

# مدل سازی آماری زبان فارسی

محمد مهدی حاجی

بخش مهندسی و علوم کامپیوتر  
دانشکده مهندسی، دانشگاه شیراز

بهار ۱۳۸۴

## چکیده

این گزارش قسمتی از پژوهش‌های انجام شده در زمینه طرح تحقیقاتی "مدلسازی آماری زبان فارسی" را ارائه می‌دهد. در این گزارش پس از ارائه مقدمه و تاریخچه به توصیف چگونگی طراحی و تنظیم مدل‌های آماری برای کلمات و جملات فارسی پرداخته شده است. در هر مساله مدلسازی در ابتدا دو عامل اصلی می‌باید ابداع و تدوین شوند- ساختار مناسب و جامع برای مدل و سپس چگونگی تنظیم و محاسبه پارامترهای آن. در این تحقیق بر اساس مطالعات و بررسی‌های همه جانبه و با توجه به سوابق و کاربردهای متعدد، مدل پنهان مارکوف برای کلمات و مدل‌های N-gram برای جملات در نظر گرفته شده‌اند. برای تنظیم و محاسبه پارامترهای مدل روش‌های یادگیری هوشمند در ارائه خواهند شد. از آنجا که هر روش یادگیری نیاز به مجموعه‌ای از داده‌های آموزشی دارد لذا روشی کارآمد و جامع برای نمایش متن باید در نظر گرفته شود. در این تحقیق از استاندارد‌های بین‌المللی جهت ذخیره و پردازش فایل‌های متن که شرح آن در این گزارش آمده است استفاده شده است. سپس به چگونگی به‌هنگار سازی و ایجاد یک لغتنامه جامع پرداخته شده است. برای داده‌های آموزشی طیف وسیعی از متون زبان فارسی در زمینه‌های مختلف ادبیات فارسی در نظر گرفته شده است طوری که مدل یاد گرفته شده مبتنی بر توزیع‌های مختلف و وسیعی از داده‌های آموزشی باشد. بدیهی است که در گزارش‌های بعدی چگونگی پیاده‌سازی و آموزش مدل کلمات و جملات و نهایتاً نرم‌افزاری که بتواند در کاربردهای مدلسازی آماری زبان مورد استفاده قرار گیرد ارائه خواهد شد.

## قسمت اول

### مقدمه

یک مدل آماری زبان، احتمال جملات را توصیف می‌کند. یعنی با استفاده از آن میتوان با توجه به اطلاعات جاری، احتمال کلمات بعد را پیش‌بینی کرد. هرچند که مدل‌های دیگری از جمله مدل‌های توانی (Exponential Models) (Rosenfeld et al., 2001) یا گرامرهای مستقل از متن (Context Free Grammars) (Corazza et al., 1993) برای مدل سازی زبانهای طبیعی پیشنهاد شده‌اند، اما در عمل مدل‌های آماری N-gram (Rosenfeld, 2000) از سایر روشها عملی‌تر و موثرتر بوده‌اند. زیرا علاوه بر سادگی پیاده‌سازی، نسبت به نويز هم مقاوم هستند.

کاربردهای عمده مدل آماری زبان در سیستم‌های اتوماتیک تشخیص گفتار (Automatic Speech Recognition) (Katz, 1987) و تشخیص متون ماشینی (OCR) (Lu et al., 1999) یا دستنویس (Handwritten Recognition) (Vinciarelli et al., 2004)

است. با توجه به حجم روزافزون اطلاعات و اهمیت پردازش خودکار گفتار و اسناد توسط کامپیوتر، میتوان به اهمیت مدلسازی زبان پی برد. از کاربردهای دیگر مدل زبان میتوان به این موارد اشاره کرد: استفاده در سیستم های ترجمه متون، پردازش زبانهای طبیعی (NLP) (Manning and Schütze, 1999)، پردازش و تصحیح خودکار لغات در پردازنده های متن، تبدیل متن به صدا (TTS) (Werner et al., 2004) و بهبود نتایج بازیابی اطلاعات متنی در موتورهای جستجو. در همه این کاربردها، عملکرد سیستم با استفاده از مدل آماری زبان بهتر میشود اما در سیستم های بازشناسی گفتار وجود مدل زبان حیاتی تر است؛ در واقع یک سیستم تشخیص گفتار بدون استفاده از مدل های زبان به هیچ وجه نمی تواند دقت قابل قبولی داشته باشد.

با توجه به کاربردهای گسترده مدل های زبان، در سالهای گذشته تحقیقات بسیاری برای مدل سازی زبانهای پراستفاده جهان و بیش از همه زبان انگلیسی انجام شده است و این تحقیقات همچنان ادامه دارند. جالب توجه است که مدل های آماری زبان انگلیسی تقریباً از سال ۱۹۸۰ در سیستم های واقعی به کار رفته اند (Rosenfeld, 2000) اما برای مدل سازی زبان فارسی متاسفانه هنوز هیچ تحقیق گسترده و قابل اعتنایی صورت نگرفته است. البته با توجه به اینکه در کشور ما روند اتوماسیون عمر زیادی ندارد و عمده کاربردهای مدل آماری زبان در سیستم های اتوماتیک پردازش و تشخیص گفتار و متن است، کمبود تحقیق در این زمینه کاملاً قابل انتظار است. هدف اصلی این تحقیق و اولین گزارش آن که در پی می آید، جبران این کاستی و پایه ریزی اصول اولیه است. امید است که این پژوهش انگیزه و راهگشای کارهای بزرگتر بعدی باشد.

## محتوای گزارش

در این گزارش پس از مقدمه و تاریخچه به بررسی چگونگی تهیه و تدارک داده های آموزشی برای یادگیری و تنظیم پارامترهای مدل مورد نظر پرداخته شده است. همانگونه که در شرح آمده است داده های آموزشی از طیف وسیعی از متون فارسی در زمینه ها و موضوعات مختلف جمع آوری شده اند. سپس روش به هنجارسازی این انبوه داده ها و روش استاندارد بین المللی یونی کد برای نمایش متون مورد بحث قرار گرفته است. در قسمت بعدی چگونگی ایجاد یک لغتنامه جامع و کامل از دیدگاه آماری پرداخته شده است. مدل های آماری مورد نظر در این تحقیق، مبتنی بر مدل های پنهان مارکوف برای کلمات و مدل های N-gram برای جملات میباشد که پارامترهای آنها بر اساس الگوریتم های یادگیری که شرح داده خواهد شد محاسبه و تنظیم خواهند شد. در این گزارش مدل آماری مارکوف به طور کلی و مدل پنهان مارکوف به خصوص معرفی و شرح داده شده است.

## داده های آموزشی

همانند هر مدل آماری دیگر در این سیستم نیز به داده هایی برای آموزش یا در واقع تنظیم پارامترهای آزاد مدل احتیاج داریم. برای تهیه یک مدل آماری، داده های آموزشی باید دارای توزیع احتمالی باشند که بعداً مدل با آنها سروکار دارد؛ یعنی به طور کلی داده های آموزشی باید فضای احتمال را پوشش دهند. یعنی اگر بخواهیم مدل آماری به دست آمده در همه جا قابل استفاده باشد، واضح است که باید انواع مختلف متون - از جمله ادبی، سیاسی، اقتصادی، طنز، عقیدتی و ... - را برای آموزش سیستم به کار گیریم.

## تهیه داده های آموزشی

برای تهیه داده های آموزشی با چنین حجم زیادی تنها راه عملی استفاده از متون موجود در شبکه اینترنت است. داده های آموزشی که در این تحقیق به کار گرفته شده اند، شامل (۱) متون عمومی: مجلات سروش (که ۹۵٪ حجم کل داده های آموزشی را به خود اختصاص داده اند) و بعضی از مقالات موجود در سایت الشیعه؛ (۲) متون تاریخی: از جمله کلیله و دمنه و (۳) اشعار: از جمله دیوان حافظ و برخی از دفاتر شعر نیما یوشیج و فروغ فرخ زاد. حجم کل نوشته جات آموزشی که در نهایت به دست آوردیم، در حدود ۱۰۰ مگابایت (با فرمت یونی کد) است شامل حدود ۷۳۰۰۰۰۰ کلمه می باشد. لازم به توضیح است که مجلات سروش شامل نشریات سروش کودکان، سروش نوجوان، سروش جوان، سروش بانوان، سروش اندیشه و هفته نامه سروش میباشد که این نشریات شامل انواع مختلف متون فارسی هستند. در زیر از هر یک از این نشریات بخشهایی آورده شده است:

### سروش کودکان:

خروس به روی بام ویرانه ای پرید و قوقولی قو سر داد. روباه از فرصت استفاده کرد. و به سر خروس پرید و خروس را به دندان گرفت و فرار کرد...  
این طوری بود که کورینلیوس تصمیم گرفت کروکودیل ها را به حال خودشان بگذارد و از کنار رودخانه دور شود اما هنوز راه زیادی نرفته بود که به یک میمون برخورد کرد...

### سروش نوجوان:

اتفاق هایی هم که در کارتون ها می افتد، خیلی وقت ها اتفاق های واقعی نیستند یا اگر اتفاق های واقعی اند، نتایجی غیرمعمول دارند؛ مثلاً گربه ای پایش به چیزی گیر می کند و به جای این که زمین بخورد، بدنش کش می آید...

سروش جوان:

زیبایی آمریکایی بیشترین جوایز اسکارامسال را برد. چیزی که از قبل میشد حدس زد در میان برندگان جوایز، نام کوین اسپیسی دیده می‌شود... گزارشی که باعث می‌شود همان روز عصر تمام آن مزرعه و کشتارگاه قرنطینه و تمام سی صدوچهل خوک و سی ودو گاو موجود در آن‌ها معدوم شوند...

سروش بانوان:

بعد مردم می‌گویند چرا زندگی‌های الان این قدر بی‌دوام است؟ چرا بین زن و شوهرها جنگ و دعواست؟ چرا فلان است، چرا بهمان است؟...

سروش اندیشه:

مکتب فلسفه جاودانه (حکمت خالده) با دیدگاهی همبسته است که مطابق آن دانشی آغازین، یا دانشی قدسی وجود دارد که بر متافیزیک عام این حقیقت‌نمایی مبتنی است...

هفته نامه سروش:

سفر قندهار، فیلمی درباره مردم افغانستان است... هرچه زیبگانه و خیل‌تواند جمله در این خانه طفیل تواند... رضا عطاران را حتما می‌شناسید، بازیگر مجموعه‌های طنز...

همانطور که در این نمونه‌ها مشاهده میشود، مجموعه متون آموزشی شامل اسم‌های خاص ایرانی (رضا عطاران) و غیر ایرانی (کوین اسپیسی)، عدد (صدوچهل) و ترکیبات تخصصی (حکمت خالده) نیز هست. بنابراین کاملاً قابل انتظار است که لغتنامه فارسی که از این متون استخراج خواهد شد، لغات بسیار زیادی داشته باشد که بسیاری از آنها ممکن است در لغتنامه‌های عادی وجود نداشته باشند.

## نرمال سازی فایل‌های متنی

فایل‌هایی متنی که از اینترنت جمع‌آوری کرده ایم، به فرمت HTML هستند. در این فایل‌ها، کاراکترهای فارسی با سیستم‌های کدگذاری متفاوتی (Unicode، Cp1256 و UTF8) نشان داده میشوند. در اولین مرحله نرمال سازی، فایل ورودی را با هر سیستم کدگذاری به فایلی یونی‌کد تبدیل میکنیم. برای این کار یک برنامه جاوا نوشته شده است. در زیر ابتدا توضیحاتی

درباره استاندارد یونی کد ارائه میکنیم و سپس نرمال سازی های بیشتری را که روی متون انجام میشود شرح خواهیم داد.

## استاندارد یونی کد

استاندارد یونی کد، استاندارد جهانی کدگذاری نویسه هاست که برای ارائه متون برای پردازش کامپیوتری به کار میرود. این استاندارد با ویرایش دوم استاندارد بین‌المللی SO/IEC 10646-1:2000 کاملاً سازگار است و همان نویسه‌ها و کدهای ISO/IEC 10646 را دارد. استاندارد یونی کد اطلاعات بیشتری نیز در مورد نویسه‌ها و کاربردهایشان فراهم کرده است، پس در واقع هر پیاده‌سازی سازگار با یونی کد، با ISO/IEC 10646 نیز سازگار است و روش رسمی پیاده‌سازی استاندارد ISO/IEC 10646 یونی کد است.

استاندارد یونی کد را پیش‌تازان صنعت کامپیوتر، شرکت‌هایی چون اوراکل، آی بی ام، اپل، سان، سای بیس، مایکروسافت، هیولت پاکارد، و بسیاری دیگر پذیرفته‌اند. استانداردهایی چون XML، جاوا، اکماسکریپت (جاواسکریپت)، LDAP، CORBA 3.0، WML، و غیره، یونی کد را ملزم می‌دانند. یونی کد در بسیاری از سیستم‌عامل‌ها، همه‌ی مرورگرهای اخیر، و بسیاری از محصولات دیگر پشتیبانی می‌شود. یونی کد مستقل از محیط، مستقل از برنامه، و مستقل از زبان به همه نویسه‌ها اعداد یکتایی اختصاص میدهد. تا قبل از اختراع یونی کد، صدها سیستم کدگذاری مختلف برای تخصیص این اعداد وجود داشت. نویسه‌های هیچ کدگذاری‌ای به تنهایی کافی نبود: مثلاً اتحادیه اروپا به چندین کدگذاری مختلف برای پوشاندن همه‌ی زبان‌هایش نیاز داشت. حتی برای زبانی مثل انگلیسی نیز هیچ کدگذاری‌ای به تنهایی برای همه حروف، علائم نقطه‌گذاری، و نمادهای فنی متداول کافی نبود. این سیستم‌ها با هم تعارض نیز داشتند. یعنی، دو کدگذاری ممکن بود از اعداد یکسان برای دو نویسه مختلف، یا از اعداد مختلف برای نویسه‌های یکسان استفاده کنند. کامپیوترها، خصوصاً سرورهای شبکه، از کدگذاری‌های مختلف و متعددی پشتیبانی می‌کردند و بنابراین هرگاه داده‌ها از کدگذاری‌ها یا محیط‌های مختلف عبور می‌کردند، در معرض خطر تحریف قرار می‌گرفتند.

یونی کد امکان کدگذاری همه نویسه‌های مورد استفاده در نوشتن زبان‌های دنیا را فراهم می‌سازد. این استاندارد از یک کدگذاری ۱۶ بیتی استفاده می‌کند که برای بیش از ۶۵۰۰۰ نویسه جا فراهم می‌کند. اگر چه ۶۵۰۰۰ نویسه برای کدگذاری اکثر نویسه‌هایی که در زبان‌های مهم دنیا استفاده می‌شود کافی است، با این وجود، یونی کد و ISO 10646 شیوه‌گسترشی به نام UTF-16 فراهم کرده‌اند که امکان اضافه کردن حدود یک میلیون نویسه دیگر را نیز فراهم

میکند. این دامنه برای کلیه نویسه‌های عالم، از جمله پوشش کامل همه خط‌های باستانی نیز کافی است.

