

Skript zur Vorlesung Bioinformatik I
gelesen Prof. Stadler

Sven Findeiß

11. Januar 2005

Inhaltsverzeichnis

1 Was ist Bioinformatik?	4
1.1 Informationsquellen	4
1.2 Problemstellung	4
2 Alignments	4
2.1 DNA – Sequenzen	4
2.2 RNA – Sequenzen	4
2.3 Proteine	4
2.4 Begriff der Homologie	4
2.4.1 Alignment	5
2.4.2 Editieroperationen	5
2.4.3 Forderung an Abbildungen	5
2.4.4 Drei Editiermöglichkeiten	5
2.5 Kostengünstiges Alignment	6
2.5.1 schlechte Möglichkeit	6
2.5.2 Bessere Möglichkeit	6
2.5.3 Needleman - Wunsch - Algorithmus	8
2.6 Wiederholung	9
2.6.1 globale Alignments und Needleman – Wunsch – Algorithmus	9
2.6.2 lokale Alignments und Smith – Waterman – Algorithmus	10
2.7 Wiederholungs – Treffen („Repeat – Matches,“)	10
2.7.1 Ziel	10
2.7.2 Beispiel	10
2.8 Overlap – Matches	11
2.8.1 Problem	11
2.9 Gaps	11
2.9.1 Gotoh – Algorithmus (Gotoh 1988)	11
2.9.2 Line – Space – Algorithmus (Hirschberg 1975, Myers Miller 1988)	12
2.10 Scoring Funktionen für paarweise Alignments	13
2.10.1 Vergleich von zwei Modellen	15
2.11 Multiple Alignments (Bei Fragen Mail an dress@mis.mpg.de)	17
2.11.1 Clustal W	17
2.11.2 DIALIGN (Burkhard Morgenstern, Thomas Werner)	18
2.12 Progressives Alignment	19
2.12.1 Dynamic Programming für $N > 2$ Sequenzen	19
2.12.2 Alignieren von 2 Alignments	20
2.12.3 “Sum of Peirs“ Score	21
2.13 Probabilistische paarweise Alignments	24
2.13.1 Stochastisches Backtracking	28
2.14 Suboptimale Alignments	29
3 Phylogenetische (Stamm-) Bäume	29
3.1 Distanz Messung	30
3.2 Genetische Distanz	30
3.2.1 Bedingungen	30

3.2.2	Ziel	30
3.3	Ziel aller Konstruktionen	34
3.4	Splits	36
3.4.1	Maximum Porsimoney Problem	39

1 Was ist Bioinformatik?

1.1 Informationsquellen

- **Thema der Bioinformatik I**

Genbank (NCBI) → fr Sequenzinformation

- **Thema der Bioinformatik II**

PDB (Brookhaven) → 3D – Strukturinformation über Proteine (u. wenige RNA – Moleküle)

- aus Röntgenstrahlungsanalysen
- aus Kernspinresonanz

- Heranziehen von bekannten Daten aus bekannten Sequenzen
- Konzept der Homologie

1.2 Problemstellung

String -to -String EDIT Problem

Nachweis bzw. Wissen über Homologie (Lösung über Algorithmen)

2 Alignments

2.1 DNA – Sequenzen

Wörter ber das Alphabet A,C,G,T

Achtung! da zweisträngig, müssen beide Stränge kontrolliert werden

$N = A/G/C/T$; $R = G/A$; $Y = T/C$

2.2 RNA – Sequenzen

Wörter über das Alphabet A,U,G, C

$Y = U/C$

2.3 Proteine

20 + 1 (Cystein; Selenocystein) + 1 Aminosäuren (Alphabet)

Definition Alphabet – A:

- Wort endlicher Länge
- ohne ϵ

2.4 Begriff der Homologie

Man spricht von homolog, wenn sich die 2 Sequenzen von einem gemeinsamen Verfahren ableiten lassen.

2.4.1 Alignment

```
A U G C G
| | | | |
A U G C G
```

2.4.2 Editieroperationen

Ein Alignment kann durch drei Editieroperationen verändert werden.

1. Deletion (Löschen einer Base aus der Sequenz)
2. Insertion (Einfügen einer Base in die Sequenz)
3. Substitution ($A \leftrightarrow U$)

2.4.3 Forderung an Abbildungen

```
k < l
|   |
u < v
```

Wenn (k,u) , (l,v) matches im Alignment und $k < l \rightarrow u < v$.

2.4.4 Drei Editiermöglichkeiten

1. Substitution ($A \leftrightarrow U$)

Kosten:

$$\sigma(i, j); i, j \in A$$

$$\sigma(i, i) = 0$$

$$\sigma(i, j) > 0; i \neq j$$

$$\sigma(i, j) = \sigma(j, i)$$

2. Insertion

$$\delta(i); i \in A$$

3. Deletion

$$\delta(i); i \in A$$

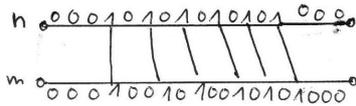
Folgende zwei Regeln müssen gelten:

1. $\sigma(i, j) \leq \delta(i) + \delta(j)$
2. $\sigma(i, j) + \sigma(j, k) \geq \sigma(i, k)$

2.5 Kostengünstiges Alignment

2.5.1 schlechte Möglichkeit

- 2 Sequenzen mit Länge n und m
 - Ausprobieren aller Alignments
 - günstigstes wird ausgewählt

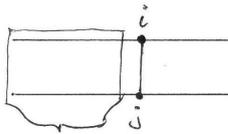


→ wieviele Alignments gibt es mit k Substitutionen

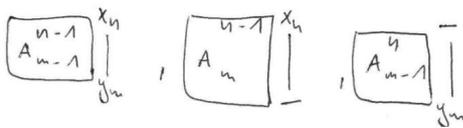
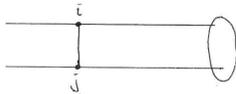
→ $\binom{n}{k} \cdot \binom{m}{k}$ → Produkt ist #Alignment mit genau k Substitutionen

$$\# \text{Alignments} = \sum_k^{\min(n,m)} \binom{n}{k} \cdot \binom{m}{k} = \binom{n+m}{n}$$

2.5.2 Bessere Möglichkeit



$$\text{cost} \left(A_{j-1}^{i-1} \right) + \text{cost} \left(A_{j-1}^{i+1} \right) + \sigma(i, j) \text{ (so dass (i,j) Substitution)}$$



$$\text{cost} \left(A_{m-1}^{n-1} \right) + \sigma(x_n, y_m), \text{cost} \left(A_m^{n-1} \right) + \delta(x_n), \text{cost} \left(A_{m-1}^n \right) + \delta(y_m)$$

$$\min \left(\text{mincost} \left(A_{m-1}^{n-1} \right) + \sigma(x_n, y_m), \text{mincost} \left(A_m^{n-1} \right) + \delta(x_n), \text{mincost} \left(A_{m-1}^n \right) + \delta(y_m) \right)$$

$$\begin{aligned}
D_{ij} &:= \text{mincost}(A_j^i) \\
&\rightarrow \text{mincost}(A_m^n) = D_{nm} \\
&\rightarrow D_{nm} = \min(D_{n-1,m-1} + \sigma(x_n, y_m); D_{n-1,m} + \delta(x_n); D_{n,m-1} + \delta(y_m)) \\
D_{0m} &= \sum_{j=1}^m \delta(y_j) \\
D_{n0} &= \sum_{i=1}^n \delta(x_i)
\end{aligned}$$

1. $D_{00} = 0$
2. for $i = 1 \dots n$ $D_{i0} = D_{i-1,0} + \delta(x_i)$
3. for $j = 1 \dots m$ $D_{0j} = D_{0,j-1} + \delta(y_j)$
4. Hamming Distanz for $i = 1 \dots n$;
for $j = 1 \dots m$;

$$D_{ij} = \min(D_{i-1,j-1} + \sigma(x_i, y_j); D_{i-1,j} + \delta(x_i); D_{i,j-1} + \delta(y_j))$$

Aufwand:

(1-3)

$O(n)$, $O(m)$ worst = best Case

(4)

$O(n \cdot m) \rightarrow$ quadratisch $\approx O(n^2)$

Vorteil:

linearer Speicheraufwand, da nur vorhergehende Zeile q als Arbeitszeile gemerkt werden muss

Verbesserung:

- unelegant
 - merken, welche der 3 Alternativen im Schritt vorher die geringsten Kosten produziert hat
 - Pfad von (0,0) zu (n,m) beschreibt Alignment \rightarrow Knoten ist match
 - ABER: Pfadspeicherung nötig \rightarrow unelegant
- eleganter
 - Start bei (n,m)

- Anschauen der 3 Möglichkeiten, ob

$$D_{i-1,j-1} + \sigma(x_i, y_j) = D_{i,j} \quad \text{oder} \quad D_{i-1,j} + \delta(x_i) = D_{i,j} \quad \text{oder} \quad D_{i,j-1} + \delta(y_j) = D_{i,j}$$

- Besser, da: i und j werden in jedem Schritt kleiner

Backtracking:

Unter Backtracking versteht man, dass wieder Herausfinden, welchen Weg man gegangen ist um zur richtigen Lösung gekommen zu sein.

