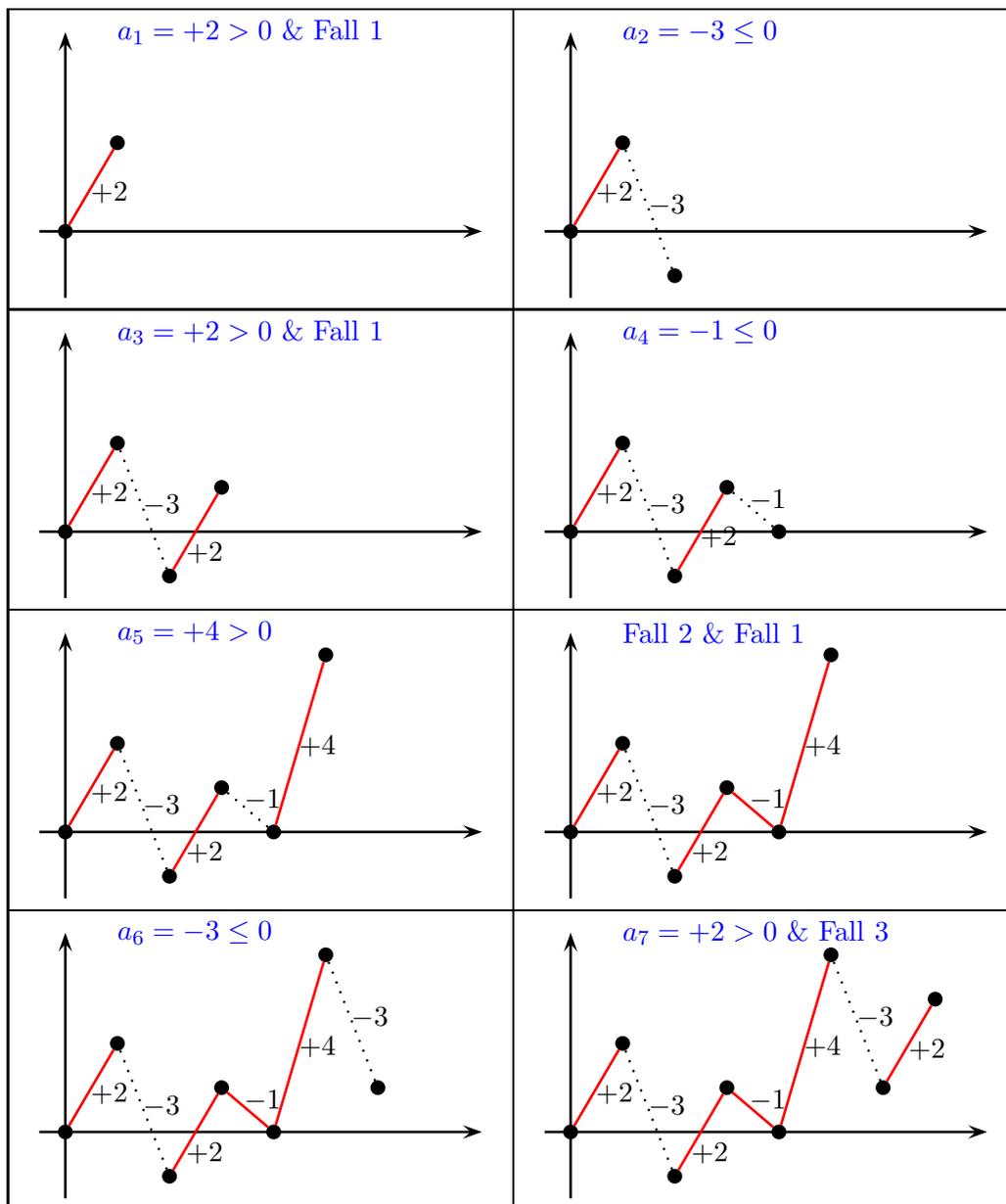


Aufgabe 1 (8 Punkte)

Ermittle mit dem in der Vorlesung angegebenen Linearzeit-Algorithmus für AMSS alle maximal bewerteten Teilfolgen von a und gib dabei alle Zwischenschritte sowie die jeweils angewendeten Fälle an. Gib auch an, welche maximal bewerteten Teilfolgen vom Algorithmus ausgegeben werden.

$$a = (+2, -3, +2, -1, +4, -3, +2)$$

Lösungsskizze (nicht ausreichend für die volle Punktzahl)



Ausgegeben wird also (1, 1), (3, 5) und (7, 7).

