

PROYECTO FINAL ANÁLISIS DE DATOS (Septiembre de 2021)

Jorge W. Alba, Juan C. Bolaños, Juan B. Dávila, Byron S. Huaraca

Resumen –

En el presente documento se detalla el proceso de recopilación, concentración, análisis y visualización de datos de varias temáticas propuestas con el objetivo de demostrar todo el conocimiento relacionado al análisis de datos adquiridos durante el semestre y de esta forma tener bases sólidas que puedan ser aplicadas en el mundo laboral.

I. DEFINICIÓN DE CASO DE ESTUDIO

Se realizará el estudio a datos actualizados sobre las temáticas de pulso político de Ecuador tanto por ciudades como por provincias, juegos online por países, ranking de felicidad por países y noticias o eventos mundiales. Los datos serán recopilados de diferentes fuentes y almacenados en diferentes bases de datos tanto SQL como NoSQL y posteriormente la información será analizada y visualizada para proporcionar conclusiones relevantes de la información previamente recopilada.

II. OBJETIVOS

• General

Recopilar datos de diferentes fuentes con el uso de herramientas y métodos para su respectivo análisis y visualización.

• Específicos

- Analizar los casos de estudio.
- Diseñar una arquitectura de solución.
- Establecer el cronograma de actividades.
- Escoger las herramientas necesarias.
- Recopilar datos de fuentes como Facebook, Twitter, TikTok, LinkedIn y con métodos como Webscraping.
- Concentrar los datos para realizar su respectivo análisis.
- Visualizar los datos con herramientas de visualización de datos.

- Documentar las conclusiones obtenidas del previo análisis de datos.

III. CRONOGRAMA DE ACTIVIDADES

El cronograma de actividades se lo visualiza en la Figura 1.

Cronograma de Actividades

Encargados	Tareas	Semana 1	Semana 2	Semana 3
• Jorge Alba • Juan Bolaños • Bernabé Dávila • Byron Huaraca	Organización del proyecto			
	Delegación de Actividades	2		
	Elección de herramientas	2		
	Estimación de tiempos	1		
• Jorge Alba • Juan Bolaños • Bernabé Dávila • Byron Huaraca	Recolección de datos			
	Búsqueda de datos		10	
	Generación de Scripts		15	
	Limpieza de datos		5	
• Jorge Alba • Juan Bolaños • Bernabé Dávila • Byron Huaraca	Concentración de datos			
	Importación de datos en las bases de datos		8	
	Concentración en elasticsearch		4	
	Análisis y visualización de datos			
• Jorge Alba • Juan Bolaños • Bernabé Dávila • Byron Huaraca	Importación de datos en herramientas de visualización			10
	Generación de gráficos			10
	Resultados			
	Interpretación de gráficos			10
• Jorge Alba • Juan Bolaños • Bernabé Dávila • Byron Huaraca	Documentación			15
	Total, de horas por semana	5	42	45
Total, de horas		92		

Figura 1 Cronograma

IV. RECURSOS Y HERRAMIENTAS

Las herramientas usadas para el proceso son:

- **Web:** Es la fuente mas grande de datos en la cual se pueden recopilar datos actualizados que son de gran ayuda para cualquier tipo de actividad y sobre todo para el análisis de datos.
- **Python:** Lenguaje de programación interpretado que ofrece varios métodos para el análisis de datos y relacionado a bases de datos.
- **SQL Server:** Es un sistema de gestión enfocado en bases de datos relacionales.
- **MySQL:** Sistema de gestión de bases relacionales de código abierto.

