# word2vec Demo

Efficient Estimation of Word Representations in Vector Space (Mikolov et al, 2013).

Andres Calderon and Hinna Shabir

University of California, Riverside

November 18, 2015

# Agenda

# Agenda

## word2vec source code

- https://code.google.com/p/word2vec/.
- Provides an efficient implementation of the continuous bag-of-words and skip-gram.
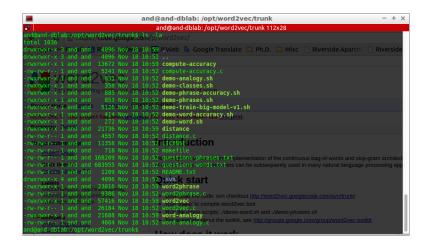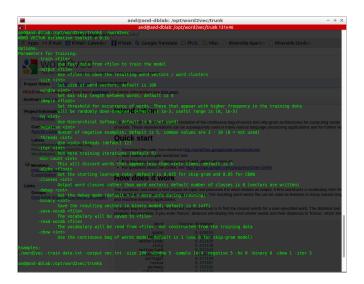- Clean and well documented code in C.

# word2vec source code

# word2vec source code

# word2vec source code

# word2vec source code

## demo-word.sh

```
## Compile the code...
make
## Download and unzip the training file...
if [ ! -e text8 ]; then
  wget http://mattmahoney.net/dc/text8.zip -O text8.gz
  gzip -d text8.gz -f
fi
## Run the model (taking time)...
time ./word2vec -train text8 -output vectors.bin -cbow 1 -size 200 -window 8 -negative 25 -hs 0
  ↪    -sample 1e-4 -threads 20 -binary 1 -iter 15
## Query word distances...
./distance vectors.bin
```

# text8 file

## demo-word.sh output

```
and@and-dblab:~/Documents/Projects/C++/word2vec/trunk$ ./demo-word.sh
make: Nothing to be done for 'all'.
--2015-11-12 18:18:10--  http://mattmahoney.net/dc/text8.zip
Resolving mattmahoney.net (mattmahoney.net)... 98.139.135.129
Connecting to mattmahoney.net (mattmahoney.net)|98.139.135.129|:80... connected.
HTTP request sent, awaiting response... 200 OK
Length: 31344016 (30M) [application/zip]
Saving to: text8.gz

text8.gz
  ↪    100%[===================================================>]   29.89M  1.74MB/s   in 18s

2015-11-12 18:18:28 (1.70 MB/s) - text8.gz saved [31344016/31344016]

Starting training using file text8
Vocab size: 71291
Words in train file: 16718843
Alpha: 0.000005  Progress: 100.10%  Words/thread/sec: 113.47k
real  10m15.450s
user  36m52.552s
sys  0m4.388s
Enter word or sentence (EXIT to break):
```

## demo-00.sh

```
## Get a small file...
head -c 5000000 text8 > text8_small
## Build the model...
./word2vec -train text8_small -output vectors_small.bin -cbow 1 -size 100 -window 5 -negative 0 -hs
 ↪    25 -threads 1 -iter 4 -min-count 2 -binary 1
## Query word distances...
./distance vectors_small.bin
```

# demo-01.sh

```
## Text model saving vocabulary
./word2vec -train text8_small -output vectors_small_50.txt -cbow 1 -size 50 -window 5 -negative 0
  ↪    -hs 25 -threads 1 -iter 4 -binary 0 -save-vocab vocab.txt
## Text model with just 3 dimensions
./word2vec -train text8_small -output vectors_small_3.txt -cbow 1 -size 3 -window 5 -negative 0 -hs
  ↪    25 -threads 1 -iter 4 -binary 0
## See the results...
echo "Text model size 50..."
head -n 5 vectors_small_50.txt
echo "Vocabulary..."
head -n 5 vocab.txt
echo "Text model size 3..."
head -n 5 vectors_small_3.txt
```

# demo-word.sh revisited

- **distance** can load a pre-trained model...
- Let's try some examples...
    1. california
    2. sciences
    3. happiness
    4. man
    5. ...

