# Seven Deadly Sins

# A diachronic analysis

with a geometric comparison of word embedding models

Data Semantics exam project
July 2022

Gaetano Chiriaco 882638
Riccardo Porcedda  886719
Gianmarco Russo 887277

# Introduction

*Project goals:*

1.  Diachronic analysis of the meanings associated with the seven deadly sins.
    Are these concept still relevant today?
    What was every sin related to in the past and what about present days?

2.  *CADE* is an effective tool to compare corpora from different contexts (topical or temporal), but it relies on only *Word2Vec* embedding model: is it sufficient?
    Would an hypothetical *GloVe* implementation produce different results?

# Data

The Corpus used for the analysis is the COHA (Corpus Of Historical American english). For the sake of our study(and due to computational limitations) we only selected a subset of it:

- *Past*: documents from 1810-1899 ~ 11000 documents
- *Post*: documents from 2000-2009 ~ 14000 documents

Leaving the 20th century as a transition period for the language to change and evolve.

# Data Loading e Preprocessing

The documents have been read and loaded into two different datasets.
Then we applied the following preprocessing operations:

- *lower case*
- *stopword removal*
- *lemmatization*
- *splitting the texts into individual sentences*

# Word Phrases

We also considered recognition of *word phrases*, words that must go together to express their real meaning: *deadly_sins, heavenly_virtue, New_York, user_friendly*, etc…

We used the gensim built-in method *Phrases* which automatically detects bigrams.

The results didn't improve the performance of our models, so we decided to not include this part in the final work.

# Word2Vec & CADE

We used *CADE*[1] for the alignment of the corpora, so we created the compass and then then trained two slices (past and post) with *Word2Vec*[2]. This time we chose the continuous bag of word method (*CBOW*). Parameters used for *w2v*:

- *ncomponents* = 100,
- *window size* = 5, we also tried with window=3 and 1 but we chose the model that gave the best results.
- *min count* = 10,
- *epochs* =10.

# LUST (CBOW)

## COHA (1810-1899)

**Most similar words:**
- "greed" (0.84)
- "avarice" (0.83)
- "insatiable" (0.78)
- "covetousness" (0.77)
- "pander" (0.75)

**Most similar to Lust (2000) in 1800 COHA:**
- "passion" (0.69)
- "lust" (0.66)
- "lustful" (0.60)
- "voluptuary" (0.60)
- "passionate" (0.60)

**Analogy example:**
- "lust" – "insatiate" + "proud"=
- prouder (0.69)
- pride (0.68)
- proudest (0.65)

## COHA (2000-2009)

**Most similar to Lust (1800) in 2000 COHA:**
- "lust" (0.66)
- "greed" (0.64)
- "vengeance" (0.58)
- "avarice" (0.55)
- "jealousy" (0.51)

**Most similar words:**
- "jealousy" (0.72)
- "yearn" (0.68)
- "greed" (0.67)
- "desire" (0.66)
- "longing" (0.66)

**Analogy example:**
- "lust" – "yearn" + "smite"=
- "wrath" (0.59)
- "almighty" (0.58)
- "sinner" (0.57)

# GREED (CBOW)

## COHA (1810-1899)

**Most similar words:**
- "avarice" (0.90)
- "lust" (0.84)
- "cupidity" (0.82)
- "rapacity" (0.82)
- "covetousness" (0.81)

**Most similar to Greed(2000) in 1800 COHA:**
- "greed" (0.70)
- "rapacity" (0.69)
- "avarice" (0.68)
- "inhumanity" (0.65)
- "self-seeking" (0.64)

**Analogy example:**
- "greed" - "desire" + "power"=
- cupidity (0.69)
- lust (0.68)
- avarice (0.65)

## COHA (2000-2009)

**Most similar to Greed(1800) in 2000 COHA:**
- "greed" (0.70)
- "avarice" (0.59)
- "greedy" (0.53)
- "lust" (0.53)
- "dishonesty" (0.52)

**Most similar words:**
- "cruelty" (0.77)
- "stupidity" (0.77)
- "cowardice" (0.76)
- "laziness" (0.76)
- "ignorance" (0.75)

**Analogy example:**
- "greed" - "bad" + "good"=
- generosity (0.70)
- virtue (0.69)
- boundless (0.65)

# GLUTTONY (CBOW)

## COHA (1810-1899)

**Most similar words:**
- "drunkenness" (0.77)
- "sensuality" (0.76)
- "debauchery" (0.74)
- "lewdness" (0.73)
- "intemperance" (0.73)

**Most similar to Gluttony(2000) in 1800 COHA:**
- "long-indulged" (0.70)
- "abnormity" (0.69)
- "self-defeating" (0.68)
- "eroticism" (0.65)
- "book-buying" (0.64)

**Analogy example:**
- "Gluttony" – "Eat" + "Wrath"=
- malignity (0.71)
- rage (0.68)
- ire (0.66)

## COHA (2000-2009)

**Most similar to Gluttony(1800) in 2000 COHA:**
- "binge" (0.60)
- "debauchery" (0.58)
- "promiscuity" (0.57)
- "overindulgence" (0.55)
- "vile" (0.55)

**Most similar words:**
- "avowal" (0.58)
- "irreverence" (0.58)
- "believe" (0.57)
- "puerile" (0.57)
- "impulsive" (0.57)

**Analogy example:**
- "Gluttony" – "food" + "desire"=
- unshakable (0.63)
- lust (0.63)
- betrayal (0.62)

# ENVY (CBOW)

## COHA (1810-1899)

**Most similar words:**
- "envious" (0.71)
- "vanity" (0.62)
- "flatters" (0.61)
- "pride" (0.61)
- "jealous" (0.61)

**Most similar to Envy(2000) in 1800 COHA:**
- "idolize" (0.54)
- "vied" (0.53)
- "infatuate" (0.48)
- "despises" (0.48)
- "envy" (0.47)

**Analogy example:**
- "envy" – "jealousy" + "pride"=
- proud(0.71)
- scorn(0.68)
- proudest(0.66)

## COHA (2000-2009)

**Most similar to Envy(1800) in 2000 COHA:**
- "jealousy" (0.48)
- "envy" (0.47)
- "swoon" (0.46)
- "lust" (0.43)
- "delight" (0.43)

**Most similar words:**
- "loathe" (0.71)
- "admiration" (0.58)
- "jealousy" (0.57)
- "jealous" (0.57)
- "despise" (0.57)

**Analogy example:**
- "envy" – "bad" + "good"=
- admire (0.68)
- admiration (0.66)
- devotion (0.66)

# SLOTH (CBOW)

## COHA (1810-1899)

**Most similar words:**
- "indolence" (0.71)
- "sensuality" (0.62)
- "self-indulgence" (0.61)
- "gluttony" (0.61)
- "sluggishness" (0.61)

**Most similar to Sloth(2000) in 1800 COHA:**
- "scourg" (0.57)
- "duta" (0.57)
- "big-nosed" (0.56)
- "likho" (0.55)
- "capra" (0.54)

**Analogy example:**
- "sloth" - "sleep" + "eat"=
- gluttony (0.64)
- glutton (0.61)
- pamper (0.60)

## COHA (2000-2009)

**Most similar to Sloth(1800) in 2000 COHA:**
- "greed" (0.48)
- "boredom" (0.47)
- "mindless" (0.46)
- "filth" (0.43)
- "laziness" (0.43)

**Most similar words:**
- "yellow-eyed" (0.63)
- "half-horse" (0.63)
- "gambol" (0.62)
- "halfman" (0.61)
- "pomposity" (0.60)

**Analogy example:**
- "sloth" - "sleep" + "eat"=
- fajitas (0.60)
- mutton (0.59)
- goat (0.58)

# WRATH (CBOW)

## COHA (1810-1899)

**Most similar words:**
- "ire" (0.80)
- "anger" (0.79)
- "rage" (0.78)
- "fury" (0.76)
- "vengeance" (0.72)

**Most similar to Wrath(2000) to 1800 COHA:**
- "korah" (0.66)
- "euph" (0.66)
- "righteous" (0.65)
- "wrath" (0.61)
- "blaspheme" (0.61)

