

In [1]:

```
from gensim.models import FastText
from gensim.models.word2vec import LineSentence
import multiprocessing

WIKI_SEG_TXT = "wiki_seg.txt"

max_cpu_counts = multiprocessing.cpu_count()
word_dim_size = 300

sentences = LineSentence(WIKI_SEG_TXT)

model = FastText(sentences, vector_size=word_dim_size, workers=max_cpu_counts)

output_model = f"fasttext.zh.{word_dim_size}.model"
model.save(output_model)

print(model.wv.vectors.shape)
model.wv.vectors

model = FastText.load('fasttext.zh.300.model')
vocab = model.wv.key_to_index

print(f"總共收錄了 {len(vocab)} 個詞彙")

print("印出 20 個收錄詞彙:")
print(list(vocab.keys())[:20])
```

(1281108, 300)

總共收錄了 1281108 個詞彙

印出 20 個收錄詞彙:

['年', '月', '日', '中', '10', '12', '11', '小行星', '中國', '時', '-', '日', '本', '美國', '20', '香港', '臺灣', '15', '位於', '30', '站']

In [2]:

```
vec = model.wv['數學家']  
print(vec.shape)  
vec
```

(300,)

Out[2]:

In [3]:

```
word = "這肯定沒見過 "  
  
# 若強行取值會報錯  
try:  
    vec = model.wv[word]  
except KeyError as e:  
    print(e)
```

In [4]:

```
model.wv.most_similar("飲料", topn=10)
```

Out[4]:

```
[('輝劍', 0.9708890914916992),  
 ('名松', 0.9569700360298157),  
 ('飲料類', 0.9305379986763),  
 ('飲料機', 0.9210789799690247),  
 ('飲料罐', 0.9001978635787964),  
 ('軟飲料', 0.887269139289856),  
 ('經米濱', 0.8780335783958435),  
 ('茶飲料', 0.8691356182098389),  
 ('飲品', 0.846047580242157),  
 ('廣慈宮', 0.797712504863739)]
```

4/5

```

-7.71618187e-01, 2.67051649e+00, 1.42840219e+00, 5.49594223e-01,
In [8]: -9.26083922e-01, -7.49866486e-01, -2.26152205e+00, 2.78727508e+00,
model.wv.most_similar("合約")
-9.70590472e-01, -1.26464534e+00, -5.52120388e-01, 6.20141923e-01,
-1.35305405e+00, -1.21958661e+00, 6.87196255e-01, 1.08689809e+00,
Out[8]: 1.08027685e+00, -6.81514263e-01, -3.49934429e-01, 1.01837456e+00,
1.27114236e+00, -8.56064260e-01, -1.88862312e+00, -3.79661947e-01,
[('德康', 9.71523624e-01, 9.911900368e-01, 2.21391845e+00, -5.95065236e-01,
('合同', 1.43144258e+00, 6.0382159640827e-01, -1.55585563e+00, 2.02197289e+00,
('綠蠅', 8.80503496e-01, 5.5822801343784e+00, 8.57517242e-01, -1.27263620e-01,
('合同額', 7.91467817e-01, 9.451733804111e+00, -5.62168479e-01, -2.58745623e+00,
('合同期', 2.02166010e-01, 5.39823455722e-01, 1.19647527e+00, 9.45923388e-01,
('合同商', 2.13526511e+00, 6.838932395e-01, -1.44336033e+00, -9.26598251e-01,
('籤合同', 2.32770368e-01, 7.338562919256e-01, -3.14604831e+00, -6.57505989e-01,
('合同制', 1.75698948e-01, 7.986431128591e-01, -9.00605321e-03, 1.76274085e+0
0], 簽約', 0.6969977617263794),
('合同條款', 0.69920885443687439)]

```

In [9]:

```
model.wv.similarity("連結", "鏈接")
```

Out[9]:

0.42695913

In [10]:

```
model.wv.similarity("連結", "陰天")
```

Out[10]:

-0.03664172

In [11]:

```
print(f"Loading {output_model}...")
new_model = FastText.load(output_model)
```

Loading fasttext.zh.300.model...

In [12]:

```
model.wv.similarity("連結", "陰天") == new_model.wv.similarity("連結", "陰天")
```

Out[12]:

True

In [ ]: