



دانشگاه الزهراء (س)

دانشکده فنی و مهندسی

پایان نامه

جهت اخذ درجه کارشناسی ارشد

رشته مهندسی کامپیوتر - گرایش هوش مصنوعی

عنوان

تخمین اطمینان خروجی ترجمه ماشینی

استاد راهنما

دکتر نوشین ریاحی

دانشجو

مرضیه صالحی

اسفند ۱۳۹۲

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه الزهراء (س)

دانشکده فنی و مهندسی

پایان نامه

جهت اخذ درجه کارشناسی ارشد

رشته مهندسی کامپیوتر - گرایش هوش مصنوعی

عنوان

تخمین اطمینان خروجی ترجمه ماشینی

استاد راهنما

دکتر نوشین ریاحی

استاد مشاور

دکتر شهرام خدیوی

دانشجو

مرضیه صالحی

اسفند ۱۳۹۲

کلیه دستاوردهای این تحقیق متعلق به دانشگاه الزهراء (س)

می باشد.

سپاس و ستایش خداوند را که مرا در این راه یاری کرد.

با تشکر از

اساتید بزرگوارم سرکار خانم دکتر ریاضی و جناب آقای دکتر خدیوی که مرا در انجام

این پژوهش راهنمایی کردند.

همسر فداکارم که در این راه همراه من بود.

و پدر و مادر عزیزم که دلسوزانه پشتیبانم بودند.

چکیده

به دلیل ابهام ذاتی موجود در زبان طبیعی، تقریباً همه فناوری‌های پردازش زبان طبیعی ناکاملند. با این حال با تخمینی از کیفیت خروجی، کاربران می‌توانند به طور مناسبی با ناکامل بودن آنها برخورد کنند. زمینه تحقیقاتی که به این مسئله می‌پردازد تخمین اطمینان نامیده می‌شود. هدف اصلی تخمین اطمینان کاربردی‌تر کردن فناوری‌های ناکامل است. در حوزه ترجمه ماشینی نیز با وجود پیشرفت‌های به دست آمده در سال‌های اخیر، این فناوری هنوز قادر به ترجمه دقیق متون نیست به طوری که گاهی ممکن است ترجمه متنی توسط ماشین و ویرایش خروجی توسط انسان، زمان بیشتری نسبت به ترجمه مستقیم توسط انسان بگیرد. در چنین حالتی، داشتن تخمینی از درستی خروجی ترجمه ماشینی برای ویرایشگران مفید است تا تلاش خود را به جملات نادرستی معطوف کنند که نیاز به تغییرات پرهزینه ندارند. علاوه بر پس‌ویرایش، تخمین اطمینان خروجی ترجمه ماشینی برای کاربردهایی که هدف آنها ارتقا کیفیت ترجمه ماشینی است، از قبیل ترکیب چند سامانه ترجمه‌گر، بازتولید خروجی و مرتب‌سازی دوباره لیست چند ترجمه برتر، مفید می‌باشد. تخمین اطمینان یا سنجش کیفیت خروجی ترجمه ماشینی یکی از موضوعات چالشی در زمینه ترجمه ماشینی محسوب می‌شود. همچنین برای جفت زبان انگلیسی-فارسی پژوهش‌های چندانی در زمینه تخمین اطمینان صورت نگرفته است. در این پژوهش مجموعه‌ای از ویژگی‌های مبتنی بر ساختار و مبتنی بر محتوای مستقل از سامانه ترجمه‌گر ارائه شده و کارایی چند روش یادگیری برای ترکیب این ویژگی‌ها بررسی شده است. مزیت ویژگی‌های ساختاری ارائه شده این است که برخلاف اکثر ویژگی‌های ارائه شده در گذشته، برای بررسی صحت ساختاری کلمه مقصد از جمله مبدا استفاده می‌کنند. همچنین برای نخستین بار از بردار زمینه برای تخمین اطمینان استفاده شده است و چالش متفاوت بودن فضای بردارهای مبدا و مقصد، با استفاده از روشی آماری حل شده است. نتایج به دست آمده از آزمون ویژگی‌های پیشنهادی در قالب جداولی ارائه گردیده است.

کلمات کلیدی: تخمین اطمینان، ترجمه ماشینی، بردار زمینه، اطلاعات متقابل، ویژگی‌های مبتنی بر ساختار، ویژگی‌های مبتنی بر محتوا.

