

Оценка ассоциативной связи слов

Word2Vec, PMI score, Collocations

Вальба О.В, Лайко Р.С

Значения forward/backward показывают ассоциативную силу (free association response probability) между исходными словами согласно работе Nelson et al¹

source	target	forward	backward
ability	able	0.077	0.014
tenor	voice	0.053	0.0
proprietor	sale	0.020	0.0
useless	worthless	0.111	0.088
...

Таблица 1: Пример исходных данных (source: Nelson et al.¹)

¹Behavior Research Methods, Instruments, Computers 2004, 36 (3), 402–40, The University of South Florida free association, rhyme, and word fragment norms

Оценка ассоциативной близости

1. Получить численное представление пар слов
2. Оценить их "близость"
3. Вычислить корреляцию с исходной ассоциативностью:
 - Коэфф. Пирсона $\rho_{X,Y} = \frac{E[(X-\mu_X)(Y-\mu_Y)]}{\sigma_X \sigma_Y}$
 - Коэфф. Кендалла $\tau = \frac{1}{n(n-1)} \sum_{i \neq j} \text{sgn}(x_i - x_j) \text{sgn}(y_i - y_j)$

Оценка ассоциативной близости

- Косинусная близость распределенных векторов слов

$$\cos(\theta) = \frac{\langle w_1, w_2 \rangle}{\|w_1\| \|w_2\|}$$

- GP-метрика²³

$$GP(w_1, w_2) = \begin{cases} \frac{P(w_1, w_2)}{\|w_2\|}, & P(w_1, w_2) < \|w_2\| \\ 1 + \frac{P(w_1, w_2)}{\|w_2\|}, & P(w_1, w_2) \geq \|w_2\| \end{cases}$$

- Коэффициент PMI (Pointwise Mutual Information)

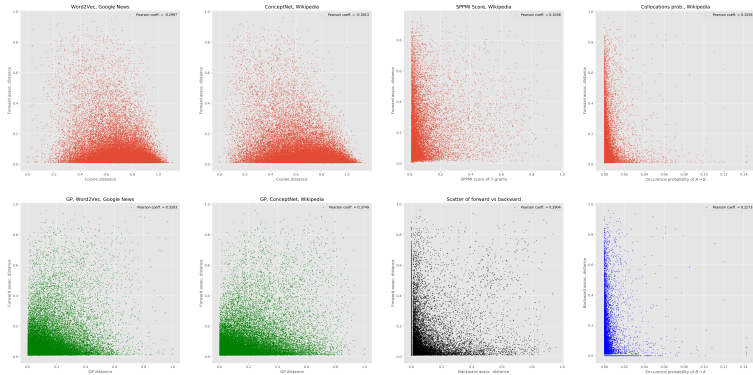
$$pmi(w_1; w_2) = \log \frac{p(w_2|w_1)}{p(w_2)}$$

- Вероятность коллокации n-грамов

²David Galea, Peter Bruza, Deriving Word Association Networks from Text Corpora

³ $P(w_1, w_2) = \frac{\langle w_1, w_2 \rangle}{\|w_2\|}$

Визуализация результатов



Модель	Корр. Пирсона	Корпус
Word2Vec	0.29	Google News
ConceptNet	0.36	Wikipedia + Wiktionary
SPPMI	0.16	Wikipedia (latest dump)
GloVe	0.21	Twitter
fastText	0.28	Wikipedia + IMDB
n-grams occ. _{forward}	0.15	Wikipedia (latest dump)
n-grams occ. _{backward}	0.23	Wikipedia (latest dump)

Пороговые значения

- Выберем порог ассоциативной близости
- Возьмем пары (source \rightarrow target) у которых оно больше этого порога
- Посчитаем корреляцию (Пирсона, Кендалла)

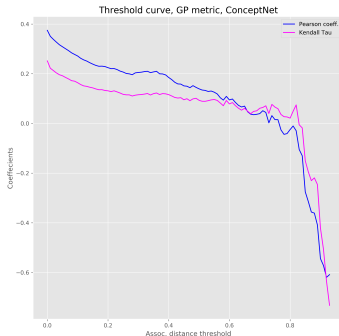


Рис. 1: Зависимость величины корреляции от порогового значения ассоциативной близости в исходном датасете

- Работа с аналогичными датасетами на русском языке
- Распределение Ципфа
- Совмещение сетевых и распределенных подходов