

APLIKASI DYNAMIC PROGRAMMING DALAM ALGORITMA NEEDLEMAN-WUNSCH UNTUK PENJAJARAN DNA DAN PROTEIN

Dian Perdhana Putra - 13507096

Teknik Informatika ITB
Jl. Ganesha 10 Bandung
e-mail: if17096@students.if.itb.ac.id

ABSTRAK

Dynamic programming adalah salah satu algoritma yang digunakan untuk menemukan nilai optimal dari suatu permasalahan. Dalam *dynamic programming*, pemecahan suatu masalah dibagi menjadi beberapa langkah (*step*) atau tahapan (*stage*) sedemikian hingga solusi dari sebuah persoalan dapat dipandang sebagai serangkaian keputusan yang saling berkaitan.

Penerapan *dynamic programming* yang amat penting dalam bidang bioinformatika adalah penjajaran *sequence DNA* dan *protein*. Salah satu penerapan *dynamic programming* adalah pada algoritma Needleman-Wunsch untuk mencari *best alignment* dari penjajaran dua buah *sequence*. Makalah ini akan membahas tentang penggunaan algoritma Needleman-Wunsch, berikut dengan contoh penggunaannya dalam penjajaran *sequence DNA* dan *protein*.

Kata kunci: *dynamic programming*, bioinformatika algoritma Needleman-Wunsch, *sequence*, *alignment*, *DNA*, *protein*.

1. PENDAHULUAN

1.1 DNA dan Asam Amino

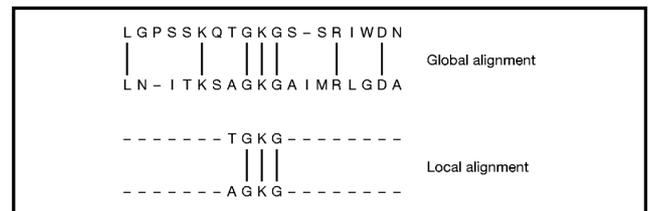
Asam deoksiribonukleat (*deoxyribonucleic acid*), atau biasa disebut *DNA*, adalah biomolekul yang berupa asam nukleat (terdapat dalam inti sel atau *nucleus*), yang berfungsi untuk menyimpan informasi genetik pada suatu organisme. *DNA* berbentuk untai ganda (*double helix*) yang disatukan oleh ikatan hidrogen antara basa-basa di dalam kedua untai tersebut. Basa-basa tersebut adalah Adenin (A), Sitosin (C), Guanin (G), dan Timin (T). Adenin berikatan dengan Timin, dan Sitosin berikatan dengan Guanin. *DNA* dapat melakukan replikasi *DNA*, di mana ketika replikasi *DNA* ini dilakukan, terbentuklah *DNA-DNA* baru yang memiliki informasi genetik yang serupa dengan induknya.

Asam amino adalah suatu senyawa organik yang tersusun dari karbon, hidrogen, oksigen, dan nitrogen. Asam amino merupakan pembentuk dari protein dalam tubuh manusia dan diperlukan untuk metabolisme. Ada 20

macam asam amino yang diperlukan untuk membentuk protein dalam tubuh manusia, yaitu Alanine (Ala,A), Arginine (Arg,R), Asparagine (Asn,N), Aspartate (Asp,D), Cysteine (Cys,C), Glutamine (Gln,Q), Glutamate (Glu,E), Glycine (Gly,G), Histidine (His,H), Isoleucin (Ile,I), Leucine (Leu,L), Lysine (Lys,K), Methionine (Met,M), Phenylalanine (Phe,F), Proline (Pro,P), Serine (Ser,S), Threonine (Thr,T), Tryptophan (Trp,W), Tyrosine (Tyr,Y), dan Valine (Val,V). *DNA* menciptakan instruksi-instruksi yang diperlukan untuk sintesis asam amino.