- **SQLite:** herramienta para la gestión de bases de datos relacionales contenida en una pequeña biblioteca que se encuentra escrita en C.
- **MongoDB:** sistema de gestión de bases de datos enfocada en bases de datos no relacionales que almacena en estructuras de datos BSON.
- **CouchDB:** Gestor para bases de datos no relacionales de código abierto enfocada en datos que asume la web de forma completa.
- **Elasticsearch:** Herramienta que ofrece un motor de búsqueda con interfaz RESTful y manejo de documentos JSON.
- **Logstash:** Herramienta que permite recolección, análisis y almacenamiento de logs para búsquedas.
- **Kibana:** Herramienta de visualización de datos que trabaja en conjunto con elasticsearch.
- **Tableau:** Software que permite la visualización de datos interactivos con enfoque en inteligencia empresarial.
- **Power BI:** Herramienta para visualizaciones interactivas y que ofrece capacidades de inteligencia empresarial con interfaz sencilla de manejar.
- **Twitter Developer Count:** Cuenta de Twitter que permite el acceso a partes del servicio ofrecido por medio de APIs que permite realizar la búsqueda de datos.
- **Kaggle:** Plataforma que ofrece recursos para Machine Learning y situaciones relacionadas a la ciencia de datos.

V. ARQUITECTURA DE SOLUCIÓN

La Arquitectura de la solución se la muestra en la Figura 2.

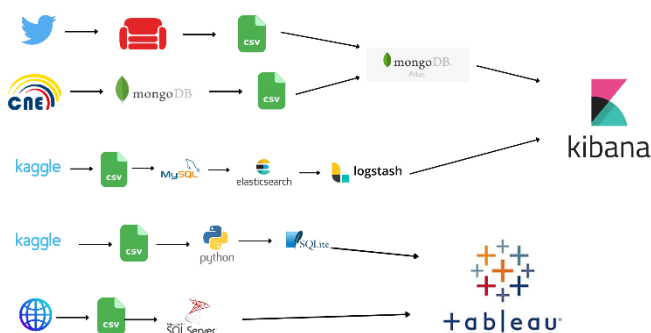


Figura 2 Arquitectura

VI. EXTRACCIÓN DE DATOS

Extracción de datos desde Twitter: La red social Twitter proporciona herramientas a los desarrolladores para facilitar la extracción de datos a través de su aplicación. Entre los principales usos que se le da a la API de Twitter se encuentran:

- Utilizar las potentes API de Twitter para ayudar a una empresa a escuchar, actuar y descubrir.

- Creación de herramientas para que las personas en Twitter integren o mejoren su experiencia en la plataforma.

- Utilizar la API de Twitter para obtener puntos de datos históricos y en tiempo real para su próximo proyecto de investigación [1].

Para poder acceder a esta cuenta de desarrollador, previamente se debe solicitar una cuenta de desarrollador a través de un formulario en el que detallas las intenciones con las que vas a usar su API. De esta forma se aseguran y evitan un mal uso de los datos que proporciona Twitter.

Cuando escribimos un tweet la información que se suele ver es el usuario que lo escribió. Pero la información que no se suele ver es mucho más potente como podemos ver en la Figura 3.



Figura 3 Anatomía de un Tweet

Extracción de datos de la Web: La web proporciona una información inmensa de cualquier temática por lo cual en algún sitio web que muestre información en tablas con la ayuda de Excel permite extraer los datos de dichas tablas y tener los datos disponibles para ser almacenados en una base de datos y posterior a ello ser analizados y presentados en gráficos. Por lo cual los datos sobre gamers por país y gamers de COD son fáciles de reconocer y entender las tendencias presentadas [2, 3].

Los datos fueron extraídos de Kaggle, en este se mostrará el informe mundial sobre la felicidad [4], la cual es una encuesta realizada por las Naciones Unidas que mide cómo ha evolucionado la felicidad de los ciudadanos en los últimos años en 156 países, al descargar de esta página se consigue un archivo .csv. Luego se procedió a utilizar código de Python para poder exportar los datos a SQLite y posterior a ello se descargó un complemento para utilizar la herramienta en la que se realizarán las visualizaciones que es Tableau [5].

Extracción de datos a través de Kaggle: Esta página web proporciona los datos que el usuario necesite, en este caso, se

lo hace en formato CSV para realizar la importación de dichos archivos a la base de datos, en este caso “MySQL”. De esta forma se facilita la creación de una tabla asignando los títulos de cada columna del CSV como los índices de la tabla [6].

