

---

# Práctica-3: Herramientas ETL PDI (Pentaho Data Integration).

---



**UNIVERSIDAD  
DE GRANADA**

Sistemas Multidimensionales (2019-2020)

Daniel Bolaños Martínez  
danibolanos@correo.ugr.es  
Grupo 2 - Viernes 15:30h

# Índice

1. Crea una BD PostgreSQL (cordoba\_danibolanos). En el esquema public de esa BD crea las tablas cuando, donde y padron. La estructura de estas tablas ha de ser similar a la de las hojas correspondientes del archivo cordoba-ETL-danibolanos.xlsx de la práctica anterior. Crea las transformaciones y jobs como se especifica en el guión. 3
2. Crea una BD PostgreSQL (prueba). En el esquema public de esa BD crea la cuando\_danibolanos. La estructura de estas tablas ha de ser similar a la de las hojas correspondientes del archivo cordoba-ETL-danibolanos.xlsx de la práctica anterior. Define el contenido de esa tabla mediante una transformación usando como origen la hoja Provincia del archivo generado mediante Power Query en la actividad anterior. 18
3. Bibliografía. 28

1. Crea una BD PosgreSQL (cordoba\_danibolanos). En el esquema public de esa BD crea las tablas cuando, donde y padron. La estructura de estas tablas ha de ser similar a la de las hojas correspondientes del archivo cordoba-ETL-danibolanos.xlsx de la práctica anterior. Crea las transformaciones y jobs como se especifica en el guión.

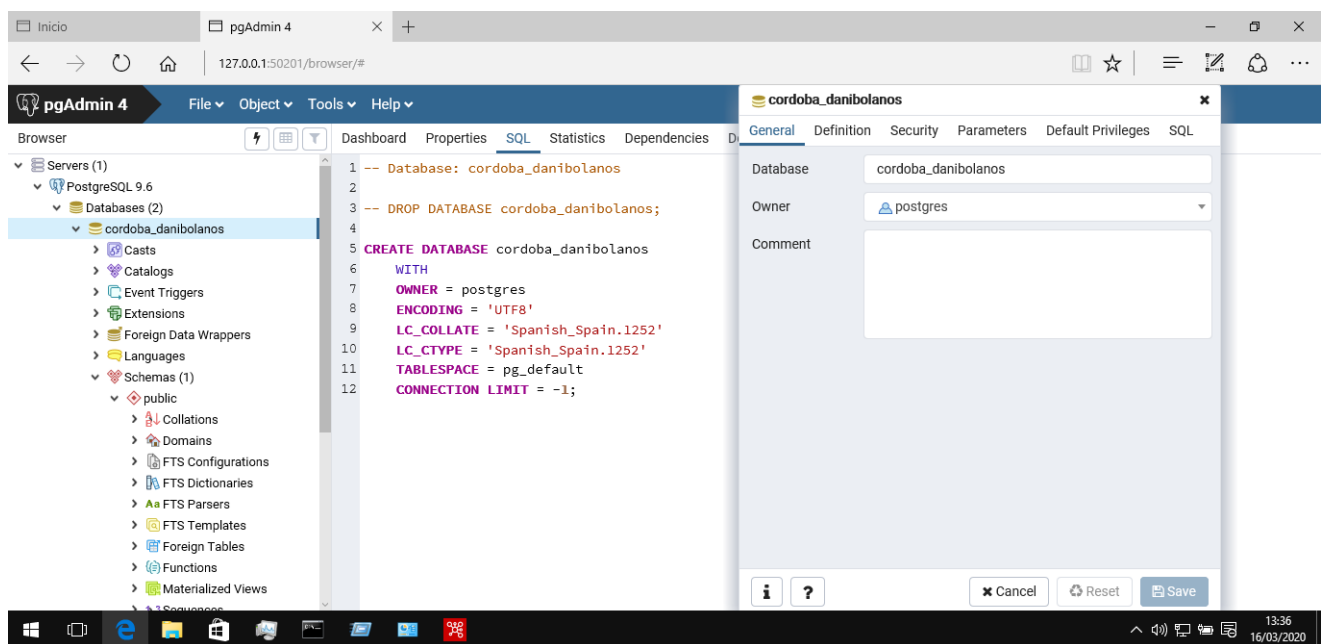


Figura 1: Creación de la BD cordoba\_danibolanos en *pgAdmin*.

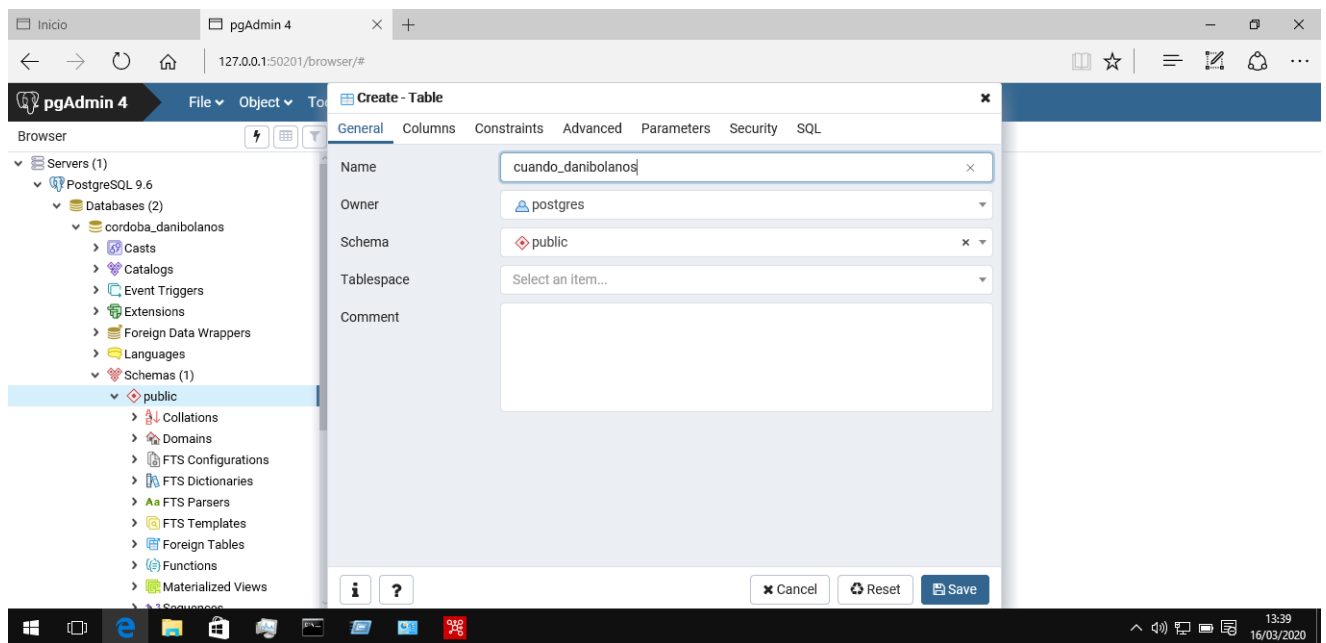


Figura 2: Creación de las tablas **cuando/donde/padron\_danibolanos**.

Creamos las tablas **cuando\_danibolanos**, **donde\_danibolanos** y **padron\_danibolanos** para la base de datos **cordoba\_danibolanos** creada anteriormente.

