

طراحی و پیاده‌سازی یک تحلیلگر صرفی برای زبان فارسی

وحید موایی

استاد: دکتر بهرام وزیرنژاد

دانشکده زبان‌ها و زبان‌شناسی، دانشگاه صنعتی شریف

چکیده

در این پروژه با استفاده از قوانین موجود در دستور زبان، سیستمی برای تجزیه صرفی کلمات طراحی نموده‌ایم. سعی شده است تا سیستم از ویژگی‌هایی چون سبک بودن، سریع بودن، مبتنی بر قاعده بودن، کمترین میزان استثناها را داشتن، توجه به ویژگی‌های زبان فارسی و چالش‌های این زبان برخوردار باشد. در این پروژه سعی شده است با استفاده از زبان برنامه نویسی Java، الگوریتمی برای تحلیل صرفی و ریشه یابی واژگان زبان فارسی ارائه و پیاده‌سازی شود. مطالعه و بررسی الگوریتم‌های موجود مانند الگوریتم porter و بررسی دشواری‌های نوشتاری فارسی و ناهماهنگی‌های آن و مشکلات نوشتار رایانه‌ای زبان فارسی از چالش‌هایی بود که با آن مواجه بودیم. همچنین روش‌های ریشه‌یابی شکل‌های صرف شده و مشتقات فعل‌های فارسی می‌توانند از موضوعات مورد بحث در آینده باشند. دلیل محدودیت‌های موجود برای نگهداری، مدل‌سازی و ایجاد مستندات در زبان‌های برنامه نویسی دیگر مانند Python, Ruby, C++, C، سادگی و در عین حال قدرتمندی زبان Java و وجود توابع بسیار برای پردازش متن و همچنین وجود بسته‌های نرم افزاری موجود برای NLP مانند Stanford NLP، زبان Java برای پیاده سازی پروژه انتخاب گردید.

کلید واژگان: تجزیه‌گر صرفی، صرف واژگان فارسی، تحلیلگر صرفی، پردازش خودکار متن، زبان Java.

1. مقدمه

انواع کلمه یا اقسام کلمه (parts of speech) به طور سنتی برای اشاره به طبقات دستوری کلمات به کار می‌رود و مهمترین انواع کلمه در دستور هابی که به پیروی از دستورنویسان روم و یونان باستان نوشته شده‌اند معمولاً عبارت است از اسم، ضمیر، فعل، قید، صفت، حرف اضافه، حرف ربط، صوت که اغلب می‌توان حرف تعریف و ادات را نیز به آن افزود. این تقسیم‌بندی از ویژگی‌های زبان‌های یونانی و لاتین است و به هیچ وجه جنبه عمومی و جهانی در همه زبانها ندارد و از طرف دیگر مفهوم "انواع کلمه" همانند مفهوم خود "کلمه" در دستورهای سنتی دارای ابهام است. بنابراین برای رفع این مشکل و محدود کردن اصطلاح "انواع کلمه" لازم است "واژه" (lexeme) ها را منظر قرار دهیم و صورت‌های تصریفی واحدهای صرفی را با عنوان مثلاً "صورت‌های اسم" یا "صورت‌های فعل" مطرح کنیم، با این کار تقسیم‌بندی کلمات، مبتنی بر ملاک‌های صوری خواهد بود که از عمومیت بیشتر برخوردار است.

واژه‌ها از نظر صرفی و نحوی رفتارهای خاصی در نظام زبان دارند و واژه‌هایی که دارای رفتارهای صرفی و نحوی و نیز مشخصه‌های معنایی مشابه باشند در گروه‌های جداگانه قرار می‌گیرند. هرکدام از این گروه با عنوان خاصی مانند اسم، صفت، فعل و غیره مشخص می‌شود. مثلاً "اسم" به گروهی از واژه‌ها اطلاق می‌شود که می‌تواند با تکواژ جمع تصریف شود، و یا به عبارت دیگر جمع بسته شود و از نظر نحوی نقش فاعل، مفعول و غیره ایفا کند. این مسئله درباره بقیه گروه‌های واژه‌ها نیز صادق است. از این رو می‌توان فرض کرد که واژه‌ها در واژگان ذهنی (mental lexicon) اهل زبان نیز طبقه‌بندی شده‌اند و به همین دلیل اهل زبان بدون آموزش و صرفاً با اتکا به توانش زبانی خود می‌توانند واژه‌های نوساخته را در طبقات خاص قرار دهند و متناسب با ویژگی صرفی و نحوی آنها را به کار برند. در ادامه به اختصار به توصیف ساختار تصریفی هرکدام از طبقات مورد اشاره در بالا می‌پردازیم. لازم به ذکر است که برخی از اقسام کلمه مانند حروف ربط صورت تصریفی ندارند و به همین دلیل به آنها نخواهیم پرداخت.

2. ساختار تصریفی انواع کلمه در زبان فارسی

2.1. ساختار تصریفی اسم

$$[\text{اسم}] + [(\text{تکواژ جمع})] + \left[\begin{array}{c} (\text{یای نکره}) \\ (\text{یای بند موصولی}) \\ (\text{واژه‌بست‌های شخصی / ضمائر متصل}) \\ (\text{کسره اضافه}) \end{array} \right] + [(\text{واژه‌بست‌های ربطی})]$$

2.2. ساختار تصریفی صفت

$$[صفت] + \left[\begin{array}{l} \text{(تکواژ صفت تفضیلی ساز)} \\ \text{(تکواژ صفت عالی ساز)} \end{array} \right] + \left[\begin{array}{l} \text{(واژه‌بست‌های ربطی)} \end{array} \right]$$

2.3. ساختار تصریفی فعل

$$\left[\begin{array}{l} \text{(تکواژ وجه امری و التزامی)} \\ \text{(تکواژ استمراری)} \end{array} \right] + [فعل] + \left[\begin{array}{l} \text{(تکواژ نمود کامل)} \\ \text{(شناسه‌ها)} \end{array} \right] + \left[\begin{array}{l} \text{(واژه‌بست‌های شخصی)} \end{array} \right]$$

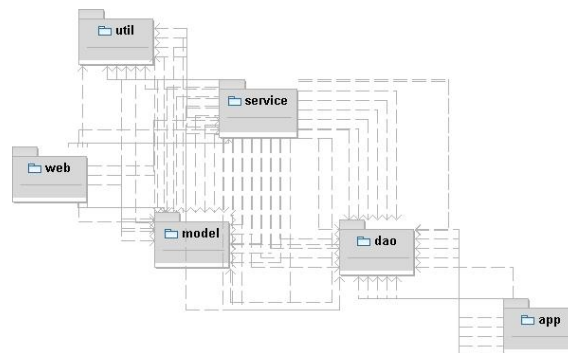
2.4. ساختار تصریفی قید

$$[قید] + \left[\begin{array}{l} \text{(تکواژ قید تفضیلی ساز)} \\ \text{(-tar/)} \end{array} \right]$$