**Analogy example:**
- "wrath" – "anger" + "pride"=
- "glory" (0.65)
- "proud" (0.64)
- "ambition" (0.63)

## COHA (2000-2009)

**Most similar to Wrath(1800) to 2000 COHA:**
- "fury" (0.71)
- "rage" (0.70)
- "anger" (0.62)
- "wrath" (0.61)
- "vengeance" (0.58)

**Most similar words:**
- "vengeance" (0.74)
- "vengeful" (0.68)
- "smite" (0.67)
- "unto" (0.66)
- "righteous" (0.66)

**Analogy example:**
- "wrath" – "sin" + "beast"=
- "ferocious" (0.60)
- "vulture" (0.58)
- "roc" (0.58)

# PRIDE (CBOW)

## COHA (1810-1899)

**Most similar words:**
- "proud" (0.75)
- "self-esteem" (0.73)
- "vanity" (0.73)
- "arrogance" (0.70)
- "ambition" (0.68)

**Most similar to Pride(2000) to 1800 COHA:**
- "pride" (0.59)
- "loyalty" (0.58)
- "disinterestedness" (0.54)
- "patriotism" (0.54)
- "self-devotion" (0.53)

**Analogy example:**
- "Pride" – "sin" + "nation"=
- patriotism (0.57)
- nationality (0.57)
- emulous (0.56)

## COHA (2000-2009)

**Most similar to Pride(1800) to 2000 COHA:**
- "pride" (0.59)
- "arrogance" (0.50)
- "jealousy" (0.48)
- "stubbornness" (0.47)
- "lust" (0.46)

**Most similar words:**
- "bravery" (0.66)
- "loyalty" (0.64)
- "gratitude" (0.61)
- "courage" (0.61)
- "indignation" (0.60)

**Analogy example:**
- "pride" – "sin" + "sport"=
- "rugby" (0.60)
- "athletic" (0.60)
- "soccer" (0.58)

# Graphical representation using t-SNE

*t-SNE algorithm*[3] is often used to produce a 2D or 3D graphical representation of k-dimensional observations.
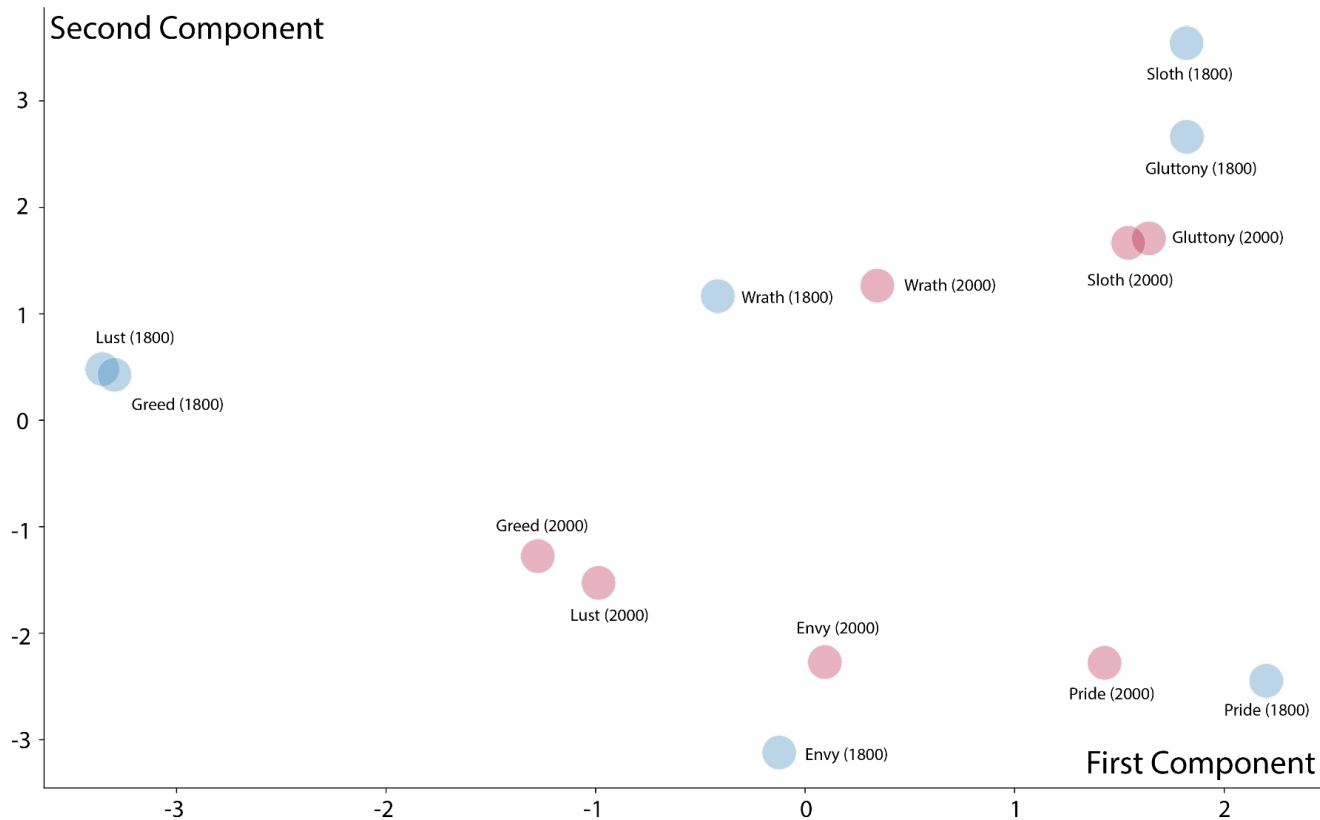
The aim of the *SNE algorithms* is to preserve the input distances in the output space.

We used this technique to obtain a scatter plot of the 14 sins word embeddings.
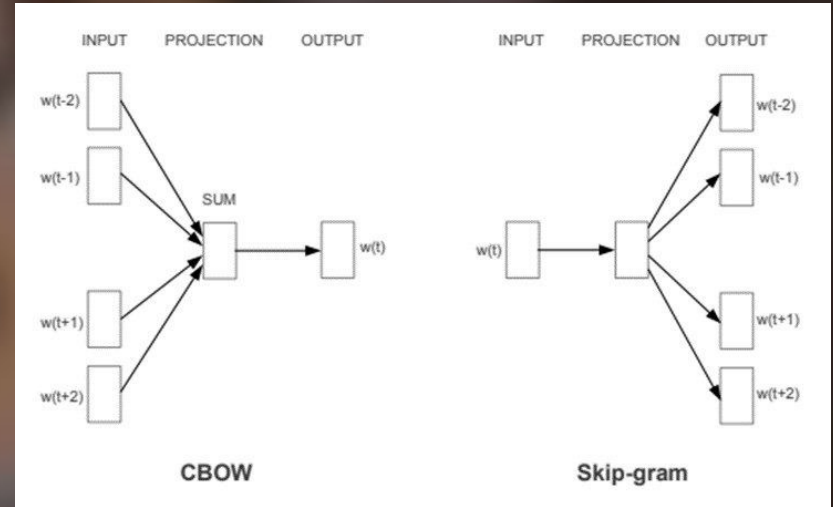
***Criticalities:***

- *perplexity* and *number of iteration* are chosen arbitrarily and can have a crucial impact of the output.
- If the intrinsic dimensionality of the input data is not 2D or 3D, this type of representation doesn't make sense.

# T-SNE visualization results on W2V CBOW

# Word2Vec Skip-gram

We repeated the experiment with the same parameters (windows size, min count, etc..) but this time we used the *Skip-gram* method. This method usually gives better representations for rarely used word compared to *CBOW*. The parameters settings are the same as *CBOW*.

