

**UNIVERSITY OF MORATUWA**

Faculty of Engineering



Department of Electronic and Telecommunication Engineering

**BM4321: Genomic Signal Processing**

**Assignment 1**

**Promoter Discovery in Bacteria**

K.G. Abeywardena – 160005C

This report is submitted in partial fulfillment of the requirements  
for the module BM4321: Genomic Signal Processing

**June 2020**

## OVERVIEW OF ASSIGNED GENBANK ACCESSION OF THE BACTERIA

01	Organism	Escherichia coli
02	Family	Enterobacteriaceae
03	Genus	Escherichia
04	Description	Typically present at the lower intestines of humans, where it is the dominant facultative anaerobe present, but it is only one minor constituent of the complete intestinal microflora.
05	Accession	NZ_AP018808.1
06	Summary of Genome	
	Base Pairs	5678205
	Coding Genes/Protein Count	(5478) 5868
	Gene Count	(5990) 6129
	Sense strand Genes	2786
	Antisense strand Genes	2692
	GC content (%)	50.54

### QUESTION 1

- a) Standard Intact Query Local Search (IQLS) to locate *Pribnow box promotor (TATAAT)* within upstream position from 5 to 30

**Results:** By accessing the NCIB Genome Website<sup>i</sup>, under *genome assemblies*, the relevant accession number chromosome *NZ\_AP018808.1* was located under the column *Replicons* and its corresponding protein table of 5868 proteins under the column *CDS* and downloaded the full chromosome as a *fasta file* and protein table as a *csv file*.

Based on the *protein table*, sequences having a length of *50 bases upstream* and *3 bases downstream* were extracted from each coding gene on both sense and antisense strand. To make sure the extracted sequences contain 53 bases after the removal of *EOL*, a *safety of 3 bases* was added and then adjusted for 53 bases. The upstream direction had to be adjusted based on from which strand the sequences are extracted.

First the presence of *Methionine* site was checked for each of the obtained sequences by checking whether the last three bases correspond to *ATG*. By performing *Standard IQLS*, *Pribnow box promotor (TATAAT)* was located within a search region of upstream positions from 5 to 30 in the selected sequences. For intact query, scores were assigned as; *match = 1*, *mismatch = -1*, *gap penalty = 2*.

# seq with mutated <i>Methionine</i>	633	11.56%
# seq in sense strand with <i>Pribnow box</i> located	2112	75.81%
# seq in antisense strand with <i>Pribnow box</i> located	2042	75.85%
Total seq with <i>Pribnow box</i> located	4154	75.83%
# seq with <i>Pribnow box</i> unlocated	691	12.61%

**Discussion:** Based on the analysis for the presence of *Methionine site*, 633 sequences were detected with mutated *methionine sites* which were discarded as it corresponded to a lower percentage compared to the available genes. Since *Pribnow Box* can easily mutate (*A -> T* or *T -> A*), using *exact alignment* in traditional *Standard IQLS* to locate *Pribnow Box* is not ideal. Instead, *W-search* was used to locate *Pribnow Box* where both *Adenine (A)* and *Thymine (T)* were considered as a common base *W*. Based on the search, *Pribnow Box* was located in more than 75% sequences within search region.

- b) Obtaining *Position Probability Matrix (PPM)* using first 1000 sequences with 10 positions for the *Pribnow box*.

**Results:** Out of the selected 4154 sequences first 1000 sequences were used to obtain the *Position Frequency Table* which was then converted to the *PPM* using the below equation. As the *k value*, 0.01 was used.

$$p_{j,N} = \frac{f_{j,N} + k}{\sum_N (f_{j,N} + k)} \quad \text{Where } f_{j,N} - \text{frequency of a base at position } j (N \in \{A, C, G, T\})$$

$$p_{j,N} - \text{probability of a base at position } j (N \in \{A, C, G, T\})$$

Table 1: PPM obtained from the first 1000 sequences for Pribnow Box with Entropy measure

Base	Position									
	1	2	3	4	5	6	7	8	9	10
A	0.554	0.438	0.456	0.423	0.343	0.294	0.273	0.290	0.359	0.305
C	0.001	0.019	0.052	0.096	0.185	0.231	0.222	0.224	0.190	0.180
G	0.001	0.003	0.029	0.102	0.116	0.195	0.265	0.268	0.272	0.328
T	0.444	0.540	0.463	0.379	0.356	0.280	0.240	0.218	0.179	0.187
S	0.988	0.864	0.599	0.284	0.129	0.018	0.005	0.010	0.059	0.052

**Discussion:** To prevent the *alignment scores* being 0.0, a small constant  $k = 0.01$  has been added prior to calculating the probabilities for each cell in PPM. Based on the *PPM*, consensus sequence was identified as ATTATAAAAG and the consensus score as -9.77. It can be seen that *Guanine* (G) has been identified as one of the bases in the consensus sequence that indicates possible mutations of G or C is also possible within *Pribnow Box* regions.

## QUESTION 02

a) Obtaining the entropy measures for each position and using suitable entropy measure eliminating the redundant position of the obtained PPM

**Results:** Using the following equation, the entropy values for the 10 positions were calculated using  $p_{0,N} = 0.25$  for each base. (as shown in Table 1)

$$I_j = \frac{1}{\ln(2)} \sum_N p_{j,N} \ln \left( \frac{p_{j,N}}{p_{0,N}} \right)$$

$I_j$  – entropy of the column  $j$   
 $p_{j,N}$  – probability of a base at position  $j$  ( $N \in \{A, C, G, T\}$ )  
 $p_{0,N}$  – initial probability of a base at position  $j$  ( $N \in \{A, C, G, T\}$ )

To eliminate the redundant positions of the PPM obtained for the Pribnow Box, *entropy threshold* = 0.02 was selected. Based on the column entropies shown in the Table 1, the columns 6,7,8 are identified as redundant as they contain lesser certainty compared to the threshold selected.

Table 2: Reduced PPM after removing redundant positions

Base	Position						
	1	2	3	4	5	6 (9)	7 (10)
A	0.554	0.438	0.456	0.423	0.343	0.359	0.305
C	0.001	0.019	0.052	0.096	0.185	0.190	0.180
G	0.001	0.003	0.029	0.102	0.116	0.272	0.328
T	0.444	0.540	0.463	0.379	0.356	0.179	0.187

**Discussion:** Using the given entropy equation shown above, entropy values for each position were calculated. Based on Table 1 we can see that column 6,7,8, the probabilities of occurrence of one of the four bases seems equi-probable. This in turns reduces the entropy value calculated as it denotes a measure of certainty. Hence based on the figures and prior knowledge on PPM reduction based on entropy values for Pribnow Box search, the entropy threshold was set to 0.02 and obtained Table 2. Based on the reduced PPM, the consensus sequence was identified as ATTATAG and consensus score was -6.01. Again we can see that *Guanine* (G) at the last position is having a considerable probability and an entropy value meaning it is likely to locate mutations of G or C for the Pribnow Bow within this genome.

### QUESTION 03

#### a) Statistical alignment of sequences in *test set* with the initial PPM

**Results:** The remaining 3154 sequences were statistically aligned with the initial PPM shown in Table 1 using a moving window of size of 10. For a given sequence (of 25 bases long), 16 windows were aligned against the Initial PPM and from the resulted array of 16 alignment scores subtracted the consensus score for the initial PPM to obtain the relative scores. If at least one of the 16 relative scores were less than or equal to the threshold for which it was tested, the particular test sequence was considered to have Pribnow Box promoter. The promoter search was done for 5 threshold values (i.e. -1, -2, -3, -4, -5) and the results were as follows.

Table 3: Results of Statistical Alignment for Pribnow Box search using Initial PPM

	Threshold				
	-1	-2	-3	-4	-5
# Genes positive for Pribnow Box	998	2179	2757	3004	3092
% (out of 3154)	31.6%	69.1%	87.4%	95.2%	98.0%

#### b) Statistical alignment of sequences in *test set* with the reduced PPM

**Results:** Similar to part (a), the 3154 sequences were statistically aligned with the *reduced PPM* shown in Table 2. Instead of using a window size of 10, a window size of 7 was used that resulted in 19 alignment scores from which 19 relative scores were obtained by subtracting consensus score for reduced PPM. If at least one of the 19 relative scores were less than or equal to the threshold for which it was tested, the particular test sequence was considered to have Pribnow Box Promoter.

