

# Texto como data

Fernanda Sobrino

6/19/2021

Introducción

Análisis cuantitativo de texto

Proceso de análisis de texto

Características de un texto

# Introducción

## Historicamente (pre 2000)

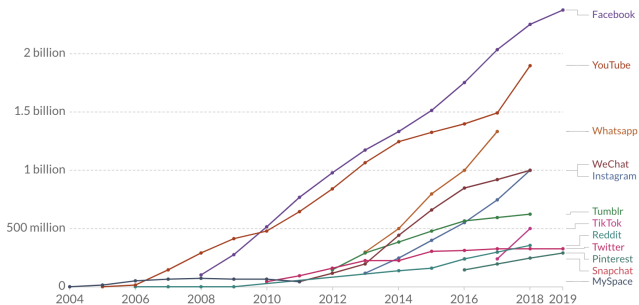
- ▶ las interacciones sociales se dan por medio de texto
- ▶ las ciencias sociales han evitado estudiar este tipo de datos
- ▶ Por qué?
  - ▶ difíciles de conseguir
  - ▶ pérdida de tiempo
  - ▶ no generalizables (formatos distintos, etc)
  - ▶ difícil de guardar y alcanzar (archivos históricos)
  - ▶ intensivos computacionalmente

# Post 2000: Redes Sociales

## Number of people using social media platforms, 2004 to 2019

Estimates correspond to monthly active users (MAUs). Facebook, for example, measures MAUs as users that have logged in during the past 30 days. See source for more details.

[+ Add data](#)



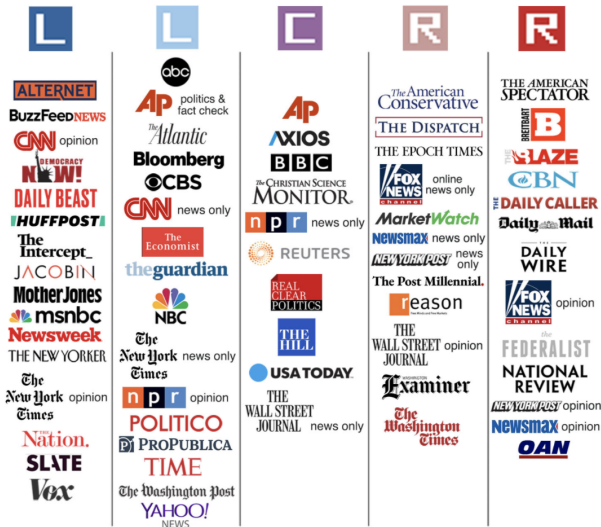
Source: Statista and TNW (2019)

CC BY

# Post 2000: Medios tradicionales

## AllSides™ Media Bias Chart

All ratings are based on online content only — not TV, print, or radio content.  
Ratings do not reflect accuracy or credibility; they reflect perspective only.



L LEFT 
 L LEAN LEFT 
 C CENTER 
 R LEAN RIGHT 
 R RIGHT

<sup>100</sup> Inquiries numbered 7, 16, and 17 are not to be asked in respect to infants. Inquiries numbered 11, 12, 15, 16, 17, 19, and 20 are to be answered (if at all) merely by an affirmative mark as /.

Post Office: New York City John Miller Ass't Marshal

[illegible]

# Por qué paso esto?

- ▶ incremento en la disponibilidad de texto no estructurado
- ▶ reducción en los costos de almacenamiento 1956 \$10,000 por megabyte, 2014 menos de 0.0001 por magabyte
- ▶ explosión en métodos, lenguajes de programación y librerías capaces de lidiar con texto
  - ▶ generalizable: un mismo método puede ser usado entre bases de datos
  - ▶ sistemático: parámetros y estadísticas comparables entre datos
  - ▶ barato: R, Python, Julia, etc pueden lidiar con estos datos y cualquiera puede hacerlo desde su casa



## Anal  s cuantitativo de texto

# Procesamiento del texto a datos estructurados

When I presented the supplementary budget to this House last April, I said we could work our way through this period of severe economic distress. Today, I can report that notwithstanding the difficulties of the past eight months, we are now on the road to economic recovery.

In this next phase of the Government's plan we must stabilise the deficit in a fair way, safeguard those worst hit by the recession, and stimulate crucial sectors of our economy to sustain and create jobs. The worst is over.

This Government has the moral authority and the well-grounded optimism rather than the cynicism of the Opposition. It has the imagination to create the new jobs in energy, agriculture, transport and construction that this green budget will

docs	words													
	made	because	had	into	get	some	through	next	where	many	irish			
t06_kenny_fg	12	11	5	4	8	4	3	4	5	7	10			
t05_cowen_ff	9	4	8	5	5	5	14	13	4	9	8			
t14_gaolainh_sf	3	3	3	4	7	3	7	2	3	5	6			
t01_lenihan_ff	12	1	5	4	2	11	9	16	14	6	9			
t11_gormley_green	0	0	0	3	0	2	0	3	1	1	2			
t04_morgan_sf	11	8	7	15	8	19	6	5	3	6	6			
t12_ryan_green	2	2	3	7	0	3	0	1	6	0	0			
t10_guinn_lab	1	4	4	2	8	4	1	0	1	2	0			
t07_odonnell_fg	5	4	2	1	5	0	1	1	0	3	0			
t09_higgins_lab	2	2	5	4	0	1	0	0	2	0	0			
t03_burton_lab	4	8	12	10	5	5	4	5	8	15	8			
t13_cuffe_green	1	2	0	0	11	0	16	3	0	3	1			
t08_gilmore_lab	4	8	7	4	3	6	4	5	1	2	11			
t02_bruton_fg	1	10	6	4	4	3	0	6	16	5	3			

Descriptive statistics  
on words

Scaling documents

Classifying documents

Extraction of topics

Vocabulary analysis

Sentiment analysis

# Métafora del pajar de Grimmer

El análisis automatizado mejoran la lectura

- ▶ analizar una pajita de heno  $\equiv$  entender el significado de una oración
  - ▶ Humanos: **Muy buenos**
  - ▶ Máquina: **Bastante malas**
- ▶ organizar todo el pajar  $\equiv$  describir, clasificar y escalar textos
  - ▶ Humanos: **Bastante malos**
  - ▶ Máquina: **Muy buenas**

# Principios del analisis cuantitativo de texto (Grimmer & Stuart 2013)

1. Todos los modelos están mal pero algunos son útiles
  - ▶ proceso de generación del texto desconocido
  - ▶ complejidad del lenguaje: temporalidad, sinónimos, ironía, sarcasmo, etc.
  - ▶ los modelos fallan al intentar capturar el lenguaje pero son buenos para tareas específicas
2. Los métodos cuantitativos aumentan las capacidades humanas pero no remplazan al humano
  - ▶ **algoritmo**: organizan, dirigen, sugieren
  - ▶ **humano**: lee e interpreta

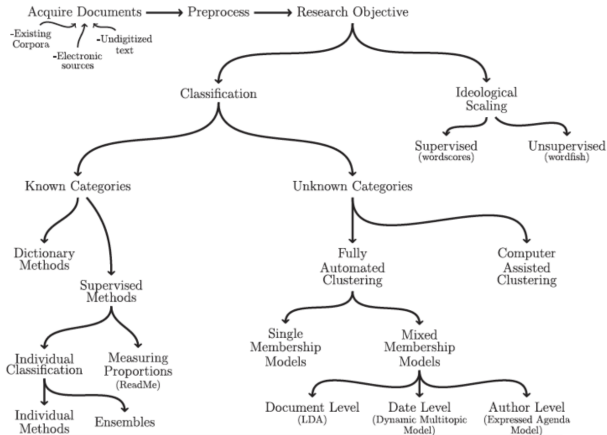
# Principios del analisis cuantitativo de texto (Grimmer & Stuart 2013)

3. No existe un un método mejor que otro para automatizar texto
  - ▶ aprendizaje supervisado → categorías pre determinadas
  - ▶ aprendizaje no supervisado → categorías por clasificar
4.  $\implies$  validar, validar, validar
  - ▶ distinto rendimiento dependiendo de la tarea en cuestión
  - ▶ poca teoría que sustente que método escoger
  - ▶ **Evitar** aplicar métodos ciegamente

# Supuestos del analisis cuantitativo del texto

1. El texto representa alguna característica de interés
  - ▶ un atributo del autor
  - ▶ sentimiento/emoción
  - ▶ relevancia de alguna cuestión política, etc
2. El texto se puede representar extrayendo sus características
  - ▶ bolsa de palabras (bag of words)
  - ▶ word embeddings, etc
3. Una matriz de características del documento puede analizarse usando métodos cuantitativos para producir estimados validos y significativos de la característica de interés

# Métodos de análisis de texto



## Proceso de análisis de texto



# Proceso de análisis de texto

1. Seleccionar el texto: definir el corpus
2. Convertir el texto a algo que la máquina pueda leer
3. Definir los documentos es decir la unidad de análisis
4. Definir las características: tokens, diccionarios, lenguaje, etc
5. Convertir las características a una matriz
6. Escoger un método cuantitativo para el análisis de las matrices
7. Presentar los resultados obtenidos

## Ejemplo: Hoberg and Phillips 2016

- ▶ nueva manera de clasificar industrias usando la descripción de sus productos
- ▶ esto permite clasificaciones de industria mas flexibles y capaces de evolucionar con el tiempo
- ▶ por qué queremos saber esto? analizar mejor los efectos de la competencia y la oferta de distintos productos

# Ejemplo: Hoberg and Phillips 2016

Cómo lo hacen?

- ▶ todas las empresas que cotizan en bolsa tienen reportes 10-K
- ▶ estos reportes describen los productos que cada empresa provee
- ▶ para cada empresa  $i$  cuentan cuantas veces aparece cada palabra y lo representan como  $c_i$
- ▶ calculan que tan diferentes son dos productos usando el angulo entre los vectores  $c_i$  y  $c_j$

## Ejemplo: Hoberg and Phillips 2016

- ▶ definen las industrias agrupando las empresas con las  $c$ 's mas cercanas
- ▶ 300 grupos distintos para que coincida con la clasificación de industrias pre existentes
- ▶ asignan una industria a cada una de las empresas por año  $\hat{v}_{it}$
- ▶ examinan el efecto de shocks a la industria militar y de software en la competencia y la oferta de productos

# Ejemplo: Hoberg and Phillips 2016

## SubMarket 1: Entertainment (Sample Focal Firm: Wanderlust Interactive)

43 rivals: Maxis, Piranha Interactive Publishing, Brilliant Digital Entertainment, Midway Games, Take Two Interactive Software, THQ, 3DO, New Frontier Media, ...

SIC codes of rivals: computer programming and data processing [sic3=737] (24 rivals), motion picture production and allied services [sic3=781] (4 rivals), misc other (13 rivals)

Core words: entertainment (42), video (42), television (38), royalties (35), internet (34), content (33), creative (31), promotional (31), copyright (31), game (30), sound (29), publishing (29), ...

## SubMarket 2: Medical services (Sample Focal Firm: Quadramed Corp)

66 rivals: IDX Systems, Medicus Systems, Hpr, Simione Central Holdings, National Wireless Holdings, HCIA, Apache Medical Systems, ...

SIC codes of rivals: computer programming and data processing [sic3=737] (45 rivals), insurance agents, brokers, and service [sic3=641] (5 rivals), miscellaneous health services [sic3=809] (4 rivals), management and public relations services [sic3=874] (3 rivals), misc other (9 rivals)

Core words: client (59), database (54), solution (49), patient (47), copyright (47), secret (47), physician (47), hospital (46), healthcare (46), server (45), resource (44), functionality (44), billing (44), ...

## SubMarket 3: Information Transmission (Sample Focal Firm: FAXSAV)

259 rivals: Omtool Ltd, Concentric Network, Premiere Technologies, International Telecommunication Data Systems, IDT Corp, Axent Technologies, Solopoint, Precision Systems, Netrix Corp, ...

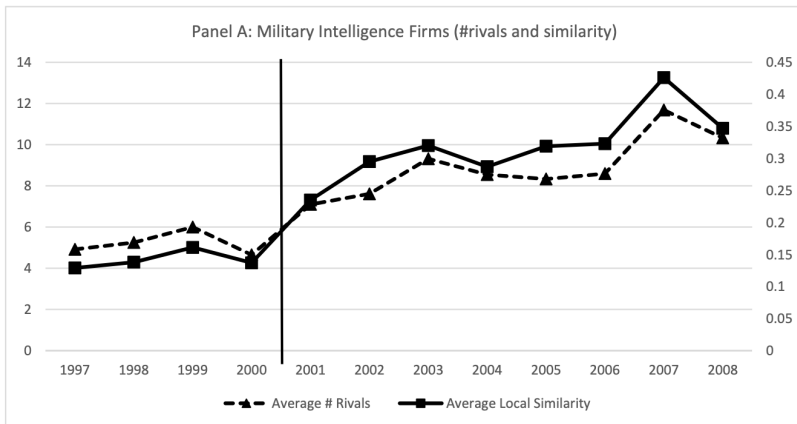
SIC codes of rivals: computer programming and data processing [sic3=737] (112 rivals), communications equipment [sic3=366] (45 rivals), telephone communications [sic3=481] (38 rivals), computer and office equipment [sic3=357] (29 rivals), communications services, other [sic3=489] (7 rivals), miscellaneous business services [sic3=738] (7 rivals), misc other (15 rivals)

Core words: internet (236), telecommunications (211), interface (194), communication (188), solution (187), platform (184), architecture (182), call (177), infrastructure (173), voice (173), functionality (173), server (173), ...

## Ejemplo: Hoberg and Phillips 2016

Resultado: después del 11 de Septiembre del 2001 hay entrada de empresas al sector militar de inteligencia

### Military Intelligence Firms: Competitor Changes



## Características de un texto

# Algunas definiciones

- ▶ **corpus:** conjunto amplio y estructurado de texto para analizar
- ▶ **documento:** cada una de las unidades del corpus
- ▶ **tokens:** partes del documento pueden ser palabras, frases, oraciones, etc.
- ▶ **stems:** reducir una palabra a su raíz por ejemplo hablar → habl
- ▶ **lemmas:** toma en consideración el análisis morfológico de la palabra niñas/niños → niño
- ▶ **stop words:** palabras que por lo general se excluyen por no aportar nada al análisis, depende de la aplicación si las dejas o no. Ejemplos: un, una, es, en, por, pero



