

Big Query

Anna Salazar

August 10, 2022

Índex

1	Què és BigQuery?	1
1.1	Per què hauríem d'utilitzar BigQuery en lloc d'altres eines?	1
2	Creació i treball amb conjunts de dades i taules	2
2.1	Configuració de la Plataforma de Google Cloud (GCP)	2
2.2	Creació d'un conjunt de dades	3
2.2.1	Limitacions	3
2.3	Definició d'una taula de BigQuery des de la interfície d'usuari	4
2.4	Consultar una taula de BigQuery senzilla	4
2.5	Càrrega de dades per crear una taula de BigQuery	4
2.6	Consulta de dades i visualització d'estadístiques de consultes	4
2.7	Creació d'una taula a partir d'un resultat de consulta	4
2.8	Creació d'una consulta a partir d'un filtre	4
3	Execució de consultes i visualització de resultats	5
3.1	Conjunts de dades públics a BigQuery	5
3.2	Configuració i ús de memòria cau de BigQuery	5
3.3	Taules externes de BigQuery	5
3.4	Integració de BigQuery amb Data Studio	5

1 Què és BigQuery?

BigQuery és un motor d'anàlisi de macrodades (Big Data) que permet executar consultes SQL al núvol sobre les dades emmagatzemades en aquest, sense importar el volum de les dades ni el tipus de consultes que es volen fer. El motor de consulta és capaç de treballar sobre terabytes de dades en qüestió de segons, i sobre petabytes en pocs minuts. Avui en dia, les empreses estan adoptant cada cop més la presa de decisions basades en dades i fomentant una cultura oberta en la qual les dades no estan aïllades dins dels departaments. BigQuery, en proporcionar els mitjans tecnològics per a promoure un canvi cultural cap a l'agilitat i l'obertura, realitza un paper molt important en l'augment del ritme de la innovació.

Treballar amb dades a BigQuery implica 3 aspectes principals: l'emmagatzemament, la incorporació de les dades i la consulta d'aquestes, Google s'encarrega de tota la resta. Com BigQuery és un servei totalment gestionat, no és necessari configurar ni instal·lar res en el nostre ordinador i, pel mateix motiu, no necessitem un administrador de la base de dades. Simplement, podem entrar en el nostre projecte de Google Cloud des del nostre navegador i començar a analitzar.

Pel que fa a l'emmagatzemament, les dades es guarden en una taula estructurada, la qual cosa significa que es pot utilitzar l'SQL estàndard per a facilitar la consulta i l'anàlisi de dades. BigQuery és perfecta pel Big Data perquè gestiona tot aquest emmagatzemament i està proveïda d'operacions d'escalabilitat que funcionen de forma automàtica sense que l'usuari s'hagi d'involucrar, per la qual cosa mai haurem de preocupar-nos per la grandària de les dades amb els quals treballem. Part de la consideració de disseny darrere de BigQuery és animar als usuaris a centrar-se en els coneixements en lloc de la infraestructura. Quan s'introdueixen les dades a BigQuery no és necessari pensar en els diferents tipus d'emmagatzemament, ni en els seus avantatges pel que fa a velocitat i cost; l'emmagatzemament està totalment gestionat.

Per a més informació sobre BigQuery, es pot consultar la pàgina de [Google Cloud](#).

1.1 Per què hauríem d'utilitzar BigQuery en lloc d'altres eines?

Una de les característiques més rellevants que presenta BigQuery és que es tracta d'una plataforma sense servidor, és a dir, que els servidors s'executen en segon pla, sense l'intervenció de l'usuari. A més, presenta una alta disponibilitat, la qual cosa es tradueix en que no cal preocupar-se per la caiguda dels servidors, ja que el servei s'encarrega d'això. Per últim, BigQuery també té propietats d'escalabilitat automàtica que fan possible gestionar fins a petabytes de dades. Aquestes característiques no estan disponibles a la majoria de plataformes d'emmagatzament de dades tradicionals, cosa que fa destacar BigQuery entre moltes.

Com en molts altres magatzems de dades, BigQuery és capaç de treballar amb moltes fonts de dades diferents. Pot extreure dades del seu propi sistema d'arxius, de Google Cloud Storage i de moltes fonts més. Després de fer-ho, es poden consultar aquestes dades utilitzant SQL estàndard o SQL heretat, el rendiment en qualsevol cas és excel·lent. Els resultats de les consultes solen emmagatzemar-se en la memòria cau durant 24 hores, de manera que les següents execucions d'aquesta consulta només hauran d'obtenir les dades de la cau en lloc de fer-ho del disc.

