

تبديل متن محاوره‌اي فارسي به رسمی به کمک N_gram

ناديه آرمین^۱ ، مهرتوش شمس‌فرد^۲

^۱ آزمایشگاه پردازش زبان طبیعی، گروه مهندسی کامپیوتر، دانشکده مهندسی برق و کامپیوتر، دانشگاه شهید بهشتی
na.armin@gmail.com

^۲ استادیار، آزمایشگاه پردازش زبان طبیعی، گروه مهندسی کامپیوتر، دانشکده مهندسی برق و کامپیوتر، دانشگاه شهید بهشتی
m-shams@sbu.ac.ir

چکیده

با افزایش متون محاوره، تبدیل متن محاوره‌ای به رسمی یکی از چالش‌های موجود در پردازش زبان است. در این مقاله به ارائه و آزمون، راهکاری برای تبدیل متون محاوره‌ای به متون رسمی با استفاده از ترکیب روش‌های مبتنی بر قاعده و مدل‌سازی آماری می‌پردازیم. همچنین چگونگی ساخت پیکره، نحوه توکن‌بندی و نحوه یافتن ریشه کلمات، قوانین تبدیل کلمات محاوره‌ای به رسمی، الگوریتم پیشنهادی و نحوه بکارگیری N_gram بیان می‌شود.

كلمات کلیدی

پردازش زبان، متن محاوره، پیکره، N_gram

۱ - مقدمه

تحلیل ریخت شناسی^۳ در زبان فارسی معرفی کرد. با توجه به این که در این کار روشی برای گسترش این سیستم روی متون و بلاغ‌های فارسی ارائه شده، برخی از قواعد ساختواری برای متون محاوره ای معروفی و اعمال شده است.

اسدی در سال ۲۰۰۷[۳]، تحقیقی بر روی آواهای حذف شده در محاورات فارسی ارائه کرد. پیکره استفاده شده در این تحقیق بیست دقیقه مکالمه سه فارسی زبان است. این تحقیق نشان می‌دهد حذف آواها در یک کلمه به نوع آن کلمه (کلمات دستوری که کاربرد بیشتری در کلام دارند، تغییر بیشتری نسبت به کلمات دیگر کرده‌اند) و محل قرارگیری آن آوا در کنار حروف دیگر وابسته است.

بررسی منابع نشان می‌دهد اگر چه مطالعات زبانی در حوزه متون محاوره‌ای فارسی صورت گرفته ولی سابقه‌ای از تحقیق در باره تبدیل متون محاوره‌ای به رسمی جهت اعمال سایر پردازش‌های متداول دیده نمی‌شود. از آنجاکه بیشتر پردازش‌های زبانی که صورت گرفته است متمرکز بر متون رسمی بوده است، برای استفاده از این پردازش‌ها روی متون محاوره باید متن محاوره‌ای را به رسمی تبدیل کرد. در این مقاله به ارائه و آزمون راهکاری برای تبدیل متون محاوره‌ای به متون رسمی با استفاده از ترکیب روش‌های مبتنی بر قاعده و مدل‌سازی آماری و پیکره بنیان می‌پردازیم.

این مقاله بصورت ذیل سازماندهی شده است: در بخش ۲ به چگونگی ساخت پیکره می‌پردازیم. در بخش ۳ الگوریتم پیشنهادی و نحوه توکن‌بندی و نحوه یافتن ریشه کلمات بیان می‌شود. بخش ۴

منظور از متن محاوره‌ای نوشتمن متن به شکلی است که عمولاً فارسی زبانان به آن شیوه، تکلم می‌کنند و با فارسی رسمی معیار، متفاوت است. تبدیل متن محاوره‌ای به رسمی یکی از مسائل پیش رو برای پردازش زبان طبیعی است که کاربردهای فراوانی در تحلیل و پردازش گفتار گوینده، نظرات^۱ بازدیدکنندگان سایت‌ها، بررسی و بلاغ‌ها، پیام-

های کوتاه و نظایر این‌ها می‌تواند داشته باشد.

لازم به ذکر است که زبان محاوره با زبان عامیانه متفاوت است. برخی از واژه‌ها، نه محاوره‌اند و نه رسمی بلکه صرفًا عامیانه‌اند. کلمات عامیانه از نظر معنایی مبهم و از نظر لحن جالب‌نده و بسته به آواز کلام مردم از آنها استفاده می‌کنند. بسیاری از واژه‌های عامیانه صرفاً از سوی برخی از افراد جامعه به کار می‌روند و تعداد قابل توجهی از مردم از آن بی‌اطلاع هستند. در واقع به تعداد گروه‌های اجتماعی جداگانه انواع واژه‌های عامیانه وجود دارد. به طور مثال واژه‌های عامیانه دانشجویان با دانش آموزان دیبرستان متفاوت است.