یونی کد برای کلیه‌های نویسه‌های مورد استفاده در زبانهای عمده دنیا کد تعیین کرده است. به علت فراخ بودن فضای تخصیص نویسه، این استاندارد بسیاری از نمادهای لازم برای حروفچینی با کیفیت بالا را نیز در بر گرفته است. از خط‌های مورد پشتیبانی این استاندارد میتوان به لاتین (در بر گیرنده اکثر زبان‌های اروپایی)، سیریلیک (روسی، صربی)، یونانی، عربی (شامل عربی، فارسی، اردو، کردی)، عبری، هندی، ارمنی، آسوری، چینی، کاتاکانا و هیراگانا (ژاپنی)، و هانگول (کره ای) اشاره کرد. به علاوه، تعداد زیادی نماد ریاضی، علائم نقطه گذاری و علامت‌های متفرقه در این استاندارد وجود دارد. این استاندارد برای علامت‌های ترکیب شونده یا اعراب‌ها نیز کدهایی در نظر گرفته است که از جمله آنها علامت‌هایی چون «~» هستند که در ترکیب با حروف پایه، حروف تغییر لحن یافته ای چون «ñ» را می سازند. آخرین نسخه یونی کد، در مجموع، ۴۹۱۹۴ نویسه دارد. به علاوه، ۶۴۰۰ کد نیز برای مصرف خصوصی در نظر گرفته شده است که برنامه‌نویسان میتوانند از آنها برای نویسه‌ها و نمادهای خودشان استفاده کنند.

به طور کلی، اصول یونی کد به این شرح می باشند: نویسه‌های شانزده‌بیتی، ترتیب مفهومی (در مقابل دیداری)، کارایی، یکی سازی (اختصاص یک کد به نویسه‌های مشترک در چند زبان مختلف)، نویسه نه شکل (یک «ع»، و نه چهارتا: «ع»، «ع»، «ع»، «ع»)، ترکیب پویا، بار معنایی (حرف بودن، مقدار عددی و ...)، دنباله‌های معادل (امکان ذخیره‌سازی یک متن به دو شکل مختلف ولی معادل)، متن ساده (و نه مفاهیمی مثل تغییر قلم، جدول‌بندی، و صفحه آرایی)، قابلیت تبدیل (هر متن موجود در قالب یک مجموعه‌نویسه سنتی باید بدون از بین رفتن معنا قابل تبدیل به یونی کد باشد).

یونی کد شیوه‌ای نیز برای کدگذاری ۸ بیتی متون مشخص کرده است که نویسه‌ها را پس از اعمال یک تابع خاص به کدشان، به صورت دنباله‌هایی که از یک تا چهار بایت دارند نگه می دارد. این شیوه که با نام UTF-8 شناخته می شود، به این خاطر که نویسه‌های ASCII را عیناً حفظ میکند و در نتیجه، هم برنامه‌های قدیمی راحت‌تر با آن کنار می آیند و هم طول پرونده‌های لاتین را زیاد نمی کند، بسیار محبوب است. در واقع بسیاری از سیستم‌هایی که ادعای پشتیبانی یونی کد را می کنند، تنها UTF-8 را پشتیبانی می کنند و پرونده‌های یونی کدی، اعم از کاربردهای اینترنتی یا موارد استفاده محلی، عمدتاً در قالب UTF-8 ذخیره شده اند. در استاندارد یونی کد، نویسه‌های فارسی در بلوک مربوط به خط عربی قرار دارند. این بلوک برای دربرگرفتن نویسه‌های زبان‌هایی که از خط عربی استفاده می کنند، مثل فارسی، اردو، پشتو، سندی، و کردی گسترش یافته است. این بلوک نشانه‌های قرآنی از قبیل نشانه‌های سجده و پایان آیه، و علائم وقف را نیز در بر دارد.

در یونی کد با وجود یکی سازی کدهای حروف مشترک، برای حروف فارسی ای که بار معنایی یا نمایشی متفاوت با حروف عربی دارند، نویسه های جداگانه در نظر گرفته شده است. یعنی کلیه حروف خاص فارسی (پ، چ، ژ، گ) و نیز «ک» و «ی» ی فارسی که با حرف مشابه در عربی تفاوت نمایشی دارند، مکان جداگانه ای به خود اختصاص داده اند. کلیه اعراب های متداول حضور دارند و میان شکل فارسی/ اردو و عربی ارقام نیز به علت شکل و رفتار متفاوت تفاوت هایی منظور گشته است.

از طرف دیگر، علائم نقطه گذاری ای چون نقطه و فاصله که شکل یکسانی در خط های لاتین و عربی دارند، کد یکسان دارند. علائمی چون پرانتز نیز، بسته به جهت متن، آینه ای میشوند، یعنی به طور مثال، نویسه ۰۰۲۸ نماینده «پرانتز باز» است، و نه «پرانتز سمت چپ». یونی کد اتصال مجازی و فاصله مجازی را نیز تحت نامهای «اتصال با عرض صفر» و «بی اتصالی با عرض صفر» به رسمیت می شناسد. به خاطر سازگاری با استاندارد های موجود در بعضی از کشورهای عربی، ISO 10646 و نتیجتاً یونی کد بلوک جداگانه ای را نیز به شکل های مختلف حروف عربی اختصاص داده است که استفاده از آنها شدیداً منع شده است. این بلوک معمولاً فقط برای تعیین جای شکل های مختلف حروف در قلمها به کار میرود.

همچنین این استاندارد توضیحات مفصل و دقیقی درباره شیوه های پیاده سازی، از جمله شیوه «متصل سازی حروف» و «نمایش متون راست به چپ و دوجهته» دارد که برنامه نویس را از مراجعه به راهنمای محلی بی نیاز می سازد. خواننده میتواند برای توضیحات بیشتر در این زمینه به مستندات یونی کد مراجعه کند ([www.unicode.org](http://www.unicode.org)).

## نرمال سازی های بیشتر

با توجه به مطالبی که درباره استاندارد یونی کد گفتیم، میتوان دریافت که گاهی ممکن است کلمات یکسان با مجموعه ای از نویسه های غیر یکسان نمایش داده شوند. به عنوان مثال کلمه «یک» میتواند چهار شکل نمایش مختلف داشته باشد، که از ترکیب دو «ی» (فارسی و عربی) و دو «ک» (فارسی و عربی) به دست می آیند. بنابراین، در اولین مرحله نرمال سازی، هر حرفی را که در استاندارد یونی کد با دو یا چند نویسه مختلف نمایش داده میشود، فقط با یکی از این نویسه ها نشان می دهیم؛ به عنوان مثال هر جا «ک» عربی باشد، به جای آن «ک» فارسی میگذاریم. علاوه بر این، حرف «-» (تطویل) هم که برای کنترل طول کلمات استفاده میشود و هیچ اطلاعات مفیدی را دربر ندارد برای پردازش های بعدی حذف میکنیم.

مرحله بعدی، حذف اعراب کلمات است، با وجود اینکه گاهی تفاوت دو کلمه فقط با اعراب مشخص میشود، اما چون در اکثریت نوشته جات آموزشی که از اینترنت جمع آوری شده اند، کلمات اعراب گذاری نشده اند، مجبور هستیم اعراب کلمات را به طور کلی حذف کنیم. به بیان دیگر، کلمات یا باید همواره دارای اعراب باشند، یا همواره دارای اعراب نباشند و چون کلمات



موجود در نوشته جات آموزشی همواره دارای اعراب نیستند، پس تنها راه برای نرمال سازی این است که کلمات هیچ اعرابی نداشته باشند.

باید توجه داشت که قواعد نرمال سازی که در بالا ذکر کردیم، اگرچه بسیاری از موارد را پوشش می دهند اما هنوز حالت‌هایی وجود دارند که کلمات یکسان با مجموعه ای از نویسه های غیر یکسان نمایش داده میشوند که عمدتاً به علت غیراستاندارد بودن شیوه نگارش است. به عنوان مثال، بعضی از نویسندگان ها ترجیح می دهند «ء» آخر کلمات را ننویسند و به این ترتیب مثلاً «اشیاء» و «اشیا» از نظر سیستم ما، که یک مجموعه از حروف فارسی پشت سر هم را به عنوان یک کلمه در نظر میگیرد، دو کلمه متفاوتند. به طور مشابه، در بعضی از نوشته جات آموزشی، به جای «ة» آخر کلمات از «ه» استفاده شده است و به این ترتیب مثلاً «حمیده» و «حمیده» دو کلمه متفاوت میشوند. اما مشکل بزرگتر، سر هم یا جدا نوشتن کلماتی است که از بخشهای مختلفی درست شده اند که این مشکل عمدتاً در نگارش افعال استمراری فارسی و جمع بستن کلمات با «ها» پیش می آید. به عنوان مثال، با توجه به تعریف کلمه از نظر سیستم ما، «می روم» شامل دو کلمه «می» و «روم» است، اما «میروم» یک کلمه متفاوت است. و به همین ترتیب، «درخت ها» شامل دو کلمه «درخت» و «ها» است اما «درختها» یک کلمه متفاوت است.

نکته قابل ذکر دیگر این است که سیستم ارائه شده چون رویکردی کاملاً آماری دارد، تلاشی برای ریشه یابی لغات انجام نمیدهد و بنابراین به عنوان مثال «کتابم»، «کتابت»، «کتابش» و ... از نظر این سیستم، کلمات متفاوتی هستند. این مساله ناشی از اتصال ضمائر ملکی به انتهای کلمات است.

## استخراج لغتنامه فارسی

پس از نرمال سازی مجموعه نوشته جات (داده های) آموزشی، می توان لغتنامه فارسی را استخراج کرد. بار دیگر یاد آوری می کنیم که در این مدل، یک دنباله پشت سر هم از نویسه ها را به عنوان یک کلمه در نظر می گیریم. همانطور که قابل انتظار است، مجموعه نوشته جات آموزشی شامل نویز هم می باشد، مثلاً وقتی که یک کلمه اشتباه تایپ شده باشد. گاهی اوقات هم بسته به نوع متن، کلمات محاوره ای در متن وجود دارند؛ مثلاً «بتونیم» به جای «بتوانیم». و همانطور که قبلاً ذکر شد، کلمات خاص زبانهای خارجی و اعداد هم کلماتی از مجموعه نوشته جات هستند. برای کاهش دادن این کلمات ناخواسته، یک روال ساده را در پیش می گیریم: فرکانس تکرار تمامی کلمات موجود در مجموعه نوشته جات را محاسبه میکنیم و سپس لغاتی را که فرکانس تکراری کمتر از یک حد آستانه (از پیش تعیین شده) دارند حذف میکنیم. در این تحقیق مقدار آستانه را ۳ قرار دادیم، یعنی در واقع کلماتی را که کمتر از ۳ بار پیش آمده بودند، به عنوان نویز در نظر گرفته و حذف کردیم. با این رویه، یک لغتنامه شامل

حدود ۴۱۰۰۰ کلمه به دست آمد. این مجموعه بر حسب فرکانس تکرار مرتب شده است، تا اگر برای کاربردی خاص مثلا به هزار لغت پر فرکانس تر (پر کاربردتر) فارسی احتیاج است بتوان به راحتی از هزار لغت ابتدای این مجموعه استفاده کرد. ده عضو ابتدای این لغتنامه به ترتیب اینها هستند: و، می، به، در، که، از، را، این، است، با. همانطور که قابل انتظار است پر فرکانس ترین لغت فارسی حرف ربط «و» است. بدیهی است که قرار گرفتن کلمه «می» به عنوان دومین کلمه لغتنامه به دلیل وجود تعداد تکرارهای کلمه «می» نیست بلکه به دلیل تعداد تکرار زیاد «می» در بخش ابتدایی افعال استمراری (مانند می روم) است که به شکل غیر چسبان نوشته (تایپ) شده است. حتی اگر در متون اولیه تعدادی کلمه «می» هم وجود داشته باشد، چون در مرحله نرمال سازی، اعراب کلمات حذف می شوند، اینها عملا با «می» فرقی نخواهند داشت و به افزایش تعداد تکرارهای کلمه «می» کمک خواهند کرد. نکته جالب توجه این است که پرفرکانس ترین لغات فارسی (یعنی اعضای نخست لغتنامه) حروف ربط هستند و به طور مشابه در زبان انگلیسی پرفرکانس ترین کلمه The است و کلمات پرفرکانس بعدی حروف ربط انگلیسی از جمله of, and, to هستند.

در زیر بعضی از کلمات لغتنامه ۴۱۰۰۰ کلمه ای را مشاهده می کنید:

شانگهای	کورتکس	نوزدهم	اینتراکتیو
صدوچهل	بودیسم	جدیدترین	آفساید
سیدمحمدرضا	انتگرال	یورونیوز	نوددرصد
پدرومادرها	ارتودنسی	پرتابل	بلندبلند
ازتهران	توتیا	پراهمیت	دوقسمتی
گیگز	مغیلان	اسکیزوفرنی	اوپک
کامیونیکیشن	تغار	شیشلول	مستثناست
بتونیم	هزارتو	امپایر	مستظرفه

کاملا واضح است که علیرغم تلاشی که برای حذف لغات نويز کردیم، هنوز لغات ناخواسته ای در لغتنامه وجود دارند. البته قابل انتظار است که با افزایش حجم نوشته جات آموزشی وبلا بردن حد آستانه ی حداقل فرکانس تکرار، کلمات ناخواسته کاهش یابند. اما نکته مثبت این لغتنامه، داشتن لغاتی است که در لغتنامه های عادی وجود ندارند و حتی ممکن است لغات فارسی یا عربی هم نباشند، مثلا انتگرال و اینترنت، اما به دلیل استفاده زیاد واقعا بخشی از زبان فارسی جدید شده باشند. اشکالات تایپی رایج و حتی بی معنی بودن برخی اعضای لغتنامه استخراج شده اهمیت چندانی ندارد؛ زیرا مدل آماری بعدا با متونی شامل همین لغات سروکار خواهد داشت.

## مدلهای مارکوف

چون مدلی که برای کلمات ارائه می دهیم، بر مبنای مدل‌های مارکوف است، در این بخش مقدمه ای بر مدل‌های مارکوف و سپس مدل‌های پنهان مارکوف ارائه می دهیم. خروجی یک فرآیند در جهان واقعی به شکل یک سیگنال پیوسته یا گسسته مشاهده می شود. یک مساله حیاتی در علوم ساختن مدل‌هایی برای این سیگنال‌های واقعی است. مدل سازی یک سیگنال مزایای فراوانی به همراه دارد. اولاً، مدل، پایه ای برای توصیف تئوری سیگنال فراهم میکند که میتواند برای پردازش سیگنال استفاده شود تا خروجی خواص مطلوبی داشته باشد. ثانیاً، مدل میتواند اطلاعات بسیار مفیدی درباره منبع سیگنال بدهد، بدون اینکه احتیاجی به خود منبع باشد. نهایتاً و از همه مهمتر، مدلها میتوانند در عمل به خوبی کار کنند و امکان تحقق سیستم های عملی مهمی را فراهم می آورند.

بسته به نوع و خواص سیگنال، راه های مختلفی برای مدل کردن یک سیگنال وجود دارد. به طور کلی، یک سیگنال میتواند معین یا نامعین (تصادفی یا آماری) باشد. مدل‌های معین از بعضی خواص شناخته شده سیگنال استفاده می کنند و مقادیر پارامترهای مدل را تخمین می زنند. در طرف دیگر، در مدل‌های آماری، یک فرآیند تصادفی، سیگنال را توصیف میکند. برای کاربرهایی نظیر تشخیص گفتار یا دستخط که با نویز و عدم قطعیت همراه هستند، مدل‌های آماری از کارایی بهتری برخوردارند. مدل‌های پنهان مارکوف (Rabiner, 1989)، که همچنین منابع مارکوف یا توابع آماری زنجیره های مارکوف نامیده می شوند، در تئوری مخابرات یکی از پرکاربرترین مدل‌های آماری هستند.