فهرست مطالب

۱	فصل اول: مقدمه
۲	۱-۱ ترجمه ماشینی آماری
۵	۲-۱ تخمین اطمینان
۷	۱-۲-۱ تخمین اطمینان ترجمه ماشینی در سطح کلمه
۸	۳-۱ طرح مسئله
۱۰	۴-۱ اهداف و نوآوری ها
۱۱	۵-۱ ساختار پایان نامه
۱۱	۶-۱ جمع بندی
۱۲	فصل دوم: پیشینه پژوهش
۱۳	۱-۲ مقدمه
۱۳	۲-۲ انواع تخمین اطمینان
۱۴	۳-۲ کاربردهای تخمین اطمینان خروجی ترجمه ماشینی
۱۴	۴-۲ چالش‌های پیش روی تخمین اطمینان
۱۶	۵-۲ معماری سامانه تخمین اطمینان
۱۶	۱-۵-۲ استفاده از درجه اطمینان
۱۷	۲-۵-۲ استفاده از روش‌های طبقه بندی
۱۸	۶-۲ روش‌های یادگیری ماشین استفاده شده
۱۸	۱-۶-۲ ترکیب خطی
۱۸	۲-۶-۲ رگرسیون منطقی
۱۹	۳-۶-۲ ماشین بردار پشتیبان (SVM)
۲۰	۴-۶-۲ بیز ساده
۲۱	۵-۶-۲ پرسپترون چند لایه
۲۲	۶-۶-۲ رگرسیون کمینه مربعات جزئی
۲۲	۷-۶-۲ درخت تصمیم
۲۳	۸-۶-۲ رگرسیون خطی تغییر داده شده
۲۳	۹-۶-۲ پیشینه انتروپی
۲۳	۷-۲ طبقه بندی پیشنهادی جهت معرفی ویژگی‌های تخمین اطمینان
۲۴	۱-۷-۲ ویژگی‌های پیش ترجمه
۲۶	۲-۷-۲ ویژگی‌های پس ترجمه

۳۵.....	۸-۲ جمع بندی
۳۶.....	فصل سوم: معرفی روش‌های تولید داده‌های آموزشی.....
۳۷.....	۱-۳ مقدمه
۳۷.....	۲-۳ برچسب زنی دستی
۳۸.....	۱-۲-۳ استاندارد کردن رای ها
۳۹.....	۲-۲-۳ میانگین گیری از رای ها
۳۹.....	۳-۳ برچسب زنی خودکار
۴۱.....	۱-۳-۳ روش ارزیابی برچسب زنی خودکار
۴۲.....	۲-۳-۳ ملاحظات لازم در هنگام داشتن چند مرجع.....
۴۳.....	۴-۳ روش نیمه خودکار
۴۵.....	۵-۳ مقایسه روش‌های موجود
۴۵.....	۶-۳ جمع بندی
۴۶.....	فصل چهارم: ویژگی‌های پیشنهادی برای تخمین اطمینان
۴۷.....	۱-۴ مقدمه
۴۷.....	۲-۴ سامانه تخمین اطمینان
۴۸.....	۳-۴ معرفی ویژگی‌های ارائه شده
۴۸.....	۱-۳-۴ ویژگی‌های مبتنی بر محتوا (سازگاری محتوایی)
۵۶.....	۲-۳-۴ ویژگی‌های مبتنی بر ساختار (سازگاری ساختاری)
۶۲.....	۴-۴ تخمین اطمینان به عنوان یک مسئله برچسب زنی دنباله
۶۳.....	۵-۴ جمع بندی
۶۴.....	فصل پنجم: پیاده سازی و آزمون
۶۵.....	۱-۵ مقدمه
۶۵.....	۲-۵ تهیه داده‌های آموزش و آزمون
۶۵.....	۱-۲-۵ امکانات و ابزارهای مورد نیاز
۶۵.....	۲-۲-۵ پیاده سازی
۶۶.....	۳-۵ پیش پردازش
۶۶.....	۱-۳-۵ امکانات و ابزارهای مورد نیاز
۶۶.....	۲-۳-۵ پیاده سازی
۷۰.....	۴-۵ استخراج ویژگی‌ها
۷۰.....	۱-۴-۵ امکانات و ابزارهای مورد نیاز

۷۰	۲-۴-۵ معیارهای ارزیابی
۷۱	۳-۴-۵ احتمال پسین، جمع رتبه‌ها و فرکانس مربوط
۷۳	۴-۴-۵ اطلاعات متقابل بین زبانی و درون زبانی
۷۴	۵-۴-۵ الگوی کلمات و برچسب POS مقصد
۷۵	۶-۴-۵ سازگاری فعل
۷۶	۷-۴-۵ همترازی نقش دستوری
۷۷	۸-۴-۵ عبارات دستوری
۷۸	۹-۴-۵ ویژگی‌های مبتنی بر بردار زمینه
۸۱	۱۰-۴-۵ ویژگی مبتنی بر اطلاعات متقابل بهبود یافته
۸۳	۵-۵ ترکیب ویژگی‌ها و ارزیابی نهایی
۸۴	۱-۵-۵ امکانات و ابزارهای مورد نیاز
۸۴	۲-۵-۵ نتایج
۸۸	۶-۵ جمع بندی
۸۹	فصل ششم: نتیجه گیری و توسعه‌های آتی
۹۰	۱-۶ نتیجه گیری
۹۱	۲-۶ توسعه‌های آتی
۹۲	مراجع
۹۸	پیوست: شبه کد پیاده سازی سامانه

