# Running Word2Vec with Chinese/English Wikipedia Dump

LUO Linkai

Deep Learning Research & Application Center

March 6, 2017

# Outline

# Word2Vec Revisited

# What is Word2Vec

- The corpus: "The cat sat on the mat"
- One-hot representation: "sat" = $[0, 0, 0, 1]^T$
- Distributed representation: "sat" = $[0.01, -0.02, 0.03]^T$
- Word similarity: "sat" $\approx$ "sit"
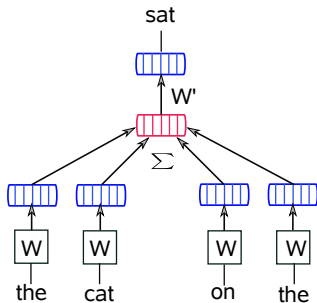- Word2Vec is the technique to obtain such meaningful word vectors given corpus

# Idea

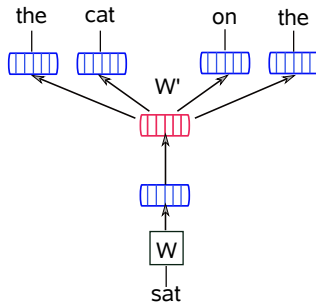<start> The cat sat on the mat <end>

<start> The cat  sat  on the  mat <end>

<start> The  cat sat on the  mat <end>

- CBOW: $p(sat|[the, cat, on, the])$  (treats an entire context as one observation)
- Skip-gram: $p(the|sat)\ p(cat|sat)\ p(on|sat)\ p(the|sat)$ (treats each context-target pair as a new observation)

# Word2Vec Algorithms
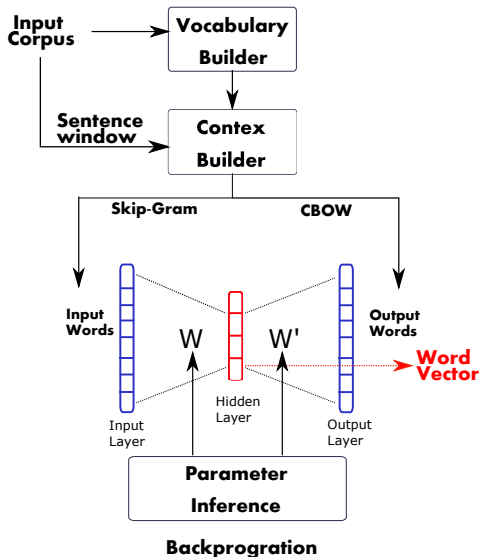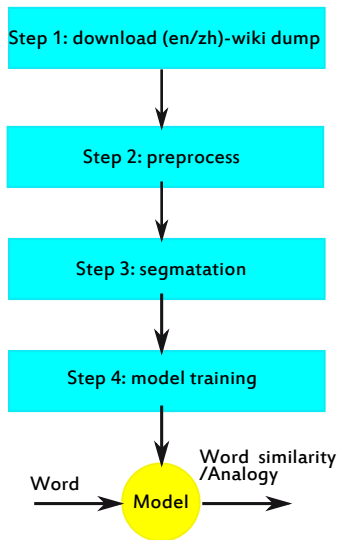


CBOW                          Skip-gram

# Word2Vec Decomposed

# Steps of Wikipedia Word2Vec

English Wikipedia Word2Vec

# Download English Wikipedia Dump

- https://dumps.wikimedia.org/enwiki/latest/
  enwiki-latest-pages-articles.xml.bz2
- 13+ G file size, 4,000,000+ articles
- Preprocess, 2G+ memory required

```
# Remove words occuring less than 20 times, and words
    occuring in more than 10% of the documents. (keep_n is
    the vocabulary size)
wiki.dictionary.filter_extremes(no_below=20, no_above
    =0.1,^^Ikeep_n=100000)
```

# Train Word2Vec Model

- Use gensim with the following settings:

```
1 sentences = SentencesIterator(wiki)
2 model = gensim.models.Word2Vec(sentences=sentences, size
    =200, min_count=3, window=5, workers=cores)
```