Table 4: Results of Statistical Alignment for Pribnow Box search using Reduced PPM

	Threshold				
	-1	-2	-3	-4	-5
# Genes positive for Pribnow Box	1519	2439	2875	3040	3115
% (out of 3154)	48.2%	77.3%	91.2%	96.4%	98.8%

#### c) Comparison between the Pribnow Box search using Initial PPM and Reduced PPM

		-1	-2	-3	-4	-5
With Initial PPM	# Genes Negative for Pribnow Box	2156	975	397	150	62
	% (out of 3376)	68.4%	30.9%	12.6%	4.8%	2.0%
With Reduced PPM	# Genes Negative for Pribnow Box	1635	715	279	114	39
	% (out of 3376)	51.8%	22.7%	8.9%	3.6%	1.2%

**Discussion:** By comparing the success rates for statistical alignment methods, alignment with *reduced PPM* shows better results compared to the alignment with initial PPM. Clear distinction can be seen for the threshold vales from -1 to -3 where the alignment with reduced PPM misses only 51.8%, 22.7% and 8.9% respectively as opposed to the 68.4%, 30.9% and 12.6% for the alignment with initial PPM. When the threshold w.r.t. to consensus scores increases, more than 95% of the sequences have been aligned for both the methods. Due to the reduction of columns having high uncertainty compared to others, alignment with reduced PPM has decreased the complexity of the calculations and improved the alignments by discarding uncommon mutations of *Guanine (G)* or *Cytosine (C)*.

As discussed earlier, since both the consensus scores contain *Guanine (G)* at the end position and that it has a considerate probability of 0.328, the success rates can be affected by the distribution differences of the test sequences and the train sequences selected. The results may improve if the selected sequences are shuffled prior to splitting them to train and test instead of selecting the first 1000 sequences.

## QUESTION 04

- a) Detection of possible mutated Pribnow Box promoters using Non-Intact Query Local Search (NIQLS) for genes that returned negative of the presence of Pribnow Box from both alignment methods

**Results:** It is highly likely that the Pribnow Box promoter to have mutations including insertions and point mutations due to substitution. Due to these possible mutations, the sequences may not have aligned statistically with either the initial PPM or reduced PPM or with both of them. In this section, we use a NIQLS with *W-search* to detect promoters with insertions possible. The resultant alignments are further analyzed for point mutations specifically for point mutations of *G* or *C*. The analysis was done for all the threshold values in terms of insertions and point mutations. For convenience, the contingency tables between number of insertion and number of point mutations at the aligned positions of the promoter for threshold -1 is shown as follows. For NIQLS, scores were assigned as; *match* = 3, *mismatch* = -3, *gap penalty* = 2.

Table 5: Contingency table for mutated Pribnow Boxes for threshold -1 related to Initial PPM

	Number of point mutations										Total		
Number of insertions		0		1		2		3		3>			
	0	347	0.17	92	0.05	86	0.04	34	0.02	10	0.01	569	0.25
	1	655	0.33	78	0.04	48	0.02	10	0.01	1	0.00	792	0.41
	2	402	0.20	47	0.02	14	0.01	1	0.00	0	0.00	464	0.24
	3	116	0.06	13	0.01	1	0.00	0	0.00	0	0.00	130	0.07
	3>	42	0.02	1	0.00	0	0.00	0	0.00	0	0.00	43	0.02
Total	1562	0.78	231	0.12	149	0.07	45	0.02	11	0.01	1998	1.00	

Table 6: Contingency table for mutated Pribnow Boxes for threshold -1 related to Reduced PPM

	Number of point mutations										Total		
Number of insertions		0		1		2		3		3>			
	0	118	0.08	65	0.04	84	0.05	34	0.02	10	0.01	311	0.20
	1	467	0.30	76	0.05	48	0.03	10	0.01	1	0.00	602	0.39
	2	390	0.25	47	0.03	14	0.01	1	0.00	0	0.00	452	0.29
	3	116	0.08	13	0.01	1	0.00	0	0.00	0	0.00	130	0.08
	3>	42	0.03	1	0.00	0	0.00	0	0.00	0	0.00	43	0.03
Total	1133	0.74	202	0.13	147	0.10	45	0.03	11	0.01	1538	1.00	

Considering a maximum of 2 insertions and 2-point mutations (i.e. up to 4 maximum mutations), the detection rates for mutated Pribnow Box promoters were calculated as follows for all the thresholds.

