



مدرس: دکتر عسگری

پردازش زبان طبیعی

## گزارش تمرین سری یک

شماره دانشجویی: ۴۰۱۱۰۰۵۲۸

نام و نام خانوادگی: امیر ملک حسینی

### چکیده

این گزارش، فرآیند ایجاد یک مجموعه داده از دستور پخت غذاهای محلی ایرانی را تشریح می‌کند. هدف اصلی، جمع‌آوری دستور پخت‌ها از منابع آنلاین مختلف، ساختارمند کردن داده‌ها، حاشیه‌نویسی (An-notation) آن‌ها با فراداده‌های مرتبط (نوع وعده، مناسبت، کیفیت) و انجام ارزیابی کیفی بود. داده‌ها از چندین وبسایت پوشش‌دهنده استان‌های مختلف با استفاده از کتابخانه‌های پایتون مانند requests و BeautifulSoup استخراج (Scrape) شدند. داده‌های جمع‌آوری‌شده تجمیع و سپس با استفاده از پلتفرم LabelStudio حاشیه‌نویسی شدند. ارزیابی کیفی بر اساس قضاوت حاشیه‌نویس در مورد وضوح و کامل بودن دستور پخت انجام شد. ورودی‌هایی که کیفیت پایینی داشتند، فیلتر شدند. مجموعه داده نهایی شامل ۱۷۰ دستور پخت با کیفیت بالا است. تحلیل‌ها نشان داد که تعداد کل کلمات در فیلدهای مرتبط در مجموعه داده نهایی ۲۶,۵۲۹ کلمه است و میانگین طول دستورالعمل‌ها تقریباً ۱۱۹ کلمه به ازای هر دستور پخت می‌باشد. این مجموعه داده، منبعی ساختاریافته برای وظایف پردازش زبان طبیعی (NLP) مرتبط با آشپزی ایرانی فراهم می‌کند.

## ۱ مقدمه

هدف این پروژه، گردآوری یک مجموعه داده ساختاریافته از دستور پخت غذاهای سنتی ایرانی از استان‌های مختلف بود. چنین مجموعه داده‌ای می‌تواند برای تحقیقات پردازش زبان طبیعی، حفظ میراث فرهنگی و کاربردهای آشپزی ارزشمند باشد. به دلیل کمبود مجموعه داده‌های ساختاریافته و آماده در این حوزه، از رویکرد استخراج وب برای جمع‌آوری اطلاعات از منابع آنلاین متفاوت استفاده شد. داده‌های جمع‌آوری شده سپس پردازش، برای ویژگی‌های خاص حاشیه‌نویسی و از نظر کیفی فیلتر شدند. استان‌های هدف شامل موارد زیر بودند:

- اصفهان
- فارس (شیراز)
- هرمزگان
- چهارمحال و بختیاری
- خوزستان
- بوشهر
- کهگیلویه و بویراحمد

این گزارش، روش‌شناسی مورد استفاده برای جمع‌آوری داده‌ها، تجمیع، حاشیه‌نویسی، پاک‌سازی و تحلیلی از مجموعه داده نهایی را ارائه می‌دهد.

## ۲ روش‌شناسی

این فرآیند شامل چندین مرحله مجزا بود: جمع‌آوری داده از طریق استخراج وب، تجمیع داده‌ها از منابع مختلف، حاشیه‌نویسی با استفاده از LabelStudio و در نهایت، پاک‌سازی و فیلتر کردن داده‌های حاشیه‌نویسی شده.

### ۱.۲ جمع‌آوری داده‌ها

استخراج وب با استفاده از پایتون انجام شد و عمدتاً از کتابخانه requests برای دریافت صفحات وب و BeautifulSoup4 برای تجزیه (Parse) محتوای HTML استفاده گردید. به دلیل تفاوت در ساختار و نحوه ارائه محتوا، اسکرپرها (Scrapers) مخصوصی برای هر وب‌سایت هدف توسعه داده شد. چالش‌ها و تکنیک‌های کلیدی شامل موارد زیر بودند:

- **مدیریت ساختارهای متنوع:** وب‌سایت‌ها از طرح‌بندی‌های HTML متفاوتی (مانند جداول، لیست‌ها، تگ‌های عنوان خاص مثل <h2>، <h3>) برای ارائه اطلاعات دستور پخت (عنوان، مواد لازم، دستورالعمل‌ها، تصاویر) استفاده می‌کردند. اسکرپرها برای شناسایی و استخراج داده‌ها از این ساختارهای خاص، سفارشی‌سازی شدند.
- **پردازش متن فارسی:** توابع کمکی (clean\_title، convert\_persian\_numbers) برای مدیریت اعداد فارسی (تبدیل '۰'-'۹' به '۰'-'۹') و کلمات عددی فارسی (مثلاً 'پنج' به ۵) ایجاد شدند. عناوین برای حذف شماره‌گذاری یا خط تیره ابتدایی، پاک‌سازی شدند.
- **تجزیه مواد لازم:** تابعی (parse\_amount\_unit) برای تجزیه خطوط مواد لازم توسعه داده شد که سعی می‌کرد نام ماده، مقدار عددی و واحد (مانند «گرم»، «پیمانه»، «قاشق») را جدا کند. این تابع فرمت‌های مختلفی از جمله کسرها ('۱/۲')، عبارات خاص ('مقدار لازم') و کلمات عددی فارسی را مدیریت می‌کرد.

- **استخراج دستورالعمل‌ها:** دستورالعمل‌ها اغلب درون پاراگراف‌ها (<p>) یا آیتم‌های لیست (<li>) پس از عناوین خاص (مانند «طرز تهیه») یافت می‌شدند. گاهی از عبارات منظم (Regular expressions) برای شناسایی مراحل شماره‌گذاری شده استفاده می‌شد.

- **استخراج تصاویر:** URL تصاویر، به‌ویژه تصویر اصلی دستور پخت، از تگ‌های <img> استخراج می‌شد که اغلب نیازمند شناسایی بر اساس نزدیکی به عنوان دستور پخت یا کلاس‌های CSS خاص بود. از های URL پایه برای تبدیل مسیرهای نسبی تصاویر استفاده شد.

هر اسکریپت یک فایل JSON حاوی لیستی از اشیاء دستور پخت برای استان مربوطه خود تولید کرد.

## ۲.۲ تجميع داده‌ها

فایل‌های JSON منفرد تولید شده از فرآیند استخراج برای هر استان (shiraz\_foods.json, isfahan\_foods.json و غیره) در یک فایل اصلی به نام Local\_Foods.json ترکیب شدند. اسکریپت تجميع به سادگی هر فایل JSON ورودی را خوانده و محتویات آن (با فرض اینکه هر فایل حاوی لیستی از دستور پخت‌ها بود) را به یک لیست اصلی اضافه می‌کرد که سپس ذخیره می‌شد. این امر منجر به یک مجموعه داده اولیه با ۱۸۵ ورودی دستور پخت شد. نمونه‌ای از ساختار یک ورودی در فایل تجميع شده اولیه در ادامه نشان داده شده است:

```

1  {
2    "title": "بندری سوسیس",
3    "location": {
4      "province": "خوزستان",
5      "city": "اهواز",
6      "coordinates": {
7        "latitude": 318327.31,
8        "longitude": 67062.48
9      }
10   },
11   "ingredients": [],
12   "instructions": [
13     "طرز پخت در این استان به خصوص در ایران، جنوبی مناطق غذاهای معروفترین این بندری سوسیس
    این اولین مواد که آنجا از می شود استفاده گوجه فرنگی رب و ادویه پیاز، سوسیس، از آن تهیه
    و نیمه آماده غذای یک می شود، سرو ساندویچی نان با معمولاً و است دسترس در و کم غذا
    ،. دارید غذا پخت برای کمی وقت که است زمان هایی مناسب می آید حساب به فست فودی
    "مخصوص بندری سوسیس تهیه طرز"
14   ],
15   "meal_type": [
16     "اصلی",
17     "دسر",
18     "غذا پیش"
19   ],
20   "occasion": [
21     "شام",
22     "ناهار",
23     "صبحانه"
24   ],
25   "images": {
26     "نهایی تصویر": "https://api2.kojaro.com/media2021/2-5440-f3474-e6546-ec-bc19-
27     ddce96cbf82667-c46109c1067c5ba76732f9?w&1920=q80="
28   }

