Лекция 3. **Языковые модели**

Языковые модели

- Что такое языковые модели?
- Общая идея работы
- Классический подход к работе с языковыми моделями
- Нейросетевой подход к работе с языковыми моделями

Что такое языковая модель?

Рассмотрим, например, модели поездов.

• имеют некоторые свойства поездов

(выглядят как поезда)

- могут вести себя подобно поездам
- хорошие модели больше похожи на настоящие

поезда





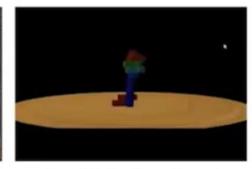
Что такое языковая модель?

Рассмотрим, например, физические модели.

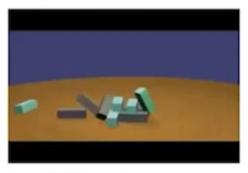
- помогут понять, какие события лучше согласуются с миром, какие более вероятны
- может предсказать, что произойдет, учитывая некоторый «контекст»



Will it fall?



In which direction?



Different masses



Complex scenes



Infer the mass



Predict fluids

Что такое языковая модель?

Отличается понятие события: для языка событие — это языковая единица (текст, предложение, токен, символ).

Языковые модели (LM) оценивают вероятность различных языковых единиц: символов, токенов, последовательностей токенов.

Примеры языковых моделей

Мы имеем дело с языковыми моделями каждый день.

Например:

• поисковая система



• переводчик, электронная почта



Примеры языковых моделей

Мы имеем дело с языковыми моделями каждый день.

Например:

• поиск опечаток

```
I saw a catt

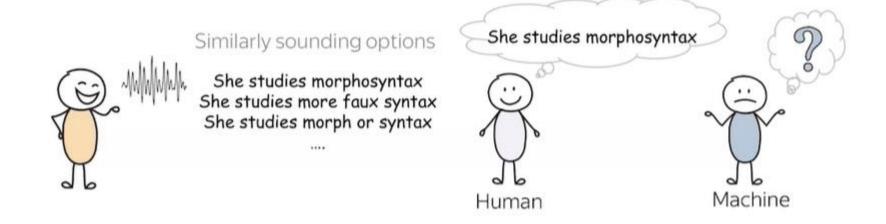
Probably you meant I saw a cat
```

```
I saw a catt

cat

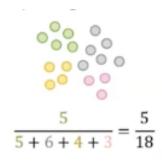
car
```

Сложности машинного восприятия

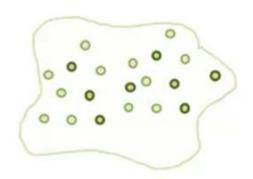


Какова вероятность появления предложения в языке?

Какова вероятность вытащить зелёный шар?



Можно ли рассуждать аналогично для предложений?



P(mut the tinming tebn is the)=
$$\frac{0}{|corpus|} = 0$$

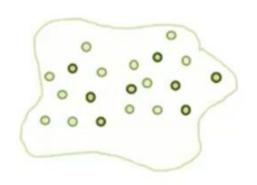
P(the mut is tinming the tebn)= $\frac{0}{|corpus|} = 0$

Text corpus

При таком подходе предложения, которые никогда не встречались в корпусе, получат нулевую вероятность.

Какова вероятность появления предложения в языке?

Можно ли рассуждать аналогично для предложений?



P(the mut is tinming the tebn)=
$$\frac{0}{|corpus|} = 0$$

P(mut the tinming tebn is the)=
$$\frac{0}{|corpus|} = 0$$

При таком подходе предложения, которые никогда не встречались в корпусе, получат нулевую вероятность.

Мы не сможем надежно оценить вероятности предложений, если будем рассматривать их как атомарные единицы!

Идея: разбивать предложения на меньшие части

- читаем предложение слово за словом
- обновляем вероятность каждый раз, когда встречается новый токен

```
P(I
                                                                                a) =
             P(I
                      saw) =
                                                                                                         P(I
                                                                                                                           a) =
                                                                                                                  saw
                                                               P(I) \cdot P(saw|I) \cdot
             P(I) \cdot P(saw|I)
                                                                                                          P(I) \cdot P(saw|I) \cdot P(a|I|saw)
                                                       Probability of I saw
 Probability of I Probability of saw given I
                                                                                                  Probability of I saw Probability of a given I saw
                                                                                                         P(I
                                                                                                                               cat on) =
                                                P(I
                                                                        cat) =
P(I
                       cat) =
                                                                                                         P(I) \cdot P(saw|I) \cdot P(a|I|saw) \cdot P(cat|I|saw|a).
                                                P(I) \cdot P(saw|I) \cdot P(a|I|saw) \cdot P(cat|I|saw|a)
P(I) \cdot P(saw|I) \cdot P(a|I|saw)
                                                                                                                  Probability of I saw a cat
                                                 Probability of I saw a Probability of cat given I saw a
 Probability of I saw a
```

Идея: разбивать предложения на меньшие части

- $(y_1, y_2, ..., y_n)$ последовательность токенов
- $P(y_1,y_2,...,y_n)$ вероятность увидеть эти токены в данной последовательности

$$P(y_1, y_2, \dots, y_n) = P(y_1) \cdot P(y_2|y_1) \cdot P(y_3|y_1, y_2) \cdot \dots \cdot P(y_n|y_1, \dots, y_{n-1}) = \prod_{t=1}^{n} P(y_t|y_{< t})$$

Структура моделирования языка с написанием слева направо

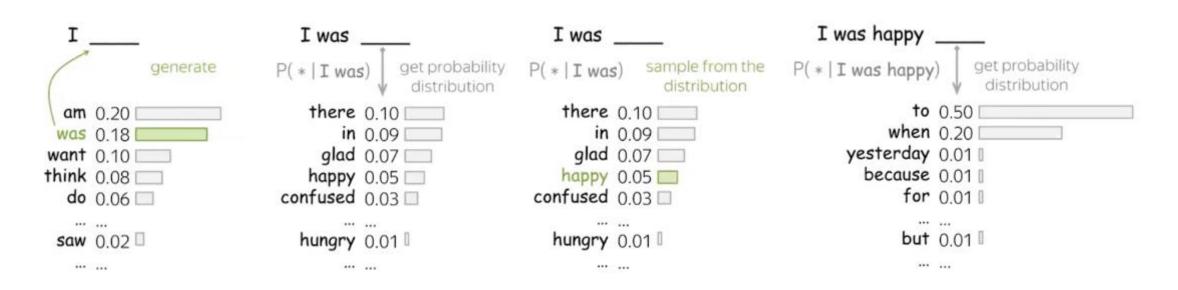
Модели различаются по способу вычисления вероятностей.