# LUST (SKIP-GRAM)

## COHA (1810-1899)

**Most similar words:**
- "avarice" (0.78)
- "covetousness" (0.75)
- "greed" (0.75)
- "bestial" (0.75)
- "insatiate" (0.75)

**Most similar to Lust (2000) in 1800 COHA:**
- "lust" (0.67)
- "passion" (0.62)
- "self-worship" (0.62)
- "concupiscence" (0.62)
- "all-consuming" (0.61)

**Analogy example:**
- "lust" - "insatiate" + "proud"=
- pride (0.67)
- prouder(0.64)
- haughty(0.63)

## COHA (2000-2009)

**Most similar to Lust (1800) in 2000 COHA:**
- "lust" (0.67)
- "greed" (0.67)
- "avarice" (0.65)
- "idolatrous" (0.65)
- "unquenchable" (0.63)

**Most similar words:**
- "erotic" (0.64)
- "arouse" (0.63)
- "seduction" (0.61)
- "bloodlust" (0.60)
- "desire" (0.60)

**Analogy example:**
- "lust" - "yearn" + "smite"=
- "smote" (0.59)
- "smites" (0.58)
- "whosoever" (0.57)

# GREED (SKIP-GRAM)

## COHA (1810-1899)

**Most similar words:**
- "avarice" (0.83)
- "lust" (0.75)
- "rapacity" (0.75)
- "covetousness" (0.75)
- "niggardliness" (0.75)

**Most similar to Greed(2000) in 1800 COHA:**
- "greed" (0.70)
- "avarice" (0.69)
- "niggardliness" (0.67)
- "lust" (0.67)
- "un-godliness" (0.67)

**Analogy example:**
- "greed" – "desire" + "power"=
- "avarice" (0.66)
- "lust" (0.64)
- "cupidity" (0.63)

## COHA (2000-2009)

**Most similar to Greed(1800) in 2000 COHA:**
- "greed" (0.70)
- "avarice" (0.63)
- "rapacious" (0.62)
- "hunger" (0.58)
- "thieve" (0.58)

**Most similar words:**
- "avarice" (0.71)
- "unutterable" (0.67)
- "epochal" (0.67)
- "cowardice" (0.66)
- "laziness" (0.66)

**Analogy example:**
- "greed" – "bad" + "good"=
- "virtue" (0.70)
- "rectitude" (0.69)
- "cleverness" (0.65)

# GLUTTONY (SKIP-GRAM)

## COHA (1810-1899)

**Most similar words:**
- "wine-bibbing" (0.76)
- "lewdness" (0.72)
- "intemperance" (0.71)
- "sensuality" (0.71)
- "greediness" (0.70)

**Most similar to Gluttony(2000) in 1800 COHA:**
- "enthrals" (0.77)
- "selfannihilation" (0.76)
- "concupiscence" (0.76)
- "unchastised" (0.75)
- "enthu" (0.75)

**Analogy example:**
- "gluttony" - "eat" + "wrath"=
- "obduracy" (0.65)
- "rage" (0.65)
- "wreaks" (0.65)

## COHA (2000-2009)

**Most similar to Gluttony(1800) in 2000 COHA:**
- "overindulgence" (0.70)
- "promiscuity" (0.70)
- "debauchery" (0.67)
- "scrofula" (0.64)
- "impotence" (0.64)

**Most similar words:**
- "abasement" (0.83)
- "largehearted" (0.82)
- "unquenchable" (0.81)
- "wholehearted" (0.81)
- "inexpressible" (0.81)

**Analogy example:**
- "gluttony" - "food" + "desire"=
- "unthinking" (0.69)
- "masochism" (0.67)
- "ufe" (0.67)

# ENVY (SKIP-GRAM)

## COHA (1810-1899)

**Most similar words:**
- "envious" (0.75)
- "jealousy" (0.72)
- "hate" (0.67)
- "vanity" (0.66)
- "malice" (0.66)

**Most similar to Envy(2000) in 1800 COHA:**
- "envy" (0.68)
- "self-exaltation" (0.64)
- "dotings" (0.64)
- "masterpassion" (0.63)
- "hard-heartedness" (0.62)

**Analogy example:**
- "envy" – "jealousy" + "pride"=
- "proud" (0.79)
- "prouder" (0.68)
- "proudest" (0.64)

## COHA (2000-2009)

**Most similar to Envy(1800) in 2000 COHA:**
- "envy" (0.68)
- "conscientiousness" (0.62)
- "inveterate" (0.62)
- "covetousness" (0.60)
- "spunk" (0.59)

**Most similar words:**
- "despise" (0.67)
- "enmity" (0.64)
- "devotion" (0.64)
- "protectiveness" (0.63)
- "love" (0.63)

**Analogy example:**
- "envy" – "bad" + "good"=
- "admire" (0.65)
- "devotion" (0.64)
- "admiration" (0.64)

# SLOTH (SKIP-GRAM)

## COHA (1810-1899)

**Most similar words:**
- "indolence" (0.70)
- "brutishness" (0.68)
- "gluttony" (0.66)
- "sensuality" (0.66)
- "slothfulness" (0.65)

**Most similar to Sloth(2000) in 1800 COHA:**
- "two-footed" (0.57)
- "ourang" (0.57)
- "unmatchable" (0.56)
- "cameleopard" (0.55)
- "diddles" (0.54)

**Analogy example:**
- "sloth" - "sleep" + "eat"=
- "glutton" (0.64)
- "greediness" (0.63)
- "flesheating" (0.63)

## COHA (2000-2009)

**Most similar to Sloth(1800) in 2000 COHA:**
- "two-toed" (0.65)
- "sloth" (0.62)
- "douc" (0.58)
- "perfectionism" (0.58)
- "endemic" (0.58)

**Most similar words:**
- "half-horse" (0.75)
- "two-toed" (0.75)
- "half-man" (0.75)
- "tapir" (0.72)
- "dungheap" (0.70)

**Analogy example:**
- "sloth" - "sleep" + "eat"=
- "mollusk" (0.61)
- "blintz" (0.61)
- "carnivorous" (0.61)

# WRATH (SKIP-GRAM)

## COHA (1810-1899)

**Most similar words:**
- "anger" (0.82)
- "rage" (0.80)
- "ire" (0.80)
- "fury" (0.78)
- "indignation" (0.77)

**Most similar to Wrath(2000) to 1800 COHA:**
- "wrath" (0.63)
- "quarequa" (0.61)
- "injur" (0.61)
- "scorneth" (0.61)
- "despoiler" (0.61)

**Analogy example:**
- "wrath" – "anger" + "pride"=
- "oermastereth" (0.66)
- "proud" (0.66)
- "scorneth" (0.65)

## COHA (2000-2009)

**Most similar to Wrath(1800) to 2000 COHA:**
- "fury" (0.70)
- "rage" (0.68)
- "outrage" (0.66)
- "anger" (0.66)
- "ire" (0.65)

**Most similar words:**
- "philistine" (0.74)
- "smites" (0.68)
- "smite" (0.67)
- "amalekites" (0.66)
- "vengeance" (0.66)

**Analogy example:**
- "wrath" – "sin" + "beast"=
- "bellowing" (0.58)
- "intulo" (0.58)
- "dervish" (0.57)

# PRIDE (SKIP-GRAM)

## COHA (1810-1899)

**Most similar words:**
- "proud" (0.77)
- "vanity" (0.73)
- "resentment" (0.70)
- "oermastereth" (0.67)
- "haughtiness" (0.67)

**Most similar to Pride(2000) to 1800 COHA:**
- "pride" (0.74)
- "implacability" (0.57)
- "favouritism" (0.57)
- "vainglory" (0.57)
- "self-abnegating" (0.56)

**Analogy example:**
- "pride" - "sin" + "nation"=
- "empire" (0.57)
- "patriotism" (0.57)
- "warlike" (0.56)

## COHA (2000-2009)

**Most similar to Pride(1800) to 2000 COHA:**
- "pride" (0.74)
- "proud" (0.63)
- "coxcomb" (0.63)
- "bashful" (0.61)
- "deject" (0.61)

**Most similar words:**
- "protectiveness" (0.64)
- "proud" (0.63)
- "joy" (0.63)
- "indignation" (0.62)
- "grandiosity" (0.62)

**Analogy example:**
- "pride" - "sin" + "sport"=
- "intramural" (0.59)
- "collegiate" (0.57)
- "rugby" (0.56)

# T-SNE visualization results on W2V Skipgram

# Embedding models comparison: *GloVe & W2V*

Since *CADE* is based on *Word2Vec*, we wondered if an hypothetical *GloVe*[4] implementation would provide different results.

*GloVe* is a bit different from *w2v*: while the latter does an incremental and 'sparse' training of a neural network by repeatedly iterating over a training corpus, *GloVe* computes a global co-occurence matrix to better gather statistical informations of the analyzed text corpus.

Furthermore, *GloVe* combines the advantages of two learning methods: global matrix factorization like <u>latent semantic analysis</u> (LSA) and local context window method like *Skip-gram* (which performs better on analogies).

The *GloVe* technique has a simpler <u>least square</u> cost or error function that reduces the computational cost of training the model.

# Embedding models comparison: first approach

The first comparison is really straightforward: evaluating the "agreements" of the three models (*CBOW*, *Skip-gram* and *GloVe*) when computing the most similar words with respect to the same token.

In this case we considered checking the top-5 most similar using *cosine similarity*.

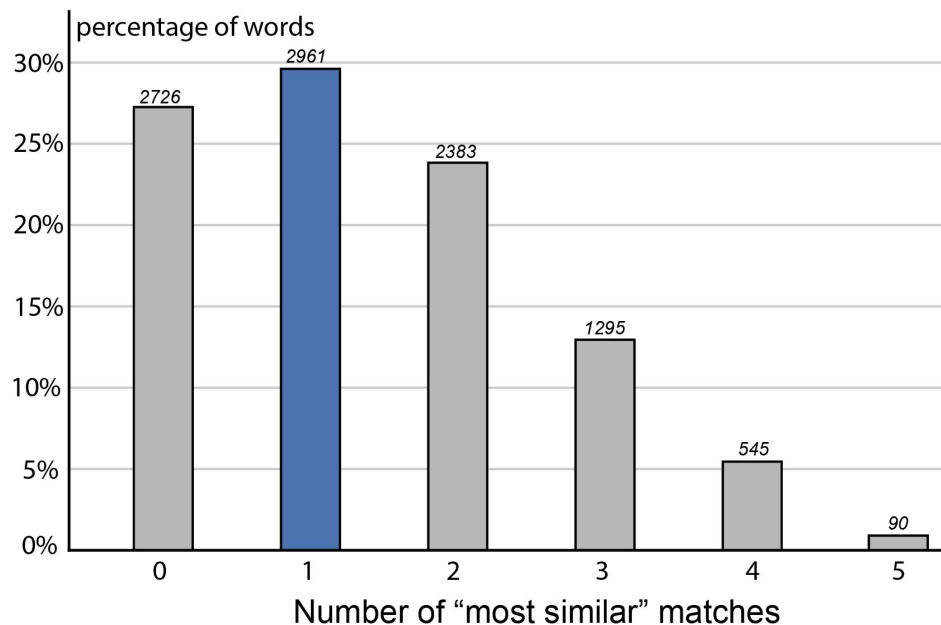# Comparison of Word2Vec CBOW word embeddings and Word2Vec Skipgram word embeddings (COHA 1810-1899)

*Tested on 10000 words of the corpus*

| Most similar word to "**uncle**" (CBOW) | ... "uncle" (Skipgram) |
|---|---|
| *1. aunt* | *1. mom* |
| *2. father* | *2. father* |
| 3. granpa | *3. aunt* |
| 4. granma | 4. son |
| *5. mom* | 5. cousin |

**3 words match**



The highest number of matches are those between *w2v CBOW* and *Skip-gram*, having just a 13% of no matches. We also got a little amount of perfect matches.

# Comparison of Word2Vec CBOW word embeddings and Word2Vec Skipgram word embeddings (COHA 2000-2009)
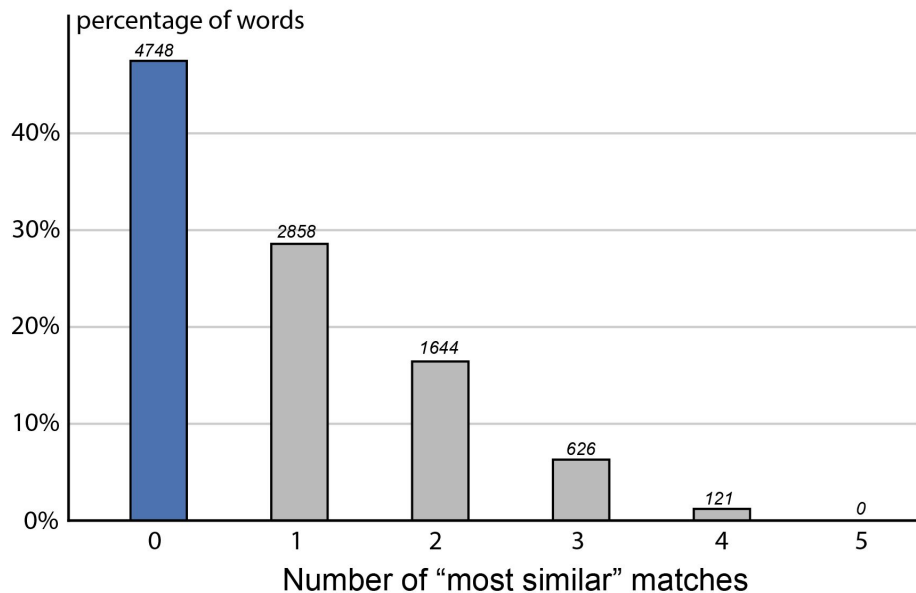
*Tested on 10000 words of the corpus*



The behaviour is similar to the "past" corpus, with the comparison between *w2v CBOW* and *Skip-gram* being the one with the highest number of agreements.

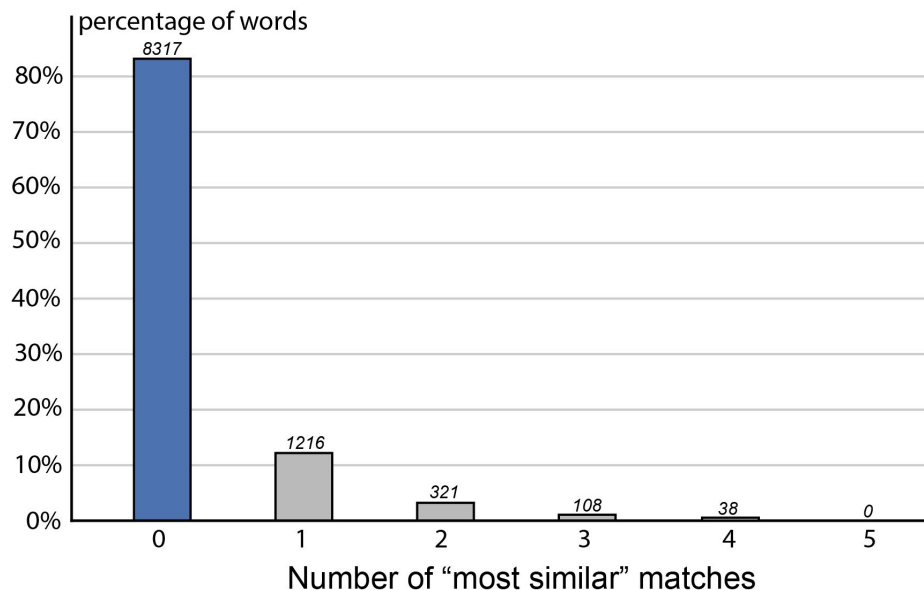**Comparison of Word2Vec CBOW word embeddings and GloVe word embeddings (COHA 1810-1899)**

*Tested on 10000 words of the corpus*

percentage of words

40%
35%
30%
25%
20%
15%
10%
5%
0%

4085
2970
2031
759
155
0

0   1   2   3   4   5

Number of "most similar" matches

**Comparison of Word2Vec CBOW word embeddings and GloVe word embeddings (COHA 2000-2009)**

*Tested on 10000 words of the corpus*

percentage of words

80%
70%
60%
50%
40%
30%
20%
10%
0%

7655
1685
495
121
44
0

0   1   2   3   4   5

Number of "most similar" matches

When we compare a *w2v* model with *GloVe* things change drastically: 5/5 matches are gone and almost half of the sample has 0 matches.

The results between *GloVe* and *w2v* models trained on 2000-2009 corporus show an even higher percentage of zero matches

**Comparison of Word2Vec Skipgram word embeddings and GloVe word embeddings (COHA 1810-1899)**

*Tested on 10000 words of the corpus*

percentage of words

(Bar chart with values: 4748, 2858, 1644, 626, 121, 0 over "Number of "most similar" matches" 0, 1, 2, 3, 4, 5)

**Comparison of Word2Vec Skipgram word embeddings and GloVe word embeddings (COHA 2000-2009)**

*Tested on 10000 words of the corpus*

percentage of words

(Bar chart with values: 8317, 1216, 321, 108, 38, 0 over "Number of "most similar" matches" 0, 1, 2, 3, 4, 5)

Again, same behaviour. We've observed a high disagreement rate between *GloVe* and *w2v* models.

# Notes on this approach

We have to point out that these high level of disagreement could be due to the fact that the *GloVe* library doesn't allow to set window size and min count as parameters, resulting in a model setup different than the one of *Word2Vec*.

For the same reason, in the second comparison approach we work with a subset of the word embeddings, considering only words that are in both models' vocabularies.

# Embedding models comparison: geometric approach

All comparisons we found in the literature are *performance-on-task based* comparisons: e.g. which model does perform better on analogies?[5,6]

We tried to better quantitatively assess how much similar these models are by considering geometric properties of the word embeddings.

For all models, we chose to fix the dimensions to 100.

# Geometric approach: first hypothesis

Both *Word2Vec* and *GloVe* are capable of finding analogies.

E.g. **king - man + woman ~ queen** <u>works in both models.</u>

Analogies are constructed as linear combinations of embedded words and many authors tried to theoretically explain the emergence of this linearity[6].

## Hypothesis n°1:
<u>Embedding spaces of *Word2Vec* and *GloVe* can be transformed one into the other via linear mapping</u>

# Linear mapping hypothesis

For any linear mapping f(<u>w</u>) of <u>w</u>, **additivity** is satisfied, so:

f(king) - f(man) + f(woman) = f(king-man+woman) ~ f(queen)

Then, we try to estimate a transformation matrix so that:

$$w_{\text{GloVe}} = M_{\text{w2v} \rightarrow \text{GloVe}} \, w_{\text{w2v}}$$

But can a linear transformation preserve similarities?
The cosine similarity is based on the angles between word vectors, so, in general, the answer is no (only orthogonal transformations preserve angles).

# Linear mapping hypothesis

To evaluate the goodness of the mapping, we decided to compute the cosine similarity of each embedded word mapped from the *w2v* model with the same word embedded in *GloVe*.

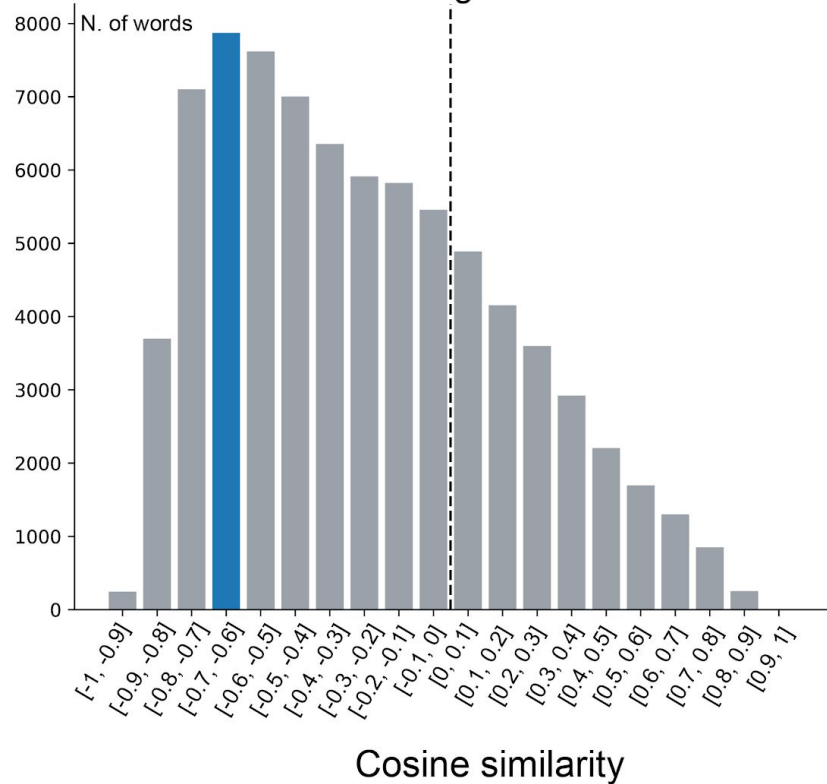Let's see the results from comparisons of *CBOW*, *Skip-gram* and *GloVe* on COHA (1810-1899) and COHA (2000-2009)

**Comparison between CBOW mapped embeddings and GloVe (COHA 1810-1899)**
*Linear transformation*

*Training set*

N. of words

Cosine similarity

*Test set*

N. of words

Cosine similarity

**Comparison between CBOW mapped embeddings and GloVe (COHA 2000-2009)**
*Linear transformation*

# Comparison between Skipgram mapped embeddings and GloVe (COHA 1810-1899)
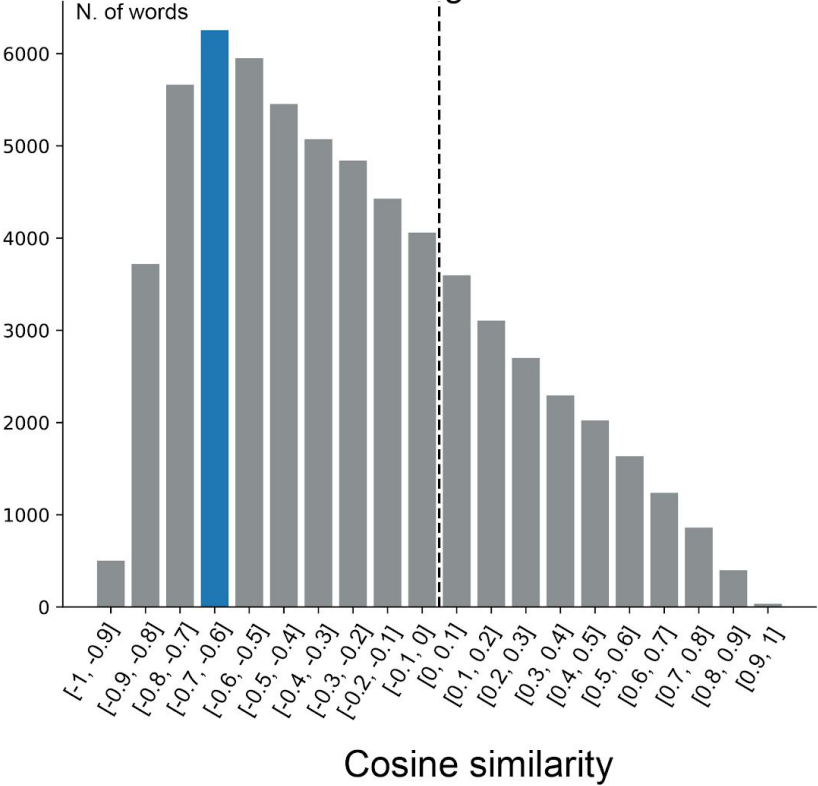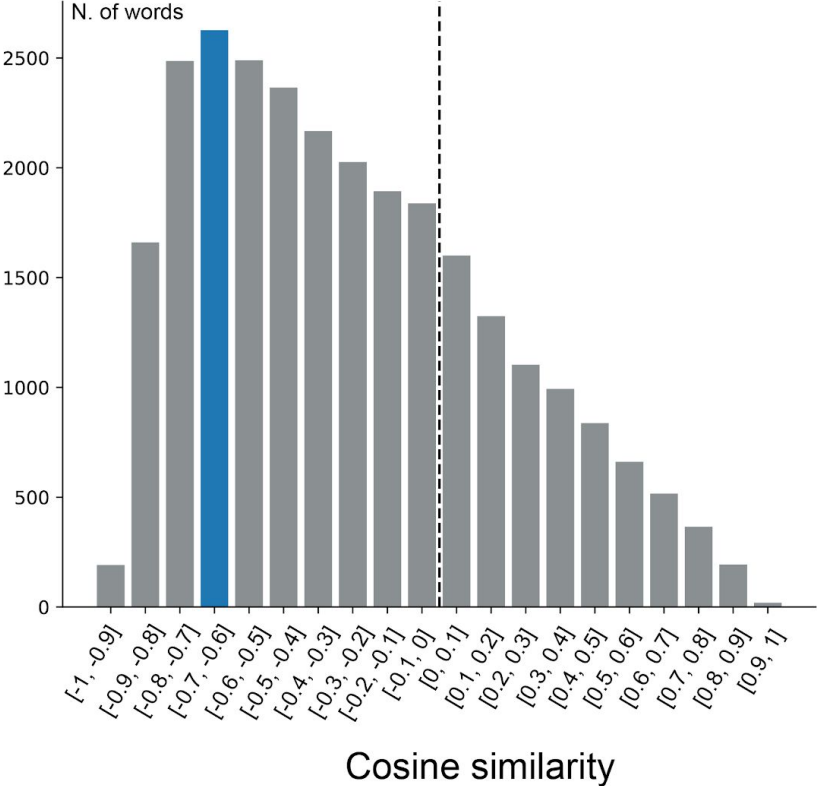## Linear transformation

# Comparison between Skipgram mapped embeddings and GloVe (COHA 2000-2009)
## Linear transformation

# Geometric approach: second hypothesis

We can look for a model which preserves some relaxed but essential properties.

Focusing on ideal perfect similarities and analogies, we only need three properties:
- Collinear points must remain collinear
- Parallel vectors must remain parallel
- Length ratios of parallel vectors must be preserved

There is a type of transformation which has the above properties and it's **affine transformation**.
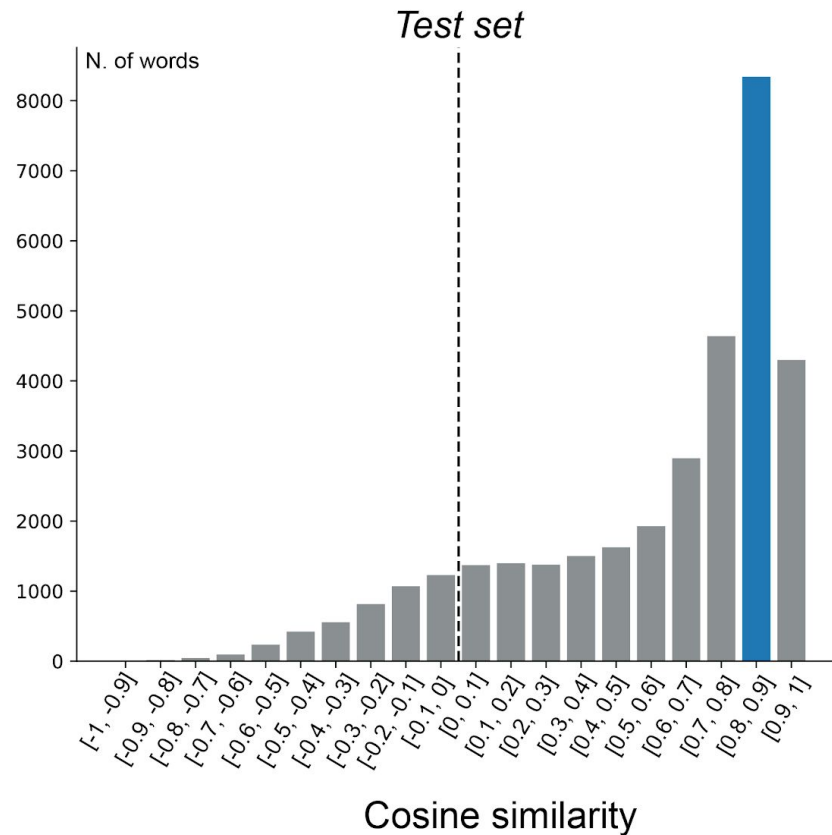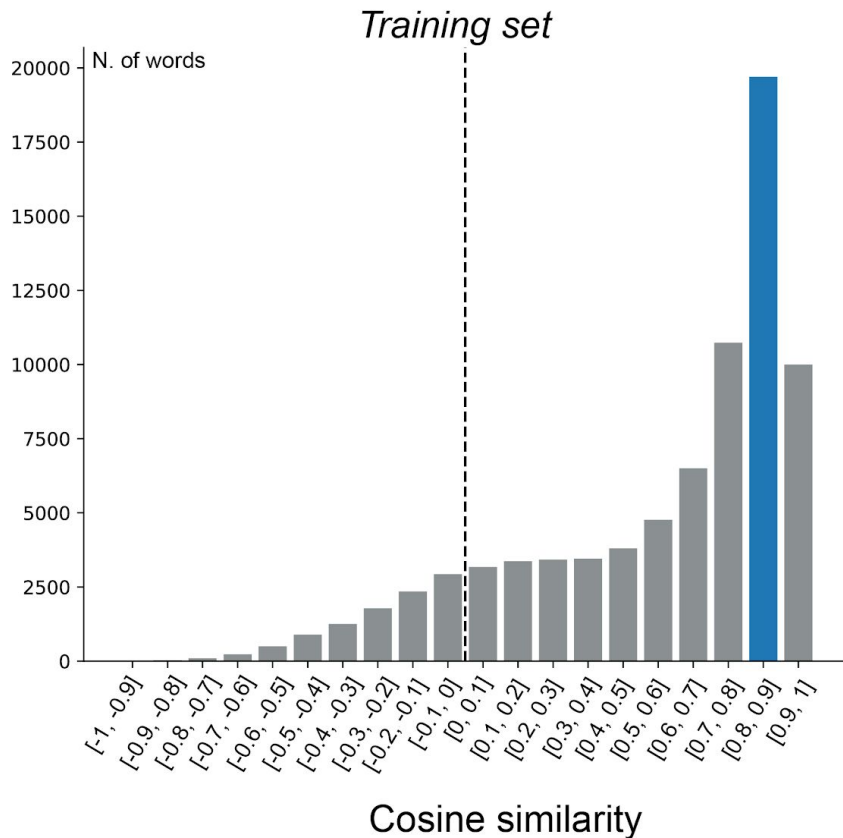
# Affine mapping hypothesis

**Hypothesis n°2:**

Embedding spaces of *Word2Vec* and *GloVe* can be transformed one into the other via affine mapping

$$w_{\text{GloVe}} = M w_{\text{w2v}} + b$$

If this hypothesis holds, then we can conclude that the compared word embedding models are **affine spaces**.
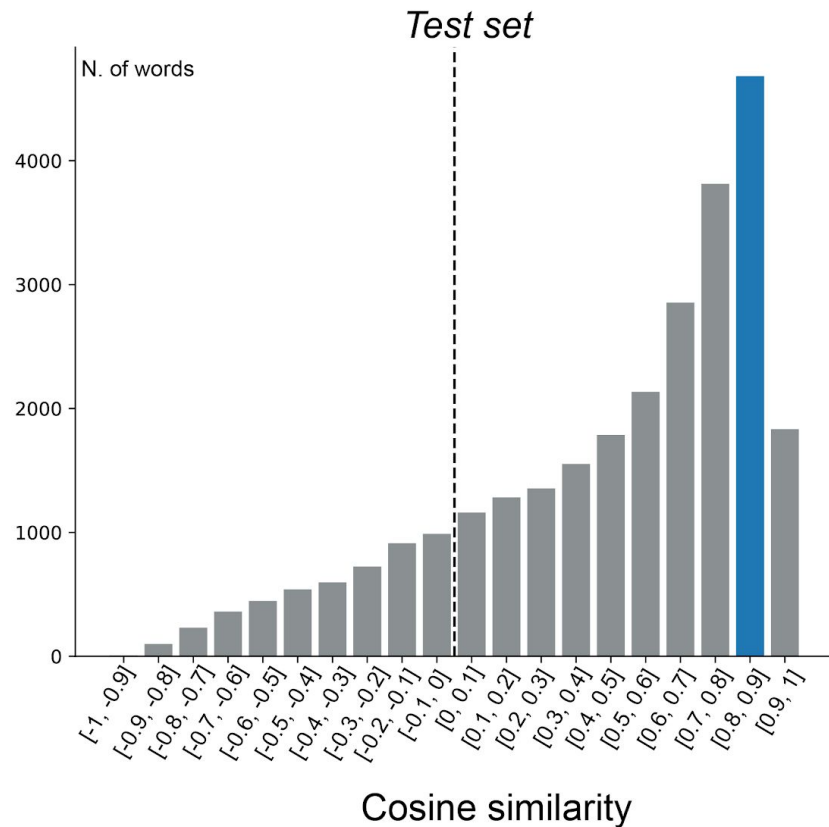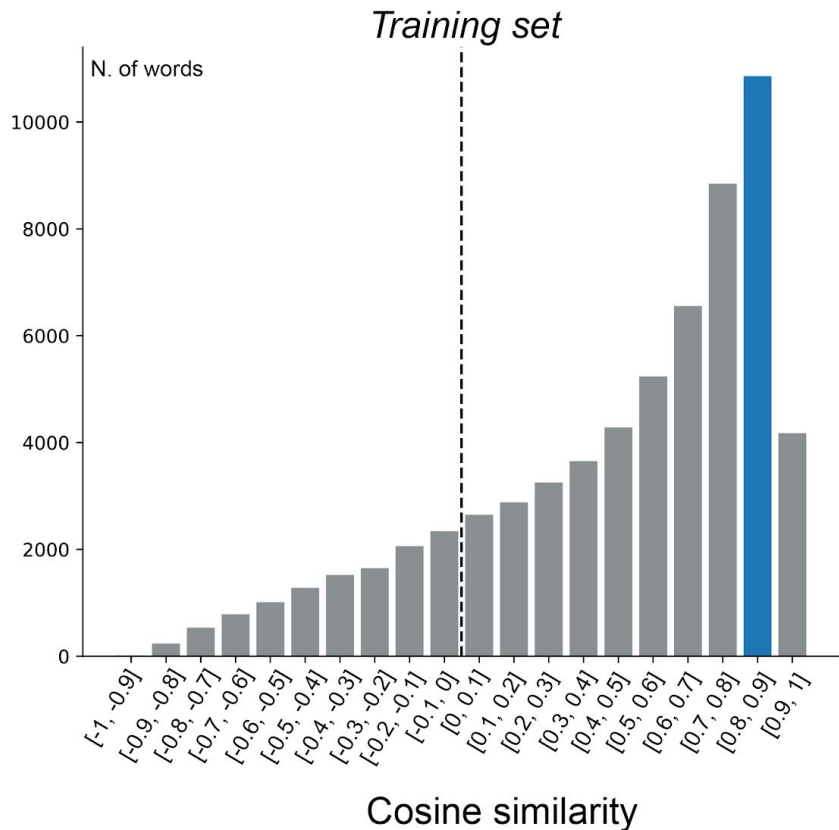
**Comparison between CBOW mapped embeddings and GloVe (COHA 1810-1899)**
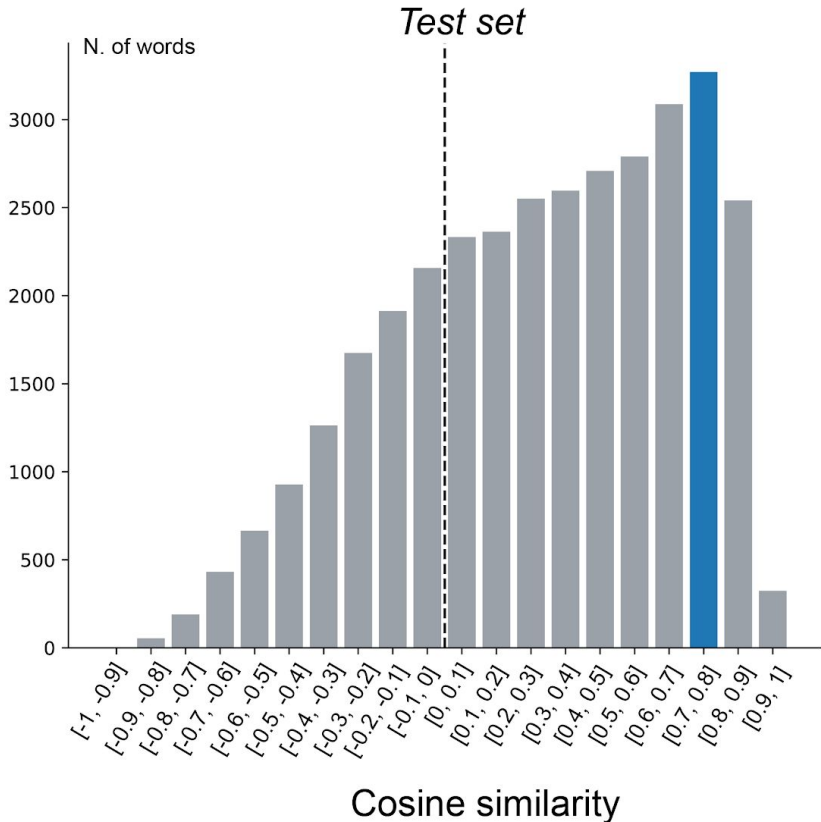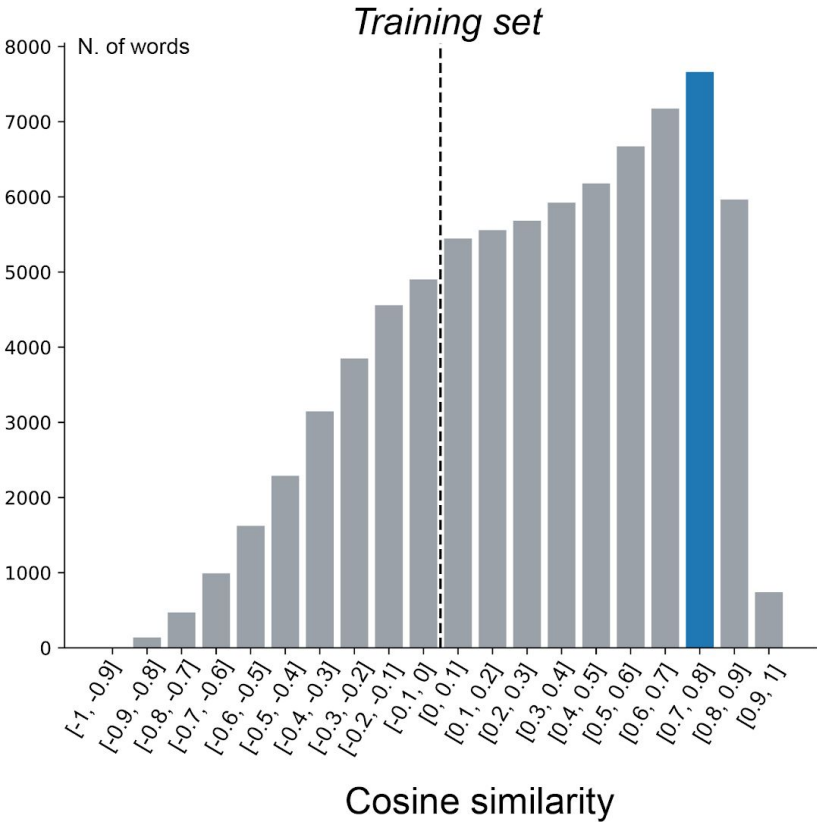*Affine transformation*

*Training set*

N. of words

Cosine similarity

*Test set*

N. of words

Cosine similarity

**Comparison between CBOW mapped embeddings and GloVe (COHA 2000-2009)**
*Affine transformation*

*Training set*

N. of words

10000

8000

6000

4000

2000

0

[-1, -0.9] [-0.9, -0.8] [-0.8, -0.7] [-0.7, -0.6] [-0.6, -0.5] [-0.5, -0.4] [-0.4, -0.3] [-0.3, -0.2] [-0.2, -0.1] [-0.1, 0] [0, 0.1] [0.1, 0.2] [0.2, 0.3] [0.3, 0.4] [0.4, 0.5] [0.5, 0.6] [0.6, 0.7] [0.7, 0.8] [0.8, 0.9] [0.9, 1]

