S&DS 355 / 555
**Introductory Machine Learning**
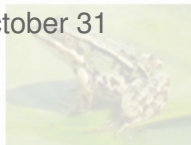
# Word Embeddings

Thursday, October 31

3. litoria        4. leptodactylidae        5. rana        7. eleutherodactylus

Yale

**For Today**

- Quiz ✓
- Class-based LMs (redux)
- Embeddings
- Embedding embeddings
- Next: Neural language models

# Class-based bigram model

- Model takes form

$$p(w_2 \mid w_1) = p(\text{class}(w_2) \mid \text{class}(w_1)) \, p(w_2 \mid \text{class}(w_2))$$
$$= p(c_2 \mid c_1) \, p(w_2 \mid c_2)$$

Brown et al., "Class-based $n$-gram models of natural language"

# Class-based bigram model

- Model takes form

$$p(w_2 \mid w_1) = p(\text{class}(w_2) \mid \text{class}(w_1))\, p(w_2 \mid \text{class}(w_2))$$
$$= p(c_2 \mid c_1)\, p(w_2 \mid c_2)$$

- Use bottom-up agglomerative clustering to group the words.

Brown et al., "Class-based *n*-gram models of natural language"

# Class-based bigram model

- Model takes form

$$p(w_2 \mid w_1) = p(\text{class}(w_2) \mid \text{class}(w_1))\, p(w_2 \mid \text{class}(w_2))$$
$$= p(c_2 \mid c_1)\, p(w_2 \mid c_2)$$

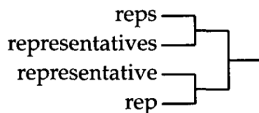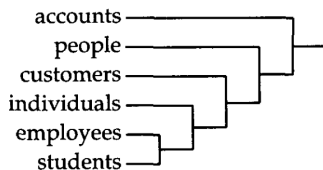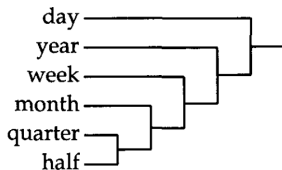- Use bottom-up agglomerative clustering to group the words.

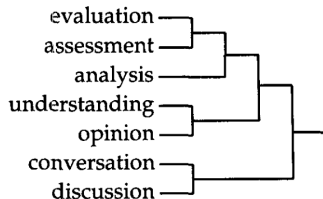- In each step, merge the pair of classes that gives the smallest reduction in likelihood of the data. (The MLE bigram model has the greatest likelihood.)

Brown et al., "Class-based *n*-gram models of natural language"

# Sample merges

# Sample clusters

Friday Monday Thursday Wednesday Tuesday Saturday Sunday weekends Sundays Saturdays

June March July April January December October November September August

people guys folks fellows CEOs chaps doubters commies unfortunates blokes

down backwards ashore sideways southward northward overboard aloft downwards adrift

water gas coal liquid acid sand carbon steam shale iron

great big vast sudden mere sheer gigantic lifelong scant colossal

man woman boy girl lawyer doctor guy farmer teacher citizen

American Indian European Japanese German African Catholic Israeli Italian Arab

pressure temperature permeability density porosity stress velocity viscosity gravity tension

mother wife father son husband brother daughter sister boss uncle

machine device controller processor CPU printer spindle subsystem compiler plotter

John George James Bob Robert Paul William Jim David Mike

anyone someone anybody somebody

feet miles pounds degrees inches barrels tons acres meters bytes

director chief professor commissioner commander treasurer founder superintendent dean custodian

liberal conservative parliamentary royal progressive Tory provisional separatist federalist PQ

had hadn't hath would've could've should've must've might've

asking telling wondering instructing informing kidding reminding bothering thanking deposing

that tha theat

head body hands eyes voice arm seat eye hair mouth

# Group globally, compute locally

- Clusters contain syntactic and semantic elements
- Surprising, since use local statistics only
- "A word is known by the company it keeps"
- Two words are similar if they appear with similar words

# Perplexity

*Perplexity* is a simple transformation of the likelihood:

$$\text{Perplexity}(\theta) = \left( \prod_{i=1}^{N} p_\theta(w_i \mid w_{1:i-1}) \right)^{-1/N}$$

Inverse of geometric mean.

If perplexity is 100, then the model predicts, on average, as if there were 100 equally likely words to follow.

# Class-based bigram model

- Clusters contain a lot of information

- Each word in a single class; partition of vocabulary

- Number of parameters: $O(C^2 + C \cdot V)$ where $C$ is number of classes.

- Reduces size of model with relatively small increase in perplexity, $244 \mapsto 271$.

