Word-embeddings and Word2Vec

Yoav Goldberg

Last Time

- Language modeling.
- 1-hot X matrix
- Embedding layer
- (pre-trained) Word Embeddings.

Neural LM $\operatorname{softmax}(\Box)$ \uparrow what's in W3? $\Box W^3 + b^3$

columns of W3 correspond to vocab items!

 $\mathbb{R}^{d_{out}}$ $\Box W^3 + b^3$ \mathbb{R}^{d_2} $g(\Box W^2 + b^2)$ \mathbb{R}^{d_1} $g(\Box W^1 + b^1)$ $\mathbb{R}^{d_{in}}$ $\mathbf{E}_{[x_{k-4}]} \circ \mathbf{E}_{[x_{k-3}]} \circ \mathbf{E}_{[x_{k-2}]} \circ \mathbf{E}_{[x_{k-1}]}$

 $\mathbb{R}^{d_{out}}$

$$\uparrow \\ encode(x_{k-4}, x_{k-3}, x_{k-2}, x_{k-1})$$

toeory as

10 Neural Worder Embeddings

тJ



000 1 15 20 30 10 few half five six two three four several some many other those^{these} all both

chief head chairman director spokesman executive trader

Country and Capital Vectors Projected by PCA 2 China ------≫Beijing 1.5 **Russi**a Japan Moscow 1 Ankara ⇒√Tokvo Turkey -0.5 Poland Germany -0 France ≫Warsaw Berlin 📉 -0.5 Italy -Paris Athens Greece Rome Spain -1 Madrid -1.5 Portugal ->>Lisbon -2 -0.5 -1.5 0.5 1.5 -1 -2 0 2 1

Target Word	BoW5
	nightwing
batman	aquaman
	catwoman
	superman
	manhunter
hogwarts	dumbledore
	hallows
	half-blood
	malfoy
	snape
turing	nondeterministic
	non-deterministic
	computability
	deterministic
	finite-state
	gainesville
	fla
florida	jacksonville
	tampa
	lauderdale
	aspect-oriented
	smalltalk
object-oriented	event-driven
	prolog
	domain-specific
	singing
dancing	dance
	dances
	dancers
	tap-dancing

[word similarity example]

- We trained a language model.
- We ended up with vector representations of words.
- These representations are useful -- they encode various aspects of word similarity.





Neural LM

towards word2vec

- Training the language model is expensive. (Why?)
- Predicting a word based on previous words is nice, but can we do better?
- If all we care about are the word vectors...

- Radically simplify the neural LM.
- Very fast training.
- Obtain good word representations.

word2vec

word2vec

- dog
 - cat, dogs, dachshund, rabbit, puppy, poodle, rottweiler, mixed-breed, doberman, pig
- sheep
 - cattle, goats, cows, chickens, sheeps, hogs, donkeys, herds, shorthorn, livestock
- november
 - october, december, april, june, february, july, september, january, august, march
- jerusalem
 - tiberias, jaffa, haifa, israel, palestine, nablus, damascus katamon, ramla, safed
- teva
 - pfizer, schering-plough, novartis, astrazeneca, glaxosmithkline, sanofi-aventis, mylan, sanofi, genzyme, pharmacia

- instead of predicting word based on previous words...
- ...predict word based on **surrounding** words.

$$P(w_5|w_1w_2w_3w_4\Box) \longrightarrow P(w_3|w_1w_2\Box w_4w_5)$$

- Getting rid of the softmax over the vocabulary:
 - probabilities --> scores
 - score the quality of a sequence, not all the words
 - ranking based loss

 $score(D, A, \Box, C, F; G) = g(\mathbf{xU})\mathbf{v}$ $\mathbf{x} = (\mathbf{E}_{[D]} \circ \mathbf{E}_{[A]} \circ \mathbf{E}_{[G]} \circ \mathbf{E}_{[C]} \circ \mathbf{E}_{[F]})$

score
$$(D, A, \Box, C, F; G) = g(\mathbf{xU})\mathbf{v}$$

 $\mathbf{x} = (\mathbf{E}_{[D]} \circ \mathbf{E}_{[A]} \circ \mathbf{E}_{[G]} \circ \mathbf{E}_{[C]} \circ \mathbf{E}_{[F]})$

score $(c_1, c_2, \Box, c_3, c_4; w) >$ score $(c_1, c_2, \Box, c_3, c_4; w') + 1$ observed word random word

score
$$(D, A, \Box, C, F; G) = g(\mathbf{xU})\mathbf{v}$$

 $\mathbf{x} = (\mathbf{E}_{[D]} \circ \mathbf{E}_{[A]} \circ \mathbf{E}_{[G]} \circ \mathbf{E}_{[C]} \circ \mathbf{E}_{[F]})$

$$score(c_1, c_2, \Box, c_3, c_4; w) > score(c_1, c_2, \Box, c_3, c_4; w') + 1$$

observed word random word

 $L(w, c, w') = max(0, 1 - (score(c_{1:k}; w) - score(c_{1:k}; w')))$

How do you train this? (discuss)

