

Práctica OpenRefine

Big Data, Máster universitario de Tecnologías del sector financiero

uc3m | Universidad **Carlos III** de Madrid

Andoni Alcelay Izarzugaza
18-10-2018

Examinar los datos con Facet/Text Facet. Por ejemplo, en Vecindad o Lugar de Muerte. Identifica posibles errores

Se empezará con examinar los datos de la Vecindad, clicando en el triangulo desplegable de esa columna y clicando en Facet->Text facet. Como se puede observar en la Ilustración 2, se agrupa toda la información y sus repeticiones en un view en la zona izquierda de pantalla.



Ilustración 2: Facet de la columna Vecindad

Estos datos podrían no ser del todo fiables por lo que se hará un análisis cluster para detectar posibles errores. Con este análisis se resolverán pequeñas diferencias entre nombres de pueblos y ciudades como se puede observar en la Ilustración 3. Aun así hay casos de pueblos como Palencia y Valencia que tienen un parecido muy grande pero que se tratan de pueblos diferentes, de ahí que se de la opción a seleccionar el cluster.

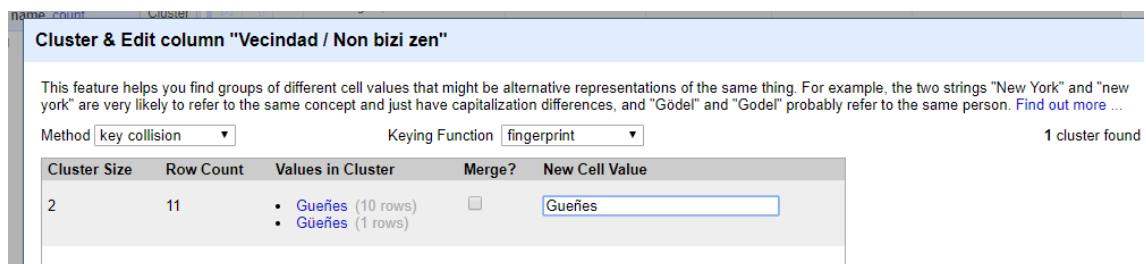


Ilustración 3: Ejemplo de cluster

Después de analizarlo con el cluster de palabras parecidas, se han encontrado parecidos entre más pueblos a los que le faltaban guiones o -a al final, algo muy común en pueblos de Euskadi como son Legazpi y Legazpia, que se tratan del mismo pueblo, uno redactado en Euskera y el otro en castellano.

A continuación se procede a hacer lo mismo con el lugar de su muerte, ya que pueden haber más diferencias.

Cluster en la columna de personas y lugares

Para hacer el cluster sobre personas se ha clicado sobre Edit cells -> Edit and cluster y nos aparece una tabla como la que se muestra en la ilustración 4.

Cluster & Edit column "Nombre y Apellidos / Izen-abizenak"

This feature helps you find groups of different cell values that might be alternative representations of the same thing. For example, the two strings "New York" and "new york" are very likely to refer to the same concept and just have capitalization differences, and "Godel" and "Godel" probably refer to the same person. [Find out more ...](#)

Method: **key collision** Keying Function: **fingerprint** 21 clusters found

Cluster Size	Row Count	Values in Cluster	Merge?	New Cell Value
2	2	Rodriguez Ruiz, Timoteo (1 rows) Ruiz Rodriguez, Timoteo (1 rows)	<input type="checkbox"/>	Rodriguez Ruiz, Timoteo
2	2	Fernandez Perez, Andres (1 rows) Perez Fernandez, Andres (1 rows)	<input type="checkbox"/>	Fernandez Perez, Andres
2	2	Fernandez Rodriguez, Jose (1 rows) Rodriguez Fernandez, Jose (1 rows)	<input type="checkbox"/>	Fernandez Rodriguez, Jose
2	2	Mariño Perez, Juan (1 rows) Perez Mariño, Juan (1 rows)	<input type="checkbox"/>	Mariño Perez, Juan
2	2	Arambarri Gallastegui, Fernando (1 rows) Gallastegui Arambarri, Fernando (1 rows)	<input type="checkbox"/>	Arambarri Gallastegui, Fernando
2	2	Blazquez Ramiro, Heliodoro (1 rows) Ramiro Blazquez, Heliodoro (1 rows)	<input type="checkbox"/>	Blazquez Ramiro, Heliodoro
2	2	Lozano Pastor, Santiago (1 rows) Pastor Lozano, Santiago (1 rows)	<input type="checkbox"/>	Lozano Pastor, Santiago

Buttons: Select All, Unselect All, Export Clusters, Merge Selected & Re-Cluster, Merge Selected & Close, Close

Ilustración 4: Ejemplo de cluster sobre la columna personas

En esta ventana se pueden ver los nombres y apellidos de personas que pueden coincidir entre orden o apellidos parecidos. Esto hay que tenerlo en cuenta ya que se ha podido meter al revés y realmente se trata de una misma persona. Ya que no se puede tener el criterio de nombre y apellidos para saber si se trata de una misma persona, en cada fila hay una opción "Browse this cluster" que dará una información más detallada sobre la persona sobre la que se hará el cluster. Al clicar en la primera fila, se obtiene la información mostrada en la ilustración 5.

