

Idee: Drei Matrizen:

Raphael
r_reitzi@cs...

- M — enden in (Mis)Match,
- I — enden in Insertion
- D — ——— Deletion

- $M_{0,0} = 0$
 $M_{i,0} = \infty$
 $M_{j,0} = \infty$
 $M_{i,j} = p(S_i, T_j)$
 $+ \min \begin{cases} M_{i-1,j-1} \\ I_{i-1,j-1} \\ D_{i-1,j-1} \end{cases}$

- $I_{0,0} = \infty$
 $I_{i,0} = \infty$
 $I_{0,j} = p + j \cdot \sigma$
 $I_{i,j} = \sigma + \min \begin{cases} M_{i,j-1} + p \\ I_{i,j-1} \\ D_{i,j-1} + p \end{cases}$

- $D_{0,0} = \infty$
 $D_{i,0} = p + i \cdot \sigma$
 $D_{0,j} = \infty$
 $D_{i,j} = \sigma + \min \begin{cases} M_{i-1,j} + p \\ I_{i-1,j} + p \\ D_{i-1,j} \end{cases}$

Algorithmus: wie bisher, aber berechne

im Schritt (i,j) $M_{i,j}, I_{i,j}, D_{i,j}$

→ Laufzeitfaktor 3

Ergebnis: $\min \{ M_{m,n}, I_{m,n}, D_{m,n} \}$

Erstmal beschrieben von

Needham : An improved algorithm
for matching biological
sequences. (1982)

Alignment wie gehabt durch Backtrac.
(durch drei Tabellen)

6] a) $(M1), (M2), (M3)$ klar.

ad $(M4)$:

Gegeben optimale Alignments für
 $(x,y), (y,z)$.

Wir konstruieren Alignment für (x,z)
mit Kosten $\leq \text{Sim}(x,y) + \text{Sim}(y,z)$.

(gzz; der optimale Alignment ist höchstens besser.)

Beobachtung: Kosten werden spaltenweise berechnet.

↳ Idee: Alignments ineinander shufflen und "zusammenschieben", bis die Worte passen.

1) Füge in Alignments $(x, y) / (y, z)$ ^{so wenig} Gap-Spalten _{wie möglich} ein, bis die y-Zeilen zusammenpassen.

wichtig: $\delta(-, -) = 0$

2) Kombinieren zu Spalten der Form

$$\begin{matrix} x \\ y \\ z \end{matrix} \begin{pmatrix} a \\ b \\ c \end{pmatrix} \in \left(\sum \cup \{-\} \right)^3$$

wichtig: keine Spalte $\begin{pmatrix} - \\ - \\ - \end{pmatrix}$

3) Spaltenweise Δ -Ungl. nachweisen, also dass

" $\delta(a, c) \leq \delta(a, b) + \delta(b, c)$ "

- $\begin{pmatrix} - \\ b \\ - \end{pmatrix} \rightsquigarrow b \neq - \rightsquigarrow \begin{pmatrix} - \\ - \end{pmatrix}$ (und dann weg)

$$\rightsquigarrow \underbrace{\delta(-, -)}_{\text{Beitrag zu Sim}(x, z)} = 0 \leq 2g = \underbrace{\delta(-, b)}_{\text{Beitrag zu Sim}(x, y)} + \underbrace{\delta(b, -)}_{\text{Beitrag zu Sim}(y, z)}$$

- $\begin{pmatrix} - \\ - \\ c \end{pmatrix} \rightsquigarrow \begin{pmatrix} - \\ c \end{pmatrix}$ ($\begin{pmatrix} a \\ - \end{pmatrix}$ symmetrisch)

$$\rightsquigarrow \delta(-, c) = g \leq 0 + g = \delta(-, -) + \delta(-, c)$$

- $\begin{pmatrix} - \\ b \\ c \end{pmatrix} \rightsquigarrow \begin{pmatrix} - \\ c \end{pmatrix}$ ($\begin{pmatrix} a \\ b \\ - \end{pmatrix}$ symmetrisch)

$$\rightsquigarrow \delta(-, c) = g \stackrel{p \geq 0}{\leq} g + p(b, c) = \delta(-, b) + \delta(b, c)$$