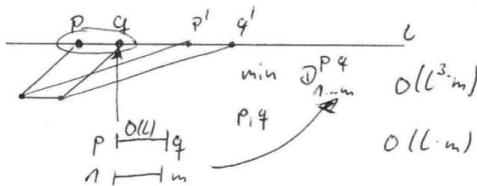
2.5.3 Needleman - Wunsch - Algorithmus

Reflexion, was bis hier gemacht wurde:

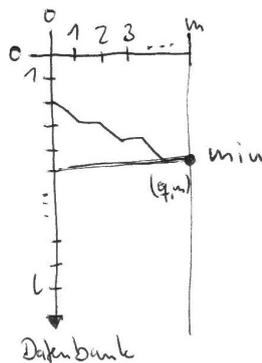
- homologe Alignments wurden verglichen

Jetzt

1. sehr langes Stück vs. sehr kuzes Stück



$D_{p-1,0} = \emptyset \rightarrow$ hat keine Kosten



$$D_{q,m} = \min D_{q',m}$$

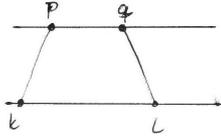
q liefert letzte Substitution, die ungleich ist

in letzter Spalte steht optimales Alignment, d.h. min. Eintrag wird gesucht

Backtracking beginnt bei min.

\rightarrow Matrix ermöglicht auch Suche nach Kosten unterhalb eines Schwellenwertes

2. Suche des besten Ausschnitt – matches



Scoring – Funktion hier besonders wichtig

→ Änderung der Kosten: $\sigma(x, x) < 0$

perfect match bekommt also negative Kosten, da sonst leeres Alignment immer beste Lösung

$$D_{p,k} =$$

normaler Alignment

$$D_{p-1,k-1} + \sigma(x_p, y_k); D_{p-1,k} + \delta(x_p); D_{p,k-1} + \delta(y_k)$$

neues Alignment

$$\sigma(x_p, y_k)$$

	A	G	C	C	A
C	1	1	-1	-1	1
C	1	1	-1	-2	0
C	1	1	-1	-2	-1

A G C C A
 | |
 C C C

- bestes Alignment durch Backtracking
 - start bei negativster Zahl
 - Rücklauf
 - bis zur ersten pos. Zahl
 - fertig

2.6 Wiederholung

2.6.1 globale Alignments und Needleman – Wunsch – Algorithmus

Sequenz $x = x_1, \dots, x_m$ $y = y_1, \dots, y_n \in \Sigma^*$

$$F(i, j) := \max(F(i-1, j-1) + s(x_i, y_j); F(i-1, j) - d; F(i, j-1) - d)$$

$F(i, j)$ = Summe für Alignments zwischen x_1, \dots, x_i und $y_1 \dots y_i$

Ausgang: (0,0)

max, da Interpretation von s und d unterschiedlich zur letzten Vorlesung

→ $F(m, n)$ groß (max) → x und sehr ähnlich (Score)

→ letztes Mal: $F(m, n)$ klein (min) → kurze Editierdistanz

2.6.2 lokale Alignments und Smith – Waterman – Algorithmus

$$F(i, j) := \max(0; F(i - 1, j - 1) + s(x_i, y_j); F(i - 1, j) - d; F(i, j - 1) - d)$$

$$s := \max_{i,j} F(i, j)$$

2.7 Wiederholungs – Treffen („Repeat – Matches,,)

2.7.1 Ziel

Partitioniere x in Regionen, die mit Teilsequenzen von y überdeckt werden.

2.7.2 Beispiel

x = HEAGAWGHEE

y = PAWHEAE

→

HEAGAWGHEE

HEA-AW-HE-

Rekurrenzen (Funktionen, die von sich selbst abhängen = Rekursion) für dynamische Programmierung

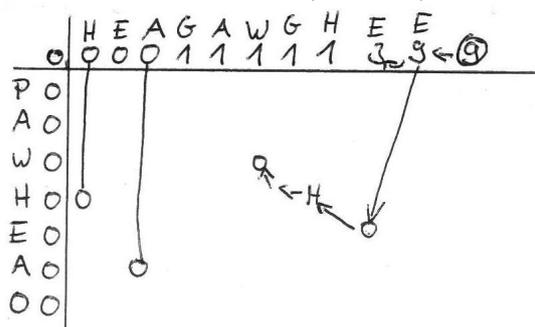
Parameter T: Threshold, der mind – Score für matchende Blöcke

$$F(0, 0) = 9$$

$$F(i, 0) = \max(F(i - 1, 0); F(i - 1, j) - T) \rightarrow j \in [1; n]$$

$$F(i, j) = \max(F(i, 0); F(i - 1, j - 1) + s(x_i, y_j); F(i - 1, j) - d; F(i, j - 1) - d)$$

NEU: Starte Traceback $F(m + 1, 0)$

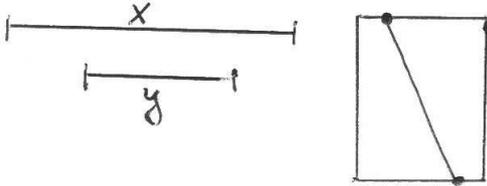


Bsp. in dem Buch besser erklärt

2.8 Overlap – Matches

2.8.1 Problem

Welcher Teilstring \bar{x} von x hat das beste Alignment mit y .



$$F(i, 0) = 0$$

$$F(i, j) = \max(F(i-1, j-1) + s(x_i, y_j); F(i-1, j) - d; F(i, j-1) - d)$$

$$F(0, j) = F(0, j-1) - d$$

2.9 Gaps

$$F(i, j) = \max(F(i-1, j-1) + s(x_i, y_j);$$

$$F(k, j) + \gamma(j-k)[k \in [0 : i-1]];$$

$$F(i, k) + \gamma(i-k)[k \in [0 : j-1]])$$

Gesamtkosten: $O(m \cdot n(n+m)) \rightarrow$ Zeit

2.9.1 Gotoh – Algorithmus (Gotoh 1988)

$O(m \cdot n)$ Zeit für affine Gap – Kosten

Verbesserung:

$$F(i, j) = \max(F(i-1, j-1) + s(x_i, y_j);$$

$$F(k, j) + \gamma(k)[k \in [0 : i-1]];$$

$$F(i, k) + \gamma(k)[k \in [0 : j-1]])$$

offene Gap – Kosten: $\gamma(k) = -d - e \cdot k$

d = Gap – open cost

e = Gap – length cost

Gotoh - Algorithmus
 3 Matrizen ($m \times n$)

- M für Matches
- I_x für Gaps in x
- I_y für Gaps in y

$$M(i, j) = \max(M(i-1, j-1) + s(x_i, y_i); I_x(i-1, j) + s(x_i, y_i); I_y(i, j-1) + s(x_i, y_j))$$

$$I_x(i, j) = \max(M(i-1, j) - d; I_x(i-1, j) - e)$$

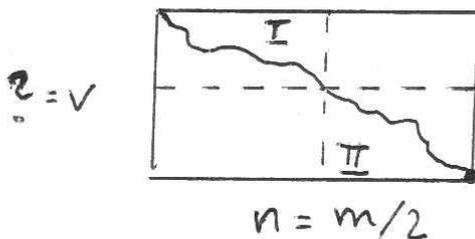
$$I_y(i, j) = \max(M(i, j-1) - d; I_y(i, j-1) - e)$$

$$(2^{18})^2 = 2^{36} = 2^{16} MB = 2^6 GB!$$

2.9.2 Line - Space - Algorithmus (Hirschberg 1975, Myas Miller 1988)

Idee:

Divide - and - Conquer



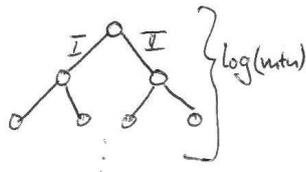
Hilfsgröße $c(i, j)$, fr $i \geq n$

Bei Ausrechnen von $F(i, j)$ bezeichne (i', j') jene Indizes, von denen das Maximum gewählt wird, d.h. $(i', j') \in \{(i-1, j-1), (i-1, j), (i, j-1)\}$

$$c(i, j) = j', \text{ falls } i = n \setminus \setminus \\ c(i', j'), \text{ sonst} \setminus \setminus$$

$c(m, n)$ liefert v . An Position (u, v) spalte Problem auf in Teilprobleme I und II und löse I und II rekursiv.

maximale Rekursive Tiefe:



$\log(m+n)$

Kosten auf jeder Ebene:

$$O(n \cdot m)$$

Gesamtaufwand:

$$O(m \cdot n \cdot \log(m + n))$$

2.10 Scoring Funktionen für paarweise Alignments

$x = x_1, \dots, x_n$ Zeitpunkt 0

↓ t

$y = y_1, \dots, y_m$ Zeitpunkt t

Annahmen:

1. Nur Substitution als Gegenstand der Veränderung
→ $y_m = y_n$, da x und y gleich lang

$x_1, x_2, x_3, \dots, x_{n-1}, x_n$

↓, ↓, ..., ↓, ↓

$y_1, y_2, y_3, \dots, y_{m-1}, y_m$

2. Sequenzpositionen sind unabhängig
3. Uniformität

Alle 3 Annahmen sind nicht in der Realität umsetzbar → max. Ignoranz der Biologie!

- es gibt Deletion und Insertion
- Mutationen treten wesentlich häufiger an gepaarten Positionen auf als an ungepaarten

Mutationsrate: Rate, mit der Mutationen versucht werden durchzusetzen.

Substitutionsrate: Rate, mit der Mutationen auch wirklich durchgesetzt werden → können tatsächlich beobachtet werden, da Selektion bereits erfolgt.

→ für uns erst einmal nur Substitutionsrate von Interesse

→ Prob steht für Wahrscheinlichkeit

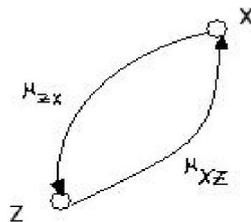
→ $p(x)$ = Wahrscheinlichkeit, dass x auftritt

→ $p(y|x)$ = Übergangswahrscheinlichkeit

$$\begin{aligned}
\text{Prob}(y_1, y_2, \dots, y_n; t | x_1, x_2, \dots, x_n; 0) &= \prod_{i=1}^n \text{Prob}(y_i; t | x_i; 0) \\
\text{Prob}(y_1, \dots, y_n | t) &= \sum_{x_1, \dots, x_n} (y_1, \dots, y_n; t | x_1, \dots, x_n; 0) \cdot p(x_1, \dots, x_n; 0) \\
&= \sum_{x_1, \dots, x_n} \prod_{i=1}^n \text{Prob}(y_i; t | x_i; 0) \cdot \prod_{i=1}^n p(x_i | 0) \\
&= \sum_{x_1} \sum_{x_2} \sum_{x_3} \dots \sum_{x_n} \cdot \prod_{i=1}^n \text{Prob}(y_i; t | x_i; 0) \cdot p(x_i | 0) \\
&= \prod_{i=1}^n \sum_{x_i} \text{Prob}(y_i; t | x_i; 0) \cdot p(x_i | 0) \\
&= \prod_{i=1}^n p(y_i | t) \\
p(y_i | t) &= \sum_{x_i} \text{Prob}(y_i; t | x_i; 0) \cdot p(x_i; 0)
\end{aligned}$$

Änderungsrate:

$$\begin{aligned}
\frac{dp(x)}{dt} &= \sum_{Z \neq X} \mu_{XZ} \cdot p(Z) - \underbrace{\sum_{Z \neq X} \mu_{ZX}}_{\mu_{XX} := -\sum_{Z \neq X} \mu_{ZX}} \cdot p(x) \\
&= \underbrace{\sum_{Z \neq X} \mu_{XZ} \cdot p(Z)}_{\text{Wahrscheinlichkeit X zu erzeugen}} \\
&\quad - \underbrace{\sum_{Z \neq X} \mu_{ZX} \cdot p(X)}_{\text{Wahrscheinlichkeit Z zu erzeugen aus X}}
\end{aligned}$$



$$\sum_{Z \neq X} \mu_{ZX} \rightarrow \mu_{XX} := \sum_{Z \neq X} \mu_{ZX} = \sum_Z \mu_{XZ} \cdot p(Z)$$

$$\frac{d}{dt} \cdot p = \mu \cdot p \text{ mit } A \cdot B; p(0) = p_0$$

Nebenrechnung zu Radioaktiven Zerfall:

$$u = a \cdot u \qquad u(t) = K \cdot e^{at}$$

$$u'(t) = K \cdot a e^{at}$$

$$u(0) = K \cdot e^0 = K = u_0$$

$$u(t) = e^{at} \cdot u_0$$

Nebenrechnung zu Taylor Polynom:

$$e^{at} = 1 + \frac{at}{1!} + \frac{1}{2!}(at)^2 + \frac{1}{3!}(at)^3 + \dots$$

$$p(t) = e^{\mu \cdot t} \cdot p_0$$

$$e^{\mu \cdot t} = I + \frac{t}{1!}M + \frac{t^2}{2!}M^2 + \frac{t^3}{3!}M^3 + \dots$$

$$\underbrace{[e^{Mt}]_{KL}} \quad := \delta_{KL} + \frac{t}{1!}M_{KL} + \frac{t^2}{2!}(M^2)_{KL} + \dots$$

Eintrag na der KL-ten Stelle

$$[e^{Mt}]_{KL} \rightarrow \text{mit } \delta = 1, \text{ fr } L = K; 0, \text{ sonst}$$

→ konvergiert für jedes beliebiges t → Nachweis, dass e^{Mt} funktioniert

$$p(t) = e^{Mt} \cdot p_0$$

$$P_X(t) = \sum_Z [e^{Mt}]_{XZ} \cdot p_0(Z)$$

→ einfachstes Modell nach dem man Sequenzen evolvieren kann

2.10.1 Vergleich von zwei Modellen

Modell I:

“evolviert“ $x \xrightarrow{t} y$

Modell II:

zufällig

Frage:

Was ist wahrscheinlicher?

$$p^{II}(x, y) = p(x) \cdot p(y)$$

Oder

$$p^I(x, y) = p(y|x) \cdot p(x)$$

$$\begin{aligned} & \frac{p^I(x, y)}{p^{II}(x, y)} \\ &= \frac{p^I(x, y)}{p(x) \cdot p(y)} \\ &= \prod_i \frac{p^I(x_i, y_i)}{p(x_i) \cdot p(y_i)} \rightarrow \text{odds ratio} \\ & \log \frac{p^I(x, y)}{p(x) \cdot p(y)} \\ &= \log \prod_i \frac{p^I(x_i, y_i)}{p(x_i) \cdot p(y_i)} \\ &= \sum_i \log \underbrace{\frac{p^I(x_i, y_i)}{p(x_i) \cdot p(y_i)}}_{\sigma(x_i, y_i)} \rightarrow \text{log odds ratio} \end{aligned}$$

$\sigma(x_i, y_i) > 0$ Wenn x_i und y_i häufiger auftreten als sie im zufälligen Versuch erwartet werden.

- positiv bei Substitution
- negativ bei Deletion, Insertion

$$\begin{aligned} & \frac{p^I(x, y)}{p(x) \cdot p(y)} \\ &= \frac{[e^{tM}]_{XY} \cdot p(y)}{p(x) \cdot p(y)} \quad \rightarrow \sigma(x, y) = \log \frac{[e^{tM}]_{XY}}{p(x)} \end{aligned}$$

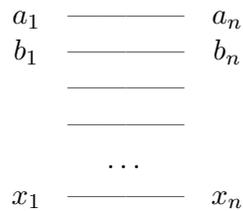
Wahrscheinlichkeit, dass an einer bestimmten Position eine Substitution statt gefunden hat:

$$\begin{aligned} d &= \sum_Y \sum_{X \neq Y} p(x|y) \\ d \cdot N &= \text{Erwartungswert fr Substitutionen} \\ d &= \sum_Y \sum_{X \neq Y} [e^{tm}]_{XY} \end{aligned}$$

2.11 Multiple Alignments (Bei Fragen Mail an dress@mis.mpg.de)

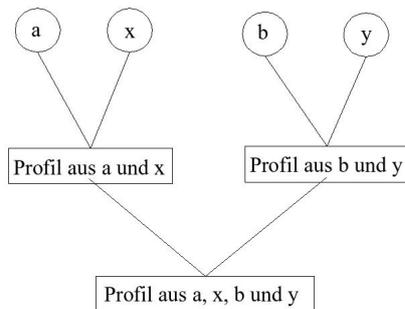
2.11.1 Clustal W

Clustal W ist das bekannteste Multiples Alignment Programm.



Funktion:

- alle Sequenzen werden erst einmal paarweise aligniert → es entsteht der **guide tree**



Entstehung eines Profils

a_1	a_2	a_3	a_4
—	x_1	—	x_2
$\frac{a_1}{2}, \frac{blank}{2}$	$\frac{a_2}{2}, \frac{x_1}{2}$	$\frac{a_3}{2}$	$\frac{a_4}{2}, \frac{x_2}{2}$

Beispiel:

20 Aminosäuren

$R[01]$

$R[1, \dots, 3] = R^n = \{x_1, \dots, x_n\}$

$(\alpha_1, \dots, \alpha_{20}, \alpha_{blank})$

$(x_1, \dots, x_{20}, x_{blank})$

Bedingungen: $(x_i \geq 0; \sum_i x_i = 1) \rightarrow$ Ähnlichkeit zu Wahrscheinlichkeiten!

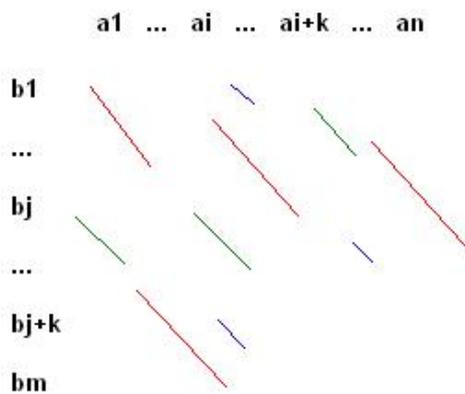
- danach werden Profile aligniert (dieses funktioniert wie bei normalen Sequenzen)

Es entstehen affine Cap – Kosten

$$\alpha + (k - 1)\beta$$

gesucht sind nun möglichst optimale α und β

Die Länge des Matches ist auch entscheidend für seinen Wert



Gesucht:

- System von Diagonalen, die aneinandergelegt eine lange durchgehende Gerade bilden
- der eine Teil der Diagonalen hört dort auf, wo der nächste beginnt
- die Summe ist am größten $-\ln(p(\dots))$ (Wahrscheinlichkeit)

$$v = \begin{pmatrix} a_i, \dots, a_{i+k} \\ b_j, \dots, b_{j+k} \end{pmatrix} \in R^\alpha$$

Sortierung der Geraden nach der Größe.

und dann nimmt man die am besten passenden zusammen und wiederholt diesen Vorgang so oft bis alle Geraden verbunden sind und ein endgültiges Alignment entstanden ist.

2.12 Progressives Alignment

2.12.1 Dynamic Programming für $N > 2$ Sequenzen

$N = 3$:

Man nimmt die drei Sequenzen und aligniert sie genau so wie man 2 Sequenzen alignieren

würde.

$$\begin{array}{l} x_1 \text{---} x_i \\ y_1 \text{---} y_j \\ z_1 \text{---} z_k \end{array}$$

Es ergeben sich folgende drei Fälle.

$$\begin{array}{ccc} \text{Fall 1:} & \text{Fall 2:} & \text{Fall 3:} \\ \begin{pmatrix} i-1 \\ j-1 \\ k-1 \end{pmatrix} \begin{array}{c} x_i \\ | \\ y_j \\ | \\ z_k \end{array} & \begin{pmatrix} i-1 \\ j-1 \\ k-1 \end{pmatrix} \begin{array}{c} x_i \\ | \\ y_j \\ | \\ z_k \end{array} & \begin{pmatrix} i-1 \\ j-1 \\ k-1 \end{pmatrix} \begin{array}{c} x_i \\ | \\ y_j \\ | \\ z_k \end{array} \end{array}$$

$$D_{ijk} = \min \left(\begin{array}{c} \text{Fall 1} \\ D_{i-1,j-1,k-1} + \sigma(x_i, y_j, z_k) \\ \text{Fall 2} \\ D_{i-1,j-1,k} + \sigma(x_i, y_i, --) \\ D_{i-1,j,k-1} + \sigma(x_i, --, z_k) \\ D_{i,j-1,k-1} + \sigma(--, y_j, z_k) \\ \text{Fall 3} \\ D_{i-1,j,k} + \sigma(x_i, --, --) \\ D_{i,j-1,k} + \sigma(--, y_j, --) \\ D_{i,j,k-1} + \sigma(--, --, z_k) \end{array} \right)$$

Allgemein:

Es gibt $2^N - 1$ Fälle

MEM $O(n^N)$

CPU $O(n^N \cdot 2^N)$

→ für alle Fälle mit $N > 3$ nicht praktisch relevant, da der Aufwand unrealistisch hoch wird.

2.12.2 Alignieren von 2 Alignments

\bar{x}_1	\bar{x}_2	\bar{x}_3	...	\bar{x}_n
x_{11}	x_{12}	x_{13}	...	x_{1n}
x_{21}	x_{22}	x_{23}	...	x_{2n}
...
x_{N1}	x_{N2}	x_{N3}	...	x_{Nn}

$N \times n; A \cup \{-\}$

\bar{y}_1	\bar{y}_2	\bar{y}_3	\dots	\bar{y}_n
y_{11}	y_{12}	y_{13}	\dots	y_{1m}
y_{21}	y_{22}	y_{23}	\dots	y_{2m}
\dots	\dots	\dots	\dots	\dots
y_{M1}	y_{M2}	y_{M3}	\dots	y_{Mm}

$M \times m; A \cup \{-\}$

\bar{x} = Vektor x
 – steht für Cap

$$\bar{x}_k \in (A \cup \{-\})^N \setminus \{(-, -, \dots)\}$$

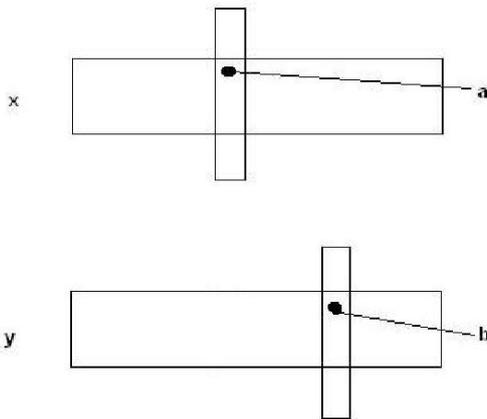
$$\bar{y}_l \in (A \cup \{-\})^M \setminus \{(-, \dots, -)\}$$

⇒ Spalten werden als Elemente über einem erweiterbaren Alphabet aufgefasst
 ⇒ man kann nun die einzelnen Vektoren matchen bzw. mismatchen
 ⇒ Problem ist die Kostenfunktion

2.12.3 “Sum of Peirs“ Score

1. match

$$\sigma(\bar{x}_k, \bar{y}_l) = \sum_{a=1}^N \sum_{b=1}^M \sigma(x_{ak}, y_{bl})$$



Mittelwert:

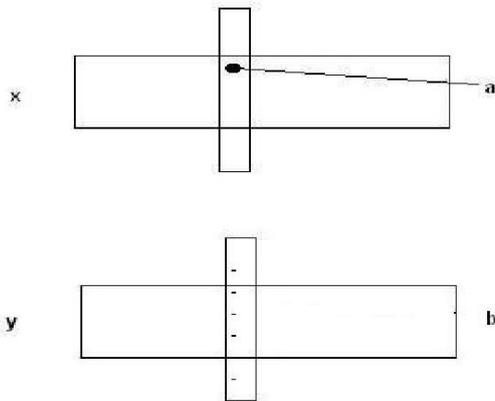
$$\sigma(x_k, y_l) = \frac{1}{N} \cdot \frac{1}{M} \cdot \sum \sigma(x_{ak}, y_{bl})$$

