



پردازش زبان طبیعی

نیم سال دوم ۰۴-۰۳

مدرس: احسان الدین عسگری

مهلت ارسال: ۲ اردیبهشت

داده

تمرین اول

- مهلت ارسال پاسخ تا ساعت ۲۳:۵۹ روز مشخص شده است.
- **انتخاب چالش در کوئرا :** در تمرین هایی که چند چالش دارند، فقط یک نفر از هر گروه در کوئرا باید چالش مورد نظر گروه را انتخاب کند. امکان تغییر چالش تا قبل از زمان ددلاین انتخاب چالش وجود دارد. البته ذکر این نکته ضروری است که هر چالش محدودیتی برای تعداد افرادی که آن را انتخاب می کنند دارد. بنابراین در اسرع وقت برای انتخاب چالش اقدام کنید.
- **امکان ارسال با تأخیر و قوانین آن:** در طول ترم امکان ارسال با تأخیر تمرین ها بدون کسر نمره تا سقف ۱۲ روز (با سقف ۵ روز برای هر تمرین) وجود دارد. همچنین، به ازای هر روز تأخیر غیر مجاز ۱۵ درصد از نمره تمرین کسر خواهد شد.
- **قابلیت باز اجرای کامل نوت بوک ها:** توجه داشته باشید که نوت بوک های شما باید قابلیت باز اجرای ۱۰۰ درصد داشته باشند و در صورت نیاز به نصب کتابخانه یا دسترسی به یک فایل، مراحل نصب و دانلود (از یک محل عمومی) در همان نوت بوک قرار داده شود.
- **نحوه آپلود فایل های پروژه و مدل ها:** تمامی فایل های مرتبط به پروژه که حجم کمی دارند، باید به شکل فایل فشرده در CW آپلود شوند. اگر حجم یک فایل زیاد بود (مانند فایل ذخیره شده یک مدل در صورتی که بیش از ۲۰۰ مگابایت باشد)، تنها همان فایل را در یک محل عمومی (نظیر Google Drive) آپلود کرده و لینک دانلود را در نوت بوک و مستندات قرار دهید.
- **کافی است که فقط یکی از اعضای گروه پروژه را آپلود کند. اما حتماً در گزارش کار، نام همه اعضای گروه همراه با شماره دانشجویی آن ها ذکر شود.**
- **اهمیت گزارش کار:** بخشی از نمره شما به گزارش کار شما اختصاص دارد. در گزارش کار لازم نیست خط به خط کاری که کرده اید را توضیح دهید؛ بلکه باید به شکل کلی ایده تان برای حل مساله را شرح دهید. چند نمونه از خروجی های مساله را در گزارش بیاورید و بر اساس آن، رفتار برنامه تان را تحلیل کنید. در صورتی که پارامتری مانند دقت، صحت یا هر معیار دیگری خواسته شده باشد، آن ها را محاسبه کرده و در گزارش خود ارائه دهید.
- **تأثیر موارد امتیازی:** دقت داشته باشید، موارد امتیازی که در این تمرین آمده است، صرفاً بر روی امتیاز همین تمرین اثر دارد و روی نمرات تمرین های دیگر تأثیرگذار نخواهد بود.
- **نحوه پرسش سؤال و رفع ابهام:** در صورت وجود هرگونه ابهام یا مشکل، در کوئرا مطرح کنید و از ارسال پیام مستقیم به تیم تدریس خودداری نمایید.

توضیحات کلی

یکی از مهم ترین موضوعات پردازش زبان طبیعی کار بر روی داده، به خصوص دادگان مخصوص هر زبان یا دامنه های خاص می باشد. امروزه بحث دادگان مخصوص هر زبان و همچنین دادگان با کیفیت مهم شده است و با استفاده از این دادگان می توان LLM مطابق با ترجیح کاربران و دنبال کننده دستورات کاربران در عین رعایت امنیت تولید کرد. در این تمرین به این موضوع می پردازیم. هدف از این تمرین جمع آوری و کراولینگ دادگان فارسی با تمرکز بر مباحث فرهنگی می باشد به این شکل که بعد از به اتمام رسیدن تمرین یک مجموعه داده با کیفیت شامل دادگان متنی و مولتی مودال مرتبط