Цель: понять, как вычислять условную вероятность $P(y_t|y_1,y_2,...,y_{t-1})$

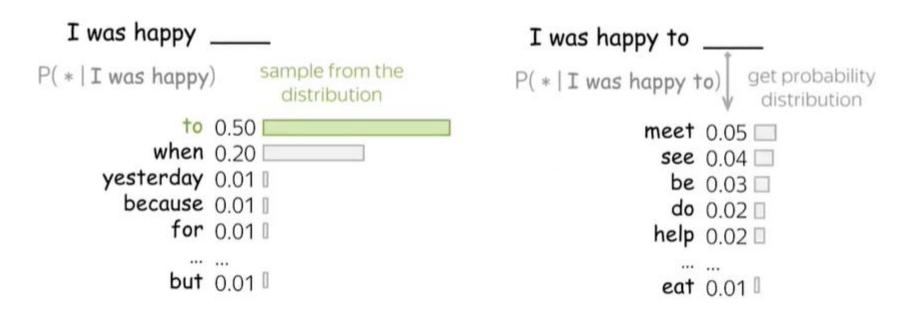
1 вариант: использование n-грамм (вычисления «вручную»)

2 вариант: нейросети

Генерация текста с помощью языковой модели



Генерация текста с помощью языковой модели



Традиционный подход к работе с языковыми моделями

Цель: необходимо вычислить условную вероятность в формуле

$$P(y_1, y_2, \dots, y_n) = P(y_1) \cdot P(y_2|y_1) \cdot P(y_3|y_1, y_2) \cdot \dots \cdot P(y_n|y_1, \dots, y_{n-1}) = \prod_{t=1}^{n} P(y_t|y_{< t})$$

Оценка на основе глобальной статистики текстового корпуса, т. е. подсчета.

Как вычислить условные вероятности?

Простой способ:

$$P(y_t | y_1, y_2, ..., y_{t-1}) = \frac{N(y_1, y_2, ..., y_t)}{N(y_1, y_2, ..., y_{t-1})},$$

где $N(y_1,y_2,\dots,y_t)$ — количество раз, которое мы встречаем (y_1,y_2,\dots,y_t) в обучающем корпусе

Марковское предположение

Вероятность слова зависит только от фиксированного числа предыдущих слов.

$$P(y_t \mid y_1, y_2, ..., y_{t-1}) = P(y_t \mid y_{t-n+1}, ..., y_{t-1})$$
all previous tokens
$$n-1 \text{ previous token}$$

Например:

- n=3 (trigram model): $P(y_t | y_1, y_2, ..., y_{t-1}) = P(y_t | y_{t-2}, y_{t-1})$
- n=2 (bigram model): $P(y_t | y_1, y_2, ..., y_{t-1}) = P(y_t | y_{t-1})$
- n=1 (unigram model): $P(y_t | y_1, y_2, ..., y_{t-1}) = P(y_t)$

Марковское предположение. Пример

P(I saw a cat on a mat) = ?Before After (3-gram) P(I saw a cat on a mat) = (I saw a cat on a mat) = P(I)P(I) P(saw | I) $\cdot P(saw \mid I) \longrightarrow \cdot P(saw \mid I)$ $\cdot P(a | I saw)$ $\cdot P(a | I saw) \longrightarrow \cdot P(a | I saw)$ $\cdot P(cat | I saw a)$ $\cdot P(cat \mid \exists saw a) \rightarrow \cdot P(cat \mid saw a)$ $\cdot P(on | T saw a cat) \longrightarrow \cdot P(on | a cat)$ · P(on | I saw a cat) · P(a | I saw a cat on) $P(a | \mathbf{I} \mathbf{saw} \mathbf{a} \mathbf{cat} \mathbf{on}) \longrightarrow P(a | \mathbf{cat} \mathbf{on})$ · P(mat | I saw a cat on a) · P(mat | I saw a cat on a) · P(mat | on a)

Возможная проблема

P(mat | I saw a cat on a) = P(mat | cat on a) =
$$\frac{N(cat on a mat)}{N(cat on a)}$$

$$P(\text{mat } | \text{ cat on a}) = \frac{N(\text{cat on a mat})}{N(\text{cat on a})} = ?$$

not good: can not compute the probability

P(mat | cat on a) =
$$\frac{N(cat \text{ on a mat})}{N(cat \text{ on a})} = ?$$

not good: zeros out probability of the sentence

Решение проблемы (Back-off)

Использовать меньше контекста в ситуациях, о которых мы мало что знаем:

- если можно, использовать триграмму
- если нет, то биграмму
- если нет, то униграмму

$$P(\text{mat } | \text{ cat on a}) = \frac{N(\text{cat on a mat})}{N(\text{cat on a})} = ?$$

$$\frac{\text{zero}}{\text{not good: can not compute the probability}}$$

N(cat on a) = 0
$$\longrightarrow$$
 try "on a"
P(mat | cat on a) \approx P(mat | on a)
N(on a) = 0 \longrightarrow try "a"
P(mat | on a) \approx P(mat | a)
N(a) = 0 \longrightarrow try unigram
P(mat | a) \approx P(mat)

Лучшее решение: интерполяция

```
P(\text{mat } | \text{ cat on a}) = \frac{N(\text{cat on a mat})}{N(\text{cat on a})} = ?
\frac{\text{zero}}{\text{not good}} : \text{can not compute the probability}
```

