



**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<sup>۱</sup> چارچوب Wav2Vec می باشد. در این چارچوب ابتدا با روش یادگیری خودنظارتی<sup>۲</sup> یادگیری اولیه با داده های بدون برچسب انجام و در گام دوم با کمک یک روش تحت نظارت<sup>۳</sup> 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<sup>۴</sup> 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

برای قسمت داده های برچسب دار از مجموعه داده ShEMO<sup>۵</sup> استفاده شده است. اندازه داده ها در حدود ۳ ساعت و ۲۵ دقیقه است که از نمایشنامه های رادیویی برخط<sup>۶</sup> استخراج شده است. این داده ها صحبت های ۸۷ نفر فارسی زبان می باشد.

<sup>۱</sup> Automatic Speech Recognition

<sup>۲</sup> Self-Supervised Learning

<sup>۳</sup> Supervised Learning

<sup>۴</sup> <https://commonvoice.mozilla.org/en>

<sup>۵</sup> <https://www.kaggle.com/datasets/mansourehk/shemo-persian-speech-emotion-detection-database>

<sup>۶</sup> Online



### ۳. پیاده سازی

تنظیم دقیق<sup>۷</sup> مدل یادگیری در زبان فارسی بر روی مجموعه داده هدف توسط افراد مختلف انجام گرفته است. برای تسریع در اجرا و ارزیابی نتایج در این تحقیق از یکی از پیاده سازی‌ها<sup>۸</sup> استفاده و اجرا شده است. پیاده سازی در دو بخش ذیل مدنظر قرار گرفته است.

#### ۳,۱ آموزش و تنظیم دقیق مدل

برای قسمت تنظیم دقیق مدل از محیط کلب<sup>۹</sup> و زبان برنامه نویسی پایتون استفاده شده است. جزئیات پیاده سازی در فایل مورد نظر<sup>۱۰</sup> وجود دارد که در ادامه مراحل اصلی پیاده‌سازی تشریح می‌شود.

#### ۳,۱,۱ تفکیک داده های آموزش از تست

داده های برچسب دار آموزش و تست به نسبت ۱ به ۱۰ و به صورت تصادفی تعیین می‌شوند. تعداد داده های آموزش ۲۵۵۴ و تعداد داده های تست ۲۸۴ می‌باشد. داده های با فرمت CSV ساخته و ذخیره می‌شوند. برای آگاهی از صحت اطلاعات برخی از داده ها به صورت تصادفی نمایش داده می‌شوند.

#### ۳,۱,۲ پیش پردازش داده ها

عملیات پیش پردازش متن به شرح ذیل بر روی داده های آموزش و تست به صورت جداگانه انجام می‌گیرد. تعداد حروف مجموعه داده‌های آموزش و تست جهت بررسی و مقایسه اطلاعات نمایش داده می‌شود.

- **حذف کاراکترهای خاص :** برخی از کاراکترهای خاص نظیر ؟ از داده های متنی حذف می‌شود. لیست مجموعه این کاراکترها در برنامه مشخص شده است.
- **تبدیل برخی از کاراکترهای خاص :** برخی از حرفهای فارسی نظیر ی و ک به جهت کدینگهای متعدد به یک کد تبدیل می‌شود. لیست این تبدیلهای در برنامه مشخص شده است.
- **عملیات نرمالیز کردن متن :** با کمک تابع کتابخانه `hazm` متن نرمالیزه می‌شود. به عنوان مثال فاصله ها در برخی از کلمات به نیم فاصله تبدیل می‌شود.

<sup>7</sup> Fine-Tuned

<sup>8</sup> <https://huggingface.co/m3hrdadfi/wav2vec2-large-xlsr-persian>

<sup>9</sup> <https://colab.research.google.com/>

<sup>10</sup> `Fine_Tune_XLSR_Wav2Vec2_on_Persian_ShEMO_ASRL_with_Transformers_ipynb.ipynb`



### ۳,۱,۳ پیش پردازش فایل های صوتی

از آنجاییکه فایل های صوتی مجموعه داده با نرخ نمونه گیری ۴۸ کیلوهرتز نمونه گیری شده اند در این گام نمونه گیری کلیه فایل های صوتی با نرخ ۱۶ کیلوهرتز انجام می شود.

### ۳,۱,۴ یادگیری

در این مرحله مدل بر اساس اطلاعات برچسب دار آموزش داده می شود. شاخص ارزیابی در طول مدت زمان آموزش میزان خطای کلمات می باشد که تابع `compute_metrics` این محاسبه را انجام می دهد.

### ۳,۲ ارزیابی

هدف از پیاده سازی این بخش تعیین میزان پیش بینی خطا بر اساس مدل آموزش داده شده در مرحله قبل می باشد. این بخش از برنامه با زبان برنامه نویسی پایتون<sup>۱۱</sup> تحت محیط توسعه نرم افزار پایچارم<sup>۱۲</sup> بر روی سیستم عامل ویندوز انجام شده است. میزان خطا بر اساس شاخص `WER`<sup>۱۳</sup> تعیین می گردد. بر این اساس خطا بر روی مجموعه داده مورد نظر ۳۴,۸۷٪ می باشد که فاصله زیادی با مقاله اصلی<sup>۱۴</sup> دارد. با عنایت به چالش نحوه بارگزاری فایل های `mp3` از تابع کتابخانه `audio2numpy` استفاده شده است.

### ۴. سایر مدل ها

افراد متعددی مدل های تنظیم شده 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 مدل ها ضرورت بررسی بیشتر در نحوه طراحی و پیاده سازی آنها را افزایش خواهد داد.

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

<sup>۱۱</sup> Python

<sup>۱۲</sup> PyCharm

<sup>۱۳</sup> Word Error Rate

<sup>۱۴</sup> <https://arxiv.org/abs/2006.11477>