

## D08 基礎語言模型：N-Gram



>
重要知識點
>
什麼是語言模型？
>
語言模型機率
>
N-Gram模型
>

### NLP自然語言學習實戰馬拉松

#### ► Day8 - 基礎語言模型：N-Gram

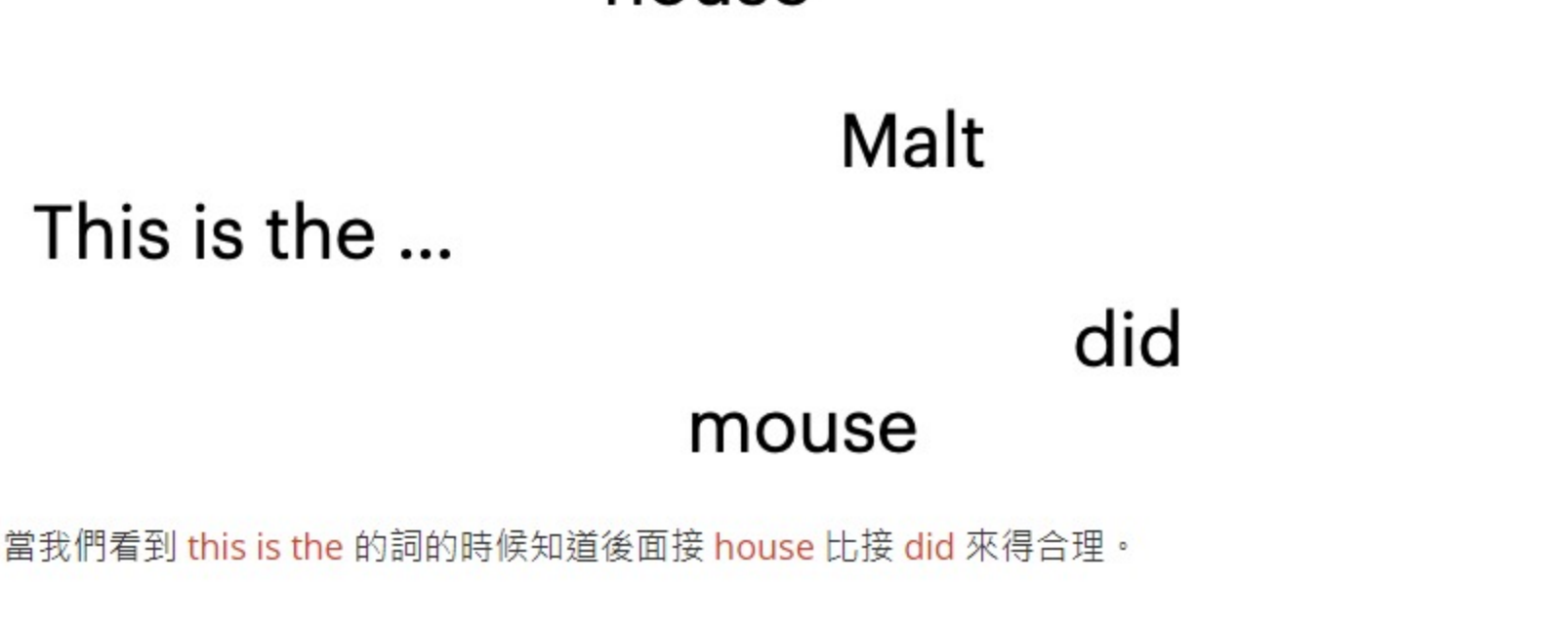
陪跑專家：Leo Liou 劉冠宏

#### 重要知識點

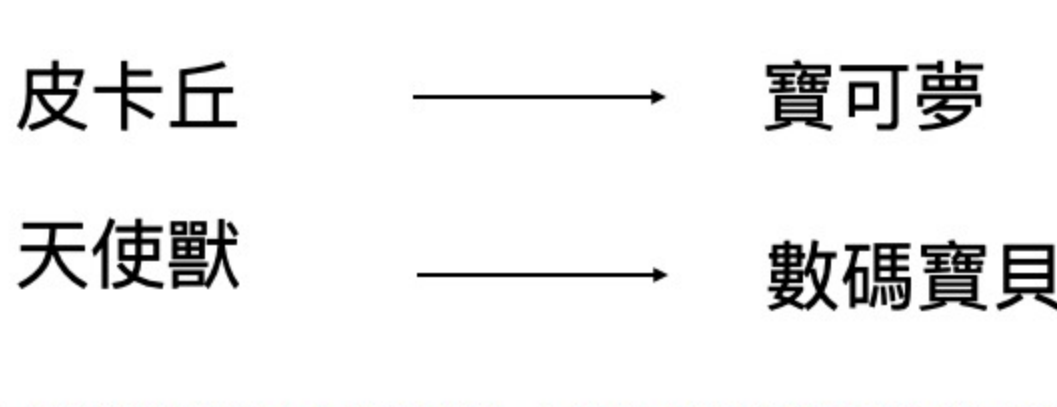


- 了解何謂語言模型
- 了解 N-Gram 語言模型的原理以及應用

#### 什麼是語言模型？



當我們看到 **this is the** 的詞的時候知道後面接 **house** 比接 **did** 來得合理。



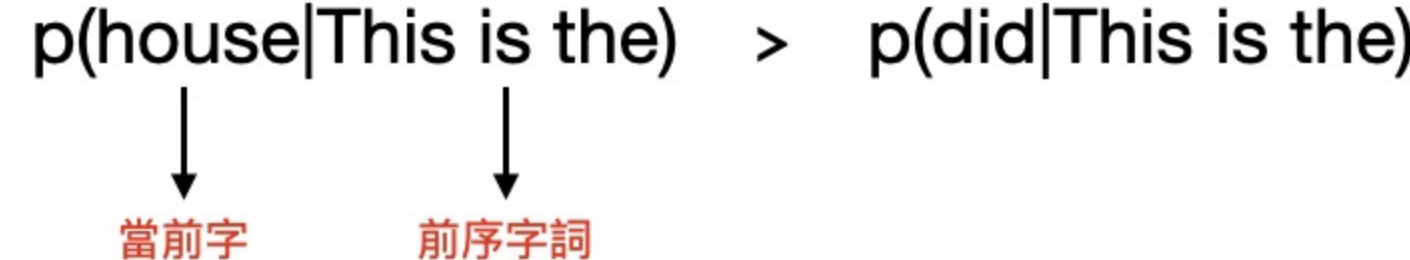
同理當我們想到可口可樂的時候，可能會聯想到百事可樂(而不是珍珠奶茶)，人們長年在學習累積語言後有判斷含不合理與聯想的能力。

那 NLP 有沒有這種判斷含不合理與聯想的能力呢？

N-gram 就是基於這種聯想的語言模型。

語言模型如何做到聯想或判斷句子的合理與不合理？

人類根據其他字詞(this is the)來判斷當前字(house)出現合理與否，而語言模型也是一樣的！根據其他字詞來判斷當前字出現的機率，因此便可以藉由機率大小的方式來判斷句子出現可能性。



#### 語言模型機率

corpus(文本資料)：

This is the house that Jack built.

This is the malt.

That lay in the house that Jack built.

This is the rat.

That ate the malt.

That lay in the house that Jack built.

This is the cat.

That killed the rat.

That ate the malt.

