#### Word Embeddings

What works, what doesn't, and how to tell the difference for applied research

Pedro L. Rodríguez Arthur Spirling

Vanderbilt University New York University

May 14, 2020

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

We spent thousands of dollars running hundreds of embedding models and performing thousands of human validation tasks to get these main takeaways so you don't have to

#### Pedro L. Rodríguez Arthur Spirling

Vanderbilt University New York University

May 14, 2020

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで



◆□▶ ◆□▶ ◆ 臣▶ ◆ 臣▶ ○ 臣 ○ の Q @





Which one of these portraits is more realistic?

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三 のへぐ

	WELFARE
dependency	reform
Select the best candidate context word for the cue word provided by clicking on the respective checkbox below the word.	

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

(4日) (個) (目) (目) (目) (の)()

Explosion of interest in word embeddings.

Explosion of interest in word embeddings. These are real valued vectors that are used for two purposes:

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ □ のへぐ

Explosion of interest in word embeddings. These are real valued vectors that are used for two purposes:

1. feature representations for downstream NLP/ML tasks.

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Explosion of interest in word embeddings. These are real valued vectors that are used for two purposes:

- 1. feature representations for downstream NLP/ML tasks.
- 2. tools for studying word usage and meaning---"semantics".

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ □ のへぐ

Explosion of interest in word embeddings. These are real valued vectors that are used for two purposes:

- 1. feature representations for downstream NLP/ML tasks.
- 2. tools for studying word usage and meaning-"semantics".

As with all such representational strategies (topic models, TF-IDF), there are (literally thousands) of modeling options available...

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・

Explosion of interest in word embeddings. These are real valued vectors that are used for two purposes:

- 1. feature representations for downstream NLP/ML tasks.
- 2. tools for studying word usage and meaning-"semantics".

As with all such representational strategies (topic models, TF-IDF), there are (literally thousands) of modeling options available...

How should political scientists choose among them?

(日) (日) (日) (日) (日) (日) (日)

## Why do we ask?

How should political scientists choose among them?

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

# Some Background

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

### Firth quote, so you know we're legit

Firth quote, so you know we're legit

#### "You shall know a word by the company it keeps."

(Firth, 1957)





I had a cup of coffee this morning

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ



#### I had a cup of coffee this morning

◆□▶ ◆□▶ ◆ 臣▶ ◆ 臣▶ ○ 臣 ○ の Q @



I had a cup of **coffee** this morning

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ



◆□▶ ◆□▶ ◆三▶ ◆三▶ ○○○



▲ロト ▲御 ト ▲ 臣 ト ▲ 臣 ・ 今 Q ()・



I had a cup of tea this afternoon

◆□▶ ◆□▶ ◆ 臣▶ ◆ 臣▶ ○ 臣 ○ の Q @



▲ロト ▲御 ト ▲ 臣 ト ▲ 臣 ・ 今 Q ()・



・ロト・西・・田・・田・・日・

# Some Foreground

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ



 $\dots$  automatically generate a dictionary from unlabeled corpora (Hamilton et al, 2016)



 $\ldots$  automatically generate a dictionary from unlabeled corpora (Hamilton et al, 2016)

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

... model ideology of parliamentarians (Rheult & Cochrane, 2019)

 $\ldots$  automatically generate a dictionary from unlabeled corpora (Hamilton et al, 2016)

... model ideology of parliamentarians (Rheult & Cochrane, 2019)

... improve performance of readme (Jerzak et al, 2018)

 $\ldots$  automatically generate a dictionary from unlabeled corpora (Hamilton et al, 2016)

... model ideology of parliamentarians (Rheult & Cochrane, 2019)

... improve performance of readme (Jerzak et al, 2018) ... understand how words have changed meaning over time (Rodman, 2019)

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

 $\ldots$  automatically generate a dictionary from unlabeled corpora (Hamilton et al, 2016)

... model ideology of parliamentarians (Rheult & Cochrane, 2019)

... improve performance of readme (Jerzak et al, 2018) ... understand how words have changed meaning over time (Rodman, 2019)

 $\rightarrow$  choose a model (Word2Vec, GloVe),

 $\ldots$  automatically generate a dictionary from unlabeled corpora (Hamilton et al, 2016)

... model ideology of parliamentarians (Rheult & Cochrane, 2019)

... improve performance of readme (Jerzak et al, 2018) ... understand how words have changed meaning over time (Rodman, 2019)

 $\rightarrow$  choose a model (Word2Vec, GloVe), an architecture within that model (CBOW, Skipgram),

 $\ldots$  automatically generate a dictionary from unlabeled corpora (Hamilton et al, 2016)

... model ideology of parliamentarians (Rheult & Cochrane, 2019)

... improve performance of readme (Jerzak et al, 2018) ... understand how words have changed meaning over time (Rodman, 2019)

 $\rightarrow$  choose a model (Word2Vec, GloVe), an architecture within that model (CBOW, Skipgram), parameters within that architecture (window size, embedding length)

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・

 $\ldots$  automatically generate a dictionary from unlabeled corpora (Hamilton et al, 2016)

... model ideology of parliamentarians (Rheult & Cochrane, 2019)

... improve performance of readme (Jerzak et al, 2018) ... understand how words have changed meaning over time (Rodman, 2019)

 $\rightarrow$  choose a model (Word2Vec, GloVe), an architecture within that model (CBOW, Skipgram), parameters within that architecture (window size, embedding length) and a training set (pretrained, locally fit).

choose a model (Word2Vec, GloVe), an architecture within that model (CBOW, Skipgram), parameters within that architecture (window size, embedding length) and a training set (pretrained, locally fit).

