

# Assignment 1, Bioinformatics

Davide Cozzi, 829827

# Capitolo 1

## esercizio 1

### 1.1 Versione 1

Si assumano stringhe su  $\Sigma = \{a, c, g, t, A, C, G, T\}$  indicizzate a partire dalla posizione 0, per comodità.

L'algoritmo si divide in due parti:

1. divisione in “token” delle due stringhe. A partire da ogni stringa si ottiene un vettore di sottostringhe di uguali caratteri
2. confronto dei due vettori di sottostringhe

Si procede specificando che:

- $len(X)$  restituisce la lunghezza di  $X$ , che sia una stringa o un vettore
- $push(X, Y)$  effettua l'operazione di **push** di  $Y$  nel vettore/stringa  $X$
- $string(Y)$  effettua il cast di  $Y$  a **string**
- $Y[i, e]$  specifica la sottostringa di  $Y$  che va dall'indice  $i$  incluso all'indice  $e$  incluso
- $lowercase(Y)$  porta  $Y$  in *lowercase*

Si ha quindi la prima parte dell'algoritmo, che effettua la divisione in “token”.

**Esempio 1.** Vediamo un esempio di input e output della funzione *Splitter*.

- **input:** "aaataaaggggccccctttttttttttcc"
- **output:** ["aaa", "t", "aaa", "gggg", "cccc", "tttttttttt", "cc"]

Si ha quindi lo pseudocodice:

---

**Algorithm 1** Algoritmo per lo split in “token” delle stringhe

---

```

function SPLITTER(s)
  result  $\leftarrow$  [ ]
  n  $\leftarrow$  len(s)
  last_mismatch  $\leftarrow$  0
  if n == 1 then
    push(result, string(s[0]))
    return result
  end if
  tmp  $\leftarrow$  ""
  for i  $\leftarrow$  0 to n do
    if i > 0 and s[i - 1]  $\neq$  s[i] then
      push(result, tmp)
      tmp  $\leftarrow$  ""
    end if
    push(tmp, s[i])
    if i == (n - 1) then
      push(result, tmp)
    end if
  end for
  return result
end function

```

---

Una volta ottenuto il vettore dei “token” delle due sequenze in input basta confrontare i due vettori.

Prima di vedere il confronto si specifica che:

- le due sequenze vengono trasformate in *lowercase* per praticità
- si assume che le due sequenze non possono essere stringhe vuote  $\varepsilon$

Fatte queste assunzioni si può fare ancora un ultima osservazione prima di procedere con l'algoritmo. Qualora i due vettori prodotti dalla funzione

**Splitter** fossero di lunghezza diversa allora l'algoritmo non potrà in nessun caso dire che ci sia stata un'infezione in quanto si assume che non ci siano né *gaps* né *inserimenti*. In altri termini qualora i due vettori siano di diversa cardinalità sicuramente, con queste premesse, non è avvenuta alcuna infezione. Vediamo ora le altre casistiche. Si itera contemporaneamente sull'elemento  $i$ -esimo dei due vettori e:

- se i due elementi  $i$ -esimi presentano il primo carattere diverso allora, date le premesse, si specifica che non può essere avvenuta un'infezione
- se i due elementi  $i$ -esimi presentano il primo carattere uguale allora si procede confrontando i due elementi  $i$ -esimi (si assuma  $X$  vettore relativo alla “sequenza originale” e  $Y$  vettore relativo alla sequenze che si vuole dimostrare l'eventuale infezione). Innanzitutto la lunghezza di  $Y[i]$  deve essere maggiore di 2 volte quella di  $X[i]$ , in modo da verificare che per ogni base ci sia stata l'aggiunta di almeno una base. Inoltre
  - se il primo carattere è “a” allora la lunghezza di  $Y[i] - X[i]$  deve essere minore o uguale a 5 volte quella di  $X[i]$ , in quanto per ogni “a” in  $X[i]$  posso avere al più 5 “a” in  $Y[i]$  in aggiunta ad ogni “a” di  $X[i]$
  - se il primo carattere è “t” allora la lunghezza di  $Y[i] - X[i]$  deve essere minore o uguale a 10 volte quella di  $X[i]$ , in quanto per ogni “t” in  $X[i]$  posso avere al più 10 “t” in  $Y[i]$  in aggiunta ad ogni “T” di  $X[i]$
  - se il primo carattere è “c” o “g” allora non ho bisogno di alcun controllo ulteriore, non avendo un “limite superiore”, al numero di basi uguali che si possono aggiungere. Mi basta quindi controllare che per ogni base ne sia stata aggiunta almeno una