VII. ANÁLISIS DE INFORMACIÓN

a. Pulso Político 20 Ciudades más importantes del Ecuador

Se recaudó un total de 1000 tweets aproximadamente, en los cuales podemos clasificar o filtrar con las herramientas la cantidad de menciones que se obtuvo de cada ciudad que analizamos y así darnos cuenta del volumen de interacciones que tuvo cada ciudad en el proceso de elecciones como se muestra en la Figura 4. También podemos filtrar la información de en cuantos de estos Tweets se mencionaba a Guillermo Lasso para darnos Cuenta de que la gran mayoría de los Tweets que recaudamos lo mencionan ya sea positiva o negativamente como se visualiza en la Figura 5.

b. Pulso en Provincias del Ecuador

A través de la página del Consejo Nacional Electoral se recaudó los datos de la primera Vuelta de las Elecciones presidenciales en donde se puede notar el porcentaje de ciudadanos que votaron por cada candidato. También se especifica que estos datos pueden tener un $\pm 3\%$ de error así que los datos no son del todo precisos. Aun así, estos datos nos sirven para analizar porque los 3 candidatos más solicitados fueron Andrés Arauz, Guillermo Lasso, y Yaku Pérez como podemos ver en la Figura 6 [7].

Para visualizar mejor esto lo dividiremos en regiones del Ecuador. Así por ejemplo tenemos los datos de la Costa Figura 7 y podemos ver que ampliamente hubo una ventaja de Andrés Arauz. Por otra parte, en las regiones de Sierra y Amazonía Figura 8 y Figura 9 podemos ver que Yaku Pérez se llevó una cantidad significativa de votos. Lo sorprendente es que Guillermo Lasso que finalmente fue el que ganó las elecciones no se mostró fuerte en ninguna región en particular. Aunque a estos datos debe incorporarse el voto de los extranjeros, que no se encuentra en este estudio.

c. Juegos Online por país

Los datos relacionados a los gamers por país muestran que Estados Unidos es el país con mayor índice de jugadores tanto en general. Por otro lado, existen varios países que cuentan con muy pocos gamers como por ejemplo Madagascar, Aruba o Bermuda, como se visualiza en la Figura 10.

En el análisis de los gamers del juego en específico COD, se puede observar que en Estados Unidos también se tiene el mayor índice de gamers de dicho juego, mientras que para este análisis también se tiene que este juego no es muy famoso en países como Bolivia, El Salvador o Omán, como se muestra en la Figura 11.

d. Ranking de felicidad

Al utilizar el código en Python para ingresar los datos que se tiene en csv y utilizando la herramienta DB Browser (SQLite) se observan todos los datos como se visualiza en la Figura 12. En la Figura 13 se muestra que tanto es el nivel de generosidad con la que se ha calificado en cada uno de los países que se encuentran en el Ranking siendo Indonesia con la gente más generosa y Grecia postula como la que es menos generosa. En la Figura 14 se nos muestra en cual de todas las regiones se encuentra un mayor score que se obtuvo en los índices de felicidad teniendo un mayor porcentaje en África y el menor en Norte América.

e. Eventos o noticias mundiales

El análisis de los datos correspondientes al número de casos de COVID-19 registrados dio como resultado que el país con mayores casos es Estados Unidos con un número cercano a los 30 millones de casos como se puede observar en la Figura 15. Por otro lado, al realizar un análisis de la cantidad de muertos registrada por COVID-19 se puede obtener que Estados Unidos sigue siendo el país con mayor registro de muertos, observando esta información en la Figura 16. Por tanto, el país más afectado durante la pandemia del COVID-19 es Estados Unidos.

El análisis realizado hacia los candidatos estadounidenses en las anteriores votaciones dio como resultado que el estado de Delaware fue el que mayores votos recaudó para ese año como se puede observar en la Figura 17.

El análisis realizado hacia el número de contagiados durante el Ébola en Europa en los años de 2014 – 2016 dio como resultados que Liberia fue el país con un número de casos posibles y casos detectados en nivel medio, apreciándose así en la Figura 18. A su vez se puede observar que este mismo número va en constante incremento hacia países bajos.

El análisis realizado a las olimpiadas de Tokio 2021 posicionó a Estados Unidos como el país con mayor número de medallas de oro obtenidas en las olimpiadas, con un total de 39 medallas, visualizándose así en la Figura 19.

