

Alignment-Verfahren zum Vergleich biologischer Sequenzen

Hans-Joachim Böckenhauer
Dennis Komm

Volkshochschule Zürich

30. April 2014

Fragestellung

Finde eine Methode zum Vergleich von DNA-Molekülen oder Proteinen

Motivation

- Suche in Genom- oder Protein-Datenbanken
- Erstellung von phylogenetischen Bäumen
- Teilproblem bei der DNA-Sequenzierung

Vorgehen:

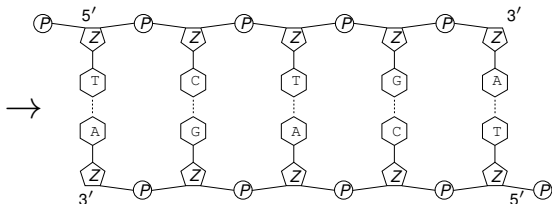
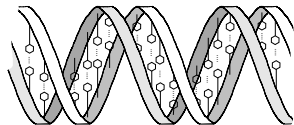
- Finde geeignete Datenstruktur für die Moleküle
- Definiere geeignetes Ähnlichkeitsmass
- Entwirf einen möglichst effizienten Algorithmus zur Berechnung der Ähnlichkeit bezüglich dieses Masses

Modellierung der Daten

Moleküle als Strings:

DNA und Proteine sind lange kettenförmige Moleküle bestehend aus wenigen oft wiederholten Grundbausteinen

⇒ Darstellung als Strings



→ TCTGA

Anforderung:

Ähnlichkeitsmass soll häufige Veränderungen in DNA- oder Proteinsequenzen widerspiegeln:

- Austausch einzelner Basen oder Aminosäuren
- Einfügen oder Löschen kurzer Teilsequenzen

Alignments

Idee:

Schreibe beide Strings buchstabenweise untereinander, füge dabei an beliebigen Stellen Lückensymbole ein

Beispiel:

Eingabe: Strings $s = \text{GACGATTATG}$ und $t = \text{GATCGAATAG}$

mögliche Alignments:

$$\begin{array}{ll} s' = \text{GA-} \color{red}{\text{CGATTATG}} & s'' = \text{GAC-} \color{red}{\text{GATTATG}} \\ t' = \text{GAT} \color{red}{\text{CGAATA}} \text{-G} & t'' = \text{GAT} \color{red}{\text{CGAATAG}} \text{-} \end{array}$$

$$\begin{array}{l} s''' = \text{GACGAT} \text{——} \text{TA-TG} \\ t''' = \text{——} \text{GAT} \color{red}{\text{CGAATAG}} \text{——} \end{array}$$

Bemerkung:

Spalten bestehend aus zwei Lücken sind sinnlos, kommen also nicht vor

Idee zur Bewertung:

- Alignment spaltenweise bewerten, dann über alle Spalten aufsummieren
- Spalte mit Lücke erhält Kosten g
- Spalte mit Buchstaben a und b erhält Kosten $p(a, b)$
- $p(a, b)$ ist Null für $a = b$ und gross für $a \neq b$
- Ziel: Minimiere die Kosten

Edit-Distanz

Beispiel für Bewertung:

Edit-Distanz (Levenshtein, 1966): Zähle Mismatches und Lücken,
d. h. $g = 1$, $p(a, a) = 0$ und $p(a, b) = 1$ für $a \neq b$

Beispiel:

Eingabe: Strings $s = \text{GACGATTATG}$ und $t = \text{GATCGAATAG}$

mögliche Alignments:

$$\begin{array}{ll} s' = \text{GA}\text{-CGATTATG} & s'' = \text{GAC}\text{-GATTATG} \\ t' = \text{GATCGA}\text{ATA-G} & t'' = \text{GATCGA}\text{ATAG-} \end{array}$$

$$\begin{array}{l} s''' = \text{GACGAT}\text{—TA-TG} \\ t''' = \text{—GATCGA}\text{ATAG—} \end{array}$$

Edit-Distanz:

$$d_{\text{edit}}(s', t') = 3 \quad d_{\text{edit}}(s'', t'') = 5 \quad d_{\text{edit}}(s''', t''') = 10$$

Frage: Wie kann man ein optimales Alignment finden?

Idee: Probiere alle Alignments durch

Abschätzung des Rechenaufwands: Wieviele verschiedene Alignments gibt es?

Theorem

Seien s und t zwei Strings der Länge n .

