# Laboratorio de datos: web scraping y Procesamiento de Lenguaje Natural

## Clase 7. Un acercamiento a los word embeddings

# Hipótesis distribucional

- "El significado deriva del uso de las palabras en el lenguaje" (Wittgenstein)

- Podemos captar el sentido de las palabras según su "compañía"

- Palabras cercanas tienen sentidos "cercanos"

- Ítems lingüísticos con distribuciones similares tienen significados similares"

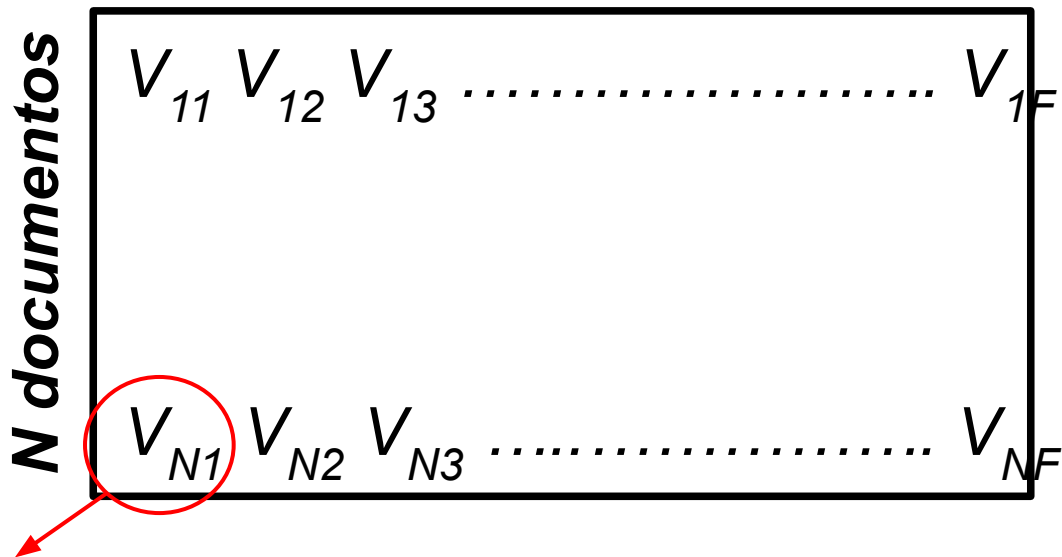- Idea de co-ocurrencia => términos que ocurren juntos

# TFM Co-ocurrencia a nivel documento

Palabras, bigramas, trigramas, lemas, solo la raíz de la palabra...

**F términos**

Matriz $M$ =

**N documentos**

$$V_{11} \; V_{12} \; V_{13} \; .................... \; V_{1F}$$

$$V_{N1} \; V_{N2} \; V_{N3} \; .................... \; V_{NF}$$

Frecuencia del término

factor~data
EIDAES_UNSAM

- La matriz de documentos-términos suele tener muchos ceros
- Problema: se hace difícil medir la relación entre los distintos documentos o términos

|  | Palabra 1 | Palabra 2 | Palabra 3 | Palabra 4 | Palabra 5 |
|---|---|---|---|---|---|
| Relato 1 | 0 | 0.12 | 0.01 | 0 | 0 |
| Relato 2 | 0 | 0 | 0.44 | 0.15 | 0.65 |
| Relato 3 | 0.11 | 0.31 | 0.28 | 0 | 0 |
| Relato 4 | 0 | 0 | 0.05 | 0.21 | 0 |
| Relato 5 | 0 | 0.13 | 0 | 0.07 | 0 |

(...)

(...)

La correlación lineal entre filas nos da una idea de la similitud del significado entre relatos

La correlación lineal entre columnas nos da una idea de la similitud del significado entre palabras

Pero hay un problema: la mayor parte de los valores son 0

"Sobre la mesa hay un florero con margaritas y jazmines"

"El vaso lleno de flores está apoyado sobre una mesada"

