

344.175 VL: Natural Language Processing

n-Gram Language Models



Navid Rekab-saz

Email: navid.rekabsaz@jku.at

Office hours: <https://navid-officehours.youcanbook.me>

Agenda

- n -gram language models
- Count-based n -gram LM
- Neural n -gram LM

Agenda

- ***n*-gram language models**
- Count-based *n*-gram LM
- Neural *n*-gram LM

Language Modeling

- Language Modeling is the task of predicting a word (or a subword or character) given a context:

$$P(v|\text{context})$$

- A Language Model (LM) can answer the questions like



$$P(v|the\ students\ opened\ their)$$

Language Modeling – formal definition

- Given a sequence of words $x^{(1)}, x^{(2)}, \dots, x^{(t)}$, a **language model** calculates the **probability distribution** of next word $x^{(t+1)}$ over all words in vocabulary

$$P(x^{(t+1)} | x^{(1)}, \dots, x^{(t-1)}, x^{(t)})$$

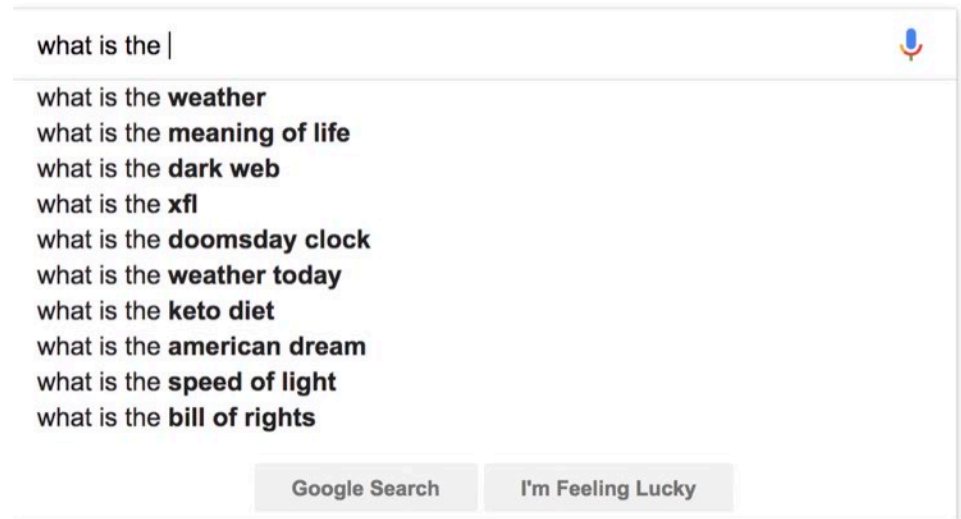
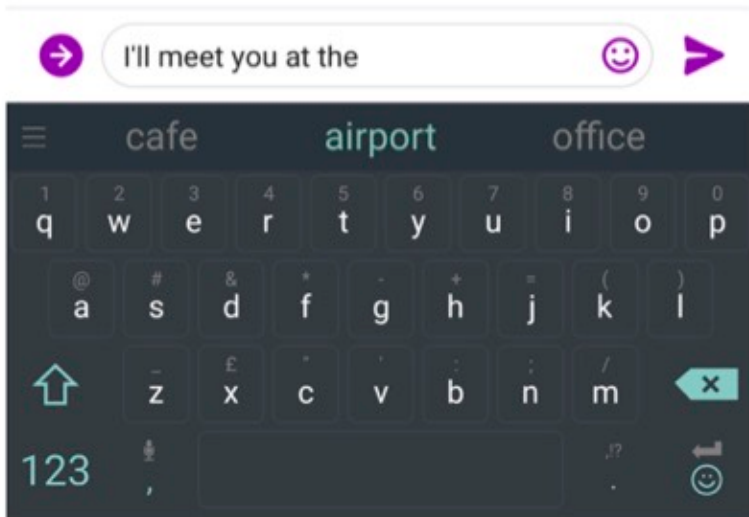
x is any word in the vocabulary $\mathbb{V} = \{v_1, v_2, \dots, v_N\}$

☞ Note: this is the definition of directed left-to-right language models.

Why Language Modeling?

- Language Modeling is a **benchmark task** that helps us measure our progress on **understanding language**
- LMs are a **subcomponent** of many NLP tasks, especially those involving **generating text** or estimating the **probability of text**:
 - Predictive typing
 - Spelling/grammar correction
 - Automatic speech recognition (ASR)
 - Handwriting recognition
 - Machine translation
 - Summarization
 - Dialogue /chatbots
 - etc.

Direct usages for next word prediction



Probability of a Text

- A Language Model can also assign probability to the validity a piece of text
 - How probable it is that a sentence appears in a language.

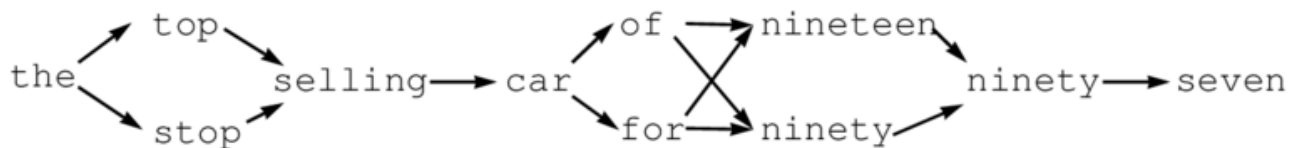
$$P(x^{(1)}, \dots, x^{(T)}) = ?$$

- According to a (directed left-to-right) Language Model, the probability of a given text is computed by:

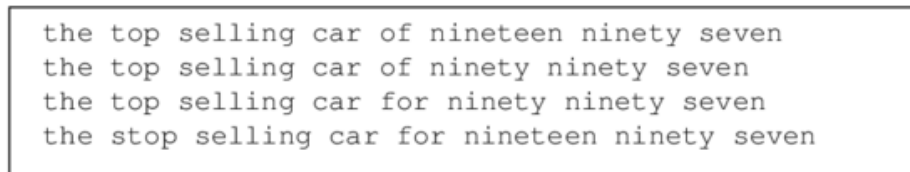
$$P(x^{(1)}, \dots, x^{(T)}) = P(x^{(1)}) \times P(x^{(2)} | x^{(1)}) \times \dots \times P(x^{(T)} | x^{(1)}, \dots, x^{(T-1)})$$

$$P(x^{(1)}, \dots, x^{(T)}) = \prod_{t=1}^T P(x^{(t)} | x^{(1)}, \dots, x^{(t-1)})$$

