

# Rozdział 1

## Odległości na przestrzeni ciągów znaków

### 1.1. Podstawowe definicje

**Definicja 1.1.** Niech  $\Sigma = \{\sigma_1, \dots, \sigma_k\}$  będzie skończonym uporządkowanym zbiorem o liczności  $|\Sigma|$ , zwanym alfabetem. Napisem nazywamy skończony ciąg znaków z  $\Sigma$ . Zbiór wszystkich napisów o długości  $n$  nad  $\Sigma$  jest oznaczony przez  $\Sigma^n$ , podczas gdy przez  $\Sigma^* = \bigcup_{n=1}^{\infty} \Sigma^n$  rozumiemy zbiór wszystkich napisów utworzonych ze znaków z  $\Sigma$  [4].

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  $\varepsilon$ . Przez  $|s|$ , dla każdego napisu  $s \in \Sigma^*$ , rozumiemy jego długość, czyli liczbę znaków w napisie. Ciąg napisów i/lub znaków oznacza ich złączenie, np.  $stu$  to napis powstały ze złączenia napisów  $s, t$  oraz  $u$ , natomiast  $abc$ , to napis powstały ze złączenia znaków  $a, b$  oraz  $c$ . 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  $i \leq j$ . Zakładamy również, że jeśli  $j < i$ , to  $s_{i:j} = \varepsilon$  [4, 17].

**Definicja 1.2.** Załóżmy, że napis  $s$  jest reprezentacją złączenia trzech, być może pustych, podnapisów  $w, x$  i  $y$ , tj.  $s = wxy$ . Wówczas podnapis  $w$  nazywamy przedrostkiem, natomiast podnapis  $y$  – przyrostkiem [4].

**Definicja 1.3.** Podnapis złożony z kolejnych, przylegających do siebie, 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 [4].

**Przykład 1.1.** Niech  $\Sigma$  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  $s_1 = \text{e}$ ,

$s_2 = 1$ ,  $s_3 = a$ . Podnapis 1:2 napisu  $s$  to  $s_{1:2} = e1$ . W napisie tym mamy do czynienia jedynie z  $q$ -gramami o  $q$  równym jeden, dwa oraz trzy, odpowiednio:  $e$ ,  $1$ ,  $a$ ;  $e1$ ,  $1a$  oraz  $e1a$ .

We wszystkich przykładach niniejszego rozdziału zakładamy, jeśli nie podano inaczej, że alfabet składa się z 32 liter polskiego alfabetu oraz liter  $q$ ,  $v$  i  $x$ .

## 1.2. Odległości na przestrzeni ciągów znaków

Określanie odległości między napisami jest ważną częścią przetwarzania danych, zwłaszcza w problemach takich jak dopasowanie statystyczne, wyszukiwanie tekstów, klasyfikacja tekstów czy sprawdzanie pisowni. Największa trudność polega na policzeniu podobieństwa między dwoma napisami w terminach metryk na przestrzeni ciągów znaków. W niniejszym podrozdziale zajmiemy się odległościami na napisach, tj. funkcjami  $d : \Sigma^* \times \Sigma^* \rightarrow [0, \infty)$ . W literaturze można znaleźć wiele różnych funkcji tego typu, które różnią się genezą powstania, podejściem do problemu oraz zastosowaniami. W pracy zajmiemy się jednak odległościami, która można podzielić na trzy grupy:

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

Pierwszy rodzaj odległości jest najczęściej używany w algorytmach zajmujących się optymalnym dopasowaniem napisów, dlatego też poświęcimy mu największą część niniejszego rozdziału. Aby wyliczyć odległość opartą na operacjach edycyjnych, trzeba określić liczbę fundamentalnych transformacji potrzebnych do przetworzenia jednego napisu w drugi. Mogą się w nich zawierać zamiany, usunięcia, wstawienia oraz transpozycje znaków. Odległości oparte na  $q$ -gramach pozwalają na określenie podobieństwa między dwoma napisami poprzez porównanie występowania  $q$ -elementowych ciągów znaków. Natomiast miary heurystyczne są rzadko stosowane, gdyż nie mają silnych matematycznych podstaw, ale zostały rozwinięte jako praktyczne narzędzie stosowane w konkretnych przypadkach.

### 1.2.1. Odległości oparte na operacjach edycyjnych

**Ścieżka edycyjna i bazowe operacje edycyjne.** *Odległość edycyjna*  $ED(s, t)$  między dwoma napisami  $s$  i  $t$  to minimalna liczba operacji edycyjnych potrzebna do przetworzenia  $s$  w  $t$  (i  $\infty$ , gdy taki ciąg nie istnieje) [10]. *Ścisłą odległością edycyjną* nazywamy minimalną liczbę nienakładających się operacji edycyjnych, które pozwalają przekształcić jeden napis w drugi, i które nie przekształcają dwa razy tego samego podnapisu [4].

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 przekształceniu napisu  $s$  w napis  $t$ , są oznaczane przez  $s \rightarrow t$ . Zbiór wszystkich bazowych operacji edycyjnych oznaczamy przez  $\mathbb{B}$  [4].

Bazowe operacje edycyjne są zazwyczaj ograniczone do:

- usunięcie znaku:  $1 \rightarrow \varepsilon$ , tj. usunięcie litery  $1$ , np.  $e1a \rightarrow ea$ ,

- 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.  $ela \rightarrow lea$ .

Często transpozycja znaków nie należy do zbioru operacji bazowych, jako że można ją zastąpić usunięciem i wstawieniem znaku. W niniejszej pracy jednak, operacja ta należy do zbioru operacji bazowych.

**Własność 1.1.** Zakładamy, że  $\mathbb{B}$  spełnia następujące własności [4]:

- 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 1.4.** 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 [4].