1.2 Penjajaran Barisan

Sequence Alignment (penjajaran barisan) adalah prosedur untuk menjajarkan dua buah barisan (*sequence*) dari *DNA* atau protein dengan tujuan mencari kesamaan di antara barisan-barisan tersebut atau untuk membuktikan bahwa kedua *sequence* yang dibandingkan berasal dari *sequence* yang sama. Ada dua macam metode *sequence alignment*, yaitu *global alignment* dan *local alignment*.



Gambar 1. Perbedaan antara *global alignment* dan *local alignment* dari dua buah *sequence* [5]

Dalam *global alignment*, penjajaran dilakukan untuk keseluruhan *sequence*, yaitu dengan menggunakan sebanyak mungkin karakter dalam *DNA*. Sedangkan *local alignment* adalah menjajarkan sebagian dari *sequence-sequence*, biasanya bagian yang memiliki tingkat kemiripan yang cukup tinggi.

Dalam bioinformatika, *dynamic programming* digunakan sebagai metode untuk menjajarkan dua *sequence* protein atau *DNA*. *Dynamic programming* digunakan karena metode ini dapat memberikan hasil *alignment* yang paling optimal di antara *sequence-sequence* yang diberikan.

$$D(i, j) = \begin{cases} D(0,0) = 0 \\ \max \begin{cases} D(i-1, j-1) + s(x_i, y_i) \\ D(i-1, j) + g_x \\ D(i, j-1) + g_y \end{cases} \end{cases} \quad (2)$$

dimana :

$s(x_i, y_i)$: elemen matriks substitusi di residu i pada *sequence* x dan residu j di *sequence* y .

$D(i-1, j-1)$: elemen matriks D di diagonal kiri atas .

$D(i, j-1)$: elemen matriks D di kiri $D(i, j)$.

$D(i-1, j)$: elemen matriks D di atas $D(i, j)$.

g_x : *gap penalty* untuk *sequence* x .

g_y : *gap penalty* untuk *sequence* y .

3. Traceback Step

Setelah matriks nilai berukuran $(M+1) \times (N+1)$ telah terisi sepenuhnya, maka *alignment score* maksimum dari dua buah *sequence* adalah nilai dari elemen paling kanan bawah dari matriks nilai, yaitu $D(M+1, N+1)$.

Lakukan *traceback step* dengan menerapkan *dynamic programming* di sini, mulai dari $D(M+1, N+1)$.

hingga ke elemen pertama tabel, $D(0,0)$.

Misalkan ada dua buah *sequence*, yaitu *sequence* A dan *sequence* B , maka prosedur untuk melakukan *traceback step* adalah :

1. Buat dua buah *string*, disebut *alignmentA* dan *alignmentB*, untuk menyimpan *alignment* yang terbentuk. Jika *sequence* A lebih panjang dari *sequence* B , maka panjang *alignmentA* dan *alignmentB* sama dengan panjang *sequence* A . Sebaliknya, jika *sequence* B lebih panjang dari *sequence* A , maka panjang *alignmentA* dan *alignmentB* sama dengan panjang *sequence* B . Inisialisasi *alignment* dan *alignmentB* dengan String kosong ("").
2. Buat variabel pencacah untuk menyimpan *alignment score*, yaitu *alignmentScore*. Inisialisasi dengan 0.
3. Buat empat buah variabel bertipe bilangan bulat, yaitu :
 $score$, untuk menyimpan $D(i, j)$
 $scoreDiag$, untuk menyimpan $D(i-1, j-1)$
 $scoreKiri$, untuk menyimpan $D(i, j-1)$
 $scoreAtas$, untuk menyimpan $D(i-1, j)$
4. Buat dua buah variabel iterasi, i dan j . Inisialisasi i dengan nilai dari panjang *sequence* A dan j dengan nilai dari panjang *sequence* B .
5. Selama $(i \geq 0)$ dan $(j \geq 0)$, lakukan pengulangan sebagai berikut :
 - Jika $score = scoreDiag + S(A_{i-1}, B_{j-1})$, maka :
 $alignmentA \leftarrow A(i-1) + alignmentA$

