

Teilstring

Teilsequenz

Overlap

String-Matching Problem: Text/Muster geg \rightarrow jede Mustervorkommen

naiver Alg. $O(m \cdot (n-m))$: $p = a^m, t = a^n$

String Matching Automat: $O(n)$ [Berechnung $O(|\Sigma| \cdot m^2)$ naive / $O(|\Sigma| \cdot m)$ opt.]

Boyer-Moore Alg: Bad-Char-Regel: Funktion für alle $a \in \Sigma$ letzte Position von links matchen, verschieben lt. Pos.

opt.: $O(n/m)$
worst: $O(n \cdot m)$

Good-Suf-Regel:

einfacher Suffix-Baum $t \rightarrow t.\$$: jedes Suffix von $t.\$$ ab Wurzel eintragen

String-Seed: Von Wurzel Pfad mit Muster suchen, ab da alle Blätter \neq Vorkommen

Größe: $O(|t|^2 \cdot |\Sigma|)$, z.B. $t = a^n b^n \$$

kompakter Suffix Baum: jedes innere Knoten mind. 2 Kinder, Knoten mit Strings oder Anfangs-/Endposition im Text
 $\Rightarrow O(n \cdot \log n)$ einfacher \Rightarrow Knoten mit einem Nachfolger eliminieren \Rightarrow Kompression

Seeds: $O(n \cdot \log n + m \cdot (|\Sigma| + k))$, $k \in$ Anzahl Vorkommen p im t , $|\Sigma| < \frac{n}{k}, k \geq 2$.

vollgenerierter Suffix Baum: für n Texte mit Textsymbolen $\$_1 \dots \$_n \notin \Sigma$

Longest-Common-Substring: vollgenerierter Suffix-Baum, mit Menge Text: an jedem Knoten \Rightarrow maximale Länge finden dessen Knoten alle Indizes enthält

$O(n \cdot (\log n + N^2 + N \cdot |\Sigma|))$

reversive Overlaps: vollgenerierter Suffix-Baum, innere Knoten mit Menge Index invertierter Texten

$O(n \cdot (\log n + |\Sigma| + N))$

Alignments: Einfügung / Deletion / Match / Mismatch: Bewertungsfunktion (Match / Mismatch)
Edit-Distanz ($a/b = 1, a/a = 0, g = 1, \min$), $g = 2 \cdot \text{Set}$ für Differenzsymbol, goal (min, max)
($a/a = 1, a/b = -1, g = \text{opt, max}$)

global: dyn. Programmierung mit Ähnlichkeitsmatrix $O(n \cdot m)$
graphisch: Edit-Graph, längster Weg finden (polynomuell)

lokal: Teilstrings betrachten, dyn. Programmierung

semi-global: Lücken am Anfang und Ende ignorieren
erste Zeile/Spalte mit 0 initialisieren / letzte Zeile/Spalte: finale max.

vollgenerierte Bewertungsfunktion: Lückentolerant berücksichtigen (Kosten für's öffnen + Kosten die proportional zur Länge sind)

Bewertungsmatrix 20×20 , PAM-Matrizen, akzeptierte Mutationen
PAM-Einheits: $S \rightarrow X$ ohne Inversion/Deletion, umkehr. Mutation je 100 AS.

1: wähle ähnliche Sequenz von selbst kopieren \rightarrow selbe Länge, Lücken best. \rightarrow PAM
2: $PAM_1 \rightarrow PAM_k$: $F = 20 \times 20$ $F(i, j) = w$ wie oft $a_i \rightarrow a_j$ in einer PAM-Einheitsmatrix

mult. Alignment: kein Lückensymbol an der selben Stelle bei allen Strings.

Position das Zeichen mit höchstem Vorkommen (ohne Differenzsymbol)

Abstand: Summe über alle vom consensus abweichenden Symbole an jeder Pos. | Bioinf 4

Mult-Consensus-Align-Problem: Optimierungsproblem: mult. Align mit kleinstm Abstand
↳ Summe über alle Positionen und Spalten mit paarweisen Edit-Fehl.
dyn. Prog mit k -dim. Array \rightarrow exa. Problem $O(k^2 \cdot 2^k \cdot n^k)$

Das-Mult-Consensus-Align-Problem (DMSPAP): Entscheidungsvariante, mit Bewertungsfkt
 $\sum (\sum_i f_i)^2 \rightarrow \mathbb{Q}$ und $d \in \mathbb{N}$. liefert true wenn mult. Align mit $SP \leq d$, bzgl. \sum

Kleinste gemeinsame Supersequenz über $\Sigma = \{0,1\}$ (D-(0,1)-SSP) $|H| \leq m \rightarrow$ NP-vollst.
↳ DMSPAP = NP-vollst. do if-time

paarweise Alignments: S Menge Strings, T Teilmenge S , $A'(A'')$ mult. Align SCT.
↳ A' komp. A'' falls $A' = A''$, wenn Spalten der Strings mit nur '-' eliminiert werden.