- Required RAM: $\mathcal{O}(\text{size} \cdot |V|)$
- $\sim$ 24 hr training on virtual server
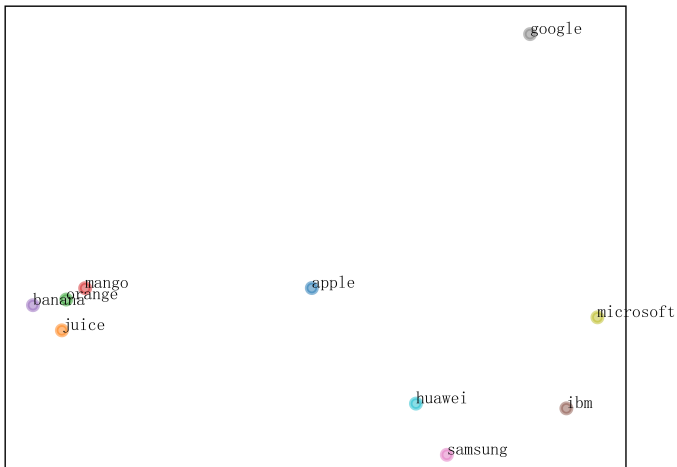
# Word2Vec Result

- Word similarity:

```
1  model = gensim.models.Word2Vec.load(os.path.join(MODEL_PATH
       , 'word2vec.model'), mmap='r')  # load large array
2  print(model.similarity('apple', 'mango'), model.similarity(
       'apple', 'ibm'))
3  0.465574139702 0.531822866913
```
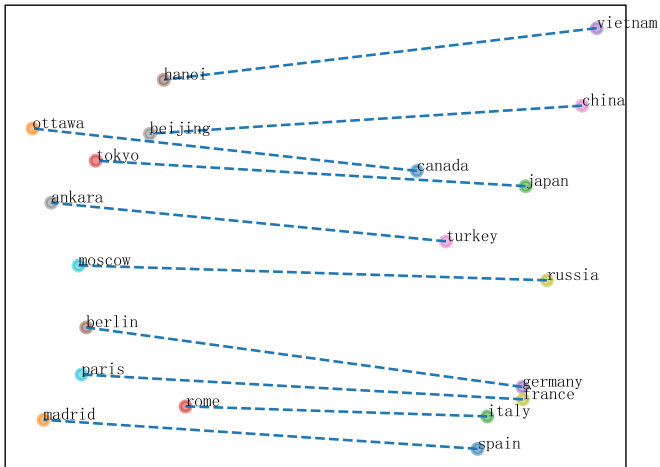
- Word Analogy:

```
1  King−man+woman
2
3  "queen" − similarity: 0.719937
4  "monarch" − similarity: 0.610968
5  "princess" − similarity: 0.608749
6  "prince" − similarity: 0.607932
7  "empress" − similarity: 0.586417
8  "throne" − similarity: 0.568553
9  "emperor" − similarity: 0.553523
10 "isabeau" − similarity: 0.540711
11 "amalasuntha" − similarity: 0.537015
```

# Word Vector: "Apple"

# Word Analogy: Country-Capital Paris

# Chinese Wikipedia Word2Vec

# Download Chinese Wikipedia Dump

- `https://dumps.wikimedia.org/zhwiki/latest/`
  `zhwiki-latest-pages-articles.xml.bz2`
- it contains traditional Chinese and simplified Chinese articles
- 1.3+ G file size, 230,000+ articles
- Preprocess, 2G+ memory required

# Preprocessing

- Use OpenCC to translate from simplified Chinese to traditional Chinese
- Developed on C++, supporting Python
-

```
1 openCC = OpenCC('s2twp')  # convert from Simplified Chinese
      to Traditional Chinese
2 converted = openCC.convert(original)
3 print(openCC.convert('这是简体中文'))
4 這是簡體中文
5 print(openCC.convert('智能手机'))
6 智慧手機
7 print(openCC.convert('一条短信'))
8 一條簡訊
```
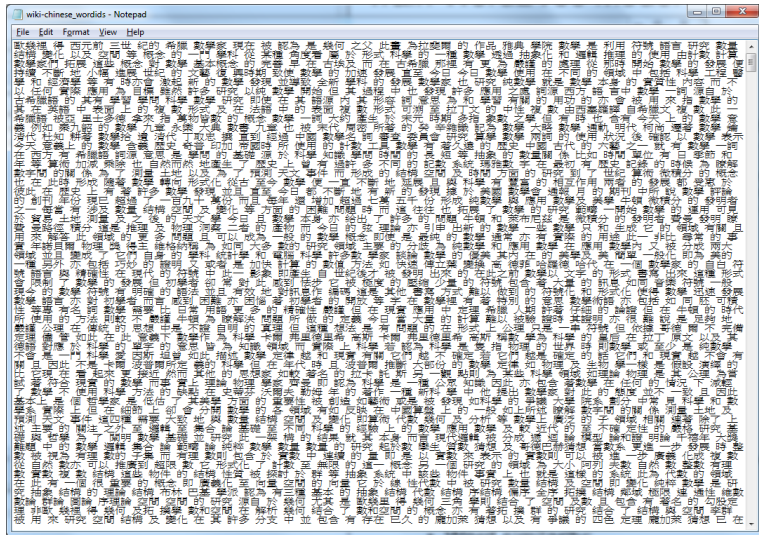
# Segmentation

- English uses whitespace, dots, etc to separate words, but Chinese does not
- 下雨天留客天留我不留
  - 下雨天留客, 天留我不留
  - 下雨天, 留客天, 留我不? 留
  - ...
- Many existing tools: stanford, fudan, jieba ...

