

Projeto: Sumarizador e Tradutor

Insper 2020
Arthur Rizzo e Vitor Liu

O projeto: problema abordado



Fontes: youmatter.world e pngwig



Fonte: medium.com



Fonte: endpoint

O projeto: aplicação e motivação

Motivação

- Avanços de técnicas de deep-learning em PLN.
- Ganho de eficiência na análise de documentos, livros, artigos, notícias, etc.

Principais aplicações

- Seleção de informação a ser traduzida/consumida.
- Categorização de documentos em língua estrangeira.

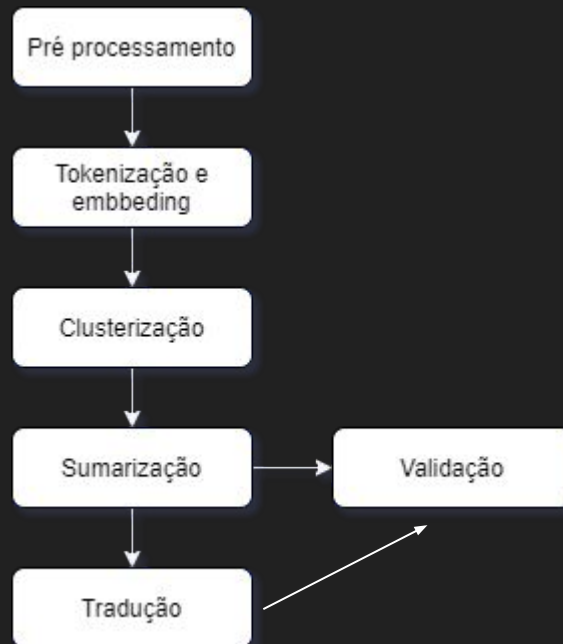
Objetivo do projeto

Realizar a sumarização de documentos e traduzir a sumarização. Com isso é possível:

- Gerar representações resumidas de documentos (melhorando o projeto anterior)
- Disponibilizar prévias de documentos em diferentes línguas
- Ganhar eficiência no mundo da tradução e sumarização

Metodologia

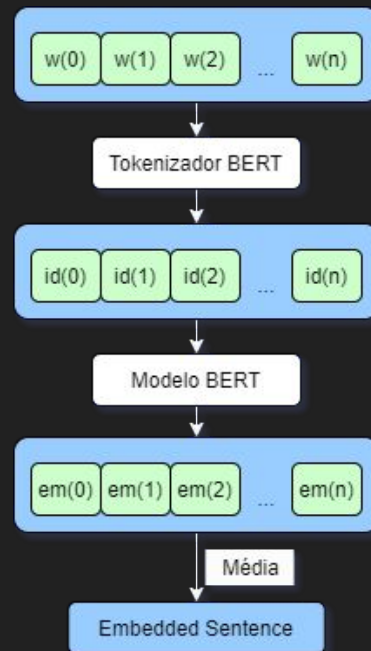
- 1 - Pré processamento
- 2 - Tokenização e embedding
- 3 - Clusterização
- 4 - Sumarização
- 5 - Resultados da sumarização
- 6 - Tradução



