

# *Cosas que aprendí mirando 200.000 anglicismos*

Elena Álvarez Mellado  
@lirondos

Prolegómenos

**TRABALENGUA**

**ANGLICISMOS**

**MACHINE  
LEARNING**

**YO**

**VOSOTROS**

*3 motivos por los que  
el anglicismo merece ser estudiado  
(y ninguno de ellos es el purismo lingüístico)*

## *3 motivos por los que el anglicismo merece ser estudiado (y ninguno de ellos es el purismo lingüístico)*

- El préstamo léxico es una manifestación de cómo las lenguas entran en contacto y cambian con el tiempo

## *3 motivos por los que el anglicismo merece ser estudiado (y ninguno de ellos es el purismo lingüístico)*

- El préstamo léxico es una manifestación de cómo las lenguas entran en contacto y cambian con el tiempo
- El préstamo léxico es una fuente de incorporación de nuevas palabras

# *El préstamo léxico como fuente de palabras nuevas*

- Nuevas realidades, nuevas palabras:

*online, software, podcast, streaming...*

# *El préstamo léxico como fuente de palabras nuevas*

- Nuevas realidades, nuevas palabras:

*online, software, podcast, streaming...*

- Viejas realidades, nuevas palabras:

*nude (color carne), low cost (barato), azeyte (olio)*



# La sociolingüística del anglicismo en la calle



Julià Guillamon  
@JuliaGuillamon

Vallcarca, avui.

[Translate Tweet](#)



4:45 PM · Oct 13, 2021 · Twitter Web App

## *3 motivos por los que el anglicismo merece ser estudiado (y ninguno de ellos es el purismo lingüístico)*

- El préstamo léxico es una manifestación de cómo las lenguas entran en contacto y cambian con el tiempo
- El préstamo es una fuente de incorporación de nuevas palabras
- El proceso de adaptación arroja luz sobre cuáles son las expectativas lingüísticas de los hablantes de una lengua

*Existe un interés social sobre la introducción  
de anglicismos (Fundéu, RAE)*

## Objetivo del proyecto

Crear un sistema capaz de detectar y monitorizar automáticamente el uso de préstamos léxicos en general y de anglicismos en particular.

Ejemplo: *Recetas de noviembre para el **batch cooking***

*Las prendas **bestsellers** se estampan con motivos florales, '**animal print**' o a retales tipo **patchwork***

*¿Cómo lo haríais?*

# *La ambigüedad contraataca: el diccionario no basta*

- *Prime time* es un anglicismo:
  - *prime* es una forma del verbo *primar*
  - *time* es una forma del verbo *timar*

# *La ambigüedad contraataca: el diccionario no basta*

- *Prime time* es un anglicismo:
  - *prime* es una forma del verbo *primar*
  - *time* es una forma del verbo *timar*
- *Social media* es un anglicismo:
  - Pero tanto *social* como *media* son también palabras en español
  - *Media social* no es un anglicismo

# *La ambigüedad contraataca: el diccionario no basta*

- *Prime time* es un anglicismo:
  - *prime* es una forma del verbo *primar*
  - *time* es una forma del verbo *timar*
- *Social media* es un anglicismo:
  - Pero tanto *social* como *media* son también palabras en español
  - *Media social* no es un anglicismo
- La palabra *primer*:
  - *Recomendamos usar un 'primer' hidratante*
  - *Ha ganado el primer premio*



## *La ambigüedad contraataca (II): nombres propios*

- *La compañía Apple triplica resultados este trimestre*

## *La ambigüedad contraataca (II): nombres propios*

- *La compañía Apple triplica resultados este trimestre*
- *Apple triplica resultados este trimestre*

## *La ambigüedad contraataca (II): nombres propios*

- *La compañía Apple triplica resultados este trimestre*
- *Apple triplica resultados este trimestre*
- *Ghosting, la tendencia que causa furor en redes*

*Cuando la casuística del sistema de reglas es demasiado complejo y no consigue dar cuenta del fenómeno que queremos abordar  
→ Aprendizaje automático*

## *El aprendizaje automático muy muy muy resumido*

- Si le damos a un ordenador una cantidad suficiente de ejemplos del fenómeno que queremos modelar, el ordenador aprenderá a reconocerlo
- Aprendizaje por imitación (¿es verdaderamente aprendizaje? ¿es inteligencia?)

# *La metáfora del seguro de automóvil*

- Cuando compras un seguro de coche, el precio de la póliza dependerá de diversas características: edad, color del coche, experiencia al volante, partes previos:
  - Una mujer de 52 años, 27 años de antigüedad, coche nuevo blanco
  - Un hombre de 23, 1 año de antigüedad, coche viejo rojo
- Imaginad que tuviésemos miles de ejemplos así, y supiésemos además si la persona dio algún parte o no. ¿Podríamos asignar pesos a cada uno de estos rasgos y calcular la probabilidad de que un conductor dé un parte?

## ¿Cómo lo hacemos?

- En vez de partes del seguro y conductores, pensad en textos.
- Necesitamos una colección de textos/un corpus rico en el fenómeno que queremos modelar
  - Necesitamos un corpus de español rico en anglicismos
- En el corpus deberemos marcar a mano qué cosas son anglicismo y cuáles no (anotación)
- Al ordenador le daremos ese corpus anotado con anglicismos
  - El sistema aprenderá cuáles son las características de los anglicismos de nuestro corpus. Cuando le demos una palabra nueva podrá intentar adivinar si es un anglicismo o no
- Es básicamente reconocimiento de patrones. El ordenador no aprende realmente nada ni tiene inteligencia (tal y como la entendemos los humanos)