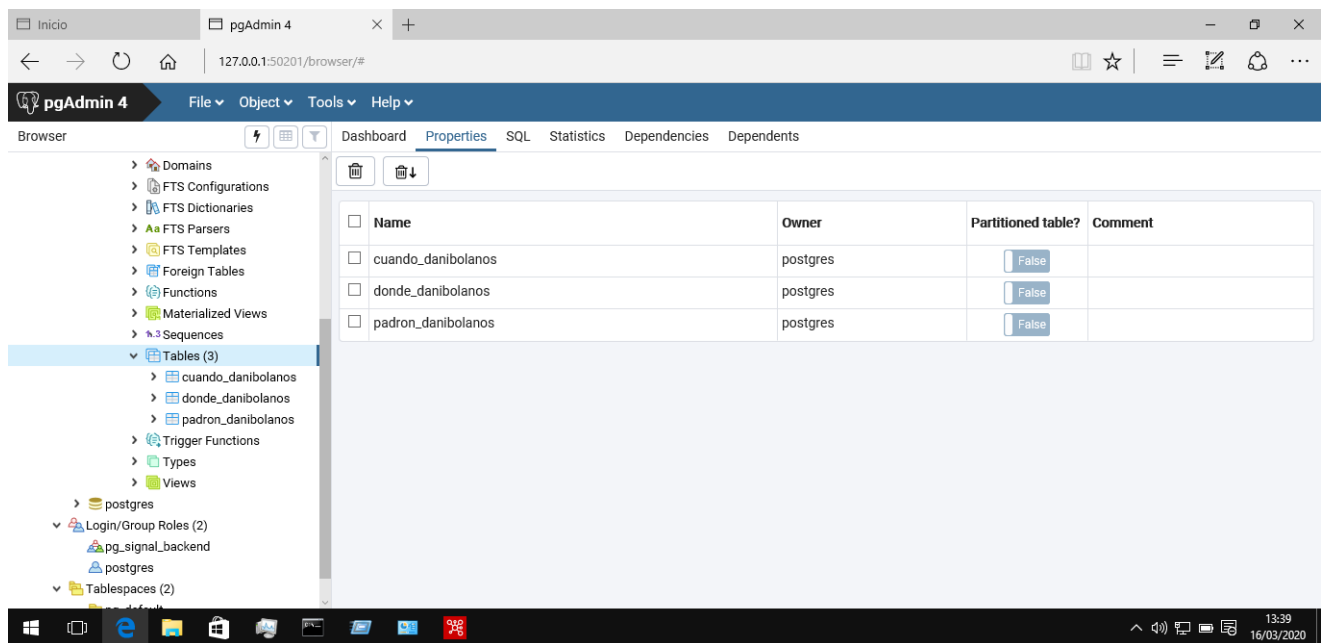


Figura 3: Tablas creadas en la BD.

Podemos observar las tres tablas creadas para la base de datos **cordoba\_danibolanos**.

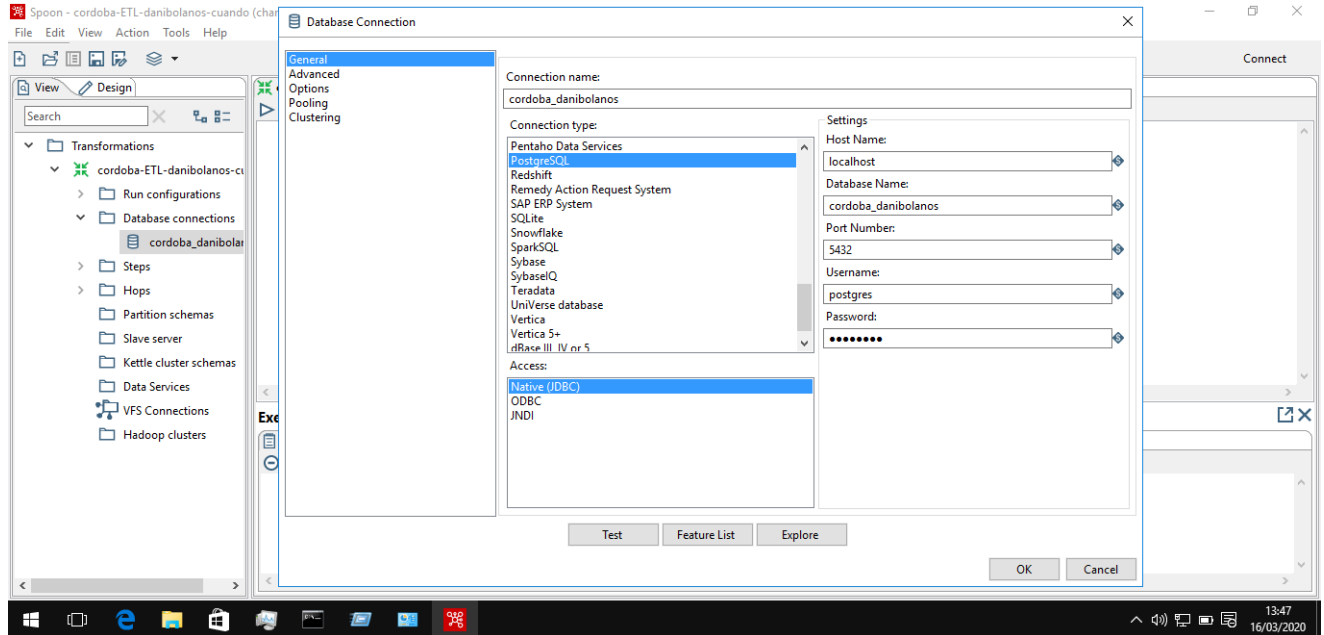


Figura 4: Creamos la conexión de la base de datos en *Spoon*.

Configuramos la conexión a la BD que hemos creados en *pgAdmin* para las transformaciones en *Spoon*. Usamos la configuración especificada en el guión.

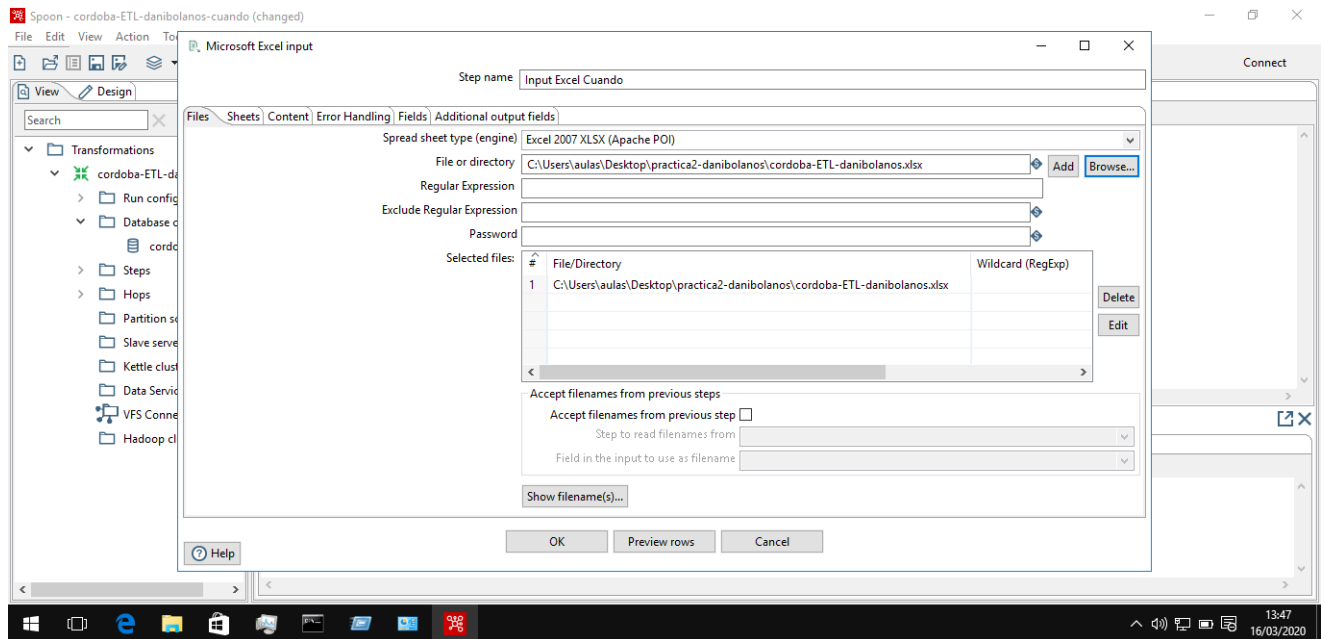


Figura 5: Extraemos los datos con el paso *Input Excel*.

Configuramos el paso *Input Excel* la obtención de los datos del archivo **cordoba-ETL-danibolanos.xlsx**.