```

```

۲۹ {,
۳۰
۳۱ }
۳۲ "title": "ماسواں آتش",
۳۳ "location": {
۳۴ "province": "خوردستان",
۳۵ "city": "اهواز",
۳۶ "coordinates": {
۳۷ "latitude": 318327.31,
۳۸ "longitude": 67062.48
۳۹ }
۴۰ },
۴۱ "ingredients": [],
۴۲ "instructions": [
۴۳ "@maria.saadat62", اینستاگرام: عکس منبع
۴۴ "تره، گشنیز، اسفناج، سبزی، غلیظ، ساییده، کشک، پیاز، گندم، بلغور، نخود، گوشت، باں ماسواں آتش
خشک، پونه، یا داغ، نعنای باں معمولاً و شده در ست زرد چوبه و قرمز، فلفل، نمک، روغن، جعفری →
",
کردن استفاده می توان نیز سفید، لوبیا، مانند، دیگری، حبوبات، از آتش این طبخ در می شود، ترین →
۴۵ "محلی، زبان در (ماسواں) دارد، زیادی، طرفداران، خوزستان، استان در، دزفول، شهر، در، ماسواں محلی، آتش
", است (آتش به معنای اینجا در) «باں و است» باں «ماست» معنی به →
۴۶ "ماسواں آتش، تهیه، طرز"
۴۷ [,
۴۸ "meal_type": [
۴۹ "اصلی",
۵۰ "دسر",
۵۱ "غذا، پیش"
۵۲ [,
۵۳ "occasion": [
۵۴ "شام",
۵۵ "ناهار",
۵۶ "صبحانه"
۵۷ [,
۵۸ "images": {
۵۹ "https://api2.kojaro.com/media2021/3--c421dbd97-e534-b19-bb05261-a2d45b07667-c46109c1067c5ba7673ebe?w&1920=q80=": "نهایی، تصویر
→
۶۰ {
۶۱ {,

```

## ۳.۲ برچسب زنی داده‌ها

مجموعه داده تجمیع شده (Local\_Foods.json) برای حاشیه نویسی به LabelStudio (<https://labelstud.io/>) وارد شد. فرآیند حاشیه نویسی بر افزودن اطلاعات معنایی و ارزیابی کیفیت داده‌ها متمرکز بود. وظایف حاشیه نویسی:

۱. طبقه‌بندی نوع وعده (meal\_type): تخصیص برچسب‌هایی مانند 'پیش غذا'، 'غذای اصلی' یا 'دسر'.
۲. طبقه‌بندی مناسبت (occasion): تخصیص برچسب‌هایی که مناسبت‌های مناسب را نشان می‌دهند، مانند 'صبحانه'، 'ناهار' یا 'شام'. انتخاب چند گزینه مجاز بود.
۳. ارزیابی کیفیت (quality): تخصیص برچسب 'خوب' یا 'بد' بر اساس کیفیت و کامل بودن درک شده.

دستور پخت، به‌ویژه دستورالعمل‌ها و لیست مواد لازم. عواملی مانند وضوح، جریان منطقی و کامل بودن ظاهری مراحل در نظر گرفته شدند.