Cosine similarity

*Test set*

N. of words

4000

3000

2000

1000

0

[-1, -0.9] [-0.9, -0.8] [-0.8, -0.7] [-0.7, -0.6] [-0.6, -0.5] [-0.5, -0.4] [-0.4, -0.3] [-0.3, -0.2] [-0.2, -0.1] [-0.1, 0] [0, 0.1] [0.1, 0.2] [0.2, 0.3] [0.3, 0.4] [0.4, 0.5] [0.5, 0.6] [0.6, 0.7] [0.7, 0.8] [0.8, 0.9] [0.9, 1]
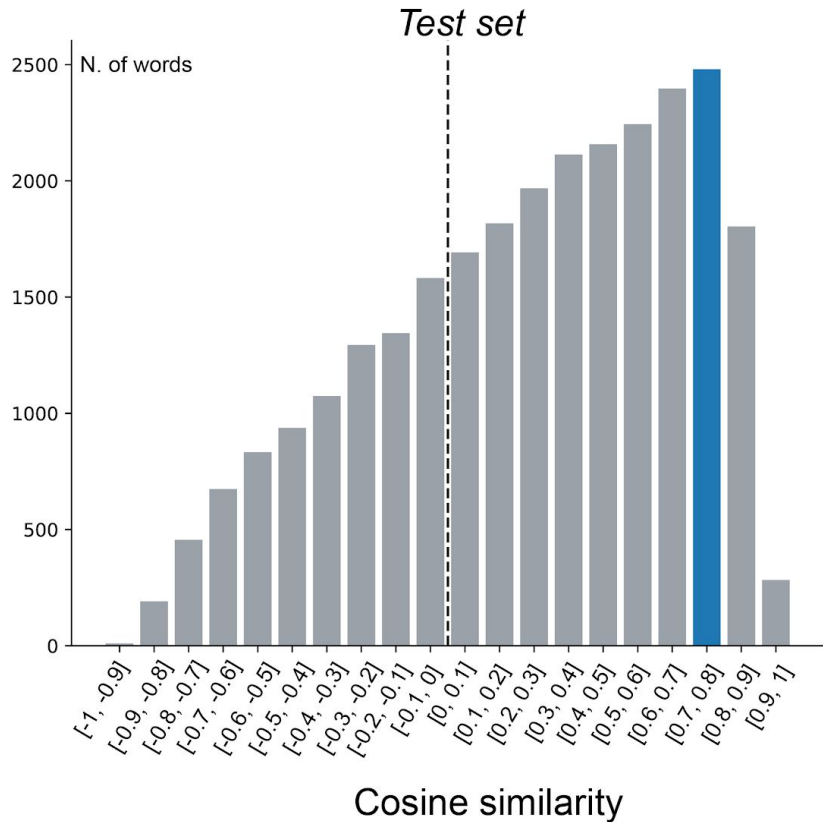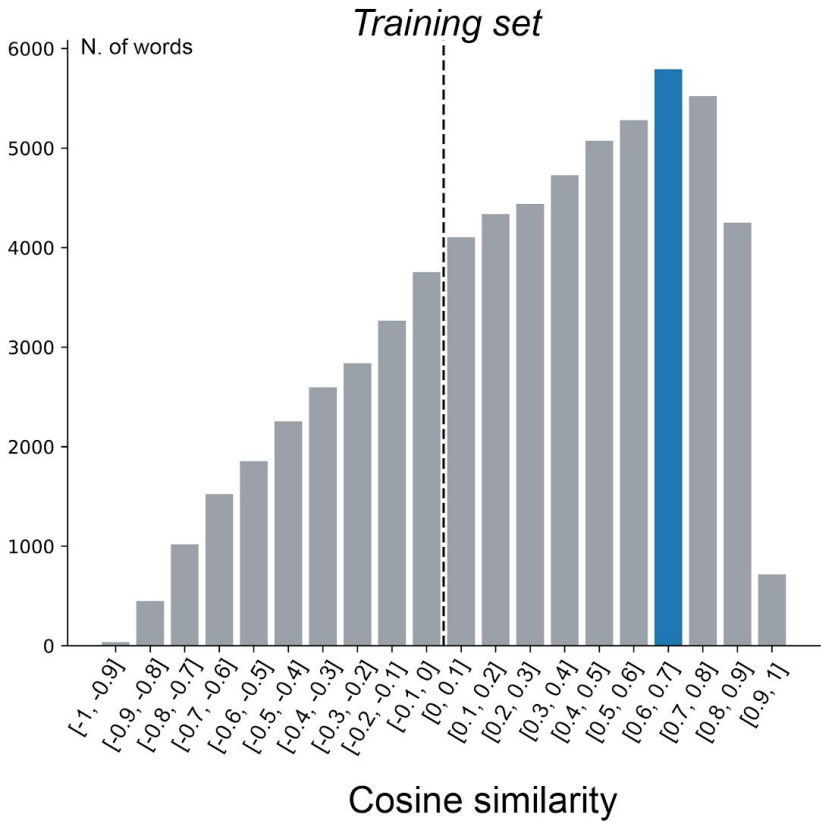
Cosine similarity

# Comparison between Skipgram mapped embeddings and GloVe (COHA 1810-1899)
## Affine transformation

# Comparison between Skipgram mapped embeddings and GloVe (COHA 2000-2009)
## *Affine transformation*

# Notes on estimation of the transformation matrices

We used the function **OT_mapping_linear** from the **POT** library (Python Optimal Transport).

The aim of this function is to find a mapping X→Y of two empirical distributions so that the expected value of |X-Y|$^2$ is minimized[7].

# Future works and improvements

Having seen some promising results in the comparison of *CBOW* and *GloVe*, we believe it would be worth to improve our method in order to better understand relationship between word embeddings model and hopefully aim to a unified theory. Our suggestions are:

- Rewrite the *GloVe* Python library and implement the possibility of setting up the same parameters of *Word2Vec*
- Evaluate if a better estimator of the mapping matrices can be implemented
- Try constraints on the estimator (e.g. the matrix must be orthogonal)
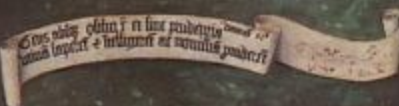- Implement a "*GloVe*" version of *CADE*

# References

- Code repository: https://github.com/grusso98/Data-Semantics

**Bibliography**:

1. Federico Bianchi, Valerio Di Carlo, Paolo Nicoli, Matteo Palmonari. 2020. Compass-aligned Distributional Embeddings for Studying Semantic Differences across Corpora: https://arxiv.org/pdf/2004.06519.pdf
2. Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Distributed Representations of Words and Phrases and their Compositionality. URL: https://arxiv.org/pdf/1310.4546.pdf
3. Laurens van der Maaten, Geoffrey Hinton. 2008. Visualizing Data using t-SNE. URL: https://www.jmlr.org/papers/volume9/vandermaaten08a/vandermaaten08a.pdf
4. Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. GloVe: Global Vectors for Word Representation. URL: https://aclanthology.org/D14-1162/
5. Bin Wang, Student Member, IEEE, Angela Wang, Fenxiao Chen, Student Member, IEEE, Yuncheng Wang and C.-C. Jay Kuo, Fellow, IEEE. Evaluating Word Embedding Models: Methods and Experimental Results. 2019. URL: https://arxiv.org/pdf/1901.09785.pdf
6. Giacomo Berardi, Andrea Esuli, Diego Marcheggiani. Word Embeddings Go to Italy: a Comparison of Models and Training Datasets. 2015. URL: http://iir2015.isti.cnr.it/slides/ITAEmbeddings.pdf
7. Knott, M. and Smith, C. S. "On the optimal mapping of distributions", Journal of Optimization Theory and Applications Vol 43, 1984. URL: https://link.springer.com/article/10.1007/BF00934745

The End

Thanks for your attention!

Data Semantics exam project
July 2022

Gaetano Chiriaco 882638
Riccardo Porcedda 886719
Gianmarco Russo 887277