

مدل کردن کشش زمانی واج برای سیستم تبدیل متن به گفتار فارسی

به کمک شبکه عصبی

مجید نم نبات، عباس کوچاری، سعید رحیمی، بهرنگ قاسمی زاده، مرتضی سالاریان¹
{namnabat,kochari,qasemizadeh,rahimi,salaraan}@digitalclone.net

چکیده

در این مقاله از یک شبکه عصبی چند لایه پرسپترون برای مدل کردن کشش زمانی استفاده شده است. بدین منظور ابتدا 60 دقیقه گفتار برچسب گذاری شد و از آن برای داده های آموزشی و آزمایشی استفاده گردید. ویژگیهای مورد استفاده برای پیشگویی شبکه عصبی فقط شامل ویژگیهای واجی و هجایی می باشند و از اطلاعات نحوی که مربوط به نقش کلمه در جمله می باشد استفاده نشده است. بعد از انتخاب ویژگیهای اولیه، بازنگری بر روی ویژگیها صورت گرفت و تعداد ویژگیها کاهش داده شد. سپس از دو روش تغییر مقیاس تقسیم بر ماکزیمم و تبدیل Z-SCORE برای نرمال کردن مقدار کشش زمانی استفاده گردید. برای ارزیابی از معیار میزان همبستگی میان کشش زمانی پیشگویی شده و مقدار واقعی استفاده شده است که برای روش تقسیم بر ماکزیمم و تبدیل Z-SCORE به ترتیب برابر 77.82% و 80.39% بدست آمد.

کلمات کلیدی: شبکه های عصبی چند لایه پرسپترون²، تبدیل Z-SCORE، کشش زمانی، تبدیل متن به گفتار.

1- مقدمه

یکی از بخشهای مهم سیستمهای تبدیل متن به گفتار³، تخمین نوای متن ورودی می باشد. اجرای اصلی نوای گفتار عبارتند از میزان کشش واجها، منحنی پیچ جملات و شدت صدا و عناصر فرعی آن ریتم صحبت، سرعت بیان و تن صدا می باشند. بطور کلی پیشگویی پارامترهای نوای گفتار بسیار دشوار و پیچیده می باشد. این پیچیدگی اولاً ناشی از نامشخص بودن تأثیرات متقابل اجزای مختلف نوا بر یکدیگر و دوماً ناشی از وقایع نادر⁴ می باشد. بررسیهای تجربی نشان دهنده ارتباط پیچیده و نامشخص میان اجزای مختلف نوا برای انتقال مفاهیم گفتار می باشند. بواقع انتقال مفاهیم گفتار توسط گویندگان با ترکیب تمامی پارامترهای مختلف نوا به شنوندگان صورت می گیرد. بطور مثال تمامی اجزای نوای گفتار در تعیین اجزای با اهمیت گفتار نسبت به بخشهای دیگر نقش دارند. مسئله دیگر وجود وقایع نادر برای مدل کردن نوای گفتار می باشد. این وقایع هر چند به ندرت اتفاق می افتند ولی مجموع فرکانس وقوع تمامی آنها به اندازه ای می باشد که همیشه حداقل یک واقعه نادر در هر جمله روی می دهد، لذا این وقایع نباید نادیده گرفته شوند[3]. در سیستمهای تبدیل متن به گفتار کنونی عمدتاً از مدل کردن اجزای پیچیده نوا همچون ریتم و سرعت بیان گفتار صرفنظر می شود و علاوه بر این ارتباط اجزای اصلی نوا بصورت مستقل در نظر گرفته می شود. عموماً در این سیستمها ابتدا میزان

¹ اعضای تیم سیستم تبدیل متن به گفتار کلونایزر

² Multi-Layer Perceptron Neural Networks (MLP)

³ Text To Speech (TTS)

⁴ Rare Events

کشش زمانی واجها تعیین و سپس با استفاده از میزان کشش زمانی واجها، منحنی پیچ جملات تعیین می شود. پیشگویی دقیق میزان کشش زمانی، علاوه بر افزایش میزان طبیعی بودن خروجی سیستم در افزایش کارایی مدل‌های تخمین پیچ نیز به دلیل وابستگی آن به کشش زمانی نقش مهمی دارد.