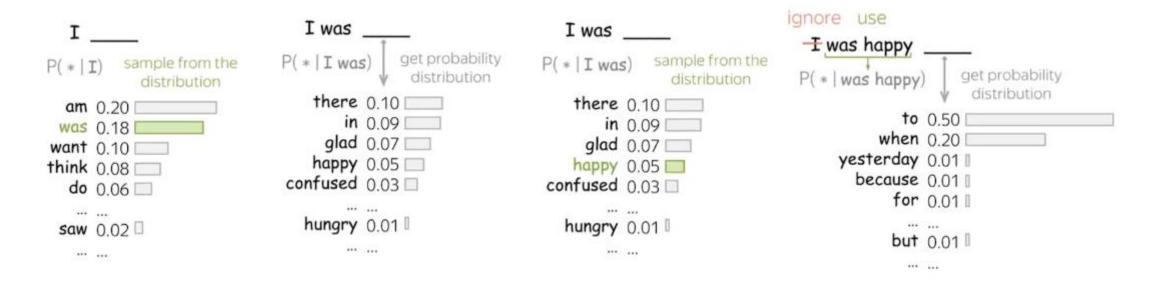
$$\widehat{\mathbf{P}}(\text{mat } \mid \text{cat on a}) \approx \lambda_3 \mathbf{P}(\text{mat } \mid \text{cat on a}) + \\ \lambda_2 \mathbf{P}(\text{mat } \mid \text{on a}) + \\ \lambda_1 \mathbf{P}(\text{mat } \mid \text{a}) + \\ \lambda_0 \mathbf{P}(\text{mat}) \\ \sum_i \lambda_i = 1$$

Сглаживание Лапласа

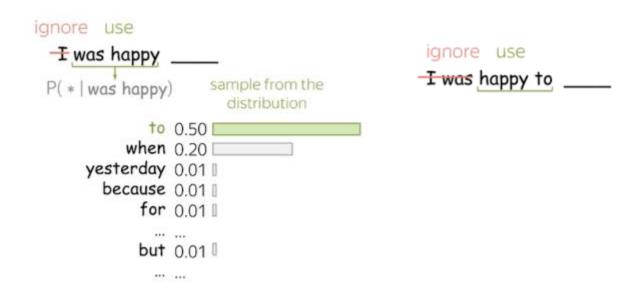
$$P(\text{mat } | \text{ cat on a}) = \frac{\frac{\text{zero}}{\text{N(cat on a mat)}}}{\text{N(cat on a)}} = ?$$
not good: zeros out probability of the sentence

$$\widehat{P}(\text{mat } | \text{ cat on a}) = \frac{\delta + N(\text{cat on a mat})}{\delta \cdot |V| + N(\text{cat on a})}$$

Генерация текста с помощью n-граммной языковой модели (триграммная модель)



Генерация текста с помощью n-граммной языковой модели (триграммная модель)



Пример сгенерированного текста

when this option may be the worst day of amnesty international delegations visited israel , and felt that his sisters , that they are reserved for zyryanovsk concentrating factory there is a member of the shire ," given as to damage the expansion of a meeting over a large health maintenance organization , smoking airconditioning , designated smoking area . eos so even when i talk a bit short , there was no easy thing to do different buffer flushing strategies in the future , due to huge list of number - one just has started production of frits in the process and has free wi - fi " operation eos

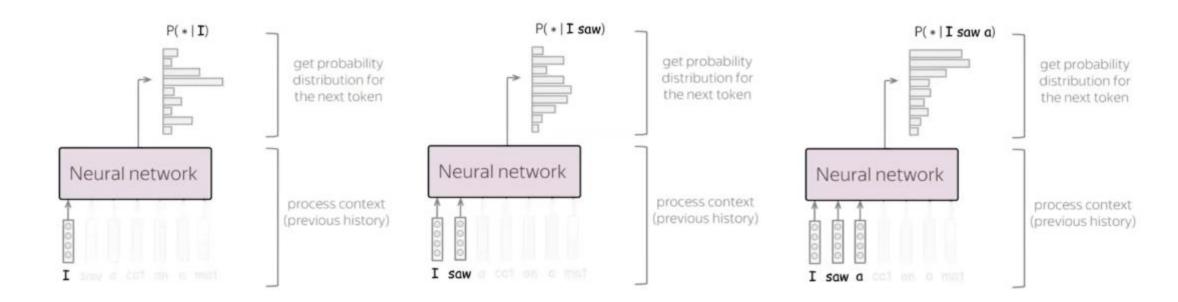
Нейросетевой подход

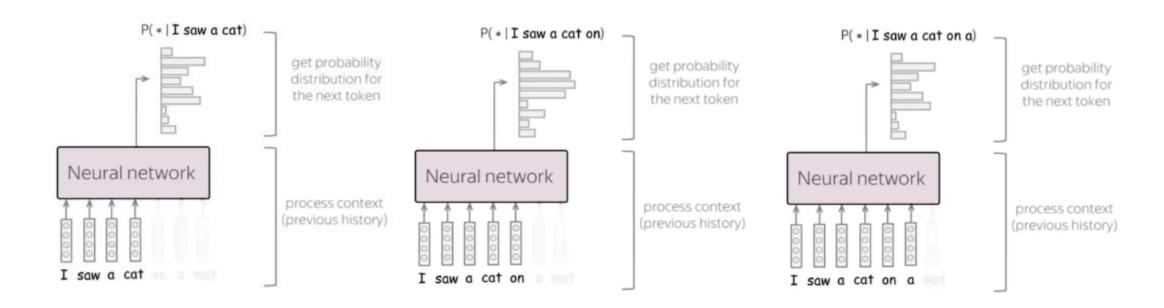
Цель: необходимо вычислить условную вероятность в формуле:

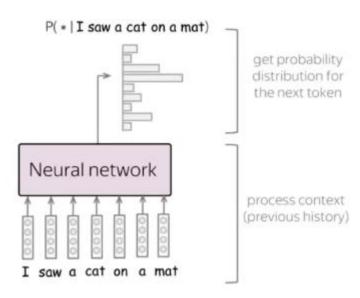
$$P(y_1, y_2, \dots, y_n) = P(y_1) \cdot P(y_2|y_1) \cdot P(y_3|y_1, y_2) \cdot \dots \cdot P(y_n|y_1, \dots, y_{n-1}) = \prod_{t=1}^n P(y_t|y_{< t})$$

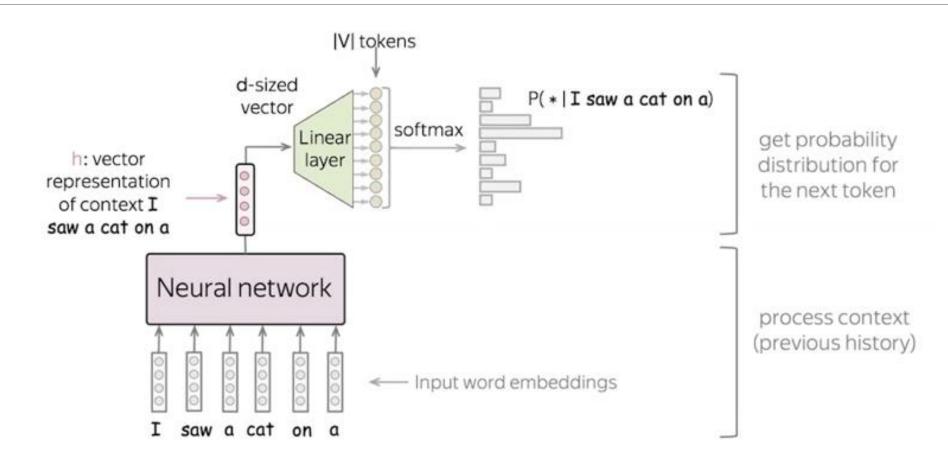
В нейронных сетях нужно обучить нейронную сеть прогнозировать вероятности.

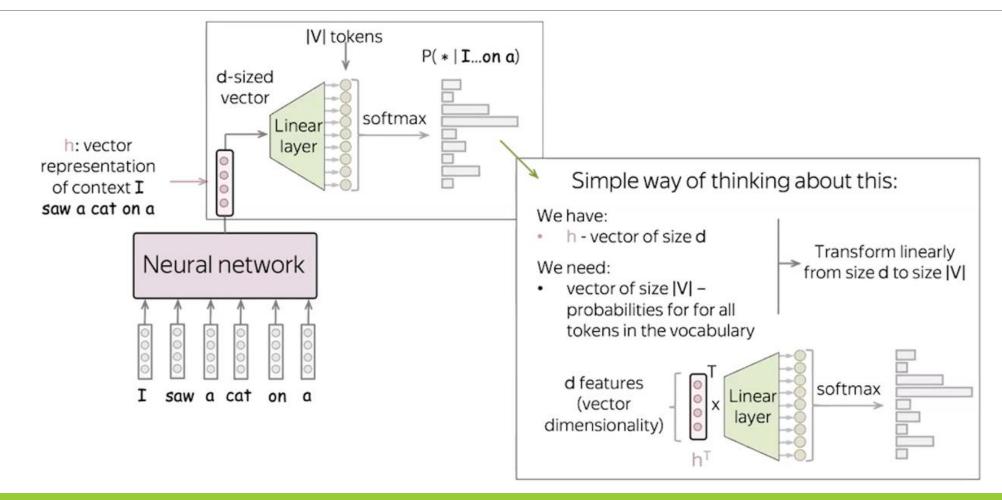
- получить векторное представление контекста
- оценить вероятности (прогноз распределения вероятностей для следующего токена)

