



دانشکده مهندسی کامپیوتر

پردازش زبان‌های طبیعی

تمرین دوم

دکتر احسان‌الدین عسگری

گروه ۱، آیین احادی نیا، محمد مهدی ابوترابی، ارشان دلیلی

۱۲ آذر ۱۴۰۱

فهرست مطالب

۱	پیش‌پردازش	۳
۲	استخراج نام	۳
۳	استخراج ایمیل	۴
۴	استخراج شهر	۴
۵	استخراج تاریخ تولد	۵
۶	استخراج شماره تلفن	۶
۷	استخراج وضعیت شغلی	۷
۸	استخراج عنوان شغلی	۷
۹	استخراج حقوق مورد انتظار	۸
۱۰	استخراج نوع شغل مورد نظر	۹
۱۱	استخراج بخش‌های اضافه	۹
۱۲	نتایج	۹

در این تمرین، ابزاری جهت استخراج اطلاعات رزومه با استفاده از عبارات منظم برای زبان فارسی توسعه می‌دهیم.

۱ پیش‌پردازش

برای پیش‌پردازش از ابزار نرمالیزیشن کتاب‌خانه hazm استفاده می‌کنیم. البته در بعضی از مازول‌ها و تسک‌ها پیش‌پردازش‌های مخصوص آن را اعمال می‌کنیم.

۲ استخراج نام

برای استخراج نام، مازول NameDetector را می‌سازیم. این مازول با دریافت اسامی کراول شده در اینترنت (شامل نام و نام‌خانوادگی) و پترن‌های مشخص نام، شروع به استخراج نام‌ها از رزومه می‌کند. تصویر این مازول را در زیر مشاهده می‌نمایید. این مازول توانایی پیدا کردن نام‌های ترکیبی و نام‌هایی که نام خانوادگی آن در دیتاست نام‌های خانوادگی وجود ندارد را نیز دارا می‌باشد. (لازم به ذکر است دیتاست نام در فولدر resources قرار دارد.)

```
import json
import re

class NameDetection:
    def __init__(self):
        with open('cv_info_extractor/resources/last_names_regex.json', 'r', encoding='utf-8') as file:
            self.last_names_reg = json.loads(file.read())

        with open('cv_info_extractor/resources/first_names_regex.json', 'r', encoding='utf-8') as file:
            self.first_names_reg = json.loads(file.read())

        self.pattern = f'^([\\W]{{{self.first_names_reg}}}[\\W]+){(self.last_names_reg)}{{{self.first_names_reg}}}[\\W]+([\\W]{{{self.last_names_reg}}}[\\W]+)$'

    def match_name(self, inp):
        matches = []
        count_pattern = self.pattern.format()
        for matched in re.finditer(count_pattern, inp):
            start, end = matched.span()
            inp = inp[start] + '#' * (end - start) + inp[end:]
            matches.append(matched)
        return matches

    def find_name(self, text):
        matched_names = self.match_name(text)
        if not matched_names:
            return ['Not Found'] * 3
        full_name = matched_names[0].group().strip()
        if matched_names[0].groups()[13]:
            first_name = matched_names[0].groups()[3].strip()
        elif len(matched_names[0].groups()) > 9 and matched_names[0].groups()[10]:
            first_name = matched_names[0].groups()[9].strip()
        else:
            first_name = 'None'
        if len(matched_names[0].groups()) > 7 and matched_names[0].groups()[7]:
            last_name = matched_names[0].groups()[7].strip()
        elif matched_names[0].groups()[4]:
            last_name = matched_names[0].groups()[4].strip()
        else:
            last_name = 'None'
        return full_name, first_name, last_name
```

شکل ۱: تصویر مازول نام

۳ استخراج ایمیل

برای استخراج ایمیل، ماژول EmailDetector را می‌سازیم. این ماژول با پترن‌های ایمیل شروع به استخراج ایمیل از رزومه می‌کند. تصویر این ماژول را در زیر مشاهده می‌نمایید.

```
class EmailDetection:
    def __init__(self):
        self.pattern = r"\b(\w+([-+.\[\]\dot\@])\w+)*@(\[\at\])\w+([-+.\[\]\dot\@])\w+*(\.\[\dot\@])\w+([-+.\[\]\dot\@])\w+*"

    def match_email(self, inp):
        matches = []
        count_pattern = self.pattern.format()
        for matched in re.finditer(count_pattern, inp):
            start, end = matched.span()
            inp = inp[:start] + '#' * (end - start) + inp[end:]
            matches.append(matched)
        return matches

    def find_email(self, text):
        matched_emails = self.match_email(text)
        if not matched_emails:
            return 'Not Found'
        return matched_emails[0].group().strip()
```

