



دانشکده مهندسی کامپیوتر

درس

تمرین

استاد

گروه ۱

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 بخش‌های اضافه	۷
۹	نتایج	۸

۱ پیش پردازش

۲ استخراج نام

```
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+))({{self.last_names_reg}}){{0,1}}$"

    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] + 'x' * (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 len(matched_names[0].groups()) == 1:
            first_name = matched_names[0].groups()[0].strip()
        elif len(matched_names[0].groups()) > 9 and matched_names[0].groups()[9]:
            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()[1-4]:
            last_name = matched_names[0].groups()[1-4].strip()
        else:
            last_name = 'None'
        return full_name, first_name, last_name
```

٣

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.

```

import re
from persiantools import digits

class DateDetection:
    def __init__(self):
        self.pattern1 = r'\d{1,2}(\s*)?(?=[\s\-\./])\d{1,2}(\s*)?(?=[\s\-\./])\d{2,4}'
        self.pattern2 = r'\d{2,4}(\s*)?(?=[\s\-\./])\d{1,2}(\s*)?(?=[\s\-\./])\d{1,2}'
        self.pattern3 = r'\d{2,4}(\s*)?(?=[\s\-\./])\d{1,2}'
        self.pattern = f'{self.pattern1}|{self.pattern2}|{self.pattern3}'
        self.keywords = [
            'والتون',
            'والتون',
            'والتون'
        ]

    def match_date(self, inp):
        matches = []
        inp = digits.fa_to_en(inp)
        for matched in re.finditer(self.pattern, inp):
            start, end = matched.span()
            inp = inp[start] + '#' * (end - start) + inp[end:]
            matches.append(matched.group().strip())
        return matches

    def find_date_number(self, text):
        for keyword in self.keywords:
            if keyword not in text:
                continue
            text1 = text[text.find(keyword):]
            matched_dates = self.match_date(text1)
            if not matched_dates:
                matched_dates = self.match_date(text1[::-1])
            if not matched_dates:
                return 'Not Found'
            return matched_dates[0]
        return 'Not Found'

```

شکل ۴: تصویر ماژول استخراج تاریخ تولد

۶ استخراج شماره تلفن

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.

```

import re
from persiantools import digits

class PhoneNumberDetection:
    def __init__(self):
        self.pattern1 = r'(^|\W)((\+98|0)79\d{9})(\W|$)'
        self.pattern2 = r'(^|\W)((\+98|0)\s)79\d{9})(\W|$)'
        self.pattern3 = r'(^|\W)((\+98|0)\s)?(9\d{2})\s)?(\d{4})(\s)?(\d{3})(\W|$)'
        self.pattern4 = r'(^|\W)((\+98|0)\s)?(9\d{2})\s)?(\d{3})(\s)?(\d{4})(\W|$)'
        self.pattern5 = r'(^|\W)((\+98|0)\s)?(9\d{2})\s)?(\d{3})(\s)?(\d{4})(\W|$)'
        self.pattern = f'{self.pattern1}|{self.pattern2}|{self.pattern3}|{self.pattern4}|{self.pattern5}'

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

    def find_phone_number(self, text):
        matched_phones = self.match_phone_number(text)
        if not matched_phones:
            matched_phones = self.match_phone_number(text[::-1])
        if not matched_phones:
            return 'Not Found'
        return matched_phones[0].group().strip()

```

شکل ۵: تصویر ماژول 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",
            "علاقه‌مند",
            "علاقه مند",
            "بیکار",
            "بشغول",
            "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 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 ۷: 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 ۸: تصویری از خروجی دو رزومه