**سیاست حاشیه‌نویسی:** برچسب‌ها برای meal\_type و occasion بر اساس مواد لازم دستور پخت، روش تهیه، تصاویر مرتبط (در صورت وجود) و دانش عمومی آشپزی ایرانی تخصیص داده شدند. برچسب quality ذهنی بود اما هدف آن شناسایی دستور پخت‌هایی با دستورالعمل‌های نامشخص، ناقص یا بی‌معنی بود. در خصوص استان‌هایی که اطلاعات کمی از آنها در دسترس بود (مانند استان مهگیلویه و بویراحمد) این سیاست‌ها با سختگیری کمتری همراه بودند. **توافق بین حاشیه‌نویس‌ها:** حاشیه‌نویسی توسط یک نفر انجام شد. بنابراین، محاسبه معیارهای استاندارد توافق بین حاشیه‌نویس‌ها مانند کاپای کوهن یا کاپای فلایس، که سازگاری بین چندین ارزیاب را اندازه‌گیری می‌کنند، در این مرحله امکان‌پذیر نبود. قابلیت اطمینان برچسب‌ها به سازگاری و پایبندی حاشیه‌نویس واحد به سیاست حاشیه‌نویسی بستگی دارد. خروجی از LabelStudio (LabelStudio\_Output.json) شامل داده‌های اصلی به همراه حاشیه‌نویسی‌های جدید و فراداده‌های مربوط به فرآیند حاشیه‌نویسی (شناسه حاشیه‌نویس، مهرهای زمانی و غیره) بود. نمونه‌ای از چنین ورودی قبل از پاک‌سازی در زیر نشان داده شده است:

```

1  {
2    "title": "بندری سوسیس",
3    "location": {
4      "province": "خوزستان",
5      "city": "اهواز",
6      "coordinates": {
7        "latitude": 318327.31,
8        "longitude": 67062.48
9      }
10   },
11   "ingredients": [],
12   "instructions": [
13     "طرز پخت در این استان به خصوص ایران، جنوبی مناطق غذاهای معروف‌ترین بندری سوسیس
14     این اولین مواد که آنجا از می شود استفاده گوجه‌فرنگی رب و ادویه پیاز، سوسیس، از آن تهیه
15     و نیمه آماده غذای یک می شود، سرو ساندویچی نان با معمولاً و است دسترس در و کم غذا
16     ،. دارید غذا پخت برای کمی وقت که است زمان‌هایی مناسب می آید حساب به فست‌فودی
17     "مخصوص بندری سوسیس تهیه طرز"
18   ],
19   "meal_type": "اصلی غذای",
20   "occasion": {
21     "choices": [
22       "ناهار",
23       "شام"
24     ]
25   },
26   "images": {
27     "https://api2.kojaro.com/media2021/2-5440-f3474-e6546-ec-bc19-
28     ddce96cbf82667-c46109c1067c5ba76732f9?w&1920=q80="
29     "تهایی تصویر"
30   },
31   {
32     "id": 80,
33     "quality": "بد",
34     "annotator": 1,
35     "annotation_id": 80,
36     "created_at": "202505-01-T18473482.44:05:Z",
37     "updated_at": "202505-01-T21324598.15:59:Z",
38     "lead_time": 091.12
39   }
40 }

```

```

۳۳ {,
۳۴
۳۵
۳۶
۳۷ }
۳۸ "title": "ماسوا آش",
۳۹ "location": {
۴۰ "province": "خوزستان",
۴۱ "city": "اهواز",
۴۲ "coordinates": {
۴۳ "latitude": 318327.31,
۴۴ "longitude": 67062.48
۴۵ }
۴۶ },
۴۷ "ingredients": [],
۴۸ "instructions": [
۴۹ "تره، گشنیز، اسفناج، سبزی غلیظ، ساییده کشک پیاز، گندم، بلغور، نخود، گوشت، با ماسوا آش
    → خشک پونه یا داغ نعنای معمولی و شده درست زردچوبه و قرمز فلفل نمک، روغن، جعفری
    → ",
    "کردن استفاده می توان نیز سفید لوبیا مانند دیگری حبوبات از آش این طبخ در می شود تزیین
    → محلی زبان در ماسوا دارد زیادی طرفداران خوزستان، استان در دز فول شهر در ماسوا محلی آش
    → است. آش به معنای اینجا در «با و است» با «ماست» معنی به
۵۱ ],
۵۲ "meal_type": "غذا پیش",
۵۳ "occasion": "صبحانه",
۵۴ "images": {
۵۵ "https://api2.kojaro.com/media2021/3--c421dbd97-e534-b19-bb05261-
    → a2d45b07667-c46109c1067c5ba7673ebe?w&1920=q80=": "نهایی تصویر"
۵۶ },
۵۷ "id": 81,
۵۸ "quality": "بد",
۵۹ "annotator": 1,
۶۰ "annotation_id": 81,
۶۱ "created_at": "202505-01-T18349729.44:11:Z",
۶۲ "updated_at": "202505-01-T21878880.16:06:Z",
۶۳ "lead_time": 509.7
۶۴ {,

```

## ۴.۲ پاک سازی و فیلتر کردن داده ها

مرحله نهایی شامل پاک سازی داده های حاشیه نویسی شده از LabelStudio بود. یک اسکریپت پایتون برای موارد زیر استفاده شد:

۱. فیلتر بر اساس کیفیت: پیمایش ورودی ها در LabelStudio\_Output.json و حذف هر ورودی که فیلد quality آن 'بد' علامت گذاری شده بود.

۲. حذف فراداده ها: حذف فراداده های خاص حاشیه نویسی که توسط LabelStudio اضافه شده بودند (مانند id, quality, annotator, annotation\_id, created\_at, updated\_at, lead\_time) از ورودی های باقی مانده با کیفیت 'خوب'.

مجموعه داده پاک سازی و فیلتر شده حاصل، که فقط شامل دستور پخت های با کیفیت بالا با اطلاعات ضروری و

برچسب‌های اضافه‌شده meal\_type و occasion بود، با نام Filtered\_Data.json ذخیره شد. این فرآیند اندازه مجموعه داده را از ۱۸۵ به ۱۷۰ ورودی کاهش داد.

## ۳ نتایج و تحلیل

مجموعه داده نهایی (Filtered\_Data.json) شامل ۱۷۰ ورودی دستور پخت است که کیفیت خوبی دارند.

### ۱.۳ مرور کلی مجموعه داده

مجموعه داده نهایی ساختار اصلی (عنوان، مکان، مواد لازم، دستورالعمل‌ها، تصاویر) را حفظ کرده و شامل فیلدهای حاشیه‌نویسی‌شده دستی meal\_type و occasion است. نمونه‌ای از یک ورودی نهایی و پاک‌سازی‌شده در زیر نشان داده شده است:

```
1 }
2 "title": "باقالی‌شوید کوفته",
3 "location": {
4   "province": "اصفهان",
5   "city": "اصفهان",
6   "coordinates": {
7     "latitude": 6539.32,
8     "longitude": 666.51
9   }
10 }
11 "ingredients": [
12   {
13     "name": "ریحان و مرزه ترخون، شوید، گشنیز، جعفری، تره، کوفته سبزی",
14     "amount": 0.1,
15     "unit": "گرم کیلو"
16   },
17   {
18     "name": "باقلا",
19     "amount": 0.1,
20     "unit": "گرم کیلو"
21   },
22   {
23     "name": "چربی بدون گوشت",
24     "amount": 0.500,
25     "unit": "گرم"
26   },
27   {
28     "name": "برنج",
29     "amount": 0.2,
30     "unit": "پیمانه"
31   },
32   {
33     "name": "جو یا جو بلغور",
34     "amount": 0.2,
35     "unit": "غذاخوری قاشق"
36   },
37   {
```

```

۳۸     "name": "مرغ تخم",
۳۹     "amount": 0.1,
۴۰     "unit": "عدد"
۴۱     },
۴۲     ],
۴۳     "name": "دارچین و زردچوبه و فلفل و نمک",
۴۴     "amount": "لازم مقدار",
۴۵     "unit": "لازم مقدار"
۴۶     {
۴۷     [
۴۸     "instructions": ]
۴۹     "و کنید آبکش را آن در آمدن پز نیم حالت به وقتی. بپزید آن داخل را باقالا و ریخته آب قابلمه در",
۵۰     "شوند خردن کاملا تا ریخته میکس در را جو و برنج سبزی، گوشت، با برنج پیمانه یک",
۵۱     "بگیرید را باقالا ها پوست",
۵۲     "و بریزید را دارچین و زردچوبه فلفل، نمک، باقالا، برنج، بقیه شدند، ترکیب خوب مواد که زمانی",
۵۳     "دهید و رز دست با",
۵۴     "شوند یک دست تا کنید مخلوط دوباره و بشکنید را مرغ تخم",
۵۵     "بجوشد بگذارید و گذاشته ملایم حرارت روی روغن پیمانه نصف و آب پیمانه دو با را قابلمه یک",
۵۶     "کنید گرد دست با و بردارید مواد از پرتقال یک اندازه به",
۵۷     "بپزند تا بگذارید آن کف در را کوفته ها و کنید کم را شعله آمد جوش آب که آن از بعد",
۵۸     [
۵۹     "meal_type": "اصلی غذای",
۶۰     "occasion": }
۶۱     "choices": ]
۶۲     "ناهار",
۶۳     "شام"
۶۴     [
۶۵     {
۶۶     "images": }
۶۷     "https://storage.jainjas.com/storage/blog/1/original1/"
۶۸     "c5152a1fb064c1b95aaa85648894194.jpg"
۶۹     {
۷۰     {
۷۱     }
۷۲     "title": "شیرازی فالوده",
۷۳     "location": }
۷۴     "province": "فارس",
۷۵     "city": "شیراز",
۷۶     "coordinates": }
۷۷     "latitude": 5926.29,
۷۸     "longitude": 5836.52
۷۹     {
۸۰     {
۸۱     "ingredients": ]
۸۲     {
۸۳     "name": "گندم نشاسته",
۸۴     "amount": 0.300,
۸۵     "unit": "گرم"