در این مقاله تمرکز بیشتر بر شناسایی و تبدیل متون محاوره‌ای است.

Popowich و همکاران در سال ۱۹۹۷ [۲]، روی ترجمه زبان انگلیسی محاوره بررسی نظری و عملی انجام دادند. آنها با روشی لغوی، یک سیستم تمام خودکار برای ترجمه زبان محاوره انگلیسی به اسپانیایی ارائه کردند.

رسمی ما پسوند "مان" را داریم در حالی که در متن محاوره‌ای پسوند "مون" جایگزین آن می‌شود (که در این صورت در متن رسمی "کوچه مان" را به عنوان یک کلمه در نظر می‌گیریم اما در متن محاوره "کوچه مون" به عنوان یک کلمه توکن بندی می‌شود) [4]. در این پروژه از شکل ساده‌ای از قطعه‌بند که تنها به به جداسازی کلمات توسط جداکننده‌هایی مانند: "فاصله" ، "!" ، ":" ، "? ، ";" و امثال آنها بسته شد.

۳ - یافتن ریشه کلمه به کمک ریشه یاب

به علت تفاوت در ساختار کلمات، ریشه‌یاب رسمی در این پروژه قابل استفاده نبود. برای ساخت ریشه‌یاب ابتدا سعی شد با تغییردادن ریشه‌یاب رسمی STeP-1 [5] برای کلمات محاوره ریشه‌یاب جدید ساخته شود. اما استفاده از ریشه‌یاب رسمی به علت پیچیده بودن، ریشه‌یابی کلمات در سطوح پایین و در نظر گرفتن جزئیات زیاد برای این پروژه کار را پیچیده‌تر و زمان اجرا را بیشتر می‌کرد. خصوصاً اینکه برخی از اعمال برای یافتن ریشه عملاً در این پروژه بدون کاربرد بودند. بنابراین با توجه به اینکه ریشه‌یاب محاوره هم از لحاظ پایگاه ریشه‌ها و وندها و هم به جهت قوانین ساختواری مورد استفاده با ریشه‌یاب رسمی متفاوت است بهتر دیدیم ریشه‌یاب تصریفی جدیدی خاص متن محاوره‌ای تهیه نماییم. این ریشه‌یاب با توجه به هدف این پژوهش که تبدیل متن محاوره‌ای به رسمی است ساده شده است. لذا این ریشه‌یاب، در سطح ساختوار تصریفی عمل می‌کند و تنها ادات جمع و "ی" نکره را در اسمای و شناسه‌ها، ضمایر مفعولی (متصل) و پیشوندهای فعلی مانند "نمی" ، "می" ، "ب" و "ن" را در افعال شناسایی و جدا می‌کند [6]. و ریشه‌یابی را تا جایی ادامه می‌دهد که ریشه یا خودش یا معادل رسمی آن در واژگان یافت شود. برای یافتن معادل رسمی از قوانین معرفی شده استفاده می‌شود. و برای برخی که قانون خاصی پیدا نشده معادل رسمی آن از واژگان ساخته شده استخراج می‌شود.

۴ - استفاده از قوانین برای ایجاد معادل رسمی کلمه

یافتن قوانین برای تبدیل کلمات محاوره‌ای به رسمی کار مشکلی است. در زبان محاوره برخی از کلمات تنها به جهت آوازی متفاوت شده‌اند (مانند: دلت (Delat) که در محاوره به دلت (Delete) تبدیل می‌شود) و برخی از کلمات هم به لحاظ آوازی و هم به لحاظ نوشتاری دگرگون گشته‌اند (مانند: خیابان (khiaban) که در محاوره به خیابون (khiabun) تبدیل می‌شود). برخی از کلمات و اصطلاحات تنها در زبان محاوره کاربرد دارند و در گفتار و نوشتار رسمی به کار نمی‌روند (مانند "واسه" "بجای" "برای") و بلعکس برخی از کلمات و اصطلاحات تنها در سبک رسمی کاربرد دارند (مانند "به همین منظور"). تغییرات در زبان محاوره بیشتر به علت راحتی تلفظ و بیان رخ می‌دهد.