Aufgabe 2 (8 Punkte)

Gib für $t\$ = t_1 \cdots t_{10}\$ = aabababaab\$$ das zugehörige Suffix-Array A inklusive der zugehörigen LCP-Tabelle L an.

Zeichne dort die LCP-Intervalle mit den Suffix-Links (ohne die Suffix-Links von einelementigen LCP-Intervallen) ein.

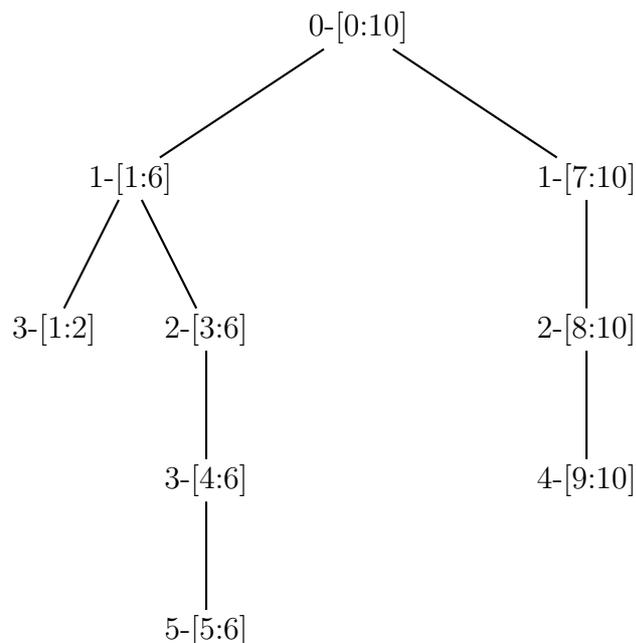
Zeichne weiterhin den zugehörigen LCP-Intervall-Baum.

Bemerkung: Die Ordnung auf dem Alphabet Σ sei $\$ < a < b$.

Lösungsskizze (nicht ausreichend für die volle Punktzahl)

i	$A[i]$	$L[i]$	$t^{A[i]}$
0	11	-1	$\$$
1	8	0	$aab\$$
2	1	3	$aabababaab\$$
3	9	1	$ab\$$
4	6	2	$abaab\$$
5	4	3	$ababaab\$$
6	2	5	$abababaab\$$
7	10	0	$b\$$
8	7	1	$baab\$$
9	5	2	$babaab\$$
10	3	4	$bababaab\$$

Im Folgenden ist der LCP-Intervall-Baum (ohne Blätter) dargestellt.



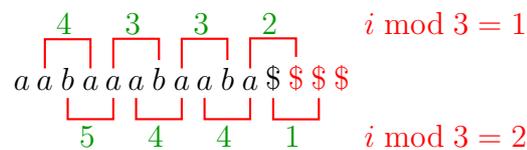
Aufgabe 3 (8 Punkte)

Erstelle für das Wort $t\$ = t_0 \cdots t_{10}\$ = aabaaabaaba\$$ ein Suffix-Array nach dem Algorithmus von Kärkkäinen und Sanders. Gib dabei alle Zwischenschritte an, wobei der rekursive Aufruf von Hand sortiert werden darf.

Hinweis: Beim Mischen von A_0 mit A_{12} soll für jede Ergebnisposition in A angegeben werden, ob 1 oder 2 Zeichenvergleiche erforderlich sind oder ob auf die Ordnung von A_{12} zurückgegriffen wird.

