

Organización de Datos (75.06): Trabajo práctico n.1

Análisis exploratorio de datos Caso de estudio: Jampp



Alumnos: Grupo 14

Apellido, nombre	Padrón	Correo electrónico
Aristegui, Facundo	90646	faristegui@fi.uba.ar
Castellanos, Cesar	81404	castellanos.cesar@gmail.com
Cozza, Fabrizio	97402	fabrizio.cozza@gmail.com
Rozanec, Matias	97404	rozanecm@gmail.com

Introducción	3
Enunciado	3
Acerca de Jampp	3
Challenge 7506	3
Objetivo	4
Procesamiento y análisis	5
Información general sobre el análisis	5
Datos utilizados	5
Lenguaje y bibliotecas utilizadas	5
Repositorio de Github	5
Inspección de las columnas	6
Inspección de datos anonimizados	11
Análisis del flujo de datos	23
Análisis de Auctions	26
¿Cuáles son las plataformas que participan de las subastas (android / ios)? y ¿Cual es la plataforma con mayor participación en las subastas?	26
¿De qué fecha obtenemos la información? y ¿Cuántas subastas por día hay en el periodo?	26
¿Cuál es el día que se producen más subastas en 1 semana?	27
Heatmap por día y cantidad de participaciones en subastas	28
¿Cual es la cantidad de subastas por hora y por día?	29
¿A qué hora se producen la mayor y la menor cantidad de subastas?	29
¿Cuales fueron los 5 dispositivos que tuvieron mayores participaciones en las subastas?	30
¿Cuántas participaciones en subastas por dispositivo hay en el periodo?	32
¿Como es la distribución del participación por día según la plataforma?	33
¿Cómo es la participación de los sources en las subastas?	34
¿Cómo es la participación de los sources en las subastas para cada plataforma?	35
Análisis de Clicks	36
¿Cómo se distribuyen los clicks en la pantalla de los teléfonos?	36
Heatmap del set de datos	38
¿De qué fecha obtenemos la información? y ¿Cuántos clicks por día hay en el periodo?	39
Heatmap de clicks por hora y días de la semana	40
¿De qué anunciante provienen los clicks?	41
¿Cuales son las 5 fuente con más clicks?	42
Análisis de Installs	43
¿Cuantas instalaciones hubo en el periodo?	43

¿A qué hora se producen la mayor y la menor cantidad de instalaciones?	44
Heatmap de hora y días de la semana	46
¿Cuántas instalaciones hay por país?	47
¿Cuáles fueron las 10 apps más instaladas?	47
¿Cuáles fueron las 10 en las cuales se instalaron más apps?	48
Análisis de Events	49
¿Cuántos eventos hubo en el periodo?	49
¿A qué hora se producen la mayor y la menor cantidad de eventos?	50
Heatmap por día de la semana y hora	52
¿Cuáles son las apps que más eventos generaron?	53
Conclusión	54

Introducción

Enunciado

El enunciado se encuentra en el siguiente link:

<https://docs.google.com/document/d/1kZRZynbiVVR-DLtGVaKf7EvpXo5REZzM2Lga7Sdz7h0/edit>

Acerca de Jampp

Jampp es una empresa centrada en una plataforma de Mobile App Marketing. La plataforma ayuda a anunciantes a promover sus apps a nivel global y también a recuperar a aquellos usuarios que ya instalaron la app pero están inactivos. Uno de los beneficios del servicio es que optimiza la compra de tráfico en base al nivel de actividad que los usuarios tienen en la app, de esa manera, sus clientes no obtienen solo instalaciones, sino que logran resultados concretos de negocio.

En palabras de su cofundador Diego Meller, *Jampp brinda una combinación única de volumen (logrado mediante la agregación de prácticamente todo el inventario de publicidad móvil disponible) y performance. La plataforma aprende qué señales producen los usuarios que mejor convierten para cada aplicación y, en función de eso, elige donde y que avisos mostrar a cada usuario. [...] El marketing de aplicaciones ya no se trata de conseguir instalaciones (más del 80% de los usuarios dejan de usar la aplicación luego de 6 meses), sino de lograr que esas instalaciones se conviertan en usuarios activos.*

Jampp hace esto posible a través de su plataforma propietaria de compra programática en tiempo real (Real Time Bidding), que está integrada a 18 exchanges de RTB y más de 150 ad networks móviles.

Challenge 7506

Según lo expuesto por Ignacio Javier Mermet, vocero de Jampp en la materia, los problemas con que enfrenta la compañía vienen definidos por la **saturación de usuarios** y saber **cuándo se podrá apostar** por un usuario.

Un usuario se maneja en varias aplicaciones durante el día. Algunas de esas aplicaciones son clientes de Jampp y permiten conocer cómo se comporta ese usuario dentro de la aplicación. Participa en varias subastas durante el día, en las que se invita a Jampp a participar.

Objetivo

El objetivo del presente trabajo práctico es obtener un entendimiento tal de los datos brindados por la empresa, que le permita al equipo sentirse cómodo con los mismos, entendiendo no solamente la información intrínseca de cada set de datos, sino también la relación existente entre los mismos, descubriendo detalles que no sean del todo evidentes mediante un método minucioso de análisis, comprobando cada hipótesis y dejando en claro los descubrimientos consiguientes para ver si efectivamente la hipótesis se cumplió según lo esperado o, si en cambio, los datos revelan un comportamiento anti-intuitivo, como suele suceder muchas veces durante este tipo de análisis.

Procesamiento y análisis

Información general sobre el análisis

Datos utilizados

La empresa proveyó los siguientes datasets:

- Datos de subastas RTB
- Instalaciones de aplicaciones
- Clicks de publicidades mostradas
- Eventos dentro de aplicaciones

Lenguaje y bibliotecas utilizadas

Para realizar el análisis exploratorio y sus correspondientes visualizaciones, se utilizó la web app *Jupyter Notebook*¹. Algunos integrantes decidieron utilizarlo en su versión con *JupyterLab*².

Para la carga y manipulación de datos se utilizó la librería *Pandas*³ (v. 0.24.1).

Para las visualizaciones se utilizaron las librerías *matplotlib*⁴ y *seaborn*⁵.

Para el cálculo de la matriz distancias se utilizó *scipy*⁶.

Al contar con algunos datasets grandes, se consideró el uso de *Modin*⁷, una librería que con cambiar una sola línea pasaría a utilizar todos los cores de la cpu obteniendo speedups muy significantes. Sin embargo, debido a que la herramienta presentó algunos comportamientos no esperados y que finalmente no hubo mayores problemas con el tamaño de los datasets, esta herramienta quedó descartada.

Repositorio de Github

Todo el desarrollo se encuentra subido a un repositorio de Github, donde se encuentran todos los análisis realizados necesarios para llevar a cabo la redacción del presente informe.

A continuación, el link a dicho repositorio: <https://github.com/faristegui/tpdatos2019>

¹ <https://jupyter.org/>

² <https://jupyterlab.readthedocs.io/en/stable/>

³ <https://pandas.pydata.org/>

⁴ <https://matplotlib.org/>

⁵ <https://seaborn.pydata.org/>

⁶ <https://www.scipy.org/>

⁷ <https://modin.readthedocs.io/en/latest/>

Inspección de las columnas

A continuación se muestra un resumen del análisis previo de las columnas de los datos para lograr una eficiente lectura de los datos, motivado principalmente por el tamaño de los archivos Auction e Installs, de 1.4 Gb y 750 Mb respectivamente.

Auctions

Columna	Tipo por defecto	Análisis
<i>auction_type_id</i>	<i>float64</i>	Se puede omitir de la lectura, sus valores son nulos.
<i>country</i>	<i>int64</i>	Se puede omitir de la lectura ya que hay un único valor.
<i>date</i>	<i>object</i>	Se tendrá que convertir a datetime durante el análisis.
<i>device_id</i>	<i>int64</i>	Contiene pocos valores pero grandes, se opta por leerlo como category ya que ocupa menos espacio en memoria.
<i>platform</i>	<i>int64</i>	Contiene valores pequeños, se opta por leerlo como un int8.
<i>ref_type_id</i>	<i>int64</i>	Contiene valores pequeños, se opta por leerlo como un int8.
<i>source_id</i>	<i>int64</i>	Contiene valores pequeños, se opta por leerlo como un int8.

Events

Columna	Tipo por defecto	Análisis
<i>date</i>	<i>object</i>	Se tendrá que convertir a datetime durante el análisis.
<i>event_id</i>	<i>int64</i>	Contiene valores pequeños, se opta por leerlo como un int8.
<i>ref_type</i>	<i>int64</i>	Contiene valores pequeños, se opta por leerlo como un int8.
<i>ref_hash</i>	<i>int64</i>	Dejarlo como está.
<i>application_id</i>	<i>int64</i>	Contiene valores pequeños, se opta por leerlo como un int8.
<i>attributed</i>	<i>bool</i>	Dejarla como está.
<i>device_countrycode</i>	<i>int64</i>	Se puede omitir de la lectura ya que hay un único valor.

<i>device_os_version</i>	<i>float64</i>	Contiene pocos valores, se opta por leerlo como una category.
<i>device_brand</i>	<i>float64</i>	Se puede omitir de la lectura, sus valores son nulos.
<i>device_model</i>	<i>float64</i>	Contiene pocos valores, se opta por leerlo como una category.
<i>device_city</i>	<i>float64</i>	Contiene pocos valores, se opta por leerlo como una category.
<i>session_user_agent</i>	<i>float64</i>	Se puede omitir de la lectura ya que hay un único valor.
<i>trans_id</i>	<i>float64</i>	Se deja con propósitos de mergeo.
<i>user_agent</i>	<i>float64</i>	Se puede omitir de la lectura, sus valores son nulos.
<i>event_uuid</i>	<i>object</i>	Se puede omitir de la lectura, son todos valores únicos y en gran cantidad.
<i>carrier</i>	<i>float64</i>	Contiene pocos valores, se opta por leerlo como una category.
<i>kind</i>	<i>float64</i>	Contiene pocos valores, se opta por leerlo como una category.
<i>device_os</i>	<i>float64</i>	Se puede omitir de la lectura ya que hay un único valor.
<i>wifi</i>	<i>float64</i>	Se puede omitir de la lectura ya que hay un único valor.
<i>connection_type</i>	<i>object</i>	Contiene pocos valores, se opta por leerlo como una category.
<i>ip_address</i>	<i>int64</i>	Dejarlo como está.
<i>device_language</i>	<i>float64</i>	Contiene pocos valores, se opta por leerlo como una category.

Installs

Columna	Tipo por defecto	Análisis
<i>created</i>	<i>object</i>	Se tendrá que convertir a datetime durante el análisis.
<i>application_id</i>	<i>int64</i>	Contiene valores pequeños, se opta por leerlo como un int8.
<i>ref_type</i>	<i>int64</i>	Dejarla como está.

<i>ref_hash</i>	<i>int64</i>	Dejarla como está.
<i>click_hash</i>	<i>float64</i>	Se puede omitir de la lectura, sus valores son nulos.
<i>attributed</i>	<i>bool</i>	Se puede omitir de la lectura ya que hay un único valor.
<i>implicit</i>	<i>bool</i>	
<i>device_countrycode</i>	<i>int64</i>	Contiene pocos valores, se opta por leerlo como una category.
<i>device_brand</i>	<i>float64</i>	Contiene pocos valores, se opta por leerlo como una category.
<i>device_model</i>	<i>float64</i>	Dejarla como está.
<i>session_user_agent</i>	<i>object</i>	Contiene pocos valores, se opta por leerlo como una category.
<i>user_agent</i>	<i>object</i>	Contiene valores tipo string.
<i>event_uuid</i>	<i>object</i>	Si no existe la posibilidad de recuperarle datos, se puede omitir.
<i>kind</i>	<i>object</i>	Si no existe la posibilidad de recuperarle datos, se puede omitir.
<i>wifi</i>	<i>object</i>	Contiene valores tipo bool.
<i>trans_id</i>	<i>object</i>	Se deja con propósitos de mergeo.
<i>ip_address</i>	<i>int64</i>	Dejarla como está.
<i>device_language</i>	<i>float64</i>	Contiene pocos valores, se opta por leerlo como una category.

Clicks

Columna	Tipo por defecto	Análisis
<i>advertiser_id</i>	<i>int64</i>	Contiene pocos valores, se opta por leerlo como una category.
<i>action_id</i>	<i>float64</i>	Se puede omitir de la lectura, sus valores son nulos.
<i>source_id</i>	<i>int64</i>	Contiene pocos valores, se opta por leerlo como una category.
<i>created</i>	<i>object</i>	Se tendrá que convertir a datetime durante el análisis.

<i>country_code</i>	<i>int64</i>	Se puede omitir de la lectura ya que hay un único valor.
<i>latitude</i>	<i>float64</i>	Dejarla como está.
<i>longitude</i>	<i>float64</i>	Dejarla como está.
<i>wifi_connection</i>	<i>bool</i>	Se puede omitir de la lectura ya que hay un único valor.
<i>carrier_id</i>	<i>float64</i>	Contiene pocos valores, se opta por leerlo como una category.
<i>trans_id</i>	<i>object</i>	Se deja con propósitos de mergeo.
<i>os_minor</i>	<i>float64</i>	Contiene pocos valores, se opta por leerlo como una category.
<i>agent_device</i>	<i>float64</i>	Dejarla como está.
<i>os_major</i>	<i>float64</i>	Contiene pocos valores, se opta por leerlo como una category.
<i>specs_brand</i>	<i>int64</i>	Contiene pocos valores, se opta por leerlo como una category.
<i>brand</i>	<i>float64</i>	Contiene pocos valores, se opta por leerlo como una category.
<i>timeToClick</i>	<i>float64</i>	No hace falta tanta precisión, se opta por leerlo como un float16.
<i>touchX</i>	<i>float64</i>	No hace falta tanta precisión, se opta por leerlo como un float16.
<i>touchY</i>	<i>float64</i>	No hace falta tanta precisión, se opta por leerlo como un float16.
<i>ref_type</i>	<i>int64</i>	Contiene pocos valores, se opta por leerlo como una category.
<i>ref_hash</i>	<i>int64</i>	Contiene valores pequeños, se opta por leerlo como un int8.

Algunas observaciones a priori

Llamó la atención que no haya ningún valor en el campo *auction_id* de Auctions. No tanto por el valor que agregaría al análisis dicho campo, sino por la presencia del mismo. Se maneja la hipótesis de que en principio se iba a incluir, y que en la versión entregada al curso finalmente se omitió.

En cuanto al campo correspondiente al país, llamó la atención que 3 de los 4 datasets contengan solamente un valor, repetido en más de la mitad de los registros en el dataset restante. No parecería ser un campo que aporte mucha información a priori.

Similarmemente a lo ocurrido en Auctions con *auction_id*, en el dataset Events se detectan todos valores nulos en el campo de *device_brand*. Ídem en *click_hash* de Installs y *action_id* de Clicks.

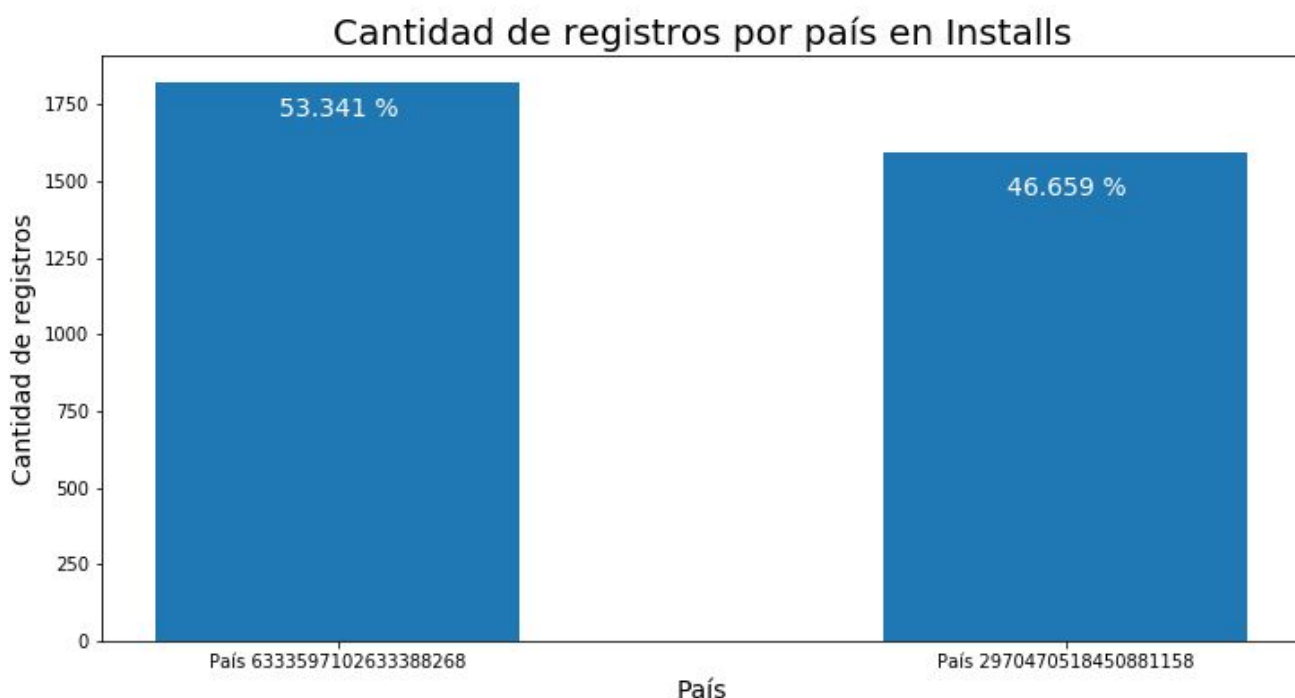
Lo que más poderosamente llamó la atención es que no existan registros con instalaciones atribuidas a Jampp.

Inspección de datos anonimizados

Lo primero que se observó en los datasets es que, tal como ya había sido anunciado, los datos fueron anonimizados. Si bien en principio esto no pareció ser de importancia, cambió un poco la forma de tratar los datos. Se listan a continuación los datos que fueron transformados durante dicho proceso, comentando en algunos casos qué impacto tuvieron durante el análisis.

Auctions

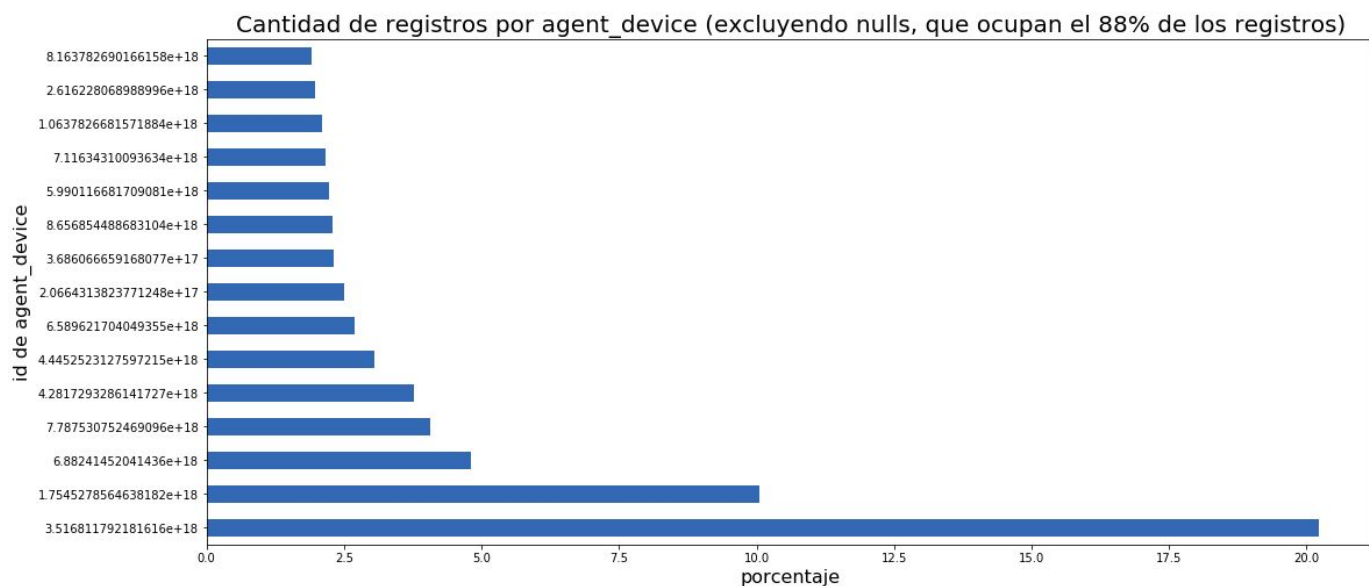
- country: si bien este dato pasó por *string hashing*, en el dataset toma un solo valor. Lo mismo pasa con los campos *country_code* y *device_countrycode* de Clicks y Events respectivamente. En los tres casos el valor existente es el mismo. No es así en el dataset Installs, donde en el campo *device_countrycode* aparecen dos valores con una distribución muy pareja, como se muestra en el gráfico a continuación. El valor de country de los otros tres datasets se corresponde con el valor mayoritario.