VIII. VISUALIZACIÓN DE INFORMACIÓN

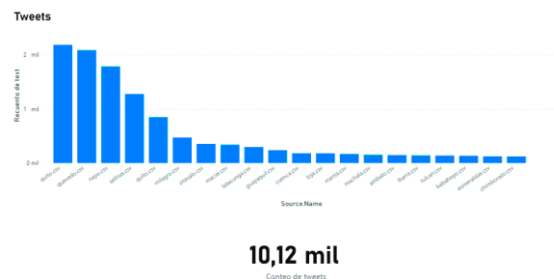


Figura 4 Tweets por ciudades



Gráfico de área que muestra la presencia de Andrés Arauz, Guillermo Lasso y Yaku Pérez por provincias. El eje Y representa la presencia (0.0 a 0.5). El eje X muestra las provincias: Sucumbios, Orellana, Zamora Chinchipe, Morona Santiago, Pastaza y Napo. La leyenda indica: Andrés Arauz (azul), Guillermo Lasso (naranja) y Yaku Pérez (rojo).

Provincia	Andrés Arauz	Guillermo Lasso	Yaku Pérez
Sucumbios	0.35	0.00	0.00
Orellana	0.25	0.00	0.00
Zamora Chinchipe	0.15	0.00	0.00
Morona Santiago	0.15	0.00	0.00
Pastaza	0.15	0.00	0.00
Napo	0.15	0.00	0.00

[illegible]

Gráfico de área que muestra la preferencia por los candidatos Andrés Arauz, Guillermo Lasso y Yaku Pérez por provincias. El eje Y representa la preferencia (0.0 a 0.5). El eje X muestra las provincias: Manabí, Los Ríos, Esmeraldas, Santa Elena, Guayas, El Oro, Santo Domingo de Guacaya y Guayas. La leyenda indica: Andrés Arauz (azul), Guillermo Lasso (naranja) y Yaku Pérez (verde).

Provincia	Andrés Arauz	Guillermo Lasso	Yaku Pérez
Manabí	0.40	0.10	0.50
Los Ríos	0.35	0.15	0.50
Esmeraldas	0.30	0.20	0.50
Santa Elena	0.25	0.15	0.60
Guayas	0.20	0.25	0.55
El Oro	0.15	0.30	0.55
Santo Domingo de Guacaya	0.15	0.35	0.50
Guayas	0.10	0.40	0.50

