# UNIT 12

## 語言搜尋引擎

December 9, 2014

## 目標

- 建立一搜尋引擎用於搜尋英文詞語用法。
- 可輔助英語學習與文章寫作。

### 搜尋例子

- `adj. beach`: 即代表搜尋 beach 前面出現過的形容詞。
- `play * role`: 搜尋 play 與 role 中間最常出現的字詞組合。
- `go ?to home`: go 與 home 之間是否要放 to。
- `go * movie`: go 與 role 中間最常出現的字詞組合。
- `kill the _`: 最常被 kill 的東西是。

# 用 Google 查英文



Google | "a * beach" | 🔍

網頁　圖片　影片　新聞　地圖　更多▾　搜尋工具

約有 1,720,000,000 項結果 (搜尋時間：0.80 秒)

**Beach Vacations - TripAdvisor**
www.tripadvisor.com/Inspiration-g1-c1-World.html ▾ 翻譯這個網頁
Koh Samui's myriad beaches present everything you could ever dream of in a tropical beach. Picture-perfect sands, coconut trees and palm fronds adorn each ...

**Hyatt Siesta Key Beach Resort, A Hyatt Residence Club ...**
www.tripadvisor.com › ... › Siesta Key › Siesta Key Hotels ▾ 翻譯這個網頁
★★★★★ 評分：5 - 221 則評論 - 消費：$
DaydreaminginDenver has 1 more review of Hyatt Siesta Key Beach Resort, A Hyatt Residence Club. "Absolute luxury on a beautiful beach". 5 of 5 stars ...

**Beach Vacations and Family Beach Resorts | Travel + Leisure**
www.travelandleisure.com › Trip Ideas ▾ 翻譯這個網頁
ooking for beach vacations? Travel + Leisure features top beach family resorts in popular destinations, plus insider tips and guides. Search for beach vacation ...

**Cheap hotels Myrtle Beach - Hotels.com**
www.hotels.com › ... › Myrtle Beach Hotels ▾ 翻譯這個網頁
A Myrtle Beach travel guide – championship golf and 60 miles of beaches. Take a vacation in Myrtle Beach and you've entered one of the South's great ...

**Myrtle Beach Hotels - Hotels.com**
www.hotels.com › ... › South Carolina hotels ▾ 翻譯這個網頁
A Myrtle Beach travel guide – championship golf and 60 miles of beaches. Take a vacation in Myrtle Beach and you've entered one of the South's great ...

**City of Huntington Beach, CA - Beach Wedding Information**
www.huntingtonbeachca.gov › residents › beach info ▾ 翻譯這個網頁
A: The beach is open until 10:00 pm every night. Most of the ... A: Beach Operations at

# 用 Google 查英文

# 語法設計

| 語法 | 說明 |
|---|---|
| _ | 單一任意字詞 |
| * | 零到多個任意字詞 |
| ?term | term 可有可無 |
| term1 \| term2 | term1 或 term2 |
| adj. det. n. v. prep | 形容詞、冠詞、名詞、動詞、介繫詞 |

## 搜尋例子

- `adj. beach`: 即代表搜尋 beach 前面出現過的形容詞。
- `play * role`: 搜尋 play 與 role 中間最常出現的字詞組合。
- `go ?to home`: go 與 home 之間是否要放 to。
- `go * movie`: go 與 role 中間最常出現的字詞組合。
- `kill the _`: 最常被 kill 的東西是。

# Lab 12

- 目標：完成語法第一項 _
  - 任意位置置入 _
  - 最長 4-gram

## Query 範例

- `play _ _ role`
- `kill the _`
- `a _ beach`

- 輸入資料：citeseerx 的許多句子
- 輸出結果：
  - key: 所有會有結果的 query
  - value: 符合 query 的前 100 名 ngram 與 count。

# Lab 12 - 輸出

- key: 所有會有結果的 query
- value: 符合 query 的前 100 名 ngram 與 count。

## 輸出範例

| Key | Ngrams | Counts |
|---|---|---|
| a _ beach | a sandy beach | 486 |
| | a private beach | 416 |
| | a beautiful beach | 314 |
| | a small beach | 175 |
| | ... | |
| kill the _ | kill the people | 189 |
| | kill the other | 174 |
| | kill the process | 163 |
| | kill the enemy | 160 |
| | ... | |

# 隨堂測驗

目標

- 依 MapReduce 架構，設計每階段 mapper, reduce 的輸入輸出來完成 Lab 12
- 在紙寫撰寫簡單輸入、輸出的 key-value 範例表達概念即可

小提示

- 可有 1 至多個 map, reduce 流程
- 考慮 mapper 的輸入資料切割影響
- mapper 輸入為 value 或 key-value，輸出為 key-value
- reducer 輸入為 grouped key-values，輸出為 key-value

