

INSTRUKCJA DO PROJEKTU

PRZETWARZANIE JĘZYKA NATURALNEGO

Ćwiczenie 1: Wstępne przetwarzanie tekstu

Jan Daciuk



Unia Europejska
Europejski Fundusz
Rozwoju Regionalnego



Projekt współfinansowany ze środków Unii Europejskiej w ramach Europejskiego Funduszu Rozwoju Regionalnego
Program Operacyjny Polska Cyfrowa na lata 2014–2020

Oś Priorytetowa nr 3 „Cyfrowe kompetencje społeczeństwa” Działanie nr 3.2 „Innowacyjne rozwiązania na rzecz aktywizacji cyfrowej”
Tytuł projektu: „Akademia Innowacyjnych Zastosowań technologii Cyfrowych (AI Tech)”

1 Cel

Zanim zacznie się właściwe przetwarzanie tekstu związane bezpośrednio z celem tego przetwarzania, należy tekst odpowiednio przygotować. Tekst może występować w różnych formatach, z których najpopularniejsze to:

- czysty tekst
- HTML
- XML
- PDF

Nawet czysty tekst może wymagać np. zmiany kodowania liter. Wydobycie czystego tekstu z innych formatów wymaga pewnego wysiłku. Celem tego ćwiczenia jest zapoznanie studentów z problemami związanymi z pozyskiwaniem tekstu i możliwymi ich rozwiązaniami. Spodziewanym wynikiem ma być nabycie umiejętności wydobycia tekstu z plików nie tylko tekstowych, ale także w formatach takich jak HTML, XML i PDF. Nabycie tych umiejętności nastąpić ma przez stopniowe rozwijanie przykładowych programów.

2 Używane narzędzia

Językiem świetnie nadającym się do przetwarzania tekstu jest Python. Książka „*Natural Language Processing with Python*” Stevena Birda, Ewana Kleina i Edwarda Lopera dostępna jest pod adresem: <http://www.nltk.org/book>. Opisuje ona bardzo przydatną bibliotekę `nltk`.

Do przetwarzania tekstów w HTML-u przydatna jest biblioteka BeautifulSoup.

Do wydobywania tekstu z formatu PDF najlepiej skorzystać z bibliotek pdfplumber lub pdfminer. Do odmiany słów można użyć programu `morfeusz2` (z interfejsem w Pythonie) dostępnego na stronie <http://morfeusz.sgjp.pl/>. Do wyszukiwania wzorców w tekście służy biblioteka `re`, chociaż można użyć bardziej rozbudowanej biblioteki `pyparsing`.

3 Źródła danych

W ćwiczeniu będziemy używać następujących źródeł danych:

1. podręczniki systemowe Linuksa dostępne za pomocą polecenia `man`
2. Portal:Current Events dostępny pod adresem https://en.wikipedia.org/wiki/Portal:Current_events (część Wikipedii). Jego polskojęzyczny odpowiednik Portal:Aktualności <https://pl.wikipedia.org/wiki/Portal:Aktualności> niestety (jeszcze) nie zawiera wiele materiału.
3. Narodowy Korpus Języka Polskiego dostępny pod adresem <http://nkjp.pl/index.php?page=14&lang=0>, a konkretnie jego milionowy podkorpus
4. Akty prawne z Internetowego Systemu Aktów Prawnych dostępnego pod adresem <http://isap.sejm.gov.pl/>

4 Przetwarzanie czystych plików tekstowych

Na zajęciach wprowadzających wykonywane ćwiczenia zmierzały do wytworzenia narzędzia do pokazywania użycia słów w kontekście. Zadanie oceniane: należy napisać program realizujący KWIC: *key word in context*. Program powinien akceptować jako parametry słowo oraz listę plików tekstowych i drukować w kolejnych wierszach szukane słowo w środku wiersza począwszy od tej samej kolumny w każdym wierszu, a po jego lewej i prawej stronie lewy i prawy kontekst — słowa, cyfry i znaki przestankowe — w jakim znajduje się kolejne wystąpienie szukanego słowa w danym pliku. Przykład wyszukiwania w Kodeksie Postępowania Cywilnego podany jest poniżej.

./kwic.py uczestnik D19640296Lj.txt

lub 3, wskazany jako wytwórca lub uczestnik niowym: 1) od dnia czynności, gdy uczestnik ów prawnych. Art. 520. § 1. Każdy uczestnik y stosuje się odpowiednio, jeżeli uczestnik ienia o zniesieniu współwłasności uczestnik iu wniosku wnioskodawca bądź inny uczestnik zenie akt uprawniona jest strona, uczestnik zenia może żądać każda strona lub uczestnik ci egzekucyjnych. § 4. Strona lub uczestnik zeniu lub zawieszeniu albo jeżeli uczestnik od dnia jego doręczenia albo gdy uczestnik eprocesowego, jeżeli pozwany albo uczestnik

procesu wprowadzenia do obrotu ta był przy niej obecny lub był o je ponosi koszty postępowania związa postępował niesumienne lub oczyw nie może dochodzić roszczeń przew postępowania zmarł lub został poz postępowania lub interwenient. § postępowania, jeżeli uprawdopodob postępowania zawiadamiają sądowy nie otrzymał zawiadomienia o licy , który takiego żądania nie zgłosi postępowania nieprocesowego podni

W wersji rozszerzonej program powinien mieć opcję włączającą poszukiwanie form odmienionych podanej formy hasłowej słowa. Np. po podaniu formy hasłowej „ołówek”, program znajdzie także formy odmienione „ołówka”, „ołówkowi”, „ołówkiem”, „ołówku”, „ołówki”, „ołówków”, „ołówkom”, „ołówkami” i „ołówkach”, każdą pokazując w kontekście użycia w tekście w podanych plikach.