☆	7786	Rodriguez Ruiz, Timoteo	Sestao	Lekeitio	1936-10-20T00:00:00Z	Muerto frente / Frontean hila
☆	7996	Ruiz Rodriguez, Timoteo	Sestao	Lekeitio	1936-10-19T00:00:00Z	Muerto frente / Frontean hila

Ilustración 5: Información detallada sobre el primer dato del cluster

Al analizar estos datos, se puede observar que definitivamente sí se trata de esa misma persona, por lo que se le pondrá el nombre en común a los dos.

Hay casos en los que, aunque la persona sea la misma, aparentemente, hay datos ausentes, como el caso de Fernando Arambarri, al que le faltan los datos como la Vecindad, aun así, ya que todos los demás campos de las filas coinciden, se tomarán como una misma persona.

Buscar errores ortográficos y registros duplicados

Para eliminar errores ortográficos y registros duplicados se ordena por nombre de las personas de la tabla y se empiezan a clusterizar para evitar duplicados, todo esto se ha realizado en el punto anterior [Cluster en la columna de personas y lugares](#)

Reconciliar la columna lugares con Wikidata

Si se quiere reconciliar una columna con wikidata, se clicla en el desplegable de esa columna-> Reconcile->Start reconciling...

Una vez clicado ahí se abrirá la una ventana como la que se muestra en la Ilustración 6.

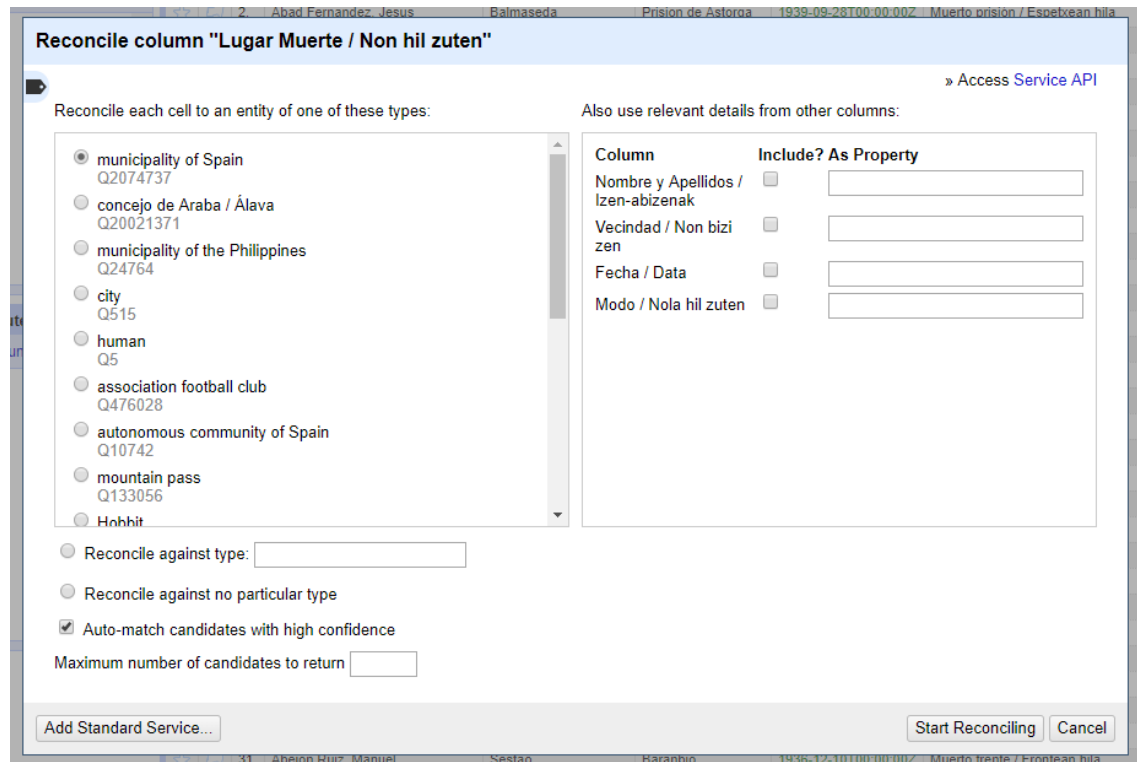


Ilustración 6: Reconcilio de datos con Wikidata

En este caso reconciliaremos la información por la columna municipios de España. El resultado es el mostrado en la Ilustración 7.

All	Nombre y Apellidos / Izen-abizenak	Vecindad / Non bizi	Lugar Muerte / Non hil zuten	Fecha / Data	Modo / Nola hil zuten
1.	Abad Angulo, Florentino				Muerto frente / Frontean hila
2.	Abad Fernandez, Jesus	Balmaseda	Prision de Astorga	1939-09-28T00:00:00Z	Muerto prisión / Espetxean hila
3.	Abad Huerta, Victor	Burgos	Getxo	1937-04-08T00:00:00Z	Muerto frente / Frontean hila
4.	Abad Ruiz, Luis	Bilbo	Larrauri	1937-05-14T00:00:00Z	Muerto frente / Frontean hila
5.	Abad Torre, Francisco				Muerto frente / Frontean hila
6.	Abad Torres, Mariano	Santurtzi	Legutio	1936-12-02T00:00:00Z	Muerto frente / Frontean hila
7.	Abadia Anaut, Vicente	Isaba			Muerto frente / Frontean hila
8.	Abaitua Arizmendi, Antonio				Muerto frente / Frontean hila
9.	Abaitua Arizmendi, Severiano				Muerto frente / Frontean hila
10.	Abaitua Perez, Pedro Jose Luis	Gasteiz	Azazeta	1937-04-01T00:00:00Z	Fusilado / Fusilatua

Ilustración 7: Resultado de reconciliación de pueblos de España

Como se puede observar, 1284 entradas de han quedado en blanco y 5549 se han matcheado. El resto, 2767 no se han podido reconciliar debido a que algunos de ellos, como por ejemplo “Prisión de Astorga”, que no corresponde a un municipio, sino que, a un establecimiento, en este caso una prisión. Podría especificarse que esa prisión se ubica en Astorga, pero sería pérdida de información. Por otro lado, hay otros que no se informan en Wikidata por lo que no se han podido matchear, como por ejemplo Larrauri o Azazeta.

Enriquecer el conjunto de datos a partir de Wikidata

Ya que hay datos en los lugares de nacimiento que, como se ha visto en el punto anterior, [Reconciliar la columna lugares con Wikidata](#), se han reconciliado. A partir de estos datos se puede sacar otra información y así enriquecer la información actual. Para ello se clic en el desplegable de Lugar muerte->Edit Column->Add columns from reconciled values. Se enriquecerán los datos de sus coordenadas y su área, como se muestra en la Ilustración 8.

