



## گزارش پروژه پایانی - فاز نهایی

استخراج رابطه در زبان فارسی با بررسی وابستگی‌های جهانی Seraji و PerDT

مبانی پردازش زبان و گفتار

استاد: دکتر بهروز مینایی بیدگلی

اعضای گروه:

دانیال بازمانده ۹۷۵۲۱۱۳۵

محمدحسین کریمیان ۹۷۵۲۱۴۶۸

**Style Definition:** TOC 1: Tab stops: 3.32 cm, Right + 3.96 cm, Right + 17.51 cm, Left, Leader: ...

**Style Definition:** TOC 2: Tab stops: 5.37 cm, Right + 17.51 cm, Left, Leader: ...

## فهرست مطالب

۱	هدف پروژه .....	۳
۲	مفهوم استفراخ رابطه .....	۳
۳	وابستگی‌های جهانی .....	۴
۳.۱	وابستگی جهانی <i>Seraji</i> .....	۶
۳.۲	وابستگی جهانی <i>PerDT</i> .....	۷
۳.۳	مقایسه‌ی وابستگی‌های جهانی .....	۸
۴	فریم‌ورک <i>Predpatt</i> .....	۹
۵	آماده‌سازی دیتاست .....	۱۱
۵.۱	مرحله اول: جمع‌آوری داده‌ها .....	۱۱
۵.۲	مرحله دوم: پاکسازی داده‌ها .....	۱۴
۵.۳	مرحله سوم: پردازش آماری داده‌ها .....	۱۴
۶	پیاپی‌سازی فرایند استفراخ رابطه .....	۱۶
6.1	استخراج آرگومان‌ها با استفاده از وابستگی جهانی <i>Seraji</i> .....	۱۷
6.2	استخراج آرگومان‌ها با استفاده از وابستگی جهانی <i>PerDT</i> .....	۱۹
6.3	به دست آوردن روابط و اجرای وابستگی <i>Seraji</i> روی تمامی جملات .....	۲۰

## ۱ هدف پروژه

هدف اصلی و غایی این پروژه، پیاده‌سازی یک سیستم هوشمند با استفاده از ابزارهای پردازش زبان‌های طبیعی و یادگیری ماشین، در جهت استخراج رابطه از جملات زبان فارسی است. موضوع پروژه انتقابی در واقع، ترکیبی از دو پروژه از لیست پروژه‌های داده شده است. چرا که در این پروژه، استخراج رابطه با استفاده از هر دو وابستگی جهانی Seraji و PerDT انجام شده است.

## ۲ مفهوم استخراج رابطه

در ابتدا، به توضیح و تفصیل مفهوم استخراج رابطه (Relation Extraction) در پردازش زبان‌های طبیعی می‌پردازیم.

عملیات استخراج رابطه در آموزش مدل‌های هوش مصنوعی مبتنی بر پردازش زبان طبیعی برای یادگیری روابط موجودیت‌های مختلف در یک متن برای تجزیه و تحلیل داده‌های استخراج‌شده انجام می‌شود. تکنیک‌های مختلفی برای انجام استخراج موجودیت وجود دارد، از تطبیق رشته ساده تا رویکردهای خودکار پیچیده‌تر.

منظور از استخراج رابطه مجموعه تسک‌هایی است که به استخراج روابط معنایی موجود در متون می‌پردازد. این روابط معمولاً بین دو یا چند نهاد از یک نوع خاص (مانند شخص، سازمان، مکان و...) رخ می‌دهند و در تعدادی از دسته‌بندی‌های معنایی قرار می‌گیرند. (مانند متاهل، مجرد، استخدام‌شده و...) به عبارت بهتر، استخراج رابطه وظیفه‌ی پیش‌بینی صفات و روابط را برای موجودات و انتیتی‌های درون یک جمله دارد و مولفه‌ای کلیدی برای سافت گراف‌های دانش است.

به عنوان نمونه، در زبان انگلیسی جمله‌ی «Barack Obama was born in Honolulu, Hawaii» را در نظر بگیرید. در این جمله هدف پیدا کردن رابطه‌ی «bornInCity» است. این رابطه بین دو موجودیت (Entity) «Barak Obama» و «Honolulu» وجود دارد.

البته در این پروژه ما با جملات و دادگان فارسی کار خواهیم داشت و به پیاده‌سازی استخراج رابطه در زبان فارسی پرداخته‌ایم.

به عنوان نمونه در زبان فارسی، جمله‌ی «دانیال بازمانده در شهر مشهد به دنیا آمد.» را در نظر بگیرید. هدف از تست استخراج رابطه پیدا کردن رابطه‌ی «به دنیا آمدن» است. ما انتظار داریم پس از انجام موفق عملیات، آرگومان‌های این رابطه به دست آیند. در واقع، «دانیال بازمانده» و «شهر مشهد» آرگومان‌ها و «به دنیا آمد.» رابطه‌ای است که بین این دو آرگومان وجود دارد.

به عنوان نمونه‌ی دیگر می‌توان جمله‌ی «پاریس در فرانسه است.» را بیان کرد. در این جمله آرگومان‌ها «پاریس» و «فرانسه» هستند و «در» رابطه‌ای است که بین این دو موجودیت وجود دارد. همانطور که واضح است، الزامی به فعل بودن و معنی فعل دادن برای رابطه وجود ندارد.

انواع روش‌های مختلفی برای انجام عملیات استخراج رابطه وجود دارد:

- ۱) Rule-based RE
- ۲) Weakly Supervised RE
- ۳) Supervised RE
- ۴) Distantly Supervised RE
- ۵) Unsupervised RE

هر کدام از روش‌های فوق، مزایا و معایبی دارند که تفصیل آنها از موصله‌ی این گزارش خارج است.

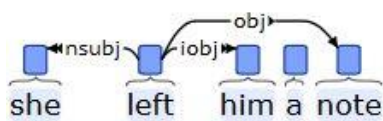
از کاربردهای مهم استخراج رابطه می‌توان به جستجوی سافت‌وار یافته، تجزیه و تحلیل احساسات (Sentiment Analysis)، سیستم‌های پرسش و پاسخ (Question Answering) و خلاصه‌سازی متن می‌توان نام برد.

