



Exame de qualificação de mestrado

Aprendizagem de representação através do uso de redes neurais convolucionais na recuperação de trecho de código-fonte

Marcelo de Rezende Martins
sob orientação do Prof. Dr. Marco Aurélio Gerosa

Instituto de Pesquisas Tecnológicas do Estado de São Paulo - IPT

Intro

Recuperação de trecho de código-fonte consiste em recuperar um trecho de código a partir de um repositório de códigos-fontes, de modo a atender a intenção do desenvolvedor, expressa em linguagem natural ¹².

¹Jose Cambroner, Hongyu Li, Seohyun Kim, Koushik Sen, and Satish Chandra. 2019. When deep learning met code search.

²Xiaodong Gu, Hongyu Zhang, and Sunghun Kim. 2018. Deep code search.

Code Retrieval: Dada uma questão em linguagem natural $q \in \mathbb{Q}$, um modelo F_r será treinado a recuperar os trechos $\mathbb{C}^+ \subset \mathbb{C}_a$ com a maior pontuação:

$$\mathbb{C}^+ = \underset{c \in \mathbb{C}_a}{\operatorname{argmax}} F_r(q, c) \quad (1)$$

Abordagem

Sejam \mathbb{Q} e \mathbb{C} conjuntos de dados heterogêneos. *Joint embedding* pode ser formulado como:

$$f: q \rightarrow t_q \rightarrow h_{\theta}(t_q, t_c) \leftarrow t_c \leftarrow c: g \quad (2)$$

Joint Embedding

Como representar as palavras e os tokens das questões e trechos de código-fonte?	<i>Word2Vec</i>
Como representar as sentenças?	CNN
Como aproximá-los?	Função de custo <i>hinge</i>

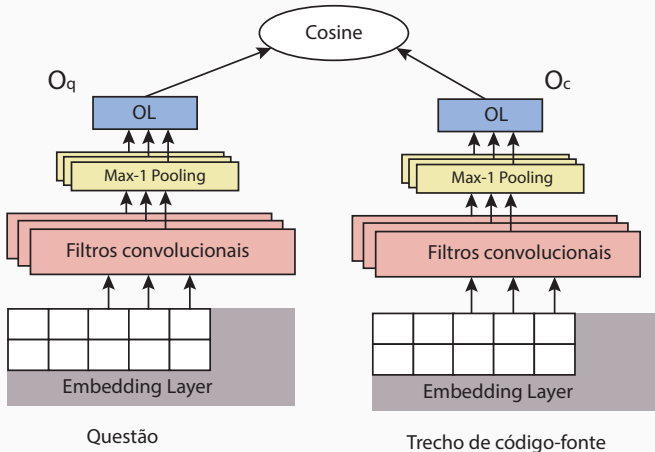


Figura 1: Arquitetura CNN proposta para recuperação de trecho de código-fonte.

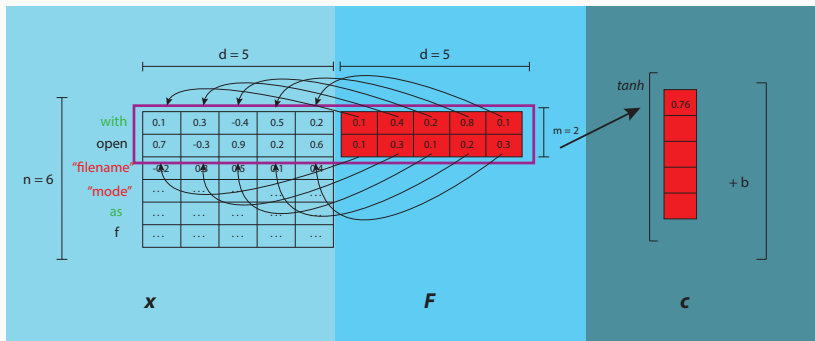


Figura 2: Primeiro passo da operação de convolução em um vetor de entrada x composto por vetores de representação distribuída de cada palavra da sentença.

Questões

- A aprendizagem de representação através do CNN auxilia na recuperação de trecho de código-fonte?
- O CNN é capaz de extrair as características mais relevantes de modo a facilitar o modelo a encontrar uma correlação entre as questões e os trechos de código-fonte?

Indiretamente:

- As interações locais auxiliam na aproximação das intenções aos trechos de código?

Dados de treinamento	Conjunto de pares de questões e trechos de código-fonte em Python coletados do Stack Overflow por Yao et al. (2018) ³
Dados para avaliação final	Conjunto de dados anotados manualmente e disponibilizados por Yao et al. (2018) ³
Métrica de desempenho	<i>MRR</i>

³Yao, Ziyu and Weld, Daniel S. and Chen, Wei-Peng and Sun, Huan. 2018. StaQC: A Systematically Mined Question-Code Dataset from Stack Overflow

Arquiteturas de referência para comparação	<ul style="list-style-type: none">• Embedding• Rede neural com mecanismo de atenção proposto por Cambronero et al. (2019)⁴
--	--

⁴Jose Cambronero, Hongyu Li, Seohyun Kim, Koushik Sen, and Satish Chandra. 2019. When deep learning met code search.

Análise dos resultados	<ul style="list-style-type: none">• Inspeção manual• Análise dos piores casos• Patologia das redes neurais (Feng et al. , 2018)⁵
------------------------	---

⁵Shi Feng, Eric Wallace, Alvin Grissom II, Mohit Iyyer, Pedro Rodriguez, Jordan Boyd-Graber. 2018. Pathologies of Neural Models Make Interpretations Difficult.

Experimento piloto

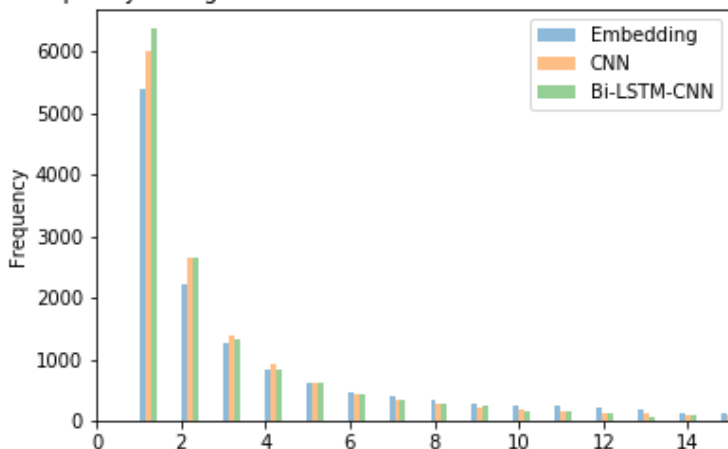
Modelos	Resultados (MRR)
Embedding	0,52 \pm 0,01
CNN	0,58 \pm 0,01
bi-LSTM-CNN	0,60 \pm 0,02

Tabela 1: Resultado preliminar do modelo CNN em comparação com outras duas arquiteturas (bi-LSTM com CNN e Embedding). Estes resultados foram obtidos a partir da amostra EVAL.

⁶Marcelo de Rezende Martins e Marco Aurélio Gerosa. 2019. Um estudo preliminar sobre o uso de uma arquitetura deep learning para seleção de respostas no problema de recuperação de trecho de código-fonte.

Histograma das posições dos trechos de código-fonte relevantes

Frequency Histogram of Rank Position of First Relevant Code Snippet



Python and appending items to text and excel file⁷

BiLSTM-CNN

```
Yvalues = [1, 2, 3, 4, 5]
file_out = open('file.csv', 'wb')
mywriter = csv.writer(file_out, delimiter = '\n')
mywriter.writerow(Yvalues)
file_out.close()
```

CNN

```
import csv

with open("output.csv", "wb") as f:
    writer = csv.writer(f)
    writer.writerow(a)
```

⁷<https://stackoverflow.com/questions/24593478/python-and-appending-items-to-text-and-excel-file>

Próximos passos

- Implementação da arquitetura proposta por Cambroner et al. (2019)⁸
- Adição de regularização aos modelos
- Coleta e análise dos resultados

⁸Jose Cambroner, Hongyu Li, Seohyun Kim, Koushik Sen, and Satish Chandra. 2019. When deep learning met code search.

Perguntas?