

Rozdział 1

Odległości na przestrzeni ciągów znaków

1.1. Podstawowe definicje

CZY TO JAKO CIĄGŁY TEKST CZY WSZYSTKO PISAĆ W DEFINICJI???

Niech $\Sigma = \{\Sigma_i\}$ będzie skończonym uporządkowanym alfabetem o wielkości $|\Sigma|$. *Napisem* nazywamy skończony ciąg znaków z Σ . Zbiór wszystkich napisów o długości n nad Σ jest oznaczony przez Σ^n , podczas gdy przez $\Sigma^* = \bigcup_{n=1}^{\infty} \Sigma^n$ rozumiemy zbiór wszystkich napisów utworzonych ze znaków z Σ [1].

O ile nie podano inaczej, używamy zmiennych s, t, u, v, w, x, y jako oznaczenie napisów oraz a, b, c do oznaczenia napisów jednoznakowych albo po prostu *znaków*. Pusty napis jest oznaczany przez ε . Przez $|s|$, dla każdego napisu $s \in \Sigma^*$, rozumiemy jego długość, czyli liczbę znaków w napisie. Ciąg zmiennych oznaczających napisy i/lub znaki oznaczają ich złączenie [1]. Dla rozróżnienia napisów od zmiennych reprezentujących napis, te pierwsze oznaczamy pismem maszynowym, np. **napis**.

Poprzez s_i rozumiemy i -ty znak z napisu s , dla każdego $i \in \{1, \dots, |s|\}$. Podciąg kolejnych przylegających do siebie znaków z napisu nazywamy *podnapisem*. Podnapisem napisu s , który zaczyna się od i -tego znaku, a kończy na j -tym znaku, oznaczamy przez $s_{i:j}$, tj. $s_{i:j} = s_i s_{i+1} \dots s_j$ dla $i < j$. Zakładamy również, że jeśli $j < i$, to $s_{i:j} = \varepsilon$ [1, 6].

Założmy, że napis s jest reprezentacją złączenia trzech, być może pustych, podnapisów w, x i y , tj. $s = wxy$. Wówczas podnapis w nazywamy *prefiksem*, natomiast podnapis y – *sufiksem* [1].

Podnapis złożony z kolejnych znaków napisu, o ustalonej długości q jest nazywany *q -gramem*. q -gramy o q równym jeden, dwa lub trzy mają specjalne nazwy: *unigram*, *bigram* i *trigram*. Jeśli $q > |s|$, to q -gramy napisu s są napisami pustymi [1].

Przykład 1.1. Niech Σ będzie alfabetem złożonym z 26 małych liter alfabetu łacińskiego oraz niech $s = \text{ela}$. Wówczas mamy $|s| = 3$, $s \in \Sigma^3$ oraz $s \in \Sigma$. Co więcej, mamy $s_1 = \text{e}$, $s_2 = \text{l}$, $s_3 = \text{a}$. Podnapis 1 : 2 napisu s to $s_{1:2} = \text{el}$. W napisie tym mamy do czynienia jedynie z q -gramami o q równym jeden, dwa oraz trzy: **e**, **l**, **a**; **el**, **la** oraz **ela** odpowiednio.

1.2. Odległości na przestrzeni ciągów znaków

W tym podrozdziale zajmiemy się odległościami na przestrzeni ciągów znaków. Można je podzielić na trzy grupy:

- oparte na operacjach edycyjnych (*edit operations*),
- oparte na q -gramach,
- miary heurystyczne.

BLA BLA JAKIEŚ LANIE WODY O METRYKACH Pierwszy rodzaj odległości jest najczęściej używany w algorytmach zajmujących się optymalnym dopasowaniem, dlatego też poświęcimy mu największą część niniejszego rozdziału. Odległości oparte na q -gramach ... (TUTAJ COŚ O NICH). Natomiast miary heurystyczne są rzadko stosowane, będąc zazwyczaj używane w konkretnych przypadkach. Miary te miały jednak swój wkład w historię optymalizacji napisów, zatem pokrótce opiszemy je pod koniec tego rozdziału.

1.2.1. Odległości oparte na operacjach edycyjnych

HISTORIA ODLEGŁOŚCI EDYCYJNYCH?

Ścieżka edycyjna i bazowe operacje edycyjne. *Odległość edycyjna* $ED(s, t)$ pomiędzy dwoma napisami s i t to minimalna liczba operacji edycyjnych potrzebna do przetworzenia s w t (i ∞ , gdy taki ciąg nie istnieje) [3]. *Ścieżką odległości edycyjnej* nazywamy minimalną liczbę nie nakładających się operacji edycyjnych, które pozwalają przekształcić jeden napis w drugi, i które nie przekształcają dwa razy tego samego podnapisu [1].

Napis może zostać przetworzony w drugi poprzez wykonanie na nim ciągu przekształceń jego podnapisów. Ten ciąg nazywany jest *ścieżką edycyjną (śladem edycji?)*, podczas gdy przekształcenia są nazywane *bazowymi operacjami edycyjnymi*. Bazowe operacje edycyjne, które polegają na mapowaniu napisu s w napis t , są oznaczane przez $s \rightarrow t$. Zbiór wszystkich bazowych operacji edycyjnych oznaczamy przez \mathbb{B} [1].

Bazowe operacje edycyjne są zazwyczaj ograniczone do:

- usunięcie znaku: $l \rightarrow \varepsilon$, tj. usunięcie litery l , np. $ela \rightarrow ea$,
- wstawienie znaku: $\varepsilon \rightarrow k$, tj. wstawienie litery k , np. $ela \rightarrow elka$,
- zamiana znaku: $e \rightarrow a$, tj. zamiana litery e na a , np. $ala \rightarrow ela$,
- transpozycja: $el \rightarrow le$, tj. przestawienie dwóch przylegających liter e i l , np. $ela \rightarrow lea$.