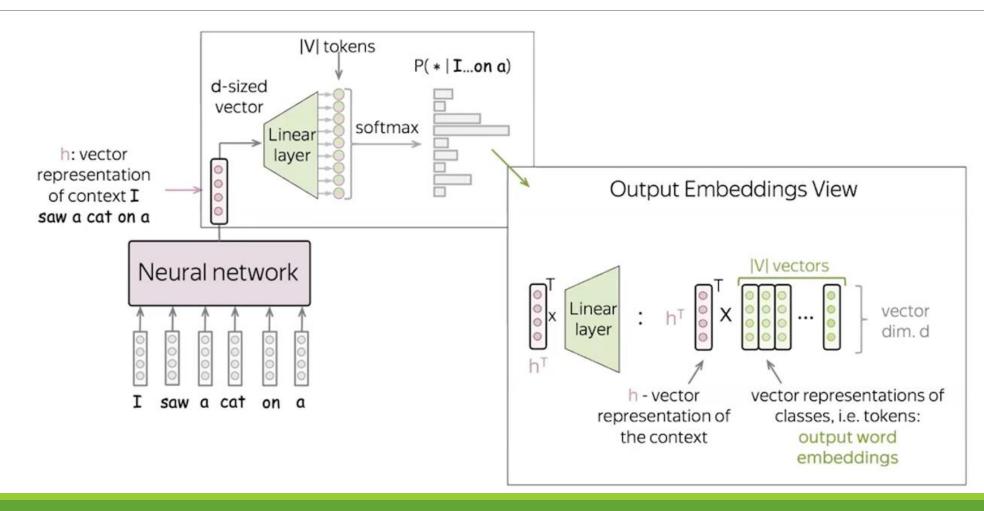




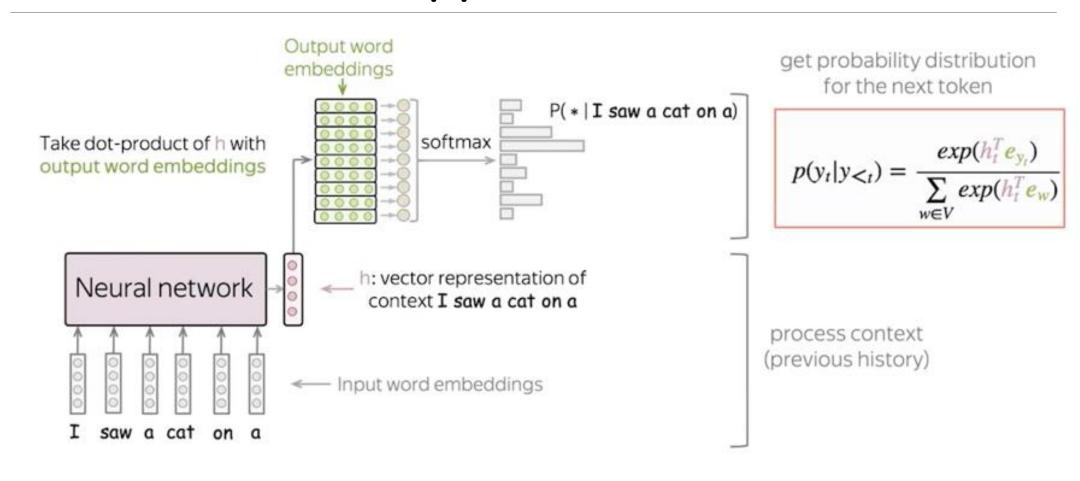








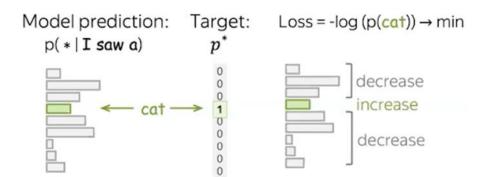
Представление вложений выходных данных



Обучение: кросс-энтропия



Training example: I saw a cat on a mat <eos>



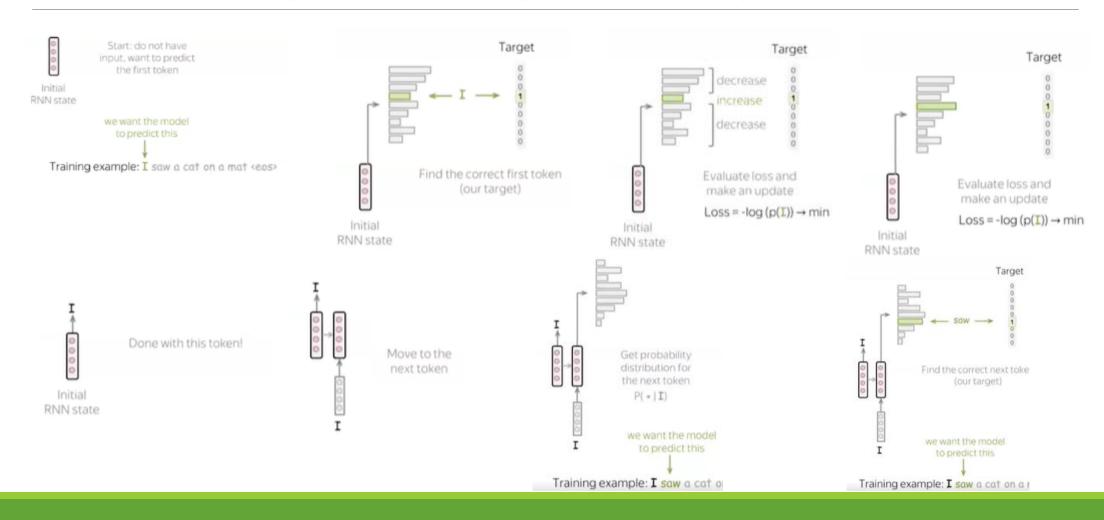
Cross-entropy loss:

$$-\sum_{i=1}^{|V|} p_i^* \cdot \log P(y_t = i|x) \to \min (p_k^* = 1, p_i^* = 0, i \neq k)$$

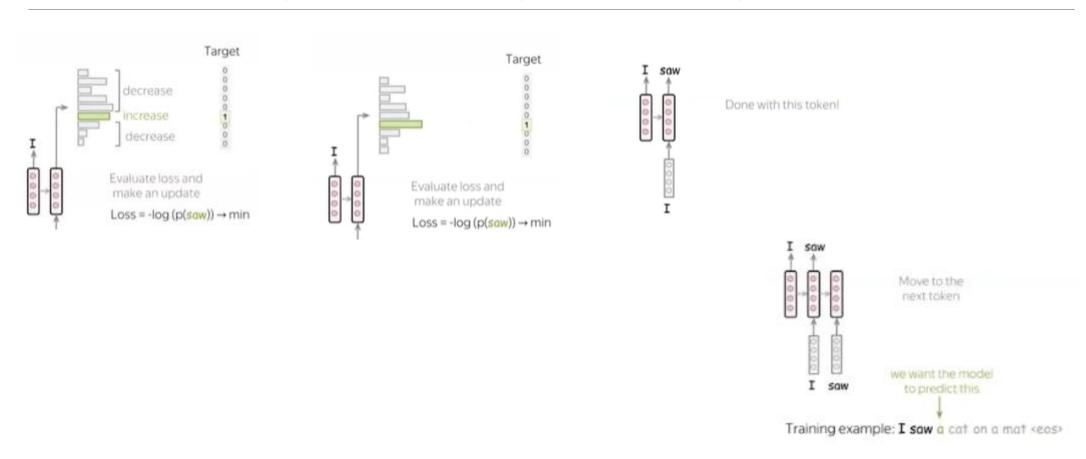
For one-hot targets, this is equivalent to