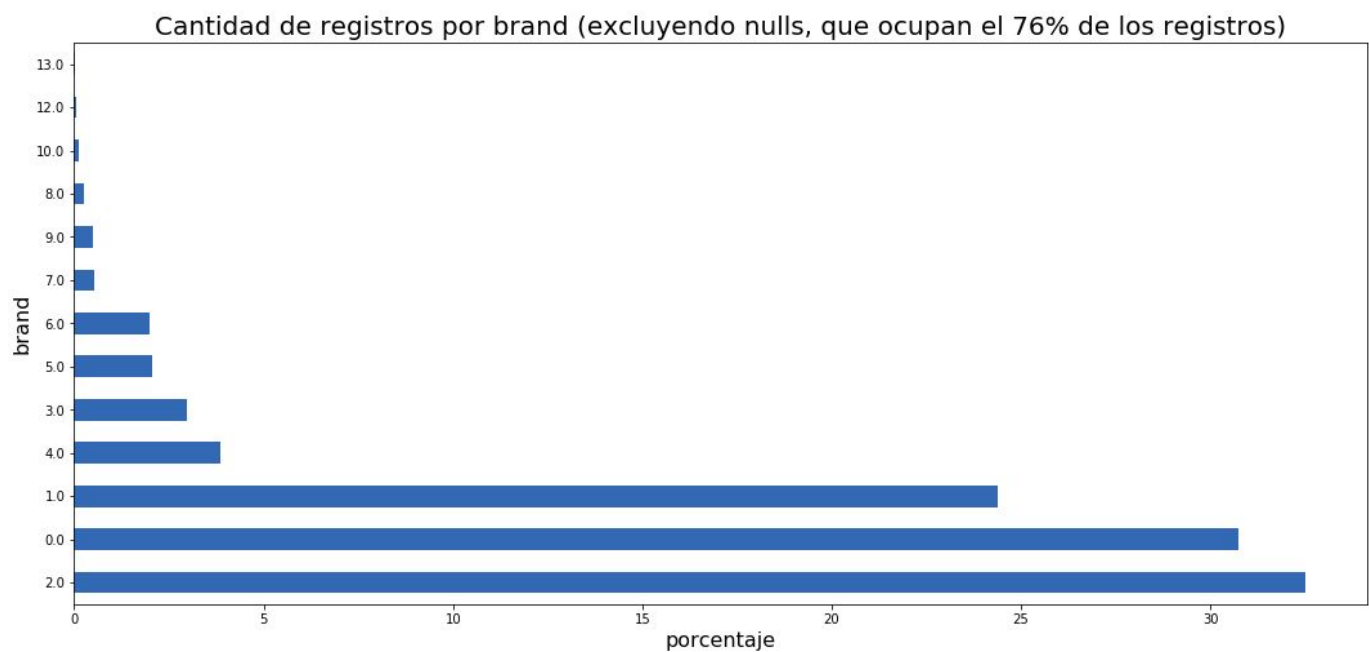
- device_id: como los ids no tienen información intrínseca, el hecho de que hayan sido hashados no modificó sustancialmente el análisis.
- source_id:

Clicks

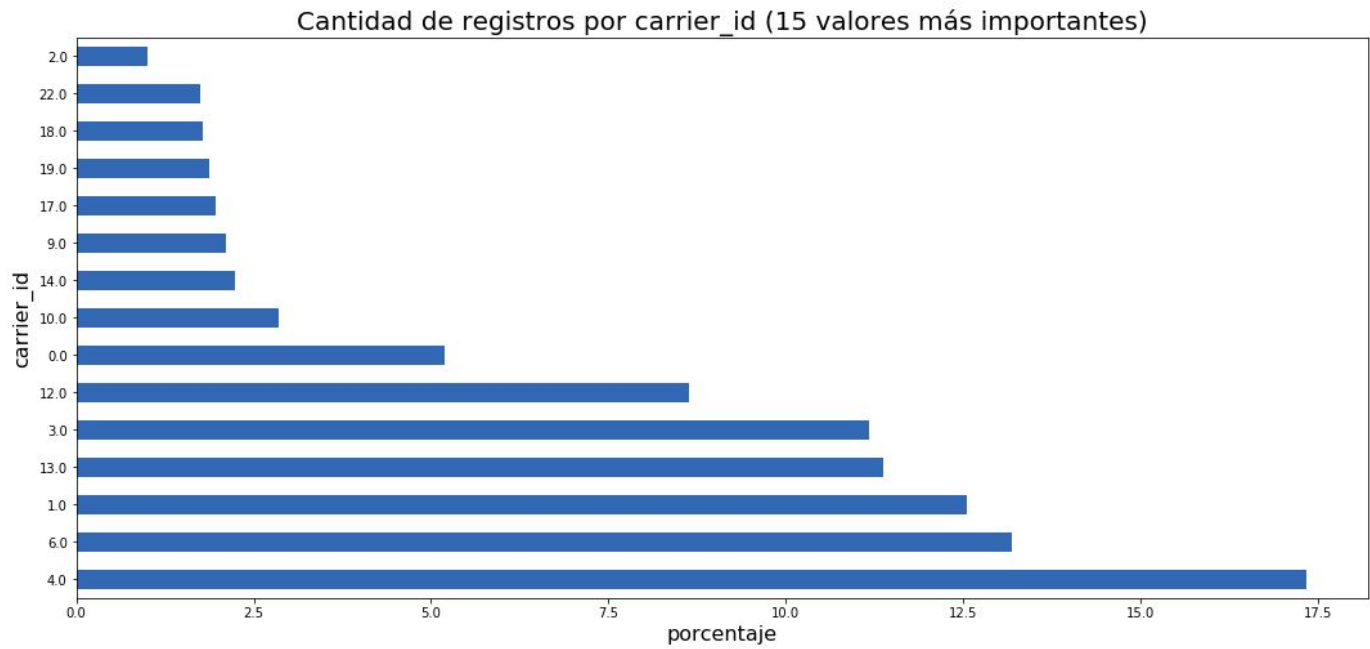
- advertiser_id: ídem caso anterior de id. En este caso se vio que no hay ningún campo nulo, y que se cuenta con un total de 7 ids únicos.
- agent_device: el encoding propiamente dicho no molesta en el análisis. Lo que sí se puede notar es que el 88% de los registros contiene valores nulos, y el resto sigue una ley de potencias, según se puede observar en el siguiente gráfico que muestra los 15 valores más frecuentes (de los 190 valores únicos que hay).



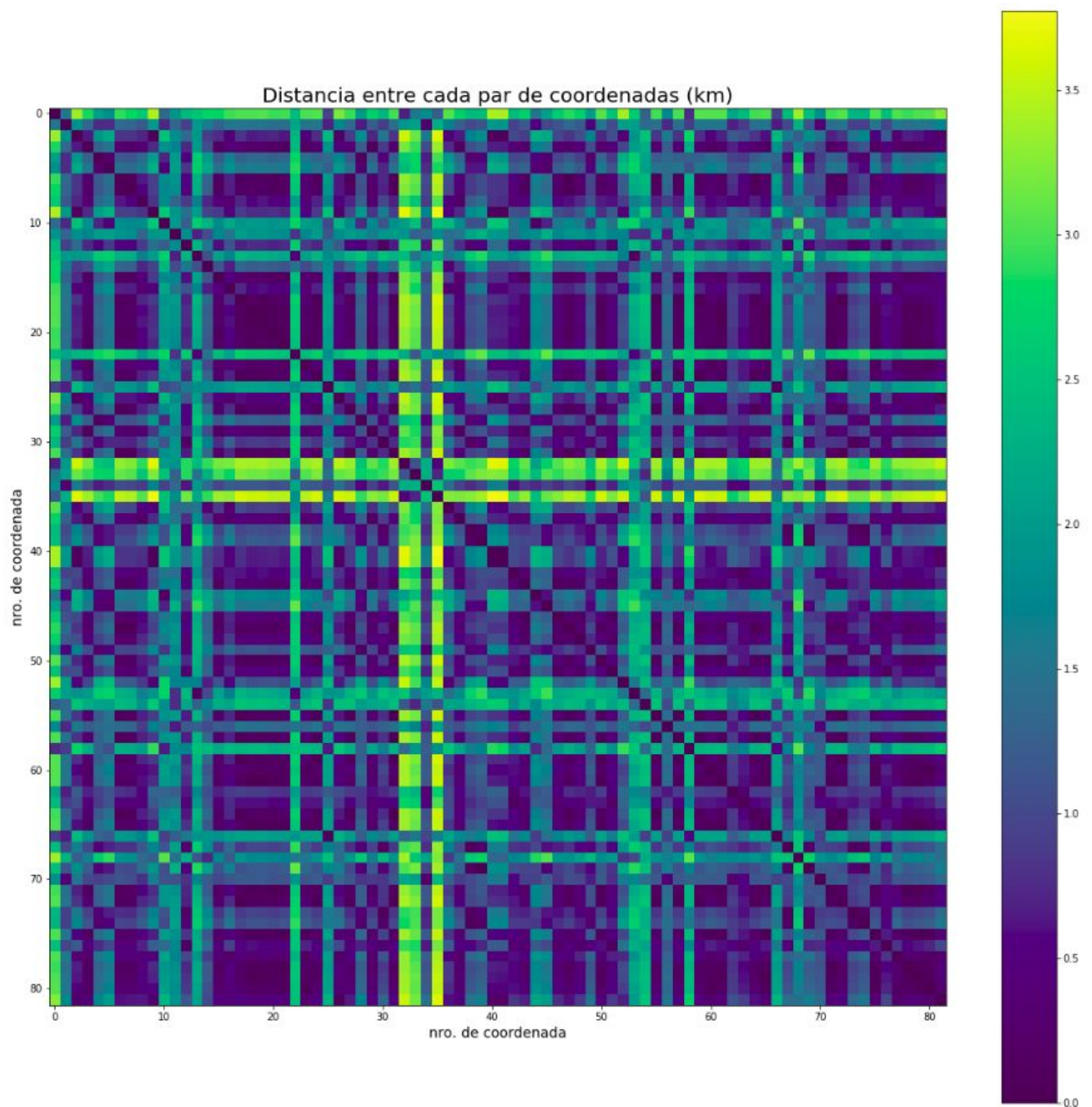
- brand: En un principio, desconocer explícitamente la marca de un dispositivo no debería traer problemas, aunque sí serviría percibir las marcas reales para poder poner un foco mayor en ellas, ya que aquellos usuarios que poseen esta marca serán los que más interactúen (por ejemplo la 1.0, 0.0 y 2.0 en el gráfico). Nuevamente se observa una gran cantidad de nulos, y luego una distribución similar al caso anterior.



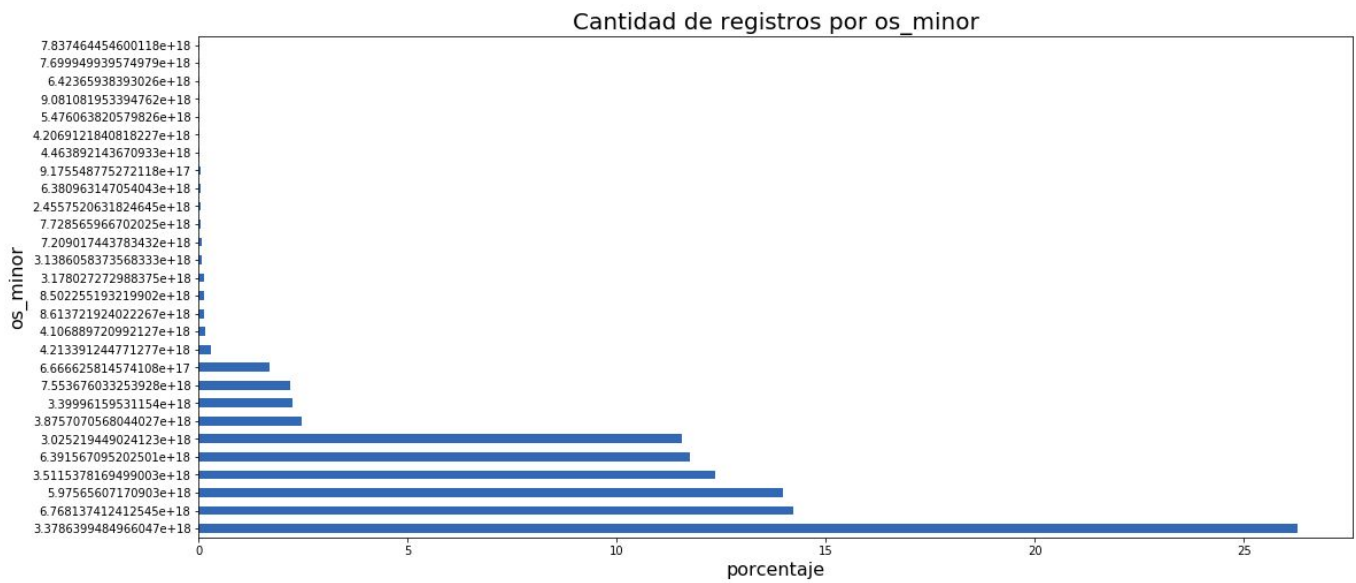
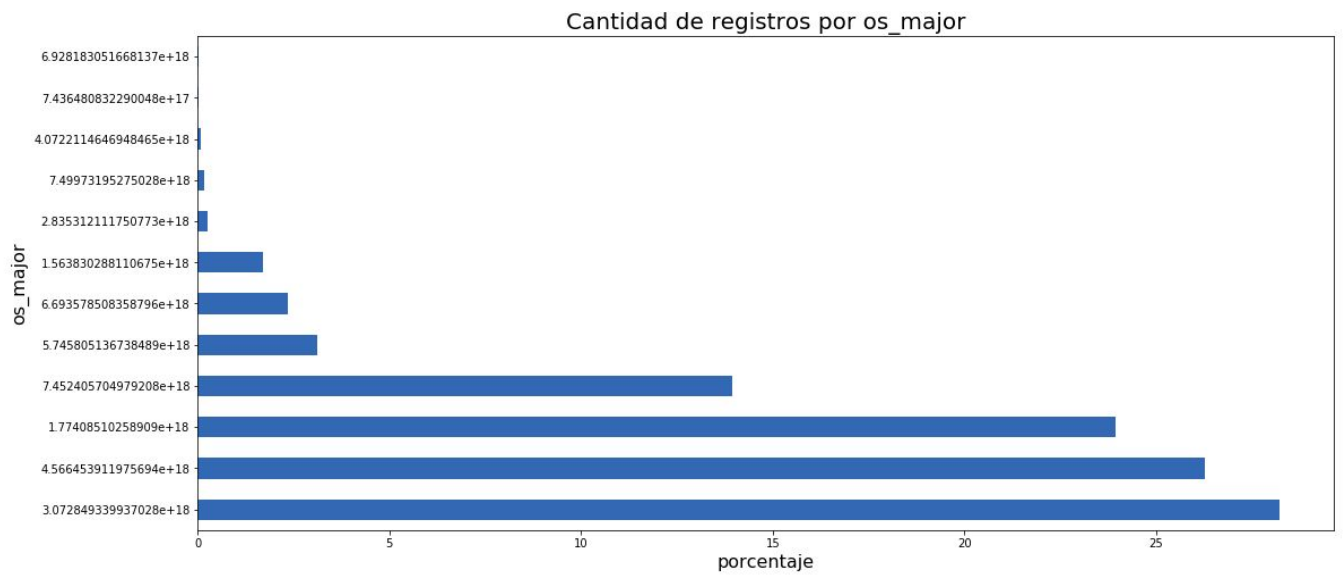
- carrier_id: en cuanto a la anonimización, sucede como en el caso anterior, que no trae mayores problemas. La diferencia reside en la cantidad de campos nulos, que en este caso está por debajo del 0.1%, por lo que se termina obteniendo un panorama mucho más real.



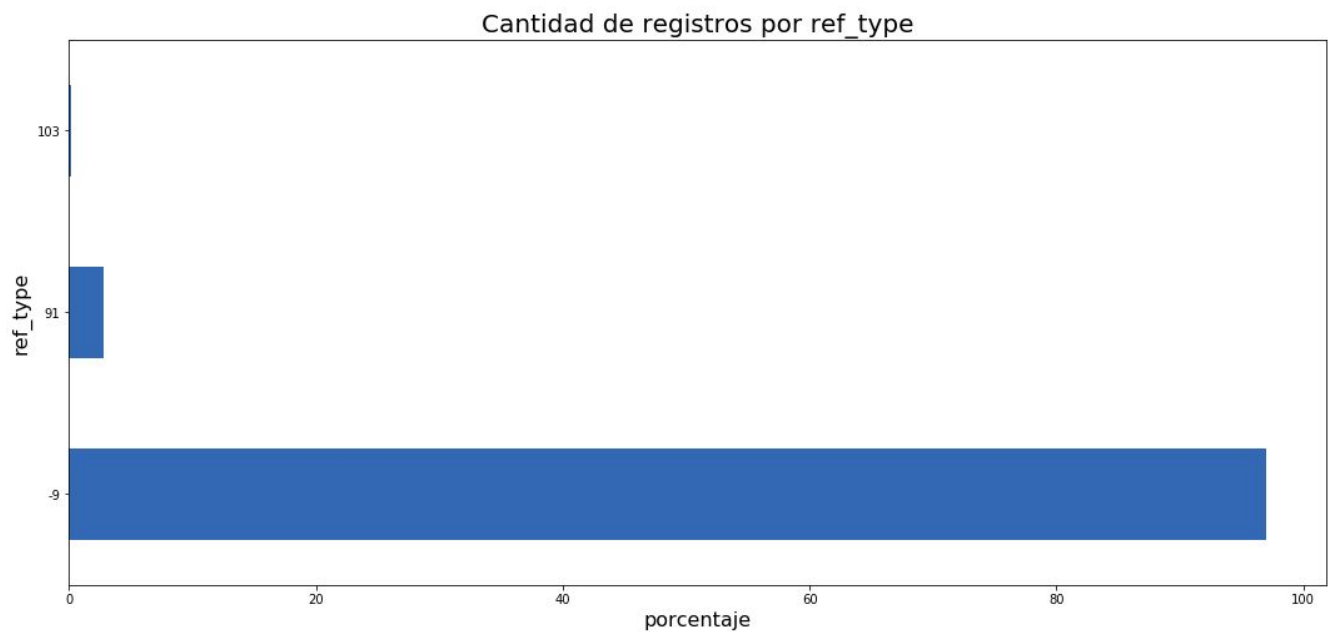
- latitud - longitud: la idea principal fue mostrar estos datos en un mapa, pero como se realizó una transformación lineal sobre estos datos, esto ya no será posible. Sin embargo, al ser la transformación lineal, se pudo llevar a cabo un análisis de distancias. A continuación se muestra un heatmap mostrando las distancias entre cada par de las 82 coordenadas únicas presentes en el dataset. Notar que en general las distancias son muy chicas, dándose las máximas entre un par muy reducido de coordenadas y con valores que rodean los 3.5 km.



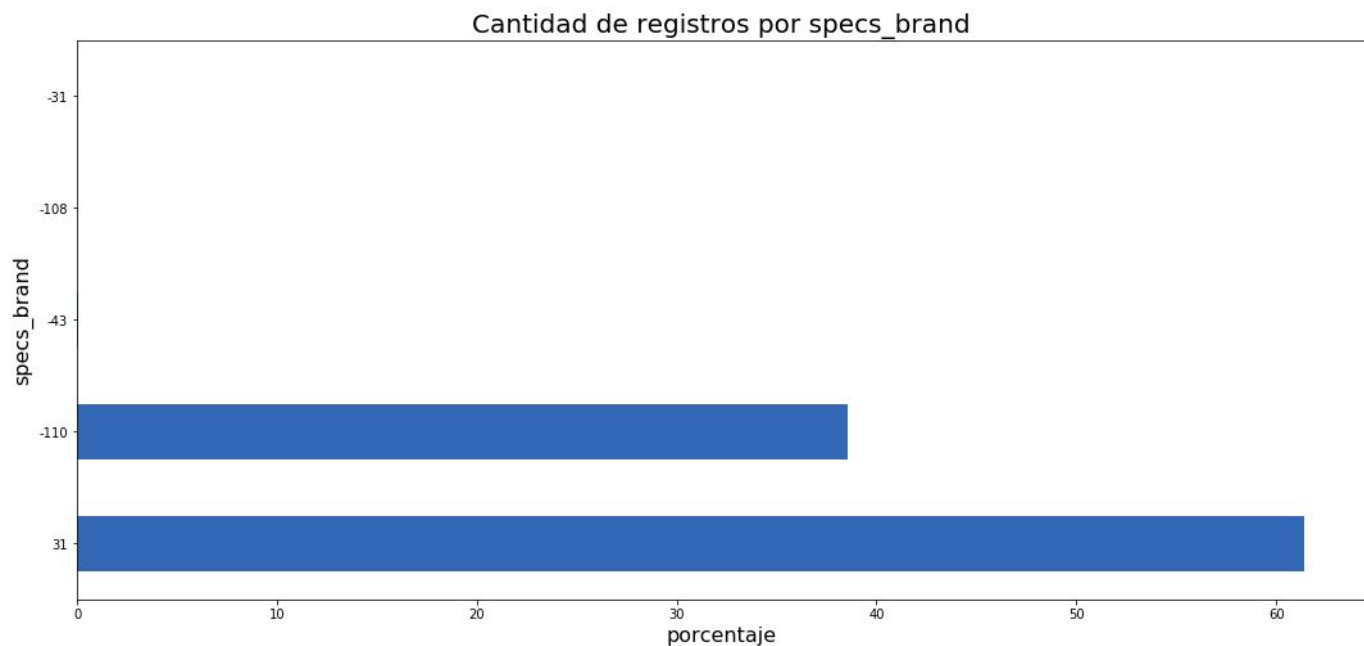
- `os_mayor` y `os_minor`: prácticamente no hay valores nulos. En `os_mayor`, si bien hay 12 valores únicos, 3 de ellos concentran el 75% de los registros; si se toman los 4 valores con mayor cantidad de registros ya se cubre el 90% de los registros.
En `os_minor` hay 28 valores posibles, entre los cuales hay un valor que abarca el 26% de los registros no nulos; le siguen 5 valores que abarcan una cantidad aproximadamente pareja de registros; entre estos 6 valores ya se abarca el 90% del dataset.



- ref_hash: este registro no importa que esté hashado, ya que el campo en sí no tiene valor intrínseco. El propósito de este campo es poder cruzar datos entre los distintos datasets.
- ref_type: casi todos los registros contienen el mismo valor (no hay ninguno con valor nulo).

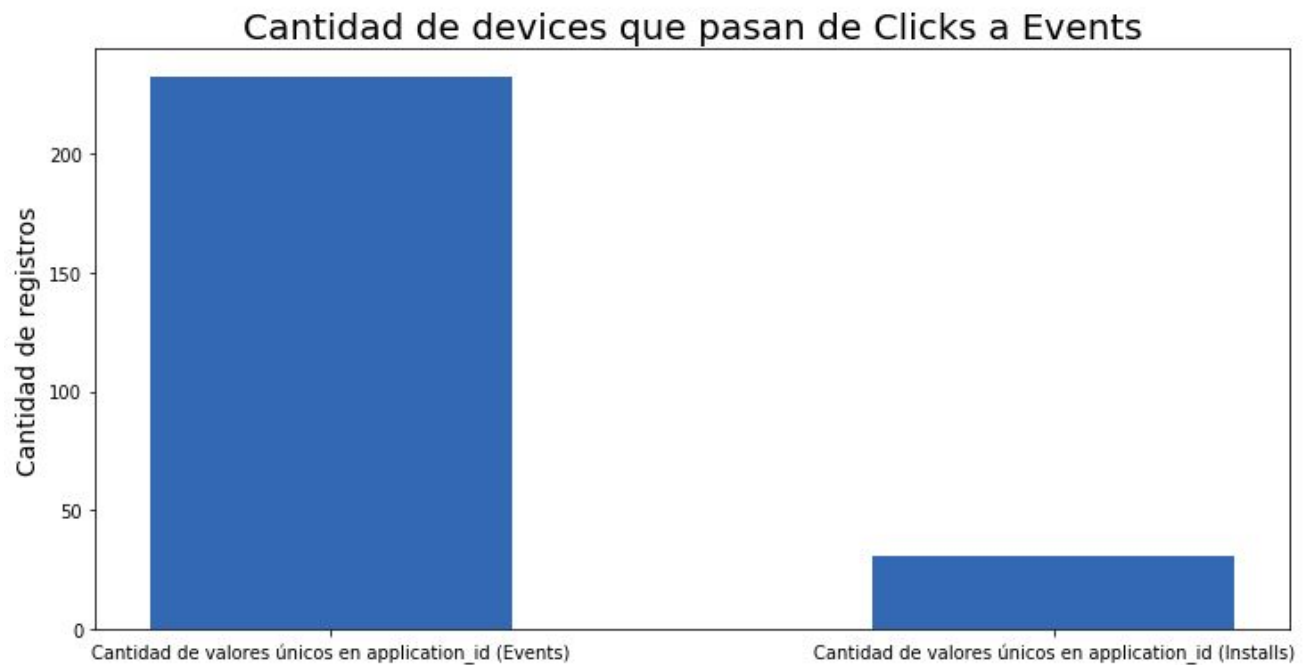


- specs_brand: no se llega a deducir qué es, y tampoco se proveyó una descripción del campo. Este campo toma solamente 5 valores posibles, pero casi la totalidad de los registros está dividida entre dos valores posibles.



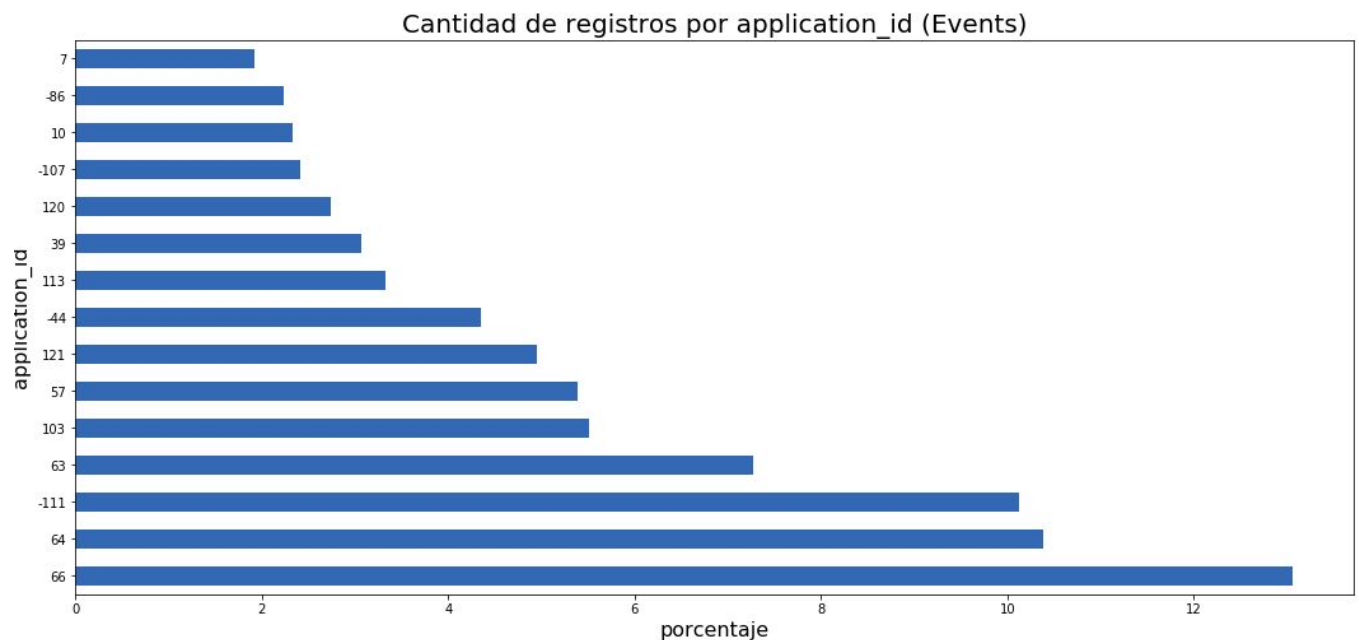
Events

- application id: registro compartido con Installs. Al no tener valor intrínseco, no cambia en nada la anonimización.
En ambos datasets la cantidad de valores únicos presentes varía muchísimo.

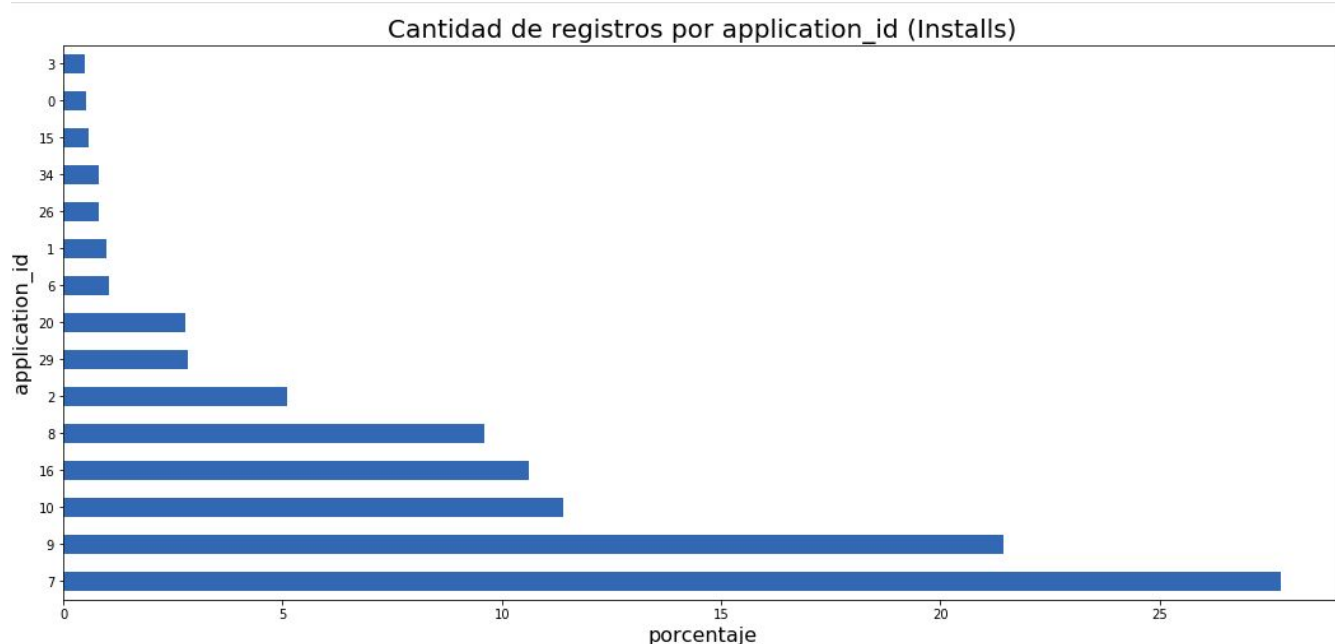


Por esta razón se dificulta bastante hacer una comparación a priori de ambos campos.

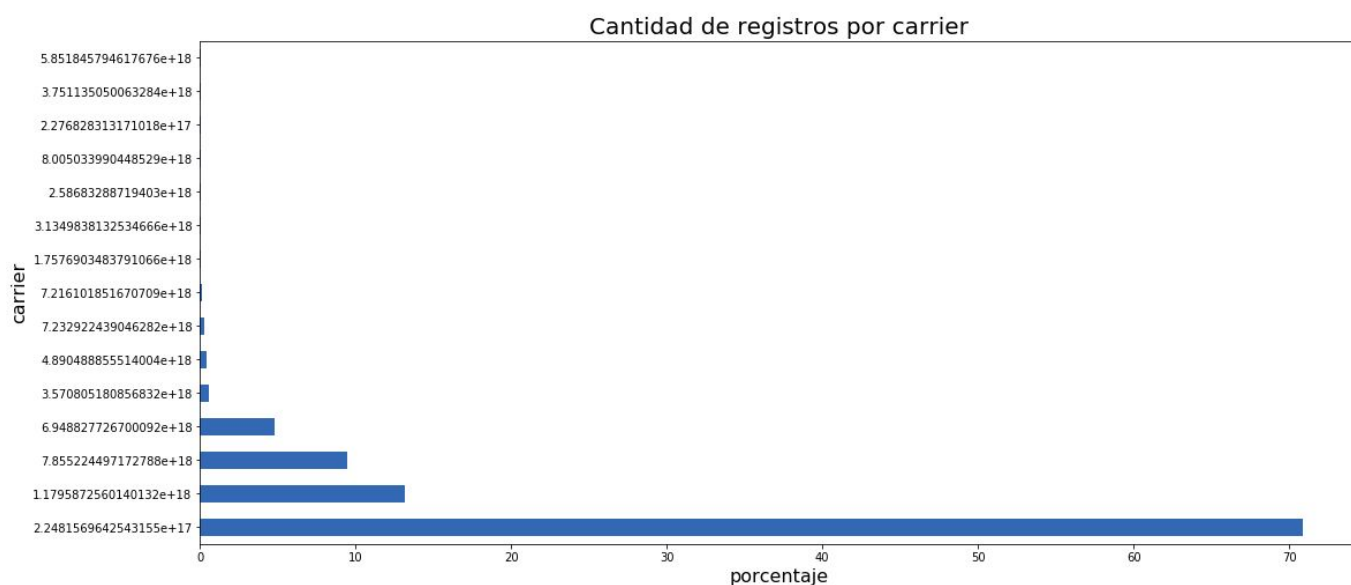
En el caso de Events hay 3 ids que toman por encima del 10% de los registros de Events cada uno.



En Installs termina sucediendo algo similar: hay 4 valores que abarcan más del 10% de los valores no nulos cada uno.

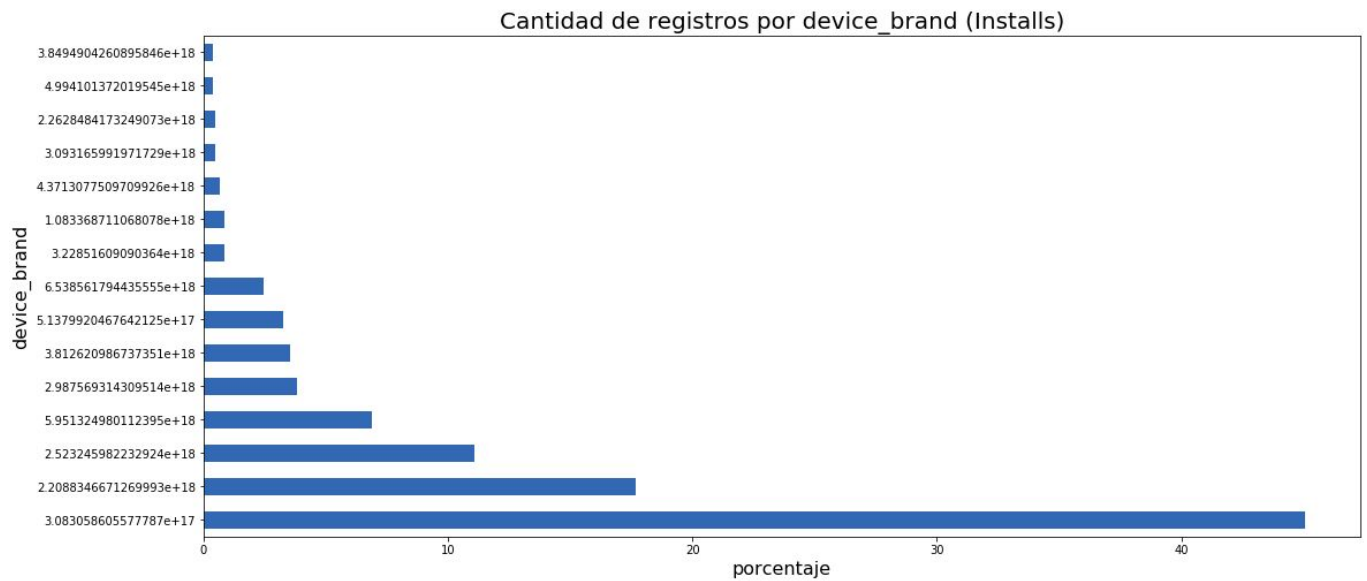
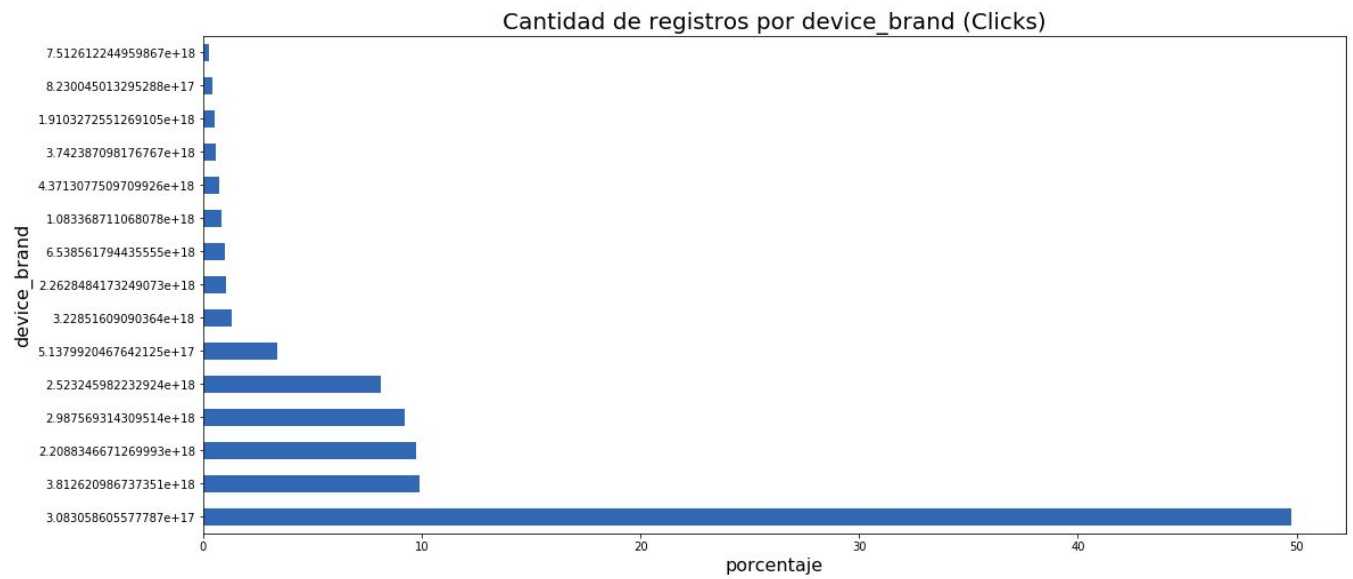


- carrier: se llegó a pensar que el hash sería compartido con el carrier_id de Clicks, sin embargo los valores no tienen nada que ver. En este caso, del 25% de registros que toman valores no nulos en este campo, el 75% toma un mismo valor (de 85 posibles).

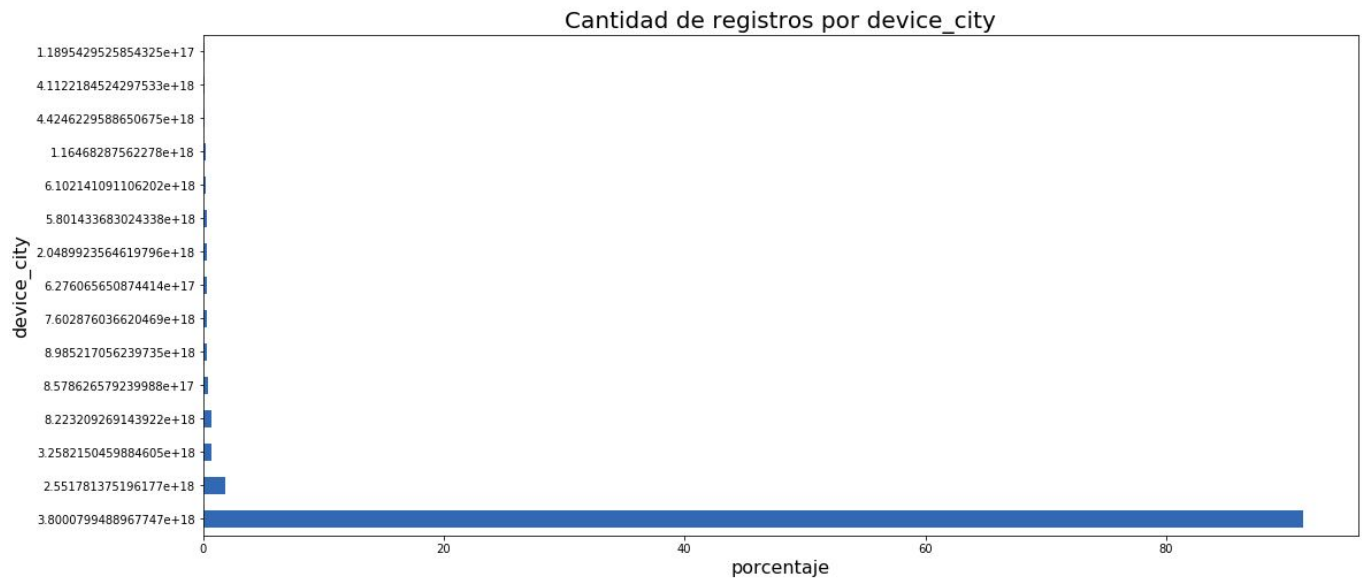


- device_brand: nuevamente se pensaría que habría valores en común con el dataset Clicks. Sin embargo tampoco tienen nada que ver, a pesar de que en ambos casos se trata de la marca del dispositivo. La distribución difiere del caso anterior también: en este caso, del 50% de campos no nulos, el 50% tiene un mismo valor, mientras que el resto de los valores ocupan cada uno menos del 10% del 50% restante.

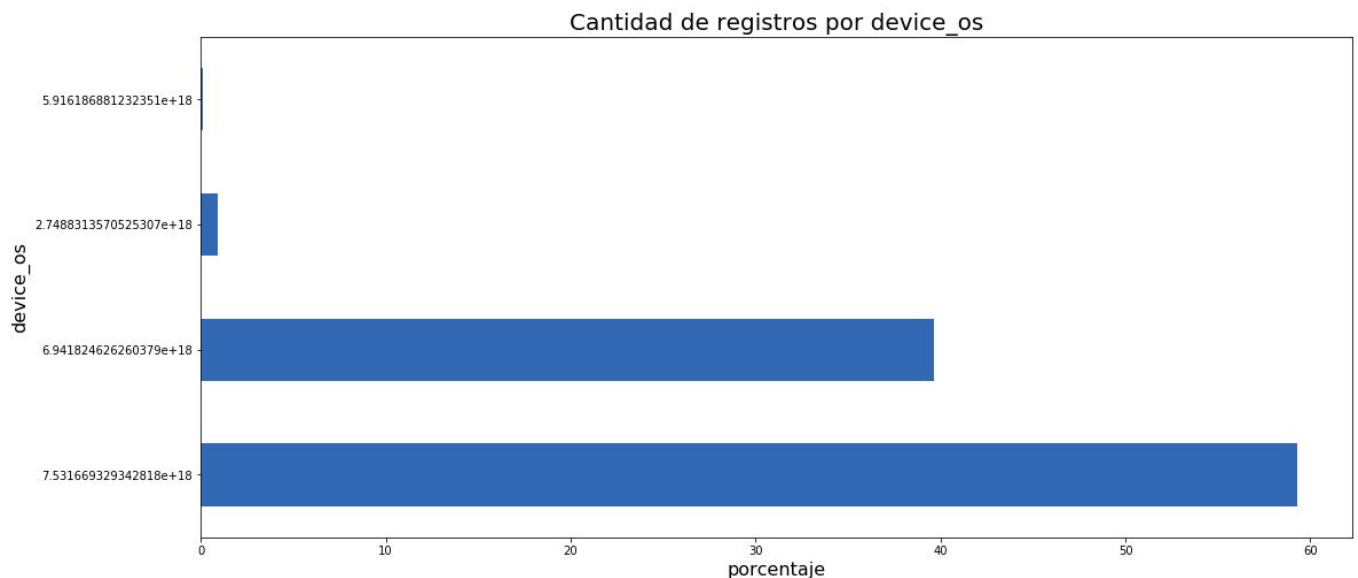
Hay que notar que sin embargo, este campo es compartido con Installs. En ambos casos se obtiene una distribución similar, aunque en un orden levemente cambiado, a excepción del caso mayoritario que en ambos casos ronda el 50% de los registros no nulos.



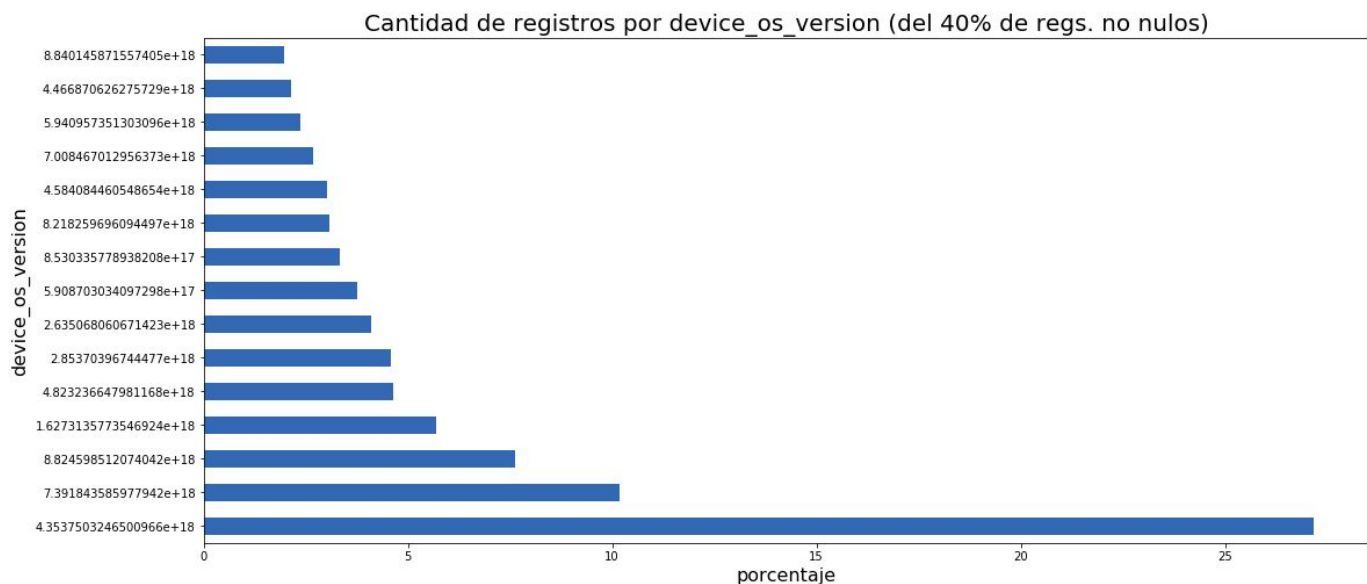
- device_city: al igual que pasó con las coordenadas, sería bueno poder hacer una representación visual de la distribución de ciudades. Así, lo que se puede rescatar es que del 25% de campos no nulos, más del 90% llevan el mismo valor.



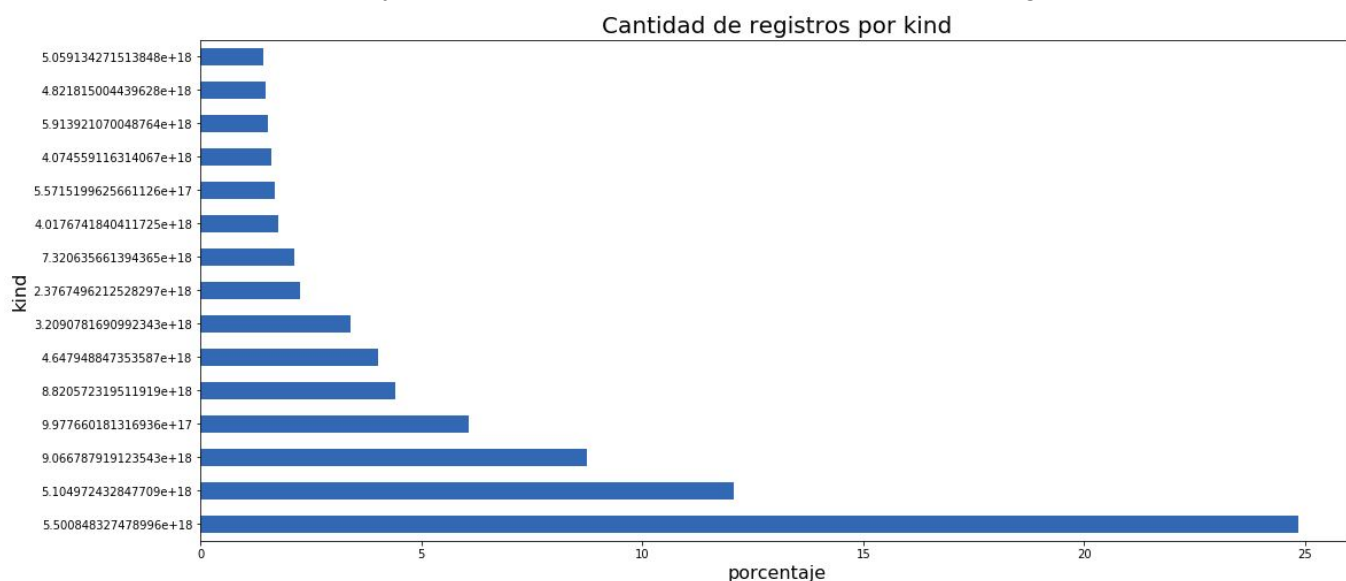
- device_os: casi el 75% de los registros no llevan valor. De los restantes, se supuso en principio habrá una gran cantidad con Android, una menor cantidad con iOS y una cantidad todavía menor con Windows Phone. Si bien no hay forma de verificar que efectivamente se trate de esos valores con las respectivas cantidades, el plot obtenido refleja lo esperado, agregando un 4to OS.



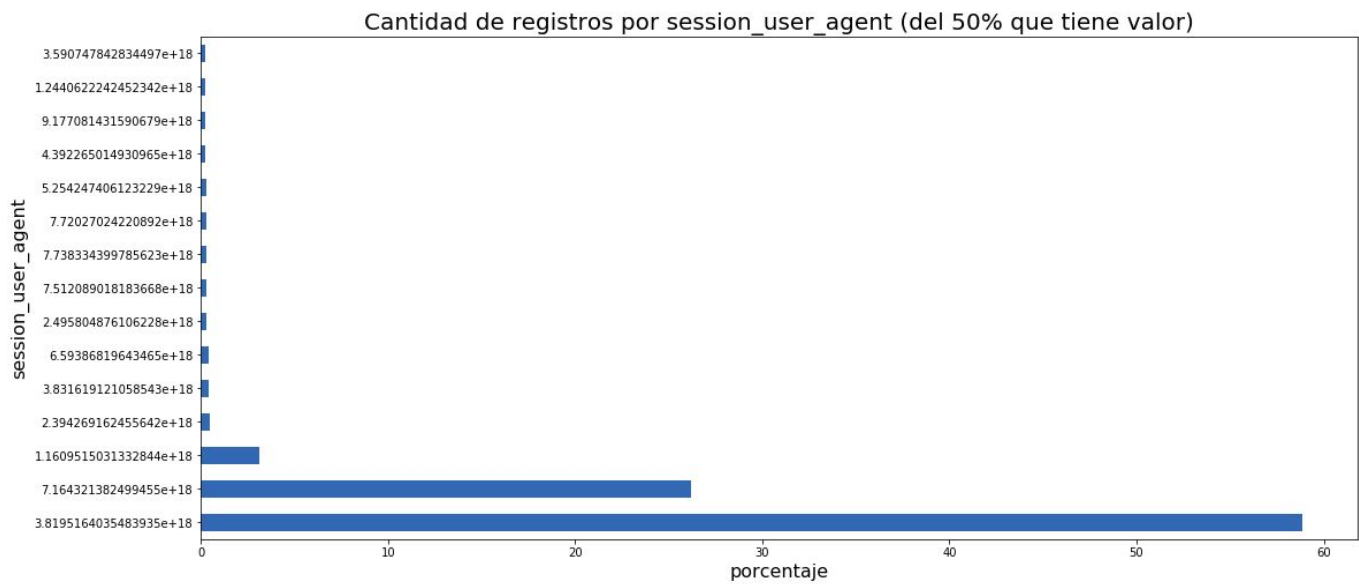
- device_os_version: la distribución sigue una ley de potencias.



- event_id: como el campo no contiene información intrínseca y además se cuenta con un identificador único de evento (event_uuid) que además es compartido entre datasets y no ha sufrido transformaciones, se decidió ignorar este campo.
- ip_address: es un dato muy ligado a la ubicación.
- kind: nuevamente la información sigue una ley de potencias. Casi no hay valores nulos, y el valor más abarcativo abarca el 25% de los registros.



- session_user_agent: aproximadamente la mitad de los registros no tiene información. De lo que resta, toman uno de los 1460 valores posibles. Sin embargo, hay un valor que abarca el (casi) 60% de los registros no nulos, y otro que toma otro (casi) 30% de los registros no nulos.

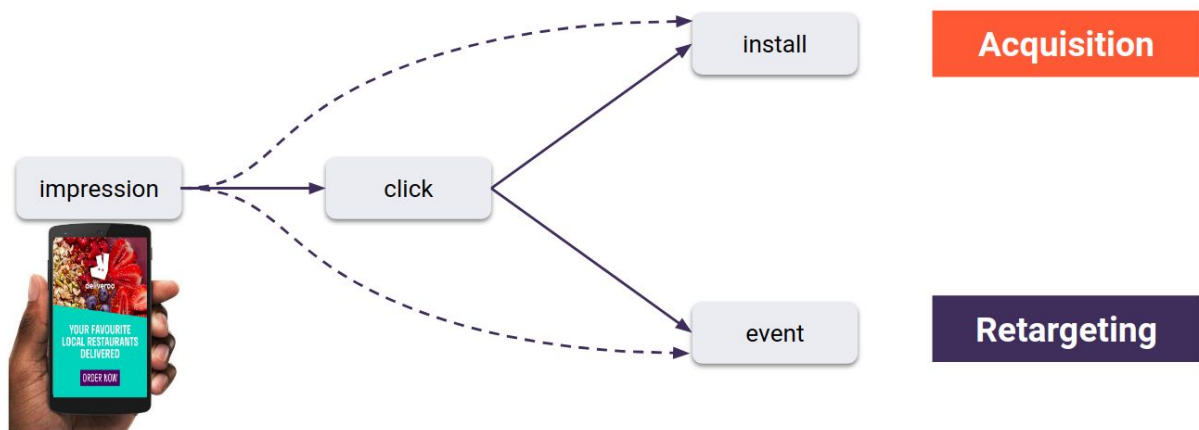


- user_agent: si bien también hay una completitud del 50% aprox., esta vez los campos toman uno de 5000 valores, y la distribución es mucho más pareja que del campo anterior. No hay ningún valor presente en más del 3% de los registros no nulos.

Análisis del flujo de datos

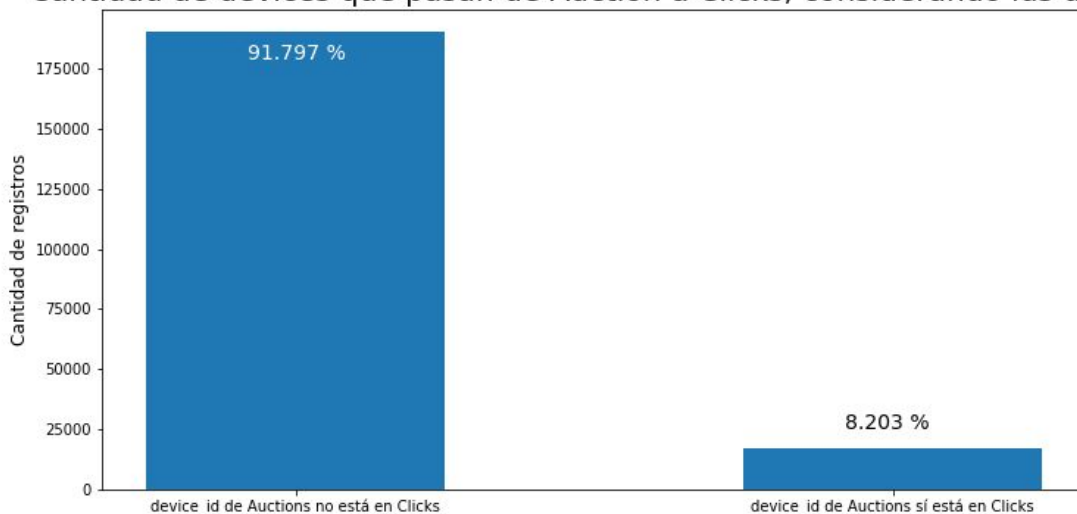
En la imagen a continuación se puede apreciar el flujo de datos según la presentación oficial de Jampp.

Transaction lifecycle: possibilities after impression

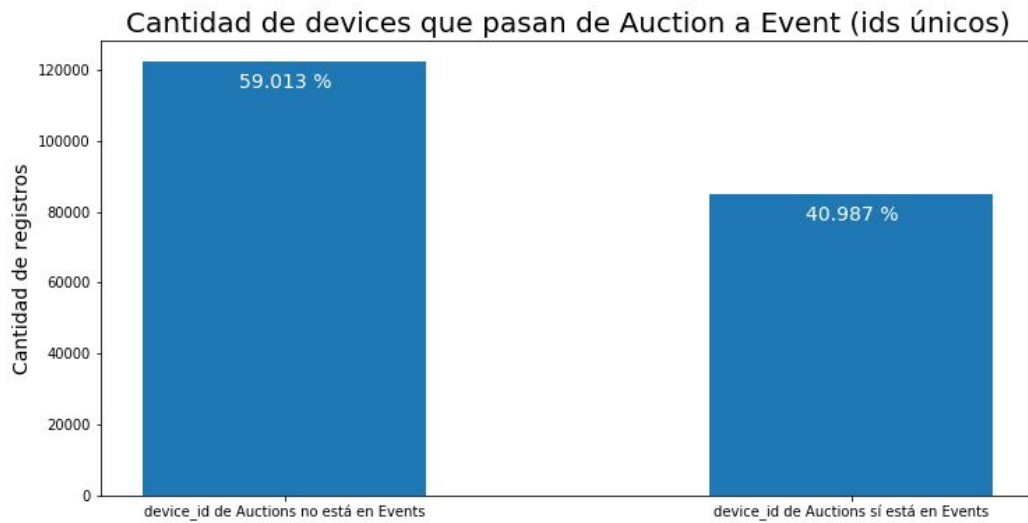


En primer lugar se observa la cantidad de dispositivos que pasaron de tener una impresión (visto en el dataset Auctions bajo la columna device_id) a un click (visto en el dataset Clicks bajo la columna ref_hash). En el gráfico siguiente se puede ver la cantidad y el porcentaje respectivo de los ids en Auctions que se encuentran también en Clicks. Como se supuso, pasa un menor porcentaje. De todas formas se esperaba una diferencia muchísimo más marcada.

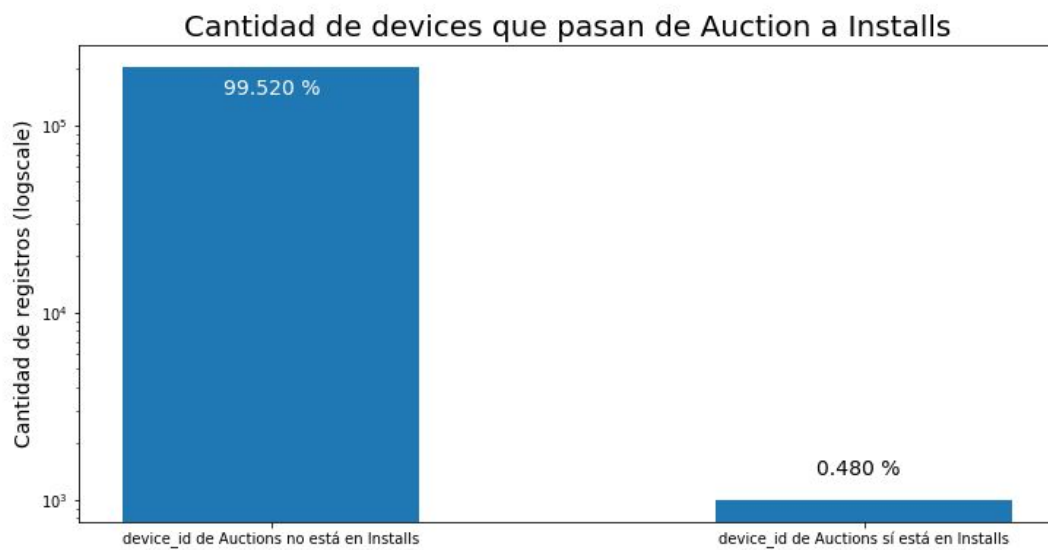
Cantidad de devices que pasan de Auction a Clicks, considerando ids únicos



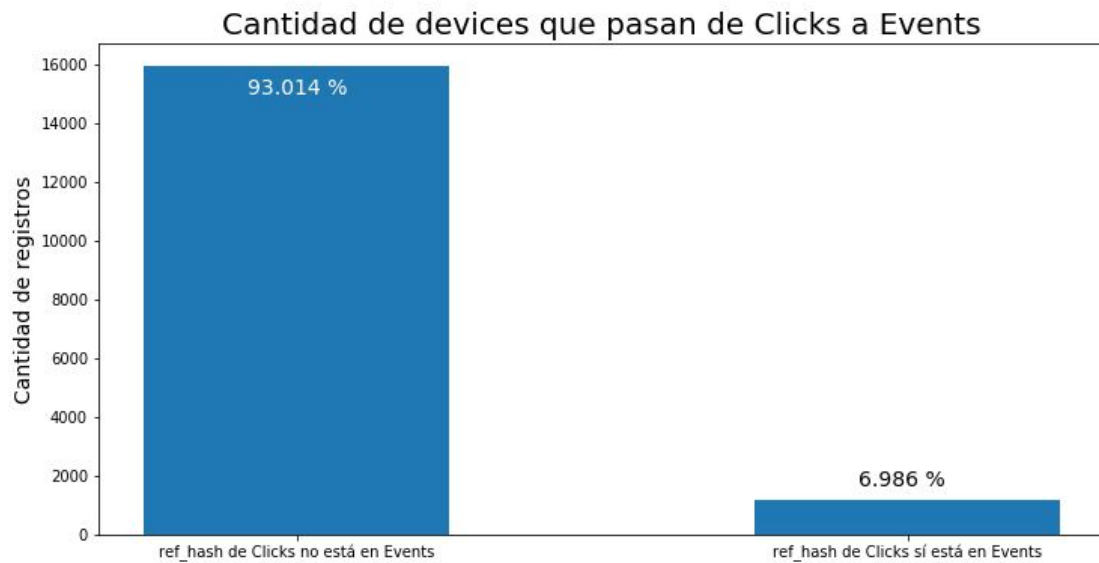
Como era de esperar, hay muchos más datos que pasan de Auctions a Events, considerando que un click es *un posible* evento.



Finalmente, de los dispositivos presentes en Auctions, no pasa casi ninguno a Installs, según se esperaba.



Devices de Clicks a Events:

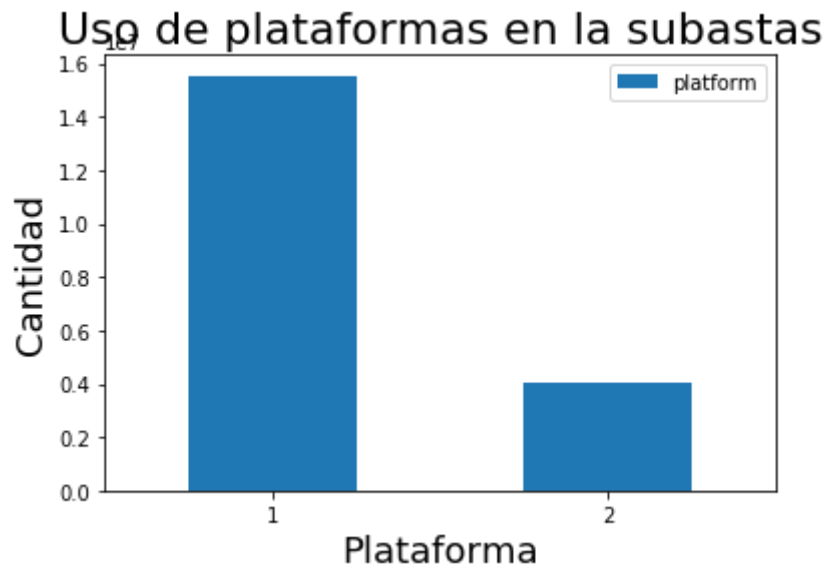


Devices que pasan de Clicks a Installs: no se grafica debido a que de 17119 valores únicos, pasan solamente 7, lo cual no se llega a apreciar en ningún gráfico.

De los más de 2 millones valores de uuid presentes en Events, solamente 860 se vuelven a repetir en Installs.

Análisis de Auctions

¿Cuáles son las plataformas que participan de las subastas (android / ios)? y ¿Cual es la plataforma con mayor participación en las subastas?



Vemos cuáles son las plataformas que participan de las subastas. Recordando que los datos se encuentran anonimizados (únicamente sabemos que existen ios y android) observamos dos plataformas participantes (1 y 2 respectivamente) siendo la 1 la que posee una mayor participación.

¿De qué fecha obtenemos la información? y ¿Cuántas subastas por día hay en el periodo?

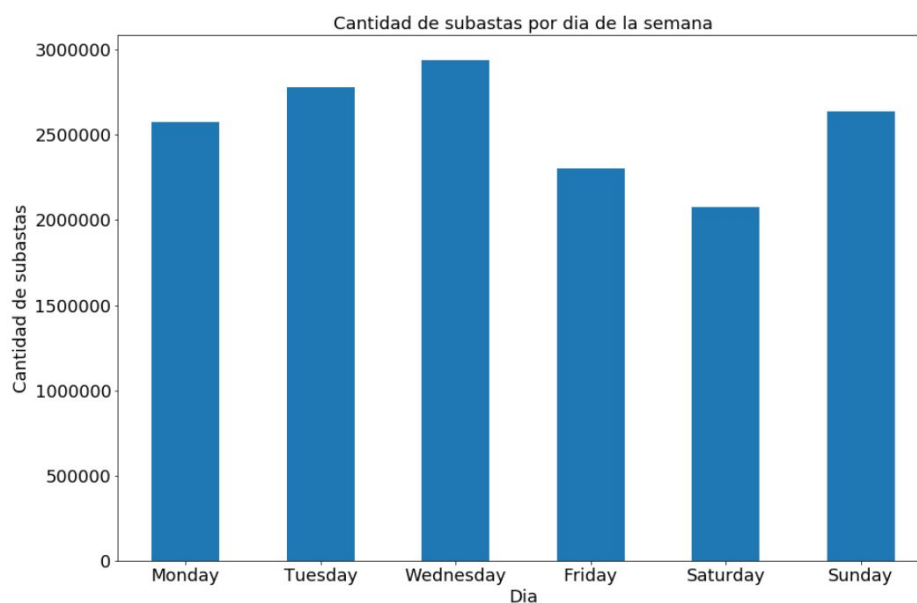
Observamos que tenemos información del 3 al 13 de marzo.

		total
dayofyear	date_single	
64	2019-03-05	1182401
65	2019-03-06	1032970
66	2019-03-07	2047661
67	2019-03-08	2303002
68	2019-03-09	2074552
69	2019-03-10	2637534
70	2019-03-11	2574916
71	2019-03-12	2779910
72	2019-03-13	2938373

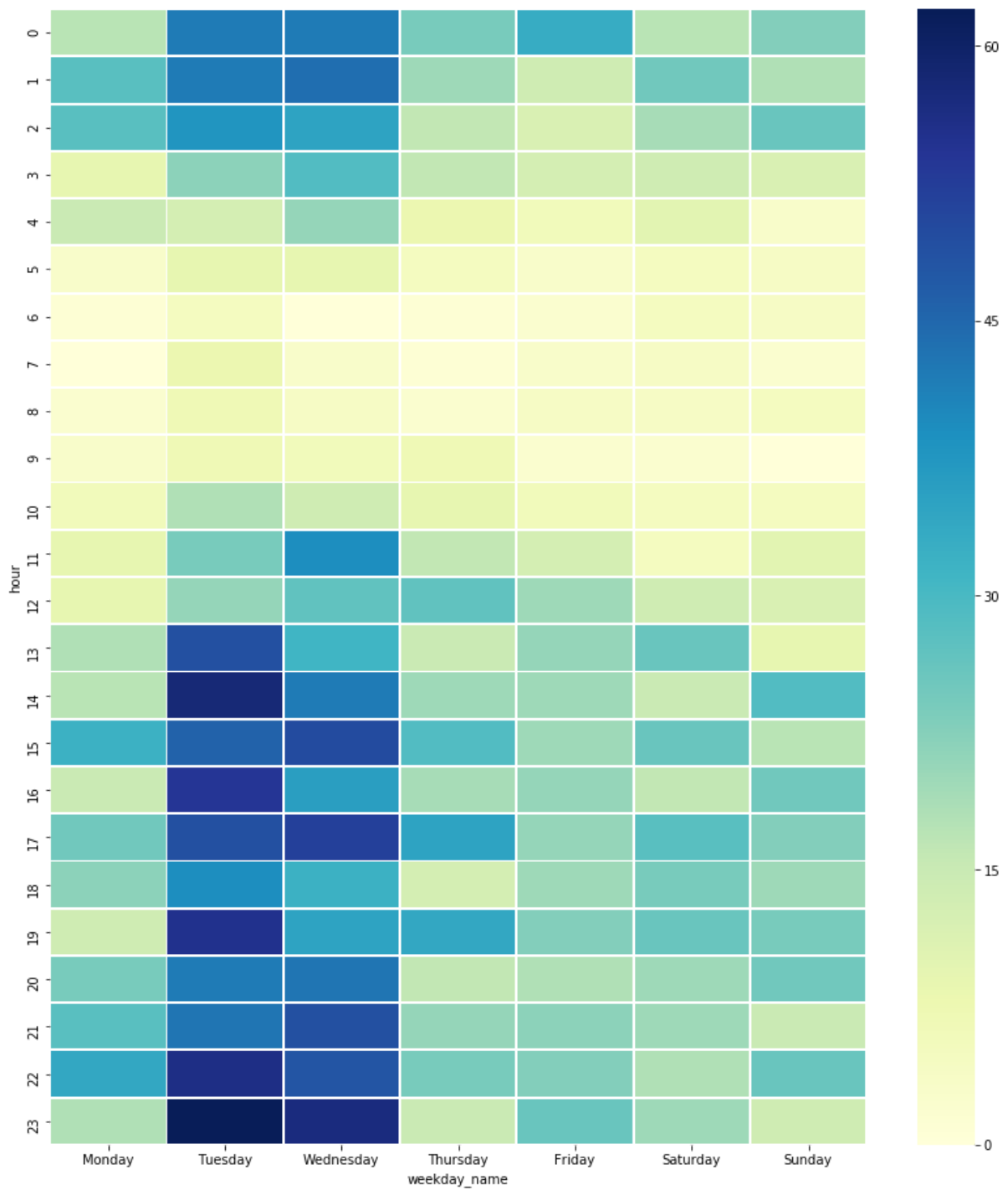


Observamos que se proyecta sobre la cantidad de subastas una tendencia alcista. Pero el periodo de fechas sobre el cual se trabaja es de una semana para decidir si esta tendencia permanecerá a lo largo del tiempo o sufrirá alguna fluctuación.

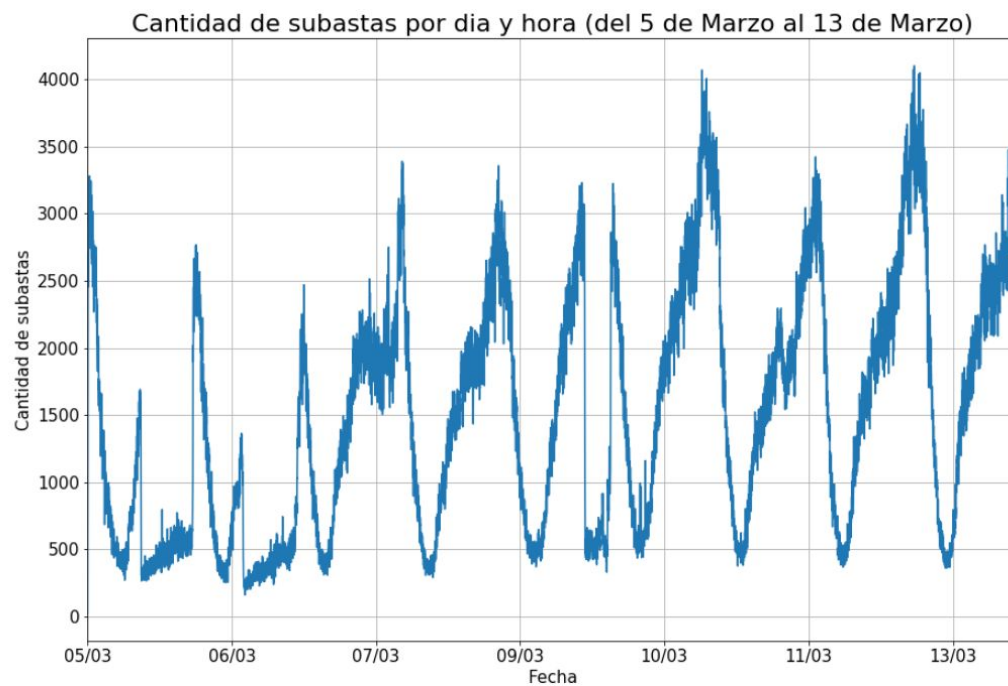
¿Cuál es el día que se producen más subastas en 1 semana?



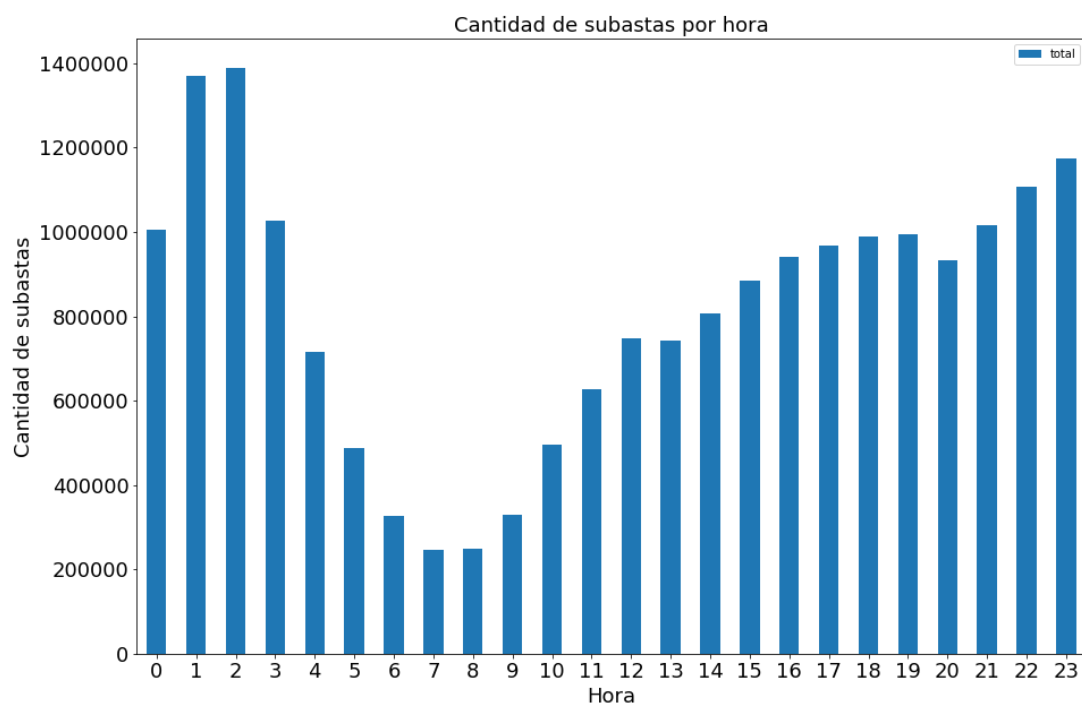
Heatmap por día y cantidad de participaciones en subastas



¿Cual es la cantidad de subastas por hora y por día?

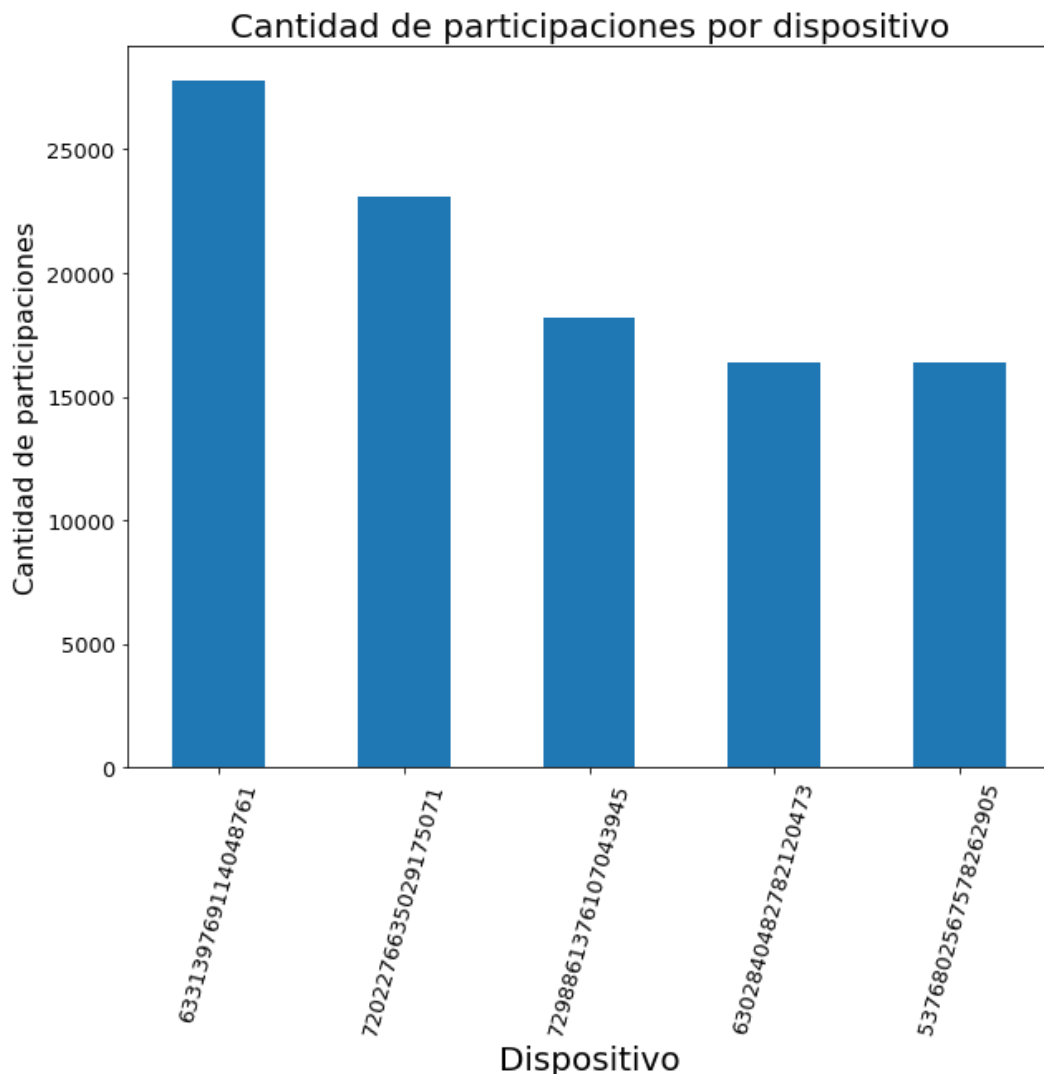


¿A qué hora se producen la mayor y la menor cantidad de subastas?

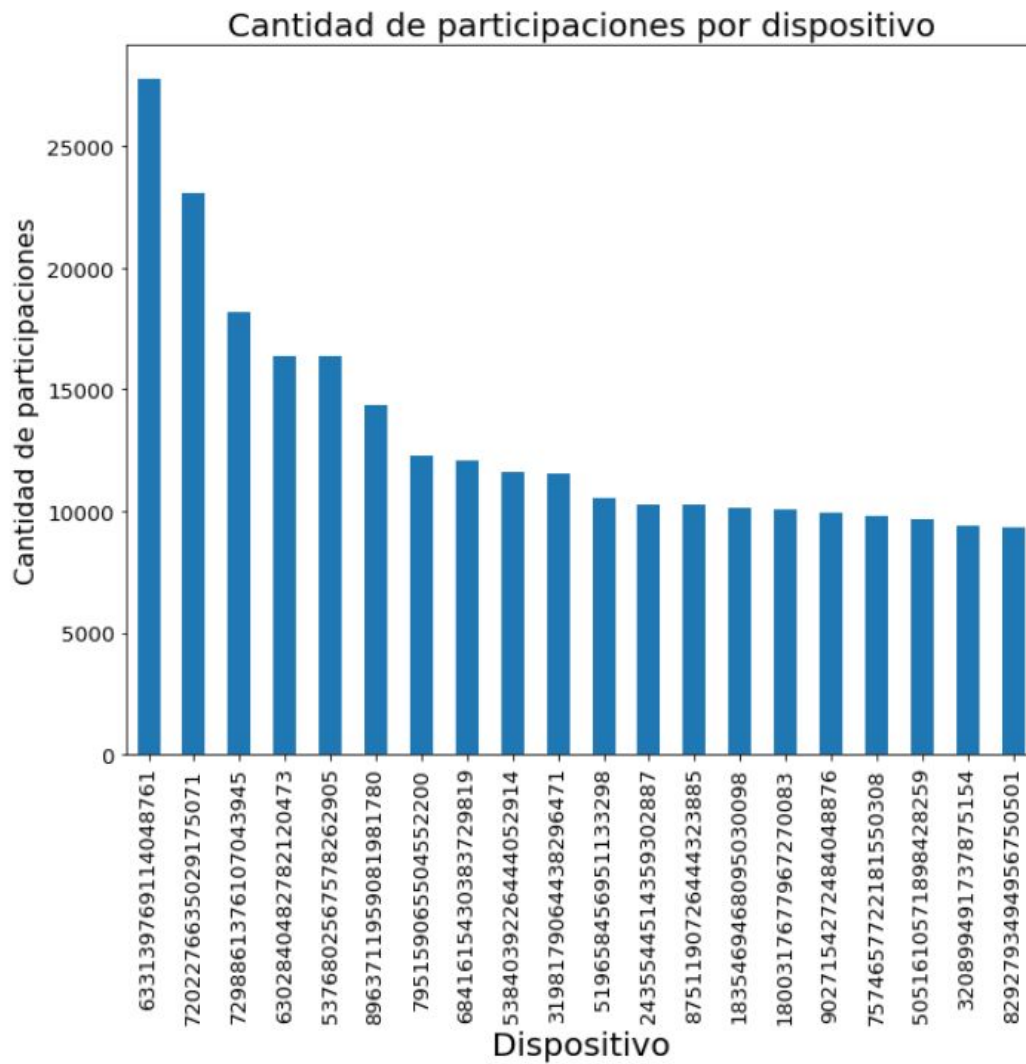


Observamos que las horas donde se encontraron mayor cantidad de subastas son entre la 1 y las 2 de la mañana y el periodo entre las 6 las 9 de la mañana donde se produjo la menor cantidad de subastas. Como observación vemos que desde la mañana a medida que suceden las horas aumenta la cantidad de subastas, llegando al pico máximo al inicio de la madrugada y luego decayendo hasta las primeras horas de la mañana.

¿Cuales fueron los 5 dispositivos que tuvieron mayores participaciones en las subastas?

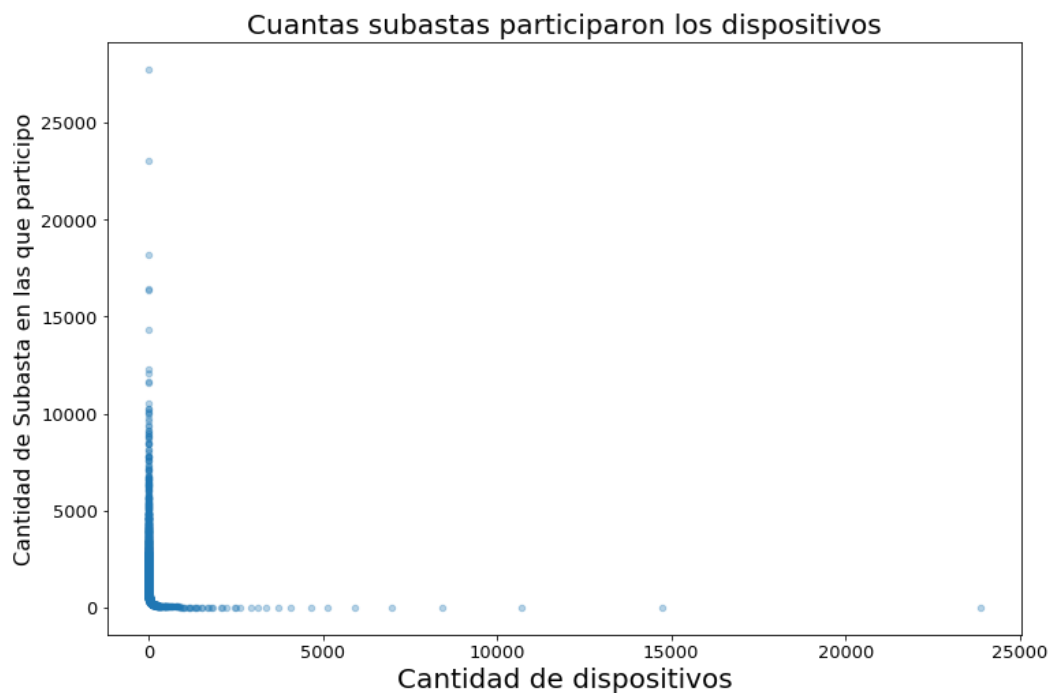


Aquí se observan los dispositivos que tuvieron mayor participación en la subastas. Estos números de participación son realmente muy altos. Estos nos podría indicar algún tipo de fraude o granja de instalación de apps.

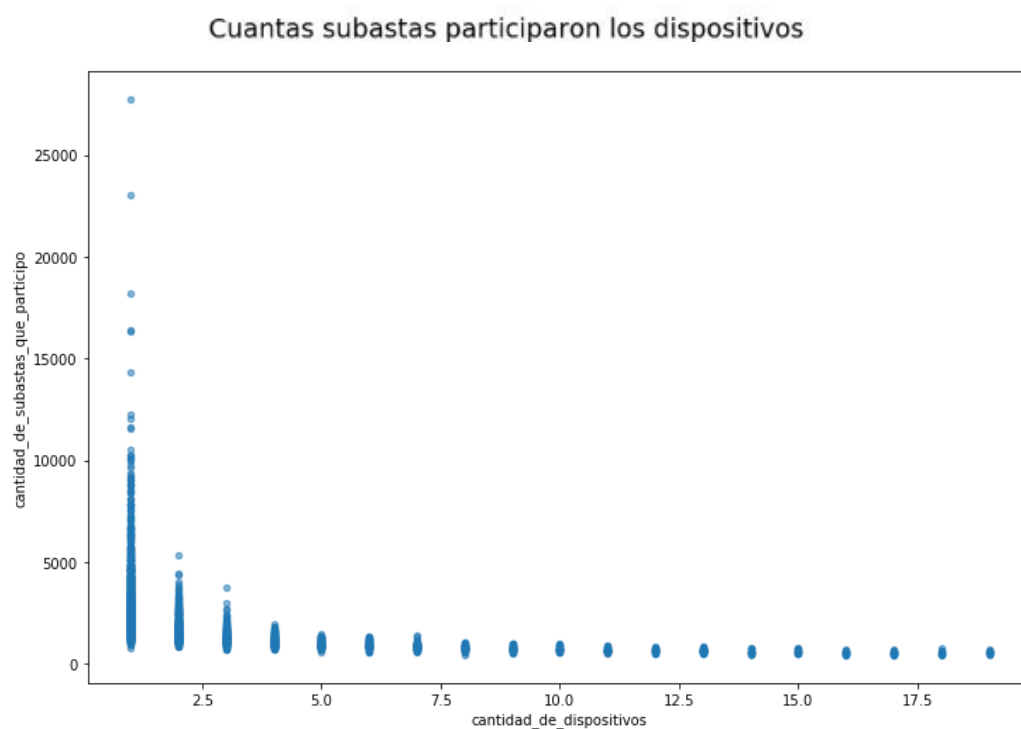


Viendo la gráfica anterior decidimos añadir un gráfico con los primeros 20 y vemos su participación en las subastas. Continúan siendo una participación excesiva en las subastas.

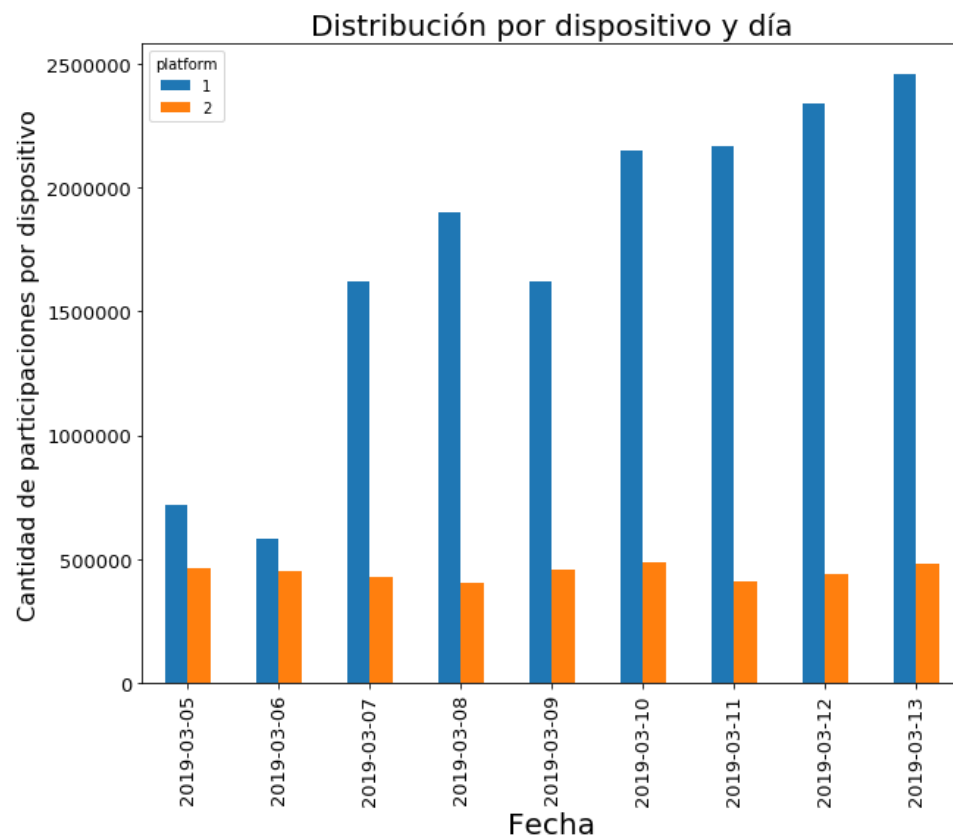
¿Cuántas participaciones en subastas por dispositivo hay en el periodo?



Aquí observamos que hay dispositivos que participan de subastas de manera excesiva para el periodo de tiempo del set de datos. Lo vimos en los gráficos anteriores que su cantidad de participación es muy alta. Reducimos la cantidad de dispositivos observados a 20.

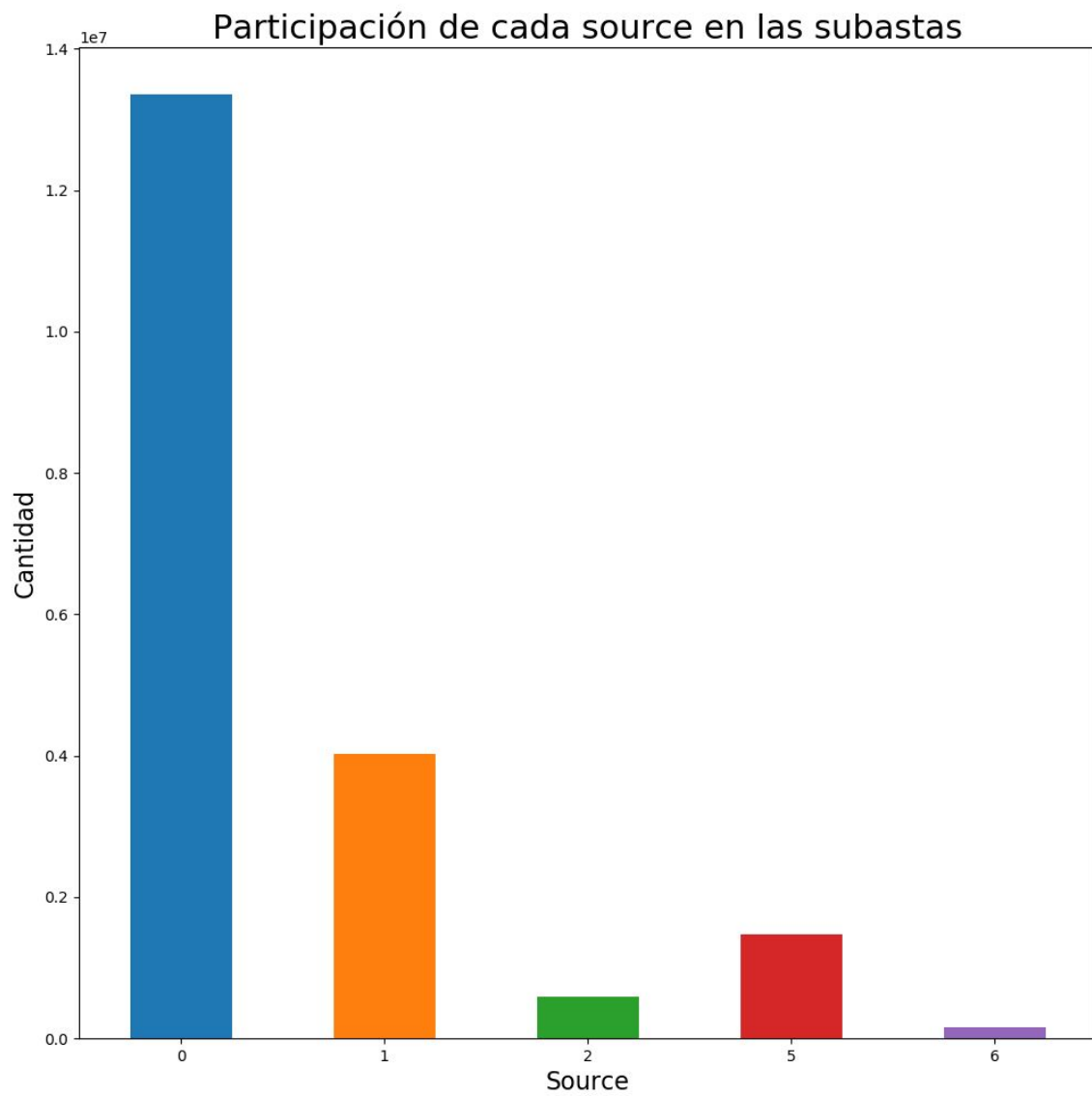


¿Como es la distribución del participación por día según la plataforma?

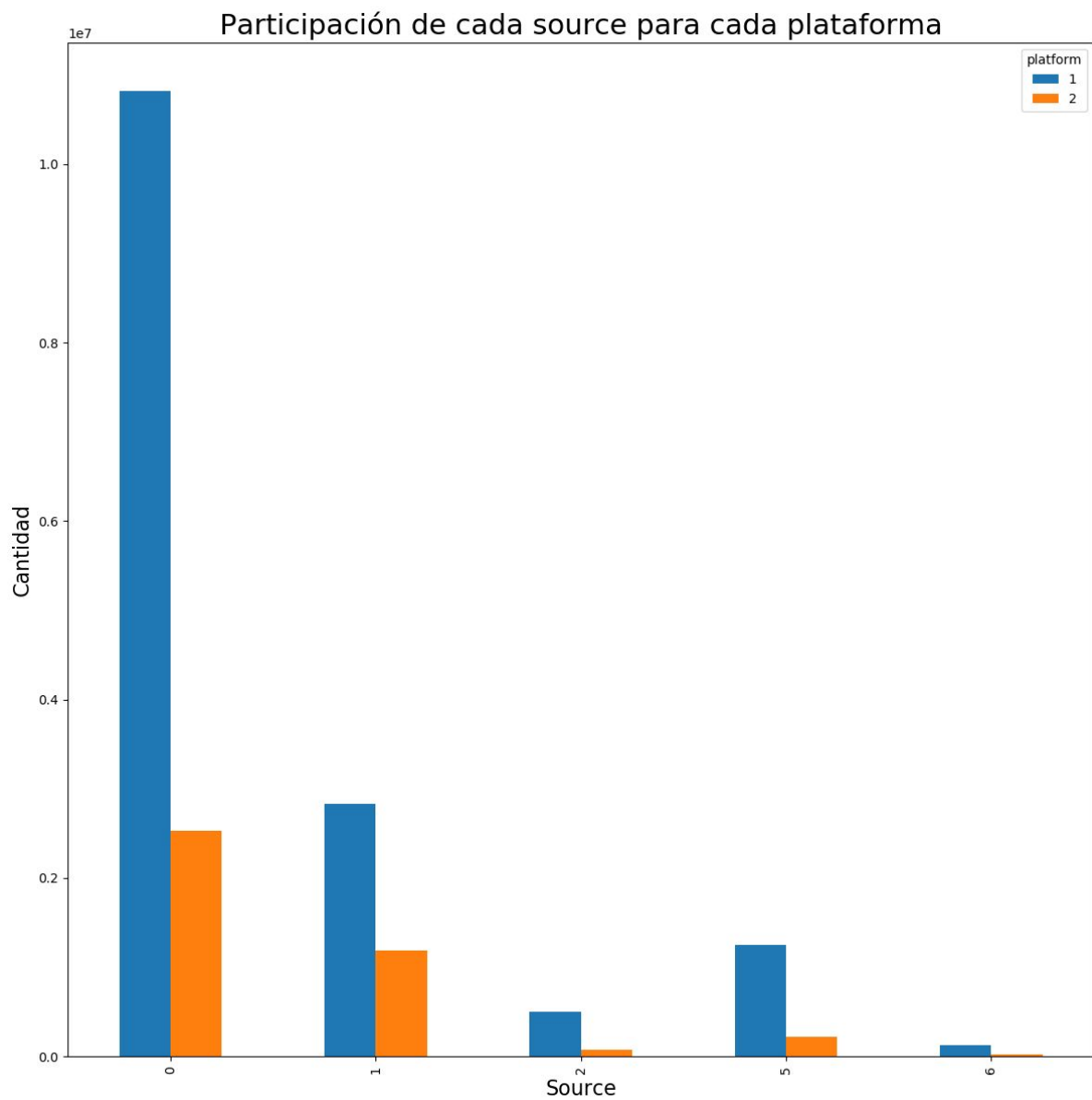


Aquí se ve que la plataforma 1 se observa un incremento diario de participación. En cuanto a plataforma 2 mantuvo una participación pareja durante el periodo siendo casi la misma participación por día.

¿Cómo es la participación de los sources en las subastas?

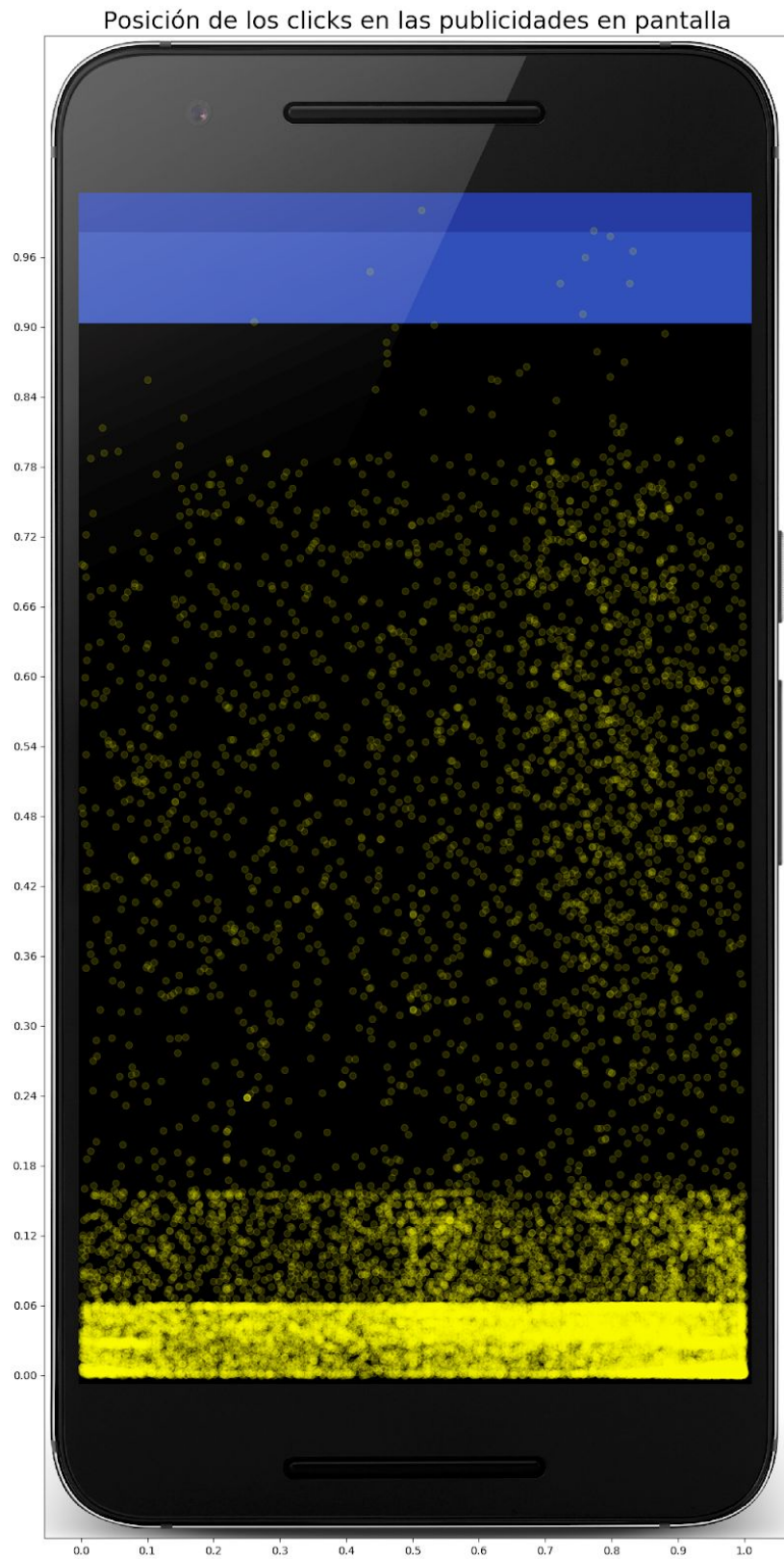


¿Cómo es la participación de los sources en las subastas para cada plataforma?



Análisis de Clicks

¿Cómo se distribuyen los clicks en la pantalla de los teléfonos?



Como se puede ver claramente en el gráfico, las publicidades suelen aparecer en la parte inferior de la pantalla de los teléfonos.

Investigando un poco acerca de anuncios en aplicaciones (App advertising⁸), en general existen dos tipos de anuncios: los de *pantalla completa* y los *banners*.

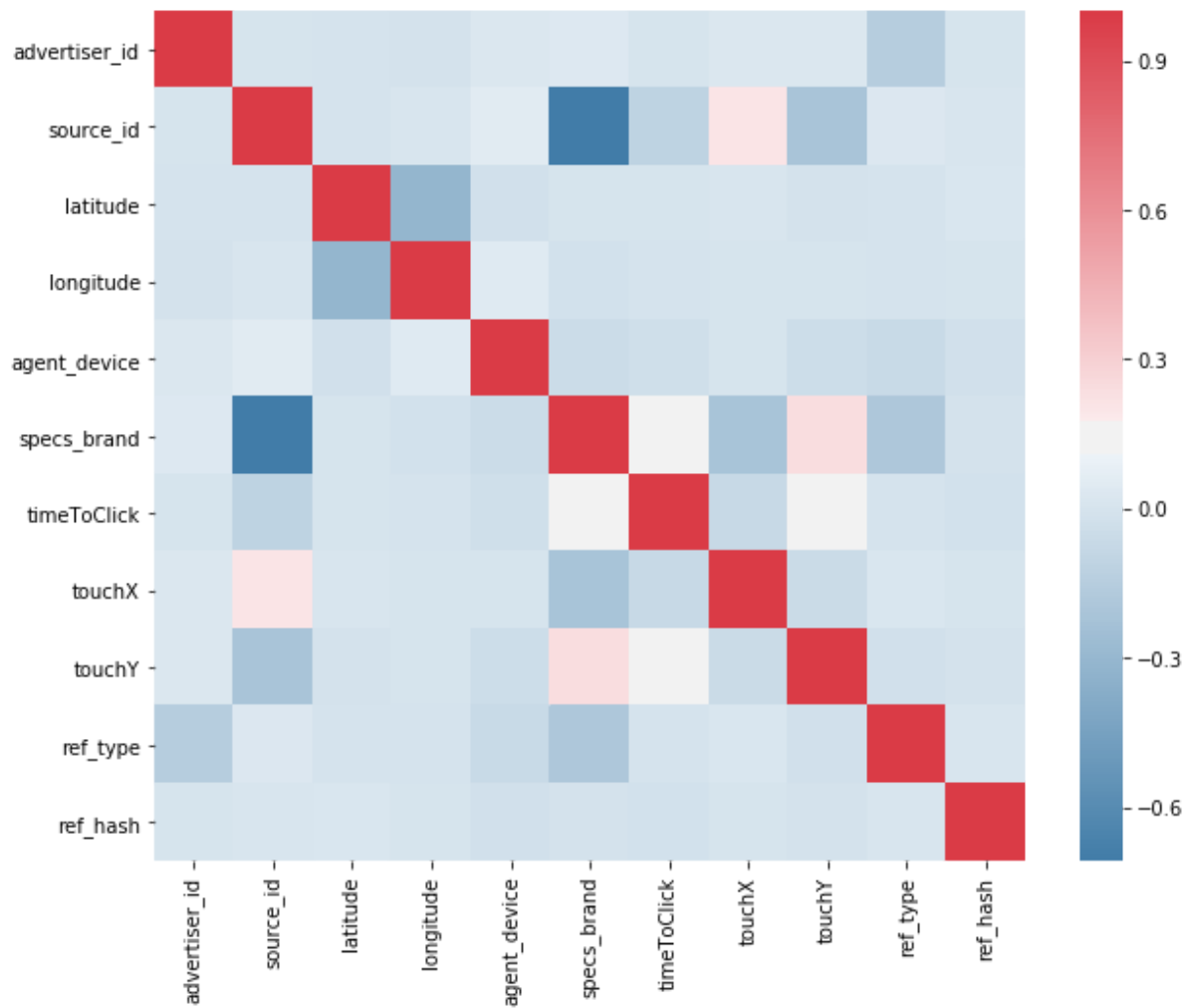
Los *anuncios de pantalla completa* brindan más espacio para un mensaje más amplio. Llamam la atención gracias al contenido de medios interactivos y visuales. Sin embargo, hay un inconveniente, este formato de publicidad en la aplicación es un poco intrusivo.

En cambio, el *banner* puede integrarse sin esfuerzos particulares. Es la opción más barata que puede usarse en cualquier pantalla. La única contra de este formato de anuncio puede ser presa del síndrome de la "ceguera del banner"; esto ocurre cuando el usuario ignora la información del mismo.

Dicho esto, si bien cada uno tiene sus problemas, posiblemente la opción del banner sea la más económica y fácil de integrar, y por este motivo, la más usada (en relación a los usuarios que los clientes de Jampp manejan).

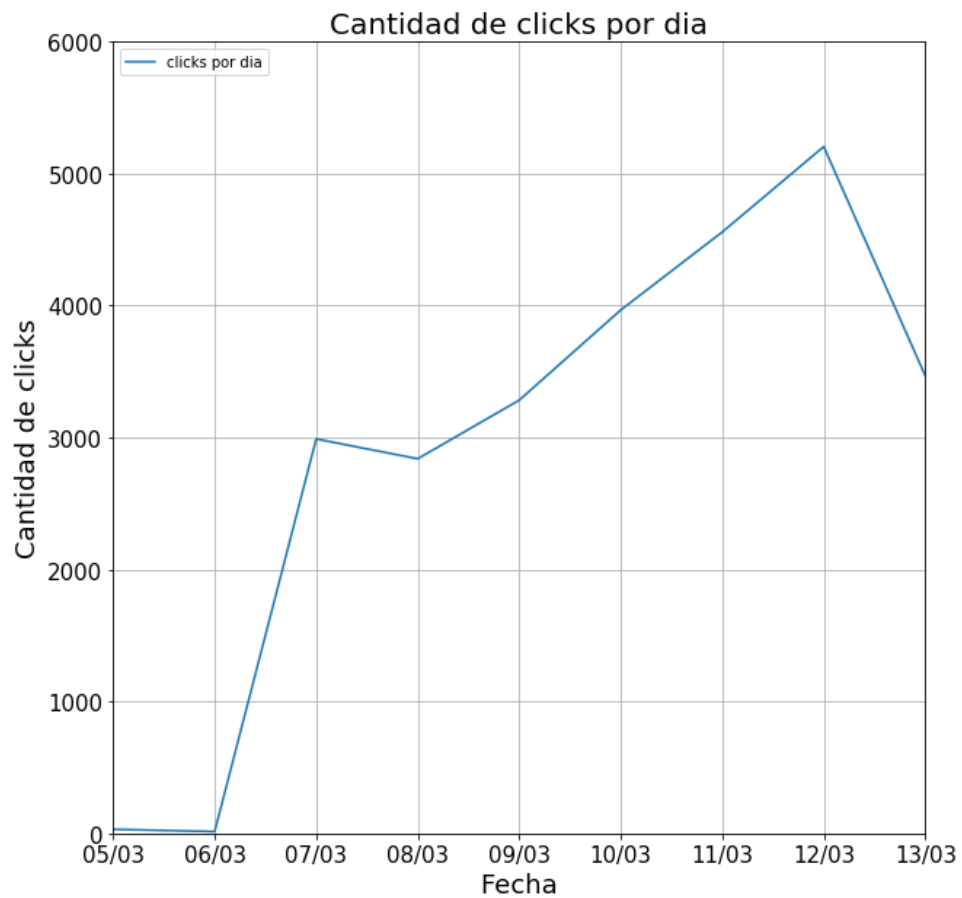
⁸ <https://splitmetrics.com/blog/guide-to-mobile-app-advertising/>

Heatmap del set de datos



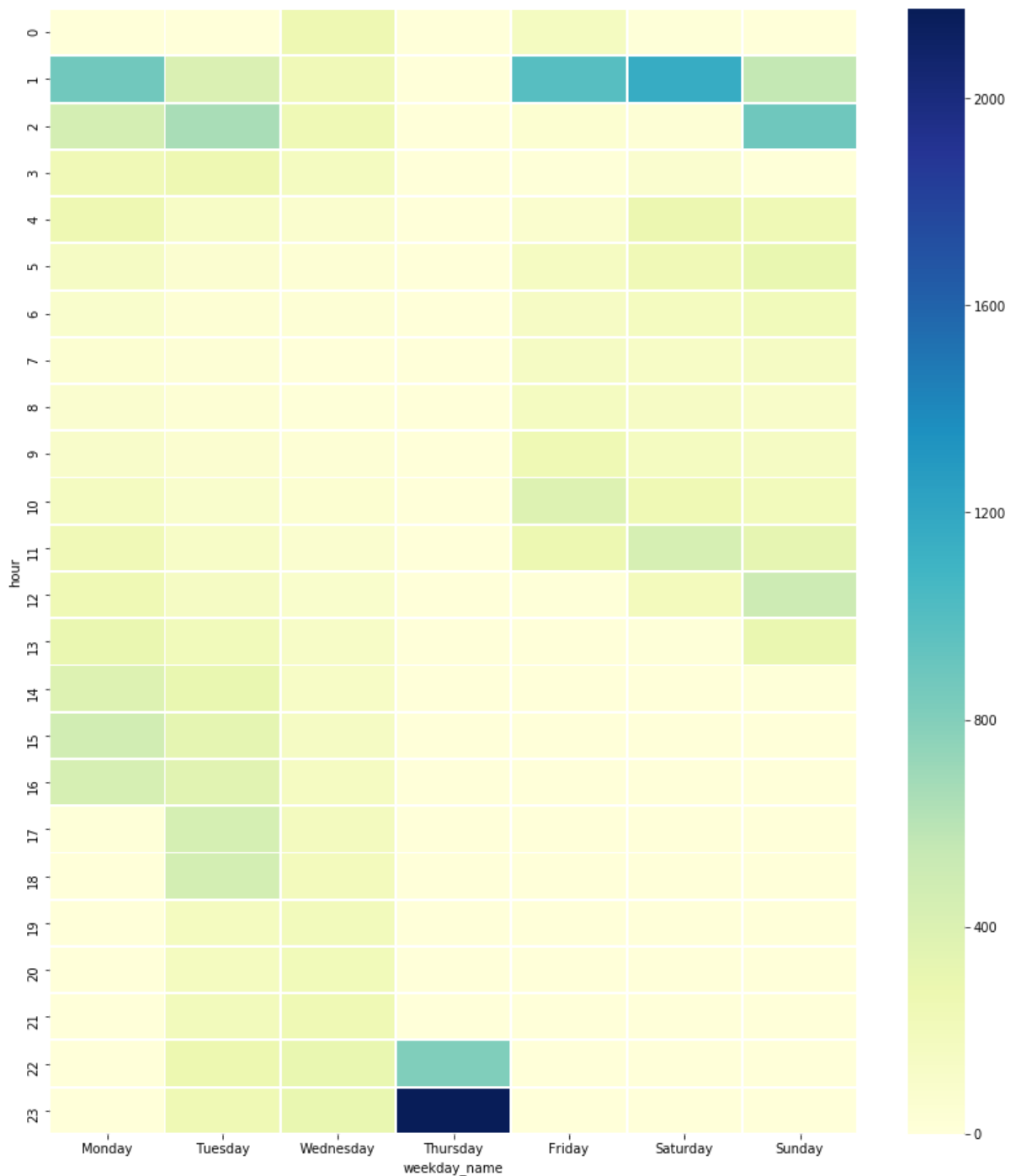
No se encuentran pares de features con correlación significativa. La mayor correlación (negativa) se encuentra entre specs_brand y source_id.

¿De qué fecha obtenemos la información? y ¿Cuántos clicks por día hay en el periodo?



El día Martes es el que se produce más clicks. Se mantiene la tendencia observada en los otros datasets. Crece a medida que surge el periodo pero en este caso el ultimo día se produce un descenso.

Heatmap de clicks por hora y días de la semana

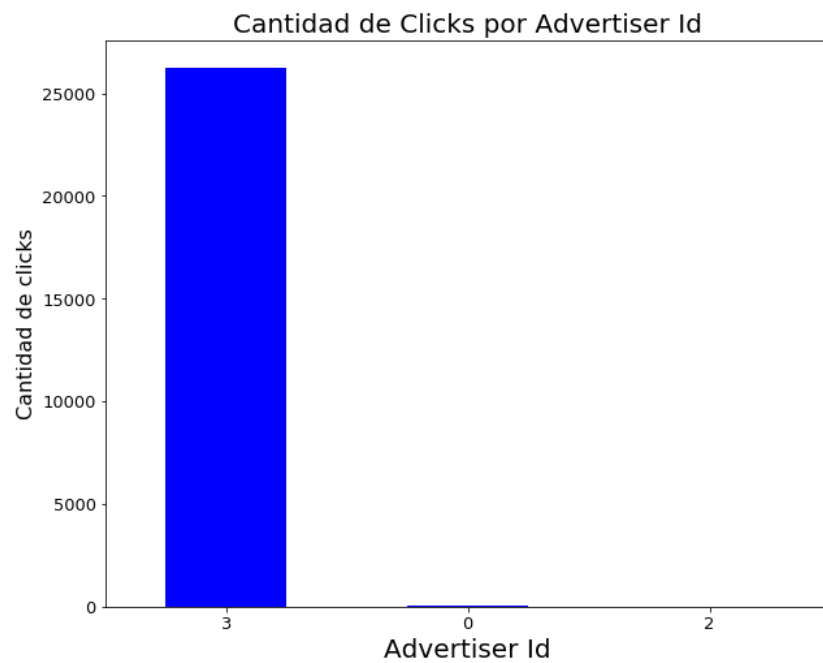


Se puede ver que los Lunes, Martes, Viernes y sábados, los horarios más frecuentados son a las 1 y 2 am. Los jueves a las 22 y más a las 23 están especialmente marcados. Habrá sido circunstancial de los días elegidos para formar el dataset, o será un comportamiento que se extiende en el tiempo?

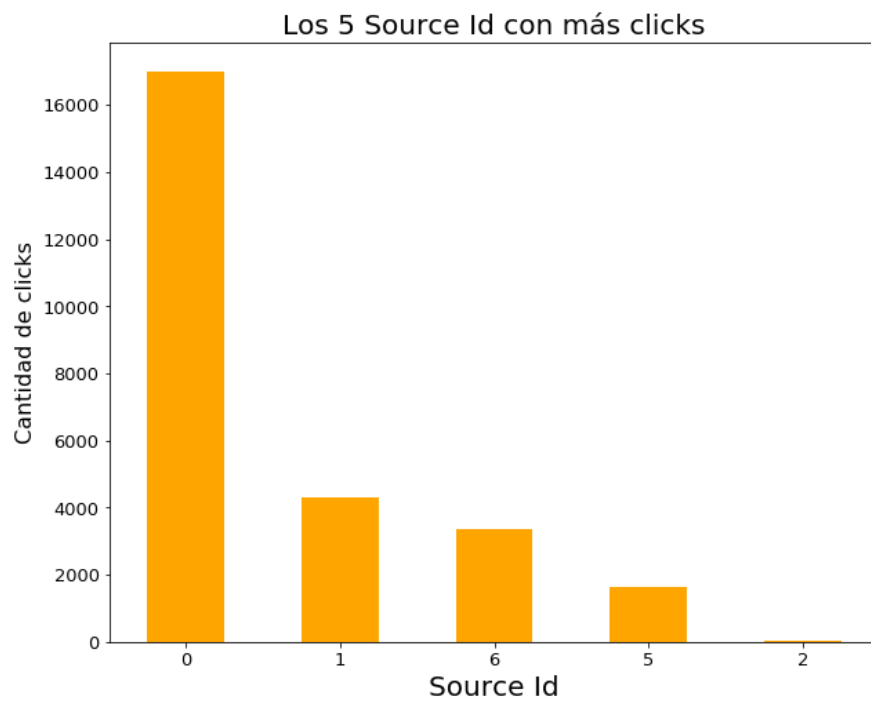
¿De qué anunciante provienen los clicks?

En este caso nos encontramos con que casi la totalidad provienen del mismo anunciante.

total	
advertiser_id	
3	26263
0	70
2	12



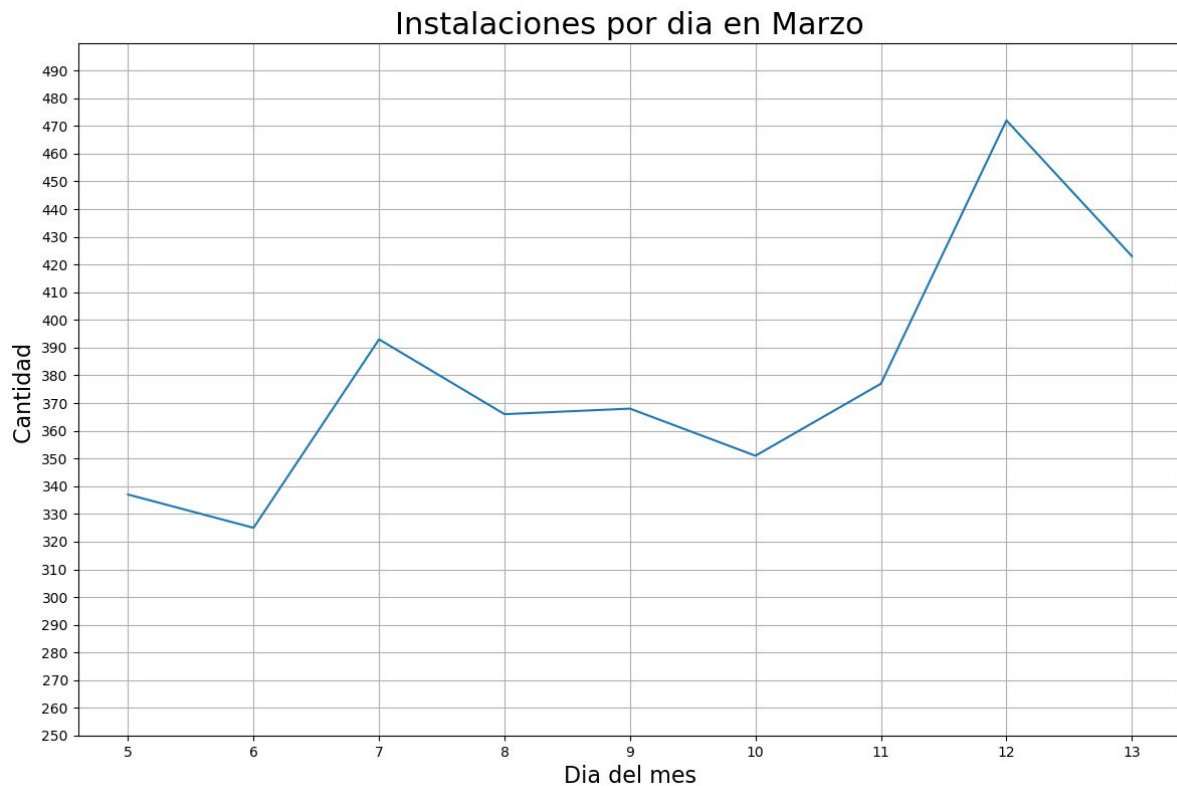
¿Cuales son las 5 fuente con más clicks?



Aquí se observa una predominancia del source 0 sobre los restantes.

Análisis de Installs

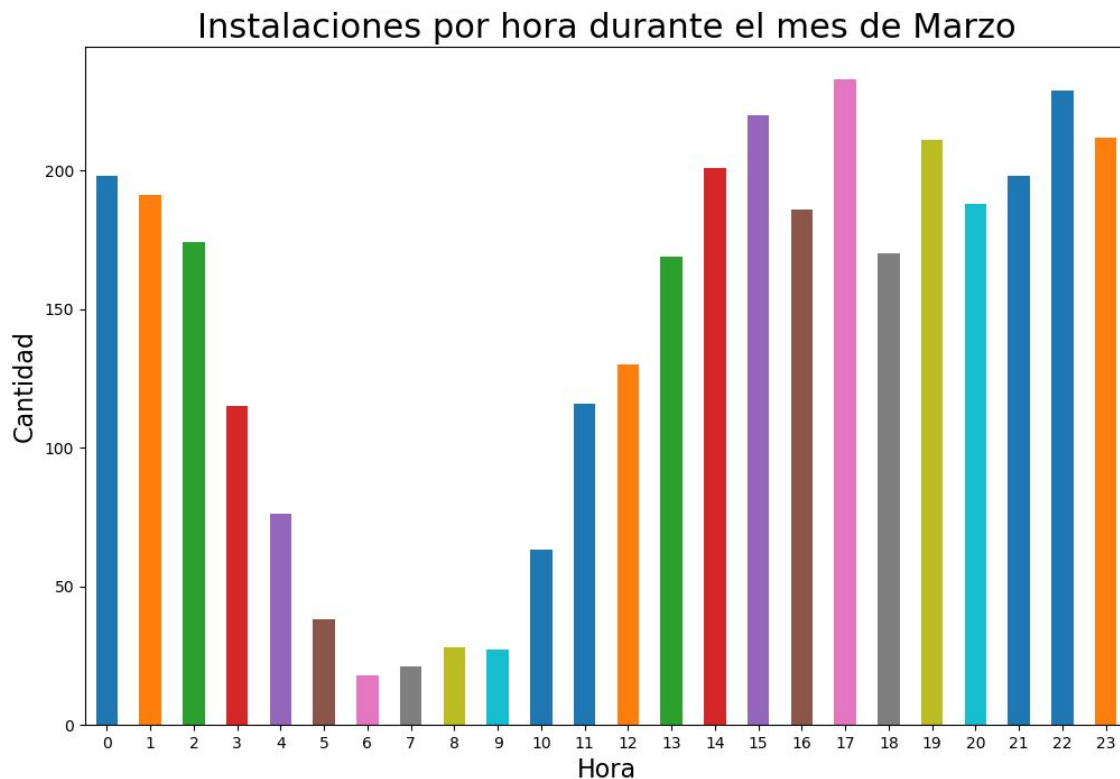
¿Cuántas instalaciones hubo en el periodo?



Aquí podemos ver que en el día doce es cuando más instalaciones se producen, esto no quiere decir que el próximo mes vaya a suceder lo mismo.

Si esto fuera en un lapso diferente y más amplio, podríamos darnos cuenta de alguna tendencia en la cual usuarios tienden a instalar más por ejemplo, pero lamentablemente aquí no es posible este tipo de análisis.

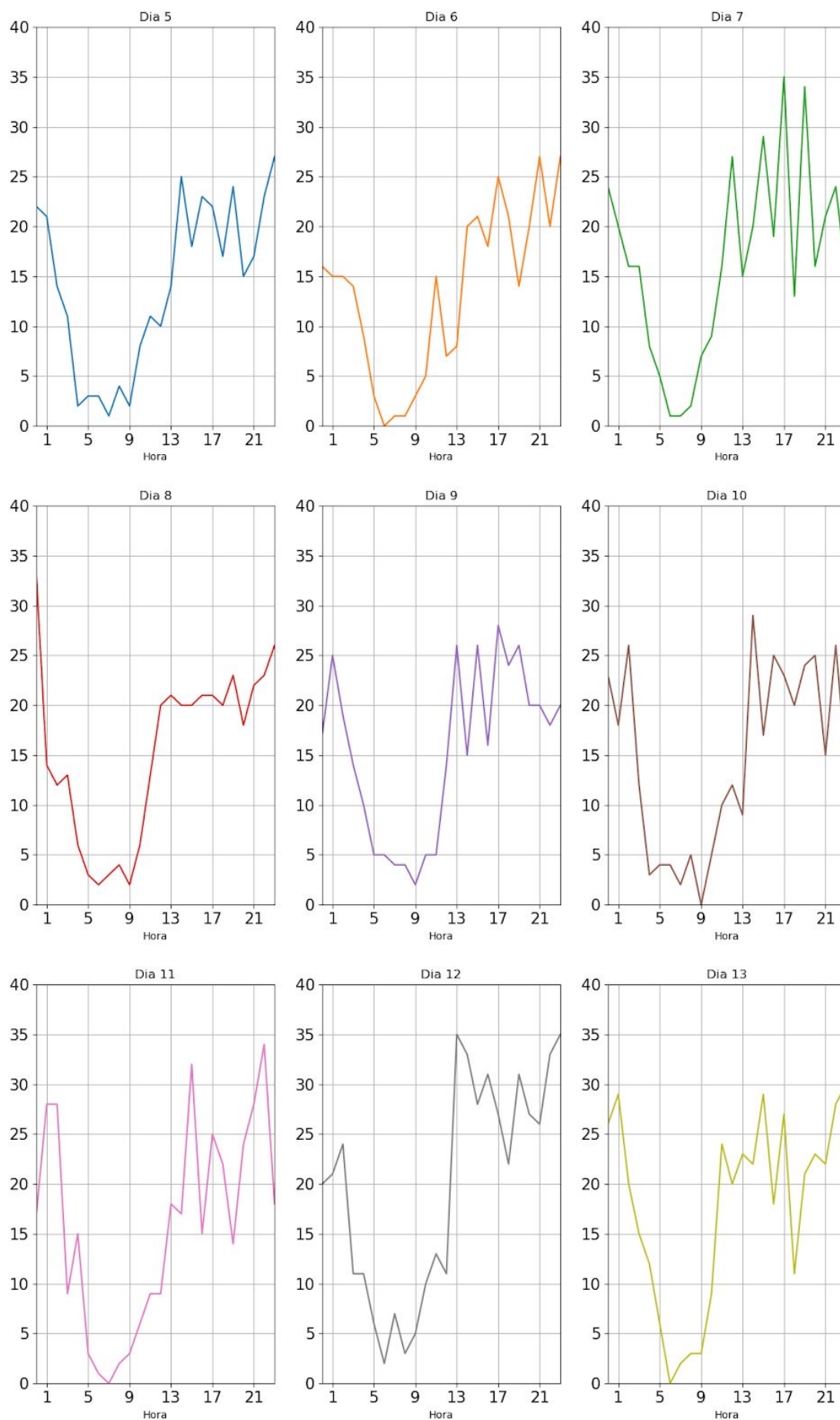
¿A qué hora se producen la mayor y la menor cantidad de instalaciones?



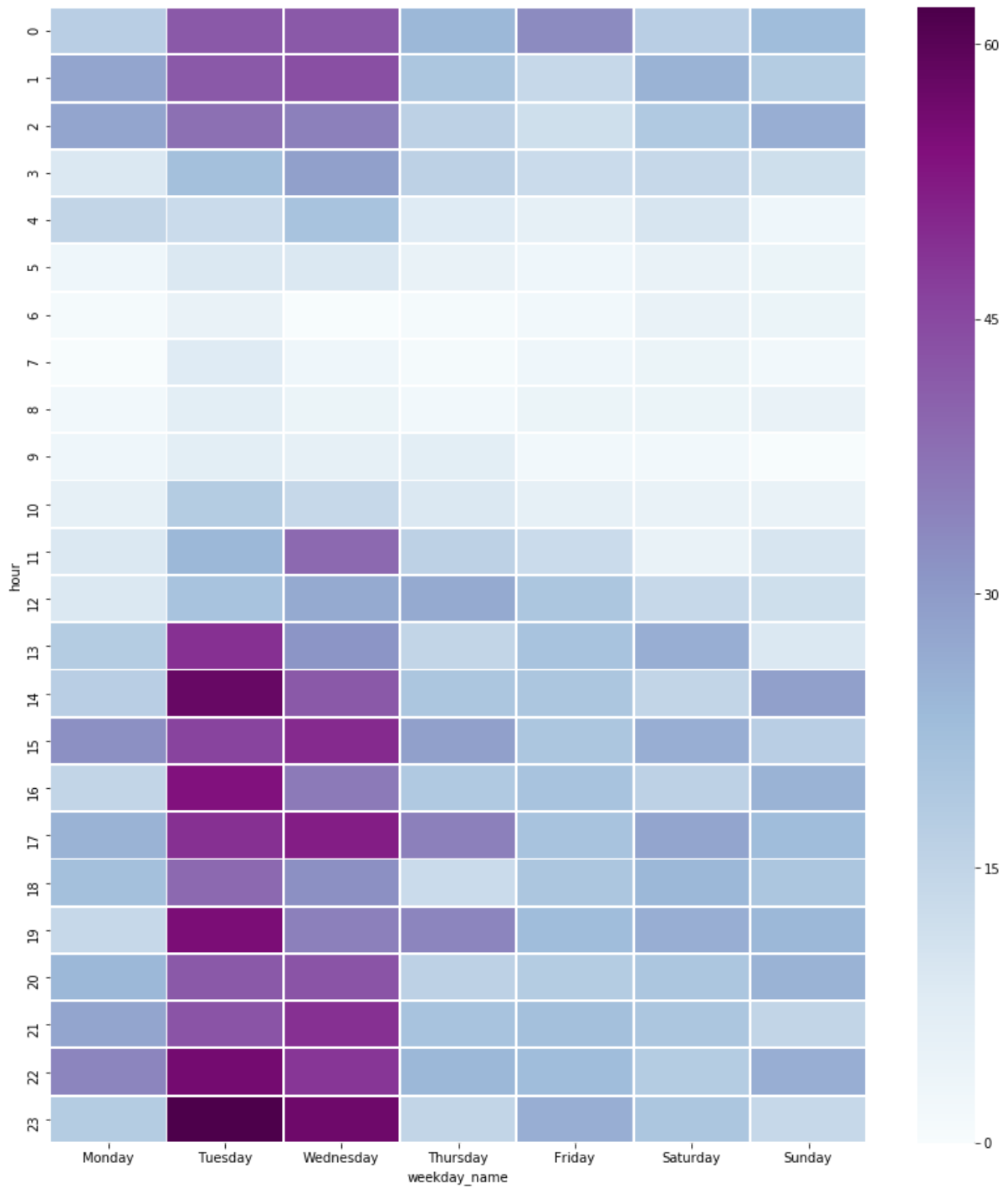
Volviendo un poco hacia lo mencionado de la representación de la realidad, en este gráfico, tanto como en aquellos a continuación (horas para cada día de marzo), se ve que durante la noche la interacción de usuarios disminuye.

Lo único quizá notable aquí es que hasta las 2 AM, la interacción de usuarios continúa de manera no muy distinta a lo que puede ser por ejemplo la interacción a las 4 PM, lo que se podría uno imaginar que no debería suceder ya que es un horario en el que la gente generalmente duerme.

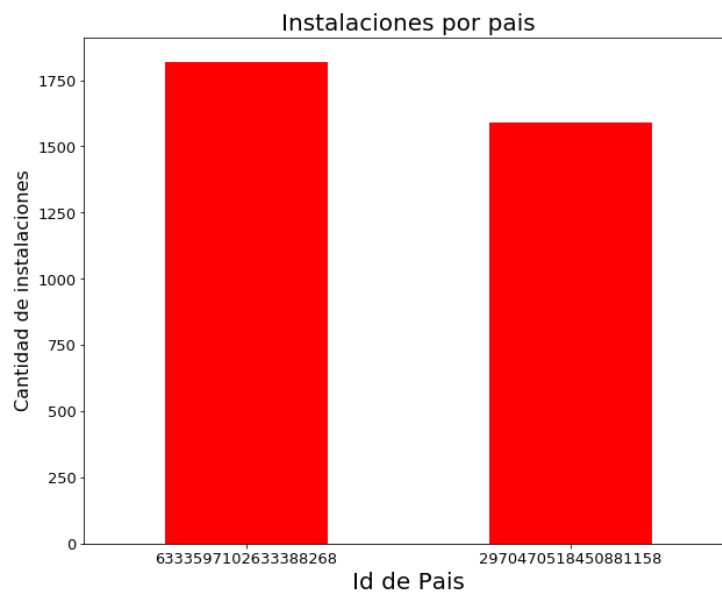
Además, vale la mención de que a las 5 PM es cuando más instalaciones se realizan.



Heatmap de hora y días de la semana

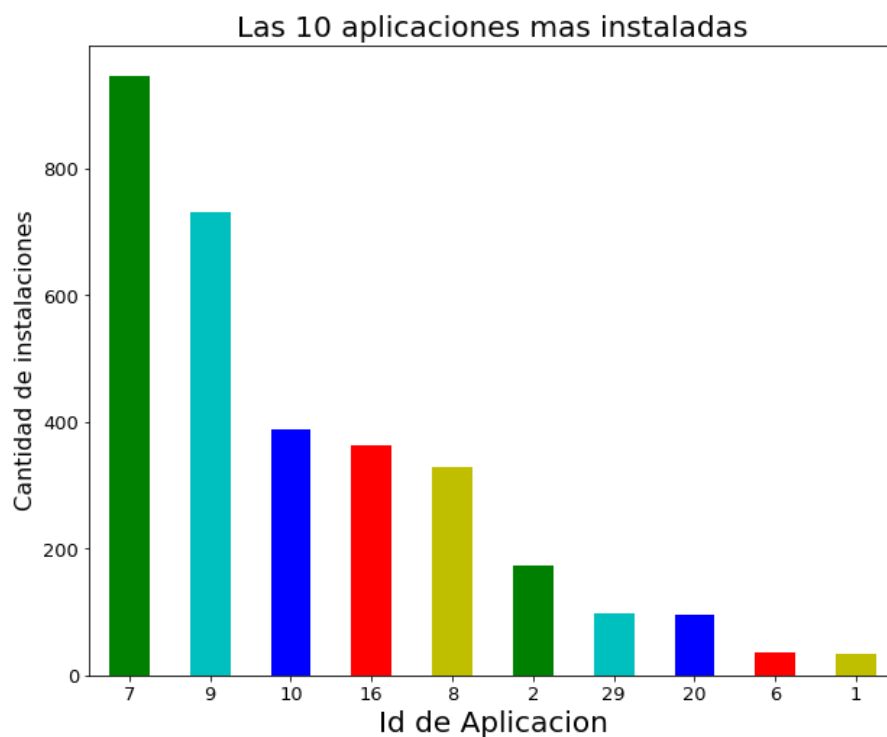


¿Cuántas instalaciones hay por país?



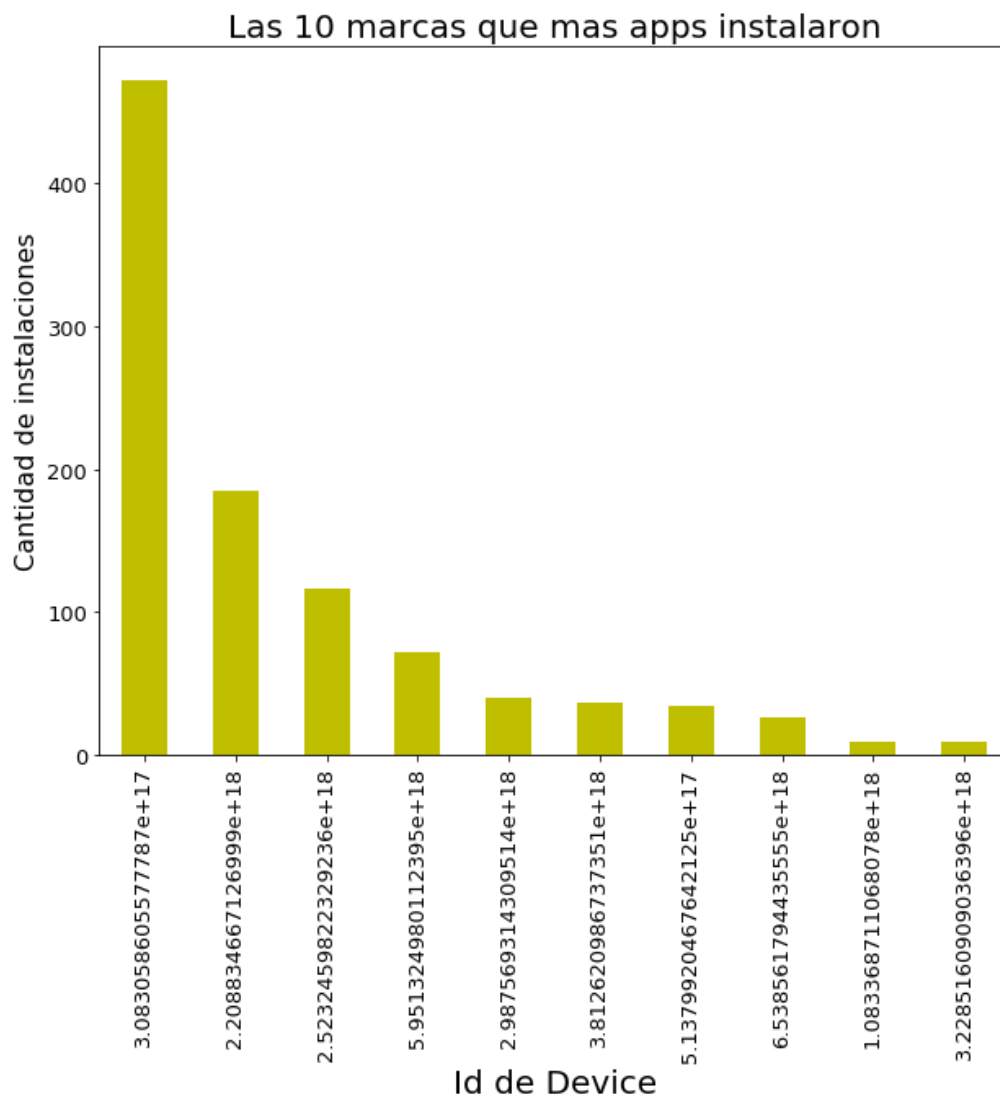
Aquí vemos que el set de datos nos muestra instalaciones que fueron realizadas desde 2 países, distinto por ejemplo a lo encontrado en el de Subastas que todas provienen del mismo país.

¿Cuales fueron las 10 apps más instaladas?



Obtenemos las apps que más instalaciones tuvieron. No nos es posible determinar su nombre o origen ya que están ocultos los datos pero observamos que de las 2 primeras hay una buena fuente de instalaciones.

¿Cuáles fueron las 10 en las cuales se instalaron más apps?



Definitivamente no nos es posible obtener la marca. Pero observamos el resultado y vemos que de la primera marca se produjeron muchas instalaciones con respecto a las siguientes.

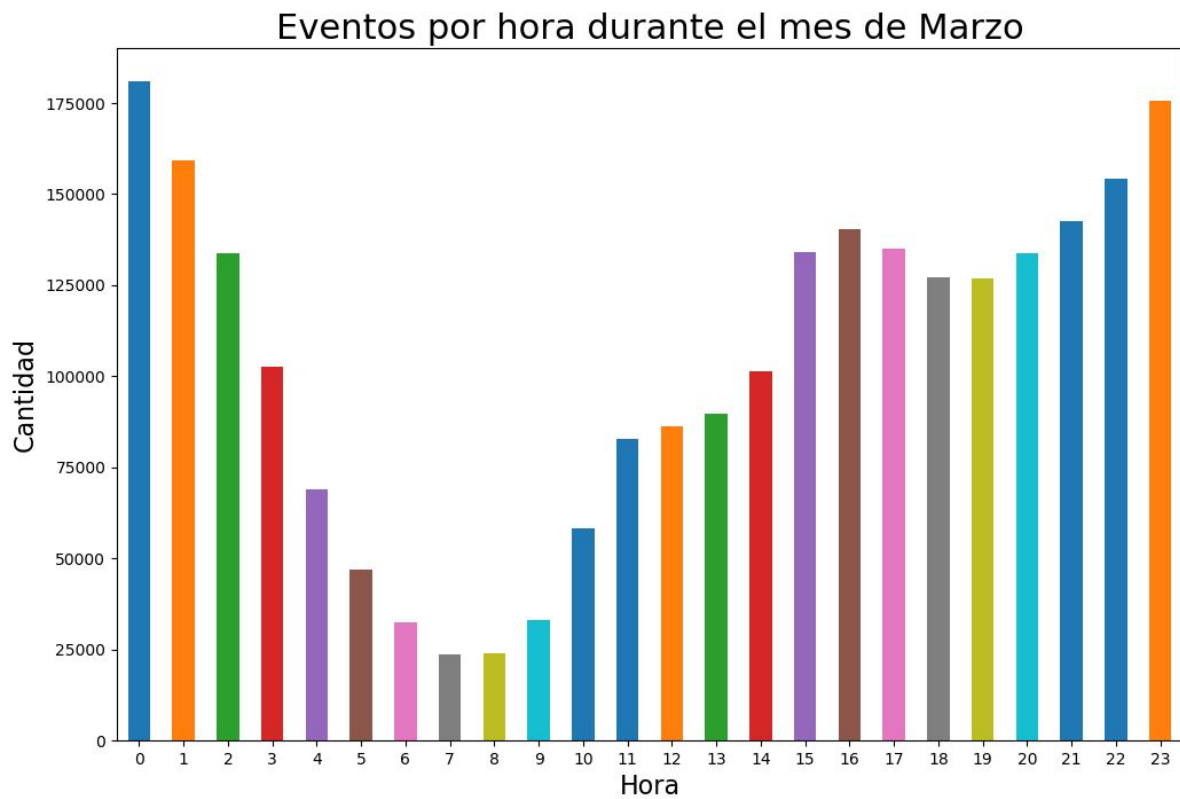
Análisis de Events

¿Cuántos eventos hubo en el periodo?

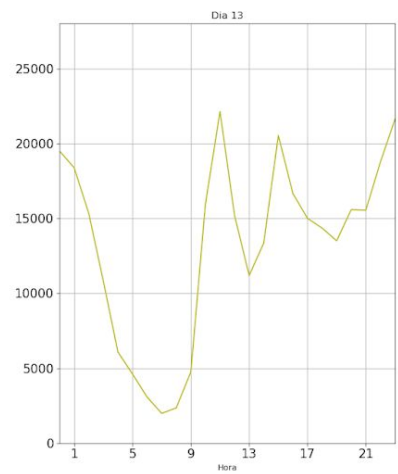
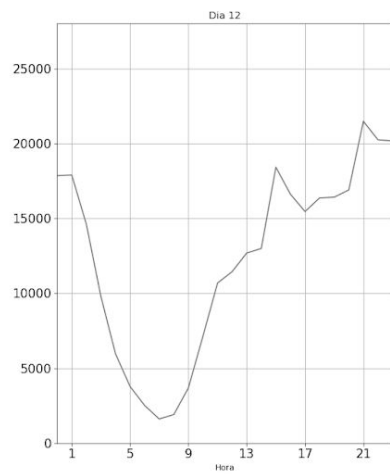
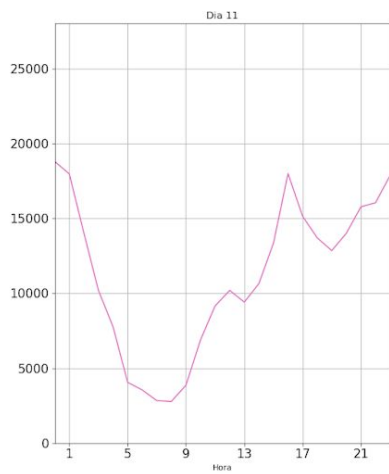
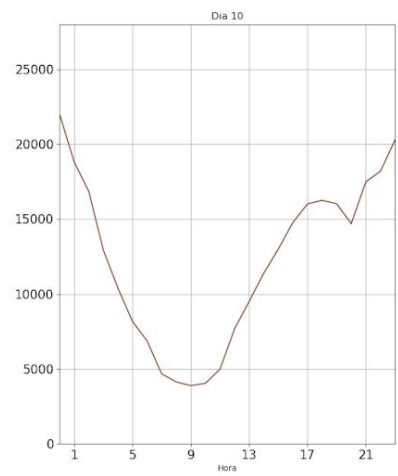
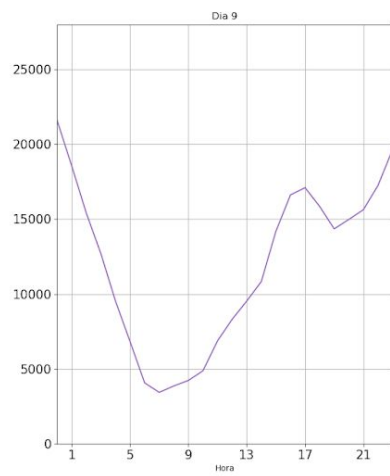
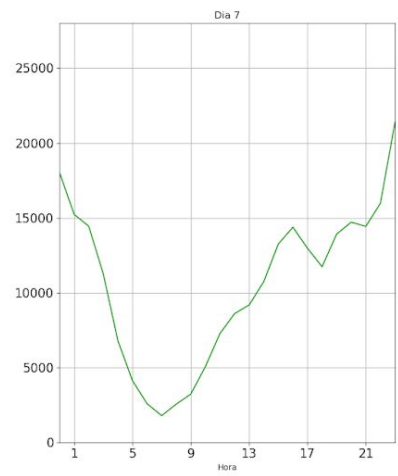
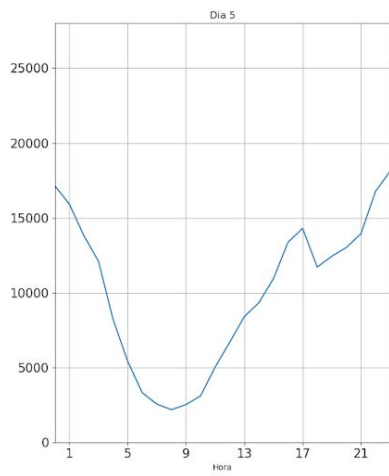


Aquí el análisis es análogo a las instalaciones, en el día trece es en el que más eventos se realiza pero no tiene mucho significado, incluso tiene menos significado por la naturaleza de lo que puede significar realizar un evento.

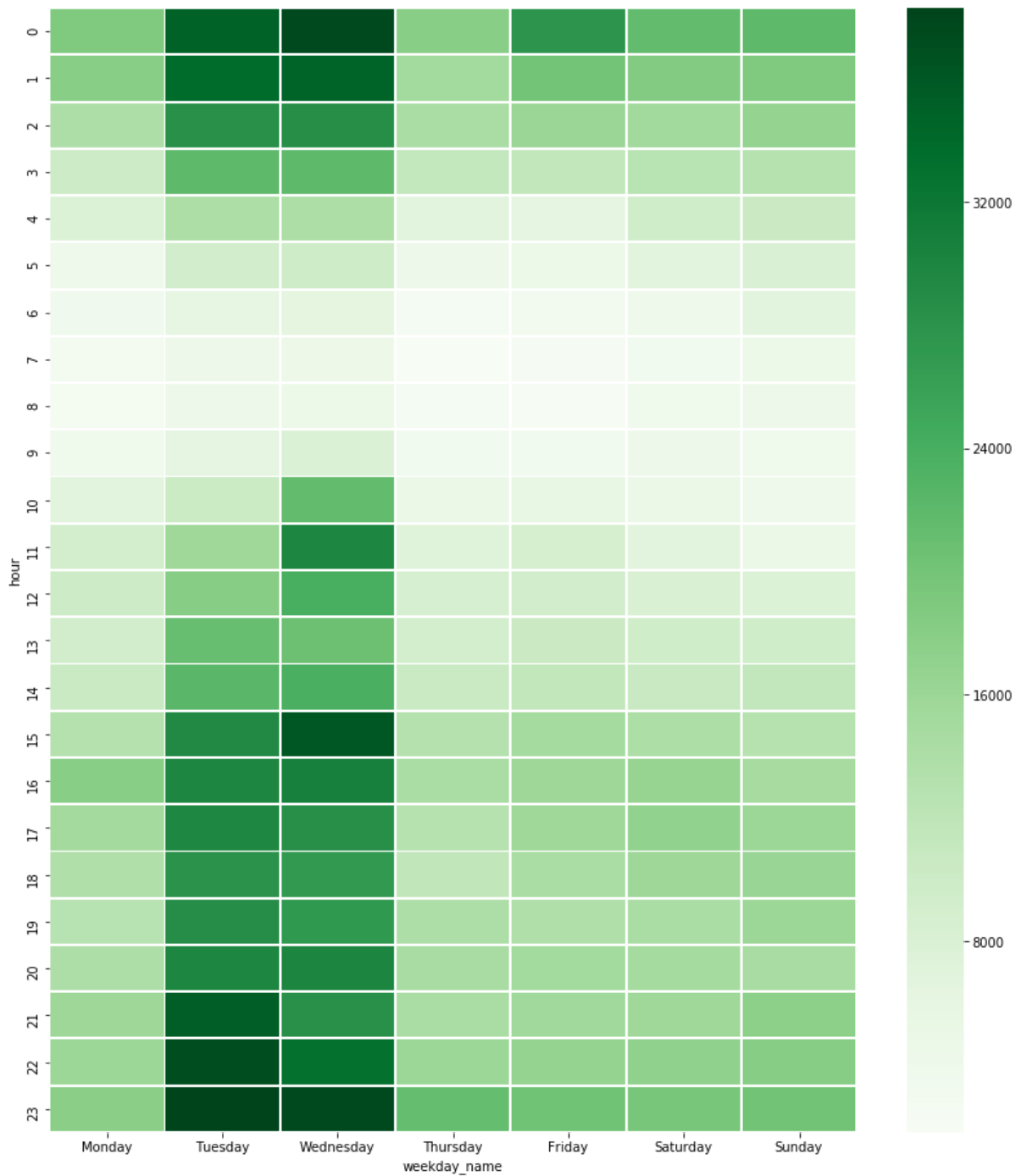
¿A qué hora se producen la mayor y la menor cantidad de eventos?



Análisis análogo a instalaciones por hora y por día, aquí a las 0 HS es cuando más eventos se realizan.

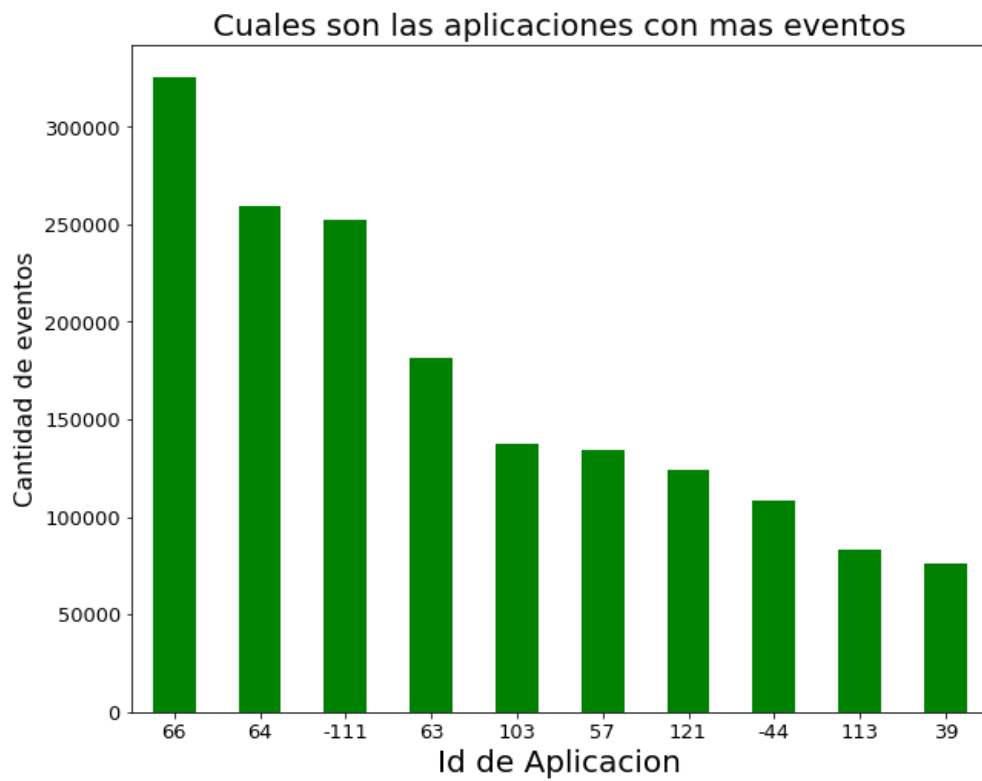


Heatmap por día de la semana y hora



Se observa una distribución mayor durante martes y miércoles.

¿Cuales son las apps que más eventos generaron?



Conclusión

A lo largo de este informe se realizaron visualizaciones y análisis partiendo de distintos puntos de vista con el fin de comprender, analizar y proponer mejoras al dominio en cuestión.

Tanto el modelo de negocio como los datos no resultaron amigables ni fácilmente entendibles a primera vista (algunos tampoco luego de recibir varias aclaraciones de la cátedra), por lo que requirieron también investigación previa al análisis de la información expuesta: algunos de sus resultados no son triviales.

El objetivo principal planteado al comienzo de la exploración de datos, más allá de analizar su información en sí, consistió en proponer alguna mejora de manera tal que Jampp pueda lograr una mejora en cuanto a las instalaciones atribuidas. Este objetivo hizo que el descubrimiento de la ausencia de registros con instalaciones atribuidas resultara un tanto chocante.

Fuera de la cantidad contada de inconvenientes, el grupo puede ahora afirmar que el entendimiento del set de datos es adecuado como para poder pasar a la segunda parte del tp.