Własność 1.1. Zakładamy, że \mathbb{B} spełnia następujące własności [1]:

- jeśli $s \rightarrow t \in \mathbb{B}$, to odwrotna operacja $t \rightarrow s$ również należy do \mathbb{B} ;
- $a \rightarrow a \in \mathbb{B}$ (operacja identycznościowa dla jednego znaku należy do \mathbb{B});
- zbiór \mathbb{B} jest zupełny: dla dwóch dowolnych napisów s i t , istnieje ślad edycji, który przekształca s w t .

Zauważmy, że zbiór \mathbb{B} nie musi być skończony.

Odległość edycyjna. Podobieństwo dwóch napisów może być wyrażone jako długość ścieżki edycyjnej, dzięki której jeden napis zostaje przekształcony w drugi:

Definicja 1.1. *Mając dany zbiór bazowych operacji edycyjnych, odległość edycyjna $ED(s, t)$ jest równa długości najkrótszej ścieżki edycyjnej, która przekształca napis s w napis t . Najkrótsza ścieżka, która przekształca napis s w napis t jest nazywana optymalną ścieżką edycyjną [1].*

Przykładowe odległości edycyjne: Hamminga, najdłuższego wspólnego podnapisu (*longest common substring*), Levenshteina, optymalnego dopasowania napisów (*optimal string alignment*), Damareu-Levenshteina. Odległości te różnią się zbiorem bazowych operacji edycyjnych. Jeśli w zbiorze tym znajduje się tylko zamiana znaków, to mamy do czynienia z odległością Hamminga. Gdy zbiór bazowych operacji edycyjnych zawiera wstawienia i usunięcia znaków, to jest to odległość najdłuższego wspólnego podnapisu. Gdyby \mathbb{B} powiększyć o zamianę znaków, to otrzymamy odległość Levenshteina. Dwie ostatnie odległości, tj. optymalnego dopasowania napisów oraz Damareu-Levenshteina, mają w zbiorze bazowych operacji edycyjnych usunięcie, wstawienie, zamianę oraz transpozycję znaków.

Przykład 1.2. JAKIŚ PRZYKŁĄD SCIEZKI I OPTYMALNEJ SCIEZKI.

Definicja odległości edycyjnej może być również interpretowana jako minimalny koszt, dzięki któremu przekształcamy jeden napis w drugi. Definicję można uogólnić na dwa sposoby. Po pierwsze, bazowe operacje edycyjne mogą mieć przydzielone koszty (wagi) $\delta(a \rightarrow b)$ [7]. Zazwyczaj koszt każdej operacji wynosi jeden, jednak można, na przykład, nadać transpozycji mniejszy koszt niż operacji wstawienia znaku. Dalej, można rozszerzyć funkcję kosztu δ na ścieżkę edycyjną $E = a_1 \rightarrow b_1, a_2 \rightarrow b_2, \dots, a_{|E|} \rightarrow b_{|E|}$ poprzez $\delta(E) = \sum_{i=1}^{|E|} \delta(a_i \rightarrow b_i)$ [1].

Odtąd poprzez odległość między napisem s a napisem t będziemy rozumieć minimalny ze wszystkich możliwych kosztów ścieżek przekształcających s w t . Odległości zdefiniowane w ten sposób zazwyczaj są nazywane *uogólnionymi* odległościami edycyjnymi.

Po drugie, zbiór operacji edycyjnych \mathbb{B} może zostać rozszerzony o ważone zamiany (substytucje) (pod)napisów, zamiast operacji edycyjnych wykonywanych na pojedynczych znakach [4]. Odległości zdefiniowane w ten sposób zazwyczaj są nazywane *rozszerzonymi* odległościami edycyjnymi. Przykładowo, \mathbb{B} może zawierać operację $x \rightarrow ks$ o koszcie jednostkowym. Wówczas rozszerzona odległość pomiędzy napisami *xero* i *ksero* wynosi jeden, podczas gdy standardowa (zwykła, nierozszerzona) odległość wyniosłaby dwa [1].

Definicja 1.2. *Mając dany zbiór bazowych operacji edycyjnych \mathbb{B} oraz funkcję δ , która nadaje koszt wszystkim bazowym operacjom edycyjnym z \mathbb{B} , ogólna(?) odległość edycyjna pomiędzy napisami s i t jest zdefiniowana jako minimalny koszt ścieżki edycyjnej, która przekształca s w t [1].*

Zazwyczaj koszt pojedynczej operacji z \mathbb{B} jest równy jeden. Czasem jednak nadaje się poszczególnym operacjom różne koszty, dając np. transpozycji mniejszą wagę niż wstawieniu znaku. Gdy koszt wszystkich operacji jest równy jeden, to mówimy po prostu o odległości edycyjnej, natomiast gdy różne operacje mają różne wagi, to mówimy o *uogólnionej* odległości edycyjnej.

Własność 1.2. Zakładamy, że funkcja kosztu $\delta(s \rightarrow t)$ ma następujące własności [1]:

- $\delta(s \rightarrow t) \in \mathbb{R}$ (koszt operacji jest liczbą rzeczywistą),
- $\delta(s \rightarrow t) = \delta(t \rightarrow s)$ (symetria),
- $\delta(s \rightarrow t) \geq 0$, $\delta(s \rightarrow s) = 0$ i $\delta(s \rightarrow t) = 0 \Rightarrow s = t$ (dodatnia określoność ??),
- $\forall \gamma > 0$ zbiór bazowych operacji $\{s \rightarrow t \in \mathbb{B} \mid \delta(s \rightarrow t) < \gamma\}$ jest skończony (skończoność podzbioru bazowych operacji, których koszt jest ograniczony z góry).

