

# Lecture 05: Language modelling

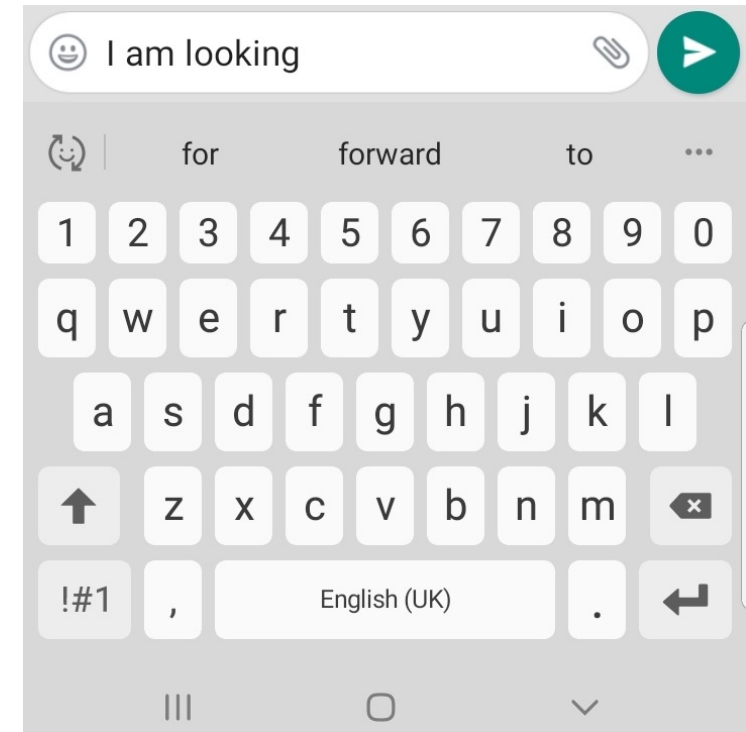
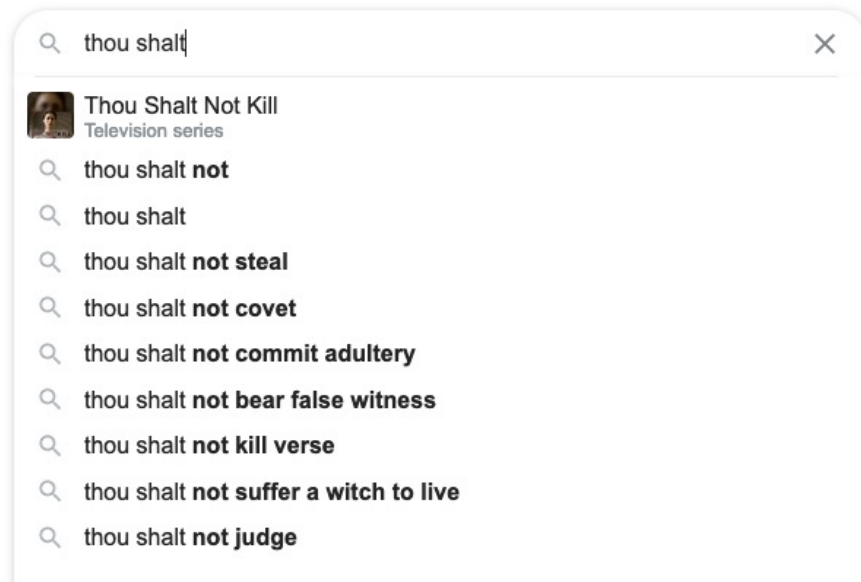
# Overview

---

- Language modelling
- Probabilistic Language modelling with N-gram
- Neural language model
- Evaluating language model