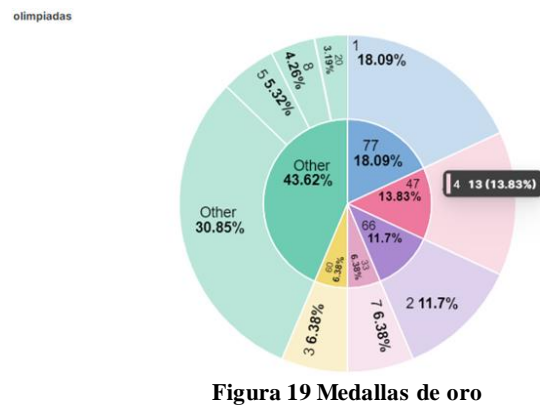
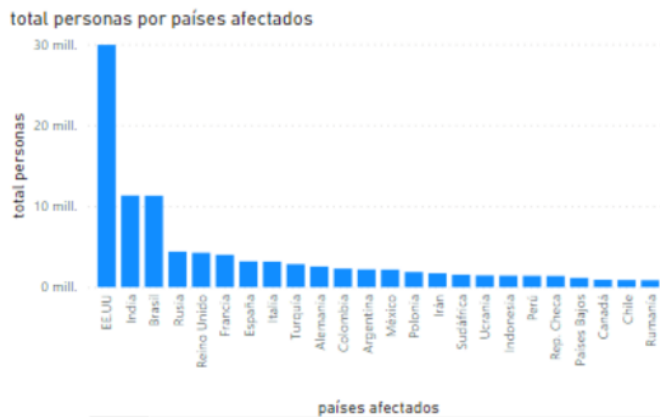
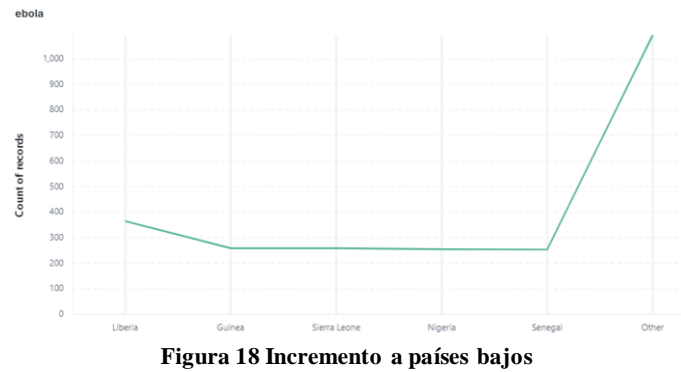
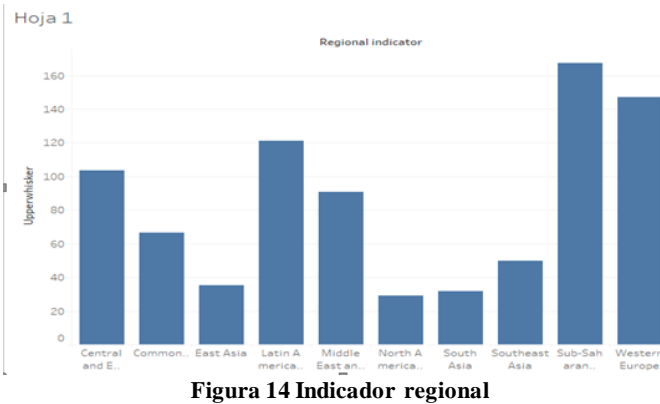
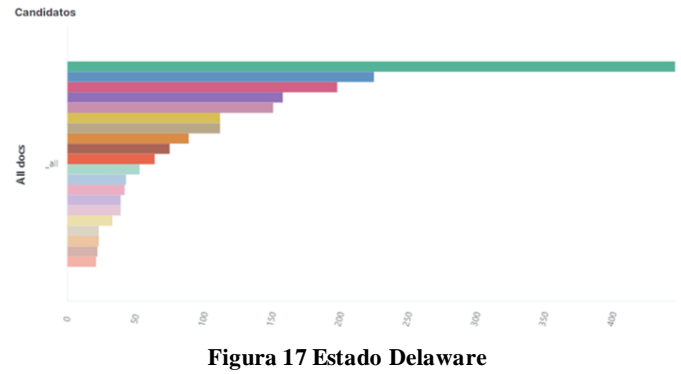
Gráfico de área que muestra la actividad de Andrés Arauz, Guillermo Lasso y Yaku Pérez en las redes sociales (Twitter, Facebook, YouTube) por provincia. El eje Y representa la actividad (0.1 a 0.5). El eje X muestra las provincias: Imbabura, Cacha, Pichincha, Azuay, Cotacachi, Loja, Cotacachi, Chimborazo, Bolívar y Tungurahua. La leyenda indica: Andrés Arauz (azul), Guillermo Lasso (naranja) y Yaku Pérez (rojo).

Country	Games
South Africa	1
Spain	2
Sri Lanka	1
Sweden	3
Switzerland	1
Syrian Arab Republic	1
Taiwan, Republic of China	1
Thailand	1
Trinidad and Tobago	1
Tunisia	1
Turkey	1
Turkmenistan	1
Ukraine	1
United Arab Emirates	4
United Kingdom	4
United States	21
United States Minor Outlying Islands	1
Uruguay	1
Uzbekistan	1
Venezuela	1
Viet Nam	1
Yemen	1
Zambia	1

País	Researcher
United States	32
Argentina	27
Japan	25
Mexico	25
Saudi Arabia	25
Australia	19
Netherlands	17
Belgium	12
Chile	12
Emirates	11
Colombia	9
Honduras	8
Poland	8
China	7
France	7
Hong Kong	7
Portugal	7
Russia	7
Ireland Kingdom	6
Spain	6
New Zealand	5
Slovenia	5
Turkey	5
Brazil	5
Cyprus	5
Greece	5
Sweden	5
Iceland	5
South Africa	5
Austria	4
Canada	4
Denmark	4
Germany	4
Italy	4
Luxembourg	4
Norway	4
Peru	4
Switzerland	4
Costa Rica	3
Czech Republic	3
Ecuador	2
Kuwait	2

Country name		Regional indicator	Ladder score	Standard error of ladder score	upperwhisker	lowerwhisker	Logged GDP
Filter	Filter	Filter	Filter	Filter	Filter	Filter	Filter
83	Congo (Brazzaville)	Sub-Saharan Africa	5.342	0.097	5.533	5.151	
84	China	East Asia	5.339	0.029	5.397	5.281	
85	Ivory Coast	Sub-Saharan Africa	5.306	0.078	5.46	5.152	
86	Armenia	Commonwealth of Independent States	5.283	0.058	5.397	5.168	
87	Nepal	South Asia	5.269	0.07	5.406	5.132	
88	Bulgaria	Central and Eastern Europe	5.266	0.054	5.371	5.16	
89	Maldives	South Asia	5.198	0.072	5.339	5.057	
90	Azerbaijan	Commonwealth of Independent States	5.171	0.04	5.25	5.091	
91	Cameroon	Sub-Saharan Africa	5.142	0.074	5.288	4.996	
92	Senegal	Sub-Saharan Africa	5.132	0.068	5.266	4.998	
93	Albania	Central and Eastern Europe	5.117	0.059	5.234	5.001	
94	North Macedonia	Central and Eastern Europe	5.101	0.051	5.202	5.001	
95	Ghana	Sub-Saharan Africa	5.088	0.067	5.219	4.958	
96	Niger	Sub-Saharan Africa	5.074	0.102	5.273	4.875	
97	Turkmenistan	Commonwealth of Independent States	5.066	0.036	5.136	4.996	
98	Gambia	Sub-Saharan Africa	5.051	0.089	5.225	4.877	
99	Benin	Sub-Saharan Africa	5.045	0.073	5.189	4.901	
100	South Asia	South Asia	5.039	0.046	5.110	4.961	

Figura 12 SQLite

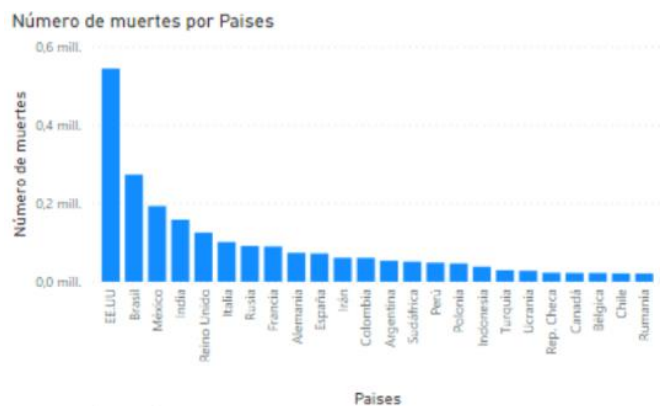


IX. RESULTADOS

Los datos que se recolectaron nos muestran el fuerte impacto que el candidato Guillermo Lasso tuvo en la opinión de la gente en las redes sociales. Lo que justifica su elección como presidente del Ecuador. También como podía esperarse la mayor cantidad de interacción vino de la ciudad de Quito.

A través de la información que se obtuvo del CNE por cada provincia del Ecuador pudimos realizar el análisis y determinar qué candidato tuvo ventaja en cada región del Ecuador y el resultado es curioso ya que en la primera vuelta electoral Guillermo Lasso no fue el favorito en ninguna de las regiones del Ecuador.

Los datos se encuentran en muchos lugares y para conocer el número de jugadores por países se pueden encontrar varias fuentes las cuales se van actualizando constantemente, como se mostró en la sección de visualizaciones los resultados de los jugadores arrojó que en Estados Unidos se encuentra la mayoría



de estos esto se debe a que en ese país existe un avance tecnológico notable, y el estilo de vida es muy bueno puesto que su economía es estable. Al observar el otro lado de la situación existen países donde no hay gran incidencia de jugadores esto se puede deber a que el nivel de vida de sus habitantes no permite que muchas familias tengan en su hogar consolas de videojuegos o incluso internet.