2. mismatch

$$\sigma(\bar{x}_k, -) = M \cdot \sum_{a=1}^N \underbrace{\sigma(x_{ak}, -)}_{\delta(x_{nk})}$$

→ man fügt in das Alignment einen Vektor von Gaps ein

$$\sigma(-, \bar{y}_l) = M \cdot \sum_{b=1}^M \underbrace{\sigma(-, y_{bl})}_{\delta(y_{mk})}$$



$$\sigma(-, -) = \emptyset (\text{keine Kosten})$$

Es gibt zwei Arten von Gaps bei affinen Gap – Kosten:

A U G - - - - -

- extention gap

-- gap open

$$\sigma(--, -) = \textit{konstant} = 0$$

A U G - - - C

A - - - - - C

⇒ Problem, wenn man sehr viele Alignments alignieren will

Bessere Möglichkeit?

$$\sum_{a=1}^N \sum_{b=1}^M \sigma(x_{ak}, y_{bk})$$

NR: (wie oft kommt α oder β in der k-ten oder l-ten Spalte vor?)

$$A' = A \cup \{--\}$$

$$n_k(\alpha) = |\{a | x_{ak} = \alpha \in A'\}|$$

$$m_l(\beta) = |\{b | y_{bl} = \beta \in A'\}|$$

$$= \sum_{\alpha \in A'} \sum_{\beta \in A'} n_k(\alpha) \cdot n_l(\beta) \cdot \sigma(\alpha, \beta)$$

$\overline{x_k}$

| bildet Spalte im fertigen Ergebnis $\overline{z_q}$

$\overline{y_l}$

$$n_q''(\alpha) = n_k(\alpha) + m_l(\alpha) \dots$$

wie oft kommt ein Buchstabe im Alignment vor \Rightarrow Profile

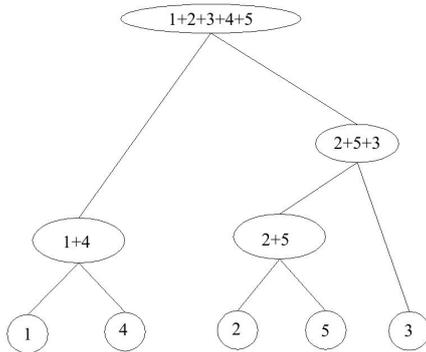
$$O(n) \cdot O(|A|) \cdot N$$

Anmerkung: Ein Profil bekommt nie mehr Zeilen, es wird nur länger.

\Rightarrow man fasst nun an immer paarweise Alignments (Sequenzen später Profile) zusammen.

\Rightarrow Frage ist aber welche Alignments sollten zusammen gefasst werden.

Guide -- Tree:



möglichst gleiche/ ähnliche Alignments werden zusammengefasst

$$O(n) \cdot N \cdot O(|A|^2) \cdot O(N)$$

2.13 Probabilistische paarweise Alignments

\mathring{A} = paarweises Alignment

$$\begin{aligned}
 Score(\mathring{A}) &= \sum_{i=1}^n Score\left(\begin{matrix} x_{1i} \\ y_{2i} \end{matrix}\right) \\
 Sc\left(\begin{matrix} x_{1i} \\ y_{2i} \end{matrix}\right) &= K \cdot \log \frac{p(x_{1i}, x_{2i})}{p(x_{1i}) \cdot p(x_{2i})} \\
 &= K \cdot \log \prod_i \frac{p(x_{1i}, x_{2i})}{p(x_{1i}) \cdot p(x_{2i})} \\
 &= K \cdot \log \frac{p(x_1, x_2)}{p(x_1) \cdot p(x_2)} \\
 p(x_1, x_2) &= Prob(\mathring{A}) \\
 &= K \cdot \log Prob(\mathring{A}) - K \cdot \log(p(x_1)) - K \cdot \log(p(x_2)) \\
 K \cdot \log(p(x_1)) - K \cdot \log(p(x_2)) &= C_0 \\
 &= -C_0 + K \cdot \log Prob(\mathring{A})
 \end{aligned}$$

$$\begin{aligned}
 \log Prob(\mathring{A}) &= \frac{Score(\mathring{A}) + C_0}{K} \\
 C'_0 &= \frac{C_0}{K} \\
 \frac{1}{k} \cdot Score(\mathring{A}) + C'_0 & \\
 Prob(\mathring{A}) &= e^{C'_0} \cdot e^{\frac{1}{K} \cdot Score(\mathring{A})}
 \end{aligned}$$

Wenn man alle Möglichkeiten zu alignieren aufsummiert erhält man den Wert 1 (Wahrscheinlichkeiten eines Ereignis).

$$\sum_{\mathring{A}} Prob(\mathring{A}) = 1$$

$$\begin{aligned}
 Prob(\mathring{A}) &= e^{C'_0} \sum_{\mathring{A}} e^{\frac{Score(\mathring{A})}{K}} \\
 e^{C'_0} &= \frac{1}{Z} \\
 \sum_{\mathring{A}} e^{\frac{Score(\mathring{A})}{K}} &= Z
 \end{aligned}$$

$Z :=$ Zustandsmenge des Alignment

$$\underline{Prob(\dot{A}) = \frac{1}{Z} \cdot \exp\left(\frac{Score(\dot{A})}{K}\right)}$$

$\Omega \subseteq \dot{A}$ | Alignments von $x_1 \dots x_n$

$\Omega^{(p,q)} \subseteq \dot{A}$ | Alignments von x und y , so dass x_p und y_q gematcht ist

$$\begin{array}{cccc} x_1 \dots x_{p-1} & x_p & x_{p+1} \dots x_n & \\ & | & & \\ y_1 \dots y_{q-1} & y_q & y_{q+1} \dots y_m & \end{array}$$

$$\begin{aligned} Prob(\dot{A} \in \Omega) &= \sum_{\dot{A} \in \Omega} Prob(\dot{A}) \\ &= \frac{1}{Z} \sum_{\dot{A} \in \Omega} \exp\left(\frac{Score(\dot{A})}{K}\right) \\ &= \frac{\sum_{\dot{A} \in \Omega} \exp\left(\frac{Score(\dot{A})}{K}\right)}{Z} = \frac{Z(\Omega)}{Z} \end{aligned}$$

Anmerkung:

$Prob(\dot{A} \in \Omega^{(p,q)})$

- sehr nahe 1 \Rightarrow p und q sind sehr wahrscheinlich gematcht wurden
- sehr nahe 0 \Rightarrow p und q sind sehr wahrscheinlich nicht gematcht wurden

$$Sc \left(\begin{array}{ccc} x_1 \dots x_{p-1} & x_p & x_{p+1} \dots x_n \\ & | & \\ y_1 \dots y_{q-1} & y_q & y_{q+1} \dots y_m \end{array} \right) = Sc \left(\begin{array}{c} x_1 \dots x_{p-1} \\ y_1 \dots y_{p-1} \end{array} \right) + \sigma(x_p, y_q) + Sc \left(\begin{array}{c} x_{p+1} \dots x_n \\ y_{q+1} \dots y_m \end{array} \right)$$

$$\begin{aligned}
\sum_{\mathring{A} \in \Omega^{p,q}} e^{\text{Score}(\mathring{A})} &= \sum_{\mathring{A}'} \sum_{\mathring{A}''} e^{\frac{\text{Score}(\mathring{A}') + \text{Score}(\mathring{A}'') + \sigma(x_p, y_q)}{K}} \\
&= \text{Alignment vom ersten Teil} \\
&= \text{Alignment vom zweiten Teil} \\
&= e^{\frac{\sigma(x_p, y_q)}{K}} \cdot \sum_{\mathring{A}'} e^{\frac{\text{Score}(\mathring{A}')}{K}} \cdot \sum_{\mathring{A}''} e^{\frac{\text{Score}(\mathring{A}'')}{K}} \\
&= e^{\frac{\sigma(x_p, y_q)}{K}} \cdot Z_{p-1, q-1} \\
&= e^{\frac{\sigma(x_p, y_q)}{K}} \cdot \tilde{Z}_{p+1, q+1}
\end{aligned}$$

$$Z(\Omega^{(p,q)}) = Z_{p-1, q-1} \cdot \tilde{Z}_{p+1, q+1} \cdot e^{\frac{\sigma(x_p, y_q)}{K}}$$