Lösungsskizze (nicht ausreichend für die volle Punktzahl)

Zuerst konstruieren wir alle Tripel, die an Position $i \bmod 3 \neq 0$ beginnen, und sortieren diese



Wir erhalten nun

$$\begin{aligned}
 t^{(1)} &= 4 \cdot 3 \cdot 3 \cdot 2 \\
 t^{(2)} &= 5 \cdot 4 \cdot 4 \cdot 1 \\
 \tilde{t} &= 4 \cdot 3 \cdot 3 \cdot 2 \cdot 5 \cdot 4 \cdot 4 \cdot 1 \\
 \tilde{t} \cdot 0 &= 4 \cdot 3 \cdot 3 \cdot 2 \cdot 5 \cdot 4 \cdot 4 \cdot 1 \cdot 0
 \end{aligned}$$

Wir erhalten nach dem rekursiven Sortieren der Suffixe von $\tilde{t} \cdot 0 = 2 \cdot 2 \cdot 6 \cdot 5 \cdot 3 \cdot 4 \cdot 2 \cdot 1 \cdot 0$:

$$\begin{aligned}
 A'[0] &= 8 \hat{=} 0 && \hat{=} \epsilon \\
 A'[1] &= 7 \hat{=} 10 && \hat{=} \$\$\$ \\
 A'[2] &= 3 \hat{=} 254410 && \hat{=} a\$ \$ \dots \\
 A'[3] &= 2 \hat{=} 3254410 && \hat{=} aab a\$ \$ \dots \\
 A'[4] &= 1 \hat{=} 33254410 && \hat{=} aab aab a\$ \$ \dots \\
 A'[5] &= 6 \hat{=} 410 && \hat{=} aba \$ \$ \$ \\
 A'[6] &= 0 \hat{=} 433254410 && \hat{=} aba aab aab, a\$ \$ \dots \\
 A'[7] &= 5 \hat{=} 4410 && \hat{=} aba aba \$ \$ \$ \\
 A'[8] &= 4 \hat{=} 54410 && \hat{=} baa aba aba \$ \$
 \end{aligned}$$

Damit erhalten wir A_{12} :

$$\begin{aligned}
 A_{12}[0] &= 2 + 3 \cdot (7 - 4) = 11 \hat{=} \$ \\
 A_{12}[1] &= 1 + 3 \cdot (3) = 10 \hat{=} a\$ \\
 A_{12}[2] &= 1 + 3 \cdot (2) = 7 \hat{=} aaba\$ \\
 A_{12}[3] &= 1 + 3 \cdot (1) = 4 \hat{=} aabaaba\$ \\
 A_{12}[4] &= 2 + 3 \cdot (6 - 4) = 8 \hat{=} aba\$ \\
 A_{12}[5] &= 1 + 3 \cdot (0) = 1 \hat{=} abaaabaaba\$ \\
 A_{12}[6] &= 2 + 3 \cdot (5 - 4) = 5 \hat{=} abaaba\$ \\
 A_{12}[7] &= 2 + 3 \cdot (4 - 4) = 2 \hat{=} baaabaaba\$
 \end{aligned}$$

Nun bestimmen wir A_0 :

A_0	0	1	2	3
	0	3	6	9
11				
10			9	
7			9	6
4	3			
8				
1	3	0		
5				
2				
	3	0	9	6
	a	a	b	b
	a	a	a	a
	a	b	$\$$	a

Nun müssen wir noch A_0 und A_{12} mischen (hierbei sind in Blau bzw. Rot die Werte mit $i \bmod 3 = 1$ bzw. $i \bmod 3 = 2$ in A_{12}):

A_0	0	1	2	3
	3	0	9	6
	a	a	b	b
	a	a	a	a

A_{12}	0	1	2	3	4	5	6	7
	11	10	7	4	8	1	5	2
	$\$$	a	a	a	a	a	a	b
		$\$$	a	a	b	b	b	a

A	0	1	2	3	4	5	6	7	8	9	10	11
	11 ¹	10*	3*	7*	0*	4 ¹	8 ¹	1 ¹	5 ¹	9*	2*	6

Bei den Vergleichen gibt \cdot^1 bzw. \cdot^2 an, dass ein Vergleich des ersten bzw. zweiten Zeichens der entsprechenden Suffixe zur Bestimmung der Reihenfolge beim Mischen ausreichend war, \cdot^* zeigt an, dass die relative Ordnung der Suffixe in A_{12} erforderlich war. Bei den Werten ohne Zusatzinformation war kein Vergleich mehr erforderlich.

