

Implementación de la tesis: “Multiple Sequence Alignment using Particle Swarm Optimization”

Vicente Machaca Arceda¹

¹Universidad Nacional de San Agustín de Arequipa. Email: vicente.machaca.a@gmail.com

Resumen

Este reporte presenta la implementación y análisis de la tesis: “Multiple Sequence Alignment using Particle Swarm Optimization”. Dicha tesis propone el uso de *Particle Swarm Optimization* (PSO) para solucionar el problema de *Multiple Sequence Alignment* (MSA). El autor propone la posición de cada *gap*, como vector para la representación de las partículas, luego define un *crossover* para simular el movimiento de una partícula hacia la partícula líder. Además, también se propone una mutación para diversificar la solución. En las pruebas realizadas, se comprobó que PSO obtiene muy buenos resultados e incluso similares a CLUSTALW.

1 Introducción

El área de Bioinformática ha tenido un auge en las últimas décadas. Por ejemplo, el proyecto: *The Human Genome Project* (HGP) que inició en 1990 y fue completado en 2003, donde participaron varios países como EE.UU., Reino Unido, Japón, Francia, Alemania, España y China (NIH 2021); el proyecto tenía como objetivo secuenciar todo el genoma humano, el cual resultó ser una secuencia de aproximadamente 3.2 billones de pares de bases (Archibald 2018). En la actualidad existen diversas áreas de investigación como: la predicción de estructura de proteínas, predicción de la función de una proteína a partir de estructuras de redes de proteínas, descubrimiento de medicamentos, predicción de enfermedades a partir del genoma, análisis de virus, etc. Es tan grande este campo de estudio que incluso se ha dividido en otras áreas como *metagenomics*, *proteomics*, *chemical informatics*, etc.

1.1 Problema

Uno de los tantos problemas que existen en bioinformática, es el alineamiento múltiple de secuencias. Se puede definir el alineamiento de secuencias como un método que permite determinar el grado de similitud entre dos o más secuencias (Xiong 2006), además el método puede insertar *gaps* dentro de las secuencias de consulta, con el objetivo de lograr la mayor cantidad de bases alineadas. Por ejemplo, en la Figura 1, se muestra el resultado luego de alinear varias secuencias de aminoácidos. Como podemos ver, se ha insertado *gaps* (-) para así maximizar la cantidad de aminoácidos que coinciden en la misma posición.

*Multiple Sequence Alignment
using Particle Swarm
Optimization* (2009)

Tipo de tesis:
Tesis de maestría

Autor de la tesis:
Fabien Bernard

Autor del reporte:
V.E. Machaca Arceda

© The Author(s) 2009. Published
by University of Pretoria Press
on behalf of the Society for
Political Methodology.

El alineamiento de secuencias puede dividirse en dos grupos: *pair-wise sequence alignment* (PSA) y *multiple sequence alignment* (MSA) (Xiong 2006). La diferencia radica en que el primero alinea solo dos secuencias y el segundo puede alinear dos a más secuencias. Además, el problema de MSA es considerado un problema NP completo (Wang y Jiang 1994). Debido a esto es que se han planteado heurísticas que logran obtener una solución local; el algoritmo más utilizado es CLUSTAL, fue propuesto por Higgins y Sharp (1988), este algoritmo ha sido mejorado con CLUSTALV (Higgins, Bleasby y Fuchs 1992), CLUSTALW (Thompson, Higgins y Gibson 1994) y CLUSTALX (Jeanmougin

```

RLA0_METVA --MIDAKSEHKIAPWKIEEVNALKLLKSANVIALIDMMEVPVAVOLQEIRDK
RLA0_METJA ---METKVKAHVAPWKIEEVKTLKGLIKSKPVVAIVDMMDVPAPOLQEIRDK
RLA0_PYRAB -----MAHVAEWKKKEVEELANLIKSYPVIALVDVSSMPAYPLSQMRRL
RLA0_PYRHO -----MAHVAEWKKKEVEELAKLIKSYPVIALVDVSSMPAYPLSQMRRL
RLA0_PYRFU -----MAHVAEWKKKEVEELANLIKSYPVVALVDVSSMPAYPLSQMRRL
RLA0_PYRKO -----MAHVAEWKKKEVEELANLIKSYPVIALVDVAGVPAYPLSKMRDK
RLA0_HALMA MSAESERKTETIPEWKQEEVDATVFMIESYESVGVVNIAGIPSRQLQDMRRD
RLA0_HALVO MSESEVRQTEVIPQWKREEVDLVDFIESYESVGVVGVAGIPSRQLQSMRRE
RLA0_HALSA MSAEEQRTTEEVPEWKRQEVAEVLDLLETYDSVGVVNVGTGIPSKQLQDMRRG
RLA0_THEAC -----MKEVSQKKELVNEITRIKASRSVAIVDTAGIRTRQIODIRCK
RLA0_THEVO -----MRKINPKKKEIVSELAODITKSKAVAIVDIKGVRTROMODIRAK
RLA0_PICTO -----MTEPAQWKIDFVKNLENEINSRKVAAIVSIKGLRNNEFQKIRNS