$alignmentB \leftarrow B(j-1) + alignmentB$
 $alignmentScore \leftarrow alignmentScore + S(A_{i-1}, B_{j-1})$
 $i \leftarrow i - 1$
 $j \leftarrow j - 1$

- Jika $score = scoreAtas + g_x$ maka :
 $alignmentA \leftarrow A(i-1) + alignmentA$
 $alignmentB \leftarrow "-" + alignmentB$
 $alignmentScore \leftarrow alignmentScore + S(A_{i-1}, B_{j-1})$
 $i \leftarrow i - 1$
- Jika $score = scoreKiri + g_y$ maka :
 $alignmentA \leftarrow "-" + alignmentA$
 $alignmentB \leftarrow B(j-1) + alignmentB$
 $alignmentScore \leftarrow alignmentScore + S(A_{i-1}, B_{j-1})$
 $j \leftarrow j - 1$

6. Tuliskan *alignment* yang didapatkan.

3.2.1 Contoh penerapan algoritma Needleman-Wunsch

Misalkan, terdapat :

Sequence A : IDEA

Sequence B : WIKIPEDIA

Gap penalty untuk *sequence* A : -10

Gap penalty untuk *sequence* B : -10

Matriks substitusi, S (dibuat dengan mengacu pada *BLOSUM62 substitution matrix*) :

Tabel 1: Matriks substitusi S

	D	E	W	A	P	I	K
D	6	2	-4	-2	-1	-3	-1
E	2	5	-3	-1	-1	-3	-1
W	-4	-3	11	-3	-4	-3	-3
A	-2	-1	-3	4	-1	-1	-1
P	-1	-1	-4	-1	7	-3	-1
I	-3	-3	-3	-1	-3	4	-3
K	-1	1	-3	-1	-1	-3	5

1. Inisialisasi

Panjang *sequence* A : 4

Panjang *sequence* B : 9

Buat matriks nilai D berukuran 5×10 , indeks dimulai dari 0, isi dengan *gap penalties* :

Tabel 2 : Matriks nilai D awal

	W	I	K	I	P	E	D	I	A	
I	0	-10	-20	-30	-40	-50	-60	-70	-80	-90
D	-10									
E	-20									
A	-30									
A	-40									

2. Mengisi matriks

Untuk $D(1,1)$:

	W
I	-10

$$D(0,0) = 0$$

$$D(0,1) = -10$$

$$D(1,0) = -10$$

$$S(I, W) = -3$$

Maka :

$$D(1,1) = \max \left((D(0,0) + (-3)), (D(0,1) + (-10)), (D(1,0) + (-10)) \right)$$

$$D(1,1) = \max (-3, -20, -20)$$

$$D(1,1) = -3$$

Untuk $D(1,2)$:

	W	I
I	-10	-20
I	-3	

$$D(0,1) = -10$$

$$D(0,2) = -20$$

$$D(1,1) = -3$$

$$S(I, I) = 4$$

Maka :

$$D(1,2) = \max \left((D(0,1) + 4), (D(0,2) + (-10)), (D(1,1) + (-10)) \right)$$

$$D(1,2) = \max (-6, -30, -13)$$

$$D(1,2) = -6$$

Pengisian matriks dilakukan hingga seluruh elemen matriks nilai terisi. Matriks yang terbentuk adalah sebagai berikut :

Tabel 3 : Matriks nilai D yang telah terisi penuh

	W	I	K	I	P	E	D	I	A	
I	0	-10	-20	-30	-40	-50	-60	-70	-80	-90
D	-10	-3	-6	-16	-26	-36	-46	-56	-66	-76
E	-20	-13	-6	-7	-17	-27	-34	-40	-50	-60
A	-30	-23	-16	-5	-10	-18	-22	-32	-42	-51
A	-40	-33	-24	-15	-6	-11	-19	-24	-33	-38

3. Traceback Step

Setelah matriks nilai D terisi penuh, maka *alignment score* optimal adalah -38, yaitu nilai dari $D(4,9)$. *Traceback step* dilakukan untuk mendapatkan *alignment* dari *sequence A* dan *sequence B* yang memiliki *alignment score* sama dengan -38.

