# Statistical Machine Learning
# for Large and Unstructured Data

## Word Embeddings

Stephen Hansen
University College London

# Word Embeddings

A word embedding is a low-dimensional vector representation of a word.

Ideally in this low-dimensional vector space words with similar meanings will lie close together.

The construction of word embeddings has been a major topic in NLP in the past decade.

Embedding vectors can be of interest in their own right, or else form part of the representation of a document for other tasks.

# Local Contexts and Word Embeddings

Recall that $w_{d,n}$ is the $n$th word in document $d$.

The *context* of $w_{d,n}$ is a length-$2L$ window of words around $w_{d,n}$:

$$C(w_{d,n}) = [w_{d,n-L}, w_{d,n-L+1}, \ldots, w_{d,n+L-1}, w_{d,n+L}]$$

Can truncate context appropriately if window stretches past beginning or end of text.

In line with Firth's distributional hypothesis, modern word embedding models seek to generate similar embeddings for words that share similar contexts.

# GloVe

The GloVe model [Pennington et al., 2014] begins with a $V \times V$ matrix **W** of local word co-occurrences.

$W_{ij}$ is the number of times term $j$ appears within the context of $i$.

Assign to each term $v$ an embedding vector $\boldsymbol{\rho}_v \in \mathbb{R}^V$.

$$\min \sum_{i,j} f(W_{i,j}) \left( \boldsymbol{\rho}_i^T \boldsymbol{\rho}_j - \log(W_{i,j}) \right)^2$$

Terms that co-occur frequently will have more highly correlated embedding vectors.

# Word2Vec

Word2vec [Mikolov et al., 2013a, Mikolov et al., 2013b] is another well-known model for word embeddings that incorporates local context.

In addition to an embedding vector, each term is assigned a context vector $\boldsymbol{\alpha}_v \in \mathbb{R}^V$.

Word vectors are chosen to solve word-prediction tasks:

$$\Pr\left[w_{d,n} = v \mid C(w_{d,n})\right] = \frac{\exp(\overline{\boldsymbol{\alpha}}_{d,n}^T \boldsymbol{\rho}_v)}{\sum_{v'} \exp(\overline{\boldsymbol{\alpha}}_{d,n}^T \boldsymbol{\rho}_{v'})} \text{ where } \overline{\boldsymbol{\alpha}}_{d,n} = \frac{1}{2L} \sum_{w \in C(w_{d,n})} \boldsymbol{\alpha}_w$$

Example of self-supervised learning.

The version of Word2Vec is the continuous-bag-of-words model; the skip-gram model instead predicts context words given a center word.

# Terms Close to Uncertainty in FOMC Transcripts

| term | sim | term | sim |
|---|---|---|---|
| uncertainties | 0.741 | challenges | 0.415 |
| anxiety | 0.48 | fragility | 0.405 |
| pessimism | 0.479 | clarity | 0.401 |
| skepticism | 0.465 | concerns | 0.4 |
| optimism | 0.445 | risks | 0.397 |
| caution | 0.442 | disagreement | 0.387 |
| gloom | 0.437 | volatility | 0.384 |
| uncertain | 0.433 | tension | 0.383 |
| sensitivity | 0.427 | certainty | 0.382 |
| angst | 0.426 | skepticism | 0.38 |

# Terms Close to Risk

| term | sim | term | sim |
|---|---|---|---|
| risks | 0.737 | misdirected | 0.385 |
| threat | 0.609 | odds | 0.379 |
| danger | 0.541 | uncertainty | 0.375 |
| dangers | 0.463 | concern | 0.371 |
| vulnerability | 0.457 | prospect | 0.37 |
| chances | 0.451 | instability | 0.363 |
| breakout | 0.433 | potentially | 0.352 |
| probability | 0.426 | concerns | 0.352 |
| possibility | 0.409 | challenges | 0.346 |
| likelihood | 0.406 | risking | 0.342 |

# Concept Detection

# Expanding Dictionaries

One application of word embeddings is to augment human judgment in the construction of dictionaries.

Motivation is that economists are experts in which concept might be most important in a particular setting, but not in which words relate to that concept.

