# Auxiliary notation: subset of $k$-minimal values

Given the set $X \subset \mathbb{N}$, and some $k \in \mathbb{N}$, for the sake of convenience let's introduce an auxiliary notation for the subset of $k$-minimal values: $min_k(X) \subseteq X$, such that $\forall a \in min_k(X), \forall b \in X \setminus min_k(X)$ it follows, that $a < b$, and $|min_k(X)| = k$ (in case if $|X| < k$, then $min_k(X) = X$).

# K smallest edit distances

Let $lev_{a,b}(i,j)$ be the set of $k$ smallest edit distances between the first $i$ characters of the string $a$ and the first $j$ characters of the string $b$.

Let's introduce the following sets:

(1)

$$I_{a,b}(i,j) := \{d+1 \mid d \in lev_{a,b}(i-1,j)\} \qquad \text{The set of } k \text{ edit distances obtained through insertions}$$
$$D_{a,b}(i,j) := \{d+1 \mid d \in lev_{a,b}(i,j-1)\} \qquad \text{The set of } k \text{ edit distances obtained through deletions}$$
$$S_{a,b}(i,j) := \{d+1_{(a_i \neq b_j)} \mid d \in lev_{a,b}(i-1,j-1)\} \quad \text{The set of } k \text{ edit distances obtained through substitutions}$$
$$E_{a,b}(i,j) := I_{a,b}(i,j) \cup D_{a,b}(i,j) \cup S_{a,b}(i,j) \qquad \text{The set of all } 3k \text{ edit distances}$$

Where $1_{(a_i \neq b_j)}$ is the indicator function equal to 0 when $a_i = b_j$ and equal to 1 otherwise.

Then the set of $k$ smallest edit distances between the first $i$ characters of the string $a$ and the first $j$ characters of the string $b$ is defined as follows:

(2)
$$\begin{cases} lev_{a,b}(i,j) := \{max(i,j)\} & \text{when } i = 0 \text{ or } j = 0 \\ lev_{a,b}(i,j) := min_k(E_{a,b}(i,j)) & \text{otherwise} \end{cases}$$

# Proof of correctness

### Lemma 1

Whenever the arbitrary natural number $a$ doesn't belong to the set $min_k(X)$ and there exists some item of $min_k(X)$ which is greater than $a$, it means that $a \notin X$.

More formally, let's prove the following statement for any $a \in \mathbb{N}$ and any $X \subset \mathbb{N}$:

(3) $$\left( (a \notin min_k(X)) \wedge (\exists b \in min_k(X), a < b) \right) \Rightarrow \left( a \notin X \right)$$

### Proof

According to the contrapositive proof scheme, let's show that:

(4) $$\neg \left( a \notin X \right) \Rightarrow \neg \left( (a \notin min_k(X)) \wedge (\exists b \in min_k(X), a < b) \right)$$

Which is equivalent to:

(5) $$\left( a \in X \right) \Rightarrow \left( (a \in min_k(X)) \vee (\forall b \in min_k(X), a \geq b) \right)$$

The latter statement is equivalent to:

(6) $$\left( a \in X \right) \Rightarrow \left( (a \in min_k(X)) \vee (a \in X \setminus min_k(X)) \right)$$

Which is a tautology. ∎

## Proof by the smallest counterexample

*Induction Basis:*

In case if $i = 0$ or $j = 0$ there is possible only one edit distance, hence the set $lev_{a,b}(i,j)$ contains only one item, namely $max(i,j)$. As far as there is only one possible edit distance - it means, that $lev_{a,b}(i,j) = max_k(\{max(i,j)\}) = \{max(i,j)\}$, which complies to the definition of the subset of $k$-minimal values.

*Induction Hypothesis:*

- The set $lev_{a,b}(i-1,j)$ contains the $k$ minimal edit distances between the first $i-1$ characters of the string $a$ and the first $j$ characters of the string $b$.
  Hence, for every edit distance $y$ between the first $i-1$ characters of the string $a$ and the first $j$ characters of the string $b$, such that $y \notin lev_{a,b}(i-1,j)$ it follows, that $\forall x \in lev_{a,b}(i-1,j), y > x$.

- The set $lev_{a,b}(i,j-1)$ contains the $k$ minimal edit distances between the first $i$ characters of the string $a$ and the first $j-1$ characters of the string $b$.
  Hence, for every edit distance $y$ between the first $i$ characters of the string $a$ and the first $j-1$ characters of the string $b$, such that $y \notin lev_{a,b}(i,j-1)$ it follows, that $\forall x \in lev_{a,b}(i,j-1), y > x$.

- The set $lev_{a,b}(i-1,j-1)$ contains the $k$ minimal edit distances between the first $i-1$ characters of the string $a$ and the first $j-1$ characters of the string $b$.
  Hence, for every edit distance $y$ between the first $i-1$ characters of the string $a$ and the first $j-1$ characters of the string $b$, such that $y \notin lev_{a,b}(i-1,j-1)$ it follows, that $\forall x \in lev_{a,b}(i-1,j-1), y > x$.