Usage in Automatic Speech Recognition (ASR)



Input Lattice



N-Best List

n-gram Language Model

- Recall: a *n*-gram is a chunk of *n* consecutive words.

the students opened their _____

- unigrams: “*the*”, “*students*”, “*opened*”, “*their*”
 - bigrams: “*the students*”, “*students opened*”, “*opened their*”
 - trigrams: “*the students opened*”, “*students opened their*”
 - 4-grams: “*the students opened their*”
-
- A *n*-gram Language Model collects frequency statistics of different *n*-grams in a corpus, and use these to calculate probabilities

***N*-gram LM as a conditional probability**

- **Markov assumption:** decision at time t depends only on the current state
- In n -gram Language Model: predicting $x^{(t+1)}$ depends on preceding $n-1$ words
- Without Markovian assumption:

$$P(x^{(t+1)} | x^{(1)}, \dots, x^{(t-1)}, x^{(t)})$$

- n -gram Language Model:

$$P(x^{(t+1)} | x^{(t-n+2)}, \dots, x^{(t-1)}, x^{(t)})$$


 $n-1$ words



Agenda

- n -gram language models
- **Count-based n -gram LM**
- Neural n -gram LM

n-gram LM using term counts

- Based on definition of conditional probability:

$$P(x^{(t+1)} | x^{(t-n+2)}, \dots, x^{(t)}) = \frac{P(x^{(t-n+2)}, \dots, x^{(t)}, x^{(t+1)})}{P(x^{(t-n+2)}, \dots, x^{(t)})}$$

- The *n*-gram probability is calculated by counting *n*-grams and [*n*−1]-grams in a large corpus of text:

$$P(x^{(t+1)} | x^{(t)}, \dots, x^{(t-n+2)}) \approx \frac{\text{count}(x^{(t-n+2)}, \dots, x^{(t)}, x^{(t+1)})}{\text{count}(x^{(t-n+2)}, \dots, x^{(t)})}$$

Example

- Example: learning a 4-gram Language Model

~~as the exam clerk started the clock, the students opened their _____~~



condition on this

$$P(v | \text{students opened their}) = \frac{P(\text{students opened their } v)}{P(\text{students opened their})}$$

- For example, suppose that in the corpus:
 - “students opened their” occurred 1000 times
 - “students opened their books” occurred 400 times
 - $P(\text{books} | \text{students opened their}) = 0.4$
 - “students opened their exams” occurred 100 times
 - $P(\text{exams} | \text{students opened their}) = 0.1$

Example – a bigram LM

- Trained on the data of a restaurant dialogue system

Bigram counts:

	i	want	to	eat	chinese	food	lunch	spend
i	5	827	0	9	0	0	0	2
want	2	0	608	1	6	6	5	1
to	2	0	4	686	2	0	6	211
eat	0	0	2	0	16	2	42	0
chinese	1	0	0	0	0	82	1	0
food	15	0	15	0	1	4	0	0
lunch	2	0	0	0	0	1	0	0
spend	1	0	1	0	0	0	0	0

Bigram LM:

	i	want	to	eat	chinese	food	lunch	spend
i	0.002	0.33	0	0.0036	0	0	0	0.00079
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0054	0.0011
to	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087
eat	0	0	0.0027	0	0.021	0.0027	0.056	0
chinese	0.0063	0	0	0	0	0.52	0.0063	0
food	0.014	0	0.014	0	0.00092	0.0037	0	0
lunch	0.0059	0	0	0	0	0.0029	0	0
spend	0.0036	0	0.0036	0	0	0	0	0

Count-based n -gram LMs – limitations

■ Sparsity

- What if the denominator never occurred in corpus?
 - **Backoff**: Probability is calculated for lower n -grams.
 - E.g., in „*students opened their v*“, if „*students opened their*“ does not exist, language model probability is calculated for “*opened their*”
 - Trigram is backed off to let a bigram do the job!
- What if the nominator never occurred in corpus?
 - Approached by various **smoothing** methods

Laplace smoothing

- Add a small number like $\delta = 1$ to the count of all words:

$$P(x^{(t+1)} | x^{(t-n+2)}, \dots, x^{(t)}) = \frac{\#(x^{(t-n+2)}, \dots, x^{(t)}, x^{(t+1)})}{\#(x^{(t-n+2)}, \dots, x^{(t)})}$$

$$P(x^{(t+1)} | x^{(t-n+2)}, \dots, x^{(t)}) = \frac{\#(x^{(t-n+2)}, \dots, x^{(t)}, x^{(t+1)})}{\sum_{v \in V} \#(x^{(t-n+2)}, \dots, x^{(t)}, v)}$$

$$P_{Laplace}(x^{(t+1)} | x^{(t-n+2)}, \dots, x^{(t)}) = \frac{\#(x^{(t-n+2)}, \dots, x^{(t)}, x^{(t+1)}) + 1}{\sum_{v \in V} (\#(x^{(t-n+2)}, \dots, x^{(t)}, v) + 1)}$$

$$P_{Laplace}(x^{(t+1)} | x^{(t-n+2)}, \dots, x^{(t)}) = \frac{\#(x^{(t-n+2)}, \dots, x^{(t)}, x^{(t+1)}) + 1}{\sum_{v \in V} (\#(x^{(t-n+2)}, \dots, x^{(t)}, v)) + |V|}$$

$$P_{Laplace}(x^{(t+1)} | x^{(t-n+2)}, \dots, x^{(t)}) = \frac{\#(x^{(t-n+2)}, \dots, x^{(t)}, x^{(t+1)}) + 1}{\#(x^{(t-n+2)}, \dots, x^{(t)}) + |V|}$$

Laplace smoothing example

Original bigram counts:

	i	want	to	eat	chinese	food	lunch	spend
i	5	827	0	9	0	0	0	2
want	2	0	608	1	6	6	5	1
to	2	0	4	686	2	0	6	211
eat	0	0	2	0	16	2	42	0
chinese	1	0	0	0	0	82	1	0
food	15	0	15	0	1	4	0	0
lunch	2	0	0	0	0	1	0	0
spend	1	0	1	0	0	0	0	0