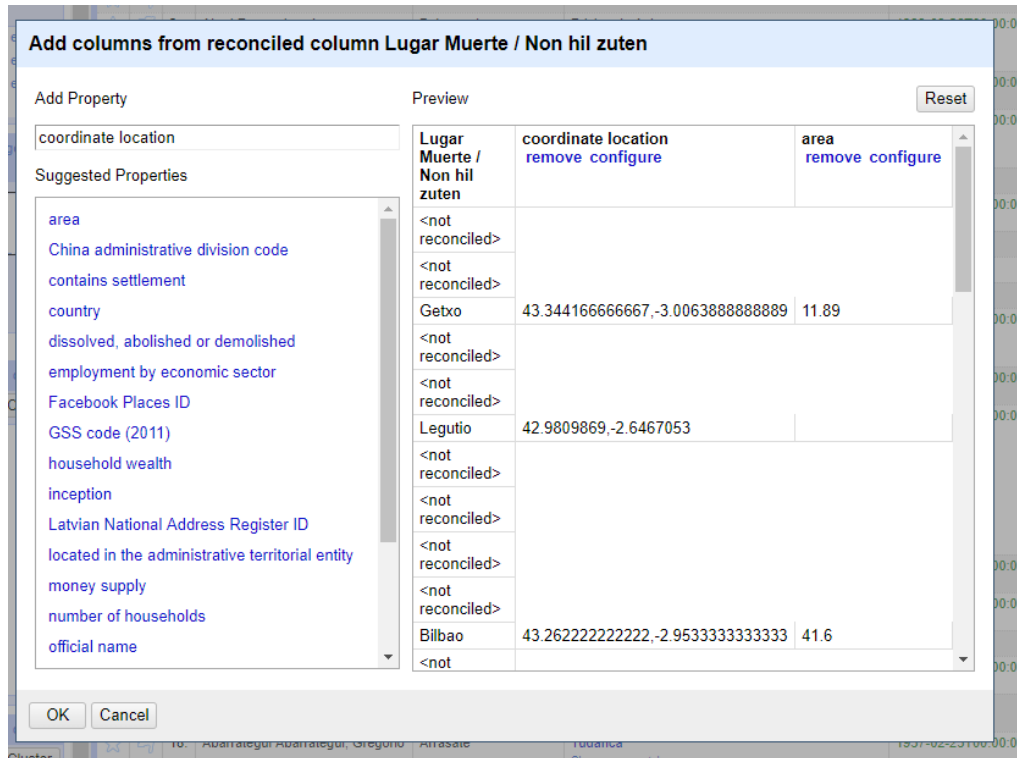


Ilustración 8: Añadir columnas de coordenadas y área

Una vez terminado el proceso, se puede observar como a la tabla actual se le han añadido dos columnas nuevas, que son las coordenadas de los municipios y su área. Esta información se ha recibido solo para las filas en las que se ha conseguido reconciliar con la información de Wikidata. Todas las demás se quedan en blanco, ya que no hay información acerca de ellas. El resultado final se puede ver en la ilustración 9, filtradas también por municipio, cogiendo el municipio Durango, Bizkaia.

Lugar Muerte / Non hil zuten	Nombre y Apellidos / Izena	Vecindad / Non I	Lugar Muerte / Non hil zuten	coordinate location	area	Fecha / Data
Arota Zaramona, Florencio	Arota Zaramona, Florencio	Durango	Durango, Bizkaia	43.166974,-2.6321025	10.91	1937-04-26T00:00:00Z
Arota Zaramona, Victor	Arota Zaramona, Victor	Durango	Durango, Bizkaia	43.166974,-2.6321025	10.91	1937-03-31T00:00:00Z
Arota Zaramona, Maria Dolores	Arota Zaramona, Maria Dolores	Durango	Durango, Bizkaia	43.166974,-2.6321025	10.91	1937-03-31T00:00:00Z
Arota Zaramona, Maria Cruz	Arota Zaramona, Maria Cruz	Durango	Durango, Bizkaia	43.166974,-2.6321025	10.91	1937-03-31T00:00:00Z
Arota Zaramona, Estefania	Arota Zaramona, Estefania	Durango	Durango, Bizkaia	43.166974,-2.6321025	10.91	1937-03-31T00:00:00Z
Arota Zaramona, Jose	Arota Zaramona, Jose	Durango	Durango, Bizkaia	43.166974,-2.6321025	10.91	1937-03-31T00:00:00Z
Arota Zaramona, Angel	Arota Zaramona, Angel	Durango	Durango, Bizkaia	43.166974,-2.6321025	10.91	1937-03-31T00:00:00Z
Arota Zaramona, Nicolas	Arota Zaramona, Nicolas	Durango	Durango, Bizkaia	43.166974,-2.6321025	10.91	1937-03-31T00:00:00Z
Arota Zaramona, Ana Maria	Arota Zaramona, Ana Maria	Durango	Durango, Bizkaia	43.166974,-2.6321025	10.91	1937-03-31T00:00:00Z
Arota Zaramona, Josefa	Arota Zaramona, Josefa	Durango	Durango, Bizkaia	43.166974,-2.6321025	10.91	1937-03-31T00:00:00Z
Arota Zaramona, Bea Balza	Arota Zaramona, Bea Balza	Durango	Durango, Bizkaia	43.166974,-2.6321025	10.91	1937-03-31T00:00:00Z
Arota Zaramona, Beatriz	Arota Zaramona, Beatriz	Durango	Durango, Bizkaia	43.166974,-2.6321025	10.91	1937-03-31T00:00:00Z
Arota Zaramona, Marcos	Arota Zaramona, Marcos	Durango	Durango, Bizkaia	43.166974,-2.6321025	10.91	1937-03-31T00:00:00Z
Arota Zaramona, Maria Josefa	Arota Zaramona, Maria Josefa	Durango	Durango, Bizkaia	43.166974,-2.6321025	10.91	1937-03-31T00:00:00Z
Arota Zaramona, Pedro	Arota Zaramona, Pedro	Durango	Durango, Bizkaia	43.166974,-2.6321025	10.91	1937-03-31T00:00:00Z
Arota Zaramona, Florencio	Arota Zaramona, Florencio	Durango	Durango, Bizkaia	43.166974,-2.6321025	10.91	1937-03-31T00:00:00Z
Arota Zaramona, Teresa	Arota Zaramona, Teresa	Durango	Durango, Bizkaia	43.166974,-2.6321025	10.91	1937-03-31T00:00:00Z
Arota Zaramona, Mariategui	Arota Zaramona, Mariategui	Durango	Durango, Bizkaia	43.166974,-2.6321025	10.91	1937-03-31T00:00:00Z

Ilustración 9: Ejemplo del enriquecimiento con wikidata