Vector semantics and embeddings unit

- ▶ Lexical semantics, words as vectors (last week Monday)
- ▶ Catch-up day (last week Wednesday)
- ▶ Word2Vec (last week Friday)
- ▶ Finish Word2Vec (**today**)
- ▶ Overview of neural nets (Wednesday)
- ▶ Neural-net language models, in lab (Friday)

Today:

- ▶ Review of premise and elements
- ▶ Details of the training algorithm
- ▶ Observing results
- ▶ Problems with data-driven methods

**Goal:** Find word embeddings, vectors that represent words in a semantic space.

**Word2vec premise:** Train a classifier on a "fake" task and use that classifier's weights/parameters as word embeddings.

**Word2vec algorithm outline:** Given an corpus,
- ▶ Find the vocabulary of the corpus
- ▶ Collect training data from the corpus
- ▶ Train a classifier on that data
- ▶ Return the classifier's weights

**The classification task:** Given *target word w* and potential *context word c*, is *c* likely (or, *how likely is c*) to be used near *w*? That is, is *c* a true context word for *w*?

Alice   looked   all   around   her   at   the   flowers
              $c_0$      $c_1$    $w$      $c_2$    $c_3$

**The training data:** For every token *w* in the corpus, pair it with the tokens found *L* positions before it and *L* positions after it (positive examples). For every positive example, find *k* randomly chosen negative examples.

| $w$ | $c_{pos}$ | $c_{neg_0}$ | $c_{neg_1}$ | $c_{neg_2}$ |
|---|---|---|---|---|
| around | looked | pineapple | earnestly | asleep |
| around | all | . . . | | |
| around | her | . . . | | |
| around | at | . . . | | |

Algorithm parameters:

*L*, window size ($L = 2$ above)

*k*, number of negative samples ($k = 3$ above)

**Choosing negative samples:**

*For training a binary classifier, we also need negative examples. . . . [For each training instance,] we'll create k negative samples, each consisting of a target w plus a "noise word" $c_{neg}$. A noise word is a random word from the lexicon, constrained not to be the target word w. . . .*

*The noise words are chosen according to their weighted unigram frequency $p_\alpha(w)$, where $\alpha$ is a weight. . . .*

$$p_\alpha(w) = \frac{count(w)^\alpha}{\sum_{w'} count(w')^\alpha}$$

*[Weighting p, for example at $\alpha = .75$] gives better performance because it gives rare noise words slightly higher probability: for rare words, $P_\alpha(w) > P(w)$.*

Jurafsky and Martin, 6.8.2, pg 20–21

**The classifier:** Logistic regression

$$P(+ \mid w, c) = \sigma(\mathbf{w} \cdot \mathbf{c})$$

where

- $P(+ \mid w, c)$ is the probability $c$ is a context word for $w$.
- $\mathbf{w}$ is a vector for $w$ as a target word.
- $\mathbf{c}$ is a vector for $c$ as a context word.
- $\mathbf{w} \cdot \mathbf{c}$ is the dot product
- $\sigma x = \frac{1}{1+exp(-x)}$ is the logistic (sigmoid) function.

Parameters to the classifier: **W** and **C**, each $V \times D$ matrices.

Parameter to the algorithm: $D$, length of embedding

**The loss function:** The cross-entropy (negative log likelihood) loss. For a given data point $(w, c_{pos}, c_{neg_0}, \ldots c_{neg_{k-1}})$, the loss is

$$-\left( \log \sigma(\boldsymbol{c_{pos}} \cdot \boldsymbol{w}) + \sum_{i=0}^{k-1} \log \sigma(-\boldsymbol{c_{neg_i}} \cdot \boldsymbol{w}) \right)$$

**The training algorithm:** Stochastic gradient descent. For each data point $(w, c_{pos}, c_{neg_0}, \ldots c_{neg_{k-1}})$,

$$\boldsymbol{c_{pos}} \ -= \ \eta(\sigma(\boldsymbol{c_{pos}} \cdot \boldsymbol{w}) - 1)\boldsymbol{w}$$

$$\boldsymbol{c_{neg_i}} \ -= \ \eta(\sigma(\boldsymbol{c_{neg_i}} \cdot \boldsymbol{w}))\boldsymbol{w}$$

$$\boldsymbol{w} \ -= \ \eta\left( (\sigma(\boldsymbol{c_{pos}} \cdot \boldsymbol{w}) - 1)\boldsymbol{c_{pos}} + \sum_{i=0}^{k-1}(\sigma(\boldsymbol{c_{neg_i}} \cdot \boldsymbol{w}))\boldsymbol{c_{neg_i}} \right)$$

Parameter to the algorithm: $\eta$, the learning rate

Coming up:

- ▶ Read J&M chapter7 (Mon, Nov 13)

- ▶ Work on stylo assignment

Word2Vec assignment coming...