فهرست شکل‌ها

۳	شکل ۱-۱: مدل مبتنی بر کلمه [1]
۳	شکل ۱-۲: مدل مبتنی بر عبارت [1]
۶	شکل ۱-۳: مدل مبتنی بر درخت [1]
۱۷	شکل ۲-۱: استفاده از درجه اطمینان برای طبقه بندی [13]
۱۸	شکل ۲-۲: استفاده از روش‌های یادگیری ماشین
۲۴	شکل ۲-۳: طبقه بندی پیشنهادی بر اساس ویژگی‌های استفاده شده
۲۵	شکل ۲-۴: مثالی از ویژگی عبارت مبدا

- شکل ۲-۵: مثالی از ویژگی برچسب POS عبارت مبدا ۲۵
- شکل ۲-۶: مثالی از ویژگی زمینه عبارت ۲۶
- شکل ۳-۱: هیستوگرام رأی‌های دو رأی دهنده متفاوت [۲] ۳۸
- شکل ۳-۲: مثالی از معیارهای خطای مختلف: رشته ABCBDB با ABBCE مقایسه شده است ۴۰
- شکل ۳-۳: همبستگی معیارهای ارزیابی مختلف با قضاوت انسان [10] ۴۲
- شکل ۴-۱: فاز آموزش سامانه تخمین اطمینان خروجی ترجمه ماشینی ۴۷
- شکل ۴-۲: فاز آزمون سامانه تخمین اطمینان خروجی ترجمه ماشینی ۴۸
- شکل ۴-۳: وزن دهی ویژگی مبتنی بر اطلاعات متقابل با مدل نمایی ۵۱
- شکل ۴-۴: تاثیر اندازه بردارها در فاصله اقلیدسی ۵۳
- شکل ۴-۵: فاصله اقلیدسی نرمال سازی شده ۵۴
- شکل ۴-۶: ویژگی‌های سازگاری فعل ۵۸
- شکل ۴-۷: مراحل استخراج زمان و شخص جمله مبدا ۵۹
- شکل ۴-۸: مراحل استخراج عبارات دستوری جمله مقصد ۶۱
- شکل ۴-۹: فلوجارت تولید مجموعه کلمات به جا و نابجا ۶۲
- شکل ۵-۱: مراحل پیش پردازش ۶۷

فهرست جدول ها

- جدول ۲-۱: معرفی ویژگی‌های استفاده شده در تخمین اطمینان و مزایا و معایب آن ها ۳۴
- جدول ۳-۱: مقایسه روش‌های تولید داده برای تخمین اطمینان ۴۵
- جدول ۵-۱: مشخصه‌های آماری پیکره ها ۶۶
- جدول ۵-۲: نتایج ویژگی‌های احتمال پسین، جمع رتبه‌ها و فرکانس مربوط ۷۳
- جدول ۵-۳: نتایج مربوط به ویژگی‌های مبتنی بر اطلاعات متقابل ۷۴
- جدول ۵-۴: نتایج ویژگی‌های لغوی ۷۵

- جدول ۵-۵: نتایج ویژگی‌های سازگاری فعل ۷۵
- جدول ۵-۶: نتایج ویژگی‌های احتمال همترازی نقش دستوری ۷۷
- جدول ۵-۷: نتایج ویژگی‌های عبارات دستوری ۷۷
- جدول ۵-۸: مقایسه روش‌های محاسبه فاصله/شباهت بردارهای زمینه ۸۰
- جدول ۵-۹: نتایج مربوط به ترکیب ویژگی‌های مبتنی بر بردار زمینه کسینوسی ۸۰
- جدول ۵-۱۰: مقایسه ویژگی‌های مبتنی بر بردار محتوای کسینوسی با سامانه‌های پایه ۸۱
- جدول ۵-۱۱: مقایسه مقادیر مختلف پارامترهای مدل نمایی و خطی ۸۲
- جدول ۵-۱۲: بررسی مفید بودن وزن فعل ۸۳
- جدول ۵-۱۳: مقایسه مدل نمایی با مدل یکنواخت ۸۳
- جدول ۵-۱۴: مقایسه ویژگی مبتنی بر اطلاعات متقابل بهبود یافته با سامانه‌های پایه ۸۳
- جدول ۵-۱۵: مقایسه روش برجسب زنی دنباله با برجسب زنی به کلمات ۸۴
- جدول ۵-۱۶: نتایج سامانه‌های پایه ۸۶
- جدول ۵-۱۷: نتایج مربوط به ترکیب ویژگی‌های مبتنی بر ساختار ۸۶
- جدول ۵-۱۸: نتایج مربوط به ترکیب ویژگی‌های مبتنی بر محتوا ۸۷
- جدول ۵-۱۹: نتایج مربوط به ترکیب ویژگی‌های مبتنی بر محتوا با ویژگی‌های مبتنی بر ساختار ۸۷
- جدول ۵-۲۰: نتایج مربوط به ترکیب ویژگی‌های مبتنی بر ساختار با ویژگی‌های مبتنی بر بردار زمینه کسینوسی ۸۷
- جدول ۵-۲۱: نتایج مربوط به ترکیب ویژگی مبتنی بر اطلاعات متقابل بهبود یافته با ویژگی‌های ساختاری ۸۸

فصل اول:

مقدمه

۱-۱ ترجمه ماشینی آماری

ترجمه ماشینی آماری^۱ را می‌توان به صورت یادگیری عمل ترجمه کردن از روی داده‌های آموزشی توسط ماشین، تعریف کرد. این کار به صورت طراحی مدل و یادگیری پارامترهای مدل از روی داده‌ها انجام می‌شود. در ترجمه ماشینی آماری، ترجمه به صورت یک فرآیند تصمیم‌گیری مدل می‌شود. با داشتن یک جمله مبدا $f_1^J = f_1 \dots f_j \dots f_J$ ، جمله مقصد $e_1^I = e_1 \dots e_i \dots e_I$ با بیشینه کردن احتمال پسین به دست می‌آید:

$$\hat{e}_1^I = \underset{I, e_1^I}{\operatorname{argmax}} \{ \Pr(e_1^I | f_1^J) \} = \underset{I, e_1^I}{\operatorname{argmax}} \{ \Pr(f_1^J | e_1^I) \cdot \Pr(e_1^I) \} \quad (1-1)$$

بنابراین به دو مدل نیاز است. مدل ترجمه $\Pr(f_1^J | e_1^I)$ و مدل زبانی $\Pr(e_1^I)$ که آنها را می‌توان به طور مستقل مدل کرد. مدل زبانی جمله مقصد خوش‌ترتیبی جمله مقصد را در بر می‌گیرد. مدل‌های ترجمه را می‌توان به صورت زیر تقسیم بندی کرد [1]:

۱. مدل‌های مبتنی بر کلمه

۲. مدل‌های مبتنی بر عبارت

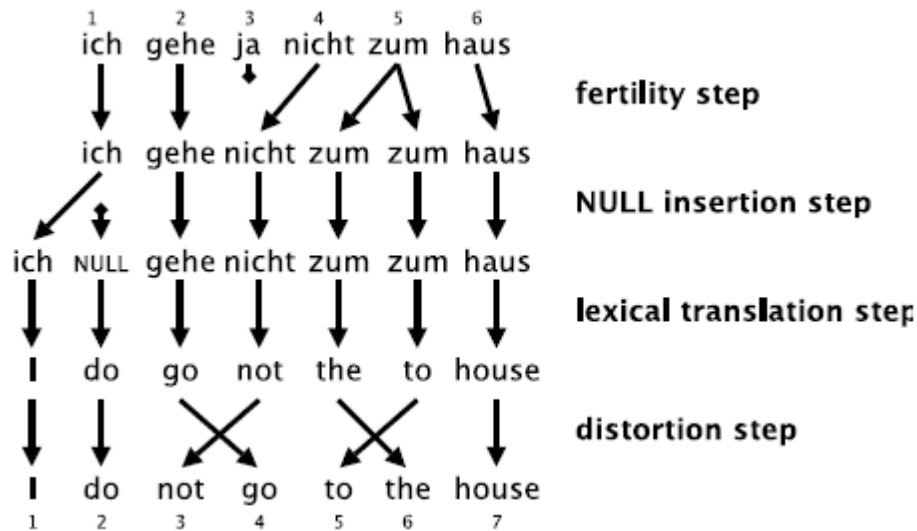
۳. مدل‌های مبتنی بر درخت

امروزه بیشتر سامانه‌های ترجمه‌گر آماری بر مبنای عبارات دوزبانه هستند [2, 3, 4, 5, 6, 7, 8]. از اینرو این مدل با جزئیات بیشتری بیان می‌شود.

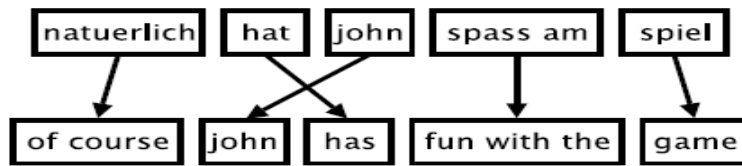
الف) مدل‌های مبتنی بر کلمه: مدل‌هایی هستند که کوچکترین واحدی که مدل‌ها بر اساس آنها ساخته می‌شود کلمه است. مثلاً مدل‌هایی که شامل احتمال‌هایی از این قبیلند: احتمال اینکه یک کلمه خاص، ترجمه کلمه خاص دیگری باشد. در شکل ۱-۱ شمایی از این مدل آمده است. این شکل مربوط به مدل IBM-3 است.

ب) مدل‌های مبتنی بر عبارت: در این مدل‌ها احتمالات بر روی واحد بزرگتری به نام عبارت تعریف می‌شوند. عبارت به معنی دنباله‌ای از کلمات است و لزوماً به معنی عبارتی که از لحاظ زبان معتبر باشد نیست. عبارت‌های دوزبانه از پیکره دوزبانه آموزش که همترازی کلمات روی آن انجام شده‌است به دست می‌آیند. در شکل ۱-۲ شمایی از مدل مبتنی بر عبارت دیده می‌شود.