Zauważmy, że ostatnia własność jest zawsze spełniona dla skończonego zbioru \mathbb{B} .

Twierdzenie 1.3. Z własności 1.1 i 1.2 wynika, że:

- dla każdych dwóch napisów s i t , istnieje ścieżka o minimalnym koszcie, tj. właściwie (prawdziwie, odpowiednio?) zdefiniowana odległość edycyjna z s do t [1],
- ogólna odległość edycyjna z definicji 1.2 jest metryką [7].

Dowód. Żeby udowodnić, że $ED(s, t)$ jest metryką, musimy pokazać, że $ED(s, t)$ istnieje, jest dodatnio określona, symetryczna oraz subaddytywna (tj. spełnia nierówność trójkąta).

Z własności 1.2 wynika, że funkcja kosztu jest nieujemna (JA TEGO NIW WIDZE) i że tylko identyczność ma koszt równy zero. Stąd, bez utraty ogólności, możemy rozważyć jedynie takie ścieżki edycyjne, które nie zawierają operacji identycznościowych. Zatem, jeśli $s = t$, to jedyna optymalna ścieżka (która nie zawiera operacji identycznościowych) jest pusta i ma zerowy koszt. Jeśli $s \neq t$, to z zupełności zbioru bazowych operacji edycyjnych wynika, że istnieje jedna lub więcej ścieżek edycyjnych, które przekształcają s w t . Wszystkie te ścieżki składają się z operacji edycyjnych o ściśle dodatnim koszcie.

Niech γ będzie kosztem ścieżki przekształcającej s w t . Rozważmy zbiór A ścieżek edycyjnych, które przekształcają s w t i których koszt jest ograniczony z góry przez γ . Zbiór A jest niepusty i składa się z operacji edycyjnych o dodatnim koszcie mniejszym niż γ . Zbiór operacji bazowych, których koszt jest ograniczony z góry przez γ jest skończony, co dowodzi, że zbiór A jest również skończony. Ponieważ A jest niepusty i skończony, to ścieżki edycyjne o minimalnym (dodatnim) koszcie istnieją i należą do A . Stąd, $ED(s, t) > 0$ dla $s \neq t$, tj. odległość edycyjna jest dodatnio określona.

Aby udowodnić symetrię odległości edycyjnej, rozważmy optymalną ścieżkę E , która przekształca s w t , oraz odpowiadającą jej odwrotną ścieżkę E_r , która przekształca t w s . Równość ich kosztów $\delta(E) = \delta(E_r)$ wynika z symetrii funkcji kosztu i symetrii zbioru operacji bazowych \mathbb{B} .

Aby pokazać subaddytywność, rozważmy optymalną ścieżkę E_1 , która przekształca s w t , optymalną ścieżkę E_2 , która przekształca t w u , oraz złożenie ścieżek $E_1 E_2$, które przekształca s w u . Z tego, że $\delta(E_1 E_2) = \delta(E_1) + \delta(E_2) = ED(s, t) + ED(t, u)$ oraz $\delta(E_1 E_2) \geq ED(s, u)$ wynika, że $ED(s, t) + ED(t, u) \geq ED(s, u)$. ■

Odległość edycyjna jest metryką, nawet gdy funkcja kosztu δ nie jest subaddytywna. Co więcej, ponieważ ciąg nakładających się operacji, które przekształcają s w t , mogą mieć mniejszy koszt niż $\delta(s \rightarrow t)$, $\delta(s \rightarrow t)$ może być większe niż $ED(s, t)$. Rozważmy, na przykład, następujący alfabet: $\{a, b, c\}$, gdzie symetria i brak subaddytywności funkcji δ jest zdefiniowana

następująco:

$$\begin{aligned}\delta(\mathbf{a} \rightarrow \mathbf{c}) &= \delta(\mathbf{b} \rightarrow \mathbf{c}) = 1 \\ \delta(\mathbf{a} \rightarrow \varepsilon) &= \delta(\mathbf{b} \rightarrow \varepsilon) = \delta(\mathbf{c} \rightarrow \varepsilon) = 2 \\ \delta(\mathbf{a} \rightarrow \mathbf{b}) &= 3\end{aligned}$$

Można zobaczyć, że $3 = \delta(\mathbf{a} \rightarrow \mathbf{b}) > \delta(\mathbf{a} \rightarrow \mathbf{c}) + \delta(\mathbf{c} \rightarrow \mathbf{b}) = 2$. Stąd optymalna ścieżka edycyjna $(\mathbf{a} \rightarrow \mathbf{c}, \mathbf{c} \rightarrow \mathbf{b})$ przekształca \mathbf{a} w \mathbf{b} z kosztem równym 2.

Ścisła odległość edycyjna. Subaddytywność odległości edycyjnej pozwala używać metod właściwych przestrzeniom metrycznym. Niemniej jednak, problem minimalizacji zbioru nakładających się operacji edycyjnych, może być trudny. Aby zrównoważyć złożoność obliczeniową, zazwyczaj używana jest funkcja podobieństwa, zdefiniowana jako minimum kosztu *ściślej ścieżki edycyjnej*. Ta ostatnia nie zawiera nakładających się na siebie operacji edycyjnych i nie modyfikuje tego samego podnapisu więcej niż raz. Odpowiadająca jej odległość edycyjna nazywana jest *ściłą odległością edycyjną* [1].

Lemat 1.4. *Dowolna nieściśła odległość edycyjna ogranicza z dołu odpowiadającą jej ściłą odległość edycyjną.*

Lemat 1.5. *Ścisła odległość Levenshteina o jednostkowym koszcie operacji jest równa nieściślej odległości Levenshteina o jednostkowym koszcie operacji.*

Powyższe wynika natychmiast z obserwacji, że optymalna ścieżka edycyjna zawiera jednoznakowe usunięcia, wstawienia oraz zamiany, które nigdy nie modyfikują znaku więcej niż raz.