Aufgabe 4 (8 Punkte)

Betrachte das Wort $t = t_1 \cdots t_{12} \in \{A, I, N, S\}^*$ und die zugehörige Burrows-Wheeler-Transformierte $\hat{t} = \hat{t}_0 \cdots \hat{t}_{12} = \text{SNSN\$NANAAAAI} \in \{\$, A, I, N, S\}$.

- Gib die Tabellen $C(\cdot)$ und $Occ(\cdot, \cdot)$ für \hat{t} an.
- Rekonstruiere t aus \hat{t} .
- Konstruiere den Wavelet-Tree zu \hat{t} .
- Bestimme die Werte von $Occ(A, 11)$ und $Occ(N, 7)$ nach der in der Vorlesung vorgestellten Methode aus dem Wavelet-Tree aus c).

Hinweis: Das zu betrachtende Alphabet ist $\Sigma \cup \{\$\} = \{\$, A, I, N, S\}$, wobei die Ordnung auf dem Alphabet durch die hier angegebene Reihenfolge gegeben ist.

Lösungsskizze (nicht ausreichend für die volle Punktzahl)

Wir erhalten die folgenden Tabellen:

$Occ(\cdot, \cdot)$	$\$$	A	I	N	S
0	0	0	0	0	1
1	0	0	0	1	1
2	0	0	0	1	2
3	0	0	0	2	2
4	1	0	0	2	2
5	1	0	0	3	2
6	1	1	0	3	2
7	1	1	0	4	2
8	1	2	0	4	2
9	1	3	0	4	2
10	1	4	0	4	2
11	1	5	0	4	2
12	1	5	1	4	2

$C(\cdot)$	i
$\$$	0
A	1
I	6
N	7
S	11

Der letzte Buchstabe von t ist immer \hat{t}_0 , also $t_{12} = \hat{t}_0 = S$. Der Buchstabe unmittelbar vor dem zuletzt berechneten Buchstaben \hat{t}_i in t ist $\hat{t}_{i'}$ und berechnet sich mittels

$$i' = \text{LF}(i) = C(\hat{t}_i) + Occ(\hat{t}_i, i - 1).$$

Somit gilt:

i	\hat{t}_i	$C(\hat{t}_i) + Occ(\hat{t}_i, i - 1) =$	i'
0	S	$C(S) + Occ(S, -1) = 11 + 0 = 11$	
11	A	$C(A) + Occ(A, 10) = 1 + 4 = 5$	
5	N	$C(N) + Occ(N, 4) = 7 + 2 = 9$	
9	A	$C(A) + Occ(A, 8) = 1 + 2 = 3$	
3	N	$C(N) + Occ(N, 2) = 7 + 1 = 8$	
8	A	$C(A) + Occ(A, 7) = 1 + 1 = 2$	
2	S	$C(S) + Occ(S, 1) = 11 + 1 = 12$	
12	I	$C(I) + Occ(I, 11) = 6 + 0 = 6$	
6	A	$C(A) + Occ(A, 5) = 1 + 0 = 1$	
1	N	$C(N) + Occ(N, 0) = 7 + 0 = 7$	
7	N	$C(N) + Occ(N, 6) = 7 + 3 = 10$	
10	A	$C(A) + Occ(A, 9) = 1 + 3 = 4$	
4	$\$$	Fertig! =	

Aufgabe 5 (8 Punkte)

Sei Σ ein Alphabet, $t = t_1 \cdots t_n \in \Sigma^n$ und $1 < k \in \mathbb{N}$. Konstruiere einen Algorithmus mit optimaler Laufzeit, der alle kürzesten Wörter $w \in \Sigma^*$ findet, die in t mindestens k -mal auftreten, aber das zugehörige gespiegelte Wort w^R nicht in t auftritt.

