

```
In [1]: import spacy
from random import random
from bisect import bisect
```

## Lab 04 - Edit distance

### Implementacja algorytmu wyznaczania odległości edycyjnej wraz z odtworzonymi krokami

```
In [2]: def delta(a, b):
        return 0 if a == b else 1

def edit_distance(x, y, delta, show_answer = False):
    m = len(x)
    n = len(y)
    table = [[0]*(n+1) for _ in range(m+1)]

    for i in range(m+1):
        table[i][0] = i

    for i in range(n+1):
        table[0][i] = i

    for i in range(1, m+1):
        for j in range(1, n+1):
            table[i][j] = min(table[i-1][j] + 1, table[i][j-1]+1, table[i-1][j-1] + delta(x[i-1], y[j-1]))

    if show_answer:
        return f'x + ' -> ' + y, '-'*(len(x) + len(y) + 4) + get_solution(table, x, y, m, n)
    return table[m][n]

def get_solution(table, s1, s2, i, j):
    if table[i][j] == 0:
        return []

    if i == 0:
        step = s1 + ' -> ' + '|' + s2[j-1] + '|' + s1
        s1 = s2[j-1] + s1
        return [step] + get_solution(table, s1, s2, i, j-1)
    elif i == 0:
        step = s1 + '->' + s1[:i-1] + ' _'
        s1 = s1[:i-1]
        return [step] + get_solution(table, s1, s2, i-1, j)
    else:
        if table[i][j] == table[i-1][j] + 1:
            step = s1 + ' -> ' + s1[:i-1] + '_' + s1[i:]
            s1 = s1[:i-1] + s1[i:]
            return [step] + get_solution(table, s1, s2, i-1, j)
        elif table[i][j] == table[i][j-1] + 1:
            step = s1 + '|' + s1[i:] + '|' + s2[j-1] + '|' + s1[i:]
            s1 = s1[:i] + s2[j-1] + s1[i:]
            return [step] + get_solution(table, s1, s2, i, j-1)
        elif table[i][j] == table[i-1][j-1] + 1:
            return get_solution(table, s1, s2, i-1, j-1)

        step = s1 + ' -> ' + s1[:i-1] + '*' + s2[j-1] + '*' + s1[i:]
        s1 = s1[:i-1] + s2[j-1] + s1[i:]
        return [step] + get_solution(table, s1, s2, i-1, j-1)
```

```
In [3]: for line in edit_distance('los', 'kloc', delta, True):
        print(line)
```

```
los -> kloc
-----
los -> lo*c*
loc -> |k|loc
```

```
In [4]: for line in edit_distance('Łódź', 'Łodz', delta, True):
        print(line)
```

```
Łódź -> Łodz
-----
Łódź -> Łód*z*
Łódź -> |Ł|odz
Łódź -> *L*odz
```

```
In [5]: for line in edit_distance('kwintesencja', 'quintessence', delta, True):
        print(line)
```

```
kwintesencja -> quintessence
-----
kwintesencja -> kwintesencj
kwintesencj -> kwintesencj*e
kwintessence -> kwintes|s|ence
kwintessence -> *q*uintessence
kwintessence -> *q*uintessence
```

```
In [6]: for line in edit_distance('ATGAATCTTACCGCCTCG', 'ATGAGGCTCTGGCCCCCTG', delta, True):
        print(line)
```

```
ATGAATCTTACCGCCTCG -> ATGAGGCTCTGGCCCCCTG
-----
ATGAATCTTACCGCCTCG -> ATGAATCTTACCGCCTT_G
ATGAATCTTACCGCCTCG -> ATGAATCTTACCG_CCTG
ATGAATCTTACCGCCTCG -> ATGAATCTTA|G|CCCCCTG
ATGAATCTTACCGCCTCG -> ATGAATCTT*G*GCCCTG
ATGAATCTTGGCCCCCTG -> ATGAATCT|C|TGGCCCCCTG
ATGAATCTCTGGCCCCCTG -> ATGAA*G*CTCTGGCCCCCTG
ATGAAGCTCTGGCCCCCTG -> ATGA*G*GCTCTGGCCCCCTG
```

### Implementacja algorytmu znajdowania maksymalnego wspólnego podciągu

```