./kwic.py -m uczestnik D19640296Lj.txt

Sejmu s. 2/524 Art. 3. Strony i uczestnicy rzepisach postępowania stronom i uczestnikom zeby sąd może udzielić stronom i uczestnikom yjaśnienia sprawy. § 2. Strony i uczestnicy ególny stanowi inaczej. Strony i uczestnicy otokółów lub pism. § 2. Strony i uczestnicy zy drzwiach zamkniętych strony i uczestnicy u na utrwalanie przez strony lub uczestników strującego dźwięk. § 2. Strony i uczestnicy k. § 3. Sąd zakazuje stronie lub uczestnikowi unku zaś do przeciwnika i innych uczestników zacji gwarancji procesowych jego uczestników na odległość. W takim przypadku uczestnicy o postępowanie do miejsca pobytu uczestników Sejmu s. 56/524 z miejsca pobytu uczestników pełnomocników, świadków i innych uczestników słuchania stron w sprawie między uczestnikami awną może być dopuszczony między uczestnikami azki są dla nich wspólne. Współ- uczestników kopiowane przez strony i innych uczestników ub 3, wskazany jako wytwórca lub uczestnik Jeżeli weźmie udział, staje się uczestnikiem się, że zainteresowany nie jest uczestnikiem , sąd wezwie go do udziału w spra iału w sprawie wezwany staje się uczestnikiem . W razie potrzeby wyznaczenia ku złożeniu przez któregośkolwiek z uczestników skuteczne tylko wtedy, gdy inni uczestnicy zędu. Art. 513. Niestawiennictwo uczestników zygnięciem sprawy może wysłuchać uczestników rzyrzeczenia oraz w nieobecności uczestników eż zażądać od osób, które nie są uczestnikami iowym: 1) od dnia czynności, gdy uczestnik od dnia doręczenia zawiadomienia uczestnika Krajowego Rejestru Sądowego dla uczestników w prawnych. Art. 520. § 1. Każdy uczestnik em w sprawie. § 2. Jeżeli jednak uczestnicy sztów lub włożyć go na jednego z uczestników ów postępowania wyłożonych przez uczestników czestników. § 3. Jeżeli interesy uczestników są sprzeczne, sąd może włożyć na uczestnika owania poniesionych przez innego uczestnika stosuje się odpowiednio, jeżeli uczestnik 2. Zainteresowany, który nie był uczestnikiem 525. Akta sprawy dostępne są dla uczestników . Sąd wzywa do udziału w sprawie uczestników ełnoletniej, sąd może na wniosek uczestnika aniem postanowienia sąd wysłucha uczestników zędu doręcza odpis postanowienia uczestnikom dopuszczalne, sąd może skierować uczestników ia władzy rodzicielskiej. Jeżeli uczestnicy osobistego stawiennictwa innych uczestników z 1980 r. obowiązuje zastępstwo uczestników lub radcy prawnego, a także gdy uczestnikiem wezwaniem na rozprawę poucza się uczestnika oraz o tym, że w razie działania uczestnika asadnieniem doręcza się z urzędu uczestnikom mowa w art. 5182 § 1, zawiadamia uczestników ddziale niniejszym, sąd wysłucha uczestników enia o zniesieniu współwłasności uczestnik ść rzeczy własność przechodzi na uczestników

postępowania obowiązani są dokon postępowania nie wolno czynić uż postępowania występującym w spra postępowania obowiązani są przyt postępowania mają prawo przegląd postępowania mają prawo do otrzy postępowania mają prawo do otrzy postępowania przebiegu posiedzeń postępowania uprzedzają sąd o za postępowania utrwalenia przebieg - z chwilą doręczenia im tego za , ochronę praw osób, którym pisma postępowania mogą brać udział w postępowania oraz 2022-08-24 ©Ka postępowania do sali sądowej sąd postępowania, wymaga przedstawie tej czynności na fakt jej dokona tej czynności tylko w wypadkach, tych należy zawiadomić o rozpraw postępowania jedynie, w przypadk procesu wprowadzenia do obrotu t . Na odmowę dopuszczenia do wzięc , sąd wezwie go do udziału w spra . W razie potrzeby wyznaczenia ku oświadczenia na piśmie cofnięcie nie sprzeciwili się temu w termi nie tamuje rozpoznania sprawy. P na posiedzeniu sądowym lub zażąd , może również zażądać od osób, k , złożenia wyjaśnień na piśmie. A był przy niej obecny lub był o j o dokonaniu czynności; 3) w przy postępowania, którym postanowien ponosi koszty postępowania związ są w różnym stopniu zainteresowa w całości. To samo dotyczy zwrot . § 3. Jeżeli interesy uczestnikó są sprzeczne, sąd może włożyć na , którego wnioski zostały oddalon . Przepis powyższy stosuje się od postępował niesumienne lub oczy postępowania zakończonych prawom postępowania oraz za zezwoleniem postępowania, w którym zapadło p postępowania lub z urzędu, przy postępowania. § 3. Na postanowie postępowania, prokuratorowi, Pol do mediacji. Przedmiotem mediacj postępowania nie uzgodnili osoby postępowania stosuje się w spraw postępowania przez adwokatów lub postępowania, jego organem, jego o treści przepisów art. 211a § 2 w celu osiągnięcia korzyści mają postępowania oraz prokuratorowi. postępowania, prokuratora, podmi postępowania. Art. 59819. § 1. D nie może dochodzić roszczeń prze wskazanych w postanowieniu. Jeże

