

اللَّهُمَّ صَلِّ وَسَلِّمْ وَبَارِكْ عَلَى سَيِّدِنَا مُحَمَّدٍ



محققین:

حمیدرضا لطفی نیا
محبوبه صید محمد خانی

استاد گرامی:

جناب آقای دکتر عبدالهی

موضوع پایان نامه:

پردازش زبان طبیعی (NLP)

نام دانشگاه:

علمی کاربردی جهاد دانشگاهی مشهد

دوره تحصیلی:

کاردانی

رشته تحصیلی:

فناوری اطلاعات و ارتباطات

سال

تحصیلی: 1394

انسان‌ها از توانایی‌های بصری و زبانی خود برای مشاهده اطراف خود و ایجاد ارتباط با یکدیگر استفاده می‌کنند. آن‌ها با داشتن یک تصویر براحتی می‌توانند یک توصیف زبانی از آن را ارائه داده و بطور مشابه وقتی چیزی بصورت زبانی توصیف شود، می‌توانند تصویری از آن را مجسم کنند. برای انسان‌ها این یک عمل طبیعی است، ولی داشتن این توانایی‌ها در ماشین، یک وظیفه چالش‌برانگیز است. آن به دانش حوزه‌های تحقیقاتی مانند بینایی کامپیوتر و پردازش زبان نیاز دارد. اگرچه، هر دو رشته فوق از هوش مصنوعی گرفته شده و امروزه حوزه‌های تحقیقاتی فعالی را در خود دارند، کار کردن بصورت ایزوله و بدون ارتباط با حوزه دیگر، برای محققان این دو حوزه یک تهدید است. تحقیقات انجام شده در یکی از این دو رشته به نظر می‌رسد که کارایی مناسبی در رشته دیگر نداشته باشد. هر دوی این رشته‌ها بصورت مجزا رونق پیدا کرده و اطلاعات بصری و متنی را جدا از هم استخراج می‌کنند. اما، هنوز بسیاری از مسائل دنیای واقعی وجود دارند که نیازمند دانش و تخصصی از هر دو رشته هستند، چراکه حضور همزمان داده‌های متنی و بصری، امری طبیعی و چشم‌گیر بوده و براحتی در دسترس می‌باشد؛ برای مثال، زیرنویس در ویدیوها، تصاویر تگ شده در سایت‌های شبکه‌های اجتماعی و غیره. همچنین رشد فاحش داده‌های بصری و متنی در وب و انبارداده‌های خصوصی، نیاز به جست‌وجو، ساماندهی و استخراج این داده‌ها به منظورهای مختلف را ایجاد کرده است. بنابراین، کاوش برای امکان‌سنجی و این که این ادغام دو رشته چگونه می‌تواند انجام شده کاربردهای حاصل شده کجا می‌توانند بکار رفته و مفید باشند، امری مهم و اجباری است.