In [7]: def longest_common_sequence(x, y):
        ranges = [len(y)]
        y_letters = list(y)
        for i in range(len(x)):
            positions = [j for j, l in enumerate(y_letters) if l == x[i]]
            positions.reverse()
            for p in positions:
                k = bisect(ranges, p)
                if k == bisect(ranges, p - 1):
                    if k < len(ranges) - 1:
                        ranges[k] = p
                    else:
                        ranges[k:k] = [p]
            return len(ranges) - 1

In [8]: print('Najdłuższy wspólny podciąg słów los i kloc ma długość: {}'.format(longest_common_sequence('los', 'kloc')))
```

Najdłuższy wspólny podciąg słów los i kloc ma długość: 2

```
In [9]: print('Najdłuższy wspólny podciąg słów Łódź i Łodz ma długość: {}'.format(longest_common_sequence('Łódź', 'Łodz')))
```

Najdłuższy wspólny podciąg słów Łódź i Łodz ma długość: 1

```
In [10]: print('Najdłuższy wspólny podciąg słów kwintesencja i quintessence ma długość: {}'.format(longest_common_sequence('kwintesencja', 'quintessence')))
```

Najdłuższy wspólny podciąg słów kwintesencja i quintessence ma długość: 8

```
In [11]: print('Najdłuższy wspólny podciąg słów ATGAATCTTACCGCCTCG i ATGAGGCTCTGGCCCCCTG ma długość: {}'.format(longest_common_sequence('ATGAATCTTACCGCCTCG', 'ATGAGGCTCTGGCCCCCTG')))
```

Najdłuższy wspólny podciąg słów ATGAATCTTACCGCCTCG i ATGAGGCTCTGGCCCCCTG ma długość: 13

### Podział tekstu na tokeny i usunięcie 3% tokenów

```
In [12]: def text_with_removed_tokens(doc : spacy.tokens.Doc, tokens_to_delete : float):
        reduced_text = []

        for token in doc:
            if random() > tokens_to_delete :
                reduced_text += [token.text_with_ws]
            else:
                cnt = token.text_with_ws.count('\n')
                for _ in range(cnt):
                    Reduced_text += ['\n']

        return reduced_text

In [13]: with open('romeo-i-julia-700.txt', mode='r', encoding='utf-8') as file:
        data = file.read()

        nlp = spacy.load("pl_core_news_sm")
        doc = nlp(data)
        reduced_text1 = text_with_removed_tokens(doc, 0.03)
        reduced_text2 = text_with_removed_tokens(doc, 0.03)

        with open('variant1.txt', mode='w', encoding='utf-8') as save_file:
            for token_text in reduced_text1:
                save_file.write(token_text)

        with open('variant2.txt', mode='w', encoding='utf-8') as save_file:
            for token_text in reduced_text2:
                save_file.write(token_text)

        lcs = longest_common_sequence(reduced_text1, reduced_text2)

        print('Ilość tokenów na które został podzielony tekst: {}'.format(len(doc)))
        print('Ilość tokenów usuniętych z pierwszej wersji podziału: {}'.format(len(doc) - len(reduced_text1)))
        print('Ilość tokenów usuniętych z drugiej wersji podziału: {}'.format(len(doc) - len(reduced_text2)))
        print('Najdłuższy wspólny podciąg tokenów: {}'.format(lcs))
```

