Memoire de stage de Maîtrise IGIS ITA

K-mers et k-mers espacés

Gabriel Toublanc

26 janvier 2017

Université de Rouen, U.F.R des Sciences et Techniques de Saint-Etienne-du-Rouvray, Equipe LITIS TIBS

Introduction

Introduction

Comptage des k-mers

Les k-mers sont des facteurs de séquences ADN.

Les k-mers sont des facteurs de séquences ADN.

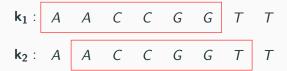
Ex: Avec une séquence x = AACCGGTT, on a les k-mers de taille 6 (6-mers) suivants:

Les k-mers sont des facteurs de séquences ADN.

Ex: Avec une séquence x = AACCGGTT, on a les k-mers de taille 6 (6-mers) suivants:

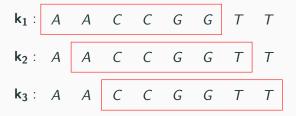
Les k-mers sont des facteurs de séquences ADN.

Ex: Avec une séquence x = AACCGGTT, on a les k-mers de taille 6 (6-mers) suivants:



Les k-mers sont des facteurs de séquences ADN.

Ex: Avec une séquence x = AACCGGTT, on a les k-mers de taille 6 (6-mers) suivants :



2

Utilisé afin de simuler les erreurs d'insertions sur les lectures longues.

Utilisé afin de simuler les erreurs d'insertions sur les lectures longues.

Ex : Avec une séquence x = AACCGGTT et le motif m = 111011, on a les 5-mers espacés suivants :

Utilisé afin de simuler les erreurs d'insertions sur les lectures longues.

Ex : Avec une séquence x = AACCGGTT et le motif m = 111011, on a les 5-mers espacés suivants :

$$\mathbf{k_1}: A \quad A \quad C \quad X \quad G \quad G \quad T \quad T$$

Utilisé afin de simuler les erreurs d'insertions sur les lectures longues.

Ex : Avec une séquence x = AACCGGTT et le motif m = 111011, on a les 5-mers espacés suivants :



Utilisé afin de simuler les erreurs d'insertions sur les lectures longues.

Ex : Avec une séquence x = AACCGGTT et le motif m = 111011, on a les 5-mers espacés suivants :



3

Utilisé afin de simuler les erreurs de délétion sur les lectures longues.

Utilisé afin de simuler les erreurs de délétion sur les lectures longues.

 $\left(\textbf{L}-\textbf{k}+1\right)*\textbf{4}^{t}$ k-mers espacés à insertion possibles.

Utilisé afin de simuler les erreurs de délétion sur les lectures longues.

 $(L-k+1)*4^t$ k-mers espacés à insertion possibles.

Ex : Avec une séquence <math>x = AACCGGTT et le motif

 $\mathbf{m} = \mathbf{111011}$, on a les 6-mers espacés suivants :

Utilisé afin de simuler les erreurs de délétion sur les lectures longues.

 $(\mathbf{L} - \mathbf{k} + \mathbf{1}) * \mathbf{4}^{t}$ k-mers espacés à insertion possibles.

Ex : Avec une séquence <math>x = AACCGGTT et le motif

 $\mathbf{m}=\mathbf{111011}$, on a les 6-mers espacés suivants :

 $\mathbf{k_1}: A \quad A \quad C \xrightarrow{A, C, G, T} C \quad G \quad G \quad T \quad T$

Utilisé afin de simuler les erreurs de délétion sur les lectures longues.

 $(L-k+1)*4^t$ k-mers espacés à insertion possibles.

 $\mathbf{E}\mathbf{x}$: Avec une séquence $\mathbf{x} = \mathbf{AACCGGTT}$ et le motif

 $\mathbf{m}=\mathbf{111011}$, on a les 6-mers espacés suivants :

$$\mathbf{k_1}: A \quad A \quad C \xrightarrow{A, C, G, T} C \quad G \quad G \quad T \quad T$$
 $\mathbf{k_2}: A \quad A \quad C \quad C \xrightarrow{A, C, G, T} G \quad G \quad T \quad T$

Utilisé afin de simuler les erreurs de délétion sur les lectures longues.