**Przykład 1.2.** Weźmy napisy **foczka** oraz **kozak**. Ścieżka edycyjna między nimi może mieć następującą postać:

$$\begin{array}{ccccccc} \text{foczka} & \xrightarrow{\text{trans. } z \ i \ k} & \text{fockza} & \xrightarrow{\text{trans. } c \ i \ k} & \text{fokcza} & \xrightarrow{\text{trans. } o \ i \ k} & \text{fkocza} \\ \xrightarrow{\text{wst. } k} & & & & & & \xrightarrow{\text{us. } f} & \text{kocza} & \xrightarrow{\text{us. } c} & \text{koza} \\ & & & & & & & & & \xrightarrow{\text{us. } c} & \text{kozak} \end{array}$$

Optymalna ścieżka natomiast ma następującą postać:

$$\text{foczka} \xrightarrow{\text{zm. } f \ na \ k} \text{koczka} \xrightarrow{\text{us. } c} \text{kozka} \xrightarrow{\text{trans. } k \ i \ a} \text{kozak}.$$

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. Formalne definicje powyższych funkcji znajdują się w dalszej części niniejszego rozdziału.

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)  $\delta(a \rightarrow b)$  [18]. Zazwyczaj koszt każdej operacji wynosi jeden, jednak można, na przykład, nadać transpozycji mniejszy koszt niż operacji wstawienia znaku. Dalej, można rozszerzyć funkcję kosztu  $\delta$  na ścieżkę edycyjną  $E = a_1 \rightarrow b_1, a_2 \rightarrow b_2, \dots, a_{|E|} \rightarrow b_{|E|}$  przez  $\delta(E) = \sum_{i=1}^{|E|} \delta(a_i \rightarrow b_i)$  [4]. Odtąd przez 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 ważne zamiany (substytucje) (pod)napisów, zamiast operacji edycyjnych wykonywanych na pojedynczych znakach [15]. Odległości zdefiniowane w ten sposób zazwyczaj są nazywane *rozszerzonymi* odległościami edycyjnymi. Przykładowo,  $\mathbb{B}$  może zawierać operację  $\mathbf{x} \rightarrow \mathbf{ks}$  o koszcie jednostkowym. Wówczas rozszerzona odległość pomiędzy napisami  $\mathbf{xero}$  i  $\mathbf{ksero}$  wynosi jeden, podczas gdy standardowa (zwykła, nierozszerzona) odległość wyniosłaby dwa [4].

**Definicja 1.5.** *Mając dany zbiór bazowych operacji edycyjnych  $\mathbb{B}$  oraz funkcję  $\delta$ , która nadaje koszt wszystkim bazowym operacjom edycyjnym z  $\mathbb{B}$ , uogólniona odległość edycyjna między napisami  $s$  i  $t$  jest zdefiniowana jako minimalny spośród kosztów wszystkich możliwych ścieżek edycyjnych, które przekształcają  $s$  w  $t$  [4].*

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 *ważonej* odległości edycyjnej.

**Własność 1.2.** *Zakładamy, że funkcja kosztu  $\delta(s \rightarrow t)$  ma następujące własności [4]:*

- $\delta(s \rightarrow t) \geq 0$  (koszt operacji jest liczbą nieujemny),
- $\delta(s \rightarrow t) = \delta(t \rightarrow s)$  (symetria),
- $\delta(s \rightarrow s) = 0$  i  $\delta(s \rightarrow t) = 0 \Rightarrow s = t$  (identyczność),
- $\forall \gamma > 0$  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. dobrze zdefiniowana odległość edycyjna z  $s$  do  $t$  [4],
- ogólna odległość edycyjna z definicji 1.5 jest metryką [18].

*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 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  $\gamma$  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  $\gamma$ . Zbiór  $A$  jest niepusty i składa się z operacji edycyjnych o dodatnim koszcie mniejszym niż  $\gamma$ . Zbiór operacji bazowych, których koszt jest ograniczony z góry przez  $\gamma$  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  $E_1E_2$ , które przekształca  $s$  w  $u$ . Z tego, że  $\delta(E_1E_2) = \delta(E_1) + \delta(E_2) = ED(s, t) + ED(t, u)$  oraz  $\delta(E_1E_2) \geq ED(s, u)$  (gdyż  $E_1E_2$  nie musi być optymalną ścieżką, przekształcającą  $s$  w  $u$ ) wynika, że  $ED(s, t) + ED(t, u) \geq ED(s, u)$ . ■

Odległość edycyjna jest metryką, nawet gdy funkcja kosztu  $\delta$  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:  $\{\mathbf{a}, \mathbf{b}, \mathbf{c}\}$ , gdzie symetria i brak subaddytywności funkcji  $\delta$  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 *ścistej 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 *ścisłą odległością edycyjną* [4]:

**Definicja 1.6.** Niech napisy  $s$  i  $t$  zostaną podzielone na tę samą liczbę, być może pustych, podnapisów:  $s = s_1s_2 \dots s_l$  i  $t = t_1t_2 \dots t_l$ , takich, że  $s_i \rightarrow t_i \in \mathbb{B}$ . Ścisłą ścieżką edycyjną nazywamy taką ścieżkę, że nie występują w niej następujące operacje:

- $s_i \rightarrow s_{i_j} \rightarrow t_i$  (modyfikacja tego samego podnapisu więcej niż raz),
- $s_i \rightarrow t_i, s_{i+1} \rightarrow t_{i+1}, t_i t_{i+1} \rightarrow t_k$  (nakładające się operacje).

**Lemat 1.4.** *Dowolna nieściśła odległość edycyjna ogranicza z dołu odpowiadającą jej ściśłą odległość edycyjną [4].*

**Lemat 1.5.** *Ściśła odległość Levenshteina o jednostkowym koszcie operacji bazowych jest równa nieściśłej odległości Levenshteina o jednostkowym koszcie operacji bazowych [4].*

Powyższe wynika natychmiast z obserwacji, że optymalna ścieżka edycyjna zawiera jednoznaczne usunięcia, wstawienia oraz zamiany, które nigdy nie modyfikują podnapisu 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 [4].*

*Dowód.* Ściśła 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 ściśłej odległości Damerau-Levenshteina [4]. ■

Rozważmy napisy  $ab$ ,  $ba$  oraz  $acb$ . Z jednej strony, najkrótsza nieściśła ścieżka edycyjna, która przekształca  $ba$  w  $acb$ , tj.  $(ba \rightarrow ab, \varepsilon \rightarrow c)$  zawiera dwie operacje: najpierws zamienia znaki  $a$  i  $b$ , a następnie wstawia  $c$  pomiędzy nie. Zauważmy, że wstawienie przekształca już transformowany napis. Jednakowoż, jeśli kolejne przekształcenia tego samego podnapisu są wykluczone, to najkrótsza ścieżka edycyjna, która przekształca  $ba$  w  $acb$ , składa się z trzech operacji edycyjnych, np.  $(b \rightarrow \varepsilon, \varepsilon \rightarrow c, \varepsilon \rightarrow b)$ . Stąd, nieściśła odległość edycyjna jest równa dwa, podczas gdy ściśła odległość wynosi trzy [4].