Lemat 1.6. *Nieściśła odległość Damerau-Levenshteina oraz ściśła odległość Damerau-Levenshteina są różnymi funkcjami. Co więcej ściśła odległość Damerau-Levenshteina nie jest metryką, gdyż nie jest subaddytywna [1].*

Dowód. Ścisła odległość Damerau-Levenshteina traktuje transpozycję (tj. zamianę dwóch przylegających do siebie znaków) jako bazową operację edycyjną. Aby udowodnić lemat podamy przykład, w którym zakaz modyfikacji znaków już stransponowanych odróżnia odległość Damerau-Levenshteina od ściślej odległości Damerau-Levenshteina [1]. ■

Rozważmy napisy \mathbf{ab} , \mathbf{ba} oraz \mathbf{acb} . Z jednej strony, najkrótsza nieściśła ścieżka edycyjna, która przekształca \mathbf{ba} w \mathbf{acb} , tj. $(\mathbf{ba} \rightarrow \mathbf{ab}, \varepsilon \rightarrow \mathbf{c})$ zawiera dwie operacje: najpierw, zamienia znaki \mathbf{a} i \mathbf{b} , a następnie wstawia \mathbf{c} pomiędzy nie. Zauważmy, że wstawienie przekształca już transformowany napis. Jednakowoż, jeśli kolejne przekształcenia tego samego podnapisu są wykluczone, to najkrótsza ścieżka edycyjna, która przekształca \mathbf{ba} w \mathbf{acb} , składa się z trzech operacji edycyjnych, np. $(\mathbf{ba} \rightarrow \varepsilon, \varepsilon \rightarrow \mathbf{b})$. Stąd, nieściśła odległość edycyjna jest równa dwa, podczas gdy ściśła odległość wynosi trzy [1].

Ścisła odległość Damerau-Levenshteina nie spełnia nierówności trójkąta, gdyż

$$\mathbf{ba} \xrightarrow[1]{\text{transp. } b \text{ i } a} \mathbf{ab} + \mathbf{ab} \xrightarrow[1]{\text{wst. } c} \mathbf{acb},$$

natomiast

$$\mathbf{ba} \xrightarrow[1]{\text{us. } b} \mathbf{a} \xrightarrow[1]{\text{wst. } c} \mathbf{ac} \xrightarrow[1]{\text{wst. } b} \mathbf{acb}$$

zatem

$$2 = ED(\mathbf{ba}, \mathbf{ab}) + ED(\mathbf{ab}, \mathbf{acb}) \leq ED(\mathbf{ba}, \mathbf{acb}) = 3.$$

Ponieważ ścisła i nieścisła odległość Damerau-Levenshteina są różnymi funkcjami, tę pierwszą nazywa się często *odległością optymalnego dopasowania napisów*. Od tego momentu w niniejszej pracy ścisłą odległość Damerau-Levenshteina nazywamy odległością optymalnego dopasowania napisów, natomiast nieścisłą odległość Damerau-Levenshteina nazywamy odległością Damerau-Levenshteina [6].

Optymalne dopasowanie.

Niech napisy s i t zostaną podzielone na tę samą liczbę, być może pustych, podnapisów: $s = s_1 s_2 \dots s_l$ i $t = t_1 t_2 \dots t_l$, takich, że $s_i \rightarrow t_i \in \mathbb{B}$. Co więcej, zakładamy, że s_i i t_j nie mogą być puste dla $i = j$. Mówimy, że ten podział definiuje *dopasowanie* $A = (s_1 s_2 \dots s_l, t_1 t_2 \dots t_l)$ pomiędzy napisami s i t , w którym podnapis s_i jest dopasowane do podnapisu t_i [1].

Dopasowanie reprezentuje ścisłą ścieżkę edycyjną $E = s_1 \rightarrow t_1, s_2 \rightarrow t_2, \dots, s_l \rightarrow t_l$. Definiujemy *koszt dopasowania* A jako koszt odpowiadającej mu ścieżki edycyjnej i oznaczamy przez $\delta(A)$:

$$\delta(A) = \sum_{i=1}^l \delta(s_i \rightarrow t_i) \quad (1.1)$$

Optymalne dopasowanie to dopasowanie o najmniejszym koszcie.

PIEKNY RYSUNEK

Przykład 1.3. Przykład optymalnego dopasowania pomiędzy słowami **foczka** i **kozak** prezentuje rys. TUTAJNUMER. Odpowiadająca mu ścieżka edycyjna składa się z zamiany $f \rightarrow k$, usunięcia $c \rightarrow \varepsilon$ oraz transpozycji $ka \rightarrow ak$.

Warto zauważyć, że istnieje różnowartościowe (1-1) mapowanie pomiędzy zbiorem ścisłych ścieżek edycyjnych i zbiorem dopasowań: każda ścisła ścieżka edycyjna o minimalnym koszcie reprezentuje dopasowanie o najmniejszym koszcie i odwrotnie. Stąd można zastąpić problem znalezienia optymalnej ścisłej odległości edycyjnej poprzez problem znalezienia optymalnego dopasowania.

Obliczanie odległości edycyjnej. Główną zasadą dynamicznego algorytmu, liczącego koszt optymalnego dopasowania, jest wyrażenie kosztu dopasowania pomiędzy napisami s i t , używając kosztu dopasowania ich prefiksów. Rozważmy prefiks $s_{1:i}$ o długości i i prefiks $t_{1:j}$ o długości j napisów s i t odpowiednio. Załóżmy, że $A = (s_1 s_2 \dots s_l, t_1 t_2 \dots t_l)$ jest optymalnym dopasowaniem pomiędzy $s_{1:i}$ i $t_{1:j}$, którego koszt oznaczamy przez $C_{i,j}$ [1].