Z_n, m := Zustandsmenge über die Summe aller Sequenzen
 \tilde{Z} := Postfixe
 $Z_{n,m} = Z = \tilde{Z}_{1,1}$

$$\text{Prob} \begin{pmatrix} p \\ \vdots \\ q \end{pmatrix} = \frac{Z_{p-1, q-1} \cdot e^{\frac{\sigma(x_p, y_q)}{K}} \cdot \tilde{Z}_{p+1, q+1}}{Z}$$

$$\begin{aligned}
Z_{kl} &= \begin{array}{c} x_1 \text{---} x_k \\ y_1 \text{---} y_l \end{array} \\
&= \begin{array}{c} \text{Fall 1:} \\ x_1 \text{---} x_{k-1} \quad x_k \\ y_1 \text{---} y_{l-1} \quad y_l \\ \text{Fall 2:} \\ x_1 \text{---} x_{k-1} \quad x_k \\ y_1 \text{---} y_l \quad - \\ \text{Fall 3:} \\ x_1 \text{---} x_k \quad - \\ y_1 \text{---} y_{l-1} \quad y_l \end{array} = \begin{array}{c} \text{Fall 1:} \\ \text{Score}(\dot{A}_{k-1,l-1}) + \sigma(x_k, y_l) \\ \text{Fall 2:} \\ \text{Score}(\dot{A}_{k-1,l}) + \sigma(x_k, -) \\ \text{Fall 3:} \\ \text{Score}(\dot{A}_{k,l-1}) + \sigma(-, y_l) \end{array} \\
&= \sum_{\text{Alignments } \dot{A}} e^{\frac{\text{Score}(\dot{A})}{K}} \\
&= \sum_{\text{Alignments } \dot{A} \text{ mit match}} e^{\frac{\text{Score}(\dot{A}_{k-1,l-1})}{K}} \cdot e^{\frac{\sigma(x_k, y_l)}{K}} \\
&+ \sum_{\text{Alignments } \dot{A} \text{ mit delet von } x_k} e^{\frac{\text{Score}(\dot{A}_{k-1,l})}{K}} \cdot e^{\frac{\sigma(x_k, -)}{K}} \\
&+ \sum_{\text{Alignments } \dot{A} \text{ mit insert von } y_l} e^{\frac{\text{Score}(\dot{A}_{k,l-1})}{K}} \cdot e^{\frac{\sigma(-, y_l)}{K}} \\
&= Z_{k-1, l-1} \cdot e^{\frac{\sigma(x_k, y_l)}{K}} + Z_{k-1, l} \cdot e^{\frac{\sigma(x_k, -)}{K}} + Z_{k, l-1} \cdot e^{\frac{\sigma(-, y_l)}{K}}
\end{aligned}$$

Needlman Wunsch Algorithmus

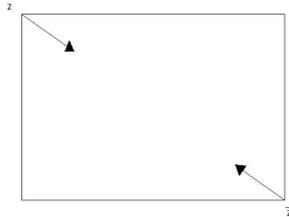
$$M_{kl} = \max \begin{pmatrix} M_{k-1, l-1} + \sigma(x_k, y_l) \\ M_{k-1, l} + \sigma(x_k, -) \\ M_{k, l-1} + \sigma(-, y_l) \end{pmatrix}$$

$$\begin{aligned}
M_{00} &= 0 & Z_{00} &= 1 \\
M_{0l} &= \sum_{j=1}^l \sigma(-, y_j) & e^{\frac{\text{Score}([y_1 y_2 \dots y_l])}{K}} \\
& & &= e^{\sum_{j=1}^l \frac{\sigma(-, y_j)}{K}} \\
M_{k0} &= \sum_{i=1}^k \sigma(x_i, -) & e^{\frac{\text{Score}([x_1 x_2 \dots x_k])}{K}} \\
& & &= e^{\sum_{i=1}^k \frac{\sigma(x_i, -)}{K}}
\end{aligned}$$

$$f_{kl} := \frac{Z_{kl}}{\exp\left(\frac{M_{kl}}{K}\right)}$$

$$\begin{aligned}
\tilde{Z}_{i,j} &= \begin{array}{c} i \quad i+1 \text{ --- --- } -n; \quad i \quad i+1 \text{ --- --- } -n; \quad - \quad i \text{ --- --- } -n \\ | \quad \quad \quad \quad \quad \quad \quad | \quad \quad \quad \quad \quad \quad \quad | \\ j \quad j+1 \text{ --- --- } -m; \quad - \quad j \text{ --- --- } -m; \quad j \quad j+1 \text{ --- --- } -m \end{array} \\
&= \begin{pmatrix} \sigma(x_i, y_j) + \text{Score}(\mathring{A}^{i+1, j+1}) \\ \sigma(x_i, -) + \text{Score}(\mathring{A}^{i+1, j}) \\ \sigma(-, y_j) + \text{Score}(\mathring{A}^{i, j+1}) \end{pmatrix} \\
&= \tilde{Z}_{i+1, j+1} \cdot e^{\frac{\sigma(x_i, y_j)}{K}} + \\
&= \tilde{Z}_{i+1, j} \cdot e^{\frac{\sigma(x_i, -)}{K}} + \\
&\quad \tilde{Z}_{i, j+1} \cdot e^{\frac{\sigma(-, y_j)}{K}}
\end{aligned}$$

$$\tilde{Z}_{n+1, m+1} = 1$$



gesucht:

Wahrscheinlichkeiten die beagen wie wahrscheinlich man aus einer Richtung gekommen ist.

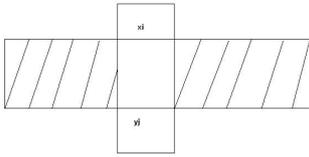
$$\begin{aligned}
Z_{kl} &= Z_{k-1, l-1} e^{\frac{\sigma(x_k, y_l)}{K}} \dots \text{match } \searrow \\
&\quad + Z_{k-1, l} e^{\frac{\sigma(x_k, -)}{K}} \dots \text{del } \downarrow \\
&\quad + Z_{k, l-1} e^{\frac{\sigma(-, y_l)}{K}} \dots \text{insert } \rightarrow
\end{aligned}$$

2.13.1 Stochastisches Backtracking

$$\begin{aligned}
\text{Prob}(\searrow | k, l) &= \frac{Z_{k-1, l-1} e^{\frac{\sigma(x_k, y_l)}{K}}}{Z_{kl}} \rightarrow k-1, l-1 \\
\text{Prob}(\downarrow | k, l) &= \frac{Z_{k-1, l} e^{\frac{\sigma(x_k, -)}{K}}}{Z_{kl}} \rightarrow k-1, l \\
\text{Prob}(\rightarrow | k, l) &= \frac{Z_{k, l-1} e^{\frac{\sigma(-, y_l)}{K}}}{Z_{kl}} \rightarrow k, l-1
\end{aligned}$$

2.14 Suboptimale Alignments

U_{ij} = Score des optimalen Alignments in dem x_i und y_j matchen.

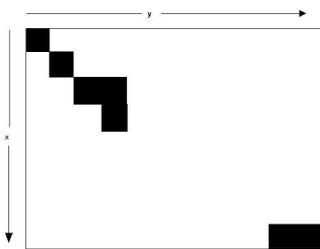


$$M_{i-1,j-1} + \sigma(x_i, y_j) + \tilde{M}_{i+1,j+1}$$

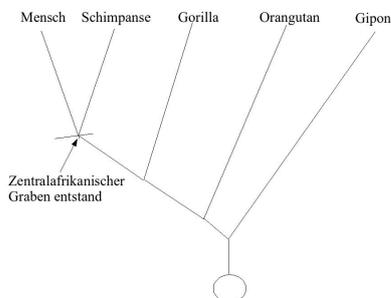
$$\tilde{M}_{i,j} = \max \begin{pmatrix} \tilde{M}_{i+1,j+1} + \sigma(x_i, y_j), \\ \tilde{M}_{i+1,j} + \sigma(x_i, -), \\ \tilde{M}_{i,j+1} + \sigma(-, y_j), \end{pmatrix}$$

Score des optimalen Alignments:

$$M_{nm} = \tilde{M}_{11}$$



3 Phylogenetische (Stamm-) Bäume



D ist die Distanz zwischen Lebewesen.