Ściśła odległość Damerau-Levenshteina nie spełnia nierówności trójkąta, gdyż

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

natomiast

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

zatem

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

Ponieważ ściśła i nieściśła odległość Damerau-Levenshteina są różnymi funkcjami, tę pierwszą nazywa się często *odległością optymalnego dopasowania napisów*. Od tego momentu w niniejszej pracy ściśłą odległość Damerau-Levenshteina nazywamy odległością optymalnego dopasowania napisów, natomiast nieściśłą odległość Damerau-Levenshteina nazywamy odległością Damerau-Levenshteina [17].

### Optymalne dopasowanie.

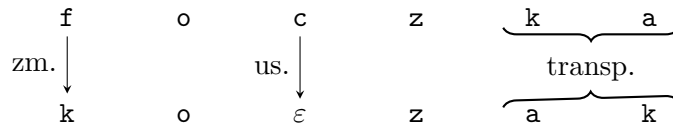
**Definicja 1.7.** *Niech napisy  $s$  i  $t$  zostaną podzielone na tę samą liczbę, być może pustych, podnapisów:  $s = s_1 s_2 \dots s_l$  i  $t = t_1 t_2 \dots t_l$ , takich, że  $s_i \rightarrow t_i \in \mathbb{B}$ . Co więcej, zakładamy, że  $s_i$  i  $t_j$  nie mogą być puste dla  $i = j$ . Mówimy, 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 dopasowany do podnapisu  $t_i$  [4].*

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 go 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 [4].

**Przykład 1.3.** Przykład optymalnego dopasowania pomiędzy słowami **foczka** i **kozak** prezentuje rys. 1.1. Odpowiadająca mu ścieżka edycyjna składa się z zamiany **f** → **k**, usunięcia **c** →  $\varepsilon$  oraz transpozycji **ka** → **ak**.



Rysunek 1.1: Przykład optymalnego dopasowania między napisami **foczka** i **kozak**.

Warto zauważyć, że istnieje różnowartościowe (1-1) mapowanie między zbiorem ścisłych ścieżek edycyjnych i zbiorem optymalnych dopasowań: każda ścisła ścieżka edycyjna o minimalnym koszcie reprezentuje dopasowanie o najmniejszym koszcie i odwrotnie. Stąd można zastąpić problem znalezienia optymalnej ścisłej odległości edycyjnej przez problem znalezienia optymalnego dopasowania, co też zastosujemy dalej [4].

**Obliczanie odległości edycyjnej.** Główną zasadą dynamicznego algorytmu, liczącego koszt optymalnego dopasowania, jest wyrażenie kosztu dopasowania pomiędzy napisami  $s$  i  $t$ , używając kosztu dopasowania ich przedrostków. Rozważmy przedrostek  $s_{1:i}$  o długości  $i$  i przedrostek  $t_{1:j}$  o długości  $j$ , odpowiednio napisów  $s$  i  $t$ . Załóżmy, że  $A = (s_1 s_2 \dots s_l, t_1 t_2 \dots t_l)$  jest optymalnym dopasowaniem między  $s_{1:i}$  i  $t_{1:j}$ , którego koszt oznaczamy przez  $C_{i,j}$  [4].

Używając równania 1.1 oraz definicji optymalnego dopasowania, łatwo pokazać, że  $C_{i,j}$  może zostać policzone przy użyciu następującej ogólnej rekurencji [15]:

$$\begin{aligned} C_{0,0} &= 0 \\ C_{i,j} &= \min\{\delta(s_{i':i} \rightarrow t_{j':j}) + C_{i'-1,j'-1} \mid s_{i':i} \rightarrow t_{j':j} \in \mathbb{B}\}. \end{aligned} \quad (1.2)$$

Można zauważyć, że:

- koszt dopasowania napisów  $s$  i  $t$  jest równy  $C_{|s|,|t|}$ ;
- wszystkie optymalne dopasowania mogą zostać wyznaczone przez odwracanie rekurencji 1.2 (przechodzenie od tyłu), tj. obliczanie najpierw  $C_{0,0}$ , następnie  $C_{1,1}$  itd.

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

$$\delta(s_{i':i} \rightarrow t_{j':j}) = [s_{i':i} \neq t_{j':j}] \quad (1.3)$$

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  $j' = j$ . Dalej, 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 rekurencyjna i można ją zapisać następująco:

**Definicja 1.8.** Odległością Hamminga nazywamy [5]:

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

Odległość Hamminga zlicza liczbę indeksów (p. rys. 1.2), 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.

**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  $d_{\text{hamming}}(\text{kozak}, \text{foczka}) = \infty$ , gdyż  $|\text{kozak}| \neq |\text{foczka}|$ .

|     |   |   |     |   |
|-----|---|---|-----|---|
|     | k | o | z   | a |
| zm. | ↓ |   | zm. | ↓ |
|     | f | o | k   | a |
|     | 1 | 2 | 3   | 4 |

Rysunek 1.2: Przykład dopasowania przy pomocy odległości Hamming między napisami **koza** i **foka**.

Rozważmy teraz odległość najdłuższego wspólnego podnapisu (ang. *longest common substring*), 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 ogólnej rekurencji 1.2, 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:

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

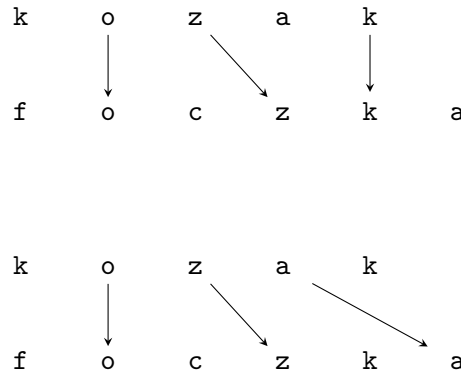
Odległość najdłuższego wspólnego podnapisu przyjmuje wartości ze zbioru  $\{0, |s| + |t|\}$ , przy czym maksimum jest osiąganę, gdy  $s$  i  $t$  nie mają ani jednego wspólnego znaku. Odległość tę oznaczamy przez  $d_{\text{lcs}}$ .



