

Análisis de Comunidades Virtuales con R



Felipe Ortega
GSyC/Libresoft
Universidad Rey Juan Carlos

[@identica](https://twitter.com/identica)
[@jfelipe](https://twitter.com/jfelipe)



© 2011 Felipe Ortega.

Algunos derechos reservados.

Este documento se distribuye bajo una licencia
Creative Commons Reconocimiento-CompartirIgual 3.0,
disponible en

<http://creativecommons.org/licenses/by-sa/3.0/es/>

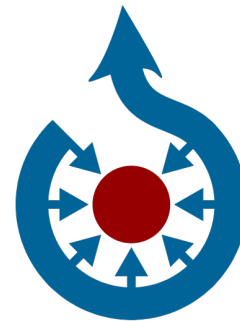
Comunidades virtuales

- Grupos de personas que interactúan y colaboran usando sistemas de información y comunicación (normalmente basados en Internet).

- Blogs.
- Microblogs.
- Wikis.
- Redes sociales.
- Software libre.
- Foros...

The Twitter logo, featuring the word "twitter" in a light blue, lowercase, sans-serif font with a white outline and a soft drop shadow.

debian

The Flickr logo, with the word "flickr" in a blue, lowercase, sans-serif font, where the "i" and "k" are pink.

WIKIMEDIA
COMMONS



WIKIPEDIA
La enciclopedia libre



¿Por qué R?

- Software libre.
- Comunidad amplia y muy activa.
- Extensa documentación.
- Bibliotecas disponibles.
- Adecuado para automatizar.
- Buen rendimiento.
- Soporte para reproducibilidad.



Caso de estudio: Wikipedia

- Análisis cuantitativo macroscópico.
 - “*Wikipedia: A quantitative analysis*” (2009).
 - Comparativa de las 10 mayores Wikipedias.
 - <http://felipeortega.net/sites/default/files/thesis-jfelipe.pdf>
- Gran cantidad de datos para analizar.
 - +300 millones de cambios, sólo para la Wikipedia en inglés.
 - Casi 10 millones de artículos enciclopédicos.
- WikiXRay: Python + MySQL + R.
 - <http://meta.wikimedia.org/wiki/WikiXRay>

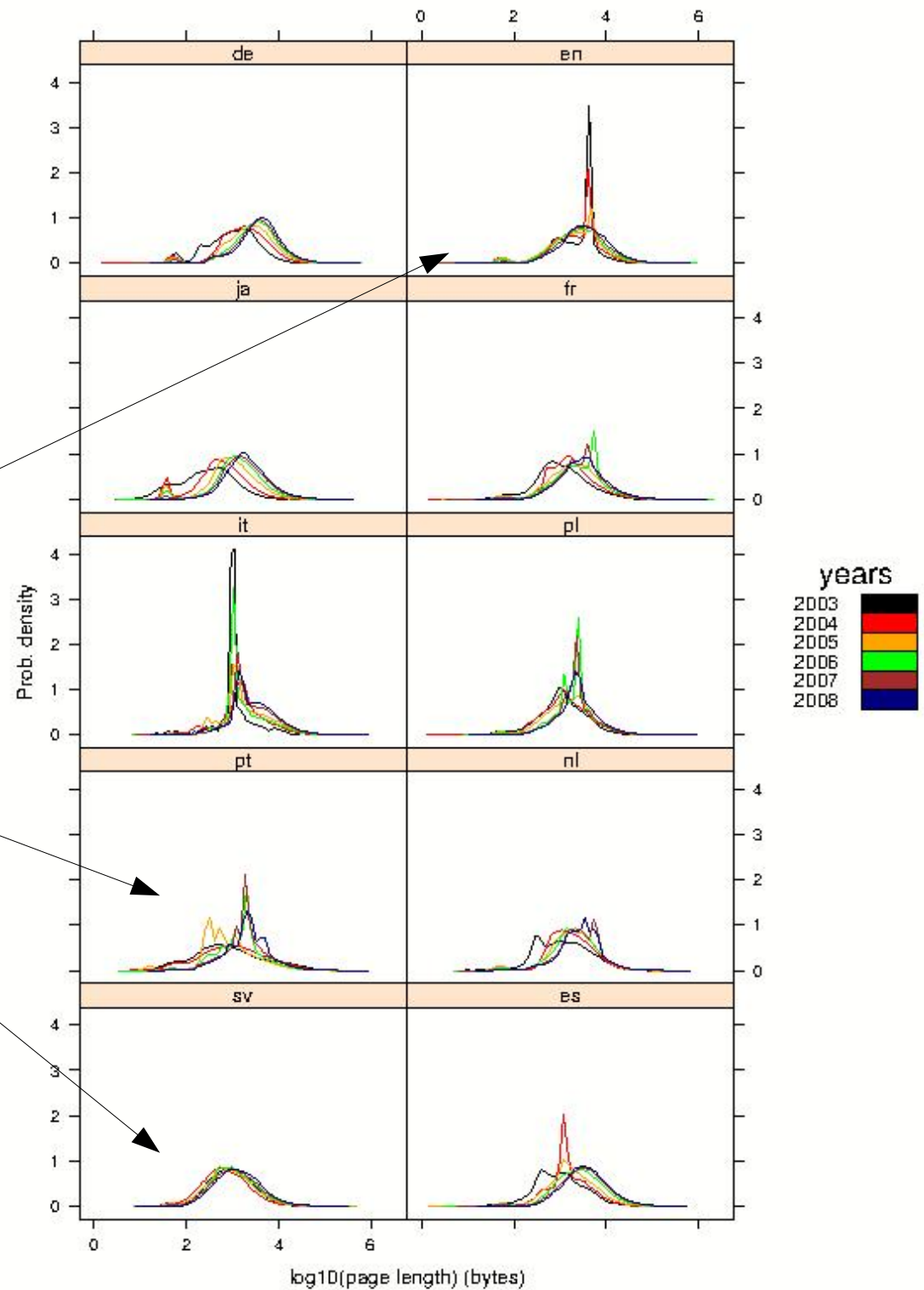
Condiciones de estudio

- Similares a muchos otros proyectos *open data*.
- Análisis del conjunto completo de datos disponibles.
 - No se realiza muestreo.
- Investigación *reproducibile*
 - Y *replicable*.
- *Automatización*.



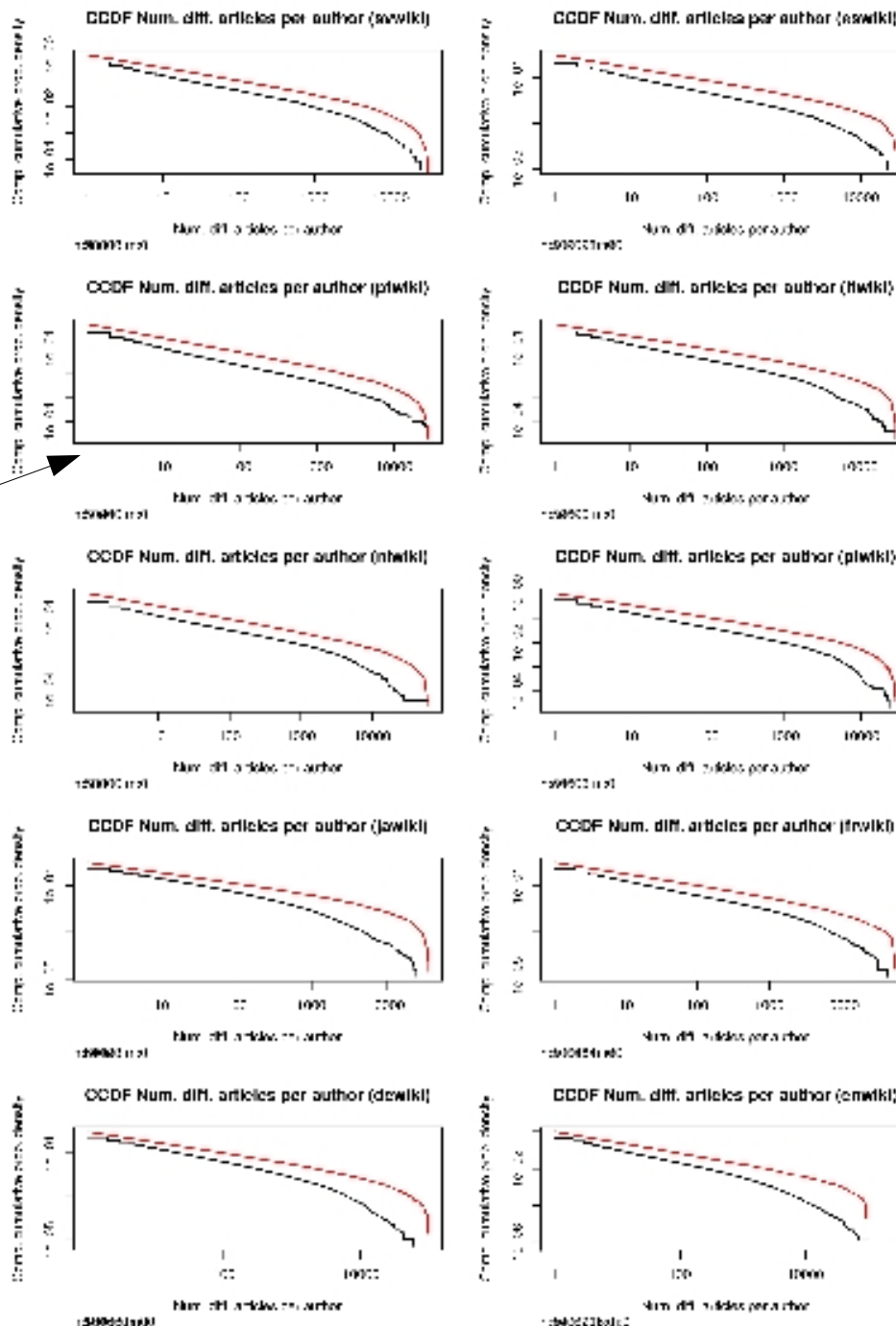
Contenidos

Diferentes
Patrones
evolutivos



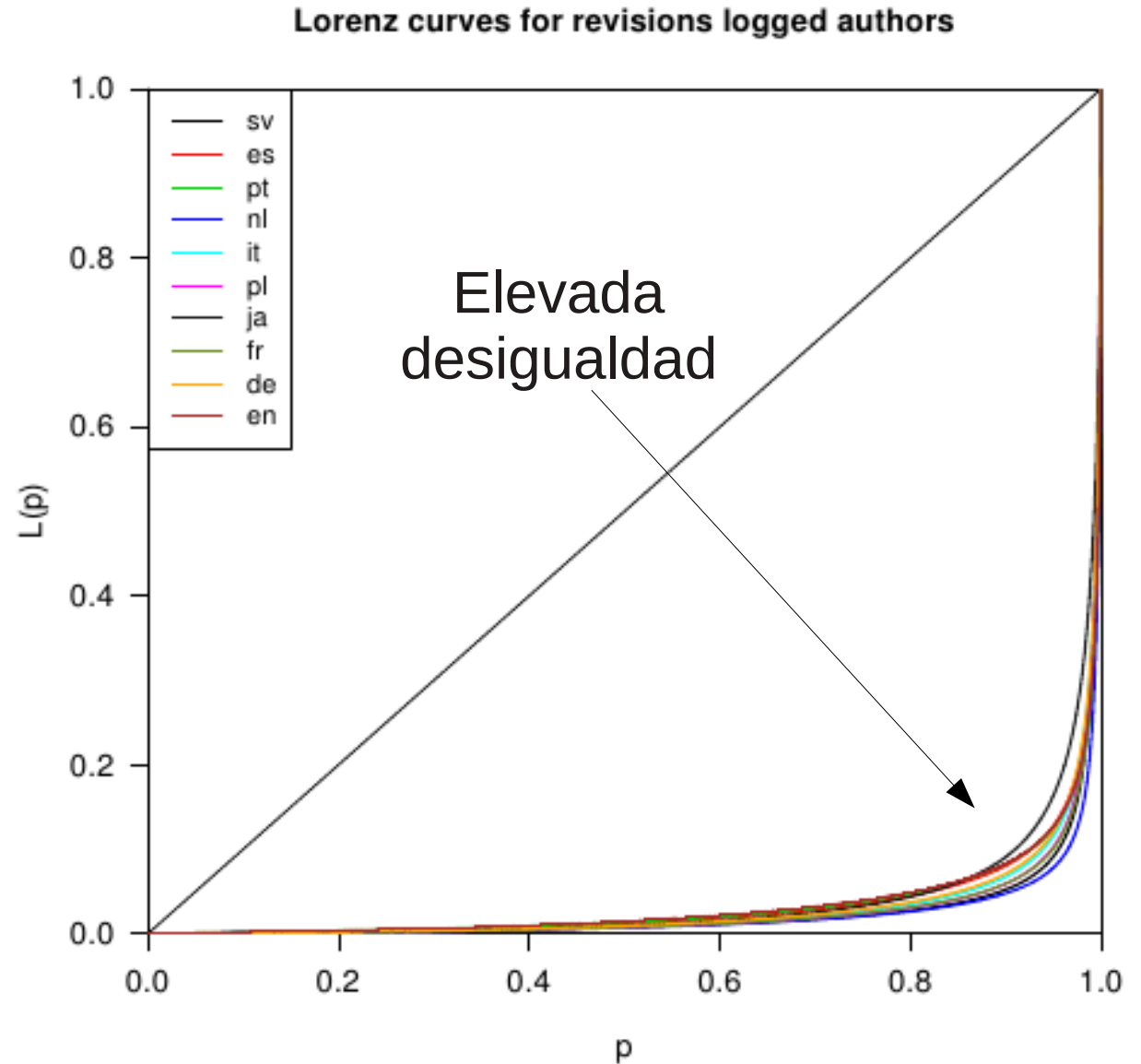
Editores

Distribución
de Pareto
truncada



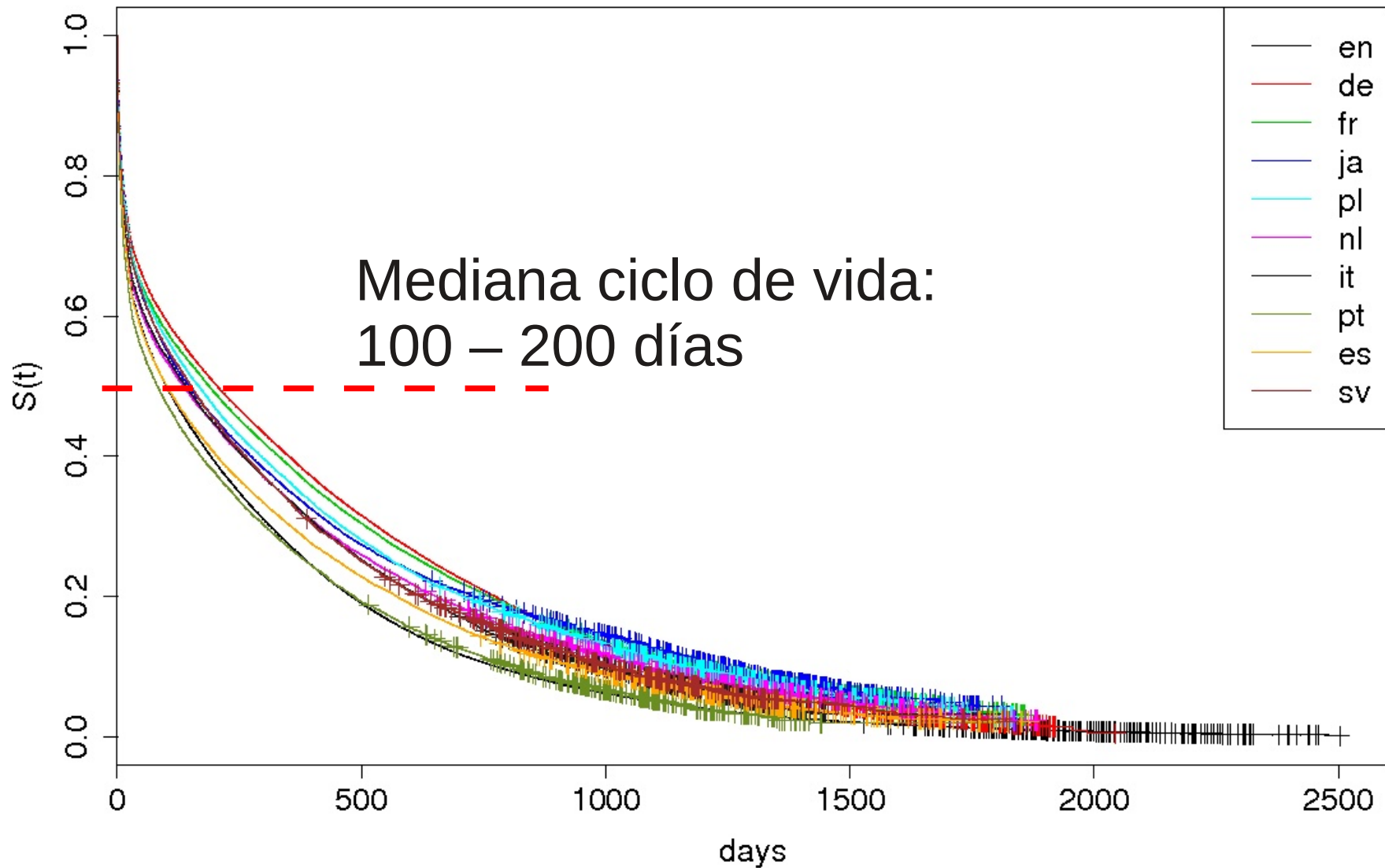
Reparto de esfuerzo

Rango de
Coeficientes
de Gini:
92% - 95%



Ciclo de vida

$S(t)$ for logged authors



Conclusiones

- **Motivos** para utilizar R.
 - Potencia, flexibilidad, rendimiento.
 - Fácil automatización del proceso de análisis.
 - Favorecer reproducibilidad de estudios.
 - Software libre: incorporación de mejoras y aportaciones.
- Próximos pasos.
 - Análisis de **redes sociales**.
 - Análisis **longitudinal**: predicción y *clustering*.

Créditos

- Portada:
 - Water drops on spider web, por [Tttrung](#) [Public domain], via [Wikimedia Commons](#).
- Página 3:
 - Twitter-logo, via [Wikimedia Commons](#).
 - Wikimedia Commons logo, (c) Wikimedia Foundationi, via [Wikimedia Commons](#).
 - Wikipedia-logo-es, por [Nohat](#) (c) Wikimedia Foundationi, via [Wikimedia Commons](#).
 - Flickr-wordmark, por Brands of the Worldi, via [Wikimedia Commons](#).
 - Debian open logo, por Debiani, via [Wikimedia Commons](#).
 - KDE logo, por KDEi, via [Wikimedia Commons](#).
 - Gnome logo, por David Vignoni, via [Wikimedia Commons](#).
- Página 4:
 - R logo, por R Foundation, via [Wikimedia Commons](#).
- Página 6:
 - R-square-logo, por Rsquaremca (Own work) [Dominio Público], via [Wikimedia Commons](#) (modificada).