

## خلاصه تحلیلی مقاله

Similarity Textual Semantic Cross-Lingual for Transformer Ensemble

### ۱ مسئله اصلی و اهمیت آن

شماحت معنایی متون (STS - Semantic Textual Similarity) به مسئله‌ای گفته می‌شود که در آن هدف، تعیین میزان نزدیکی معنایی بین دو متن به صورت یک مقدار پیوسته، معمولاً در بازه  $[0, 5]$  است. در حالت Cross-Lingual، دو متن مورد مقایسه به زبان‌های متفاوتی نوشته شده‌اند.

این مسئله نقش کلیدی در بسیاری از کاربردهای مهم پردازش زبان طبیعی ایفا می‌کند؛ از جمله بازیابی اطلاعات چندزبانه، سیستم‌های پرسش‌پاسخ، ترجمه ماشینی، موتورهای جستجو و سامانه‌های تولید پاسخ مبتنی بر بازیابی (Retrieval-Augmented Generation). چالش اصلی در این حوزه، ایجاد یک فضای برداری مشترک است که بتواند جملات زبان‌های مختلف را به‌گونه‌ای نمایش دهد که نزدیکی بردارها بازتاب‌دهنده شماحت معنایی واقعی آن‌ها باشد.

مقاله حاضر با تمرکز ویژه بر جفت‌زبان فارسی-انگلیسی، روشی مبتنی بر ترکیب (ensemble) چند مدل ترانسفورمری چندزبانه ارائه می‌دهد که بدون استفاده از ترجمه ماشینی، شماحت معنایی بین‌زبانی را با دقت بالاتری نسبت به روش‌های پیشین محاسبه می‌کند.

### ۲ ورودی‌ها و خروجی‌های مدل

#### ورودی

ة یک جفت جمله  $(s_1, s_2)$

ة جملات می‌توانند به زبان‌های متفاوت باشند (در این مقاله، فارسی و انگلیسی)

#### خروجی

ة یک نمره پیوسته شماحت معنایی:

$$\text{Similarity}(s_1, s_2) \in [0, 5]$$

ة این نمره میزان همپوشانی مفهومی دو جمله از دیدگاه انسانی را نشان می‌دهد.

### ۳ داده‌های مورد استفاده

برای آموزش و ارزیابی مدل‌ها، از مجموعه داده PESTS استفاده شده است. این دیتاست شامل ۵۳۷۵ جفت جمله فارسی-انگلیسی است که هر جفت جمله دارای یک نمره شماحت معنایی تعیین شده توسط ارزیاب انسانی می‌باشد. داده‌ها به سه بخش آموزش، اعتبارسنجی و آزمون تقسیم شده‌اند و به طور خاص برای مسئله شماحت معنایی بین‌زبانی طراحی شده‌اند. استفاده از این دیتاست یکن از نقاط قوت مقاله است، زیرا منابع استاندارد برای ارزیابی STS در زبان فارسی بسیار محدود هستند. علاوه بر این، برای بررسی تعمیم‌پذیری روش پیشنهادی، آزمایش‌هایی روی دیتاست SemEval ۲۰۱۷ برای جفت‌زبان عربی-انگلیسی نیز انجام شده است.

## ۴ روش پیشنهادی مقاله

روش ارائه شده مبتنی بر استفاده هم زمان از چند مدل ترنسفورمری چند زبانه و ترکیب خروجی آنها به منظور بهبود دقت تخمين شbahت معنایی است. بخلاف بسیاری از روش های مبتنی بر ترجمه ماشینی، این روش مستقیماً در یک فضای برداری مشترک عمل می کند.

### مدل های مورد استفاده

سه مدل ترنسفورمری زیر برای ترکیب انتخاب شده اند:

paraphrase-xlm-r-multilingual-v1

paraphrase-multilingual-mpnet-base-v2

mstsbs-paraphrase-multilingual-mpnet-base

این مدل ها با استفاده از دیتاست PESTS برای وظیفه STS ریزنظمی شده اند.

### روش های Ensemble

۱. میانگین گیری از نمرات شbahت: میانگین شbahت کسینوسی خروجی مدل ها به عنوان نمره نهایی استفاده می شود.
۲. الحق بردارها و کاهش بعد: بردارهای خروجی مدل ها الحق شده و پس از کاهش بعد با، PCA شbahت محاسبه می شود.
۳. شبکه عصبی پیش خور: در این روش، یک شبکه عصبی ساده برای تخمين نمره شbahت بر اساس ترکیب بردارها آموزش داده می شود.

### شبه کد کلی

```
for each transformer model:  
    e1 = encode(sentence1)  
    e2 = encode(sentence2)  
    store embeddings  
  
combine embeddings  
output similarity score
```

## ۵ Baseline و باز تولید نتایج مقاله

به عنوان baseline از یک مدل ترنسفورمری منفرد چند زبانه (paraphrase-xlm-r-multilingual-v1) استفاده شد. در این حالت، شbahت معنایی دو جمله با استفاده از شbahت کسینوسی بین های embedding آنها محاسبه گردید. این مدل روی دیتاست PESTS اجرا شد و ضریب همبستگی پیرسون بین نمرات پیش بینی شده مدل و نمرات انسانی دیتاست به عنوان معیار ارزیابی استفاده گردید. این نتیجه مبنای مقایسه برای روش های ensemble قرار گرفت.

Correlation Pearson	روش
۸۲.۰	(Baseline) Model Single
۸۶.۰	(Average) Ensemble
۸۸.۰	dims(۲۰۰) PCA + Ensemble

جدول ۱: مقایسه عملکرد مدل پایه و روش های ensemble روی دیتاست PESTS

## ۶ تحلیل اثر کاهش بُعد (PCA)

برای بررسی اثر کاهش بُعد روی های ensemble embedding از روش تحلیل مؤلفه های اصلی (PCA) استفاده شد. ابعاد مختلفی از فضای ویژگی (۵۰ تا ۴۰۰ بُعد) مورد آزمایش قرار گرفت.

نتایج نشان داد که کاهش بُعد تا یک مقدار بهینه می تواند باعث حذف نویز و بهبود همبستگی با نمرات انسانی شود، در حالی که کاهش بیش از حد بُعد منجر به افت عملکرد می گردد. این تحلیل نشان می دهد که ترکیب ensemble با تکنیک های کاهش بُعد می تواند یک راهکار عملی برای بهبود دقت و کاهش هزینه محاسباتی باشد.

## ۷ نتایج اصلی

نتایج تجربی نشان می دهد که روش های ensemble به طور معناداری عملکرد بهتری نسبت به مدل های منفرد دارند. بهترین عملکرد مربوط به ترکیب ensemble همراه با PCA بوده است. همچنین آزمایش روی جفت زبان عربی-قاچالیسی نشان داد که روش پیشنهادی قابلیت تعمیم مناسبی به زبان های دیگر دارد، هرچند کیفیت و حجم داده نقش مهمی در میزان بهبود ایفا می کند.

## ۸ محدودیت ها و ایده های ادامه

### محدودیت ها

ة افزایش هزینه محاسباتی به دلیل استفاده از چند مدل ترنسفورمیر

ة زمان آموزش و استنتاج بالاتر در روش های ensemble

ة وابستگی عملکرد به کیفیت دیتاست ریز تنظیم

### ایده های ادامه

ة استفاده از knowledge distillation برای فشرده سازی ensemble

ة توسعه دیتاست های بزرگ تر برای زبان فارسی

ة به کارگیری روش در سامانه های عملیاتی چند زبانه