

Busca Semântica com Sentence-Transformers

Continuação dos Experimentos

Juliana Resplande

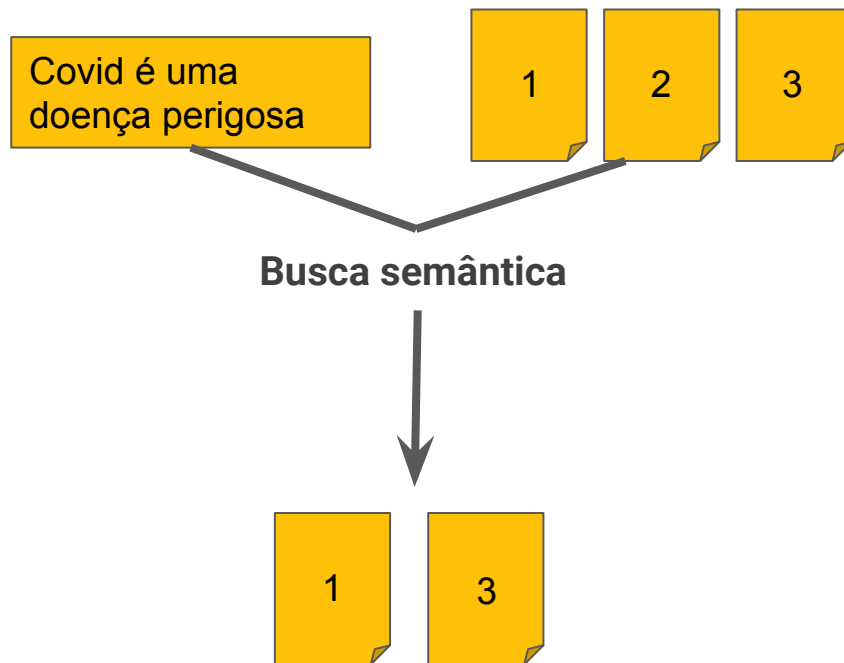
Agenda

1. Busca
2. Sentence Transformers
3. Novos experimentos
4. Conclusões

Busca semântica

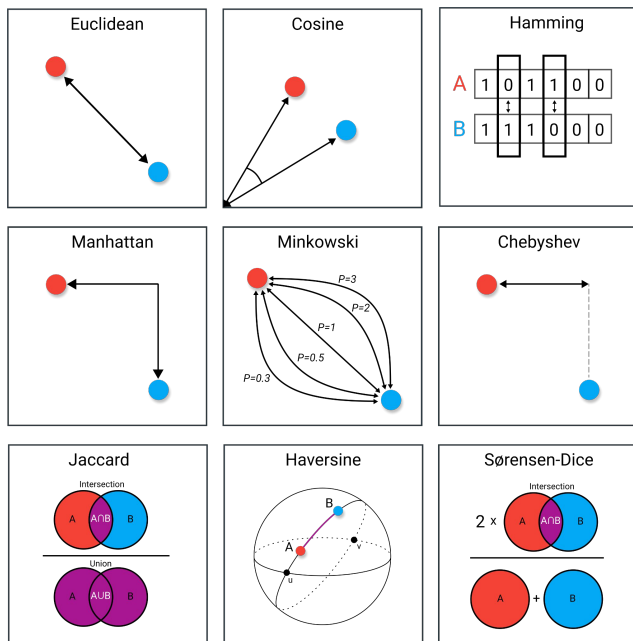
Busca semântica

- Information Retrieval (IR)
- Funcionamento:
 - a. Gera as **representações semânticas**
 - b. Compara o vetor da representação pela **distância**
- Dataset: MS-Marco



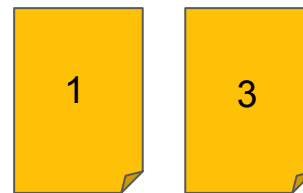
Considerações

- Como medir a **distância**?



- Considerar o **ranking** ou não?
 - Tarefa de similaridade
 - Tarefa de Ranqueamento

**“O documento 1 é mais importante
que o documento 3?”**



Introdução

1. Extrair embeddings de parágrafos
2. Extrair embeddings do termo de busca
3. Distância de cosseno para obtenção de dados mais semelhantes

Sentence Transformers

Cross-Encoder

X

Bi-encoder

- **BERT**
- **Entrada:**
 - Concatenação das sentenças
- Função de regressão depende do modelo
- Devagar

- **SBERT**
- **Entrada:**
 - Sentenças paralelamente
- Representações independentes
- Pode usar similaridade cosseno
- Rápido

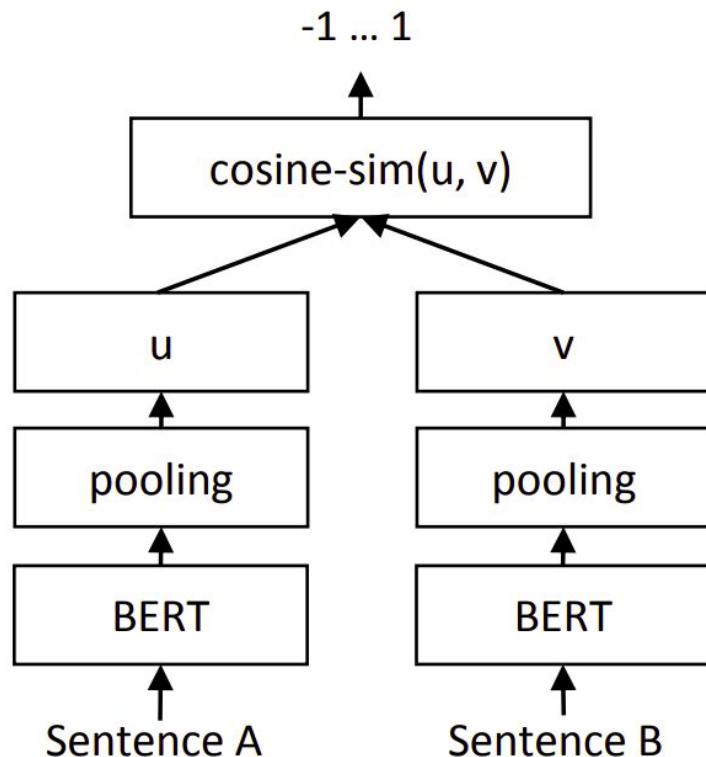
Sentence embeddings

Práticas comuns:

- Embeddings **individuais** usando BERT
- **Inicialização aleatória**
- Tamanho fixo por meio de média ou CLS

Sentence BERT

- BERT “**siamês**” para a tarefa
- Carrega o **BERT pré-treinado**
- Tamanho fixo por meio de **média**, máximo, CLS



Inferência Textual (NLI)

- **Classificação** de pares de sentenças
- Exemplos: QQP, QNLI, SNLI, MNLI, RTE
- **Treino em NLI** gera boas representações semânticas de texto

“Se consegue correlacionar bem quaisquer dois pares de sentenças, provavelmente correlacionará bem a query com o documento”

A soccer game with multiple males playing.

Some men are playing a sport.

entailment

SNLI

What programming languages do what?

What do all programming languages do?

duplicate

QQP

Sentence-Transformers

- Repositório oficial do Sentence-BERT
- Suporte aos transformers do HuggingFace
- Possui busca semântica implementada

Modos de treino iniciais:

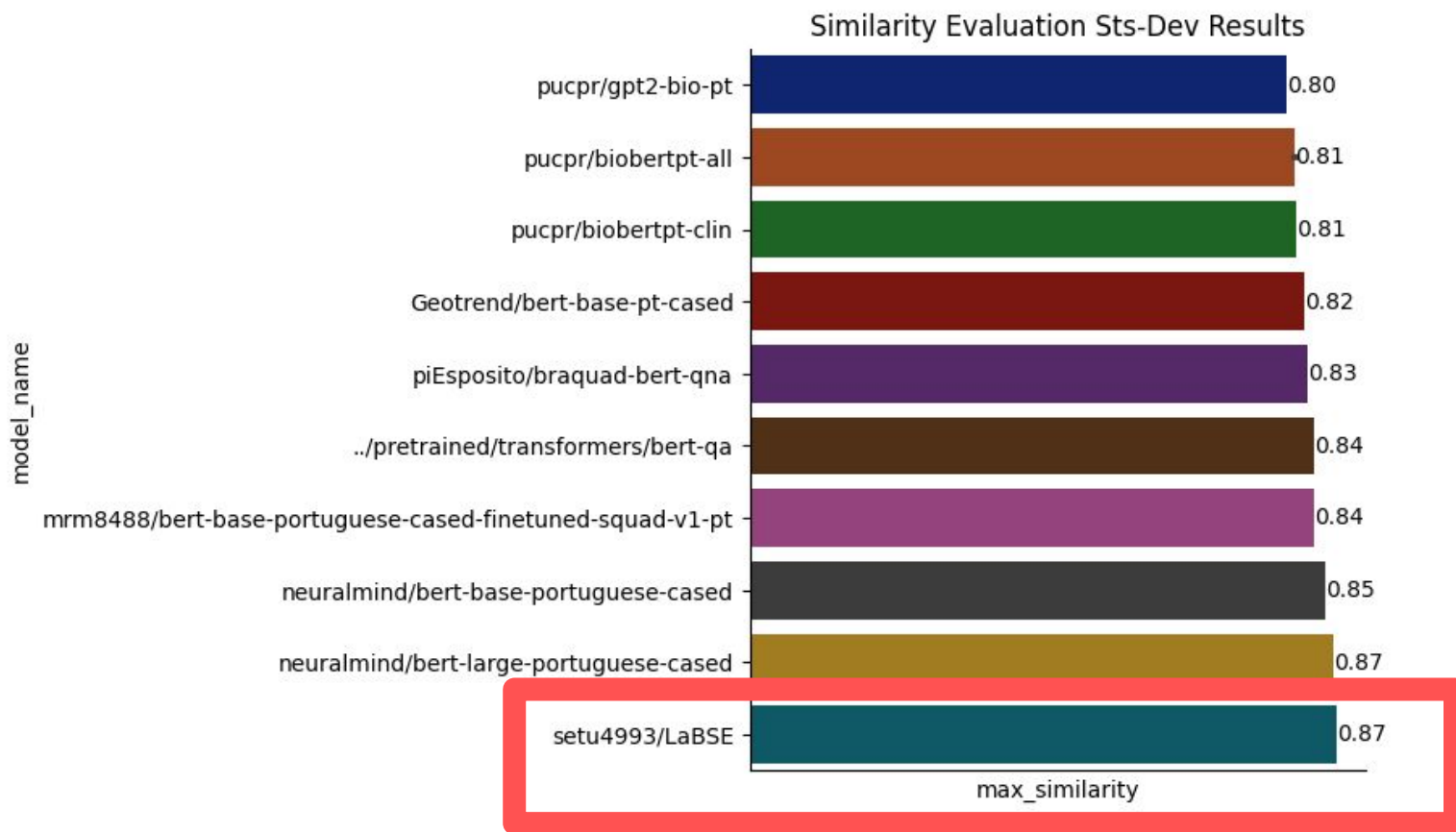
- Treinar no NLI
- Treinar no NLI e depois no STS-B
- Treinar em multi-tasking NLI e STS-B
- *Knowledge Distillation* em tradução

Novos experimentos

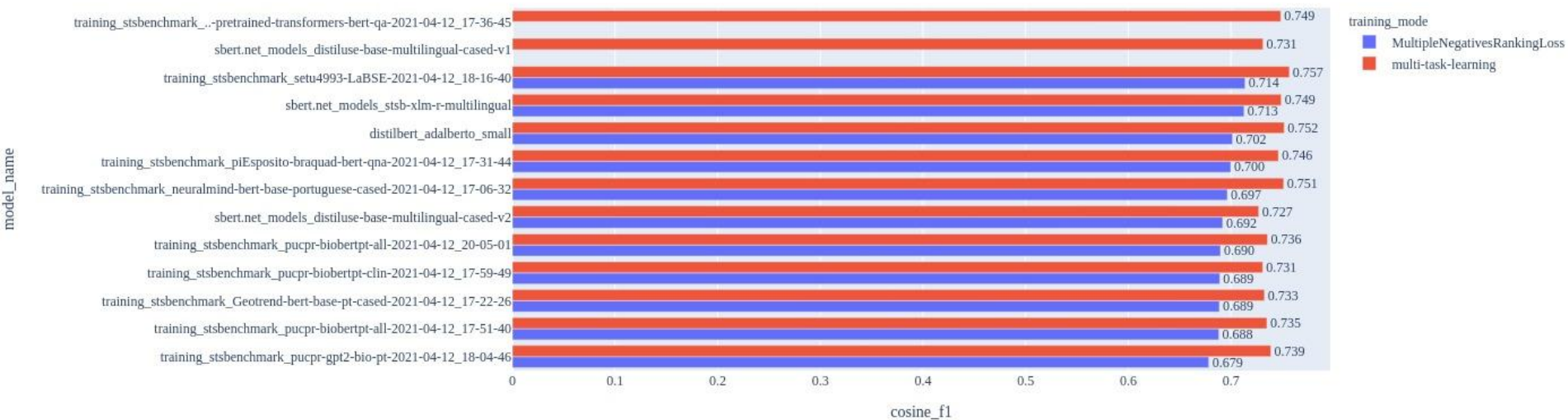
Treinos mais voltados para busca semântica

1. Treino pipeline STS + Quora Duplicate Questions
 - a. Treino normal STS
 - b. Treino no dataset do Quora modificado para busca semântica
 - i. Duplicate Questions Classification: Given two questions, are these questions duplicates?
 - ii. Duplicate Questions Mining: Given a large set (like 100k) of questions, identify all question pairs that are duplicates
 - iii. Duplicate Questions Information Retrieval: Given a large corpus (350k+) of questions. For a new, unseen question, find the most related (i.e. duplicate) questions in this corpus.
2. Treino Ms-Marco
3. Treino DPR
4. Treino não-supervisionado

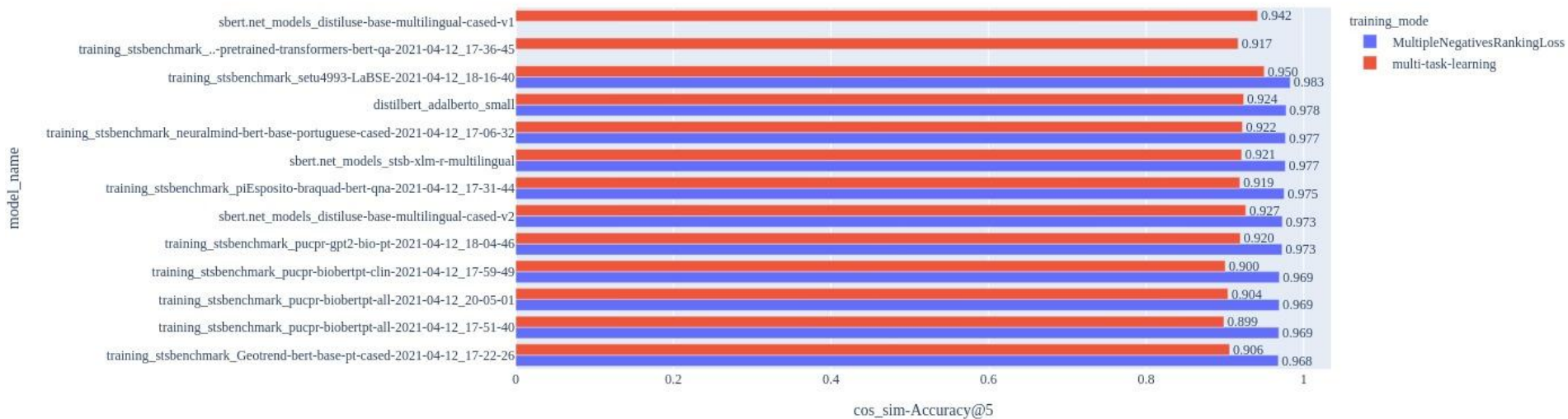
Resultados



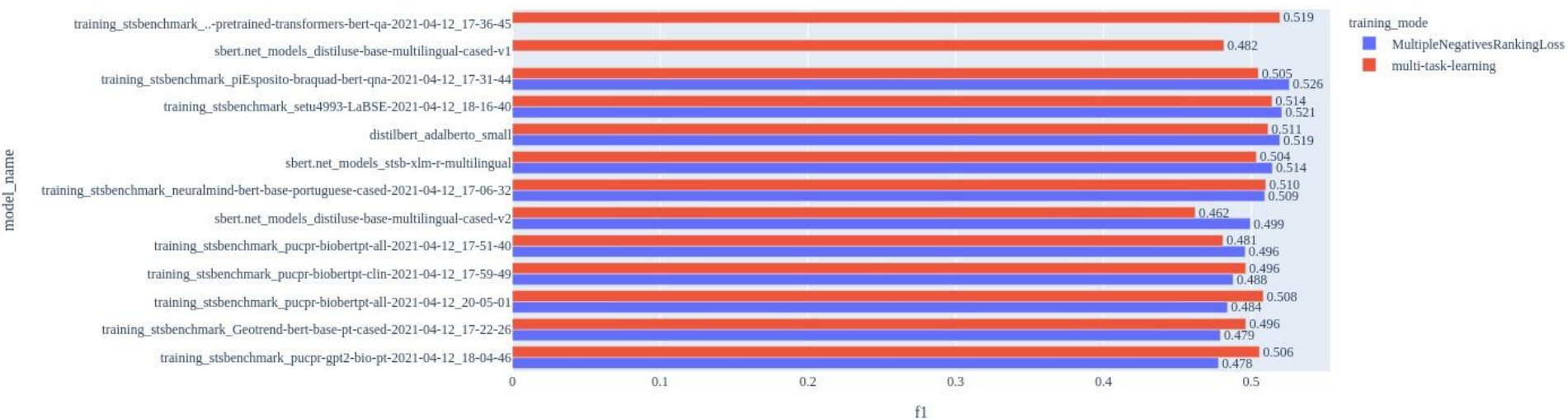
Binary Classification Evaluation Results



Information-Retrieval Evaluation Results



Paraphrase Mining Evaluation Dev Results



Conclusões

Conclusões/Trabalhos Futuros

- Importância de usar tarefas e métricas de ranqueamento
- LABSE: modelos multilíngues em detrimento de modelos monolíngues
- Finalizar o treino com o MS-Marco
- Treinar DPR
- Treino Não-Supervisionado
- Melhorar treino DeCLUTR