#### The Problem

(ロ)、(型)、(E)、(E)、 E) の(()


There are **no** generally accepted downstream tasks in political science:



There are **no** generally accepted downstream tasks in political science: 'extrinsic' evaluation criteria make no sense (e.g. analogy banks, learner accuracy).

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三 のへぐ

There are **no** generally accepted downstream tasks in political science: 'extrinsic' evaluation criteria make no sense (e.g. analogy banks, learner accuracy).

So, embeddings are **only** useful to the extent they capture semantically meaningful information about politics:

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

There are **no** generally accepted downstream tasks in political science: 'extrinsic' evaluation criteria make no sense (e.g. analogy banks, learner accuracy).

So, embeddings are **only** useful to the extent they capture semantically meaningful information about politics: focus on 'intrinsic' evaluation criteria.

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

There are **no** generally accepted downstream tasks in political science: 'extrinsic' evaluation criteria make no sense (e.g. analogy banks, learner accuracy).

So, embeddings are **only** useful to the extent they capture semantically meaningful information about politics: focus on 'intrinsic' evaluation criteria.

But how can we evaluate this?

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三 のへぐ

# The Solution



We propose a "Turing test":





We propose a "Turing test": ask crowdworkers whether output from humans or machine (model) fits a cue better.

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ



We propose a "Turing test": ask crowdworkers whether output from humans or machine (model) fits a cue better.

We get remarkable, human-like performance from embeddings models in terms of meaning.  $\checkmark$ 

Embeddings models have multiple parameters,

Embeddings models have multiple parameters, especially window-size and embedding dimensions.

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ □ のへぐ

Embeddings models have multiple parameters, especially window-size and embedding dimensions.

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

And should you use pretrained or locally fit?

Embeddings models have multiple parameters, especially window-size and embedding dimensions.

And should you use pretrained or locally fit?

We use our technical critera on fit and stability and our Turing test to provide advice.

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Embeddings models have multiple parameters, especially window-size and embedding dimensions.

And should you use pretrained or locally fit?

We use our technical critera on fit and stability and our Turing test to provide advice.

Avoid small windows, few dimensions but otherwise results are robust to these parameter choices.  $\checkmark$ 

Embeddings models have multiple parameters, especially window-size and embedding dimensions.

And should you use pretrained or locally fit?

We use our technical critera on fit and stability and our Turing test to provide advice.

Avoid small windows, few dimensions but otherwise results are robust to these parameter choices.  $\checkmark$ 

Pretrained embeddings work about as well as anything else.  $\checkmark$ 

# Implementing Choices

▲ロト ▲御 ▶ ▲ 臣 ▶ ▲ 臣 ▶ ● 臣 ● のへで

# Implementing Choices

We train a series of (local) models.

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

We focus on two hyperparameter choices (25 combinations):

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

We focus on two hyperparameter choices (25 combinations):

window-size —1, 6, 12, 24 and 48



We focus on two hyperparameter choices (25 combinations):

window-size —1, 6, 12, 24 and 48

We focus on two hyperparameter choices (25 combinations):

window-size —1, 6, 12, 24 and 48

For each pair we estimate 10 sets of embeddings (250 in all).

We focus on two hyperparameter choices (25 combinations):

window-size —1, 6, 12, 24 and 48

embedding dimension —50, 100, 200, 300 and 450

For each pair we estimate 10 sets of embeddings (250 in all).

Also compare to GloVe and Word2Vec (skip-gram) pretrained.

technical criteria: model loss and computation time.

model variance (stability): within-model Pearson correlation of nearest neighbor rankings across multiple initializations.

query search ranking correlation: Pearson and rank correlations of cosine similarities.

human preference: a "Turing test" assessment and rank deviations from human generated lists.

100 200 300 400 Embedding Dimensions

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ



・ロト・西・・田・・田・・日・



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ ○臣 - の々ぐ



▲ロト ▲園 ト ▲ 臣 ト ▲ 臣 ト 一臣 - のへ(で)



◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ



**Embedding Dimensions** 

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで



**Embedding Dimensions** 

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ● □ ● ● ● ●

## Technical Criteria: Computation Time



◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

technical criteria: model loss and computation time.

model variance (stability): within-model Pearson correlation of nearest neighbor rankings across multiple initializations.

query search correlations: Pearson and rank correlations of cosine similarities.

human preference: a "Turing test" assessment and rank deviations from human generated lists.

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三 のへぐ

technical criteria: model loss and computation time.

model variance (stability): within-model Pearson correlation of nearest neighbor rankings across multiple initializations.

query search correlations: Pearson and rank correlations of cosine similarities.

human preference: a "Turing test" assessment and rank deviations from human generated lists.

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00



◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

technical criteria: model loss and computation time.

model variance (stability): within-model Pearson correlation of nearest neighbor rankings across multiple initializations.

query search correlations: Pearson and rank correlations of cosine similarities.

human preference: a "Turing test" assessment and rank deviations from human generated lists.

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00


▲□▶ ▲□▶ ▲ □▶ ▲ □ ▶ ▲ □ ● ● の Q ()



▲□▶ ▲□▶ ▲ □▶ ▲ □▶ □ のへぐ



▲□▶ ▲□▶ ▲ □▶ ▲ □▶ □ のへぐ



▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

technical criteria: model loss and computation time.

model variance (stability): within-model Pearson correlation of nearest neighbor rankings across multiple initializations.

query search correlations: Pearson and rank correlations of cosine similarities.

human preference: a "Turing test" assessment and rank deviations from human generated lists.

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

# RShiny App: Definitions

#### **Context Words**

A famous maxim in the study of linguistics states that:

You shall know a word by the company it keeps. (Firth, 1957)

This task is designed to help us understand the nature of the "company" that words "keep": that is, their CONTEXT.

Specifically, for a CUE WORD, its CONTEXT WORDS include words that:

 Tend to occur in the vicinity of the CUE WORD. That is, they are words that appear close to the CUE WORD in written or spoken language.

#### AND/OR

Tend to occur in similar situations to the CUE WORD in spoken and written language. That is, they are words that
regularly appear with other words that are closely related to the CUE WORD.

For example, CONTEXT WORDS for the cue word COFFEE include:

- 1. cup (tends to occur in the vicinity of COFFEE).
- 2. tea (tends to occur in similar situations to COFFEE, for example when discussing drinks).

Click "Next" to continue

Next

#### RShiny App-1: Context Word Generation

#### Task 3 of 10

#### welfare

Click here to enter text

Press enter to save entry.

- reform - help - poor

Number of unique words entered: 3 Number of words required to satisfy minimum: 7 Time remaining: 156 secs

Please input at least 10 context words before clicking "Next".

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ のQ@

#### Next



Select the best candidate context word for the cue word provided by clicking on the respective checkbox below the word.

Click "Next" to continue





Select the best candidate context word for the cue word provided by clicking on the respective checkbox below the word.





▲□▶▲□▶▲≡▶▲≡▶ ≡ めぬぐ



Select the best candidate context word for the cue word provided by clicking on the respective checkbox below the word.



▲□▶▲□▶▲≡▶▲≡▶ ≡ めぬぐ



Select the best candidate context word for the cue word provided by clicking on the respective checkbox below the word.



▲□▶▲□▶▲≡▶▲≡▶ ≡ めぬぐ

Machine: 1 - Human: 0

immigration equality taxes freedom democrat justice welfare democracy republican



◆□ > ◆□ > ◆豆 > ◆豆 > ̄豆 = のへで



▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三 のへぐ



▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで







### Candidate: GloVe pretrained Baseline: Human



### Candidate: GloVe pretrained Baseline: W2V pretrained



▲□▶ ▲圖▶ ▲臣▶ ▲臣▶ 三臣 - のへの



We get remarkable, human-like performance from embeddings models in terms of meaning.  $\checkmark$ 

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ □ のへぐ

We get remarkable, human-like performance from embeddings models in terms of meaning.  $\checkmark$ 

Avoid small windows, few dimensions but otherwise results are robust to these parameter choices.  $\checkmark$ 

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

We get remarkable, human-like performance from embeddings models in terms of meaning.  $\checkmark$ 

Avoid small windows, few dimensions but otherwise results are robust to these parameter choices.  $\checkmark$ 

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

Pretrained, 'default' embeddings work about as well as anything else.  $\checkmark$ 

We get remarkable, human-like performance from embeddings models in terms of meaning.  $\checkmark$ 

Avoid small windows, few dimensions but otherwise results are robust to these parameter choices.  $\checkmark$ 

Pretrained, 'default' embeddings work about as well as anything else.  $\checkmark$ 

GloVe appears to be more robust than Word2vec, but both are equally liked by humans.  $\checkmark$ 

#### Thank you!



R software package to apply our proposed metrics and framework: coming soon.

GitHub: http://github.com/ArthurSpirling/EmbeddingsPaper

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへの

#### Paper:

../Paper/Embeddings\_SpirlingRodriguez.pdf

FAQ:
../Project\_FAQ/faq.md