شکل ۲: تصویر ماژول ایمیل

۴ استخراج شهر

برای استخراج شهر، ابتدا اسامی را که بین شهرها و استان‌ها مشترک هستند (مانند تهران و یزد) کراول می‌کنیم و در فایل cities.csv ذخیره می‌کنیم. سپس کل متن را پیمایش می‌کنیم و اسامی شهرها، استان‌ها و اسامی مشترک را mask می‌کنیم. سپس با استفاده از الگوهای این ماژول، نام شهرها و استان‌ها را استخراج می‌کنیم.

```

import pandas as pd
import re

class CityProvinceExtractor:
    def __init__(self):
        df = pd.read_csv("cv_info_extractor/resources/cities.csv")
        self.cities = set(df["city"].to_list())
        self.provinces = set(df["province"].to_list())
        self.common_names = set.intersection(self.cities, self.provinces)

    def find(self, original_text):
        original_text = re.sub(r"[\s]", "", original_text)
        original_text = re.sub(r"\\s", " ", original_text)
        original_text = re.sub(r"\\n", " ", original_text)
        text = original_text

        for common_name in self.common_names:
            text = text.replace(common_name, "#" * len(common_name))
        for city in self.cities:
            text = text.replace(city, "#" * len(city))
        for province in self.provinces:
            text = text.replace(province, "#" * len(province))

        matches = []

        for match in re.finditer(r"([#6]+) ([#6]+) ", text):
            city = original_text[match.span(1)[0]: match.span(1)[1]]
            province = original_text[match.span(2)[0]: match.span(2)[1]]
            text = {
                "text": match.start(),
                "city": match.end() - match.start(),
                "province": match.end(),
            }
            matches.append({"city": city, "province": province})

        for match in re.finditer(r"([#6]+) استان ([#6]+) ", text):
            city = original_text[match.span(1)[0]: match.span(1)[1]]
            province = original_text[match.span(2)[0]: match.span(2)[1]]
            text = {
                "text": match.start(),
                "city": match.end() - match.start(),
                "province": match.end(),
            }
            matches.append({"city": city, "province": province})

        for match in re.finditer(r"([#6]+) ", text):
            province = original_text[match.span(1)[0]: match.span(1)[1]]
            text = {
                "text": match.start(),
                "city": match.end() - match.start(),
                "province": match.end(),
            }
            matches.append({"city": None, "province": province})

        if not matches:
            return [{"city": "None", "province": "None"}]

```

شکل ۳: تصویر مازول استخراج شهر

۵ استخراج تاریخ تولد

برای استخراج تاریخ تولد، ابتدا با استفاده از الگوهای تاریخ شروع به استخرا تاریخ می‌کنیم. سپس با چک کردن مکان کلمات مرتبط با تاریخ تولد، نزدیک‌ترین تاریخ (که تا ۱۰۰ کاراکتر بعد می‌تواند باشد) را برمی‌گردانیم.

شکل ۴: تصویر مازول استخراج تاریخ تولد

برای استخراج شماره تلفن، ابتدا الگوهای متداول شماره تلفن و انواع مرسوم نوشتن شماره را به دست آورده و با استفاده از این الگوها، شماره تلفن را به دست می‌آوریم.

شکل ۵: تصویر مازول استخراج شماره تلفن

۷ استخراج وضعیت شغلی

برای استخراج وضعیت شغلی، از تعدادی کلمات کلیدی استفاده می‌کنیم. سپس با استفاده از تعدادی حالت پیش‌فرض برای وضعیت شغلی، در صورت وجود آن‌ها، برگردانده می‌شوند و در غیر این صورت، در سوابق شغلی شروع به بررسی تعدادی کلمات کلیدی با مضمون “تا کنون” می‌کند و در صورت وجود آن‌ها را برمی‌گرداند. هم‌چنین برای حقوق مورد انتظار و نوع شغل مورد نظر نیز ماژول‌های جداگانه ساخته شده که از روی کلمات کلیدی و مقادیر مشخص اقدام به استخراج این موارد می‌نماید که ساختاری بسیار شبیه به وضعیت اشتغال دارد.

```
class EmploymentStatusExtractor:
    def __init__(self) -> None:
        keys = [
            "وضعیت اشتغال",
            "وضعیت اشتغالی",
            "وضعیت استخدام",
            "وضعیت استخدامی",
            "وضعیت کار",
            "وضعیت کاری",
        ]
        values = [
            "علاقه مند",
            "علاقه‌مند",
            "علاقه مند",
            "بیکار",
            "بفول",
            "فعال",
        ]
        up_to_now = [
            "تاکنون",
            "تا اکنون",
            "تا زمان حال",
            "تا حال",
            "تا الان",
            "تا الان",
        ]
        self.pattern_keys = r"({})\s*:".format("|".join(keys))
        self.pattern_values = r"({})".format("|".join(values))
        self.pattern_up_to_now = r"({})".format("|".join(up_to_now))

    def find(self, text1, job_text):
```

