

Ejercicios2_algoritmos

Tania Gonzalo y Daniel Parra

WORKSHEET

Ejercicio 1 (4 puntos)

En este ejercicio probarás el algoritmo Needleman-Wunsch en una secuencia corta de partes de hemoglobina (código PDB 1AOW) y mioglobina 1 (código PDB 1AZI). Aquí alineará la secuencia HGSAQVKGHG con la secuencia KTEAEMKASEDLKKHGT.

Las dos secuencias están dispuestas en una matriz en la Tabla 1. Las secuencias comienzan en la esquina superior derecha, y las penalizaciones por desfase inicial se enumeran en cada posición inicial de desfase. La penalización por desfase se considera -8. Las puntuaciones de similitud $S_{i,j}$ procedentes de la búsqueda de coincidencias proceden de la tabla BLOSUM40.

Para resolver este ejercicio, hemos aplicado el algoritmo de Needleman-Wunsch para alinear dos secuencias de proteínas que nos indican: una parte de la hemoglobina (HGSAQVKGHG) y una parte de la mioglobina (KTEAEMKASEDLKKHGT). Este algoritmo de alineamiento global nos permite encontrar la mejor correspondencia entre las secuencias considerando inserciones, eliminaciones y sustituciones de aminoácidos.

Para llevar a cabo la alineación, seguimos los siguientes pasos:

1. **Inicialización de la matriz:** Se construyó una matriz en la que las secuencias se dispusieron en los ejes horizontal y vertical. Se asignaron penalizaciones de apertura de brecha (-8) en la primera fila y la primera columna.
2. **Cálculo de los valores en la matriz:** Cada celda se llenó utilizando la fórmula de Needleman-Wunsch:

$$F(i, j) = \max \begin{cases} F(i-1, j-1) + S(i, j) & \text{(match o mismatch según la matriz BLOSUM40)} \\ F(i-1, j) + \text{gap penalty} & \text{(inserción)} \\ F(i, j-1) + \text{gap penalty} & \text{(eliminación)} \end{cases}$$

Donde $S(i, j)$ es el puntaje de sustitución obtenido de la matriz BLOSUM40.

La tabla obtenida se muestra en Figura 1 donde aparecen todas las celdas rellenas además de las flechas indicativas de la procedencia de cada uno de los números para posteriormente poder llevar a cabo la reconstrucción del alineamiento.

3. **Obtención del score final:** Una vez completada la matriz, el score final de alineación se encuentra en la casilla inferior derecha (Figura 1). Como podemos observar, en nuestro caso el score obtenido es de -21.
4. **Retroceso para encontrar las alineaciones óptimas:** Siguiendo las flechas que indicaban la mejor elección en cada celda, realizamos un recorrido inverso para reconstruir los alineamientos óptimos. Encontramos dos posibles caminos (alternativas), lo que dio lugar a dos alineaciones diferentes, ambos mostrados en la parte inferior de la Figura 1.

Este procedimiento nos permitió obtener una alineación óptima de las secuencias, teniendo en cuenta tanto coincidencias como penalizaciones por inserciones y eliminaciones.

Algoritmos

Ejercicio 1

	H	G	S	A	Q	V	K	G	H	G	
Q	0	-8	-16	-24	-32	-40	-48	-56	-64	-72	-80
K	-8	-1	-9	-16	-24	-31	-39	-42	-50	-58	-66
T	-16	-9	-3	-7	-15	-23	-30	-38	-44	-52	-60
E	-24	-16	-11	-3	-8	-13	-21	-29	-37	-44	-52
A	-32	-24	-15	-10	-6	-13	-21	-28	-36	-43	-50
E	-40	-32	-23	-15	-6	-4	-12	-20	-28	-36	-43
M	-48	-39	-31	-23	-14	-4	-3	-11	-19	-27	-34
K	-56	-47	-39	-31	-22	-12	-3	-11	-19	-27	-34
A	-64	-55	-46	-38	-26	-20	-11	-3	-11	-19	-27
S	-72	-63	-54	-41	-34	-25	-19	-5	-4	-11	-19
E	-80	-71	-62	-49	-42	-32	-27	-13	-4	-11	-19
D	-88	-79	-70	-57	-50	-40	-35	-21	-12	-4	-11
L	-96	-87	-78	-65	-58	-48	-38	-29	-20	-12	-4
K	-104	-95	-86	-73	-66	-56	-46	-32	-28	-20	-14
K	-112	-103	-94	-81	-74	-64	-54	-40	-34	-28	-22
H	-120	-111	-102	-89	-82	-72	-62	-48	-42	-34	-29
G	-128	-119	-110	-97	-90	-80	-70	-56	-50	-42	-37
T	-136	-127	-118	-105	-98	-88	-78	-64	-58	-50	-45

Obtenemos dos alineamientos con puntuación -21

* Alineamiento 1

-	-	H	G	S	-	-	A	-	Q	-	V	K	G	H	G	-
K	T	E	A	E	M	K	A	S	E	D	L	K	K	H	G	T

* Alineamiento 2

-	-	H	G	-	-	S	A	-	Q	-	V	K	G	H	G	-
K	T	E	A	E	M	K	A	S	E	D	L	K	K	H	G	T

Figura 1: Solución manual del ejercicio 1. En la parte superior se muestra la matriz obtenida tras la aplicación del algoritmo de Needleman-Wunsch. En la parte inferior se muestran los dos alineamientos óptimos obtenidos.

Ejercicio 2 (6 puntos)

Dado el conjunto de secuencias múltiples:

- *S1: PPGVKSDCAS*
- *S2: PADGVKDCAS*
- *S3: PPDGKSDS*
- *S4: GADGKDCCS*
- *S5: GADGKDCAS*

Utilice el popular método de alineación progresiva para alinear globalmente el conjunto anterior de secuencias. Genere el árbol guía por unión de vecinos. Compare su resultado (alineamiento) con el de Clustal-Omega.