Alignment-Raum: Knoten sind Strings, Kanten mit opt. Alignment zw. zwei Strings
 T ist Stern \Rightarrow Star-Alignment

Bewertungsfkt. $\delta: (\Sigma \cup \{-\})^2 \rightarrow \mathbb{R}$ ist gut: $\delta(a,a) = 0$, $\delta(a,c) = \delta(a,b) + \delta(b,c) \Rightarrow \delta(a,b) \geq 0$

Fluoridieren: FASTA: suche Muster in DB mit Sequenzen

1. Hot-Spots ($k=6$ DNA / $k=2$ Protein) = exakte Matches der Länge k
2. Zusammenfassung: Matrix $M_{i,j} = 1$ wenn $p_i = t_j$ für Muster p (Search) t
Suche 10 best. Stäpfe (Diagonal) mit HS am Anfang/Ende
Bewertung: sowohl HS positiv / Länge drehen negativ
3. Teilalignments zu längeren zusammensetzen
4. Iterativlösung

BLAST: 1. Zets (Teilstrings ähnlich sind) mit Länge w (11 DNA / 3 Proteine)
2. Paare von Zets mit Abstand $\leq d$
3. Paare an den Enden erweitern bis Alignment-Bewertung nicht mehr steigt
wenn Bewertung $>$ Schwellenwert $S \rightarrow$ High-Score-Pair.

physikalische Kation: Marker und Funktion die Positionen angibt
Double-Digest: zwei Enzyme $A+B$, beide Vollkoden $A, B, AB \rightarrow$ Anordnung finden
nach: $|A|! \cdot |B|!$

Reduktion auf Set-Partition \rightarrow NP-vollständig
Partial-Digest: teilweise Verdau: $\binom{k}{2}$ Elemente (ideell)
Anordnung mit Reaktoren (jeweils größtes verbleibendes Element an
linken oder rechten Rand) finden.

Worst case: $O(2^k \cdot k \log k) \rightarrow$ Exponentiell, in der Regel aber gute Laufzeit
DDP einfacher Experiment / schwere Berechnung \Leftrightarrow PDP aufwändiges Experiment / einfache Berechnung

Hybridisierung mit DNA-chips, Anordnung für nxm Matrix: Consecutive Ones Property (COP)
 \rightarrow PA-Bäume: jedes Blatt ein Element, P-Knoten Kinder beliebig permutierbar,
Q-Knoten Kinder invertierbar, entsteht aus universellem Baum durch
schrittweise Anwendung der Restriktionen (Zeilen). $O(n+m+k)$

Felder: 0 nuzigler Block von 1 (jedes negativ), 1-wertig Block von 0 (jedes positiv),
Block von 0 (diagonal), dicker $\hat{=}$ Block von Nullen \Rightarrow Permutation zur Minimierung
 \Rightarrow Problem der minimalen dickenwand (durch Spaltenpermutation)
Spaltenabstandsgraph: Knoten $\hat{=}$ Spalten, ungerichtete Kanten mit Gewichtung
Planung-Diagramm (jeweils Zeilen an denen Spalten unterschiedlich sind)
 \rightarrow finde Pfad durch Graphen mit geringsten Kosten. Von unterschiedlichen Lösungen
mit beginnenden / endenden HS/AB Blöcken vermeiden: Einziges Spalte $\hat{=}$ mit
nur 0ern und 1ern suchen \rightarrow MinLM $\hat{=}$ TSP (NP-vollständig)

da Dreiecksungleichung: $\Delta - TSF \rightarrow 1,5$ -Approximation, bzw. 2 mit Span Biolinf 5
 \rightarrow minimaler Spannbaum \Rightarrow Tiefensuch

Bestimmung Basenreihenfolge: Gabelknotenpunkte (~ 1000 GP)
 Shotgun-Sequenzierung \rightarrow Fragment-Assembly Problem (~ 100 KGP)
 1500 Fragmente ≈ 500 bp, Überdeckung 7,5
 \rightarrow overlap-Bestimmung (paarweise), nicht exakt \rightarrow Layout (Anordnung) \rightarrow
 Consensus (Majority rule)

