

بسم الله الرحمن الرحيم

موضوع پروژه :

سیستم استخراج اصطلاحات تخصصی حوزه اقتصاد (مبتنی بر BERT)

ارائه دهنگان :

خانم نرگس علی حیدری

خانم فاطمه قیصری

آقای علیرضا فرزانه

استاد پژوهه:

سرکار خانم دکتر پیشگو

موضوع درس:

یادگیری ماشین

مقدمه

در عصر انفجار اطلاعات، دسترسی سریع به مفاهیم کلیدی در یک حوزه تخصصی اهمیت فراوانی دارد ساخت واژه نامه‌های تخصصی به صورت دستی، فرآیندی زمان بر و پرهزینه خواهد بود. هدف این پروژه، خودکارسازی این فرآیند برای متون اقتصادی است تا بتوان با دریافت یک متن خام، لیستی از واژگان کلیدی آن را استخراج کرد. استخراج خودکار واژگان تخصصی (Automatic Glossary Extraction) یکی از چالش‌های کلیدی در پردازش زبان طبیعی است که هدف آن شناسایی اصطلاحات دامنه محور از متون غیرساختاریافته می‌باشد. در این پژوهش، یک (Pipeline) مبتنی بر مدل زبانی DistilBERT و الگوریتم‌های خوش بندی برای استخراج واژگان حوزه اقتصاد طراحی شده است که برخلاف روش‌های سنتی مبتنی بر آمار (مانند TF-IDF)، روش glossex بر درک معنایی و برداری کلمات تمرکز دارد.

هدف پروژه

استخراج خودکار اصطلاحات تخصصی حوزه اقتصاد از یک متن (Corpus) است؛ به طوری که خروجی نهایی یک لیست رتبه بندی شده از "ترم‌های اقتصادی" باشد (شامل تک واژه و عبارت‌های چندواژه‌ای)، بدون این‌که دیتاست برچسب‌خوردهی بزرگ داشته باشیم.

Weakly Supervised (نظرارت ضعیف)

این پروژه weakly supervised است زیرا به جای داشتن برچسب‌های کامل، از seed words (چند کلمه‌ی نمونه‌ی اقتصادی و چند کلمه‌ی عمومی) برای هدایت الگوریتم استفاده می‌کنیم و با سیگнал ضعیف (weak signal) دامنه را تشخیص می‌دهیم.

(Dataset) داده‌ها

برای استخراج واژگان، یک corpus در حوزه‌ی اقتصاد استفاده شده corpus شامل تنها چند جمله در حوزه‌ی اقتصاد (مانند مفاهیم عرضه و تقاضا، تورم، رشد اقتصادی و بازارها)، که می‌توان اندازه‌ی corpus را افزایش داد در واقع افزایش اندازه‌ی corpus نقش مهمی در بهبود پوشش واژگان و پایداری مراحل بعدی pipeline خواهد داشت.

در این پروژه پیکرهای متنی اقتصاد در فایل زیر قرار دارد:

data/raw/economics_sample.txt •

این فایل ابتدا یک نمونه‌ی کوچک بود و سپس برای بهتر شدن خروجی، متن دامنه‌ای بزرگ‌تر به آن اضافه شد (افزایش حجم متن → افزایش تنوع واژگانی → کاندیدهای بهتر).

معماری کلی سیستم (Pipeline Overview)

پروژه به صورت یک خط لوله (Pipeline) طراحی شد که از خامترین داده تا خروجی نهایی را تولید می‌کند:

Preprocess: ۱ پاکسازی، نرمال‌سازی، توکن‌سازی

Baseline (TF-IDF): ۲ استخراج کلمات مهم بر اساس فراوانی/تمایز

Embedding: ۳ ساخت بردار معنایی برای واژگان/لمّاها با مدل زبانی

Clustering: ۴ خوشبندی بردارها برای گروه‌بندی ترم‌های نزدیک معنایی

Filtering (Seed-guided): ۵ انتخاب خوش‌های اقتصادی با مقایسه شباهت به seed های اقتصادی و seed های عمومی

Phrase Extraction + Ranking .6 (فاز ۲): استخراج عبارت‌های چندکلمه‌ای- n (TF-IDF + Embedding similarity) و رتبه‌بندی هیبرید gram

Evaluation: ۷ Gold Glossary ارزیابی در K های مختلف (50/100/200) نسبت به

Phase1

در فاز ۱ یک اسکریپت EDA ساخته شد و خروجی نمودارها ذخیره شد مثل توزیع طول جمله، توزیع طول کلمات، Zipf، رشد واژگان، top frequent words و.... هدف EDA در گزارش:

- نشان دهد متن ورودی چقدر بزرگ است،
- تنوع واژگان چقدر است،

- آیا متن واقعاً دامنه‌ای است یا نه،
- آیا داده برای استخراج واژگان کافی هست یا باید غنیتر شود.