Table 7: Summary of the mutated Pribnow Box detection for both the cases

		-1	-2	-3	-4	-5
With Initial PPM	Total Number of Detections (% of negative genes from Question 3)	1998 (92.7)	921 (94.5)	369 (92.9)	133 (88.7)	51 (82.3)
	# Genes with possible mutated Pribnow Boxes	1769	692	196	40	14
	% mutations from Total Number of Detections	88.54	75.14	53.12	30.08	27.45
With Reduced PPM	Total Number of Detections (% of negative genes from Question 3)	1538 (94.1)	678 (94.8)	256 (91.8)	99 (86.8)	32 (82.1)
	# Genes with possible mutated Pribnow Boxes	1309	462	108	28	7
	% mutations from Total Number of Detections	85.11	68.14	42.19	28.28	21.88

**Discussion:** When searching for the mutated Pribnow boxes, we searched for mainly point mutations and insertions to the promoter. Since, Pribnow Boxes are likely to have mutations of *A* or *T* without compromising its functionality, we used *W-search* with NIQLS to detect possible Pribnow Boxes with insertions present. Based on the analysis shown in Table 5 and Table 6 for threshold -1, we can see that

majority of the aligned sequences have either 1 or 2 insertions present. Since we did the alignment using *W-search*, these insertions represent either an insertion of *Guanine(G)* or an insertions of *Cytosine(C)*. Further it shows that number of insertions in aligned sequences rarely exceed 3. (with 2% and 3% for the two cases)

Further by analyzing the percentages of point mutations of *G* or *C*, probability of having such mutations is below 25% in both cases. These point mutations hinder the functionality of Pribnow Box promoter as it changes the number of H-bonds present in the promoter. As can be seen, the point mutations more than 3 is highly unlikely having a 1% of total detected genes for both cases.

Based on the Table 7, out of the negatively detected genes in Question 03, more than 90% have been aligned using NIQLS for first three thresholds. Based on the criteria imposed that only a maximum of 4 mutations possible, we have selected only intersection of the first 3 rows and first 3 columns from the respective contingency tables. It can be seen that for the first three thresholds more than 50% out of the total aligned genes have been detected as positive for having a mutated Pribnow Box in both the cases. With a threshold of -4 and -5, it is unlikely that statistical alignment missed many of the possible Pribnow boxes intact even with point mutations present hence they have a lower percentile in both the cases.

## QUESTION 05

a) Standard IQLS to locate *Sigma Binding site (TTGACA)* within upstream position from 30 to 50

**Results:** As described in *Question 01*, after following similar procedure and performing Standard IQLS, within a search region of upstream positions from 30 to 50 we selected the genes with TTGACA Box present.

# seq with mutated <i>Methionine</i>	633	11.56%
# seq in sense strand with <i>TTGACA</i> box located	1976	70.93%
# seq in antisense strand with <i>TTGACA</i> box located	1931	71.73%
Total seq with <i>TTGACA</i> box located	3907	71.32%
# seq with <i>TTGACA</i> box unlocated	938	17.12%

**Discussion:** It can be seen that using the *Standard IQLS* more than 70% are identified with a *TTGACA* Box within the search region of 30 to 50 bases upstream. Since it is not possible to use *W-search* to detect the presence of the TTGACA box with the presence of *Guanine (G)* and *Cytosine(C)*, the only method suitable is by performing a traditional *Standard IQLS* that looks for *exact alignment*.

b) Obtaining *Position Probability Matrix (PPM)* and the entropy measures for each position

**Results:** Using the first 1000 sequences of the selected sequences, the PPM was obtained for 10 positions following the same procedure as in *part(b) of Question 01*. Then the respective column entropies were calculated and using an entropy threshold of 0.01 the redundant positions of 5, 6, 8, 9, 10 were eliminated (columns with blue outline) to derive the reduced PPM.