```



```

۸۶ {
۸۷ }
۸۸ "name": "شکر",
۸۹ "amount": 0.1,
۹۰ "unit": "دسته داران فرانسوی لیوان"
۹۱ {
۹۲ }
۹۳ "name": "آب",
۹۴ "amount": 0.4,
۹۵ "unit": "دسته داران فرانسوی لیوان"
۹۶ {
۹۷ }
۹۸ "name": "گلاب",
۹۹ "amount": 0.4,
۱۰۰ "unit": "غذاخوری قاشق"
۱۰۱ {
۱۰۲ }
۱۰۳ "name": "یا آلبیمو شیرازی ترش لیمو",
۱۰۴ "amount": 0.1,
۱۰۵ "unit": "عدد"
۱۰۶ {
۱۰۷ }
۱۰۸ "name": "آلبالو شربت",
۱۰۹ "amount": "لازم مقدار",
۱۱۰ "unit": "لازم مقدار"
۱۱۱ {
۱۱۲ [
۱۱۳ "instructions": ]
۱۱۴ "کنید درست را شیرازی فالوده شربت: اول مرحله",
۱۱۵ "کنید اضافه را گلاب: دوم مرحله",
۱۱۶ "بزنید هم را فالوده شربت: سوم مرحله",
۱۱۷ "بجوشانید را آب و نشاسته: چهارم مرحله",
۱۱۸ "دهید تفت را آب و نشاسته: پنجم مرحله",
۱۱۹ "کنید آماده را نشاسته دادن فرم و سایل: ششم مرحله",
۱۲۰ "دهید فرم را فالوده نشاسته: هفتم مرحله",
۱۲۱ "بریزید شربت در را نشاسته ها: هشتم مرحله",
۱۲۲ "کنید سرو را شیرازی فالوده: آخر مرحله"
۱۲۳ [
۱۲۴ "meal_type": "دسر",
۱۲۵ "occasion": }
۱۲۶ "choices": ]
۱۲۷ "ناهار",
۱۲۸ "شام",
۱۲۹ "صبحانه"
۱۳۰ [
۱۳۱ {
۱۳۲ "images": }
۱۳۳ "فالوده: https://blog.okcs.com/wp-content/uploads/2024/07/jpg"
۱۳۴ {
۱۳۵ }