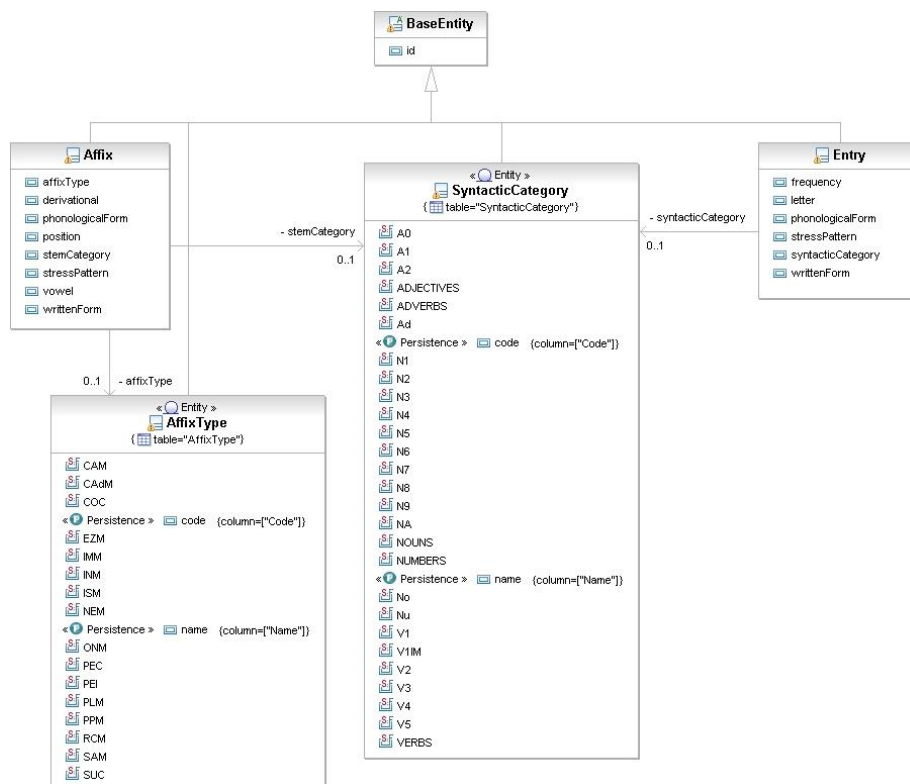
3. ساختار کلی برنامه

با استفاده از مبانی نظری و زیرساخت‌های تئوریک که در بالا به اختصار به آن اشاره شد (و توضیح بیشتر در باره آن، خود یک مقاله جداگانه خواهد بود)، یک برنامه مبتنی بر وب و با معماری چندلایه برای پیاده‌سازی تحلیلگر صرفی نوشته شد. در لایه مدل برنامه از Hibernate و 3.0 و Annotation‌های معرفی شده در آن برای دسترسی به داده‌ها در پایگاه داده استفاده می‌شود. پایگاه داده مبتنی بر MySQL® 5.0 است که داده‌ها در آن پایگاه داده Flexicon پر شده است. در لایه وب از Struts 2.0 و Jsp و برای زیرساخت برنامه هم از Spring 2.5 استفاده شده است. به جای توضیح زبانی الگوریتم‌های استفاده شده، نمودارهای UML آنها از جمله Class Diagram ها و Sequence Diagram ها با استفاده از ابزار Eclipse Europa® و Omondo® با استفاده از روش مهندسی معکوس، از روی برنامه ساخته شده‌اند. بنابراین این نمودارها را در ادامه می‌آوریم و از توضیح بیشتر پرهیز می‌کنیم.