Bigram counts added by 1:

	i	want	to	eat	chinese	food	lunch	spend
i	6	828	1	10	1	1	1	3
want	3	1	609	2	7	7	6	2
to	3	1	5	687	3	1	7	212
eat	1	1	3	1	17	3	43	1
chinese	2	1	1	1	1	83	2	1
food	16	1	16	1	2	5	1	1
lunch	3	1	1	1	1	2	1	1
spend	2	1	2	1	1	1	1	1

Laplace smoothing example

LM without smoothing:

	i	want	to	eat	chinese	food	lunch	spend
i	0.002	0.33	0	0.0036	0	0	0	0.00079
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0054	0.0011
to	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087
eat	0	0	0.0027	0	0.021	0.0027	0.056	0
chinese	0.0063	0	0	0	0	0.52	0.0063	0
food	0.014	0	0.014	0	0.00092	0.0037	0	0
lunch	0.0059	0	0	0	0	0.0029	0	0
spend	0.0036	0	0.0036	0	0	0	0	0

LM with Laplace smoothing:

	i	want	to	eat	chinese	food	lunch	spend
i	0.0015	0.21	0.00025	0.0025	0.00025	0.00025	0.00025	0.00075
want	0.0013	0.00042	0.26	0.00084	0.0029	0.0029	0.0025	0.00084
to	0.00078	0.00026	0.0013	0.18	0.00078	0.00026	0.0018	0.055
eat	0.00046	0.00046	0.0014	0.00046	0.0078	0.0014	0.02	0.00046
chinese	0.0012	0.00062	0.00062	0.00062	0.00062	0.052	0.0012	0.00062
food	0.0063	0.00039	0.0063	0.00039	0.00079	0.002	0.00039	0.00039
lunch	0.0017	0.00056	0.00056	0.00056	0.00056	0.0011	0.00056	0.00056
spend	0.0012	0.00058	0.0012	0.00058	0.00058	0.00058	0.00058	0.00058

Count-based n -gram LMs – limitations

■ Sparsity

- What if the denominator never occurred in corpus?
 - **Backoff**: Probability is calculated for lower n -grams.
 - E.g., in „*students opened their v*“, if „*students opened their*“ does not exist, language model probability is calculated for “*opened their*”
 - Trigram is backed off to let a bigram do the job!
- What if the nominator never occurred in corpus?
 - Approached by various **smoothing** methods
- Sparsity issue becomes even more prominent in higher n -gram!

Count-based n -gram LMs – limitations

■ Storage

- An n -gram language model needs to store all levels of n -grams, from uni- to n -gram, observed in the corpus
- Increasing n radically worsens the storage problem!

■ No understanding of tokens relations

- Semantic and syntactic relations between words are fully ignored
 - “*book*” and “*books*” or “*car*” and “*automobile*” are treated completely separately

Generating text

- A trigram LM trained on Reuters corpus (1.7 M words)

today the _____

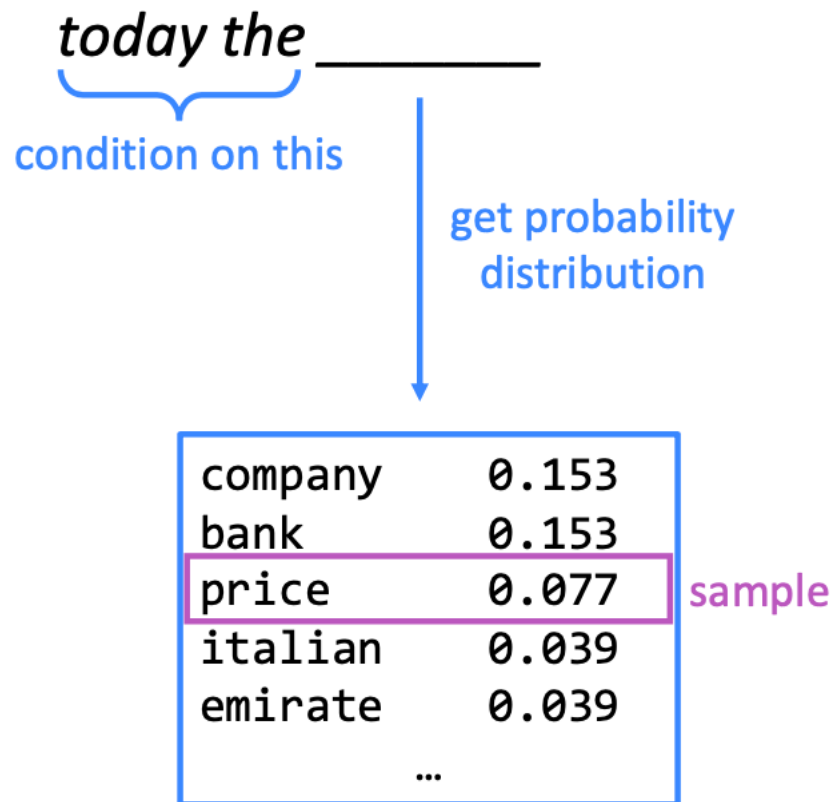
get probability
distribution

company	0.153
bank	0.153
price	0.077
italian	0.039
emirate	0.039
...	