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

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 (rys. 1.3):



Rysunek 1.3: Przykład odległości najdłuższego wspólnego podnapisu między napisami **kozak** i **foczka**.

Jak widać, 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 ogólnej rekurencji 1.3:

- $i' = i - 1$  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:

$$C_{i,j} = \min \begin{cases} 0, & \text{gdy } i = j = 0 \\ C_{i-1,j} + 1, & \text{gdy } i > 0 \\ C_{i,j-1} + 1, & \text{gdy } j > 0 \\ C_{i-1,j-1} + [s_i \neq t_j], & \text{gdy } i, j > 0 \end{cases} \quad (1.4)$$

Odległość Levenshteina oznaczamy przez  $d_{\text{lv}}$ .

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

Powyższy przykład ilustruje dodatkową elastyczność w porównaniu do odległości najdłuższego wspólnego podnapisu, bowiem daje ona mniejszą wartość odległości między napisami, jako że w przypadku pierwszego znaku potrzebujemy jedynie zamiany, zamiast wstawienia i usunięcia [17]. Co więcej, ścieżka edycyjna między tymi słowami może być inna i zawierać usunięcie i wstawienie zamiast dwóch ostatnich zamian znaków.

Przypomnijmy, że mówimy o ważonej 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 **kozak** i **foczka** wynosi:  $d_{lv}(\text{kozak}, \text{foczka}) = 1.9$ , bo  $\text{kozak} \xrightarrow[0.3]{zm. k na f} \text{fozak} \xrightarrow[1]{wst. c} \text{foczak} \xrightarrow[0.3]{zm. a na k} \text{foczkk} \xrightarrow[0.3]{zm. k na a} \text{foczka}$ .

Ważona 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 transformacji napisu  $t$  w napis  $s$  [17]. Dobrze obrazuje to następujący przykład:

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

$$d_{lv}(\text{kozak}, \text{foczka}) = 1.9, \quad (1.5)$$

gdyż

$$\text{kozak} \xrightarrow[0.3]{zm. k na f} \text{fozak} \xrightarrow[1]{wst. c} \text{foczak} \xrightarrow[0.3]{zm. a na k} \text{foczkk} \xrightarrow[0.3]{zm. k na a} \text{foczka},$$

natomiast

$$d_{lv}(\text{foczka}, \text{kozak}) = 1, \quad (1.6)$$

gdyż

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

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

$$d_{lv}(\text{kozak}, \text{foczka}) = 1,$$

gdyż

$$\text{kozak} \xrightarrow[0.3]{zm. k na f} \text{fozak} \xrightarrow[0.1]{wst. c} \text{foczak} \xrightarrow[0.3]{zm. a na k} \text{foczkk} \xrightarrow[0.3]{zm. k na a} \text{foczka},$$

czyli analogicznie, jak w przypadku 1.6. Natomiast

$$d_{lv}(\text{foczka}, \text{kozak}) = 1.9,$$

bo

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

czyli analogicznie, jak w przypadku 1.5.

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. Spójrzmy na następujące bezpośrednie uogólnienie rekurencji 1.4, dodające transpozycję do zbioru bazowych operacji edycyjnych [4]:

$$C_{i,j} = \min \begin{cases} 0, & \text{gdy } i = j = 0 \\ C_{i-1,j} + 1, & \text{gdy } i > 0 \\ C_{i,j-1} + 1, & \text{gdy } j > 0 \\ C_{i-1,j-1} + [s_i \neq t_j], & \text{gdy } i, j > 0 \\ C_{i-2,j-2} + 1, & \text{gdy } s_i = t_{j-1}, s_{i-1} = t_j \text{ oraz } i, j > 1 \end{cases} \quad (1.7)$$

Rekurencja 1.7 liczy odległość optymalnego dopasowania napisów (ozn.  $d_{\text{osa}}$ ), czyli ściśłą odległość Damerau-Levenshteina, która nie zawsze jest równa odległości Damerau-Levenshteina. Dla przykładu, odległość między napisami **ba** i **acb** wyliczona przy pomocy rekurencji 1.7 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 [19]. W ich definicji zamiana zostaje zastąpiona przez minimalizację po możliwych transpozycjach między danym znakiem a wszystkimi nie przetransformowanymi znakami, przy czym koszt transpozycji wzrasta wraz z odległością między transponowanymi znakami [17]. Innymi słowy, do  $\mathbb{B}$  należą wstawienia, usunięcia, zamiany oraz operacje  $axb \rightarrow bya$  o koszcie równym  $|x| + |y| + 1$  [4]. Mając tak zdefiniowane  $\mathbb{B}$  ogólna rekurencja 1.2 dla odległości Damerau-Levenshteina przedstawia się następująco:

$$C_{i,j} = \min \begin{cases} 0, & \text{gdy } i = j = 0 \\ C_{i-1,j} + 1, & \text{gdy } i > 0 \\ C_{i,j-1} + 1, & \text{gdy } j > 0 \\ C_{i-1,j-1} + [s_i \neq t_j], & \text{gdy } i, j > 0 \\ \min_{\substack{0 < i' < i, 0 < j' < j \\ s_i = t_{j'}, s_{i'} = t_j}} C_{i'-1,j'-1} + (i - i') + (j - j') - 1 & \end{cases} \quad (1.8)$$

Co więcej, Lowrance i Wagner wykazali, że wewnętrzne minimum w rekurencji 1.8 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**  $\xrightarrow[1]{zm. k \text{ na } f}$  **fozak**  $\xrightarrow[1]{wst. c}$  **foczak**  $\xrightarrow[1]{transp. a \text{ i } k}$  **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 [17]:

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

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 na $q$ -gramach

Przypomnijmy, że  $q$ -gramem nazywamy napis składający się z  $q$  kolejnych (przylegających) znaków.  $q$ -gramy związane z napisem  $s$  są otrzymywane przez przesuwanie przez napis  $s$  „okna” o szerokości  $q$  znaków i zapisaniu występujących  $q$ -gramów. Przykładowo digramy napisu **ela** to **e****l** i **l****a**. Oczywiście taka procedura nie ma sensu, gdy  $q > |s|$  lub gdy  $q = 0$ . Z tego powodu definiujemy następujące przypadki brzegowe dla wszystkich odległości  $d_q(s, t)$  opartych na  $q$ -gramach:

$$\begin{aligned} d_q(s, t) &= \infty, \text{ gdy } q > \min\{|s|, |t|\}, \\ d_0(s, t) &= \infty, \text{ gdy } |s| + |t| > 0, \\ d_0(\varepsilon, \varepsilon) &= 0. \end{aligned}$$

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  $\mathcal{Q}(s, q)$  oznaczmy zbiór unikalnych  $q$ -gramów występujących w napisie  $s$ , to możemy zdefiniować odległość Jaccarda [17]:

**Definicja 1.9.** Niech  $\mathcal{Q}(s, q)$  oznacza zbiór unikalnych  $q$ -gramów występujących w napisie  $s$ . Wówczas odległość Jaccarda,  $d_{\text{jac}}$ , między napisami  $s$  i  $t$  definiuje się jako

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

gdzie  $|\cdot|$  oznacza licznosc zbioru.

Odległość Jaccarda przyjmuje wartości z przedziału  $[0, 1]$ , gdzie 0 odpowiada pełnemu pokryciu 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 **papaja** i **japa** dla  $q = 2$  wynosi:  $d_{\text{jac}}(\text{papaja}, \text{japa}, 2) = 0.25$ , bo  $\mathcal{Q}(\text{papaja}, 2) = \{\text{pa}, \text{ap}, \text{aj}, \text{ja}\}$ , a  $\mathcal{Q}(\text{japa}, 2) = \{\text{ja}, \text{ap}, \text{pa}\}$ , więc odległość wynosi  $1 - \frac{3}{4} = 0.25$ .

Inną odległością opartą na  $q$ -gramach jest odległość  $q$ -gramowa. Otrzymuje się ją przez wylistowanie  $q$ -gramów występujących w obu napisach i policzenie  $q$ -gramów, które nie są wspólne dla obu napisów [17]. Formalnie można to zapisać następująco:

**Definicja 1.10.** Niech  $s = s_1 s_2 \dots s_n$  będzie napisem z  $\Sigma^*$  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 wystąpień wszystkich możliwych  $q$ -gramów z  $\Sigma^q$  w  $s$ . Niech  $s, t \in \Sigma^*$  oraz  $q > 0$  będzie liczbą naturalną. Odległość  $q$ -gramową między napisami  $s$  i  $t$  definiuje się następująco [16]:

$$d_{\text{qgram}}(s, t, q) = \|\mathbf{v}(s, q) - \mathbf{v}(t, q)\|_1 = \sum_{i=1}^{|\Sigma|^q} |v_i(s, q) - v_i(t, q)|. \quad (1.9)$$

Wzór 1.9 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, zamiast sprawdzać wystąpienie wszystkich możliwych  $q$ -gramów z  $\Sigma^q$  w napisach  $s$  i  $t$ , wystarczy policzyć jedynie liczbę faktycznie występujących  $q$ -gramów w obu napisach, by obliczyć odległość  $q$ -gramową [17].

**Przykład 1.10.** Niech  $\Sigma = \{a, j, p\}$ . Wówczas odległość  $q$ -gramowa między napisami **papaja** i **japa** dla  $q = 2$  wynosi:  $d_{\text{qgram}}(\text{papaja}, \text{japa}, 2) = 2$ . 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  $d_{\text{qgram}}(\text{papaja}, \text{japa}, 2) = \|(1, 1, 1, 2) - (0, 1, 1, 1)\|_1 = 2$ .

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 [17].

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

$$d_{\text{cos}}(s, t, q) = 1 - \frac{\mathbf{v}(s, q) \cdot \mathbf{v}(t, q)}{\|\mathbf{v}(s, q)\|_2 \|\mathbf{v}(t, q)\|_2}, \quad (1.10)$$

gdzie  $\|\cdot\|_2$  oznacza zwykłą normę Euklidesową. Odległość cosinusowa wynosi zero, gdy  $s = t$  oraz jeden, gdy  $s$  i  $t$  nie mają wspólnych  $q$ -gramów. Odległość ta powinna być interpretowana jako kąt pomiędzy  $\mathbf{v}(s, q)$  i  $\mathbf{v}(t, q)$ , jako że drugie wyrażenie równania 1.10 przedstawia cosinus kąta między dwoma wektorami.

**Przykład 1.11.** Niech  $\Sigma = \{a, j, p\}$ . Wówczas odległość cosinusowa między napisami **papaja** i **japa** dla  $q = 2$  wynosi:  $d_{\text{cos}}(\text{papaja}, \text{japa}, 2) \approx 0.127$ , bo  $\mathbf{v}(\text{papaja}, 2) = (1, 1, 1, 2)$ , a  $\mathbf{v}(\text{japa}, 2) = (0, 1, 1, 1)$  (p. przykład 1.10), więc  $d_{\text{cos}}(\text{papaja}, \text{japa}, 2) = 1 - \frac{4}{\sqrt{3} \cdot \sqrt{7}} \approx 0.127$ .

Wszystkie trzy odległości oparte na  $q$ -gramach są nieujemne i symetryczne. Odległości Jaccarda i  $q$ -gramowa spełniają również nierówność trójkąta (dowód dla tej pierwszej poniżej), w odróżnieniu od odległości cosinusowej. Żadna z powyższych miar nie spełnia warunku identyczności, ponieważ zarówno  $Q(s, q)$ , jak i  $\mathbf{v}(s, q)$  jest funkcją wiele-do-jednego. Jako przykład, zauważmy, że  $Q(\text{abaca}, 2) = Q(\text{acaba}, 2)$  oraz  $\mathbf{v}(\text{abaca}, 2) = \mathbf{v}(\text{acaba}, 2)$ , więc  $d_{\text{jac}}(\text{abaca}, \text{acaba}, 2) = d_{\text{qgram}}(\text{abaca}, \text{acaba}, 2) = d_{\text{cos}}(\text{abaca}, \text{acaba}, 2) = 0$ . Innymi słowy, odległość oparta na  $q$ -gramach równa zero, nie gwarantuje, że  $s = t$ . Inne własności  $\mathbf{v}(s, q)$  można znaleźć w [16].