Untuk $D(4,9)$:

	W	I
A	-42	-51
A	-33	-38

$$score = -38$$

$$score \text{ diagonal} = -42 + S(A, A) = -42 + 4 = -38$$

$$score \text{ atas} = -51 + gap \text{ sequence II} = -51 + (-10) = -61$$

$$score \text{ kiri} = -33 + gap \text{ sequence I} = -33 + (-10) = -43$$

Karena $score = score \text{ diagonal} = -38$, maka
 $i = i - 1 = 4 - 1 = 3$
 $j = j - 1 = 9 - 1 = 8$
 Sehingga sel selanjutnya adalah $D(3,8)$.

Untuk $D(3,8)$:

	W	I
E	-40	-50
E	-32	-42

$$score = -42$$

$$score \text{ diagonal} = -40 + S(I, E) = -40 + 3 = -37$$

$$score \text{ atas} = -50 + gap \text{ sequence II} = -50 + (-10) = -60$$

$$score \text{ kiri} = -32 + gap \text{ sequence I} = -32 + (-10) = -42$$

Karena $score = score \text{ kiri} = -42$, maka
 $j = j - 1 = 8 - 1 = 7$
 Sehingga sel selanjutnya adalah $D(3,7)$.

Traceback step dilakukan hingga elemen terakhir yaitu $D(0,0)$.

Tabel 4 : *best alignment* dari *sequence A* dan *sequence B* pada matriks nilai D

	W	I	K	I	P	E	D	I	A
0	-10	-20	-30	-40	-50	-60	-70	-80	-90
I	-10	-3	-6	-16	-26	-36	-46	-56	-66
D	-20	-13	-6	-7	-17	-27	-34	-40	-50
E	-30	-23	-16	-5	-10	-18	-22	-32	-42
A	-40	-33	-24	-15	-6	-11	-19	-24	-33

Hasil Alignment

Sequence A : ---I-DE-A
 Sequence B : | |
 WIKIPEDIA

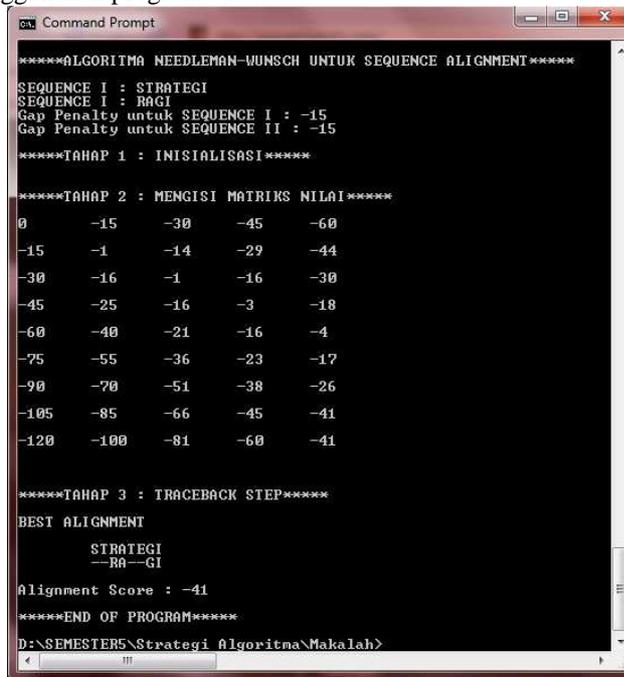
Alignment score =

$$\begin{aligned}
 &S(-, W) + S(-, I) + S(-, K) + S(I, I) + S(-, P) \\
 &+ S(D, E) + S(E, D) + S(-, I) + S(A, A) \\
 &= (-10) + (-10) + (-10) + 4 + (-10) + 2 + 2 + (-10) + 4 \\
 &= -38.
 \end{aligned}$$