u wniosku wnioskodawca bądź inny uczestnik dokonany wpisie sąd zawiadamia uczestników postępowania. Nie zawiadamia się uczestnika wiąże się wpis. § 12. Na wniosek uczestnika enia zawiadomienia o wpisie. Dla uczestnika nik zawiadamia wnioskodawcę oraz uczestników enie, na które wezwie wszystkich uczestników e lub wykaz, sąd zawiadomi o tym uczestników onym przez spadkobiercę sąd oraz uczestnicy nne osoby niż te, które wskazali uczestnicy rozdziału. Jednakże ten, kto był uczestnikiem cywilnego. Art. 683. Na żądanie uczestnika go na podstawie zgodnego wniosku uczestników Krajowego Rejestru Sądowego jest uczestnikiem której wnioskodawca jest jedynym uczestnikiem enie akt uprawniona jest strona, uczestnik enia może żądać każda strona lub uczestnik sąd uwzględni interesy stron lub uczestników ować prawa i obowiązki stron lub uczestników enia o udzieleniu zabezpieczenia uczestnikom Organ egzekucyjny może żądać od uczestników Przepis § 1 stosuje się także do uczestnika . 7631. Komornik poucza strony i uczestników ancję bezpieczeństwa komornika i uczestników ości egzekucyjnej oraz stronom i uczestnikom i egzekucyjnych. § 4. Strona lub uczestnik § 1. Komornik zawiadamia strony, uczestników § 1. Komornik zawiadamia strony, uczestników ienie możliwości obrony praw jej uczestników 5. Komornik umożliwia stronom i uczestnikom celarii komorniczej. Stronom ani uczestnikom 2) wskazać miejsce zamieszkania uczestników a komornik zawiadamia znanych mu uczestników etowej Krajowej Rady Komorniczej uczestników oszacowania zawiadomi znanych mu uczestników nia liczy się od dnia doręczenia uczestnikowi tórym mowa w art. 945 § 4, a dla uczestników ancelaria Sejmu s. 421/524 1) 2) uczestnikom zaofiarował najwyższą cenę, oraz uczestników uchaniu tak jego, jak i obecnych uczestników po wysłuchaniu wnioskodawcy oraz uczestników eniu lub zawieszeniu albo jeżeli uczestnik dzającym, na wniosek nabywcy lub uczestnika 0 sprzedaży komornik zawiadamia uczestników ie licytacji komornik zawiadamia uczestników 4) sumę, jaka przypada każdemu z uczestników cego sprzedaż zarządca zawiadomi uczestników haskiej z 1996 r. i wyznacza dla uczestników w terminie, o którym mowa w § 3, uczestnicy nnego organu państwa obcego albo uczestników nnego organu państwa obcego albo uczestników nnego organu państwa obcego albo uczestników y (WE) nr 2201/2003; 4) wezwania uczestników nnego organu państwa obcego albo uczestników a sporządza się tylko na żądanie uczestnika od dnia jego doręczenia albo gdy uczestnik uzasadnienie na wniosek strony, uczestnika procesowego, jeżeli pozwany albo uczestnik

postępowania zmarł lub został po postępowania. Nie zawiadamia się , który na piśmie zrzekł się zawi postępowania zawarty w akcie not , który zrzekł się zawiadomienia, postępowania o sporządzenie spis sprawy. Jeżeli przed tym posiedz . Art. 659. Spadkobierca, który s mogą zadawać spadkobiercy pytani . W postanowieniu o stwierdzeniu postępowania o stwierdzenie naby działu, zgłoszone nie później ni , dział spadku będzie rozpoznany postępowania, chociażby nie był postępowania, przepis § 1 stosuj postępowania lub interwenient. § postępowania, jeżeli uprawdopodo postępowania w takiej mierze, ab postępowania na czas trwania pos postępowania, wraz z pouczeniem postępowania złożenia wyjaśnień. postępowania, którego dotyczy cz postępowania niezastępowanych pr postępowania, poszanowania godno obecnym podczas czynności dokony postępowania zawiadamiają sądowy postępowania i administracyjny o postępowania, z wyjątkiem dłużni . § 3. Postanowienie o nadaniu kl postępowania zapoznanie się z za postępowania nie wydaje się zapi postępowania. § 2. Jeżeli nieruc . § 2. Komornik wzywa ponadto prz , o których nie ma wiadomości, or oraz dokona obwieszczenia stosow zawiadomienia, o którym mowa w a , którym nie doręczono zawiadomie postępowania; organowi gminy, ur przetargu dokonuje się jedynie w . Art. 988. § 1. Postanowienie o , jeżeli stawia się na posiedzeni nie otrzymał zawiadomienia o lic postępowania. Nabywca może być n stosownie do art. 954. § 3. Jeże postępowania stosownie do art. 9 podziału; sumy, które mają być w stosownie do art. 954. Art. 1064 postępowania termin sześciu tygo postępowania nie wystąpią z wnio postępowania o stwierdzenie jury postępowania o stwierdzenie jury postępowania o stwierdzenie jury postępowania do złożenia wniosku postępowania. Art. 1107. § 1. Do postępowania zgłoszone w termini , który takiego żądania nie zgłos postępowania lub osoby ubiegając postępowania nieprocesowego podn