### ۳ وابستگی‌های جهانی

وابستگی جهانی (UD) پروژه‌ای است که در حال توسعه‌ی ماشینی‌نویسی از نظر زبانی سازگار برای بسیاری از زبان‌ها، با هدف تسهیل توسعه‌ی تمیزه‌کننده چندزبانه، یادگیری بین زبانی و تجزیه و تحلیل تمقیقات از منظر گونه‌شناسی زبان است. فلسفه‌ی کلی این است که فهرستی جهانی از مقوله‌ها و دستورالعمل‌ها برای تسهیل ماشینی‌نویسی ثابت سافت‌وارهای مشابه در بین زبان‌ها فراهم کند.

به عبارت بهتر، وابستگی‌های جهانی (UD) چارچوبی برای ماشینه‌نویسی ثابت دستور زبان (بفش‌هایی از گفتار، ویژگی‌های صرفی و وابستگی‌های نحوی) در زبان‌های مختلف انسانی است. UD یک تلاش با بیش از ۳۰۰ مشارکت‌کننده است که نزدیک به ۲۰۰ درخت را از بیش از ۱۰۰ زبان تولید می‌کنند.

این طرح، تجزیه و تحلیل نحوی جملات را از نظر وابستگی‌های دستور زبان مربوطه تولید می‌کند. هر وابستگی از یک تابع نحوی مشخص می‌شود که روی لبه‌ی وابستگی (فلش در شکل) نشان داده می‌شود.



برای موفقیت یک UD فوب، موارد زیر باید رعایت شوند:

- در زمینه‌ی تحلیل زبانی برای هر زبان رضایت‌بفش و مناسب باشد.
- برای گونه‌شناسی زبانی فوب باشد. به‌عنوان مثال، زمینه‌ای مناسب برای ایجاد موازی زبانی بین زبان‌ها و خانواده‌های زبانی فراهم کند.
- برای ماشینه‌نویسی سریع و ثابت توسط یک ماشینه‌نویس انسانی مناسب باشد.
- به راحتی توسط یک غیر زبان‌شناس قابل درک و استفاده باشد.
- برای تجزیه کامپیوتری با دقت بالا مناسب باشد.
- به فوبی از وظایف درک زبان پایین دستی پشتیبانی کند. (استفراخ رابطه، درک مطلب، ترجمه ماشینی و...)

در زبان فارسی، وابستگی جهانی ماوی داده‌هایی از چندین بانک درختی (Treebank) است که توسط تیم‌های مختلف در زمان‌های مختلف و با ابزارهای تبدیل متفاوت ایجاد شده‌اند. در حال حاضر دو وابستگی جهانی مهم در زبان فارسی وجود دارند: Seraji و PerDT.

### ۳/۱ وابستگی جهانی Seraji

بانک درفتی وابستگی جهانی فارسی (Persian UD) نسخه تبدیل‌شده‌ی بانک درفتی وابستگی فارسی اوپسالا (UPDT) است. (سراجی، ۲۰۱۵). این بانک درفتی، طرح ماشیه‌نویسی اصلی خود را بر اساس وابستگی‌های تایپ‌شده‌ی استنفورد دارد. (دی مارنر و همکاران، ۲۰۰۶؛ د مارنر و منینگ، ۲۰۰۸)

این طرح برای زبان فارسی گسترش یافت تا روابط نموی خاص زبان را در بر گیرد که نمی‌توانست توسط طرح اولیه توسعه‌یافته برای انگلیسی پوشش داده شود. این بانک درفتی شامل ۶۰۰۰ جمله‌ی متون مکتوب با تنوع دامنه بزرگ، از نظر ژانرهای مختلف (شامل مقالات روزنامه، داستان‌ها، توضیحات فنی و اسناد در مورد فرهنگ و هنر) و نشانه‌گذاری است. تغییرات در نشانه‌گذاری به دلیل تغییرات املائی کلمات مرکب و عبارات ثابت در زبان است.

UPDT اصلی توسط مؤگان سراجی و زیر نظر جواکیم نیور و کارینا جهانی در دانشگاه اوپسالا ساخته شد. جدای از طرح ماشیه‌نویسی جهانی و قوانین کلی در UD، UD فارسی و UPDT در نشانه‌گذاری تفاوت بیشتری دارند. تمام کلمات ماوی کلیده‌ای تقسیم‌نشده (قطعات ضمائر و کوپولا) که با برچسب‌های پیچیده در UPDT ماشیه‌نویسی شده‌اند، از کلیدواژه‌ها جدا شده‌اند و برچسب‌های متمایز در UD فارسی دریافت کرده‌اند.

تبدیل UPDT به وابستگی‌های جهانی به صورت نیمه خودکار انجام شده است. در این فرآیند، از یک اسکریپت تبدیل برای معکوس کردن سر و روابط وابسته در تبدیل‌کننده مرف اضافه (prep) و مفعول یک مرف اضافه (pobj) استفاده کردیم. علاوه بر این، ما از اسکریپت‌های دیگری که برای فارسی طراحی شده‌اند، برای جدا کردن انواع کلیتیک‌ها از میزبان خود استفاده کرده‌ایم. متعاقباً قوانین مختلفی را برای بازنویسی تگ‌های درشت قسمت گفتار و برچسب‌های وابستگی اضافه کردیم.

## ۳/۲ وابستگی جهانی PerDT

بانک درختی وابستگی جهانی فارسی (PerUDT) حاصل پوشش خودکار بانک درخت وابستگی فارسی (PerDT) با اصلاحات دستی گسترده است که توسط ممدصادق رسولی، پگاه صفری، امیرسعید مولودی و علیرضا نوریان در سال ۲۰۱۳ پیاده‌سازی شده است.

بانک درختی اصلی شامل ۲۹ هزار جمله است که از متون فارسی معاصر در ژانرهای مختلف از جمله: اخبار، مقالات دانشگاهی، مقالات مجلات و داستان‌ها نمونه برداری شده است. این بانک درختی بر اساس یک طرح زبان فاص ماشیه‌نویسی شد و تبدیل خودکار آن شامل سه مرحله اصلی بود:

(۱) بازبینی توکن‌سازی (Revising Tokenization)

(۲) نگاشت POS (POS Mapping)

(۳) نگاشت وابستگی (Dependency Mapping)

در مرحله‌ی نشانه‌گذاری، به منظور جداسازی عطف‌های چند کلمه‌ای افعال ساده که به صورت یک نشانه در PerDT گروه بندی شده اند، از دستورالعمل‌های موجود برای یافتن خودکار افعال اصلی پیروی کردیم. همچنین به طور خودکار کلیدهای ضمیری را جدا کردیم.

در مرحله تبدیل POS، ما از آخرین تگر فارسی NER مبتنی بر (Taher et al., 2020) BERT با اصلاحات دستی برای گسترش فراخوان استفاده کردیم. از طریق هفت موجودیت مختلف شناسایی‌شده توسط برپسب‌گذار، از Person و Location برای علامت‌گذاری برپسب‌های PROPON استفاده کردیم. در مجموع، PerDT شامل ۴۳ رابطه‌ی نمودی بدون نگاشت مستقیم است. مروف ربط از ابتدای جمله تا انتها مرتب شده اند و مهمتر از آن، مروف اضافه به عنوان سر عبارات اضافه و افعال کمکی به عنوان سر جملات قرار می‌گیرند.

بنابراین ما ترتیب مروف ربط را از انتها به ابتدا از طریق یک اسکریپت و قواعد مناسب برای تبدیل هر نوع رابطه به نسخه UD آن به درستی مرتب کردیم. در طول کل فرآیند و در پایان هر مرحله، نتایج را بررسی کردیم و در صورت نیاز اصلاحات دستی را اعمال کردیم.

### ۳/۳ مقایسه‌ی وابستگی‌های جهانی

در تصویر زیر، جدولی از مقایسهٔ «درفت بانک وابستگی فارسی اوپسالا» با هریک از وابستگی‌های جهانی Seraji و PerDT ارائه داده است. مبنای این وابستگی‌های جهانی می‌گذاریم و نشان می‌دهد هر برچسب اجزای سنف در وابستگی جهانی چه معادلی در پروژه‌ی دانشگاه اوپسالا (Seraji) و «پروژه‌ی دادگان وابستگی زبان فارسی (PerDT)» دارد.

[illegible]



#### ۴ فریم‌ورک *Predpatt*

Predpatt مخفف عبارت Predicate-Argument Extraction from Universal Dependencies

یک کتابخانه‌ی قدرتمند توسعه‌یافته است که به عنوان چارچوبی از الگوهای استخراج گزاره (Predicate)-

آرگومان قابل توسعه، قابل تفسیر و زبان فنتی ارائه می‌شود. PredPatt عملیات تجزیه نموی عمیق پروژه‌ی وابستگی جهانی را به یک لایه معنایی گم‌عمق اولیه پیوند می‌دهد. این می‌تواند پایه‌ای برای لایه‌بندی آینده‌ای از ماشین‌نویسی‌های معنایی در بالای درخت‌های وابستگی جهانی باشد. همچنین به‌طور جداگانه می‌تواند به عنوان یک مؤلفه زبان‌شناختی به‌فوبی از سازوکار «IE Universal» در نظر گرفته شود.

این کتابخانه در سراسر زبان‌های دنیا کار می‌کند. *Predpatt* با سایر سیستم‌های برجسته‌ی Open IE مقایسه شده است و نتایج به‌دست‌آمده نشان می‌دهد که PredPatt به بهترین دقت و یادآوری دست می‌یابد.

کاربرد اصلی *PredPatt* استخراج روابط از جملات ورودی به آن است. به این صورت که تلاش می‌کند با استفاده از الگوریتم‌های قوی پیاده‌سازی‌شده، آرگومان‌ها و رابطه‌های موجود در یک جمله را به دست آورد. با استفاده از این آرگومان‌ها و رابطه‌ها می‌توان گراف معنایی (*Semantic Graph*) مرتبط با جملات را ترسیم کرد.

به عنوان نمونه می‌توان کاربرد آن را در زبان انگلیسی مشاهده کرد. همانطور که مشخص است، آرگومان‌های *a* و *b* از جمله به دست آمده است که هرکدام از آرگومان‌ها نقش *POS* مفصوص به خود را دارند.

*?a extracts ?b from ?c*

*?a: PredPatt*

*?b: predicates*

*?c: text*

*?a gave ?b to ?c*

*?a: Chris*

*?b: the book*

*?c: Pat*

به عنوان نمونه‌ی دیگر در جملات زبان فارسی داریم:

علی دایی بازیکن فوتبال بازنشسته تیم ملی ایران و باشگاه پرسپولیس است.

$a$ ? بازیکن  $b$ ? است.

$a$ ? علی دایی

$b$ ? فوتبال (بازنشسته تیم ملی ایران و باشگاه پرسپولیس)

در مثال بالا،  $a$  و  $b$  آرگومان‌های ما هستند. یعنی در این جمله، علی دایی و فوتبال آرگومان‌های رابطه ممسوب می‌شوند. هم‌پنین کلمه‌ی «بازیکن» به عنوان رابطه‌ی بین این دو آرگومان ممسوب می‌شود.

**نکته مهم:** برای تشخیص اینکه کدام کلمه به عنوان رابطه در نظر گرفته شود، ما همیشه تمام کلماتی که بین دو آرگومان می‌آیند را به عنوان رابطه‌ی آن دو آرگومان در نظر می‌گیریم. برای استناد به درستی این کار، از مقاله‌ی Farsbase bkp نوشته دکتر مینایی بهره گرفته‌ایم.

سهراب سپهری شاعر، نویسنده و نقاش اهل ایران بود.

$a$ ? شاعر بود

$a$ ? سهراب سپهری

$a$ ? نویسنده

$a$ ? سهراب سپهری

$a$ ? نقاش  $b$ ?

$a$ ? سهراب سپهری

$b$ ? اهل ایران

مطابق توضیحات قبل، دو آرگومان «سهراب سپهری» و «اهل ایران» داریم و کلمات «شاعر»، «نویسنده» و «نقاش» نیز به عنوان روابط بین این آرگومان‌ها در نظر گرفته می‌شوند. پس سه رابطه در مجموع داریم.

