# Similaridad Semántica

¿25?

0 ——————————————————————— 100

¿felicidad?

amor ———————————————————— comida

# Similaridad Semántica

Corpus de textos



felicidad

0.1

0.3

0.05

0.15

amor

0.05

0.05

comida

0.2

flan

0.4

0.01

0.15

licuadora

# La semántica de una palabra puede deducirse de su contexto

- Sarabaraban aparece de noche
- Cuando hay un peligro aparece sarabaraban
- Sarabaraban puede ayudarte
- Sarabaraban es un científico

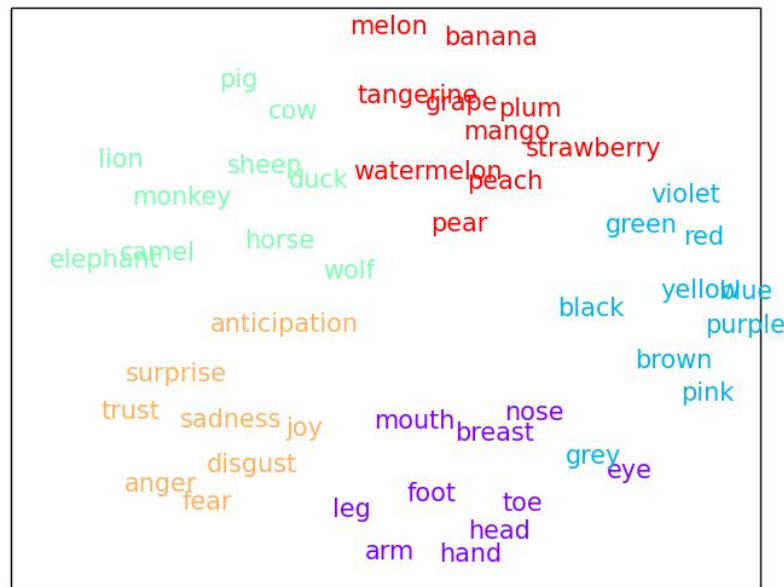La semántica de una pala[bra puede] deducirse de su con[texto]
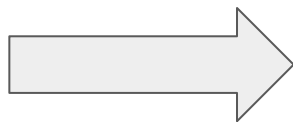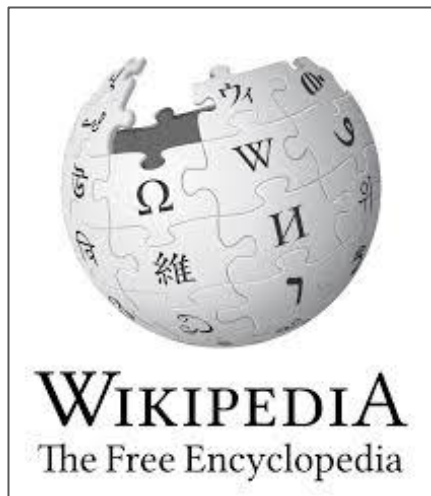
- **Sarabaraban** aparece de noch[e]
- Cuando hay un peligro apare[ce]
- **Sarabaraban** puede ayudarte
- **Sarabaraban** es un científico



J. R. Firth 1957

# Word-embeddings

Corpus de textos

Podemos calcular **cercanías** entre palabras

Vector Space Models



felicidad

amor  0.45        0.18  comida

# Vector Space Model

## Term-Document matrix

|  | Doc 1 | Doc2 | Doc 3 | Doc 4 | Doc 5 | Doc 6 | ... | Doc N |
|---|---|---|---|---|---|---|---|---|
| choripan | 1 | 0 | 3 | 2 | 0 | 0 | ... | 0 |
| vino | 0 | 2 | 0 | 0 | 0 | 0 | ... | 0 |
| chimichurri | 0 | 1 | 2 | 1 | 0 | 0 | ... | 0 |
| uva | 1 | 0 | 0 | 0 | 0 | 2 | ... | 1 |
| pera | 0 | 0 | 0 | 0 | 3 | 0 | ... | 1 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| kiwi | 0 | 0 | 0 | 0 | 1 | 2 | ... | 0 |

# Vector Space Model

## Term-Document matrix

| | Doc 1 | Doc2 | Doc 3 | Doc 4 | Doc 5 | Doc 6 | ... | Doc N |
|---|---|---|---|---|---|---|---|---|
| choripan | 1 | 0 | 3 | 2 | 0 | 0 | ... | 0 |
| vino | 0 | 2 | 0 | 0 | 0 | 0 | ... | 0 |
| chimichurri | 0 | 1 | 2 | 1 | 0 | 0 | ... | 0 |
| uva | 1 | 0 | 0 | 0 | 0 | 2 | ... | 1 |
| pera | 0 | 0 | 0 | 0 | 3 | 0 | ... | 1 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| kiwi | 0 | 0 | 0 | 0 | 1 | 2 | ... | 0 |

# Transformación TF-IDF

$$\text{tf-idf}(t, d) = \text{tf}(t, d) \cdot \text{idf}(t)$$

$$\text{idf}(t) = \log \frac{1 + |D|}{1 + |\{d : t \in d\}|} + 1$$

Número de documentos en el set de entrenamiento

Número de documentos en los que aparece el término $t$

# Vector Space Model

## Term-Document matrix

|             | Doc 1 | Doc2 | Doc 3 | Doc 4 | Doc 5 | Doc 6 | ... | Doc N |
|-------------|-------|------|-------|-------|-------|-------|-----|-------|
| choripan    | 0.83  | 0    | 0.37  | 0.72  | 0     | 0     | ... | 0     |
| vino        | 0     | 0.02 | 0     | 0     | 0     | 0     | ... | 0     |
| chimichurri | 0     | 0.91 | 0.22  | 0.31  | 0     | 0     | ... | 0     |
| uva         | 0.01  | 0    | 0     | 0     | 0     | 0.55  | ... | 0.18  |
| pera        | 0     | 0    | 0     | 0     | 0.13  | 0     | ... | 0.11  |
| ...         | ...   | ...  | ...   | ...   | ...   | ...   | ... | ...   |
| kiwi        | 0     | 0    | 0     | 0     | 0.41  | 0.22  | ... | 0     |

# Vector Space Model

## Term-Document matrix

|  | Doc 1 | Doc2 | Doc 3 | Doc 4 | Doc 5 | Doc 6 | ... | Doc N |
|---|---|---|---|---|---|---|---|---|
| choripan | 0.83 | 0 | 0.37 | 0.72 | 0 | 0 | ... | 0 |
| vino | 0 | 0.02 | 0 | 0 | 0 | 0 | ... | 0 |
| chimichurri | 0 | 0.91 | 0.22 | 0.31 | 0 | 0 | ... | 0 |
| uva | 0.01 | 0 | 0 | 0 | 0 | 0.55 | ... | 0.18 |
| pera | 0 | 0 | 0 | 0 | 0.13 | 0 | ... | 0.11 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| kiwi | 0 | 0 | 0 | 0 | 0.41 | 0.22 | ... | 0 |

# Similaridad coseno

$$\text{cossim}\left(\vec{v}_1, \vec{v}_2\right) = \cos(\alpha) = \frac{\vec{v}_1 . \vec{v}_2}{|\vec{v}_1| . |\vec{v}_2|}$$

Doc 2

vino

choripan

Doc 1

chimichurri

# Similaridad coseno

$$\text{cossim}\left(\vec{v}_1, \vec{v}_2\right) = \cos(\alpha) = \frac{\vec{v}_1 . \vec{v}_2}{|\vec{v}_1| . |\vec{v}_2|}$$

$$\text{cossim}\left(\vec{v}_1, \vec{v}_2\right) = \frac{\sum_{i=1}^{N} v_{1,i} . v_{2,i}}{\sqrt{\sum_{i=1}^{N} v_{1,i}^2} \sqrt{\sum_{i=1}^{N} v_{2,i}^2}}$$
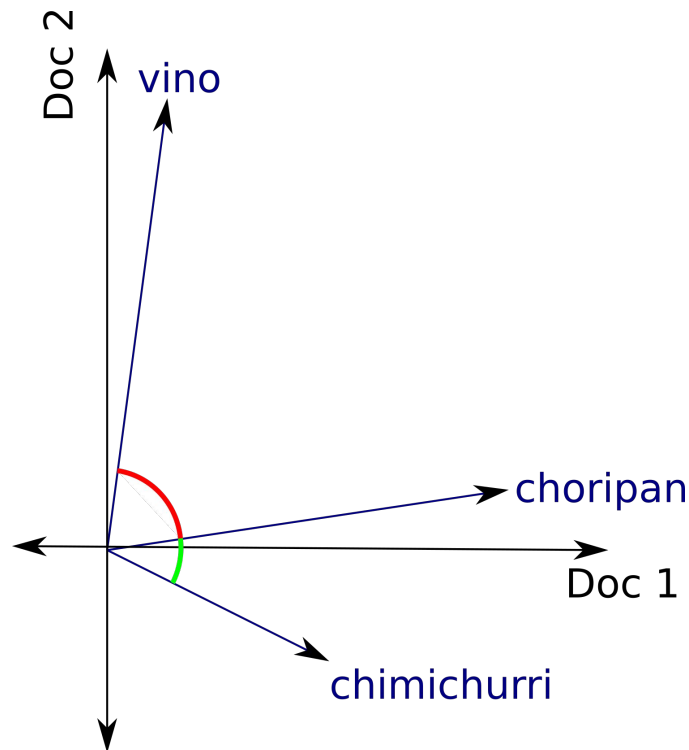
$$\text{cossim}\left(\vec{v}_1, \vec{v}_2\right) \in [-1, 1]$$

Ejemplo

$$
\begin{aligned}
\vec{v}_1 &= (0, 5, 1) \\
\vec{v}_2 &= (1, 0, 2) \\
\text{cossim}\left(\vec{v}_1 . \vec{v}_2\right) &= \frac{0x1 + 5x0 + 1x2}{\sqrt{0^2 + 5^2 + 1^2}\sqrt{1^2 + 0^2 + 2^2}} \approx 0.175
\end{aligned}
$$



15

# Vector Space Model

## Term-Document matrix

| | Doc 1 | Doc2 | Doc 3 | Doc 4 | Doc 5 | Doc 6 | ... | Doc N |
|---|---|---|---|---|---|---|---|---|
| choripan | 0.83 | 0 | 0.37 | 0.72 | 0 | 0 | ... | 0 |
| vino | 0 | 0.02 | 0 | 0 | 0 | 0 | ... | 0 |
| chimichurri | 0 | 0.91 | 0.22 | 0.31 | 0 | 0 | ... | 0 |
| uva | 0.01 | 0 | 0 | 0 | 0 | 0.55 | ... | 0.18 |
| pera | 0 | 0 | 0 | 0 | 0.13 | 0 | ... | 0.11 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| kiwi | 0 | 0 | 0 | 0 | 0.41 | 0.22 | ... | 0 |

# Information-Retrieval

## Query: "El chimichurri es un condimento típico de Argentina..."

Entreno el TF-ID en el dataset

Aplico los pesos entrenados

| | Doc 1 | Doc2 | Doc 3 | Doc 4 | Doc 5 | Doc 6 | ... | Doc N | | Query |
|---|---|---|---|---|---|---|---|---|---|---|
| choripan | 0.83 | 0 | 0.37 | 0.72 | 0 | 0 | ... | 0 | | 0.15 |
| vino | 0 | 0.02 | 0 | 0 | 0 | 0 | ... | 0 | | 0 |
| chimichurri | 0 | 0.91 | 0.22 | 0.31 | 0 | 0 | ... | 0 | | 0.74 |
| uva | 0.01 | 0 | 0 | 0 | 0 | 0.55 | ... | 0.18 | | 0 |
| pera | 0 | 0 | 0 | 0 | 0.13 | 0 | ... | 0.11 | | 0.09 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | | ... |
| kiwi | 0 | 0 | 0 | 0 | 0.41 | 0.22 | ... | 0 | | 0 |

# Term-context matrix

Choripan. The Argentine *choripán* consists of a sausage made out of beef and pork, hot off the grill, split down the middle, and served on a roll. The *chorizo* may be used whole or cut in half lengthwise, in which case it is called a *mariposa* (butterfly). It is customary to add sauces on the bread, most likely *chimichurri*.

# Term-context matrix

|  | choripan | vino | chimichurri | uva | pera | ... | kiwi |
|---|---|---|---|---|---|---|---|
| choripan | 47 | 54 | 23 | 5 | 2 | ... | 1 |
| vino | 54 | 354 | 17 | 21 | 3 | ... | 4 |
| chimichurri | 23 | 17 | 59 | 1 | 1 | ... | 0 |
| uva | 5 | 21 | 1 | 203 | 20 | ... | 19 |
| pera | 2 | 3 | 1 | 20 | 399 | ... | 11 |
| ... | ... | ... | ... | ... | ... | ... | ... |
| kiwi | 1 | 4 | 0 | 19 | 11 | ... | 61 |

# Term-context matrix

| | choripan | vino | chimichurri | uva | pera | ... | kiwi |
|---|---|---|---|---|---|---|---|
| choripan | 47 | 54 | 23 | 5 | 2 | ... | 1 |
| vino | 54 | 354 | 17 | 21 | 3 | ... | 4 |
| chimichurri | 23 | 17 | 59 | 1 | 1 | ... | 0 |
| uva | 5 | 21 | 1 | 203 | 20 | ... | 19 |
| pera | 2 | 3 | 1 | 20 | 399 | ... | 11 |
| ... | ... | ... | ... | ... | ... | ... | ... |
| kiwi | 1 | 4 | 0 | 19 | 11 | ... | 61 |

# Term-Document matrix

- First-order co-occurrence

- Syntagmatic associations

- Ejemplo:
  mostaza - hamburguesa

# Term-Context matrix

- Second-order co-occurrence

- Paradigmatic associations

- Ejemplo:
  mostaza - ketchup

# Similaridad semántica

Qué está más asociado al **choripan**, el **vino** o el **chimichurri**?

## Term-context matrix

$f$ =

| W \ C | choripan | vino | chimichurri | en |
|---|---|---|---|---|
| choripan | 30 | 10 | 10 | 30 |
| vino | 10 | 55 | 5 | 50 |
| chimichurri | 10 | 5 | 15 | 10 |
| en | 30 | 50 | 10 | 500 |

# Pointwise Mutual Information (PMI)

$$\mathrm{PMI}(w, c) = \log \left( \frac{P(w,c)}{P(w) . P(c)} \right)$$

| W \ C | choripan | vino | chimichurri | en |
|-------|----------|------|-------------|-----|
| choripan | 30 | 10 | 10 | 30 |
| vino | 10 | 55 | 5 | 50 |
| chimichurri | 10 | 5 | 15 | 10 |
| en | 30 | 50 | 10 | 500 |

Church and Hanks (1990) Word Association Norms, Mutual Information, and Lexicography

# Ejemplo

$$f = $$

| W \ C | choripan | vino | chimichurri | en | $N_C$ |
|---|---|---|---|---|---|
| choripan | 30 | 10 | 10 | 30 | 80 |
| vino | 10 | 55 | 5 | 50 | 120 |
| chimichurri | 10 | 5 | 15 | 10 | 40 |
| en | 30 | 50 | 10 | 500 | 590 |
| $N_W$ | 80 | 120 | 40 | 590 | 830 |

$$p_{ij} = \frac{f_{ij}}{\sum_{i=1}^{W} \sum_{j=1}^{C} f_{ij}}$$

$$p_{i*} = \frac{\sum_{j=1}^{C} f_{ij}}{\sum_{i=1}^{W} \sum_{j=1}^{C} f_{ij}}$$

$$p_{*j} = \frac{\sum_{i=1}^{W} f_{ij}}{\sum_{i=1}^{W} \sum_{j=1}^{C} f_{ij}}$$

$$\mathrm{PMI}(w, c) = \log\left(\frac{p_{ij}}{p_{i*} p_{*j}}\right)$$

$$\mathrm{PMI}(w = chimi, c = chori) = \log\left(\frac{\frac{10}{830}}{\frac{80}{830}\frac{40}{830}}\right) \approx \log\left(\frac{0.01204}{0.00465}\right) = 0.9514$$

# Ejemplo

$$\text{PMI}(w, c) = \log\left(\frac{P(w,c)}{P(w) \cdot P(c)}\right)$$

PMI =

| W \ C | choripan | vino | chimichurri | en |
|---|---|---|---|---|
| choripan | **1.359** | -0.146 | 0.951 | -0.639 |
| vino | -0.146 | **1.154** | -0.146 | -0.534 |
| chimichurri | **0.951** | -0.146 | **2.052** | -1.045 |
| en | -0.639 | -0.534 | -1.045 | **0.176** |

# Otra forma de verlo

$$\text{PMI}(w, c) = \log\left(\frac{P(w,c)}{P(w) \cdot P(c)}\right)$$

$$\text{PMI}(w, c) = \log\left(\frac{P(c|w)P(w)}{P(w) \cdot P(c)}\right)$$

$$= \log\left(\frac{P(c|w)}{P(c)}\right)$$

$$= \log\left(\frac{P(w|c)}{P(w)}\right)$$

$$= \log\left(\frac{P(c=chori|w=chimi)}{P(c=chori)}\right) = \log\left(\frac{10/40}{80/830}\right) = 0.9514$$

| W \ C | choripan | vino | chimichurri | en | $N_C$ |
|---|---|---|---|---|---|
| choripan | 30 | 10 | 10 | 30 | 80 |
| vino | 10 | 55 | 5 | 50 | 120 |
| chimichurri | 10 | 5 | 15 | 10 | 40 |
| en | 30 | 50 | 10 | 500 | 590 |
| $N_W$ | 80 | 120 | 40 | 590 | 830 |

# ¿Que significa un PMI negativo?

$$\mathrm{PMI}(w, c) = \log\left(\frac{P(w,c)}{P(w) \cdot P(c)}\right)$$

Si $P(w) = P(c) = 10^{-6}$
Para que PMI<0 $\longrightarrow$ $P(w,c) < 10^{-12}$

# Se suele usar el Positive-PMI (PPMI)

$$\text{PPMI}(w, c) = \max\left(0, \log\left(\frac{P(w,c)}{P(w) \cdot P(c)}\right)\right)$$

PPMI =

| W \ C | choripan | vino | chimichurri | en |
|-------|----------|------|-------------|-----|
| choripan | 1.359 | 0 | 0.951 | 0 |
| vino | 0 | 1.154 | 0 | 0 |
| chimichurri | 0.951 | 0 | 2.052 | 0 |
| en | 0 | 0 | 0 | 0.176 |

# Cuidado!

$$\text{PPMI}(w, c) = \max\left(0, \log\left(\frac{P(w,c)}{P(w) \cdot P(c)}\right)\right)$$

En el caso límite donde 2 palabras están totalmente correlacionadas como por ejemplo (hocus pocus):
P(w,c)=P(w)=P(c):

$$\text{PMI}(w, c) = \log\left(\frac{P(w,c)}{P(w) \cdot P(c)}\right) = \log\left(\frac{1}{P(w,c)}\right) = -\log(P(w,c))$$

P(w,c) chicos dan PPMIs mas grandes!

Jurafsky and Martin (2017) Speech and Language Processing, 3rd editions
Bouma (2009) Normalized Pointwise Mutual Information in Collocation Extraction

# Soluciones

$$\text{PPMI}(w, c) = \max\left(0, \log\left(\frac{P(w,c)}{P(w) \cdot P(c)}\right)\right)$$

➢ Filtrar $f_{ij} < k$

Levy et al (2015) Improving distributional similarity with lessons learned from word-embeddings
Jurafsky and Martin (2017) Speech and Language Processing, 3rd editions
Bouma (2009) Normalized Pointwise Mutual Information in Collocation Extraction

# Soluciones

$$\text{PPMI}(w, c) = \max\left(0, \log\left(\frac{P(w,c)}{P(w) \cdot P(c)}\right)\right)$$

➢ Filtrar $f_{ij} < k$

➢ Aumentar P(c), Levy et al. (2015)   $P_\alpha(c) = \dfrac{count(c)^\alpha}{\sum_c counts(c)^\alpha}$

Levy et al (2015) Improving distributional similarity with lessons learned from word-embeddings
Jurafsky and Martin (2017) Speech and Language Processing, 3rd editions
Bouma (2009) Normalized Pointwise Mutual Information in Collocation Extraction

# Soluciones

$$\text{PPMI}(w, c) = \max\left(0, \log\left(\frac{P(w,c)}{P(w) \cdot P(c)}\right)\right)$$

➢ Filtrar $f_{ij} < k$

➢ Aumentar P(c), Levy et al. (2015)     $P_\alpha(c) = \frac{count(c)^\alpha}{\sum_c counts(c)^\alpha}$

➢ Add-k smoothing

|            | choripan | vino  | chimichurri | en     |
|------------|----------|-------|-------------|--------|
| choripan   | 30+k     | 10+k  | 10+k        | 30+k   |
| vino       | 10+k     | 45+k  | 5+k         | 50+k   |
| chimichurri| 10+k     | 5+k   | 20+k        | 10+k   |
| en         | 30+k     | 50+k  | 10+k        | 500+k  |

Levy et al (2015) Improving distributional similarity with lessons learned from word-embeddings
Jurafsky and Martin (2017) Speech and Language Processing, 3rd editions
Bouma (2009) Normalized Pointwise Mutual Information in Collocation Extraction

# Soluciones

$$\text{PPMI}\left(w, c\right) = \max\left(0, \log\left(\frac{P(w,c)}{P(w) \cdot P(c)}\right)\right)$$

➢ Filtrar $f_{ij}$ <k

➢ Aumentar P(c), Levy et al. (2015) $\qquad P_\alpha(c) = \dfrac{count(c)^\alpha}{\sum_c counts(c)^\alpha}$

➢ Add-k smoothing

➢ Normalized-PMI, Bouma (2009)

|  | choripan | vino | chimichurri | en |
|---|---|---|---|---|
| choripan | 30+k | 10+k | 10+k | 30+k |
| vino | 10+k | 45+k | 5+k | 50+k |
| chimichurri | 10+k | 5+k | 20+k | 10+k |
| en | 30+k | 50+k | 10+k | 500+k |

$$\text{NPMI}\left(w, c\right) = \frac{\log\left(\frac{P(w,c)}{P(w) \cdot P(c)}\right)}{-\log(P(w,c))}$$

Levy et al (2015) Improving distributional similarity with lessons learned from word-embeddings
Jurafsky and Martin (2017) Speech and Language Processing, 3rd editions
Bouma (2009) Normalized Pointwise Mutual Information in Collocation Extraction

# Aplicaciones: Collocations (expresiones)

Tokenización:
➢ "Las tardecitas de **Buenos Aires** tiene ese qué sé yo, viste?
➢ "Spinetta era de **villa urquiza**"

# Aplicaciones: Collocations (expresiones)

Tokenización:
➢ "Las tardecitas de **Buenos Aires** tiene ese qué sé yo, viste?
➢ "Spinetta era de **villa urquiza**"

Más ejemplos:
- martin fierro
- salud mental
- ping pong
- dar paja
- susana gimenez

# Identificación de collocations

Opciones:
➢ Usar listas de collocation

➢ Buscar bigramas en un corpus con un alto PPMI

Choripan. The Argentine *choripán* consists of a sausage made out of beef and pork, hot off the grill, split down the middle, and served on a roll. The *chorizo* may be used whole or cut in half lengthwise, in which case it is called a *mariposa* (butterfly). It is customary to add sauces on the bread, most likely *chimichurri*.

# Otras opciones:

➢ Implementación de Gensim: NPMI

➢ Implementación de NLTK: PMI

Choripan**.** The Argentine *choripán* consists of a sausage made out of beef and pork, hot off the grill, split down the middle, and served on a roll. The *chorizo* may be used whole or cut in half lengthwise, in which case it is called a *mariposa* (butterfly). It is customary to add sauces on the bread, most likely *chimichurri*.

# Observatorio del cine

# "Brilliance = males" stereotype

A) Story about a really really **smart** person

B) Story about a really really **nice** person



Bian et al. (2017) Gender stereotypes about intellectual ability emerge early and influence children's interests

# Cuales son las fuentes de estereotipos?

## Hipótesis:

➢ Tratamiento diferencial de padres y maestras/os

➢ Falta de roles modelos

➢ Exposición a productos culturales que refuerzan el estereotipo

# "Brilliance = males" stereotype  in films subtitles



opensubtitles.org

NLP

Más de 11.000 subtítulos

# Análisis

"Brilliance" related words:

*ingenious, genius, ingeniousness, ingeniously, bright, brightness, brightly, brilliant, brilliance, brilliantly, clever, cleverness, cleverly, intelligent, intelligence, intelligently*.

Pronombres femeninos:

*she, hers, her, herself*.

Pronombres masculinos:

*he, his, he, himself*.

# Quantifying "brilliance=male" stereotype in films

$$PMI(w, c) = log\left(\frac{p(w, c)}{p(w)p(c)}\right)$$

*w = pronombres*

*c = "brilliance" related words*

# Quantifying "brilliance=male" stereotype in films

$$PMI(w, c) = log\left(\frac{p(w, c)}{p(w)p(c)}\right)$$

*w = pronombres*
*c = "brilliance" related words*

➢ Asociación sintagmática: 1st order co-occurrences

➢ Facil interpretacion

# Quantifying "brilliance=male" stereotype in films

$$PMI(w, c) = log\left(\frac{p(w, c)}{p(w)p(c)}\right)$$

*w = pronombres*
*c = "brilliance" related words*

➢ Asociación sintagmática: 1ˢᵗ order co-occurrences

➢ Facil interpretacion

Gender bias    brilliance-male association    brilliance-female association

$$\Delta PMI = PMI(w_m, c) - PMI(w_f, c)$$

# Sliding time-window to compute co-occurrence



## SubRip File

```
279
00:24:09,973 --> 00:24:11,665
I want you to rest well, and a
month from now...

280
00:24:11,839 --> 00:24:15,203
This Hollywood big shot's
gonna give you what you want.

281
00:24:15,373 --> 00:24:18,032
Too late. They start shooting in
a week.

282
00:24:18,573 --> 00:24:21,561
I'm gonna make him an offer he
can't refuse.

283
00:24:24,606 --> 00:24:26,334
Now just go outside, enjoy
yourself...

284
00:24:26,572 --> 00:24:30,505
and forget about all this
nonsense.
```

1) Given a target word in a SubRip file ("him" in this case), all frames falling within a time-window are identified.

| 279 | 280 | 281 | 282 | 283 | 284 |

Δt       Δt       Time

2) Text contained in all context frames is cleaned, tokenized and lemmatized.

### Tokens in the current context

this, hollywood, big, shot, 's, gon, na, give, you, what, you, want, too, late, they, start, shoot, in, a, week, i, 'm, gon, na, make, an, offer, he, ca, n't, refuse, now, just, go, outside, enjoy, yourself
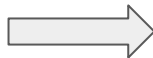
3) For each token, its number of appearances in the current context is added to the row corresponding to the target word in the co-occurrence matrix.

## Co-occurrence matrix

**Context Tokens**

| | her | him | offer | smart |
|---|---|---|---|---|
| her | 7 | 5 | 2 | 2 |
| him | 5 | 12 | … 9 | 11 |
| | | ⋮ | ⋮ | |
| offer | 3 | 9 | … 0 | 3 |
| smart | 5 | 11 | 3 | 3 |

**Target Tokens**

The process is repeated for every word in every subtitle under analysis.
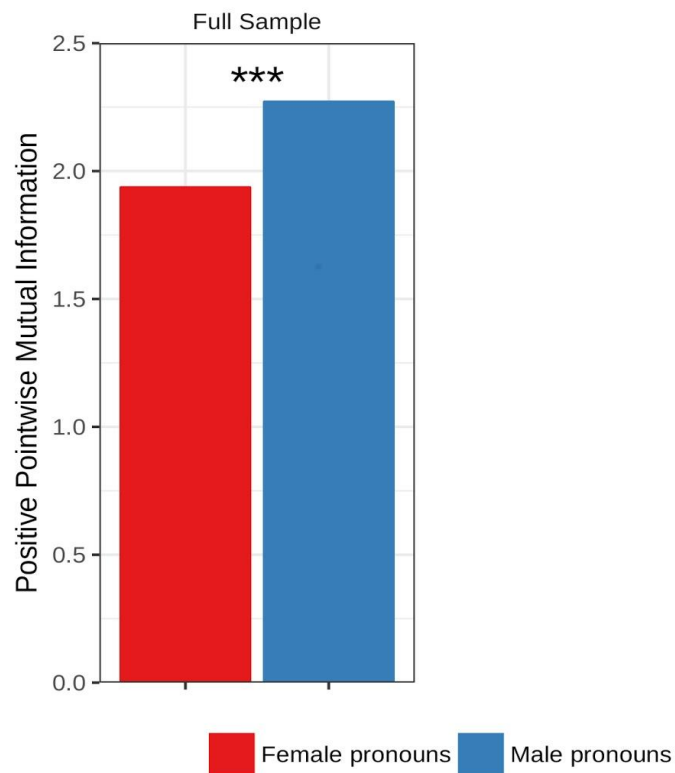
# Unifico filas y columnas de pronombres y estereotipos

➢ Female pronouns = {she, hers, her, herself }

➢ Male pronouns = {he, his, him, himself }

| | he | dancer | pilot | ... | table |
|---|---|---|---|---|---|
| she | 100 | 40 | 15 | ... | 50 |
| he | 200 | 10 | 300 | ... | 50 |
| ... | ... | ... | ... | ... | ... |
| bus | 40 | 5 | 25 | ... | 5 |
| cup | 20 | 9 | 5 | ... | 45 |

| | he | dancer | estereo | ... | table |
|---|---|---|---|---|---|
| F pron | 200 | 125 | 50 | ... | 80 |
| M pron | 400 | 30 | 500 | ... | 90 |
| ... | ... | ... | ... | ... | ... |
| bus | 40 | 5 | 1 | ... | 5 |
| cup | 20 | 9 | 2 | ... | 45 |

# *"brilliance = male"* stereotype



$$\Delta PMI = log\left(\frac{p(c|w_m)}{p(c|w_f)}\right) = 0.33$$

$$\frac{p(c|w_m)}{p(c|w_f)} = 1.26$$

**Test de significancia**
Ejemplo: Odd ratio

**Contingency table**

|       | c     | not c  | total       |
|-------|-------|--------|-------------|
| $w_f$ | $c_f$ | $nc_f$ | $c_f+nc_f$  |
| $w_m$ | $c_m$ | $nc_m$ | $c_m+nc_m$  |

# "brilliance = male" stereotype



$$\Delta PMI \quad = \quad log\left(\frac{p(c|w_m)}{p(c|w_f)}\right) = 0.33$$

$$\frac{p(c|w_m)}{p(c|w_f)} \quad = \quad 1.26$$

$$\Delta PMI \quad = \quad 0.4$$

$$\frac{p(c|w_m)}{p(c|w_f)} \quad = \quad 1.32$$
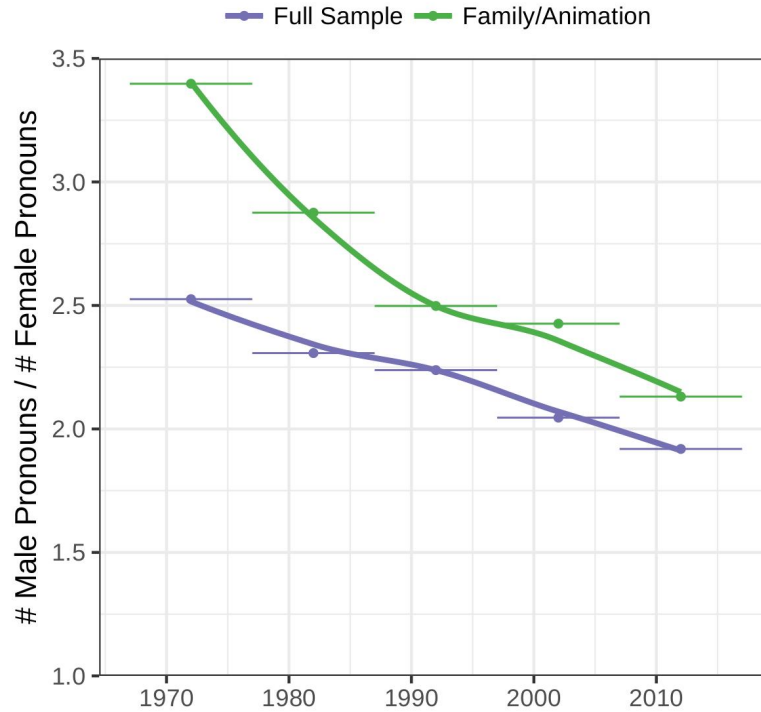
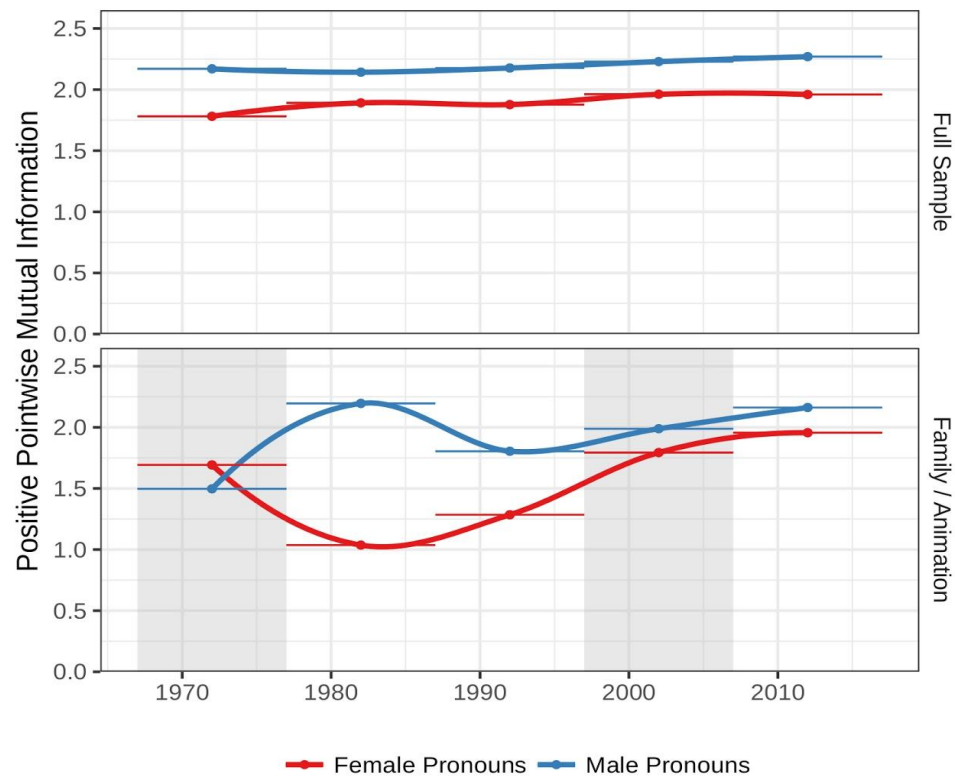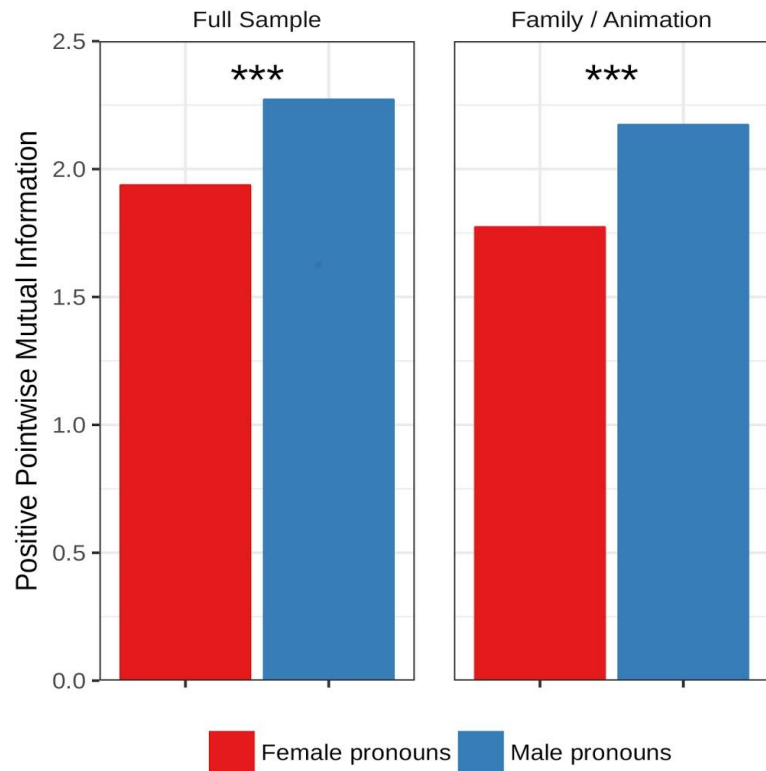# *"brilliance = male"* stereotype

# Roles estereotipados

- **Feminine** stereotyping roles
  dancer, decorator, designer, dietician, florist, homemaker, housekeeper, model, nanny, typist...

- **Masculine** stereotyping roles
  engineer, programmer, physicist, architect, detective, pilot, firefighter, inventor, mechanic, officer…

- **non-**stereotyping roles
  assistant, cashier, editor, poet, reporter, worker, doctor, lawyer, servant...

# Stereotyping roles

# Frecuencias:



Full Sample — Family/Animation

Son las películas las que realmente han estado moviendo todo en Estados Unidos desde que fueron inventadas. Te muestran qué hacer, cómo hacerlo, cuándo hacerlo y cómo sentirse al respecto.

Andy Warhol