خروجی‌ها و مستندسازی فاز ۱

در فاز ۱ موارد زیر در پروژه قرار دارد:

- گزارش فاز ۱ docs/phase-1-report.md
- ثبت نتایج baseline و برنامه‌ی فاز ۲ (experiment plan)
- ذخیره نمودارهای EDA در مسیر نتایج

فاز ۲: (Phase-2) روش پیشنهادی Weakly Supervised + ارزیابی

فاز ۲ روی "روش اصلی" و "بهبود و مقایسه" تمرکز دارد.

ایده‌ی weak supervision در این پروژه از طریق لیست seed‌ها پیاده‌سازی شد:

. (economics, market) نمونه‌های اقتصادی مثل data/seeds/economics.txt

(.....education, school) نمونه‌های عمومی مثل data/seeds/general.txt

مراحل فاز ۲ :

preprocess

در مرحله اول

- متن به توکن‌ها شکسته شد
- نرمال‌سازی اولیه انجام شد
- توکن‌های نامعتبر حذف شدند

پس از نرمال‌سازی و حذف توکن‌های نامعتبر خروجی این قسمت وارد بخش embedding می‌شود.

پس از پیشپردازش، مشخصات داده به صورت زیر بود:

• تعداد توکن‌های نهایی: 8981 token

data\processed\preprocess.json : خروجی

saliency scoring

در این مرحله میزان تخصصی بودن هر کلمه سنجیده می‌شود. برای این کار از مقیاس zipf استفاده شده است تا فرکانس کلمه در متن پژوهش با فرکانس آن در زبان انگلیسی عمومی مقایسه شود کلماتی که در متن ما تکرار بالا و در زبان عمومی تکرار پایینی دارند امتیاز بیشتری می‌گیرند

```
score = math.log(domain_prob / general_prob + 1e-9)
```

demo\output\top_saliency_terms.csv : خروجی

Embedding

در این قسمت از تکنولوژی ترانسفورمر استفاده شده هر کلمه ورودی رو به یک (vector) یک لیست طولانی از اعداد در فضای ریاضی تبدیل می‌کند در مثل آورده شده هر 8918 توکن حالا در فضای ریاضی یک موقعیت و معنایی دارند ما 2787 embedding برداریم. این کاهش امبدینگ‌ها به دلیل حذف کلمات تکراری - حذف stopword- ریشه‌یابی یا lemmatization است. استفاده از مدل embeddings.py برای تبدیل کلمات به بردارهای 768 بعدی است این مرحله باعث می‌شود کلماتی مثل "Inflation" "economic" در فضای ریاضی به هم نزدیک شوند. برخلاف مدل‌های word2vec این بردارها حاوی اطلاعات context هستند که دقت خوشبندی رو به شدت افزایش می‌دهد.

تعداد امبدینگ‌های نهایی: 2784 embedding

data\processed\lemma_embeddings.json : خروجی

Clustering

توكن‌ها بر اساس شباهت برداری خوش‌بندی شدند تا واژگان هم معنی در کنار یکدیگر قرار گیرند. یعنی کلماتی که بردارهای آنها نزدیک به هم بوده در یک خوش‌بندی قرار گرفتند مثل "inflation" در نهایت در این پروژه 696 خوش‌بندی تشكیل شد که هر کدام نماینده یک مفهوم بالقوه بودند در این مرحله از Agglomerative

تعداد خوش‌بندی نهایی : cluster 696

خروجی: data\processed\clusters.json

Filtering

برای تشخیص خوش‌بندی‌های مرتبط با اقتصاد، از یک روش مبتنی بر seed words استفاده شد:

در این قسمت از دو seed استفاده کردیم

• واژگان مرتبط با اقتصاد Economic seeds:

• واژگان عمومی و غیرتخصصی General seeds:

این دقیقاً همان جایی است که پروژه weakly supervised می‌شود: ما برچسب دقیق نداریم، ولی با یک سیگنال کم‌هزینه (seed list) جهتدهی می‌کنیم.

در filtering با مقایسه شباهت خوش‌بندی‌ها به Seed‌های اقتصادی، واژگان نهایی استخراج می‌شوند.

اگر میانگین شباهت اقتصادی > عمومی باشد، خوش‌بندی "اقتصادی" در نظر گرفته می‌شود.