## ۵ آماده‌سازی دیتاست

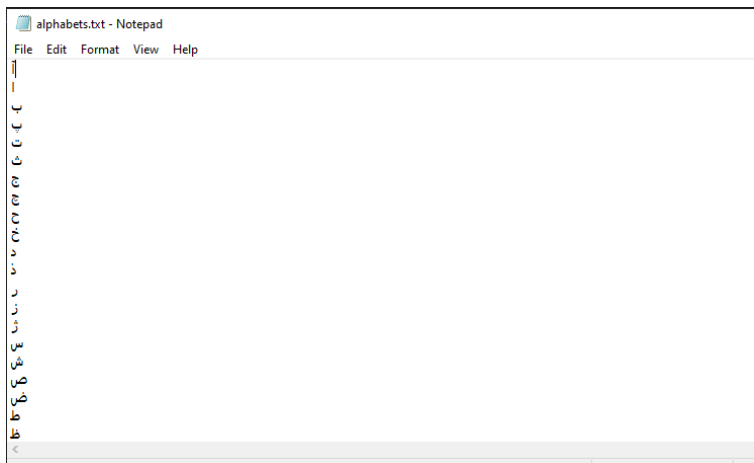
### ۵/۱ مرحله اول: جمع آوری داده‌ها

به عنوان اولین فاز، امتیاج به جمع‌آوری مقدار زیادی داده داریم. زیرا برای دستیابی به هدف نهایی، نیاز به تعداد جملات زیادی داریم تا عملیات آموزش و (آزمون را روی آنها انجام دهیم و مدل خود را بهبود دهیم. به عنوان سروس و مرجع از دیتاست سایت جامع ویکی‌پدیا استفاده می‌کنیم.

برای دستیابی به داده‌ها و جملات موجود در این وبسایت بزرگ، به Crawl کردن آنها می‌پردازیم. برای این کار، با استفاده از کتابخانه‌ی BeautifulSoup اقدام به فرز در فایل‌های html می‌کنیم و اطلاعات را مرملة به مرملة استخراج می‌کنیم. توضیحات مربوط به این بخش به طور کامل در گزارش فاز قبلی آمده است و اینجا به طور فاصله به مرامل انجام کار اشاره می‌کنیم و به همین بسنده می‌کنیم.

در گام اول، ابتدا تمامی الفباهای موجود در این وبسایت را به دست می‌آوریم:

[illegible]



در گام بعد، اقدام به دستیابی به تمام موضوعاتی که با آن مروف شروع می‌شوند، می‌کنیم. به این خاطر که داده‌های مورد نیاز ما حداکثر سی هزار جمله می باشد، فقط موضوعاتی که در صفحه اول مروف هستند را در نظر می‌گیریم و همچنین مروفی از فهرست که بیش از یک حرف دارند را هم حذف می‌کنیم.

صفحة قبلی (Zevtera) | صفحة بعد (آئومیلیستا)

- /
- 10-آ
- 10-آ ناندیروئت
- 10-آ ناندیروئت 2
- 4-آ اسکای هاوک
- 4-آ اسکای هاوک
- 4-آ اسکای هاوک
- 6-آ اینترودر
- 6-آ اینترودر
- 7-آ کرسیر
- 7-آ کرسیر
- 86929-آ
- 7-آ سائالاول
- 7-آ کس
- 7-آ لیست
- 7-آ کس
- 7-آ لیست
- 7-آ می
- 7-آ می
- 7-آ می
- 7-آ ها
- 7-آ
- 7-آ ناندیروئت
- 7-آ ناندیروئت 2
- 7-آ اسکای هاوک
- 7-آ اسکای هاوک
- 7-آ کس

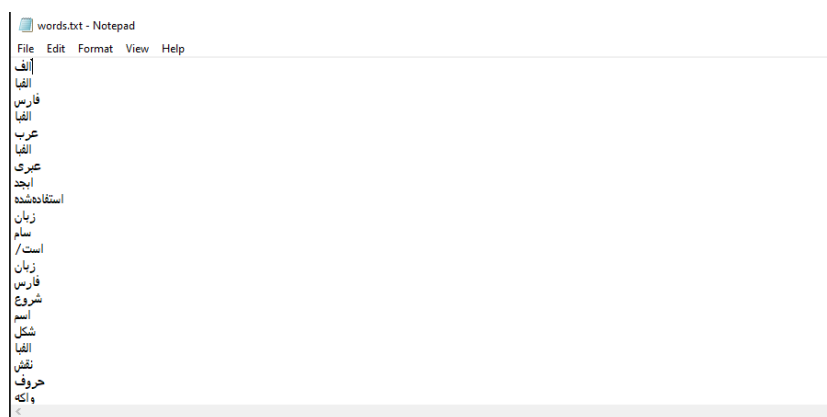
```
hrefst.txt - Notepad
File Edit Format View Help
https://fa.wikipedia.org/wiki/ا
https://fa.wikipedia.org/wiki/ا_آلبر
https://fa.wikipedia.org/wiki/ا_اس_رم
https://fa.wikipedia.org/wiki/ا_اس_پ
https://fa.wikipedia.org/wiki/ا_ب_ث_آفریقا
https://fa.wikipedia.org/wiki/ا_ب_ث_آفریقا
https://fa.wikipedia.org/wiki/ا_ب_س_آفریقا
https://fa.wikipedia.org/wiki/ا_ب_س_آفریقا
https://fa.wikipedia.org/wiki/ا_ث_میلان
https://fa.wikipedia.org/wiki/ا_ث_میلان
https://fa.wikipedia.org/wiki/ا_ج_آربری
https://fa.wikipedia.org/wiki/ا_دوبوا
https://fa.wikipedia.org/wiki/ا_والثون_لیتر
https://fa.wikipedia.org/wiki/ا_والثون_لیتر
https://fa.wikipedia.org/wiki/ا_کونهابالی
https://fa.wikipedia.org/wiki/ا_کیتین_هو
https://fa.wikipedia.org/wiki/ا_ا_گ
https://fa.wikipedia.org/wiki/ا_ارد
https://fa.wikipedia.org/wiki/ا_اس_رم
https://fa.wikipedia.org/wiki/ا_اس_ا
https://fa.wikipedia.org/wiki/ا_اس_ترگو_مورش
https://fa.wikipedia.org/wiki/ا_اس_پ
https://fa.wikipedia.org/wiki/ا_اس_رم
https://fa.wikipedia.org/wiki/ا_اس_موناکو
https://fa.wikipedia.org/wiki/ا_ث_میلان
https://fa.wikipedia.org/wiki/ا_ث_میلان_وتیر_ملی_فوتبال_ایتالیا
https://fa.wikipedia.org/wiki/ا_ث_میلان
```