5 Przetwarzanie stron w HTML-u.

Skrypt `getnews.py` pobiera wiadomości z portalu Current Events Wikipedii i drukuje znalezione odnośniki do Wikipedii z podziałem na dni i kategorie wydarzeń. Korzysta z biblioteki BeautifulSoup. Należy rozbudować go tak, aby w przypadku, gdy zdarzenia zaszły w konkretnych, wymienionych bezpośrednio w tekście państwach, drukował informacje o polu powierzchni tych państw. Informacje te można znaleźć na stronach Wikipedii korzystając z wyluskanych przez skrypt odnośników. Warto wcześniej obejrzeć takie strony i skorzystać z funkcji „Inspect element”. W wersji rozszerzonej powinien także drukować liczbę mieszkańców wymienionych miast. Uwaga: istnieje wiele bibliotek Pythona, które ułatwiają to zadanie, np. `pycountry`, `countrylist`, `geography3`, `flashgeotext`.

getnews.py

```
#!/usr/bin/python
from bs4 import BeautifulSoup
import re
import urllib.request
```

```

def parse_ul_tree(el, top_level = True):
    """
    Parse events of a specific theme at a specific date.
    :param el: unordered list containing themes or individual events.
    :return: a list of all links in that element.
    """
    all_links = []
    for ll in el.find_all("li", recursive=False):
        # Each line contains also subordinate ULs, but up to \n is the current line text.
        tt = ll.text.split("\n")[0]
        tt = re.sub("\([^\)]*\)$", "", tt, count=1) # remove news agency in parentheses if present (.*? didn't work)
        links = [[lnks["href"], lnks.text] for lnks in ll.find_all("a", recursive=False) if lnks["href"][:6] == "/wiki/"]
        all_links.extend(links)
        for uu in ll.find_all("ul", recursive=False):
            all_links.extend(parse_ul_tree(uu, False))
    return all_links

def parse_month():
    """
    Parse a description of events in one month of Wikipedia events.
    It can also be used to parse the current events.
    :return: nothing
    """
    base = "https://en.wikipedia.org/wiki/Portal:Current_events"
    with urllib.request.urlopen(base) as fp:
        soup = BeautifulSoup(fp, "html.parser")
        # The div with region role is deep in the body hierarchy. It represents one day.
        for d in soup.find_all("div", role="region"):
            date = []
            category = ""
            if "id" in d.attrs:
                year, month, day = d["id"].split("_")
                date = [day, month, year]
                print("{0} {1} {2}".format(day, month, year))
            for t1 in d.find_all(recursive=False):
                if t1.name == "div" and "class" in t1.attrs and "description" in t1["class"]:
                    for t in t1.find_all(recursive=False):
                        if t.name == "p":
                            # An opening paragraph contains category name
                            category = t.b.contents[0]
                            print(category)
                        if t.name == "ul":
                            # The following unordered list contains events grouped in themes.
                            # Each theme is one line in the top-level list
                            # Parse those events
                            events = parse_ul_tree(t)
                            print(events)

if (__name__ == "__main__"):
    parse_month()

```

6 Przetwarzanie plików XML

Narodowy Korpus Języka Polskiego jest zbiorem oznaczonych tekstów w języku polskim. Jego stumilionowy podzbiór jest publicznie dostępny. Każdy tekst opisany jest w katalogu zawierającym pliki:

- `ann_groups.xml` — grupy składniowe
- `ann_morphosyntax.xml` — interpretacje morfoskładniowe poszczególnych segmentów
- `ann_named` — jednostki nazewnicze
- `ann_segmentation` — reprezentacja segmentacji tekstu na zdania i na segmenty (w przybliżeniu: słowa ortograficzne)
- `ann_senses` — informacja o znaczeniach wybranych słów wieloznacznych
- `ann_words` — słowa składniowe i ich interpretacje morfoskładniowe
- `header.xml` — nagówek, zawiera m.in. informacje o pochodzeniu tekstu
- `text.xml` — właściwy tekst

Skrypt `proc_nkjp.py` wydobywa i drukuje z katalogu opisu tekstu o danej nazwie informacje z pliku `ann_morphosyntax.xml` o słowie i jego możliwych interpretacjach morfoskładniowych.