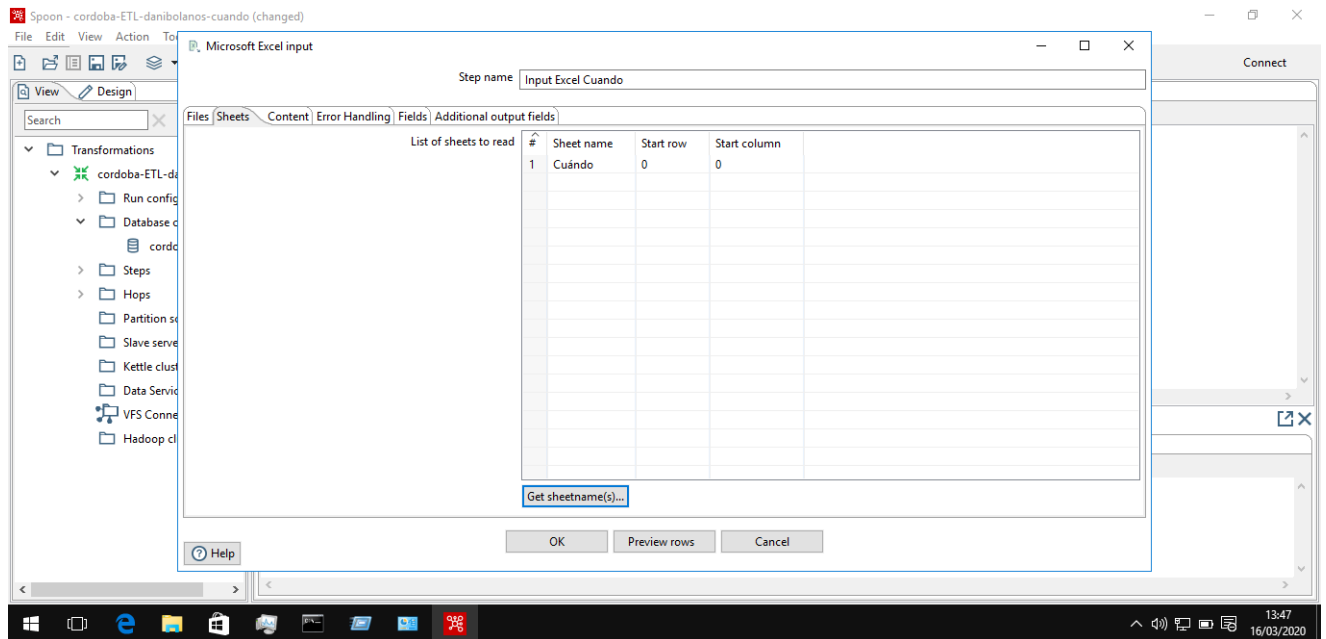


Figura 6: Extraemos los datos de la hoja **Cuándo**.



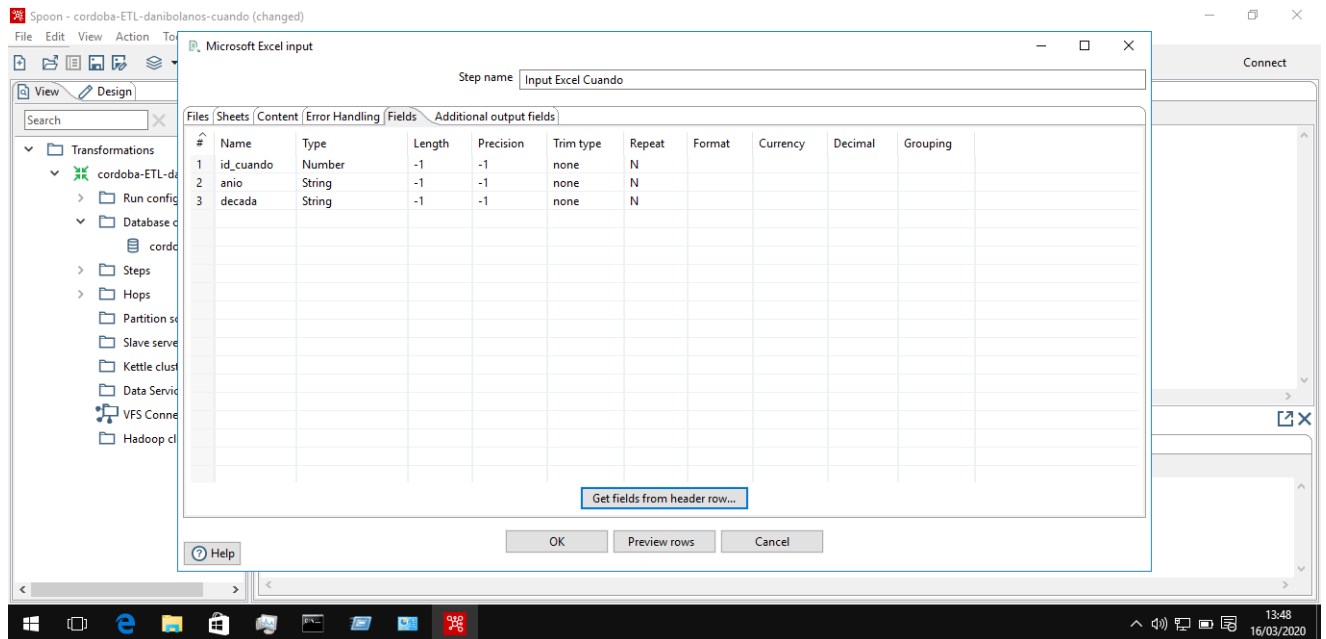


Figura 7: Campos de la hoja **Cuándo** con los criterios especificados.

Para cada hoja (**Cuándo**, **Dónde**, **Padrón**) del archivo excel generado en la práctica anterior, seleccionamos los campos y ponemos los nombres en minúscula, sin tildes y utilizando el criterio *camel\_case*.

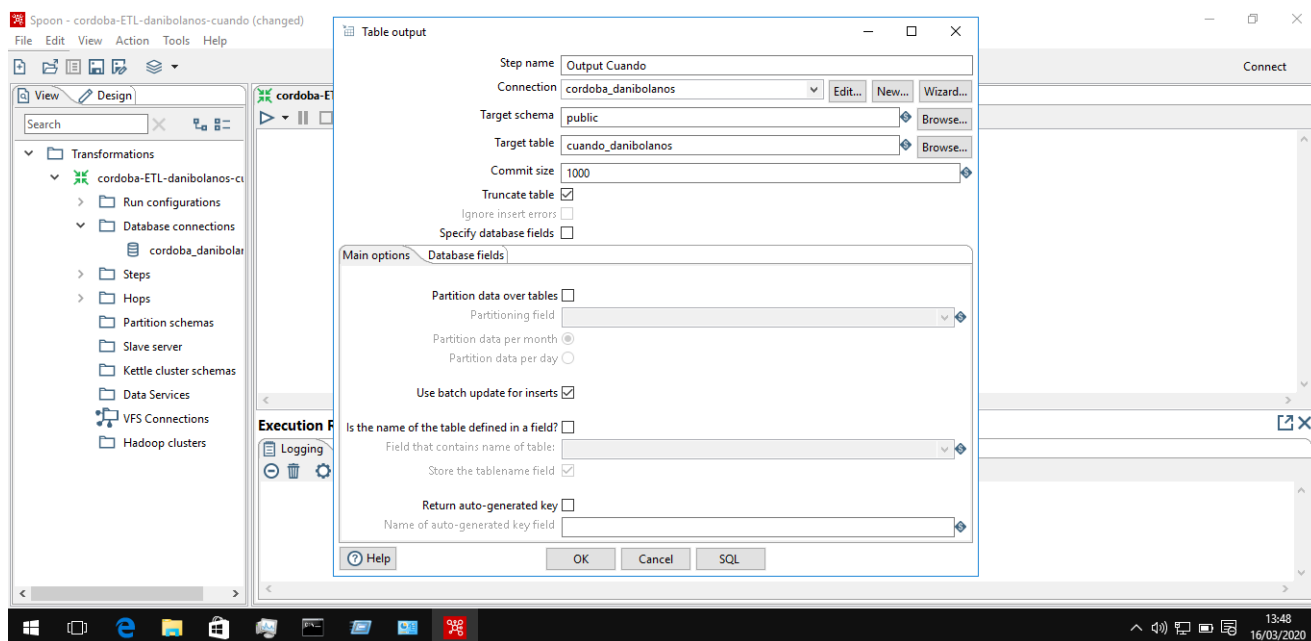


Figura 8: Configuración del paso *Table Output*.

Configuramos tal y como se indica en el guión para la hoja **Cuándo** el paso *Table Output* que se encargará de crear la tabla con los datos leídos para añadirla a la BD **cordoba\_danibolanos** de *pgAdmin*.