One can specify a set of 'seed' words and then find nearest neighbors of those words to populate a dictionary.

Strategy adopted by several recent papers:

1. [Hanley and Hoberg, 2019]
2. [Li et al., 2021]
3. [Bloom et al., 2021]
4. [Davis et al., 2020]

# Embedding Dictionaries

Dictionaries provide a coarse representation of concepts in that some relevant terms might be missing altogether, and strength of association with concept isn't accounted for.

One strategy is to measure the association between documents and word lists in an embedding space rather than the bag-of-words space.

Recent example is [Gennaro and Ash, 2022] which studies emotional language in politics using the Congressional Record corpus.

Set $A$ of words represents emotion, and set $C$ of words represents cognition (both from LIWC).

Emotionality of speech $i$ is

$$Y_i = \frac{\text{sim}(\boldsymbol{d_i}, \boldsymbol{A}) + b}{\text{sim}(\boldsymbol{d_i}, \boldsymbol{C}) + b}$$
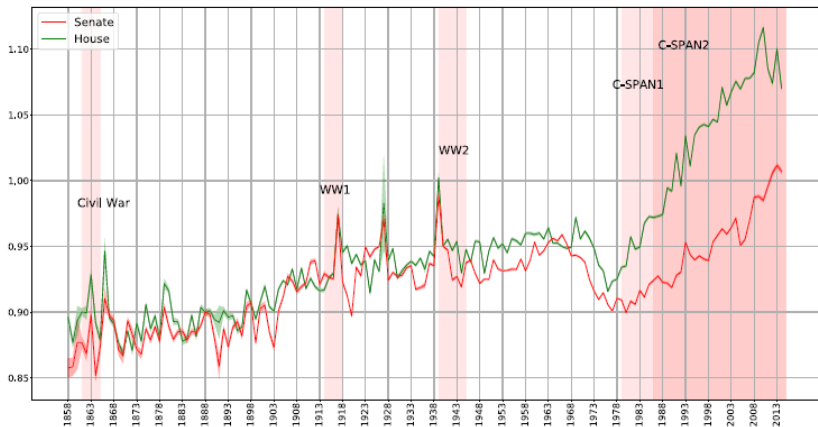
# Results



Fig. 2. *Emotionality in U.S. Congress by Chamber, 1858–2014.*
*Notes:* Time series of emotionality in the Senate (red) and the House of Representatives (green).

# Transfer Learning

In the above examples, Word2vec is fit directly to the data of interest.

In many use cases, one instead uses an estimated model from an auxiliary corpus for word embeddings, or as starting values in embedding estimation.

This is known as transfer learning and becomes particularly important for large-scale language models.

# Importance of Training Corpus

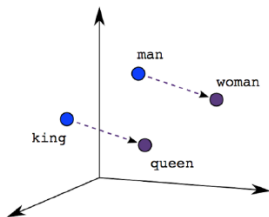Relationships among words can vary depending on the training corpus.

Example of training word embeddings on Wiki/Newswire text and on Harvard Business Review.
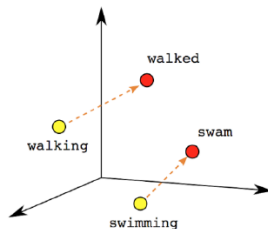
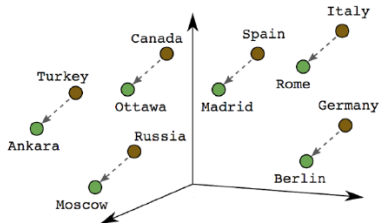| team | | leader | |
| --- | --- | --- | --- |
| HBR | Generic | HBR | Generic |
| teams | teams | leadership | leaders |
| project_team | squad | leaders | leadership |
| management_team | players | manager | party |
| executive_team | football | person | opposition |
| group | coach | strong_leader | led |
| staff | league | chief_executive | rebel |

