

Błędna pisownia może być powodowana przez różne czynniki:

- błędy w przekazywaniu sygnałów do mięśni odpowiedzialnych za ruch palców (literówki, np. *litreówki*);
- nieznanomość pisowni słów i zasad ortografii (wymowa błędnej formy jest taka sama lub podobna do wymowy formy poprawnej, np. *Stary człowiek i może* jako tytuł znanego utworu literackiego a nie reklama leku);
- nieznanomość morfologii (np. *upartosc* zamiast *upór*).

Metody poprawiania pisowni powinny uwzględniać model błędów użytkownika.

Odległość Levenshteina dwóch łańcuchów znaków  $x$  i  $y$  to minimalna liczba prostych operacji edycyjnych, które przekształcają łańcuch  $x$  w łańcuch  $y$  (lub odwrotnie).

**Proste operacje edycyjne to:**

- wstawienie znaku, np. kota  $\rightarrow$  kwota
- usunięcie znaku, np. kwota  $\rightarrow$  kota
- zamiana znaku, np. kota  $\rightarrow$  koza
- zamiana miejscami dwóch sąsiadujących znaków (tzw. czeski błąd), np. pisk  $\rightarrow$  psik

# Odległość Levenshteina

$$\begin{aligned}ed(X_{1\dots i+1}, Y_{1\dots j+1}) &= ed(X_{1\dots i}, Y_{1\dots j}) && \text{jeśli } x_{i+1} = y_{j+1} \\ & && \text{(ostatnie znaki} \\ & && \text{takie same)} \\ &= 1 + \min\{ed(X_{1\dots i-1}, Y_{1\dots j-1}), && \text{jeśli } x_i = y_{j+1} \\ & \quad ed(X_{1\dots i+1}, Y_{1\dots j}), && \text{i } x_{i+1} = y_j \\ & \quad ed(X_{1\dots i}, Y_{1\dots j+1})\} && \text{(dwa ostatnie znaki} \\ & && \text{zamienione miejscami)} \\ &= 1 + \min\{ed(X_{1\dots i}, Y_{1\dots j}), && \text{w pozost. przypadkach} \\ & \quad ed(X_{1\dots i+1}, Y_{1\dots j}), \\ & \quad ed(X_{1\dots i}, Y_{1\dots j+1})\} \\ ed(X_{1\dots 0}, Y_{1\dots j}) &= j && 0 \leq j \leq n \\ ed(X_{1\dots i}, Y_{1\dots 0}) &= i && 0 \leq i \leq m \\ ed(X_{1\dots -1}, Y_{1\dots j}) &= ed(X_{1\dots i}, Y_{1\dots -1}) = \max(m, n) && \text{(Definicje brzegowe)}\end{aligned}$$

# Odległość Levenshteina

przestawienie: (*kula*, *kual*)

$$\begin{aligned} ed(X_{1\dots i+1}, Y_{1\dots j+1}) &= ed(X_{1\dots i}, Y_{1\dots j}) && \text{jeśli } x_{i+1} = y_{j+1} \\ & && \text{(ostatnie znaki} \\ & && \text{takie same)} \\ &= 1 + \min\{ed(X_{1\dots i-1}, Y_{1\dots j-1}), && \text{jeśli } x_i = y_{j+1} \\ & \quad ed(X_{1\dots i+1}, Y_{1\dots j}), && \text{i } x_{i+1} = y_j \\ & \quad ed(X_{1\dots i}, Y_{1\dots j+1})\} && \text{(dwa ostatnie znaki} \\ & && \text{zamienione miejscami)} \\ &= 1 + \min\{ed(X_{1\dots i}, Y_{1\dots j}), && \text{w pozost. przypadkach} \\ & \quad ed(X_{1\dots i+1}, Y_{1\dots j}), \\ & \quad ed(X_{1\dots i}, Y_{1\dots j+1})\} \\ ed(X_{1\dots 0}, Y_{1\dots j}) &= j && 0 \leq j \leq n \\ ed(X_{1\dots i}, Y_{1\dots 0}) &= i && 0 \leq i \leq m \\ ed(X_{1\dots -1}, Y_{1\dots j}) &= ed(X_{1\dots i}, Y_{1\dots -1}) = \max(m, n) && \text{(Definicje brzegowe)} \end{aligned}$$