# Next word prediction

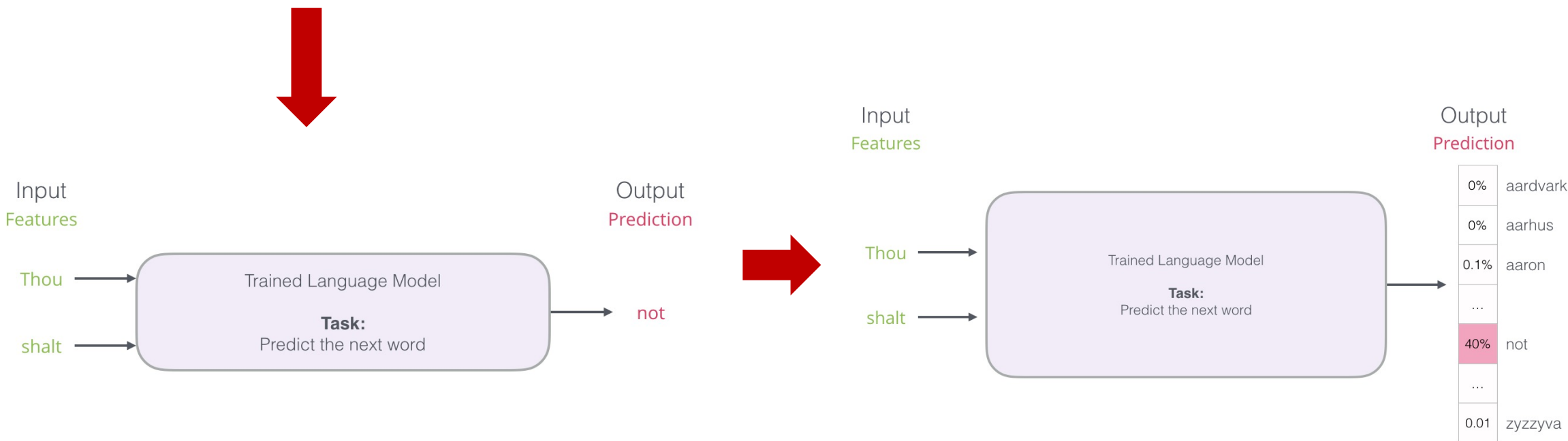
---



# Next word prediction

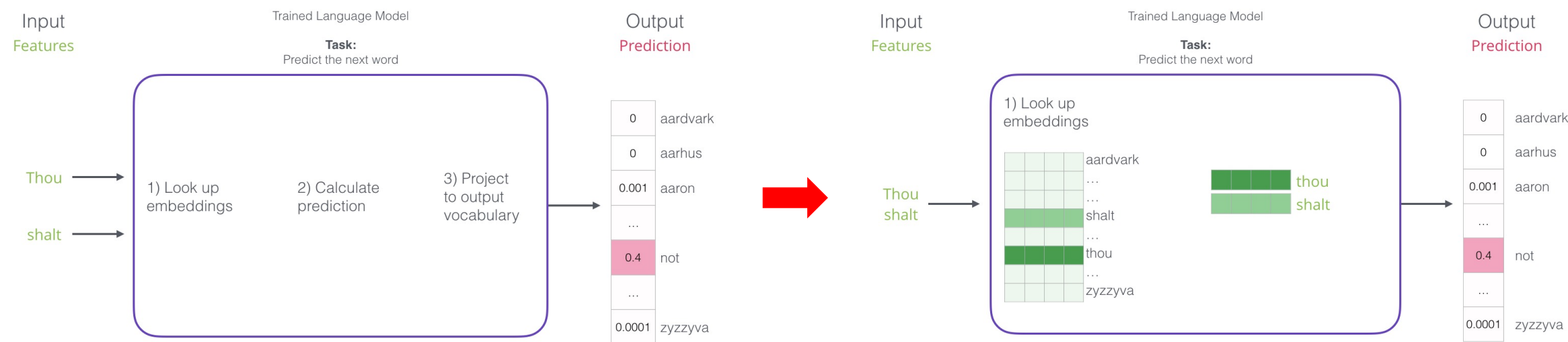
input/feature #1      input/feature #2      output/label

Thou    shalt    \_\_\_\_\_



# Language modeling

## Three main steps in a neural language models



# Language modeling

---

- Language modelling is the task of predicting the next word given an input sequence of words over a vocabulary  $V$ .
  - a system that assigns probabilities to a piece of text from a vocabulary  $V$ .
- Given a sequence of words  $w^{(1)}, w^{(2)}, \dots, w^{(t)}$ , a language model compute the probability distribution of the next word  $w^{(t+1)}$

$$P(w^{(t+1)} | w^{(t)}, \dots, w^{(1)}), \quad \forall w \in V = \{w^1, \dots, w^{|V|}\}$$

- Model that assigns a probability to the sequence of words

$$P(w^{(1)}, w^{(2)}, \dots, w^{(t)})$$

# Language modelling frameworks

---

Currently there are two general techniques of training a language model these can be classed as

- Traditional or probabilistic language model
  - N-gram model
    - Chain rule of probability
    - Markov assumption
- Neural language model
  - Based on neural networks (RNN)
  - LSTM, GRU,
    - Bidirectional RNN
    - Seq2seq (Encoder-decoder)

# Language model – n-grams

---

- N-gram language models are based on probabilities of chunks of word.