- Mismo sentido pero ninguna palabra en común
- Una solución ya la vimos: LDA, STM => detección de tópicos
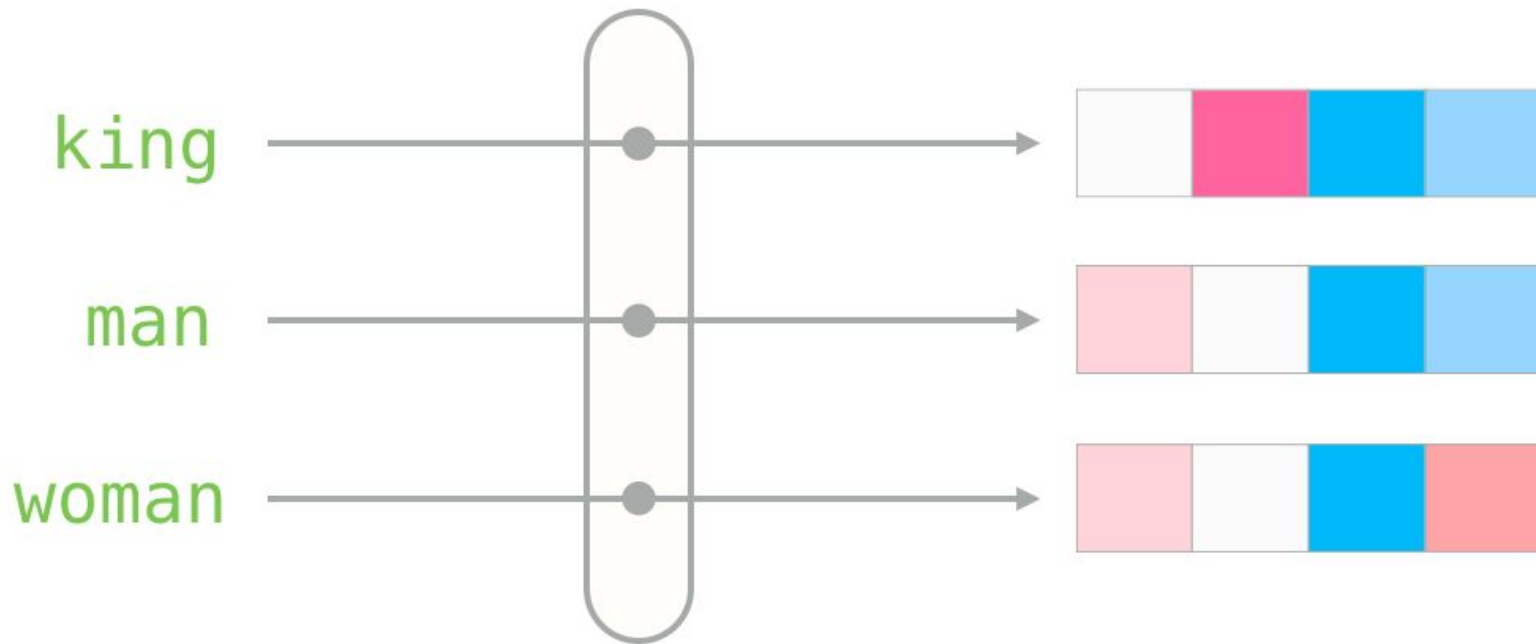
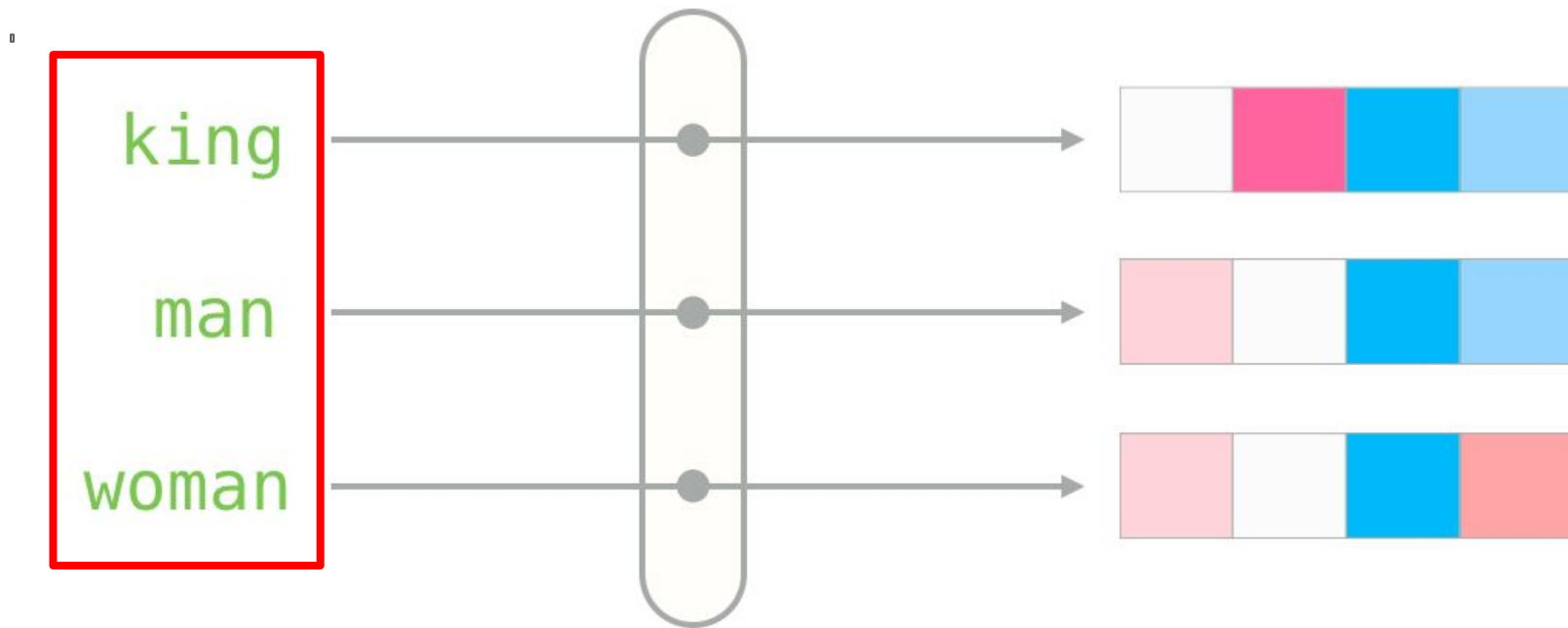- **Otra solución: word embeddings**

# Word embeddings => idea general

- Reducir la dimensión del vocabulario
  - ~50.000 palabras a ~100 => representación no "esparsa" sino densa

- Flexibilizar supuestos de BoW: cada columna/término/dimensión es un término y se asume independencia

- Hay interacción entre palabras => es esperable que la dimensionalidad sea menor