# Odległość Levenshteina

usunięcie: (bab,baba)

$$\begin{aligned}ed(X_{1\dots i+1}, Y_{1\dots j+1}) &= ed(X_{1\dots i}, Y_{1\dots j}) \\ &= 1 + \min\{ed(X_{1\dots i-1}, Y_{1\dots j-1}), \\ &\quad ed(X_{1\dots i+1}, Y_{1\dots j}), \\ &\quad ed(X_{1\dots i}, Y_{1\dots j+1})\} \\ &= 1 + \min\{ed(X_{1\dots i}, Y_{1\dots j}), \\ &\quad ed(X_{1\dots i+1}, Y_{1\dots j}), \\ &\quad ed(X_{1\dots i}, Y_{1\dots j+1})\}\end{aligned}$$

jeśli  $x_{i+1} = y_{j+1}$   
(ostatnie znaki  
takie same)  
jeśli  $x_i = y_{j+1}$   
i  $x_{i+1} = y_j$   
(dwa ostatnie znaki  
zamienione miejscami)  
w pozost. przypadkach

$$\begin{aligned}ed(X_{1\dots 0}, Y_{1\dots j}) &= j & 0 \leq j \leq n \\ ed(X_{1\dots i}, Y_{1\dots 0}) &= i & 0 \leq i \leq m \\ ed(X_{1\dots -1}, Y_{1\dots j}) &= ed(X_{1\dots i}, Y_{1\dots -1}) = \max(m, n) & \text{(Definicje brzegowe)}\end{aligned}$$

# Odległość Levenshteina

$$\begin{aligned}ed(X_{1\dots i+1}, Y_{1\dots j+1}) &= ed(X_{1\dots i}, Y_{1\dots j}) \text{ wstawienie: } (babab, baba)_{-1} \\ &= 1 + \min\{ed(X_{1\dots i-1}, Y_{1\dots j-1}), ed(X_{1\dots i+1}, Y_{1\dots j}), ed(X_{1\dots i}, Y_{1\dots j+1})\} \\ &= 1 + \min\{ed(X_{1\dots i}, Y_{1\dots j}), ed(X_{1\dots i+1}, Y_{1\dots j}), ed(X_{1\dots i}, Y_{1\dots j+1})\}\end{aligned}$$

(ostatnie znaki takie same)  
jeśli  $x_i = y_{j+1}$   
i  $x_{i+1} = y_j$   
(dwa ostatnie znaki zamienione miejscami)  
w pozost. przypadkach

$$\begin{aligned}ed(X_{1\dots 0}, Y_{1\dots j}) &= j & 0 \leq j \leq n \\ ed(X_{1\dots i}, Y_{1\dots 0}) &= i & 0 \leq i \leq m \\ ed(X_{1\dots -1}, Y_{1\dots j}) &= ed(X_{1\dots i}, Y_{1\dots -1}) = \max(m, n) & \text{(Definicje brzegowe)}\end{aligned}$$