Al obtener los datos del informe mundial de felicidad se investigó que estos fueron obtenidos de la encuesta mundial de Gallup el cual se pidió a todos los encuestados que puntúen aspectos de su vida dando un valor de 0 a 10, siendo cero la peor vida posible, además, de datos como la generosidad, la corrupción existente en el país, el apoyo social, la dystopia, la libertad entre otros. Así se pudo construir un score de acuerdo para puntuar a todas las ciudades del mundo, además, que en herramientas de visualización como Tableau se puede obtener un mejor sondeo de estos resultados.

Los datos analizados sobre los casos y muertes registradas posicionan a Estados Unidos como el país más afectado, más aún esto supone que aún siendo un país de primer mundo, no se encuentra libre de ser uno de los que más estragos sufrió. Además, al observar los análisis de las olimpiadas de Tokio 2021 Estados Unidos destaca, por lo que se puede concluir que, aun siendo el país más afectado durante la pandemia, no deja de ser uno de los que más destaca en actividades deportivas.

X. CONCLUSIONES Y RECOMENDACIONES

La Api de Twitter tiene múltiples ventajas para extraer información y utilizando las herramientas que se nos enseñó como Kibana o Power Bi podemos filtrar la información que más nos sirva para realizar un análisis en profundidad.

La web ofrece un sinnúmero de información la cual puede ser recopilada fácilmente con la ayuda de herramientas como Excel y posteriormente ser tratada y analizada para diferentes propósitos.

Cuando se necesita recolectar información de algún ámbito ya sea político, económico o de cualquier índole paginas como Kaggle nos proporcionan datos específicos, auténticos y utilizados para realizar cualquier tipo de mediciones.

El realizar el análisis de información sobre un determinado campo o tema nos permite recolectar datos precisos, esto facilita la toma de decisiones al momento de implementar o desarrollar algún proyecto.

XI. DESAFÍOS Y PROBLEMAS

Para recopilación de datos desde la web con la ayuda de Excel es necesario realizar búsquedas en sitios web en donde se presente la información en tablas puesto que si la información requerida no se encuentra de esta forma su recopilación es muy complicada si se desea usar la herramienta de Excel.

La herramienta de Tableau aunque intuitiva con la mayoría de conexiones a una determinada base de datos, en lo que se tomo en algo complicado fue conectarlo a una base de datos SQLite puesto que era necesario descargar un complemento necesario en la pagina oficial para poder conectarlo a dicha base.

Para recaudar la información desde Twitter hubo la condición de no extraer información demasiado rápido y seguido ya que la API de twitter consideraba que no estaba permitido a realizar esta acción de forma repetida y constante.

XII. REFERENCIAS

- [1] «Anatomía de un Tweet,» DIGITALTROUPE, 17 marzo 2019. [En línea]. Available: <http://www.digitaltroupe.com/anatomia-de-un-tweet/>. [Último acceso: 11 septiembre 2021].
- [2] «Call of Duty (CoD) - statistics & facts,» Statista, 13 Agosto 2021. [En línea]. Available: <https://www.statista.com/topics/8300/call-of-duty-cod/>. [Último acceso: 11 Septiembre 2021].
- [3] «Highest Earnings By Country,» Sportsearnings, 2021. [En línea]. Available: <https://www.esportsearnings.com/countries>. [Último acceso: 11 Septiembre 2021].
- [4] «World Happiness Report 2021,» Kaggle, [En línea]. Available: <https://www.kaggle.com/ajaypalsinghlo/world-happiness-report-2021>. [Último acceso: 11 septiembre 2021].
- [5] «SQLite ODBC Driver,» ch-werner, 20 junio 2020. [En línea]. Available: <http://www.ch-werner.de/sqliteodbc/>. [Último acceso: 11 septiembre 2021].
- [6] «Datasets,» Kaggle, [En línea]. Available: <https://www.kaggle.com/datasets>. [Último acceso: 11 septiembre 2021].
- [7] «Consejo Nacional Electoral,» Gobierno del Ecuador, 14 febrero 2021. [En línea]. Available: <https://www.gob.ec/>. [Último acceso: 11 septiembre 2021].