روشهای تخمین کشش زمانی را می توان به دو دسته عمده روشهای قانون-گرا⁵ و روشهای آماری تقسیم بندی نمود. کلات اولین سیستم تخمین کشش زمانی را بصورت قانون-گرا در سیستم تبدیل متن به گفتار MITalk ارائه نمود [1,2]. این سیستم برای بسیاری از زبانهای دیگر پیاده سازی و مورد استفاده قرار گرفت. امروزه با بزرگتر شدن دادگانهای گفتار، روشهای آماری مورد توجه قرار گرفته اند. روشهای آماری را می توان به دو دسته مدل‌های پارامتری و غیر پارامتری تقسیم بندی نمود. در مدل‌های پارامتری ساختار پردازش ویژگیهای ورودی مشخص است. بطور مثال می توان به مدل مجموعی از حاصلضربها⁶ اشاره نمود [13]. در این روش تاثیر ویژگیهای مختلف در میزان کشش زمانی بصورت یک حاصل جمع از تاثیر ویژگیهای مختلف در نظر گرفته می شود. تاثیر هر ویژگی نیز بصورت حاصلضرب تاثیر این ویژگی نسبت به ویژگیهای دیگر مدل می شود [10]. عمده مدل‌های غیرپارامتری شبکه های عصبی و درختهای تصمیم گیری [10] و منحنی های رگرسیون متعدد تطبیقی⁷ [6] می باشند که برای پیشگویی کشش زمانی مورد استفاده قرار گرفته اند. بررسیها نشان می دهد که شبکه های عصبی و منحنی های رگرسیون متعدد تطبیقی کارایی بهتری نسبت به درختهای تصمیم گیری دارند [7,8]. کمپل واحد تخمین کشش را هجا بدلیل پایدارتر بودن آن نسبت به واج در نظر گرفت. وی برای محاسبه میزان کشش زمانی واجها بر اساس کشش زمانی هجا از مفهوم z-score استفاده نمود [4]. استفاده از z-score برای مدل کردن کشش زمانی امکان مدل کردن سرعت بیان گفتار را نیز براحتی فراهم می نماید [5]. روش MARS روش جدیدی می باشد که برای پیشگویی کشش زمانی کارایی مناسبی در حدود 90% برای زبان آلمانی داشته است. این روش برای مدل کردن داده ها از مجموع یکسری منحنی های رگرسیون استفاده می کند. آموزش این روش نیازمند حجم محاسباتی بالا و زمانبر می باشد [7]. چانگ برای مدل کردن کشش زمانی از درختهای تصمیم گیری و رگرسیون استفاده کرده است. وی بردار ویژگی مربوط به حروف صدا دار را از صامت جدا ساخته و برای هر یک درخت تصمیم گیری جداگانه ای ساخته است [9].

در این پژوهش از شبکه های عصبی چند لایه پرسپترون برای تخمین کشش زمانی واجها استفاده شده است و دو روش مختلف تغییر مقیاس برای تخمین کشش زمانی بررسی شده است. در بخش 2 این مقاله دادگان مورد استفاده، در بخش 3 ویژگیهای انتخاب شده شرح داده می شوند. در بخش 4 مدل شبکه عصبی شرح داده می شوند. بخش 5 به ارزیابی سیستم اختصاص یافته است و در انتها بخش 6 به نتیجه گیری و ارائه پیشنهادات می پردازد.