 $score(D, A, \Box, C, F; G) = g(\mathbf{xU})\mathbf{v}$ $\mathbf{x} = (\mathbf{E}_{[D]} \circ \mathbf{E}_{[A]} \circ \mathbf{E}_{[G]} \circ \mathbf{E}_{[C]} \circ \mathbf{E}_{[F]})$

$$score(c_1, c_2, \Box, c_3, c_4; w) > score(c_1, c_2, \Box, c_3, c_4; w') + 1$$

observed word random word

 $L(w, c, w') = max(0, 1 - (score(c_{1:k}; w) - score(c_{1:k}; w')))$

from Collobert and Weston to Word2Vec

- Back to probabilities
- Simplify the scoring function

Back to probabilities

 $L(w, c, w') = max(0, 1 - (score(c_{1:k}; w) - score(c_{1:k}; w')))$ \downarrow $P(D = 1 | w, c) = \sigma(s(w, c)) = \frac{1}{1 + e^{-s(w, c)}}$

Back to probabilities

$$P(D = 1 | w, c) = \sigma(s(w, c)) = \frac{1}{1 + e^{-s(w, c)}}$$

$$\mathcal{L}(\Theta; D, \bar{D}) = \sum_{(w,c)\in D} \log P(D=1|w,c) + \sum_{(w',c)\in \bar{D}} \log P(D=0|w',c)$$

Back to probabilities

good word+context pairs **bad** word+context pairs $P(D = 1 | w, c) = \sigma(s(w, c)) = \frac{1}{1 + e^{-s(w,c)}}$ $\mathcal{L}(\Theta; D, \overline{D}) = \sum_{(w,c)\in D} \log P(D = 1 | w, c) + \sum_{(w',c)\in \overline{D}} \log P(D = 0 | w', c)$

Back to probabilities

$$P(D = 0|w, c) = 1 - P(D = 1|w, c)$$
$$P(D = 1|w, c) = \sigma(s(w, c)) = \frac{1}{1 + e^{-s(w, c)}}$$
$$\mathcal{L}(\Theta; D, \bar{D}) = \sum_{(w,c)\in D} \log P(D = 1|w, c) + \sum_{(w',c)\in \bar{D}} \log P(D = 0|w', c)$$

• Simplify the score function:

CBOW Skip-grams

• Simplify the score function:

CBOW Skip-grams

$$P(D = 1 | w, c) = \sigma(s(w, c)) = \frac{1}{1 + e^{-s(w, c)}}$$

what is c?

• Simplify the score function:

CBOW

$$\operatorname{score}(w; c_1, ..., c_k) = \left(\sum_{i=1}^k \mathbf{E}_{[c_i]}\right) \cdot \mathbf{E}'_{[w]}$$

$$P(D=1|w,c_{1:k}) = \frac{1}{1+e^{-(\mathbf{w}\cdot\mathbf{c_1}+\mathbf{w}\cdot\mathbf{c_2}+\ldots+\mathbf{w}\cdot\mathbf{c_k})}}$$

• Simplify the score function:

Skip-grams: $P(D = 1|w, c_i) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{c_i}}}$ $P(D = 1|w, c_{1:k}) = \prod_{i=1}^{k} P(D = 1|w, c_i) = \prod_{1=i}^{k} \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{c_i}}}$ $\log P(D = 1|w, c_{1:k}) = \log \sum_{i=1}^{k} \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{c_i}}}$

• Simplify the score function:

Skip-grams:

$$\mathcal{L}(\Theta; D, \bar{D}) = \sum_{(w,c)\in D} \log P(D=1|w,c) + \sum_{(w',c)\in \bar{D}} \log P(D=0|w',c)$$

$$= \arg \max_{\theta} \sum_{(w,c)\in D} \log \frac{1}{1 + e^{-v_c \cdot v_w}} + \sum_{(w,c)\in D'} \log(\frac{1}{1 + e^{v_c \cdot v_w}})$$

• How to train a word2vec model? (discuss)

Enriching Word Vectors with Subword Information

Piotr Bojanowski* and Edouard Grave* and Armand Joulin and Tomas Mikolov Facebook AI Research {bojanowski,egrave,ajoulin,tmikolov}@fb.com

Represent a word by the sum of its char ngrams.

Represent a word by the sum of its char ngrams.

Where = "<where>" "<where", "where>", "<wher", "where", "where", "here>", "<whe", "wher", "here", "ere>", "<wh",

 $v_{where} = E["<where>"] + E["<where"] +$

$$= \arg\max_{\theta} \sum_{(w,c)\in D} \log \frac{1}{1 + e^{-v_c(v_w)}} + \sum_{(w,c)\in D'} \log(\frac{1}{1 + e^{v_d(v_w)}})$$

Represent a word by the sum of its char ngrams.

"<where>" "<where", "where>", "<wher", "where", "here>", "<whe", "wher", "here", "ere>", "<wh",

- What are the benefits of FastText? Where do we expect it to be **better**?
- When do we expect it to work **worse**?
- Can you think of a generalization of it to non-human-language tasks?
- Can you think of non-human-language-tasks where it will hurt and not help?
- Can you think of a human languages where it may hurt?

Represent a word by the sum of its char ngrams.

"<where>" "<where", "where>", "<wher", "where", "here>", "<whe", "wher", "here", "ere>", "<wh",

$v_{where} = E["<where>"] + E["<where"] +$

Think:

Isn't this very expensive to run at inference time? Why or why not?

Other word vectors

- Other **contexts** are also possible.
- Other **algorithms** are also possible.
 - GloVe
 - NCE
 - SVD
- Implementations:
 - word2vec, word2vecf, GloVe, gensim, your own...

Using Word Vectors

- Use for initializing the word embeddings in other models. ("pre-training". more soon.)
- Use for finding similar words. (how?)
- Find a word similar to a group of words. (how?)
- Find the word that does not belong to a group. (how?)

Word Similarity

Similarity is calculated using cosine similarity:

$$sim(d ec{og}, c ec{a}t) = rac{d ec{og} \cdot c ec{a}t}{||d ec{og}|| \, ||c ec{a}t||}$$

For normalized vectors (||x|| = 1), this is equivalent to a dot product:

$$sim(d \vec{o} g, c \vec{a} t) = d \vec{o} g \cdot c \vec{a} t$$

Normalize the vectors when loading them.

Finding the most similar words to dog

- Compute the similarity from word \vec{v} to all other words.
- This is a single matrix-vector product: $W \cdot \vec{v}^{\top}$

- Result is a |V| sized vector of similarities.
- Take the indices of the k-highest values.
- FAST! for 180k words, d=300: ~30ms

Most Similar Words, in python+numpy code

W,words = load_and_normalize_vectors("vecs.txt")
W and words are numpy arrays.
w2i = {w:i for i,w in enumerate(words)}

dog = W[w2i['dog']] # get the dog vector

sims = W.dot(dog) # compute similarities

most_similar_ids = sims.argsort()[-1:-10:-1]
sim_words = words[most_similar_ids]

Similarity to a group of words

- "Find me words most similar to cat, dog and cow".
- Calculate the pairwise similarities and sum them:

$$W \cdot \vec{cat} + W \cdot \vec{dog} + W \cdot \vec{cow}$$

- Now find the indices of the highest values as before.
- Matrix-vector products are wasteful. Better option:

$$W \cdot (\vec{cat} + \vec{dog} + \vec{cow})$$

tagging + pre-training

we can use the **E** we got from LM training to initialize **E** for the POS tagging task. NOUN

The brown **fox** jumped over

VERB

brown fox jumped over the

(why is that helpful?)

PREP

fox jumped over the lazy

Pre-training

- A large part of the success of feed-forward networks in NLP comes from the use of pre-trained word embeddings.
- Pre-trained embeddings are an easy way to perform semi-supervised learning (or transfer learning).
- But notice: fine-tuning the pre-trained embeddings means that some features change, while most stay the same...

Pre-training

- Define an auxiliary task that you suspect is correlated with your prediction problem.
- Train a model to perform this task.
- Take features representations from this model as inputs to another model.

Pre-training

- In word2vec, auxiliary tasks are "predict word based on a window of size k around it" (CBOW) or "predict neighboring words in a window of k around a focus word" (skipgram).
- More generally, "predict word based on some context of the word".
- This is useful, as we can get tons of training data for free.
- The choice of contexts determines the resulting word representations.

Some possible contexts

Window of k around a word.

(smaller k: more syntactic. Larger k: more semantic)

- **Positional window around a word.** (more syntactic)
- Aligned words in a different language. Requires parallel corpus (synonyms, paraphrases)
- Neighbors in a dependency tree. Requires parser. (functional similarity)

Effect of Context

Target Word	Bag of Words (k=5)	Dependencies
	Dumbledore	Sunnydale
	hallows	Collinwood
Hogwarts	half-blood	Calarts
(Harry Potter's school)	Malfoy	Greendale
	Snape	Millfield
	Related to Harry Potter	Schools

Levy and Goldberg 2014 Dependency-based word embeddings

Effect of Context

Target Word	Bag of Words (k=5)	Dependencies
	nondeterministic	Pauling
	non-deterministic	Hotelling
Turing	computability	Heting
(computer scientist)	deterministic	Lessing
	finite-state	Hamming
	Related to computability	Scientists

Levy and Goldberg 2014 Dependency-based word embeddings

Effect of Context

Target Word	Bag of Words (k=5)	Dependencies
	singing	singing
	dance	rapping
dancing	dances	breakdancing
(dance gerund)	dancers	miming
	tap-dancing	busking
	Related to dance	Gerunds

Levy and Goldberg 2014 Dependency-based word embeddings

- no layers.
- no softmax. (so negative sampling.)
- cbow or skipgrams objectives.
- many implementations available.
- more in the NLP course.