- Lograr introducir una métrica de distancia para que palabras "cerca" en el nuevo espacio estén "cerca" semánticamente estén cerca.
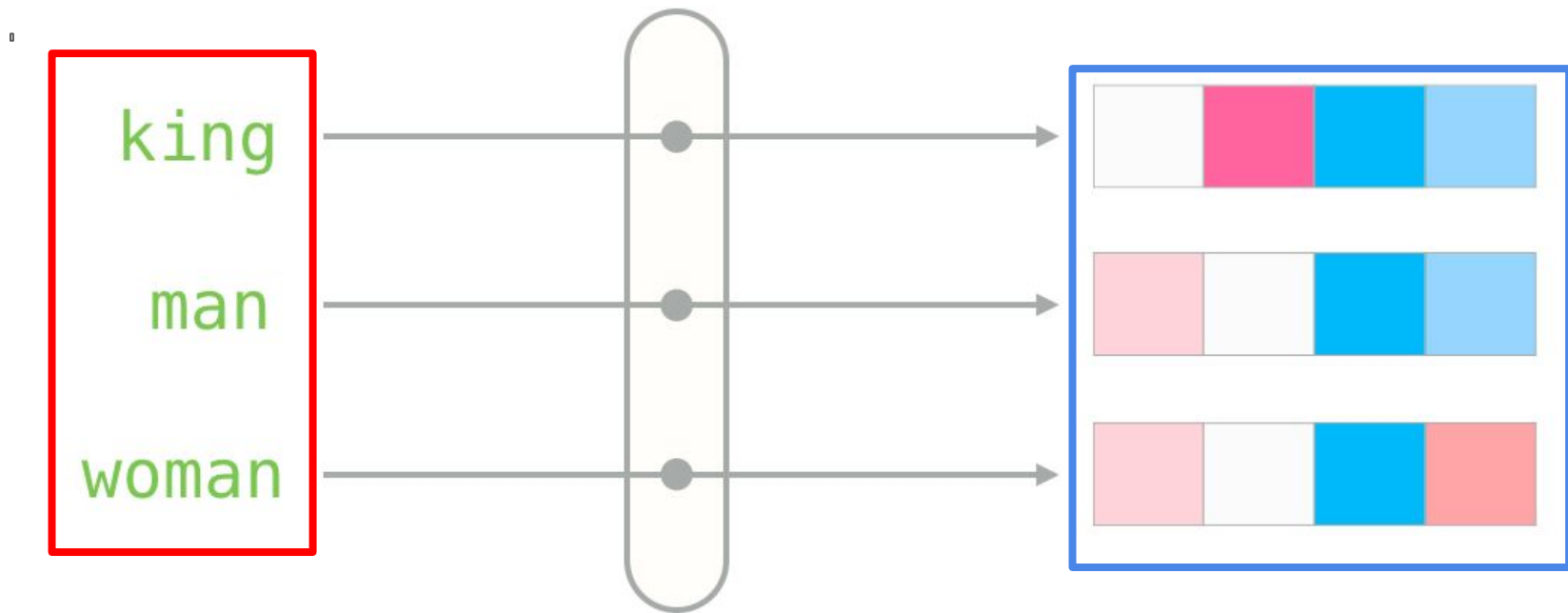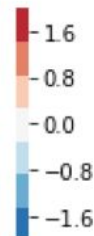
# word2vec

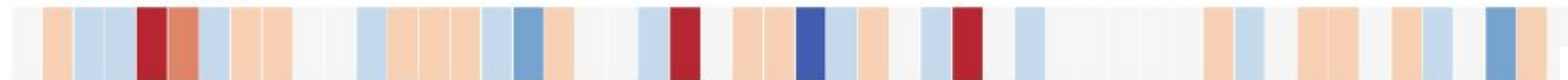**Fuente:** https://jalammar.github.io/illustrated-word2vec/

# word2vec

**Fuente:** https://jalammar.github.io/illustrated-word2vec/

# word2vec

**Fuente:** https://jalammar.github.io/illustrated-word2vec/

# word2vec



"king"

"Man"

"Woman"