¹Statistical machine translation



شکل ۱-۱: مدل مبتنی بر کلمه [1]



شکل ۱-۲: مدل مبتنی بر عبارت [1]

احتمال پسین $P(f_1^I | e_1^I)$ مستقیماً با استفاده از ترکیب لگاریتم-خطی وزن دار مدل زبانی، مدل ترجمه عبارت و مدل لغوی مبتنی بر کلمه به دست می‌آید.

در روش مبتنی بر عبارت، سه مرحله ضروری است: قطعه‌بندی جمله ورودی، ترجمه هر عبارت به عبارت معادل زبان مقصد و تعیین ترتیب درست عبارات ترجمه شده در جمله مقصد.

ابتدا جمله ورودی f به دنباله‌ای از I عبارت f_1^I تقسیم می‌شود. سپس هر عبارت f_i^I به عبارت معادل \bar{e}_i ترجمه می‌شود. در نهایت ممکن است عبارات ساخته شده جابجا شوند. برای مثال در یکی از سامانه‌های مبتنی بر عبارت شناخته شده به نام موزز [3]، از ۴ مدل برای ترجمه استفاده می‌شود:

- جدول عبارات $\phi(\bar{f}_i | \bar{e}_i)$
- مدل زبانی $P_{LM}(e)$

- مدل جابجایی^۱
- جریمه کلمه^۲

هر یک از موارد فوق، بخشی از کیفیت ترجمه را کنترل می‌کند. مدل جابجایی، ترتیب عبارات در جمله را کنترل می‌کند و با هر جابجایی هزینه‌ای را به سامانه تحمیل می‌نماید.

جریمه کلمه، وظیفه کنترل تعداد کلمات در جمله را بر عهده دارد و اجازه نمی‌دهد که جملات با طول بسیار بزرگ و یا بسیار کوچک تولید شوند. به هر یک از اجزاء فوق وزنی نسبت داده می‌شود و احتمال ترجمه دو عبارت به شکل زیر محاسبه می‌گردد:

$$P(e|f) = \phi(f|e)^{weight_\phi} \cdot LM^{weight_{LM}} \cdot D(e, f)^{weight_d} \cdot W(e)^{weight_w} \quad (2-1)$$

بطوریکه $P(f|e)$ در پیاده‌سازی موزز به شکل زیر تجزیه می‌شود:

$$P(e|f) = P(f|e)P_{LM}(e)\omega^{length(e)} \quad (3-1)$$

بر اساس ترجمه مبتنی بر عبارت ابتدا جمله ورودی f به I عبارت \bar{f}_1^I افزای می‌شود سپس هر عبارت \bar{f}_i به عبارت معادل \bar{e}_i از زبان مقصد ترجمه می‌شود. در نهایت عبارات در جمله معادل جابجا خواهند شد. احتمال ترجمه عبارت با تابع $\Phi(\bar{f}_i|\bar{e}_i)$ بیان می‌شود. جابجایی عبارات خروجی در حالت پیش فرض با استفاده از مدل ساده $d(start_i, end_{i-1})$ محاسبه می‌شود. $start_i$ ابتدای عبارتی است که به i امین عبارت از زبان مقصد ترجمه شده و end_{i-1} انتهای عبارتی است که به $i-1$ امین عبارت از زبان مقصد ترجمه شده است. ساده ترین روش تعریف مدل جابجایی به صورت $d(start_i, end_{i-1}) = \alpha^{|start_i - end_{i-1} - 1|}$ ، با انتخاب مقدار مناسب برای α است. در نهایت بهترین ترجمه اینگونه محاسبه می‌شود:

$$\begin{aligned} e_{best} &= \arg \max_e P(e|f) \\ &= \arg \max_e P(f|e)P_{LM}(e)\omega^{length(e)} \end{aligned} \quad (4-1)$$

بطوریکه:

$$P(\bar{f}_1^I | \bar{e}_1^I) = \Phi(\bar{f}_i | \bar{e}_i) d(start_i, end_{i-1}) \quad (5-1)$$

¹Distortion Model

²Word Penalty

برای یافتن بهترین دنباله‌ای که بیشترین احتمال ترجمه را دارا باشد، لازم است فضای حالات ممکن جستجو گردد. به دلیل گستردگی فرضیات خروجی، جستجوی این فضا از مرتبه نمایی بوده و حل مسئله یافتن جمله با حداکثر احتمال، یک مسئله NP-Complete است. بنابراین برای حل بهینه این مسئله از الگوریتم‌های جستجوی بهینه مانند جستجوی شعاعی^۱ یا A^* استفاده می‌شود. ارتقاء کیفیت ترجمه ابتدا با جدول عبارات دقیق‌تر و سپس تنظیم مناسب‌تر پارامترهای رمزگشایی حاصل می‌شود. برای این منظور ابتدا مدل بهتری تخمین زده می‌شود. مثلاً مدل زبانی مرتبه بالاتر استفاده می‌گردد. مدل زبانی از دو جنبه موثر است. اول در محاسبه هزینه آتی^۲ برای امتیاز دهی عبارات، که با بهبود آن عبارات بهتری تولید خواهد شد. از طرف دیگر برای اعتبارسنجی رشته خروجی تولید شده و برای کنترل کیفیت ترجمه نهایی استفاده می‌گردد و بنابراین در دو مرحله می‌تواند به بهبود سامانه منجر شود. اما نکته این است که با ارتقاء و بهبود مدل به سمت مدل بالاتر که لاجرم پیچیدگی بیشتری را نیز به همراه خواهد داشت، ضروری است که فضای جستجو نیز متناسب با آن افزایش داده شود.