# Odległość Levenshteina

$$\begin{aligned}ed(X_{1\dots i+1}, Y_{1\dots j+1}) &= ed(X_{1\dots i}, Y_{1\dots j}) && \text{jeśli } x_{i+1} = y_{j+1} \\ &&& \text{(ostatnie znaki)} \\ &= 1 + \min\{ed(X_{1\dots i-1}, Y_{1\dots j-1}), && \text{zamiana: (kod,kot)} \\ &\quad ed(X_{1\dots i+1}, Y_{1\dots j}), && \text{jeśli } x_i = y_{j+1} \\ &\quad ed(X_{1\dots i}, Y_{1\dots j+1})\} && \text{i } x_{i+1} = y_j \\ &&& \text{(dwa ostatnie znaki} \\ &&& \text{zamienione miejscami)} \\ &= 1 + \min\{ed(X_{1\dots i}, Y_{1\dots j}), && \text{w pozost. przypadkach} \\ &\quad ed(X_{1\dots i+1}, Y_{1\dots j}), \\ &\quad ed(X_{1\dots i}, Y_{1\dots j+1})\} \\ ed(X_{1\dots 0}, Y_{1\dots j}) &= j && 0 \leq j \leq n \\ ed(X_{1\dots i}, Y_{1\dots 0}) &= i && 0 \leq i \leq m \\ ed(X_{1\dots -1}, Y_{1\dots j}) &= ed(X_{1\dots i}, Y_{1\dots -1}) = \max(m, n) && \text{(Definicje brzegowe)}\end{aligned}$$

# Odległość Levenshteina

$$\begin{aligned}ed(X_{1\dots i+1}, Y_{1\dots j+1}) &= ed(X_{1\dots i}, Y_{1\dots j}) && \text{jeśli } x_{i+1} = y_{j+1} \\ &&& \text{(ostatnie znaki} \\ &&& \text{takie same)} \\ &= 1 + \min\{ed(X_{1\dots i-1}, Y_{1\dots j-1}), && \text{usunięcie: (kot, kota)} \\ &\quad ed(X_{1\dots i+1}, Y_{1\dots j}), && \text{↑ } x_{i+1} = y_j \\ &\quad ed(X_{1\dots i}, Y_{1\dots j+1})\} && \text{(dwa ostatnie znaki} \\ &&& \text{zamienione miejscami)} \\ &= 1 + \min\{ed(X_{1\dots i}, Y_{1\dots j}), && \text{w pozost. przypadkach} \\ &\quad ed(X_{1\dots i+1}, Y_{1\dots j}), \\ &\quad ed(X_{1\dots i}, Y_{1\dots j+1})\} \\ ed(X_{1\dots 0}, Y_{1\dots j}) &= j && 0 \leq j \leq n \\ ed(X_{1\dots i}, Y_{1\dots 0}) &= i && 0 \leq i \leq m \\ ed(X_{1\dots -1}, Y_{1\dots j}) &= ed(X_{1\dots i}, Y_{1\dots -1}) = \max(m, n) && \text{(Definicje brzegowe)}\end{aligned}$$



# Odległość Levenshteina

$$\begin{aligned}ed(X_{1\dots i+1}, Y_{1\dots j+1}) &= ed(X_{1\dots i}, Y_{1\dots j}) && \text{jeśli } x_{i+1} = y_{j+1} \\ &&& \text{(ostatnie znaki} \\ &&& \text{takie same)} \\ &= 1 + \min\{ed(X_{1\dots i-1}, Y_{1\dots j-1}), && \text{jeśli } x_i \neq y_{j+1} \\ &\quad ed(X_{1\dots i+1}, Y_{1\dots j}), && \text{wstawienie: (kota, kot)} \\ &\quad ed(X_{1\dots i}, Y_{1\dots j+1})\} && \text{(dwa ostatnie znaki} \\ &&& \text{zamienione miejscami)} \\ &= 1 + \min\{ed(X_{1\dots i}, Y_{1\dots j}), && \text{w pozost. przypadkach} \\ &\quad ed(X_{1\dots i+1}, Y_{1\dots j}), \\ &\quad ed(X_{1\dots i}, Y_{1\dots j+1})\} \\ ed(X_{1\dots 0}, Y_{1\dots j}) &= j && 0 \leq j \leq n \\ ed(X_{1\dots i}, Y_{1\dots 0}) &= i && 0 \leq i \leq m \\ ed(X_{1\dots -1}, Y_{1\dots j}) &= ed(X_{1\dots i}, Y_{1\dots -1}) = \max(m, n) && \text{(Definicje brzegowe)}\end{aligned}$$