Ilość tokenów na które został podzielony tekst: 2699  
Ilość tokenów usuniętych z pierwszej wersji podziału: 51  
Ilość tokenów usuniętych z drugiej wersji podziału: 53  
Najdłuższy wspólny podciąg tokenów: 2564

```
In [14]: def show_diff(x, y):
        L = [[0 for _ in range(len(y) + 1)] for _ in range(len(x) + 1)]
        for i in range(1, len(x) + 1):
            for j in range(1, len(y) + 1):
                if x[i - 1] == y[j - 1]:
                    L[i][j] = L[i - 1][j - 1] + 1
                else:
                    L[i][j] = max(L[i - 1][j], L[i][j - 1])

        lines = []
        i = len(x) - 1
        j = len(y) - 1

        while i >= 0 and j >= 0:
            if x[i] == y[j]:
                i, j = i-1, j-1
            elif L[i][j-1] >= L[i-1][j]:
                lines.append(f"> |{j}| {y[j]}")
                j -= 1
            elif L[i][j-1] < L[i-1][j]:
                lines.append(f"< |{i}| {x[i]}")
                i -= 1

            while j >= 0:
                lines.append(f"> |{j}| {y[j]}")
                j -= 1

            while i >= 0:
                lines.append(f"< |{i}| {x[i]}")
                i -= 1
            lines.reverse()
            for line in lines:
                print(line)

In [15]: with open('variant1.txt', mode='r', encoding='utf-8') as file:
        lines1 = file.read().splitlines()
        with open('variant2.txt', mode='r', encoding='utf-8') as file:
            lines2 = file.read().splitlines()

        show_diff(lines1, lines2)
```

```
< [2] Romeo Julia
> [2] Romeo i Julia
> [11] * PARYS – młody Weronieńczyk szlachetnego rodu, krewny księcia
> [11] * PARYS – młody Weronieńczyk szlachetnego rodukrewny księcia
> [12] * MONTEKI, KAPULET – naczelnicy dwóch domów nieprzyjecznych sobie
> [13] * STARZEC – stryjeczny brat Kapuleta
> [14] * ROMEO – syn Montekiego
> [14] * MONTEKI, KAPULET – naczelnicy dwóch domów nieprzyjecznych sobie
> [13] * STARZEC – stryjeczny brat Kapuleta
> [14] * ROMEO – Montekiego
> [18] * LAURENTY – ojciec franciszkanin
> [18] * LAURENTY – ojciec franciszkanin
> [19] * JAN – brat z tegoż zromadzenia
> [20] * BALTAZAR – słuźacy Rómea
> [19] * JAN – brat z tegoż zromadzenia
> [20] * BALTAZAR – Rómea
> [25] * PAŻ PARYSA
> [25] * PAŻ PARYSA
> [37] Rzec odbywa się przez większa część sztuki w Weronie, przez część piatego aktu Mantui.
> [45] Rzec odbywa się przez większa część sztuki w Weronie, przez część piatego w Mantui.
> [45] Dwa rody, zacne jednakó 1 sławne
> [45] Dwa rody, zacne jednakó 1 sławne –
> [50] Z łon tych dwu wrogów wzięło źwie,
> [50] Z łon tych dwu wrogów wzięło bowiem źwie,
> [53] Śmierć ich stiumia rodzicielskie boje.
> [53] Śmierć ich rodzicielskie boje.
> [54]
> [55] Tej ich miłości przebieg zbył bolesny
> [54]
> [55] Tej ich miłości przebieg zbył
> [56] I jak się ojców nienawid nie zmienia,
> [57] Aż ja zakończy dzieł zgon przedwczesny,
> [56] I jak się ojców nienawid nie zmienia,
> [57] Aż ja zakończy dzieł przedwczesny,
> [61] Jest w nim co złego, mi usuniem błędy..
> [61] Jest w nim co , my usuniem błędy..
> [72] Plac publiczny. Wchodza Samson i Grzegorz uzbrojeni w tarcze i miecze. /
> [72] / Plac publiczny, Samson i Grzegorz uzbrojeni w tarcze i miecze. /
> [78]
> [80] GRZEGORZ
> [88]
> [90] GRZEGORZ
> [97] Mam zwyczaj drapać zaraz, jak mię kto rozrucha.
> [97] zwyczaj drapać zaraz, jak mię kto rozrucha.
> [102] Tak, ale nie zwykłeś się dać rozruchać.
> [102] Tak, ale nie zwykłeś się dać rozruchać.
> [112] Rozruchać się tyle znaczy co ruszyć się z / być walecznym jest to stać nieporuszenie: pojmuję , że skut
> [112] Rozruchać się tyle znaczy co ruszyć się z miejsca; być walecznym jest to stać nieporuszenie: pojmuję ,
> [112] Rozruchać się tyle znaczy co ruszyć się z miejsca; być walecznym jest to stać nieporuszenie: pojmuję ,
> [117] Te psy z domu Montekich usuniętych z drugiej wersji podziału: 53
> [117] Te psy z domu Montekich rozruchać mię mogą tylko stania na miejscu. Będę jak dla każdego mężczyzny i ka
> [117] Te psy z domu Montekich rozruchać mię mogą tylko do stania na miejscu. Będę jak mur dla każdego mężczyzn
> [127] Prawda, dlatego to kobiety, jako najsłabsze, tulą się zawsze do muru. Ja odtrące od muru ludzi Monteki
> [127] Prawda, dlatego to kobiety, jako najsłabsze, tulą się zawsze do muru. Ja też odtrące od muru ludzi Mont
> [152] Tym lepiej, że się liczysz do zwierząt; bo gdybyś się liczył do ryb, to byłbyś pewnie sztokfiszem. Weź
> [152] Tym lepiej, że się liczysz do zwierząt; bo gdybyś się liczył do , to byłbyś pewnie sztokfiszem. Weź no
> [152] Tym lepiej, że się liczysz do zwierząt; bo gdybyś się liczył do , to byłbyś pewnie sztokfiszem. Weź no
> [159] Mój giwer już dobyty: zaczeć ich, ja stanę z tyłu.
> [159] Mój już dobyty: zaczeć ich, ja stanę z tyłu.
> [174] Ja bym się miał bać z twojej przyczyny!
> [174] Ja bym się miał bać z twojej przyczyny!
> [189] Nie chca, ale jak śmia. Ja im gębę wykrzywie; hanba im, jeśli to ścierpia.
> [189] Nie chca, ale jak śmia. Ja im gębę ; hanba im, jeśli to ścierpia.
> [200]
> [202] ABRAHAM
> [204] Czy na nas się skrzywiłeś, panie?
> [204] Czy na nas się skrzywiłeś, panie?
> [211] Będziemy-ż mieli prawo sobą, jak powiem: tak jest?
> [211] Będziemy-ż mieli prawo za sobą, jak powiem: tak jest?
> [216] .
> [216] Nie.
> [221] Nie, mości panie; nie skrzywiłem na was, tylko skrzywiłem się tak sobie.
> [221] Nie, mości panie; nie skrzywiłem się na was, tylko skrzywiłem się tak sobie.
> [228] Waść szukasz?
> [228] Waść szukasz?
> [233] Zaczepki szukasz?
> [233] Zaczepki nie.
> [233] Zaczepkinie.
> [238] Jeżeli jej szukasz, jestem na waszine usługi. Mój pan tak dobry jak i wasz.
> [238] Jeżeli jej szukasz, to jestem na waszine usługi. Mój pan tak dobry jak i wasz.
> [257] Powiedz: lepszy. Oto nadchodzi z krewnych mego pana.
> [257] Powiedz: lepszy. Oto nadchodzi jeden z krewnych mego pana.
> [272] Dobądźcie mieczów, macie serca. Grzegorz, pamiętaj o swoim pchnięciu.
> [272] Dobądźcie mieczów, jeśli macie serca. Grzegorz, pamiętaj o swoim pchnięciu.