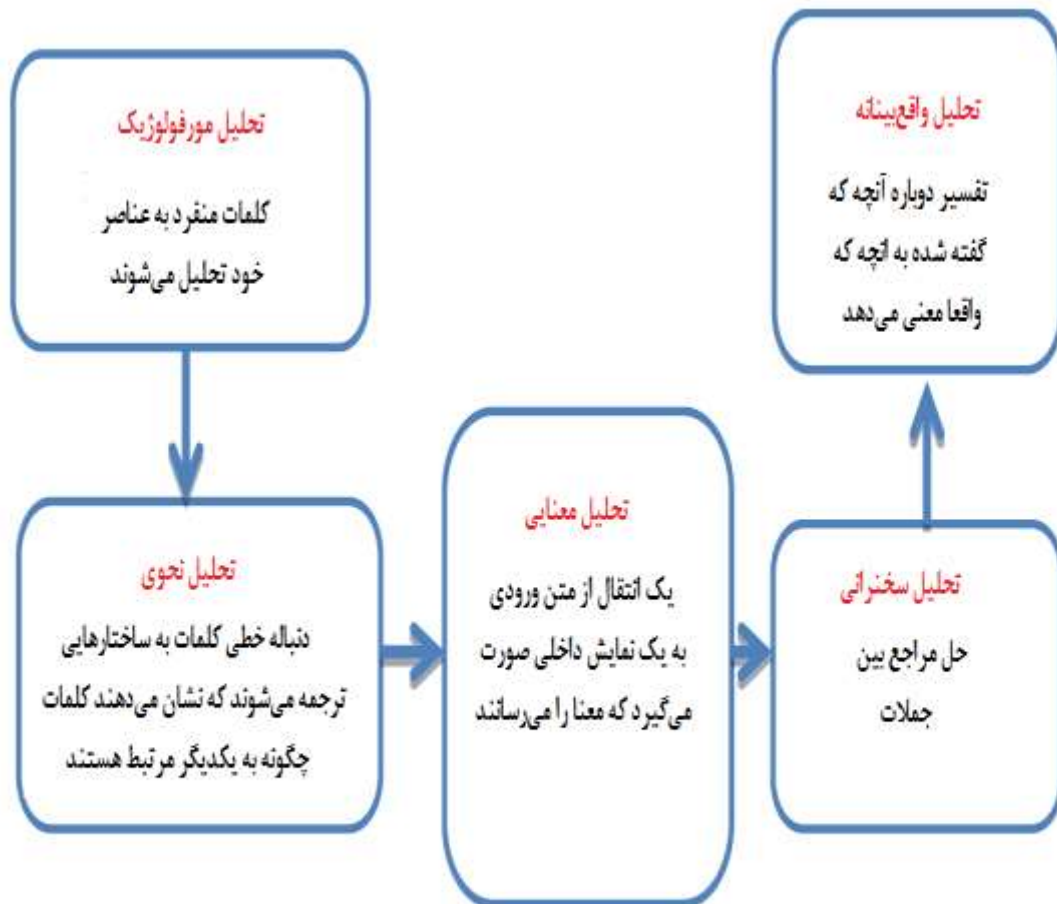
7	پردازش زبان طبیعی
8	محتویات زبان طبیعی
8	تاریخچه زبان طبیعی
8	محدودیت ها
9	موانع اساسی
9	پردازش زبان طبیعی آمارگرا
9	کارکرد های مهم پردازش زبان طبیعی
9	خلاصه سازی
10	محتویات خلاصه سازی
10	معنی خلاصه سازی
10	اهمیت خلاصه سازی
10	سیستم خلاصه سازی چیست؟
11	انواع خلاصه سازی
11	انواع مخاطب و متن
12	انواع اطلاعات ورودی
12	سبک متن
12	نوع کاربر
12	تاریخچه خلاصه سازی خودکار
13	مدل سازی اطلاعات

14	مدل دودویی
15	مدل برداری
15	مدل احتمالی
16	تعیین میزان ربط هر سند
16	تفاوت بازیابی داده و بازیابی اطلاعات
16	معیارهای ارزیابی
16	محتویات معیارهای ارزیابی
18	فرآیند ترجمه
18	روش ها
19	روش قاعده مندی
19	روش بین زبانی
19	روش مبتنی بر فرهنگ لغت
20	روش آماری
20	روش مبتنی بر مثال
20	ترجمه ماشینی پیوندی
21	مسئله های اصلی
21	ابهام زدایی
21	گفتارهای غیر استاندارد
22	واحدهای اسمی
22	نرم افزارهای کاربردی

23	ارزیابی ترجمه خودکار در حوزه های مختلف
24	ویراستاری
24	محتویات ویراستاری
24	کار ویراستاری
24	انواع ویراستاری
25	ریشه واژه و مترادفها
25	نتیجه گیری
26	منابع

پردازش زبان طبیعی

حوزه دیگری که در نظر گرفته شده، NLP است. اگر خیلی هسته‌ای نگاه کنیم، NLP به معنای توانمند ساختن کامپیوترها برای استخراج معنا از صحبت‌ها یا متن نوشتاری موجود در یک ورودی زبان طبیعی است. همچنین طیف NLP شامل تولید جملاتی به زبان طبیعی توسط کامپیوترها است که شبیه به جملاتی باشند که ما انسان‌ها بکار می‌بندیم. شکل 2 لسیتی از وظایف را که در یک کاربرد معمول از زبان طبیعی انجام می‌شوند، نشان می‌دهد.



شکل ۲: مراحل پردازش زبان طبیعی (NLP)

محتویات زبان طبیعی

۱. تاریخچه

۲. محدودیت‌ها

۳. موانع اساسی

۴. پردازش زبان‌های طبیعی آمارگرا

۵. کارکردهای مهم پردازش زبان‌های طبیعی

تاریخچه زبان طبیعی

به طور کلی تاریخچه پردازش زبان طبیعی از دهه ۱۹۵۰ میلادی شروع می‌شود. در ۱۹۵۰ آلن تورینگ مقاله معروف خود را که امروزه به عنوان ملاک هوشمندی شناخته می‌شود، منتشر ساخت.

نخستین تلاش‌ها برای ترجمه توسط رایانه ناموفق بودند، به طوری که ناامیدی بنگاه‌های تأمین بودجه پژوهش از این حوزه را نیز در پی داشتند. پس از اولین تلاش‌ها آشکار شد که پیچیدگی زبان بسیار بیشتر از چیزی است که پژوهشگران در ابتدا پنداشته بودند. بی‌گمان حوزه‌ای که پس از آن برای استعانت مورد توجه قرار گرفت زبان‌شناسی بود. اما در آن دوران نظریه زبان‌شناسی وجود نداشت که بتواند کمک شایانی به پردازش زبان‌ها بکند. در سال ۱۹۵۷ کتاب ساختارهای نحوی اثر نوام چامسکی زبان‌شناس جوان آمریکایی که از آن پس به شناخته‌شده‌ترین چهره زبان‌شناسی نظری تبدیل شد به چاپ رسید از آن پس پردازش زبان با حرکت‌های تازه‌ای دنبال شد اما هرگز قادر به حل کلی مسئله نشد.

محدودیت‌ها

پردازش زبان‌های طبیعی بسیار جذاب است و برای ارتباط بین انسان و ماشین محسوب می‌شود و در صورت عملی شدنش به طور کامل می‌تواند تحولات شگفت‌انگیزی را در پی داشته‌باشد. سیستم‌های قدیمی محدودی مانند SHRDLU که با واژه‌های محدود و مشخصی سر و کار داشتند، بسیار عالی عمل می‌کردند، به طوری که پژوهشگران را به شدت نسبت به این حوزه امیدوار کرده بودند. اما در تقابل با چالش‌های جدی‌تر زبانی و پیچیدگی‌ها و ابهام‌های زبان‌ها، این امیدها کم‌رنگ شدند. مسئله پردازش زبان‌های طبیعی معمولاً یک مسئله AI-Complete محسوب می‌شود، چرا که محقق شدن آن به طور کامل مستلزم سطح بالایی از درک جهان خارج و حالات انسان برای ماشین است.

موانع اساسی

رایانه برای آن که بتواند برداشت درستی از جمله‌ای داشته باشد و اطلاعات نهفته در آن جمله را درک کند، گاهی لازم است که برداشتی از معنای کلمات موجود در جمله داشته باشد و تنها آشنایی با دستور زبانکافی نباشد. مثلاً جمله حسن سيب را نخورد برای این که کال بود. و جمله حسن سيب را نخورد برای این که سیر بود. ساختار دستوری کاملاً یکسانی دارند و تشخیص این که کلمات «کال» و «سیر» به «حسن» برمی‌گردند یا به «سب»، بدون داشتن اطلاعات قبلی درباره ماهیت «حسن» و «سب» ممکن نیست.

دستور هیچ زبانی آن قدر دقیق نیست که با استفاده از قواعد دستوری همیشه بتوان به نقش هریک از اجزای جمله‌های آن زبان پی برد.

پردازش زبان‌های طبیعی آمارگرا

پردازش زبان‌های طبیعی به شکل آمارگرا عبارت است از استفاده از روش‌های تصادفی، احتمالاتی و آماری برای حل مسائلی مانند آنچه در بالا ذکر شد. به‌ویژه از این روش‌ها برای حل مسائلی استفاده می‌کنند که ناشی از طولانی بودن جملات و زیاد بودن تعداد حالات ممکن برای نقش کلمات هستند. این روش‌ها معمولاً مبتنی بر نمونه‌های متنی و مدل‌های مارکف هستند.

کارکردهای مهم پردازش زبان‌های طبیعی

خلاصه‌سازی خودکار

استخراج اطلاعات

بازیابی اطلاعات

تشخیص گفتار

ویرایش

زبان‌شناسی محاسباتی

خلاصه‌سازی

خلاصه‌سازی خودکار به عملیاتی گفته می‌شود که به وسیله یک برنامه کامپیوتری، باعث کاهش حجم متن و تولید خلاصه‌ای حاوی مهمترین نکات متن می‌شود.