| | 1.6 |
| | 0.8 |
| | 0.0 |
| | -0.8 |
| | -1.6 |

**Fuente:** https://jalammar.github.io/illustrated-word2vec/

# word2vec



**Fuente:** https://jalammar.github.io/illustrated-word2vec/

factor~data
EIDAES_UNSAM

# word2vec
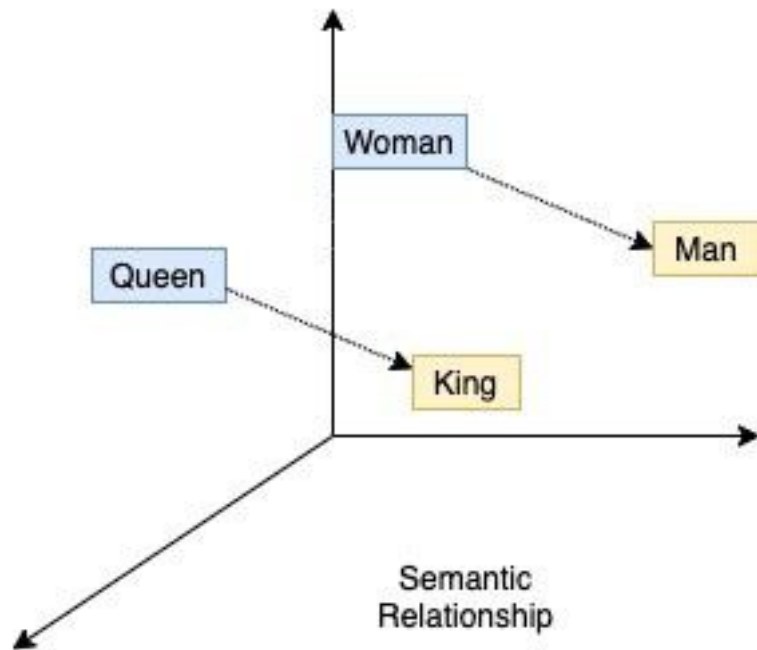


king − man + woman ~= queen

**Fuente:** https://jalammar.github.io/illustrated-word2vec/

# word2vec

# Evaluación de embeddings

Table 1: *Examples of five types of semantic and nine types of syntactic questions in the Semantic-Syntactic Word Relationship test set.*

| 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 |

factor~data
EIDAES_UNSAM

# Evaluación de embeddings

Table 4: *Comparison of publicly available word vectors on the Semantic-Syntactic Word Relationship test set, and word vectors from our models. Full vocabularies are used.*

| Model | Vector Dimensionality | Training words | Accuracy [%] | | |
|---|---|---|---|---|---|
| | | | Semantic | Syntactic | Total |
| Collobert-Weston NNLM | 50 | 660M | 9.3 | 12.3 | 11.0 |
| Turian NNLM | 50 | 37M | 1.4 | 2.6 | 2.1 |
| Turian NNLM | 200 | 37M | 1.4 | 2.2 | 1.8 |
| Mnih NNLM | 50 | 37M | 1.8 | 9.1 | 5.8 |
| Mnih NNLM | 100 | 37M | 3.3 | 13.2 | 8.8 |
| Mikolov RNNLM | 80 | 320M | 4.9 | 18.4 | 12.7 |
| Mikolov RNNLM | 640 | 320M | 8.6 | 36.5 | 24.6 |
| Huang NNLM | 50 | 990M | 13.3 | 11.6 | 12.3 |
| Our NNLM | 20 | 6B | 12.9 | 26.4 | 20.3 |
| Our NNLM | 50 | 6B | 27.9 | 55.8 | 43.2 |
| Our NNLM | 100 | 6B | 34.2 | **64.5** | 50.8 |
| CBOW | 300 | 783M | 15.5 | 53.1 | 36.1 |
| Skip-gram | 300 | 783M | **50.0** | 55.9 | **53.3** |

**Fuente:** Efficient Estimation of Word Representations in Vector Space (2013) https://arxiv.org/abs/1301.3781 Tomas Mikolov, Kai Chen, Greg Corrado, Jeffrey Dean

factor~data
EIDAES_UNSAM

# Evaluación de embeddings

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 |

factor~data
EIDAES_UNSAM

# Usos posibles

- Similitud entre palabras y documentos
- Similitud entre palabras "target" y palabras de contexto al resultado

- Autocompletado
- Traducción automática
- Encontrar clusters de palabras con significados similares
- Buscar analogías entre palabras

- Modelo semántico del lenguaje para comparar con procesamiento del lenguaje hecho por humanos

# Aplicaciones en Ciencias Sociales - Estereotipos

## The Geometry of Culture: Analyzing the Meanings of Class through Word Embeddings

Austin C. Kozlowski,[a] Matt Taddy,[b] and James A. Evans[a,c]



**Figure 2.** Conceptual Diagram of (A) the Construction of a Cultural Dimension; (B) the Projection of Words onto That Dimension; and (C) the Simultaneous Projection of Words onto Multiple Dimensions

factor~data
EIDAES_UNSAM

# Aplicaciones en Ciencias Sociales - Estereotipos



*Measuring Cultural Dimensions*

To identify cultural dimensions in word embedding models, we average numerous pairs of antonym words. Cultural dimensions are calculated by simply taking the mean of all word pair differences that approximate a given dimension, $\dfrac{\sum_p^{|P|} \overrightarrow{p_1} - \overrightarrow{p_2}}{|P|}$ , where $p$ are all antonym word pairs in relevant set $P$, and $\overrightarrow{p_1}$ and $\overrightarrow{p_2}$ are the first and second word vectors of each pair.[17] The projection of a normalized word vector onto a cultural dimension is calculated with cosine similarity, as is the angle between cultural dimensions.

# Aplicaciones en Ciencias Sociales - Estereotipos

# Aplicaciones en Ciencias Sociales - Estereotipos



**Figure 2.** Conceptual Diagram of (A) the Construction of a Cultural Dimension; (B) the Projection of Words onto That Dimension; and (C) the Simultaneous Projection of Words onto Multiple Dimensions

# Aplicaciones en Ciencias Sociales - Estereotipos



**Figure 3.** Projection of Music Genres onto Race and Class Dimensions of the Google News Word Embedding (Gray) and Average Survey Ratings for Race and Class Associations (Black)

**Figure 10.** Words That Project High and Low on the Employment Dimension of Word Embedding Models Trained on Texts Published at the Beginning and End of the Twentieth Century; 1900–1919 and 1980–1999 Google Ngrams Corpus

factor~data
EIDAES_UNSAM

# Aplicaciones en Ciencias Sociales

**ARTICLE**  **OPEN**

## Automated analysis of free speech predicts psychosis onset in high-risk youths

Gillinder Bedi[1,2,9], Facundo Carrillo[3,9], Guillermo A Cecchi[4], Diego Fernández Slezak[3], Mariano Sigman[5], Natália B Mota[6], Sidarta Ribeiro[6], Daniel C Javitt[1,7], Mauro Copelli[8] and Cheryl M Corcoran[1,7]

**BACKGROUND/OBJECTIVES:** Psychiatry lacks the objective clinical tests routinely used in other specializations. Novel computerized methods to characterize complex behaviors such as speech could be used to identify and predict psychiatric illness in individuals.

**AIMS:** In this proof-of-principle study, our aim was to test automated speech analyses combined with Machine Learning to predict later psychosis onset in youths at clinical high-risk (CHR) for psychosis.

**METHODS:** Thirty-four CHR youths (11 females) had baseline interviews and were assessed quarterly for up to 2.5 years; five transitioned to psychosis. Using automated analysis, transcripts of interviews were evaluated for semantic and syntactic features predicting later psychosis onset. Speech features were fed into a convex hull classification algorithm with leave-one-subject-out cross-validation to assess their predictive value for psychosis outcome. The canonical correlation between the speech features and prodromal symptom ratings was computed.

**RESULTS:** Derived speech features included a Latent Semantic Analysis measure of semantic coherence and two syntactic markers of speech complexity: maximum phrase length and use of determiners (e.g., *which*). These speech features predicted later psychosis development with 100% accuracy, outperforming classification from clinical interviews. Speech features were significantly correlated with prodromal symptoms.

**CONCLUSIONS:** Findings support the utility of automated speech analysis to measure subtle, clinically relevant mental state changes in emergent psychosis. Recent developments in computer science, including natural language processing, could provide the foundation for future development of objective clinical tests for psychiatry.

factor~data
EIDAES_UNSAM

# Aplicaciones en Ciencias Sociales - Estereotipos

## Semantics derived automatically from language corpora contain human-like biases

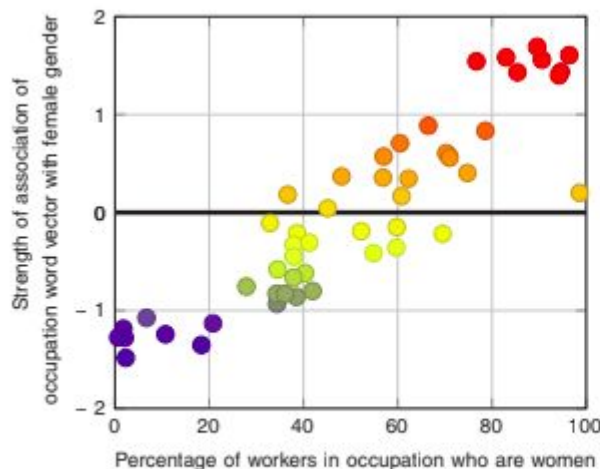Aylin Caliskan,[1*] Joanna J. Bryson,[1,2*] Arvind Narayanan[1*]



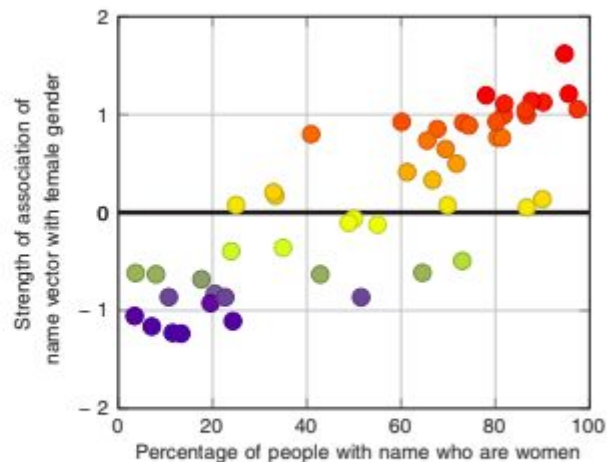Fig. 1. Occupation-gender association. Pearson's correlation coefficient $\rho = 0.90$ with $P < 10^{-18}$.



Fig. 2. Name-gender association. Pearson's correlation coefficient $\rho = 0.84$ with $P < 10^{-13}$.

factor~data
EIDAES_UNSAM

# Aplicaciones en Ciencias Sociales - Trayectorias

- Arquitectura de embeddings / transformers como forma de generar representaciones comprimidas de trayectorias en múltiples dimensiones.
- Life2Vec
  - Dataset n~3.000.000 de habitantes
    - Historias laborales / ingresos
    - Historias migratorias
    - Historias de salud
- Tarea: predicción de fallecimiento
- Usan algo similar a lo que funciona por detrás de chatGPT

Germans Savcisens [1], Tina Eliassi-Rad [2,3], Lars Kai Hansen[1], Laust Hvas Mortensen [4,5], Lau Lilleholt [6,7], Anna Rogers[8], Ingo Zettler [6,7] & Sune Lehmann [1,7] ✉

Here we represent human lives in a way that shares structural similarity to language, and we exploit this similarity to adapt natural language processing techniques to examine the evolution and predictability of human lives based on detailed event sequences. We do this by drawing on a comprehensive registry dataset, which is available for Denmark across several years, and that includes information about life-events related to health, education, occupation, income, address and working hours, recorded with day-to-day resolution. We create embeddings of life-events in a single vector space, showing that this embedding space is robust and highly structured. Our models allow us to predict diverse outcomes ranging from early mortality to personality nuances, outperforming state-of-the-art models by a wide margin. Using methods for interpreting deep learning models, we probe the algorithm to understand the factors that enable our predictions. Our framework allows researchers to discover potential mechanisms that impact life outcomes as well as the associated possibilities for personalized interventions.