2 Creació i treball amb conjunts de dades i taules

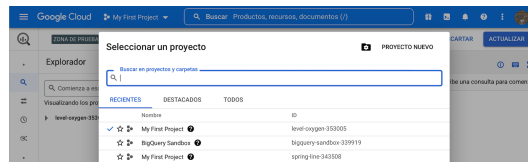
2.1 Configuració de la Plataforma de Google Cloud (GCP)

Per utilitzar aquesta eina d'anàlisi només ens caldrà crear un compte a Google Cloud i treballar a la zona de proves que ofereix Google per treballar de forma gratuïta. Per fer servir la zona de proves (Sandbox) seguirem els passos següents:

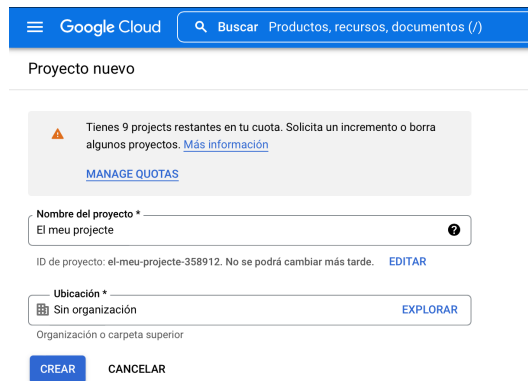
1. En primer lloc, ens dirigim a la interfície d'usuari de [BigQuery](#). Des d'aquesta interfície es poden realitzar la majoria de les operacions.

2. Accedeix al teu compte de Google o crea un nou compte si encara no en tens cap. Si és el primer cop que iniciis sessió a Google Cloud, hauràs de marcar el país i acceptar les condicions de servei.

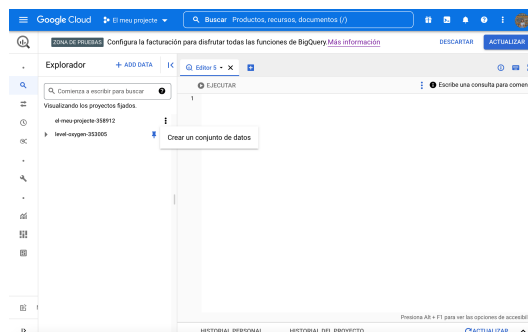
3. Un cop dins, podem veure com és l'espai de treball SQL. Hi ha una secció de l'Explorador a l'esquerra que ens permet navegar en projectes, conjunts de dades i taules. Per tal de fer servir la zona de proves, haurem de crear un projecte.



Introdueix un nom al teu projecte i fes clic a Create. En el nostre cas, hem anomenat el projecte *El meu projecte*, i treballarem sobre aquest per il·lustrar el funcionament de la plataforma.



4. Un cop creat el projecte, seràs redirigit a la interfície web de BigQuery.
5. Ara ja pots carregar o consultar dades en el teu projecte sense cap compte de facturació adjunta.



2.1.1 Limitacions

Per a l'ús de la zona de proves gratuïta que ofereix Google, haurem de tenir en compte un seguit de limitacions.

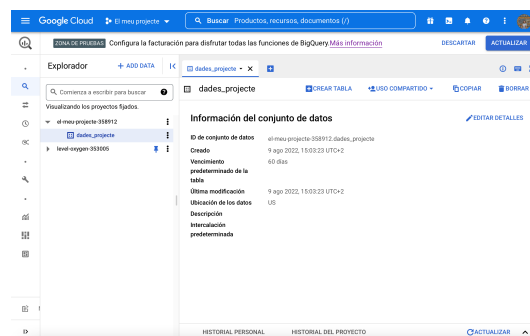
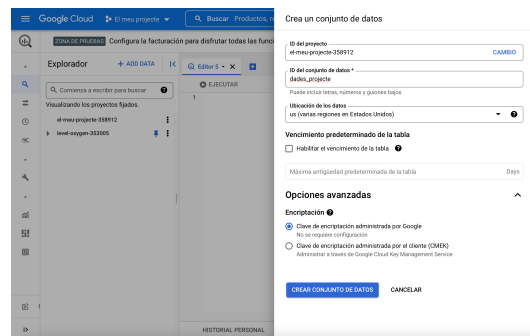
En primer lloc, ens trobem amb un màxim de 10 Gb d'emmagatzemament i 10 Tb de consulta al mes. Al llarg d'aquest projecte no utilitzarem un volum de dades més gran ni sobrepassarem el límit d'espai de consulta, però s'han de tenir en compte aquestes limitacions si l'objectiu és treballar amb el format gratuït.

A més, ens trobem que tots els conjunts de dades tenen el temps de caducitat de la taula per defecte establerta en 60 dies. Per tant, totes les taules, vistes o particions de les taules caducaran automàticament passats els 60 dies.

Una altra característica destacable és que els projectes de la zona de proves no són compatibles amb:

- La transmissió de dades
- Sentència de llenguatge de manipulació de dades (DML)
- Servei de transferència de dades de BigQuery

2.2 Creació d'un conjunt de dades



- 2.3 Definició d'una taula de BigQuery des de la interfície d'usuari
- 2.4 Consultar una taula de BigQuery senzilla
- 2.5 Càrrega de dades per crear una taula de BigQuery
- 2.6 Consulta de dades i visualització d'estadístiques de consultes
- 2.7 Creació d'una taula a partir d'un resultat de consulta
- 2.8 Creació d'una consulta a partir d'un filtre

3 Execució de consultes i visualització de resultats

3.1 Conjunts de dades públics a BigQuery

3.2 Configuració i ús de memòria cau de BigQuery

3.3 Taules externes de BigQuery

3.4 Integració de BigQuery amb Data Studio