Sparsity problem:
not much granularity
in the probability
distribution

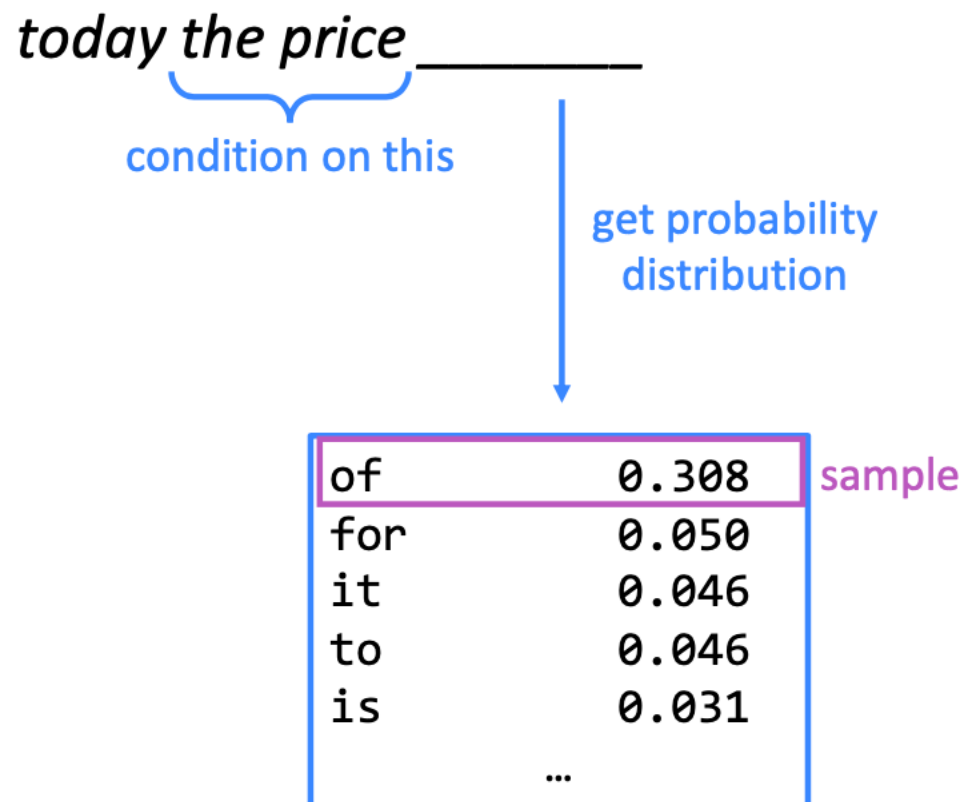
Generating text

- Generating text by sampling from the probability distributions



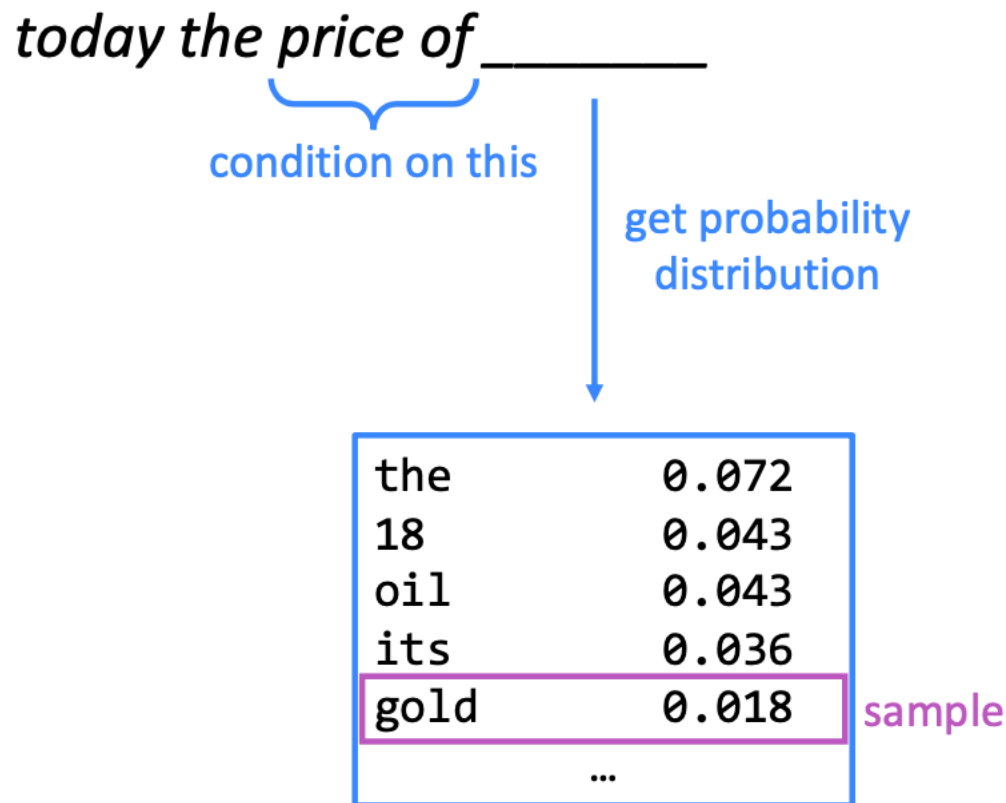
Generating text

- Generating text by sampling from the probability distributions



Generating text

- Generating text by sampling from the probability distributions



Generating text

- Generating text by sampling from the probability distributions

today the price of gold per ton , while production of shoe lasts and shoe industry , the bank intervened just after it considered and rejected an imf demand to rebuild depleted european stocks , sept 30 end primary 76 cts a share .

- Decently good in **syntax** ... but **incoherent!**
- Increasing n makes the text more coherent but also intensifies the discussed issues

Generating text

- n -gram LMs trained on Shakespeare's works

1

gram

–To him swallowed confess hear both. Which. Of save on trail for are ay device and rote life have

–Hill he late speaks; or! a more to leg less first you enter

2

gram

–Why dost stand forth thy canopy, forsooth; he is this palpable hit the King Henry. Live king. Follow.

–What means, sir. I confess she? then all sorts, he is trim, captain.

3

gram

–Fly, and will rid me these news of price. Therefore the sadness of parting, as they say, 'tis done.

–This shall forbid it should be branded, if renown made it empty.

4

gram

–King Henry. What! I will go seek the traitor Gloucester. Exeunt some of the watch. A great banquet serv'd in;

–It cannot be but so.

Generating text

- n -gram LMs trained on Wall Street Journal

1
gram

Months the my and issue of year foreign new exchange's september were recession exchange new endorsed a acquire to six executives

2
gram

Last December through the way to preserve the Hudson corporation N. B. E. C. Taylor would seem to complete the major central planners one point five percent of U. S. E. has already old M. X. corporation of living on information such as more frequently fishing to keep her

3
gram

They also point to ninety nine point six billion dollars from two hundred four oh six three percent of the rates of interest stores as Mexico and Brazil on market conditions

Agenda

- n -gram language models
- Count-based n -gram LM
- **Neural n -gram LM**

n-gram language modeling with neural networks

Recall

- The aim of a *n*-gram Language Model is to calculate:

$$P(x^{(t+1)} | x^{(t-n+2)}, \dots, x^{(t)})$$

- We can use a feed forward neural network to estimate this probability
- Immediate benefits:
 - Smooth probability estimation
 - Exploiting the semantic space of word embeddings (probably better generalization)

Neural n -gram LM – preparing training data

- Preparing training data for a neural 4-gram Language Model in the form of (context, next word), namely $(x^{(t-2)}x^{(t-1)}x^{(t)}, x^{(t+1)})$:

- For a given text corpus:

a fluffy cat sunbathes on the bank of river ...

- Training data items would be:

(<bos> <bos> <bos>, a)

(<bos> <bos> a, fluffy)

(<bos> a fluffy, cat)

(a fluffy cat, sunbathes)

(fluffy cat sunbathes, on)

(cat sunbathes on, the)

(sunbathes on the, bank)

...

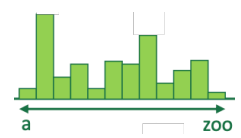
<bos> is a special token added to dictionary, referring to *beginning of sentence*

Also, dot is commonly replaced with <eos> token – *end of sentence*

Neural n -gram Language Model – architecture

$$P(x^{(t+1)} | \text{a fluffy cat}) = \hat{y}_{x^{(t+1)}}^{(t)}$$

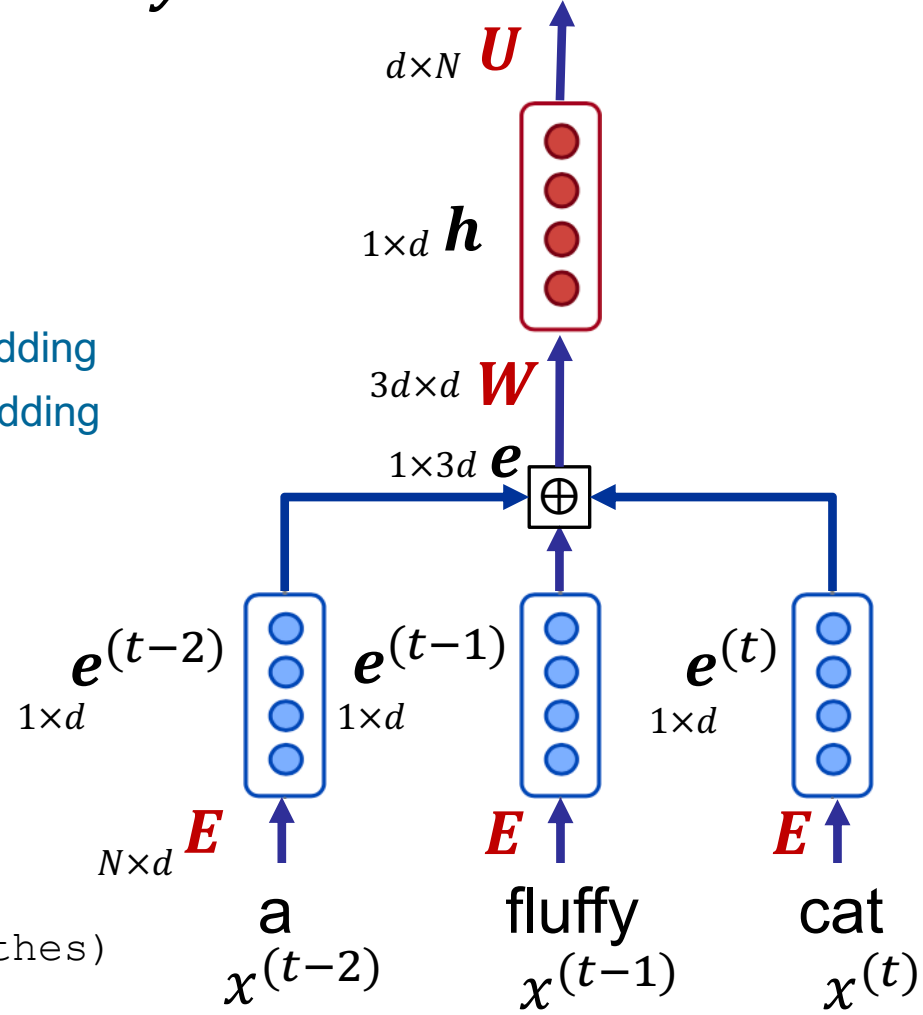
$$P(\text{sunbathes} | \text{a fluffy cat}) = \hat{y}_{\text{sunbathes}}^{(t)}$$



- $E \rightarrow N \times d$
- $W \rightarrow d \times ((n - 1) \times d)$
- $U \rightarrow d \times N$

E is called encoder embedding
 U is called decoder embedding or output projection

N size of vocabulary
 d embeddings dimension



A data item in training data:
(a fluffy cat, sunbathes)

N size of vocabulary
 d embeddings dimension
 $(n - 1)$ number of preceding words
Parameters are shown in red

Formulation

Encoder

- From word to word embedding:
 - One-hot vector of word $x^{(t)}$ is $\mathbf{x}^{(t)}$ vector: $\mathbf{x}^{(t)} \rightarrow 1 \times N$
 - In $\mathbf{x}^{(t)}$, all values are 0 and only the value corresponding to the word $x^{(t)}$ is set to 1
 - Fetching word embedding: $\mathbf{e}^{(t)} = \mathbf{x}^{(t)} \mathbf{E}$ $\mathbf{e}^{(t)} \rightarrow 1 \times d$
 - In practice, $\mathbf{e}^{(t)}$ is achieved by fetching the vector of $x^{(t)}$ from \mathbf{E} . No need for $\mathbf{x}^{(t)}$ in practice
- Concatenation of $(n - 1)$ word embeddings:

$$\mathbf{e} = [\mathbf{e}^{(t-2)}, \mathbf{e}^{(t-1)}, \mathbf{e}^{(t)}] \quad \mathbf{e} \rightarrow 1 \times (n - 1)d$$

- Hidden layer: $\mathbf{h} = \tanh(\mathbf{e}\mathbf{W} + \mathbf{b}^w)$ $\mathbf{h} \rightarrow 1 \times d$

$$\mathbf{E} \rightarrow N \times d$$

$$\mathbf{W} \rightarrow (n - 1)d \times d$$

$$\mathbf{b}^w \rightarrow 1 \times d$$

Formulation

N size of vocabulary

d embeddings dimension

$(n - 1)$ number of preceding words

Parameters are shown in red

Decoder

- Predicted logits:

$$\mathbf{z} = \mathbf{h}\mathbf{U} + \mathbf{b}^u \quad \mathbf{z} \rightarrow 1 \times N$$