مال با ایجاد یک درخواست (*request*) به هرکدام از لینک‌های موجود، اقدام به دستیابی به جملات موجود در هر صفحه می‌کنیم. برای هر صفحه یک موجودیت *json* شامل id, title و subject درست می‌کنیم.

```
"id": 1,
"title": "الف",
"text": "الف اولین حرف الفبای فارسی و الفبای عربی و الفبای عبری و سایر ابجد های استفاده شده توسط زبان های سامی است. این الفبای فارسی هم در نقش یکی از حروف واکنه بلند (مصوت بلند) است که گاه الف در زبان فارسی شروع کننده خیلی از اسم هاست به صورت (ا) در کلماتی چون: آب، باد و دریا است. و هم در نقش تلفظ همزه به کمک حرف های واکنه به کار می رود، مانند: آبز انسان و استاد (با حروف واکنه کوتاه فتحه - کسره و ضمه) و مانند: آب (الف اول در اصل=ااب)، او، این (با حروف واکنه بلند) آ یا الف مددار (به عربی: ألف ممدودة) است که آن را با همزه یا الف یکی الف در لغت نامه ها «ا» متخف اسم است الف. این نشانه با نام های مختلفی الف می شناسند، اما در حقیقت حرف جداگانه الف است. این حرف تغییر شکلیافته حرف عبری א است فرم های معین در آغاز فصل «آ» نوشته شده است: «ا» و «آ» و «ا» را در الفبای فارسی یک حرف به حساب الف می رسی شده است. از جمله آورند، اما در حقیقت دو حرف جداگانه اند...» در لغت نامه دهخدا نیز آمده است: «آ. (حرف) الف لیثنه، مقابل همزه یا الف متحرکه، حرف اول است از حروف مجاء و در حساب چثل آن را به «یک» دارند.» اما این نشانه در حقیقت ترکیبی از دو حرف از الفبای فارسی است: «ا» به عنوان صامت که شکلی از اشکال همزه است، و «آ» که مصوت بدون نامی است از مصوت های شگانه خط فارسی. اما به اشتباه در دستور خط فارسی مصوب فرهنگستان زبان و ادب فارسی، در «جدول الفبای فارسی» نشانه های اصلی، سی و سه نشانه معرفی شده است و این نشان «آ» جزو نشانه های سی و سه گانه ای که الفبای فارسی را تشکیل می دهند آورده نشده است. تمام مطالب این بخش از لغت نامه دهخدا، ماده «آ»، نقل شده است؛ مگر آن که از منبع دیگری نام برده شده باشد
```

## ۵/۲ مرحله دوم: پاکسازی داده‌ها

در این مرحله با استفاده از کتابخانه‌ی هضم که همانند *nlk* در زبان انگلیسی است، اقدام به نرمال‌سازی هر کدام از جملات می‌کنیم. سپس برای به دست آوردن کلمات، با *stemming* ریشه اصلی آن‌ها را به دست می‌آوریم و در نهایت چک می‌کنیم که این کلمات بین *stop words* های زبان فارسی هستند یا خیر و اگر در فایل *stop words* نبود، آن را در لیست کلمات در فایل *words.txt* می‌نویسیم که در شکل زیر این فایل را می‌بینیم.



### ۵/۳ مرحله سوم: پردازش آماری داده‌ها

در نهایت تمامی جملات به دست آمده را جدا می‌کنیم و در فایل *sentences.json*، همراه با آیدی و موضوع آن می‌نویسیم. از این جملات به عنوان دادگان اصلی در مهت استفاده در بخش‌های بعدی استفاده می‌کنیم.

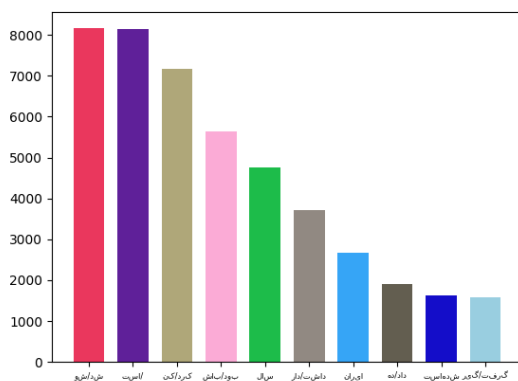
```
{
  "id": "30501",
  "text": "ی‌ما در رده «پیروان حزب» طبقه‌بندی شد (MSDAP) پدز مایرما از سال ۱۹۳۳ عضو حزب ناسیونال سوسیالیست";
  "subject": "پوزگن مایرما";
},
{
  "id": "30502",
  "text": "ن تجویز شده، عضو از سازمان جوانان (جوتگولک) شد که شاعره‌ای از سازمان شبه‌نظامی جوانان میلیتری بشمار می‌رفت";
  "subject": "پوزگن مایرما";
},
{
  "id": "30503",
  "text": "( معینان نراثر از محدودیت سنی در آن سازمان ماند و مجبور نشد در ۱۴ سالگی به سازمان جوانان میلیتری بپیوندد ";
  "subject": "پوزگن مایرما";
},
{
  "id": "30504",
  "text": "ن نیز دوباره داوطلب شده بود، خود را نا زمانی که آمریکایی‌ها منطقه را اشغال کردند از افران پلیس پنهان کرد ";
  "subject": "پوزگن مایرما";
},
}
```

تعداد جملات نهایی ۳۱۳۶۳ عدد می‌باشند که برای استفاده در وابستگی‌های جهانی مورد استفاده خواهند بود.