*The student open their*

- Definition: An n-gram is a chunk of n consecutive words.
  - **u**nigram: unit of single word – “the”, “student”, “opened”, “their”
  - **b**igrams: unit of double words – “the student”, “student opened”, “opened their”
  - **t**riagram: unit of triple words – “the student opened”, “student opened their”
  - **4**-gram: unit of 4 words – “the student opened their”
- The main idea behind *n-gram* models is to collect statistics about the frequency of different *n-grams* and use this to predict the next word.



# N-gram language models

---

- Given a sequence of words  $w^1, w^2, \dots, w^T$ , the probability of this sequence occurring according to the language model is obtained using the chain rule:

$$\begin{aligned} P(w^1, w^2, \dots, w^T) &= P(w^1) \times P(w^2 | w^1) \times \dots \times P(w^T | w^{T-1}, \dots, w^1) \\ &= \prod_{t=1}^T P(w^t | w^{t-1}, \dots, w^1) \end{aligned}$$

# n-gram Language Models

---

## The chain rule

$$P(w_1, w_2, \dots, w_{T-1}, w_T) = \prod_{t=1}^T P(w_t | w_{t-1}, w_{t-2}, \dots, w_1)$$

<b>the</b>	cat	sat	on	the	mat	$P(w_1)$
the	<b>cat</b>	sat	on	the	mat	$P(w_2   w_1)$
the	cat	<b>sat</b>	on	the	mat	$P(w_3   w_2, w_1)$
the	cat	sat	<b>on</b>	the	mat	$P(w_4   w_3, w_2, w_1)$
the	cat	sat	on	<b>the</b>	mat	$P(w_5   w_4, w_3, w_2, w_1)$
the	cat	sat	on	the	<b>mat</b>	$P(w_6   w_5, w_4, w_3, w_2, w_1)$

# n-gram language model

---

## *Markov assumptions:*

- $w^{(t+1)}$  depends only on the preceding  $n-1$  words.

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

- By conditional probability theory we get

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

Prob of n-gram

Prob of (n-1)-gram

# n-gram Language Models

---

Markov assumption on n-gram

$$P(w_1, w_2, \dots, w_{T-1}, w_T) \approx \prod_{t=1}^T P(w_t | w_{t-1}, \dots, w_{t-n+1})$$

<b>the</b>	cat	sat	on	the	mat	$P(w_1)$
the	<b>cat</b>	sat	on	the	mat	$P(w_2   w_1)$
the	cat	<b>sat</b>	on	the	mat	$P(w_3   w_2, w_1)$
the	cat	sat	<b>on</b>	the	mat	$P(w_4   w_3, w_2)$
the	cat	sat	on	<b>the</b>	mat	$P(w_5   w_4, w_3)$
the	cat	sat	on	the	<b>mat</b>	$P(w_6   w_5, w_4)$

# n-gram Language Models

---

How do we get all these probabilities?

- The n-gram and *(n-1)-gram* probabilities are estimated by counting the frequency of occurrence in the vocabulary:

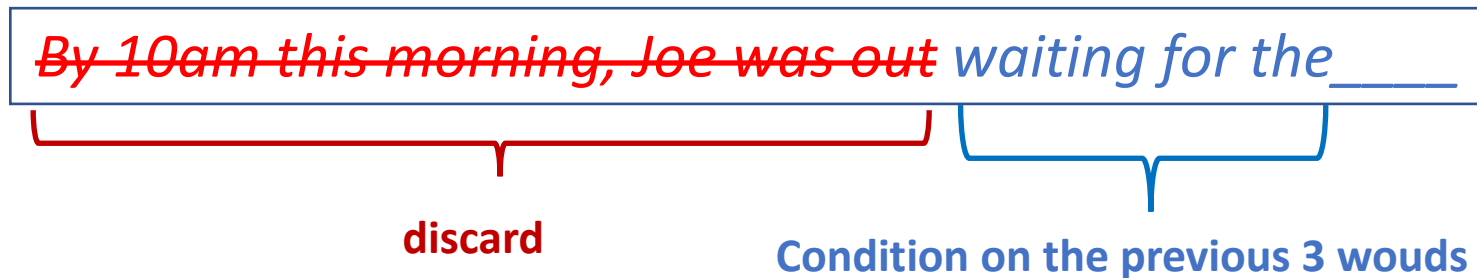
$$\approx \frac{\text{count}(w^{t+1}, w^t, \dots, w^{t-n+2})}{\text{count}(w^t, \dots, w^{t-n+2})}$$

- Sequence frequency of occurrence is a reasonable estimate of the probability

# n-gram Language Models: Example

---

Assume we are learning a *trigram* language model.



Example: given we have a corpus

- “waiting for the” occurred 1000 times
- “waiting for the bus” occurred 350 times
- “waiting for the car” occurred 100 times
- Then we get that
  - $P(\text{bus} \mid \text{waiting for the}) = 0.35$
  - $P(\text{car} \mid \text{waiting for the}) = 0.10$

$$P(\text{w} \mid \text{waiting for the}) = \frac{\text{count}(\text{waiting for the } \text{w})}{\text{count}(\text{waiting for the})}$$

# From n-gram probabilities to language model

---

- A language  $L \subseteq V^*$  is a (possibly infinite) set of strings over a (finite) vocabulary  $V$ .
- $P(w^t | w^{t-1})$  defines a distribution over all **2-grams** in  $V$ :
$$\forall w \in V : \sum_{w' \in V} P(w^t = w' | w^{t-1} = w) = 1$$
- By multiplying this distribution  $N$  times we get one distribution over all strings of the same length  $N$  ( $V^N$ ):
  - Probability of one  $N$ -word string:  $P(w_1 \dots w_N) = \prod_{i=1}^N P(w^{(i)} = w_i | w^{(i-1)} = w_{i-1})$
  - Probability of all  $N$ -word string:  $P(V^N) = \sum_{w, w' \in V} \prod_{i=1}^N P(w^{(i)} = w | w^{(i-1)} = w')$

# From n-gram probabilities to language model

---

A language model  $P(L) = P(V^*)$  should define one distribution  $P(V^*)$  that sums to one over *all* strings in  $L \subseteq V^*$ , *regardless of their length*:

$$P(L) = P(V^1) + P(V^2) + P(V^3) + \dots + P(V^n)$$

## Solution:

Add *end-of-Sentence (EOS)* or *beginning-of-sentence (BOS)* token to vocabulary  $V$ , assumptions;

- Each string ends with EOS (or BOS if focus is on start of string)
- EOS can only appear at the end of a string (or BOS at start of string)



# From n-gram probabilities to language model

---

In a trigram model

$$P(w^1 w^2 w^3) = P(w^1) P(w^2 | w^1) P(w^3 | w^2, w^1)$$

- The only trigram is  $P(w^3 | w^2, w^1)$
- $P(w^1)$  and  $P(w^2 | w^1)$  are not trigrams
- Add  $n-1$  BOS symbols to the sentence for an n-gram model

$$BOS_1 BOS_2 w^1 w^2, \dots, w^n$$

# From n-gram probabilities to language model

---

- A language model can be regarded as a *stochastic process*:
  - At each time step, randomly pick one token.
  - Stop when the word you picked is a special *EOS* token.
- Add an *EOS* to all sentences in your corpus, thus our vocabulary is now defined as

$$V^{EOS} = V \cup \{EOS\} \text{ or} \\ V^{BOS} = V \cup \{BOS\}$$

- With the new vocabulary we can get a single distribution over strings of any length because  $P(EOS | \dots)$  will be high enough that we are always guaranteed to stop after generating a finite number of words.