# Pointwise mutual information (PMI)

Average mutual information

$$I(W_1, W_2) = \sum_{w_1, w_2} p(w_1, w_2) \log \frac{p(w_1, w_2)}{p(w_1)p(w_2)}$$

Pointwise mutual information (PMI)

$$\log \left( \frac{p_{\text{near}}(w_1, w_2)}{p(w_1)p(w_2)} \right)$$

- How likely are specific words/clusters to co-occur together within some window, compared to if they were independent?

# Example clusters from PMI

we our us ourselves ours
question questions asking answer answers answering
performance performed perform performs performing
tie jacket suit
write writes writing written wrote pen
morning noon evening night nights midnight bed
attorney counsel trial court judge
problems problem solution solve analyzed solved solving
letter addressed enclosed letters correspondence
large size small larger smaller
operations operations operating operate operated
school classroom teaching grade math
street block avenue corner blocks
table tables dining chairs plate
published publication author publish writer titled
wall ceiling walls enclosure roof
sell buy selling buying sold

# Shortcomings of word clusters

- Can't use vector space operations
- Doesn't give "features" of words
- These are addressed with "distributed representations" (next)

# Core idea of embeddings

- Form a language model but replace classes by vectors, one for each word

- Use PMI-like scores to fit the vectors

- Can be applied whenever have cooccurrence data.

## Embedding LM

Model is

$$p(w_2 \mid w_1) = \frac{\exp(\phi(w_2)^T \phi(w_1)}{\sum_w \exp(\phi(w)^T \phi(w_1))}.$$

As before,

$$\begin{aligned} \ell(\phi) &= \sum_{w_1, w_2} p(w_1, w_2) \log p(\phi_2 \mid \phi_1) p(w_2 \mid \phi_2) \\ &= I(\Phi_1, \Phi_2) - H(W) \end{aligned}$$

Thus, we want embedding vectors with high mutual information.

# Constructing embeddings

Carry out stochastic gradient descent over the embedding vectors $\phi \in \mathbb{R}^d$ (where $d \approx$ 50–100 is chosen by trial and error)

This is what Mikolov et al. (2014, 2015) did at Google. With a couple of heuristics:

"Distributed representations of words," (2014) "Efficient representations of words in vector space" (2015)

# Constructing embeddings

Heuristics used:

- Skip-gram: predict surrounding words from current word

# Constructing embeddings

Heuristics used:

- Skip-gram: predict surrounding words from current word
- An issue with this is that it "over generates" the data. With text `the lazy brown fox jumped` we will have $p(\text{brown}|\text{lazy})$ and $p(\text{brown}|\text{fox})$

# Constructing embeddings

Heuristics used:

- Skip-gram: predict surrounding words from current word
- An issue with this is that it "over generates" the data. With text `the lazy brown fox jumped` we will have $p(\text{brown}|\text{lazy})$ and $p(\text{brown}|\text{fox})$
- Second is computational. The bottleneck is computing the denominator in the logistic (softmax) probability.

# Constructing embeddings

Heuristics used:

- Skip-gram: predict surrounding words from current word

- An issue with this is that it "over generates" the data. With text `the lazy brown fox jumped` we will have $p(\text{brown}|\text{lazy})$ and $p(\text{brown}|\text{fox})$

- Second is computational. The bottleneck is computing the denominator in the logistic (softmax) probability.

- Use "negative sampling": Approximation

$$\sum_w \exp(\phi(w)^T \phi(w_1))$$
$$\approx \exp(\phi(w_2)^T \phi(w_1)) + \sum_{\text{random } w} \exp(\phi(w)^T \phi(w_1))$$

# Analogies

These heuristics enable training on very large text collections. Leads to vector representations of words with interesting properties.

For example, analogies:

`king` is to `man` as ? is to `woman`

## Analogies

These heuristics enable training on very large text collections. Leads to vector representations of words with interesting properties.

For example, analogies:

`king` is to `man` as ? is to `woman`

`Paris` is to `France` as ? is to `Germany`

## Analogies

These heuristics enable training on very large text collections. Leads to vector representations of words with interesting properties.

For example, analogies:

<div align="center">

king is to man as ? is to woman

Paris is to France as ? is to Germany

</div>

$$\phi(\text{king}) - \phi(\text{man}) \stackrel{?}{\approx} \phi(\text{queen}) - \phi(\text{woman})$$

$$\widehat{w} = \arg\min_{w} \|\phi(\text{king}) - \phi(\text{man}) + \phi(\text{woman}) - \phi(w)\|^2$$

Does $\widehat{w} = \text{queen}$?