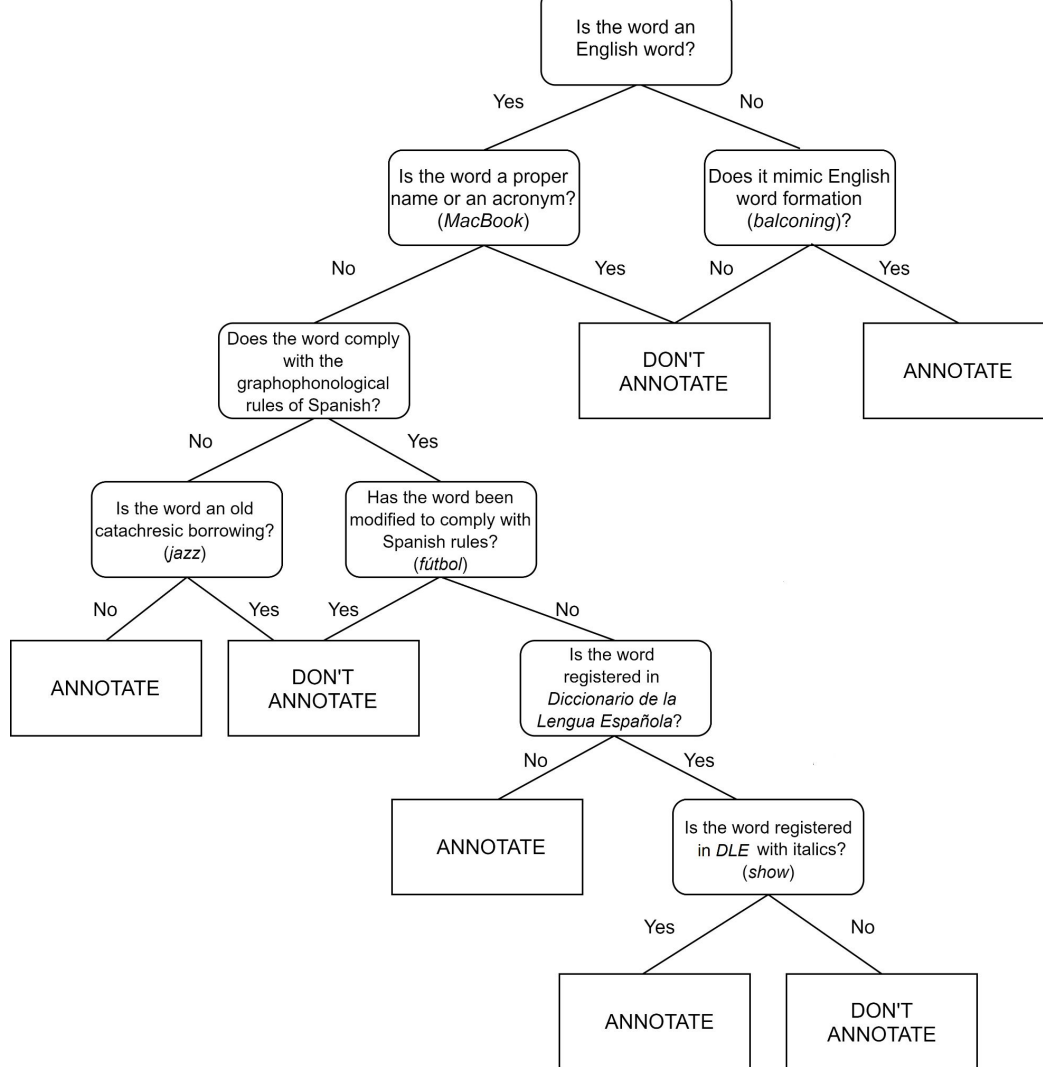
## *La metáfora del examen de matemáticas*

- Le damos al ordenador un montón de problemas resueltos (corpus de entrenamiento)
- El ordenador tiene acceso tanto al problema como a la solución
- El ordenador “aprenderá” a partir de esos problemas resueltos
- Cuando el modelo ya está entrenado (= el ordenador ha terminado de “estudiar”) le damos problemas nuevos (el examen). Esta vez SIN darle la solución (datos de evaluación)
- Evaluamos qué tal ha hecho el examen nuestro modelo



# 1. *Construir un corpus anotado con anglicismos*

- Recopilar 21,570 titulares de prensa en español de elDiario.es
- Marcar a mano cuáles de esas palabras son anglicismos y cuáles no. Esto es más fácil decirlo que hacerlo → guía de anotación (ponga un lingüista en su vida)



En O  
este O  
mes O  
especialmente O  
puede O  
ser O  
de O  
utilidad O  
apuntarnos O  
al O  
batch B-ENG  
cooking I-ENG

Benching B-ENG  
, O  
estar O  
en O  
el O  
banquillo O  
de O  
tu O  
crush B-ENG  
mientras O  
otro O  
juega O  
de O  
titular O

## 2. Transformar los datos en algo que el ordenador entienda

La fiebre de los podcasts llega a España

○ ○ ○ ○ ENG ○ ○ ○

¿Os acordáis del ejemplo del seguro de coche? Pensad en palabras:

- Cada palabra estará representada por un conjunto de rasgos: qué caracteres que conforman la palabra, si va en mayúscula o no, si lleva signos de puntuación o no, etc.
- La selección de rasgos es una tarea fundamental y muy lingüística (tanto la observación experimental como la intuición lingüística pueden ayudar)

## *De palabras a rasgos*

Cada palabra estará representada por los siguientes rasgos:

- la palabra en sí: podcast
- terminación: 'ast'
- trigramas de caracteres: 'pod', 'odc', 'dca', 'cas', 'ast'
- mayúscula/minúscula: minúscula
- palabras adyacentes: antecedita por comilla (modelo secuencial)

Durante el entrenamiento, el modelo “aprenderá” cuánto de importantes son estos rasgos y cuáles son los indicios de que algo sea o no un anglicismo.

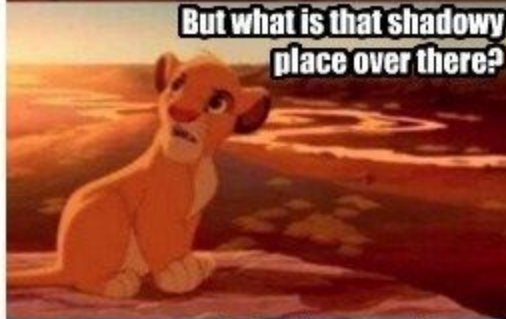
## *Más allá de la forma*

- Nuestros rasgos giran todos en torno a la forma de la palabra
- La forma puede representar aspectos morfológicos, pero no relativos al significado de las palabras
- ¿Es posible representar la semántica de una forma numérica que un ordenador pueda procesar?

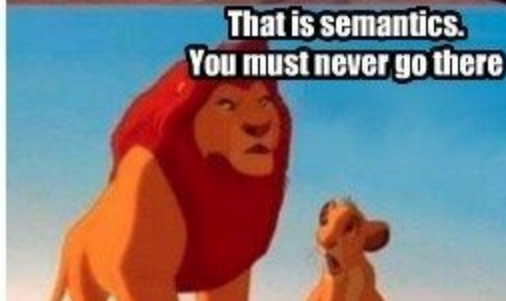
**This is linguistics. You can study any  
aspect of language you want, and you  
will always find structure**



**But what is that shadowy  
place over there?**



**That is semantics.  
You must never go there**



## *Palabras como vectores (aka word embeddings)*

El agua hierve a 100°C

He pedido un vaso de limonada

Luisa no bebe limonada

Prefiero beber un vaso de agua

Siempre come helado de  
chocolate de postre



## *Palabras como vectores (aka word embeddings)*

El agua hierve a 100°C

He pedido un vaso de limonada

Luisa no bebe limonada

Prefiero beber un vaso de agua

Siempre come helado de  
chocolate de postre

	hervir	beber	vaso	chocolate
agua				
limonada				
helado				

## *Palabras como vectores (aka word embeddings)*

El agua hierve a 100°C

He pedido un vaso de limonada

Luisa no bebe limonada

Prefiero beber un vaso de agua

Siempre come helado de  
chocolate de postre

	hervir	beber	vaso	chocolate
agua	1	1	1	0
limonada	0	1	1	0
helado	0	0	0	1

agua = (1, 1, 1, 0)

limonada = (0, 1, 1, 0)

helado = (0, 0, 0, 1)

## Palabras como vectores (aka word embeddings)

El agua hierve a 100°C

He pedido un vaso de limonada

Luisa no bebe limonada

Prefiero beber un vaso de agua

Siempre come helado de chocolate de postre

	hervir	beber	vaso	chocolate
agua	1	1	1	0
limonada	0	1	1	0
helado	0	0	0	1

agua = (1, 1, 1, 0)

limonada = (0, 1, 1, 0)

helado = (0, 0, 0, 1)

Estos 2 son  
más parecidos  
=  
espacialmente  
más cerca

### 3. *Cómo entrenar a tu ~~dragón~~ modelo*

- Tenemos nuestro corpus anotado
- Hemos transformado los datos en algo que el ordenador pueda beber (rasgos, embeddings)
- Ahora, a entrenar (a estudiarse los exámenes resueltos)
  - El entrenamiento como dardos

## 4. *Hora de evaluar*

- Cuando el entrenamiento ha terminado, evaluamos
- Proporcionamos ejemplos nuevos al modelo (sin solución). El modelo analizará cada palabra y a la luz de los rasgos que observe intentará adivinar qué palabras son anglicismos
- Es decir, nuestro modelo está ahora haciendo predicciones

# *La evaluación*

- ¿Cómo de bien le sale el examen?
- Para evaluar un modelo tenemos en cuenta 2 aspectos:
  - Cuántas palabras marcadas por el modelo como “anglicismo” eran verdaderamente anglicismos
  - Cuántos anglicismos se dejó el modelo fuera (eran anglicismos y modelo los marcó como no anglicismo)
  - La combinación de estas dos se llama valor F1

Para el modelo de Lázaro,  $F1 = 87$  (100 sería un modelo perfecto). No es perfecto, ¡pero no está mal!

## 5. *Creando Observatorio Lázaro*

- Ya tenemos un modelo capaz de detectar automáticamente anglicismos en prensa española. ¡Ya podemos crear el observatorio!
- El observatorio es un sistema que:
  - Extrae los artículos publicados cada día en 8 medios españoles de prensa escrita: El País, La Vanguardia, ABC, El Mundo, elDiario, 20minutos, El Confidencial, EFE
  - Los artículos se mandan al modelo
  - El modelo extrae los anglicismos que aparecen en los artículos
  - Cada anglicismo se almacena en una base de datos junto a la información del avistamiento (fecha, URL, medio, contexto)

## 6. *Visualizando, que es gerundio*

Con la información almacenada podemos:

- ver cómo cambia la frecuencia de uso a lo largo del tiempo
- hacer resumen con los anglicismos nuevos del día



## Observatorio Lázar

Observatorio del anglicismo en la prensa española



8

Medios observados



129 105 365

Palabras analizadas



258 021

Anglicismos totales



12 527

Anglicismos únicos

Observatorio Lázar analiza y extrae automáticamente los anglicismos aparecidos en las noticias del día de ocho medios de comunicación españoles: elDiario.es, El País, El Mundo, ABC, La Vanguardia, El Confidencial, 20minutos y EFE. Las gráficas que aparecen en esta página se actualizan semanalmente. [Más sobre Observatorio Lázar](#).

## Anglicismos más frecuentes

Evolución de los 10 anglicismos más frecuentes (frecuencia por cada millón de palabras)

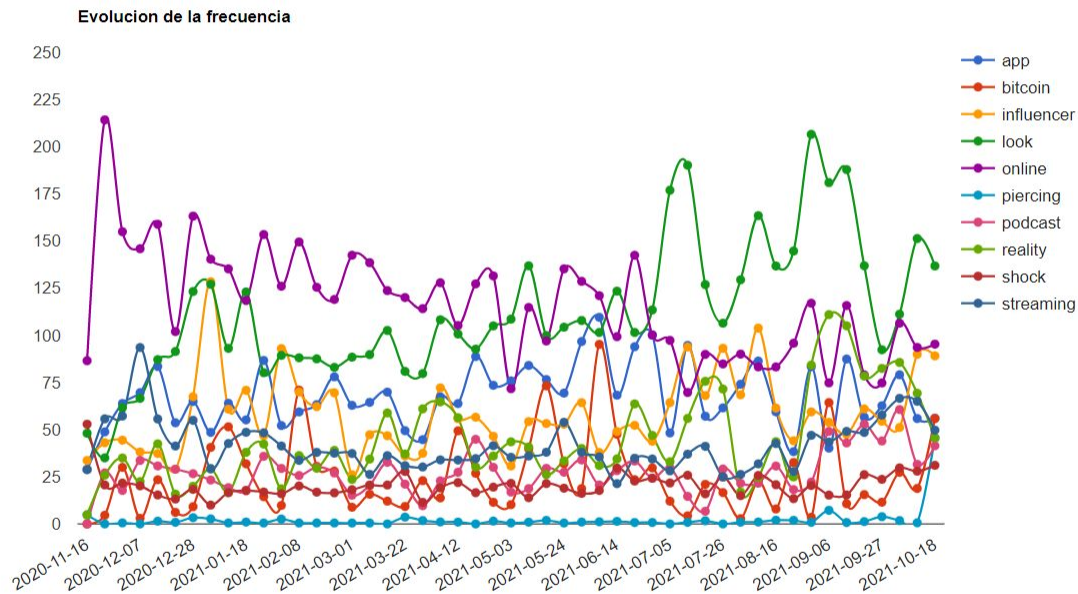
Evolucion de la frecuencia

250

2010 2011 2012 2013 2014 2015 2016 2017 2018 2019 2020 2021 2022 2023 2024 2025

## Anglicismos más frecuentes

Evolución de los 10 anglicismos más frecuentes (frecuencia por cada millón de palabras)



# podcast

Anglicism: yes

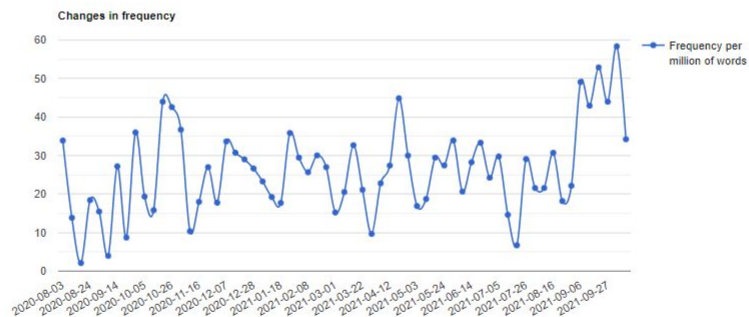
Forms: podcast podcasts

Average frequency\*: 26.313

Frequency during the last 30 days\*: 47.818

Newspaper sections: Portada Televisión Cultura Gente Economía

\* Frequency per million of words (measured from August 2020 onwards).

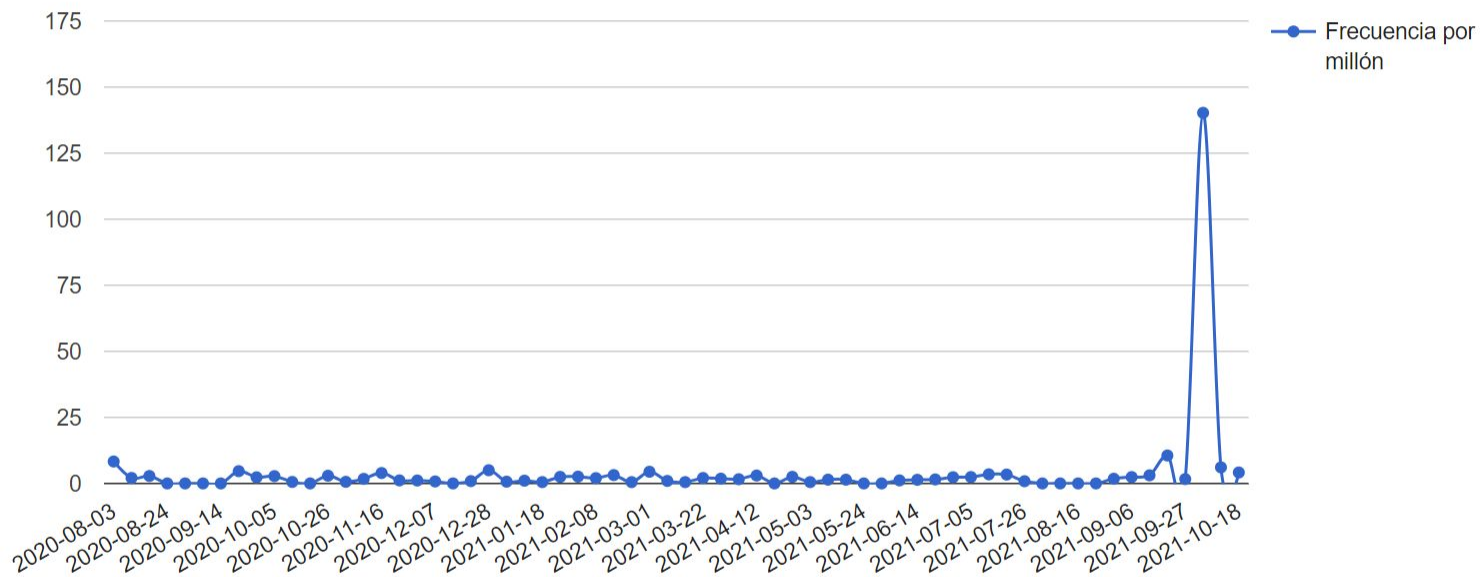


Show  entries

Search:

Borrowing <sup>*</sup>	Context	Newspaper	Date
podcast	Grito Sordo y que guioniza el contenido que propone Ignatius Farray, como el 'podcast' 'Payasos y fuego', En	elpais	18-10-2021

### Evolucion de la frecuencia



# offshore

Anglicismo: sí

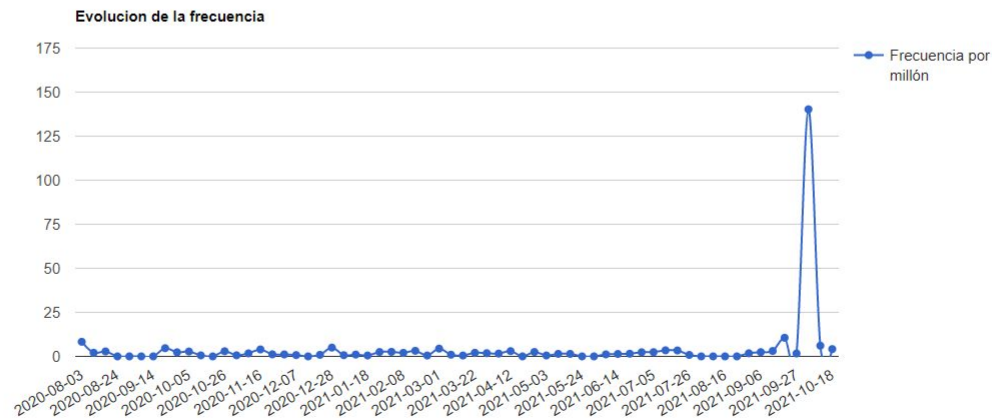
Formas: offshore   offshores

Frecuencia media de aparición\*: 4.143

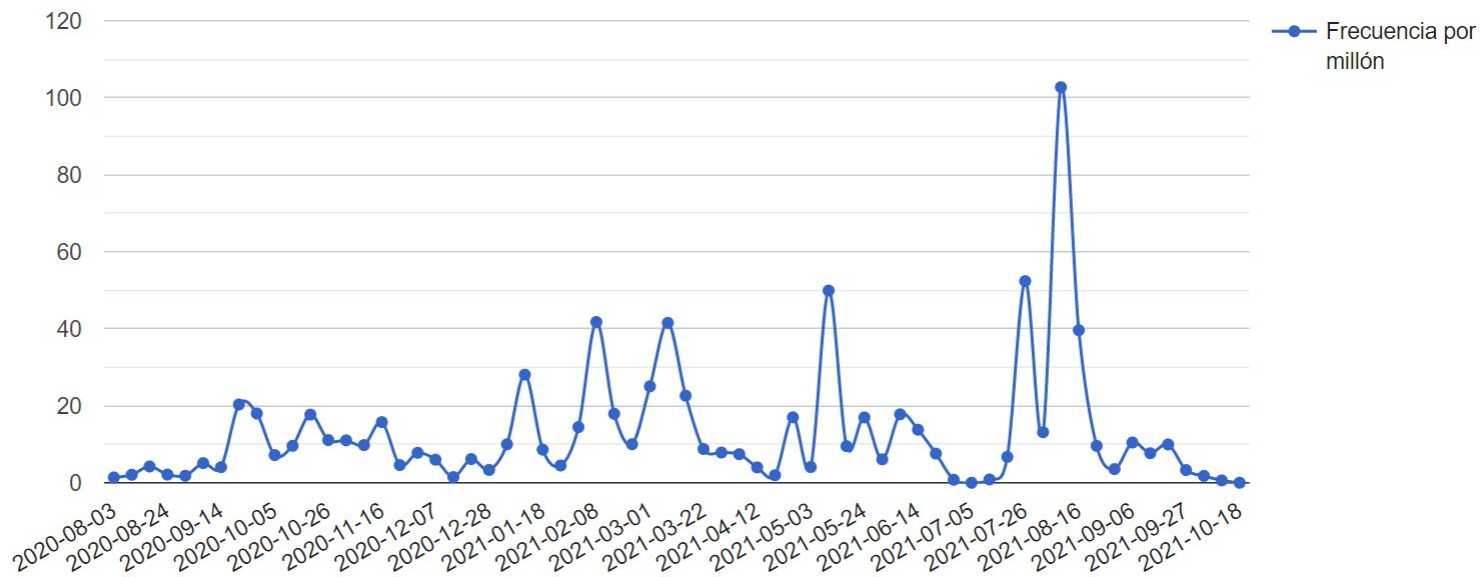
Frecuencia de aparición en el último mes\*: 35.18

Secciones habituales: Economía   Portada   Internacional   Política   España

\* Frecuencia por cada millón de palabras medida desde agosto de 2020.



### Evolucion de la frecuencia



# rider

Anglicismo: sí

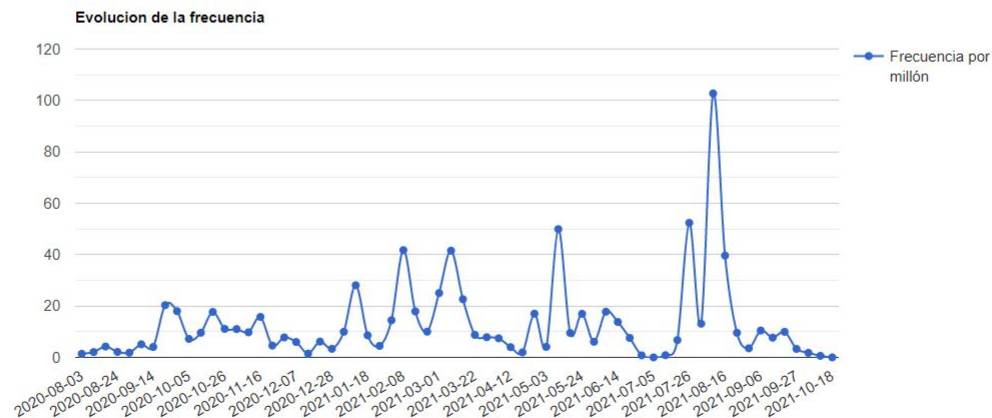
Formas: rider riders

Frecuencia media de aparición\*: 12.965

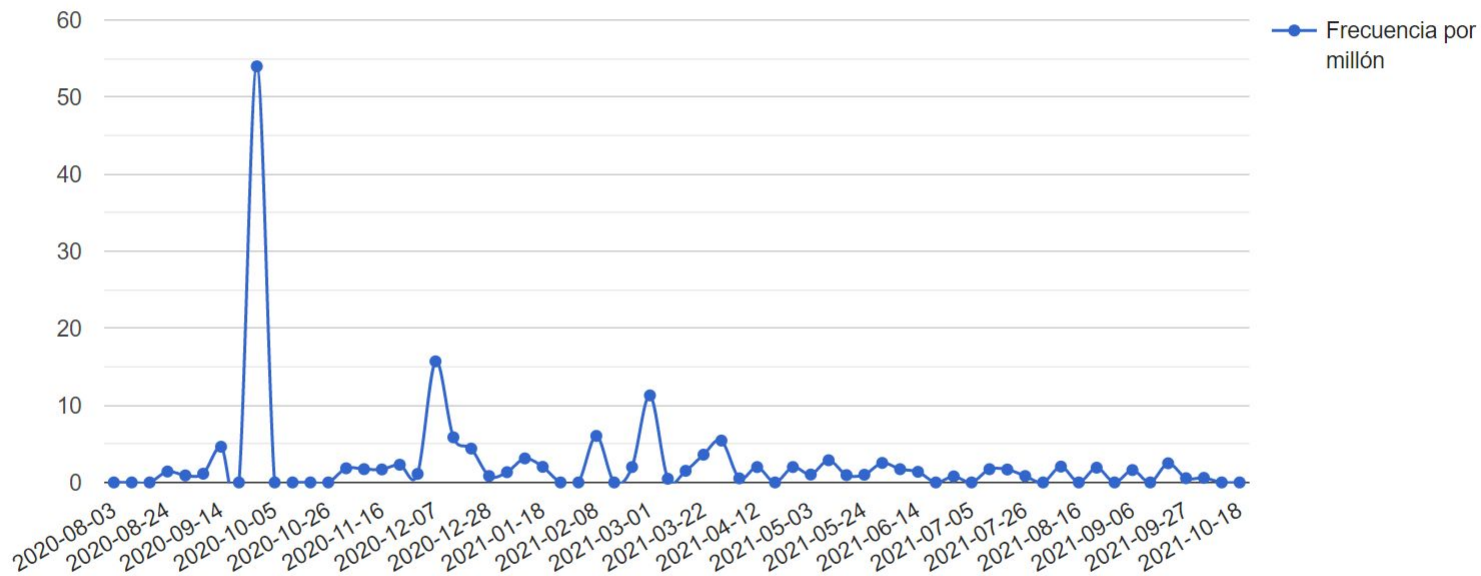
Frecuencia de aparición en el último mes\*: 2.39

Secciones habituales: Economía Portada Tecnología España Medio Ambiente

\* Frecuencia por cada millón de palabras medida desde agosto de 2020.



Evolucion de la frecuencia





# black

Anglicismo: sí

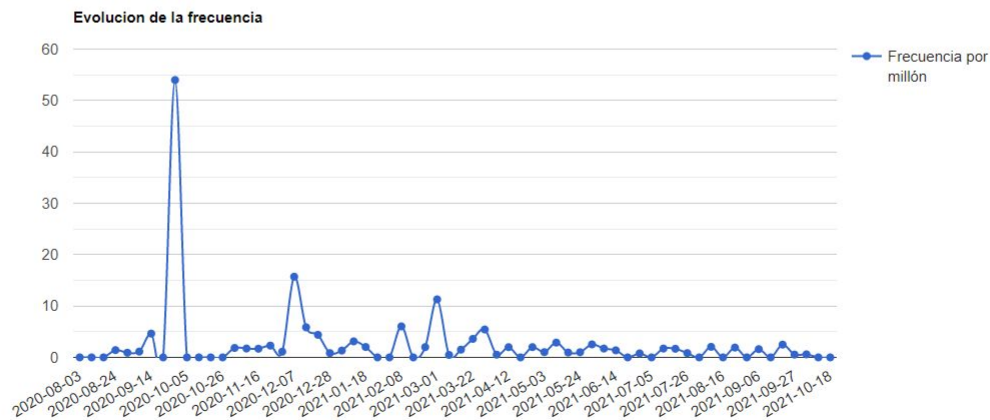
Formas: black

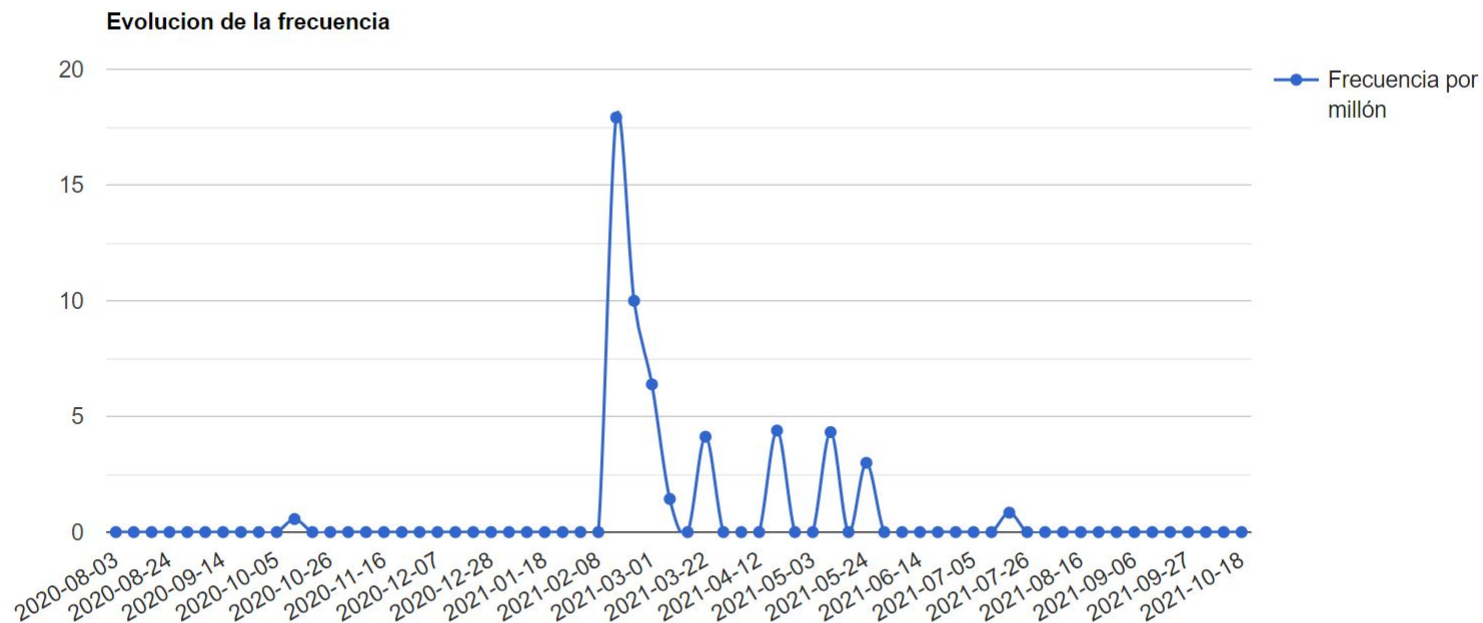
Frecuencia media de aparición\*: 2,647

Frecuencia de aparición en el último mes\*: 0,531

Secciones habituales: Portada Economía España Política Cultura

\* Frecuencia por cada millón de palabras medida desde agosto de 2020.





# foam

Anglicismo: sí

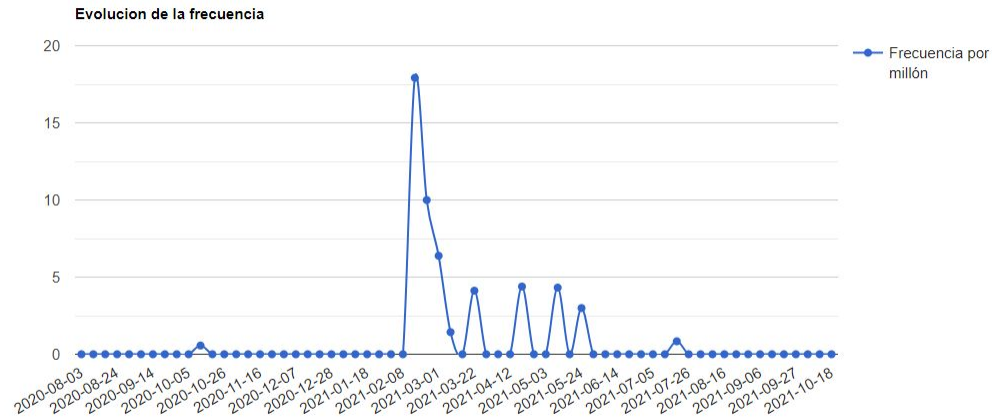
Formas: foam

Frecuencia media de aparición\*: 0.997

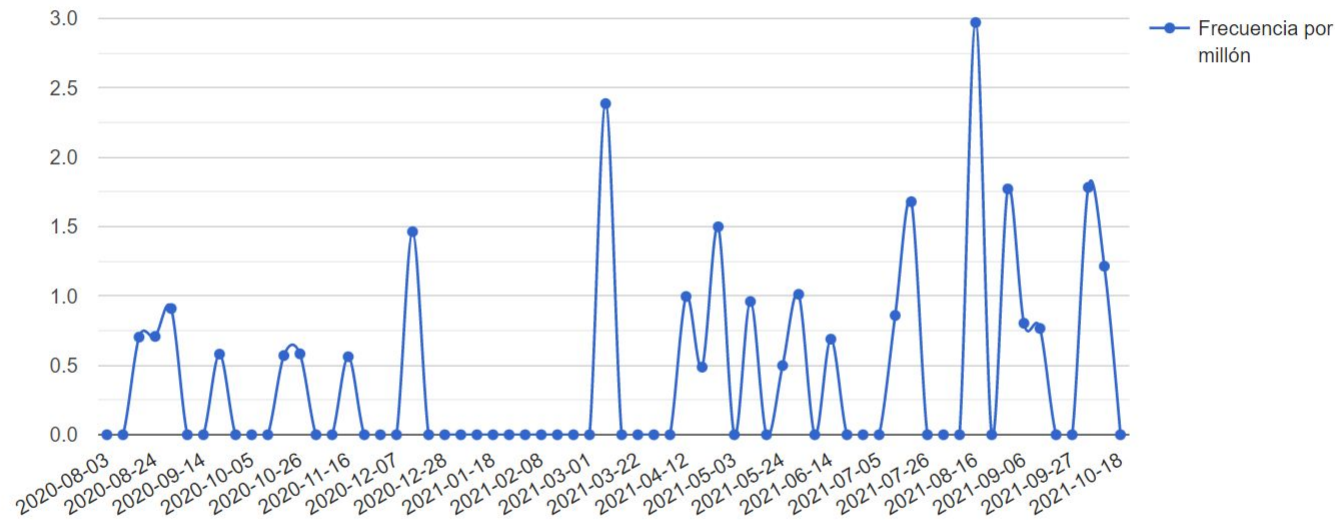
Frecuencia de aparición en el último mes\*: 0

Secciones habituales: Portada Política España Opinión Moda

\* Frecuencia por cada millón de palabras medida desde agosto de 2020.



Evolucion de la frecuencia



# booster

Anglicismo: sí

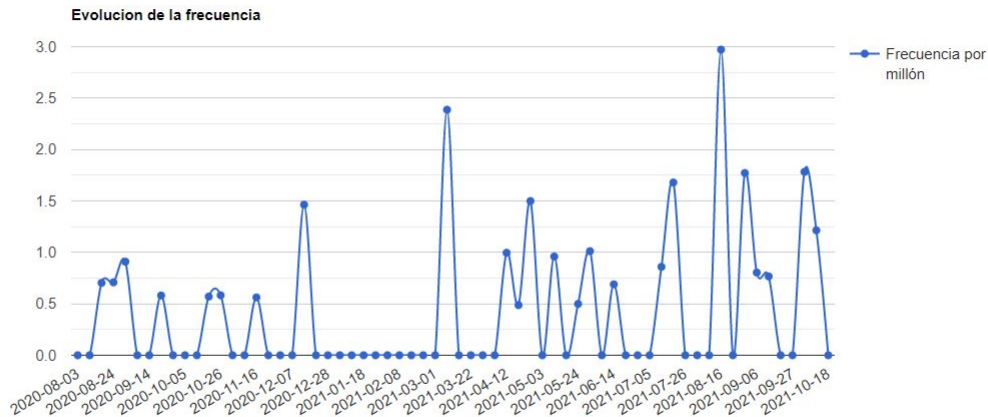
Formas: booster boosters

Frecuencia media de aparición\*: 0.393

Frecuencia de aparición en el último mes\*: 0.664

Secciones habituales: Moda Salud Sociedad Portada España

\* Frecuencia por cada millón de palabras medida desde agosto de 2020.



anxiety baking

"...milénicos son especialmente dados al anxiety baking, la práctica de preparar..."

[Translate Tweet](#)



Horneamos por encima de nuestras posibilidades  
El frenesí repostero de la cuarentena se explica por la necesidad de llenar horas muertas, las ansias de aplausos en las redes sociales y la búsqueda de un ...  
[elpais.com](#)

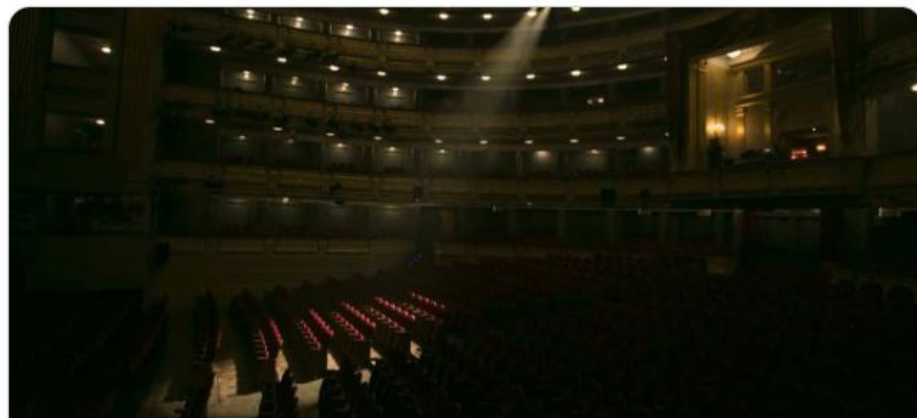
6:35 PM · May 3, 2020 · [lazarobot](#)

3 Retweets 9 Likes

old school

"...en su diseño básico muy old school, pero también decisivo en..."

[Translate Tweet](#)



## EL PAÍS

La función no puede continuar

El dramaturgo alemán Roland Schimmelpfennig llora el cierre de los teatros en este artículo escrito durante el confinamiento por el coronavirus

[elpais.com](#)



Observatorio Lázaro

@lazarobot



male tears

"...y a formatos clásicos de meme) para tratar algunos temas, como las male tears vertidas tras el estreno del videojuego Last of us 2, cuya..."

[Translate Tweet](#)



verne.elpais.com

Videojuegos y el machismo que siempre pica

Hablamos con la creadora de Feminismoen8bits, una cuenta de Instagram en la que se habla de videojuegos desde una ...

4:12 PM · Nov 23, 2020 · lazarobot





Observatorio Lázarro

@lazarobot



red carpet

"...Laguna solo hace vestidos largos de red carpet, este look de Clara Courel..."

[Translate Tweet](#)



vanitatis.elconfidencial.com

Señoras y señores, comienza la MBFWM: la crónica de las ...  
Repasamos los momentos más señalados de la pasarela  
madrileña

9:53 PM · Sep 12, 2020 · lazarobot

Crying

"...: . ORIGINAL:. . Cryings not for me. . ..."



Esta es la playlist que necesitas para salir a la calle y disfrutar del aire fresco  
Si buscas temazos para practicar ejercicio al aire libre o, simplemente, quieres dar  
un paseo y mantener el ánimo arriba, aquí tienes lo que necesitas.

 [20minutos.es](https://20minutos.es)

brilli-brilli

"...glitter o brilli-brilli de la artista..."

[Translate Tweet](#)



## DESIGN

Tres españoles recopilan más de 300 obras de 'Arte Covid' de todo el mundo. ¿H...  
Cuando esto acabe será fundamental recoger un testimonio emocional y artístico  
de cómo el virus nos ha afectado , dice uno de los impulsores de Covid Art ...

[elpais.com](#)



Observatorio Lázarobot

@lazarobot



shosho

"...eso sería como el shosho de la Bernarda...."

[Translate Tweet](#)



elpais.com

El 'divorcio' obligado de los Estopa

David y José Muñoz, 21 años después de su primer éxito, llevan una vida corriente y gozan de un sólido patrimonio qu...

8:26 PM · May 5, 2020 · lazarobot

1 Quote Tweet   2 Likes



**Observatorio Lázarobot**

@lazarobot



date prisa

"...Si te decides por estas, date prisa y fíchalas por el mismo precio..."

[vanitatis.elconfidencial.com/estilo/moda/20...](https://vanitatis.elconfidencial.com/estilo/moda/20...)

[Translate Tweet](#)

4:46 PM · May 4, 2020 · lazarobot

# ¿Preguntas?

Elena Álvarez Mellado  
<http://observatoriolazaro.es/>  
@lirondos @lazarobot