2- دادگان مورد استفاده

برای تهیه داده های مورد نیاز بمنظور آموزش و تست سیستم از دادگان فارس دات بزرگ استفاده شده است. در این دادگان 100 گویشور با تنوع سن، جنس، میزان تحصیلات مختلف و همچنین با لهجه های متعدد صحبت کرده اند. هر گویشور بطور متوسط حدود 4000 کلمه از متون منتخب از روزنامه های کثیرالانتشار در داخل کشور را با توجه به تنوع موضوعی بطور رسمی خوانده است. ضبط صدا در یک محیط آرام و با فرکانس نمونه برداری 22050 هرتز انجام شده است. برای استخراج داده ها از 60 دقیقه گفتار یک گویشور مرد از این دادگان استفاده شده است. سیگنال صحبت این گوینده بطور دستی در سطح واج به دقت تقطیع شده است. برای استخراج ویژگیهای ورودی لازم است که متن این گفتار برچسب گذاری شود. برای اینکار ابتدا کلمات سازنده متن استخراج و با رشته آوانگاری صحبت آن تراز شده اند. سپس این متن و کلمات آن در سطح واج، هجا، کلمه برچسب گذاری شده اند. لازم به ذکر است که اطلاعات نحوی مورد استفاده قرار نگرفته اند و تخمین کشش زمانی مستقل از نقش کلمات در جمله و براساس اطلاعات واجی و هجاهای کلمه صورت گرفته است.

⁵ Rule-Based

⁶ Sum-of-Products

⁷ Multivariate Adaptive Regression Splines (MARS)

3- ویژگیها

برای ارزیابی کارایی ویژگیهای مختلف، 35 ویژگی از دادگان متنی استخراج گردید. سپس تاثیر هر یک از این ویژگیها در کارایی مدل شبکه عصبی مورد ارزیابی قرار می گیرد تا در نهایت از یک زیر مجموعه بهینه از این ویژگیها برای آموزش هر یک از مدلها استفاده شود. این ویژگیها عبارتند از:

- شماره شناسه واج مورد نظر و واجهای همسایه آن: در زبان فارسی 30 واج وجود دارد که با یک شناسه عددی می توان آنها را از هم تفکیک کرد. طول همسایگی برای هر واج در کلمه، دو در نظر گرفته شده است. این بدین معنی است که علاوه بر شناسه واج مورد نظر، از شناسه های دو واج بعدی و قبلی آن نیز بعنوان ویژگی استفاده می شود.
- ویژگیهای آوایی واج مورد نظر و واجهای همسایه: این ویژگیها عبارتند از کلاس آوایی هر واج، محل تولید آن، بیواک، واکنار یا صدا دار بودن آن. برای صدا دارها، ویژگی کلاس آوایی متمایز کننده واجهای صدا دار کشیده و کوتاه از یکدیگر در نظر گرفته شده است. علاوه بر این ویژگی محل تولید برای صدا دارها آنها را نظر مشابهت میان ویژگیهای ثانویه آنها از یکدیگر متمایز می نماید.
- موقعیت واج مورد نظر در کلمه و هجا و انتهای هجا: ویژگی موقعیت واج در کلمه برای متمایز ساختن واجهای اول، دوم، آخر و یکی مانده به آخر و بقیه موقعیتهای در کلمه از یکدیگر به کار می رود. هر هجا را می توان به سه بخش ابتدایی⁸، هسته⁹ و انتهای¹⁰ تقسیم بندی نمود. ویژگی موقعیت واج در هجا مشخص کننده بخشی از هجا می باشد که واج مورد نظر به آن تعلق دارد. با توجه به ساختار هجابندی زبان فارسی، دو بخش ابتدایی و هسته همواره شامل یک واج می باشند ولی بخش انتهای می تواند حداقل صفر و حداکثر دو واج داشته باشد. در صورتیکه واج مورد نظر متعلق به بخش انتهای باشد، ویژگی موقعیت در انتهای هجا، جایگاه آن را در این بخش مشخص می نماید.
- ویژگیهای آوایی هجای شامل واج مورد نظر و هجاهای همسایه: در زبان فارسی نوع هسته هجا و طول بخش انتهای آن مشخص کننده ویژگیهای هجا می باشند. برای تعیین هسته هجا دو ویژگی نوع و شماره شناسه هسته در نظر گرفته شده است. نوع هسته کشیده یا کوتاه بودن هسته و شماره شناسه متمایز کننده واجهای صدا دار با یک نوع می باشد. این ویژگی علاوه بر هجای شامل واج مورد نظر، برای هجای قبلی و بعدی آن نیز استخراج شده است.
- تعداد هجاها و موقعیت هجای شامل واج در کلمه: برای موقعیت هجا در کلمه سه حالت ابتدا، انتها و وسط کلمه بودن در نظر گرفته شده است. علاوه بر این حالتها، بعضی کلمات می توانند تک هجایی باشند که برای این هجاها نیز یک حالت در نظر گرفته شده است. برای مشخص کردن تعداد هجاهای کلمه شامل واج مورد نظر نیز پنج حالت در نظر گرفته شده است. کلمات با تعداد هجای کمتر از چهار تا با حالتهای جداگانه و کلمات با بیش از چهار هجا نیز با یک حالت مشخص می شوند.

3-1- بازنگری ویژگیها

بعد از آموزش با داده های آموزشی و تست کردن، نتایج حاکی از این نکته داشت که استفاده از بعضی از ویژگیها باعث کندي سیستم و حتي کاهش کارایی می گردد. در این مرحله برای داشتن ویژگیهای مناسب که از آنها بتوان بهترین مدل را آموزش داد، با تغییر ویژگیهای آموزشی بهترین آنها که نتایج بهتری و هزینه کمتری دارند انتخاب گردید. ویژگیهای انتخابی 15 ویژگی می باشند که بعضی از آنها از ترکیب ویژگیها بدست آمده اند و بیانگر ویژگیهای دیگر با صرف هزینه کمتر می باشند. بررسیها نشان می دهند که در میان واجهای همسایه، خصوصیات واج دو تا قبل و همچنین ویژگیهای هجای بعدی دارای تاثیر بسیار کم و قابل اغماض در کارایی مدلها می باشند.

⁸ Onset

⁹ Nuclear

¹⁰ Coda

جدول 1- ویژگیهای استفاده شده برای شبکه عصبی

تعداد نورون	ویژگی	سطح واجشناسی
8	کلاس آوایی واج قبلی	واج
12	محل تولید واج فعلی	
3	ترکیب بیواک، واکدار یا صدادار بودن واج فعلی	
8	کلاس آوایی واج قبلی	
12	محل تولید واج قبلی	
3	ترکیب بیواک، واکدار یا صدادار بودن واج قبلی	
8	کلاس آوایی واج بعدی	
12	محل تولید واج بعدی	
3	ترکیب بیواک، واکدار یا صدادار بودن واج بعدی	
8	کلاس آوایی واج دو تا بعد	
12	محل تولید واج دو تا بعد	
3	ترکیب بیواک، واکدار یا صدادار بودن واج دو تا بعد	
2	مکان واج در انتها	
2	نوع هسته هجا	
3	نوع انتهای هجا	

4- شبکه عصبی MLP

در ابتدا از یک شبکه عصبی چند لایه پرسپترون برای تعیین کشش زمانی واجها استفاده شد. این شبکه عصبی یک شبکه پیشرو با چهار لایه می باشد که دارای یک لایه ورودی، یک لایه خروجی و دو لایه مخفی است. دو لایه مخفی به ترتیب دارای 30 و 6 نورون در نظر گرفته شده است. در لایه ورودی 104 نورون و در لایه خروجی برای تعیین میزان کشش واجها یک نورون در نظر گرفته شده است. تابع فعالیت لایه های ورودی و میانی، تانژانت هایپربولیک و لایه خروجی به صورت خطی انتخاب شده است. برای مقدار دهی اولیه وزنها و مقادیر بایاس لایه اول از الگوریتم گوین-ویدرو¹¹ استفاده شده است [11]. در این الگوریتم از میانگین داده های آموزشی بعنوان مقادیر اولیه وزنها و بایاس استفاده میشود. در جدول 1 ویژگیهای بهینه که برای آموزش شبکه انتخاب شده است همراه با تعداد نورونهای اختصاص یافته به هر ویژگی نشان داده شده است. تمامی این ویژگیها گسسته می باشند و برای ارائه این ویژگیها به شبکه های عصبی بهتر است برای هر مقدار از هر ویژگی یک عضو از بردار ورودی اختصاص داده شود. در اینحالت با توجه به مقدار هر ویژگی، یک عضو از اعضای اختصاص یافته از بردار ورودی به آن ویژگی یک و بقیه صفر در نظر گرفته میشوند.