# Bi-gram Count

## Bi-gram Count Mapper 範例

| Input(value) | Output(key => value) |
|---|---|
| C D C D | C D => 2 |
| | D C => 1 |
| B C D A | B C => 1 |
| | C D => 1 |
| | D A => 1 |
| C D A B | C D => 1 |
| | D A => 1 |
| | A B => 1 |

## Reducer 範例

| Input(key => value) | Output(key => value) |
|---|---|
| A B => 1 | A B => 1 |
| B C => 1 | B C => 1 |
| C D => 2 | C D => 4 |
| C D => 1 | |
| C D => 1 | |
| D A => 1 | D A => 2 |
| D A => 1 | |
| D C => 2 | C C => 2 |

# Lab12

## Lab12 Mapper 範例

| Input(value) | Output(key => value) |
|---|---|
| A B C 200 | A B C => A B C 200 |
|  | _ B C => A B C 200 |
|  | A _ C => A B C 200 |
|  | A B _ => A B C 200 |
|  | _ _ C => A B C 200 |
|  | _ B _ => A B C 200 |
|  | A _ _ => A B C 200 |
| A D C 300 | _ D C => A D C 300 |
|  | A _ C => A D C 300 |
|  | ... |
| A E C 100 | _ E C => A E C 100 |
|  | A _ C => A E C 100 |
|  | ... |

# Lab12

## Lab12 Reducer 範例

| Input(value) | Output(key => value) |
|---|---|
| A _ C => A B C 200 | A _ C =>  A D C 300, |
| A _ C => A D C 300 |            A B C 200, |
| A _ C => A E C 100 |            A E C 100 |
| A B _ => A B C 200 | A B _ =>  A B C 200 |
| A D _ => A D C 300 | A D _ =>  A D C 300 |
| A E _ => A E C 100 | A E _ =>  A E C 100 |
| A _ _ => A B C 200 | A _ _ =>  A D C 300, |
| A _ _ => A D C 300 |            A B C 200, |
| A _ _ => A E C 100 |            A E C 100 |
| _ B C => A B C 200 | _ B C =>  A B C 200 |
| _ D C => A D C 300 | _ D C =>  A D C 300 |
| _ E C => A E C 100 | _ E C =>  A E C 100 |
| … | … |

需完成六支程式

- 產生 ngram count 的 mapper, reducer
- 產生 query result 的 mapper, reducer
- 將 query result 轉為 database
  (試試 python 內建的 shelve 或 sqlite3 套件)
- Database 介面程式，讓使用者輸入 query ， 即時取得
  result

# python shelve

```python
import shelve
d = shelve.open('data.shelve')
d['odds'] = [1, 3, 5, 7, 9]
print d['odds']
d['evens'] = [2, 4, 6, 8, 10]
d['hello'] = 'world'
del d['hello']
d['zipcodes'] = {'hsinchu': 300, 'zhongli': 320}
print d.keys()
d.close()
```

Google "python shelve" for official documents

```python
#!/usr/bin/env python
# -*- coding: utf-8 -*-


def ngrams(words):
    for length in range(1, 5 + 1):
        for ngram in zip(*(words[i:] for i in range(length))):
            yield ngram


def mapper(files):
    import fileinput
    from nltk.tokenize import word_tokenize
    from collections import Counter
    ngram_counter = Counter()
    for line in fileinput.input(files):
        line = line.decode('iso-8859-1')
```

# Ngram Count II

```
18            words = word_tokenize(line.lower())
19            ngram_counter.update(ngrams(words))
20
21        for ngram, count in ngram_counter.iteritems():
22            print (u' '.join(ngram) + u'\t' + unicode(count)).encode('utf-8')
23
24
25    def line_to_ngram(line):
26        line = line.decode('iso-8859-1')
27        return line.split(u'\t', 1)[0]
28
29
30    def line_to_count(line):
31        line = line.decode('iso-8859-1')
32        return int(line.split(u'\t', 1)[1])
33
34
35    def reducer(files):
36        import fileinput
```

## Ngram Count III

```
37        from itertools import groupby, imap
38
39        for ngram, lines in groupby(fileinput.input(files), key=line_to_ngram):
40            count = sum(imap(line_to_count, lines))
41            print (ngram + u'\t' + unicode(count)).encode('utf-8')
42
43
44   if __name__ == '__main__':
45        import argparse
46        import sys
47        parser = argparse.ArgumentParser(description='N-gram counter')
48        parser.add_argument(
49            '-r', '--reducer', action='store_true', help='reducer mode')
50        parser.add_argument(
51            '-m', '--mapper', action='store_true', help='mapper mode')
52        parser.add_argument('files', metavar='FILE', type=str, nargs='*',
53                            help='input files')
54
55        args = parser.parse_args()
```

# Ngram Count IV

```
56
57    if (args.mapper and args.reducer
58            or
59            not args.mapper and not args.reducer):
60        parser.print_help()
61        sys.exit(1)
62
63    if args.mapper:
64        mapper(args.files)
65    elif args.reducer:
66        reducer(args.files)
```