根據貝氏定理：

$$p(\text{house}|\text{This is the}) = \frac{\text{count}(\text{this is the house})}{\text{count}(\text{this is the...})}$$
$$= \frac{1}{4}$$

由上述的說明可以知道，若能得到一個句子出現的機率，如：P(This is a house)，可轉化為下列表示。

取的下列句子機率：

$$W = (W1W2W3W4,...,Wm)$$
$$P(W1,W2,W3,W4,...,Wm) = ?$$

根據鏈乘率：

$$P(W1,W2,W3,W4,...,Wm) = P(W1) \cdot P(W2|W1) \cdot P(W3|W1,W2) \cdot ... \cdot P(Wm|W1,...,Wm-1)$$

因此若想知道「我喜歡深度学习更喜歡NLP」含不合理(機率大小)可以藉由下面方法得到

$$P(\text{我喜歡深度学习更喜歡NLP}) = P(\text{我}) \cdot P(\text{喜}) | \text{我}) \cdot P(\text{歡}) | \text{我,喜}) \cdot ... \cdot P(\text{我喜歡深度学习更喜歡NLP})$$

為了簡化問題，我們可以假設當前字出現只跟前一個字有關

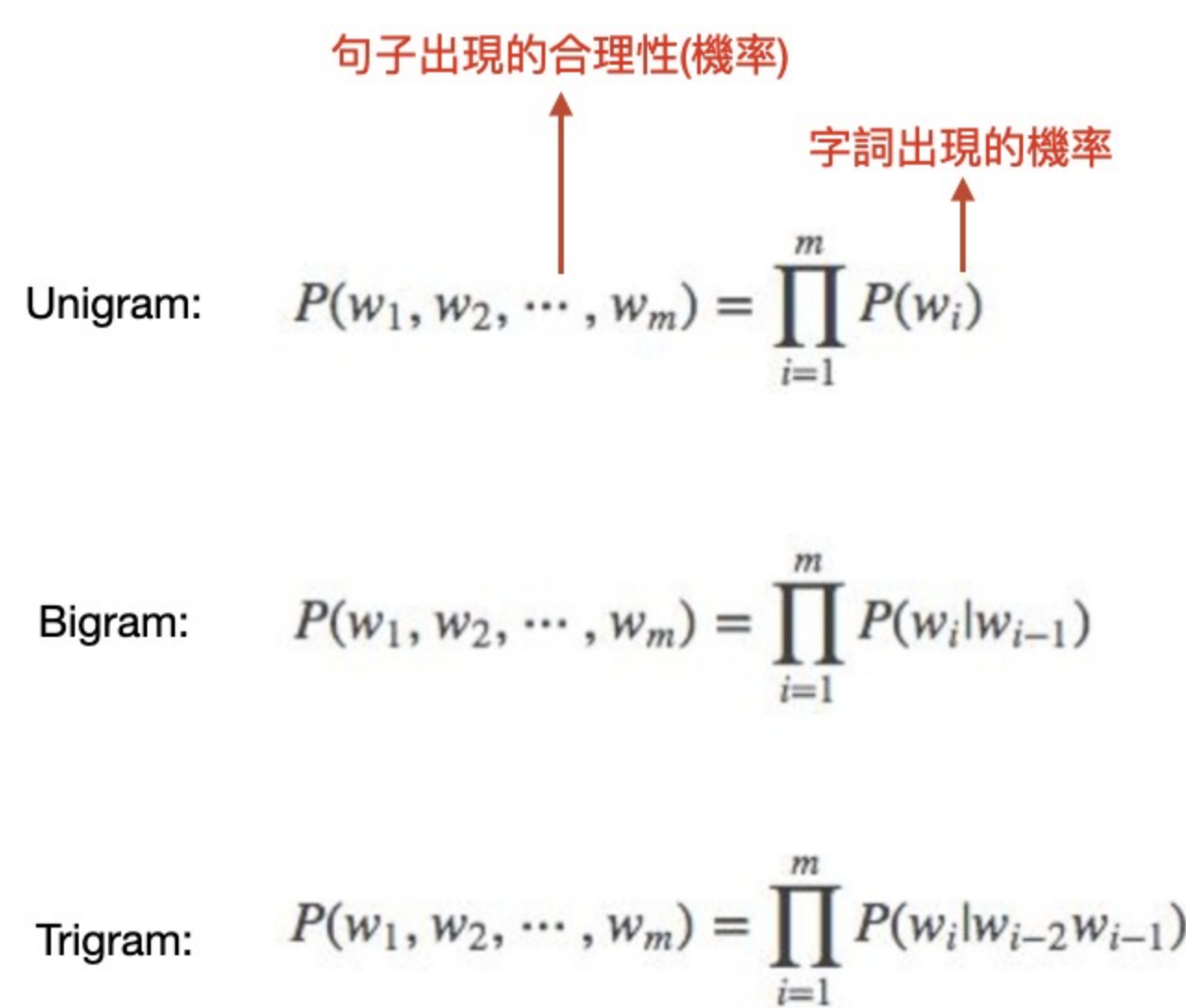
$$P(\text{我喜歡深度学习更喜歡NLP}) = P(\text{我}) \cdot P(\text{喜}) | \text{我}) \cdot P(\text{歡}) | \text{喜}) \cdot ... \cdot P(\text{我}) | \text{我})$$

這樣的假設算法即是 N-Gram 模型中 n=2 的情形，又稱作 Bigram

PS：試著思考為何需要做這樣的簡化呢？(詳細請見作業)

#### N-Gram模型

N-Gram 語言模型是基於統計的語言模型算法，主要是將文本中的內容取最靠近的 N 個字當作條件概率計算的先驗條件，形成長度是 N 的字詞片段序列，每個字詞片段及稱為 gram。常見的 N-Gram 模型有 Unigram(1-gram)、Bigram(2-gram)、Trigram(3-gram)。



$$\text{Unigram: } P(w_1, w_2, \dots, w_m) = \prod_{i=1}^m P(w_i)$$

$$\text{Bigram: } P(w_1, w_2, \dots, w_m) = \prod_{i=1}^m P(w_i|w_{i-1})$$

$$\text{Trigram: } P(w_1, w_2, \dots, w_m) = \prod_{i=1}^m P(w_i|w_{i-2}w_{i-1})$$

#### Bigram語言模型

我們知道了 Bigram 語言模型，現在我們來看看 Bigram 模型如何判斷句子的合理性，為了考慮字詞在開頭與結尾的價值，可以加入 start 的 token 與 end 的 token 在文本中。

Bigram (n=2)：

$$P(W1,W2,W3,W4,...,Wm) = P(W1|<start>) \cdot P(W2|W1) \cdot P(W3|W2) \cdot ... \cdot P(Wm|Wm-1) \cdot P(<end>|Wm)$$

corpus：

That is the house that Jack built.

This is the malt.

It is a fancy house.

$$\text{Probabilities:}$$
$$P(\text{this is the house}) = P(\text{this}<start>) \cdot P(\text{is}|\text{this}) \cdot P(\text{the}|\text{is}) \cdot P(\text{house}|\text{the}) \cdot P(<end>|\text{house})$$
$$= (1/3) \cdot (1) \cdot (2/3) \cdot (1/2) \cdot (1/2)$$

#### N-Gram應用場景

語言模型可以應用的場景很多，N-Gram 常用的應用場景像是「搜尋推薦」、「關鍵字正」、「分類系統」等。

搜尋推薦 (根據部分輸入，推薦接下來可能的字詞)



分類系統(目前分類算法通常以自然語言課程中介紹的為大宗)

$$P(Y_1) = p(\text{我})p(\text{喜欢}|\text{我})p(\text{自然语言处理}|\text{喜欢})$$

$$P(Y_2) = p(\text{我})p(\text{喜}|\text{我})p(\text{欢}|\text{喜})p(\text{自然}|\text{欢})p(\text{语言}|\text{自然})p(\text{处理}|\text{语言})$$

$$P(Y_3) = p(\text{我})p(\text{喜}|\text{我})p(\text{欢}|\text{喜})p(\text{喜}|\text{自})p(\text{然语言}|\text{自})p(\text{处理}|\text{然语言})$$