- Predicted probability distribution:

$$\hat{\mathbf{y}}^{(t)} = \text{softmax}(\mathbf{z}) \quad \hat{\mathbf{y}}^{(t)} \rightarrow 1 \times N$$

- Probability of any next word v at step t :

$$P(v|x^{(t)}, \dots, x^{(t-n+2)}) = \hat{y}_v^{(t)}$$

$$\mathbf{U} \rightarrow d \times N$$

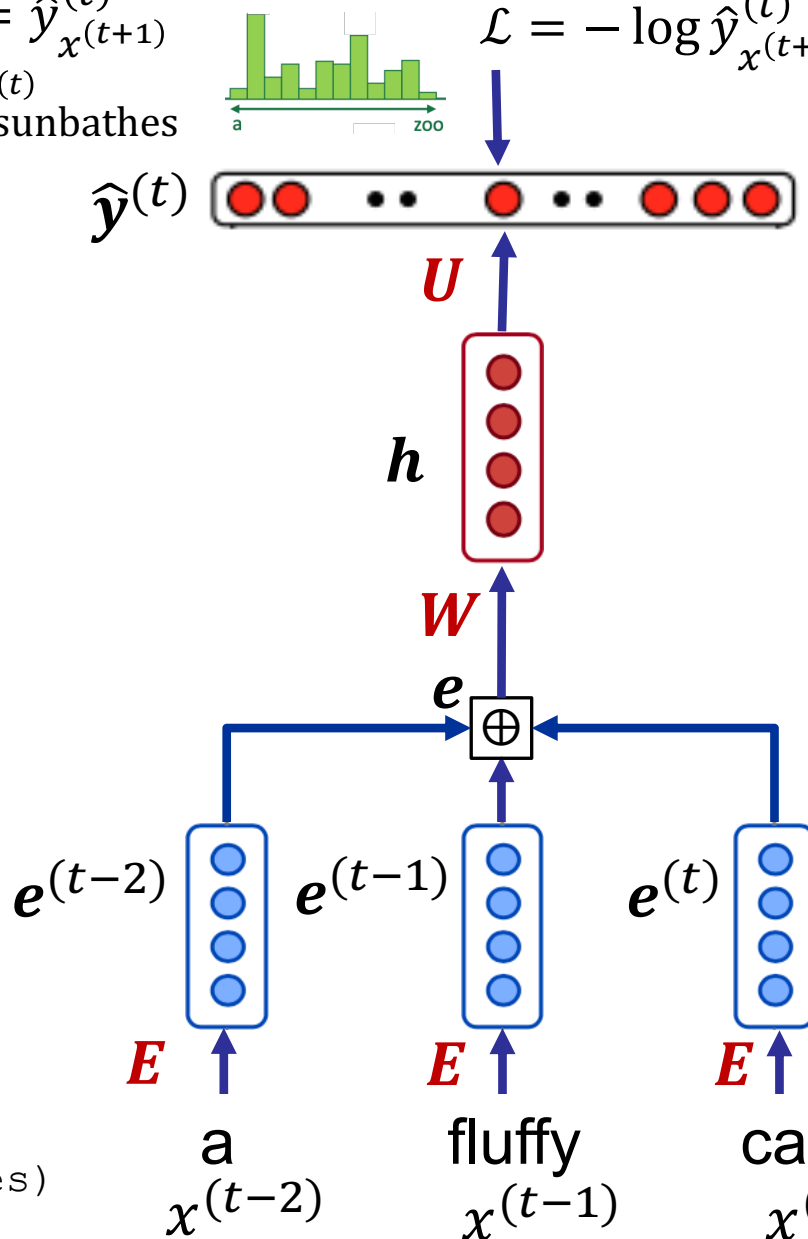
$$\mathbf{b}^u \rightarrow 1 \times N$$

Loss function

$$P(x^{(t+1)} | \text{a fluffy cat}) = \hat{y}_{x^{(t+1)}}^{(t)}$$

$$P(\text{sunbathes} | \text{a fluffy cat}) = \hat{y}_{\text{sunbathes}}^{(t)}$$

$$\mathcal{L} = -\log \hat{y}_{x^{(t+1)}}^{(t)} = -\log \hat{y}_{\text{sunbathes}}^{(t)}$$



A data item in training data:
(a fluffy cat, sunbathes)

Training n -gram neural LM

- Start with a large text corpus: $x^{(1)}, \dots, x^{(T)}$
- In every step t , give $n-1$ previous words as input and output the predicted probability distribution of the next words $\hat{y}^{(t)}$
- Loss function at t is Negative Log Likelihood of the predicted probability for the actual next word $x^{(t+1)}$

$$\mathcal{L}^{(t)} = -\log P(x^{(t+1)} | x^{(t-n+2)}, \dots, x^{(t)})$$

$$\mathcal{L}^{(t)} = -\log \hat{y}_{x^{(t+1)}}^{(t)}$$

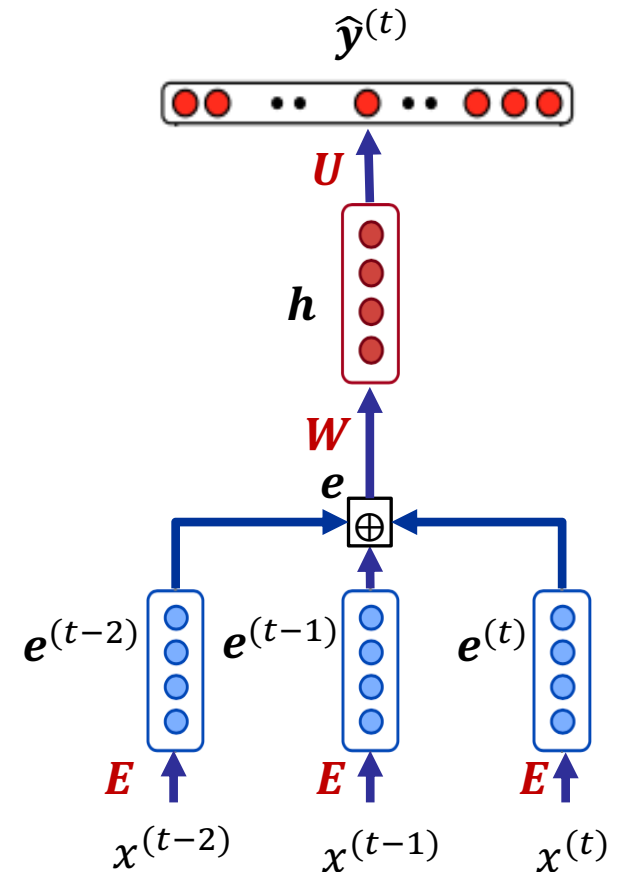
- Overall loss is the average over all time steps:
 - In practice, loss is calculated over batches of text chunks

$$\mathcal{L} = \frac{1}{T} \sum_{t=1}^T \mathcal{L}^{(t)}$$

About transfer learning

- E provides a vector for each word in dictionary
- We can initialize (some) vectors of E with pre-trained word embeddings like GloVe or word2vec
 - In this case, for every word in E we fetch the corresponding vector from a pre-trained word embedding model
 - If the word doesn't exist, as before, its vector is randomly initialized
- We can also do the same for U
- This *better* initialization of parameters is a form of **transfer learning**

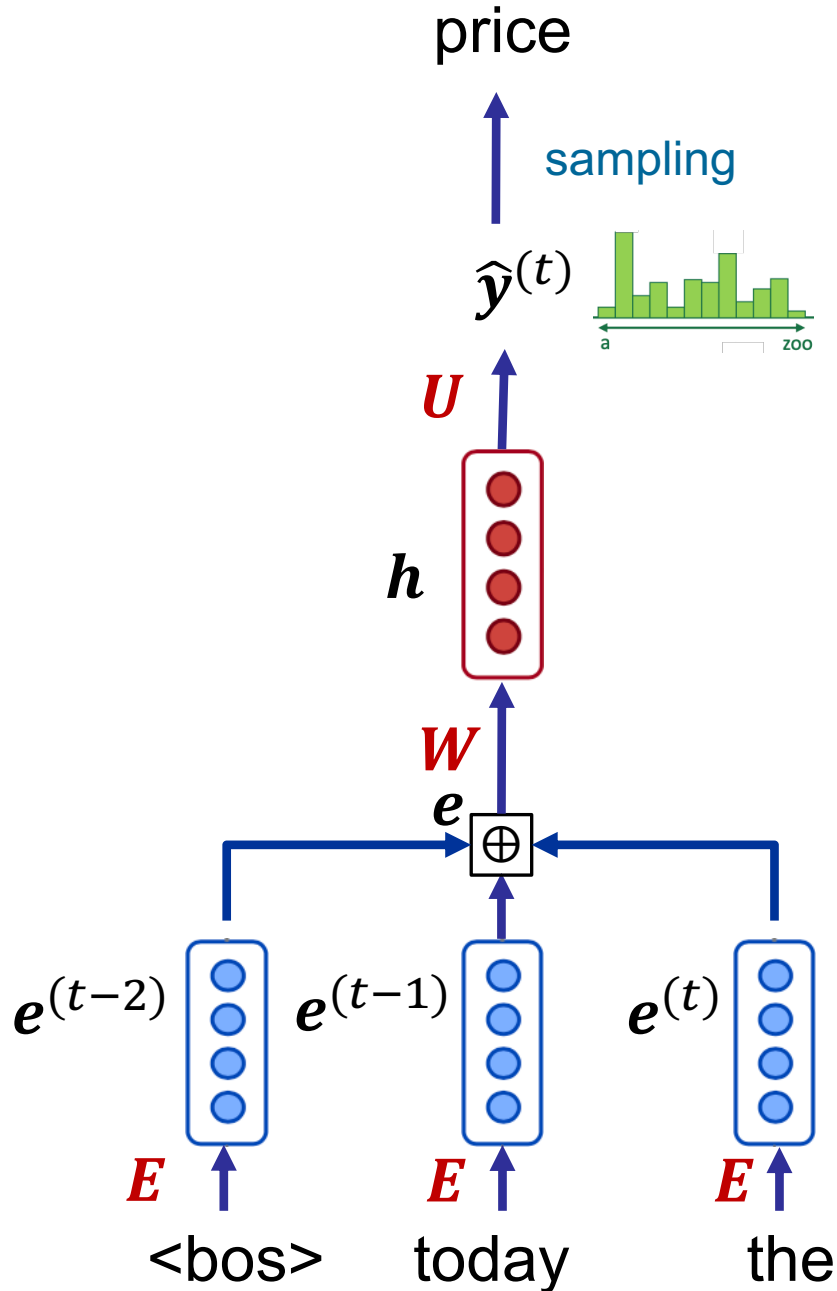
👉 Even without transfer learning, after training the LM, the E and U matrices provide *proper* word embeddings which can be used independently



