

Übungen zur Vorlesung Sequenzanalyse

Universität Bielefeld, SS 2020

Prof. Dr. Jens Stoye · Dr. Marília D. V. Braga · Leonie R. Brockmann · Rebecca K. Pfeil

<https://gi.cebitec.uni-bielefeld.de/teaching/2020summer/sa>

Übungsblatt 3 vom 7.5.2020

Abgabe am 14.5.2020 bis 12:00 Uhr (mittags)

Aufgabe 1 (Rank und Unrank)

(3 Punkte)

Gegeben sind das Alphabet $\Sigma = \{A, C, G, T\}$ mit $r_\Sigma(A) = 0$, $r_\Sigma(C) = 1$, $r_\Sigma(G) = 2$, $r_\Sigma(T) = 3$ und die Wortlänge $q = 5$. Verwende die absteigende Variante der Codierung von $q - 1$ nach 0.

1. Berechne den Rang des Wortes **AGACT**. Gib alle Zwischenschritte an.
2. Berechne den Rang des Wortes **GACTT** ohne vollständige Neuberechnung, sondern durch ein Update in konstanter Zeit (aus dem Rang des Wortes **AGACT**). Gib alle Zwischenschritte an.
3. Welches Wort hat den Rang 559?

Aufgabe 2 (Worte mit gleichem q -Gramm-Profil)

(7 Punkte)

Gegeben sei der String $x = \text{CATGCATATGCA}$; finde alle Strings, die von x unterschiedlich sind, aber das gleiche q -Gramm-Profil haben; und zwar:

1. für $q = 5$
2. für $q = 4$
3. für $q = 3$

Aufgabe 3 (Maximal-Matches-Distanz)

(3 Punkte)

Gegeben seien die Sequenzen $x = \text{montreal}$ und $y = \text{motremotreaal}$. Berechne $\delta(x||y)$ und $\delta(y||x)$. Gib jeweils die links-nach-rechts- und rechts-nach-links-Partitionen an.

Welche Beobachtung machst du? Was kannst du über die Maximal-Matches-Distanz aussagen?

Aufgabe 4 (q -Gramm- und Maximal-Matches-Distanzen als Filter)

(7 Punkte)

Gegeben seien die Sequenzen:

$x = \text{AATCGAGGTAC}$

$y_1 = \text{AAGATCGGACC}$

$y_2 = \text{AATCGCGGTAC}$

$y_3 = \text{AGGTACAATCG}$

$y_4 = \text{CTGAACGTCTG}$

Wir wollen entscheiden, ob die Sequenzen y_1, \dots, y_4 eine Edit-Distanz von max. 1 zur Sequenz x haben können, ohne alle Edit-Distanzen zu berechnen.

1. Berechne die 2-Gramm-Profile des Wortes x und y_1, \dots, y_4 . Filtere die Sequenzen y_1, \dots, y_4 mit Hilfe der 2-Gramm-Distanz. Welche Sequenzen können ausgeschlossen werden?
2. Filtere die übrigen Sequenzen mit Hilfe der Maximal-Matches-Distanz. Welche Sequenzen bleiben als Kandidaten übrig?
3. Nenne einen weiteren Filter, den man auch noch verwenden könnte.