Si ha quindi il seguente pseudocodice dell'algoritmo, che ritorna  $\top$  o  $\perp$  a seconda che sia possibile che  $seq2$  sia una versione infettata di  $seq1$ :

---

**Algorithm 2** algoritmo di verifica dell'infezione

---

```

function CHECKINFECTION( $seq1, seq2$ )
  if  $len(seq1) == 0$  or  $len(seq2) == 0$  then
    return  $\perp$ 
  end if
   $vseq1 \leftarrow Splitter(lowercase(seq1))$ 
   $vseq2 \leftarrow Splitter(lowercase(seq2))$ 
   $check \leftarrow \top$ 
   $check1 \leftarrow \top$ 
   $check2 \leftarrow \top$ 
  if  $len(vseq1) \neq len(vseq2)$  then
    return  $\perp$ 
  end if
   $j \leftarrow len(vseq1)$ 
  for  $i \leftarrow 0$  to  $j$  do
    if  $vseq1[i][0] \neq vseq2[i][0]$  or  $\neg check$  then
      return  $\perp$ 
    end if
    if  $vseq1[i][0] == 'a'$  then
       $check1 \leftarrow len(vseq2[i]) \geq 2 \cdot len(vseq1[i])$ 
       $check2 \leftarrow (len(vseq2[i]) - len(vseq1[i]) \leq 5 \cdot len(vseq1[i]))$ 
       $check \leftarrow check1$  and  $check2$ 
    else if  $vseq1[i][0] == 't'$  then
       $check1 \leftarrow len(vseq2[i]) \geq 2 \cdot len(vseq1[i])$ 
       $check2 \leftarrow (len(vseq2[i]) - len(vseq1[i]) \leq 10 \cdot len(vseq1[i]))$ 
       $check \leftarrow check1$  and  $check2$ 
    else
       $check \leftarrow len(vseq2[i]) \geq 2 \cdot len(vseq1[i])$ 
    end if
  end for
  return  $check$ 
end function

```

---

**Esempio 2.** Vediamo un esempio di infezione:

- **input:** "ATAGCTC" e "AAATTAAAGGGCCCCCTTTTTTCC"
- **output:**  $\top$

ATAGCTC  
 AAATTAAAGGGGCCCTTTTTTCC

- **input:** "ATAGCTC"  $e$   
"AAATAAAGGGCCCCCTTTTTTCC"
- **output:**  $\perp$

ATAGCTC  
 AAATAAAGGGGCCCCCTTTTTTCC

- **input:** "ATAGCTC"  $\epsilon$   
"AAATTAAAAAAGGGGCCCCCTTTTTTTCC"
- **output:**  $\perp$

ATAGCTC  
 AAATTAAAAAAAGGGGCCCTTTTTTCC

- **input:** "ATAGCTC"  $e$   
"AAACCTTAAAAAGGGCCCCCTTTTTT"
- **output:**  $\perp$

ATAGCTC  
AAACCTTAA AAAAGGGGCCCCCTTTTTT

5

## 1.2 Versione 2

Si assumano stringhe su  $\Sigma = \{a, c, g, t, A, C, G, T\}$  indicizzate a partire dalla posizione 0, per comodità.

Indico anche una versione intuitivamente più semplice. In questo caso si scorre la sequenze che si suppone essere originale, tenendo conto dei caratteri uguali ripetuti. Appena si ha un cambio di carattere si verifica se nella sequenza che si vuole dimostrare essere una mutazione di quella originale si ha una corretta sequenza di caratteri secondo la specifica.

Si ha che:

- per comodità le sequenze sono riportate in lowercase e si assume che, in presenza di anche solo una sequenza nulla, l'algoritmo restituisce  $\perp$
- **seq1** e **seq2** sono rispettivamente la sequenza originale al sequenza che si vuole dimostrare essere la mutazione
- **co** e **cm** sono le variabili che ogni volta accumulano il conteggio dei caratteri uguali consecutivi rispettivamente sulla sequenza originale e su quella che si suppone mutata
- *i* e *j* sono rispettivamente gli indici per la prima e per la seconda sequenza

In poche parole si itera sulla sequenza originale, aggiornando di volta in volta il contatore relativo finché si ha lo stesso carattere. Nel momento in cui si ha un cambiamento o si è arrivato all'ultimo carattere della sequenza si ferma il conteggio e si verifica il conteggio sulla seconda sequenza (anche in questo caso facendo attenzione a non andare *out of bounds*). Qualora la seconda sequenza presenti, all'indice a cui si è arrivati con l'algoritmo, un carattere diverso da quello della prima si può restituire  $\perp$ . Controllando la seconda sequenza si aggiorna il rispettivo contatore e l'indice. Una volta che anche sulla seconda sequenza si ha un cambio di carattere o si è arrivati alla fine si confrontano i due contatori secondo le specifiche, confrontandoli in modo analogo a quanto fatto con le lunghezze dei "token" nella *versione 1*. Finito questo controllo si azzerano i contatori e si verifica che, qualora la sequenza originale sia stata visitata interamente, anche la seconda sia conclusa. In caso contrario si restituisce  $\perp$ .

Si ha quindi il seguente pseudocodice.

**Algorithm 3** algoritmo di verifica dell'infezione, seconda versione

---

```

function CHECKINFECTION(seq1, seq2)
  m  $\leftarrow$  len(seq1)
  n  $\leftarrow$  len(seq2)
  if m == 0 or n == 0 then
    return  $\perp$ 
  end if
  seq1  $\leftarrow$  lowercase(seq1)
  seq2  $\leftarrow$  lowercase(seq2)
  co  $\leftarrow$  0
  cm  $\leftarrow$  0
  j  $\leftarrow$  0
  for i  $\leftarrow$  0 to m do
    co  $\leftarrow$  co + 1
    if i == n - 1 or seq1[i]  $\neq$  seq[i + 1] then
      if seq1[i]  $\neq$  seq2[j] then
        return  $\perp$ 
      end if
      while j  $\neq$  n - 1 and seq2[j] == seq2[j + 1] do
        cm  $\leftarrow$  cm + 1
        j  $\leftarrow$  j + 1
      end while
      cm  $\leftarrow$  cm + 1
      j  $\leftarrow$  j + 1
      if seq1[i] == 'a' and cm  $\geq$  2 · co and cm - co  $\leq$  5 · co then
        check  $\leftarrow$   $\top$ 
      else if seq1[i] == 't' and cm  $\geq$  2 · co and cm - co  $\leq$  10 · co then
        check  $\leftarrow$   $\top$ 
      else if (seq1[i] == 'c' or seq1[i] == 'g') and cm  $\geq$  2 · co then
        check  $\leftarrow$   $\top$ 
      else
        return  $\perp$ 
      end if
      co  $\leftarrow$  0
      cm  $\leftarrow$  0
      if  $\neg$ check or (i == m - 1 and j  $\neq$  n) then
        return  $\perp$ 
      end if
    end if
  end for
  return check
end function

```

---



Dal punto di vista pratico questo secondo algoritmo evita la fase di pre-processing vista nel primo algoritmo con la funzione **Splitter**. Anziché avere a priori le sottostringhe di cui confrontare le lunghezze, la seconda versione tiene conto di volta in volta dei caratteri uguali consecutivi, confrontando i contatori (motivo per cui anche gli esempi possono essere riadattati anche a questa seconda versione, sostituendo le lunghezze dei token coi due contatori).

## Capitolo 2

### Esercizio 2

# Capitolo 3

## Esercizio 3

La **distanza di Hamming** tra due stringhe, che si assumono di uguale lunghezza, è il numero di indici per i quali i due caratteri associati sulle due stringhe sono diversi. È quindi un conteggio delle sostituzioni necessarie per passare da una stringa all'altra.