`proc.nkjp.py`

```
#!/usr/bin/python3
import xml.etree.ElementTree as ET

def tag_uri_and_name(elem):
    """
    Divide element name with a prefix being namespace
    into namespace uri and tag.
    :param elem: element in an xml tree
    :return: a 2-tuple: namespace uri and tag
    """
    if elem.tag[0] == "{":
        uri, ignore, tag = elem.tag[1:].partition("}")
    else:
        uri = None
        tag = elem.tag
    return uri, tag

def get_next_morph(filename):
    """
    Get the next word with all its annotations.
    :param filename: file named ann_morphosyntax.xml in NKJP
    :return: a 3-tuple with the word, its possible interpretations,
    and the correct interpretation
    """
    events = ["start-ns", "start", "end"]
    xmlns = ""
    path = []
    ctags = []
    msds = []
    interps = []
    disamb = []
    base = ""
    orth = ""
    for (event, elem) in ET.iterparse(filename, events=events):
        # React to opening and closing of tags
        tag = ""
        if event == "start-ns":
            if elem[0] == "{":
                xmlns = elem[1]
        elif event == "start":
            # For tag openings, construct a path to the current tag.
            # The path is a list.
            # Use either tags, or tags concatenated with some attribute
            # values (e.g. "type", "name", or "value") as path items.
            # Store path current leaves in variable `branch`
            ns, tag = tag_uri_and_name(elem)
            branch = tag
            if tag == "fs":
                branch = tag + ":" + elem.attrib["type"]
            elif tag == "f":
                branch = tag + ":" + elem.attrib["name"]
            elif tag == "symbol":
                if path[-1] == "f:ctag":
                    # The part-of-speech tag for the word
                    ctags.append(elem.attrib["value"])
                elif (path[-1] == "f:msd"
                      or path[-1] == "vAlt" and path[-2] == "f:msd"):
                    # Possible morphosyntactic descriptions of the word,
                    # as obtained from a morphological analyzer
                    msd = elem.attrib["value"]
                    # They are appended to the part-of-speech
                    msd_suffix = msd
                    if msd != "":
                        msd_suffix = ":" + msd
                    for c in ctags:
                        # Append msds to each part-of-speech unless empty
                        if c == "ign":
                            interps.append("ign")
                        else:
                            interps.append(base + ":" + c + msd_suffix)
                    msds.append(elem.attrib["value"])
            path.append(branch)
        elif event == "end":
            branch = path.pop()
            if branch == "f:interps":
```

Dostarczony skrypt należy rozszerzyć w ten sposób, aby drukowana była też informacja o właściwym oznaczeniu (w skrypcie przypisywana do zmiennej `disamb`). Aby była właściwie wydrukowana, należy oczywiście wpisać do niej odpowiednią wartość. Nazwa zmiennej sugeruje, gdzie w drzewie XML znaleźć pożądaną informację. Poniżej przedstawiono początek wzorcowego wyjścia skryptu dla pierwszego podkatalogu NKJP.

7

W wersji rozszerzonej skrypt powinien jako pierwszy argument przyjąć słowo, dla którego informacje (jak w powyższym przykładzie) mają być wypisane, razem z analogiczną informacją dla słowa poprzedzającego i następującego. Przykład poniżej („<” oznacza poprzednie słowo, „=” — bieżące, a „>” — następne słowo):

```
<=orth=stanowi, interps=['stan:subst:sg:dat:m3', 'stanowić:subst:sg:ter:imperf', 'stanowić:fin:sg:ter:imperf', 'stanowy:subst:pl:nom:m1:pos', 'stanowy:fin:pl:nom:m1:pos', 'stanowię:subst:pl:voc:m1:pos', 'stanowię:fin:pl:voc:m1:pos', 'stanowię:adj:pl:voc:m1:pos'], disamb=stanowić:fin:sg:ter:imperf  
==orth=o, interps=['o:interj', 'o:interj:acc', 'o:prep:acc', 'o:interj:loc', 'o:prep:loc', 'ojciec:interj:pun', 'ojciec:prep:pun', 'ojciec:brev:pun'], disamb=o:prep:loc  
>=orth=wymiarze, interps=['wymiar:subst:sg:loc:m3', 'wymiar:subst:sg:voc:m3'], disamb=wymiar:subst:sg:loc:m3  
<=orth=decydująca, interps=['decydować:pact:sg:acc:f:imperf:aff', 'decydować:pact:sg:inst:f:imperf:aff', 'decydujący:pact:sg:acc:f:pos', 'decydujący:adj:sg:acc:f:pos', 'decydujący:pact:sg:inst:f:pos', 'decydujący:adj:sg:inst:f:pos'], disamb=decydować:pact:sg:inst:f:imperf:aff  
==orth=o, interps=['o:interj', 'o:interj:acc', 'o:prep:acc', 'o:interj:loc', 'o:prep:loc', 'ojciec:interj:pun', 'ojciec:prep:pun', 'ojciec:brev:pun'], disamb=o:prep:loc  
>=orth=losie, interps=['los:subst:sg:loc:m3', 'los:subst:sg:voc:m3'], disamb=los:subst:sg:loc:m3  
<=orth=chodzi, interps=['chodzić:fin:sg:ter:imperf'], disamb=chodzić:fin:sg:ter:imperf  
==orth=o, interps=['o:interj', 'o:interj:acc', 'o:prep:acc', 'o:interj:loc', 'o:prep:loc', 'ojciec:interj:pun', 'ojciec:prep:pun', 'ojciec:brev:pun'], disamb=o:prep:acc  
>=orth=sensowność, interps=['sensowność:subst:sg:nom:f', 'sensowność:subst:sg:acc:f'], disamb=sensowność:subst:sg:acc:f  
<=orth=już, interps=['już:qub'], disamb=już:qub  
==orth=o, interps=['o:interj', 'o:interj:acc', 'o:prep:acc', 'o:interj:loc', 'o:prep:loc', 'ojciec:interj:pun', 'ojciec:prep:pun', 'ojciec:brev:pun'], disamb=o:prep:acc  
>=orth=religię, interps=['religia:subst:sg:acc:f'], disamb=religia:subst:sg:acc:f  
<=orth=, interps=[':interp'], disamb=:interp  
==orth=o, interps=['o:interj', 'o:interj:acc', 'o:prep:acc', 'o:interj:loc', 'o:prep:loc', 'ojciec:interj:pun', 'ojciec:prep:pun', 'ojciec:brev:pun'], disamb=o:prep:acc  
>=orth=tajemnicę, interps=['tajemnica:subst:sg:acc:f'], disamb=tajemnica:subst:sg:acc:f  
<=orth=aale, interps=['Al:depr:pl:nom:m2', 'Al:depr:pl:voc:m2', 'Ala:depr:pl:nom:f', 'Ala:subst:pl:nom:f', 'Ala:depr:pl:acc:f', 'Ala:subst:pl:acc:f', 'Ala:depr:pl:voc:f', 'Ala:subst:pl:voc:f', 'Ali:depr:pl:nom:m2', 'Ali:subst:pl:nom:m2', 'Ali:depr:pl:nom:m2', 'Ali:depr:pl:voc:m2', 'Ali:subst:pl:voc:m2', 'Ali:depr:pl:voc:m2', 'Alo:depr:pl:nom:m2', 'Alo:subst:pl:nom:m2', 'Alo:depr:pl:nom:m2', 'Alo:depr:pl:nom:m2', 'Alo:subst:pl:nom:m2', 'Alo:depr:pl:nom:m2', 'Alo:subst:pl:nom:m1', 'Alo:depr:pl:nom:m1', 'Alo:depr:pl:nom:m1', 'Alo:subst:pl:nom:m1', 'Alo:depr:pl:voc:m1', 'Alo:subst:pl:voc:m1', 'Alo:depr:pl:voc:m1', 'Alo:subst:pl:voc:m1', 'Ale:depr', 'Ale:subst', 'Ale:depr', 'Ale:subst', 'Ale:conj', 'Ale:depr', 'Ale:conj', 'Ale:qub'], disamb=aale:conj  
==orth=o, interps=['o:interj', 'o:interj:acc', 'o:prep:acc', 'o:interj:loc', 'o:prep:loc', 'ojciec:interj:pun', 'ojciec:prep:pun', 'ojciec:brev:pun'], disamb=o:prep:acc  
>=orth=normę, interps=['norma:subst:sg:acc:f'], disamb=norma:subst:sg:acc:f
```