با فرهنگ‌های مختلف ایران و استان‌های مختلف خواهیم داشت که گامی مهم در جهت اضافه کردن ترجیحات فرهنگی ایرانی به مدل‌های زبانی می‌باشد. این تمرین به شکل **گروهی** می‌باشد و تیم شما باید ابتدا یکی از موضوعات مانند غذاها، تاریخ و غیره که توضیحات آن‌ها در ادامه آورده شده است را انتخاب کند، سپس برای هر کدام از موضوعات باید یکی از زیر بخش‌های آن را انتخاب کنید به عنوان مثال یک تیم ممکن است موضوع غذا را علاقمند و آشنا باشد همچنین به استان‌های گروه ۲ آگاهی بیشتری داشته باشد و به منابع آن دسترسی داشته باشد بنابراین موضوع تمرین آن‌ها غذاها در گروه ۲ استانی می‌شود و باید به حل این مسئله در دامنه استان‌های گروه ۲ پردازند. توجه داشته باشید که برای انتخاب موضوع تمرین یک فرم منتشر می‌شود که باید در آن علاقمندی خود را وارد کنید و براساس ظرفیت هر موضوع (حدود ۲ تا ۳ تیم) یکی از اولویت‌های شما اعلام می‌شود. پس از مشخص کردن موضوع باید منابع مختلف برای آن را گردآوری کنید که تعدادی از آن‌ها در متن تمرین به عنوان راهنمایی آورده شده است، سپس با استفاده از مطالب آموخته شده در ورکشاپ کراولینگ باید یک کد برای جمع‌آوری دادگان مورد نظر خود پیاده‌سازی کنید سپس با استفاده از موارد آموزش داده شده در درس برای استخراج اطلاعات و پیش‌پردازش، کد فریمورک خود را به شکلی گسترش دهید که اطلاعات مورد نظر را تمیز و استخراج کند. لازم به ذکر است برای هر موضوع ساختار دقیق اطلاعاتی که باید استخراج شوند آورده شده است و خروجی نهایی شما باید طبق این ساختار برای هر موضوع باشد. در نهایت باید دادگان خود را در یک چارچوب برچسب‌گذاری **labelstudio** قرار دهید و از نظر کیفیت برچسب بزنید. برای اینکار نیاز به یک سیاست برچسب‌زنی دارید که باید درون تیم مشخص کنید و این سیاست را در گزارش تمرین خود وارد کنید، سپس شروع به برچسب‌زنی دادگان از نظر کیفی کنید به این شکل که حداقل دو نفر کیفیت هر نمونه را بسنجند و در صورتی که نظرات مخالف داشتند نفر سوم داده را برچسب بزند و در صورتی که داده کیفیت مناسبی نداشت به شکل دستی درست شود. در نهایت بعد از تولید دیتاست خود باید آن را در **هاگینگ فیس** که در اختیارتان قرار می‌گیرد آپلود کنید و کدها و گزارش کار خود را در CW به همراه گزارش آماری از دادگان (مشابه کاری که در نوت‌بوک‌های درس انجام شد) آپلود کنید.

لطفاً به موارد زیر توجه فرمایید:

- **استفاده از چارچوب برچسب‌گذاری:** در این پروژه لازم است داده‌ها در یک چارچوب استاندارد برچسب زده شوند. سیستم پیشنهادی برای اینکار **labelstudio** می‌باشد. این قسمت باعث می‌شود که با فرآیند و ابزارهای برچسب‌زنی داده آشنا شوید که موضوع مهمی در صنعت و تحقیقات پردازش زبان طبیعی و به شکل کلی هوش مصنوعی می‌باشد.
- **مدل‌های زبانی بزرگ:** در صورت نیاز و زمانی که استخراج اطلاعات به شکل قاعده‌محور ممکن نبود، شما می‌توانید و حتی توصیه می‌شود که از مدل‌های زبانی بزرگ استفاده کنید، همچنین می‌توانید از این مدل‌ها در فرآیند برچسب‌زنی خود استفاده کنید. در صورتی که نیاز به استفاده از مدل‌های منبع باز LLM داشتید می‌توانید از مدل **Gemma3** استفاده کنید که نسخه‌های کم پارامتر آن قابل اجرا بر روی منابع موجود مانند کولب یا کگل می‌باشد. همچنین اگر نیاز به مدل OCR برای جمع‌آوری داده داشتید می‌توانید از مدل **Olmo** استفاده کنید. در صورتی که مشکل در اجرای این مدل‌ها داشتید از تیم تدریس کمک بگیرید.
- **کنترل کیفیت و گزارش متریک‌های برچسب‌گذاری:** پس از گردآوری داده و انجام برچسب‌گذاری، باید فرایند کیفیت‌سنجی روی داده‌ها انجام شود. در این فرآیند باید معیارهای برچسب‌زنی و توافق برچسب‌زن‌ها را با متریک ارزیابی مانند امتیاز کاپا گزارش کنید. همچنین جزئیات فرآیند برچسب‌زنی مانند برچسب‌هایی که هر داده گرفته است و همچنین داده قبل و بعد از اصلاح را در CW آپلود کنید.
- **تحویل و بارگذاری داده‌ها روی مخازن هاگینگ فیس:** تمامی داده‌ها و مدل‌های نهایی (در صورت نیاز) باید روی مخازن هاگینگ فیس بارگذاری شوند. دسترسی به هاگینگ فیس به مسئول گروه داده خواهد شد.
- یک فعالیت امتیازی برای این تمرین در نظر گرفته شده است که تاثیر مثبت بالایی خواهد داشت. این فعالیت به این صورت هست که به عنوان مثال تیم x موضوع مشاهیر هنری را انتخاب کرده است و یکی از اعضای این تیم منابع یا دادگان خوبی در خصوص غذاهای استان خاصی دارد در اینصورت می‌تواند با پرکردن فرم مخصوص این فعالیت، اطلاعات و دادگان خود را در اختیار تیم‌هایی که روی موضوع غذا کار می‌کنند قرار دهد و باعث بهبود