```

Figura 1. Ejemplo de *Multiple Sequence Alignment* (MSA). El método ha insertado *gaps* (-) en las secuencias necesarias para maximizar la cantidad de letras (aminoácidos) que concidan en la misma posición.

y col. 1998). Otro algoritmo importante es MUSCLE (Edgar 2004).

El problema de los algoritmos mencionados anteriormente, es que a pesar de ser heurísticas, el tiempo de procesamiento es muy alto. Cada segundo los datos genómicos crecen exponencialmente (Archibald 2018) y los algoritmos utilizados para buscar información en estas bases de datos, son basados en alineamiento. Entonces, mientras mas crecen los datos, mas lentos se vuelven estos algoritmos (Zablocki y col. 2009).

1.2 Objetivo

Zablocki y col. (2009), proponen en su tesis aplicar *Particle Swarm Optimization* (PSO) como alternativa a CLUSTAL para solucionar el problema de MSA.

2 Metodología aplicada

En esta sección describimos la metodología utilizada por Zablocki y col. (2009) para solucionar el problema de MSA utilizando PSO.

2.1 Representación de las partículas

Un de los primeros pasos para solucionar un problema utilizando PSO, es la representación de cada partícula. Zablocki y col. (2009) propone utilizar la posición de cada *gap* como un vector. Por ejemplo en al Figura 2, tenemos la representación del la partícula *leader* y una partícula cualquiera (ambas representan una posible solución al problema). Ahora en la Figura 3, tenemos otra forma de representar dichas partículas, en este caso solo estamos considerando la posición de cada *gap* insertado. La opción mas facil de controlar es la correspondiente a la Figura 3.

| | | | | | | | | | |
|-----------|---|---|---|---|---|---|---|---|---|
| Leader: | W | G | K | V | - | - | N | V | D |
| | W | D | K | V | - | - | N | - | - |
| | S | - | K | V | G | G | N | - | - |
| Particle: | W | G | K | - | - | V | N | V | D |
| | - | W | D | K | - | - | - | V | N |
| | S | - | K | V | G | G | - | N | - |

Figura 2. Ejemplo de la partícula *leader* y una partícula cualquiera. Ambas representan una posible solución a un problema de alineamiento de 3 secuencias de aminoácidos.

| | | | | |
|-----------|---|---|---|---|
| | 5 | 6 | | |
| Leader: | 5 | 6 | 8 | 9 |
| | 2 | 8 | 9 | |
| | 4 | 5 | | |
| Particle: | 1 | 5 | 6 | 7 |
| | 2 | 7 | 9 | |

Figura 3. Ejemplo de la partícula *leader* y una partícula cualquiera. En este caso solo estamos registrando la posición de cada *gap* insertado.

2.2 Movimiento de las partículas

Según el algoritmo de PSO, cada partícula debe acercarse al *leader* en cada iteración. Este acercamiento será implementado haciendo un *crossover*. El *crossover* permitirá que la nueva partícula (partícula en movimiento) tenga información de ambas partículas (*leader* y la partícula en movimiento).

