

Московский Физико-Технический Институт
(Государственный Университет)

Факультет Управления и Прикладной Математики
Кафедра «Интеллектуальные Системы»

ДИПЛОМНАЯ РАБОТА СТУДЕНТА 374 ГРУППЫ

«Поиск оптимальной модели детектирования нечетких дубликатов коротких текстов»

Выполнил:

студент 4 курса 374 группы

Сафин Камиль Фанисович

Научный руководитель:

к.ф.-м.н.

Чехович Юрий Викторович

Содержание

1 Введение	3
2 Основная часть	3

Аннотация

Аннотация.

1 Введение

Актуальность темы. В последнее время активно развиваются технологии, позволяющие создавать модели языков для различных задач. Самые частые задачи это: машинный перевод, тематическое моделирование, поиск заимствований, семантический анализ и многие другие.

2 Основная часть

Скрытая переменная $z \in Z$. Ее распределение: $z \sim P(z)$. Набор детерминированных функций, параметризованных вектором $\theta \in \Theta$: $f(z, \theta) : Z \times \Theta \rightarrow \mathcal{X}$. Если z — случайная переменная, то и $f(z, \theta)$ — случайная переменная в пространстве \mathcal{X} .

Правдоподобие данных:

$$P(X) = \int P(X|z, \theta)P(z)dz.$$

Здесь $f(z, \theta)$ заменено на распределение $P(X|z, \theta)$. Причем берется $P(X|z, \theta) = \mathcal{N}(X|f(z, \theta), \sigma^2 \cdot I)$. Введем $Q(z|X)$ — функцию, которая по данным X выдает распределение над переменными z , из которых вероятнее сгенерируется X .

Дивергенция Кульбака-Лейблера: $\mathcal{D}[Q(z)||P(z|X)] = E_{z \sim Q}[\log Q(z) - \log P(z|X)]$. Из правила Байеса можно получить: $\log P(z|X) = \frac{P(X|z) \cdot P(z)}{P(X)} = \log P(X|z) + \log P(z) - \log P(X)$. Тогда для дивергенции получаем:

$$\mathcal{D}[Q(z)||P(z|X)] = E_{z \sim Q}[\log Q(z) - \log P(z|X)] = E_{z \sim Q}[\log Q(z) - \log P(X|z) - \log P(z)] + \log P(X).$$

И далее:

$$\log P(X) - \mathcal{D}[Q(z)||P(z|X)] = E_{z \sim Q}[\log P(X|z) - \mathcal{D}[Q(z)||P(z)]].$$

В целом, Q — любое распределение над z , но чтобы минимизировать $\mathcal{D}[Q(z)||P(z|X)]$, логично, чтобы оно зависело от X : $Q(z) = Q(z|X)$. В итоге получаем:

$$\log P(X) - \mathcal{D}[Q(z|X)||P(z|X)] = E_{z \sim Q}[\log P(X|z)] - \mathcal{D}[Q(z|X)||P(z)]$$

. ELBO:

$$\log P(X) \geq E_{z \sim Q}[\log P(X|z)] - \mathcal{D}[Q(z|X)||P(z)] = \mathcal{L}.$$

Variational attention:

$$\mathcal{L} = E_{z, a \sim Q(z, a|X)}[\log P(X|z, a)] - \mathcal{D}[Q(z, a|X)||P(z, a)]$$

Т.к. z и a независимы друг от друга, то их совместное распределение факторизуется:

$$\mathcal{L} = E_{z \sim Q^{(z)}(z|X), a \sim Q^{(a)}(a|X)}[\log P(X|z, a)] - \mathcal{D}[Q^{(z)}(z|X)||P(z)] - \mathcal{D}[Q^{(a)}(a|X)||P(a)].$$

И финальная функция ошибки:

$$\mathcal{I} = \mathcal{I}_{rec}(X) + \lambda_{KL} \left[\mathcal{D}[Q^{(z)}(z|X)||P(z)] + \gamma_a \sum_{j=1}^{|X|} \mathcal{D}[Q^{(a)}(a_j|X)||P(a_j)] \right]$$