



FREAK DATA EXPERIENCE

Christian López - [@christianlrcalo](https://twitter.com/christianlrcalo) - info@christianlr.es

29-10-2022

SOBRE CHRISTIAN

- Ingeniero Informático - USC
- Master of Applied Data Science - Michigan University
- Desarrollador Drupal
- Devops - SysAdmin



FASE1 - ORIGEN DEL PROYECTO

Carrera terminada en 2016 (ATPSM)

Máster terminado en 2018

FASE1 - ORIGEN DEL PROYECTO

Carrera terminada en 2016 (ATPSM)

Máster terminado en 2018

Comienzo verano 2017

FASE1 - ORIGEN DEL PROYECTO

Carrera terminada en 2016 (ATPSM)

Máster terminado en 2018

Comienzo verano 2017

"Tenemos que montar algo para recuperar datos de las redes y que lo vean..."

FASE1 - ORIGEN DEL PROYECTO

Carrera terminada en 2016 (ATPSM)

Máster terminado en 2018

Comienzo verano 2017

"Tenemos que montar algo para recuperar datos de las redes y que lo vean..."

... y tenemos sobre un mes!

FIFA TV

BUENO, VAMOH A ELLO!

PRIMEROS PASOS

- Análisis - Diseño - Desarrollo - Pruebas => NOP!

PRIMEROS PASOS

- Análisis - Diseño - Desarrollo - Pruebas => NOP!
- Java - Php - Python? => Python - Tweepy

PRIMEROS PASOS

- Análisis - Diseño - Desarrollo - Pruebas => NOP!
- Java - Php - Python? => Python - Tweepy
- Facebook - Twitter - Instagram? => Twitter

PRIMEROS PASOS

- Análisis - Diseño - Desarrollo - Pruebas => NOP!
- Java - Php - Python? => Python - Tweepy
- Facebook - Twitter - Instagram? => Twitter
- SQL - NoSQL? => PostgreSQL

PRIMEROS PASOS

- Análisis - Diseño - Desarrollo - Pruebas => NOP!
- Java - Php - Python? => Python - Tweepy
- Facebook - Twitter - Instagram? => Twitter
- SQL - NoSQL? => PostgreSQL
- Entrega? => PgAdmin con consultas pre-cargadas

PRIMERA ENTREGA

Crawler de datos (cuenta gratuita - límites)

Relación Tweet - Usuario

Búsqueda según hashtags y palabras clave

Optimización de almacenamiento (no repetidos,
exclusión)

Principio de automatización - independencia

FEEDBACK?

FEEDBACK?



FEEDBACK?



Quiero más!

FEEDBACK?



Quiero más!

Y podemos?, y podemos?, y si?

FASE2 - MÁS DIMENSIONES

Geolocalización: tweets y usuarios

Datos macro y socioeconómicos

Datos meteorológicos

Datos de género

Configuración de búsqueda dinámica

Diferentes grupos de búsqueda

FASE2 - MÁS DIMENSIONES

Geolocalización: tweets y usuarios

Datos macro y socioeconómicos

Datos meteorológicos

Datos de género

Configuración de búsqueda dinámica

Diferentes grupos de búsqueda

GIPD (Gráficas - Informes - Pdfs y Dashboard)



VENGA CHAO!

PROBLEMAS

- Ejecución estática => servicios/automática
- Servidor por grupo de búsqueda => gestión de recursos (APIs)
- Relación/fusión de datos de todas las dimensiones
- Actualización versiones Tweepy - Api Twitter
- Generador informes + estadísticas + API + Dashboard



101894730

TWEETS

29038

AVERAGE LAST WEEK

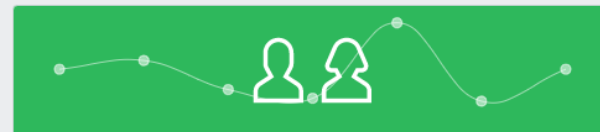


18042650

USERS

8839

AVERAGE LAST WEEK



22462455

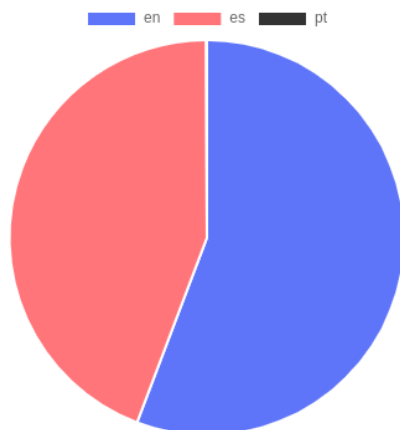
GENDERS

201

AVERAGE LAST WEEK

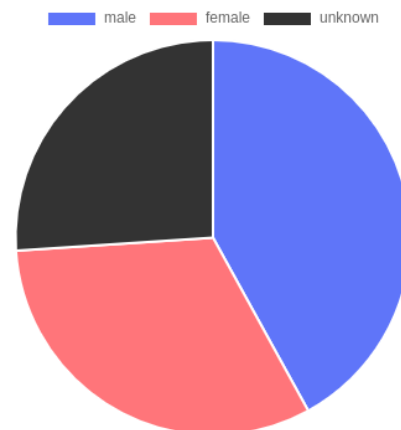
% of languages

% of top3 languages stored



% of genders

% of genders stored



Languages

0

User places

12510135

Keywords

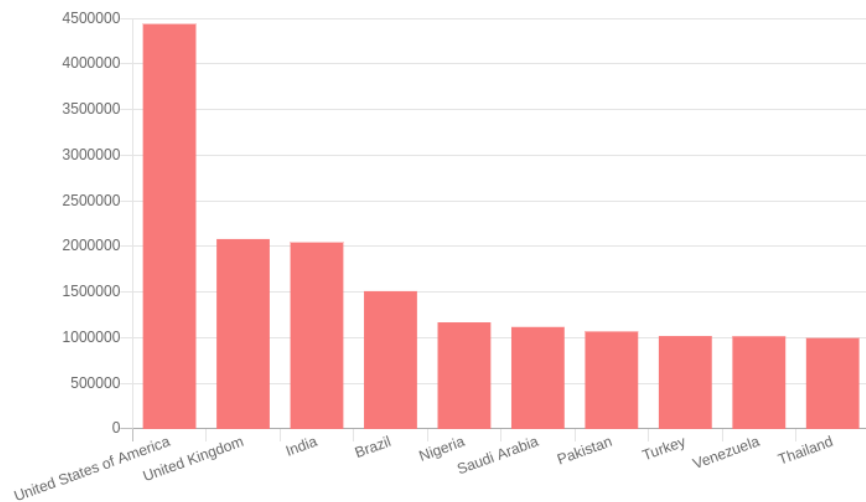
18

Hashtags

27

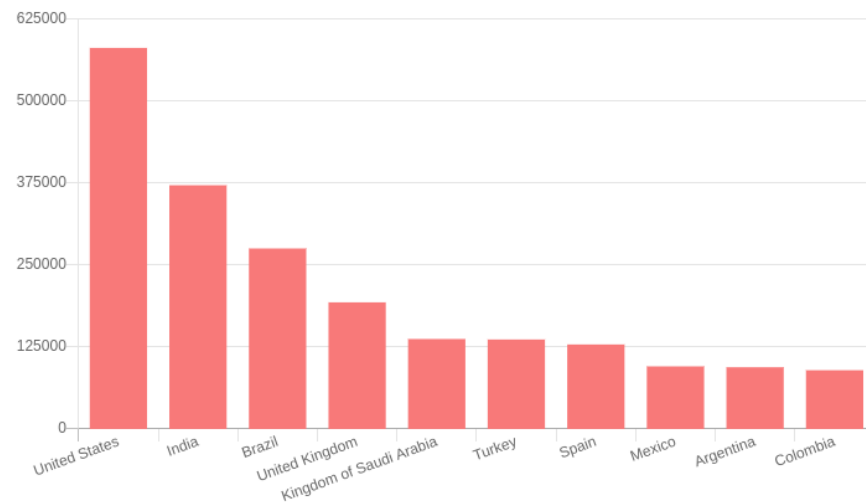
Top 10 user places

Top 10 of tweets by user location



Top 10 tweet places

Top 10 of tweets by tweet location



FASE3 - ML Y SENTIMENTANALYSIS

"Saber sentimiento/preocupación sobre un X en este dataset"

Si puede ser tanto en inglés como en español

Y además si puede ser dividido en regiones

Y claro, relacionando todas las dimesiones

FASE3 - ML Y SENTIMENTANALYSIS

"Saber sentimiento/preocupación sobre un X en este dataset"

Si puede ser tanto en inglés como en español

Y además si puede ser dividido en regiones

Y claro, relacionando todas las dimesiones

AHH! y tenemos un par de meses



NOOP!, ESTA VEZ NOP!

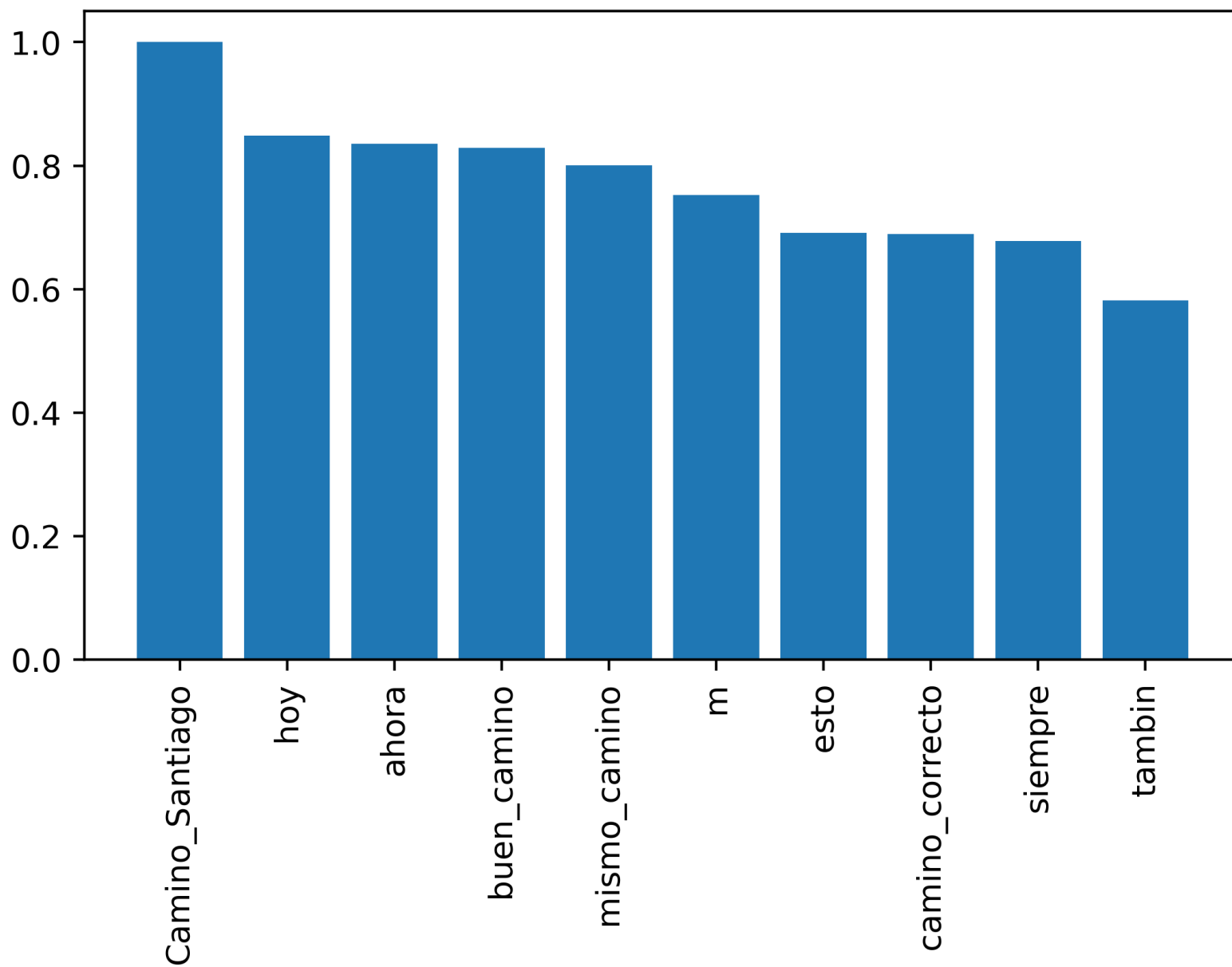
PROBLEMAS

- No hay nada formateado/ limpio
- Máquinas de crawler sin conexión entre ellas
- Necesidad de clasificación de datos => mejora de reporting
- Gran dispersidad de datos (culturas, idiomas, contextos, etc...)
- Diccionarios/corpus orientados a felicidad (+/-)

ENTREGA

- Reestructuración de las máquinas => una sola
- Reorganización de las ddbb y relaciones => una ddbb
- Restricción a principales países
- Restricción a 2 idiomas
- Reorientación de preocupación hacia felicidad
- Apoyar el proceso mediante datos estadísticos

Frequency words of all dataset



idea familia historia calle camino camino final creo
mismo paso tambien meno casa decir ahora
sigue camino gran alguien camino vida qu seguir camino
vez gente mano esto nada hoy m vida
equipo camino ms grande momento
camino parece correcto peor siempre
sigue persona aunque poco partido despu gobierno nuevo luego
cmo ah importante tema hombre nadie mal camino
mucho mundo gracia problema politica slo son lugar punto
pa camino seguir falta tena nico camino pue
aqu hecho buen llegar camino dice claro est
estn mitad camino sera otros adem mañana est camino hizo otra



BONUS - REQUISITOS EXTRA

BONUS - REQUISITOS EXTRA

Quiero un poco de todo en Twitter

BONUS - REQUISITOS EXTRA

Quiero un poco de todo en Twitter

ML con video

BONUS - REQUISITOS EXTRA

Quiero un poco de todo en Twitter

ML con video

Interactuación con mapas

BONUS - REQUISITOS EXTRA

Quiero un poco de todo en Twitter

ML con video

Interactuación con mapas

Encuestas embebidas dinámicas

¿CONCLUSIONES?

¿CONCLUSIONES?

- Hay demasiado humo (foros, documentación, ...)

¿CONCLUSIONES?

- Hay demasiado humo (foros, documentación, ...)
- La planificación es imposible

¿CONCLUSIONES?

- Hay demasiado humo (foros, documentación, ...)
- La planificación es imposible
- El rendimiento siempre queda en 2º plano

¿CONCLUSIONES?

- Hay demasiado humo (foros, documentación, ...)
- La planificación es imposible
- El rendimiento siempre queda en 2º plano
- Lo que importa es la visualización

¿CONCLUSIONES?

- Hay demasiado humo (foros, documentación, ...)
- La planificación es imposible
- El rendimiento siempre queda en 2º plano
- Lo que importa es la visualización
- SIEMPRE FALTAN DATOS

¡GRACIÑAS!

Dudas, ideas, mejoras, etc..

info@christianlr.es

[@christianlrcalo](#)