همچنین برای این که هر کلمه را با تعداد تکرارش داشته باشیم، یک دیکشنری تعریف می‌کنیم و هر بار که کلمه ای تکرار شد یکی به تعداد تکرار آن اضافه می‌کنیم و در نهایت آن را پس از *sort* کردن بر اساس تکرار بیشتر، در فایل *count.txt* ذخیره می‌کنیم.

```
count.txt - Notepad
File Edit Format View Help
[شد/شو: 8162
است: /8139
کرد/اش: 7173
بود/باش: 5690
سال: 4763
داشت/دار: 3722
ایران: 2665
داد/داده: 1914
شده/است: 1632
گرفت/گیر: 1587
شهر: 1194
هست: /1177
کار: 1149
کشور: 1121
انگلیسی: 1086
توانست/توان: 1052
تاریخ: 1042
گفت/گو: 1017
کتاب: 932
من: 916
زمین: 882
دست: 807
یافت/یاب: 805
رفت/رو: 766
آمد/آ: 766
جهان: 741
گروه: 734
شرکت: 722
نظر: 714
آلمان: 711
آب: 710
تولید: 698
مدت/ان: 678
دوران: 659
گشت/گردد: 657
عطی:
```

نمودار فراوانی کلمات در زیر قابل مشاهده است:



## ۶ پیاده‌سازی فرایند استخراج رابطه

در این مرحله، اقدام به پیاده‌سازی سیستمی در جهت استخراج روابط موجود در جملات می‌کنیم. به نحوی که سیستم مربوط با دریافت یک یا تعدادی از جملات و همچنین انتخاب نوع وابستگی جهانی (*Seraji* یا *PerDT*)، آرگومان‌های مربوطه و رابطه‌ی بین آرگومان‌ها را استخراج کند.

**توجه:** لازم به ذکر است در جهت توسعه‌ی یک سیستم نرم‌افزاری مطلوب در جهت راحتی مشتریان و *UI/UX* بهتر، اقدام به پیاده‌سازی این سیستم در قالب یک سیستم *client/server* کردیم و تلاش کردیم صفحه‌ای به عنوان رابط کاربری برای کاربر نمایش دهیم تا به واسطه‌ی *API* های سرور بتواند به استخراج آرگومان‌های مربوطه از جملات دلفواه خود بپردازد. پیاده‌سازی این بخش نرم‌افزاری جزو قسمت‌های پروژه نبوده است. برای بخش سرور از فریم‌ورک *Django* مربوط به زبان پایتون استفاده کرده‌ایم. با استفاده از این فریم‌ورک، یک صفحه طراحی کرده‌ایم که کاربر بتواند با ورودی متن مربوطه و انتخاب نوع وابستگی جهانی و انتخاب گزینه *Submit*، نتیجه را مشاهده کند.

### RelationExtractionInPersian

RelationExtractionInPersian is a relation extraction framework for Persian texts

Choose your desired Universal Dependency:

☒ Seraji ☐ PerDT

Enter your text in Persian:

علی دایی بازیکن فوتبال بازنشسته تیم ملی ایران و باشگاه پرسپولیس است.

SUBMIT

Conll-Format:

```
# generator = UDPipe 2, https://lindat.mff.cuni.cz/services/udpipe
# udpipe_model = persian-seraji-ud-2.6-200830
# udpipe_model_licence = CC BY-NC-SA
# newdoc
# newpar
# sent_id = 1
```

Output:

ه7 بازیکن b7 است  
ا7: علی دایی  
b7: فوتبال بازنشسته تیم ملی ایران و باشگاه پرسپولیس



همچنین اقدام به پیاده‌سازی *API* های مربوطه برای هر دو وابستگی جهانی پرداختیم که بتوان از دیوایس‌ها و برنامه‌های دیگر نیز با استفاده از این *API* ها ارتباط برقرار کرد.

برای استخراج آرگومان‌های مربوطه به پیاده‌سازی هر دو وابستگی جهانی پرداختیم:

### ۶/۱ استخراج آرگومان‌ها با استفاده از وابستگی جهانی Seraji

در این قسمت، پس از دریافت *ud\_type* و *text* از کاربر، در صورتی که کاربر وابستگی جهانی سرایی را انتخاب کرده باشد، اقدام به تولید یک فایل موقتی به نام *file.txt* می‌کنیم. تمامی جملات ورودی را در این فایل می‌ریزیم. سپس اقدام به اجرا کردن یک کامند از طریق *bash* می‌کنیم.

دلیل استفاده از این فایل موقتی به دلیل نیاز این کامند به خواندن اطلاعات از فایل دارد. (*data=@file.txt*) مقدار پارامتر *model* را برابر *persian* قرار دادیم تا از مدل‌های مربوط به زبان فارسی برای استخراج رابطه استفاده شود. از *normalized\_space* به عنوان *tokenizer* استفاده می‌کنیم و همچنین برای تجزیه‌گر (پارسر) از *api* ای که برای *udpipe* تعبیه شده است، استفاده می‌کنیم.

در صورتی که اجرای کامند فوق دچار خطا نشود، از کلید *result* در نتایج این دستور می‌توانیم فرمت *conllu* موردنظر را به دست آوریم.

#### منظور از Conllu Format چیست؟

*CoNLL* نامی معمولی برای قالب‌های *TSV* (مقادیر جدا شده با تب مانند *csv*) در *NLP* است و در واقع فرمتی خاص برای نمایش مقادیر است که در این حوزه بسیار پر کاربرد است. بسیاری از فرمت‌های مختلف *CoNLL* وجود دارد زیرا *CoNLL* هر سال یک کار مشترک متفاوت است که در کنفرانس سالانه *NLP* تصمیم‌گیری می‌شود. فرمت مشهور آن به این صورت است که:

- هر توکن در یک خط نشان داده می‌شود.
- هر جمله با یک خط خالی (*NewLine*) از جمله بعدی جدا می‌شود.
- هر ستون نشان‌دهنده یک *Annotation* است.

- هر کلمه در یک جمله دارای تعداد ستون های یکسانی است. (در برفی قالب ها: هر کلمه در مجموعه تعداد ستون های یکسانی دارد).
- و... که در این گزارش نمی‌گنجد.

در ادامه، با استفاده از تابع *load\_conllu* اقدام به خواندن فرمت *Conll* تولید شده می‌کنیم و تمامی تگ‌های *POS* به دست آمده و اطلاعات لازم را در یک لیست نگه‌داری می‌کنیم.

به عنوان مرملة نهایی و در قدم آخر، با ورودی دادن این لیست به کلاس *PredPatt* روابط مورد نظر و درفت رابطه ایجاد می‌شود. برای اینکه آرگومان‌های مربوطه را ببینیم، متود *pprint()* را از این کلاس صدا می‌زنیم.