Con el alineamiento final representa el logo. Para este proposito los caracteres nulos o gap son ignorados y no cuentan para el número de observaciones de una columna.

Para la resolución de este ejercicio nos piden llevar a cabo un alineación global de secuencias mediante alineación progresiva, donde el conjunto de secuencias con las que se trabaja es:

Secuencia	Cadena
S1	PPGVKSDCAS
S2	PADGVKDCAS
S3	PPDGKSDS
S4	GADGKDCCS
S5	GADGKDCAS

Metodología

Para este ejercicio se utilizó el método de alineación progresiva para realizar el alineamiento global de las secuencias. Se llevaron a cabo los siguientes pasos:

- **Alineamientos por pares:** Se utilizaron las herramientas EMBOSS Needle y Pairwise Sequence Alignment (PSA) para generar alineamientos dos a dos entre las secuencias.
- **Construcción de la matriz de distancias:** A partir de los alineamientos, se calculó una matriz de distancias basada en los valores de similitud entre las secuencias de forma manual.
- **Generación del árbol guía de forma manual:** Se empleó el método de unión de vecinos (*Neighbor Joining*) para construir el árbol guía a partir de la matriz de distancias.

- **Alineamiento progresivo manual:** Utilizando el árbol guía, se realizó el alineamiento progresivo de las secuencias.
- **Comparación con Clustal Omega:** Se comparó el alineamiento obtenido con el generado por Clustal Omega (1.2.4).
- **Generación del logo de secuencias:** Finalmente, se representó el alineamiento en formato de logo ignorando los caracteres nulos o gaps mediante el uso de las siguientes librerías en Python para el análisis y visualización de secuencias biológicas: **Biopython** (módulo AlignIO) para la manipulación de alineamientos de secuencias, **Logomaker** para la generación de logotipos de secuencia, **Matplotlib** para la visualización gráfica y **NumPy** para el procesamiento numérico.

1. Alineamientos por pares

Para los alineamientos por pares se utilizaron los siguientes parámetros:

- **Matriz de sustitución:** BLOSUM40
- **Penalización por gap:** -10 (aproximado a -8, usado en el ejercicio anterior)

Estos también se muestran en la figura Figura 2.

OUTPUT FORMAT ⓘ pair ▼	MATRIX ⓘ BLOSUM40 ▼	GAP OPEN ⓘ 10 ▼
GAP EXTEND ⓘ 10.0 ▼	END GAP ⓘ false ▼	END GAP OPEN ⓘ 10 ▼
END GAP EXTEND ⓘ 0.5 ▼		

Figura 2: Parámetros para el alineamiento dos a dos.

Se calcularon los alineamientos y sus respectivos scores, definidos como la relación entre el número de coincidencias y la cantidad total de residuos (excluyendo gaps). Los resultados obtenidos fueron los siguientes:

Secuencias comparadas	Longitud 1	Longitud 2	Coincidencias	Score
S1 vs S2	10	10	8	8/10
S1 vs S3	10	8	5	5/9
S1 vs S4	10	9	4	8/19

Secuencias comparadas	Longitud 1	Longitud 2	Coincidencias	Score
S1 vs S5	10	9	5	10/19
S2 vs S3	10	8	4	4/9
S2 vs S4	10	9	7	14/19
S2 vs S5	10	9	8	16/19
S3 vs S4	8	9	4	8/17
S3 vs S5	8	9	3	6/17
S4 vs S5	9	9	8	8/9

Los resultados de los alineamientos obtenidos los podemos observar a continuación:

```

S1          1 P-PGVKSDCAS      10
             | .||| |||
S2          1 PADGVK-DCAS      10

```

Figura 3: Alineamiento dos a dos de S1 vs S2, con 10 aa cada secuencia y 8 coincidencias, tiene un score 8/10.

```

S1          1 PPGVKSDCAS      10
             ||..|||.
S3          1 PPDGKSDS--      8

```

Figura 4: Alineamiento dos a dos de S1 vs S3, con 10 aa en S1 y 8 aa en S3; y 5 coincidencias, tiene un score 5/9.

```

S1          1 PPGVK-SDCAS      10
             |.. .||.
S4          1 --GADGKDCCS      9

```

Figura 5: Alineamiento dos a dos de S1 vs S4, con 10 aa en S1 y 9 aa en S4; y 4 coincidencias, tiene un score 8/19.

S1	1 PPGVK-SDCAS	10
	.. .	
S5	1 --GADGKDCAS	9

Figura 6: Alineamiento dos a dos de S1 vs S5, con 10 aa en S1 y 9 aa en S5; y 5 coincidencias, tiene un score 10/19.

S2	1 PADGVKDCAS	10
	
S3	1 PPDGKSDS--	8

Figura 7: Alineamiento dos a dos de S2 vs S3, con 10 aa en S2 y 8 aa en S3; y 4 coincidencias, tiene un score 4/9.

S2	1 PADGVKDCAS	10
	. .	
S4	1 GADG-KDCCS	9

Figura 8: Alineamiento dos a dos de S2 vs S4, con 10 aa en S2 y 9 aa en S4; y 7 coincidencias, tiene un score 14/19.

S2	1 PADGVKDCAS	10
	.	
S5	1 GADG-KDCAS	9

Figura 9: Alineamiento dos a dos de S2 vs S5, con 10 aa en S2 y 9 aa en S5; y 8 coincidencias, tiene un score 16/19.

S3	1 PPDGKSDS--	8
	.. .	
S4	1 GADGK-DCCS	9

Figura 10: Alineamiento dos a dos de S3 vs S4, con 8 aa en S3 y 9 aa en S4; y 4 coincidencias, tiene un score 8/17.