و گستردگی دادگان شود. هدف از این فعالیت استفاده از اطلاعات خاص فرهنگی ساکنان شهرها و استان‌های مختلف می‌باشد که حضور دانشجویان از شهرهای مختلف این امکان را فراهم می‌کند و دادگانی از این جنس که مستقیم توسط منابع هر استان معرفی شود کمک به غنای دادگان تیم‌های دیگر خواهد کرد.

برای موضوعات غذاهای محلی، جغرافیا و سایت‌های گردشگری و آداب و رسوم محلی، تقسیم‌بندی به شکل استانی می‌باشد بنابراین باید هر تیم یک گروه از استان‌ها را انتخاب کند. انتخاب دسته‌بندی استان‌ها توسط گروه‌ها در فرم مربوطه به همراه انتخاب چالش انجام خواهد شد. دسته‌بندی استان‌ها عبارتند از:

- **گروه ۱:** تهران، قزوین، مازندران، سمنان، گلستان، البرز، قم
 - **گروه ۲:** اصفهان، فارس، بوشهر، چهارمحال و بختیاری، هرمزگان، کهگیلویه و بویراحمد
 - **گروه ۳:** آذربایجان شرقی، کردستان، آذربایجان غربی، زنجان، گیلان، اردبیل
 - **گروه ۴:** لرستان، ایلام، کرمانشاه، همدان، مرکزی، خوزستان
 - **گروه ۵:** خراسان رضوی، خراسان جنوبی، خراسان شمالی، کرمان، یزد، سیستان و بلوچستان
- برای موضوع تاریخ دسته‌بندی وجود ندارد و هدف جمع‌آوری وقایع تاریخی مانند جنگ‌ها می‌باشد. برای موضوع مشاهیر دسته‌بندی براساس موارد زیر است و نیاز هست هر تیم که موضوع مشاهیر را انتخاب می‌کند یکی از دسته‌های زیر را انتخاب کند.

- مشاهیر علمی و فناوری و فلسفی
- مشاهیر ادبی
- مشاهیر هنری و موسیقی
- مشاهیر سیاسی
- مشاهیر ورزشی

در این بخش، هدف ما جمع‌آوری و **ساختاردهی اطلاعات** مربوط به غذاهای محلی ایران از منابع متنی مختلف در سطح وب است. برای این منظور، داده‌های متنی مرتبط با غذاهای استان‌ها را از اینترنت استخراج می‌کنیم (دسته‌بندی استان‌ها در بخش ابتدایی تمرین مشخص شده است). منابع پیشنهادی برای این کار شامل وبلاگ‌های آشپزی، سایت‌های آموزش آشپزی، مقالات مرتبط با گردشگری غذایی و **ویکی‌پدیا** است. این منابع اطلاعاتی درباره **نام غذاها، مواد اولیه و مراحل تهیه** ارائه می‌دهند که می‌توانند در ساختاردهی داده‌های مرتبط مفید باشند.

برای استخراج اطلاعات از این منابع، می‌توان از کتابخانه‌های BeautifulSoup، Scrapy و Selenium استفاده کرد. این ابزارها امکان پردازش و تحلیل داده‌های متنی را فراهم می‌کنند. بسته به منبع داده، ممکن است لازم باشد که **مراحل پیش‌پردازش متن** شامل حذف نویزهای متنی، یکسان‌سازی فرمت‌ها و استانداردسازی اطلاعات انجام شود. این کار باعث افزایش کیفیت داده‌های استخراج‌شده خواهد شد.

پس از اتمام مراحل پردازش داده، اطلاعات مرتبط با هر غذای محلی در یک **ساختار استاندارد** ذخیره می‌شوند. این ساختار شامل **نام غذا، استان و شهر مربوطه، مواد اولیه همراه با مقادیر آن‌ها، دستور تهیه و نوع وعده غذایی** خواهد بود. نمونه‌ای از این ساختار در قالب JSON به شکل زیر ارائه شده است:

```
1 {
2   "title": "قلیه ماهی",
3   "location": {
4     "province": "خوزستان",
5     "city": "آبادان",
6     "coordinates": {
7       "latitude": 30.3473,
8       "longitude": 48.2934
9     }
10  },
11  "ingredients": [
12    {
13      "name": "ماهی",
14      "amount": 500,
15      "unit": "گرم"
16    },
17    {
18      "name": "سبزی (گشنیز و شنبلیله)",
19      "amount": 250,
20      "unit": "گرم"
21    }
22  ],
23  "instructions": [
24    "پیاز را خرد کرده و در روغن تفت دهید تا طلایی شود.",
25    "تمبر هندی را در آب حل کرده و از صافی رد کنید، سپس به قابلمه اضافه کنید.",
26    "ماهی را به قطعات متوسط برش داده و به خورش اضافه کنید."
27  ],
28  "meal_type": [
29    "غذای اصلی",
30    "دریایی"
31  ],
32  "occasion": [
```

```

33     "ناهار",
34     "شام"
35 ],
36 "images": {
37     "تصویر نهایی": "https://example.com/images/ghalieh_mahi.jpg",
38     "مرحله ۱": "https://example.com/images/step1.jpg",
39     "مرحله ۲": "https://example.com/images/step2.jpg"
40 }
41 }