```

## ۲.۳ تحلیل کمی

یک اسکرپت برای انجام تحلیل کمی پایه روی مجموعه داده نهایی و پاک‌سازی شده (Filered\_Data.json) استفاده شد. معیارهای کلیدی به‌دست‌آمده عبارتند از:

- **تعداد کل ورودی‌ها:** ۱۷۰
- **تعداد کل کلمات (فیلدهای مرتبط):** ۲۶,۵۲۹ کلمه. این شمارش شامل کلمات فیلدهایی مانند عنوان، مکان (استان، شهر)، مواد لازم (نام، واحد)، دستورالعمل‌ها، نوع وعده و انتخاب‌های مناسب می‌شود.
- **میانگین طول دستورالعمل‌ها:** ۱۱۹/۴۶ کلمه به ازای هر دستور پخت. این میانگین بر اساس تمام ۱۷۰ دستور پخت در مجموعه داده نهایی محاسبه شد، زیرا همه آن‌ها از فیلتر کیفیت عبور کرده و حاوی دستورالعمل بودند.

این معیارها یک نمای کلی از اندازه و حجم محتوای متنی مجموعه داده ارائه می‌دهند. میانگین طول دستورالعمل‌ها نشان می‌دهد که دستور پخت‌ها به طور کلی سطح جزئیات معقولی دارند.

## ۳.۳ تحلیل کیفی

مرحله فیلتر کردن بر اساس حاشیه‌نویسی 'کیفیت'، قابلیت اطمینان مجموعه داده را با حذف ورودی‌هایی با دستورالعمل‌های ضعیف توصیف‌شده یا ناقص، به طور قابل توجهی بهبود بخشید. ۱۷۰ ورودی باقی‌مانده به طور کلی دستور پخت‌هایی با مراحل منسجم را نشان می‌دهند. با این حال، به دلیل ماهیت استخراج وب از منابع متنوع با سطوح مختلف جزئیات و سبک‌های نوشتاری، ممکن است ناسازگاری‌های بالقوه‌ای همچنان وجود داشته باشد. تجزیه مواد لازم، اگرچه کاربردی است، ممکن است منجر به برخی تغییرات در نمایش واحد یا طبقه‌بندی نادرست گاه‌به‌گاه مقدار در مقابل واحد، به‌ویژه برای اندازه‌های کمتر رایج، شده باشد. حاشیه‌نویسی دستی meal\_type و occasion اطلاعات معنایی ارزشمندی را اضافه می‌کند که همیشه به طور صریح در متون منبع وجود ندارد.

## ۴ نتیجه‌گیری

این پروژه با موفقیت یک مجموعه داده از ۱۷۰ دستور پخت غذای سنتی ایرانی از هفت استان را ایجاد کرد. این فرآیند شامل توسعه اسکرپت‌های وب هدفمند، جمع‌آوری داده‌های جمع‌آوری‌شده، حاشیه‌نویسی آن برای نوع وعده، مناسب و کیفیت با استفاده از LabelStudio و فیلتر کردن ورودی‌های با کیفیت پایین بود. مجموعه داده نهایی اطلاعات ساختاریافته‌ای شامل عنوان، مکان، مواد لازم با مقادیر/واحدهای تجزیه‌شده، دستورالعمل‌های گام‌به‌گام، تصاویر و برچسب‌های معنایی را فراهم می‌کند. تحلیل کمی نشان داد که پایه متنی قابل توجهی (بیش از ۲۶,۰۰۰ کلمه) و میانگین طول دستورالعمل تقریباً ۱۱۹ کلمه وجود دارد. در حالی که رویکرد تک‌حاشیه‌نویس ارزیابی قابلیت اطمینان برچسب را از طریق معیارهای استاندارد مانند کاپا محدود می‌کند، مرحله فیلتر کردن کیفیت، قابلیت استفاده مجموعه داده را افزایش می‌دهد. این مجموعه داده به عنوان یک منبع ارزشمند برای وظایف آتی پردازش زبان طبیعی مرتبط با تحلیل دستور پخت، تولید متن یا انفورماتیک فرهنگی متمرکز بر آشپزی ایرانی عمل می‌کند.