دسته مهمی از فرآیند های تصادفی، فرآیندهای مارکوف است که دارای خواصی است که مطالعه ریاضی آنها را امکان پذیر میکند. در جهان واقعی، معمولاً مطلوب است که یک دنباله متغیرهای تصادفی وابسته را - که مقدار هر متغیر به عنصر (یا عناصر) قبلی در دنباله بستگی دارد - بررسی کنیم. در یک فرآیند مارکوف مقدار متغیر تصادفی جاری برای پیش بینی مقدار متغیرهای تصادفی آینده کافی است. به بیان دیگر، وقتی عنصر جاری را داشته باشیم، عناصر آینده از عناصر گذشته مستقل شرطی اند. فرض کنید  $X = (X_1, \dots, X_T)$  دنباله متغیرهای تصادفی باشد که مقادیری از فضای محدود  $S = \{s_1, \dots, s_N\}$  میگیرند. خواص مارکوف با دو رابطه زیر بیان میشوند:

$$(1) \quad P(X_{t+1} = s_k | X_1, X_2, \dots, X_t) = P(X_{t+1} = s_k | X_t)$$

$$(2) \quad P(X_{t+1} = s_k | X_t) = P(X_2 = s_k | X_1)$$

خاصیت دوم را تغییر ناپذیری با زمان می نامیم. اگر دنباله  $X$  هر دو خاصیت مارکوف را داشته باشد، آن را یک زنجیره مارکوف می نامیم.

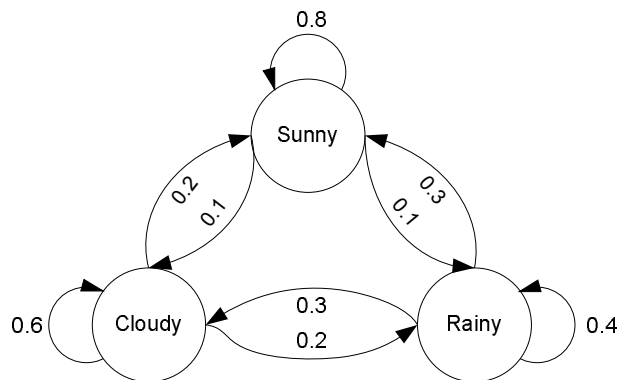
یک زنجیره مارکوف را می توان با بردار تصادفی وضعیت اولیه  $\Pi$  و ماتریس تصادفی انتقال  $A$  به طور کامل توصیف کرد:

$$(۳) \quad \pi_i = P(X_1 = s_i)$$

$$(۴) \quad a_{ij} = P(X_{t+1} = s_j | X_t = s_i)$$

که  $\pi_i \geq 0, \forall i$  و  $\sum_{i=1}^N \pi_i = 1$  و  $a_{ij} \geq 0, \forall i, j$  و  $\sum_{j=1}^N a_{ij} = 1, \forall i$ .

برای روشن کردن این مطالب، مثالی درباره پیش بینی وضع هوا را در نظر بگیرید که در آن با استفاده از تاریخچه مشاهدات وضع هوا در گذشته می خواهیم هوای فردا را حدس بزنیم. برای سادگی فرض می کنیم هوا سه حالت بیشتر نداشته باشد: خورشیدی، ابری و بارانی و وضعیت هوا در طول یک روز یکسان است؛ یعنی تغییر حالتی در وسط روز اتفاق نمی افتد. اگر فرض کنیم زنجیره وضعیت هوا در روزهای متوالی مارکوف است (که البته در جهان واقعی فرض درستی نیست)، آنگاه ماشین حالت محدود شکل ۱ با احتمالات انتقال (تغییر وضعیت) دلخواهی که روی پیکان ها داده شده است این زنجیره مارکوف را نشان می دهد. توجه کنید که مجموع احتمالات پیکانهای خروجی در هر یک از سه حالت ۱ است. از شکل ۱ واضح است که مدل مارکوف را می توان به عنوان یک ماشین حالت محدود نامعین (تصادفی) به شمار آورد که در آن احتمالات روی پیکانهای تغییر وضعیت قرار گرفته اند.



شکل ۱ - مدل مارکوف مثال پیش بینی وضع هوا

فرض کنید  $s_1 = Sunny$ ،  $s_2 = Cloudy$  و  $s_3 = Rainy$  باشد و در اولین روز هوا خورشیدی باشد. آنگاه:

$$\Pi = (1.0, 0.0, 0.0)$$

$$A = \begin{bmatrix} 0.8 & 0.1 & 0.1 \\ 0.2 & 0.6 & 0.2 \\ 0.3 & 0.3 & 0.4 \end{bmatrix}$$

محاسبه احتمال یک دنباله از وضعیتها  $X_1, \dots, X_K$  برای یک زنجیره مارکوف به سادگی با رابطه زیر انجام میشود:

$$\begin{aligned} P(X_1, \dots, X_K) &= P(X_1) P(X_2 | X_1) P(X_3 | X_1, X_2) \dots P(X_K | X_1, \dots, X_{K-1}) \\ &= P(X_1) P(X_2 | X_1) P(X_3 | X_2) \dots P(X_K | X_{K-1}) \\ (5) \quad &= \pi_{X_1} \prod_{t=1}^{K-1} a_{X_t, X_{t+1}} \end{aligned}$$

بنابراین در مثال بالا احتمال اینکه هوا در هفت روز آینده به ترتیب خورشیدی، خورشیدی، بارانی، بارانی خورشیدی، ابری و خورشیدی باشد یا در واقع احتمال  $O = s_1, s_1, s_3, s_3, s_1, s_2$   $s_1$  میتواند به شکل زیر محاسبه شود:

$$\begin{aligned} P(O | \text{Model}) &= \pi_{s_1} P(s_1 | s_1) P(s_1 | s_1) P(s_3 | s_1) P(s_3 | s_3) P(s_1 | s_3) P(s_2 | s_1) P(s_1 | s_2) \\ &= \pi_1 a_{11} a_{11} a_{13} a_{33} a_{31} a_{12} a_{21} \\ &= 1.0 (0.8) (0.8) (0.1) (0.4) (0.3) (0.1) (0.2) \\ (6) \quad &= 1.536 \times 10^{-4} \end{aligned}$$

به طور کلی وقتی از مدل‌های مارکوف صحبت می‌کنیم، منظور ما مدل‌های مارکوف مرتبه اول است که در آن تاریخچه ای به طول ۱ (یعنی یک عنصر قبلی) برای پیش بینی رفتار آینده استفاده میشود. اما گاهی اوقات برای پیش بینی حالت‌های آینده تاریخچه بزرگتری لازم است. در یک مدل مارکوف مرتبه  $n$ ، برای پیش بینی حالت بعدی، از  $n$  حالت قبلی استفاده می‌شود. اما همواره می‌توان با تغییر شکل نمایش فضای حالت، هر مدل مارکوف مرتبه  $n$  را به یک مدل مارکوف مرتبه یک تبدیل کرد. بنابراین از نظر تئوری، فرض مارکوف مرتبه اول، محدود کننده نیست.

## قسمت دوم

### مدل‌های پنهان مارکوف

مدل‌های پنهان مارکوف (Hidden Markov Models)، که به اختصار HMM نامیده میشوند، یکی از قوی‌ترین ابزارها برای پردازش سیگنال‌ها میباشند. انواع مختلف HMM علیرغم محدودیتهایی که دارند، هنوز پر استفاده‌ترین تکنیک در سیستم‌های مدرن بازشناسی گفتار و تشخیص متون هستند. مدل پنهان مارکوف، کل الگوی ورودی را به عنوان یک بردار

ویژگی تکی مدل نمیکند، بلکه رابطه بین بخشهای متوالی یک الگو را استخراج میکند، زیرا هر بخش نسبت به کل ورودی کوچکتر و بنابراین مدل سازی آن ساده تر است. یک HMM را در واقع میتوان یک ماشین حالت محدود (Finite State Machine) احتمالاتی به حساب آورد که هر حالت با یک تابع تصادفی مرتبط است. فرض میشود که در یک دوره گسسته از زمان  $t$ ، مدل در یک حالت است و با یک تابع تصادفی از آن حالت خروجی ای (مشاهده ای) تولید میکند. بر مبنای تابع احتمال انتقال حالت جاری، مدل مارکوف در زمان  $t+1$  تغییر حالت میدهد. دنباله حالتی که مدل از آن میگذرد معمولاً پنهان است، و تنها یک تابع احتمالاتی از آن آشکار است، که مشاهدات تولید شده به وسیله تابع تصادفی مربوط به حالتها است، به همین دلیل در نام گذاری این مدلها از صفت "پنهان" استفاده شده است. یک HMM را در واقع میتوان یک فرآیند تصادفی به طور ناتمام (جزئی) مشاهده شده در نظر آورد. یک HMM با عناصر زیر توصیف میشود:

(۷)  $N$ : تعداد حالتهای مدل

(۸)  $S = \{s_1, s_2, \dots, s_N\}$ : مجموعه حالتها

(۹)  $\Pi = \{\pi_i = P(s_i \text{ at } t = 1)\}$ : احتمالات حالت اولیه

(۱۰)  $A = \{a_{ij} = P(s_j \text{ at } t+1 | s_i \text{ at } t)\}$ : احتمالات تغییر حالت

(۱۱)  $M$ : تعداد علائم قابل مشاهده (تولید شده)

(۱۲)  $V = \{v_1, v_2, \dots, v_M\}$ : مجموعه علائم قابل مشاهده

(۱۳)  $B = \{b_i(v_k) = P(v_k \text{ at } t | s_i \text{ at } t)\}$ : احتمالات تولید (انتشار) علائم قابل مشاهده

(۱۴)  $O_i$ : علامت مشاهده شده در زمان  $t$

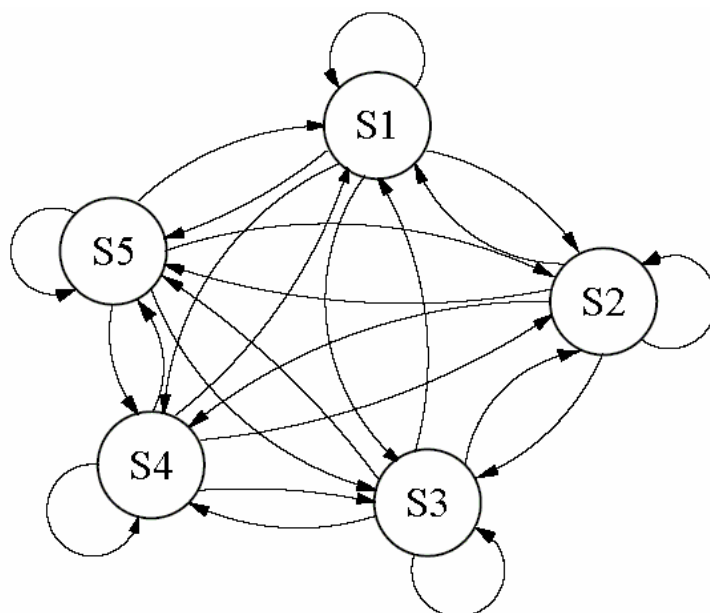
(۱۵)  $T$ : طول دنباله مشاهدات

(۱۶)  $\lambda = (A, B, \Pi)$ : نماد خلاصه برای مدل پنهان مارکوف

واضح است که بر احتمالات  $\Pi$ ،  $A$  و  $B$  سه قید وجود دارد:  $\sum_{i=1}^N \pi_i = 1$ ،  $\sum_{j=1}^N a_{ij} = 1$ ،  $\forall i$  و

$$\sum_{k=1}^M b_i(v_k) = 1, \forall i$$

ساختار ماتریس احتمالات تغییر حالت  $A$  توپولوژی HMM را تعیین میکند. اگر  $a_{ij} \neq 0 \forall i, j$  یعنی هر حالت مدل از هر حالت دیگر با یک گذر قابل رسیدن باشد، مدل "کاملاً متصل" یا "ارگودیک" (Ergodic) نامیده میشود (شکل ۲). برای مدل سازی کلمات زبان که بعداً شرح داده خواهد شد، از مدل ارگودیک استفاده خواهد شد.



شکل ۲- یک مدل مارکوف ارگودیک پنج حالت

یک توپولوژی دیگر که در کاربردهای بازشناسی گفتار و متن بسیار پر استفاده است، توپولوژی "چپ به راست" (LR) یا "بکیس" (Bakis) میباشد که در آن حالت‌های با شماره پایین تر، مشاهدات نخستین را تولید میکنند. ترتیب زمانی در HMM های چپ به راست با قرار دادن صفرهای ساختاری در مدل به شکل قیدهای  $\Pi = \{1, 0, \dots, 0\}$  و  $a_{ij} = 0, i > j$  اعمال میشود، که یعنی مدل از اولین حالت (حالت با کمترین شماره و در واقع سمت چپ ترین حالت) شروع میکند و در هر حالت تغییر فقط میتواند به حالت‌های با شماره بالاتر صورت گیرد. به عنوان یک محدودیت دیگر، بیشتر اوقات در HMM های چپ به راست، اندازه پرش‌های رو به جلو در هر حالت محدود میشود و به این ترتیب از تغییر حالت زیاد جلوگیری میشود، یعنی برای یک  $\Delta$  ثابت،  $a_{ij} = 0, j > i + \Delta$  در نظر گرفته میشود.

مثال زیر به درک کاربرد مدل‌های پنهان مارکوف کمک میکند. فرض کنید یک نفر برای مدتی در اتاقی زندانی است و میخواهد از وضعیت هوای بیرون اطلاع داشته باشد. تنها اطلاعاتی که او از دنیای بیرون دارد این است که آیا کسی که هر روز برای او وعده غذایی را می آورد با چتر وارد میشود یا نه؛ بنابراین برای مشاهده حمل چتر  $V = \{True, False\}$  برای سادگی فرض میشود که هوا تنها سه وضعیت دارد: آفتابی، ابری و بارانی، و یک روز معادل یک بازه زمانی است، یعنی وضعیت هوا در طول روز تغییر نمیکند. فرض کنید احتمال حمل چتر در یک روز به شرط اینکه هوا آفتابی باشد 0.1 باشد، همچنین احتمال حمل چتر به شرط ابری بودن 0.3 و احتمال حمل چتر به شرط بارانی بودن 0.7 باشد. هدف این است که فرد از مشاهدات (حمل کردن یا نکردن چتر) درباره وضعیت هوای بیرون (که از او پنهان است) نتیجه ای بگیرد. فرض

کنید  $w_i$  وضعیت هوا در روز  $i$  باشد، و متغیر بولی  $u_i$  به این معنی باشد که در آن روز چتر مشاهده میشود یا خیر. با استفاده از قانون بیس (Bayes):

$$(۱۷) \quad P(w_1, \dots, w_n | u_1, \dots, u_n) = \frac{P(u_1, \dots, u_n | w_1, \dots, w_n) P(w_1, \dots, w_n)}{P(u_1, \dots, u_n)}$$

احتمال  $P(w_1, \dots, w_n)$  معادل مدل مارکوف مثال قبل است، و  $P(u_1, \dots, u_n)$  احتمال از پیش معلوم دیدن دنباله ای خاص از رخدادهاى حمل کردن یا نکردن چتر است. اگر فرض شود که به ازای همه  $i$  ها، به شرط  $w_i, u_i$  از هر  $u_j$  و  $w_j$  به ازای هر  $j \neq i$  مستقل باشد، آنگاه احتمال  $P(u_1, \dots, u_n | w_1, \dots, w_n)$  این گونه محاسبه میشود  $\prod_{i=1}^n P(u_i | w_i)$ .