# Odległość edycyjna – przykład

◇	g	u	p	c	h	i	s	
◇	0	1	2	3	4	5	6	7
g	1	0	1	2	3	4	5	6
ł	2	1	1	2	3	4	5	6
u	3	2	1	2	3	4	5	6
p	4	3	2	1	2	3	4	5
s	5	5	3	2	2	3	4	4
i	6	5	4	3	3	3	3	4

# Odległość edycyjna – przykład

$$\begin{aligned} ed(\mathbf{g}, \mathbf{g}) &= ed(\epsilon, \epsilon) \\ &= 0 \end{aligned}$$

◇	g	u	p	c	h	i	s	
◇	0	1	2	3	4	5	6	7
g	1	0	1	2	3	4	5	6
ł	2	1	1	2	3	4	5	6
u	3	2	1	2	3	4	5	6
p	4	3	2	1	2	3	4	5
s	5	5	3	2	2	3	4	4
i	6	5	4	3	3	3	3	4

gg  
gg

# Odległość edycyjna – przykład

$$\begin{aligned}
 ed(g, g\acute{t}) &= 1 + \min\{ed(\epsilon, g), ed(g, g), ed(\epsilon, g\acute{t})\} \\
 &= 1 + \min\{1, 0, 2\} = 1
 \end{aligned}$$

◇	g	u	p	c	h	i	s	
◇	0	1	2	3	4	5	6	7
g	1	0	1	2	3	4	5	6
ł	2	1	1	2	3	4	5	6
u	3	2	1	2	3	4	5	6
p	4	3	2	1	2	3	4	5
s	5	5	3	2	2	3	4	4
i	6	5	4	3	3	3	3	4

g  
g   ł

# Odległość edycyjna – przykład

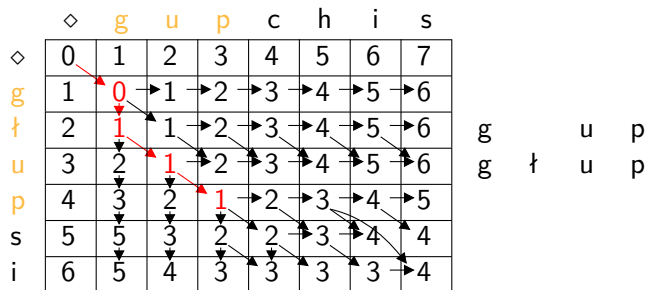
$$\begin{aligned} ed(\mathbf{gu}, \mathbf{g\acute{t}u}) &= ed(\mathbf{g}, \mathbf{g\acute{t}}) \\ &= 1 \end{aligned}$$

◇	g	u	p	c	h	i	s	
◇	0	1	2	3	4	5	6	7
g	1	0	1	2	3	4	5	6
ł	2	1	1	2	3	4	5	6
u	3	2	1	2	3	4	5	6
p	4	3	2	1	2	3	4	5
s	5	5	3	2	2	3	4	4
i	6	5	4	3	3	3	3	4

g                    u  
g                    ł                    u

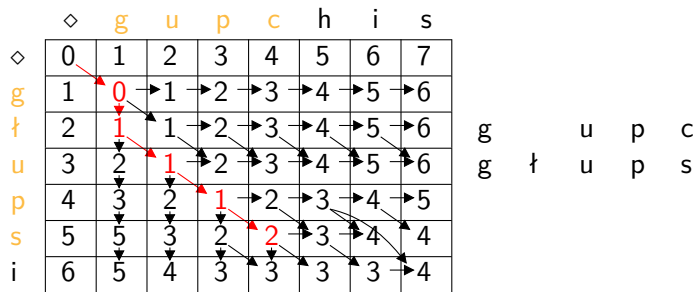
# Odległość edycyjna – przykład

$$\begin{aligned} ed(\text{gup}, \text{głup}) &= ed(\text{gu}, \text{głu}) \\ &= 1 \end{aligned}$$