- (*) • $\begin{pmatrix} a \\ b \\ c \end{pmatrix} \rightsquigarrow \begin{pmatrix} a \\ c \end{pmatrix}$

$$\rightsquigarrow \delta(a, c) = p(a, c) \stackrel{\Delta\text{-Ungl. für } p}{\leq} p(a, b) + p(b, c)$$

$$= d(a,b) + d(b,c)$$

$$\bullet \begin{pmatrix} a \\ - \\ c \end{pmatrix} \rightsquigarrow \begin{pmatrix} a \\ - \\ - \\ c \end{pmatrix}$$

$$\rightsquigarrow d(a,-) + d(-,c) = 2g$$

Kosten vor- wie nachher. \parallel

b) Vorschlag: Nehme "Fastmetrik"
(\rightsquigarrow Beweis wie in a) + ϵ) \square

$$\Sigma = \{a, b, c\}$$

$$0 < g < 1$$

$$P = \begin{pmatrix} 0 & 1 & 3 \\ 1 & 0 & 1 \\ 3 & 1 & 0 \end{pmatrix}$$

Klar: (M1-3)
erfüllt.

Aber:

$$p(a,b) + p(b,c) = 2$$

$$p(a,c) = 3$$

$\rightsquigarrow \Delta$ -Ungl. verletzt.

Beweis, dass Metrik, wie in a)
außer für (*) ; sonst bräuchten
wir Δ -Ungl nicht.

ad (*)

$$\begin{pmatrix} A \\ B \\ C \end{pmatrix} \rightsquigarrow \begin{cases} \begin{pmatrix} A & - \\ - & C \end{pmatrix} & \text{, sonst } \text{i)} \\ \begin{pmatrix} A \\ C \end{pmatrix} & \text{, } B \in \{A, C\} \text{ ii)} \end{cases}$$

$$\rightsquigarrow \text{i)} \delta(A, -) + \delta(-, C)$$

$$= 2g < 2 \stackrel{\text{i)}}{\leq} \delta(A, B) + \delta(B, C)$$

$$\text{ii)} \delta(A, C) \stackrel{p \geq 0}{\leq} \delta(A, B) + \delta(B, C)$$

$\in \{\delta(A, B), \delta(B, C)\}$

□

c) (M1-3) — s.o.

(M4) ähnlich, können aber nicht alles zerlegen.

Neues Symbol: \vdash — Start eines Gap

\rightsquigarrow Schreiben $\begin{pmatrix} a \\ + \end{pmatrix} \rightarrow \text{cost } p+g$ für die ersten (linken) Gaps, $\begin{pmatrix} a \\ - \end{pmatrix}$ sonst. (insertions analog.)

$\hookrightarrow \text{cost } g$

- Behandle viele Fälle wie oben; genauer. $\rightsquigarrow \begin{pmatrix} a \\ c \end{pmatrix}$ (Δ -Ungl. gilt)



Werden \vdash zu $-$ und umgekehrt?

für:

$$\begin{pmatrix} a \\ b \\ c \end{pmatrix}, \begin{pmatrix} a \\ b \\ - \end{pmatrix}, \begin{pmatrix} a \\ - \\ - \end{pmatrix}, \begin{pmatrix} - \\ b \\ c \end{pmatrix}, \begin{pmatrix} - \\ b \\ - \end{pmatrix},$$

$$\begin{pmatrix} - \\ - \\ c \end{pmatrix}, \begin{pmatrix} - \\ - \\ - \end{pmatrix}, \begin{pmatrix} a \\ - \\ - \end{pmatrix}, \begin{pmatrix} a \\ - \\ - \end{pmatrix},$$

$$\begin{pmatrix} - \\ b \\ - \end{pmatrix}, \begin{pmatrix} - \\ b \\ - \end{pmatrix}, \begin{pmatrix} - \\ b \\ - \end{pmatrix}, \begin{pmatrix} - \\ b \\ - \end{pmatrix},$$

- Übrig bleiben Blöcke der Form

z.B.

$$\begin{pmatrix} a & a & - & a \\ - & - & - & - \\ + & c & c & + \end{pmatrix}$$

$$\begin{pmatrix} u \\ - \\ w \end{pmatrix}, \quad u, w$$

$$\rightsquigarrow \begin{pmatrix} u & - \\ - & w \end{pmatrix}$$

, dann passt es.

Siehe: □
Waterman, Smith, Beyer (1976)
Some biological sequence metrics.

d) (M2) geht kaputt:

- für semi-global (x, xy) bzw. (x, yx)

• für local (x, y, z, y)

haben Score 0, aber die Worte sind
i.A. nicht gleich.

(Ann: $f(a, a) = 0$ f.a. a ; sonst ist aber
 $\text{Sim}(x, x) \neq 0 \leadsto f(\mu_2)$)