# 2019 0416 第九週 蒐集詞彙語意研究資料 ——網路爬蟲

## Pepper 有幾種意思?

salt and pepper roast red pepper in the oven

會先想到:查劍橋英中辭典

https://dictionary.cambridge.org/us/
dictionary/english-chinese-traditional/pepper

#### 詞,原形,詞意,定義



#### 詞,原形,詞意,定義

#### 詞意 2

#### pepper noun (VEGETABLE)

ALSO sweet pepper, US ALSO bell pepper [C] a vegetable that is usually green, red, or yellow, has a rounded shape, and is hollow with seeds in the middle

甜椒,燈籠椒

a red/green pepper

#### 紅/青椒

Peppers are usually cooked with other vegetables or eaten raw in salads.

甜椒通常與其他蔬菜一同烹調,或放在沙拉裡 生吃。

Red peppers are ideal for roasting in the oven.

紅椒很適合用烤箱烘烤來吃。

#### 定義

翻譯

等級

例句

#### HTML 基本概念 <a href="https://en.wikipedia.org/wiki/HTML">https://en.wikipedia.org/wiki/HTML</a>

- Texts + nested elements
  - start tag, end tag
- Start tag may contain
  - attributes
    - value (one value (sometimes more) per attribute

```
<tag attribute1="value1" attribute2="value2">
<tag attribute1="value1 value2">
```

- Use class attribute to classify similar elements (semantics)
- Example

```
Paragraph 1 Paragraph 2
<a href="https://www.wikipedia.org/">A link to Wikipedia!</a>
<input type="text" /> <!-- This is for text input -->
<input type="file" /> <!-- This is for uploading files -->
<input type="checkbox" /> <!-- This is for checkboxes -->
```

## 相對網頁的格式(1)

### headword

<span class="headword"><span class="hw">pepper</span>

pepper

pepper noun (POWDER)

noun ⋅ UK (1) /'pep.ər/ US (1) /'pep.ə-/



### headword + pos + guideword

<div class="txt-block txt-block--alt2"><span class="hw">pepper</span> <span class="pos" >noun/ span>

<span class="guideword" >(<span>POWDER</span>)</span></span>





### epp-xref: <span class="epp-xref A2">A2</span>

### gram: <span class="gram">[<span class="gcs"> <span class="gc">U</span></span>]</span>

### def + trans

<b class="def">a grey or white powder produced by crushing dry peppercorns, used to give a spicy, hot

A2 [U] a grey or white powder produced by crushing dry peppercorns, used to give a spicy, hot taste to food

胡椒粉

taste to food.<span class="trans" lang="zh-Hant">胡椒粉</span>

### eg and trans

#1: <span class="eg">freshly ground black pepper <span class="trans" lang="zh-Hant">現磨的黑胡椒粉</span>

#2: <span class="eg"><span class="b">salt and pepper</span> <span class="trans" lang="zh-Hant">鹽和胡椒</span>

freshly ground black pepper

現磨的黑胡椒粉

salt and pepper

鹽和胡椒

### 相對網頁的格式(2)

### headword + pos + guideword

pepper noun (VEGETABLE)

<div class="txt-block txt-block--alt2"><span class="hw">pepper</span> <span class="pos">noun</span>
<span class="guideword"> (<span>VEGETABLE</span>)</span> </span>

### epp-xref <span class="epp-xref B1">B1</span>





ALSO sweet pepper, US ALSO bell pepper

### usage+v <span class="usage">also</span></span> <span class="v">sweet pepper</span>, </span><span class="region">US</span> <span class="v">bell pepper</span></span>

### gram <span class="gram">[C]</span>

[C]

vegetable that is usually green, red, or yellow, has a rounded shape, and is hollow with seeds in the middle

### def + trans

甜椒,燈籠椒

<b class="def"> a vegetable that is usually green, red, or yellow, has a rounded shape, and is hollow with seeds in the middle <span class="trans" lang="zh-Hant"> 甜椒,燈籠椒</span>

a red/green pepper

a reargreen peppe

pper P

Peppers are usually cooked with other vegetables or eaten raw in salads.

紅/青椒

甜椒通常與其他蔬菜一同烹調,或放在沙拉裡 生吃。

Red peppers are ideal for roasting in the oven.