Prawie żadne akty prawne publikowane przez Sejm RP i rząd nie zawierają spisów treści w metadanych. Nie zawierają ich też, ze względu na charakter aktów prawnych, w treści. Czasami okazują się nawet skanami bez bezpośredniego dostępu do treści jako ciągu znaków. Skrypt `ustawa.py` drukuje bardzo uproszczony spis treści podając numery rozdziałów i strony, na których się zaczynają, o ile w treści dokumentu są takowe wyróżnione i dokument nie jest skanem.

8


```

"""
Print chapters in a PDF file (in Polish).
:param filename: PDF file name
:return: nothing
Processing is simplified. The file may contain no chapters.
"""
with pdfplumber.open(filename) as pdf:
    page_number = 1
    for p in pdf.pages:
        t = p.extract_text()
        for r in re.findall(r'^(Rozdział)\W(\w*)', t,
                           re.M | re.I):
            print(r[0], r[1], 'strona', page_number)
            page_number = page_number + 1

if (__name__ == "__main__"):
    import sys
    for filename in sys.argv[1:]:
        print(filename)
        print_chapters(filename)

```

Sam numer rozdziału pozwala się szybciej poruszać po dokumencie, ale nie informuje o znaczeniu poszczególnych części dokumentu. Należy podany skrypt rozszerzyć, tak aby drukował tytuły części, ksiąg, działów, rozdziałów i numery artykułów, o ile takie występują. Poniżej podany jest przykład początkowych wierszy wyjścia dla kodeksu postępowania cywilnego.

```
./ustawa-base.py D20211805Lj.pdf
```

```

D20211805Lj.pdf
TYTUŁ WSTĘPNY strona 1
CZĘŚĆ PIERWSZA strona 4
KSIĘGA PIERWSZA strona 4
TYTUŁ I strona 4
DZIAŁ I strona 4
Rozdział 1 strona 4
Oddział 1 strona 4
Oddział 2 strona 5
Rozdział 2 strona 7
Oddział 1 strona 7
Oddział 2 strona 7
Oddział 3 strona 9
Oddział 4 strona 10
DZIAŁ II strona 11
DZIAŁ III strona 13
TYTUŁ II strona 15
TYTUŁ III strona 16
TYTUŁ IIIa strona 17
TYTUŁ IIIb strona 18
TYTUŁ IIIc strona 18
TYTUŁ IV strona 18
DZIAŁ I strona 18
...

```

W wersji rozszerzonej należy także drukować także tytuł aktu prawnego. Można to osiągnąć na różne sposoby, np. inspirując się rozwiązaniem podanym na stronie: https://pdfminersix.readthedocs.io/en/latest/howto/character_properties.html. Przykładowe wyjście podane jest poniżej. Swoboda prawodawcy w ustalaniu formatu jednostek redakcyjnych aktu czyni dokładne wykonanie tego zadania dość karkołomnym; w przykładzie widać pewne uproszczenia.

```
./ustawa-ext.py D20211805Lj.pdf
```

```

D20211805Lj.pdf
Ustawa z dnia 17 listopada 1964 r. Kodeks postępowania cywilnego
TYTUŁ WSTĘPNY strona 1
CZĘŚĆ PIERWSZA strona 4
KSIĘGA PIERWSZA strona 4
TYTUŁ I strona 4
DZIAŁ I strona 4
Rozdział 1 strona 4
Oddział 1 strona 4
Oddział 2 strona 5
Rozdział 2 strona 7
Oddział 1 strona 7
Oddział 2 strona 7
Oddział 3 strona 9
Oddział 4 strona 10