Generating text

text seed:
today the

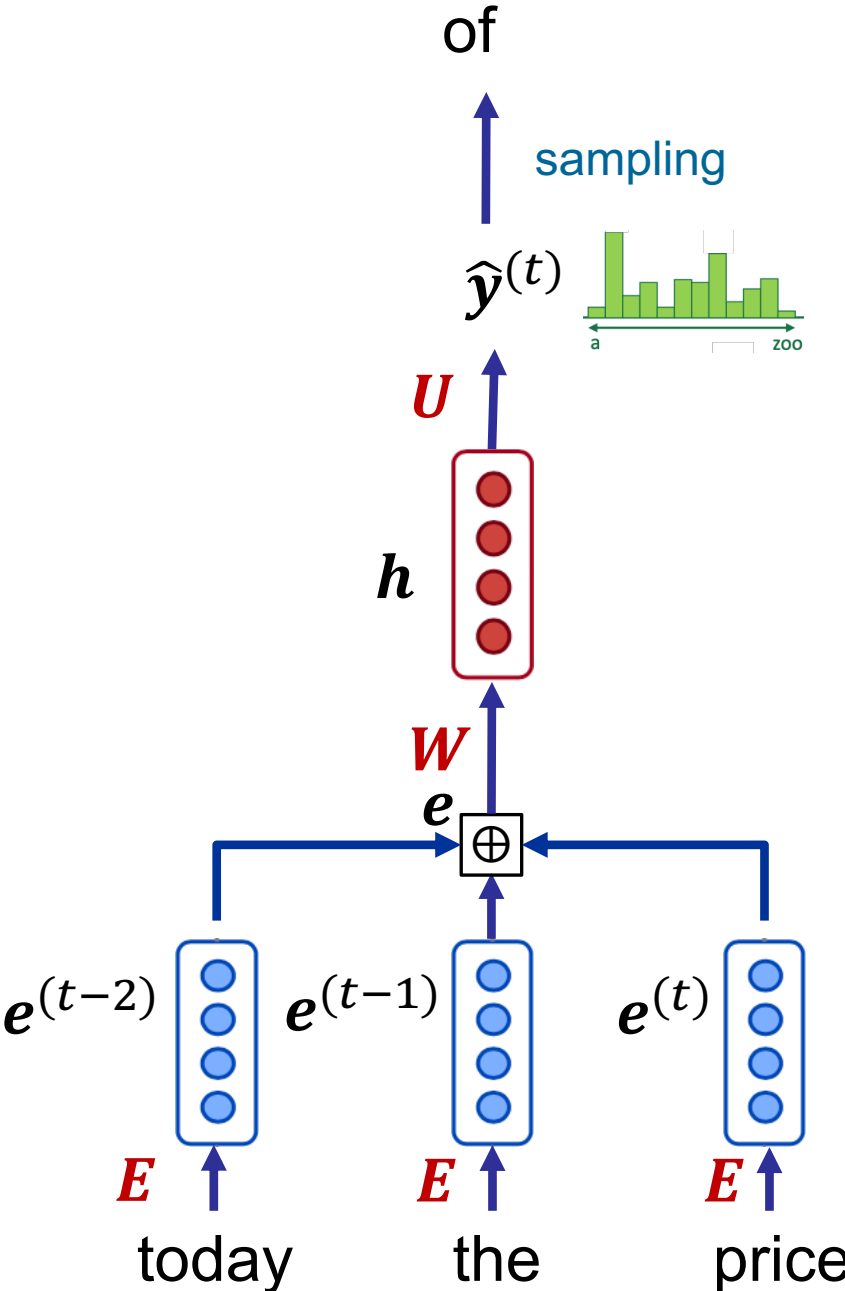
after generation:
today the price



Generating text

(before):
today the price

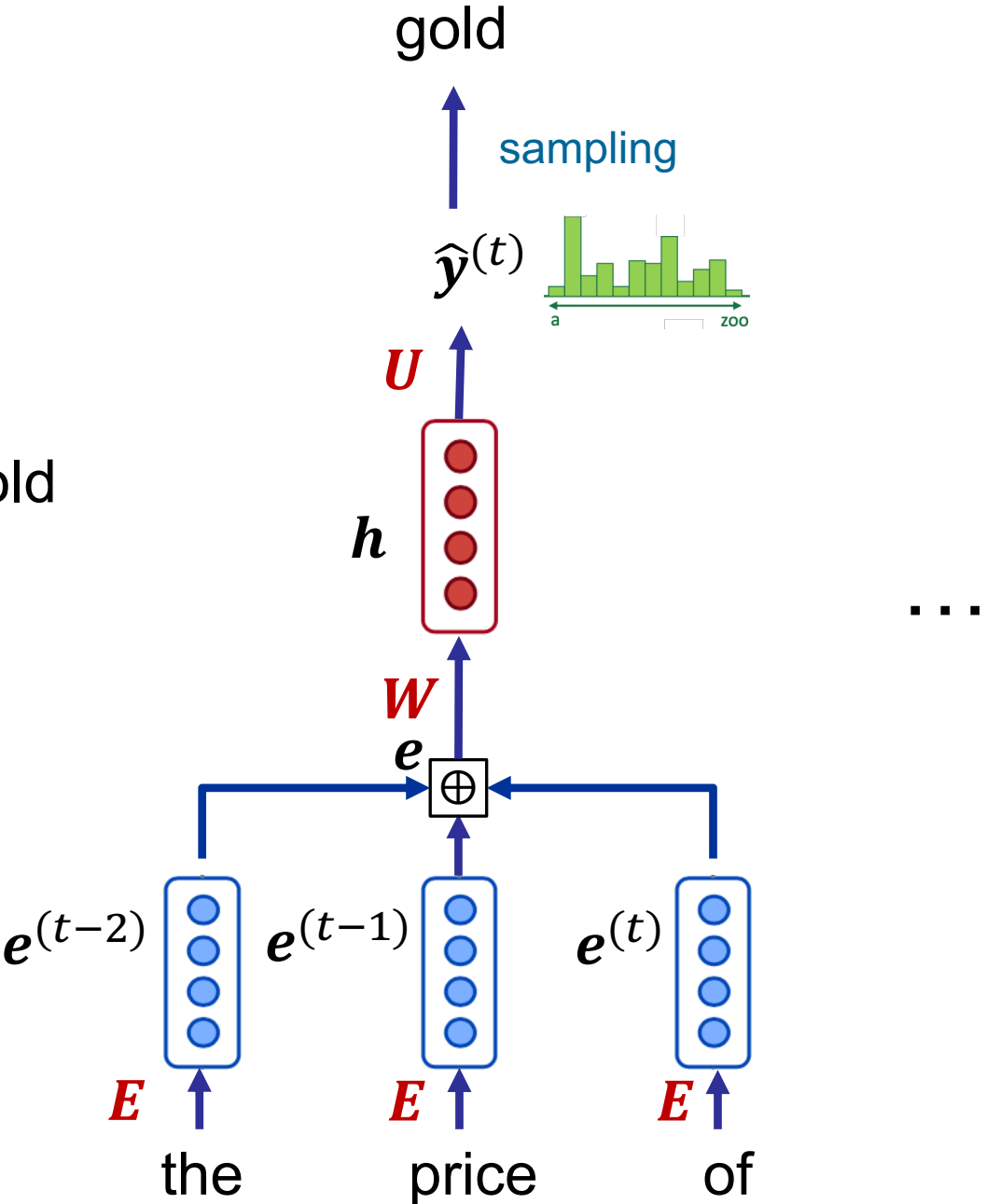
(after):
today the price of



Generating text

(before):
today the price of

(after):
today the price of gold



Neural n -gram LMs – summary

- Neural n -gram LMs predict next word probabilities
- Neural n -gram LMs benefit from semantic relations of words, provided by the encoder and decoder embeddings
- Neural n -gram LMs provide a smooth probability distribution
- At inference time, neural n -gram LMs require a forward pass
 - Count-based n -gram LMs might be more convenient at inference time in practice, since they calculate probabilities simply from the stored counts