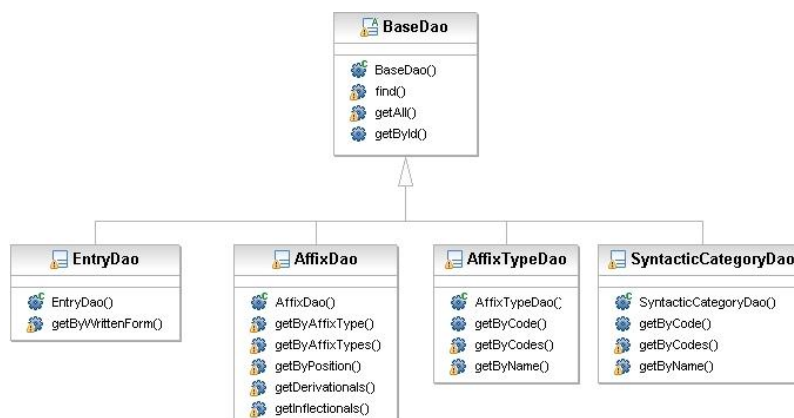
3.1. ساختار بسته‌های برنامه



3.2. ساختار کلاس‌های لایه مدل



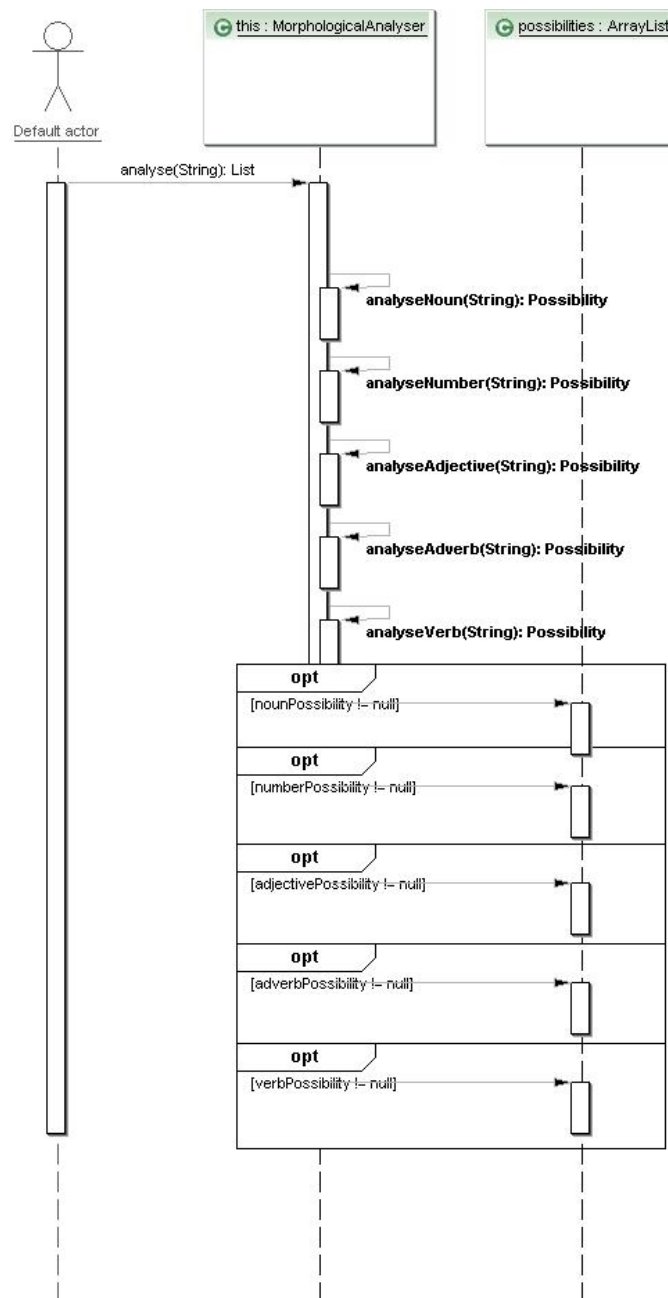
3.3. ساختار کلاس‌های لایه DAO



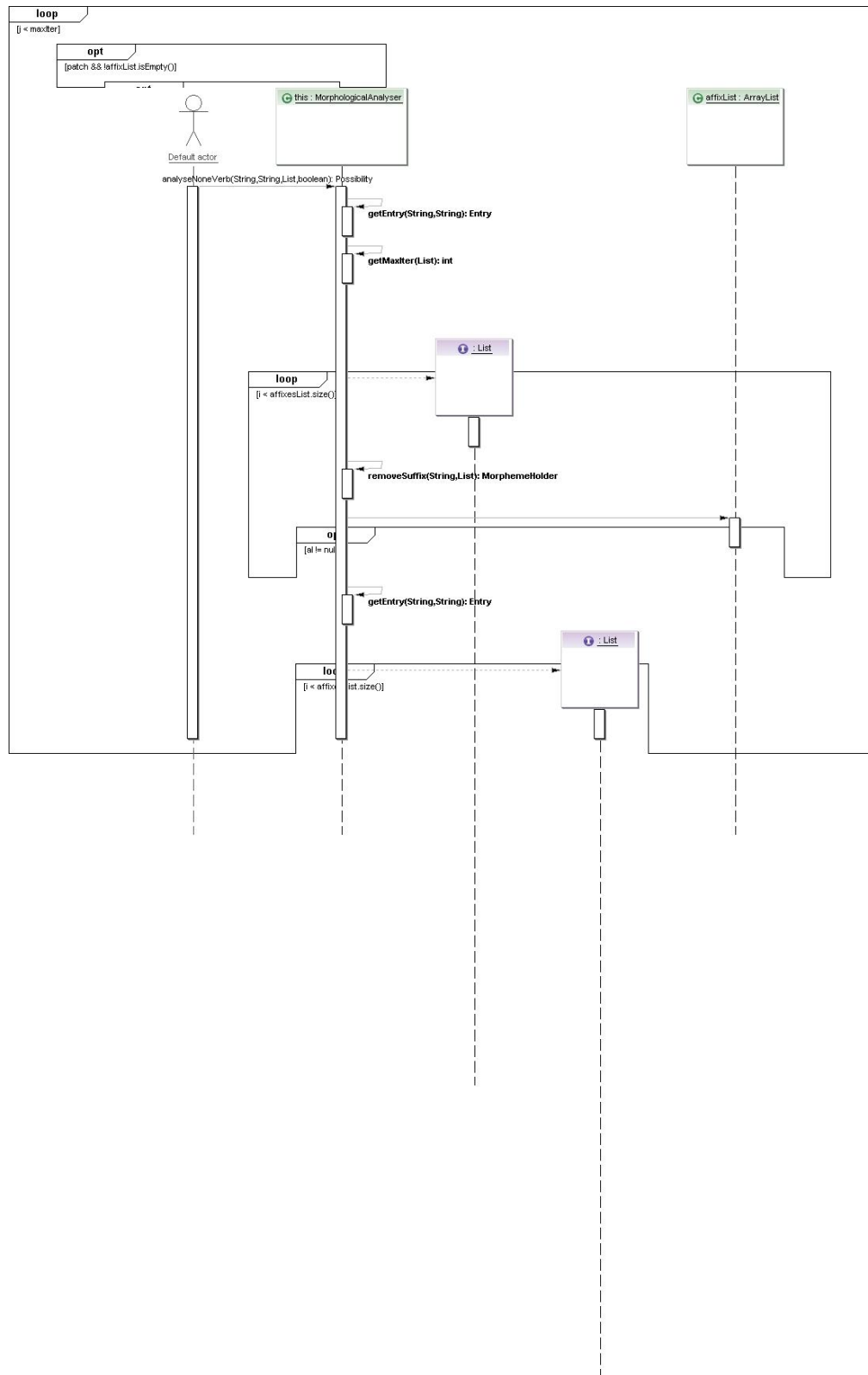