$$-\log P(y_t = cat|x) \rightarrow min$$

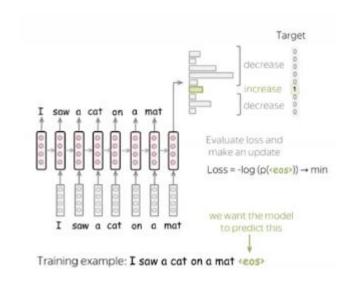
Обучение: кросс-энтропия

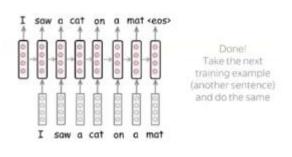


Обучение: кросс-энтропия



Обучение: кросс-энтропия

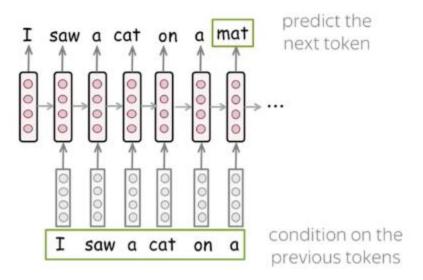




Training example: I saw a cat on a mat <eos>

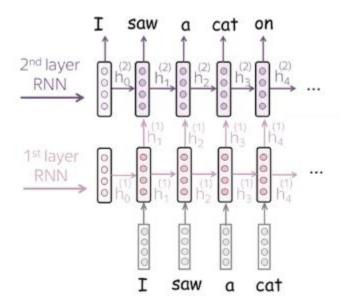
Рекуррентные модели

Читаем текст, на каждом шаге предсказываем следующий токен



Рекуррентные модели

Многослойность: передача состояний из одной RNN в другую

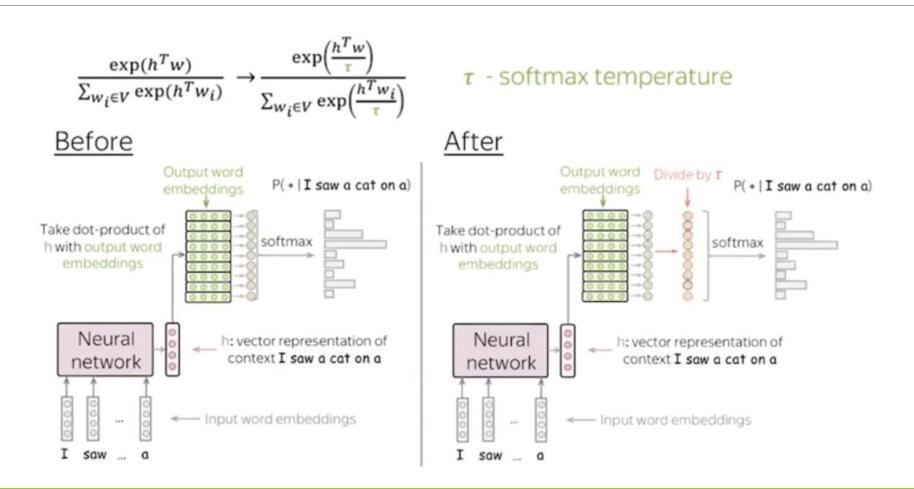


Связность и разнообразие

Сгенерированные тексты должны быть:

- связными и осмысленными— сгенерированный текст должен иметь смысл;
- разнообразными модель должна быть способна создавать очень разные образцы.

Использование softmax temperature



Пример: temperature = 2

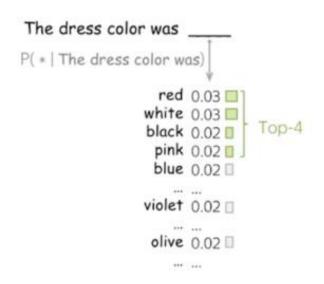
paradise sits farms started paint hollow almost unprecedented decisions, care using withdrawal from rebel cis (, saying graphics mongolia official line, greeted agenda victor is exploring anger :) draw testify liberalization decay productive 2 went exchanges of marketing drawing enabling challenging systematic crisis influencing the executive arrangement performs designs

Пример: temperature = 0,2

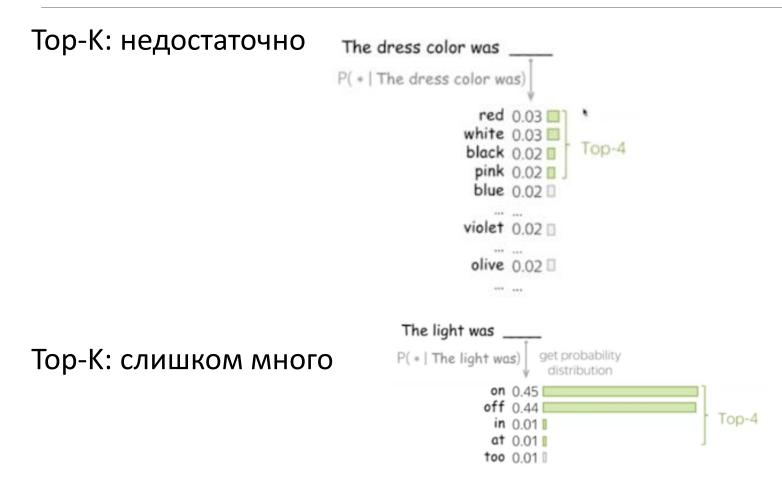
the first time the two - year - old - old girl with a new version of the new

Связность и разнообразие

• выбирать из первых К наиболее вероятных токенов (обычно К составляет около 10-50)

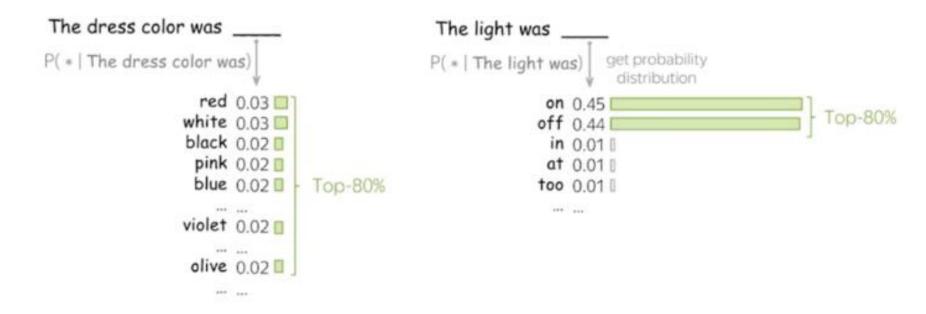


Как выбрать К?