關鍵字正

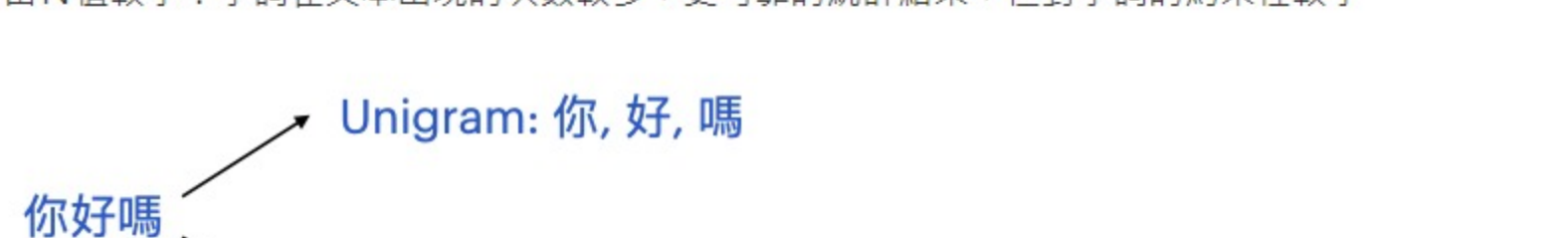
今天天氣很好，適合出門玩

$$P(\text{天}|\text{今}) \cdot P(\text{天}|\text{天}) \cdot P(\text{氣}|\text{天}) \cdot P(\text{很}|\text{好}) < \text{threshold}$$

