



دانشگاه اصفهان

تمرین چهارم درس پردازش زبان و گفتار
استاد درس: دکتر حمیدرضا برادران کاشانی
دستیاران آموزشی: آیین کوپایی – هاجر مظاهری

تاریخ بارگذاری تمرین: ۱۴۰۳/۰۳/۲۳

تاریخ تحویل تمرین: ۱۴۰۳/۰۴/۱۲

هدف تمرین:

- آشنایی با آموزش نیمه نظارتی (semi-supervised training)
- آشنایی با مدل wav2vec 2.0 از پیش آموزش دیده برای تشخیص خودکار گفتار
- تنظیم دقیق مدل wav2vec 2.0 با داده‌های گفتاری برچسب‌گذاری شده بسیار کم
- ساخت یک سیستم بازشناسی گفتار فارسی

مقاله خوانی:

با توجه به مقاله [3] و [4] به سوالات زیر پاسخ دهید:

- معماری wav2vec 2.0 و فرآیند پیش آموزش را با جزئیات توضیح دهید.
- تفاوت مدل wav2vec2.0 و wav2vec xlsr-53 در چیست؟
- رمزگشایی در مدل wav2vec2.0 با چه الگوریتمی انجام می‌شود؟ روش را توضیح دهید.
- برای بهبود نتایج بدست آمده چه روش یا تکنیکی را پیشنهاد می‌کنید.

مراحل تمرین:

۱- دسترسی به مجموعه داده

برای تنظیم دقیق مدل wav2vec2 از مجموعه داده موزیلا (common voice) استفاده می‌کنیم. جدول ۱ مشخصات مجموعه داده را نشان می‌دهد.

جدول ۱: مجموعه داده

Dataset	Version	language
common_voice	۶.۱	Persian

روش دسترسی به مجموعه داده موزیلا:

- روش آفلاین: به سایت common_voice به آدرس [1] رفته و در قسمت زبان، زبان فارسی (Persian) را انتخاب کنید و با پذیرفتن تعهدات و ثبت ایمیل خود مجموعه داده را دانلود کنید.

۷- فیلتر زمانی

برای آموزش مدل از فایل های صوتی با طول ۴ تا ۶ ثانیه استفاده کنید و برای تست از فایل های با طول کمتر ۱۵ ثانیه استفاده کنید.

(تعداد فایل های آموزش و تست را گزارش کنید.)

۸- padding و Data Collator

جمع آوری داده را تعریف کنید جمع آوری داده برای مدل های گفتاری منحصربه فرد است زیرا ویژگی های ورودی و برچسبها به صورت مستقل بررسی می شوند. ویژگی های ورودی باید توسط استخراج کننده ویژگی و برچسبها توسط توکنایزر مدیریت شوند. همچنین برخلاف اکثر مدل های NLP در مدل Wav2Vec2 طول ورودی بسیار بیشتری از طول خروجی است. به عنوان مثال، یک نمونه با طول ورودی ۵۰۰۰۰ دارای طول خروجی حدود ۱۰۰ است. بنابراین آموزش مدل Wav2Vec2 به padding ویژه نیاز دارد.

۹- معیار ارزیابی

برای ارزیابی سیستم بازشناسی گفتار از معیار نرخ خطای کلمه (WER) استفاده کنید.

۱۰- بارگیری مدل

مدل facebook/wav2vec2-large-xlsr-53 را با استفاده از Wav2Vec2ForCTC بارگیری کنید.

۱۱- تنظیم مدل wav2vec2 با مجموعه داده موزیلا

ابتدا بخش استخراج کننده ویژگی را فریز کنید و سپس با توجه به جدول ۲ زیر مدل را آموزش دهید.

جدول ۲: تنظیمات اجرا

num_train_epochs	5
fp16	True
attention_dropout	0.1
learning_rate	1e-4
warmup_steps	1000

۱۲- پکیج های مورد نیاز برای تنظیم مدل wav2vec2

Datasets >= 1.18.3

Transformers == 4.11.3

Librosa

Jiwer

Torchaudio

Hazm

Num2fawords

۱۳- تنظیم batch_size

با توجه به جدول زیر می‌توانید اندازه batch_size را تنظیم کنید.

جدول ۳: تنظیمات batch_size

Batch-size	سخت افزار
Batch-size \leq 12	T4 16G (Colab)
Batch-size \leq 6	RTX 8 G

نکات تحویل

۱- پاسخ خود را در پوشه ای به اسم NLP_NAME_FAMILY_HW4 و در قالب zip بارگذاری نمایید.

۲- این پوشه باید حاوی موارد زیر باشد:

- کد نوشته شده در قالب یک فایل jupyter notebook
- تشریح و تحلیل نتایج بدست آمده از نظر شما
- آپلود مدل در یک فضای ابری یا مخزن هاینک فیس و ارسال لینک
- پاسخ سوالات مقاله خوانی
- گزارش نتایج بدست آمده از تنظیم دقیق مدل wav2vec مطابق با جدول ۴

جدول ۴: نتایج اجرا

Batch Size
Num Epoch
تعداد فایل‌های مورد استفاده برای آموزش مدل
تعداد فایل‌های مورد استفاده برای تست مدل
مدت زمان آموزش
تعداد پارامترهای مدل
تعداد پارامترهای فعال در زمان آموزش مدل
WER (Word Error Rate)

۳- لازم به ذکر است که رعایت قوانین نگارشی حائز اهمیت است.

- [1] <https://commonvoice.mozilla.org/en/datasets>
- [2] https://huggingface.co/datasets/mozilla-foundation/common_voice_6_1
- [3] Baevski, A., Zhou, Y., Mohamed, A., & Auli, M. (2020). wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems*, 33, 12449-12460.
- [4] Conneau, A., Baevski, A., Collobert, R., Mohamed, A., & Auli, M. (2020). Unsupervised cross-lingual representation learning for speech recognition. *arXiv preprint arXiv:2006.13979*.