We live in the age of algorithm-driven prediction of human behavior. The predictions range from those at the global and population level, with societies allocating vast resources to predicting phenomena such as global warming[1] or the spread of infectious diseases[2], all the way to the constant flow of individual micro-predictions that shape our reality and behavior as we use social media[3]. When it comes to individual life outcomes, however, the picture is more complex. Sociodemographic decade interval, we show that accurate individual predictions are indeed possible. Our dataset includes a host of indicators, such as health, professional occupation and affiliation, income level, residency, working hours and education (Dataset section).

The main reason why we are currently experiencing this 'age of human prediction' is the advent of massive datasets and powerful machine learning algorithms[8,9]. Over the past decade, machine learning

factor~data
EIDAES_UNSAM
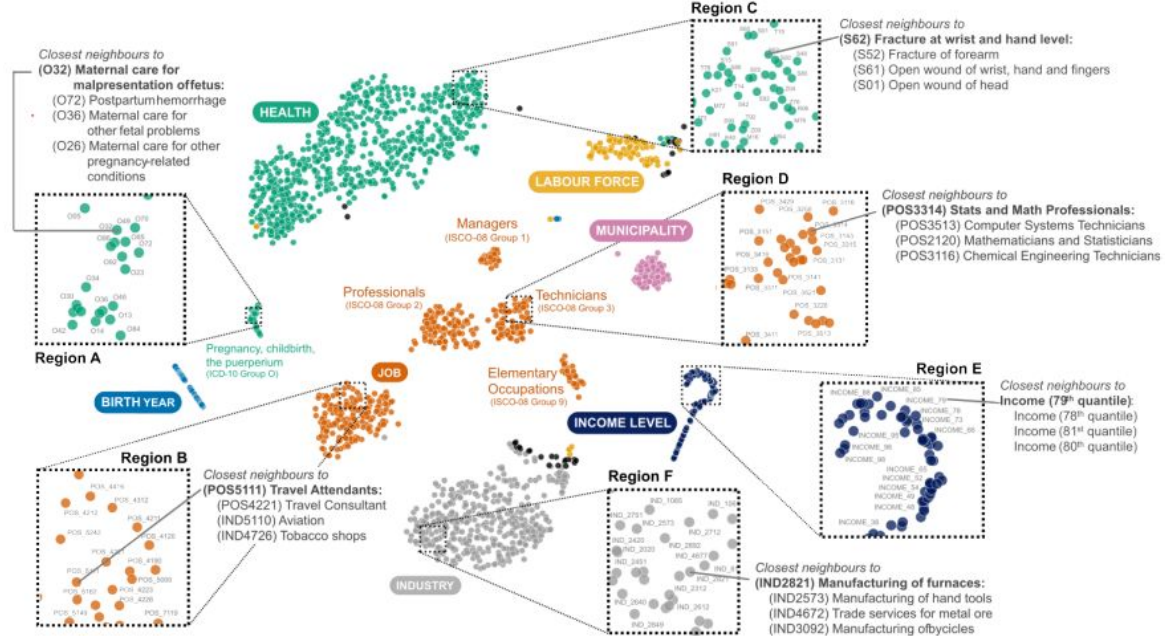
# Aplicaciones en Ciencias Sociales - Trayectorias

**Figure 4:** Two-dimensional projection of the concept space (using the PaCMAP [72]). Each point corresponds to a concept token in the vocabulary. Points are colored based on the concept types (several types are omitted - black points). Each region provides a closer look at several parts of the concept space. You can also see the top three closest neighbors for selected tokens (based on the cosine distance). (**A**) Diagnoses related to Pregnancy, childbirth, and the puerperium in ICD-10 [40]. (**B**) Job concepts related to Service and Sales Workers (corresponds to Job Category 5 of ISCO-08 [38]). (**C**) Injury-related diagnoses in ICD-10 [40]. (**D**) Job concepts related to Technicians and Associate Professionals (corresponds to Job Category 3 of ISCO-08 [38]). (**E**) Income-related concepts. `life2vec` arranges these concepts in increasing ordinal order. (**F**) Concepts related to the manufacturing industry in DB07 [39].

# Aplicaciones en otras disciplinas

# Aplicaciones en otras disciplinas



factor~data
EIDAES_UNSAM

# ¿Cómo sucede la magia?

# One hot encoding

- Eje Y = tiempo
- Eje X = vocabulario
- Celdas: 1 si la palabra aparece en ese "momento"; 0 si no aparece



factor~data
EIDAES_UNSAM

# Skip-gram

Cambia la unidad

Ahora el corpus es visto como un todo continuo…

No se ven los documentos por separado

Un parámetro importante: el tamaño de la ventana…

Otro metodo: CBOW (al revés)

## Source Text

## Training Samples

The quick brown fox jumps over the lazy dog. ➡

(the, quick)
(the, brown)

The quick brown fox jumps over the lazy dog. ➡

(quick, the)
(quick, brown)
(quick, fox)

The quick brown fox jumps over the lazy dog. ➡

(brown, the)
(brown, quick)
(brown, fox)
(brown, jumps)