> [277] Odstąpcie, głupcy: schowajcie miecze . Sami nie wiecie, robiecie.
> [277] Odstąpcie, głupcy: schowajcie miecze do pochew. nie wiecie, co robiecie.
> [281] Wchodzi Tybalt. /
> [281] Wchodzi Tybalt. /
> [290] BENWOLIO
> [290]
> [299] Tego wyrazu, tak jak nienawidzę
> [299] wyrazu, tak jak nienawidzę
> [300] Szatana, wszystkich Montekich i ciebie.
> [301] Broń się, nikczemny tchórz.
> [302]
> [303] / Walczą. Nadchodzi kilku przyjaciół obu partii mieszają się do zwady; wkrótce potem wchodzą mieszczani
> [303] e z pałkami.
> [300] Szatanawszystkich Montekich i ciebie.
> [301] Broń sięnikczemny tchórz.
> [302]
> [303] Walczą. Nadchodzi kilku przyjaciół obu partii i mieszają się zwady; wkrótce potem wchodzą mieszczanie z
> [303] pałkami. /
> [308] Hola! berdyszów! palek! Dalej po !
> [308] Hola! berdyszów! palek! Dalej po nich!
> [317] Mój mieczhej!
> [317] Mój miecz! hej!
> [351] Zapamiętali niesforni poddani,
> [351] Zapamiętali niesforni poddani,
> [356] Pod kara tortur wypuścić natychmiast
> [356] kara tortur wypuścić natychmiast
> [365] Porzucić swoje wygodne przybory
> [365] Porzucić wygodne przybory
> [368] Niecheci wasze przeczynaJeżeli
> [368] Niecheci wasze przeczynaJeżeli
> [375] Dalsza ma wola oznajmiona będzie.
> [375] Dalsza ma wola oznajmiona będzie
> [376] Jeszcze raz wzywam tu obecnych
> [377] Pod karą śmierci, się rozeszli.
> [378]
> [379] / Książę z orszakiem wychodzi. Podobnie Kapulet, Pani Kapulet, Tybalt, obywateli i słuźy/
> [376] Jeszcze raz wzywam wszystkich tu obecnych
> [377] Pod karą śmierci, aby się rozeszli.
> [378]
> [379] / Książę z orszakiem wychodzi. Podobnie Kapulet, Pani Kapulet, Tybalt, obywatele i słuźy. /
> [384] Kto wszczął tę nowa zwadę? Mów, synowcze,
> [384] Kto wszczął tę nowa zwadę? , synowcze,
> [394] I harde zionąc mi w uszy wyzwanie,
> [394] I harde zionąc mi w uszy wyzwanie,
> [395] Jał się wywijać nim i sieć powietrze,
> [395] Jał się nim i sieć powietrze,
> [398] Cięcia i pchnięcia zamieniali, zbiegl się
> [398] Cięcia i pchnięcia zamieniali, zbiegl się
> [416] Tam, już tak rano, syn wasz się przechadzał.
> [416] Tam, już tak rano, syn wasz się przechadzał.
> [421] Ledwie go ujrzał, pobiegłem ku niemu.
> [421] Ledwie go ujrzał, pobiegłem ku niemu.
> [418] Lecz on, spostrzegłszy skrył natychmiast
> [419] I w najciemniejszej ukrył się gestwinie.
> [420] Pociąg ten do odosobnienia
> [421] Mierząc mym własnym (serce nasze bowiem
> [422] Jest najczynnniejsze, kiedyśmy samotni),
> [422] Mierząc mym własnym (serce nasze bowiem
> [425] Z dion skrawionych tę broń buntowniczą
> [425] Z dion skrawionych tę broń buntowniczą
> [434] słońce sprzed łoża
> [434] Wesołe słońce sprzed łoża Aurory
> [435] Zaczęło ściągać cienista kotarę,
> [436] On, uciekając od widoku światła,
> [435] Zaczęło ściągać cienista kotarę,
> [436] Onuciekając widoku światła,
> [441] Jeśli się na to lekarstwo nie .