#### N 的大小

當N值較大，對字詞的約束性更高，具有更高的辨識力，複雜度較高

當N值較小，字詞在文本出現的次數較多，更可靠的統計結果，但對字詞的約束性較小



※ Bigram 的可能的 gram 數較多

#### 知識點回顧

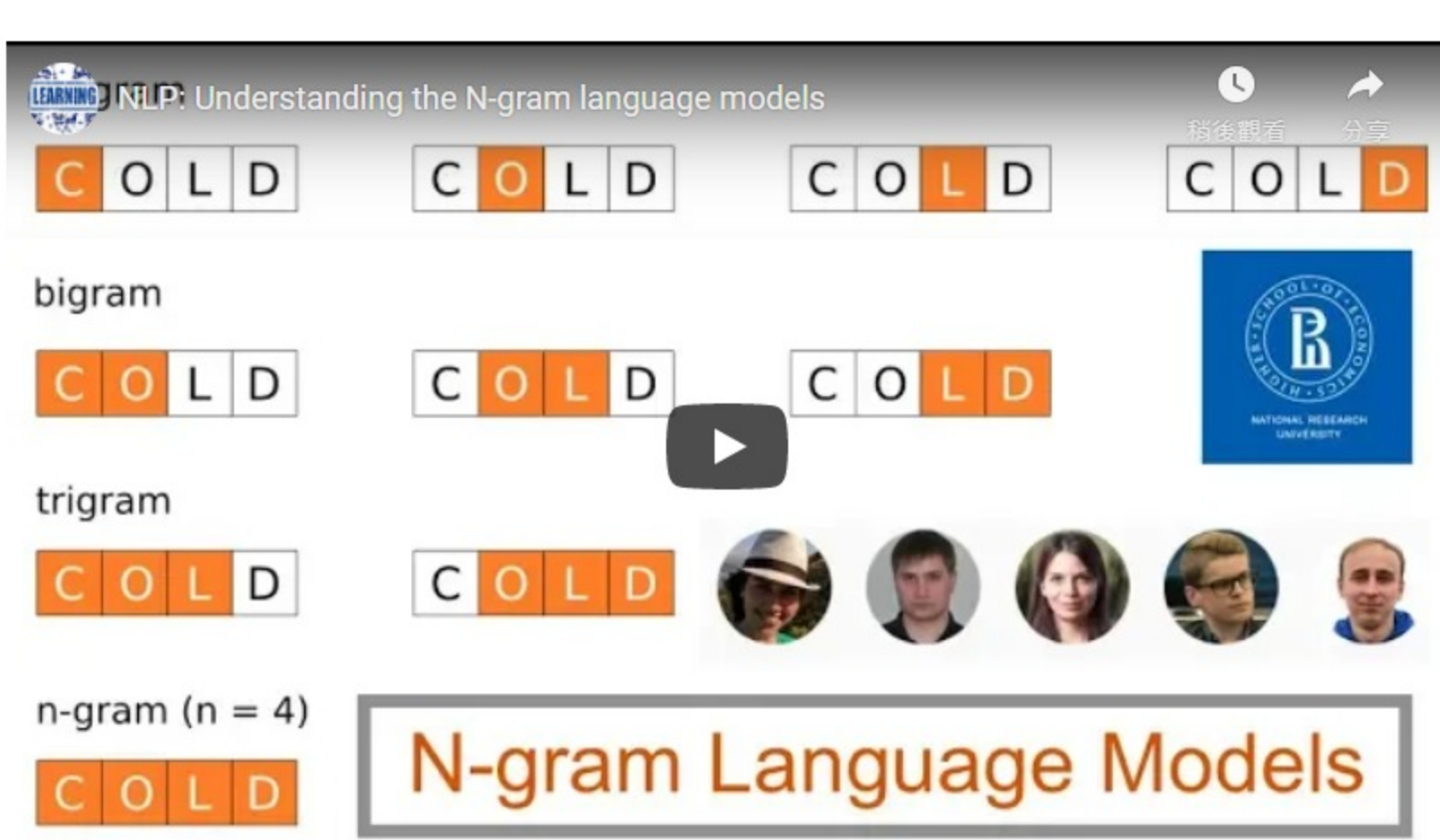
在這章給我們學理到了

- 何謂語言模型
- 了解 N-Gram 語言模型的原理以及應用

#### 延伸閱讀

網站：[基礎語言模型介紹](#)

基礎語言介紹 YouTube 影片



網站：[語言模型推薦](#)

介紹常見的語言模型應用



網站：[N-Gram語言模型補充介紹](#)

文中提到語言模型機率計算的挑戰

#### INTRODUCTION

### Language Models: N-Gram

#### A step into statistical language modeling

Shashank Kapadia

Mar 26, 2019 · 4 min read

#### Introduction

Statistical language models, in its essence, are the type of models that assign probabilities to the sequences of words. In this article, we'll understand the simplest model that assigns probabilities to sentences and sequences of words, the **n-gram**

You can think of an N-gram as the sequence of N words, by that notion, a 2-gram (or bigram) is a two-word sequence of words like "please turn", "turn you", or "your homework", and a 3-gram (or trigram) is a three-word sequence of words like "please turn you", or "turn your homework"

#### Intuitive Formulation

Let's start with equation  $P(w|h)$ , the probability of word  $w$ , given some history,  $h$ . For example,

$$P(\text{the} | \text{its water is so transparent that})$$

網站：[N-Gram中文練習檢查](#)

使用 N-gram 實踐中文打字機查



#### 中文拼寫糾正

最基本的思想，將所有的常見錯別字替換為字庫。

但是這個字庫的數據實際上非常有限，所以這是動態糾正算法。

本文主要介紹如何利用 N-gram 模型結合字詞在中文語料中的頻率，然後介紹如何根據 N-gram 模型在中文語料中計算 N-gram 模型的機率。

最後從如何輸入字詞解如何中文文本長度限制 (基於 N-gram)，並從如何方法中推導出一種公式，利用 N-gram 模型的特點結合 N-gram 模型及 N-gram 模型的機率。

#### n-gram模型

在中文語料字庫模型中，我們判斷一個字詞是否合適可以通過計算它的機率來得到，假設一個字詞  $s = [w_1, w_2, \dots, w_n]$ ，則問題可以轉換成如下形式：

$$p(s) = p(w_1, w_2, \dots, w_n) = p(w_1) \cdot p(w_2|w_1) \cdot p(w_3|w_1, w_2) \cdot \dots \cdot p(w_n|w_1, w_2, \dots, w_{n-1})$$

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

假設我們使用 N-gram 模型，則問題計算一個字詞在語料中的機率。

&lt;