Zeige die Korrektheit des Algorithmus und analysiere dessen Laufzeit.

Hinweis: Für $w = w_1 \cdots w_\ell \in \Sigma^\ell$ ist das zugehörige gespiegelte Wort $w^R = w_\ell \cdots w_1$.

Lösungsskizze (nicht ausreichend für die volle Punktzahl)

Konstruiere einen Suffix-Baum T für $t' = t'_1 \cdots t'_{2n+2} = t\#t^R\$$. Mit Hilfe einer Tiefensuche in T zähle getrennt für jeden inneren Knoten die Anzahl der nachfolgenden Blätter mit einem Suffix in t (also Blätter mit Index aus $[1 : n]$) und in t^R (also Blätter mit Index aus $[n+2 : 2n+1]$) und bestimme ebenfalls die Worttiefe jedes Knotens in T .

Das zu einem inneren Knoten korrespondierende Wort w , dessen Zähler für Blätter in t größer gleich k und dessen Zähler für t^R gleich 0 ist, kommen eben in t mindestens k -mal vor und treten in t^R nicht auf. Beachte, dass ein Wort w in t^R genau dann nicht auftritt, wenn w^R nicht in t auftritt. Dies gilt auch für alle Wörter, die auf der eingehenden Kante zu diesem Knoten enden. Markiere also in T (mit einer weiteren Tiefensuche oder parallel zur ersten) alle inneren Knoten mit den oben genannten Eigenschaften.

Für die Ausgabe kürzester solcher Wörter kommen nun zunächst nur die Wörter in Frage, die zu Kanten gehören, deren Endknoten markiert sind und deren Anfangsknoten eben nicht markiert ist. Das korrespondierende Teilwort zum Anfangsknoten ist nämlich kürzer als die korrespondierenden Wörter zu Positionen auf der Kante zum oder am Endknoten. Da wir nur kürzeste Teilwörter ausgeben wollen, sind nur Teilwörter relevant, die zur ersten Position auf einer Kante korrespondieren, für die der zugehörigen Anfangsknoten nicht markiert und der Endknoten markiert ist. Des Weiteren muss der Anfangsknoten unter all solchen eine niedrigste Worttiefe besitzen. Mit einer weiteren Tiefensuche können wir nun die minimale Worttiefe von inneren unmarkierten Knoten ermitteln, die mindestens ein markiertes Kind besitzen. Solche inneren Anfangsknoten mit minimaler Worttiefe bezeichnen wir im Folgenden als *relevant*.

Für die Ausgabe ermitteln wir mit einer weiteren Tiefensuche für die Kinder von relevanten Knoten die zugehörigen Blattlisten beschränkt auf Teilwörter aus t , also mit Blattmarkierungen aus $[1 : n]$. Tatsächlich müssen wir keine vollständigen Blattlisten speichern, sondern nur einen Index aus dieser Blattliste, z.B. den kleinsten. Für die als relevant markierten Knoten mit Worttiefe d geben wir nun zu jedem markierten Kind die zugehörige Referenz (s, e) des Wortes aus, das nach dem ersten Zeichen auf der Kante zu dem markierten Kind endet, wobei s ein (und nur ein) Knoten aus der Blattliste des Kindes ist und $e = s + d$ (das Wort entspricht also $t_s \cdots t_e$).

Die Laufzeit für Ukkonens Algorithmus zur Konstruktions von T ist linear in der Länge der Zeichenreihe $t\#t^R\$$, also $O(2n) = O(n)$. Auch die oben genannten Tiefensuchen in T können jeweils in Zeit $O(n)$ ausgeführt werden. Die Konstruktion der Blattlisten an den Kindern von relevanten Knoten ist nach Vorlesung ebenfalls in linearer Zeit möglich. Da relevante Knoten nicht Vorfahren eines anderen relevanten Knotens sein können, kann maximal an jedem markierten Kind eines relevanten Knoten eine Referenz ausgegeben werden, also insgesamt maximal $O(n)$. Somit ist die Laufzeit $O(n)$, was optimal ist.