Falloquellen: Sequenzierungsfehler, Oligonucleotid, unvollst. Überdeckung, Orientierung, Repeats

Shortest-Common-Superstring Problem (SCS): Optimizing problem \Rightarrow fordern Teilstringfraktion
 statt Längemaß (Minimierung) \rightarrow Kompressionsmaß (Maximierung): Maximum-Compression (MCC)
 Zeichen Overlap und Distanzgraphen für alle Strings, suchst kürzesten / längsten
 Pfad \Rightarrow NP-vollst. (Reduktion auf HP-Problem)
 Greedy-Superstring (bsp. Overlap): Greedy, overlap suchen, Knoten löschen, etc.
 \rightarrow damit (N. Ollsh) 4-Approx., mind. 2-fach.

Bspgl. MCC sogar 2-Approx. Beweis 3-Approx.: Greedy-Neige verhindert max. 3-Approx da od.
 kein Cycle Cover: gerichtete Kreise, sodass jeder Knoten in genau einem Zyklus \rightarrow Polynomzeit
 bestimme ein Cycle Cover \rightarrow wähle in jedem Zyklus Repräsentanten \rightarrow erweitere Cycle-Cover \rightarrow
 Kante mit geringstem Overlap löschen \rightarrow Kontraktion \rightarrow 3-Approx für SCS

Reconstruction-Modell: Teilstring-Edit-Distanz: direkte Fehlerberücksichtigung (jedoch
 nicht unvollst. Überdeckung, Oligonucleotid, Repeats, iinkl. Orientierung \rightarrow NP-schwer

SBH: Sequenzierung durch Hybridisierung: Spektren, DNA bestimmen für Länge $l \approx$ alle Teilstrings der
 Länge l , die in DNA vorkommen. Spektrum-Graph: Knoten für alle Teilstrings der Länge $l-1$,
 Kanten (gerichtet) mit Beschriftung des Strichs um hinzukommen \rightarrow finit. Endspfad, dann
 einfach kompatibler String (enthält jeden Teilstring genau einmal) \rightarrow linearzeit
 Problem: nur etwa 2000bp rekonstruierbar; was nicht durch Kettenabbruchmethode erklärbar ist.

Signale in DNA-Sequenzen: finde "interessante" Regionen (Restriktions-/Bindungsstellen / Gene)
 Finde das längste gemeinsame ähnl. Teilstring (Bewertung durch Planung-Distanz)
 \rightarrow Consensus-String Problem: NP-schwer
 ≈ 3 Approx. Algorithmen (mit Eingabe l (Länge): $1 + O(\frac{l^2}{n})$ Approx, $O(n \cdot l + n^2 \cdot l)$)

Finde häufige/seltene Strings. Autokorrelationspolynom: Potenz um die Länge von Suffix \neq Präfix
 Aufstellung $S = s_1 \dots s_m$, $X = \sum_{i=1}^m x_i$, $f = \frac{1}{k^2}$, $E[X^2] = m \cdot f$, $Var[X^2] = f \cdot m \cdot (2 \cdot \text{corr}(f, f) - (2 \cdot l - 1) \cdot f + 1)$

Hidden-Markov-Modell (HMM): W'keiten für bestimmte Ereignisse und Transitionswahrscheinlichkeit
 in anderen Modell; Berechnung W'keit eines bestimmten Pfades zu einem Angebots
 Suche besten Pfad zu einer Eingabe: HMM-dekodier-Problem = polynom. mit dyn. Prog.
 \rightarrow Viterbi-Alg.: String mit Länge n für k Parameter: $O(n \cdot k)$.
 Finde CG-Inseln mit Hilfe altern. Verfahrens

Phylogenetische Räume: Metrik (Entfernungen zw. unterschiedlichen Knoten), Symmetrie Dreieck
 Ultrametrik ($\forall a, b, c \in A: (a, b), (a, c), (b, c) \rightarrow$ zwei gleich, dritte kleiner)
 UPGMA-Alg.: Ultrametric \rightarrow phylog. Baum in $O(n^2)$: immer Knoten mit kleinstem Abstand
 verbinden.

additiver Baum: metrisches Distanzmaß gegeben, kompakt: nur Knoten aus Metrik
 Konstr. durch minimalen Spannbaum: nehme jeweils billigste Kante, falls nicht: Fehler
 Parsimony: bin. Baum mit Blattbeschriftung DNA mit Länge k .
 Fitch: Folge Wurzel ein, von den Blättern nur Wurzel: Schrittlänge der Querlinien
 auf inneren Knoten, falls 0 Verzweigung sind Kostenmaß $\neq 1$ für jede Verzweigung
 von der Wurzel zu den Blättern jeweils 1 bit. Spindel wählen, falls Wahl nicht in Knoten
 \rightarrow minimales Kosten lt. Pfadlänge - Abstand $O(n \cdot k)$
 Min-Pars-Seq-Problem: gep. Strings \rightarrow finde ultramet. phylog. Baum für Parsimony minimal
 \rightarrow Exponentiell viele \rightarrow daher Repräsent. mit Ansatz, jeweils 4 Strings wählen.

Quartett: ung. phylog. Baum für 4 Taxa, optimal min. Parsimony

Wenn T ung. phylog. Baum für S (Menge von Taxa) und Q Quartett (a,b,c,d)
Q ist Konsistenz von P₁ und P₂ distinkt, mit P₁ Pfad a-b in T und P₂ Pfad c-d.
aus Menge aller konsistenten Quartette → konsistenter Baum eindeutig mit polytom. Ästern.
↳ nicht bestimmbar, daher Parsimony für alle Teilmengen (NP O(n⁴))
finde ung. phylog. Baum mit max. dist. konsistenter Quartette → Max-Quartett-Konsist → NP

Quartett-Parsimony: optimale Quartette für jede 4-elem. Teilmenge berechnen, dann zufällige Reihenfolge der Taxa wählen mit opt. Quartett der ersten 4 Elemente beginnen, für jedes weitere Element die Knoten mit min. Kosten finden (jeweils für alle drei Elementen Mengen des Baums + neues Element die Knoten a₁, a₂ um eins erhöhen, wenn Q (a₁, a₂, a₃, b_i) und dort Element einfügen.

Genom-Duplex: Reversals → Permutation die einen Abschnitt umkehren
M, SA: Permutation gegeben, ursprüngliche Identität mit min. Reversals wieder herstellen: NP
Breakpoints: Unterbrechungen von denen die Permutation nicht zusammengesetzt, max. zwei Breakpoints je Reversal möglich → opt. braucht min $\frac{1}{2} \lfloor \log_2 \text{Reversals} \rfloor$
2-Approx: in jedem Schritt durch ein Breakpoint einfügen.
Erweiterung S₀, S₁ aufsteigend, an eine ein-elementige Absteigend
ex. absteigender Stij: eliminieren einen Breakpoint / sonst zwei Breakpoints O(n³)

Sekundärstrukturen: Minimierung freier Energie: primär → Sekundärstruktur, finde homologe Basenpaare (Stems).
Alg. von Mulschiner (Dyck-Programmierung) O(n³) nur durch Basenpaare
Erweiterung: freie Energie betrachten O(n³), Rekurrenz: keine Paarung, Paarung, Stack
neben betrachteten Positionen
Verbesserung durch Alg. von Zuker: beidseitige Stems, Bulge, Hairpin, Interior loop
durch weitere Notation: O(n⁴)

Strukturvorhersage
Protein-Folding: betrachte HP-Modell (hydrophob 1/hydrophil, polar 0) und Einleitung in 2 oder 3 dimensionales Gitter, maximale 1-1 Paare ⇒ 1-1-Problem
→ NP für 2 und 3 Dimensionen
verwenden / topologische Überbrücken (1,1)-Paar (i,j) ⇒ i gerade, j ungerade oder umgekehrt: Faltungswahrscheinlichkeit

4-Approx: teile in Mengen mit 1-ern auf, aufgesplittet nach geraden/ungeraden Positionen. X die Menge mit ungeraden Elementen, Y die Menge mit geraden. opt. Lösung hat max. $\frac{1}{2} |X| |Y|$ Paare. XA-Partitionierung

Finde Faltungspunkt f₀, so daß mind $\frac{|X|}{2}$ von X auf der linken und $\frac{|Y|}{2}$ von Y auf der rechten Seite liegen, Einsetzen von x' mit y' paaren, (x') den Rest links/rechts → y'

Protein-Threading: Primärstruktur auf bekannte Tertiärstruktur zurückführen. Dabei Datenbank mit Xernen von Tertiärstrukturen (Faltblatt/Polices) und minimale/maximale Abstände von Schleifen darstellen (als 5-Tupel) und Alignment finden, so daß Lücken nur in Schleifen auftreten und keine Tarnen zerstören.