مربوط به قوانین تبدیل کلمات محاوره‌ای به رسمی است. در بخش ۵ توضیحی در مورد N_gram و نحوه استفاده از آن در این مقاله بیان شده است و در آخر هم در بخش ۶ به بیان نتایج می‌پردازیم.

۲ - ساخت پیکره

ابتدا برای ساخت پیکره و متن رسمی است، از کتاب‌های با متن محاوره و زیرنویس فارسی چندین فیلم استفاده شد. همچنین از ده‌ها وبلاگ مختلف با حجم بالای ۱۰۰۰ کلمه (همانطور که می‌دانید زبان فارسی چهارمین زبان از لحاظ حجم وبلاگ در اینترنت محاسب می‌شود) متن استخراج شدند. به این ترتیب پیکره‌ای با بیش از ۴۴۰۰ کلمه گرد آوری شد.

سپس به بررسی کلمات محاوره در پیکره موجود و یافتن قواعد احتمالی پرداخته شد. طی بررسی‌های انجام شده بیش از ۵۰ کلمه با حجم بالای ۱۰۰۰ کلمه (همانطور که می‌دانید زبان فارسی چهارمین زبان از لحاظ حجم وبلاگ در اینترنت محاسب می‌شود) متن استخراج شد و برای کلماتی که قانون خاصی برای تبدیل آنها یافت نشد، پایگاه داده‌ای از کلمات محاوره و معادل رسمی آن با بیش از ۵۰ کلمه با تگ POS ساخته شد.

۳ - الگوریتم پیشنهادی

هدف در این الگوریتم یافتن لیست مرتبی از کلمات رسمی است که می‌تواند جایگزین یک کلمه محاوره ای در متن شود (مثال: خونه تبدیل می‌شود به خانه، خون است، خواند). انتظار می‌رود این لیست به ترتیب احتمال وقوع در بافتار متنی مورد نظر مرتب شده باشد. لذا برای پیشنهاد مناسب‌ترین کلمه در هر جمله از N_gram برای اولویت‌دهی به کلمات لیست استفاده می‌شود. الگوریتم پیشنهادی به صورت زیر است:

۱. توکن‌بندی متن محاوره‌ای
۲. جستجوی کلمه در واژگان^۳ زبان رسمی
۳. یافتن ریشه کلمه به کمک ریشه‌یاب^۴ و جستجوی ریشه در واژگان رسمی و واژگان محاوره
۴. استفاده از قوانین تبدیل محاوره به رسمی برای رسمی کردن ریشه کلمه و جستجوی ریشه در واژگان رسمی
۵. جستجو در پایگاه داده‌ی محاوره_رسمی
۶. افزودن وندهای رسمی شده به ریشه
۷. رتبه‌بندی کلمات در لیست پیشنهادی با استفاده از N_gram کلمات متن.

در ادامه برای هریک از مراحل فوق توضیحاتی ارائه خواهد شد.

۴ - توکن‌بندی متن محاوره

به علت آنکه در متن محاوره‌ای ساختار کلمات دارای تفاوت‌هایی با شکل رسمی است، استفاده از قطعه‌بندهای متن رسمی به صورت مستقیم امکان‌پذیر نمی‌باشد و نیاز به ساخت قطعه‌بند جدید برای متن محاوره‌ای یا تغییر قطعه‌بندهای موجود داریم، مثلاً (در متن