3.1 Distanz Messung

- man greift sich ein Gen heraus, welches alle zu betrachtenden Lebewesen besitzen (Cytochrom C sehr kurzes Protein) und vergleicht dieses

A	B	C	D	E	F	G
—	—	—	—	—	—	—
—	0	—	—	0	—	—

- man schreibt das Genom des Lebewesens auf und zählt statistisch wie oft welche Base vorkommt
- allgemein:
 - alignieren von Genen
 - Chromosomen oder Genome ungeeignet (Lage der Gene unterschiedlich, Chromosomensatz unterschiedlich)

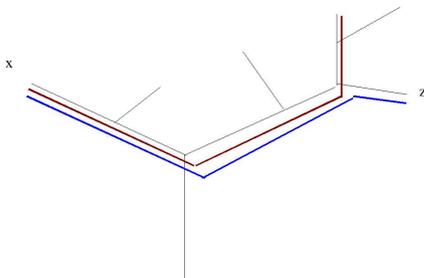
3.2 Genetische Distanz

3.2.1 Bedingungen

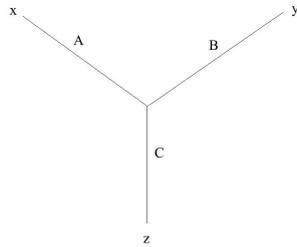
1. $D(x, x) = 0$
2. $D(x, y) > 0$
3. $D(x, y) = D(y, x)$ (Symmetrie)
4. $D(x, y) + D(y, z) \text{ rot} \geq D(x, z) \text{ blau}$ (Dreiecksungleichung)

3.2.2 Ziel

Rekonstruktion eines Baumes, wobei die Kantenlänge die genetische Distanz symbolisiert. Wenn ein Baum existiert so muss die Dreiecksungleichung anwendbar sein.



- bei 3 Punkten gibt es eine eindeutige Lösung (Dreiecksungleichung)

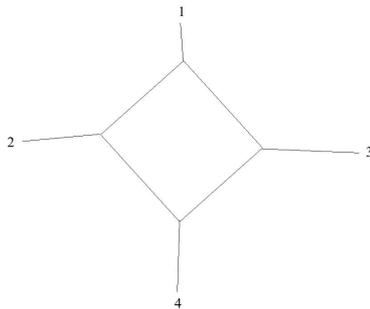


$$A = \frac{D(x, y) + D(x, z) - D(y, z)}{2}$$

$$B = \frac{D(y, x) + D(y, z) - D(x, z)}{2}$$

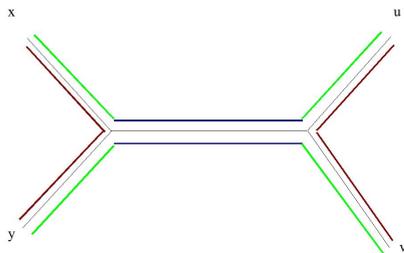
$$C = \frac{D(z, x) + D(z, y) - D(x, y)}{2}$$

- Problem ergibt sich bei 4 Punkten, da es hier keinen eindeutigen Baum gibt



Annahme:

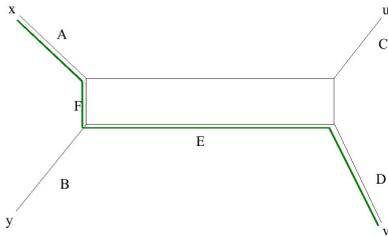
Es gäbe einen Baum x,y,u und v



So würden sich folgende Gleichungen für diesen ergeben:

$$\begin{aligned}
D(x, y) + D(u, v) &= 1 \cdot \text{rot} \\
D(x, u) + D(y, v) &= 2 \cdot \text{blau} + 1 \cdot \text{grün} > 1 \cdot \text{rot} \\
D(v, x) + D(y, u) &= 2 \cdot \text{blau} + 1 \cdot \text{grün}
\end{aligned}$$

Die größten zwei dieser Summen müssen gleich sein



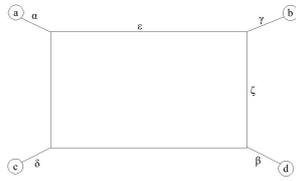
$$\begin{aligned}
A &= \frac{D(x, u) + D(x, y) - D(u, y)}{2} \\
B &= \frac{D(y, v) + D(y, x) - D(x, v)}{2} \\
C &= \frac{D(u, x) + D(u, v) - D(x, v)}{2} \\
D &= \frac{D(v, y) + D(v, u) - D(y, u)}{2} \\
E &= \frac{D(x, v) + D(y, u) - D(x, u) - D(y, v)}{2} \\
E &= \frac{D(x, v) + D(y, u) - D(x, y) - D(u, v)}{2}
\end{aligned}$$

$$A + D + E + F = 2xy$$

Bei Boxstruktur wäre es:

$$(D(x, y) + D(u, v)) \leq (D(x, u) + D(y, v)) < (D(x, v) + D(y, u))$$

$$\begin{aligned}
X &= \{a, b, c, d\} \rightarrow 4 \text{ elementige Metrik} \\
D &: X \times X \rightarrow \mathcal{R} \\
D(x, y) &= xy
\end{aligned}$$



oBdA :

a, b, c, d mit

$$ab + cd \geq$$

$$ac + bd \geq$$

$$ad + bc$$

$$\exists! \alpha, \beta, \gamma, \delta, \epsilon, \zeta \in R$$

→ so dass $ac = \alpha + \epsilon + \gamma$

$$ab = \alpha + \epsilon + \zeta + \beta$$

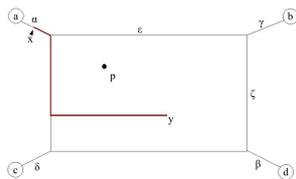
$$\alpha = ac + ad - cd$$

$$\epsilon = ab + cd - (ad + bc)$$

$$\zeta = ab + cd + (ac + bd)$$

Für mehr als 4 Elemente

Manhattan – Matrix/ City – Block – Matrix



Beispiel, dass es für mehr als 4 Punkte immer min. 2 minimale Wege gibt → darum $T(X,D)$

$$xy = D_{cb}((x_1, x_2), (y_1, y_2)) = |y_1 - x_1| + |y_2 - x_2|$$

$$f_p : X \rightarrow R$$

$$x \rightarrow D(p, x)$$

$$p = p' \leftrightarrow f_p = f_{p'}, f_p(x) + f_p(y) \geq xy$$

$f : X \rightarrow R$ ist von der Form $f = f_p$, für die obige Figur

$$1. f(x) + f(y) \geq D(x, y)$$

2. Für alle $x \in X$
 existiert $y \in X$ mit $f(x) + f(y) = D(x, y)$

Sei $T(X, D)$ die Menge aller dieser f

Annahme:

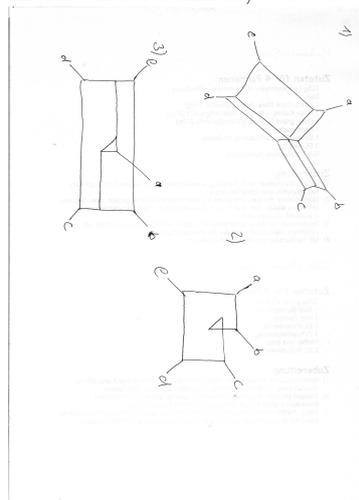
D genügt der 4 – Punkte Bedingung (eine der Strecken ist 0, d.h. so groß wie die größte \rightarrow kann weggelassen werden \rightarrow lässt sich als Baum beschreiben) \rightarrow wenn D ein Baum sein soll, dann muss diese Bedingung für 4 Punkte gelten $\Rightarrow T(X, D)$ ist „Baum,, für D .