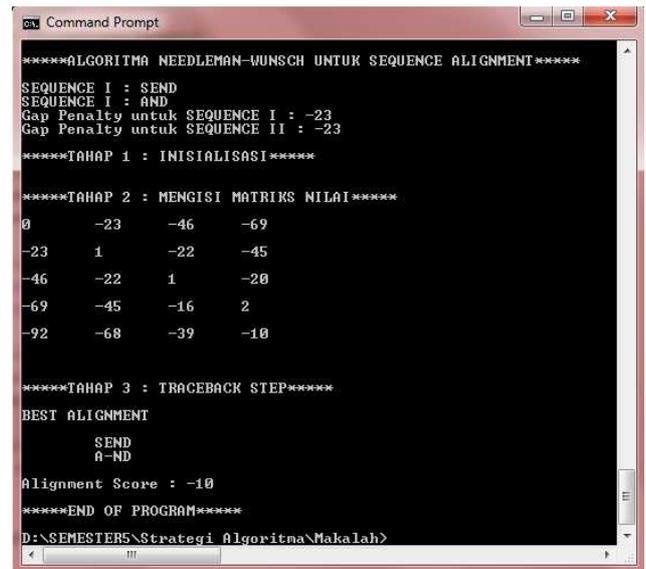
4. HASIL EKSPERIMEN

Penulis membuat program sederhana dalam bahasa Java yang mendemonstrasikan *alignment* untuk dua buah *sequence* protein. *Source code* program ini dapat diunduh di <http://students.if.itb.ac.id/~if17096/ProgramStrategiAlgoritma.zip>. Dalam program ini, matriks substitusi yang digunakan adalah *BLOSUM62 substitution matrix*.

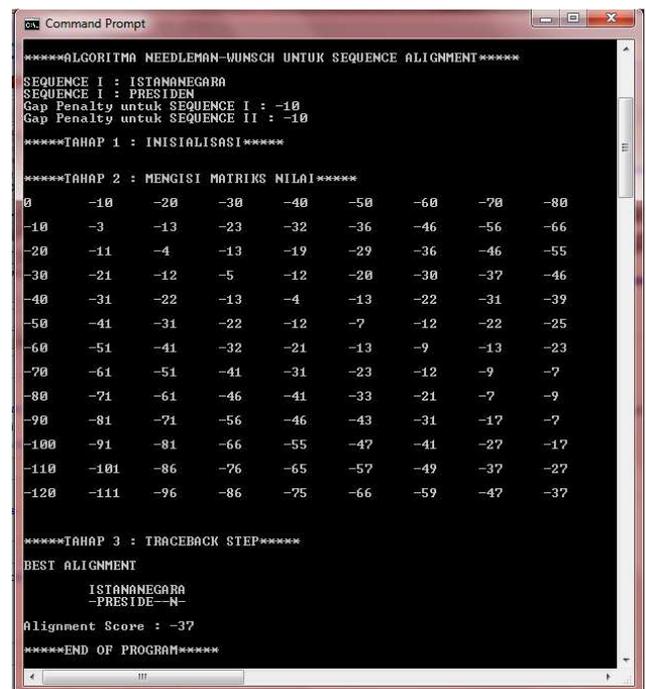
Berikut adalah beberapa hasil percobaan dengan menggunakan program tersebut :



Gambar 3 : Hasil Percobaan I



Gambar 4 : Hasil Percobaan II



Gambar 5 : Hasil Percobaan III

5. ANALISIS

Dari hasil percobaan, terbukti bahwa algoritma NeedlemanWunsch dapat digunakan untuk menemukan *best alignment*. Penyelesaian masalah dalam *traceback step* dibagi menjadi beberapa tahap atau langkah. Jumlah langkah yang terjadi sesuai dengan panjang dari *sequence* yang memiliki panjang terbesar.

Jika *sequence I* memiliki panjang m dan *sequence II* memiliki panjang n , maka untuk tahap inialisasi matriks nilai, kita perlu memasukkan nilai dari kolom pertama sebanyak m kali dan baris pertama sebanyak n kali. Pada tahap inialisasi ini, kompleksitasnya adalah $O(m + n)$. Pada tahap mengisi matriks nilai, terdapat kalang bersarang, sehingga untuk mengisi matriks nilai, kompleksitas waktunya adalah $O(mn)$. Untuk *traceback step*, ada dua langkah (*step*) yang dilakukan yaitu menandai elemen dan menemukan *final path*. Untuk tahap menandai elemen, perpindahan yang mungkin dilakukan adalah sebanyak m baris dan n kolom, sehingga kompleksitasnya adalah $O(m + n)$. Sedangkan untuk menemukan *final path*, langkah ini melibatkan maksimum n langkah (dengan asumsi $n \geq m$). Oleh karena itu kompleksitasnya adalah $O(n)$. Jadi kompleksitas total dari algoritma Needleman Wunsch adalah :

$$O(m + n) + O(mn) + O(m + n) + O(n) = O(mn).$$

Pencarian *best alignment* dengan menggunakan algoritma brute force memiliki kompleksitas sebesar $O(m!n!)$ sehingga dengan demikian algoritma Needleman-Wunsch yang memiliki kompleksitas $O(mn)$ akan jauh lebih mangkus untuk m dan n yang sangat besar.

6. KESIMPULAN

1. *Dynamic programming* memiliki penerapan yang amat beragam dalam berbagai bidang keilmuan, salah satunya dalam bidang bioinformatika, yaitu pada algoritma Needleman-Wunsch untuk menemukan *best alignment* pada penjajaran *sequence* DNA dan protein.
2. Algoritma Needleman-Wunsch dapat digunakan untuk mencari *best alignment* dari dua buah *sequence*, yaitu *alignment* yang memiliki *alignment score* maksimum.
3. Penggunaan algoritma Needleman-Wunsch untuk mencari *best alignment* dari dua buah *sequence* DNA atau protein akan jauh lebih efektif dan efisien daripada menggunakan *brute force*, terutama untuk *sequence-sequence* yang memiliki panjang sangat besar.
4. Algoritma Needleman-Wunsch memiliki kompleksitas $O(mn)$, sedangkan brute force memiliki kompleksitas $O(m!n!)$.
5. Algoritma Needleman-Wunsch adalah algoritma yang cukup bagus untuk menemukan *best alignment* dari dua buah *sequence*. Di masa yang akan datang, algoritma ini dapat dikembangkan lagi sehingga memberikan hasil yang lebih baik.

7. REFERENSI

- [1] Munir, Rinaldi. "Strategi Algoritma". Teknik Informatika, Bandung, 2009.
- [2] Böckenhauer, Hans-Joachim and Dirk Bongartz. "Algorithmic Aspects of Bioinformatics", Springer, New York, 2007.
- [3] Baxevanis, Andreas D. and B.F. Francis Oullette. "BIOINFORMATICS : A Practical Guide to the Analysis of Genes and Proteins", John Wiley & Sons, New York, 2001.
- [4] Lesk, Arthur. "Introduction to Bioinformatics", Oxford University Press Inc, New York, 2002.
- [5] Mount, David. "Bioinformatics : Sequence and Genome Analysis", Cold Spring Harbor Laboratory Press, New York, 2001.
- [6] Barnes, Michael and Ian C. Gray. "Bioinformatics for Geneticists", John Wiley & Sons, New York, 2003.
- [7] Jonassen, Inge and Kim Junhyong. "Algorithms in Bioinformatics", Springer, New York, 2005.
- [8] Cormen, Thomas H. et al. "Introduction to Algorithms 2nd ed", The MIT Press, Cambridge, 2001

<http://www.avatar.se/molbioinfo2001/dynprog/dynamic.html> diakses pada 19 Desember 2009, pukul 20.13.

http://www.avatar.se/molbioinfo2001/dynprog/adv_dynamic.html diakses pada 19 Desember 2009, pukul 20.35.

http://en.wikipedia.org/wiki/Dynamic_programming diakses pada 18 Desember 2009 pukul 21.00

http://en.wikipedia.org/wiki/Needleman-Wunsch_algorithm diakses pada 18 Desember 2009 pukul 21.15