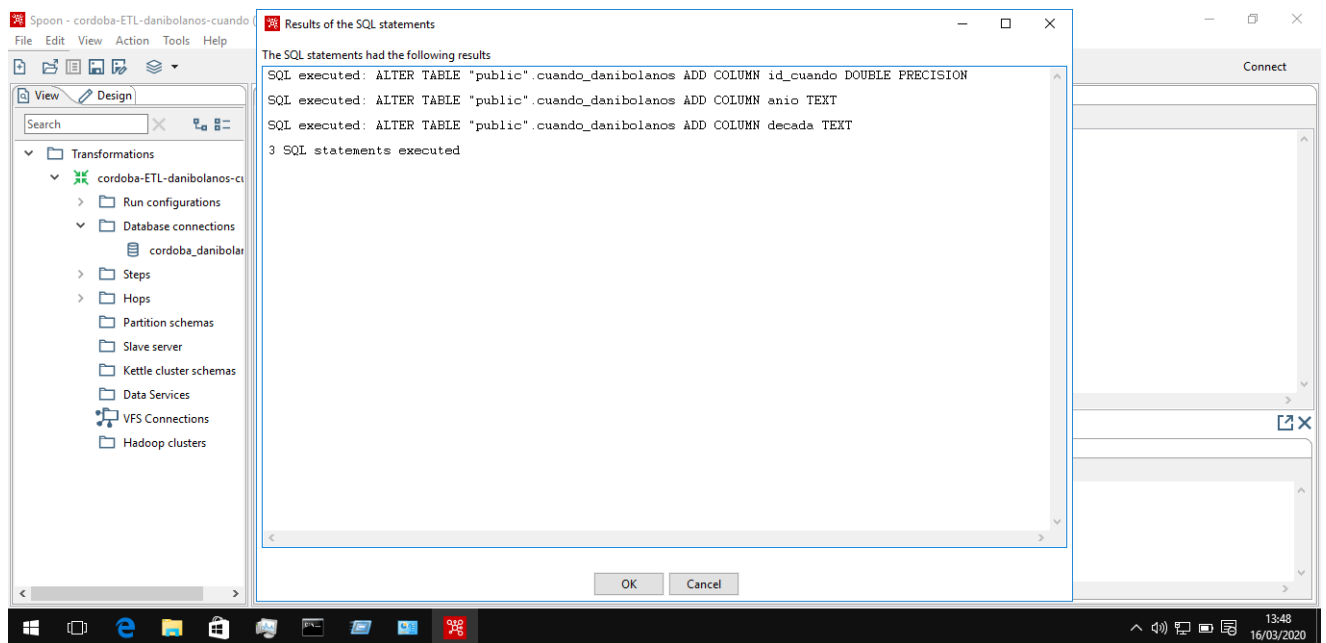


Figura 9: Generamos las sentencias SQL.

Al pulsar sobre el botón *SQL* del paso Table Output generamos las sentencias *SQL* encargadas de crear las tablas en la BD.

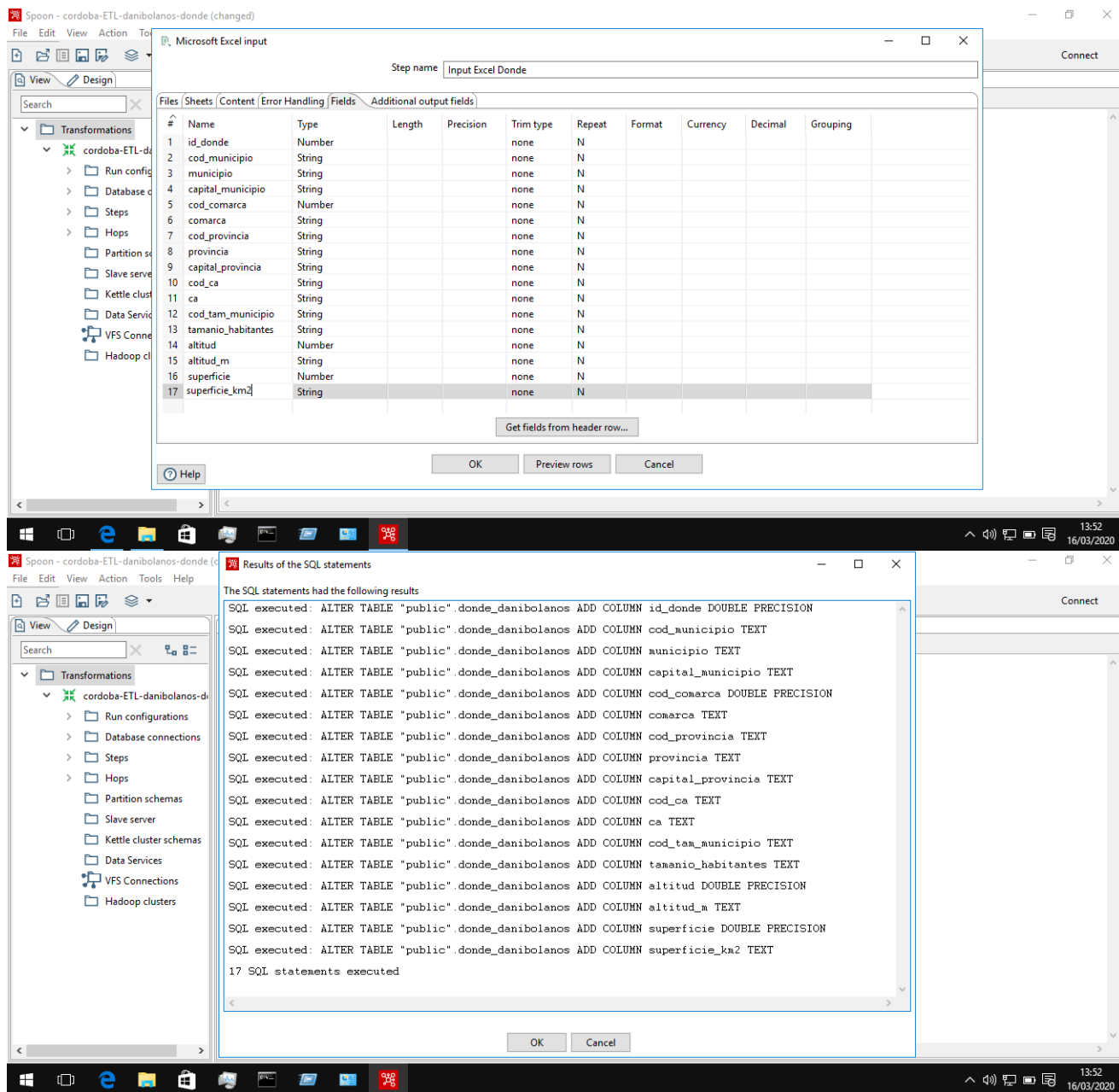


Figura 10: Realizamos el mismo proceso para **Dónde**.

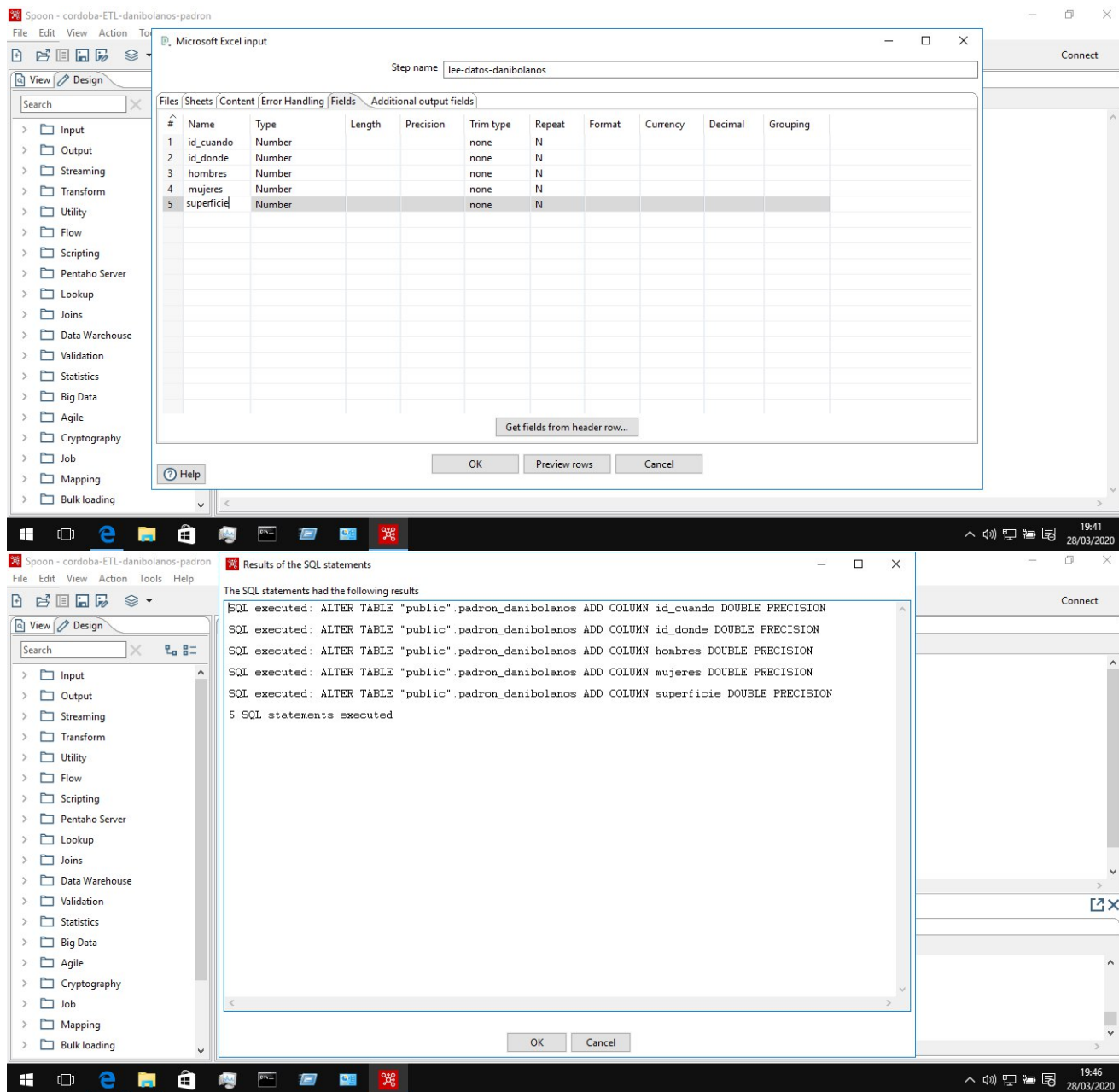


Figura 11: Realizamos el mismo proceso para **Padrón**.

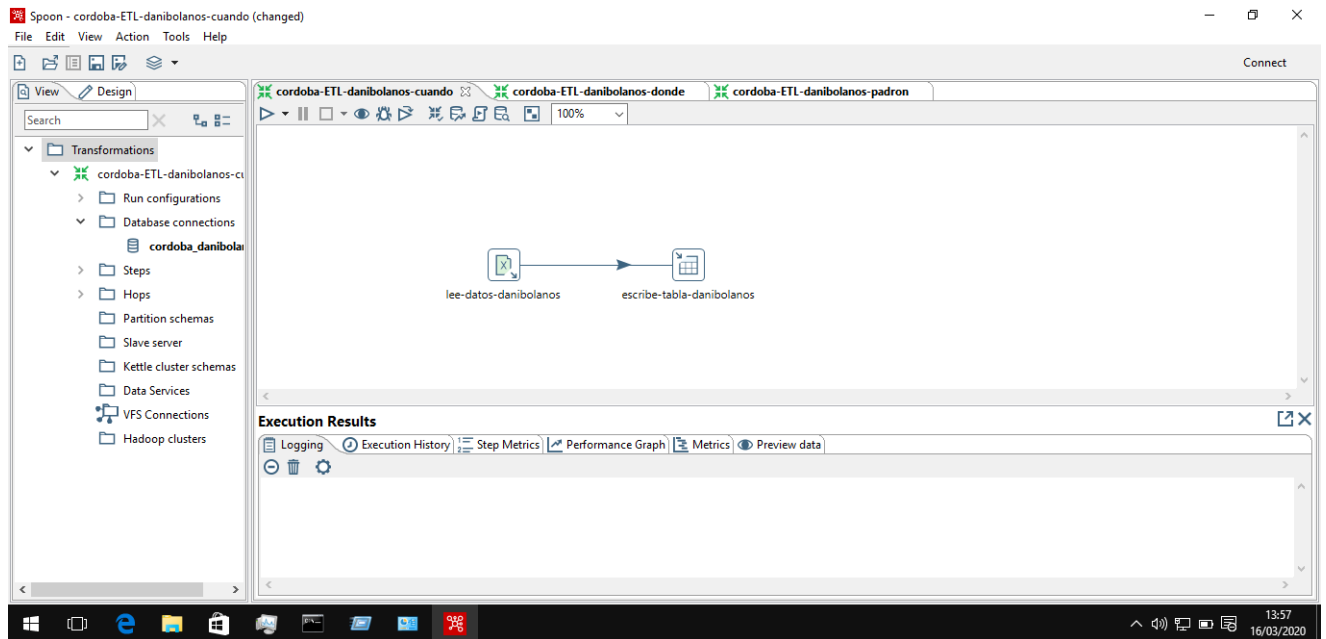


Figura 12: Ejemplo de la transformación **cordoba-ETL-danibolanos-cuando**.

Para las hojas **Dónde** y **Padrón** se han creado las transformaciones (**cordoba-ETL-danibolanos-donde** y **cordoba-ETL-danibolanos-padron**) de forma similar a la explicada.

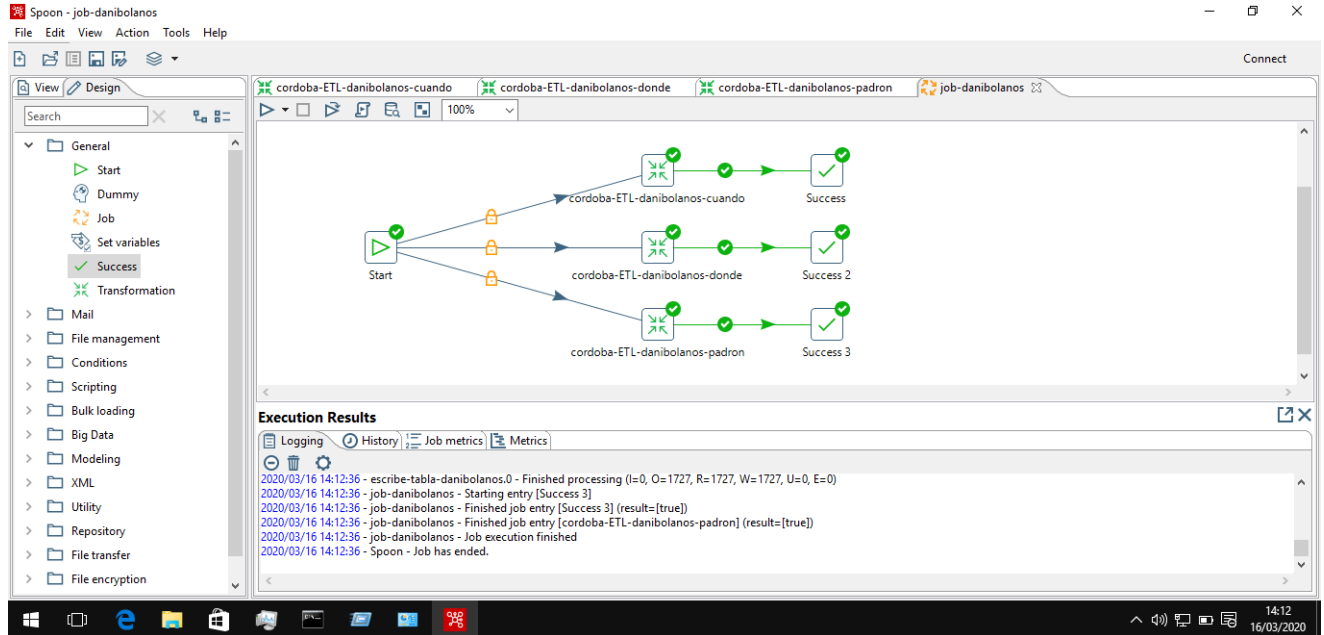


Figura 13: Job generado para las transformaciones.

Creamos una tarea 'Job' para ejecutar las tres transformaciones encargadas de crear las tablas en la BD para las tres tablas especificadas.

Creamos un paso **Start** para empezar el 'Job' y utilizaremos un paso **Success** para cada transformación para controlar si funciona todo bien o si se generan errores.

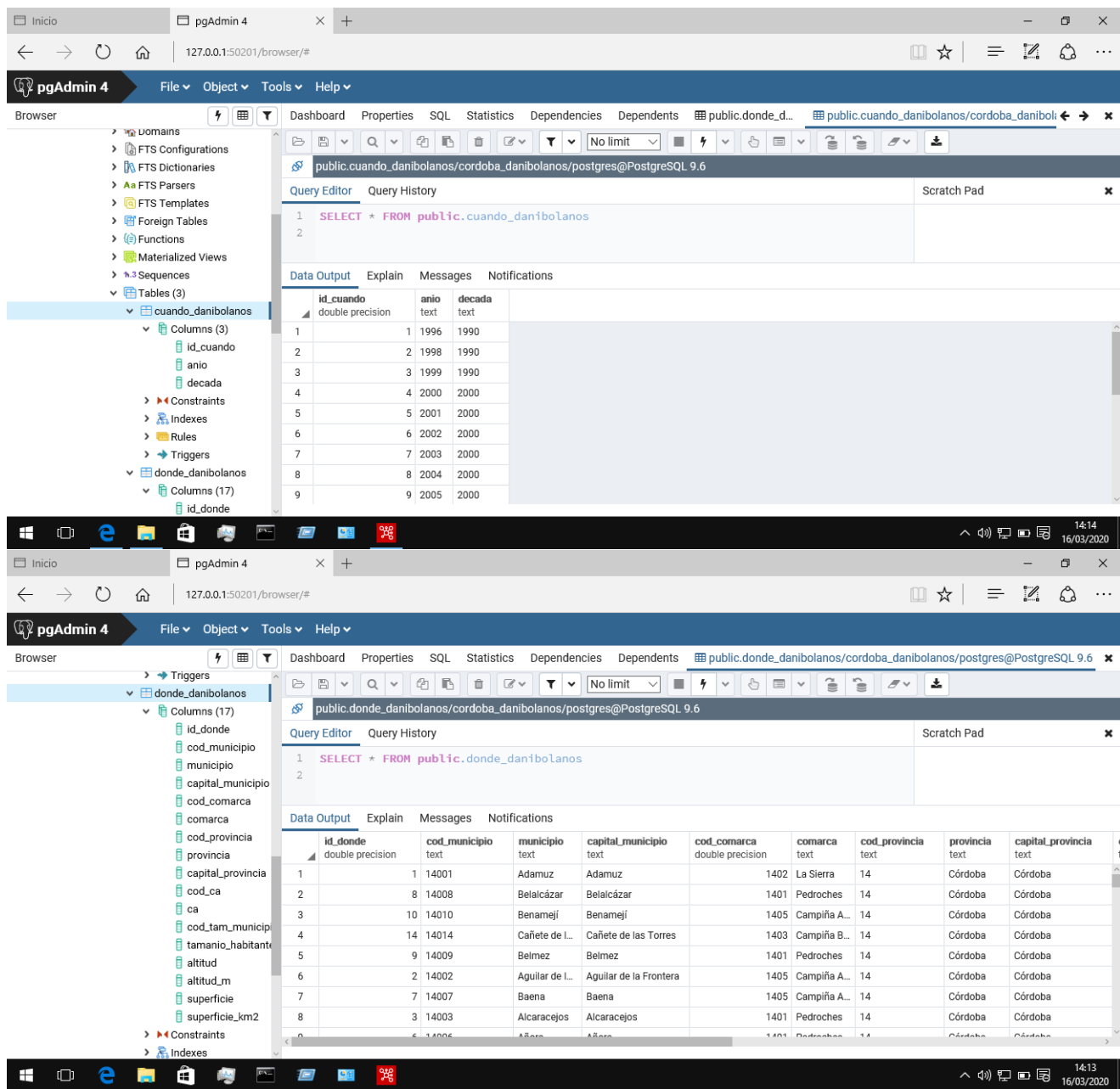


Figura 14: Datos para tablas **cuando** y **donde** en *pgAdmin*.



The screenshot shows the pgAdmin 4 web interface. On the left, the 'padron' database is selected, and the 'Columns (5)' folder is expanded. The main panel displays the 'Data Output' tab for the 'padron' table. The table contains 9 rows of data. The columns are: id\_cuando (double precision), id\_donde (double precision), hombres (double precision), mujeres (double precision), and superficie (double precision).

	id_cuando double precision	id_donde double precision	hombres double precision	mujeres double precision	superficie double precision
1		1	2241	2225	335.0024
2		2	2257	2183	335.0024
3		1	6596	6738	166.0574
4		2	6631	6766	166.0574
5		1	1906	2037	355.986746
6		2	1870	2009	355.986746
7		1	2338	2344	53.3505
8		2	2378	2357	53.3505
9		1	733	745	175.8709

Figura 15: Datos para tabla **padron** en *pgAdmin*.

Podemos observar como se han añadido correctamente los datos de las tres hojas a la BD que habíamos creado previamente en *pgAdmin*.

2. Crea una BD PostgreSQL (prueba). En el esquema public de esa BD crea la cuando\_danibolanos. La estructura de estas tablas ha de ser similar a la de las hojas correspondientes del archivo cordoba-ETL-danibolanos.xlsx de la práctica anterior. Define el contenido de esa tabla mediante una transformación usando como origen la hoja Provincia del archivo generado mediante Power Query en la actividad anterior.

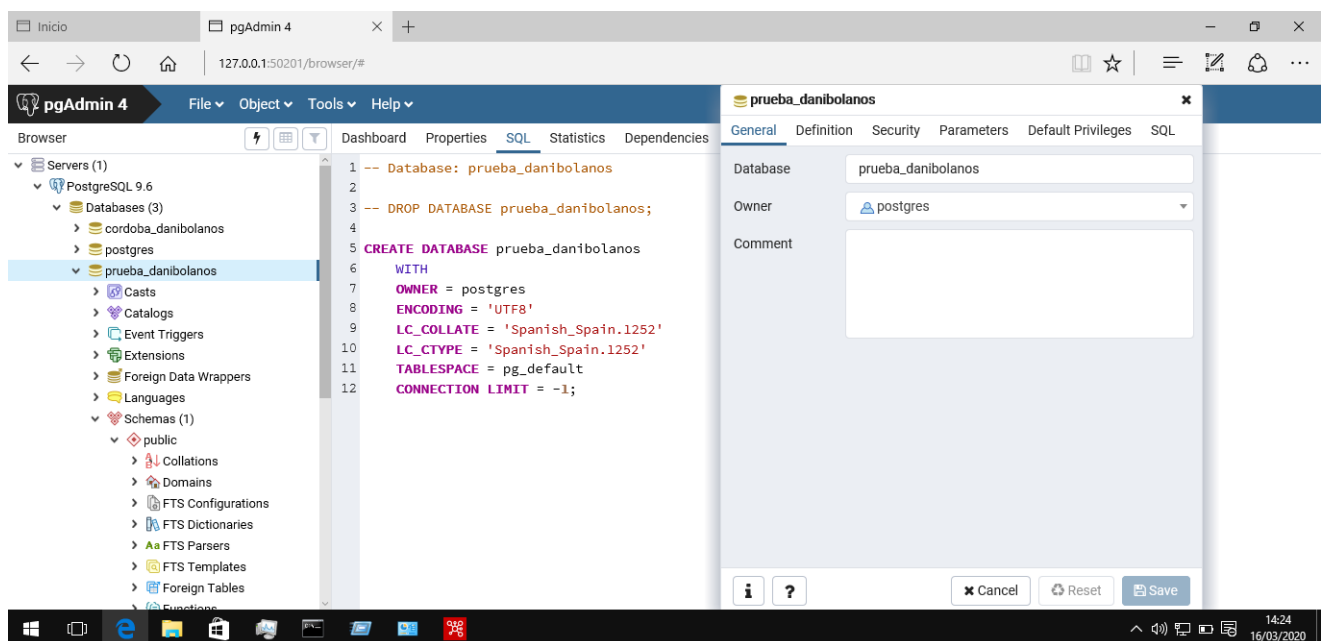


Figura 16: Creamos la BD `prueba_danibolanos` en *pgAdmin*.

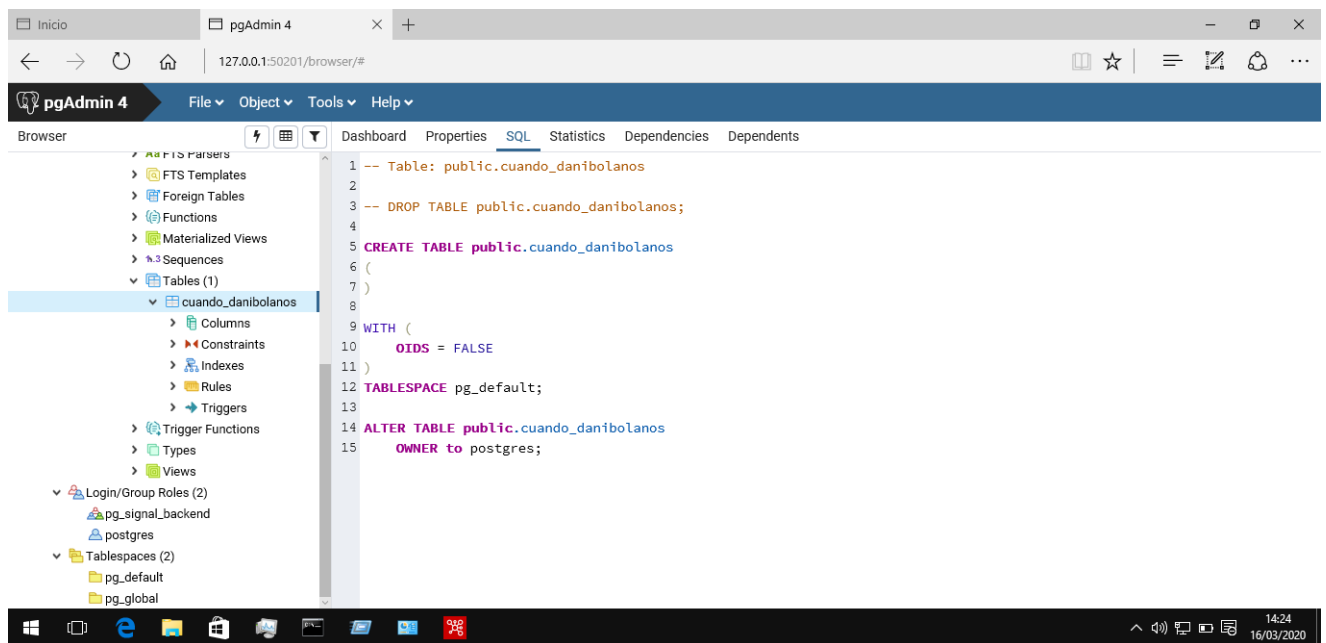


Figura 17: Creamos la tabla **cuando\_danibolanos** en la BD.

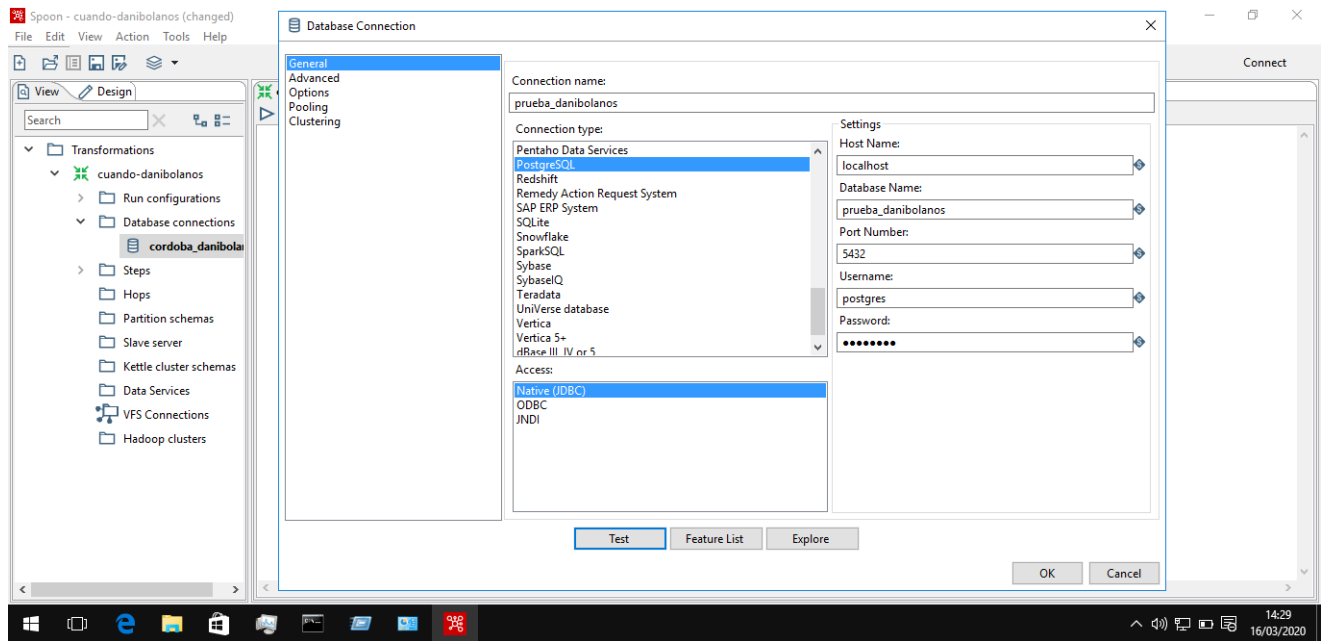


Figura 18: Configuramos la conexión para la BD prueba\_danibolanos

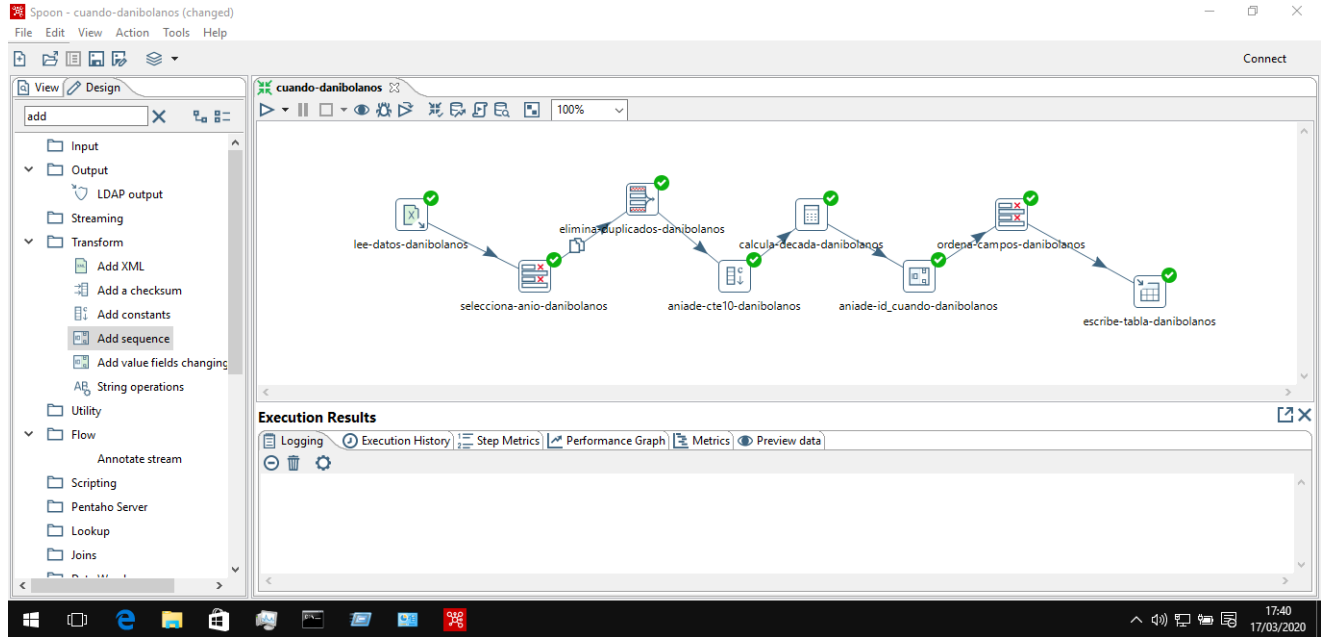


Figura 19: Configuramos los diferentes pasos en la transformación.

Para la transformación **cuando-danibolanos** creamos los siguientes pasos para pasar de la hoja **Provincia** a la tabla **cuando**.

Para ello necesitamos quedarnos sólo con el campo **anio** y completar la tabla con los campos **id\_cuando** y **decada** a partir de la información que aporta el año.

De izquierda a derecha se han usado los pasos: *Excel input file*, *Select values*, *Unique rows*, *Add constants*, *Calculator*, *Add sequence*, *Select values* y *Table output*.