مدل‌های مبتنی بر درخت: در مدل‌های مبتنی بر درخت ساختار دستور زبانی جملات نیز مورد توجه قرار می‌گیرد و شرط اعتبار عبارات مدل قبل موجود بودن تمام کلمات عبارت در یک زیر درخت دستوری از جمله است. شکل ۱-۳ شمایی از این مدل را نمایش می‌دهد.

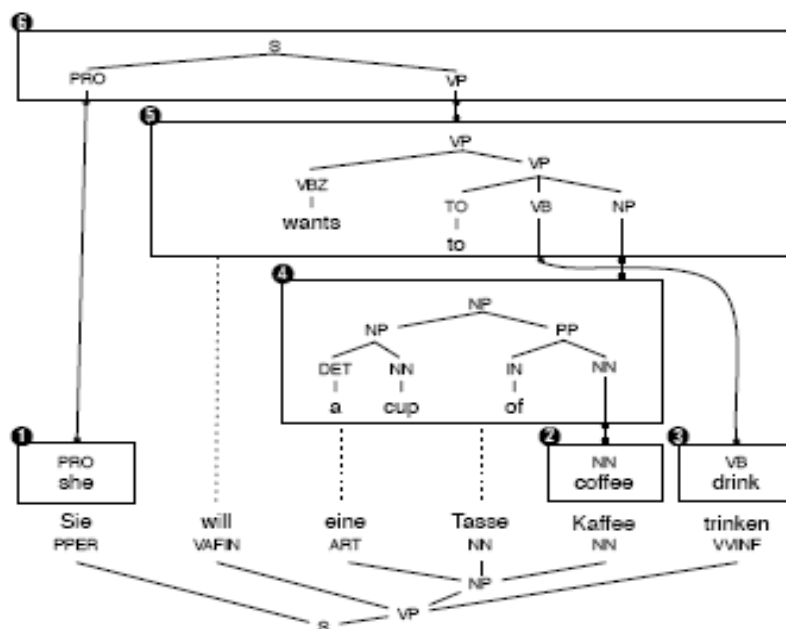
۱-۲ تخمین اطمینان

تخمین اطمینان (CE) یک رویکرد یادگیری ماشین برای محاسبه درجه درست بودن خروجی (درجه اطمینان) در یک کاربرد مبتنی بر فناوری‌های ناکامل، مانند ترجمه ماشینی است. اگر $x \in X$ ورودی ترجمه‌گر و $y \in Y$ خروجی مربوط به آن در نظر گرفته شود، درجه اطمینان تابعی به صورت $c : X \times Y \times K \rightarrow I$ است، به طوری که X دامنه ورودی‌ها، Y دامنه خروجی‌ها، K دامنه دانش اضافه مورد نیاز و I مجموعه مقادیری که نماینده درستی خروجی هستند، می‌باشند. I می‌تواند بازه‌ای از اعداد حقیقی یا مجموعه‌ای از سطوح مختلف مانند {«افزوده شده»، «جابه‌جا شده»، «جانشین شده»، «خوب»} [9] و یا مجموعه دوتایی {«درست»، «نادرست»} [10, 11, 12, 13, 14] باشد. یک حالت خاص از تابع c به صورت $P(C = 1|x, y, k)$ تعریف می‌شود که $C \in \{0, 1\}$. $C = 1$ نشان دهنده درستی y است. تخمین اطمینان به طور گسترده در حوزه تشخیص

¹ Beam search

² Future cost

گفتار استفاده شده است، اما در حوزه‌های دیگر پردازش زبان طبیعی پیشینه کوتاه‌تری دارد. تخمین اطمینان را می‌توان به صورت لایه‌ای در نظر گرفت که بر روی سامانه پایه قرار می‌گیرد، برای مثال در تشخیص گفتار خروجی‌های سامانه تشخیص گفتار به این لایه داده می‌شود و در صورتی که درجه اطمینان پایینی برای آن تخمین زده شود، از کاربر تقاضا می‌شود جمله را دوباره تکرار کند. تخمین اطمینان می‌تواند در سطوح مختلفی از جمله سطح کلمه، سطح زیرجمله [15]، سطح جمله [16] و سطح سند (به ندرت) انجام شود. هدف تخمین اطمینان در هریک از این سطوح زدن برچسب به واحد مربوط به آن سطح است.



شکل ۱-۳: مدل مبتنی بر درخت [1]

اگرچه تخمین اطمینان و ارزیابی خودکار ترجمه ماشینی با استفاده از معیارهایی از قبیل BLEU [17] و NIST [18]، هر دو کیفیت خروجی ترجمه ماشینی را تعیین می‌کنند، ولی با هم متفاوتند. مهم‌ترین تفاوت بین آنها ناشی از این واقعیت است که تخمین اطمینان بدون مقایسه خروجی ترجمه ماشینی با ترجمه‌های درست و یا پرسیدن نظر یک مترجم انسانی در مورد خروجی انجام می‌شود.