```

این ساختار امکان پردازش و تحلیل داده‌ها را برای کاربردهای مختلف، از جمله توسعه سامانه‌های معرفی غذاهای محلی و تحلیل داده‌های تغذیه‌ای فراهم می‌کند.

در این بخش، هدف ما جمع‌آوری و **ساختاردهی اطلاعات** مربوط به مهمترین رویدادهای تاریخی ایران از منابع متنی مختلف در سطح وب است که می‌تواند شامل دسته‌های زیر باشد:

- جنگ‌ها
- قراردادهای
- انقلاب‌ها
- مهمترین دستاوردهای دوره‌ها و پادشاهان
- ...

و اطلاعات مربوط به آن‌ها را با ساختار ارائه شده جمع‌آوری کنید. منابع پیشنهادی برای این کار:

- ویکی پدیا
- بلاگ‌های سایت کجارو و بیتوته
- سایت تاریخ ما
- <https://tarikh.inoor.ir/>
- <https://tarikhirani.ir/>
- <https://historydocuments.ir/>
- <https://tarikhirani.ir/>
- کتب تاریخی
- ...

برای استخراج اطلاعات از این منابع، می‌توان از کتابخانه‌های BeautifulSoup، Scrapy، Selenium برای کراول استفاده کرد. بسته به منبع داده، ممکن است لازم باشد که مراحل پیش‌پردازش متن شامل حذف نویزهای متنی، یکسان‌سازی فرمت‌ها و استانداردسازی اطلاعات انجام شود. این کار باعث افزایش کیفیت داده‌های استخراج‌شده خواهد شد.

پس از اتمام مراحل پردازش داده، اطلاعات مرتبط با هر رویداد در یک **ساختار استاندارد JSON** ذخیره می‌شوند. این ساختار شامل **منبع، نام رویداد، زمان رخداد، مکان رویداد، نتیجه، علت، طرفین، اهمیت تاریخی و توضیحات** خواهد بود. نمونه‌ای از این ساختار در قالب JSON به شکل زیر ارائه شده است:

```
1 {
2   "title": "نبرد چالدران: برخورد ایران و عثمانی",
3   "description": "نبرد چالدران میان صفویان ایران و امپراتوری عثمانی رخ داد. این جنگ به دلیل اختلافات مذهبی و "سرزمینی به وقوع پیوست و عثمانی‌ها به علت برتری تسلیحاتی، پیروز شدند، اما موفق به اشغال ایران نشدند.",
4   "period": {
5     "start_year": "۱۵۱۴",
6     "end_year": "۱۵۱۴"
7   },
8   "location": {
```

```

9     "province": "آذربایجان غربی",
10    "city": "چالدران",
11    "coordinates": {
12        "latitude": 36.2695,
13        "longitude": 59.5863
14    }
15 },
16 "causes": [
17     "رقابت مذهبی بین شیعه و سنی",
18     "اختلافات سرزمینی و توسعه طلبی عثمانی",
19 ],
20 "belligerents": [
21     {
22         "name": "صفویان",
23         "leader": "شاه اسماعیل اول"
24     },
25     {
26         "name": "امپراتوری عثمانی",
27         "leader": "سلطان سلیم اول"
28     }
29 ],
30 "result": "پیروزی عثمانی",
31 "casualties": {
32     "تلفات صفویان": "۵۰۰۰ نفر",
33     "تلفات عثمانی ها": "۲۰۰۰ نفر"
34 },
35 "impact": [
36     "کاهش نفوذ صفویان در آناتولی",
37     "اثبات برتری سلاح های گرم عثمانی"
38 ],
39 "historical_significance": "این جنگ نقطه عطفی در تاریخ ایران بود که نشان داد برتری سلاح های گرم",
40 "references": [
41     {
42         "title": "تاریخ ایران در دوران صفوی",
43         "author": "دکتر حسن نصر",
44         "year": "۱۹۹۸"
45     },
46     {
47         "title": "تاریخ نظامی عثمانی",
48         "author": "الیزابت جونز",
49         "year": "۲۰۰۵"
50     }
51 ],
52 "source": {
53     "title": "نبرد چالدران: برخورد ایران و عثمانی",
54     "author": "محمد رضایی",
55     "publication_date": "۱۵ مه ۲۰۲۳",
56     "url": "https://example.com/article"
57 }
58 }

