



پروژه درس پردازش زبان های طبیعی:

## Sparse Coding متون با الگوریتم

دکتر ممتازی

ترم بهار ۱۳۹۶-۱۳۹۷

دانشکده مهندسی کامپیوتر، دانشگاه صنعتی امیرکبیر

زمان تحویل: ۲۵ تا ۲۷ تیر ۱۳۹۷

هدف از این پروژه استفاده از بازنمایی Term Frequency و Word2Vec برای پیاده سازی Summarization با کمک Sparse Coding می باشد.

فرض کنید که پیکره ای از تعدادی جمله داریم. در صورتی که هر جمله را به صورت یک بردار مانند Term-Frequency و یا میانگین Word2Vec کلمات آن جمله نشان دهیم، می توان پیکره را به صورت یک ماتریس که هر سطر آن بردار یک جمله است در نظر گرفت. به منظور به دست آوردن خلاصه ای از این جملات می توان با کمک Sparse Coding یک مجموعه از سطرهای ماتریس را انتخاب کرد که با کمک آنها بتوان جملات دیگر را تقریب زد. به صورت دقیق تر اگر مجموعه جملات به دست آمده توسط Sparse Coding را  $S^i = [s_1^i, s_2^i, \dots, s_k^i]$  بنامیم در آن صورت هر جمله  $s_i$  (سطر مربوط به جمله  $i$  ام در ماتریس اولیه) را می توان به صورت زیر تقریب زد.

$$s_i = \sum_{j=1}^k a_{ji} s_j^i$$

که در آنها مقادیر نامنفی هستند. در Sparse Coding به دنبال این هستیم که این تقریب کمترین خطا را داشته باشد.

برای جزئیات بیشتر می توانید به مقاله [1] مراجعه کنید.

در این پروژه قصد داریم از پیکره «پاسخ» استفاده کنیم. پیکره «پاسخ» (پیکره استاندارد سامانه های خلاصه ساز) در دو مدل تک سندی و چند سندی ارائه گردیده است. پیکره تک سندی شامل 100 موضوع مختلف از انواع گونه های خبری بوده که از خبرگزاری های پربیننده ایران انتخاب شده اند. هر کدام از این موضوعات دارای 5 خلاصه چکیده ای و استخراجی می باشند که توسط کارشناسان آموزش دیده تولید شده اند. پیکره چند سندی "پاسخ" نیز شامل 50 موضوع می باشد که هر موضوع حاوی 20 سند بوده و همچنین هر موضوع شامل 5 خلاصه انسانی و چکیده ای می باشد.

برای این پروژه از هر دو پیکره تک سندی و چند سندی استفاده کنید. یک بار خلاصه سازی بر روی پیکره تک سندی اجرا شود. در حالت تک سندی پیکره جملات ورودی را مجموعه جملات یک سند در نظر بگیرید. یک بار هم به صورت مجزا خلاصه سازی را بر روی پیکره چند سندی اجرا کنید. در حالت چند سندی پیکره جملات را مجموعه تمامی جملات در تمامی سندها در نظر بگیرید. برای کار با این پیکره

مرز جملات را نقطه، علامت سوال و علامت تعجب در نظر بگیرید. شایان ذکر است در برخی از خلاصه‌های گلد این پیکره مرز جملات هیچ یک از علامت‌های ذکر شده نیست و به صورت دستی مرز جملات مشخص شده اما برای پروژه نیاز نیست شما به علامتی غیر از علامت‌های گفته شده برای مرزبندی جملات توجه کنید.

مراحل انجام پروژه:

1. دیتاست پاسخ را از لینک زیر دانلود کنید:

<https://goo.gl/hwiULH>

2. برای داکيومنت‌های هر مجموعه، با استفاده از Term Frequency برداری تولید کنید.

3. الگوریتم Sparse Coding را بر روی بردارهای تولید شده اجرا نمایید.

4. خلاصه تولید شده، جمله‌های به دست آمده توسط Sparse Coding می‌باشد.

5. ارزیابی را با استفاده از معیار ROUGE1 و ROUGE2 بدست آورید و سپس بر اساس هر یک از دو معیار Precision، Recall و F-measure را محاسبه کنید.

6. با استفاده از روش Word2Vec برای واژگان موجود در دیتاست، یک بردار تولید کنید.

7. با توجه به بردارهای به دست آمده برای واژگان، بردار هر جمله را با میانگین‌گیری از مجموع بردارهای کلمات آن به دست آورید.

8. روش گفته شده برای Term Frequency را بر روی بردارهای جدید هم اجرا کرده و به طور مشابه ارزیابی نمایید.

منابع

[1] [Multi-Document Summarization Based on Two-Level Sparse Representation Model](#)