# Odległość edycyjna – przykład

$$\begin{aligned}
 ed(\text{gupc}, \text{głups}) &= 1 + \min\{ed(\text{gup}, \text{głup}), ed(\text{gupc}, \text{głup}), \\
 &\quad ed(\text{gup}, \text{głups})\} \\
 &= 1 + \min\{1, 2, 2\} = 2
 \end{aligned}$$



# Odległość edycyjna – przykład

$$\begin{aligned}
 ed(\text{gupch}, \text{głups}) &= 1 + \min\{ed(\text{gupc}, \text{głup}), ed(\text{gupch}, \text{głup}), \\
 &\quad ed(\text{gupc}, \text{głups})\} \\
 &= 1 + \min\{2, 3, 2\} = 3
 \end{aligned}$$

◇	g	u	p	c	h	i	s	
◇	0	1	2	3	4	5	6	7
g	1	0	1	2	3	4	5	6
ł	2	1	1	2	3	4	5	6
u	3	2	1	2	3	4	5	6
p	4	3	2	1	2	3	4	5
s	5	5	3	2	2	3	4	4
i	6	5	4	3	3	3	3	4

g            u   p   c   h  
g   ł   u   p   s



# Odległość edycyjna – przykład

$$\begin{aligned} ed(\text{gupchi}, \text{głupsi}) &= ed(\text{gupch}, \text{głups}) \\ &= 3 \end{aligned}$$

◇	g	u	p	c	h	i	s
◇	0	1	2	3	4	5	6
g	1	0	1	2	3	4	5
ł	2	1	1	2	3	4	5
u	3	2	1	2	3	4	5
p	4	3	2	1	2	3	4
s	5	5	3	2	2	3	4
i	6	5	4	3	3	3	3

g            u   p   c   h   i  
g   ł   u   p   s        i

# Odległość edycyjna – przykład

$$\begin{aligned}ed(\text{gupchis}, \text{głupsi}) &= 1 + \min\{ed(\text{gupch}, \text{głup}), \\ &\quad ed(\text{gupchi}, \text{głupsi}), ed(\text{gupchis}, \text{głups})\} \\ &= 1 + \min\{3, 3, 4\} = 4\end{aligned}$$

◇	g	u	p	c	h	i	s	
◇	0	1	2	3	4	5	6	7
g	1	0	1	2	3	4	5	6
ł	2	1	1	2	3	4	5	6
u	3	2	1	2	3	4	5	6
p	4	3	2	1	2	3	4	5
s	5	5	3	2	2	3	4	4
i	6	5	4	3	3	3	3	4

g            u   p   c   h   i   s  
g   ł   u   p   s            i

# Odległość edycyjna – przykład

◇	g	u	p	c	h	i	s	
◇	0	1	2	3	4	5	6	7
g	1	0	1	2	3	4	5	6
ł	2	1	1	2	3	4	5	6
u	3	2	1	2	3	4	5	6
p	4	3	2	1	2	3	4	5
s	5	5	3	2	2	3	4	4
i	6	5	4	3	3	3	3	4

g            u   p   c   h   i   s  
g   ł   u   p            s   i

# Odległość edycyjna – przykład

◇	g	u	p	c	h	i	s	
◇	0	1	2	3	4	5	6	7
g	1	0	1	2	3	4	5	6
ł	2	1	1	2	3	4	5	6
u	3	2	1	2	3	4	5	6
p	4	3	2	1	2	3	4	5
s	5	5	3	2	2	3	4	4
i	6	5	4	3	3	3	3	4

