step to reach the
global optimum

* Basically we start with First kmer of seq1 generate its profile and move to second sequence → we send that sequence and this profile to profile → most probable - kmer , which gives probability to every kmer of second sequence and we return one with highest probability of being generated with that profile. we add new motif or returned one to list of motifs , and create new profile → there two motifs and move to seq 3 then \Rightarrow same process again until we reach end of all sequences \Rightarrow then score these new collection of motifs to the best motifs if score is less than new motifs become best motifs.

* This process ends when we have gone through all the kmers of seq2.

```

def score(motifs):
    consensus = ""
    rows = len(motifs) # rows 3
    cols = len(motifs[0]) # 4
    for j in range(0, cols):
        freq = {} # freqency
        for i in range(0, rows):
            motif[j][i]
            freq[motif[i][j]] = freq.get(motif[i][j], 0)
            +1
        consensus += max(freq, key=freq.get)

```

score = 0

for motif in motifs:

score += hamming_distance(motif, consensus)

return score