در مثال پیش بینی وضع هوا (و خیلی از مسائل دیگر) میتوان از احتمال پیشین  $P(u_1, \dots, u_n)$  صرف نظر کرد زیرا از وضع هوا مستقل است. بر مبنای فرض مارکوف مرتبه اول، معیار محتمل بودن (Likelihood)، که با احتمال متناسب است، اینگونه محاسبه میشود:

$$(۱۸) \quad \begin{aligned} P(w_1, \dots, w_n | u_1, \dots, u_n) &\propto \\ L(w_1, \dots, w_n | u_1, \dots, u_n) &= P(u_1, \dots, u_n | w_1, \dots, w_n) P(w_1, \dots, w_n) \\ &= \prod_{i=1}^n P(u_i | w_i) \prod_{i=1}^n P(w_i | w_{i-1}) \end{aligned}$$

فرض کنید روز حبس شدن فرد آفتابی بوده باشد، و روز بعد مسوول غذا با چتر وارد شده باشد. مطلوب است پیش بینی وضعیت هوای روز بعد:

در ابتدا با این فرض که روز بعد آفتابی بوده باشد معیار محتمل بودن را حساب میکنیم:

$$L(w_2 = \text{Sunny} | w_1 = \text{Sunny}, u_2 = \text{True}) = P(u_2 = \text{True} | w_2 = \text{Sunny}).$$

$$(۱۹) \quad P(w_2 = \text{Sunny} | w_1 = \text{Sunny}) = 0.1 (0.8) = 0.08$$

سپس با این فرض که روز بعد ابری بوده باشد معیار محتمل بودن را حساب میکنیم:

$$L(w_2 = \text{Cloudy} | w_1 = \text{Sunny}, u_2 = \text{True}) = P(u_2 = \text{True} | w_2 = \text{Cloudy}).$$

$$(۲۰) \quad P(w_2 = \text{Cloudy} | w_1 = \text{Sunny}) = 0.3 (0.1) = 0.03$$

و سرانجام با این فرض که روز بعد بارانی بوده باشد معیار محتمل بودن را حساب میکنیم:

$$L(w_2 = \text{Rainy} | w_1 = \text{Sunny}, u_2 = \text{True}) = P(u_2 = \text{True} | w_2 = \text{Rainy}).$$

$$(۲۱) \quad P(w_2 = \text{Rainy} | w_1 = \text{Sunny}) = 0.7 (0.1) = 0.07$$

بنابراین پر احتمال ترین حالت این است که روز بعد آفتابی بوده باشد.



## سه مساله اساسی مدل‌های پنهان مارکوف

در کاربردهای HMM، نیازمند حل حداقل یکی از سه مساله زیر هستیم:

**مساله ۱.** اگر مدل  $\lambda = (A, B, \Pi)$  را داشته باشیم، چطور میتوانیم احتمال  $P(O | \lambda)$  یعنی احتمال وقوع دنباله مشاهده  $O = O_1, O_2, \dots, O_T$  به شرط مدل را به شکل موثری حساب کنیم؟

**مساله ۲.** اگر مدل  $\lambda$  و بردار مشاهده  $O$  را داشته باشیم، چطور دنباله حالت  $S = s_1, s_2, \dots, s_T$  را انتخاب کنیم که  $P(O, S | \lambda)$  یعنی احتمال مشترک دنباله مشاهده  $O = O_1, O_2, \dots, O_T$  و دنباله حالت  $S$  به شرط مدل، ماکزیمم شود؟ به بیان دیگر، هدف این است که دنباله حالتی پیدا کنیم که مشاهدات را به بهترین نحو توصیف کند.

**مساله ۳.** اگر بردار مشاهده  $O$  را داشته باشیم، چطور پارامترهای مدل  $\lambda = (A, B, \Pi)$  را تنظیم کنیم که احتمال  $P(O | \lambda)$  یا  $P(O, S | \lambda)$  ماکزیمم شود. به بیان دیگر، هدف پیدا کردن (یا در واقع آموزش) مدلی است که داده های مشاهده شده را به بهترین نحو توصیف کند.

### حل مساله ۱

در این مساله هدف محاسبه احتمالی است که مدل  $\lambda$  دنباله مشاهده  $O$  را تولید کند. ساده ترین روش برای محاسبه  $P(O | \lambda)$  پیدا کردن  $P(O | S, \lambda)$  برای یک دنباله حالت ثابت  $S$  و ضرب آن در  $P(S | \lambda)$  و جمع برای تمام دنباله های حالت ممکن به طول  $T$  است:

$$(22) \quad P(O | \lambda) = \sum_S P(O | S, \lambda) \cdot P(S | \lambda)$$

چون  $P(S | \lambda) = \pi_{s_1} a_{s_1 s_2} a_{s_2 s_3} \dots a_{s_{T-1} s_T}$  و  $P(O | S, \lambda) = b_{s_1}(O_1) b_{s_2}(O_2) \dots b_{s_T}(O_T)$  معادله (۲۲) میتواند اینگونه نوشته شود:

$$(23) \quad P(O | \lambda) = \sum_S \pi_{s_1} b_{s_1}(O_1) a_{s_1 s_2} b_{s_2}(O_2) \dots a_{s_{T-1} s_T} b_{s_T}(O_T)$$

محاسبه احتمال با استفاده از رابطه (۲۳) عملی نیست زیرا  $N^T$  دنباله حالت وجود دارد و بنابراین  $(2T-1)N^T$  عمل ضرب و  $N^T-1$  عمل جمع لازم است. بنابراین استفاده از یک روش عملی و کارا ضروری است. خوشبختانه دو روش برای این کار وجود دارد: الگوریتم رو به جلو (Forward) و الگوریتم رو به عقب (Backward) که شرح آنها در زیر می آید:

الگوریتم رو به جلو برای هر حالت  $s$  متغیر رو به جلو  $\alpha_t(s)$  را که به شکل زیر تعریف میشود محاسبه میکند:

$$(24) \quad \alpha_t(s) = P(O_1, O_2, \dots, O_t, s_t = s | \lambda)$$

که احتمال دیدن جزئی دنباله مشاهده تا زمان  $t$  و بودن در حالت  $s$  به شرط مدل  $\lambda$  است. روال سه مرحله ای زیر  $\alpha_t(s)$  را برای تمام حالات و زمانها حساب میکند:  
 ۱- مقدار دهی نخستین:

$$(25) \quad \alpha_1(s) = \pi_s b_s(O_1), 1 \leq s \leq N$$

۲- استقرا:

$$(26) \quad \alpha_{t+1}(r) = \left[ \sum_{s=1}^N \alpha_t(s) a_{sr} \right] b_r(O_{t+1}), 1 \leq r \leq N, 1 \leq t \leq T-1$$

این احتمال رو به جلوی قرار گرفتن در حالت  $r$  در زمان  $t+1$  بر مبنای احتمال مشترک متغیرهای رو به جلوی همه حالتها در زمان  $t$  و احتمالات انتقال از هر یک از آنها به حالت  $r$  است. به این دلیل میتوان احتمال را به این شکل محاسبه کرد که از هر یک از  $N$  حالت در زمان  $t$  میتوان به طور مستقل (در زمان  $t+1$ ) به حالت  $r$  (با احتمال  $a_{sr}$ ) رسید.

۳- خاتمه:

$$(27) \quad P(O | \lambda) = \sum_{s=1}^N \alpha_T(s)$$

محاسبه متغیرهای رو به جلو برای تمام حالتها و در تمام زمانها به  $N(N-1)(T-1) + (N-1)$  عمل جمع و  $N + N(N+1)(T-1)$  عمل ضرب احتیاج دارد؛ که از مرتبه  $N^2T$  است، در مقایسه با مرتبه محاسبات در روش مستقیم که  $TN^T$  بود.

الگوریتم رو به عقب کاملاً شبیه به الگوریتم رو به جلو است با این تفاوت که در اینجا متغیر رو به عقب  $\beta_t(s)$  که در زیر تعریف میشود برای هر حالت در جهت عکس محاسبه میشود:

$$(28) \quad \beta_t(s) = P(O_{t+1}, O_{t+2}, \dots, O_T | s_t = s, \lambda)$$

این احتمال دنباله مشاهده از  $t+1$  تا  $T$  به شرط مدل  $\lambda$  و قرار گرفتن در حالت  $s$  در زمان  $t$  است. شبیه  $\alpha_t(s)$  هم میتواند برای تمام حالات و زمانها با یک الگوریتم سه مرحله ای محاسبه شود:

۱- مقدار دهی نخستین:

$$(29) \quad \beta_T(s) = 1, 1 \leq s \leq N$$

۲- استقرا:

$$(۳۰) \quad \beta_t(s) = \sum_{r=1}^N a_{sr} b_r(O_{t+1}) \beta_{t+1}(r), \quad 1 \leq s \leq N, t = T-1, T-2, \dots, 1$$

۳- خاتمه:

$$(۳۱) \quad P(O | \lambda) = \sum_{s=1}^N \pi_s b_s(O_1) \beta_1(s)$$

محاسبه  $P(O | \lambda)$  با استفاده از متغیرهای رو به عقب نیز شامل محاسباتی از مرتبه  $N^2 T$  است.

## حل مساله ۲

در این مساله میخواهیم پراحتمال ترین دنباله حالت (یعنی قسمت پنهان مدل) متناظر با یک دنباله مشاهده را پیدا کنیم. الگوریتم معروف "ویتربی" (Viterbi, 1967) یک الگوریتم برنامه نویسی پویا (Dynamic Programming) برای پیدا کردن این دنباله (مسیر) بهینه است. این الگوریتم بهترین دنباله حالت را در هر لحظه از زمان برای هر یک از  $N$  حالت نگه میدارد، و در نهایت بهترین مسیر را برای هر یک از  $N$  حالت به عنوان حالتی که آخرین مشاهده در آن صورت گرفته است محاسبه میکند، و از بین این مسیرها، مسیر دارای بیشترین احتمال را به عنوان مسیر بهینه کلی انتخاب می کند.

الگوریتم چهار مرحله ای ویتربی همان استراتژی الگوریتم رو به جلو را دنبال میکند، اما با این تفاوت که عمل جمع با ماکزیمم (یا مینیمم) جایگزین میشود. انتخاب ماکزیمم یا مینیمم بستگی به این دارد که معیار بهینگی چه انتخاب شده باشد (بدیهی است که اگر معیار انتخاب شده احتمال باشد، باید ماکزیمم شود و اگر هزینه باشد باید مینیمم شود). برای دنباله مشاهده  $O = O_1, O_2, \dots, O_T$  و مدل  $\lambda$  الگوریتم به شکل زیر بیان میشود:

۱- مقدار دهی نخستین:

$$(۳۲) \quad \delta_1(s) = \pi_s b_s(O_1)$$

$$(۳۳) \quad \psi_1(s) = 0, \quad 1 \leq s \leq N$$

که  $\delta_t(s)$  نشان دهنده وزنهای انباشته شده است وقتی که مدل در زمان  $t$  در حالت  $s$  قرار دارد، و  $\psi_t(s)$  نمایانگر حالتی است که در زمان  $t-1$  پایین ترین هزینه (بالاترین احتمال) مربوط به انتقال حالت به  $s$  در زمان  $t$  را داشته است.

۲- استقرا:

$$(۳۴) \quad \delta_t(s) = \max_{1 \leq r \leq N} [\delta_{t-1}(r) a_{rs}] b_s(O_t)$$

$$(۳۵) \quad \psi_t(s) = \arg \max_{1 \leq r \leq N} [\delta_{t-1}(r) a_{rs}], \quad 1 \leq s \leq N, 2 \leq t \leq T$$

۳- خاتمه:

$$(۳۶) \quad P^* = \max_{1 \leq s \leq N} [\delta_T(s)]$$

$$(۳۷) \quad q_T^* = \arg \max_{1 \leq s \leq N} [\delta_T(s)]$$

۴- برگشت به عقب برای پیدا کردن مسیر بهینه:

$$(۳۸) \quad q_t^* = \psi_{t+1}(q_{t+1}^*), \quad t = T-1, T-2, \dots, 1$$

و حالا  $Q^* = \{q_1^*, q_2^*, \dots, q_T^*\}$  دنباله حالت بهینه است و  $P^*$  احتمال مشترک دنباله مشاهده  $O$  و دنباله حالت بهینه  $Q^*$  است.

شبهه به الگوریتم های رو به جلو و رو به عقب، پیچیدگی محاسباتی الگوریتم ویتربی از مرتبه  $N^2T$  میباشد.

البته توجه به این نکته ضروری است که پیاده سازی مستقیم الگوریتم بالا میتواند باعث به وجود آمدن مشکل پاریز (underflow) شود، زیرا احتمالاتی که حساب میشوند شامل ضرب اعداد کوچک است که به سرعت باعث میشود محدوده اعداد تولید شده از محدوده پیش بینی شده اعداد ممیز شناور کامپیوتر خارج شود. برای حل این مشکل، الگوریتم ویتربی طوری تغییر داده میشود تا به جای احتمال با لگاریتم احتمال کار کند. این تکنیک نه تنها مشکل پاریز را حل میکند بلکه روند محاسبات را نیز سرعت میدهد، زیرا عمل جمع خیلی سریعتر از ضرب است. باید توجه کرد که الگوریتم ویتربی یک الگوریتم زمان اجراست نه یک الگوریتم آموزش که معمولا به شکل آفلاین انجام میشود، بنابراین پیاده سازی سریع این الگوریتم بسیار مطلوب و ضروری است. نسخه کارا و عملی الگوریتم ویتربی در زیر می آید:

۰- پیش پردازش

$$(۳۹) \quad \tilde{\pi}_s = \log(\pi_s), \quad 1 \leq s \leq N$$

$$(۴۰) \quad \tilde{a}_{rs} = \log(a_{rs}), \quad 1 \leq r, s \leq N$$

$$(۴۱) \quad \tilde{b}_s(O_t) = \log(b_s(O_t)), \quad 1 \leq s \leq N, 1 \leq t \leq T$$

۱- مقدار دهی نخستین:

$$(۴۲) \quad \tilde{\delta}_1(s) = \tilde{\pi}_s + \tilde{b}_s(O_1)$$

$$(۴۳) \quad \psi_1(s) = 0, \quad 1 \leq s \leq N$$

۲- استقرا:

$$(۴۴) \quad \tilde{\delta}_t(s) = \max_{1 \leq r \leq N} [\tilde{\delta}_{t-1}(r) + \tilde{a}_{rs}] + \tilde{b}_s(O_t)$$

$$(۴۵) \quad \psi_t(s) = \arg \max_{1 \leq r \leq N} [\tilde{\delta}_{t-1}(r) + \tilde{a}_{rs}], \quad 1 \leq s \leq N, 2 \leq t \leq T$$

۳- خاتمه:

$$(۴۶) \quad P^* = \max_{1 \leq s \leq N} [\tilde{\delta}_T(s)]$$

$$(۴۷) \quad q_T^* = \arg \max_{1 \leq s \leq N} [\tilde{\delta}_T(s)]$$

۴- برگشت به عقب برای پیدا کردن مسیر بهینه:

$$(۴۸) \quad q_t^* = \psi_{t+1}(q_{t+1}^*), \quad t = T-1, T-2, \dots, 1$$

و حالا  $Q^* = \{q_1^*, q_2^*, \dots, q_T^*\}$  دنباله حالت بهینه است و  $\exp(P^*)$  احتمال مشترک دنباله مشاهده  $O$  و دنباله حالت بهینه  $Q^*$  است.

شاید متوجه شده باشید که الگوریتم ویتربی تنها شامل عملیات ضرب است که با لگاریتم گیری تبدیل این اعمال ضرب به سادگی تبدیل به جمع میشوند. اما الگوریتمهای رو به جلو و رو به عقب شامل عملیات جمع هستند. در اینجا اگر لازم باشد مشکل پاریز اعداد ممیز شناور میتواند با همان لگاریتم گیری حل شود، تنها لازم است  $\log(x+y)$  حساب شود که با استفاده از تکنیک زیر قابل انجام است (Manning and Schütze, 1999):

```

if  $y - x > \log \text{big}$ 
    return  $y$ ;
else if  $x - y > \log \text{big}$ 
    return  $x$ ;
else
    return  $\min(x, y) + \log(\exp(x - \min(x, y)) + \exp(y - \min(x, y)))$ ;

```

که  $\text{big}$  یک ثابت به حد کافی بزرگ مثلا  $10^{30}$  است.

### حل مساله ۳

بر مبنای اینکه چه احتمالی برای ماکزیمم کردن انتخاب شود، دو روش کلی برای تخمین پارامترهای مدل (آموزش) وجود دارد: الگوریتم segmental k-means (Juang and Rabiner, 1990) که پارامترهای مدل را تنظیم میکند تا  $P(O, Q^* | \lambda)$  ماکزیمم شود، که

$Q^*$  دنباله حالت بهینه متناظر با دنباله مشاهده  $O$  است. و الگوریتم Baum-Welch (Rabiner, 1989) که پارامترهای مدل را تنظیم میکند تا  $P(O|\lambda)$  به یک ماکزیمم برسد، در اینجا  $P(O|\lambda)$  میتواند جمع احتمالات  $P(O, S|\lambda)$  روی تمام دنباله حالت‌های ممکن  $S$  تعبیر شود، یعنی این الگوریتم روی دنباله حالت خاصی متمرکز نمیشود. معمولا الگوریتم segmental k-means برای آموزش ترجیح داده میشود زیرا در مقایسه با الگوریتم Baum-Welch به محاسبات خیلی کمتری احتیاج دارد، و همچنین برای خیلی از کاربردهای مدلسازی و کدگشایی، معیار  $P(O, Q^*|\lambda)$  کاملا طبیعی به نظر میرسد. در دو قسمت زیر شرح الگوریتم‌ها آورده میشود.

### الگوریتم Segmental K-Means

مانند هر الگوریتم آموزش، الگوریتم segmental k-means نیز به تعدادی دنباله مشاهده (آموزشی) احتیاج دارد. فرض کنید به تعداد  $w$  از چنین دنباله‌هایی موجود باشد. هر دنباله  $O = O_1, O_2, \dots, O_{T_i}$  شامل  $T_i$  بردار مشاهده است، بنابراین در مجموع  $\sum_{i=1}^w T_i$  بردار مشاهده وجود دارد. اگر تنها یک دنباله دراز موجود باشد، میتوان آن را به تعداد اختیاری دنباله‌های کوچک تقسیم کرد و از آنها برای آموزش استفاده کرد. فرض میشود هر علامت (Symbol) مشاهده یک بردار با اندازه (بعد) یک با بیشتر باشد و بدیهی است که بردارهای مشاهده باید دارای بعد یکسان باشند. همچنین در مدل‌های پنهان مارکوف گسسته که تا اینجا مورد بحث قرار گرفته اند لازم است که این علامتها متعلق به یک مجموعه محدود باشند. الگوریتم شامل مراحل زیر است:

$N-1$  بردار آموزشی  $C_1, C_2, \dots, C_N$  را به طور تصادفی انتخاب کنید و هر یک از بردارهای آموزشی باقیمانده را به یکی از این  $N$  بردار که کمترین فاصله (مثلا اقلیدسی) را تا آن دارد منسوب کنید. بنابراین  $N$  دسته (خوشه) شکل میگیرد که هر کدام یک حالت (با شماره ای بین ۱ تا  $N$ ) نامیده میشود. نماد  $O_i \in S$  یعنی  $t$  امین علامت از یک دنباله مشاهده به حالت (دسته)  $S$  منسوب شده است. انتخاب اولیه دسته‌ها بر HMM نهایی تاثیری نمیگذارد، اما میتواند تعداد تکرارهای لازم برای آموزش مدل را تعیین کند. برای اینکه انتخاب اولیه دسته‌ها تا حد امکان توزیع پراکنده ای داشته باشد، یک استراتژی خوب برای وقتی که  $w \geq N$  این است که  $C_1$  اولین بردار مشاهده اولین دنباله انتخاب شود،  $C_2$  دومین بردار مشاهده دومین دنباله انتخاب شود و ... (Dugad and Desai, 1996). این مرحله مقدار دهی نخستین مناسبی برای کل روال آموزش آماده میکند.

۲- احتمالات نخستین و انتقال را به شکل زیر محاسبه کنید:

$$(۴۹) \quad \hat{\pi}_s = \frac{\text{number of occurrences of } \{O_1 \in s\}}{\text{total number of occurrences of } O_1}, \quad 1 \leq s \leq N$$

$$(۵۰) \quad \hat{a}_{rs} = \frac{\text{number of occurrences of } \{O_t \in r \text{ and } O_{t+1} \in s\}}{\text{total number of occurrences of } \{O_t \in r\}}, \quad 1 \leq r, s \leq N, \quad 1 \leq t \leq T_i - 1$$

۳- ماتریس میانگین و کوواریانس را به شکل زیر محاسبه کنید:

$$(۵۱) \quad \hat{\mu}_s = \frac{1}{N_s} \sum_{O_t \in s} O_t, \quad 1 \leq s \leq N$$

$$(۵۲) \quad \hat{V}_s = \frac{1}{N_s} \sum_{O_t \in s} (O_t - \hat{\mu}_s)^T (O_t - \hat{\mu}_s), \quad 1 \leq s \leq N$$

۴- توزیع احتمال را برای هر بردار مشاهده و برای هر حالت محاسبه کنید:

$$(۵۳) \quad \hat{b}_s(O_t) = \frac{1}{(2\pi)^{\frac{D}{2}} |\hat{V}_s|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (O_t - \hat{\mu}_s)^T \hat{V}_s^{-1} (O_t - \hat{\mu}_s)\right)$$

ثابت شده است که برای دسته وسیعی از توابع چگالی احتمال از جمله گاوسی، الگوریتم همگرا میشود. در اینجا تابع گاوسی به طور اختیاری انتخاب شده است.

۵- با احتمالات جدید با استفاده از الگوریتم ویتربی برای هر دنباله آموزشی، دنباله حالت بهینه  $Q^*$  را پیدا کنید. اگر انتساب حالت قبلی یک بردار مشاهده متفاوت از حالت بهینه تخمین زده شده متناظر باشد آن را به حالت جدید منسوب کنید؛ یعنی اگر  $q_t^* = s$  آنگاه  $O_t$  را به  $s$  منسوب کنید.

۶- اگر در مرحله ۵ هر کدام از بردارهای مشاهده به حالتی جدید (یعنی حالتی مخالف با حالت انتساب یافته قبلی اش) منسوب شد، آنگاه مراحل ۲ تا ۶ باید تکرار شوند و در غیر این صورت الگوریتم متوقف میشود.

### الگوریتم Baum-Welch

این الگوریتم جزء دسته الگوریتم های موسوم به ماکزیمم کردن توقع (EM) میباشد. الگوریتم EM روش بسیار پر استفاده ای برای یادگیری در مسائلی است که در آنها متغیرهای مشاهده نشده (پنهان) وجود دارند. یک الگوریتم EM با تخمین زدن مکرر مقدار مورد انتظار (Expected Value) هر متغیر پنهان با استفاده از فرضهای جاری و سپس محاسبه مجدد متحمل ترین فرضها (maximum likelihood hypothesis) با استفاده از مقادیر مورد

انتظار متغیرها، محتمل ترین فرضیه ها را جستجو میکند. به بیان دیگر، در مرحله اول، فرضیه‌های جاری برای تخمین مقدار متغیرهای پنهان استفاده میشوند و در مرحله دوم مقادیر این متغیرها برای بهبود دادن فرضیه ها استفاده میشود. میتوان اثبات کرد که چنین روندی نهایتاً فرضیه‌هایی دارای بیشترین احتمال محلی را پیدا میکند. یعنی ممکن است الگوریتم EM همانند بسیاری از الگوریتمهای تکراری ممکن است نتواند جواب بهینه سراسری را پیدا کند، اما همواره یک بهینه محلی را پیدا میکند (Mitchell, 1997).

تخمین اولیه فرضها (پارامترها)، یعنی ساخت HMM اولیه، میتواند با هر روشی انجام شود، اما یک تخمین اولیه معقول میتواند با استفاده از چهار مرحله اول الگوریتم  $k$ -segmental means به دست آید.

قبل از ارائه فرمولهای اصلی باید چند مفهوم و نماد را معرفی کنیم.  $\gamma_t(s) = P(s_t = s | O, \lambda)$  را در نظر بگیرید، که احتمال قرار داشتن در حالت  $s$  در زمان  $t$  به شرط مدل  $\lambda$  و دنباله مشاهده  $O$  است. با استفاده از قانون بیس داریم:

$$(54) \quad \gamma_s(t) = \frac{P(s_t = s, O | \lambda)}{P(O | \lambda)} = \frac{\alpha_t(s)\beta_t(s)}{P(O | \lambda)}, \quad 1 \leq s \leq N$$

که  $\alpha_t(s)$  و  $\beta_t(s)$  متغیرهای رو به جلو و رو به عقب هستند. همچنین  $\xi_t(r, s) = P(s_t = r, s_{t+1} = s | O, \lambda)$  را تعریف میکنیم که احتمال قرار داشتن در حالت  $r$  در زمان  $t$  و انتقال به حالت  $s$  در زمان  $t+1$  به شرط مدل  $\lambda$  و دنباله مشاهده  $O$  است. با استفاده از قانون بیس و خاصیت کزال بودن زنجیره مارکوف (رابطه (۱)) میتوان نشان داد:

$$(55) \quad \xi_t(r, s) = \frac{\alpha_t(r)a_{rs}b_s(O_{t+1})\beta_{t+1}(s)}{P(O | \lambda)}, \quad 1 \leq r, s \leq N$$

اگر  $\gamma_t(s)$  از  $t = 1$  تا  $T$  جمع زده شود، تعداد دفعات قابل انتظاری که حالت  $s$  دیده شده است به دست می آید، و اگر تنها تا  $T-1$  جمع زده شود، تعداد دفعات قابل انتظاری که از حالت  $s$  انتقالی انجام شده است به دست می آید. به طور مشابه، اگر  $\xi_t(r, s)$  از  $t = 1$  تا  $T$  جمع زده شود، تعداد دفعات قابل انتظاری که انتقالی از حالت  $r$  به حالت  $s$  انجام گرفته است به دست می آید:

$$(56) \quad \sum_{t=1}^{T-1} \gamma_t(s) = \text{expected number of transitions from state } s, \quad 1 \leq s \leq N$$

$$(57) \quad \sum_{t=1}^{T-1} \xi_t(r, s) = \text{expected number of transitions from state } r \text{ to state } s, \quad 1 \leq r, s \leq N$$

رابطه بین  $\gamma_t(r)$  و  $\xi_t(r, s)$  میتواند با جمع کردن  $\xi_t(r, s)$  روی  $s$  به دست آید:



$$(58) \quad \gamma_t(r) = \sum_{s=1}^N \xi_t(r, s), \quad 1 \leq r \leq N$$

و حالا فرمولهای تخمین مجدد Baum-Welch به شکل زیر تعریف میشوند:

$$(59) \quad \hat{\pi}_s = \gamma_1(s), \quad 1 \leq s \leq N$$

$$(60) \quad \hat{a}_{rs} = \frac{\sum_{t=1}^{T-1} \xi_t(r, s)}{\sum_{t=1}^{T-1} \gamma_t(r)}, \quad 1 \leq r, s \leq N$$

$$(61) \quad \hat{b}_s(v_k) = \frac{\sum_{t=1, O_t=v_k}^T \gamma_t(s)}{\sum_{t=1}^T \gamma_t(s)}, \quad 1 \leq s \leq N$$

فرمول تخمین مجدد برای  $\pi_s$  به سادگی برابر است با احتمال قرار داشتن در حالت  $s$  در زمان ۱. فرمول تخمین مجدد برای  $a_{rs}$  برابر است با نسبت تعداد دفعات قابل انتظار انتقال حالت از  $r$  به  $s$  به تعداد دفعات قابل انتظاری که از  $r$  انتقال حالتی انجام شده است. و سرانجام، فرمول تخمین مجدد برای  $b_s(v_k)$  برابر است با نسبت تعداد دفعات قابل انتظار قرار داشتن در حالت  $s$  و مشاهده علامت  $v_k$  به تعداد دفعات قابل انتظار قرار داشتن در حالت  $s$ .

### مدل (پنهان) مارکوف برای کلمات فارسی

در استفاده از HMM برای مدلسازی، یک مساله اساسی تعیین تعداد حالت‌های مدل است؛ در بیشتر اوقات، ارتباط از پیش معلومی بین تعداد حالتها و تعداد علائم قابل مشاهده وجود ندارد. مساله تعیین تعداد حالتها در HMM ها مشابه مساله تعیین تعداد لایه های مخفی و تعداد نرونهای آنها در شبکه های عصبی است. چون هنوز تئوری یا روشی اصولی برای انتخاب تعداد حالت‌های HMM وجود ندارد، به عنوان یک راه حل ساده، معمولا تعداد حالت‌های مدل را متناسب با تعداد علائم قابل مشاهده، که کاملا مشخص است، در نظر میگیرند. مثلا در کاربردهای تشخیص متون ماشینی، معمولا پنج (تا ده) حالت به هر حرف (کاراکتر) تخصیص داده میشود و بنابراین تعداد حالت‌های مدل برابر میشود با پنج (تا ده) برابر اندازه مجموعه کاراکترها (علائم قابل مشاهده). اما در مدلی که برای کلمات فارسی پیشنهاد میدهیم، هر حالت مدل را با یک حرف الفبای فارسی متناظر میکنیم، بنابراین مدل قسمت پنهانی ندارد؛ اینگونه مدلها را "مدلهای آشکار مارکوف" (Visible Markov Models) یا به اختصار VMM مینامند. چون ترتیب زمانی در مدل کلمات اهمیتی ندارد و به طور طبیعی هر حالت مدل (یعنی هر حرف) باید با یک انتقال از هر حالت دیگر قابل رسیدن باشد، معقول است که برای این مدل توپولوژی کاملا متصل (ارگودیک) را برگزینیم.