۱-۲-۱ تخمین اطمینان ترجمه ماشینی در سطح کلمه

هدف تخمین اطمینان خروجی ترجمه ماشینی در سطح کلمه نسبت دادن احتمال درستی یا امتیاز درستی (تخمین اطمینان قوی) و یا برچسب‌های گسسته (تخمین اطمینان ضعیف) به کلمات خروجی ترجمه ماشینی است که بدون استفاده از ترجمه مرجع انجام می‌شود. بر اساس نوع تخمین اطمینانی که تولید می‌شود (ضعیف یا قوی)، الگوریتم‌های یادگیری ماشین متفاوتی برای این کار استفاده می‌شوند. برای تولید تخمین اطمینان ضعیف، الگوریتم‌های طبقه بندی یادگیری ماشین برای نسبت دادن برچسب دوتایی یا چندتایی به هر کلمه، مورد استفاده قرار می‌گیرند. در حالت قوی، تکنیک‌های رگرسیون به کار می‌روند. در این بخش جزئیاتی از تخمین اطمینان سطح کلمه با استفاده از یک روش یادگیری نظارتی بیان می‌شود.

برای یک کلمه در جایگاه j از یک جمله خروجی ترجمه ماشینی، برچسب تخمین اطمینان آن به صورت زیر نشان داده می‌شود که t جمله مقصد و s جمله مبدا می‌باشد:

$$C_{j,t,s} \in \{0,1\} \quad (6-1)$$

برای تولید تخمین اطمینان تابع توزیع احتمال درستی یک کلمه در حالی که جمله مبدا و مقصد داده شده‌اند، یادگیری می‌شود. از نظر ریاضی، با داشتن جمله مبدا و مقصد، هدف یادگیری، تخمین احتمال زیر برای هر کلمه است:

$$p(C_{j,t,s} = 1 | j, t, s) \quad (7-1)$$

اما چون s و t ممکن است هر جفت جمله‌ای باشند، احتمال بالا به صورت مستقیم به دست نمی‌آید. به همین دلیل هر کلمه به صورت یک بردار از ویژگی‌های عددی نشان داده می‌شود. چنین بردار ویژگی با استفاده از یک تابع x به دست می‌آید. تابع x اطلاعات استخراج شده از جملات مبدا و مقصد و خروجی ترجمه ماشینی را به بردار ویژگی نظیر می‌کند:

$$x: (j, t, s) \in j \times l_s \times l_t \rightarrow x(j, t, s) \in R^{dw} \quad (8-1)$$

در واقع اطلاعات استخراج شده از جفت جمله مبدا و مقصد و خروجی ترجمه ماشینی، برای استخراج ویژگی‌هایی که جنبه‌های مختلف ترجمه را در بر می‌گیرند به کار می‌رود. با نمایش کلمات به این شیوه، تابع توزیع احتمال یادگیری شده به صورت زیر در می‌آید:

$$p(C_{j,t,s}; j, t, s) = P(C_{j,t,s} | x(j, t, s)) \quad (9-1)$$

که $C_{j,t,s} \in \{0,1\}$.

تابع توزیع احتمال فوق از یک مجموعه داده آموزشی یادگیری می‌شود. داده‌های آموزشی شامل نمونه‌های آموزشی برچسب خورده هستند. هر نمونه آموزشی به صورت یک جفت از بردار ویژگی (نماینده هر کلمه) و برچسب مربوط به آن است. بعد از آموزش، توزیع احتمال یادگیری شده برای انجام طبقه بندی به کار می‌رود.

$$\hat{c} : (j, t, s) \rightarrow \hat{c}(j, t, s) \in \{0,1\} \quad (10-1)$$

یک برچسب بر اساس فرمول زیر به یک بردار نسبت داده می‌شود:

$$\hat{c}(j, t, s) = \operatorname{argmax}_{c \in \{0,1\}} p(c; j, t, s) \quad (11-1)$$

در عبارت بالا برچسبی به کلمه اختصاص می‌یابد که احتمال آن از $0/5$ بیشتر باشد ولی این مقدار حد آستانه، گاهی نامناسب است چرا که هزینه دادن برچسب نادرست به کلمه درست بسیار کمتر از هزینه دادن برچسب درست به کلمه نادرست است. برای همین یک حد آستانه τ تعریف می‌شود:

$$\hat{c}(j, t, s; \tau) = \begin{cases} 1 & \text{if } p(1; j, t, s) > \tau \\ 0 & \text{otherwise} \end{cases} \quad (12-1)$$

از عبارت بالا برمی‌آید که اگر احتمال درستی یک کلمه از مقدار حد آستانه بیشتر باشد برچسب درست می‌خورد و از آنجا که در برخی کاربردها تشخیص کلمات نادرست مهمتر است، در این موارد مقدار حد آستانه بیشتر از $0/5$ در نظر گرفته می‌شود. مقدار حد آستانه می‌تواند با استفاده از معیارهای ارزیابی روی مجموعه آموزشی بهینه شود.