# Learned Analogies

Table 8: *Examples of the word pair relationships, using the best word vectors from Table 4 (Skip-gram model trained on 783M words with 300 dimensionality).*

| Relationship | Example 1 | Example 2 | Example 3 |
|---|---|---|---|
| France - Paris | Italy: Rome | Japan: Tokyo | Florida: Tallahassee |
| big - bigger | small: larger | cold: colder | quick: quicker |
| Miami - Florida | Baltimore: Maryland | Dallas: Texas | Kona: Hawaii |
| Einstein - scientist | Messi: midfielder | Mozart: violinist | Picasso: painter |
| Sarkozy - France | Berlusconi: Italy | Merkel: Germany | Koizumi: Japan |
| copper - Cu | zinc: Zn | gold: Au | uranium: plutonium |
| Berlusconi - Silvio | Sarkozy: Nicolas | Putin: Medvedev | Obama: Barack |
| Microsoft - Windows | Google: Android | IBM: Linux | Apple: iPhone |
| Microsoft - Ballmer | Google: Yahoo | IBM: McNealy | Apple: Jobs |
| Japan - sushi | Germany: bratwurst | France: tapas | USA: pizza |

# Evaluation Analogies

| Type of relationship | Word Pair 1 | | Word Pair 2 | |
|---|---|---|---|---|
| Common capital city | Athens | Greece | Oslo | Norway |
| All capital cities | Astana | Kazakhstan | Harare | Zimbabwe |
| Currency | Angola | kwanza | Iran | rial |
| City-in-state | Chicago | Illinois | Stockton | California |
| Man-Woman | brother | sister | grandson | granddaughter |
| Adjective to adverb | apparent | apparently | rapid | rapidly |
| Opposite | possibly | impossibly | ethical | unethical |
| Comparative | great | greater | tough | tougher |
| Superlative | easy | easiest | lucky | luckiest |
| Present Participle | think | thinking | read | reading |
| Nationality adjective | Switzerland | Swiss | Cambodia | Cambodian |
| Past tense | walking | walked | swimming | swam |
| Plural nouns | mouse | mice | dollar | dollars |
| Plural verbs | work | works | speak | speaks |

Mikolov et al., "Distributed representations of words," (2014); "Efficient representations of words in vector space" (2015)

# GloVe

Shortly after: Stanford group introduced a computational expedient
(with attempt to give a "principled" motivation)

$$\mathcal{O}(\phi) = \sum_{w_1, w_2} f(c_{w_1, w_2}) \left( \phi(w_1)^T \phi(w_2) - \log c_{w_1, w_2} \right)^2$$

where $c_{w, w'}$ are cooccurrence counts.

- A type of regression estimator. Can interpret/relate this to other objectives.
- Main advantage is that SGD can be carried out much more efficiently

Pennington et al., "GloVe: Global vectors for word representation," (2015)

# GloVe

$$\mathcal{O}(\phi) = \sum_{w_1, w_2} f(c_{w_1, w_2}) \left( \phi(w_1)^T \phi(w_2) - \log c_{w_1, w_2} \right)^2$$

where $c_{w,w'}$ are cooccurrence counts.

- Heuristic weighting function

$$f(x) = \left( \frac{x}{x_{\max}} \right)^\alpha$$

where $\alpha = 3/4$ set empirically.

- So $10^{-4} \mapsto 10^{-3}$. Each order of magnitude down gets "boosted" by 1/4-magnitude.

Pennington et al., "GloVe: Global vectors for word representation," (2015)

# GloVe site and code

# Recommendation via Embedding

**Notebook**

Let's go to the Python notebook!

# Embedding embeddings: t-SNE

- How can we visualize the embeddings?

# Embedding embeddings: t-SNE

- How can we visualize the embeddings?
- We're in a very high dimensional space

# Embedding embeddings: t-SNE

- How can we visualize the embeddings?
- We're in a very high dimensional space
- Can do PCA, but this will introduce an additional projection/approximation step

# Embedding embeddings: t-SNE

- How can we visualize the embeddings?

- We're in a very high dimensional space

- Can do PCA, but this will introduce an additional projection/approximation step

- Many visualization techniques exist. A currently popular one is t-SNE: "Student-t Stochastic Neighborhood Embedding"

# t-SNE

Here's the idea behind t-SNE:

- Form a language model using the embeddings

# t-SNE

Here's the idea behind t-SNE:

- Form a language model using the embeddings
- Scale and symmetrize, giving a matrix $P = [P_{ij}]$

---

Pronounced: **tee**-snee

# t-SNE

Here's the idea behind t-SNE:

- Form a language model using the embeddings
- Scale and symmetrize, giving a matrix $P = [P_{ij}]$
- Represent word $i$ by $y_i \in \mathbb{R}^2$. Use a heavy-tailed distribution (Student-t with one degree of freedom)

# t-SNE

Here's the idea behind t-SNE:

- Form a language model using the embeddings

- Scale and symmetrize, giving a matrix $P = [P_{ij}]$

- Represent word $i$ by $y_i \in \mathbb{R}^2$. Use a heavy-tailed distribution (Student-t with one degree of freedom)