بردارهای ورودی و خروجی شبکه های عصبی بهتر است که در محدود [0,1] و یا [-1,1] تغییر نمایند. برای تغییر مقیاس مقدار کشش واجها می توان به دو روش عمل کرد. در روش اول می توان بطور ساده مقدار کشش واجها را بر ماکزیمم مقدار کشش تقسیم نمود تا محدوده تغییر آنها [0,1] شود. در روش دوم از تبدیل Z-score برای تغییر محدوده استفاده می شود. برای انجام این تبدیل لازم است که لگاریتم میانگین (μ_i) و لگاریتم انحراف معیار (σ_i) کشش برای هر واج محاسبه گردد. سپس مقدار Z-score کشش زمانی (dur) بصورت زیر بدست می آید:

$$Dur_{z-score(\log)} = (\log(dur) - \mu_{\log(i)}) / \sigma_{\log(i)} \quad (1)$$

¹¹ Nguyen-Widrow

در فرمول فوق i بیانگر یک واج است. در صورتیکه توزیع لگاریتم کشش زمانی نرمال باشد با استفاده از تبدیل فوق در حدود 98% داده ها در محدوده $[-2, 2]$ قرار می گیرند. استفاده از لگاریتم باعث می شود که توزیع نرمال بهتری برای کشش زمانی بدست آید. استفاده از تبدیل Z-score برای تخمین میزان کشش زمانی، امکان مدل کردن سرعت بیان و حالت‌های بیان مختلف را با تغییر میانگین و انحراف معیار هر واج را براحتی فراهم می سازد. با استفاده از این تبدیل تاثیر میزان کشش ذاتی هر واج حذف و امکان تخمین میزان کشش واجها بهتر فراهم می شود.

برای آموزش شبکه از الگوریتم انتشارخطا به عقب جهنده¹² استفاده شده است. الگوریتم انتشار خطا به عقب جهنده علاوه بر کارآیی مناسب در حد الگوریتم لوبنبرگ-مارکواریت¹³ نسبت به آن حافظه کمتری مصرف میکند و برای شبکه های با داده های آموزشی زیاد مناسب می باشد، در صورتیکه برای لایه های میانی شبکه MLP تابع فعالیت سیگموئید انتخاب شود. در روند آموزش آن با الگوریتم انتشار خطا به عقب به دلیل کوچک بودن شیب این تابع برای ورودیهای بزرگ، دامنه گرادیان بسیار کوچک میشود و مقدار تغییر وزنها و بایاس در هر بار آموزش بسیار کوچک می شود. در این الگوریتم دامنه گرادیان یک مقدار ثابت انتخاب میشود و برای تعیین آن از شیب تابع فعالیت استفاده نمی شود [12]. برای جلوگیری از یادگیری بیش از حد داده های آموزشی¹⁴ توسط شبکه ها از 15% داده های آموزشی برای ارزیابی¹⁵ سیستم استفاده شده است.

5- پیاده سازی و نتایج

برای پیاده سازی شبکه های عصبی از نرم افزار مطلب استفاده شده است. معیار ارزیابی تعیین دقت میزان کشش زمانی تخمین زده شده، مقدار همبستگی¹⁶ می باشد. فرمول محاسبه این مقدار در زیر دیده می شود.

$$Correlation = \frac{\sum_i (x - \mu_x)(y - \mu_y)}{\sqrt{\sum_i (x - \mu_x)^2} \sqrt{\sum_i (y - \mu_y)^2}} \quad (2)$$

در جدول 2 مقدار همبستگی بدست آمده براساس روش تغییر مقیاس و Z-score دیده می شود. ویژگیهای استفاده شده در هر دو روش تغییر مقیاس، ویژگیهای جدول 1 می باشد.