```

DZIAŁ II strona 11
DZIAŁ III strona 13
TYTUŁ II strona 15
TYTUŁ III strona 16
TYTUŁ IIIa strona 17
TYTUŁ IIIb strona 18
TYTUŁ IIIc strona 18
TYTUŁ IV strona 18
DZIAŁ I strona 18
DZIAŁ II strona 21
DZIAŁ III strona 22
DZIAŁ IV strona 23
DZIAŁ V strona 24
TYTUŁ V strona 28
DZIAŁ I strona 28
DZIAŁ II strona 32
TYTUŁ VI strona 39
DZIAŁ I strona 39
Rozdział 1 strona 39
Rozdział 2 strona 47
Rozdział 3 strona 56
Rozdział 4 strona 62
Rozdział 5 strona 63
Rozdział 6 strona 64
DZIAŁ II strona 69
Rozdział 1 strona 69
Oddział 1 strona 69
Oddział 2 strona 75
Rozdział 2 strona 75
Rozdział 2a strona 83
Rozdział 3 strona 89
DZIAŁ III strona 95
Rozdział 1 strona 95
Rozdział 2 strona 96
Oddział 1 strona 96
Oddział 2 strona 98
Oddział 3 strona 101
Oddział 4 strona 104
Oddział 5 strona 107
Oddział 6 strona 108
Oddział 7 strona 109
Rozdział 3 strona 110
DZIAŁ IV strona 111
Rozdział 1 strona 111
Oddział 1 strona 111
Oddział 2 strona 116
Oddział 3 strona 118
Oddział 4 strona 120
Rozdział 1a strona 121
Rozdział 2 strona 121
Rozdział 3 strona 124
DZIAŁ V strona 125
Rozdział 1 strona 125
Rozdział 11 strona 133
Rozdział 2 strona 133
DZIAŁ Va strona 138
DZIAŁ Vb strona 144
DZIAŁ VI strona 145
DZIAŁ VII strona 149
DZIAŁ VIII strona 149
TYTUŁ VII strona 152
DZIAŁ I strona 152
Rozdział 1 strona 152
Rozdział 2 strona 153
Rozdział 3 strona 156
DZIAŁ II strona 157
DZIAŁ IIA strona 159
DZIAŁ IIB strona 163
DZIAŁ III strona 165
Rozdział 1 strona 165
Rozdział 2 strona 171
Rozdział 3 strona 173
DZIAŁ IV strona 179
DZIAŁ IVa strona 179
Rozdział 1 strona 179
DZIAŁ IVb strona 185
DZIAŁ IVc strona 185
DZIAŁ IVd strona 187
DZIAŁ IVe strona 189
DZIAŁ IVf strona 191
DZIAŁ IVg strona 192
Rozdział 1 strona 193
Rozdział 2 strona 194

Rozdział 3	strona 197
Rozdział 4	strona 198
Rozdział 5	strona 201
DZIAŁ V	strona 204
Rozdział 1	strona 204
Rozdział 2	strona 206
Rozdział 3	strona 209
DZIAŁ VI	strona 209
DZIAŁ VII	strona 213
Rozdział 1	strona 213
Rozdział 2	strona 214
DZIAŁ VIII	strona 216
Rozdział 1	strona 216
KSIEGA DRUGA	strona 219
TYTUŁ I	strona 219
TYTUŁ II	strona 226
DZIAŁ I	strona 226
Rozdział 1	strona 226
Oddział 1	strona 226
Oddział 2	strona 227
Oddział 3	strona 228
Rozdział 2	strona 229
Oddział 1	strona 229
Oddział 2	strona 230
Oddział 3	strona 231
DZIAŁ Ia	strona 234
DZIAŁ II	strona 240
Rozdział 1	strona 240
Rozdział 2	strona 243
Oddział 1	strona 243
Oddział 2	strona 248
Oddział 3	strona 253
Oddział 4	strona 256
Oddział 5	strona 257
Oddział 6	strona 262
Rozdział 3	strona 264
DZIAŁ III	strona 266
Rozdział 1	strona 266
Rozdział 2	strona 267
Rozdział 2a	strona 267
Oddział 1	strona 267
Oddział 2	strona 268
Rozdział 3	strona 269
Rozdział 4	strona 269
Rozdział 5	strona 271
Rozdział 6	strona 272
DZIAŁ IV	strona 277
Rozdział 1	strona 277
Rozdział 2	strona 284
Rozdział 3	strona 285
Rozdział 4	strona 287
Rozdział 5	strona 288
Rozdział 6	strona 288
Rozdział 6a	strona 289
Rozdział 7	strona 290
Rozdział 8	strona 291
Rozdział 9	strona 294
Dział spadku	strona 294
Rozdział 10	strona 295
DZIAŁ IVa	strona 295
DZIAŁ IVb	strona 297
DZIAŁ V	strona 298
Rozdział 1	strona 298
Rozdział 2	strona 300
Rozdział 3	strona 301
DZIAŁ VI	strona 302
KSIEGA TRZECIA	strona 306
KSIEGA CZWARTA	strona 306
CZĘŚĆ DRUGA	strona 308
TYTUŁ I	strona 309
TYTUŁ II	strona 315
TYTUŁ III	strona 322
CZĘŚĆ TRZECIA	strona 328
TYTUŁ I	strona 328
DZIAŁ I	strona 328
DZIAŁ II	strona 343
DZIAŁ IIa	strona 354
DZIAŁ IIb	strona 356
DZIAŁ IIc	strona 356
DZIAŁ IIId	strona 356
DZIAŁ IIIda	strona 357
DZIAŁ IIe	strona 358
DZIAŁ IIIf	strona 359