# Relationship Among Concepts

# Directions Encode Meaning



Male-Female

Verb Tense

Country-Capital

# Word Embeddings and Cultural Attitudes

Because word embeddings appear to capture semantically meaningful relationships among words, there is interest in using them to measure cultural attitudes.
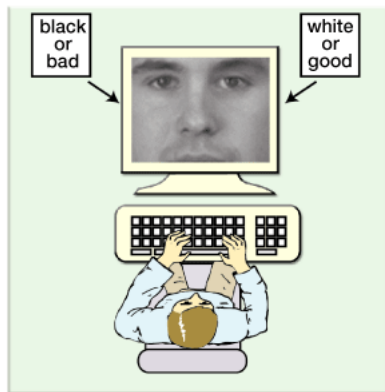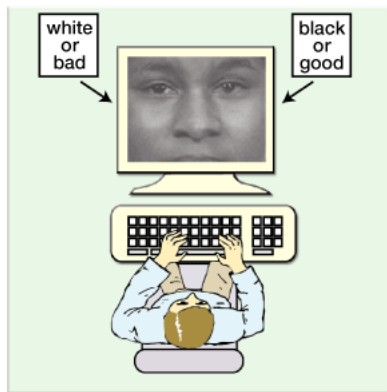
In psychology there is a long-standing Implicit Association Test that measures participants' time to correctly classify images depending on word combinations.

The hypothesis is that reaction times are shorter when word combinations more naturally belong together, which allows a measure of bias.

[Caliskan et al., 2017] have use word embeddings to ask whether similar biases exist in natural language.

# Implicit Association Test

# Word-Embedding Association Test

The Word-Embedding Association Test (WEAT) measures whether two sets of target words $X, Y$ (e.g. male, female words) differ in their relative similarity to two sets of attribute words $A, B$ (e.g. career, family words).

Let $\cos(\mathbf{x}, \mathbf{y})$ be cosine similarity between vectors $\mathbf{x}$ and $\mathbf{y}$.
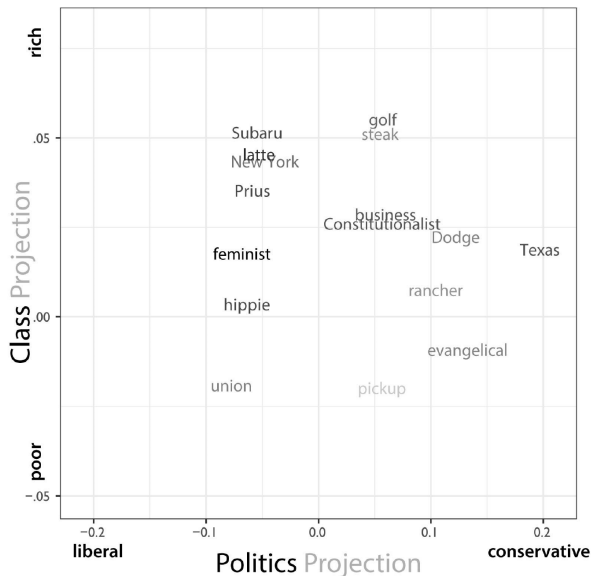
Let $s(\mathbf{w}, A, B) = \text{mean}_{\mathbf{a} \in A} \cos(\mathbf{w}, \mathbf{a}) - \text{mean}_{\mathbf{b} \in B} \cos(\mathbf{w}, \mathbf{b})$.

$$\text{WEAT} = \frac{\sum\limits_{\mathbf{x} \in X} s(\mathbf{x}, A, B) - \sum\limits_{\mathbf{y} \in Y} s(\mathbf{y}, A, B)}{\text{std}_{\mathbf{x} \in X \cup Y}\, s(\mathbf{x}, A, B)}$$

# IAT vs WEAT

| Target words | Attribute words | Original finding | | | | Our finding | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Ref. | N | d | P | $N_T$ | $N_A$ | d | P |
| Flowers vs. insects | Pleasant vs. unpleasant | (5) | 32 | 1.35 | $10^{-8}$ | 25 × 2 | 25 × 2 | 1.50 | $10^{-7}$ |
| Instruments vs. weapons | Pleasant vs. unpleasant | (5) | 32 | 1.66 | $10^{-10}$ | 25 × 2 | 25 × 2 | 1.53 | $10^{-7}$ |
| European-American vs. African-American names | Pleasant vs. unpleasant | (5) | 26 | 1.17 | $10^{-5}$ | 32 × 2 | 25 × 2 | 1.41 | $10^{-8}$ |
| European-American vs. African-American names | Pleasant vs. unpleasant from (5) | (7) | Not applicable | | | 16 × 2 | 25 × 2 | 1.50 | $10^{-4}$ |
| European-American vs. African-American names | Pleasant vs. unpleasant from (9) | (7) | Not applicable | | | 16 × 2 | 8 × 2 | 1.28 | $10^{-3}$ |
| Male vs. female names | Career vs. family | (9) | 39k | 0.72 | $<10^{-2}$ | 8 × 2 | 8 × 2 | 1.81 | $10^{-3}$ |
| Math vs. arts | Male vs. female terms | (9) | 28k | 0.82 | $<10^{-2}$ | 8 × 2 | 8 × 2 | 1.06 | .018 |
| Science vs. arts | Male vs. female terms | (10) | 91 | 1.47 | $10^{-24}$ | 8 × 2 | 8 × 2 | 1.24 | $10^{-2}$ |
| Mental vs. physical disease | Temporary vs. permanent | (23) | 135 | 1.01 | $10^{-3}$ | 6 × 2 | 7 × 2 | 1.38 | $10^{-2}$ |
| Young vs. old people's names | Pleasant vs. unpleasant | (9) | 43k | 1.42 | $<10^{-2}$ | 8 × 2 | 8 × 2 | 1.21 | $10^{-2}$ |

# Language and Culture [Kozlowski et al., 2019]
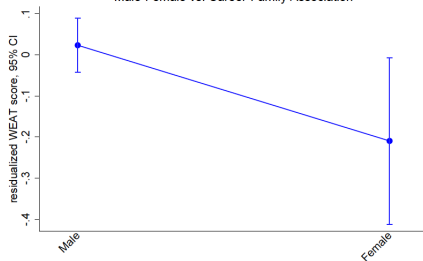
# Does Language affect Decisions?

[Ash et al., 2020] use a measure similar to WEAT to measure linguistic gender bias among judges using written opinions.

They then match judge-specific bias scores with individual judge decisions to see whether the two are related.
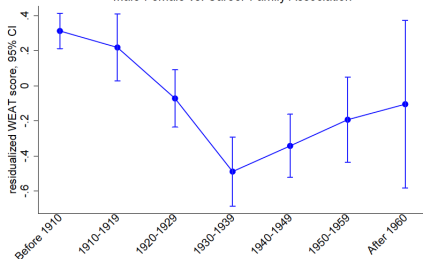
Data is the universe of US appellate court decisions from 1890-2013.

# WEAT and Judge Characteristics



WEAT Effect Size by Judge Gender
Male-Female vs. Career-Family Association

WEAT Effect Size by Judge Cohort
Male-Female vs. Career-Family Association

# Effects of WEAT

Judges with higher lexical bias are:

▶ Less likely to cast vote in favor of women's interests

▶ More likely to vote more conservatively across all issues

▶ Less likely to cite women in their opinions

▶ More likely to reverse female district judges

# Topic Modeling for Embeddings

# Incorporating Document Heterogeneity

Topic models express heterogeneity across documents within the context of the bag-of-words model.

Word embeddings exploit local context to construct a semantically meaningful vector space.

Can we combine the strengths of the two approaches?

[Dieng et al., 2020] present the embedded topic model.

# Statistical Model

1. Draw topic proportions $\boldsymbol{\theta}_d \sim \mathcal{LN}(0, 1)$.
2. For each word $w_{d,n}$:
   2.1 Draw a topic assignment $z_{dn} \sim \mathrm{Mult}(\boldsymbol{\theta}_d, 1)$.
   2.2 $w_{d,n} \sim \mathrm{Mult}\left(\mathrm{softmax}(\mathbf{P}^T \boldsymbol{\alpha}_{z_{dn}}), 1\right)$.

Where $\mathbf{P}$ is a $K \times V$ matrix whose $v$th column is the word embedding $\boldsymbol{\rho}_v$.