ازتون: از شما	ازت: از تو	برای تبدیل زبان محاوره به رسمی قانون و قاعده خاصی وجود ندارد. و چنانچه قاعده‌ای هم بیاپیم، تنها برای برخی از کلمات صدق می‌کند. افعال تغییرات بیشتری نسبت به کلمات دیگر دارند. به ویژه افعال مضارع که در زبان محاوره بیشتر به کار می‌روند.
ازشون: از آنها، از ایشان	ازش: از او، از آن	برای یافتن معادل رسمی یک ریشه از تمام قوانین موجود استفاده می‌شود و لیستی از تبدیلات ممکن برای هر کلمه بدست می‌آید. برای انتخاب درست بین حالات محتمل نیاز به بررسی کلمه در متن هست که بواسیله مدل‌های زبانی مانند N_gram این کار صورت می‌گیرد. برخی از قوانین به کار رفته در این برنامه در ذیل آمده است.
بعدمون: بعد ما	بعدم: بعد من	● ۴ - قواعد تبدیل
بعدتون: بعد شما	بعدت: بعد تو	● برخی از کلمات که دارای "ون" و "وم" هستند. و "آنها را باید به "ا" تبدیل کرد.
بعدشون: بعد او، بعد آن	بعدش: بعد او، بعد آنها، بعد آنها	● کدونه: خانه آسمون: آسمان شونه: شانه آسون: آسان لوونه: لانه
گاهی از اوقات برای اتصال ضمایر متصل بعضی از حروف کلمه حذف و بعضی اضافه می‌شوند.	برامون: برای من، برایمان براتون: برای شما، برایتان براش: برای او، برای آن، برایش براشان، برایشان	● برخی از کلمات به "ا" ختم می‌شوند باید به "ها" و یا "ان" جمع تبدیل کرد.
باهاomon: با ما	باهاam: با من	● درختا: درختها
باهاhtون: با شما	باهاat: با تو	● "و" ای که به آخر کلمات اضافه می‌شود و "رو" در بیشتر مواقع باید به "را" تبدیل شود.
باهاšون: با آنها، با ایشان	باهاš: با او، با آن	● کتاب رو: کتاب را خودشو: خودش را کتابмо: کتابم را ماهانو: ماهان را
"م" در آخر برخی از کلمات به "هم" تبدیل می‌شود.	بازم: باز هم منم: من هم	● حرف "ه" بعد از "ا" در آخر بیشتر کلمات حذف شده است.
● افعال	● همینه: همین است	● نگا: نگاه سیا: سیاه
برخی از افعال حذف می‌شوند.	خوبه: خوب است مبارکه: مبارک باشد خوبی: خوب هستی	● کلا: کلاه
● افعال ماضی	● همینه: همین است	● ضمایر
ماضی ساده و استمراری و بعيد و ملموس: اگر تغییرات در بن فعل رخداده باشد، تغییرات در دوم شخص جمع و سوم شخص جمع رخداده است (فقط اشکال تغییر یافته در زیر آمده).	رفتیں: می‌رفتید رفتیں: رفتید رفتن: رفتند	● تغییر ضمایر متصل:
● می‌رفتیں: می‌رفتید ● می‌رفتن: می‌رفتند	● رفته بودین: رفته بودید ● رفته بودن: رفته بودند	● من: من شما: شما شون: شان
● داشتین می‌رفتیں: داشتید می‌رفتید ● داشتن می‌رفتن: داشتند می‌رفتند	● داشتین می‌رفتیں: داشتید می‌رفتید ● داشتن می‌رفتن: داشتند می‌رفتند	● تو: تو او، اون: او، آن ایشون، اونها، اونا: ایشان، آنها

تغییرات در بن فعل:

ازمون: از ما	ازم: از من	● ضمایر متصل بعد از حروف اضافه به ضمایر منفصل تبدیل می‌شوند.
بهم: به من	بهمون: به ما	● بهم: به من
بهت: به تو	بهتون: به شما	● بهت: به تو
بهشون: به آنها، به ایشان		● بهشون: به آنها، به ایشان

می خواین برین: می خواهید بروید	خوندن: خواندن
می خواد بره: می خواهد برود	ماضی نقلی: ماضی نقلی در محاوره معمولاً به ماضی ساده
می خوان بن: می خواهند بروند	تبدیل می شود، البته در سوم شخص مفرد "است" حذف می شود.
بدون قاعدها	
● "ه" آخر برخی از کلمات به "ر" تبدیل می شود.	رفتیم: رفته ایم
اگه: اگر مگه: مگر دیگه: دیگر	رفتی: رفته ای
در بعضی از کلمات تعدادی از حروف اضافه می شوند.	رفتن: رفته اند
چار: چهار چل: چهل مثه: مثله	ماضی التزامی: تغییرات در سوم شخص مفرد، دوم شخص جمع و سوم شخص جمع رخ داده است.
در برخی از کلمات "ی" ای که برای تلفظ راحتتر اضافه شده است، حذف می شود.	رفته باشم
هیفده: هفده شیش: شش	رفته باشی
کوچیک: کوچک هیجده: هجده	رفته باشیم: رفته باشید
● "چی" برخی اوقات به "چه" و برخی اوقات به "چیز"	رفته باشند: رفته باشد
تبديل می شود.	
همه چی: همه چیز چی کار: چه کار	

۵ اولویت دهی به کمک N_grams

برای ساخت متن رسمی لازم است بین لیست کلمات رسمی محتمل برای یک کلمه محاوره انتخاب صحیح صورت گیرد. برای این کار استفاده از N_gram ها برای یافتن احتمال رویداد کلمه در کنار کلمات دیگر و اولویت دادن به یک کلمه پیشنهاد می شود.