# Segmentation

```
1  import jieba
2  print('/'.join(jieba.cut('下雨天留客天留我不留',cut_all=False)))
3  下雨天 / 留客 / 天留 / 我 / 不留
4  '我畢業於香港科技大學化工系'
5  我 / 畢業 / 於 / 香港科技大學 / 化工系
6  '見了他,她變得很低很低,低到塵埃里,但她心里是歡喜的,從塵埃里開出
      花來'
7  見 / 了 / 他 / , / 她 / 變得 / 很 / 低 / 很 / 低 / , / 低到 / 塵埃 / 里 / , / 但 / 她 / 心里 / 是 / 歡
      喜 / 的 / , / 從 / 塵埃 / 里 / 開出 / 花來
```

# File Format

# Train Chinese Word2Vec Model
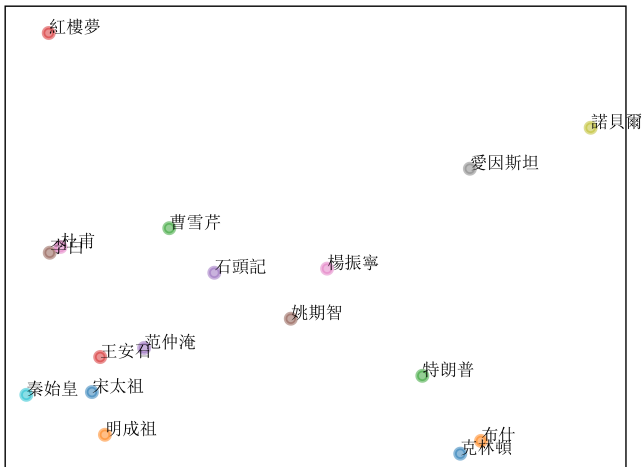
- Use gensim with the following settings:

```
sentences = SentencesIterator(wiki)
model = gensim.models.Word2Vec(sentences=sentences, size
    =300, min_count=3, window=5, workers=cores)
```
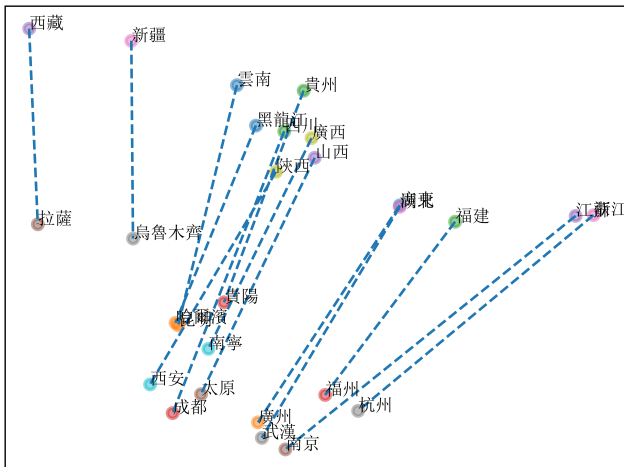
- $\sim$ 12 hr training

# Chinese Word2Vec Result

```
model = gensim.models.Word2Vec.load(os.path.join(MODEL_PATH, '
    word2vec.model'), mmap='r')
print(model.similarity('蘇軾', '東坡'))
0.612037855344

pprint(model.most_similar('奧斯卡'))
[('奧斯卡金像獎', 0.6850196123123169),
('奧斯卡獎', 0.6607037782669067),
('金球獎', 0.6433079242706299),
('東尼獎', 0.5794373750686646),
('金棕櫚獎', 0.5717018842697144),
('艾美獎', 0.5554369688034058),
('斷背山', 0.5518748760223389),
('泰坦尼克號', 0.5482996702194214),
('布克獎', 0.546217679977417),
('學院獎', 0.5445188879966736)]
```

# Word Similarity: Names

# Word Analogy: Province-Capital Paris

# Concluding Remarks

- Running Word2Vec with English/Chinese Wikipedia dump have acceptable performance in terms of word similarity/analogy
- Judge of Word2Vec's performance should also take machine translation into consideration
- Next …
  - sentence pairs
  - constructing basic LSTMs