

→ 10,000 → $\boxed{PCs} \rightleftharpoons \boxed{PCs}$

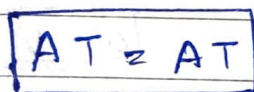
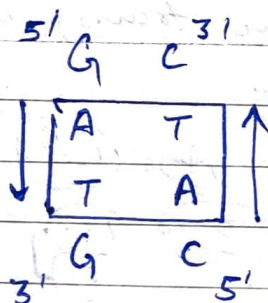
②

feature

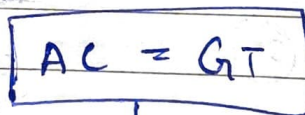
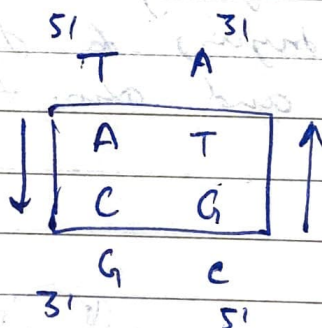
→ features :- 28 structural feature + 3 energy features

values of these 28 str. features obtained by analysis of B-DNA crystal structure and values of energy features obtained from in-house software.

- values for all 10 unique dinucleotides were obtained and averaged to get a particular value.
- Unique according to us need to be 16, however they are 10 because :-

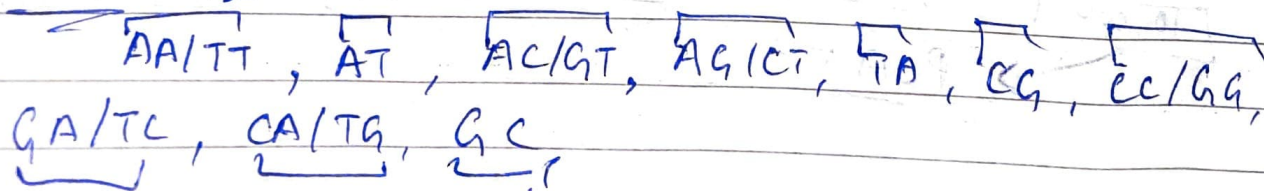


↓
AT is one of the 10 unique dinucleotides.



↓
So instead of AC & GT to be 2 different dinucleotides they both are equivalent and hence together as 1 of the unique dinucleotide step.

all 10 are :-

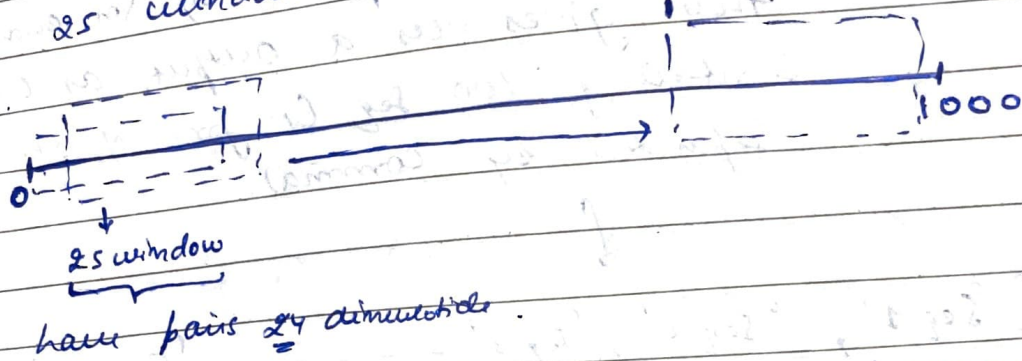


→ Training ^{Test} dataset for
TSS :- 12880 seq.

for
CDS :- 6281 seq.

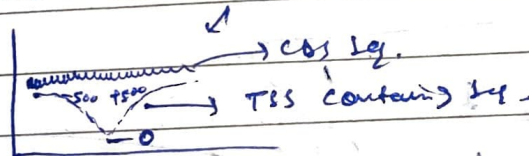
all these seq. were converted to 31 numeric sequences
based upon the values of 10 dinucleotide steps.

→ each sequence of both TSS & CDS were of length
1001 nt. on this 1001 nt a window of 25
was set at index starting 0 and ending at 974
this 25 window now have 24 different dinucleotide
steps.



eg S1 & S2

LISI Shows that all 31 parameters ^{in TSS sequence} have some
changes while CDS shows no change.



* all 31 parameters have discriminative power, however
their scale were not same and also some had
increasing effect other had decreasing (which might
skew/bias our observation) and also they were some-
what correlated, hence a PCA need to be done.
(after scaling/normalisation).

Code language :-

→ we have say a input of 100 seq.

→ from these 100 seq., we choose 1st.

→ in this 1st seq. we set a window of 25,
for this window 31 parameters are calculated for 24

dinucleotide steps.

- These 31 parameters are then normalised (0 to 1).

- Now window shift towards right (so that can cover all 1000 nucleotide). i.e. the shifting & step 3 to end is repeated 975 times.

main.py module

① → Read Seq. # It input a .txt file containing 100 sequences, separated by \n character this gives us a output as list containing 100 Seq. (in form of strings) separated by comma.

i.e.

['Seq 1', 'Seq 2', 'Seq 3' - - - - - 'Seq 100']

↓
Now this list is then converted to a map/dict whose structure is:

```
{ 0: 'Seq 1'
  1: 'Seq 2'
  2: 'Seq 3'
  ...
  99: 'Seq 100' }
```

possible optimisation :-

* python might have a inbuilt function that can directly convert an input file to a dict.

FINAL OUTPUT: SEQUENCE MAP

NCE

(2) get Parameters details :-

→ takes sequence map as input

→ it goes to each seq. in sequence map and perform following functions :-

(a) Calculate-parameter

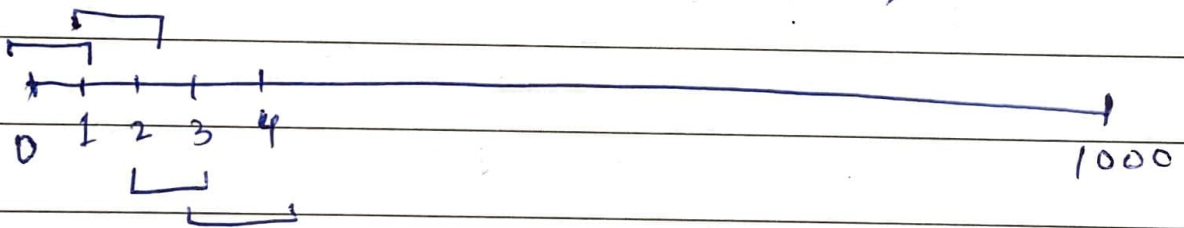
DOUBT :- Why we are considering a 25 length window. Why can't we just proceed with the whole sequence taking diminishing steps.

if we take 25 windows then total values are $975 \times 25 \times 31$, but if we take diminishing step then we have 998×31 .

① calculate parameters

- is run for each sequence (100)
- first it assigns scores/values to each ~~dist~~ dinucleotide & base pair.

for example: (A+T) will have diff score.



in this way, keeping in mind the 10 different ~~dist~~ dinucleotide steps, we assign the scores to 31 parameters.

the result of this function is a param_map of the structure: (1000) for each 1000 bp.

{ 'a': [- , - , - , - , - , - , - , - , - , -]
 'b':

ae :

}

② calculate Moving Averages (param map)

for every parameter in params-map:
 say 'a': [1000 values]
 for each bp.

we take a window of 25 and
 take average of 1st 25 values for
 parameter 'a'.

then we move right by 1 step
 & repeat the above step.

for 'a': $\left[\frac{\text{sum of 0 to 24}}{25}, \frac{\text{sum of 1 to 25}}{25}, \right.$

$\left. \frac{\text{sum of 2 to 26}}{25}, \dots, \frac{\text{sum of 975 to 999}}{25} \right]$

'a': [975 values (averaged for
 a window of 25 upto 1000)]

Doubt : why is this step needed?

this function returns moving-params-map
 of the structure:

975 x 25 x 31, but if we take
we have 998 x 31.

[① calculate parameters, ② calculate moving averages],
↳ on other

③ normalise moving averages

input: { a: [975 values] → Moving params
map.
a: [975 values] }

→ It is just normalizing all 975 values for each parameter.