```

در این بخش، هدف ما جمع‌آوری و **ساختاردهی اطلاعات** مربوط به مشاهیر ایران از منابع متنی مختلف در سطح وب است. برای این منظور، ابتدا یک دسته از مشاهیر ادبی، علمی، هنری و ... را انتخاب کنید و اطلاعات مربوط به آن‌ها را با ساختار ارائه شده جمع‌آوری کنید. منابع پیشنهادی برای این کار وبسایت ویکی‌پدیا، وبسایت گنجور و سایر وبسایت‌ها و وبلاگ‌های مرتبط با مشاهیر و آثار ادبی و علمی است.

برای استخراج اطلاعات از این منابع، می‌توان از کتابخانه‌های BeautifulSoup، Scrapy و Selenium استفاده کرد. این ابزارها امکان پردازش و تحلیل داده‌های متنی را فراهم می‌کنند. بسته به منبع داده، ممکن است لازم باشد که مراحل پیش‌پردازش متن شامل حذف نویزهای متنی، یکسان‌سازی فرمت‌ها و استانداردسازی اطلاعات انجام شود. این کار باعث افزایش کیفیت داده‌های استخراج‌شده خواهد شد.

پس از اتمام مراحل پردازش داده، اطلاعات مرتبط با هر شخص در یک **ساختار استاندارد JSON** ذخیره می‌شوند. این ساختار شامل **نام فرد**، **جنسیت**، **محل و تاریخ تولد و فوت و همچنین مقبره آن**، **دوره زندگی**، **مشاغل**، **آثار**، **رویدادهای مهم در طول زندگی آن فرد و تصاویر مربوطه** خواهد بود. نمونه‌ای از این ساختار در قالب JSON به شکل زیر ارائه شده است:

```
1 {
2   "name": "ابوالقاسم فردوسی توسی",
3   "sex": "مرد",
4   "nick-names": [
5     "حکیم توس",
6     "حکیم سخن"
7   ],
8   "birth": {
9     "date": "940 AD",
10    "location": {
11      "province": "خراسان",
12      "city": "توس",
13      "coordinates": {
14        "latitude": 36.2695,
15        "longitude": 59.5863
16      }
17    }
18  },
19  "death": {
20    "date": "1025 AD",
21    "location": {
22      "province": "خراسان",
23      "city": "توس",
24      "coordinates": {
25        "latitude": 36.2695,
26        "longitude": 59.5863
27      }
28    }
29  },
30  "tomb_location": {
31    "province": "خراسان رضوی",
32    "city": "توس",
33    "coordinates": {
34      "latitude": 36.2695,
```



```

34     "longitude": 59.5863
35   }
36 }
37 },
38 "era": "سامانیان، غزنویان",
39 "occupation": [
40   "شاعر",
41   "دهقان"
42 ],
43 "works": [
44   "شاهنامه"
45 ],
46 "events": [
47   {
48     "title": "سرودن شاهنامه",
49     "start_date": "974 AD",
50     "end_date": "1004 AD",
51     "location": {
52       "province": "خراسان",
53       "city": "توس",
54       "coordinates": {
55         "latitude": 36.2695,
56         "longitude": 59.5863
57       }
58     },
59     "related_people": [
60       "منشور ابو منصور"
61     ],
62     "description": "فردوسی در جوانی به مطالعه تاریخ ایران علاقه‌مند شد. وقتی دید که شاهنامه‌ی منشور...
63   }
64 ],
65 "image": {
66   "young": "https://example.com/images/ferdowsi_young.jpg",
67   "adult": "https://example.com/images/ferdowsi_adult.jpg",
68   "tomb": "https://example.com/images/ferdowsi_tomb.jpg"
69 }
70 }

```

تغییرات اقلیمی و منابع طبیعی و سایت‌های گردشگری در ایران

در این بخش، هدف ما جمع‌آوری و ساختاردهی اطلاعات مربوط به منابع طبیعی، جاذبه‌های گردشگری و تأثیرات تغییرات اقلیمی در استان‌های ایران از منابع متنی مختلف در سطح وب است. برای این منظور، ابتدا یک گروه از استان‌ها را مطابق توضیحات ابتدای تمرین انتخاب کرده و داده‌های متنی مرتبط با این موضوعات را از منابع مختلف اینترنتی استخراج کنید. منابع پیشنهادی وبسایت‌های گردشگری، کتاب‌های جغرافیای استانی، پایگاه‌های علمی، مقالات محیط‌زیستی، گزارش‌های تغییرات اقلیمی و ویکی‌پدیا می‌باشد. این منابع اطلاعاتی درباره ویژگی‌های طبیعی، جاذبه‌های گردشگری و تغییرات اقلیمی ارائه می‌دهند که می‌توانند در ساختاردهی داده‌های مرتبط مفید باشند. استخراج اطلاعات و پیش‌پردازش از کتاب‌های جغرافیای استانی باید بخشی از سیستم شما باشد.