紅椒很適合用烤箱烘烤來吃。

### eg + trans

#1: <span class="eg">a red/green pepper</span> <span class="trans" lang="zh-Hant">紅/青椒</span>

#2: <span class="eg">Peppers are usually cooked with other vegetables or eaten raw in salads. <span class="trans" lang="zh-Hant">甜椒通常與其他蔬菜一同烹調,或放在沙拉裡生吃。</span>

#3: <span class="eg">Red peppers are ideal for roasting in the oven.</span> <span class="trans" lang="zh-Hant">紅椒很適合用烤箱烘烤來吃。</span>

#### 相關模組

- 下載網頁模組
  - import urllib
  - o page = urllib.request.urlopen(<url>)
- 分析網頁模組
  - from bs4 import BeautifulSoup
  - soup = BeautifulSoup(page\_source, 'html.parser')
  - for div in soup.find\_all([<tag>, ...])
  - check div[<attribute>]

#### 用下載資料研究詞彙語意解岐

- 利用詞典英文、中文定義詞、引導分類詞,做語意解岐
  - 輸入:句子 (問題詞+文脈詞) presidential chopper
  - 輸出:詞典的語意條目 AIRCRAFT; TOOL; MOTORCYCLE
  - chopper

noun · UK (1) /'tfpp.ər/ US (1) /'tfar.pa/

#### chopper noun (AIRCRAFT)

[○] INFORMAL FOR helicopter
直升機 (helicopter的非正式說法)

#### chopper noun (TOOL)

[○] a small axe held in one hand
小斧頭

#### chopper noun (MOTORCYCLE)

- ② a type of motorcycle with tall handlebars (= parts that you hold) and a long front end (高重把、長前叉的)美式機車
- ② a type of bicycle that is designed to look similar to a chopper motorcycle 設計得看上去像是美式機車的自行車

## 研究議題 (1)

- 利用詞典英文、中文定義詞、引導分類詞,做語意解岐
  - 輸入: presidential chopper
  - 輸出:AIRCRAFT; TOOL; MOTORCYCLE
  - 建立關係: <文脈詞, 語意定義詞+引導分類詞>
    - (presidential, #1) —> (presidential, AIRCRAFT)
       (presidential, #2) —> (presidential, TOOL)
       (presidential, #3) —> (presidential, MOTORCYCLE)
  - 研究方法
    - Bayesian models (supervised, unsupervised)
    - EM models
    - Graph models

### 研究議題 (2)

- 自動調整語意條目的引導詞分類
  - 問題:引導詞分類不一致 -> 解答:自動化的多重分類

saffron#1 = SPICE —> #spice, #food

chopper#1 = AIRCRAFT —> #aircraft

- 問題:非歧義詞缺乏引導詞分類 --> 解答:自動分類

aircraft#1 = ? —> #aircraft

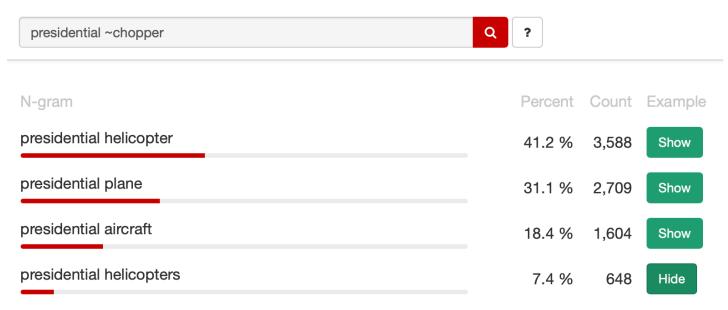
### 參考文獻

- Bayesian models (unsupervised): Yarowsky (1992)
  - Word-sense disambiguation using statistical models of Roget's categories trained on large corpora
- EM-like models: Yarowsky (1995)
  - Unsupervised word sense disambiguation rivaling supervised methods
- Graph models: Agirre et al. (2014)
  - Random walks for knowledge-based word sense disambiguation
- Topic models: Chaplot and Salakhutdinov (2018)
  - Knowledge-Based Word Sense Disambiguation Using Topic Models
- Neural models: Yuan et al. (2016)
  - Semi-supervised word sense disambiguation with neural models

## 研究議題 (cont.)

利用 word embeddings

#### linggle1012



- A six-year-old project to build state-of-the-art presidential helicopters has bogged down in a contracting quagmire that will challenge Mr. Obama 's desire to rein in military contracting expenses.
- For Mr. Obama, the program is one more inheritance from the Bush administration, which began the effort after the Sept. 11 attacks generated concern about whether presidential helicopters from the 1970s were up to the challenge of terrorist threats.