

استفاده از بازآرایی نحوی جهت بهبود ترجمه ماشینی آماری انگلیسی به فارسی

رضا سعیدی^۱، محسن کاهانی^۲، سیداحمد جکیان طوسی^۳ و بهداد بهمدی مقدس^۴

reza.saeedi@stu-mail.um.ac.ir^۱

mohesn.kahani@um.ac.ir^۲

ahmad.toosi@stu-mail.um.ac.ir^۳

behdad.behmadi@stu-mail.um.ac.ir^۴

آزمایشگاه فناوری وب دانشگاه فردوسی مشهد

چکیده

ترجمه ماشینی آماری به عنوان یکی از بهترین روش ها برای ترجمه از یک زبان به زبان دیگر شناخته می شود. برای زبان هایی که از لحاظ ساختار دارای شباهت زیادی به یکدیگر هستند خروجی این مترجم بسیار مناسب می باشد. تفاوت های ساختاری میان زبان انگلیسی و فارسی و همچنین عدم وجود پیکره دوزبانه بزرگ باعث شده است که این روش برای زبان ترجمه انگلیسی به فارسی ترجمه های مطلوبی را تولید نکند. ما در این مقاله سعی کرده ایم با استفاده از رهیافت بازآرایی کلمات، تا حد ممکن شباهت ساختاری میان عبارات انگلیسی و فارسی را افزایش دهیم. در ادامه تاثیر این عمل را بر روی بهبود نتایج خروجی مورد بررسی قرار داده ایم. به همین منظور ابتدا با کمک درخت تجزیه، مجموعه ای از قوانین بازآرایی استخراج شده است. سپس این قوانین به عنوان یک عمل پیش پردازشی بر روی عبارات انگلیسی اعمال گردیده است. نتایج بررسی ها نشان می دهد که خروجی مترجم پس از اعمال این روش منجر به بهبود کیفیت ترجمه در معیار BLEU شده است.

کلمات کلیدی

ترجمه ماشینی آماری، بازآرایی نحوی، مدل زبانی، مدل ترجمه، درخت تجزیه

۱ - مقدمه

اساس رابطه (۱) با ورود یک جمله در زبان مبدا (e) شبیه ترین ترجمه آن در زبان مقصد (f) یافت می شود [۲].

$$e' = \arg \max_e P(e|f) \quad (1)$$

با توجه به رابطه (۲) $P(e|f)$ مبتنی بر دو فاکتور است. میزان احتمال بودن یک جمله در زبان e ، که این فاکتور بعنوان مدل زبانی ($P(e)$) شناخته می شود و روشی که جملات موجود در زبان e به جملات موجود در زبان f تبدیل می شوند که این فاکتور نیز بعنوان مدل ترجمه ($P(f|e)$) شناخته می شود.

$$e' = \arg \max_e \frac{P(e)P(f|e)}{P(f)} \quad (2)$$

همانطور که در رابطه (۳) نشان داده شده است، این دو فاکتور از اعمال قانون بیز بر روی معادله اولیه حاصل می گردند. از آنجایی که مقدار $P(f)$ همواره ثابت است، آن را از بیشینه سازی حذف می کنیم. لذا داریم:

$$e' = \arg \max_e P(e)P(f|e) \quad (3)$$

ترجمه ماشینی به معنای استفاده از کامپیوتر برای ترجمه از یک زبان به زبان دیگر به صورت خودکار می باشد [۱]. تفاوت ذاتی میان برخی از زبانها کار ترجمه را بسیار مشکل می سازد. روشهای سنتی برای ترجمه ماشینی، بر روی اطلاعات زبان شناختی بشر، بصورت قوانین تبدیل متن از زبانی به زبان دیگر تکیه داشت. با توجه به وسعت دامنه اطلاعات موجود در زبان ها، استفاده از این روشها، بسیار پیچیده و نیازمند دانش بسیار می باشد. ترجمه ماشینی آماری، یک رویکرد متفاوت است که سعی دارد بطور خودکار از روی حجم زیادی از داده های آموزشی، به دانش گفته شده در ترجمه، دست پیدا کند. این دانش که عموماً بصورت احتمالاتی مبتنی بر قابلیت های گوناگون زبانی می باشد و به منظور راهنمایی در فرایند ترجمه مورد استفاده قرار می گیرد.

از این رو ترجمه ماشینی آماری می تواند به عنوان یک رویه به منظور یافتن جمله ای با بالاترین احتمال در زبان مقصد تعریف شود. بر