Let $\Sigma$ an alphabeth of events $\Sigma = \{\sigma_1 \ldots \sigma_n\}$
a transaction $w$ over $\sigma$ is an element of

$\Sigma^*$ | **example** $\Sigma = \{A, T, C, G\}$ and a word $w$ is AACCTTG

given $w, \overline{w} \in \overline{\Sigma}$ we say:

$w$ is a subsequence of $w'$ (written $w \sqsubseteq w'$)
if and only if there exists a sequence of
indexes $i_1 < \ldots < i_n$ with $n = |w'|$
such that $\overline{w}[i_1] \ldots \overline{w}[i_n] = w$

**example:** $\quad$ ATT $\sqsubseteq$ C<u>AT</u>G<u>T</u> (take $1, 2, 4$ as indexes)

$\qquad\qquad$ while

$\qquad\qquad$ TAT $\not\sqsubseteq$ CATGT ( there is no A following a T in CATGT)

Given a word $w$ we define its <u>preceding subsequent</u> as
the words $w'$ such that $\underline{w' \sqsubseteq w}$ and $\underline{|w'| = |w| - 1}$
we denote them with $w' < w$

Apriori_for_seq $(\gamma, \varepsilon) \longrightarrow \left\{ \overline{w} : \dfrac{\sum\limits_{w \in Dom(T),\ \overline{w} \in w} w}{\sum\limits_{w \in Dom(T)} \gamma(w)} \geq \varepsilon \right\}$

$\qquad\qquad\qquad\qquad\qquad$ support
$\qquad\qquad\qquad\qquad\qquad$ multiset of sequences $\qquad\qquad Sup_\gamma(w)$

$R_1 = \{\sigma : Sup_\gamma(\sigma) \geq \varepsilon\}$
$k = 1$
<u>while</u> $R_k \neq \emptyset$ <u>do</u>:

$\quad$ $R_{k+1} = \emptyset$
$\quad$ <u>for each</u> $w \in R_k,\ \overline{w} \in R_1$ <u>do</u>:

$\qquad$ if $\left( \begin{array}{l} \{w' . \overline{w} : w' < w\} \subseteq R_k \\ \underline{and}\ Sup_\gamma(w . \overline{w}) \geq \varepsilon \end{array} \right)$ <u>then</u>: $R_{k+1} = R_{k+1} \cup \{w . \overline{w}\}$

$\quad$ $k = k + 1$

<u>return</u> $\bigcup\limits_{i=1}^{k-1} R_k$


**Exercise 1**
Using the Notebook Apriori Sepsis
$\quad$ (a) Implement the Apriori_for_seq explained above,
$\quad$ (b) Test the code (a) on the sepsis sequence
$\qquad$ extraction provided in the notebook. ($\varepsilon = 0.05$)
$\quad$ (c) Extract the association rules based on sequences
$\qquad$ on the result of algorithm (a) with
$\qquad$ confidence $\delta = 0.8$
$\qquad$ (hint: join the mined patterns $X$ and $Y$ with
$\qquad\qquad X \sqsubseteq Y$ and $X \neq Y$)