Para aplicar un *crossover* entre dos partículas, debemos calcular primero la distancia entre ellas aplicando la formula de la Ecuación 1. Luego debemos calcular el punto de cruce de ambas partículas, para esto aplicamos la formula de la Ecuación 2, en este caso *length* representa la longitud de las secuencias.

$$distance = \frac{matchingGaps}{totalGaps} \quad (1)$$

$$crossPoint = rand(1, distance * length) \quad (2)$$

Por ejemplo, dada las partículas de la Figura 3 y suponiendo que hemos obtenido un *crossPoint* = 5, utilizando la Ecuación 2. La partícula en movimiento se desplazaría y su nueva representación sería el resultado de aplicar un *crossover*. En la Figura 4 mostramos el resultado de aplicar el *crossover*. Para lograr esto, por cada secuencia del *leader* insertamos sus *gaps* en la nueva partícula que son menores o iguales al *crossPoint*, luego completamos insertando los *gaps* de la partícula en movimiento que tengan *gaps* mayores al *crossPoint*.

| | | | |
|----------|---|---|---|
| | 5 | | |
| Particle | 5 | 6 | 7 |
| | 2 | 7 | 9 |

Figura 4. Resultado del movimiento de una partícula luego de aplicar *crossover*.

2.3 Función objetivo

La función objetivo propuesta por Zablocki y col. (2009) es la misma utilizada para medir el desempeño de los algoritmos basados en alineamiento como CLUSTAL para el problema de MSA. Por ejemplo en la Figura 5, se presenta un ejemplo del cálculo del *score* del alineamiento de dos secuencias. En este ejemplo se utilizó un *match score* : +1, *mismatch score*: 0, *opening gap*: -1 y *gap penalty*: -1. Para el caso de MSA, se procesa este *score* por cada posible combinación de dos secuencias. Para los experimentos se usaron estos mismos valores.

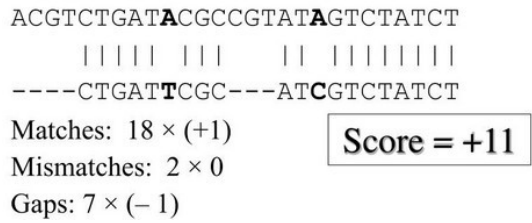


Figura 5. Ejemplo del cálculo del *score* del alineamiento de dos secuencias de ADN. En este caso se utilizó: *match score* : +1, *mismatch score* : 0, *opening gap* : -1 y *gap penalty* : -1

2.4 Mutaciones

Para diversificar al algoritmo de PSO, Zablocki y col. (2009) propuso hacer mutaciones a las partículas. El proceso de mutación consistía en escoger de manera aleatoria una partícula, luego dentro de dicha partícula se escogía una secuencia y se insertaba un *gap* en una posición aleatoria. Luego, para que todas las secuencias dentro de la partícula tengan la misma longitud, se insertaba un *gap* al principio o al final según un valor generado aleatoriamente.

3 Experimentos

En esta sección detallaremos las bases de datos utilizadas y los parametros del algoritmo PSO para poder replicar los resultados.

3.1 Bases de datos

Zablocki y col. (2009) propone un conjunto de 7 bases de datos. S1, S2, S3, S4, S5, S6 y S7 que el construyo a partir de un conjunto de secuencias de ADN. En la Tabla 1, presentamos el *ascension code* de cada secuencia utilizada. Además, el autor propuso un conjunto pequeño de secuencias para hacer pruebas rápidas, a este conjunto lo llamo S8. En nuestro caso, al ser una tesis un poco antigua, algunas secuencias ya no estaban disponibles en NCBI, y solo logramos obtener las secuencias de S6, S7 y S8.

Cuadro 1. Bases de datos utilizados por Zablocki y col. (2009). S1, S2, S3, S4, S5, S6 y S7 es el nombre que el autor definio para cada conjunto de secuencias. La segunda columna representa el *ascension code* de cada secuencia.

| Dataset | Secuencias |
|---------|---|
| S1 | HCV2L1A10 HCV2L3A5 HCV2L3C1 HCV2L3C8 HCV2L3D4 HCV2L3E6 HCV2L3A7 HCV2L3A9 HCV2L3B2 HCV2L3B1 |
| S2 | HS06674 HS06675 HS06676 HS06677 HS06679 |
| S3 | TPAHISIN TNIHISIN TNHISIN TMIHISIN TMHISIN THHISIN TFHISIN TEHISIN TCUHISIN TCHISIN TBHISIN TAUHISIN TAHISIN TTHISIN TSHISIN TRHISIN TPYHISIN TPIHISIN TPHISIN TCAHISIN TLHISIN |
| S4 | HI1U16764 HI1U16766 HI1U16768 HI1U16776 HI1U16778 HI1U16770 HI1U16774 HI1U16772 |
| S5 | HI1U16765 HI1U16767 HI1U16769 HI1U16771 HI1U16773 HI1U16775 HI1U16777 HI1U16779 |
| S6 | PP59651 PP59652 PP59653 PP59654 PP59655 PP59656 |
| S7 | AB023287 AB023286 AB023285 AB023284 AB023283 AB023279 AB023278 AB023276 |

Como se menciona anteriormente, se utilizaron las bases de datos S6, S7 y S8. En la Tabla 2, detallamos datos adicionales de las bases de datos utilizadas. Por ejemplo la base de datos mas

sencilla es S8, en la Figura 6, mostramos este conjunto de secuencias.