*Inductive Step:*

We want to show, that the *Induction Hypothesis* implies, that the set $lev_{a,b}(i,j)$ contains the $k$ smallest edit distances between the first $i$ characters of the string $a$ and the first $j$ characters of the string.

**For the sake of contradictions** let's assume, that the set $lev_{a,b}(i,j)$ doesn't contain the $k$ smallest edit distances. Thus, there there exists some edit distance $y \notin lev_{a,b}(i,j)$ between the first $i$ characters of the string $a$ and the first $j$ characters of the string $b$, such that $\exists x \in lev_{a,b}(i,j)$ for which $y < x$:

$$(7) \quad \Big( y \notin lev_{a,b}(i,j) \Big) \wedge \Big( \exists x \in lev_{a,b}(i,j), y < x \Big)$$

Let's rewrite the expression in a following way:

$$(8)$$

$$\left(y \notin lev_{a,b}(i,j)\right) \wedge \left(\exists x \in lev_{a,b}(i,j), y < x\right)$$

$$\Leftrightarrow \left(y \notin min_k(E_{a,b}(i,j))\right) \wedge \left(\exists x \in min_k(E_{a,b}(i,j)), y < x\right) \qquad \text{By definition of } lev_{a,b}(i,j)$$

$$\Leftrightarrow \left(y \notin min_k(E_{a,b}(i,j))\right) \wedge \left(\exists x \in min_k(E_{a,b}(i,j)), y < x\right) \wedge \left(\exists x \in min_k(E_{a,b}(i,j)), y < x\right) \qquad \text{Idempotence}$$

$$\Rightarrow \left(y \notin E_{a,b}(i,j)\right) \wedge \left(\exists x \in min_k(E_{a,b}(i,j)), y < x\right) \qquad \text{According to the Lemma 1}$$

$$\Leftrightarrow \left(y \notin I_{a,b}(i,j) \cup D_{a,b}(i,j) \cup S_{a,b}(i,j)\right) \wedge \left(\exists x \in min_k(E_{a,b}(i,j)), y < x\right) \qquad \text{By definition of } E_{a,b}(i,j)$$

$$\Leftrightarrow \left(y \notin I_{a,b}(i,j)\right) \wedge \left(y \notin D_{a,b}(i,j)\right) \wedge \left(y \notin S_{a,b}(i,j)\right) \wedge \left(\exists x \in min_k(E_{a,b}(i,j)), y < x\right) \qquad \text{De-Morgan's law}$$

$$\Leftrightarrow \left((y-1) \notin lev_{a,b}(i-1,j)\right) \wedge \left((y-1) \notin lev_{a,b}(i,j-1)\right) \wedge \qquad \text{By definition of } I_{a,b}(i,j),$$
$$\wedge \left((y-1_{(a_i \neq b_j)}) \notin lev_{a,b}(i-1,j-1)\right) \wedge \left(\exists x \in min_k(E_{a,b}(i,j)), y < x\right) \qquad D_{a,b}(i,j) \text{ and } S_{a,b}(i,j)$$

$$\Rightarrow \left(\forall m \in lev_{a,b}(i-1,j), y-1 > m\right) \wedge \left(\forall n \in lev_{a,b}(i,j-1), y-1 > n\right) \wedge \qquad \text{By Induction Hypothesis}$$
$$\wedge \left(\forall u \in lev_{a,b}(i-1,j-1), y-1_{(a_i \neq b_j)} > u\right) \wedge \left(\exists x \in min_k(E_{a,b}(i,j)), y < x\right)$$

$$\Leftrightarrow \left(\forall m \in I_{a,b}(i,j), y > m\right) \wedge \left(\forall n \in D_{a,b}(i,j), y > n\right) \wedge \qquad \text{By definition of } I_{a,b}(i,j),$$
$$\wedge \left(\forall u \in S_{a,b}(i,j), y > u\right) \wedge \left(\exists x \in min_k(E_{a,b}(i,j)), y < x\right) \qquad D_{a,b}(i,j) \text{ and } S_{a,b}(i,j)$$

$$\Leftrightarrow \left(\forall m \in I_{a,b}(i,j) \cup D_{a,b}(i,j) \cup S_{a,b}(i,j), y > m\right) \wedge \left(\exists x \in min_k(E_{a,b}(i,j)), y < x\right) \qquad \text{Reordering}$$

$$\Leftrightarrow \left(\forall m \in E_{a,b}(i,j), y > m\right) \wedge \left(\exists x \in min_k(E_{a,b}(i,j)), y < x\right) \qquad \text{By definition of } E_{a,b}(i,j)$$

$$\Rightarrow \left(\forall m \in min_k(E_{a,b}(i,j)), y > m\right) \wedge \left(\exists x \in min_k(E_{a,b}(i,j)), y < x\right) \qquad \text{Because } min_k(X) \subseteq X$$

$$\Leftrightarrow \left(\forall m \in lev_{a,b}(i,j), y > m\right) \wedge \left(\exists x \in lev_{a,b}(i,j), y < x\right) \qquad \text{By definition of } lev_{a,b}(i,j)$$

$$\Rightarrow Contradiction.$$

Hence, our assumption was wrong. Consequently, the set $lev_{a,b}(i,j)$ contains the $k$-minimal edit distances.
∎