فرمول محاسبه شده در پروژه cosine similarity می‌باشد

$$\text{Cosine} = \frac{\mathbf{a} \cdot \mathbf{b}}{|\mathbf{a}| \cdot |\mathbf{b}|}$$

خروجی: data\processed\final_terms.json

phrases

برای نزدیک شدن به glossary واقعی که معمولاً multi-word term است در واقع ما همیشه با تک واژه ها رو به رو نیستیم می خواهیم اصطلاحات تخصصی رو استخراج کنیم در این صورت به وسیله n-gram ها که ما در این پروژه bigram , trigram داریم عبارات رو استخراج می کنیم.

- استخراج کandidهای عبارتی با TF-IDF + n-gram انجام شد

خرожی : data\processed\phrase_candidates.json

ranked_phrases

میانگین بردار embedding های مرحله phrases را بدست آورده و میانگین شbahت آن ها با seed general و cosine از طریق شbahت انجام شده و score بروجی بررسی می شود اگر میانگین شbahت اقتصادی > عمومی باشد، خوش "اقتصادی" در نظر گرفته می شود

خرожی : data\processed\rank_phrases.json

Hybrid Ranking

ترکیب TF-IDF و Embedding

برای بهتر شدن دقت نسبت به TF-IDF :

- یک امتیاز هیرید تعریف شد که هم "اهمیت آماری (TF-IDF)" و هم "شbahت معنایی به seed های اقتصادی" را لحاظ می کرد.
- این باعث شد خروجی ها از حالت صرفاً "کلمات پرتکرار عمومی" فاصله بگیرند و اقتصادی تر شوند.
- این ترکیب طبق فرمول زیر محاسبه می شود

$$\text{Hybrid} = \alpha * \text{tfidf_n} + (1 - \alpha) * \text{emb_n}$$

الف = 0-1

الف : 0.8 به دلیل قابلیت اطمینان بالاتر

خرожی : data\processde\Ranked_Hybrid.json

Baseline

به عنوان معیار مقایسه baseline ، از روش TF-IDF برای استخراج واژگان پر تکرار استفاده شد. این روش بدون درک معنایی embedding، صرفاً بر اساس فراوانی و تکرار واژگان عمل می کند.

demo\outputs\tfidf_baseline_top_terms.csv خروجی:

Evaluation

در فاز ۲ اسکریپت ارزیابی طوری به روزرسانی شد که در چند K گزارش بدهد:

K=200	K=100	K=50	Glossex model
0.045	0.040	0.060	precision
0.090	0.040	0.030	recall
0.060	0.040	0.040	F1

K=200	K=100	K=50	Tfidf model
0.010	0.020	0.040	precision
0.020	0.020	0.020	recall
0.013	0.020	0.027	F1

تحلیل نتایج:

- در K های بزرگتر، Recall روش پیشنهادی افزایش پیدا می کند (یعنی اصطلاحات مرجع بیشتری پوشش داده می شود).
- روش پیشنهادی نسبت به baseline در اکثر K ها بهتر است خصوصاً در Recall.
- این با هدف glossary extraction هم سازگار است: در بسیاری از سناریوهای Recall مهمتر است چون حذف موارد بد ساده تر از پیدا کردن اصطلاحات جا افتاده است.

ساختار نهایی (Project Structure)

- داده خام data
 - کد اصلی پروژه src pipeline ها، baseline
 - اسکریپت‌های اجرایی تولید خروجی scripts
 - ارزیابی دو مدل پایه و مدل فاز 2 پروژه evaluation
 - گزارش‌های فاز ۱ و فاز ۲ docs
 - خروجی‌های نمایشی demo/outputs
 - تحلیل‌ها و EDA تعاملی notebooks
- جمع‌بندی نهایی

در این پروژه یک سیستم استخراج واژگان اقتصاد ساخته شد که:

- از متن دامنه‌ای اقتصاد به عنوان ورودی استفاده می‌کند،
- با baseline کلاسیک TF-IDF شروع می‌کند،
- سپس با weak hybrid ranking و embedding + seed- filtering سمت supervision می‌رود،
- و با معیارهای Precision / Recall در چند K ارزیابی می‌شود.

نتیجه‌ی نهایی نشان داد روش پیشنهادی در سناریوهای مختلف بهخصوص از نظر Recall نسبت به baseline بهبود دارد و به هدف glossary extraction نزدیک تر است.

References

- . (Encyclopedia Wikipedia – Economics Britannica) برای تعریف و توضیح کلی حوزه
- برای Encyclopaedia Britannica – *The unintended effects of markets* محتواهای و غنی‌سازی متن اقتصاد (Wikipedia).
- Velardi et al., IEEE Intelligent Systems (2008) درباره اهمیت و معماری کلی Glossary Extraction و چرایی نیاز به اتوماسیون.