مواردی که در بالا توضیح داده شد، در تصویر زیر قابل مشاهده است:

```
if request.method == "POST":
    ud_type, text = request.POST.get("ud_type"), request.POST.get("input_text")

    if ud_type == 'seraji':
        with open('file.txt', 'w', encoding='utf-8') as f:
            f.write(text.replace('\r', '').strip())
            f.close()

        bash_command = """curl -F data=@file.txt -F model=persian -F tokenizer=normalized_spaces -F tagger= -F
        parser= http://lindat.mff.cuni.cz/services/udpipe/api/process """
        try:
            send_order = subprocess.check_output(bash_command.split(), shell=True)
        except subprocess.CalledProcessError as e:
            os.remove("file.txt")
            raise RuntimeError("command '{}' return with error (code {}): {}".format(e.cmd, e.returncode, e.output))

        os.remove("file.txt")
        conll = json.loads(send_order)['result']
        ppatt = PredPatt([ud_parse for sent_id, ud_parse in load_conllu(conll)][0])
        return render(request, 'service.html',
            {'input_text': text, 'ud_type': ud_type, 'conll_text': conll, 'output_text': ppatt.pprint()})
```

## ۶/۲ استخراج آرگومان‌ها با استفاده از وابستگی جهانی *PerDT*

پیاپی‌سازی این قسمت در بعضی بخش‌ها مشابه وابستگی جهانی سرامی است. در ابتدا اقدام به دانلود مدل‌های شبکه‌های عصبی موازی برای زبان فارسی می‌کنیم. برای دانلود این مدل‌ها از پکیج *stanza* استفاده می‌کنیم. با استفاده از متود *Pipeline* اقدام به انتفاع شبکه عصبی پیش‌فرض در زبان فارسی می‌کنیم. مال که مدل مربوطه را به دست آوردیم، جمله ورودی را به عنوان ورودی به مدل می‌دهیم.

فروبی مدل تعدادی جملات برای ما فواهد بود که این جملات شامل تعدادی کلمات هستند. در این مرحله بر فلاف روش سرامی که فودش همه نقش‌ها را به ما می‌داد، فودمان باید یکی یکی ویژگی‌های مورد نیاز را به ازای هر کلمه به دست آوریم. ویژگی‌هایی از جمله شناسه، متن کلمه، لم، جایگاه آن در متن، *xpos*، مرف شروع و پایانی و...

تمامی ویژگی‌های مورد نیاز را فودمان دستی تولید می‌کنیم و در قالب یک رشته به همدیگر می‌پسبانییم تا مانند سرامی به *PredPatt* دهیم تا فروبی را برایمان تولید کند. در واقع، فرمت *conll* را در اینجا فودمان تولید کردیم و به *Predpatt* دادیم. بر فلاف روش سرامی که با ان کردن کامند مربوطه این کار را انجام می‌داد.

مواردی که در بالا توضیح داده شد، در تصویر زیر قابل مشاهده است:

```
elif ud_type == "perdt":
    stanza.download('fa') # This downloads the Persian models for the neural pipeline
    nlp = stanza.Pipeline('fa') # This sets up a default neural pipeline in Persian
    doc = nlp(text)
    for sentence in doc.sentences:
        res = "# text = " + text.strip()
        res = res.strip() + '\n'
        for word in sentence.words:
            res = res + str(word.id) + "\t" + str(word.text) + "\t" + str(word.lemma) + "\t" + \
                str(word.pos) + "\t" + str(word.xpos) + "\t" + str(word.feats) + "\t" + \
                str(word.head) + "\t" + str(word.deprel) + "\t" + str(word.start_char) + "\t" + \
                str(word.end_char)
            res = res + '\n'

    ppatt = PredPatt([ud_parse for sent_id, ud_parse in load_conllu(res)][0])
    return render(request, 'service.html',
                  {'input_text': text, 'ud_type': ud_type, 'conll_text': res, 'output_text': ppatt.pprint()})
```

### ۶/۳ به دست آوردن روابط و اجرای وابستگی Seraji روی تمامی جملات

در این مرحله، از *API* هایی که در مرحله ی قبل برای استخراج آرگومان ها با استفاده از وابستگی جهانی سرایی پیاده سازی کردیم، استفاده می کنیم تا هر یک از آرگومان ها را در تمامی جملاتی که داریم، به همراه رابطه ی بین این آرگومان ها به دست آوریم. در نهایت، دیتاست ممدودتری به دست می آوریم که شامل تمامی جملاتی است که حداقل یک رابطه در درون خود دارند. (ممکن است تعدادی از جملات طبق مدل *Seraji* اصلاً شامل رابطه نباشند.)

ابتدا اقدام به باز کردن فایل جملات و همچنین فایل جملات بعلاوه روابط آنها (برای نوشتن) می کنیم. سپس، با استفاده از آدرسی که برای *API* مدل سرایی تعیین کردیم، در حالتی که سرور در حال اجرا باشد، اقدام به ریکوئست (دین) می کنیم.

```
sentences_file = open('../data/sentence_broken/sentences.json', 'r', encoding='utf8')
results = open('../data/sentence_broken/sentences_with_relations.json', 'w', encoding='utf8')

seraji_request_url = 'http://127.0.0.1:8000/api/seraji'

sentences_json = json.load(sentences_file)

counter = 29000
subject_counter = 2159
error_counter = 0
dataset = []
```

در واقع، روی تک تک جملات دیتاست خود، آنها را به عنوان ورودی *input\_text* به *API* مربوطه تمت ریکوئست می دهیم. سپس، از فریمی به دست آمده با استفاده از عملیات مربوط به کار با رشته ها در پایتون، ابتدا اطمینان پیدا می کنیم که متما دو آرگومان وجود داشته باشند. (در صورت وجود تک آرگومان، رابطه تعریف نمی شود.)