S3	1 PPDGKSDS-	8
:	
S5	1 GADGKDCAS	9

Figura 11: Alineamiento dos a dos de S3 vs S5, con 8 aa en S3 y 9 aa en S5; y 3 coincidencias, tiene un score 6/17.

S4	1 GADGKDCCS	9
	.I	
S5	1 GADGKDCAS	9

Figura 12: Alineamiento dos a dos de S4 vs S5, con 9 aa cada secuencia y 8 coincidencias, tiene un score 8/9.

2. Construcción de la matriz de distancias

La matriz de distancias se construyó a partir del score de similitud, aplicando la transformación:

$$d = 1 - \text{score}$$

Obteniendo la siguiente matriz de distancias que se puede observar en la parte superior de Figura 13:

	S1	S2	S3	S4	S5
S1	0.00	-	-	-	-
S2	0.20	0.00	-	-	-
S3	0.44	0.56	0.00	-	-
S4	0.58	0.26	0.53	0.00	-
S5	0.47	0.16	0.65	0.11	0.00

3. Construcción del árbol guía de forma manual

A partir de la matriz de distancias, se generó el árbol guía utilizando el método de unión de vecinos (*Neighbor Joining*). El árbol guía obtenido muestra la relación evolutiva entre las secuencias y nos sirve como base para la alineación progresiva. El proceso de generación del árbol guía se hizo de forma manual y se muestra en Figura 13.

El árbol guía obtenido finalmente fue el mostrado en Figura 14.

Ejercicio 2

• Matriz de distancias tras los alineamientos Pánuire

S1	0				
S2	0,2	0			
S3	0,44	0,56	0		
S4	0,58	0,26	0,53	0	
S5	0,47	0,16	0,65	0,11	0
	S1	S2	S3	S4	S5

$$S2-S5 = 1 - \frac{16}{19} = 0,16$$

$$S3-S5 = 1 - \frac{6}{17} = 0,65$$

$$\left[\begin{array}{c} S4 \\ S5 \end{array} \right] A$$

$$A-S1 = \frac{0,58 + 0,47}{2} = 0,525$$

$$A-S2 = \frac{0,26 + 0,16}{2} = 0,21$$

$$A-S3 = \frac{0,53 + 0,65}{2} = 0,59$$

B	0		
S3	0,5	0	
A	0,3675	0,59	0
	B	S3	A

$$\left[\begin{array}{c} A \\ B \\ S3 \end{array} \right]$$

$$S1-S2 = 1 - \frac{8}{10} = 0,2$$

$$S1-S3 = 1 - \frac{5}{9} = 0,44$$

$$S1-S4 = 1 - \frac{8}{14} = 0,58$$

$$S1-S5 = 1 - \frac{10}{19} = 0,47$$

$$S2-S3 = 1 - \frac{4}{9} = 0,56$$

$$S2-S4 = 1 - \frac{14}{19} = 0,26$$

$$S3-S4 = 1 - \frac{8}{17} = 0,53$$

$$S4-S5 = 1 - \frac{8}{9} = 0,11$$

S1	0			
S2	0,2	0		
S3	0,44	0,56	0	
A	0,525	0,21	0,59	0
	S1	S2	S3	A

$$\left[\begin{array}{c} S1 \\ S2 \end{array} \right] B$$

$$B-S3 = \frac{0,44 + 0,56}{2} = 0,5$$

$$B-A = \frac{0,525 + 0,21}{2} = 0,3675$$

Figura 13: Construcción de matrices de distancias.

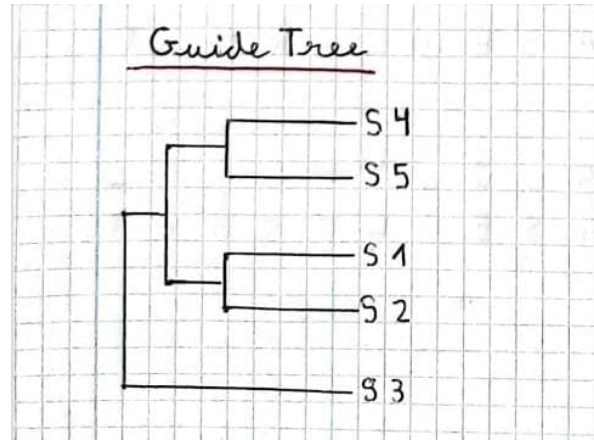


Figura 14: Guide Tree obtenido a partir de alineamientos 2 a 2 con parámetros similares al ejercicio 1.

4. Alineamiento progresivo manual

Utilizando el árbol guía, se llevó a cabo el alineamiento progresivo de las secuencias. Este alineamiento se obtuvo manualmente y posteriormente se comparará con el resultado de Clustal Omega.

El alineamiento final obtenido se muestra en Figura 15.

Alineamiento progresivo

A	G	A	D	G	K	D	C	C	S		G	A	D	G	-	K	-	D	C	C	S	
	G	A	D	G	K	D	C	A	S		G	A	D	G	-	K	-	D	C	A	S	
+										⇒												
B	P	-	P	G	V	K	S	D	C	A	S	P	A	D	G	V	K	-	D	C	A	S
	P	A	D	G	V	K	-	D	C	A	S	P	-	P	G	V	K	S	D	C	A	S
												+										
												S3	P	P	D	G	-	K	S	D	-	S

Alineamiento final

S4:	G	A	D	G	-	K	-	D	C	C	S
S5:	G	A	D	G	-	K	-	D	C	A	S
S2:	P	A	D	G	V	K	-	D	C	A	S
S1:	P	-	P	G	V	K	S	D	C	A	S
S3:	P	P	D	G	-	K	S	D	-	-	S

Figura 15: Alineamiento progresivo obtenido a partir de alineamientos 2 a 2 con parámetros similares al ejercicio 1.