*Dann gibt es **mehr als 3^n** mögliche Alignments von s und t .*

Beweisidee:

- Alignment ist eindeutig bestimmt durch die Position der eingefügten Lücken
- Zähle nur Alignments einer bestimmten einfachen Form

Anzahl möglicher Alignments

- Drei Möglichkeiten für zwei Strings a und b der Länge 1:

$$\begin{array}{ccc} a & a- & -a \\ b & -b & b- \end{array}$$

- Sei $s = s_1 s_2 \dots s_n$ und $t = t_1 t_2 \dots t_n$
- \Rightarrow Die drei Möglichkeiten anwenden auf alle Paare s_j und t_j
- **Beispiel** für $n = 3$:

$$\begin{array}{c} s_1 & - & | & - & s_2 & | & s_3 \\ - & t_1 & | & t_2 & - & | & t_3 \end{array} \quad \text{oder} \quad \begin{array}{c} s_1 & | & s_2 & - & | & - & s_3 \\ t_1 & | & - & t_2 & | & t_3 & - \end{array}$$

- \Rightarrow Daraus lassen sich 3^n Alignments zusammensetzen
- **Beachte:** Dies sind nicht alle möglichen Alignments, es gibt noch mehr, zum Beispiel

$$\begin{array}{c} - & - & s_1 & s_2 & s_3 \\ t_1 & t_2 & t_3 & - & - \end{array} \quad \text{oder} \quad \begin{array}{c} s_1 & s_2 & - & - & s_3 \\ - & t_1 & t_2 & t_3 & - \end{array}$$

Exponentielle Laufzeit

n	10	50	100	300	10 000
$10n$	100	500	1 000	3 000	100 000
$4n^2$	400	10 000	40 000	360 000	400 000 000
n^3	1 000	125 000	1 000 000	27 000 000	13 Ziffern
3^n	59 049	24 Ziffern	48 Ziffern	143 Ziffern	4 772 Ziffern

- ⇒ Vollständige Suche ist viel zu langsam
- ⇒ Intelligenter Algorithmenentwurf nötig

Prinzip der dynamischen Programmierung:

Lösung für die gesamte Eingabe zusammensetzen aus Teillösungen für Teilprobleme, beginnend mit den kleinsten Teilproblemen

Problem: Finde geeignete Teilprobleme

Idee (Needleman und Wunsch, 1970):

- Alle Paare von Anfangsstücken (Präfixen) der gegebenen Strings als Teilprobleme
- Berechne Alignments für längere Präfixe aus den optimalen Alignments für kürzere Präfixe

Beispiel für das Alignment von Präfixen

Ziel: Berechne optimales Alignment von $s = \text{ATG}$ und $t = \text{TAG}$
Unterscheide **drei Fälle** bezüglich der letzten Spalte des Alignments:

AT	G
TA	G

$d_{\text{edit}}(\text{AT}, \text{TA}) \quad +0$

AT	G
TAG	-

$d_{\text{edit}}(\text{AT}, \text{TAG}) \quad +1$

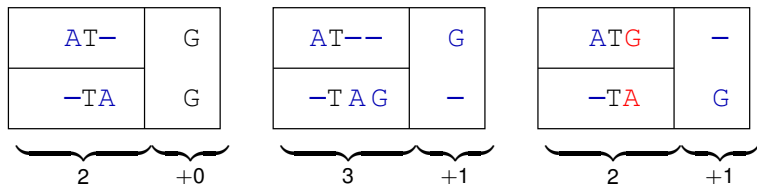
ATG	-
TA	G

$d_{\text{edit}}(\text{ATG}, \text{TA}) \quad +1$

⇒ Berechnung von $d_{\text{edit}}(\text{ATG}, \text{TAG})$ zurückgeführt auf
Berechnung der Edit-Distanz für drei Paare von Präfixen

Beispiel für das Alignment von Präfixen

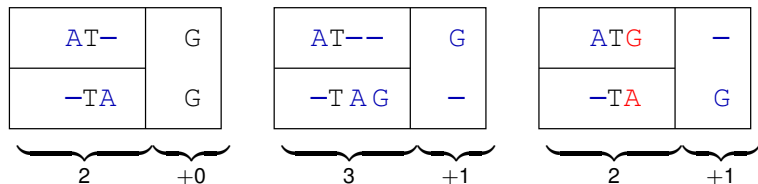
Ziel: Berechne optimales Alignment von $s = \text{ATG}$ und $t = \text{TAG}$
Unterscheide **drei Fälle** bezüglich der letzten Spalte des Alignments:



\Rightarrow Berechnung von $d_{\text{edit}}(\text{ATG}, \text{TAG})$ zurückgeführt auf Berechnung der Edit-Distanz für drei Paare von Präfixen

Beispiel für das Alignment von Präfixen

Ziel: Berechne optimales Alignment von $s = \text{ATG}$ und $t = \text{TAG}$
Unterscheide **drei Fälle** bezüglich der letzten Spalte des Alignments:



\Rightarrow **optimales Alignment** mit Edit-Distanz $d_{\text{edit}}(s', t') = 2$ ist

$$s' = \text{AT-G}$$

$$t' = \text{-TAG}$$

Definition: Der **leere String** λ ist ein String der Länge 0, er ist Präfix von jedem anderen String.

Initialisierung der Berechnung:

Alignment eines nichtleeren Präfixes mit dem leeren String ist eindeutig:

$$\begin{array}{ccc} s_1 & s_2 & \dots & s_j \\ - & - & \dots & - \end{array} \quad \text{oder} \quad \begin{array}{ccc} - & - & \dots & - \\ t_1 & t_2 & \dots & t_j \end{array}$$

Kosten: $d_{\text{edit}}(s_1 \dots s_j, \lambda) = d_{\text{edit}}(\lambda, t_1 \dots t_j) = j$

Beispiel zur Berechnung der Edit-Distanz

s \ t	0	C 1	C 2	T 3	G 4
0	0	1	2	3	4
A 1	1				
C 2	2				
T 3	3				
T 4	4				
G 5	5				

Initialisierung

Beispiel zur Berechnung der Edit-Distanz

s \ t	0	C	C	T	G
0	0	1	2	3	4
A 1	1	1			
C 2	2				
T 3	3				
T 4	4				
G 5	5				

Bestimme $d_{\text{edit}}(s_1, t_1)$:
Lücke in t einfügen

$$\begin{array}{c|c}
 - & A \\
 \underbrace{C} & \underbrace{-} \\
 d_{\text{edit}}(\lambda, t_1) & +1
 \end{array}$$

Lücke in s einfügen

$$\begin{array}{c|c}
 A & - \\
 \underbrace{-} & \underbrace{C} \\
 d_{\text{edit}}(s_1, \lambda) & +1
 \end{array}$$

Mismatch einfügen

$$\begin{array}{c|c}
 \lambda & A \\
 \underbrace{\lambda} & \underbrace{C} \\
 d_{\text{edit}}(\lambda, \lambda) & +1
 \end{array}$$

Minimum bilden

Beispiel zur Berechnung der Edit-Distanz

s \ t	0	C 1	C 2	T 3	G 4
0	0	1	2	3	4
A 1	1	1			
C 2	2	1			
T 3	3				
T 4	4				
G 5	5				

Berechne $d_{\text{edit}}(s_2, t_1)$

Beispiel zur Berechnung der Edit-Distanz

s \ t	0	C 1	C 2	T 3	G 4
0	0	1	2	3	4
A 1	1	1			
C 2	2	1			
T 3	3	2			
T 4	4	3			
G 5	5	4			

Berechne den Rest
von Spalte 1

Beispiel zur Berechnung der Edit-Distanz

s \ t	0	C 1	C 2	T 3	G 4
0	0	1	2	3	4
A 1	1	1	2	3	4
C 2	2	1	1	2	3
T 3	3	2	2	1	2
T 4	4	3	3	2	2
G 5	5	4	4	3	2

$$d_{\text{edit}}(s, t) = 2$$

Darstellung als Matrix

$s \backslash t$	0	1	...	$j-1$	j	...	n
0							
1							
2							
\vdots							
$i-1$							
i							
\vdots							
m							

Letzte Spalte des
Alignments ist

- ↓ Lücke in t
- Lücke in s
- ↘ Match/Mismatch

$$d_{\text{edit}}(s_1 \dots s_i, t_1 \dots t_j) = \min \left\{ \begin{aligned} & d_{\text{edit}}(s_1 \dots s_{i-1}, t_1 \dots t_j) + 1, \\ & d_{\text{edit}}(s_1 \dots s_i, t_1 \dots t_{j-1}) + 1, \\ & d_{\text{edit}}(s_1 \dots s_{i-1}, t_1 \dots t_{j-1}) + p(s_i, t_j) \end{aligned} \right\}$$

Analyse der Laufzeit

Eingabe: Zwei Strings $s = s_1 \dots s_m$ und $t = t_1 \dots t_n$

Initialisierung der Ränder: $m + n + 2$ Operationen

```
for  $i = 0$  to  $m$  do  
     $M(i, 0) = i$   
for  $j = 0$  to  $n$  do  
     $M(0, j) = j$ 
```