برای استخراج اطلاعات از این منابع، می‌توان از کتابخانه‌های BeautifulSoup، Scrapy و Selenium استفاده کرد. این ابزارها امکان پردازش و تحلیل داده‌های متنی را فراهم می‌کنند. بسته به منبع داده، ممکن است لازم باشد که مراحل پیش‌پردازش متن شامل حذف نویزهای متنی، یکسان‌سازی فرمت‌ها و استانداردسازی اطلاعات انجام شود. این کار باعث افزایش کیفیت داده‌های استخراج‌شده خواهد شد.

پس از اتمام مراحل پردازش داده، اطلاعات مرتبط با هر استان در یک ساختار استاندارد ذخیره می‌شوند. این ساختار شامل نام استان، منابع طبیعی، جاذبه‌های گردشگری، و تأثیرات تغییرات اقلیمی به همراه تصاویر مختلف از آنها در طول سال‌های مختلف خواهد بود. نمونه‌ای از این ساختار در قالب JSON به شکل زیر ارائه شده است:

```
1 {
2   "title": "ویژگی‌های جغرافیایی اصفهان",
3   "location": {
4     "province": "اصفهان",
5     "city": "اصفهان",
6   },
7   "geographical_features": [
8     {
9       "name": "رودخانه‌ها",
10      "description": [
11        {"name": "زاینده‌رود", "images": []},
12        {"name": "رودخانه کوهپایه", "images": []},
13        {"name": "رودخانه چم‌آباد", "images": []},
14        {"name": "رودخانه دالانکوه", "images": []},
15        {"name": "رودخانه‌های غرب اصفهان", "images": []}
16      ]
17    },
18    {
19      "name": "کوه‌ها",
20      "description": [
21        {"name": "کوه صفه", "images": []},
22        {"name": "کوه‌های کرکس", "images": []},
23        {"name": "کوه دالانکوه", "images": []},
24        {"name": "کوه‌های زیاران", "images": []},
25        {"name": "کوه‌های آران و بیدگل", "images": []}
26      ]
27    },
28    {
29      "name": "دریاچه‌ها",
```

```

30     "description": [
31         {"name": "تالاب گاوخونی", "images": []},
32         {"name": "دریاچه نمک", "images": []}
33     ]
34 },
35 {
36     "name": "پوشش گیاهی",
37     "description": [
38         "بیابان‌های مرکزی ایران",
39         "گونه‌های مقاوم به کم‌آبی",
40         "گیاگان نادر مناطق بیابانی",
41         "پوشش گیاهی مراتع کوهستانی",
42         "گونه‌های گیاهی منطقه‌ای (پالیز و بادام وحشی)"
43     ]
44 }
45 ],
46 "topography": [
47     {"name": "نیمه‌بیابانی", "description": ["بیابانی"]},
48     {"name": "کوه صفه", "description": ["کوه‌های کرکس"]},
49     {"name": "دشت‌های وسیع", "description": ["ترکیب کوه‌های مرتفع",
50         "توپوگرافی", "مناطق بیابانی و کویری"]},
51     {"name": "گسل‌های فعال مانند گسل اصفهان-کاشان", "description": ["ویژگی‌های زمین‌شناسی",
52         "وجود گسل‌های زمین‌شناسی مهم", "منطقه کوهستانی"]},
53 ],
54 "natural_resources": [
55     {"name": "چاه‌های کشاورزی", "description": ["آب‌های زیرزمینی",
56         "رودخانه‌های فصلی", "چاه‌های عمیق", "سد زاینده‌رود"]},
57     {"name": "گچ", "description": ["سنگ آهن", "منابع معدنی",
58         "معدن سنگ مرمر", "معدن مس", "گاز طبیعی"]},
59 ],
60 "tourist_attractions": [
61     {"name": "میدان نقش جهان",
62         "images": [], "year_built": "1598",
63         "architect": "علی‌اکبر اصفهانی - محمدرضا ابن حسین بنای",
64         "description": "یکی از بزرگ‌ترین میدان‌های جهان و از آثار ثبت‌شده در یونسکو."},
65     {"name": "پل خواجه", "images": [],
66         "year_built": "1650", "architect": "نامشخص",
67         "Constructor": "شاه عباس دوم",
68         "description": "یکی از زیباترین پل‌های تاریخی اصفهان با معماری بی‌نظیر."},
69     {"name": "چهل ستون", "images": [], "year_built": "1647",
70         "architect": "شیخ بهایی", "Constructor": "شاه عباس دوم",
71         "description": "کاخ‌های با شکوه با نقاشی‌های تاریخی و معماری صفوی."},
72     {"name": "کاخ عالی‌قاپو", "images": [], "year_built": "1597",
73         "architect": "علی‌اکبر اصفهانی - محمدرضا ابن حسین بنای اصفهانی",
74         "Constructor": "شاه عباس اول",
75         "description": "کاخ تاریخی با تزیینات زیبا و نمایی فوق‌العاده از میدان نقش جهان."},
76     {"name": "مسجد شیخ لطف‌الله",
77         "images": [], "year_built": "1619", "architect": "محمدرضا اصفهانی",
78         "Constructor": "شاه عباس اول",
79         "description": "مسجدی منحصر به فرد با کاشی‌کاری‌های بی‌نظیر در میدان نقش جهان."}