در مدل‌های VMM، چون دنباله حالت بهینه از پیش مشخص است، سه مساله مدل‌های HMM به شکل خیلی ساده تری حل میشود. در مدل پیشنهادی برای کلمات، دنباله‌های مشاهده (آموزشی) همان کلمات لغتنامه استخراج شده هستند، هرچند که میتوان برای ساخت (آموزش) مدل، یا در واقع محاسبه احتمالات اولیه، انتقال و انتشار (مشاهده)، از الگوریتم‌های استاندارد Segmental K-Means یا Baum-Welch استفاده کرد، اما این احتمالات را به سادگی میتوان از روابط زیر به دست آورد:

$$\pi_i = \frac{\text{No. of words starting with letter } s_i}{\text{Lexicon size}} \quad (62)$$

$$a_{ij} = \frac{\text{No. of transitions from letter } s_i \text{ to } s_j}{\text{Total number of transitions from letter } s_i} \quad (63)$$

$$b_i(v_k) \approx \begin{cases} 1 & \text{if } s_i = v_k \\ 0 & \text{otherwise} \end{cases} \quad (64)$$

حالا با داشتن این احتمالات، برای محاسبه  $P(O | \lambda)$ ، یعنی حل مساله ۱، بجای استفاده از الگوریتم رو به جلو یا رو به عقب میتوان به سادگی از رابطه زیر استفاده کرد:

$$P(O | \lambda) = \pi_{o_1} a_{o_1 o_2} \dots a_{o_{T-1} o_T} \quad (65)$$

### نتایج آزمایشی مدل ارائه شده کلمات فارسی

در زیر چند نمونه از احتمالات اولیه و انتقال که با استفاده از فرمول‌های (۶۲) و (۶۳) محاسبه شده اند آورده شده است:

$$\pi_b = 3946 / 41000 \approx 0.096$$

$$\pi_z = 81 / 41000 \approx 0.002$$

$$\pi_p = 4590 / 41000 \approx 0.112$$

$$a_{b \rightarrow a} = 1795 / 8029 = 0.224$$

$$a_{b \rightarrow b} = 86 / 8029 = 0.011$$

$$a_{b \rightarrow y} = 1035 / 8029 = 0.129$$

$$a_{z \rightarrow x} = 0.003$$

$$a_{f \rightarrow f} = 0.0005$$

$$a_{z \rightarrow z} = 0.0$$

چون مدل پنهان (یا آشکار) مارکوف یک مدل مولد است، با استفاده از مدل ارائه شده برای کلمات میتوان کلماتی را نیز تولید کرد، که در واقع محتمل ترین دنباله‌های مشاهده از دید مدل هستند. در جدول ۱ نمونه‌هایی از کلمات تولید شده توسط مدل را مشاهده میکنید:

سه حرفی	چهار حرفی	پنج حرفی	شش حرفی
بکس	جاشو	حدیده	رزماند
رزک	فسور	میری	زراندر
حیر	ادهر	سنارا	پرتزاد
ساج	جمین	ستمد	شیرداد
فرا	کافک	نانده	آیکابر
صیت	صلو	جهاعا	ستزوج

جدول ۱- نمونه هایی از کلمات تولید شده به وسیله مدل پیشنهاد شده کلمات فارسی

به منظور بررسی کیفی کلمات تولید شده توسط مدل، تعدادی از کلماتی را که بدون استفاده از مدل و به شکل کاملاً تصادفی تولید شده اند می آوریم: ضمچ، چضض، بططپ، ذضجژ، دلزندی، اژخذ، رصمژن و ظمشخجض. ملاحظه میشود که این کلمات هیچ شباهتی به کلمات زبان فارسی ندارند و هیچ یک به سادگی قابل تلفظ نیستند. اما در مقابل، مدل پیشنهاد شده توانسته است کلمات خوش آهنگی را تولید کند و حتی در برخی موارد (سطر اول جدول) موفق به تولید کلماتی شده است که واقعا عضوی از مجموعه کلمات زبان فارسی میباشند اما در لغتنامه ۴۱۰۰۰ کلمه ای، که برای آموزش مدل نیز به کار رفت، موجود نمیباشند.

به عنوان کاربردی دیگر، با استفاده از مدل ارائه شده لگاریتم احتمال (Logprob) تعدادی کلمه (دنباله مشاهده) درست و نادرست محاسبه میشود:

$$\text{Logprob}(\text{تعالی}) = -12.9$$

$$\text{Logprob}(\text{نگالی}) = -16.9$$

$$\text{Logprob}(\text{مهندسی}) = -15.4$$

$$\text{Logprob}(\text{موندشی}) = -14.1$$

$$\text{Logprob}(\text{مهندسی}) = -36.7$$

$$\text{Logprob}(\text{ستزوج}) = -22.95$$

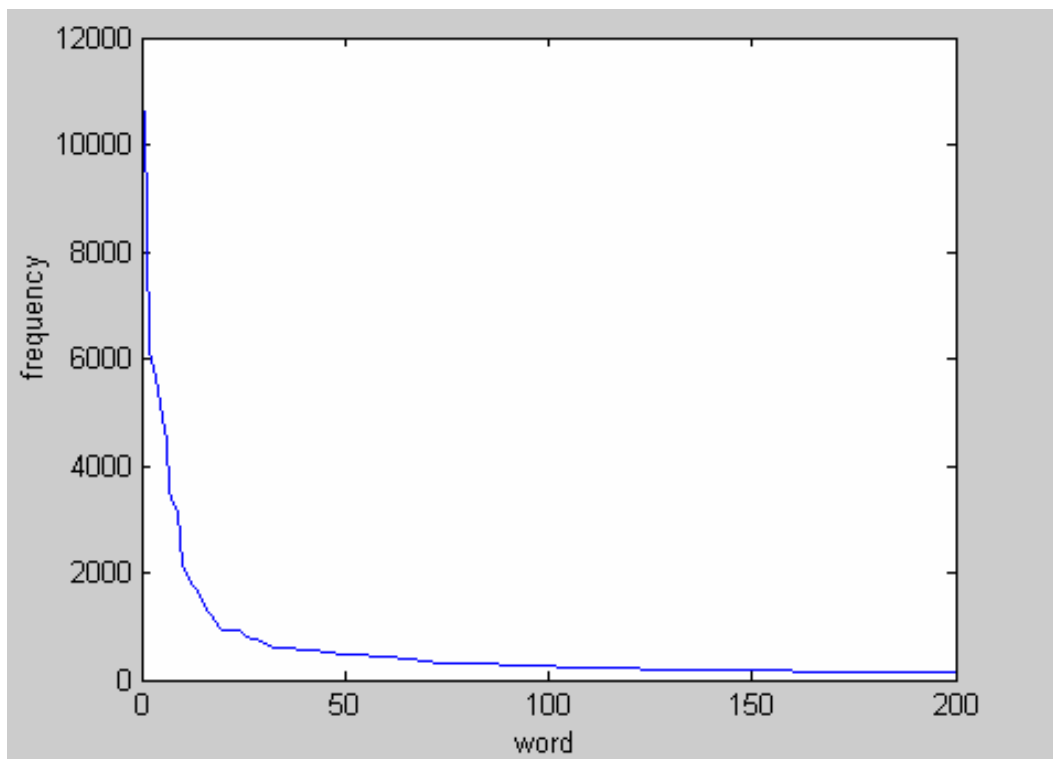
$$\text{Logprob}(\text{ظمشخجض}) = -44.7$$

واضح است که معیار Logprob برای یک کلمه هرچه نزدیک تر به صفر باشد، یعنی آن کلمه به مدل نزدیکتر است. همانطور که قبلاً گفته شد، به دلیل ماهیت تغییر پذیر کلمات زبان، این امکان وجود ندارد که برای تشخیص کلمات نادرست (غیر متعلق به زبان) از یک لغتنامه ثابت استفاده کرد؛ یعنی تشخیص کلمات نادرست بر خلاف آنچه که ممکن است در ابتدا به نظر برسد مساله ساده ای نیست. اما با توجه به مثالهای بالا به نظر میرسد که با استفاده از مدل ارائه شده میتوان حتی بدون استفاده از لغتنامه، کلمات بسیار نادرست (بسیار نامحتمل) را تشخیص داد.

## قسمت سوم

### مدل آماری زبان

همانطور که بخشهای قبل ذکر شد، فرکانس یا احتمال وقوع کلمات یک زبان بسیار متفاوت است. این احتمالات برای زبانهای طبیعی از توزیعی به نام zipf تبعیت میکند. این توزیع برای ۲۰۰ کلمه اول (پر استفاده تر) زبان فارسی با استفاده از نوشته جات آموزشی ما در نمودار شکل ۳ نشان داده شده است.



شکل ۳- نمودار فراوانی ۲۰۰ کلمه اول زبان فارسی

با شمارش دفعات وقوع یک کلمه در یک مجموعه نوشته جات و تقسیم آن بر کل تعداد کلمات میتوان یک تخمین ML برای احتمال وقوع هر کلمه به دست آورد. این ایده را میتوان به زوج کلمات (و ترکیب های طولانی تر) تعمیم داد با این هدف که بعد از دیدن یک یا چند کلمه بتوان کلمه بعدی را حدس زد. پیش بینی یا حدس زدن کلمه بعدی یکی از اعمال ضروری در کاربردهای بازشناسی گفتار، بازشناسی متون (دستنویس یا تایپی) و تشخیص و تصحیح غلطهای املائی است که در آنها تشخیص کلمه کاری دشوار است زیرا ورودی مبهم و همراه با نویز زیاد است. و بنابراین در نظر گرفتن ورودی های قبلی میتواند اطلاعات مهمی درباره گزینه های ممکن برای ورودی جاری در اختیار بگذارد. در مدل های آماری زبان هدف پیش بینی کلمه بعدی یا به بیان دیگر محاسبه احتمال دنباله کلمات (و جملات) است. روش های محاسبه

(انتساب) احتمال به یک جمله میتوانند برای محاسبه احتمال کلمه بعدی در یک جمله ناقص (و بر عکس) به کار روند. دانستن احتمال یک جمله همچنین در کاربردهایی از جمله مشخص کردن نقش کلمات، رفع ابهام معنای کلمات (در ترجمه) و تجزیه احتمالی بسیار موثر است. در مساله تصحیح غلطهای املائی ممکن است که یک یا چند اشتباه تایپی، کلمه مورد نظر را تبدیل به یک کلمه دیگر از زبان (یعنی یک کلمه درست و مجاز) کنند؛ واضح است که در این حالت استفاده از یک لغتنامه به تنهایی نمیتواند کمکی کند. اما با استفاده از اطلاعات کلمات مجاور و گرامرهای زبان در بیشتر موارد میتوان غلط را تشخیص داد و حتی تصحیح کرد. به عنوان مثال کلمه «دوستش» که در جمله «کتاب را به درستش داد» به اشتباه «درستش» تایپ شده است، ترکیب بسیار غیر محتمل «به درستش» را نتیجه داده است که با استفاده از یک گرامر آماری ساده نیز قابل تشخیص است. البته باید توجه کرد که بسیاری از کلمات و جملات کاملاً درست و بامعنی زبان ممکن است احتمال وقوع پایینی داشته باشند، که این مساله تشخی خطا را دشوار میکند.

مدل زبان یا به بیان دیگر مدل پیش بینی کلمه ای که برای فارسی استفاده کرده ایم و در اینجا شرح میدهم یک مدل آماری بسیار پر استفاده به نام N-gram است که در آن احتمال کلمه N ام با استفاده از N-1 کلمه قبلی تخمین زده میشود. موفقیت این مدلها اولین بار در آزمایشگاه های بازشناسی گفتار IBM به اثبات رسید و پس از آن در بسیار از زمینه ها مورد توجه و استفاده محققین قرار گرفتند.

### شمارش کلمات

برای تخمین احتمال نیاز به شمارش ترکیب ها داریم. همانگونه که قبلاً ذکر شد، پردازش آماری زبان بر مبنای یک مجموعه بزرگ نوشته جات است. البته باید توجه کرد که در بعضی کاربردها - از جمله بررسی گرامر، تولید صدا یا شناسایی نویسنده - لازم است تا علائمی همچون '!' و '!' را نیز به عنوان کلمه در نظر گرفت و در شمارش کلمات و ترکیبات متوالی، آنها را نیز به حساب آورد. مساله دیگری که باید به آن توجه کرد اشکال مختلف کلمات است. در بیشتر سیستم های مبتنی بر N-gram اشکال مختلف یک کلمه به عنوان کلمات مجزا در نظر گرفته میشوند؛ یعنی هیچ تلاشی برای ریشه یابی انجام نمیشود که البته این ساده سازی در بعضی کاربردها مناسب نیست.

### N-gram ساده

در این بخش مدل N-gram ساده یا هموار نشده (unsmoothed) را شرح میدهم. یاد آوری میکنیم که مدلهایی که در اینجا برای دنباله کلمات در نظر میگیریم مدلهای احتمالاتی هستند یعنی روشهایی برای انتساب احتمال به رشته های کلمات که هم میتوانند برای

محاسبه احتمال یک دنباله کلی استفاده شوند و هم برای پیش بینی احتمالاتی کلمه بعدی در یک دنباله.

در ساده ترین مدل ممکن برای دنباله کلمات، وقوع هر کلمه پس از هر کلمه دیگر مجاز است. به بیان احتمالاتی، هر کلمه احتمال وقوع یکسانی پس از هر کلمه دلخواه دارد. به عنوان مثال با داشتن لغتنامه 41000 کلمه ای، این احتمال برابر است با  $1 / 41000$  یا تقریباً 0.000024. در یک مدل نسبتاً پیچیده تر فرض میکنیم باز هم هر کلمه میتواند پس از هر کلمه دلخواه قرار گیرد اما احتمال وقوع کلمه بعدی (کلمه دوم) برابر است با فراوانی تکرار نرمال (احتمال وقوع) آن. به عنوان مثال کلمه «این» در نوشته جات آموزشی ما فراوانی تکرار نرمال 0.0182 دارد، پس احتمال اینکه «این» بعد از هر کلمه بیاید 0.0182 است. یعنی در این مدل ساده که unigram نام دارد، کلمات گذشته در پیش بینی کلمه بعدی تاثیری ندارند. در واقع با استفاده از فراوانی های نسبی میتوان یک توزیع احتمال روی کلمات بعدی به دست آورد. اما در نظر نگرفتن تاریخچه گذشته در بیشتر موارد نمیتواند پیش بینی های به درد خوری را نتیجه دهد. به عنوان مثال میدانیم که پس از عبارت «خیلی دیر» قرار گرفتن «شده» بسیار معقول تر از «این» است، اما مدل ساده unigram، «این» را محتمل تر از «شده» میداند، زیرا احتمال وقوعش به تنهایی بیشتر است. این مثال نشان میدهد که احتمال وقوع یک کلمه را باید به شرط کلمات قبل از آن تخمین زد نه به طور مجزا. اگر در مثال ذکر شده حتی تاریخچه ای به اندازه یک - یعنی یک کلمه گذشته - را در نظر بگیریم، مدلی به دست می آید که ترکیب «خیلی دیر شده» را به «خیلی دیر این» ترجیح میدهد زیرا (دیر | شده)  $P$  بیشتر از (دیر | این)  $P$  است.

با توجه به مطالب گفته شده، حالا روش محاسبه احتمال یک رشته کلمات (که به شکل  $w_1 \dots w_n$  یا  $w_1^n$  نمایش داده میشوند) را بررسی میکنیم. اگر فرض کنیم هر کلمه در موقعیت درستش یک رویداد مستقل است، میتوانیم این احتمال را به شکل زیر نشان دهیم:

$$P(w_1, w_2, \dots, w_{n-1}, w_n) \quad (66)$$

با استفاده از رابطه زنجیری احتمال این رابطه را میتوان به شکل زیر تجزیه کرد:

$$P(w_1^n) = P(w_1)P(w_2 | w_1)P(w_3 | w_1^2) \dots P(w_n | w_1^{n-1}) = \prod_{k=1}^n P(w_k | w_1^{k-1}) \quad (67)$$

حالا سوال این است که چطور احتمالاتی همچون  $P(w_n | w_1^n)$  محاسبه کنیم. به نظر میرسد هیچ راه سر راستی برای محاسبه احتمال یک کلمه به شرط دنباله ای طولانی از کلمات قبلی وجود نداشته باشد، مگر اینکه یک مجموعه نوشته جات بسیار بسیار بزرگ در دسترس باشد که عملاً این طور نیست. این مساله را با یک ساده سازی مفید حل میکنیم. یعنی احتمال یک کلمه را به شرط کلمات قبلی را تقریب میزنیم. در ساده ترین تقریب احتمال یک کلمه را به

شرط تنها یک کلمه قبلی به دست می آوریم. این مدل ساده **bigram** نام دارد که در آن  $P(w_n | w_1^n)$  با  $P(w_n | w_{n-1})$  تقریب زده میشود.