Udowodnimy teraz, że odległość Jaccarda spełnia nierówność trójkąta.

**Lemat 1.7.** Niech odległość Jaccarda będzie dana następującym wzorem:

$$d(A, B) = 1 - \frac{|A \cap B|}{|A \cup B|}.$$

Wówczas  $d$  spełnia nierówność trójkąta:

$$\forall A, B, C \quad d(A, C) \geq d(A, B) + d(B, C).$$

*Dowód.* Przypuśćmy, że nierówność nie jest spełniona, tj. istnieją zbiory  $A, B$  oraz  $C$ , takie, że  $d(A, C) > d(A, B) + d(B, C)$ . Wówczas

$$1 - \frac{|A \cap C|}{|A \cup C|} > 1 - \frac{|A \cap B|}{|A \cup B|} + 1 - \frac{|B \cap C|}{|B \cup C|}$$

lub równoważnie

$$\frac{|A \cap B|}{|A \cup B|} + \frac{|B \cap C|}{|B \cup C|} - \frac{|A \cap C|}{|A \cup C|} > 1. \quad (1.11)$$

Ponieważ postulujemy, że nierówność trójkąta nie jest spełniona, chcemy aby lewa strona ostatniej nierówności była jak najmniejsza. Stąd wystarczy rozważyć przypadki, gdy  $A \subseteq B$  lub  $B \subseteq A$ , gdyż w przeciwnym przypadku, dla  $B' = A \cap B$  mamy

$$\frac{|A \cap B'|}{|A \cup B'|} = \frac{|A \cap B|}{|A \cup (A \cap B)|} \geq \frac{|A \cap B|}{|A \cup B|}.$$

Zastępujemy zatem  $B$  przez  $B'$  i dostajemy:

$$\frac{|A \cap B'|}{|A \cup B'|} + \frac{|B \cap C|}{|B \cup C|} - \frac{|A \cap C|}{|A \cup C|} \geq \frac{|A \cap B|}{|A \cup B|} + \frac{|B \cap C|}{|B \cup C|} - \frac{|A \cap C|}{|A \cup C|}.$$

Analogicznie możemy rozważyć jedynie przypadki gdy  $C \subseteq B$  lub  $B \subseteq C$ . Przeanalizujmy teraz cztery przypadki.

Przypadek 1.  $A \subseteq B$  oraz  $C \subseteq B$ . Wówczas

$$\begin{aligned} \frac{|A \cap B|}{|A \cup B|} + \frac{|B \cap C|}{|B \cup C|} - \frac{|A \cap C|}{|A \cup C|} &= \frac{|A|}{|B|} + \frac{|C|}{|B|} - \frac{|A \cap C|}{|A \cup C|} = \frac{|A| + |C|}{|B|} - \frac{|A \cap C|}{|A \cup C|} = \\ &= \frac{|A \cup C| + |A \cap C|}{|B|} - \frac{|A \cap C|}{|A \cup C|} \leq \frac{|A \cup C|}{|B|} + \frac{|A \cap C|}{|B|} - \frac{|A \cap C|}{|B|} = \frac{|A \cup C|}{|B|} \leq 1 \end{aligned}$$

Przypadek 2.  $A \subseteq B$  oraz  $B \subseteq C$ . Wówczas

$$\frac{|A \cap B|}{|A \cup B|} + \frac{|B \cap C|}{|B \cup C|} - \frac{|A \cap C|}{|A \cup C|} = \frac{|A|}{|B|} + \frac{|B|}{|C|} - \frac{|A|}{|C|} \leq \frac{|A|}{|C|} + \frac{|B|}{|C|} - \frac{|A|}{|C|} = \frac{|B|}{|C|} \leq 1$$

Przypadek 3.  $C \subseteq B$  oraz  $B \subseteq A$ . Wówczas

$$\frac{|A \cap B|}{|A \cup B|} + \frac{|B \cap C|}{|B \cup C|} - \frac{|A \cap C|}{|A \cup C|} = \frac{|B|}{|A|} + \frac{|C|}{|B|} - \frac{|C|}{|A|} \leq \frac{|B|}{|A|} + \frac{|C|}{|A|} - \frac{|C|}{|A|} = \frac{|B|}{|A|} \leq 1$$

Przypadek 4.  $B \subseteq A$  oraz  $B \subseteq C$ . Wówczas

$$\frac{|A \cap B|}{|A \cup B|} + \frac{|B \cap C|}{|B \cup C|} - \frac{|A \cap C|}{|A \cup C|} = \frac{|B|}{|A|} + \frac{|B|}{|C|} - \frac{|A \cap C|}{|A \cup C|} \leq \frac{|A \cap C|}{|A|} + \frac{|B|}{|C|} - \frac{|A \cap C|}{|A|} = \frac{|B|}{|C|} \leq 1$$

Wszystkie cztery powyższe przypadki są w sprzeczności ze stwierdzeniem, że wyrażenie 1.11 jest ostro większe od 1. Stąd założenie jest nieprawdziwe, więc nierówność trójkąta jest spełniona [20]. ■

### 1.2.3. Miary heurystyczne

Odległość Jaro została stworzona w amerykańskim Bureau of the Census (rządowa agencja, która jest odpowiedzialna m.in. za spis ludności Stanów Zjednoczonych) w celu połączenia rekordów, które były wpisane w niewłaściwe pola formularza oraz zlikwidowaniu literówek. Pierwszy publiczny opis tej odległości pojawił się w instrukcji obsługi [8], co może wyjaśniać dlaczego nie jest rozpowszechniona w literaturze informatycznej. Jednak odległość ta została skutecznie zastosowana w statystycznych problemach dopasowania w przypadku dość krótkich napisów, głównie imion, nazwisk oraz danych adresowych [17].