Topics are now represented by vectors in the embedding space.
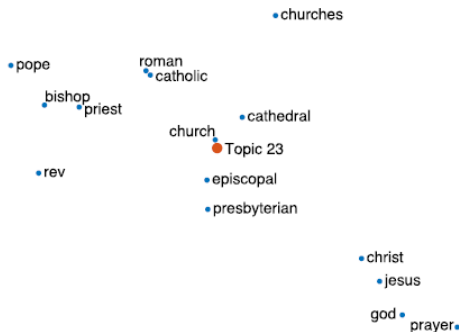
# Topic in Embedding Space



Figure 2: A topic about Christianity found by the ETM on *The New York Times*. The topic is a point in the word embedding space.

# Most Common Topics

| ETM | | | | | | |
|---|---|---|---|---|---|---|
| game | music | united | wine | company | yankees | art |
| team | mr | israel | food | stock | game | museum |
| season | dance | government | sauce | million | baseball | show |
| coach | opera | israeli | minutes | companies | mets | work |
| play | band | mr | restaurant | billion | season | artist |

Table 3: Top five words of seven most used topics from different document models on 1.8M documents of the *New York Times* corpus with vocabulary size 212,237 and $K = 300$ topics.
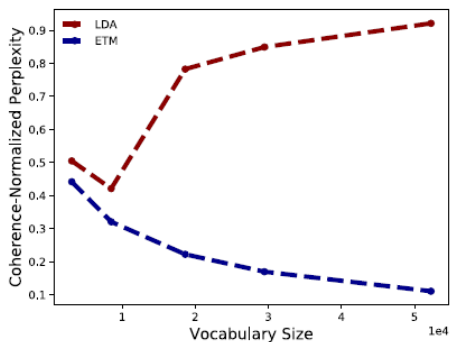
# Performance vs LDA



Figure 1: Ratio of the held-out perplexity on a document completion task and the topic coherence as a function of the vocabulary size for the ETM and LDA on the *20NewsGroup* corpus. The perplexity is normalized by the size of the vocabulary. While the performance of LDA deteriorates for large vocabularies, the ETM maintains good performance.

# Embedding Items from Consumer Choice

# Contextual Data

More generally, embeddings model data given its context.

This idea extends well beyond text.

One recent area of application is to consumer choice data.

'word' is replaced by an item that a consumer is observed to purchase.

'context' is replaced by the other items a consumer has bought.

See auxiliary slides for illustration to movie data.

# SHOPPER Model

[Ruiz et al., 2020] builds a probabilistic model of consumer choice that incorporates context.

Basic choice probability is

$$\Pr\left[\text{item } c \mid \text{items in basket}\right] \propto \exp\left\{\boldsymbol{\theta}_u^T \boldsymbol{\alpha}_c + \boldsymbol{\rho}_c \left(\frac{1}{i-1} \sum_{j=1}^{i-1} \boldsymbol{\alpha}_{y_j}\right)\right\}$$

Extended choice probability incorporates prices, seasonal effects.

Likelihood formed by summing over unobserved ordering of product choices.

Model competitive for predicting how consumer choice reacts to price changes.
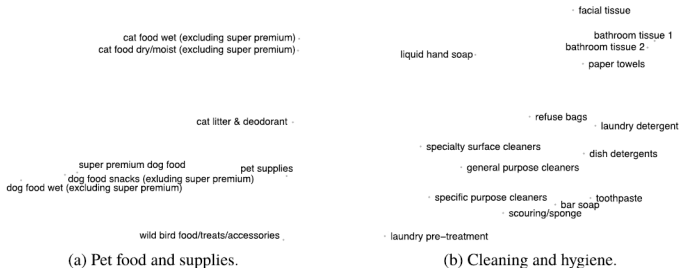
(a) Pet food and supplies.

(b) Cleaning and hygiene.