Ipotizzando di voler studiare la distanza di Hamming tra due sequenze tramite la griglia estesa si può ipotizzare di ragionare solo in termini di una delle due sequenze, ovvero considerare i pesi degli archi in ottica di una sola delle due sequenze. Si associano quindi i seguenti pesi:

- $w(diagonale) = 0$
- $w(verticale) = 0$
- $w(orizzontale) = 1$

o, in modo speculare se si vuole ragionare sull'altra sequenza:

- $w(diagonale) = 0$
- $w(verticale) = 1$
- $w(orizzontale) = 0$

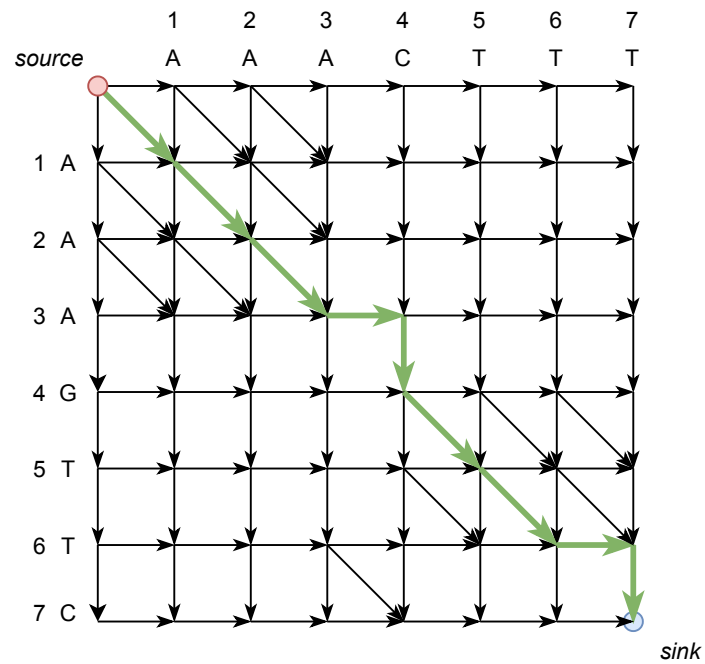
Dove gli archi diagonali rappresentano un match.

A questo punto si cerca il cammino di peso minimo che parte dal nodo *source*, posto in alto a sinistra, e arriva al nodo *sink*, posto in basso a destra.

Prendendo, ad esempio, in input:

- AAAGTTC
- AAAC TTT

Si ha un cammino minimo (non l'unico), assumendo costo degli archi in orizzontale pari a 1 e in verticale pari a 0, del tipo:



Dove in verde è segnato il cammino minimo scelto. Avendo solo due archi orizzontali, che ricordiamo avere peso 1 mentre gli altri hanno peso nullo, possiamo concludere che la distanza di Hamming tra le due stringhe è pari a 2.

# Capitolo 4

## Esercizio 4

Il problema della **longest common substring** si pone l'obiettivo di estrarre, a partire da due stringhe di lunghezza arbitraria, anche non uguale, la sottostringa, comune ad entrambe, più lunga.

Dal punto di vista della griglia si assegnano i seguenti pesi:

- $w(\textit{diagonale}) = 1$
- $w(\textit{verticale}) = -i$
- $w(\textit{orizzontale}) = -j$

Dove gli archi diagonali rappresentano un match ed  $i$  e  $j$  sono gli indici che scorrono le due stringhe, ovvero rispettivamente sulle righe e sulle colonne.

**Qualora si ottenga un valore negativo viene messo 0.**

Si ragiona quindi tenendo traccia del valore massimo raggiunto, tenendo traccia, per una maggior efficienza, anche delle coordinate. Una volta completata la griglia si avrà che il valore massimo corrisponde al nodo *sink* del cammino composto dalla più lunga sequenza possibile di archi diagonali consecutivi. Qualsiasi altra “operazione”, con archi orizzontali o diagonali, comporta infatti una perdita di punteggio e ogni volta che un cammino diagonale termina viene azzerato il punteggio, tramite pesi negativi pesati sugli indici. Azzerando ogni volta si permette di avere la costruzione di un eventuale cammino diagonale che assegna man mano il valore corrispondente alla lunghezza del cammino stesso, avendo che la diagonale pesa 1. Il valore massimo quindi altro non è che la lunghezza della *longest common substring*. Sapendo che archi diagonali corrispondono a match tra le due sequenze si ottiene che tale cammino corrisponde alla *longest common substring*.