این فرض که احتمال یک کلمه تنها به کلمه قبلی وابسته است، فرض مارکوف (مرتبه اول) است. همانطور که در مدل کلمات دیدیم، مدل های مارکوف دسته ای از مدل های احتمالاتی هستند که در آنها فرض میکنیم که میتوانیم احتمال رویدادی در آینده را بدون در نظر گرفتن تاریخچه ای طولانی از رویدادهای گذشته پیش بینی کنیم. بنابراین مدل **bigram** پایه را میتوان به عنوان نوعی زنجیره مارکوف که یک حالت به ازای هر کلمه دارد در نظر گرفت.

بدیهی است که میتوانیم مدل **bigram** را، که تنها یک کلمه گذشته را در نظر میگیرد، به مدل **trigram** که دو کلمه گذشته را در نظر میگیرد و در حالت کلی به مدل **N-gram** که **N-1** کلمه گذشته را در نظر میگیرد تعمیم دهیم. پس مدل **bigram** یک مدل مارکوف مرتبه اول، **trigram** مدل مارکوف مرتبه دوم و **N-gram** مدل مارکوف مرتبه **N-1** ام است. (در عمل مدل های **bigram** یا حداکثر **trigram** استفاده میشوند).

معادله کلی برای تقریب **N-gram** - احتمال شرطی کلمه بعدی در یک دنباله - برابر است با:

$$(68) \quad P(w_n | w_1^{n-1}) \approx P(w_n | w_{n-N+1}^{n-1})$$

این رابطه نشان میدهد که احتمال کلمه  $w_n$  به شرط همه کلمات قبلی میتواند با احتمال این کلمه به شرط  $N$  کلمه قبلی تقریب زده شود.

بنابراین برای یک گرامر **bigram** میتوانیم احتمال یک رشته کامل را با جایگزین کردن رابطه 68 در رابطه 67 محاسبه کنیم:

$$(69) \quad P(w_1^n) \approx \prod_{k=1}^n P(w_k | w_{k-1})$$

البته باید توجه کرد که پیاده سازی مستقیم این روابط ممکن است منجر به خطای پاریز (**underflow**) شود زیرا هرکدام از احتمال ها معمولا عددی بسیار کوچکتر از 1 است و حاصلضرب تعداد کمی از این اعداد هم از محدوده قابل نمایش در کامپیوتر تجاوز میکند. برای حل این مشکل، همانطور که قبلا نیز دیده اید، بجای ضرب احتمالات، لگاریتم آنها را با هم جمع میکنیم و در انتها معکوس لگاریتم نتیجه را حساب میکنیم. یعنی برای جلوگیری از خطای پاریز، بجای ذخیره و کار با احتمال، با لگاریتم آن کار میکنیم.

نکته دیگری که در اینجا لازم به توضیح است محاسبه احتمال یک کلمه در ابتدای جمله است. یعنی وقتی که کلمات گذشته ای وجود نداشته باشند. در این حالت میتوانیم به سادگی از مدل **unigram** استفاده کنیم اما برای به دست آوردن کارایی بهتر، فرض میکنیم شبه کلمه  $\langle S \rangle$  در ابتدای هر جمله (آموزشی / آزمایشی) وجود دارد که برای محاسبه احتمال ها مانند کلمات عادی شمرده میشود. بنابراین احتمال کلمه  $w$  در ابتدای یک جمله با مدل **bigram** برابر

است با  $P(w | \langle s \rangle)$  و اگر از مدل trigram استفاده کنیم، احتمال کلمه اول جمله برابر است با  $P(w | \langle s \rangle, \langle s \rangle)$  و به همین ترتیب.

آموزش مدل‌های N-gram از طریق شمارش و نرمال کردن انجام میشود. در مدل‌های احتمالاتی نرمال کردن به معنای تقسیم به یک مجموع است به طوری که احتمالات نتیجه در محدوده ۰ و ۱ قرار گیرند. به عنوان مثال برای محاسبه (آموزش) یک مدل bigram ساده، با در دست داشتن یک مجموعه نوشته جات، تعداد بارهای تکرار هر bigram (زوج کلمه) را می‌شماریم و سپس آن را بر مجموع دفعات تکرار تمام bigram هایی که با کلمه اول شروع میشوند تقسیم میکنیم. یعنی:

$$(۷۰) \quad P(w_n | w_{n-1}) = \frac{C(w_{n-1}w_n)}{\sum_w C(w_{n-1}w)}$$

این رابطه را میتوان ساده کرد، زیرا مجموع همه bigram هایی که با کلمه  $w_{n-1}$  شروع میشوند برابر است با دفعات تکرار این کلمه (یعنی تعداد unigram های  $w_{n-1}$ ). بنابراین:

$$(۷۱) \quad P(w_n | w_{n-1}) = \frac{C(w_{n-1}w_n)}{C(w_{n-1})}$$

و برای حالت کلی تخمین N-gram از رابطه زیر استفاده میشود:

$$(۷۲) \quad P(w_n | w_{n-N+1}^{n-1}) = \frac{C(w_{n-N+1}^{n-1}w_n)}{C(w_{n-N+1}^{n-1})}$$

این رابطه احتمال N-gram را با تقسیم کردن فراوانی مشاهده شده یک دنباله خاص بر فراوانی مشاهده شده پیش دنباله تخمین میزند. این نسبت فراوانی نسبی نامیده میشود. استفاده از فراوانی نسبی به عنوان روشی برای تخمین احتمال مثالی از کاربرد تکنیک MLE است؛ زیرا مجموعه پارامترهای نتیجه شده، احتمال مجموعه آموزشی  $T$  به شرط مدل  $M$  (یعنی  $P(T | M)$ ) را ماکزیمم میکند.

البته به جای استفاده از فراوانی های نسبی روشهای بهتری برای تخمین احتمالات N-gram وجود دارد که در بخش بعد آنها را بررسی میکنیم. اما آن روشهای پیشرفته نیز به نوعی از ایده فراوانی نسبی استفاده میکنند.

## مثال

جدول ۲ فراوانی تکرار bigram های زیر مجموعه ای از نوشته جات آموزشی ما را نشان میدهد که شامل ۱۶۱۶ نوع کلمه و تقریباً ۱۰۰۰۰ جمله است. همانطور که ملاحظه میکنید بیشتر مقادیر صفر هستند. در واقع چون تنها هفت کلمه را انتخاب کرده ایم، این ماتریس



خلوت تر نیز شده است. فراوانی تکرار این هفت کلمه به این شرح است: او: ۳۴۳۷، میخواهد: ۱۲۱۵، در: ۳۲۵۶، مسابقات: ۹۳۸، کشتی: ۲۱۳، شرکت: ۱۵۰۶، کند: ۴۵۹.

کند	شرکت	کشتی	مسابقات	در	میخواهد	او	
0	0	0	13	0	1087	8	او
6	8	6	0	786	0	3	میخواهد
12	0	3	860	10	0	3	در
52	2	19	0	2	0	0	مسابقات
1	120	0	0	0	0	2	کشتی
0	0	0	0	17	0	19	شرکت
0	1	0	0	0	0	4	کند

جدول ۲- فراوانی تکرار bigram های زیر مجموعه ای از نوشته جات آموزشی

جدول ۳ احتمالات bigram ها را پس از نرمال سازی - یعنی تقسیم هر سطر بر فراوانی تکرار unigram مناسب - نشان میدهد.

کند	شرکت	کشتی	مسابقات	در	میخواهد	او	
0	0	0	.0038	0	.32	.0023	او
.0049	.0066	.0049	0	.65	0	.0025	میخواهد
.0037	0	.00092	.26	.0031	0	.00092	در
.055	.0021	.020	0	.0021	0	0	مسابقات
.0047	.56	0	0	0	0	.0094	کشتی
0	0	0	0	.011	0	.013	شرکت
0	.0022	0	0	0	0	.0087	کند

جدول ۲- فراوانی تکرار bigram های زیر مجموعه ای از نوشته جات آموزشی

### چند ویژگی مدل‌های N-gram

از ویژگی های قابل انتظار مدل‌های N-gram این است که دقت (کارایی) مدل با افزایش مقدار N افزایش می یابد. علیرغم این واقعیت، عملاً در بیشتر کاربردها از مدل‌های bigram یا حداکثر trigram استفاده میشود؛ زیرا مدل‌های مرتبه بالاتر از ۳ برای آموزش مناسب، احتیاج به مجموعه نوشته جات بسیار بزرگی دارند و در غیر این صورت نمیتوان تخمین های مناسبی برای احتمالات به دست آورد. علاوه بر این، مدل‌های مرتبه بالاتر از ۳ بسیار بزرگ هستند و ذخیره و استفاده از آنها احتیاج به حافظه زیادی دارد. مرتبه حافظه مورد نیاز یک مدل N-gram برابر است با تعداد ترکیب‌های مختلف N کلمه ای که دارای حد بالای  $V^N$  است (V اندازه لغتنامه مورد استفاده است).

ویژگی دیگر مدل‌های N-gram وابستگی زیاد آنها به نوشته جات (به خصوص نوع و اندازه) آموزشی است. یکی از روش‌های معمول برای مشاهده عملکرد کیفی یک مدل N-gram، تولید رشته‌های تصادفی کلمات با استفاده از مدل است. به این ترتیب که اولین کلمه بر اساس احتمال unigram انتخاب می‌شود و سپس دومین کلمه با مدل bigram و پس از آن کلمات بعدی می‌توانند با مدل bigram یا trigram تولید شوند.

نظر به اینکه احتمالات در یک مدل آماری همچون N-gram از یک مجموعه آموزشی استخراج می‌شوند، این مجموعه را باید با دقت انتخاب کرد. اگر مجموعه آموزشی تنها مربوط به یک زمینه (موضوع) خاص باشد، احتمالات نتیجه شده نمی‌توانند برای جملات جدید به شکلی مناسب تعمیم یابند. از طرفی اگر مجموعه نوشته جات آموزشی بسیار عمومی و کلی باشد، ممکن است احتمالات نتیجه شده برای آن کاربرد خاص مناسب نباشند.

برای آموزش و سپس محاسبه کارایی یک مدل، همانند سایر مسائل یادگیری محاسباتی، مجموعه نوشته جات موجود را به دو مجموعه مجزای آموزشی و آزمایشی تقسیم می‌کنیم؛ مدل را بر اساس مجموعه آموزشی می‌سازیم و سپس کارایی آن را با استفاده از معیاری به نام سرگشتگی (Perplexity) - که معیاری مشابه با آنتروپی است - روی مجموعه آزمایشی محاسبه می‌کنیم. البته در برخی موارد به بیش از یک مجموعه آزمایشی احتیاج داریم. مثلاً فرض کنید که چند مدل مختلف را در اختیار داریم و هدف این است که بهترین را انتخاب و سپس کارایی آن را محاسبه کنیم. البته باید توجه کرد که برای مقایسه کارایی این مدل‌ها لازم است از تست‌های آماری استفاده شود تا معلوم شود تفاوت بین دو مدل تا چه حد قابل توجه است.

## هموارسازی

مسئله اصلی مدل‌های استاندارد N-gram که تا حالا بررسی کردیم، این است که تخمین‌های احتمال آنها با استفاده از یک مجموعه نوشته جات آموزشی "محدود" به دست می‌آیند که به هر حال نمی‌تواند شامل تمام ترکیب‌های (N-gram های) مجاز و با معنی زبان باشد. یعنی مثلاً ماتریس bigram با هر مجموعه آموزشی که محاسبه شود، ماتریسی خلوت است. یعنی شامل تعداد بسیار زیادی bigram با احتمال صفر است که بخشی از این bigram واقعا مجاز و بامعنی هستند و لذا باید احتمال غیر صفری داشته باشند. این مشکل را در مثال جدول ۳ ملاحظه می‌کنید، به عنوان نمونه در سطر اول، ترکیب «او شرکت» علی‌رغم معقول بودن، احتمالی صفر دارد (یعنی از دید مدل غیر ممکن تلقی می‌شود).

حتی اگر مشکل فراوانی‌های صفر وجود نداشته باشد، روش MLE وقتی مقدار فراوانی‌ها کوچک باشد، باز هم تخمین‌های ضعیفی را نتیجه می‌دهد. برای حل این مشکلات از تکنیک‌های هموارسازی (smoothing) استفاده می‌شود که در آنها به احتمالات صفر (و پایین)

مقادیر غیر صفری (و بیشتری) منسوب میشود و به همین نسبت از احتمالات زیاد کاسته میشود (تا جمع احتمالات ۱ باقی بماند). در دو بخش زیر، دو الگوریتم (ساده و پیشرفته) هموارسازی را به همراه مثال شرح میدهیم.

### هموارسازی جمع با یک

اولین و ساده ترین روشی که برای هموارسازی - یا اجتناب از فراوانی های صفر - به ذهن میرسد، روش «جمع با یک» است و اگرچه که کارایی خوبی ندارد و در عمل چندان مورد استفاده قرار نمیگیرد اما مفاهیم کلی و اساسی هموارسازی را در بر دارد و بررسی آن به درک و توصیف الگوریتم های پیچیده تر کمک میکند.

در ابتدا به منظور سادگی، کاربرد روش جمع با یک را برای هموارسازی احتمالات unigram شرح میدهیم. تخمین ML معمولی ناهموار، با تقسیم کردن دفعات تکرار یک کلمه بر تعداد کل کلمات موجود در مجموعه نوشته جات آموزشی (که با  $N$  نمایش داده میشود) به دست می آید:

$$P(w_x) = \frac{c(w_x)}{\sum_i c(w_i)} = \frac{c(w_x)}{N}$$

روشهای مختلف هموارسازی بر مبنای یک فراوانی تصحیح شده  $c^*$  هستند. تصحیح فراوانی برای روش جمع با یک اینطور تعریف میشود: مقدار فراوانی پس از جمع با یک، در ضریب نرمال سازی  $\frac{N}{N+V}$  ضرب میشود که  $V$  اندازه لغتنامه (یا تعداد نوع کلمات، یا به عبارتی دیگر تعداد unigram های منحصر به فرد) است. چون به هر کلمه (نوع کلمه) یکی اضافه میکنیم، به اندازه کل کلمات موجود به اندازه  $V$  اضافه میشود و بنابراین فراوانی تصحیح شده برای روش هموارسازی جمع با یک اینطور تعریف میشود:

$$c_i^* = (c_i + 1) \frac{N}{N + V} \quad (۷۳)$$

و این فراوانی ها با تقسیم (نرمال) کردن بر  $N$  تبدیل به احتمال میشوند. در واقع ایده اصلی همه الگوریتم های هموارسازی کاستن از فراوانی های غیر صفر (و دارای مقدار زیاد) با هدف انتساب فراوانی (و در نتیجه احتمال) غیر صفر به فراوانی های صفر است. بنابراین، بجای استفاده از فراوانی های کاهش یافته  $c^*$ ، برخی محققین ترجیح میدهند که الگوریتم های همواره سازی را بر حسب یک ضریب کاهش  $d_c$  تعریف کنند که برابر است با نسبت فراوانی کاهش یافته به فراوانی اصلی:  $d_c = \frac{c^*}{c}$ .