Table 8: PPM obtained from the first 1000 sequences for TTGACA box with Entropy measures

Base	Position									
	1	2	3	4	5	6	7	8	9	10
A	0.177	0.176	0.214	0.322	0.254	0.284	0.311	0.280	0.282	0.279
C	0.139	0.161	0.207	0.273	0.270	0.252	0.228	0.240	0.231	0.233
G	0.166	0.198	0.368	0.200	0.225	0.244	0.220	0.229	0.253	0.228
T	0.518	0.465	0.211	0.205	0.251	0.220	0.241	0.251	0.234	0.260
S	0.240	0.158	0.049	0.029	0.003	0.006	0.014	0.004	0.005	0.005



**Discussion:** Using the *initial PPM*, the consensus sequence identified was *TTGACAAAAA* with a consensus score of *-11.11*. As can be seen, the first 6 positions of the consensus directly represent the TTGACA Box. Using the *reduced PPM* the consensus sequence identified was *TTGAA* with a consensus score of *-4.72*. We can see that even position 5 and 6 corresponding to *C* and *A* have been eliminated due to the fact that they are more uncertain having all the bases equi-probable at those positions. Further the entropy values for Sigma site PPM is lower compared to Pribnow box PPM.

#### c) Statistical alignment of sequences in test set with the initial PPM

**Results:** Following the same procedure described in the *part (a) of Question 03* the statistical alignment to detect the TTGACA box was performed with Initial PPM and the following results were obtained.

Table 9: Results of Statistical Alignment for TTGACA Box search using Initial PPM

	Threshold				
	-1	-2	-3	-4	-5
# Genes Positive for TTGACA Box	333	2163	2884	2907	2907
% (out of 2907)	11.5%	74.4%	99.2%	100.0%	100.0%

#### d) Statistical alignment of sequences in test set with the reduced PPM

**Results:** Following the same procedure described in the *part (b) of Question 03* the statistical alignment to detect the TTGACA box was performed with Reduced PPM and the following results were obtained.

Table 10: Results of Statistical Alignment for TTGACA Box search using Reduced PPM

	Threshold				
	-1	-2	-3	-4	-5
# Genes Positive for TTGACA Box	1319	2794	2907	2907	2907
% (out of 2907)	45.4%	96.1%	100.0%	100.0%	100.0%

#### e) Comparison between the Pribnow Box search using Initial PPM and Reduced PPM

		-1	-2	-3	-4	-5
With Initial PPM	# Genes Negative for TTGACA Box	2574	744	23	0	0
	% (out of 2907)	88.5%	25.6%	0.8%	0.0%	0.0%
With Reduced PPM	# Genes Negative for TTGACA Box	1588	113	0	0	0
	% (out of 2907)	54.6%	3.9%	0.0%	0.0%	0.0%

**Discussion:** By comparing the success rates for statistical alignment methods, alignment with *reduced PPM* shows better results compared to the alignment with initial PPM. Clear distinction can be seen for the threshold values from -1 to -2 where the alignment with reduced PPM misses only *54.6%* and *3.9%* respectively as opposed to the *88.5%* and *25.6%* for the alignment with initial PPM. When the threshold w.r.t. to consensus scores increases, *all the sequences* have been aligned for both the methods. Due to the reduction of columns having high uncertainty compared to others, alignment with reduced PPM has decreased the complexity of the calculations and improved the alignments by discarding uncommon mutations specifically at position 5 and 6 in initial PPM as all the bases have equal probability.

#### f) Detection of possible mutated TTGACA Box promoters using Non-Intact Query Local Search (NIQLS) for genes that returned negative of the presence of TTGACA Box from both alignment methods

**Results:** As similar to the Pribnow Box, it is likely to have mutation in the *TTGACA Box promoter* in the genes. Having all the sequences positive for the presence of *TTGACA Box* for thresholds -3 to -5, it is a fair doubt that the undetected genes for thresholds -1 and -2 should be having mutated *TTGACA*

boxes. Hence as in *Question 04*, we look for the mutations (i.e. possible insertions and point mutations) using traditional NIQLS. In this we look at each position of the aligned sequences and check for possible point mutations after checking for the possible insertions. As in *Question 04*, for the convenience the results will be based on threshold -1. For NIQLS, scores were same as in *Question 04*.

