

Machine Learning in Computational Biology (Fall 2023)

Assignment #2

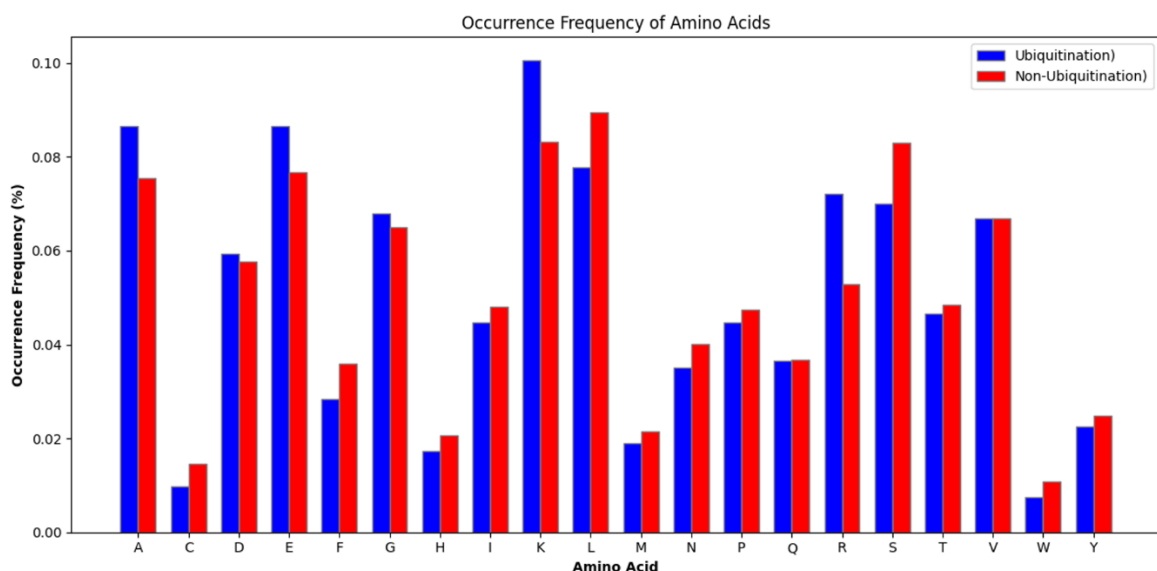
Deadline: 23:59:59 17th November, 2023 (Delayed submission is not allowed for any reason)

Purpose: to enhance the learning outcomes for the topics in “**Features Encoding and Investigation**” and “**Protein Sequence Analysis**”.

After the removal of homologous sequences in both positive and negative dataset using CD-Hit, in order to carry out a binary classification between ubiquitination and non-ubiquitination sites, please accomplish the tasks described as follows.

1. [Amino acid composition] (20 points)

Amino acid composition (AAC) is a common method used to transform protein sequences into 20-dimensional numeric vectors. As you learned from Chapter 2 in this class, please calculate AAC for each sequences in both positive and negative datasets; then, the comparison of AAC between positive (ubiquitination) and negative data (non-ubiquitination) can be displayed as the histogram shown below.



Please provide the histogram plot of AAC based on your training dataset using CD-Hit with 50% similarity.

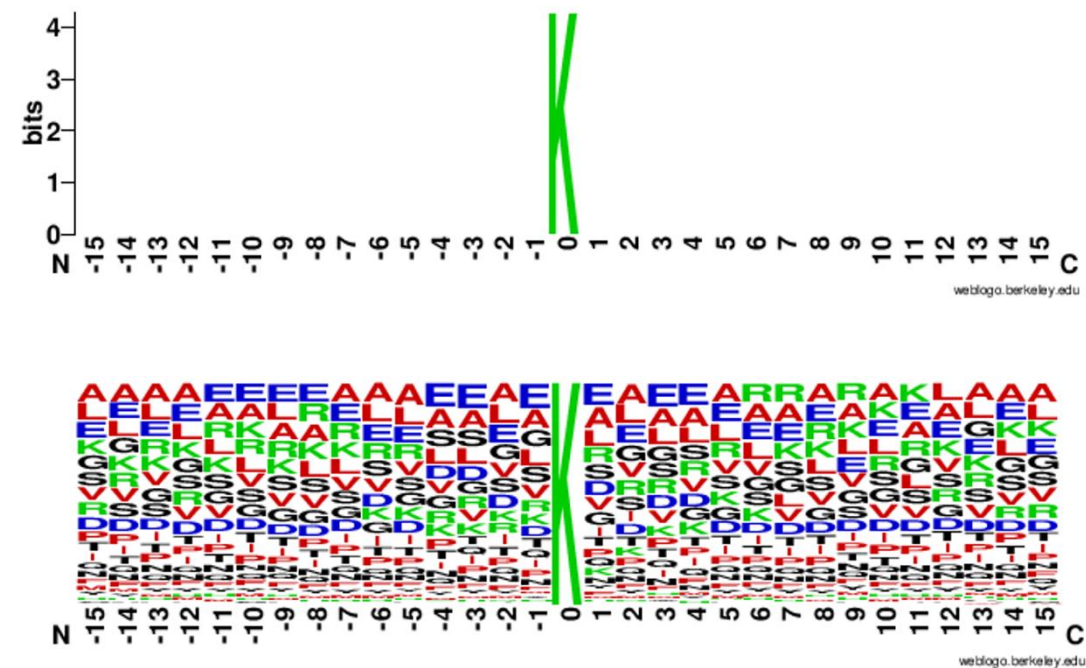
50 % CD-HIT has uploaded

2. [Sequence log and TwoSampleLogo] (20 points)

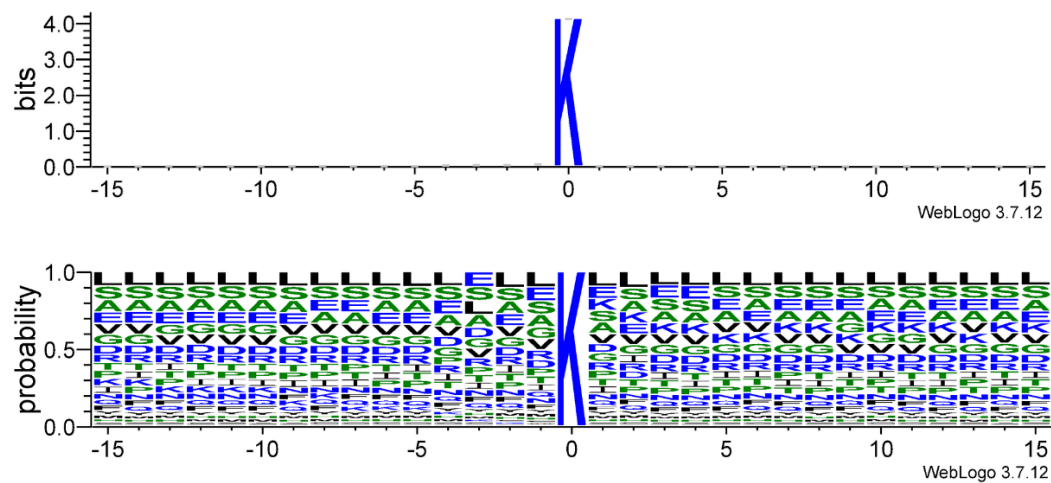
2.1 In order to observe the position-specific AAC for protein ubiquitination sites, the **WebLogo** tool can be utilized to create **frequency** and **entropy plots** of sequence logos for visualizing the potential amino acid motifs surrounding the modification sites (centered as

position 0). Please provide the **frequency** and **entropy** plots of sequence logos on both positive and negative datasets.