3.4. ساختار کلاس‌های لایه سرویس



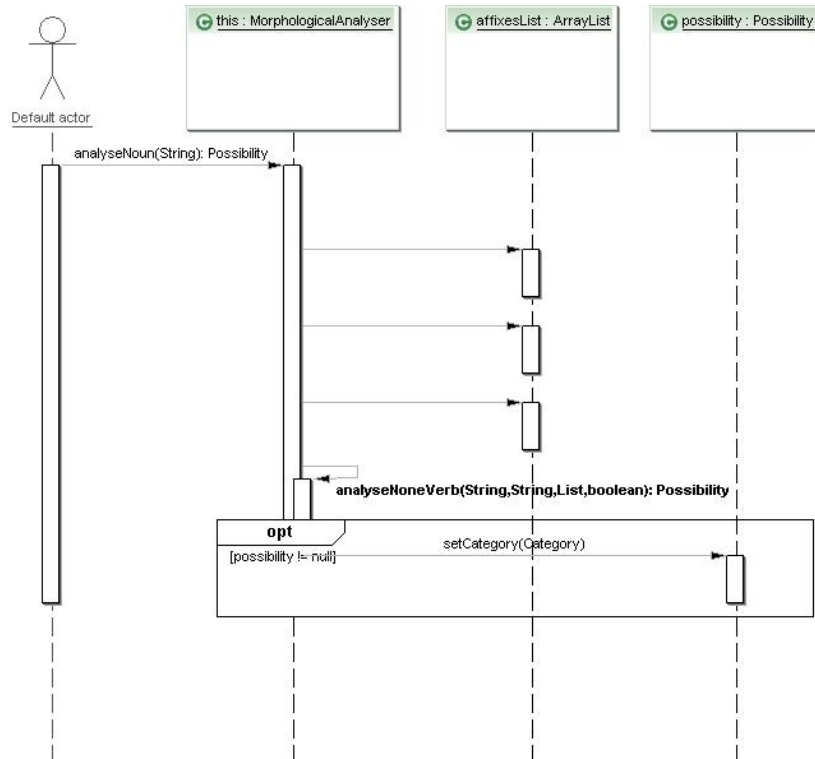
3.5. نمودار الگوریتم آنالیز کلی



3.6. نمودار الگوریتم آنالیز انواع کلمه غیر از فعل



3.7. نمودار الگوریتم آنالیز اسم



3.8. نموذج الگوریتم آنالیز فعل