g            u   p   c   h   i   s  
g   ł   u   p                    s   i

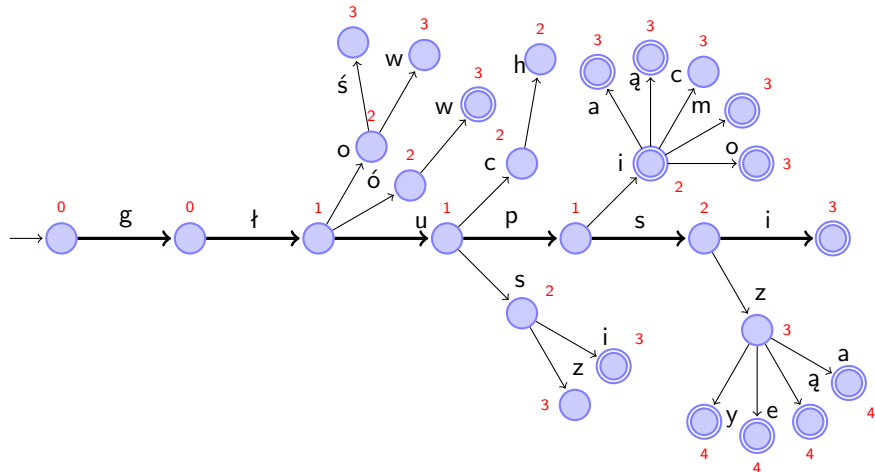
Porównywanie w całości wszystkich słów w słowniku z formą niepoprawną  $X$  (o długości  $m$ ) byłoby zbyt kosztowne. Wystarczy sprawdzić tylko część potencjalnego zamiennika  $Y$  (o długości  $n$ ) z częścią niepoprawnej formy  $X$  i zdecydować, czy sprawdzać następne litery formy:

$$\text{cuted}(X_{1\dots m}, Y_{1\dots n}) = \min_{l \leq i \leq u} \text{ed}(X_{1\dots i}, Y_{1\dots n}),$$

$$l = \min(1, n - t), u = \max(m, n + t)$$

$t$  jest maksymalną dopuszczalną odległością, operacje  $\min$  i  $\max$  zapobiegają wyjściu indeksów poza granice słowa.

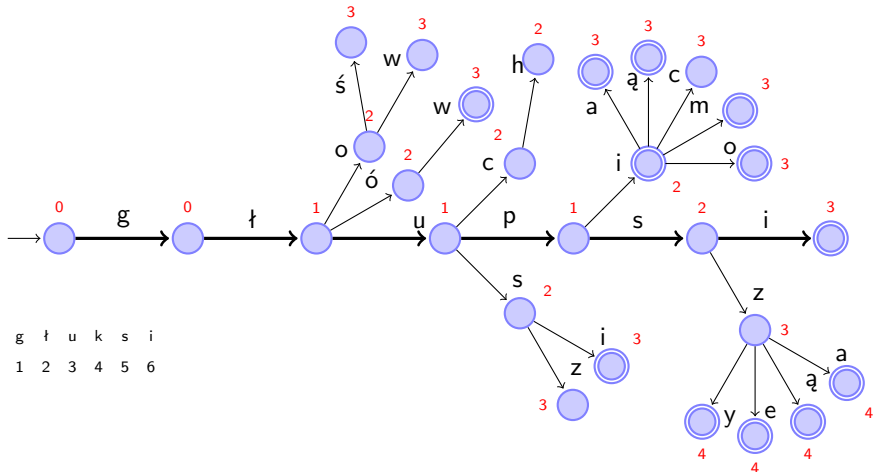
# Edycyjna odległość odcięcia



Dla  $t = 4$ ,

$$\text{cuted}(gupchis, głuch) = \min\{\text{ed}(g, głuch), \text{ed}(gu, głuch), \text{ed}(gup, głuch), \text{ed}(gupc, głuch), \text{ed}(gupch, głuch), \text{ed}(gupchi, głuch), \text{ed}(gupchis, głuch)\} = 2$$

# Edycyjna odległość odcięcia



Dla  $t = 4$ ,

$$\text{cuted}(gupchis, głuch) = \min\{\text{ed}(g, głuch), \text{ed}(gu, głuch), \text{ed}(gup, głuch), \text{ed}(gupc, głuch), \text{ed}(gupch, głuch), \text{ed}(gupchi, głuch), \text{ed}(gupchis, głuch)\} = 2$$