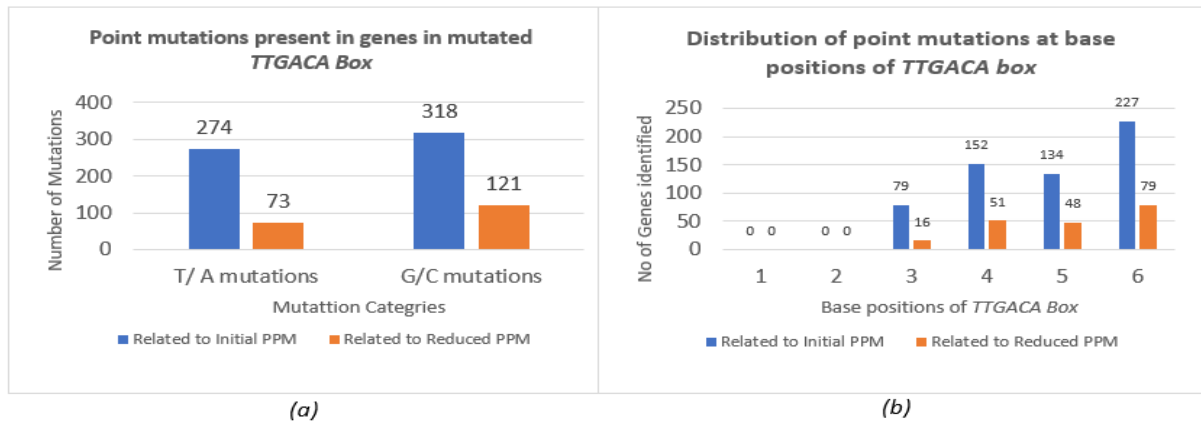


Figure 1: (a) Types of point mutations vs frequency of occurrence (b) distribution of mutations at each base position of TTGACA box (threshold -1)

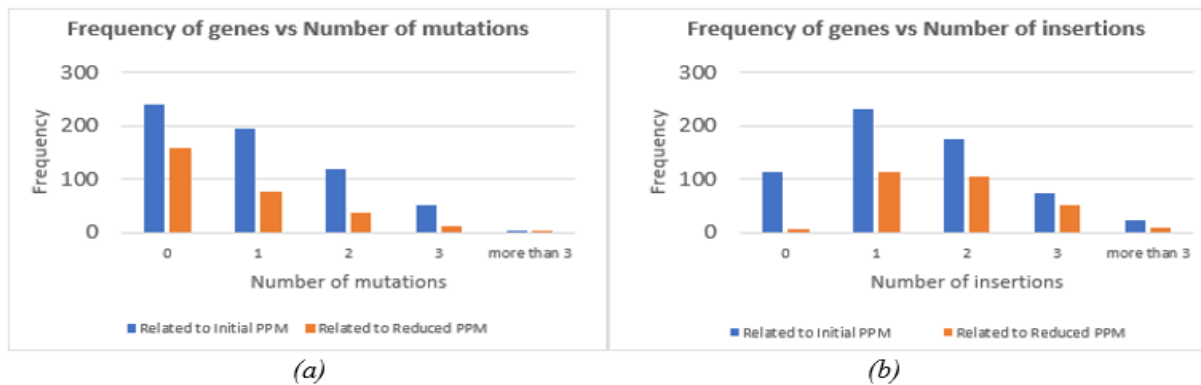


Figure 2: (a) The distribution of number of mutations detected (b) Distribution of number of insertions detected (for threshold -1)

Considering a maximum of 2 insertions and 2-point mutations (i.e. up to 4 maximum mutations), the detection rates for mutated TTGACA Box promoters were calculated as follows the thresholds -1 and -2.

Table 11: Summary of the mutated TTGACA Box detection for both the cases

		-1	-2
With Initial PPM	Total Number of Detections	608	82
	(% of negative genes from part e)	(23.6)	(11.0)
	# Genes with possible mutated TTGACA Boxes	468	65
With Reduced PPM	% mutations from Total Number of Detections	77.0	79.3
	Total Number of Detections	285	3
	(% of negative genes from Question 3)	(17.9)	(2.7)
With Reduced PPM	# Genes with possible mutated TTGACA Boxes	211	2
	% mutations from Total Number of Detections	74.0	66.7

**Discussion:** When searching for the mutated TTGACA boxes, we searched for mainly point mutations and insertions to the promoter. We used *traditional* NIQLS to detect possible TTGACA Boxes with insertions present as the sigma binding sites contains *G* and *C* that fails the effectiveness of the *W-search*. During the search, for both the statistical alignment scenarios, it was found that the returned aligned x-sequences had gaps when aligned with TTGACA box which suggests deletion of some bases



from possible sigma binding sites of the genes. Since having gaps in x-sequences during *Local search* is not common, these sequences were considered as genes without *TTGACA* boxes. Due to the high number of deletions received (514 and 49 for -1, -2 of Initial PPM) and the partial alignments received from NIQLS, we only get a small portion of sequences that detects mutated *TTGACA* box promoters.