Cálculo con parámetros por defecto

Además del cálculo utilizando BLOSUM40, se decidió repetir el proceso usando la matriz BLOSUM62 y los parámetros por defecto, para evaluar si el alineamiento final variaba en función de los parámetros. Los parámetros utilizados fueron los que se muestran en Figura 16.

OUTPUT FORMAT ⓘ	MATRIX ⓘ	GAP OPEN ⓘ
pair ▼	BLOSUM62 ▼	10 ▼
GAP EXTEND ⓘ	END GAP ⓘ	END GAP OPEN ⓘ
0.5 ▼	false ▼	10 ▼
END GAP EXTEND ⓘ		
0.5 ▼		

Figura 16: Parámetros para el alineamiento dos a dos por defecto.

Se realizaron los alineamientos por pares con estos parámetros y se generó una nueva matriz de distancias, un nuevo árbol guía y el alineamiento progresivo correspondiente.

Los scores obtenidos con estos parámetros fueron se muestran en la siguiente tabla:

Secuencias comparadas	Longitud 1	Longitud 2	Coincidencias	Score
S1 vs S2	10	10	7	7/10
S1 vs S3	10	8	5	5/9
S1 vs S4	10	9	4	8/19
S1 vs S5	10	9	5	10/19
S2 vs S3	10	8	4	4/9
S2 vs S4	10	9	7	14/19
S2 vs S5	10	9	8	16/19
S3 vs S4	8	9	3	6/17
S3 vs S5	8	9	3	6/17
S4 vs S5	9	9	8	8/9

Y los alineamientos obtenidos se muestran a continuación:

S1	1 -PPGVKSDCAS	10
	..	
S2	1 PADGVK-DCAS	10

Figura 17: Alineamiento dos a dos con parámetros por defecto, de S1 vs S2, con 10 aa cada secuencia y 7 coincidencias, tiene un score 7/10.

S1	1 PPGVKSDCAS	10
	.. .	
S3	1 PPDGKSDS--	8

Figura 18: Alineamiento dos a dos con parámetros por defecto, de S1 vs S3, con 10 aa en S1 y 8 aa en S3; y 5 coincidencias, tiene un score 5/9.

S1	1 PPGVK-SDCAS	10
	.. .	
S4	1 --GADGKDCCS	9

Figura 19: Alineamiento dos a dos con parámetros por defecto, de S1 vs S4, con 10 aa en S1 y 9 aa en S4; y 4 coincidencias, tiene un score 8/19.

S1	1 PPGVK-SDCAS	10
	.. .	
S5	1 --GADGKDCAS	9

Figura 20: Alineamiento dos a dos con parámetros por defecto, de S1 vs S5, con 10 aa en S1 y 9 aa en S5; y 5 coincidencias, tiene un score 10/19.

S2	1 PADGVKDCAS	10
	. . .	
S3	1 PPDGKSDS--	8

Figura 21: Alineamiento dos a dos con parámetros por defecto, de S2 vs S3, con 10 aa en S2 y 8 aa en S3; y 4 coincidencias, tiene un score 4/9.

S2	1 PADGVKDCAS	10
	. .	
S4	1 GADG-KDCCS	9

Figura 22: Alineamiento dos a dos con parámetros por defecto, de S2 vs S4, con 10 aa en S2 y 9 aa en S4; y 7 coincidencias, tiene un score 14/19.

S2	1 PADGVKDCAS	10
	.	
S5	1 GADG-KDCAS	9

Figura 23: Alineamiento dos a dos con parámetros por defecto, de S2 vs S5, con 10 aa en S2 y 9 aa en S5; y 8 coincidencias, tiene un score 16/19.

S3	1 PPDGKSDS-	8
	
S4	1 GADGKDCCS	9

Figura 24: Alineamiento dos a dos con parámetros por defecto, de S3 vs S4, con 8 aa en S3 y 9 aa en S4; y 3 coincidencias, tiene un score 6/17.

S3	1 PPDGKSDS-	8
 :	
S5	1 GADGKDCAS	9

Figura 25: Alineamiento dos a dos con parámetros por defecto, de S3 vs S5, con 8 aa en S3 y 9 aa en S5; y 3 coincidencias, tiene un score 6/17.

S4	1 GADGKDCCS	9
	.	
S5	1 GADGKDCAS	9

Figura 26: Alineamiento dos a dos con parámetros por defecto, de S4 vs S5, con 9 aa cada secuencia y 8 coincidencias, tiene un score 8/9.

Con ello, contruimos la matriz de distancias:

	S1	S2	S3	S4	S5
S1	0.00	-	-	-	-
S2	0.30	0.00	-	-	-
S3	0.44	0.56	0.00	-	-
S4	0.58	0.26	0.53	0.00	-
S5	0.47	0.16	0.65	0.11	0.00

Y se llevó a cabo el alineamiento progresivo de forma análoga al anterior (Figura 27)

Finalmente, construimos el guide tree:

Y el alineamiento progresivo se muestra en la figura Figura 29

Alineamiento múltiple

Para llevar a cabo el alineamiento múltiple de secuencias, se utilizaron los siguientes parámetros:

OUTPUT FORMAT ⓘ ClustalW with character counts ▼		DEALIGN INPUT ⓘ no ▼	
MBED-LIKE CLUSTERING GUIDE-TREE ⓘ no ▼	MBED-LIKE CLUSTERING ITERATION ⓘ yes ▼	COMBINED ITERATIONS ⓘ default(0) ▼	
MAX GUIDE TREE ⓘ default ▼	MAX HMM ITERATIONS ⓘ default ▼	ORDER ⓘ aligned ▼	DISTANCE MATRIX ⓘ yes ▼
OUTPUT GUIDE TREE ⓘ yes ▼			

Figura 30: Parámetros para el alineamiento múltiple.