و به عنوان روشی دیگر، میتوانیم احتمالات  $p^*$  را مستقیماً از فراوانی های اصلی حساب کنیم:

$$p_i^* = \frac{c_i + 1}{N + V}$$

حال که روش جمع با یک را برای unigram توضیح دادیم، آنرا به bigram مثال قبل (جدول ۲) اعمال میکنیم. جدول ۴ فراوانی های با یک جمع شده (هموار شده) و جدول ۵ احتمالات مربوطه را نشان میدهد.

کند	شرکت	کشتی	مسابقات	در	میخواهد	او	
1	1	1	14	1	1088	9	او
7	9	7	1	787	1	4	میخواهد
13	1	4	861	11	1	4	در
53	3	20	1	3	1	1	مسابقات
2	121	1	1	1	1	3	کشتی
1	1	1	1	18	1	20	شرکت
1	2	1	1	1	1	5	کند

جدول ۴- فراوانی های تکرار همواره شده مثال جدول ۲ با روش جمع با یک

کند	شرکت	کشتی	مسابقات	در	میخواهد	او	
.00020	.00020	.00020	.0028	.00020	.22	.0018	او
.0025	.0032	.0025	.00035	.28	.00035	.0014	میخواهد
.0027	.00021	.00082	.18	.0023	.00021	.00082	در
.021	.0012	.0078	.00039	.0012	.00039	.00039	مسابقات
.0011	.066	.00055	.00055	.00055	.00055	.0016	کشتی
.00032	.00032	.00032	.00032	.0058	.00032	.0064	شرکت
.00048	.00096	.00048	.00048	.00048	.00048	.0024	کند

جدول ۵- احتمالات هموار شده مثال جدول ۲ با روش جمع با یک

یادآوری میکنیم که احتمالات bigram با نرمال سازی هر سطر با فراوانی unigram محاسبه میشود (رابطه ۷۱). برای فراوانی های هموار شده (با روش جمع با یک) لازم است فراوانی هر unigram را به اندازه تعداد کلمات منحصر به فرد (اندازه لغتنامه که با  $V$  نشان داده میشود) افزایش دهیم:

$$(۷۴) \quad P(w_n | w_{n-1}) = \frac{C(w_{n-1}w_n) + 1}{C(w_{n-1}) + V}$$

چون در این مثال  $V = 1616$  است، فراوانی های unigram (تک کلمه ها) به این شکل افزایش می یابند:

$$\text{او: } ۵۰۵۳ = ۱۶۱۶ + ۳۴۳۷ \quad \text{میخواهد: } ۲۸۳۱ = ۱۶۱۶ + ۱۲۱۵$$

$$\begin{aligned} 4872 &= 1616 + 3256 & \text{در:} \\ 1829 &= 1616 + 213 & \text{کشتی:} \\ 2075 &= 1616 + 459 & \text{کند:} \end{aligned}$$

$$\begin{aligned} 2554 &= 1616 + 938 & \text{مسابقات:} \\ 3122 &= 1616 + 1506 & \text{شرکت:} \end{aligned}$$

و نتیجه که احتمالات هموار شده است در جدول ۵ نشان داده شد. برای اینکه ببینیم یک الگوریتم هموار سازی چگونه فراوانی های اصلی را تغییر داده است، معمولا ماتریس فراوانی را بازسازی میکنیم. فراوانی های تصحیح شده که با رابطه ۷۳ محاسبه میشوند برای این مثال در جدول ۶ نشان داده شده اند. توجه کنید که این روش (جمع با یک) فراوانی ها را به شدت تغییر داده است. مثلا (میخواهد در)  $C$  از ۷۸۶ به ۳۳۳ کاهش یافته است و در فضای احتمال نیز (میخواهد | در)  $P$  از ۰.۶۵ (در حالت ناهموار) به ۰.۲۸ (در حالت هموار) کم شده است. ضریب کاهش  $d_c$  نیز نشان میدهد که فراوانی برای هر کلمه اول چند برابر شده است. مثلا در اینجا bigram هایی که با «کشتی» شروع شده اند، تقریبا ۸ برابر کم شده اند.

او	میخواهد	در	مسابقات	کشتی	شرکت	کند
او	6	740	10	0.68	0.68	0.68
میخواهد	2	0.42	331	0.42	4	3
در	3	0.69	8	594	0.69	9
مسابقات	0.37	0.37	1	0.37	1	20
کشتی	0.36	0.12	0.12	0.12	15	0.24
شرکت	10	0.48	9	0.48	0.48	0.48
کند	1.1	0.22	0.22	0.22	0.44	0.22

جدول ۶- فراوانی های بازسازی شده مثال جدول ۲ با روش جمع با یک

تغییرات شدید مقادیر فراوانی و احتمال در اثر جابجایی زیاد توزیع احتمال به طرف احتمالات صفر (و نزدیک به صفر) میباشد. مشکل از اینجا ناشی میشود که مقدار ۱ برای جمع با هر فراوانی به طور اختیاری انتخاب شده است. البته میتوان با جمع کردن مقادیر کوچکتر (مثلا نیم یا یک هزارم) تا حدودی از این مشکل اجتناب کرد؛ اما برای انتخاب این عدد هم باید معیاری معقول وجود داشته باشد. به طور کلی هموارسازی جمع با یک روشی ضعیف است و در مقایسه با روشهای پیشرفته تر عملکرد بسیار بدتری دارد. همچنین تحقیقات نشان داده اند که واریانس های فراوانی های تولید شده با روش جمع با یک واقعا بدتر از چیزی است که با روش MLE معمولی (ناهموار) به دست آید.

### هموارسازی Witten-Bell

روش هموار سازی یا کاهش Witten-Bell تنها کمی پیچیده تر از روش جمع با یک است، اما کارایی بسیار بهتری دارد و یکی از روشهای پر استفاده در سیستم های بازشناسی گفتار

است. از روشهای پرستفاده دیگر میتوان از Good-Turing و Katz نام برد. روش Witten-Bell بر پایه یک ایده ساده اما زیرکانه درباره رویدادهای با فراوانی صفر (مشاهده نشده) است: یک رویداد یا N-gram با فراوانی صفر وقتی اتفاق می افتد اولین باری است که آن را مشاهده میکنیم؛ پس احتمال دیدن یک N-gram با فراوانی صفر میتواند با احتمال دیدن آن برای اولین بار مدل شود؛ که این یک مفهوم تکرار شونده در پردازش آماری زبان است. ایده اصلی این است که از فراوانی چیزهایی که تنها یک بار دیده شده اند (رخ داده اند) برای تخمین فراوانی چیزهایی که تا حالا دیده نشده اند (رخ نداده اند) استفاده کنیم. این ایده در برخی الگوریتمهای دیگر هموارسازی و روشهای تشخیص نقش کلمه نیز استفاده شده است. برای محاسبه احتمال دیدن یک N-gram برای اولین بار باید دفعاتی را که برای اولین بار در مجموعه نوشته جات مشاهده میشود حساب کرد؛ و این به سادگی قابل محاسبه است زیرا تعداد N-gram های اولین بار (جدید) یعنی تعداد N-gram های منحصر به فرد در مجموعه داده ها. بنابراین مجموع احتمال همه N-gram های دارای فراوانی صفر را با تعداد انواع N-gram ها تقسیم بر مجموع تعداد کل و تعداد انواع N-gram ها تخمین میزنیم:

$$(75) \quad \sum_{i: c_i=0} p_i^* = \frac{T}{N+T}$$

دلیل اینکه با مجموع تعداد کل و تعداد انواع N-gram ها نرمال میکنیم این است که مجموعه نوشته جات را میتوانیم به عنوان یک سری از رویدادها در نظر بگیریم: یک رویداد برای هر کلمه و یک رویداد برای هر نوع جدید. بنابراین رابطه 75 در واقع یک تخمین ML برای رخ دادن یک رویداد جدید است. باید توجه کرد که تعداد انواع مشاهده شده  $T$  با تعداد کل (یا اندازه لغتنامه)  $V$  که قبلا در روش هموارسازی جمع با یک استفاده شد تفاوت دارد.  $T$  تعداد انواعی است که تا کنون دیده ایم، در حالیکه  $V$  تعداد "همه انواع ممکن" است.

رابطه 75 مجموع احتمال N-gram های دیده نشده است که باید بین همه N-gram های دارای فراوانی صفر تقسیم شود. میتوانیم این تقسیم را به طور مساوی انجام دهیم. فرض کنید  $Z$  تعداد N-gram های منحصر به فرد با فراوانی صفر باشد، حالا هر unigram صفر سهم مساوی از توزیع احتمال میگیرد:

$$(76) \quad Z = \sum_{i: c_i=0} 1$$

$$(77) \quad p_i^* = \frac{T}{Z(N+T)}$$

احتمال مجموع N-gram های صفر، که از رابطه 75 حساب میشود، با کاهش دادن احتمال N-gram های دیده شده تامین میشود:

$$(۷۸) \quad p_i^* = \frac{c_i}{N+T} \quad \text{if } c_i > 0$$

به بیان دیگر، فراوانی های هموار شده را میتوانیم مستقیماً از رابطه زیر محاسبه کنیم:

$$(۷۹) \quad c_i^* = \begin{cases} \frac{T}{Z} \frac{N}{N+T}, & \text{if } c_i = 0 \\ c_i \frac{N}{N+T}, & \text{if } c_i > 0 \end{cases}$$

همانطور که ملاحظه میکنید، این روش هموارسازی برای unigram ها بسیار شبیه به روش جمع با یک است. اما وقتی معادلات به bigram تعمیم داده میشوند، تفاوت بزرگی می بینیم که به این دلیل است که در اینجا فراوانی انواع به تاریخچه ای وابسته است. برای محاسبه احتمال یک bigram تا حالا مشاهده نشده  $w_{n-1}w_{n-2}$  از احتمال مشاهده bigram جدیدی که با  $w_{n-1}$  شروع میشود استفاده میکنیم. این یعنی تخمینی که برای bigram های جدید میزنیم به یک تاریخچه کلمه وابسته میشود. کلماتی که در تعداد کمتری از bigram ها اتفاق می افتند، تخمین bigram مشاهده نشده کمتری نسبت بر کلمات فراوان تر ارائه میکنند. این واقعیت را با وابسته کردن  $T$  (تعداد انواع bigram ها) و  $N$  (تعداد کل bigram ها) به  $w_x$  (کلمه قبلی) نشان میدهیم:

$$(۸۰) \quad \sum_{i:c(w_x w_i)=0} p^*(w_i | w_x) = \frac{T(w_x)}{N(w_x) + T(w_x)}$$

دوباره باید این جمع احتمال را بین همه bigram های دیده نشده توزیع کنیم. فرض کنید  $Z$  تعداد bigram های منحصر به فرد با فراوانی صفر باشد، حالا هر bigram صفر سهم مساوی از توزیع احتمال میگیرد:

$$(۸۱) \quad Z(w_x) = \sum_{i:c(w_x w_i)=0} 1$$

$$(۸۲) \quad p^*(w_i | w_{i-1}) = \frac{T(w_{i-1})}{Z(w_{i-1})(N + T(w_{i-1}))} \quad \text{if } c_{w_{i-1}w_i} = 0$$

مانند bigram های غیر صفر، با وابسته کردن  $T$  به تاریخچه، آنها را به همان ترتیب کاهش میدهیم:

$$(۸۳) \quad \sum_{i:c(w_x w_i)>0} p^*(w_i | w_x) = \frac{c(w_x w_i)}{c(w_x) + T(w_x)}$$

از روش کاهش Witten-Bell میتوان به شکل دیگری نیز استفاده کرد. در رابطه ۸۰ احتمالات bigram هموار شده را به کلمه قبلی وابسته کردیم. یعنی  $T(w_x)$  (تعداد انواع

bigram ها) و  $N(w_x)$  (تعداد کل bigram ها) را به  $w_x$  (کلمه قبلی) وابسته کردیم. اما بجای این کار، میتوان هر bigram را به عنوان یک رویداد تکی در نظر گرفت بدون توجه به اینکه از دو کلمه تشکیل شده است. آنگاه  $T$  تعداد انواع همه bigram ها و  $N$  تعداد کل bigram های اتفاق افتاده است. با این کار در واقع بجای کاهش دادن با احتمال شرطی  $P(w_i|w_x)$  داریم از احتمال مشترک  $P(w_i w_x)$  استفاده میکنیم. در اینجا با  $P(w_i w_x)$  دقیقاً همانند یک احتمال unigram برخورد میکنیم. البته این روش کاهش کمتر از روش استاندارد (رابطه ۸۰) به کار رفته است.

## منابع

Corazza, A., De Mori, R., Gretter, R. and Satta, G. (October 1993). "Language modeling using stochastic context-free grammars", *Speech Communication*, vol. 13(1-2), pp. 163-170.

Dugad, R. and Desai, U. B. (May 1996). "A Tutorial on Hidden Markov Models", *Technical Report No. SPANN-96.1*, Indian Institute of Technology, Bombay, India.

Juang, B. H. and Rabiner, L. R. (Sep. 1990). "The Segmental K-Means Algorithm for Estimating the Parameters of Hidden Markov Models", *IEEE Trans. on Acoustics, Speech and Signal Processing*, vol. 38(9), pp. 1639-1641.

Katz, S. M. (March 1987). "Estimation of Probabilities from Sparse Data for the Language Model Component of a Speech Recognizer", *IEEE Trans. on ASSP*, vol. 35(3), pp. 400-401.

Klatt, D. H. (1987). "Review of text-to-speech conversion for English", *Journal of the Acoustical Society of America*, vol. 82(3), pp. 737-793.

Lu, Z., Bazzi, I., Kornai, A., Makhoul, J., Natarajan, P. and Schwartz, R. (1999). "A Robust Language-Independent OCR System", *Proceedings of 27<sup>th</sup> AIPR Workshop: Advances in Computer-Assisted Recognition, SPIE Proceedings*.

Manning, C. D. and Schütze, H. (1999). *Foundations of Statistical Natural Language Processing*, The MIT Press.

Mitchell, Tom M. (1997). *Machine Learning*, McGraw-Hill.

Rabiner, L. R. (1989). "A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition", *IEEE Proceedings*, vol. 77(2), pp. 257-286.

Rosenfeld, R. (2000). "Two decades of statistical language modeling: Where do we go from here?", *Proceedings of the IEEE*, vol. 88, pp. 1270-1278.



Rosenfeld, R., Chen, S. F. and Zhu, X. (January 2001). "Whole-sentence exponential language models: a vehicle for linguistic-statistical integration", *Computer Speech & Language*, vol. 15(1), pp. 55-73.

Vinciarelli, A., Bengio, S. and Bunke, H. (June 2004). "Offline Recognition of Unconstrained Handwritten Texts Using HMMs and Statistical Language Models", *IEEE Trans. on PAMI*, vol. 26(6), pp. 709-720.

Viterbi, A. J. (1967). "Error Bounds for Convolutional Codes and an Asymptotically Optimum Decoding Algorithm", *IEEE Trans. on Information Theory*, IT-13(2), pp. 260-269.

Werner, S., Eichner, M., Wolff, M. and Hoffmann, R. (July 2004). "Toward Spontaneous Speech Synthesis - Utilizing Language Model Information in TTS", *IEEE Trans. on Speech and Audio Processing*, vol. 12(4), pp. 436-445.