 $(L-k+1)*4^t$ k-mers espacés à insertion possibles.

Ex : Avec une séquence $\mathbf{x} = \mathbf{AACCGGTT}$ et le motif $\mathbf{m} = \mathbf{111011}$, on a les 6-mers espacés suivants :

Utilisé afin de simuler les erreurs de délétion sur les lectures longues.

 $(L-k+1)*4^t$ k-mers espacés à insertion possibles.

Ex: Avec une séquence x = AACCGGTT et le motif m = 111011, on a les 6-mers espacés suivants :

k ₁ :	Α	Α	c ²	1, C, G, T	С	G	G	T	T
k ₂ :	Α	Α	С	C A	, C, G, T	G	G	T	Т
k ₃ :	Α	Α	С	С	$G \stackrel{A}{\sim}$, C, G, T	G	Т	T
k ₄ :	Α	Α	С	С	G	G A	1, C, G, 7	ТТ	Т

 $\textbf{Entr\'ees}: \texttt{table_hachage} \ \textit{table}, \ \texttt{cha\^{i}nes} \ \textit{lectures}, \ \texttt{cha\^{i}ne} \ \textit{motif}, \ \texttt{entier} \ \textit{k}$

Entrées : table_hachage table, chaînes lectures, chaîne motif, entier k
pour chaque lecture de lectures faire

Entrées : table_hachage table, chaînes lectures, chaîne motif, entier k pour chaque lecture de lectures faire

```
\begin{array}{ll} \mathbf{pour} \ i = 0; \ i + k \leq |\mathit{lecture}|; \ i + = 1 \ \mathbf{faire} \\ & \mathit{kmerEntier} = 0; \\ & \mathit{kmer} = ""; \end{array}
```

 $\begin{table}{ll} \bf Entrées: table_hachage \ \it table, \ chaînes \ \it lectures, \ chaîne \ \it motif, \ entier \ \it k \ \it pour \ chaque \ \it lectures \ \it faire \ \it lectures \ \it lectu$

```
\begin{array}{l} \mathbf{pour} \ i = 0; \ i+k \leq |\mathit{lecture}|; \ i+=1 \ \mathbf{faire} \\ \\ k\mathit{merEntier} = 0; \\ k\mathit{mer} = ""; \\ \\ \mathbf{pour} \ j = 0; \ j < k; \ j+=1 \ \mathbf{faire} \\ \\ \\ | \ \mathbf{si} \ \mathit{motif[j]} \neq 0 \ \mathbf{alors} \ \mathit{kmer} = \mathit{kmer} + \mathit{lecture[i+j]}; \\ \\ \mathbf{fin} \end{array}
```

```
Entrées: table_hachage table, chaînes lectures, chaîne motif, entier k
pour chaque lecture de lectures faire
    pour i = 0; i + k < |lecture|; i + = 1 faire
         kmerEntier = 0:
         kmer = "";
         pour j = 0; j < k; j+ = 1 faire
              si motif[j] \neq 0 alors kmer = kmer + lecture[i + j];
         fin
         pour chaque nucleotide de kmer faire
              kmerEntier* = 4;
              suivant valeur de nucleotide faire
                  cas où A faire:
                  cas où C faire kmerEntier + = 1;
                  cas où G faire kmerEntier + = 2;
                  cas où T faire kmerEntier+=3;
              fin
         fin
```

```
Entrées: table_hachage table, chaînes lectures, chaîne motif, entier k
pour chaque lecture de lectures faire
    pour i = 0; i + k < |lecture|; i + = 1 faire
         kmerEntier = 0:
         kmer = "";
         pour j = 0; j < k; j+ = 1 faire
              si motif[j] \neq 0 alors kmer = kmer + lecture[i + j];
         fin
         pour chaque nucleotide de kmer faire
              kmerEntier* = 4;
              suivant valeur de nucleotide faire
                  cas où A faire:
                  cas où C faire kmerEntier + = 1;
                  cas où G faire kmerEntier + = 2;
                  cas où T faire kmerEntier+=3;
              fin
         fin
         table[kmerEntier] + = 1;
    fin
fin
```

kmerExpand

kmerExpand

 $\mbox{\bf Entr\'ees}$: chaînes $\it lectures$, chaîne $\it motif$, entier $\it k$

kmerExpand

Entrées : chaînes *lectures*, chaîne *motif*, entier *k* pour chaque *lecture de lectures* faire

kmerExpand

```
\label{eq:charge_entropy} \begin{split} & \textbf{Entr\'ees}: \text{cha\^nes } \textit{lectures}, \text{ cha\^ne } \textit{motif}, \text{ entier } \textit{k} \\ & \textbf{pour chaque } \textit{lecture } \textit{de lectures faire} \\ & | & \textbf{pour } i = 0; \ i + k \leq |\textit{lecture}|; \ i + = 1 \ \textbf{faire} \\ & | & \textit{kmer} = \textit{lecture}[i:k] \ \textit{kmerExpandRec(kmer, motif, 0)} \\ & | & \textbf{fin} \end{split}
```

kmerExpand

```
\label{eq:charge_entropy} \begin{split} &\textbf{Entr\'ees}: \text{cha\^{n}es } \textit{lectures}, \text{ cha\^{n}e } \textit{motif}, \text{ entier } \textit{k} \\ &\textbf{pour chaque } \textit{lecture } \textit{de lectures faire} \\ & | & \textbf{pour } i = 0; \ i + k \leq |\textit{lecture}|; \ i + = 1 \ \textbf{faire} \\ & | & \textit{kmer} = \textit{lecture}[i:k] \ \textit{kmerExpandRec(kmer, motif, 0)} \\ & & \textbf{fin} \end{split}
```

kmerExpandRec

Entrées : chaîne kmer, chaîne nvKmer, chaîne motif, entier posMotif, entier posKmer

kmerExpand

```
\label{eq:charge_entropy} \begin{split} &\textbf{Entr\'ees}: \text{cha\^{n}es } \textit{lectures}, \text{ cha\^{n}e } \textit{motif}, \text{ entier } \textit{k} \\ &\textbf{pour chaque } \textit{lecture } \textit{de lectures faire} \\ & | & \textbf{pour } i = 0; \ i + k \leq |\textit{lecture}|; \ i + = 1 \ \textbf{faire} \\ & | & \textit{kmer} = \textit{lecture}[i:k] \ \textit{kmerExpandRec(kmer, motif, 0)} \\ & & \textbf{fin} \end{split}
```

kmerExpandRec

Entrées : chaîne kmer, chaîne nvKmer, chaîne motif, entier posMotif, entier p

kmerExpand

kmerExpandRec

```
\label{eq:continuous_section} \begin{split} &\textbf{Entr\'ees}: \text{cha\^ine } \textit{kmer}, \text{ cha\^ine } \textit{nvKmer}, \text{ cha\^ine } \textit{motif}, \text{ entier } \textit{posMotif}, \text{ entier } \textit{posKmer} \\ &\textbf{si } \textit{posSeed} == |\textit{motif}| \textbf{ alors } \textit{affiche(nvKmer)}; \\ &\textbf{sinon } \textbf{si } \textit{motif[posMotif]} \neq 0 \textbf{ alors} \\ &\textit{kmersExpandRec(kmer, nvKmer + kmer[posKmer], posMotif + 1, posKmer + 1)}; \end{split}
```

kmerExpand

kmerExpandRec

Mots Minimaux absents (MMA)

Mots Minimaux absents (MMA)

Plus long sous-mot commun (PLS)

Plus long sous-mot commun (PLS)

Résultats obtenus

Résultats : KmersDel et kmersExpand

Résultats : Mots Minimaux absents (MMA)

Résultats : Plus long sous-mot commun (PLS)

