



گزارش فنی

تنظیم دقیق چارچوب Wav2Vec2 در زبان فارسی

Fine-Tuned Wav2vec2 In Persian

حمیدرضا اکبری

خردادماه ۱۴۰۱



۱. هدف

یکی از چارچوب های کارآمد در تشخیص خودکار^۱ چارچوب Wav2Vec می باشد. در این چارچوب ابتدا با روش یادگیری خودنظارتی^۲ یادگیری اولیه با داده های بدون برچسب انجام و در گام دوم با کمک یک روش تحت نظارت^۳ مدل گام قبل تنظیم دقیق تری خواهد شد. کارایی مناسب این روش حتی با حجم داده های برچسب دار اندک نکته کلیدی در این چارچوب می باشد. هدف از این گزارش بررسی عملی این چارچوب بر روی یک مجموعه داده و بررسی میزان کارایی این روش در تشخیص گفتار می باشد.

۲. مجموعه داده

Common Voice^۴ یک مجموعه ی از داده های صوتی در دسترس عموم می باشد. این مجموعه داده توسط افراد داوطلب در سطح جهانی پشتیبانی می شود و تولیدکنندگان برنامه های نرم افزار صوتی می توانند از این داده ها برای یادگیری مدل های پیش بینی استفاده نمایند. داده های مذکور بیش از ۹۰ زبان مهم جهان را پشتیبانی می کند. در زیر به برخی از شاخص های این مجموعه داده در زبان فارسی اشاره شده است.

اندازه داده	۹ گیگابایت
فرمت فایل صوتی	Mp3
تعداد فایل های صوتی	۴,۰۵۸

۳. مدل یادگیری خودنظارتی

همانطور که اشاره شد در گام اول چارچوب Wav2Vec، یادگیری مدل پایه بر اساس اطلاعات برچسب نشده و با استفاده از روش های یادگیری خودنظارتی انجام می گیرد. در این تحقیق با انتخاب یکی از مدل های خودنظارتی یادگرفته شده بر روی مجموعه داده های زبان فارسی، تنظیم دقیق مدل با داده های برچسب دار پیگیری و اجرا می شود. لذا انتخاب مدل

¹ Automatic Speech Recognition

² Self-Supervised Learning

³ Supervised Learning

⁴ <https://commonvoice.mozilla.org/en>



خودنظارتی یادگیری شده بر روی زبان فارسی یکی از مراحل انجام این تحقیق می‌باشد. در وب سایت ^۵Hugging Face لیست انواع مدل‌های آموزش داده شده بر روی چارچوب Wav2Vec که توسط کاربران مختلف طراحی و اجرا شده است قابل استفاده می‌باشد. با جستجو می‌توان آنها را شناسایی و برای تنظیم دقیق از آنها استفاده نمود.

۴. پیاده سازی

در این بخش بر اساس مدل پایه انتخاب شده جزئیات پیاده سازی و نتایج حاصل از آن تشریح شده است. کدمنبع^۶ به شرح پیوست برای این هدف طراحی و پیاده سازی شده است. این کد با زبان برنامه نویسی پایتون^۷ تولید و تحت محیط توسعه نرم افزار پایچارم^۸ بر روی سیستم عامل ویندوز پیاده سازی شده است.

در فایل Evaluation میزان خطا پیش‌بینی بر اساس شاخص ^۹WER محاسبه و نمایش داده خواهد شد. نحوه یادگیری مدل پایانه در آدرس https://colab.research.google.com/github/m3hrdadfi/notebooks/blob/main/Fine_Tune_XLSR_Wav2Vec2_on_Persian_ASR_with_%F0%9F%A4%97_Transformers_ipynb قابل دسترسی می‌باشد. ساختار کد به شرح زیر توصیف شده است.

- بارگزاری فایل های صوتی: تابع `speech_file_to_array_fn` مسئولیت بارگزاری و `resample` فایل را با اندازه ۱۶۰۰۰ بر عهده دارد. چالش‌هایی در نحوه بارگزاری فایل های صوتی با فرمت `mp3` وجود داشته که برای حل این چالش از توابع کتابخانه `audio2numpy` استفاده شده است.
- پیش پردازش بر روی داده های متنی برچسب دار: بر روی داده های متنی سه پیش پردازش ذیل انجام می‌شود.
 - حذف کاراکترهای خاص: برخی از کاراکترهای خاص نظیر ؟ از داده های متنی حذف می‌شود. لیست مجموعه این کاراکترها در برنامه مشخص شده است.

^۵ <https://huggingface.co/>

^۶ Source Coe

^۷ Python

^۸ PyCharm

^۹ Word Error Rate



○ تبدیل برخی از کارکترهای خاص : برخی از حرفهای فارسی نظیر ی و ک به جهت کدینگهای متعدد به یک کد تبدیل می شود. لیست این تبدیلهای در برنامه مشخص شده است.

○ عملیات نرمالیز کردن متن : با کمک تابع کتابخانه hazm متن نرمالیزه می شود. به عنوان مثال فاصله ها در برخی از کلمات به نیم فاصله تبدیل می شود.

• تنظیم دقیق : برای انجام عملیات تنظیم دقیق از مدل m3hrdadfi/wav2vec2-large-xlsr-persian استفاده شده است. شاخص ارزیابی استفاده شده برای ارزیابی WER می باشد که بر روی مجموعه داده مورد نظر خطا عدد ۳۴,۸۷٪ می باشد که فاصله زیادی با مقاله اصلی^{۱۰} دارد.

۵. سایر مدل ها

افراد متعددی مدلهای خودنظارتی مختلفی را طراحی و اجرا نموده اند که از منظر کارایی و بر اساس نتایج اعلامی به شرح جدول ذیل می باشد. میزان تمایز برخی از این مدلهای ضرورت بررسی بیشتر در نحوه طراحی و پیاده سازی آنها را افزایش خواهد داد.

مدل	کارایی (WER)
jonatasgrosmann/wav2vec2-large-xlsr-53-persian	۳۰,۱۲
m3hrdadfi/wav2vec2-large-xlsr-persian	۳۴,۸۷
m3hrdadfi/wav2vec2-large-xlsr-persian-v3	۱۰,۳۶

¹⁰ <https://arxiv.org/abs/2006.11477>