DZIAŁ III strona 360
DZIAŁ IV strona 372
DZIAŁ V strona 378
DZIAŁ VI strona 388
TYTUŁ II strona 391
DZIAŁ I strona 391
Rozdział 1 strona 391
Rozdział 2 strona 397
Rozdział 3 strona 403
DZIAŁ II strona 407
DZIAŁ III strona 411
DZIAŁ IV strona 417
DZIAŁ IVa strona 424
DZIAŁ V strona 431
DZIAŁ VI strona 434
Rozdział 1 strona 434
Rozdział 2 strona 434
Rozdział 3 strona 440
Rozdział 4 strona 444
Rozdział 5 strona 447
Rozdział 6 strona 448
Rozdział 6a strona 451
Rozdział 7 strona 455
Rozdział 8 strona 457
Rozdział 9 strona 459
DZIAŁ VIa strona 461
DZIAŁ VII strona 463
DZIAŁ VIII strona 465
Rozdział 1 strona 465
Rozdział 2 strona 471
Rozdział 3 strona 472
Rozdział 4 strona 472
TYTUŁ III strona 474
DZIAŁ I strona 474
DZIAŁ II strona 483
Rozdział 1 strona 483
Rozdział 2 strona 485
Rozdział 3 strona 488
DZIAŁ III strona 491
DZIAŁ IV strona 492
DZIAŁ V strona 493
DZIAŁ VI strona 495
CZEŚĆ CZWARTA strona 495
KSIĘGA PIERWSZA strona 495
TYTUŁ I strona 495
TYTUŁ III strona 497
TYTUŁ IV strona 501
TYTUŁ IVa strona 506
KSIĘGA PIERWSZA strona 507
KSIĘGA DRUGA strona 510
TYTUŁ I strona 510
TYTUŁ II strona 510
TYTUŁ III strona 512
TYTUŁ IV strona 512
TYTUŁ V strona 517
TYTUŁ VI strona 517
TYTUŁ VII strona 517
TYTUŁ VIIA strona 518
TYTUŁ VIII strona 520
TYTUŁ IX strona 520
TYTUŁ X strona 520
TYTUŁ XI strona 521
KSIĘGA TRZECIA strona 523
TYTUŁ I strona 523
TYTUŁ II strona 525
TYTUŁ III strona 527
TYTUŁ IV strona 527
TYTUŁ V strona 527
TYTUŁ VI strona 527
KSIĘGA CZWARTA strona 527
TYTUŁ I strona 528
TYTUŁ II strona 529
TYTUŁ III strona 533
CZEŚĆ PIĄTA strona 534
TYTUŁ I strona 534
TYTUŁ II strona 536
TYTUŁ III strona 539
TYTUŁ IV strona 542
TYTUŁ V strona 544
TYTUŁ VI strona 547
TYTUŁ VII strona 549
TYTUŁ VIII strona 552

8 Kryteria oceny

Za poszczególne zadania można dostać 2 punkty w wersji podstawowej. Rozszerzenie pozwoli zdobyć dodatkowe pół punktu.

- KWIC: wersja podstawowa — bez odmiany, wersja rozszerzona — możliwość odmiany
- informacje o miejscach wydarzeń: wersja podstawowa — kraje z podanym polem powierzchni, wersja rozszerzona — także miejscowości z liczbą mieszkańców
- właściwe oznaczenie segmentu w NKJP: wersja podstawowa — wypisanie realizowane w programie przykładowym, wersja rozszerzona — umożliwienie wyszukiwania określonego słowa i wypisanie informacji jak w wersji podstawowej, ale tylko o poprzednim, danym i następnym słowie
- spis treści aktu prawnego: wersja podstawowa — tylko spis treści, wersja rozszerzona — także tytuł ustawy.

Dodatek — Instalacja oprogramowania

Przykładowe programy wymagają Pythona i jego bibliotek. Najlepiej używać nowszych ich wersji, w szczególności Pythona w wersji 3. W najnowszych dystrybucjach Linuksa jest on dostępny, ale w innych, aby domyślnie Python był używany wersji 3, konieczna może być instalacja pakietu `python-is-python3`. Biblioteki można instalować korzystając z polecenia `pip3` (na Fedorze itp. zapewne po prostu `pip`) lub `conda`. To ostatnie wymaga instalacji `miniconda` (lub pełnej `anacondy`). Można to zrobić w przestrzeni użytkownika. Opis na stronie <https://docs.conda.io/en/latest/miniconda.html>. Jeżeli okaże się, że na danym komputerze w laboratorium `conda` nie jest zainstalowana, można ją doinstalować w przestrzeni użytkownika korzystając z opisu na podanej stronie, ale należy skorzystać z archiwum ze strony <https://repo.anaconda.com/miniconda/> do ściągnięcia wersji działającej z Pythonem 2.7; potem można zainstalować Python 3 i biblioteki dla Pythona 3.