Rozumowanie stojące za odległością Jaro jest następujące: błędny znak oraz transpozycje znaków są spowodowane błędem przy wpisywaniu, ale mało prawdopodobne jest znalezienie błędnego znaku w miejscu odległym od zamierzonego, żeby mogło to być spowodowane błędem przy wpisywaniu. Stąd odległość Jaro mierzy liczbę wspólnych znaków w dwóch napisach, które nie są zbyt odległe od siebie i dodaje karę za dopasowanie znaków, które są stransponowane. Formalna definicja wygląda następująco [17]:

**Definicja 1.11.** Niech  $s$  i  $t$  będą napisami z  $\Sigma^*$ . Niech  $m$  oznacza liczbę wspólnych znaków z  $s$  i  $t$ , przy czym zakładając, że  $s_i = t_j$ , to znak ten jest wspólny dla obu napisów, jeśli:

$$|i - j| < \left\lfloor \frac{\max\{|s|, |t|\}}{2} \right\rfloor$$

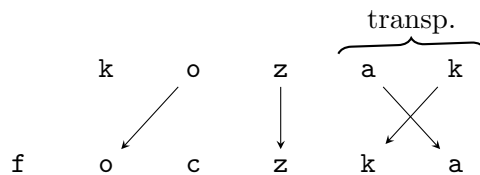
i każdy znak z  $s$  może być wspólny ze znakiem z  $t$  tylko raz. W końcu, jeśli  $s'$  i  $t'$  są podnapisami utworzonymi z  $s$  i  $t$  poprzez usunięcie znaków, które nie są wspólne dla obu napisów, to  $T$  jest liczbą transpozycji potrzebnych do otrzymania  $t'$  z  $s'$ . Transpozycje znaków nieprzylegających są dozwolone.

Wówczas odległość Jaro definiuje się jako:

$$d_{\text{jaro}}(s, t) = \begin{cases} 0, & \text{gdy } s = t = \varepsilon \\ 1, & \text{gdy } m = 0 \text{ i } |s| + |t| > 0 \\ 1 - \frac{1}{3} \left( \frac{m}{|s|} + \frac{m}{|t|} + \frac{m-T}{m} \right) & \text{w przeciwnym przypadku} \end{cases} \quad (1.12)$$

Odległość Jaro przyjmuje wartości z przedziału  $[0, 1]$ , gdzie zero oznacza, że  $s = t$ , natomiast jeden wskazuje na kompletną odmiennność napisów z  $m = T = 0$ .

**Przykład 1.12.** Odległość Jaro między napisami **kozak** i **foczka** wynosi:  $d_{\text{jaro}}(\text{kozak}, \text{foczka}) \approx 0.261$ , bo liczba wspólnych znaków wynosi  $m = 4$ , a liczba potrzebnych transpozycji wynosi  $T = 1$  (p. rys. 1.4), co daje odległość równą  $d_{\text{jaro}}(\text{kozak}, \text{foczka}) = 1 - \frac{1}{3} \left( \frac{3}{5} + \frac{4}{6} + \frac{3}{4} \right) = \frac{47}{180} \approx 0.261$ .



Rysunek 1.4: Przykład odległości Jaro między napisami **kozak** i **foczka**.

Winkler rozszerzył odległość Jaro przez włączenie dodatkowej kary za błędny znak wśród pierwszych czterech znaków napisu [17]:

**Definicja 1.12.** Niech  $s$  i  $t$  będą napisami z  $\Sigma^*$ ,  $\ell(s, t)$  oznacza długość najdłuższego wspólnego przedrostka, mającego maksymalnie cztery znaki i niech  $p$  będzie liczbą z przedziału  $[0, \frac{1}{4}]$ . Wówczas odległość Jaro-Winklera dana jest wzorem [21]:

$$d_{jw}(s, t, p) = d_{jaro}(s, t)[1 - p\ell(s, t)] \quad (1.13)$$

Czynnik  $p$  określa jak bardzo różnice w czterech pierwszych znakach w obu napisach wpływają na odległość między nimi. Zmienna  $p$  jest liczbą z przedziału  $[0, \frac{1}{4}]$ , by mieć pewność, że odległość Jaro-Winklera miała wartości w przedziale  $[0, 1]$  ( $0 \leq d_{jw}(s, t) \leq 1$ ). Jeśli  $p = 0$ , to odległość ta redukuje się do odległości Jaro i wszystkie znaki wnoszą taki sam wkład do funkcji odległości. Jeśli  $p = \frac{1}{4}$ , to odległość Jaro-Winklera jest równa zero nawet wówczas gdy tylko cztery pierwsze znaki w obu napisach pokrywają się. Powód jest taki, że podobno ludzie są mniej skłonni do popełniania błędów w czterech pierwszych znakach lub też są one lepiej zauważalne, więc różnice w pierwszych czterech znakach wskazują na większe prawdopodobieństwo, że dwa napisy są rzeczywiście różne [17]. Winkler [21] używał w swoich badaniach  $p = 0.1$  i zauważył lepsze rezultaty niż dla  $p = 0$ .

**Przykład 1.13.** Odległość Jaro-Winklera między napisami **faktura** i **faktyczny** dla  $p = 0$ ,  $p = 0.1$  oraz  $p = 0.25$  wynosi odpowiednio:

$$\begin{aligned} d_{jw}(\text{faktura}, \text{faktyczny}, p = 0.00) &\approx 0.328 = d_{jaro}(\text{faktura}, \text{faktyczny}) \\ d_{jw}(\text{faktura}, \text{faktyczny}, p = 0.10) &\approx 0.197 \\ d_{jw}(\text{faktura}, \text{faktyczny}, p = 0.25) &= 0 \end{aligned}$$

Łatwo zauważyć z równań 1.12 i 1.13, że odległości Jaro i Jaro-Winklera, dla  $p \neq \frac{a}{4}$ , są nieujemne, symetryczne i spełniają warunek identycznościowy. Nierówność trójkąta w obu przypadkach nie jest jednak spełniona. Rozważmy następujący przykład:  $s = \mathbf{ab}, t = \mathbf{cb}, u = \mathbf{cd}$ . Jako że napisy  $s$  i  $u$  nie mają wspólnych znaków, to odległość Jaro między nimi wynosi  $d_{jaro}(s, u) = 1$ , podczas gdy  $d_{jaro}(s, t) = d_{jaro}(t, u) = \frac{1}{3}$ , więc w tym przypadku  $d_{jaro}(s, u)$  jest większe od  $d_{jaro}(s, t) + d_{jaro}(t, u)$ . Z tego łatwo zauważyć, że odległość Jaro-Winklera nie spełnia nierówności trójkąta dla tego samego przykładu dla  $p \in [0, \frac{1}{4}]$  [17].