# Summary: training bi-gram model

---

- Replace all words not in training vocabulary by the unknown token *UNK*
- Enclose each sentence by special start and stop symbols

*< s > Alice was beginning to get very tired ....</s >*

- Define new vocabulary  $V' = \{V, UNK, < s >, </s >\}$
- Count the frequency of each bigram

$$C(< s >, Alice) = 1, C(Alice, was) = 1, \dots$$

- and normalize these frequencies to get probabilities

$$P(was|Alice) = \sum_{w_i \in V'} \frac{C(Alice\ was)}{C(Alice\ w_i)}$$

# Evaluating Language model

---

There are two ways to evaluate language models:

- **Intrinsic evaluation** measures how well the model captures what it is supposed to capture (e.g. probabilities)
  - **Extrinsic (task-based) evaluation** measures the language model improves the performance of the particular task.
- 
- Both cases require an evaluation metric that allows us to measure and compare the performance of different models

# Intrinsic evaluation

---

Intrinsic evaluation follows the normal procedure for every other machine learning model

- Define an evaluation metric (scoring function).
- Train the model on a training set
- Test the model on an unseen test set
- Compare language models by their scores.

# Intrinsic evaluation - perplexity

---

Perplexity is the probability of the test set normalized by the number of words:

$$PP(w_1 \dots w_N) =_{def} \sqrt[N]{\prod_{i=1}^N \frac{1}{P(w_i | w_{i-1}, \dots, w_{i-n+1})}}$$

with

$$PP(w_1 \dots w_N) =_{def} \exp \left( -\frac{1}{N} \sum_{i=1}^N \log P(w_i | w_{i-1}, \dots, w_{i-n+1}) \right)$$

## Practical issues

- Since language model probabilities are very small, multiplying them together often yields to arithmetic underflow.
- It is often better to use logarithms scale instead, so replace

# Intrinsic evaluation- perplexity

---

- **Perplexity:** most common intrinsic metric
  - based on Shannon principle
  - Perplexity is a bad approximation
    - Only use if the test data is very similar to the training set
    - Generally useful in pilot experiments
- A better model
  - Is one that assigns the higher probability to the word that occur

# Extrinsic evaluation

---

Best evaluation from comparing language model A and B

- Put A and B in a task
  - Spelling corrector, machine translator system, speech recognizer
- Run the task, get an accuracy when A is in the system
- Run the task again with B in the model, get the accuracy
  - How many misspelled words corrected properly
  - How many words were translated correctly
- Compare accuracy for A and B