Estos parámetros fueron configurados en Clustal Omega, una herramienta bioinformática para la alineación de secuencias.

Como resultado del proceso, se obtuvo el siguiente árbol guía, que representa la relación filogenética entre las secuencias alineadas (Figura 31).

Finalmente, se generó el alineamiento múltiple de las secuencias, donde se pueden observar las similitudes y diferencias entre ellas (Figura 32).

• Matriz de distancias tras los alineamientos Pairwise (BLOSUM 62):

S1	0				
S2	0,3	0			
S3	0,44	0,56	0		
S4	0,58	0,26	0,65	0	
S5	0,47	0,16	0,65	0,11	0
	S1	S2	S3	S4	S5

$$S2-S5 = 1 - \frac{16}{19} = 0,16$$

$$S3-S5 = 1 - \frac{6}{17} = 0,65$$

$$\left\{ \begin{matrix} S4 \\ S5 \end{matrix} \right\} A \rightarrow$$

$$A-S1 = \frac{0,58+0,43}{2} = 0,525$$

$$A-S2 = \frac{0,26+0,16}{2} = 0,21$$

$$A-S3 = \frac{0,65+0,65}{2} = 0,65$$

S1	0				
S2	0,3	0			
S3	0,44	0,56	0		
A	0,525	0,21	0,65	0	
	S1	S2	S3	A	

S1	0			
S3	0,44	0		
B	0,4125	0,605	0	
	S1	S3	B	

$$\left\{ \begin{matrix} A \\ S2 \end{matrix} \right\} B \leftarrow$$

$$B-S1 = \frac{0,3+0,525}{2} = 0,4125$$

$$B-S3 = \frac{0,56+0,65}{2} = 0,605$$

$$\left\{ \begin{matrix} B \\ S1 \\ S3 \end{matrix} \right\}$$

Figura 27: Construcción de matrices de distancias.

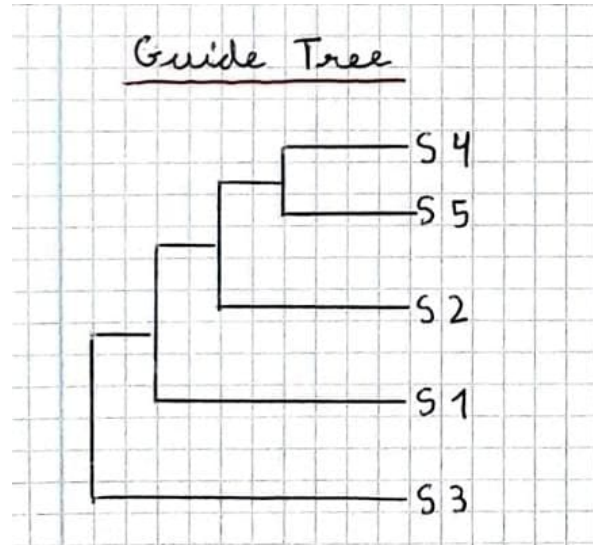


Figura 28: Guide Tree obtenido a partir de alineamientos 2 a 2

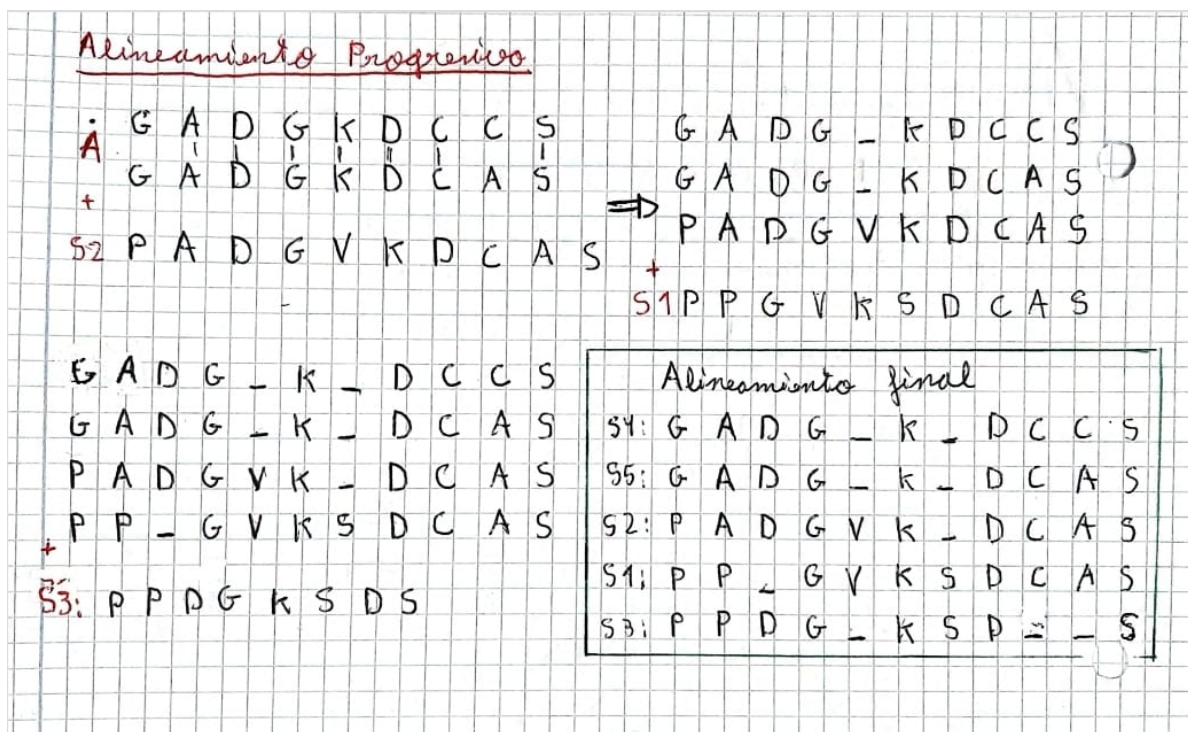


Figura 29: Alineamiento progresivo obtenido a partir de alineamientos 2 a 2 con matriz BLO-SUM62

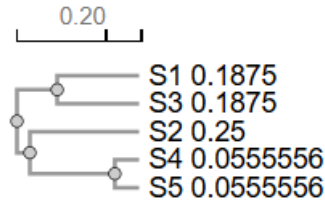


Figura 31: Guide tree del alineamiento múltiple por Clustal Omega.

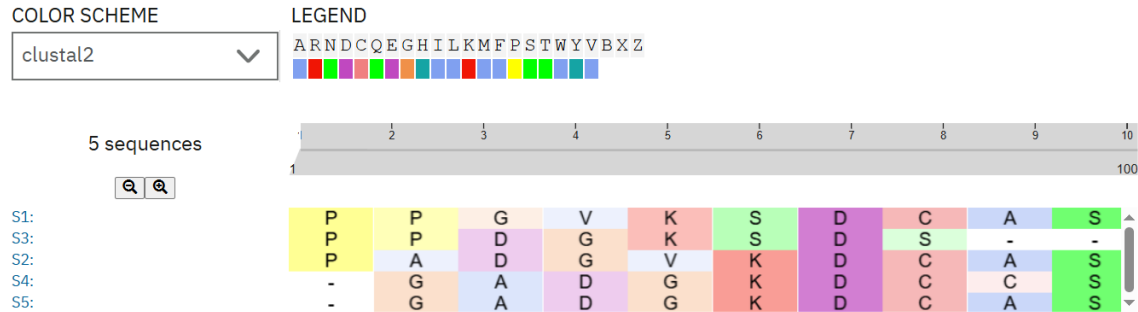


Figura 32: Alineamiento múltiple por Clustal Omega.

Este alineamiento nos permite analizar la conservación de regiones específicas, identificar mutaciones y establecer relaciones evolutivas entre las secuencias comparadas.

Resultados

Comparación alineamientos

1. Primer alineamiento (manual, con BLOSUM40)(Figura 15): Existe una brecha significativa en S3 (representada por guiones), lo que indica una menor similitud con las demás secuencias. Las secuencias S4 y S5 muestran una alta conservación en sus posiciones (similitud cercana en aminoácidos).
2. Segundo alineamiento (manual, con BLOSUM62) (Figura 29): Aparece una ligera variación en S1 (gap introducido en tras la segunda prolina (P) en lugar de tras la primera prolina (P)). El resto de las secuencias mantiene el mismo patrón que en el primer alineamiento. Esto es debido a la utilización de distintas matrices de sustitución donde la matriz BLOSUM62 favorece la detección de similitudes más estrechas en secuencias altamente relacionadas y también cómo la matriz penaliza los gaps o prioriza substituciones específicas.

3. Tercer alineamiento (Clustal Omega)(Figura 32): En este caso podemos ver que el alineamiento difiere bastante a los obtenidos de forma manual. Podemos observar menos gaps y, en caso de haberlos, estos se encuentran en los extremos (menor número de brechas y una mayor continuidad), es decir, presenta un mayor número de aminoácidos alineados, lo que puede ser el resultado de la automatización del proceso.

Con todo ello, podemos concluir que en todos los casos, S4 y S5 presentan una alta conservación, consolidando que están estrechamente relacionadas. La mayor diferencia se encuentra en la cadena S3, donde los métodos manuales marcan a S3 como más divergente (más gaps), pero Clustal Omega lo integra más en el alineamiento (mismatch en lugar de gaps).

Esto puede deberse a que Clustal Omega se basa en un algoritmo heurístico y, por tanto, no garantiza encontrar el alineamiento globalmente óptimo, teniendo dificultad frente a delecciones e inserciones, en las que incluye gaps. Además, debemos tener en cuenta la importancia de los parámetros introducidos para que esta se ajuste a la realidad. En el peor de los casos se utilizan valores predeterminados (como ocurre con Clustal Omega (1.2.4)), que pueden no ser ideales para todo el conjunto de datos. También hay que considerar que Clustal Omega usa el método HAlign, que es un algoritmo basado en modelos de Markov ocultos (HMM), para realizar alineamientos más precisos entre perfiles de secuencias.

Comparación de árboles guía

Como hemos mostrado anteriormente, se realizaron finalmente tres árboles diferentes (Figura 14, Figura 28, Figura 31):

1. Árbol basado en BLOSUM40 (manual, alineamiento progresivo y vecinos cercanos):
 - S4 y S5 son los más cercanos, agrupándose como un clado separado.
 - S1 y S2 son también los más cercanos, agrupándose como un clado separado.
 - Estos dos clados formados se unen en un antecesor común.
 - S3 es el más distante, uniéndose finalmente al resto.
2. Árbol basado en BLOSUM62 (manual, alineamiento progresivo y vecinos cercanos). La estructura difiere de la obtenida con la matriz BLOSUM40. En este caso:
 - S4 y S5 están más estrechamente relacionados.
 - S2 se incorpora al grupo S4-S5.
 - S1 se une después al conjunto principal.
 - S3 es el más distante, uniéndose finalmente al resto.

Esto refleja que cambiar la matriz de sustitución (BLOSUM40 vs. BLOSUM62) afecta a las relaciones globales en la construcción del árbol guía mediante este método.

3. Árbol generado por Clustal Omega: Los resultados en este caso se presentan en un formato numérico con distancias específicas:

- S4 y S5 siguen siendo los más cercanos, con una distancia de 0.055556 .
- S2 se sitúa más lejos de S4 y S5, con una distancia de 0.25 .
- S1 y S3 forman un subgrupo con una distancia de 0.1875 , diferente a los árboles manuales. Podemos observar que este método introduce diferencias en cómo se agrupan S1 y S3 en comparación con los anteriores.

Con todo ello, podemos concluir una relación S4-S5 consistente ya que en todos los métodos y matrices utilizadas, S4 y S5 son siempre los más cercanos, lo que indica una fuerte similitud entre estas secuencias. Además, hemos observado el impacto que tiene escoger la penalización por gaps y entre una u otra matriz BLOSUM al comparar los resultados entre BLOSUM40 y BLOSUM62 y observando que el árbol guía obtenido es diferente. Por último, ninguno de los árboles obtenidos manualmente coincide con el obtenido por el algoritmo Clustal Omega, que introduce distancias precisas e incluye la relación entre S1 y S3, obtenido por Clustal Omega (1.2.4).

Generación del logo a partir de los alineamientos obtenidos

Para representar el alineamiento final en formato de logo, se utilizó el lenguaje de programación Python y las librerías Biopython, Logosmaker, Matplotlib y NumPy.

En la construcción del logo, se ignoraron los gaps y se consideraron únicamente los residuos presentes en cada posición del alineamiento. Por ello, el número de aminoácido presentes es de 8 y, por tanto, el máximo valor que se puede obtener de cantidad de información (bits, representados en el eje y) es de 3.

1. Logo del alineamiento obtenido usando BLOSUM40

Cargamos las librerías necesarias, cargamos los documentos FASTA de los alineamientos progresivos y obtuvimos el logo correspondiente.

```
#Importación de librerías
from Bio import AlignIO
import logomaker
import matplotlib.pyplot as plt
import numpy as np

# Cargar y procesar el alineamiento
alignment = AlignIO.read("Alineamiento_final_BLOSUM40.fa", "fasta")
sequences = [str(record.seq) for record in alignment]

# Crear matriz de conteos (sin gaps)
counts_matrix = logomaker.alignment_to_matrix(sequences,\
    to_type='counts',\
```

```

    characters_to_ignore='-') # Excluir gaps para evitar sesgos

# Convertir a frecuencias relativas (ignorando gaps)
frequency_matrix = counts_matrix.div(counts_matrix.sum(axis=1), axis=0)

# Normalizar a frecuencias relativas y convertir a información (bits)
information_matrix = logomaker.transform_matrix(frequency_matrix, \
    from_type='probability', to_type='information')

# Generar el logo
logo = logomaker.Logo(
    information_matrix,
    color_scheme='chemistry', #Colores por propiedades químicas
    font_name='Arial',
    show_spines=True,
    stack_order='big_on_top' #Aminoácidos más frecuentes arriba
)

# Ajustar el estilo del gráfico
logo.ax.set_title("Logo de alineamiento progresivo con la BLOSUM40")
logo.ax.set_xlabel("Posición")
logo.ax.set_ylabel("Bits")
logo.ax.set_ylim(0, 3)

plt.tight_layout()
plt.show()

```



Podemos observar la representación gráfica del logo de alineamiento progresivo basado en la matriz de sustitución BLOSUM40. Este tipo de visualización permite evaluar la conservación y variabilidad de secuencias de aminoácidos en alineaciones múltiples, proporcionando información relevante sobre la evolución y estructura de las proteínas.

En el eje vertical aparecen representados los bits que señalan la cantidad de información en cada posición de la secuencia. Valores altos indican una mayor conservación de los aminoácidos

en esa posición. En el eje horizontal tenemos las posiciones de la secuencia donde podemos observar la ubicación específica de los aminoácidos en la alineación. En el logo en sí aparecen representadas las letras correspondientes a un aminoácido concreto, y su tamaño refleja la frecuencia con la que aparece en una posición determinada.

Cada color indica propiedades de los aminoácidos de forma que los aminoácidos azules son los básicos; los rojos los aminoácidos ácidos; los verdes polares no cargados; los morados, aromáticos; los verde oscuro, residuos especiales y los negro, hidrofóbicos.

Así, podemos observar que algunos aminoácidos como “G”, “V”, “K”, “S” y “D” tienen una alta conservación en posiciones clave, mientras que otros presentan más variabilidad.

2. Logo del alineamiento obtenido usando BLOSUM62

```
# Cargar y procesar el alineamiento
alignment = AlignIO.read("Alineamiento_final_BLOSUM62.fa", "fasta")
sequences = [str(record.seq) for record in alignment]

# Crear matriz de conteos (sin gaps)
counts_matrix = logomaker.alignment_to_matrix(sequences,\
    to_type='counts',\
    characters_to_ignore='-') # Excluir gaps para evitar sesgos

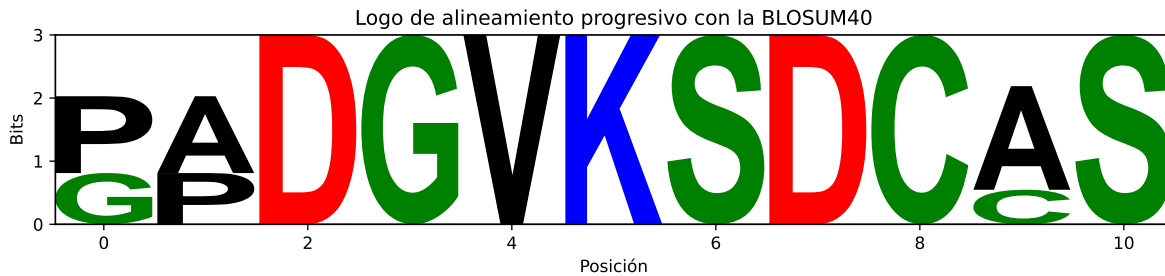
# Convertir a frecuencias relativas (ignorando gaps)
frequency_matrix = counts_matrix.div(counts_matrix.sum(axis=1), axis=0)

# Normalizar a frecuencias relativas y convertir a información (bits)
information_matrix = logomaker.transform_matrix(frequency_matrix, \
    from_type='probability', to_type='information')

# Generar el logo
logo = logomaker.Logo(
    information_matrix,
    color_scheme='chemistry', #Colores por propiedades químicas
    font_name='Arial',
    show_spines=True,
    stack_order='big_on_top' #Aminoácidos más frecuentes arriba
)

# Ajustar el estilo del gráfico
logo.ax.set_title("Logo de alineamiento progresivo con la BLOSUM40")
logo.ax.set_xlabel("Posición")
logo.ax.set_ylabel("Bits")
logo.ax.set_ylim(0, 3)
```

```
plt.tight_layout()
plt.show()
```



En este caso, el alineamiento progresivo se hizo con la matriz BLOSUM62, que tiene un mayor poder discriminativo para secuencias más cercanas evolutivamente. Podemos observar similitudes con el primer logo, pero algunas posiciones presentan diferencias en los aminoácidos predominantes o en su grado de conservación. Por ejemplo, la posición de “A” y “S” en el extremo derecho muestra una ligera diferencia en la proporción de letras, lo que sugiere un efecto de la matriz en la ponderación de las sustituciones. También existen discrepancias en la posición 2, donde podemos observar que en este caso solo tenemos un aminoácido (“D”) mientras que en el primer caso también teníamos “P”

3. Logo del alineamiento obtenido usando Clustal Omega

```
# Cargar y procesar el alineamiento
alignment = AlignIO.read("clustalo-I20250321-165928-0743-92365021-p1m.fa", \
    "fasta")
sequences = [str(record.seq) for record in alignment]

# Crear matriz de conteos (sin gaps)
counts_matrix = logomaker.alignment_to_matrix(sequences, \
    to_type='counts', \
    characters_to_ignore='-') # Excluir gaps para evitar sesgos

# Convertir a frecuencias relativas (ignorando gaps)
frequency_matrix = counts_matrix.div(counts_matrix.sum(axis=1), axis=0)

# Normalizar a frecuencias relativas y convertir a información (bits)
information_matrix = logomaker.transform_matrix(frequency_matrix, \
    from_type='probability', to_type='information')

# Generar el logo
logo = logomaker.Logo(
```

```

    information_matrix,
    color_scheme='chemistry', #Colores por propiedades químicas
    font_name='Arial',
    show_spines=True,
    stack_order='big_on_top' #Aminoácidos más frecuentes arriba
)

# Ajustar el estilo del gráfico
logo.ax.set_title("Logo de alineamiento progresivo con la BLOSUM40")
logo.ax.set_xlabel("Posición")
logo.ax.set_ylabel("Bits")
logo.ax.set_ylim(0, 3)

plt.tight_layout()
plt.show()

```



En este caso podemos observar una mayor variabilidad en las posiciones iniciales del alineamiento, con diferencias notables en la conservación de ciertas posiciones clave en comparación con los otros dos métodos. La distribución de letras es distinta en varias posiciones, como ya habíamos comentado anteriormente (apartado de comparación de alineamientos y árbol guía).

Comparación de logos

Los alineamientos manuales muestran una mayor conservación en ciertas posiciones, lo que puede deberse a que el árbol guía fue construido a mano con una estrategia específica, cambiando únicamente el uso de una matriz u otra mientras que Clustal Omega parece permitir una mayor variabilidad en algunas posiciones, lo que puede reflejar una optimización diferente basada en heurísticas computacionales.

La diferencia entre los logos generados con BLOSUM40 y BLOSUM62 indica que la elección de la matriz afecta la ponderación de las sustituciones, con BLOSUM62 favoreciendo conservaciones más estrictas. BLOSUM40, al estar diseñado para secuencias más divergentes, permite más variabilidad en algunas posiciones, lo que queda reflejado en el logo (posición 2).

Con todo ello, concluimos que el logo obtenido por Clustal Omega es el que presenta mayor variabilidad en las posiciones, mientras que el logo de forma manual mediante el uso de la matriz BLOSUM62 es el que presenta mayor conservación (posiciones 2, 3, 4, 5, 6, 7, 8 y 10) en todas las secuencias comparadas.

Concluiónes final:

El uso de diferentes matrices de sustitución afecta la conservación y la variabilidad en el alineamiento. BLOSUM40 favorece secuencias más filogenéticamente distantes, mientras que BLOSUM62 es más adecuado para secuencias cercanas. Además, hemos podido observar como los alineamientos manuales tienden a reflejar las decisiones del usuario en la construcción del alineamiento (por ejemplo, la introducción de gaps puede ser en ocasiones arbitraria), lo que puede generar patrones de conservación diferentes a los obtenidos mediante métodos automáticos como Clustal Omega.

Clustal Omega introduce una variabilidad diferente debido a su método de clustering, lo que puede dar lugar a alineamientos distintos, especialmente en regiones menos conservadas. Estas discrepancias en la conservación de ciertas posiciones pueden ser el resultado de cómo se ponderan las sustituciones en cada método, destacando la importancia de elegir la matriz y el algoritmo adecuados para cada conjunto de secuencias.

Con todo ello, podemos decir que, en general, el método de alineamiento y la matriz de sustitución pueden influir en la representación final de la conservación de secuencias. A pesar de ello, este ejercicio nos ha permitido identificar regiones conservadas y representarlas mediante un logo de secuencias usando métodos alternativos; y comprender cómo funcionan este tipo de algoritmos de alineamiento múltiple progresivos.