

Controllable text generation with small data using auxiliary in-domain enrichment

Беляев Станислав

Санкт-Петербургский Академический Университет

stasbelyaev96@gmail.com

23 марта 2018

Введение

Обзор

“If a typical person can do a mental task with less than one second of thought, we can probably automate it using AI either now or in the near future.”

— Andrew Ng, 2017

Машины умеют:

- 1 Различать формы и объекты
- 2 Имитировать стиль изображений
- 3 Отвечать на простые вопросы

Машины НЕ умеют:

- 1 Хорошо подражать высшей нервной деятельности
- 2 Понимать и обобщать сложные категории
 - Этика (юмор, мораль, норма, ...)
 - Эстетика (книги, картины, ...)

Введение

Постановка задачи

МООС платформам нужна генерация контента:

- Дешево
- Быстро
- Ультимативная защита от списывания

Особенности:

- Generic характер генерации
- Примеров готового контента мало
- Набор текстовых свойств для единицы контента (курс, тема, тэги, сложность, ...)



Задача: По набору свойств $f = \{f_i \in F\}$ сгенерировать новые примеры текстовых данных из генеральной совокупности X , соответствующих f . Возьмем в качестве X условия задач по программированию.

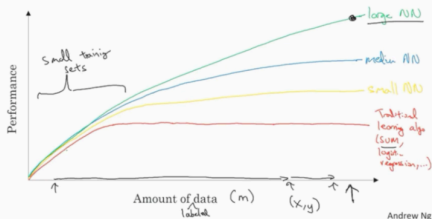
Введение

Данные в DL

“Data is the New Oil.”

— Andrew Ng, 2017

Scale drives deep learning progress



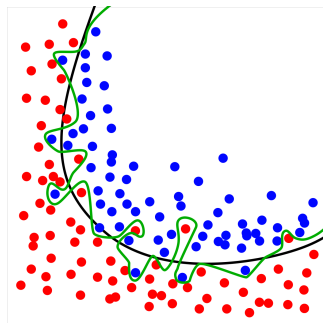
Введение

Проблема данных

Если мы не знаем паттерна для генерации и хотим уметь обобщать, то будем использовать DL и больше данных (Mikolov et al., 2010).

Но что делать, если данных мало?

- Мы не сможем обобщать
- Мы скорее всего переобучимся
- При генерации новые сэмплы будут слишком похожи на старые



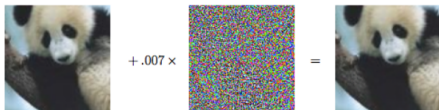
Решение: Искать похожие $X_{aux} \sim X$ in-domain данные из смежных областей.

Введение

Изображения vs текст

Изображения

$$f: \mathbb{R}^2 \rightarrow \mathbb{R}^M$$



- Непрерывное пространство
- Набор всевозможных преобразований как дифференцируемых функций
- Понятно, куда распространять градиент

Текст

...an efficient method for learning high quality distributed vector ...

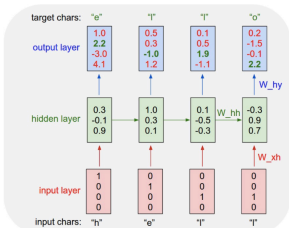
context focus word context

- Дискретное пространство
- Переменная длина
- Нет устойчивости к шуму
- Long-term зависимости
- Омонимия и контекст

Введение

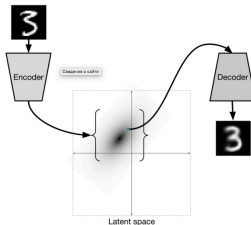
Генерация текста

RNN



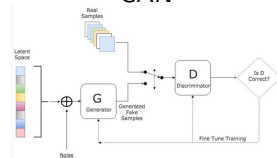
(Mikolov et al., 2011)

VAE



(Bowman et al., 2016; Hu et al., 2018)

GAN



(Yu et al., 2017; Fedus et al., 2018)

Данные

Источники

Разберемся, откуда брать данные...

RNN

Обзор

RNN

VAE

Обзор

VAE

GAN

Обзор

GAN

Оценивание

Метрики

Perplexity и смотреть глазками

Treatments	Response 1	Response 2
Treatment 1	0.0003262	0.562
Treatment 2	0.0015681	0.910
Treatment 3	0.0009271	0.296

Таблица: Table caption

Выводы

Результаты

To to и to to

Ссылки

Статьи

То то и то то



[John Smith \(2012\)](#)

Title of the publication

Journal Name 12(3), 45 – 678.