- Select $y_i$ using stochastic gradient descent

---

Pronounced: **tee**-snee

# t-SNE: Detailed algorithm

For each word $w_i$ compute a language model

$$P_{j|i} \propto \exp\left(-\frac{\|\phi(w_i) - \phi(w_j)\|^2}{2h_i^2}\right)$$

That is:

$$P_{j|i} = \frac{\exp\left(-\dfrac{\|\phi(w_i) - \phi(w_j)\|^2}{2h_i^2}\right)}{\displaystyle\sum_k \exp\left(-\dfrac{\|\phi(w_i) - \phi(w_k)\|^2}{2h_i^2}\right)}$$

---

Pronounced: **tee**-snee

# t-SNE: Detailed algorithm

For each word $w_i$ compute a language model

$$P_{j|i} \propto \exp\left(-\frac{\|\phi(w_i) - \phi(w_j)\|^2}{2h_i^2}\right)$$

That is:

$$P_{j|i} = \frac{\exp\left(-\dfrac{\|\phi(w_i) - \phi(w_j)\|^2}{2h_i^2}\right)}{\sum_k \exp\left(-\dfrac{\|\phi(w_i) - \phi(w_k)\|^2}{2h_i^2}\right)}$$

Choose the bandwith $h_i$ so that the perplexity is, say, 10. This puts the probabilities all on the same scale.

---

Pronounced: **tee**-snee

# t-SNE: Detailed algorithm

For each word $w_i$ compute a language model

$$P_{j\,|\,i} \propto \exp\left(-\frac{\|\phi(w_i) - \phi(w_j)\|^2}{2h_i^2}\right)$$

# t-SNE: Detailed algorithm

For each word $w_i$ compute a language model

$$P_{j|i} \propto \exp\left(-\frac{\|\phi(w_i) - \phi(w_j)\|^2}{2h_i^2}\right)$$

Now form

$$P_{ij} = \frac{1}{2}\left(P_{j|i} + P_{i|j}\right)$$

as a simple way of symmetrizing.

# t-SNE: Detailed algorithm

Now form Student-t distribution depending on the visualization vectors $y_i \in \mathbb{R}^2$:

$$Q_{ij} \propto \left(1 + \|y_i - y_j\|^2\right)^{-1}$$

## t-SNE: Detailed algorithm

Now form Student-t distribution depending on the visualization vectors $y_i \in \mathbb{R}^2$:

$$Q_{ij} \propto \left(1 + \|y_i - y_j\|^2\right)^{-1}$$

That is:

$$Q_{ij} = \frac{\left(1 + \|y_i - y_j\|^2\right)^{-1}}{\sum_{k \neq \ell} \left(1 + \|y_k - y_\ell\|^2\right)^{-1}}$$

This has fatter tails than a Gaussian

## t-SNE: Detailed algorithm

Finally, run stochastic gradient descent (SGD) over the vectors $y_i$ to optimize:

$$\widehat{y} = \arg\min \sum_{ij} P_{ij} \log P_{ij}/Q_{ij}$$
$$= \arg\max \sum_{ij} P_{ij} \log Q_{ij}$$

Interpretation: if $\phi(w_i)$ is very close to $\phi(w_j)$ then $y_i$ will be close to $y_j$. (long distances may be stretched further...)
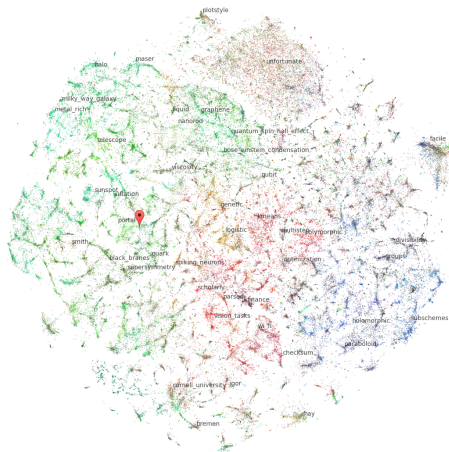
# t-SNE: More info and examples

```
https://lvdmaaten.github.io/tsne/
http://cs.stanford.edu/people/karpathy/tsnejs/
```

Note: This is just a visualization technique, to give intuition for the
high dimensional embedding

# t-SNE: Examples



`http://www.cs.cornell.edu/~ginsparg/arxiv/gmaps2.html`

# Summary: Word embeddings

- Word embeddings are vector representations of words, learned from cooccurrence statistics

- The models can be viewed in terms of logistic regression and class-based bigram models

- Surprising semantic relations are encoded in linear relations

- Various heuristics have been introduced to get scalability

- Embeddings improve with more data

- t-SNE is an algorithm for visualizing embeddings