Cuadro 2. Bases de datos utilizados en los experimentos.

| <i>Dataset</i> | <i>Longitud mínima</i> | <i>Longitud máxima</i> | <i>Cantidad de secuencias</i> |
|----------------|------------------------|------------------------|-------------------------------|
| S6 | 8 | 17801 | 153 |
| S7 | 457 | 457 | 8 |
| S8 | 7 | 10 | 5 |

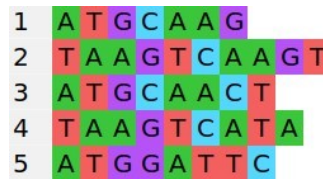


Figura 6. Conjunto de secuencias de la base de datos S8. Se utilizo Mega para colorear cada base.

3.2 Parametros

Los parametros son descritos en la Tabla 3. Para reducir el tiempo de procesamiento reducimos el valor de algunos parametros como la cantidad de iteraciones, el autor proponia en este caso 1000 iteraciones. Luego la cantidad de veces que se replico cada experimento era de 30 y nosotros solo lo replicamos 10 veces.

Cuadro 3. Parametros utilizados en los experimentos. Algunos parametros fueron distintos a los planteados por Zablocki y col. (2009) (iteraciones), con el objetivo de reducir el tiempo de procesamiento.

| <i>Parametro</i> | <i>Valor</i> |
|---|--------------|
| Iteraciones | 30 |
| Cantidad de partículas | 25 |
| Probabilidad de mutación | 0.2 |
| Gaps permitidos | 30 % |
| Cantidad de veces que se replico el experimento | 10 |

4 Resultados

En la Tabla 4, se muestra el *score* (resultado de la función objetivo) obtenido en cada base de datos. Como se puede apreciar, el *score* depende de la complejidad del problema, por ejemplo en el conjunto de secuencias S8 obtenemos un *score* bajo, pero para las secuencias de S6 (la mas compleja), obtenemos un *score* mucho mas grande. La idea es que mientras mayor sea el *score*, es mejor.

Para visualizar mejor el resultado de los alineamientos se utilizo la herramieta Mega, esta herramienta colorea cada base según su posición para ver mejor su alineamiento. En la Figura 7, mostramos el resultado luego de alinear el conjunto de secuencias de S8 utilizando PSO y CLUSTALW. Como podemos ver, existen muchas soluciones para cada problema, e incluso podemos llegar a tener un diferentes soluciones con un mismo *score*. En la Figura 8, mostramos un parte del alineamiento de las secuencias de S7, recordemos que estas secuencias tienen una longitud de 457 bases.

Cuadro 4. Resultados obtenidos en cada base de datos. Se hicieron pruebas incluyendo y excluyendo la mutación.

| <i>Dataset</i> | <i>PSO con mutación</i> | <i>PSO sin mutación</i> | <i>CLUSTALW</i> |
|----------------|-------------------------|-------------------------|-----------------|
| S6 | 12678 | 10012 | 18045 |
| S7 | 11105 | 9054 | 12564 |
| S8 | 32 | 28 | 49 |



Figura 7. **Izquierda:** Resultado luego de alinear las secuencias de la base de datos S8 utilizando PSO. **Derecha:** Resultado luego de alinear las secuencias de la base de datos S8 utilizando CLUSTALW.

5 Discusión

De la Tabla 4, comprobamos que el desempeño de PSO es bueno y casi llega a una solución similar a CLUSTAW. Además, uno de los problemas de PSO es que llega muy pronto a óptimos locales, debido a eso en los resultados comprobamos como aplicando una diversificación (mutación) se mejoran los resultados. Otro de PSO y otros algoritmos similares es en como el autor plantea el desplazamiento de cada partícula, en este trabajo se planteo utilizar un *crossover*, pero existen trabajos recientes propuestos por Zhan *y col.* (2019), Lalwani y Sharma (2019) y Moustafa *y col.* (2017) donde obtienen mejores resultados modificando la función de movimiento de cada partícula y combinando PSO con las cadenas de Markov.