جدول 2- مقدار همبستگی دو روش مختلف برای پیشگویی کشش زمانی

روش	تقسیم بر ماکزیمم	Z-score لگاریتمی
کارایی برای مجموعه آموزشی	78.31%	78.99%
کارایی برای مجموعه تست	77.82%	80.39%

6- نتیجه گیری و ارائه پیشنهادات

در این مقاله از شبکه های عصبی چندلایه برای پیشگویی کشش زمانی که از پارامترهای اصلی نوای گفتار است، استفاده شده است. برای اینکار از ویژگیهای استخراجی از دادگان فارس دات برای آموزش و آزمایش مدلها استفاده گردید. انتخاب بهینه ویژگیها نقش بسیار مهمی در تخمین مناسب کشش زمانی دارد. در اینجا مقدار کشش زمانی بصورت کلی و با نرمال کردن کشش زمانی و هم با استفاده از معیار Z-score مورد استفاده قرار گرفت که نتایج حاصل نشان از برتری استفاده از معیار Z-score نسبت به روش دیگر

¹² Resilient Backpropagation

¹³ Levenberg-Marquardt

¹⁴ Overfitting

¹⁵ Validation

¹⁶ Correlation

دارد. علت این برتری در اینست که در معیار Z-score از میانگین و واریانس هر واج برای تخمین کشش زمانی استفاده می شود که از انحراف بیش از حد مقدار پیشگویی شده جلوگیری می کند. پیشنهادی که برای ادامه کار مطرح است دسته بندی آواها به دو دسته صادر دار و بی صدا و تشکیل دو شبکه عصبی مجزا برای هر دسته می باشد. علاوه بر این می توان از دیگر روشهای یادگیری ماشین همچون درختهای تصمیم گیری و یا روش MARS برای تخمین کشش زمانی استفاده کرد. علاوه بر این اضافه کردن ویژگیهای نحوی به ویژگیهای انتخابی به احتمال زیاد باعث افزایش کارایی خواهد شد.

مراجع

- [1] D. H. Klatt, "Review of Text-to-Speech Conversion for English", J. Acoustical Society of America, vol. 82, no. 3, pp. 737-793, 1987
- [2] D. H. Klatt, "Linguistic Uses of Segmental Duration in English: Acoustic and perceptual Evidence", Journal of Acoustic Society of America, vol. 59, pp. 1209-1221, 1976.
- [3] B. Moebius, J. Van Santen, "Modeling Segmental Duration in German Text-to-Speech Synthesis", Spoken Language Conference, Vol. 4, no. 2395-2398, 1996.
- [4] W. N. Campbell, "Predicting Segmental Durations for Accommodation within a Syllable-Level Timing Framework", EuroSpeech, 1993.
- [5] R. Cordoba, J. A. Vallejo, J. M. Montero, J. Arriola-Gutierrez, M. A. Lopez, J. M. Pardo, "Selection of the Most Significant Parameters for Duration Modelling in a Spanish Text-to-Speech System Using Neural Networks", Computer Speech and Language, Vol. 16, pp. 183-203, 2002.
- [6] J. H. Friedman, "Multivariate Adaptive Regression Splines", The Annals of Statistics, vol. 19, no. 1, pp. 1-141, 1991.
- [7] M. Riedi, "Modeling Segmental Duration with Multivariate Adaptive Regression Splines", Eurospeech, vol. 5, pp. 2627-2630, 1997.
- [8] M. Riedi, "A Neural-Network Based Model of Segmental Durational for Speech Synthesis, Eurospeech, vol. 1, pp. 599-602, ESCA, 1995.
- [9] H. Chung, "Duration Models and the Perceptual Evaluation of Spoken Korean", Speech Prosody, pp. 219-222, 2002.
- [10] Lee, S. and Y.H. Oh, "Tree-based modeling of prosodic phrasing and segmental duration for Korean TTS systems", Speech Communication, vol. 28, pp. 283-300, 1999.
- [11] H. Demuth, M. Beale, "Neural Network Toolbox for Use with Matlab", Users Guide Version 3.0, 1998.
- [12] M. Riedmiller, H. Braun, "A Direct Adaptive Method for Faster Back-Propagation Learning: The RPROP algorithm", IEEE International Conference on Neural Networks, San Francisco, CA, 589-591, 1993.
- [13] V. Santen, "Assignment of Segmental Duration in Text-To-Speech Synthesis", Computer, Speech and Language, vol. 8, pp. 95-128, 1994.