Volendo si può scegliere di salvare in un vettore i valori massimi qualora coincidano, per ottenere eventuali più *longest common substring* qualora ce ne

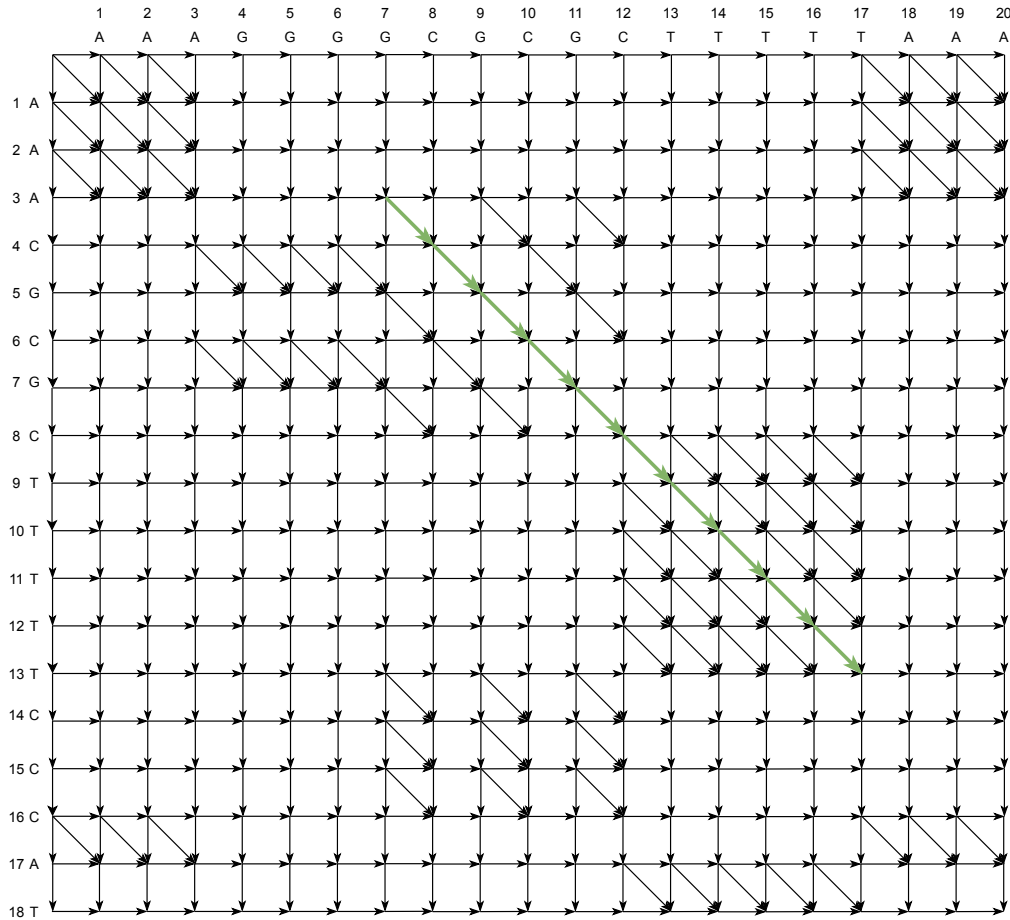
siano.

Per ricostruire la *longest common substring* si parte dal nodo *sink*, come detto definito dal valore massimo calcolato. Si aggiunge il carattere corrispondente e si risale tutto il cammino diagonale, aggiungendo di volta in volta in testa il carattere letto. Ci si ferma quindi quando non si ha più un nodo il cui punteggio è stato ottenuto tramite un arco diagonale.

Vediamo quindi un esempio pratico. Siano date:

- AAACGCGCTTTTTCCCAT
- AAAGGGGCGCGCTTTTAAA

Si costruisce quindi la griglia e si identifica il cammino diagonale più lungo:



Ricostruendo si ha che tale cammino identifica un massimo pari a 10 e si ricostruisce:

CGCGCTTTT

che è appunto la *longest common substring* delle due sequenze in input.