Fonte: O autor

Pré processamento, Tokenização e Embedding

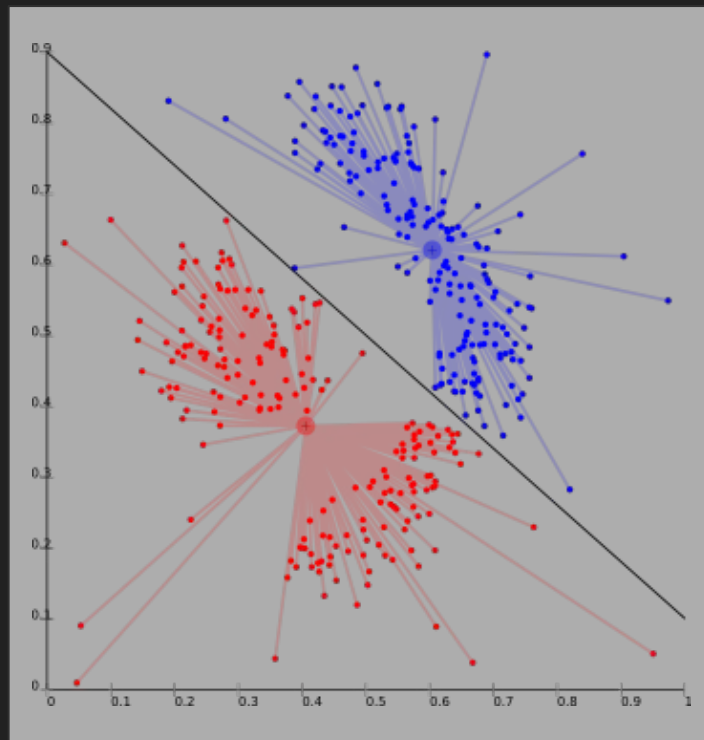
- Dataset usado: CNN stories tokenized
- Tokenizador: BERT base uncased (transformers)
- Modelo: BERT base uncased (transformers)



Fonte: O autor

Clusterização e Sumarização

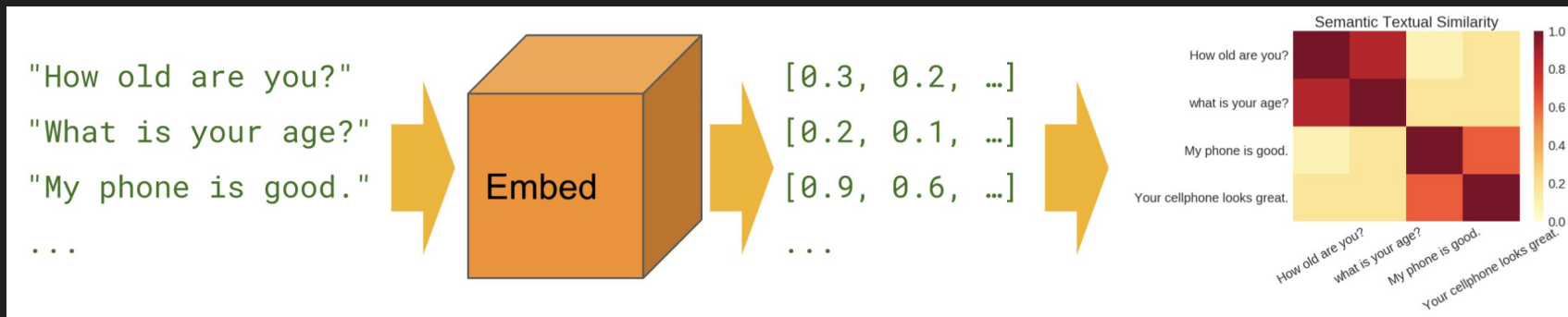
- Clusterização: MiniBatchKMeans (sklearn)
- Sumarização: Sentença mais perto do centro do cluster



Fonte: PngEgg

Validação da sumarização

- Medição da similaridade semântica entre o documento sumarizado e o “gabarito”
- Modelo: Universal Sentence Encoder (USE) - Boa performance



Fonte: TFHUB.DEV

Resultados da sumarização

Modelo usado: Universal Sentence Encode

Similaridade calculada por produto escalar 'numpy.inner()' (de 0 a 1)

Média de todos documentos:
0,3036

Média da maior correlação de
todos documentos: 0,5571

@highlight

Syrian official : Obama climbed to the top of the tree , `` does n't know how to get down ``

@highlight

Obama sends a letter to the heads of the House and Senate

@highlight

Obama to seek congressional approval on military action against Syria

@highlight

Aim is to determine whether CW were used , not by whom , says U.N. spokesman

Exemplo highlights CNN stories tokenized, o “gabarito”

Resultados da sumarização



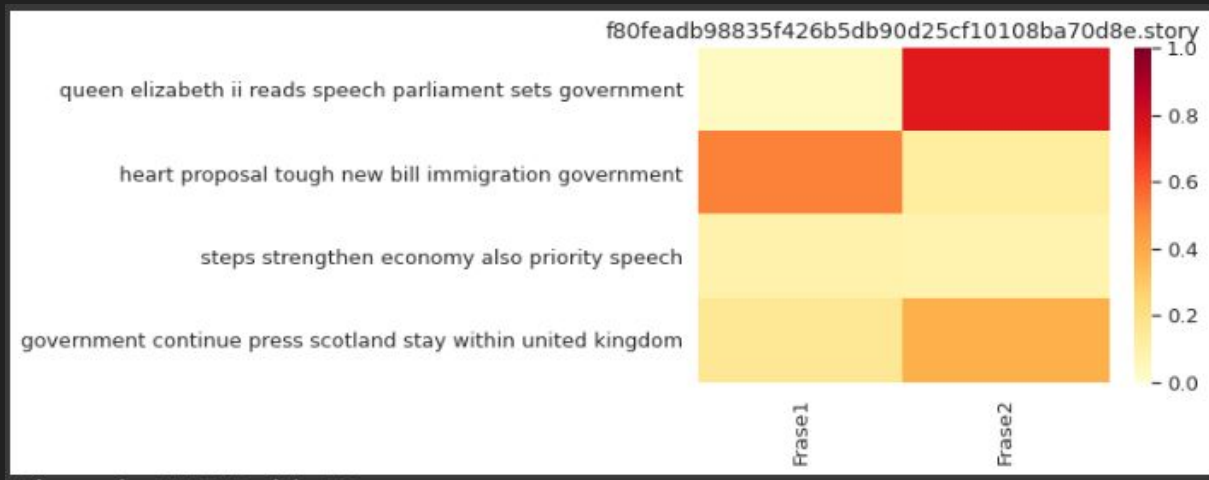
Exemplo caso médio

Resultados da sumarização

f80feadb98835f426b5db90d25cf10108ba70d8e.story

Frase 1: center government legislative agenda new bill tough new measures continue immigration reform prevent illegal immigrants accessing services entitled home office said

Frase 2: address queen elizabeth ii ceremonial state opening parliament written government although read monarch



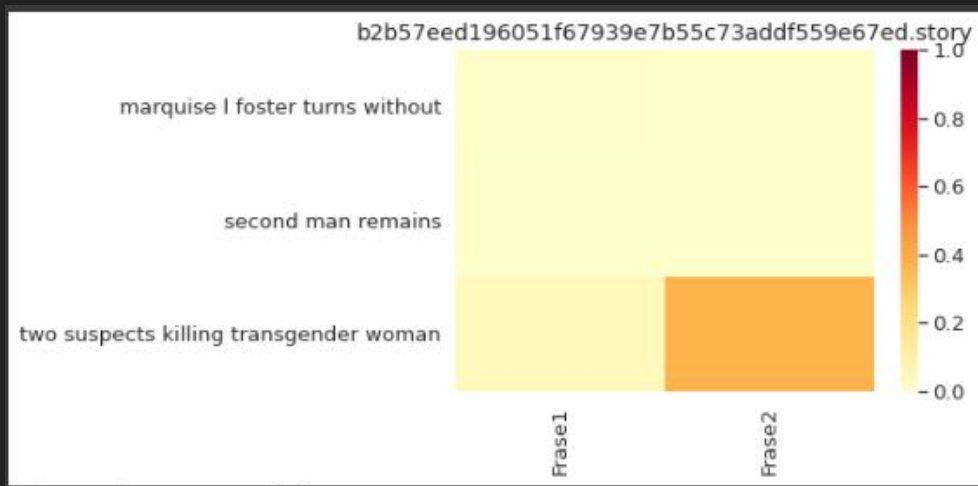
Exemplo caso bom

Resultados da sumarização

b2b57eed196051f67939e7b55c73addf559e67ed.story

Frase 1: special unit looks cases determine whether legal basis bring biascrime charges unit within office reviewing cases

Frase 2: two suspects killing victoria carmen white newark fatally shot september apartment maplewood police said statement maplewood halfhour drive west new york city



Exemplo caso ruim

Tradução - Modelo XLM

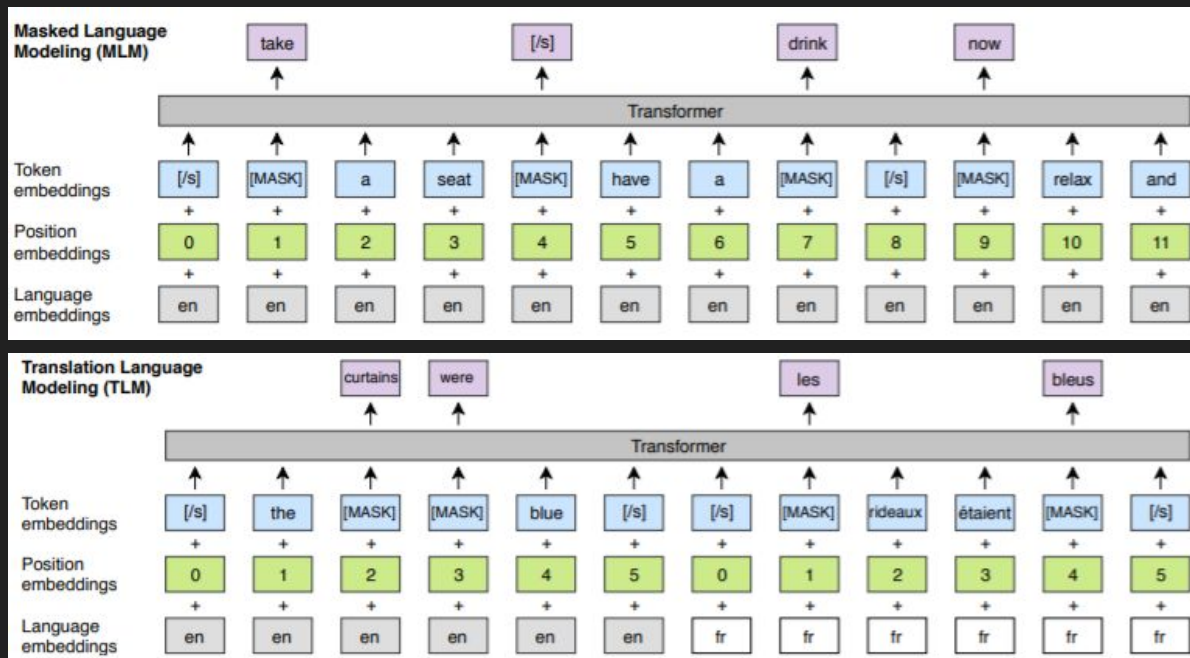
Implementação do modelo Cross-lingual Language Model Pre Training

Funcionalidades:

- Language model pre training:
 - Casual Language Model (CLM)
 - Masked Language Model (MLM)
 - Translation Language Model (TLM)
- GLUE e XNLI fine-tuning
- Supervised / Unsupervised MT training

Tradução - Modelo XLM

Treinamento MLM e TLM:



Fonte: facebookresearch/XLM

Tradução - Modelo XLM

Vantagens:

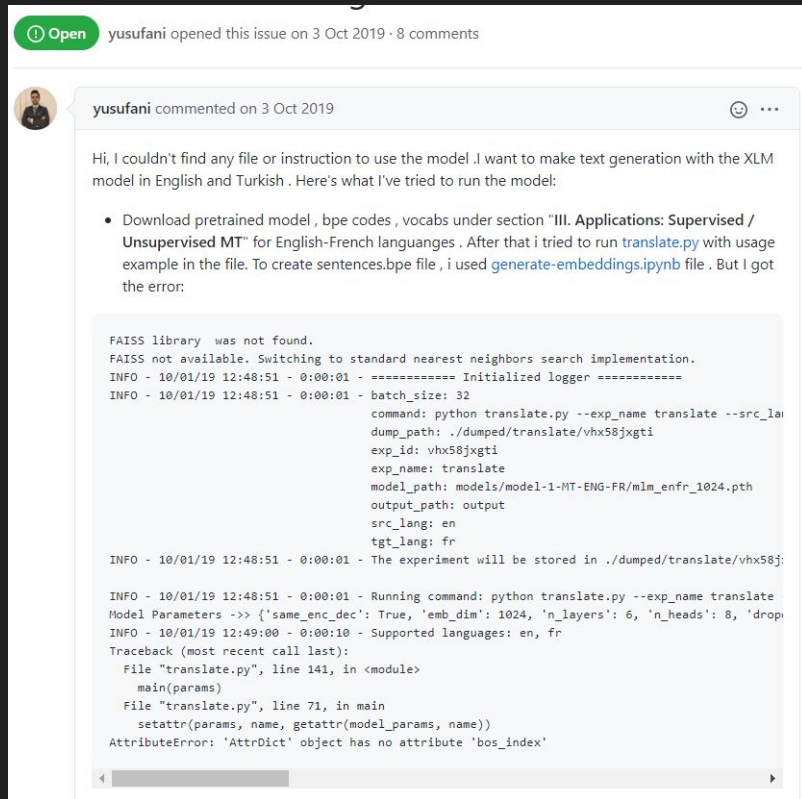
- Melhor desempenho na tradução
- Cross-lingual Language Model sem supervisão
- Possibilidade de fine-tuning

	en	fr	es	de	el	bg	ru	tr	ar	vi	th	zh	hi	sw	ur	Δ
<i>Machine translation baselines (TRANSLATE-TRAIN)</i>																
Devlin et al. (2018)	81.9	-	77.8	75.9	-	-	-	-	70.7	-	-	76.6	-	-	61.6	-
XLM (MLM+TLM)	<u>85.0</u>	<u>80.2</u>	<u>80.8</u>	<u>80.3</u>	<u>78.1</u>	<u>79.3</u>	<u>78.1</u>	<u>74.7</u>	<u>76.5</u>	<u>76.6</u>	<u>75.5</u>	<u>78.6</u>	<u>72.3</u>	<u>70.9</u>	63.2	<u>76.7</u>
<i>Machine translation baselines (TRANSLATE-TEST)</i>																
Devlin et al. (2018)	81.4	-	74.9	74.4	-	-	-	-	70.4	-	-	70.1	-	-	62.1	-
XLM (MLM+TLM)	<u>85.0</u>	79.0	79.5	78.1	77.8	77.6	75.5	73.7	73.7	70.8	70.4	73.6	69.0	64.7	65.1	74.2

Fonte: facebookresearch/XLM

Resultados da Tradução

Não foi possível realizar a tradução devido a um problema no código disponibilizado.



Fonte da imagem: XLM/issues/208

Considerações finais

Resultados resumidos

- Bons resultados na sumarização
- Falha na implementação da tradução

Para o futuro

- Implementar modelo de tradução com sucesso
- Melhorias na sumarização:
experimentalizar outras técnicas

Referências

<https://github.com/facebookresearch/XLM>

<https://arxiv.org/abs/1901.07291>

<https://tfhub.dev/google/universal-sentence-encoder/4> (USE)

<https://github.com/google-research/bert> (BERT)

<https://arxiv.org/abs/1810.04805> (BERT)

https://huggingface.co/transformers/model_doc/bert.html (BERT)

<https://github.com/facebookresearch/XLM/issues/208>

<https://www.pngegg.com/en/png-dsgwd>

Perguntas??