The quick brown fox jumps over the lazy dog. ➡

(fox, quick)
(fox, brown)
(fox, jumps)
(fox, over)

factor~data
EIDAES_UNSAM

# Skip-gram

| Contexte | | | | Mot Cible |
|---|---|---|---|---|
| The | Quick | Fox | Jump | **Brown** |
| Quick | Brown | Jumps | Over | **Fox** |
| Brown | Fox | Over | The | **Jumps** |

factor~data
EIDAES_UNSAM

# Skip-gram - Matriz de co-ocurrencias

|       | brown | dog | fox | jumps | lazy | over | quick | the |
|-------|-------|-----|-----|-------|------|------|-------|-----|
| brown | 0     | 0   | 0   | 0     | 0    | 0    | 1     | 1   |
| dog   | 0     | 0   | 0   | 0     | 1    | 0    | 0     | 1   |
| fox   | 1     | 0   | 0   | 0     | 0    | 0    | 1     | 0   |
| jumps | 1     | 0   | 1   | 0     | 0    | 0    | 0     | 0   |
| lazy  | 0     | 0   | 0   | 0     | 0    | 1    | 0     | 1   |
| over  | 0     | 0   | 1   | 1     | 0    | 0    | 0     | 0   |
| quick | 0     | 0   | 0   | 0     | 0    | 0    | 0     | 1   |
| the   | 0     | 0   | 0   | 1     | 0    | 1    | 0     | 0   |

factor~data
EIDAES_UNSAM

# Skip-gram (otro ejemplo)

Thou shalt not make a | machine in the likeness of a human mind

| thou | shalt | not | make | a | machine | in | the | ... |

Cambia la unidad

Ahora el corpus es visto como un todo continuo…

No se ven los documentos por separado

Un parámetro importante: el tamaño de la ventana…

Otro metodo: CBOW (al revés)

| input word | target word |
| --- | --- |
| not | thou |
| not | shalt |
| not | make |
| not | a |

factor~data
EIDAES_UNSAM

**Fuente:** https://jalammar.github.io/illustrated-word2vec/

# Skip-gram (otro ejemplo)



Thou [shalt not make a machine] in the likeness of a human mind

| thou | shalt | not | make | a | machine | in | the | ... |
|------|-------|-----|------|---|---------|----|----|-----|

| thou | shalt | not | make | a | machine | in | the | ... |
|------|-------|-----|------|---|---------|----|----|-----|

| input word | target word |
|------------|-------------|
| not | thou |
| not | shalt |
| not | make |
| not | a |
| make | shalt |
| make | not |
| make | a |
| make | machine |

factor~data
EIDAES_UNSAM

# Skip-gram (otro ejemplo)

Thou shalt not make a machine in the likeness of a human mind

| thou | shalt | not | make | a | machine | in | the | ... |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |

| thou | shalt | not | make | a | machine | in | the | ... |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |

| thou | shalt | not | make | a | machine | in | the | ... |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |

| thou | shalt | not | make | a | machine | in | the | ... |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |

| thou | shalt | not | make | a | machine | in | the | ... |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |

| input word | target word |
| --- | --- |
| not | thou |
| not | shalt |
| not | make |
| not | a |
| make | shalt |
| make | not |
| make | a |
| make | machine |
| a | not |
| a | make |
| a | machine |
| a | in |
| machine | make |
| machine | a |
| machine | in |
| machine | the |
| in | a |
| in | machine |
| in | the |
| in | likeness |

factor~data
EIDAES_UNSAM

**Fuente:** https://jalammar.github.io/illustrated-word2vec/

# Modelando con skipgram

| input word | target word |
|:---:|:---:|
| not | thou |
| not | shalt |
| not | make |
| not | a |
| make | shalt |
| make | not |
| make | a |
| make | machine |
| a | not |
| a | make |
| a | machine |
| a | in |
| machine | make |
| machine | a |
| machine | in |
| machine | the |
| in | a |
| in | machine |
| in | the |
| in | likeness |

not ⟶ **Untrained Model**

**Task:**
Predict neighbouring word

factor~data
EIDAES_UNSAM

**Fuente:** https://jalammar.github.io/illustrated-word2vec/

# Modelando con skipgram



not → | Untrained Model

**Task:**
Predict neighbouring word | →

| | |
|---|---|
| 0 | aardvark |
| 0 | aarhus |
| 0.001 | aaron |
| ... | |
| 0.4 | taco |
| 0.001 | thou |
| ... | |
| 0.0001 | zyzzyva |

1) Buscar palabras

2) Calcular la predicción

3) Proyectar vocabulario salida

factor~data
EIDAES_UNSAM

**Fuente:** https://jalammar.github.io/illustrated-word2vec/

# Modelando con skipgram

# Modelando con skipgram

**Fuente:** https://jalammar.github.io/illustrated-word2vec/

# Modelando con skipgram

# Modelando con skipgram => PROBLEMA



not → Untrained Model

**Task:** Predict neighbouring word

| | |
|---|---|
| 0 | aardvark |
| 0 | aarhus |
| 0.001 | aaron |
| ... | |
| 0.4 | taco |
| 0.001 | thou |
| ... | |
| 0.0001 | zyzzyva |

1) Buscar palabras

2) Calcular la predicción

**3) Proyectar vocabulario salida COSTOSO!**

factor~data
EIDAES_UNSAM

**Fuente:** https://jalammar.github.io/illustrated-word2vec/

# Modelando con skipgram => PROBLEMA



Change Task from

Untrained Model

**Task:**
Predict neighbouring word

not → thou

factor~data
EIDAES_UNSAM

# Modelando con skipgram => PROBLEMA

Change Task from

To:

not ──→ **Untrained Model** **Task:** Predict neighbouring word ──→ thou

not ──→
thou ──→
**Untrained Model** **Task:** Are the two words neighbours? ──→ 0.90