در ابتداء تعریفی کوتاه از N_gram ارائه می شود. یک N_gram یک زیر رشته به طول n از یک رشته از کلمات می باشد. از مفهوم N_gram در حیطه های وسیعی از جمله پردازش زبانهای طبیعی، استفاده می شود. می توان نشان داد که یک N_gram، یک مدل مارکوف از مرتبه $n-1$ می باشد.^[1]

فرکانس تکرار یک N_gram در یک متن تا حدود زیادی می تواند ارتباط آن متن با عبارت مورد نظر را نشان دهد. در این میان حجم سند مورد نظر نیز از اهمیت زیادی برخوردار است. به این معنی که تعداد N_gram ها در واحد حجم سند پارامتر مناسبی محاسبه می شود و تعداد آنها به تنهایی ملاک مطلوبی نیست.^[1]

در این الگوریتم ابتداء از Bi_gram استفاده می شود سپس برای تمایز بین دو احتمال مشابه از Uni_gram استفاده می گردد. برای این کار روی پیکره بیجن خان احتمال وقوع هر کلمه وابسته به کلمه قبلی آن(Bi_gram) را محاسبه کرده، سپس احتمال وقوع هر کلمه (Uni_gram) محاسبه شد. توسط احتمال محاسبه شده، لیست کلمات رسمی محتمل برای هر کلمه محاوره ای، مرتب می شود. نتایج استفاده از آن در بیشتر مواقع قابل توجه بوده است. و اولویت دهی خوبی بین کلمات ایجاد می کند.

مثالاً کلمه خونه دارای لیست احتمالی (خانه، خون است، خواند) است. برای پیدا کردن کلمه جایگزین مناسب آن ابتداء از Bi_gram

موندن: ماندن	خوندن: خواندن
ماضی نقلی: ماضی نقلی در محاوره معمولاً به ماضی ساده تبدیل می شود، البته در سوم شخص مفرد "است" حذف می شود.	●

رفتیم: رفته ایم	رفته ام
رفتی: رفته ای	رفته اید
رفتن: رفته اند	رفته است
ماضی التزامی: تغییرات در سوم شخص مفرد، دوم شخص جمع و سوم شخص جمع رخ داده است.	●
رفته باشم	رفته باشیم
رفته باشی	رفته باشید
رفته باشند: رفته باشد	رفته باشند

افعال مضارع

مضارع ساده و التزامی: اگر تغییرات در بن فعل رخ نداده باشد، تغییرات تنها در شناسه های سوم شخص مفرد، دوم شخص جمع و سوم شخص جمع وجود دارد.

بدون تغییر در بن فعل:

می خورم	می خوریم
می خوری	می خوریم: می خورید
می خورن: می خورد	می خورن: می خورند
همراه تغییر در بن فعل:	

می ریم: می رویم	می رو: می روم
می ری: می روید	
می رن: می روند	

می گم: می گوییم	می گم: می گوییم
می گی: می گویی	می گی: می گویید
می گن: می گویند	می گه: می گویند

می خونه: می خواند	می خونه: می خواند
می ذاره: می گذارد	می آرده: می آورد
بن های مضارعی که به "و" ختم می شوند. "و" در محاوره حذف می شود.	●
بروم: برم	

بگوییم: بگم	بشویم: بشم
افعال مستقبل	

ترکیب تغییرات در مضارع اخباری "خواه" و مضارع الترامی افعال است.

می خوام برم: می خواهم برم	می خوایم برم: می خواهیم برم
می خوایم برم: می خواهیم برم	می خوایم برم: می خواهیم برم

به علت جستجوهای بی‌دریی در واژگان سرعت الگوریتم پایین می‌آید. بعلت آنکه بیشتر جستجوها بر اساس نوع انجام می‌شود. جدا کردن کلمات واژگان به لحاظ نوع در جداول متفاوت می‌تواند به بهبود سرعت الگوریتم بیانجامد.

سپاسگزاری

نویسنده‌گان این مقاله از همفکری تمام اعضای آزمایشگاه پردازش زبان طبیعی و خانم سکینه دهقان برای همیاری در بررسی پیکره و مرتب کردن قوانین کمال سپاسگزاری را دارند.