Una prueba visual del buen desempeño de PSO es presentado en las Figuras 7 y 8. Como vemos en las imágenes, se aprecia que se ha obtenido un alineamiento aceptable. Cada base (letra) es representada con un color distinto y esta forma de visualización ayuda bastante para determinar el alineamiento entre dos o mas secuencias.

6 Conclusiones

La tesis propone el uso de PSO para solucionar el problema de MSA. El autor evalua los resultados en 8 bases de datos, que representan un conjunto de secuencias de diferentes longitudes. El autor obtuvo dichas secuencias de NCBI, pero al ser una tesis del 2009, algunas secuencias ya no estaban presentes.

El autor tambien propone aplicar una mutación a las partículas para evitar caer en óptimos locales, esta mutación se basaba en insertar un *gap* dentro de una secuencia de una partícula. Esta diversificación (mutación) mejoro el desempeño de PSO.

Para comprobar el desempeño de PSO, se comparó sus resultados con CLUSTALW. El desempeño de PSO fue bueno y logro un *score* similar a CLUSTALW. En la actualidad existen varias modificaciones en la forma de como se define una partícula y la función de movimiento de esta con muy buenos resultados.

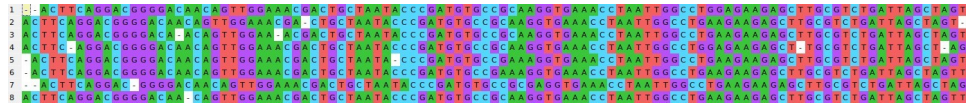


Figura 8. Resultado luego de alinear las secuencias de la base de datos S7 utilizando PSO. Solo se esta mostrando un extracto, porque las secuencias tienen una longitud de 457 bases.

Material adicional

Para material adicional a este reporte, por favor visite: <https://github.com/arceda/DOCTORADO/tree/master/ia/src/thesis>

Referencias

- Archibald, J. M. 2018. *Genomics: A Very Short Introduction*. Volumen 559. Oxford University Press.
- Edgar, R. C. 2004. "MUSCLE: multiple sequence alignment with high accuracy and high throughput". *Nucleic acids research* 32 (5): 1792-1797.
- Higgins, D. G., A. J. Bleasby y R. Fuchs. 1992. "CLUSTAL V: improved software for multiple sequence alignment". *Bioinformatics* 8 (2): 189-191.
- Higgins, D. G., y P. M. Sharp. 1988. "CLUSTAL: a package for performing multiple sequence alignment on a microcomputer". *Gene* 73 (1): 237-244.
- Jeanmougin, F., J. D. Thompson, M. Gouy, D. G. Higgins y T. J. Gibson. 1998. "Multiple sequence alignment with Clustal X". *Trends in biochemical sciences* 23 (10): 403-405.
- Lalwani, S., y H. Sharma. 2019. "Multi-objective three level parallel PSO algorithm for structural alignment of complex RNA sequences". *Evolutionary Intelligence*: 1-9.
- Moustafa, N., M. Elhosseini, T. H. Taha y M. Salem. 2017. "Fragmented protein sequence alignment using two-layer particle swarm optimization (FTLPSO)". *Journal of King Saud University-Science* 29 (2): 191-205.
- NIH. 2021. *The Human Genome Project*. Web resource. National Human Genome Research Institute.
- Thompson, J. D., D. G. Higgins y T. J. Gibson. 1994. "CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice". *Nucleic acids research* 22 (22): 4673-4680.
- Wang, L., y T. Jiang. 1994. "On the complexity of multiple sequence alignment". *Journal of computational biology* 1 (4): 337-348.
- Xiong, J. 2006. *Essential bioinformatics*. Cambridge University Press.
- Zablocki, F. B. R., y col. 2009. "Multiple sequence alignment using Particle swarm optimization". Tesis doctoral, University of Pretoria.
- Zhan, Q., N. Wang, S. Jin, R. Tan, Q. Jiang e Y. Wang. 2019. "ProbPFP: a multiple sequence alignment algorithm combining hidden Markov model optimized by particle swarm optimization with partition function". *BMC bioinformatics* 20 (18): 1-10.