Based on Figure 1 (a), we can see that out of all the possible mutations for threshold -1, bases tend to mutate for *G* or *C* rather than *T* or *A* for both the scenarios. Further, based on Figure 1 (b), it is evident that the occurrence of point mutations at position 1 and 2 of the sigma binding sites (i.e. the bases of *T* and *T*) have no mutations detected and the peaks of the distribution have occurred at position 4 and 6, both of which corresponds to the base *Adenine(A)* for both the scenarios. Further based on Figure 2 (a), the same pattern for the number of point mutations present can be observed as in Question 04, making a higher number of point mutations highly unlikely. Further we can see that based on the number of insertions present during the NIQLS alignment, most of the sequences appear to have 1 or 2 insertions.

As shown in Table 11, based on the criteria of selecting the successful detection of mutated *TTGACA* box, even though the selected proportion of genes from *part e* is considerably low (23.6% and 17.9% for the two scenarios with threshold of -1), more than 70% of the sequences have been identified successfully for the presence of mutated *TTGACA* boxes for both the scenarios.

## QUESTION 06

**Results:** To maintain the uniformity of results, elimination of redundant positions was based on entropy threshold of 0.02 to obtain Reduced PPM to detect Pribnow Box using statistical aligning.

Accession Code	Total test seq.s	Threshold									
		-1		-2		-3		-4		-5	
		#	%	#	%	#	%	#	%	#	%
NZ_CP027472.1	3044	1333	43.8	2315	76.1	2756	90.5	2944	96.7	3005	98.7
NZ_CP027352.1	3067	1464	47.7	2333	76.1	2766	90.2	2961	96.5	3020	98.5
NZ_CP027577.1	3081	1209	39.2	2275	73.8	2755	89.4	2961	96.1	3034	98.5
NZ_CP027387.1	3087	1394	45.2	2338	75.7	2772	89.8	2988	96.8	3046	98.7
NZ_CP027442.1	3081	1411	45.8	2336	75.8	2750	89.3	2980	96.7	3031	98.4
NZ_CP047378.1	3045	1384	45.5	2265	74.4	2722	89.4	2919	95.9	2994	98.3
NZ_CP028607.1	3035	1514	49.9	2382	78.5	2753	90.7	2923	96.3	2985	98.4
NZ_CP027338.1	2935	1133	38.6	2166	73.8	2640	89.9	2817	96.0	2885	98.3
NZ_CP037945.1	3033	1146	37.8	2210	72.9	2654	87.5	2909	95.9	2981	98.3
NZ_CP027555.1	2980	1425	47.8	2275	76.3	2711	91.0	2878	96.6	2942	98.7
NZ_CP032808.1	2971	1442	48.5	2321	78.1	2713	91.3	2871	96.6	2930	98.6
NZ_CP045827.1	2991	1385	46.3	2292	76.6	2722	91.0	288	96.4	2941	98.3
NZ_CP015853.1	2939	1217	41.4	2182	74.2	2646	90.0	2838	96.6	2901	98.7
NZ_CP015020.1	2948	1167	39.6	2122	72.0	2589	87.8	2794	94.8	2891	98.1
NZ_CP028592.1	2973	1483	49.9	2355	79.2	2723	91.6	2869	96.5	2928	98.5
NZ_CP032803.1	2964	1480	49.9	2332	78.7	2711	91.5	2864	96.6	2918	98.4
NZ_CP045975.1	2902	1663	57.3	2393	82.5	2702	93.1	2822	97.2	2870	98.9
NZ_CP032795.1	2973	1495	50.3	2363	79.5	2733	91.9	2878	96.8	2933	98.7
NZ_CP034806.1	2941	729	24.8	1927	65.5	2514	85.5	2763	93.9	2879	97.9
NZ_CP040305.1	2914	1412	48.5	2330	80.0	2677	91.9	2823	96.9	288	98.8
NZ_AP018808.1	3154	1519	48.2	2439	77.3	2875	91.2	3040	96.4	3115	98.8

<sup>i</sup> [https://www.ncbi.nlm.nih.gov/genome/?term=NZ\\_AP018808.1](https://www.ncbi.nlm.nih.gov/genome/?term=NZ_AP018808.1)