> [441] Jeśli się na to lekarstwo nie znajduje.
> [462] Lecz on jedyny powiernik swych smutków.
> [462] Lecz on powiernik swych smutków
> [463] Tak im jest wierny, tak zamknięty w sobie,
> [463] Tak im jest wierny, tak zamknięty w sobie,
> [469] Nie zbrakłoby nam zaradczego środka.
> [469] Nie zbrakłoby nam zaradczego środka.
> [482] Obyś w tej sprawie, co nam serce ,
> [482] Obyś w tej sprawie, co nam serce rani
> [488] BENWOLIO
> [488]
> [505] Jak nudnie
> [505] Jak nudnie
> [507] Włoka się chwile. Moi-ż to rodzice
> [507] Włoka się chwile. Moi-ż to rodzice
> [506] Włoka się chwile. Moi-ż to rodzice
> [507] Tak spiesznie w tamtą zboczyli ulicę?
> [522] Miłość więc?
> [522] Miłość ?
> [548] Niestety! Czemuż, z zasłona na skroni,
> [548] Niestety! Czemuż, z na skroni,
> [551] Jakis spór? Nie mów mi o , niem wszystkim.
> [551] Jakis spór? Nie mów mi o , niem wszystkim.
> [552] O! wy sprzecznoci niepojęte !
> [552] O! wy sprzecznoci niepojęte !
> [552] W grze tu nienawid wielka, lecz i miłość.
> [553] O! wy sprzecznoci niepojęte dźwię!
> [553] O! wy sprzecznoci niepojęte dźwię!
> [557] Szpetny chaosi wdzieków! Cieżki puchu!
> [557] Szpetny wdzieków! Cieżki puchu!
> [558] Jasna mgło! Zimny żarzeMartyw ruchu!
> [558] Jasna mgło! Zimny żarze! Martyw ruchu!
> [561] Czy się nie śmiejesz?
> [561] Czy się nie śmiejesz?
> [571] Nad czym, pocziwa duszo?
> [571] Nad czym pocziwa duszo?
> [582] A więc strzała
> [582] A więc strzała
> [583] Miłości nawet przez odbitkę działa?
> [583] Miłości nawet przez odbitkę działa?
> [584] Dość mi już ciężyli mój smutek, jego
> [584] Dość mi już ciężyli mój smutek, ty jego
> [587] Nie ulga, ale nowym jest kamieniem
> [587] ulga, nowym jest kamieniem
> [588] Dla mego serca. Miłość, przyjacielu,
> [588] Dla mego serca. Miłość, przyjacielu,
> [589] To dym, co z westchnień się unosi;
> [589] To dym, co z para westchnień się unosi;
> [592] Czymże jest więcej? Istnym amalgamem,
> [592] Czymże jest więcej? Istnym amalgamem,
> [614]
> [616] ROMEO
> [619] Mam-że wraz jęczyć mówić?
> [619] Mam-że wraz jęczyć i mówić?
> [626] Kogóż to kochasz? Powiedz.
> [626] Kogóż to kochasz? Powiedz.
> [632] Pisać testament: będzie-ż to wezwanie
> [632] Pisać testament: będzie-ż to wezwanie
> [640] Gdy to , nimes mi powierzył.
> [640] Gdy to pomysłai, nimes mi powierzył.
> [649] BENWOLIO
> [649]
> [661] Odpiera szturm spojrzeń napastniczych
> [661] Odpiera szturm spojrzeń napastniczych;
> [681] Temu skazanym - wieczne cierpieć męki.
> [681] Temu skazanym - wieczne cierpieć
> [686] Jest na to rada: przestań myśleć o niej.
> [686] Jest na to rada: przestań myśleć o niej.
```

```
In [ ]:
```