W niniejszym rozdziale przedstawiono odległości określone na przestrzeni ciągów znaków. Mając do wyboru wachlarz różnych funkcji nasuwa się pytanie której użyć. Ostateczna decyzja zależy od konkretnego przypadku, jednak istnieją pewne ogólne reguły. Wybór pomiędzy odległościami opartymi na operacjach edycyjnych i  $q$ -gramach z jednej strony, a miarami heurystycznymi z drugiej zależy w dużej mierze od długości napisów – te ostatnie są dedykowane krótszym napisom takim jak np. dane osobowe. W odróżnieniu od odległości opartych na operacjach edycyjnych i miarach heurystycznych, odległości oparte na  $q$ -gramach można łatwo policzyć dla bardzo długich tekstów, jako że liczba  $q$ -gramów możliwych do utworzenia z języka naturalnego (dla niezbyt małego  $q$ , tj.  $q \geq 3$ ) jest z reguły o wiele mniejsza niż liczba  $q$ -gramów, którą można otrzymać z całego alfabetu. Wybór spośród odległości opartych na operacjach edycyjnych zależy przede wszystkim od dokładności jaką chce się otrzymać. Przykładowo do wyszukiwania haseł w słowniku, gdzie różnice między dobranymi napisami są niewielkie, odległości pozwalające na więcej operacji edycyjnych (tak jak np. odległość Damerau-Levenshteina) mogą dać lepsze rezultaty. Odległości Jaro i Jaro-Winklera zostały



skonstruowane do krótkich, napisanych przez człowieka, napisów, więc ich zakres zastosowania powinien być jasny.



# Literatura

- [1] Wikipedia - wikipedia, wolna encyklopedia. <http://pl.wikipedia.org/wiki/Wikipedia>. Dostęp: 2015-12-01.
- [2] Léon Bottou. Stochastic gradient tricks. Grégoire Montavon, Genevieve B. Orr, Klaus-Robert Müller, redaktorzy, *Neural Networks, Tricks of the Trade, Reloaded*, Lecture Notes in Computer Science (LNCS 7700), strony 430–445. Springer, 2012.
- [3] Léon Bottou, Yoshua Bengio. Convergence properties of the k-means algorithms. *Advances in Neural Information Processing Systems 7*, strony 585–592. MIT Press, 1995.
- [4] Leonid Boytsov. Indexing methods for approximate dictionary searching: Comparative analysis. *J. Exp. Algorithmics*, 16:1.1:1.1–1.1:1.91, 2011.
- [5] R. W. Hamming. Error Detecting and Error Correcting Codes. *Bell System Technical Journal*, 29:147–160, 1950.
- [6] Trevor J. Hastie, Robert John Tibshirani, Jerome H. Friedman. *The elements of statistical learning : data mining, inference, and prediction*. Springer series in statistics. Springer, New York, 2009.
- [7] Lawrence Hubert, Phipps Arabie. Comparing partitions. *Journal of Classification*, 2(1):193–218, 1985.
- [8] Matthew A. Jaro. *UNIMATCH: A record linkage system: User manual*. United States Bureau of the Census, 1978.
- [9] J. Koronacki, J. Ćwik. *Statystyczne systemy uczące się*. Wydawnictwa Naukowo-Techniczne, 2005.
- [10] G. Navarro. A guided tour to approximate string matching. *ACM Computing Surveys*, 33(1):31–88, 2001.
- [11] William M. Rand. Objective criteria for the evaluation of clustering methods. *Journal of the American Statistical Association*, 66(336):846–850, 1971.
- [12] Andrew Rosenberg, Julia Hirschberg. V-measure: A conditional entropy-based external cluster evaluation measure. *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, strony 410–420, 2007.
- [13] Peter Rousseeuw. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *J. Comput. Appl. Math.*, 20(1):53–65, 1987.

- 
- [14] D. Sculley. Web-scale k-means clustering. *Proceedings of the 19th International Conference on World Wide Web*, WWW '10, strony 1177–1178, New York, NY, USA, 2010. ACM.
  - [15] Esko Ukkonen. Algorithms for approximate string matching. *Inf. Control*, 64(1-3):100–118, 1985.
  - [16] Esko Ukkonen. Approximate string-matching with q-grams and maximal matches. *Theoretical Computer Science*, 92(1):191 – 211, 1992.
  - [17] Mark P. J. van der Loo. The stringdist Package for Approximate String Matching. *The R Journal*, 6:111–122, 2014.
  - [18] Robert A. Wagner, Michael J. Fischer. The string-to-string correction problem. *Journal of the ACM*, 21(1):168–173, 1974.
  - [19] Robert A. Wagner, Roy Lowrance. An extension of the string-to-string correction problem. *J. ACM*, 22(2):177–183, 1975.
  - [20] Anna Wilbik, James M. Keller. A distance metric for a space of linguistic summaries. *Fuzzy Sets and Systems*, 208:79–94, 2012.
  - [21] William E. Winkler. String comparator metrics and enhanced decision rules in the fellegi-sunter model of record linkage. *Proceedings of the Section on Survey Research*, strony 354–359, 1990.
  - [22] Xindong Wu, Vipin Kumar, J. Ross Quinlan, Joydeep Ghosh, Qiang Yang, Hiroshi Motoda, Geoffrey J. McLachlan, Angus Ng, Bing Liu, Philip S. Yu, Zhi-Hua Zhou, Michael Steinbach, David J. Hand, Dan Steinberg. Top 10 algorithms in data mining. *Knowl. Inf. Syst.*, 14(1):1–37, 2007.