FIG. 3. *Two regions of the two-dimensional T-SNE projection of the features vectors $\alpha_c$ for the category-level experiment.*

TABLE 9
*Items with the highest complementarity and lowest exchangeability metrics for some query items*

| Query items | Complementarity score | | Exchangeability score | |
|---|---|---|---|---|
| Mission tortilla | 2.40 | taco bell taco seasoning mix | 0.05 | mission fajita size |
| soft taco 1 | 2.26 | mcrmck seasoning mix taco | 0.07 | mission tortilla soft taco 2 |
| | 2.24 | lawrys taco seasoning mix | 0.13 | mission tortilla fluffy gordita |
| | | | | |
| Private brand | 2.99 | bp franks meat | 0.11 | ball park buns hot dog |
| hot dog buns | 2.63 | bp franks bun size | 0.13 | private brand hotdog buns potato 1 |
| | 2.37 | bp franks beed bun length | 0.15 | private brand hotdog buns potato 2 |
| | | | | |
| Private brand mustard | 0.50 | private brand hot dog buns | 0.15 | frenchs mustard classic yellow squeeze |
| squeeze bottle | 0.41 | private brand cutlery full size forks | 0.16 | frenchs mustard classic yellow squeezed |
| | 0.24 | best foods mayonnaise squeeze | 0.21 | heinz ketchup squeeze bottle |
| | | | | |
| Private brand napkins | 0.78 | private brand selection plates 6 7/8 in | 0.09 | vnty fair napkins all occasion 1 |
| all occasion | 0.50 | private brand selection plates 8 3/4 in | 0.11 | vnty fair napkins all occasion 2 |
| | 0.49 | private brand cutlery full size forks | 0.12 | private brand selection premium napkins |

# Conclusion

Embedding models for language learn vector representations for words that depend on local context.

Word2vec was an important milestone for demonstrating how to estimate a neural language model that produced semantically coherent embeddings.

Natural next step: embedding text sequences.

# References I

Ash, E., Chen, D. L., and Ornaghi, A. (2020).
Gender attitudes in the judiciary : Evidence from U.S. circuit courts.
https://warwick.ac.uk/fac/soc/economics/research/workingpapers/2020/twerp_1256_-_ornaghi.pdf.

Bloom, N., Hassan, T. A., Kalyani, A., Lerner, J., and Tahoun, A. (2021).
The Diffusion of Disruptive Technologies.
Working Paper 28999, National Bureau of Economic Research.

Caliskan, A., Bryson, J. J., and Narayanan, A. (2017).
Semantics derived automatically from language corpora contain human-like biases.
Science, 356(6334):183–186.

Davis, S. J., Hansen, S., and Seminario-Amez, C. (2020).
Firm-Level Risk Exposures and Stock Returns in the Wake of COVID-19.
Working Paper 27867, National Bureau of Economic Research.

Dieng, A. B., Ruiz, F. J. R., and Blei, D. M. (2020).
Topic Modeling in Embedding Spaces.
Transactions of the Association for Computational Linguistics, 8:439–453.

# References II

Gennaro, G. and Ash, E. (2022).
Emotion and Reason in Political Language.
*The Economic Journal*, 132(643):1037–1059.

Hanley, K. W. and Hoberg, G. (2019).
Dynamic Interpretation of Emerging Risks in the Financial Sector.
*The Review of Financial Studies*, 32(12):4543–4603.

Kozlowski, A. C., Taddy, M., and Evans, J. A. (2019).
The Geometry of Culture: Analyzing the Meanings of Class through Word
Embeddings.
*American Sociological Review*, 84(5):905–949.

Li, K., Mai, F., Shen, R., and Yan, X. (2021).
Measuring Corporate Culture Using Machine Learning.
*The Review of Financial Studies*, 34(7):3265–3315.

Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013a).
Efficient Estimation of Word Representations in Vector Space.
arXiv:1301.3781 [cs].

# References III

Mikolov, T., Sutskever, I., Chen, K., Corrado, G., and Dean, J. (2013b).
Distributed Representations of Words and Phrases and their Compositionality.
arXiv:1310.4546 [cs, stat].

Pennington, J., Socher, R., and Manning, C. (2014).
GloVe: Global Vectors for Word Representation.
In Proceedings of the 2014 Conference on Empirical Methods in
Natural Language Processing (EMNLP), pages 1532–1543, Doha, Qatar. Association
for Computational Linguistics.

Ruiz, F. J. R., Athey, S., and Blei, D. M. (2020).
SHOPPER: A probabilistic model of consumer choice with substitutes and
complements.
The Annals of Applied Statistics, 14(1):1–27.