Top-p sampling

• на каждом шаге выбирать столько верхних токенов, чтобы покрыть р% вероятностной массы



Оценка моделей. Перплексия

Насколько модель «удивляется» при чтении нового текста?

$$L(y_{1:M}) = L(y_1, y_2, \dots, y_M) = \sum_{t=1}^{M} \log_2 p(y_t | y_{< t}) \qquad Loss(y_{1:M}) = -\sum_{t=1}^{M} \log p(y_t | y_{< t})$$

$$Loss(y_{1:M}) = -\sum_{t=1}^{M} \log p(y_t|y_{< t})$$

Log-likelihood of the text

Note: cross-entropy (our loss) is negative log-likelihood

$$Perplexity(y_{1:M}) = 2^{-\frac{1}{M}L(y_{1:M})}$$

Какие значения может принимать перплексия?

- лучшее значение это 1 (если модель идеальна и присваивает вероятность 1 правильным токенам, то логарифмические вероятности равны нулю)
- худшее значение это |V| (если модель ничего не знает о данных, она присваивает всем токенам вероятность 1/|V|, независимо от контекста).

$$Perplexity(y_{1:M}) = 2^{-\frac{1}{M}L(y_{1:M})}$$

Привязка веса (Weight Tying)

Большая часть параметров модели берется из этих матриц — вы можете уменьшить размер модели!

<u>Default</u> (no weight tying) Weight tying softmax softmax Output word embeddings Neural Neural Output word embeddings network network different the same embeddings embeddings saw ... a

Примеры сгенерированного текста

Proof. Omitted.

Lemma 0.1. Let C be a set of the construction.

Let C be a gerber covering. Let F be a quasi-coherent sheaves of O-modules. We have to show that

$$\mathcal{O}_{\mathcal{O}_X} = \mathcal{O}_X(\mathcal{L})$$

Proof. This is an algebraic space with the composition of sheaves F on $X_{étale}$ we have

$$O_X(F) = \{morph_1 \times_{O_X} (G, F)\}$$

where G defines an isomorphism $F \to F$ of O-modules.

Lemma 0.2. This is an integer Z is injective.

Proof. See Spaces, Lemma ??.

Lemma 0.3. Let S be a scheme. Let X be a scheme and X is an affine open covering. Let $U \subset X$ be a canonical and locally of finite type. Let X be a scheme. Let X be a scheme which is equal to the formal complex.

The following to the construction of the lemma follows.

Let X be a scheme. Let X be a scheme covering. Let

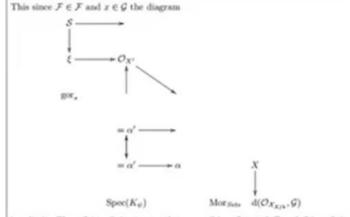
$$b: X \to Y' \to Y \to Y \to Y' \times_X Y \to X.$$

be a morphism of algebraic spaces over S and Y.

Proof. Let X be a nonzero scheme of X. Let X be an algebraic space. Let F be a quasi-coherent sheaf of O_X -modules. The following are equivalent

- F is an algebraic space over S.
- (2) If X is an affine open covering.

Consider a common structure on X and X the functor $O_X(U)$ which is locally of finite type.



is a limit. Then G is a finite type and assume S is a flat and F and G is a finite type f_* . This is of finite type diagrams, and

- the composition of G is a regular sequence,
- O_{X*} is a sheaf of rings.

Proof. We have see that $X = \operatorname{Spec}(R)$ and \mathcal{F} is a finite type representable by algebraic space. The property \mathcal{F} is a finite morphism of algebraic stacks. Then the cohomology of X is an open neighbourhood of U.

Proof. This is clear that G is a finite presentation, see Lemmas ??.

A reduced above we conclude that U is an open covering of C. The functor F is a "field

$$\mathcal{O}_{X,x} \longrightarrow \mathcal{F}_{\mathbb{F}} -1(\mathcal{O}_{X_{Out_k}}) \longrightarrow \mathcal{O}_{X_k}^{-1}\mathcal{O}_{X_k}(\mathcal{O}_{X_n}^{\mathbb{F}})$$

is an isomorphism of covering of \mathcal{O}_{X_1} . If \mathcal{F} is the unique element of \mathcal{F} such that X is an isomorphism.

The property \mathcal{F} is a disjoint union of Proposition ?? and we can filtered set of presentations of a scheme \mathcal{O}_X -algebra with \mathcal{F} are opens of finite type over S. If \mathcal{F} is a scheme theoretic image points.

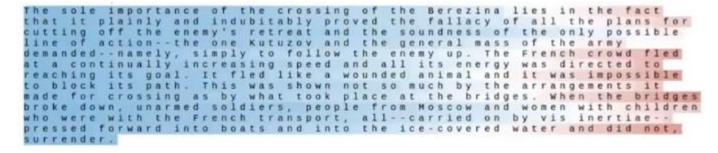
If F is a finite direct sum O_{X_h} is a closed immersion, see Lemma ??. This is a sequence of F is a similar morphism.

Примеры сгенерированного текста

```
* Increment the mise file of the new incorrect MT FILTEN group information
 " of the sire unnerstively.
static int indicate policy(void)
 int error:
 if (fd - HARN EFT) {
    . The kurnel blank will could it to userspace.
    if (se->segment < mem_total)
     unblock graph and set blocked();
   also
     ret - 1
   gote ball;
 segaddr = in_SB(in.addr);
 selector = seg / 16;
 setup works * true;
 for (1 = 0; 1 < blocks; i++) (
   seq = buf[i++];
   bpf = bd->bd.next = i = search;
   if (fd) (
     current = blocked;
 rw->name = "Getjbbregs";
 bprm self clearl(&iv->version);
 regs-bnew = blocks((SPF STATS << info->historidae)) | PFMR CLOSATHING SECONDS << 17;
 return sequable;
```

LSTM-сети на уровне символов, обученные на ядре Linux и «Войне и мире»

• Нейрон чувствителен к положению в строке



• Нейрон активируется на внутренние кавычки



• Нейрон активируется внутри операторов if

• Нейрон активируется на внутренние комментарии и цитаты

• Нейрон чувствителен к глубине кода

```
#ifdef config_AUDITSYSCALL
static inline int audit_match_class_bits(int class, u32 *mask)
{
  int i;
  if (classes[class]) {
    for (i = 0; i < AUDIT_BITMASK_SIZE; i++)
      if (mask[i] & classes[class][i])
      return 0;
}
return 1;
}</pre>
```

• Нейрон активируется перед началом новой строки

```
char *audit_unpack_string(void *bufp, size_t *remain, si
{
    char *str;
    if (I*bufp* || (len == 0) || (len > *remain))
    return ERR_PTR(-EINVAL);

/* Of the currently implemented string fields, PATH_MAX
    *defines the longest valid length.

*/

if (len > PATH_MAX)
    return ERR_PTR(-ENAMETOOLONG);

str = kmalloc(len + 1, GFP_KERNEL);

if (unlikely(!str))
    return ERR_PTR(-ENOMEM);
memcpy(str, 'bufp, len);
str[len] = 0;
    *bufp += len;
    remain -= len;
    return str;
}
```

Нейроны, следящие за настроением

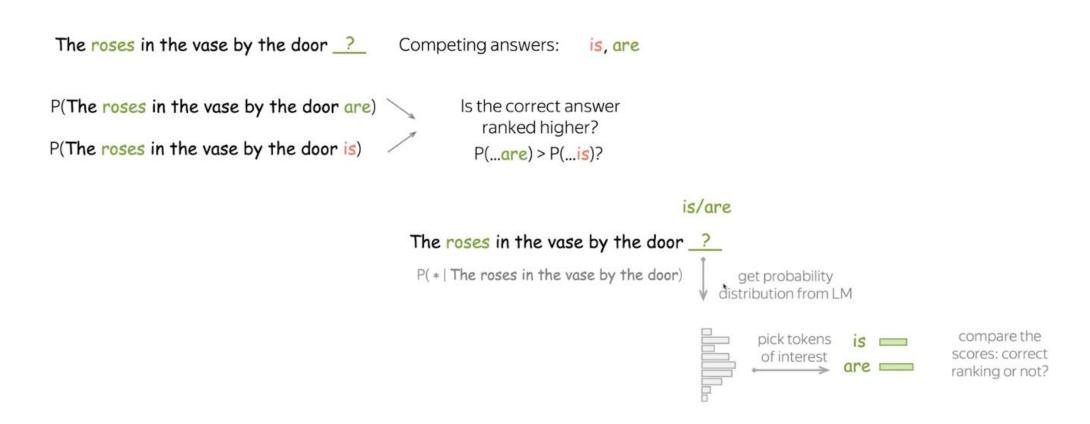
LSTM-сеть на уровне символов, обученная на отзывах Amazon

This is one of Crichton's best books. The characters of Karen Ross, Peter Elliot, Munro, and Amy are beautifully developed and their interactions are exciting, complex, and fast-paced throughout this impressive novel. And about 99.8 percent of that got lost in the film. Seriously, the screenplay AND the directing were horrendous and clearly done by people who could not fathom what was good about the novel. I can't fault the actors because frankly, they never had a chance to make this turkey live up to Crichton's original work. I know good novels, especially those with a science fiction edge, are hard to bring to the screen in a way that lives up to the original. But this may be the absolute worst disparity in quality between novel and screen adaptation ever. The book is really, really good. The movie is just dreadful.

Интерпретируемые нейроны для управления генерируемыми текстами

SENTIMENT FIXED TO POSITIVE	SENTIMENT FIXED TO NEGATIVE
I couldn't figure out the shape at first but it definitely does what it's meant to do. It's a great product and I recommend it highly	I couldn't figure out how to use the product. It did not work. At least there was no quality control; this tablet does not work. I would have given it zero stars, but that was not an option.
I couldn't figure out why this movie had been discontinued! Now I can enjoy it anytime I like. So glad to have found it again.	I couldn't figure out how to set it up being that there was no warning on the box. I wouldn't recommend this to anyone.
I couldn't figure out how to use the video or the book that goes along with it, but it is such a fantastic book on how to put it into practice!	I couldn't figure out how to use the gizmo. What a waste of time and money. Might as well through away this junk.
I couldn't figure out how to use just one and my favorite running app. I use it all the time. Good quality, You cant beat the price.	I couldn't figure out how to stop this drivel. At worst, it was going absolutely nowhere, no matter what I did.Needles to say, I skim-read the entire book. Don't waste your time.
I couldn't figure out how to attach these balls to my little portable drums, but these fit the bill and were well worth every penny.	I couldn't figure out how to play it.

Контрастная оценка: специфический текст



Контрастная оценка: специфический текст

The roses ? Simple: no attractors

The roses in the vase ? Harder: 1 attractor

The roses in the vase by the door ? Harder: 2 attractors

Attractors: nouns with different number than the subject

ЛР 2

Квантизация - это метод снижения вычислительных затрат и затрат памяти при выполнении логического вывода за счет представления весов и активаций с использованием типов данных низкой точности.

Quantization_config — это параметр (аргумент) метода from_pretrained в библиотеке Hugging Face transformers, который используется для передачи конфигурации квантования в модель.

BitsAndBytesConfig — это класс из библиотеки transformers от Hugging Face, который предназначен для настройки и управления квантованием моделей с помощью библиотеки bitsandbytes.

ЛР 3

Pydantic играет ключевую роль при работе со структурированными данными в LLM.

Самое главное применение — заставить LLM возвращать данные в строго определённом формате.

JSON string — это строка текста, которая содержит данные в формате JSON (JavaScript Object Notation).

JSON string — это **текстовое представление структурированных данных**, которое выглядит как JavaScript-объект, но является обычной строкой.