به طور خلاصه انجام تخمین اطمینان در سطح کلمه نیاز به دو کار دارد: تولید ویژگی‌های مناسب و طبقه بندی مناسب.

۱-۳ طرح مسئله

اکثر فناوری‌های زبان طبیعی ناکاملند. یکی از دلایل آن این است که به طور کلی در بر گرفتن پیچیدگی‌ها و ابهامات زبان طبیعی سخت است. دلیل دیگر به استفاده از روش‌های آماری برمی‌گردد. روش‌های آماری در دو دهه اخیر کاربرد گسترده‌ای در تشخیص گفتار و پردازش زبان طبیعی داشته‌اند و به موفقیت‌های قابل ملاحظه‌ای نیز رسیده‌اند. این موفقیت به دلیل این واقعیت است که این رویکرد مستقل از زبان است و تنها به پیکره‌های به اندازه کافی بزرگ برای تخمین چگالی‌های احتمال نیاز است. با این حال در روش‌های آماری یک مشکل وجود دارد: آنها تنها محتمل ترین نتایج را با توجه به داده‌های آموزشی و ورودی تولید می‌کنند. روشن است که این روش گاهی با توجه به انتظار انسان بهینه نیست [12]. با این وجود می‌توان با استفاده از راه

حل‌هایی خطای موجود در فناوری‌های ناکامل را کنترل کرد. برای مثال در سامانه تشخیص گفتار، درجه اطمینان پایین در آنالیز سخن یک کاربر می‌تواند موجب شود سامانه عملیات را تکرار کند. این استراتژی قابلیت بهبود دادن کارایی سامانه را دارد، به شرطی که تخمین اطمینان دقیق باشد [10]. بنابراین سنجش کیفیت نتایج، به طور خودکار مهم است. تخمین اطمینان این کار را انجام می‌دهد.

با وجود اینکه از زمان ظهور ترجمه ماشینی یعنی تقریباً ۶۰ سال پیش، بهبودهای وسیعی در این حوزه به وجود آمده‌است، کارایی سامانه‌های ترجمه ماشینی مدرن نیز هنوز فاصله زیادی تا کامل شدن دارند. ترجمه یک کلمه می‌تواند اشتباه، جابجا، یا غیر موجود باشد. همچنین کلمات اضافه‌ای می‌توانند تولید شوند. با این حال سامانه‌های ترجمه ماشینی موجود برای کاربردهایی از جمله ترجمه سریع برای جستجو و یا کمک به مترجمان انسانی برای ترجمه مفید است [19]. برای کاربردهایی که فقط نیاز به مراد جمله مبدا دارند و نیاز به ترجمه دقیق ندارند، سامانه‌های ترجمه ماشینی جدید خصوصاً سامانه‌های ترجمه ماشینی آماری جوابگو هستند. اما برای کاربردهایی از جمله ترجمه شرح اقدامات و مذاکرات رسمی که ترجمه دقیق مورد نیاز است، سامانه‌های ترجمه ماشینی فقط می‌توانند یک ترجمه ابتدایی ارائه دهند. سپس مترجم‌های انسانی این ترجمه‌های ابتدایی را با پس‌ویرایش تصحیح می‌کنند [20]. اما گاهی ممکن است پس‌ویرایش زمان بیشتری نسبت به ترجمه مستقیم توسط انسان بگیرد. به دلیل نبودن دانش در مورد کیفیت ترجمه اولیه، مترجمان انسانی ممکن است زمان خود را برای پس‌ویرایش ترجمه‌های بسیار بی ربط هدر دهند. در چنین حالتی، داشتن تخمینی از درستی یا اطمینان خروجی ترجمه ماشینی برای ویرایشگران مطلوب است تا تلاش خود را به جملاتی معطوف کنند که نیاز به تغییرات پرهزینه ندارند. همچنین این تخمین‌ها می‌توانند به ویرایشگران کمک کنند که بخش‌های درست ترجمه را پس‌ویرایش نکنند [20]. همچنین تشخیص خطا با تخصیص معیارهای اطمینان به هر کلمه، عبارت یا جمله ترجمه شده، برای بهبود کیفیت سامانه مفید است [12]. تخمین اطمینان سامانه را قادر می‌کند خطاهای محتمل را به اطلاع کاربر برساند یا تنها ترجمه‌های با اطمینان بالا را ارائه دهد [14].

در پاسخ به اینکه « چرا از امتیازاتی که خود سامانه پایه به خروجی‌ها می‌دهد برای درجه اطمینان استفاده نشود؟ » باید گفت استفاده از لایه تخمین اطمینان بر روی سامانه پایه، نسبت به استفاده از امتیازاتی که خود سامانه پایه به خروجی اختصاص می‌دهد مزایایی دارد. مزیت کلیدی آن این است که به صورت یک ماژول جدا از سامانه است. این مسئله مزایای زیر را به دنبال دارد: