

КУРСОВАЯ РАБОТА
Исследовательский проект
“Сжатие словарей для нейросетевого анализа
исходных кодов программ”

Выполнил: Андрей Гусев
Руководитель КР: Чиркова Надежда Александровна

Высшая Школа Экономики

aagusev_2@edu.hse.ru

8 июня 2020 г.

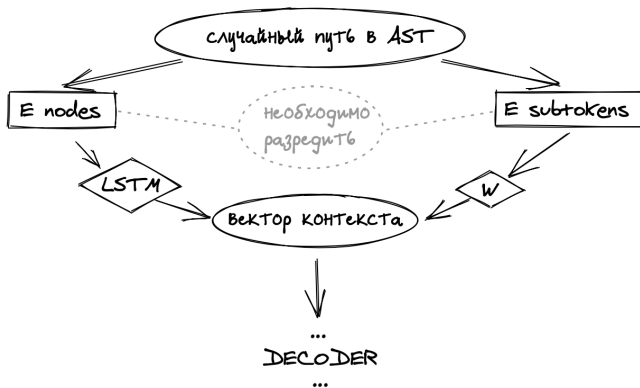
Прунинг входного слоя

Матрица эмбедингов E задает представление L элементов словаря. Сумму рассматриваемых регуляризаторов $R(E)$ можно представить в следующем виде:

$$R(E) = \underbrace{\lambda_1 \sum_{l=1}^L \|E_l\|_2}_{\text{group Lasso}} + \underbrace{\lambda_2 \|E\|_1}_{\text{Lasso}}, \text{ где } \|\cdot\|_2 \text{ — это } \ell^2\text{-норма}$$

Lasso-регуляризация стимулирует обнуление одного веса в строке, а group Lasso — к последующему обнулению всей строки.

Encoder Code2Seq



Сравнение с базовой моделью

	Val F1	Test F1	nodes	subtokens
разреженная¹	0.4239	0.4195	177	10029
ограниченная²	0.4201	0.4277	177	10029
оригинальная³	0.4092	0.4229	323	73906

¹Средние результаты трех запусков с предложенной техникой разреживания

²Результаты запуска модели с простой техникой сжатия

³Усредненные результаты исходной модели