مراجع

- [1] کاووسی، کاوه، مشیری، بهزاد، طراحی معماری یک عامل هوشمند تطبیقی برای جستجوی اطلاعات تجاری با استفاده از تئوری ترکیب اطلاعات، نشریه جهانی رسانه، دوره: ۴، ۲۰۰۷.
- [2] Popowich, F., Turcato, D., Laurens, O., and McFetridge, P., *A Lexicalist Approach to the Translation of Colloquial Text*, the 7th International Conference on Theoretical and Methodological Issues in Machine Translation, 1997.
- [3] Assadi, Sh., *Sound Deletion in colloquial Persian*, 2007.
- [4] Davarpanah, M.R., Sanji M., Aramideh M., *Farsi lexical analysis and stop word list*, Library Hi, 27 (3), 435-449, 2009.
- [5] Shamsfard M., Jafari, HS., Ilbeygi M., *STeP-I: A Set of Fundamental Tools for Persian Text Processing*, LREC 2010, Malta, 2010.
- [6] Sharifloo, A. Shamsfard, M., *A Bottom up Approach to Persian Stemming*, Proceedings of the Third International Joint Conference on Natural Language Processing, 2007.
- [7] Megerdoomian, K., *Extending a Persian Morphological Analyzer to Blogs*, In: Proceedings of the Second Workshop on Persian Language and Computers, 2006.

زیرنویس‌ها

¹ Comment

² Morphological

³ Lexicon

⁴ Stemmer

استفاده می‌شود. مثلاً اگر کلمه قبلی "می"، باشد احتمال جایگزینی کلمه "خواند" به جای "خونه" بالا است، اما احتمال وقوع کلمه "خانه" بعد از "می" پایین است.

نمونه‌هایی از ورودی و خروجی الگوریتم در جدول (۱) آمده است. همانطور که مشاهده می‌شود، لیست احتمالی هر کلمه بر اساس N_gram مرتب شده است.

جدول (۱) : نمونه خروجی

متن محاوره	تبدیلات رسمی
کتابشو خونه برد.	کتابهایش: اسم. را: حرف اضافه پسین / رو: بن مضارع. خواند: فعل مضارع / خانه: اسم / خون است. برده است / برده: اسم.
مثه اینکه نمیاد.	مثل: اسم. اینکه: قید. نمی‌آید: فعل مضارع.
می‌گوید: فعل مضارع / میگ است. آمد: فعل ماضی. بنشینید: فعل مضارع.	میگه: اومد بشینه.

۶ - نتیجه

در این مقاله به نحوه تبدیل متن محاوره‌ای به رسمی پرداخته شد. ابتدا با استفاده از برخی قوانین کلمه محاوره‌ای به لیستی از کلمات رسمی مبدل می‌شود، سپس بوسیله N_gram لیست احتمالی برای هر کلمه مرتب خواهد شد.

خطای پروژه در مورد قوانین محدود و نسبتاً قابل قبول است. به طور متوسط از هر ۱۰۰ کلمه محاوره ای داده شده به برنامه ۹۳ کلمه درست تشخیص داده شده است. یعنی برنامه دارای دقت متوسط ۹۳٪ روی پیکره آزمایشی است. برخی از خطاهای که برنامه دارد به علت کامل نبودن قطعه‌بند است. و با کامل کردن آن (استفاده از پیکره برای قطعه‌بندی) جواب‌ها بهبود خواهد یافت.

به نظر می‌رسد استفاده از یک الگوریتم یادگیر در کنار قوانینی از این نوع الزامی است. زیرا بیان تمام تغییرات در قالب قانون کاری ناممکن به نظر می‌رسد. و استفاده از یک الگوریتم یادگیر کار را به صورت نیمه خودکار پیش خواهد برد. و به علت اینکه زبان محاوره در طول زمان در حال تغییر است الگوریتم یادگیر می‌تواند قوانین را به مرور تصحیح کند. اما الگوریتم‌های یادگیر نیاز به وجود پیکره‌های موازی محاوره-رسمی و یا پیکره‌های محاوره ای برچسب خورده با برچسب‌های رسمی دارند که در حال حاضر برای زبان فارسی موجود نمی‌باشند.

بررسی بیشتر کلمات محاوره بخصوص به جهت آوایی می‌تواند برای یافتن قوانین بهتر و ساخت مبدل‌هایی از این دست راه‌گشا باشد. زیرا بیشتر تغییرات در کلمات محاوره برای بیان راحت‌تر کلمات بوجود آمداند. پس بررسی نوع تغییر در آوای کلمات قدم بعدی در این رابطه خواهد بود.