Używając równania 1.1 oraz definicji optymalnego dopasowania, łatwo pokazać, że $C_{i,j}$ może zostać policzone przy użyciu następującej ogólnej rekurencji [4]:

$$C_{0,0} = 0 \quad (1.2)$$

$$C_{i,j} = \min\{\delta(s_{i':i} \rightarrow t_{j':j}) + C_{i'-1,j'-1} \mid s_{i':i} \rightarrow t_{j':j} \in \mathbb{B}\} \quad (1.3)$$

Można zauważyć, że:

- koszt dopasowania napisów s i t jest równy $C_{|s|,|t|}$;
- wszystkie optymalne dopasowania mogą zostać wyznaczone poprzez odwracanie rekurencji 1.2 (przechodzenie od tyłu).

Rozważmy teraz odległość Hamminga, gdzie $s_{i':i} \rightarrow t_{j':j}$ to zamiany znaków o koszcie równym jeden. Stąd,

$$\delta(s_{i':i} \rightarrow t_{j':j}) = [s_{i':i} \neq t_{j':j}] \quad (1.4)$$

gdzie $[X]$ jest równe jeden, gdy warunek X jest spełniony, zero w przeciwnym przypadku. Co więcej, w tym przypadku możliwa jest tylko jedna kombinacja i' oraz j' , mianowicie $i' = i$ oraz $j' = j$. Co więcej, odległość ta jest zdefiniowana jedynie dla $|s| = |t|$, zatem $C_{i,j}$ może być policzone jedynie dla $i = j$. Wówczas definicja odległości Hamminga nie jest rekursywna (rekursyjna?) i można ją zapisać następująco:

Definicja 1.3. Odległością Hamminga *nazywamy* [2]:

$$d_{\text{hamming}}(s, t) = \begin{cases} \sum_{i=1}^{|s|} \delta(s_i \rightarrow t_i) = \sum_{i=1}^{|s|} [s_i \neq t_i], & \text{gdy } |s| = |t|, \\ \infty, & \text{w przeciwnym przypadku,} \end{cases}$$

Odległość Hamminga zlicza liczbę indeksów, na których dwa napisy mają różny znak. Odległość ta przyjmuje wartości ze zbioru $\{0, \dots, |s|\}$, gdy $|s| = |t|$, natomiast jest równa nieskończoności, gdy napisy mają różne długości.

[PIĘKNY RYSUNEK??]

Przykład 1.4. Odległość Hamminga między słowami **koza** i **foka** wynosi $d_{\text{hamming}}(\text{koza}, \text{foka}) = 2$, natomiast między słowami **kozak** i **foczka** wynosi ona $d_{\text{hamming}}(\text{kozak}, \text{foczka}) = \infty$.

Rozważmy teraz odległość najdłuższego wspólnego podnapisu, gdzie $s_{i':i} \rightarrow t_{j':j}$ to wstawienia i usunięcia znaków o koszcie równym jeden. Wówczas istnieją dwie kombinacje i' oraz j' z równania 1.4, odpowiadające usunięciu i wstawieniu, odpowiednio:

- $i' = i - 1$ oraz $j' = j$,
- $i' = i$ oraz $j' = j - 1$.

Uwzględniając powyższe uproszczenia, możemy następująco przepisać ogólną postać rekurencji 1.2 dla odległości najdłuższego wspólnego podnapisu [NIE WIEM CZY TO JEST DOBRZE!!!!]:

$$C_{i,j} = \min \begin{cases} 0, & \text{gdy } i = j = 0 \\ C_{i-1,j} + 1, & \text{gdy } i > 0 \\ C_{i,j-1} + 1, & \text{gdy } j > 0 \end{cases} \quad (1.5)$$

Odległość najdłuższego wspólnego przyjmuje wartości ze zbioru $\{0, |s| + |t|\}$, przy czym maksimum jest osiąganę, gdy s i t nie mają ani jednego wspólnego znaku. Odległość tę oznaczamy przez d_{lcs} .

Przykład 1.5. Odległość najdłuższego wspólnego podnapisu między napisami **kozak** i **foczka** wynosi: $d_{lsc}(\text{kozak}, \text{foczka}) = 5$, bo $\text{kozak} \xrightarrow[1]{us. k} \text{ozak} \xrightarrow[1]{us. a} \text{ozk} \xrightarrow[1]{wst. f} \text{fozk} \xrightarrow[1]{wst. c} \text{foczka}$.

Powyższy przykład pokazuje, że w ogólności nie ma unikalnej najkrótszej drogi transformacji jednego napisu w drugi, gdyż można zamienić kolejność usuwania (lub wstawiania) znaków i również uzyskać odległość równą 5. Można również usunąć z napisu znak **k** zamiast **a**, otrzymując taką samą odległość między napisami.

Jak sugeruje nazwa, odległość najdłuższego wspólnego podnapisu, ma też inną interpretację. Poprzez wyrażenie *najdłuższy wspólny podnapis* rozumiemy najdłuższy ciąg utworzony przez sparowanie znaków z *s* i *t* nie zmieniając ich porządku. Wówczas odległość ta jest rozumiana jako liczba niesparowanych znaków z obu napisów. W powyższym przykładzie może to być zwizualizowane następująco:

[PIĘKNY RYSUNEK]

Jak widać na rysunku, znaki **k**, **a**, **f**, **c** i **a** w pierwszym przypadku oraz **k**, **k**, **f**, **c** i **k** w drugim, pozostają bez pary, dając odległość równą 5.

Przejdźmy do odległości Levenshteina. Odległość ta dopuszcza, oprócz usunięć i wstawień, także zamiany znaków. Istnieją zatem trzy kombinacje i' oraz j' z równania 1.4:

- $i' = i - 1$ oraz $j' = j$,
- $i' = i$ oraz $j' = j - 1$,
- $i' = i - 1$ oraz $j' = j - 1$.

Stąd ogólna postać rekurencji 1.2 dla odległości Levenshteina może zostać przepisana następująco:

$$C_{i,j} = \min \begin{cases} 0, & \text{gdy } i = j = 0 \\ C_{i-1,j} + 1, & \text{gdy } i > 0 \\ C_{i,j-1} + 1, & \text{gdy } j > 0 \\ C_{i-1,j-1} + [s_i \neq t_j], & \text{gdy } i, j > 0 \end{cases} \quad (1.6)$$

Odległość Levenshteina oznaczamy przez d_{lv} .

Przykład 1.6. Odległość Levenshteina między napisami **kozak** i **foczka** wynosi: $d_{lv}(\text{kozak}, \text{foczka}) = 4$, bo $\text{kozak} \xrightarrow[1]{zm. k \text{ na } f} \text{fozak} \xrightarrow[1]{wst. c} \text{foczak} \xrightarrow[1]{zm. a \text{ na } k} \text{foczkk} \xrightarrow[1]{zm. k \text{ na } a} \text{foczka}$.

Powyższy przykład ilustruje, że dodatkowa elastyczność w porównaniu do odległości najdłuższego wspólnego podnapisu, daje mniejszą wartość odległości między napisami, jako że w przypadku pierwszego znaku potrzebujemy jedynie zamiany znaku, zamiast wstawienia i usunięcia znaków [6]. Co więcej, ścieżka edycyjna między tymi słowami może być inna i zawierać usunięcie i wstawienie zamiast dwóch ostatnich zamian znaków.

Przypomnijmy, że mówimy o uogólnionej odległości, gdy zmienimy koszty poszczególnych operacji na różne od jeden. Gdy za koszt przyjmiemy np. (0.1, 1, 0.3) dla usunięć, wstawień i

zamian znaków odpowiednio, to uogólniona odległość Levenshteina między napisami **koza** i **foka** wynosi: $d_{lv}(\text{koza}, \text{foka}) = 0.6$, bo $\text{koza} \xrightarrow[0.3]{\text{zm. } k \text{ na } f} \text{foza} \xrightarrow[0.3]{\text{zm. } z \text{ na } k} \text{foka}$.

Uogólniona odległość Levenshteina spełnia definicję metryki, gdy koszt usunięcia jest równy kosztowi wstawienia znaku. W przeciwnym przypadku nie spełnia ona założenia o symetrii. Jednakowoż, symetria zostaje zachowana przy jednoczesnej zamianie s i t oraz kosztów usunięcia i wstawienia znaku, jako że liczba usunięć znaków przy przetwarzaniu napisu s w napis t jest równa liczbie wstawień znaków przy przetwarzaniu napisu t w napis s [6]. Dobrze obrazuje to następujący przykład.

Przykład 1.7. Przyjmijmy za koszt usunięcia, wstawienia i zamiany znaku $(0.1, 1, 0.3)$ odpowiednio. Wówczas uogólniona odległość Levenshteina dla napisów **koza** i **foczka** wynosi:

$$d_{lv}(\text{koza}, \text{foczka}) = 0.5, \quad (1.7)$$

bo

$$\text{koza} \xrightarrow[0.3]{\text{zm. } k \text{ na } f} \text{foza} \xrightarrow[0.1]{\text{wst. } c} \text{focza} \xrightarrow[0.1]{\text{wst. } k} \text{foczka},$$

natomiast

$$d_{lv}(\text{foczka}, \text{koza}) = 2.3, \quad (1.8)$$

bo

$$\text{foczka} \xrightarrow[0.3]{\text{zm. } f \text{ na } k} \text{koczka} \xrightarrow[1]{\text{us. } c} \text{kozka} \xrightarrow[1]{\text{us. } k} \text{koza}.$$

Gdy za koszty przyjmiemy $(1, 0.1, 0.3)$, to uogólniona odległość Levenshteina wynosi:

$$d_{lv}(\text{koza}, \text{foczka}) = 2.3,$$

bo

$$\text{koza} \xrightarrow[0.3]{\text{zm. } k \text{ na } f} \text{foza} \xrightarrow[1]{\text{wst. } c} \text{focza} \xrightarrow[1]{\text{wst. } k} \text{foczka},$$

czyli analogicznie, jak w przypadku 1.8. Natomiast

$$d_{lv}(\text{foczka}, \text{koza}) = 0.5,$$

bo

$$\text{foczka} \xrightarrow[0.3]{\text{zm. } f \text{ na } k} \text{koczka} \xrightarrow[0.1]{\text{us. } c} \text{kozka} \xrightarrow[0.1]{\text{us. } k} \text{koza},$$

czyli analogicznie, jak w przypadku 1.7.

Zgodnie z lematem 1.5 nieściśła odległość Levenshteina jest równa ściśłej odległości Levenshteina. Z drugiej strony, ściśła odległość edycyjna jest równa kosztowi optymalnego dopasowania. Stąd rekurencja 1.2 liczy nieściśłą odległość Levenshteina. Częstym błędem jest, że następujące bezpośrednie uogólnienie, dodające transpozycję do zbioru bazowych operacji edycyjnych, rekurencji 1.6 liczy (nieściśłą) odległość Damerau-Levenshteina [1]:

$$C_{i,j} = \min \begin{cases} 0, & \text{gdy } i = j = 0 \\ C_{i-1,j} + 1, & \text{gdy } i > 0 \\ C_{i,j-1} + 1, & \text{gdy } j > 0 \\ C_{i-1,j-1} + [s_i \neq t_j], & \text{gdy } i, j > 0 \\ C_{i-2,j-2} + 1, & \text{gdy } s_i = t_{j-1}, s_{i-1} = t_j \text{ oraz } i, j > 1 \end{cases} \quad (1.9)$$

Jednak, rekurencja 1.9 liczy odległość optymalnego dopasowania napisów, czyli ściśle odległość Damerau-Levenshteina, która nie zawsze jest równa odległości Damerau-Levenshteina. Dla przykładu, odległość między napisami *ba* i *acb* wyliczona przy pomocy rekurencji 1.9 jest równa trzy, natomiast odległość Damerau-Levenshteina między tymi napisami wynosi dwa.

Rekurencyjna definicja odległości Damerau-Levenshteina została po raz pierwszy podana przez Lowrance'a i Wagnera [8]. W ich definicji zamiana zostaje zastąpiona poprzez minimalizację po możliwych transpozycjach między danym znakiem a wszystkimi nie przetransformowanymi znakami, przy czym koszt transpozycji wzrasta wraz z odległością między transponowanymi znakami. Innymi słowy, do \mathbb{B} należą wstawienia, usunięcia, zamiany oraz operacje $axb \rightarrow bya$ o koszcie równym $|x| + |y| + 1$. Mając tak zdefiniowane \mathbb{B} ogólna rekurencja 1.2 dla odległości Damerau-Levenshteina przedstawia się następująco:

$$C_{i,j} = \min \begin{cases} 0, & \text{gdy } i = j = 0 \\ C_{i-1,j} + 1, & \text{gdy } i > 0 \\ C_{i,j-1} + 1, & \text{gdy } j > 0 \\ C_{i-1,j-1} + [s_i \neq t_j], & \text{gdy } i, j > 0 \\ \min_{\substack{0 < i' < i, 0 < j' < j \\ s_i = t_{j'}, s_{i'} = t_j}} C_{i'-1,j'-1} + (i - i') + (j - j') - 1 \end{cases} \quad (1.10)$$

Co więcej, Lowrance i Wagner wykazali, że wewnętrzne minimum w rekurencji 1.10 jest osiągnięte dla największych $i' < i$ oraz $j' < j$, które spełniają $s_i = t_{j'}$ oraz $s_{i'} = t_j$. Odległość Damerau-Levenshteina oznaczamy przez d_{dl} .

Przykład 1.8. Odległość optymalnego dopasowania napisów oraz Damerau-Levenshteina między napisami *kozak* i *foczka* wynosi 3, bo *kozak* $\xrightarrow[1]{zm. k \text{ na } f}$ *f*ozak $\xrightarrow[1]{wst. c}$ *f*oczak $\xrightarrow[1]{transp. a \text{ i } k}$ *foczka*.

W przypadku odległości Levenshteina, optymalnego dopasowania napisów oraz Damerau-Levenshteina, maksymalna odległość między napisami s i t wynosi $\max\{|s|, |t|\}$. Jednak warto zauważyć, że gdy liczba dopuszczalnych operacji edycyjnych rośnie, to liczba dopuszczalnych ścieżek między napisami wzrasta, co pozwala czasem zmniejszyć odległość między napisami. Dlatego relację między zaprezentowanymi powyżej odległościami można podsumować następująco [6]:

$$\left. \begin{aligned} \infty &\geq |s| \geq d_{hamming}(s, t) \\ |s| + |t| &\geq d_{lcs}(s, t) \\ \max\{|s|, |t|\} &\end{aligned} \right\} \geq d_{lv}(s, t) \geq d_{osa}(s, t) \geq d_{dl}(s, t) \geq 0$$

Jako że odległości Hamminga i najdłuższego wspólnego podnapisu nie mają wspólnych bazowych operacji edycyjnych, to nie ma pomiędzy nimi porządku relacyjnego. Górne ograniczenie $|s|$ odległości Hamminga jest zachowane jedynie gdy $|s| = |t|$.

1.2.2. Odległości oparte q -gramach

Przypomnijmy, że q -gramem nazywamy napis składający się z q kolejnych (przylegających) znaków. q -gramy związane z napisem s są otrzymywane poprzez przesuwanie przez napis s

„okna” o szerokości q znaków i zapisaniu występujących q -gramów. Przykładowo digramy napisu `ela` to `el` i `la`. Oczywiście taka procedura nie ma sensu, gdy $q > |s|$ lub gdy $q = 0$. Z tego powodu definiujemy następujące przypadki brzegowe dla wszystkich odległości $d(s, t, q)$ opartych na q -gramach:

$$\begin{aligned} d(s, t, q) &= \infty, \text{ gdy } q > \min\{|s|, |t|\} \\ d(s, t, 0) &= \infty, \text{ gdy } |s| + |t| > 0 \\ d(\varepsilon, \varepsilon, 0) &= 0 \end{aligned}$$

Najprostszą odległością między napisami, opartą na q -gramach, otrzymuje się poprzez wypisanie unikalnych q -gramów w obu napisach i porównanie które są wspólne. Jeśli przez $\mathcal{Q}(s, q)$ oznaczymy zbiór unikalnych q -gramów występujących w napisie s , to możemy zdefiniować odległość Jaccarda [6]:

Definicja 1.4. Niech $\mathcal{Q}(s, q)$ oznacza zbiór unikalnych q -gramów występujących w napisie s . Wówczas odległość Jaccarda, $d_{jaccard}$, między napisami s i t definiuje się jako

$$d_{jaccard}(s, t, q) = 1 - \frac{|\mathcal{Q}(s, q) \cap \mathcal{Q}(t, q)|}{|\mathcal{Q}(s, q) \cup \mathcal{Q}(t, q)|}, \quad (1.11)$$

gdzie $|\cdot|$ oznacza licznosc zbioru.