Füllen der Matrix: $4 \cdot m \cdot n$ Operationen

```
for  $i = 1$  to  $m$  do  
    for  $j = 1$  to  $n$  do  
         $M(i, j) := \min\{M(i - 1, j) + 1,$   
                         $M(i, j - 1) + 1,$   
                         $M(i - 1, j - 1) + p(s_i, t_j)\}$ 
```

Ausgabe: $d_{\text{edit}}(s, t) = M(m, n)$

Analyse der Laufzeit

Laufzeit: ungefähr $4n^2$ für zwei Strings der Länge n

n	10	50	100	300	10 000
$10n$	100	500	1 000	3 000	100 000
$4n^2$	400	10 000	40 000	360 000	400 000 000
n^3	1 000	125 000	1 000 000	27 000 000	13 Ziffern
3^n	59 049	24 Ziffern	48 Ziffern	143 Ziffern	4 772 Ziffern

Anschaulich: Vergleich zweier Gene (Länge in der Grössenordnung von 10 000) braucht

- 100 MB Platz und
- < 1 Minute Zeit

Bestimmung des optimalen Alignments

t \ s	0	C 1	C 2	T 3	G 4
0	0	1	2	3	4
A 1	1	1	2	3	4
C 2	2	1	1	2	3
T 3	3	2	2	1	2
T 4	4	3	3	2	2
G 5	5	4	4	3	2

$s' = \text{A C T T G}$

$t' = \text{C C T - G}$

Problem bei Strings ungleicher Länge

Beispiel: Betrachte die Strings $s = \text{TAAGGT}$ und $t = \text{AGTTTATAGCCTGGT}$

Ein **optimales Alignment** mit Edit-Distanz $d_{\text{edit}}(s', t') = 9$ ist

$$\begin{array}{l} s' = \text{—— TA—A—— GGT} \\ t' = \text{AGTTTATAGCCTGGT} \end{array}$$

Kompakter und besser biologisch motiviert ist

$$\begin{array}{l} s'' = \text{—— TA—AGG—T——} \\ t'' = \text{AGTTTATAGCCTGGT} \end{array}$$

mit Edit-Distanz $d_{\text{edit}}(s'', t'') = 10$

Idee: Lücken am Beginn und Ende des kürzeren Strings nicht mitzählen

Umsetzung (für $|s| < |t|$):

- erste Zeile der Matrix mit Nullen initialisieren
⇒ Lücken vor dem Lesen von s_1 kostenlos
- Minimum der Werte in der letzten Zeile liefert das Ergebnis
⇒ Lücken nach dem Lesen von s_m kostenlos
- Rest des Algorithmus wie vorher

Beispiel

s \ t	0	C	C	A	C	T	T	T	G	T
	0	1	2	3	4	5	6	7	8	9
0	0	0	0	0	0	0	0	0	0	0
C 1	1	0	0	1	0	1	1	1	1	1
A 2	2	1	1	0	1	1	2	2	2	2
G 3	3	2	2	1	1	2	2	3	2	3
T 4	4	3	3	2	2	1	2	2	3	2

$s' = -CA\color{red}GT-----$

$t' = \color{blue}CCA\color{red}CTTT\color{blue}GT$

- **Lokales Alignment** (Smith und Waterman, 1981): Finde **Teilstrings** mit maximaler Ähnlichkeit, ignoriere dabei nicht passende Anfangs- und Endstücke
- Lücken stärker gewichtet als Mismatches (kommen auch in der Natur seltener vor)
- Verschiedene **Gewichtung von Mismatches**, zum Beispiel entsprechend der chemischen Ähnlichkeit der Aminosäuren beim Vergleich von Proteinsequenzen
- Zusätzliche Kosten für das **Öffnen** einer Sequenz von Lücken

⇒ leichte Abwandlungen des Algorithmus lösen auch diese Probleme

Problem: Quadratische Laufzeit ist für sehr lange Strings (zum Beispiel ganze Genome) zu gross

Ausweg: Finde Algorithmen mit **linearer Laufzeit**, die ein gutes (aber nicht notwendigerweise optimales) Alignment finden

Ideen:

- Beschränke die Anzahl der vorkommenden Lückensymbole
⇒ Berechne nur Streifen konstanter Breite der Matrix um die Mitteldiagonale herum
- Finde zunächst kurze gemeinsame Teilstrings, erweitere diese mit dynamischer Programmierung an beiden Enden, setze diese Teil-Alignments zusammen

- Effizienter Algorithmus für den Vergleich von DNA-Sequenzen
- Modellierung biologischer Fragestellungen als Informatik-Problem
- Algorithmische Technik der dynamischen Programmierung
- Laufzeitanalyse von Algorithmen