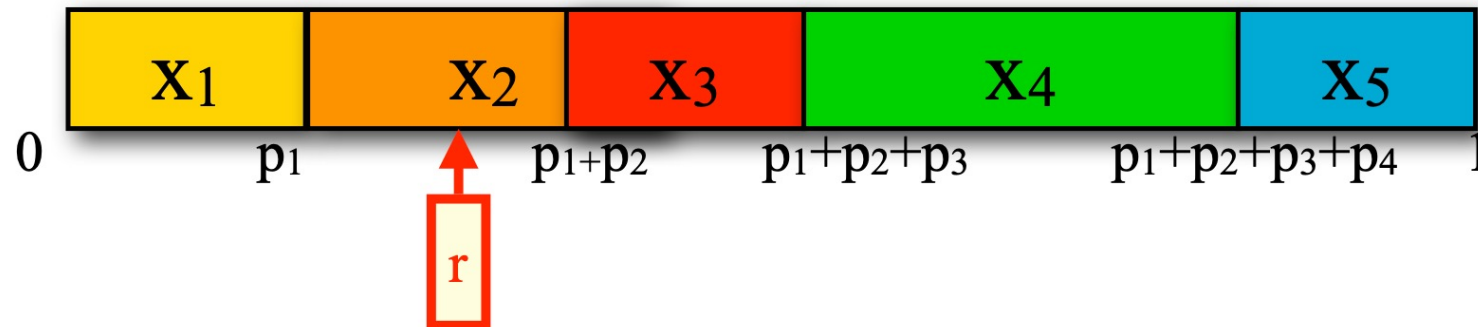


# Generating text with language models

---

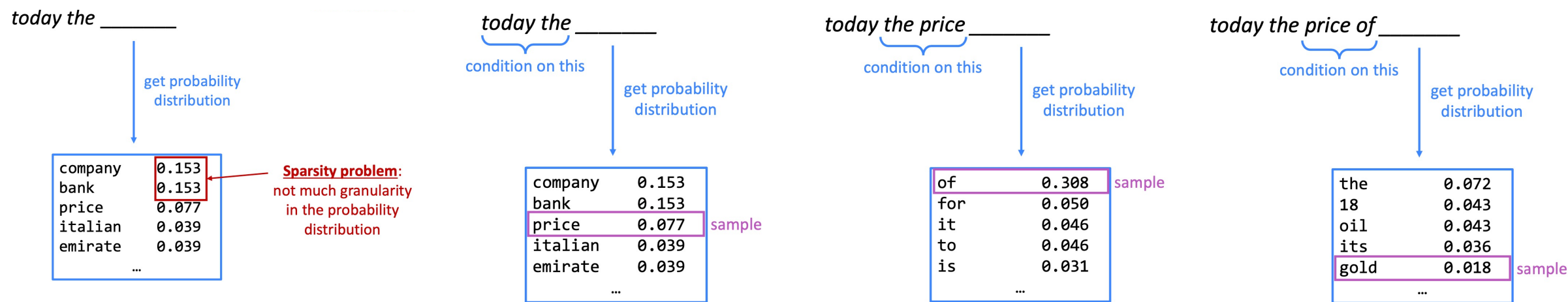
Given an  $n$ -gram language model with probability distribution  $P(X|Y = y)$

- Let  $X = \{x_1, \dots, x_N\}$  be  $N$  possible outcomes and  $P(X = x_i | Y = y) = p_i$
- Divide the interval  $[0, 1]$  into  $N$  intervals according to the probabilities of the outcomes
- Generate a random number  $r \in [0, 1]$
- Return the  $x_i$  whose interval the number  $r$  lies within



# Generating text with n-gram models

Text generation with a trigram model trained with 1.7 million words corpus (Reuters)



today the price of gold \_\_\_\_\_

*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 .*

Surprisingly grammatical!

# Limitations – number of possible parameters

---

Estimating the number of parameter per n-gram language model.

Given a vocabulary  $V$  of  $|V|$  unique tokens, where  $|V| = 10^4$

- **Unigram** model:  $|V|$  parameters  $\Leftrightarrow 10^4$  parameters
  - One distribution  $P(w^{(i)})$  with  $|V|$  outcomes, each  $w \in V$  is one outcome
- **Bigram** model:  $|V|^2$  parameters  $\Leftrightarrow 10^8$  parameters
  - $|V|$  distribution  $P(w^{(i)} | w^{(i-1)})$ , one distribution for each  $w \in V$  with  $|V|$  outcome each [each  $w \in V$  is one outcome]
- **Trigram** model:  $|V|^3$  parameters  $\Leftrightarrow 10^{12}$  parameters
  - $|V|$  distribution  $P(w^{(i)} | w^{(i-1)}, w^{(i-2)})$ , one distribution per bigram  $w'w''$  with  $|V|$  outcome each [each  $w \in V$  is one outcome]

# Limitations – sparsity problem

## Sparsity Problem 1:

- What if “*waiting for the  $w$* ” never occurred in the corpus?
- Then  $w$  has probability 0

## Solution - smoothing:

Add small  $\delta$  to the count for every  $w \in V$ .

$$P(w|\text{waiting for the}) = \frac{\text{count}(\text{waiting for the } w)}{\text{count}(\text{waiting for the})}$$

## Sparsity Problem 2:

- What if “*waiting for the*” never occurred in the corpus?
- Then we cannot calculate probability of any  $w$ .

## Solution - backoff:

condition on the  $(n-1)$ -gram, i.e. “*waiting for*”

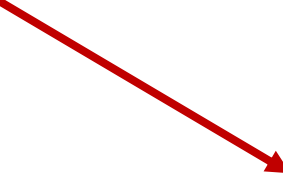
Larger  $n$  makes sparsity problem worse. Typically  $n$  should be less than or equal to 5

# Limitations – storage problems

---

## Storage Problem:

The need to store count for all n-grams you saw in the corpus.


$$P(\mathbf{w} | \text{waiting for the}) = \frac{\text{count}(\text{waiting for the } \mathbf{w})}{\text{count}(\text{waiting for the})}$$

Increasing  $n$  or increasing corpus size  $\Leftrightarrow$  increases model size