شکل ۶: تصویر ماژول وضعیت شغلی

۸ استخراج عنوان شغلی

در این بخش ما دیتاستی از عناوین شغلی مختلف که در رزومه‌های افراد پرتکرار است جمع‌آوری کردیم. سپس این عناوین شغلی را با انجام پیش‌پردازش مناسب با بخش سرآیند و بخش سوابق شغلی رزومه تطبیق می‌دهیم و اولین مورد انطباق را به عنوان خروجی برمی‌گردانیم. این کار از این جهت منطقی است که سوابق شغلی به ترتیب نزولی بر حسب تاریخ و اهمیت مرتب می‌شود بنابراین اولین مورد انطباق، آخرین سابقه شغلی فرد است.

```
import pandas as pd

class JobTitleFinder:
    def __init__(self) -> None:
        titles = list(
            pd.read_csv(
                "cv_info_extractor/resources/jobs.csv",
            )["title"]
        )
        expanded_1 = [
            text.replace("مهندس", "مهندسی") for text in titles if "مهندس" in text
        ]
        titles += expanded_1
        self.titles = sorted(titles, key=lambda s: -len(s))
        self.normalized_title = [self.normalize(text) for text in self.titles]

    def normalize(self, txt):
        txt = txt.replace(" ", "")
        txt = txt.replace("\n", "")
        txt = txt.replace("\u200C", "")
        return txt

    def find(self, text):
        normalized_text = self.normalize(text)
        min_index = float("inf")
        for idx, title in enumerate(self.normalized_title):
            if title in normalized_text:
                index = normalized_text.index(title)
                if index < min_index:
                    return self.titles[idx]
        return 'None'
```

شکل ۷: تصویر ماژول عنوان شغلی

۹ استخراج حقوق مورد انتظار

در این ماژول ما با پیش پردازش مناسب، به دنبال کلماتی مانند «حقوق مورد نظر»، «حقوق مد نظر» و موارد از این دست می‌گردیم و در فاصله مناسب از این کلمات، عدد نوشته شده را به کاربر برمی‌گردانیم.

```
import re

class ExpectedSalaryExtractor:
    def __init__(self) -> None:
        keys = [
            "حقوق مورد انتظار",
            "حقوق مورد انتظار",
            "حقوق مدنظر",
            "حقوق دلخواه",
        ]
        self.pattern_keys = r"({})\s*:".format("|".join(keys))

    def find(self, text1):
        text = text1
        for match_key in re.finditer(self.pattern_keys, text):
            text2 = text[match_key.end():]
            if text2.find('\n') != -1:
                return text2[:text2.find('\n')]
            return text2
        return "None"
```

شکل ۸: تصویر ماژول حقوق مورد نظر

۱۰ استخراج نوع شغل مورد نظر

در این مازول ابتدا پیش پردازش مناسب بر روی داده انجام می‌دهیم. سپس با جستجو کلمات کلیدی مانند «نوع شغل مورد نظر»، «نوع شغل دلخواه» به دنبال مکان‌های احتمالی که در اطراف آنها نوع شغل مورد نظر کاربر نوشته شده باشد می‌گردیم و پس از آن در اطراف هر یک از آنها با جستجو کلمات کلیدی مانند «پاره وقت»، «تمام وقت» و مواردی از این دست شغل مورد نظر کاربر را پیدا می‌کنیم.

```
import re

class JobExpectedExtractor:
    def __init__(self) -> None:
        keys = [
            "نوع شغل مورد نظر",
            "شغل مورد نظر",
            "شغل دلخواه",
            "نوع شغل مورد نظر",
        ]
        values = [
            "تمام وقت",
            "تمام وقت",
            "پاره وقت",
            "کارآموزی",
        ]
        self.pattern_keys = r"({})\s*:".format("|".join(keys))
        self.pattern_values = r"({})".format("|".join(values))

    def find(self, text1):
        text = text1
        text = re.sub(r"\n+", r" ", text)
        text = re.sub(r"\s+", r" ", text)
        for match_key in re.finditer(self.pattern_keys, text):
            for match_value in re.finditer(
                self.pattern_values, text[match_key.end() : match_key.end() + 40]
            ):
                return match_value.group(1)
        return "None"
```

شکل ۹: تصویر مازول نوع شغل

۱۱ استخراج بخش‌های اضافه

برای بخش‌های اضافی، ابتدا لیستی از عناوین این بخش‌ها (شامل سوابق تحصیلی، سوابق شغلی، دستاوردها، پروژه‌ها و غیره) را در فایل keywords.txt قرار می‌دهیم. سپس، با خواندن آن‌ها شروع به بخش کردن رزومه به هر یک از آن‌ها می‌کنیم. در نهایت این بخش‌ها را تشخیص داده و خروجی می‌دهیم.

شکل ۱۰: تصویر مازول استخراج بخش‌های اضافه

برای خروجی، کافی است که لیستی از رزومه‌هایی که می‌خواهیم را در فایل `cv_extractor.py` قرار دهیم. سپس با اجرای این فایل، خروجی به صورت CSV در فایل `output.csv` قرار می‌گیرد. در زیر نمونه‌ای از خروجی را که نتیجه حاصل از اجرای استخراج کننده بر روی دو رزومه است را مشاهده می‌کنید.

شکل ۱۱: تصویری از خروجی دو رزومه