Positive Dataset



Negative Dataset



[illegible][illegible]

3.2 Please compare the PWM of positive dataset to that of negative dataset to identify the significantly differential represented amino acids in each position.

- Difference 比較：將兩個 PWM table 互減 (Positive - Negative)

Row_Max	Column	Max Value	Row_Min	Min Value
K	-15	0.022897	S	-0.01771
K	-14	0.024052	S	-0.01807
K	-13	0.026741	S	-0.02005
K	-12	0.032707	S	-0.01682
K	-11	0.029583	L	-0.01427
K	-10	0.037498	S	-0.01974
K	-9	0.037835	S	-0.01423
K	-8	0.04471	S	-0.01692
K	-7	0.041315	S	-0.02202
K	-6	0.03085	G	-0.01175
K	-5	0.036013	S	-0.01754
K	-4	0.037397	L	-0.0134
K	-3	0.036684	L	-0.01136
K	-2	0.041962	S	-0.01352
K	-1	0.050426	L	-0.01906
A	0	0	A	0
A	1	0.026966	K	-0.03428
R	2	0.021855	K	-0.02619
E	3	0.022694	K	-0.0181
R	4	0.020598	L	-0.01266
R	5	0.023477	S	-0.0115
R	6	0.036834	S	-0.01394
R	7	0.037188	L	-0.02181
R	8	0.030564	L	-0.01579
R	9	0.032221	S	-0.01894
R	10	0.022419	L	-0.01288
R	11	0.021512	L	-0.01973
R	12	0.009148	S	-0.01527
R	13	0.018235	S	-0.01512
A	14	0.015937	L	-0.01134
R	15	0.018409	S	-0.01613

- 取 log2 的方法

$$PWM_{ij} = \log_2 \left(\frac{p_{ij}}{p_i} \right),$$

Row_Max	Column	Max Value	Row_Min	Min Value
K	-15	0.554307	S	-0.3264
K	-14	0.567166	W	-0.4545
K	-13	0.642597	S	-0.36642
K	-12	0.770633	F	-0.45847
K	-11	0.71863	C	-0.55181
K	-10	0.899197	H	-0.46987
K	-9	0.982371	C	-0.46027
K	-8	1.094827	C	-1.06241
K	-7	1.096252	C	-1.28492
K	-6	0.972668	C	-1.22821
K	-5	1.172893	W	-0.92988
K	-4	1.35903	W	-1.17046
K	-3	1.543166	C	-1.15668
K	-2	1.992372	W	-1.39236
K	-1	2.893634	W	-1.17781
K	0	0	K	0
R	1	0.47239	W	-0.92221
R	2	0.498553	C	-0.76261
R	3	0.461723	C	-1.22717
R	4	0.478413	C	-1.40441
R	5	0.522913	W	-0.68158
R	6	0.769401	C	-0.599
R	7	0.769983	W	-1.15107
R	8	0.646003	C	-0.67005
R	9	0.682533	S	-0.37348
R	10	0.519064	C	-0.73307
R	11	0.478495	C	-0.50108
R	12	0.21838	W	-0.35538
R	13	0.421493	H	-0.39687
A	14	0.282149	F	-0.34521
R	15	0.442028	S	-0.29173

3.3 Can we use the PWM to be a predictive model for the prediction of ubiquitination sites?
Explain why or why not?

可以拿去預測，只是結果好壞不可知。

因為透過兩個方法去觀察 positive data 和 negative data 的差異性

1. Difference
2. Log2 Transformation

發現，確實在特定位點會有胺基酸特異性 (就以這筆資料來說)。

比如，做完 log2 的 table 中，位置 -1 的氨基酸位點，最大和最小值是整筆資料中相差較大的。在這個位點，可以 Positive 的資料是 K 比較重要，而 Negative 則是 W 貢獻度比較大。但如果有一個統計上的顯著性衡量標準去評判何謂顯著會更好。