سپس در صورت وجود دو آرگومان، چک می کنیم که *subject* صفحه ویکی پدیا در آرگومان اول جمله باشد و بالعکس. در صورت وجود، آرگومان دوم را *object* می نامیم و تمامی کلمات بین این دو آرگومان را به عنوان رابطه ی بین این دو آرگومان در نظر می گیریم. (پیش تر توضیح داده شد که با استناد به مقاله ی *Farsbase pbk*

دکتر مینایی این فرض را کرده ایم.)

```

for sentence_dict in sentences_json[29000:]:
    try:
        response = requests.post(seraji_request_url, json={'input_text': sentence_dict['text']})
        output_text = response.json()['output_text']

        subj_index = output_text.find('?a:')
        obj_index = output_text.find('?b:')
        if subj_index != -1 and obj_index != -1:
            subject = output_text[subj_index + 3:obj_index].strip()
            if sentence_dict['subject'] in subject or subject in sentence_dict['subject']:
                object = output_text[obj_index + 3:output_text.find('\n', obj_index + 3)]
                if 'SOMETHING' not in object and len(subject) >= 2:
                    relation = output_text[output_text.find('?a') + 2: output_text.find('?b')]

                subject = remove_extra_whitespaces(subject)
                object = remove_extra_whitespaces(object)
                relation = remove_extra_whitespaces(relation)
                dataset.append(
                    {
                        'id': subject_counter,
                        'text': sentence_dict['text'],
                        'subject': subject,
                        'object': object,
                        'relation': relation
                    }
                )
                subject_counter += 1
                # if subject_counter >= 6:
                #     break

                print('Subject counter: {}'.format(subject_counter))

            counter += 1
            print(str(counter) + " from 31363 (" , end='')
            print("%.2f" % ((counter / 31363) * 100), end='')
            print("%")

            # if counter >= 29000:
            #     break
        except:
            error_counter += 1
            print('Error #{} occurred!'.format(error_counter))

results.write(json.dumps(dataset, ensure_ascii=False))
results.close()

```

در انتها، با استفاده از تابعی که فودمان تعریف کرده‌ایم (*remove\_extra\_whitespaces*)، اقدام به حذف فواصل اضافی در ابتدا و انتهای رشته‌ها می‌کنیم و تمامی روابط به‌دست‌آمده را در قالب یک انتیتی دافل لیست *dataset* ذخیره می‌کنیم.

```
def remove_extra_whitespaces(word):  
    if len(word) >= 1:  
        if word[0] == ' ':  
            word = word[1:]  
        if word[-1] == ' ':  
            word = word[:-1]  
  
    return word
```

Formatted: Centered

در انتها اقدام به ذخیره‌ی تمامی انتیتی‌های دیتاست در درون یک جیسون وامد در فایل *sentences with relations.json* می‌کنیم.

فرمت هریک از موجودیت‌های فروری به صورت زیر خواهد بود:

```
{  
  "id": 327,  
  "text": "شنا باعث کاهش آثار زیان‌بار استرس می‌شود",  
  "subject": "شنا",  
  "object": "کاهش آثار زیان‌بار استرس",  
  "relation": "باعث"  
},
```

Formatted: Centered

```
{  
  "id": 82,  
  "text": "میلان در اولین دههٔ قرن، یکی از قدرتمندترین باشگاه‌های اروپا و جهان بود",  
  "subject": "میلان",  
  "object": "اولین دههٔ قرن",  
  "relation": "در"  
},
```

```
{
  "id": 67,
  "text": ".مجموعه شعرهای احمد شاملو (دوجلدی) در آلمان غربی منتشر می‌شود",
  "subject": ".مجموعه شعرهای احمد شاملو (دوجلدی)",
  "object": ".آلمان غرب",
  "relation": "در"
},
```

همانطور که مشاهده می‌شود، توضیحات پیشین داده‌شده در مورد نتایج به‌دست‌آمده صدق می‌کنند و می‌بینیم که اکثر جملات به درستی آرگومان‌هایشان و روابط بین آن آرگومان‌ها استخراج شده است.

**توجه:** لازم به ذکر است که برای اجرا کردن این وابستگی جهانی روی جملات، با توجه به تعداد زیاد جملات و طولانی‌شدن فرایند، این کار در ۵ مرحله مجزا انجام گرفته و در نهایت، فرم‌های با همدیگر ادغام گشته‌اند و در قالب یک فایل میسون درآمده‌اند. فایل نهایی شامل مجموعه ۲۳۹۲ روابط استخراج شده از جملات می‌باشد. طبیعتاً تعداد زیادی از جملات به دلایلی از جمله عدم وجود آرگومان، تک آرگومانی، نبودن موضوع در آرگومان اول و... حذف شده‌اند و مقادیری باقی مانده‌اند که بتوانیم در بخش‌های بعدی از آنها استفاده کنیم.

تصاویر زیر، نمایانگر مراحلی از اجرای کد مربوطه روی جملات هستند. تصویر اول پیشروی ریکوئست‌ها به ازای هر جمله را نشان می‌دهد و تصویر دوم نیز اطلاعاتی در مورد سرعت آپلود و دانلود و... در ازای هر بار دستیابی به مدل *Conllu* با استفاده از دستور اجرایی در *bash* می‌باشد. (این عملیات با صدا زدن *API* مربوط به سرویس *udpipe* انجام می‌گیرد).

```
Run: connect_objects_and_relations x
21040 from 31363 (67.09%)
21041 from 31363 (67.09%)
21042 from 31363 (67.09%)
21043 from 31363 (67.09%)
21044 from 31363 (67.10%)
21045 from 31363 (67.10%)
21046 from 31363 (67.10%)
21047 from 31363 (67.11%)
21048 from 31363 (67.11%)
21049 from 31363 (67.11%)
21050 from 31363 (67.12%)
```

Formatted: Centered

```
Terminal: Local x +
[10/Jul/2022 04:12:22] "POST /api/seraji HTTP/1.1" 200 3244
% Total    % Received % Xferd  Average Speed   Time    Time     Time  Current
           Dload  Upload  Total    Spent    Left    Speed
100 4794    0 3803 100 991 6703 1746 --:--:-- --:--:-- --:--:-- 8484
[10/Jul/2022 04:12:22] "POST /api/seraji HTTP/1.1" 200 4325
% Total    % Received % Xferd  Average Speed   Time    Time     Time  Current
           Dload  Upload  Total    Spent    Left    Speed
100 2805    0 2029 100 776 4362 1668 --:--:-- --:--:-- --:--:-- 6058
[10/Jul/2022 04:12:23] "POST /api/seraji HTTP/1.1" 200 1913
% Total    % Received % Xferd  Average Speed   Time    Time     Time  Current
           Dload  Upload  Total    Spent    Left    Speed
0 0 0 0 0 0 0 0 --:--:-- --:--:-- --:--:-- 0
```