Odległość Jaccarda przyjmuje wartości z przedziału $[0, 1]$, gdzie 0 odpowiada pełne pokrycie zbiorów, tj. $\mathcal{Q}(s, q) = \mathcal{Q}(t, q)$, natomiast 1 oznacza puste przecięcie, tj. $\mathcal{Q}(s, q) \cap \mathcal{Q}(t, q) = \emptyset$.

Przykład 1.9. Odległość Jaccarda między napisami `palce` i `pałace` dla $q = 2$ wynosi: $d_{jaccard}(\text{palce}, \text{pałace}, 2) \approx 0.714$, bo $\mathcal{Q}(\text{palce}, 2) = \{\text{pa}, \text{al}, \text{lc}, \text{ce}\}$, a $\mathcal{Q}(\text{pałace}, 2) = \{\text{pa}, \text{ał}, \text{ła}, \text{ac}, \text{ce}\}$, więc odległość wynosi $1 - \frac{2}{3} \approx 0.714$.

Odległość Jaccarda nie spełnia definicji metryki, bowiem $d_{jaccard}(s, t) = 0$ nawet wówczas, gdy $s \neq t$. Przykładowo $d_{jaccard}(\text{abaca}, \text{acaba}, 2) = 0$, mimo że $\text{abaca} \neq \text{acaba}$.

Inną odległością opartą na q -gramach jest odległość q -gramowa [????]. Otrzymuje się ją poprzez wylistowanie q -gramów występujących w obu napisach i policzenie q -gramów, które nie są wspólne [dla obu napisów] [6]. Formalnie można to zapisać następująco:

Definicja 1.5. Niech $s = s_1 s_2 \dots s_n$ będzie napisem z Σ^* i niech $x \in \Sigma^q$ będzie q -gramem. Jeśli $s_i s_{i+1} \dots s_{i+q-1} = x$ dla pewnego i , to x wystąpiło w s . Niech $\mathbf{v}(s, q)$ będzie wektorem o długości $|\Sigma|^q$, którego zmienne oznaczają liczbę wszystkich możliwych wystąpień q -gramów w s . Niech $s, t \in \Sigma^*$ oraz $q > 0$ będzie liczbą naturalną. Odległość q -gramową między napisami s i t definiuje się następująco [5]:

$$d_{qgram}(s, t, q) = \|\mathbf{v}(s, q) - \mathbf{v}(t, q)\|_1 = \sum_{i=1}^{|\Sigma|^q} |v_i(s, q) - v_i(t, q)| \quad (1.12)$$

Wzór 1.12 definiuje odległość q -gramową między napisami s i t jako odległość L_1 pomiędzy $\mathbf{v}(s, q)$ i $\mathbf{v}(t, q)$. Zauważmy, że wystarczy policzyć jedynie liczbę faktycznie występujących q -gramów w obu napisach, by obliczyć odległość q -gramową [6].

Przykład 1.10. Niech $\Sigma = \{a, j, p\}$. Wówczas odległość q -gramowa między napisami **papaja** i **japa** dla $q = 2$ wynosi: $d_{qgram}(\text{papaja}, \text{japa}, 2) = 2$. Wszystkie możliwe digramy występujące w napisach **papaja** i **japa** to **aj**, **ap**, **ja** i **pa**. Zatem $\mathbf{v}(\text{papaja}, 2) = (1, 1, 1, 2)$, a $\mathbf{v}(\text{japa}, 2) = (0, 1, 1, 1)$. Stąd $d_{qgram}(\text{papaja}, \text{japa}, 2) = \|(1, 1, 1, 2) - (0, 1, 1, 1)\|_1 = 2$.

Odległość q -gramowa również nie spełnia definicji metryki, bowiem $d_{qgram}(s, t) = 0$ nawet wówczas, gdy $s \neq t$. Przykładowo $d_{qgram}(\text{abaca}, \text{acaba}, 2) = 0$, mimo że **abaca** \neq **acaba**.

Maksymalna liczba wystąpień różnych q -gramów w napisie s wynosi $|s| - q + 1$. Stąd maksymalna odległość q -gramowa między napisami s i t wynosi $|s| + |t| - 2q + 2$, osiągnięta, gdy s i t nie mają wspólnych q -gramów [6].

Skoro zdefiniowana została odległość q -gramowa w języku wektorów, każda funkcja odległościowa (distance functions) w (całkowitej) przestrzeni wektorowej może zostać zastosowana. Przykładowo można zdefiniować *odległość cosinusową* między napisami s i t :

$$d_{cos}(s, t, q) = 1 - \frac{\mathbf{v}(s, q) \cdot \mathbf{v}(t, q)}{\|\mathbf{v}(s, q)\|_2 \|\mathbf{v}(t, q)\|_2} \quad (1.13)$$

Literatura

- [1] Leonid Boytsov. Indexing methods for approximate dictionary searching: Comparative analysis. *J. Exp. Algorithmics*, 16:1.1:1.1–1.1:1.91, Maj 2011.
- [2] R. W. Hamming. Error Detecting and Error Correcting Codes. *Bell System Technical Journal*, 29:147–160, 1950.
- [3] G. Navarro. A guided tour to approximate string matching. *ACM Computing Surveys*, 33(1):31–88, 2001.
- [4] Esko Ukkonen. Algorithms for approximate string matching. *Inf. Control*, 64(1-3):100–118, Marzec 1985.
- [5] Esko Ukkonen. Approximate string-matching with q-grams and maximal matches. *Theoretical Computer Science*, 92(1):191 – 211, 1992.
- [6] Mark P. J. van der Loo. The stringdist Package for Approximate String Matching. *The R Journal*, 6:111–122, 2014.
- [7] Robert A. Wagner, Michael J. Fischer. The string-to-string correction problem. *Journal of the ACM*, 21(1):168–173, January 1974.
- [8] Robert A. Wagner, Roy Lowrance. An extension of the string-to-string correction problem. *J. ACM*, 22(2):177–183, Kwiecień 1975.