★ Optimisation: python may have a inbuilt normalising function

Output :- same as above, but all values (975) are in range 0 to 1 (after normalization)

Output is named as Normalised map

Optimisation:- In calculate-parameter function we are assigning values to every individual B-P i.e. (1000 BP for each sequence), we are extracting from values from supple table which unique codon-number step and values of them corresponding to all 31 parameters.

We can use this table in our optimisation.

→ Now we are going for PCA-regression module.

Input is the following:-

{ 0 : { a : [975 non-coding value]
 b :
 ac : [] }

975-400-1
574

100 : { a : [975 non-coding value]
 ac : [975 Non-coding value] }

0 to 100 are the sequences
a to ac are the parameters.

* Extra task :- Calculate correlation for 31 parameters

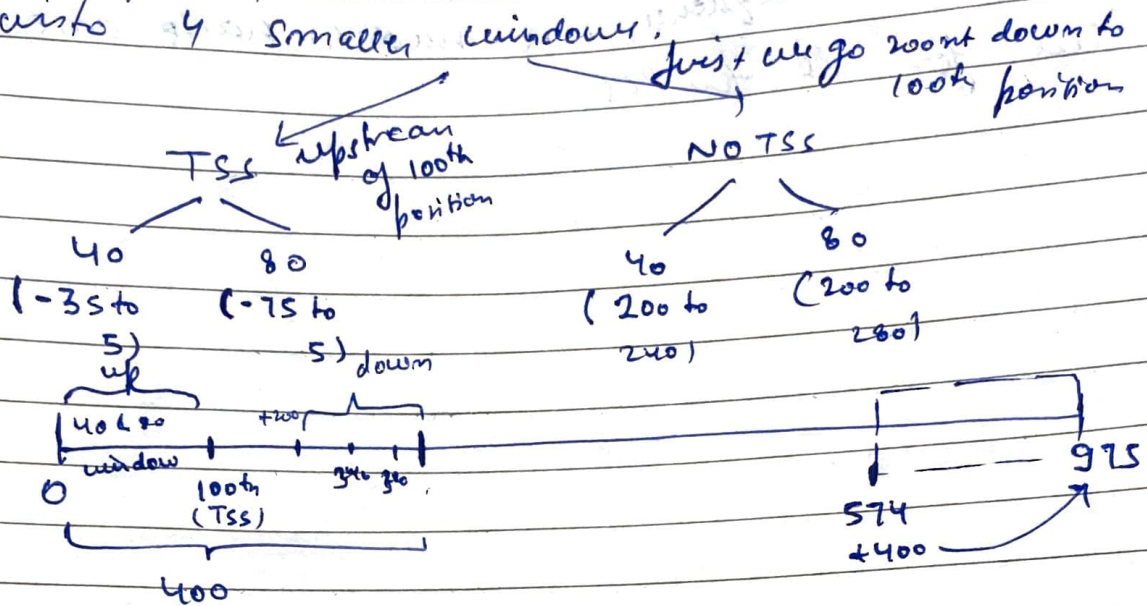
* In the PCA the main motive is to reduce the no. of features/parameters to some PCs.

for this there are 2 forkline one is with regions upstream & downstream of TSS & other is with motif combination.

① Use first discussion the 1st pipeline.

→ here we input our sequence 1st :-
 → for this sequence ~~the length is~~ all the 31 params have 975 values.

→ In this 975 values we are considering a 400 window
 → in this 400 window, our TSS is at 100th (consumption), and then we divide this 400 window into 4 smaller windows.



→ this is repeated for 574 times (i.e. the whole sequence get covered).

↳ for this we have :-

for each window exhaust window function is called, this function averages all the values of each 40/80 window, into 31 parameters.

for 975, consider a 400 window, in this window we assume that TSS lies on 100th position, from TSS.

for each sequence we get output in form of

Seq-40-map : { 0 : [31 para for TSS-40 window] - [31 para for NOT TSS 40 window] }

574 : [] [] []

SAME is for Seq-80-map & Seq-100-map

DOUBT :- why there is a leg-100-map?

⑥ Now all these 3 maps are sent for PCA with Predict-PCA.

Training

↳ Training Now for all sequence we need to have one value for all the PCs.