**Fuente:** https://jalammar.github.io/illustrated-word2vec/

# Modelando con skipgram => PROBLEMA

Change Task from                                    To:

Untrained Model

not → **Task:** Predict neighbouring word → thou

not → thou → Untrained Model **Task:** Are the two words neighbours? → 0.90

| input word | target word |
|------------|-------------|
| not | thou |
| not | shalt |
| not | make |
| not | a |
| make | shalt |
| make | not |
| make | a |
| make | machine |
| | |

| input word | output word | target |
|------------|-------------|--------|
| not | thou | **1** |
| not | shalt | **1** |
| not | make | **1** |
| not | a | **1** |
| make | shalt | **1** |
| make | not | **1** |
| make | a | **1** |
| make | machine | **1** |
| | | |

Problema!
Todos
ejemplos
positivos…

OVERFITTING

factor~data
EIDAES_UNSAM

**Fuente:** https://jalammar.github.io/illustrated-word2vec/

# Negative sampling

| input word | output word | target |
|:---:|:---:|:---:|
| not | thou | **1** |
| not | | **0** |
| not | | **0** |
| not | shalt | **1** |
| | | |
| | | |
| not | make | **1** |
| | | |
| | | |

Negative examples

**Fuente:** https://jalammar.github.io/illustrated-word2vec/

# Negative sampling

Pick randomly from vocabulary
(random sampling)

| input word | output word | target |
|------------|-------------|--------|
| not | thou | **1** |
| not | aaron | **0** |
| not | taco | **0** |
| not | shalt | **1** |
| | | |
| | | |
| not | make | **1** |
| | | |
| | | |

| Word | Count | Probability |
|------|-------|-------------|
| aardvark | | |
| aarhus | | |
| aaron | | |
| | | |
| taco | | |
| thou | | |
| zyzzyva | | |

**Fuente:** https://jalammar.github.io/illustrated-word2vec/

# La fórmula mágina de w2vec



Skipgram

| shalt | not | make | a | machine |

| input | output |
| --- | --- |
| make | shalt |
| make | not |
| make | a |
| make | machine |

Negative Sampling

| input word | output word | target |
| --- | --- | --- |
| make | shalt | 1 |
| make | aaron | 0 |
| make | taco | 0 |

**Fuente:** https://jalammar.github.io/illustrated-word2vec/

# Regresión logística en forma de red neuronal



$$z_i = \sum w_{ji} x_i$$

$$\sigma(z)_i = \frac{e^{z_i}}{\sum_{j=1}^{K} e^{z_i}}$$

$$x_i$$

factor~data
EIDAES_UNSAM

# Redes neuronales (intuición)



$$\sum w_{ji} x_i \qquad f()$$

$$\sum w_{ji} x_i \qquad f()$$

$$\sum w_{ji} x_i \qquad f()$$

$$\sigma(z)_i = \frac{e^{z_i}}{\sum_{j=1}^{K} e^{z_i}}$$

$$z_i = \sum w_{ji} x_i$$

factor~data
EIDAES_UNSAM

# Redes neuronales (intuición)

# Ahora sí... word2vec



$$\sum w_{ji} x_i$$

$$x_i$$

$$y_i = f\left(\sum w_{ji} x_i\right) \qquad z_k = f'\left(\sum w'_{kj} y_j\right)$$

# Ahora sí… word2vec

Una "unidad" por palabra en el vocabulario => One hot encoded

1 x 5



$x_i$

$$\sum w_{ji} x_i$$

$$y_i = f\left(\sum w_{ji} x_i\right) \qquad z_k = f'\left(\sum w'_{kj} y_j\right)$$

factor~data
EIDAES_UNSAM

# Ahora sí… word2vec

Una "unidad" por palabra en el vocabulario => One hot encoded 1 x 5

Una "unidad" por palabra en el vocabulario => One hot encoded

$$\sum w_{ji} x_i$$

$x_i$

$$\sum w'_{kj} y_j \quad f'()$$

$$y_i = f(\sum w_{ji} x_i) \qquad z_k = f'(\sum w'_{kj} y_j)$$

factor~data
EIDAES_UNSAM

# Ahora sí… word2vec

Este es el **embedding**. Es la representación de baja dimensionalidad de una palabra 1 x 3

Una "unidad" por palabra en el vocabulario => One hot encoded

$$\sum w_{ji} x_i \qquad f()$$

$$\sum w_{ji} x_i \qquad f()$$

$$\sum w_{ji} x_i \qquad f()$$

$$y_i = f\left(\sum w_{ji} x_i\right)$$

$$x_i$$

$$\sum w'_{kj} y_j \qquad f'()$$

$$\sum w'_{kj} y_j \qquad f'()$$

$$\sum w'_{kj} y_j \qquad f'()$$

$$\sum w'_{kj} y_j \qquad f'()$$

$$\sum w'_{kj} y_j \qquad f'()$$

$$z_k = f'\left(\sum w'_{kj} y_j\right)$$

factor~data
EIDAES_UNSAM

# Otros métodos para construir embeddings

- word2vec fue pionero (2013) pero hoy hay métodos mejores

- GloVe: trabaja directamente sobre la matriz de co-ocurrencias



factor~data
EIDAES_UNSAM