3.3 Ziel aller Konstruktionen

$$D \approx D'$$

textnormal wobei D' der 4 Punkte – – Bedingung genügt

Es gibt aber keine schnelle und gut berechenbare Lösung. Wenn man nun eine Lösung sucht merkt man schnell, dass es mehrere Möglichkeiten gibt.



zu 1)

$$\zeta = \frac{1}{2} \min \left(\max \begin{pmatrix} xu + yv \\ xv + yu \end{pmatrix} - (xy + ux) \right)$$

$$x, y \in \{a, e\}$$

$$u, v \in \{b, c, d\}$$

X, D

$A \subseteq X \rightarrow B := X - A \Rightarrow B$ ist Komplement von $A \rightarrow$ Betrachtung zweier disjunkter Teilmengen

$$A \cap B = \emptyset$$

$$\alpha_{A,B}^D := \frac{1}{2} \min_{x,y \in A; u,v \in B} \left(\max \begin{pmatrix} xu + yv \\ xv + yu \\ xy + uv \end{pmatrix} - xy - uy \right) \geq 0$$

$\alpha_{A,B}^D =$ Isolationsindex für den Split A, B

Problem:

Es gibt keinen genauen Splitpunkt. Man erhält aus einer Menge von N Spezies so ergeben sich $\frac{2^N - 2}{2}$ Splitmöglichkeiten. Überwindung nur möglich durch Wahl einer sehr kleinen (polinomialen) Teilmenge

Satz:

Gegeben X (mit $\#X = N$) und $D : X \times X \rightarrow \text{ReineMatriz.}$ ($\#X$ bedeutet $\text{card}(X)$)

$$\#\{A, B\} \in \varphi(X) : \text{EineKana} \alpha_{A,B}^D \geq 0 \leq \binom{N}{2}$$

$$\zeta(X) := \#\{A, B\} \begin{pmatrix} \emptyset \neq A, B \subseteq X \\ A \cap B = \emptyset \\ A \cup B = X \end{pmatrix}$$

wenn ein weiteres Element in einem gegebenen Split einfügen will:

gegeben: $\{A', B'\}; X' \subseteq X; X' \cup \{X\}$

folgt: 2 Möglichkeiten $\{A' + x, B'\}$ oder $\{A', B' + x\}$

Problem wenn immer etwas dazu kommt

Lösung:

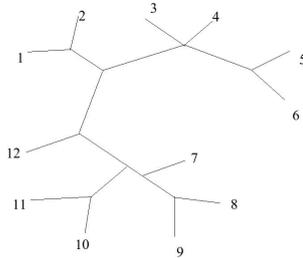
$$N = 4 : 6$$

$$N = 5 \leq 12 \text{ aber auch } \leq \binom{5}{2} = 10 \rightarrow \text{d.h. 2 mal geht etwas schief}$$

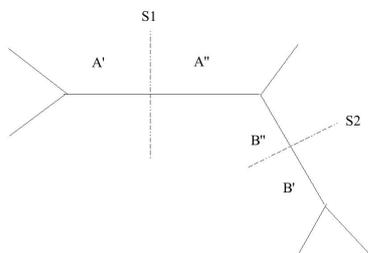
nähert sich an $O(n^5)$ im Endeffekt

3.4 Splits

Baum durch Splits beschrieben:



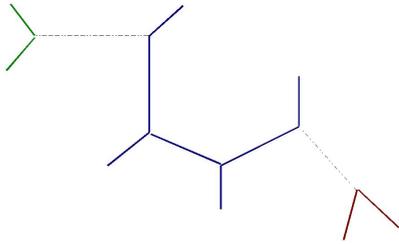
- eine Kante definiert immer einen Split
- $2N - 3$ Splits \rightarrow mehr als phylogenetisch relevant
- die Zerlegung ist eindeutig
- ein Split ist Teilung einer Menge in zwei Teilmengen
- Splitsystem ist eine Menge von Splits auf einer Menge



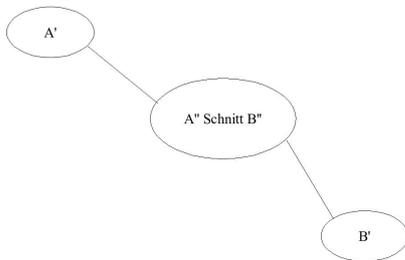
$$\begin{aligned}
 &A' \cap B' \\
 &A' \cap B'' \\
 &A'' \cap B' \\
 &A'' \cap B''
 \end{aligned}$$

Genau einer dieser 4 Durchschnitte muss leer sein.

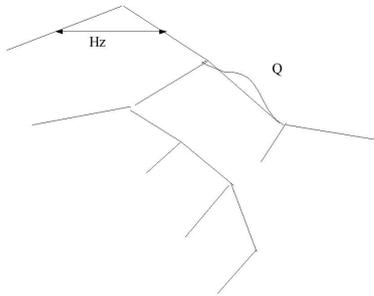
Wenn ein Baum durch zwei Splits (gestrichelte Linien) geteilt wird entstehen immer genau 3 Teilbäume. Im Extremfall sind es Punkte. System von Splits δ auf X (Menge aller Knoten) ist kompatibel, wenn genau einer der Durchschnitte leer ist für je zwei Splits (A', A'') und $(B', B'') \in \delta$.



Baum \Leftrightarrow kompatible Splitsystem



Das einfache Modell der Biologen für die Sequenzevolution:

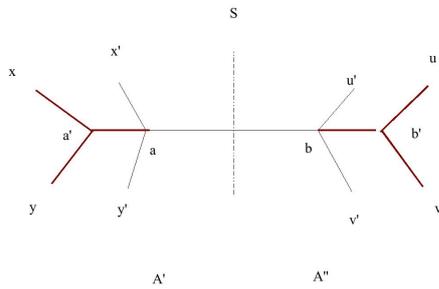


HZ : Horizontaler Gentransfer

Q : Q - Variationen

\Rightarrow es gibt viele „Störungen,, die einen solchen Baum nicht aufbauen lassen

\Rightarrow der Großteil aller Effekte kann aber zu einem solchen Baum approximiert werden



$$D(a', b') = (D(x, u) + D(y, v) - D(x, y) - D(u, b)) \cdot \frac{1}{2}$$

$$\min_{x, y \in A'; u, v \in A''} D(a', b') = D(a, b)$$

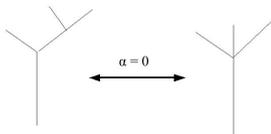
$$\alpha_S = \frac{1}{2} \min_{x, y \in A'; u, v \in A''} (xu + yv - xy - uv)$$

$\alpha_S =$ starker Isolationsindex

$$D(x, v) = \sum_S \alpha_S = \sum_S \delta_S(x, y) \cdot \alpha_S \leq D(x, y) \text{ Gleichheit gilt wenn es ein Baum ist}$$

Problem:

Wenn die Daten nicht exakt Baumartig sind, dann sind die konstruierten Distanzen immer kleiner als die eigentlichen realen Distanzen.



⇒ Splitverfahren sind sehr konservativ

- Kanten die gefunden werden sind richtig
- der Baum ist aber ungenau

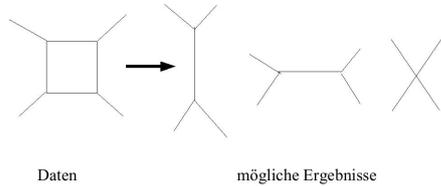
$$D(x, y) := f(T) = \min_{x, y} \sum \left(\sum_{S \in T} \delta_S(x, y) \cdot c_S - D(x, y) \right)$$

c_S bestimmt die Länge

Wie gut passt mein konstruierter Baum auf meine vorhandenen Daten???

Dieses Verfahren gibt immer einen Baum aus egal was die Daten beschreiben.

Bsp :



$B(n)$ = Anzahl der binären Bäume mit n nummerierten Blättern

$B(n + 1) = B(n)$ Anzahl der Kanten von binären Bäumen mit n Blättern

n Blätter, $n - 2$ innere Knoten

→ $2n - 2$ Knoten

→ $2n - 3$ Kanten

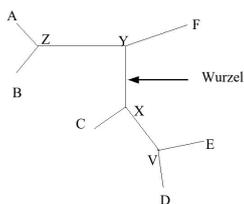
$B(n + 1) = (2n - 3) B(n)$

⇒ zu viele Bäume um sie alle durchzuprobieren

n	$2n-3$	$B(n)$
3	3	1
4	5	3
5	7	15
6	9	105
7	11	945

3.4.1 Maximum Parsimony Problem

Maximum Parsimony Problem bedeutet mit so wenig Aufwand wie möglich.



Wenn man diesen Baum so genau aufbauen kann → S : #Mutationen entlang $S \cdot \mu(S)$.

$\mu(T) = \sum_{S \in T} \mu(S)$ Gesamtzahl der Mutationen in T (Tree, Baum) mit S = Anzahl der Kanten.