```

```
80 ],
81 "additional_info": {
82   "books_source": "http://chap.sch.ir/category/%D8%AF%D9%88%D8%B1%D9%87/%D8%AF%D9
      %88%D8%B1%D9%87-%D8%A2%D9%85%D9%88%D8%B2%D8%B4-%D9%85%D8%AA%D9%88%D8%B3%D8%B7
      %D9%87/%D8%AC%D8%BA%D8%B1%D8%A7%D9%81%D9%8A%D8%A7%DE%8C-%D8%A7%D8%B3%D8%AA%D8
      %A7%D9%86-%D9%87%D8%A7"
83   }
84
85 }
```

در این بخش، هدف ما جمع‌آوری و **ساختاردهی اطلاعات** مربوط به آداب و رسوم مختص هر استان ایران، شامل مراسم مخصوص به آن منطقه در مناسبت‌های مختلف، صنایع دستی، لباس محلی و دیگر ویژگی‌های فرهنگی است. برای این منظور، ابتدا یک گروه از استان‌ها را طبق توضیحات کلی تمرین انتخاب کرده و داده‌های متنی مرتبط با این موضوعات را از منابع مختلف اینترنتی استخراج کنید. منابع پیشنهادی شامل مقالات گردشگری، وبسایت‌های فرهنگی، پژوهش‌های مردم‌شناسی و **ویکی‌پدیا** است. این منابع اطلاعاتی درباره **مراسم سنتی، لباس‌های محلی، صنایع دستی** و دیگر شاخص‌های فرهنگی ارائه می‌دهند که می‌توانند در ساختاردهی داده‌های مرتبط مفید باشند.

برای استخراج اطلاعات از این منابع، می‌توان از کتابخانه‌های BeautifulSoup، Scrapy و Selenium استفاده کرد. این ابزارها امکان پردازش و تحلیل داده‌های متنی را فراهم می‌کنند. بسته به منبع داده، ممکن است لازم باشد که **مراحل پیش‌پردازش متن** شامل حذف نویزهای متنی، یکسان‌سازی فرمت‌ها و استانداردسازی اطلاعات انجام شود. این کار باعث افزایش کیفیت داده‌های استخراج‌شده خواهد شد.

پس از اتمام مراحل پردازش داده، اطلاعات مرتبط با هر استان در یک **ساختار استاندارد** ذخیره می‌شوند. این ساختار شامل **نام استان، شهرهای مرتبط، مراسم خاص، صنایع دستی، لباس‌های محلی و سایر عناصر فرهنگی** خواهد بود. نمونه‌ای از این ساختار در قالب JSON به شکل زیر ارائه شده است:

برای مراسم‌های محلی می‌توانید از فرمت زیر استفاده کنید:

```
1 {
2   "title": "نوروز",
3   "category": "مراسم",
4   "description": "نوروز یکی از بزرگ‌ترین و قدیمی‌ترین جشن‌های ایرانی است که در آغاز بهار جشن گرفته می‌شود.",
5   "location": {
6     "province": "ایران باستان",
7     "city": "ایران باستان",
8     "coordinates": {
9       "latitude": 30.3473,
10      "longitude": 48.2934
11    }
12  },
13  "history": "نوروز از زمان‌های قدیم تا به امروز در ایران و برخی کشورهای همسایه جشن گرفته می‌شود.",
14  "elements": [
15    {
16      "name": "سبزه",
17      "description": "سبزه یکی از نمادهای نوروز است که در سفره هفت‌سین قرار می‌گیرد."
18    },
19    {
20      "name": "هفت‌سین",
21      "description": "هفت‌سین به مجموعه‌ای از هفت شیء با نام‌هایی که با حرف سین شروع می‌شود گفته می‌شود."
22    }
23  ],
24  "occasion": [
25    "نوروز",
26    "سال نو"
27  ],
28  "attire": [
29    {
30      "name": "لباس سنتی",
```

```

31     "description": "در نوروز مردم ایران معمولاً لباس‌های نو و سنتی می‌پوشند."
32   },
33 ],
34   "music": [
35     {
36       "name": "موسیقی نوروزی",
37       "instrument": "سازهای سنتی مانند تار و سه‌تار",
38       "description": "موسیقی سنتی ایرانی که در مراسم نوروز نواخته می‌شود.",
39       "audio": "https://example.com/audio/nowruz_music.mp3"
40     }
41   ],
42   "special_customs": [
43     "پذیرایی از مهمانان با شیرینی‌های مخصوص نوروز",
44     "دید و بازدید خانواده‌ها"
45   ],
46   "activities": [
47     "موسیقی زنده",
48     "رقص‌های محلی",
49     "نمایش‌های فرهنگی"
50   ],
51   "date_time": "۱ فروردین هر سال",
52   "ceremony_type": "جشنواره فرهنگی",
53 ],
54 ],
55   "images": {
56     "تصویر ۱": "https://example.com/images/nowruz1.jpg",
57     "تصویر ۲": "https://example.com/images/nowruz2.jpg",
58     "تصویر ۳": "https://example.com/images/nowruz3.jpg"
59   }
60 }
61 }

```

برای لباس‌های محلی می‌توانید از فرمت زیر استفاده کنید:

```

1 {
2   "title": "لباس محلی گیلانی",
3   "category": "لباس محلی",
4   "description": "لباس‌های سنتی مردم گیلان که در مراسم‌های مختلف، مانند عروسی‌ها پوشیده می‌شود.",
5   "materials": [
6     {
7       "name": "پارچه ابریشمی",
8       "description": "یک پارچه نرم و براق که در تولید لباس‌های محلی استفاده می‌شود."
9     },
10    {
11      "name": "نوارهای تزئینی",
12      "description": "نوارهای رنگارنگ که به عنوان تزئینات در لباس‌های گیلانی به کار می‌روند."
13    }
14  ],
15   "history": "لباس‌های محلی گیلانی در مراسمات مهم زندگی همچون عروسی‌ها و جشن‌ها مورد استفاده قرار می‌گیرد.",
16   "location": {
17     "province": "گیلان",
18     "city": "رشت",

```

```

19     "coordinates": {
20         "latitude": 30.3473,
21         "longitude": 48.2934
22     }
23 },
24 "occasion": [
25     "عروسی",
26     "جشن‌ها"
27 ],
28 "season": "تابستان و بهار",
29 "cultural_significance": "نمادهایی چون رنگ‌های خاص نشان‌دهنده‌ی وضعیت اجتماعی افراد است.",
30 "social_context": "در مراسم عروسی، لباس‌های گیلانی معمولاً نماد احترام و پذیرش مهمانان است.",
31 "gender": "مرادانه و زنانه",
32 "images": {
33     "تصویر ۱": "https://example.com/images/gilan_traditional_dress1.jpg",
34     "تصویر ۲": "https://example.com/images/gilan_traditional_dress2.jpg",
35     "تصویر ۳": "https://example.com/images/gilan_traditional_dress3.jpg",
36 }
37 }

```

برای صنایع دستی محلی می‌توانید از فرمت زیر استفاده کنید:

```

1 {
2     "title": "سفالگری لالجین",
3     "category": "صنایع دستی",
4     "description": "سفالگری یکی از هنرهای دستی معروف در شهر لالجین همدان است.",
5     "materials": [
6         {
7             "name": "خاک رس",
8             "description": "ماده اولیه برای ساخت سفال‌های لالجین که به دلیل کیفیت بالای آن استفاده می‌شود."
9         },
10        {
11            "name": "لعاب",
12            "description": "لعاب‌هایی که برای پوشش سطح سفال‌ها به کار می‌روند و باعث براق شدن آن می‌شوند."
13        }
14    ],
15    "history": "سفالگری در لالجین از دیرباز رایج بوده و امروزه یکی از مهم‌ترین صنایع دستی ایران شناخته می‌شود.",
16    "location": {
17        "province": "همدان",
18        "city": "همدان",
19        "coordinates": {
20            "latitude": 30.3473,
21            "longitude": 48.2934
22        }
23    },
24    "techniques": [
25        "چرخ‌زنی",
26        "کوزه‌گری"
27    ],
28    "production_steps": [
29        "انتخاب خاک رس",
30        "ساخت قطعات با استفاده از چرخ سفالگری",

```

```

31     "خشک کردن قطعات",
32     "لعب زدن و پخت در کوره"
33 ],
34
35     "marketplaces": [
36         "بازار همدان",
37         "نمایشگاه‌های صنایع دستی"
38     ],
39     "cultural_factors": {
40         "symbolism": "سفالگری لالچین به عنوان نمادی از هنر و ذوق مردم این منطقه شناخته می‌شود.",
41
42     },
43     "images": {
44         "تصویر ۱": "https://example.com/images/lalejin_pottery1.jpg",
45         "تصویر ۲": "https://example.com/images/lalejin_pottery2.jpg",
46         "تصویر ۳": "https://example.com/images/lalejin_pottery3.jpg" }
47 }

```