## demo-classes.sh

```
## Same than before...
make
if [ ! -e text8 ]; then
  wget http://mattmahoney.net/dc/text8.zip -O text8.gz
  gzip -d text8.gz -f
fi
## Train the model with classes rather than vectors...
time ./word2vec -train text8 -output classes.txt -cbow 1 -size 200 -window 8 -negative 25 -hs 0
  ↪    -sample 1e-4 -threads 20 -iter 15 -classes 500
## Sort the result by the second column...
sort classes.txt -k 2 -n > classes.sorted.txt
echo The word classes were saved to file classes.sorted.txt
```

# Interesting properties of the word vectors

- $\overrightarrow{paris} - \overrightarrow{france} + \overrightarrow{italy} \cong \overrightarrow{rome}$
- $\overrightarrow{king} - \overrightarrow{man} + \overrightarrow{women} \cong \overrightarrow{queen}$

# demo-analogy.sh

```
## Same that before...
make
if [ ! -e text8 ]; then
  wget http://mattmahoney.net/dc/text8.zip -O text8.gz
  gzip -d text8.gz -f
fi
echo -----------------------------------------------------------
echo Note that for the word analogy to perform well, the model
echo should be trained on much larger data set
echo Example input: paris france berlin
echo -----------------------------------------------------------
time ./word2vec -train text8 -output vectors.bin -cbow 1 -size 200 -window 8 -negative 25 -hs 0
  ↪   -sample 1e-4 -threads 20 -binary 1 -iter 15
## Call word-analogy script...
./word-analogy vectors.bin
```

## demo-analogy.sh

- Some examples...
  1. paris france bogota ...
  2. boy girl brother ...
  3. chicago illinois memphis ...
  4. poland zloty sweden ...
  5. bad worst good ...
  6. child children mouse ...
  7. going went selling ...
  8. king man queen ...
  9. mexico mexican peru ...
  10. berlin germany riyadh[1] ...
  11. woman angel man ...

---

[1] **word2phrase** will address the problem...

## demo-phrases.sh

```
## Compile...
make
## Download...
if [ ! -e news.2012.en.shuffled ]; then
  wget http://www.statmt.org/wmt14/training-monolingual-news-crawl/news.2012.en.shuffled.gz
  gzip -d news.2012.en.shuffled.gz -f
fi
## Pre-process...
sed -e "s//'/g" -e "s//'/g" -e "s/''/ /g" < news.2012.en.shuffled | tr -c "A-Za-z'_ \n" " " >
  ↪    news.2012.en.shuffled-norm0
time ./word2phrase -train news.2012.en.shuffled-norm0 -output news.2012.en.shuffled-norm0-phrase0
  ↪    -threshold 200 -debug 2
time ./word2phrase -train news.2012.en.shuffled-norm0-phrase0 -output
  ↪    news.2012.en.shuffled-norm0-phrase1 -threshold 100 -debug 2
tr A-Z a-z < news.2012.en.shuffled-norm0-phrase1 > news.2012.en.shuffled-norm1-phrase1
## Model...
time ./word2vec -train news.2012.en.shuffled-norm1-phrase1 -output vectors-phrase.bin -cbow 1 -size
  ↪    200 -window 10 -negative 25 -hs 0 -sample 1e-5 -threads 20 -binary 1 -iter 15
## Deploy...
./distance vectors-phrase.bin
```

# demo-phrases.sh output (1/3)

```
and@and-dblab:~/Documents/Projects/C++/word2vec/trunk$ ./demo-phrases.sh
make: Nothing to be done for 'all'.
--2015-11-12 18:33:08--
  ↪    http://www.statmt.org/wmt14/training-monolingual-news-crawl/news.2012.en.shuffled.gz
Resolving www.statmt.org (www.statmt.org)... 129.215.197.100
Connecting to www.statmt.org (www.statmt.org)|129.215.197.100|:80... connected.
HTTP request sent, awaiting response... 200 OK
Length: 786717767 (750M) [application/x-gzip]
Saving to: news.2012.en.shuffled.gz

news.2012.en.shuffled.gz
  ↪    100%[============================================>] 750.27M  4.25MB/s   in 3m 7s

2015-11-12 18:36:16 (4.02 MB/s) - news.2012.en.shuffled.gz saved [786717767/786717767]

...
```

# demo-phrases.sh output (2/3)

```
...
Starting training using file news.2012.en.shuffled-norm0
Words processed: 296900K     Vocab size: 33198K
Vocab size (unigrams + bigrams): 18838711
Words in train file: 296901342
Words written: 296900K
real  7m38.607s
user  7m8.592s
sys  0m15.176s

Starting training using file news.2012.en.shuffled-norm0-phrase0
Words processed: 280500K     Vocab size: 38761K
Vocab size (unigrams + bigrams): 21728781
Words in train file: 280513979
Words written: 280500K
real  7m0.022s
user  6m19.436s
sys  0m14.756s
...
```

# demo-phrases.sh output (3/3)

```
...
Starting training using file news.2012.en.shuffled-norm1-phrase1
Vocab size: 681320
Words in train file: 283545447
Alpha: 0.000005  Progress: 100.00%  Words/thread/sec: 162.97k
real  115m6.531s
user  434m57.904s
sys   1m4.464s
Enter word or sentence (EXIT to break):
```