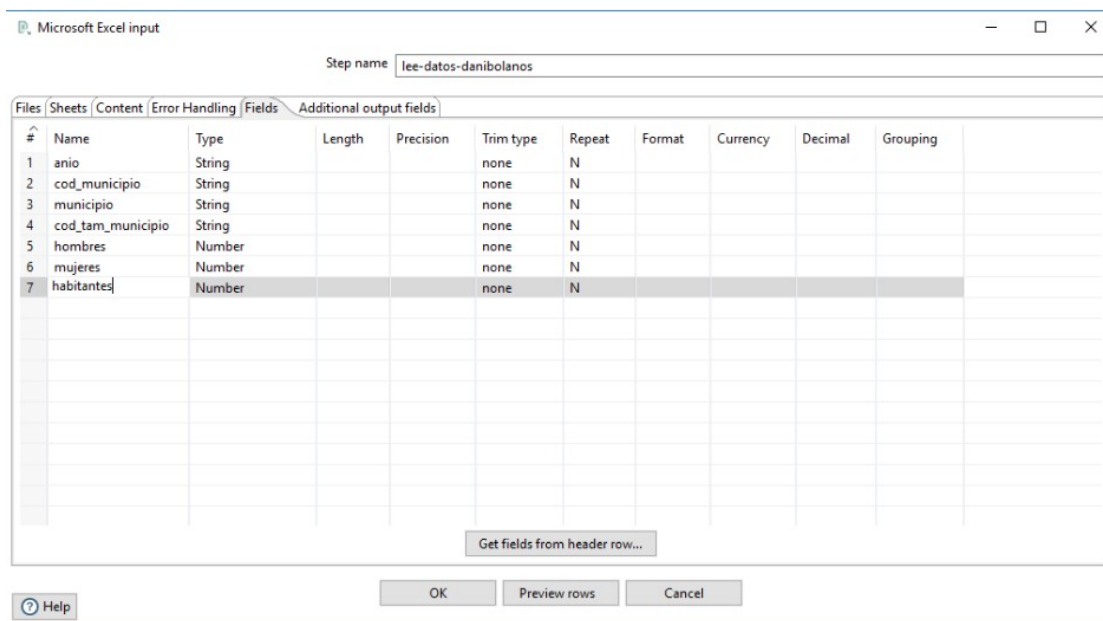


Figura 20: Leemos los datos de **Provincia** con criterio *camel\_case*.

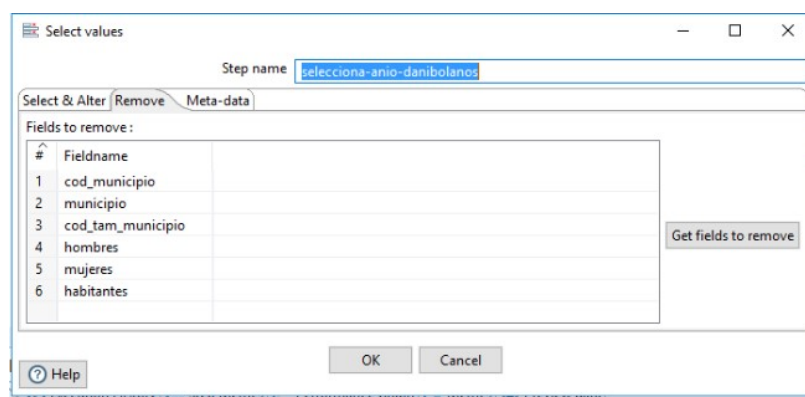


Figura 21: Eliminamos los campos innecesarios.

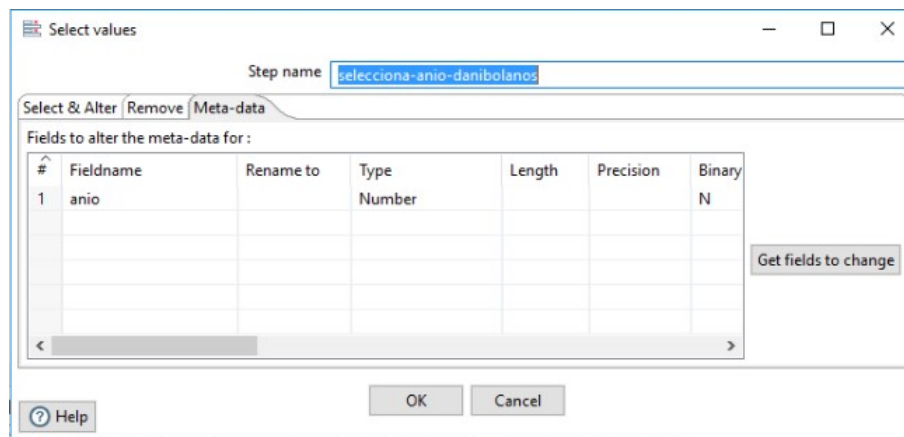


Figura 22: Cambiamos el tipo del campo **anio** a *Number*.

Cambiamos a tipo *Number* para poder hacer cálculos con la información que nos aporta para generar el campo **decada**.

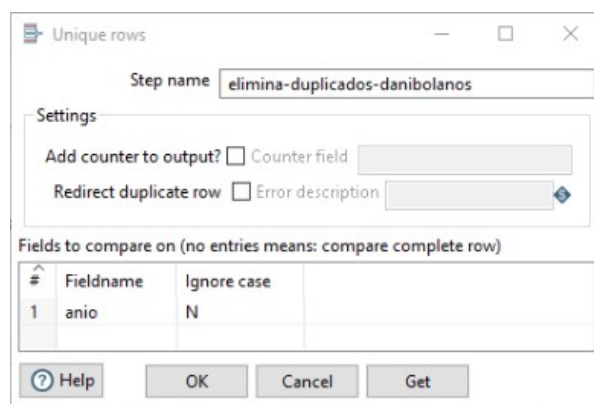


Figura 23: Eliminamos las filas duplicadas para el campo **anio**.

Add constants

Step name:

Fields:

#	Name	Type	Format	Length	Precision	Currency	Decimal	Group	Value	Set empty string?
1	cte	Number							10	N

Help OK Cancel

Figura 24: Creamos el campo **cte**.

Añadimos un campo **cte** y lo completamos con el valor 10 para poder usarlo posteriormente para calcular el campo **decada**.

Calculator

Step name:

☐ Throw an error on non existing files

Fields:

#	New field	Calculation	Field A	Field B	Field C	Value type	Length	Precision	Remove	Conversion mask	Decimal symbol	Grouping symbol	Currency symbol
1	temp1	A / B	anio	cte		None			N				
2	temp2	FLOOR( A )	temp1			None			N				
3	decada	A * B	temp2	cte		None			N				

Help OK Cancel

Figura 25: Calculamos el campo **decada** a partir del resto.

Calculamos el campo **decada** a partir de la fórmula:  $\text{floor}(\frac{\text{anio}}{10}) \cdot 10$ .



Step name: aniade-id\_cuando-danibolanos

Name of value: id\_cuando

Use a database to generate the sequence

Use DB to get sequence? ☐

Connection: prueba\_danibolanos

Schema name:

Sequence name: SEQ\_

Use a transformation counter to generate the sequence

Use counter to calculate sequence? ☒

Counter name (optional):

Start at value: 1

Increment by: 1

Maximum value: 23

Figura 26: Añadimos el campo **id\_cuando**.

Generamos los valores del campo de 1 a 23.

Step name: ordena-campos-danibolanos

Select & Alter Remove Meta-data

Fields:

#	Fieldname	Rename to	Length	Precision
1	id_cuando			
2	anio			
3	decada			

Include unspecified fields, ordered by name ☒

Figura 27: Reordenamos los valores para adecuarlos a la tabla **cuando**.

Table output

Step name: escribe-tabla-danibolanos

Connection: prueba\_danibolanos [Edit...] [New...] [Wizard...]

Target schema: public [Browse...]

Target table: cuando\_danibolanos [Browse...]

Commit size: 1000

Truncate table: ☒

Ignore insert errors: ☐

Specify database fields: ☐

Main options: Database fields

Partition data over tables: ☐

Partitioning field: [dropdown]

Partition data per month: ☒

Partition data per day: ☐

Use batch update for inserts: ☒

Is the name of the table defined in a field?: ☐

Field that contains name of table: [dropdown]

Store the tablename field: ☒

Return auto-generated key: ☐

Name of auto-generated key field: [dropdown]

[?] Help [OK] [Cancel] [SQL]

Figura 28: Escribimos los valores en la tabla.

Pulsamos el botón *SQL* para crear la tabla **cuando\_danibolanos**.

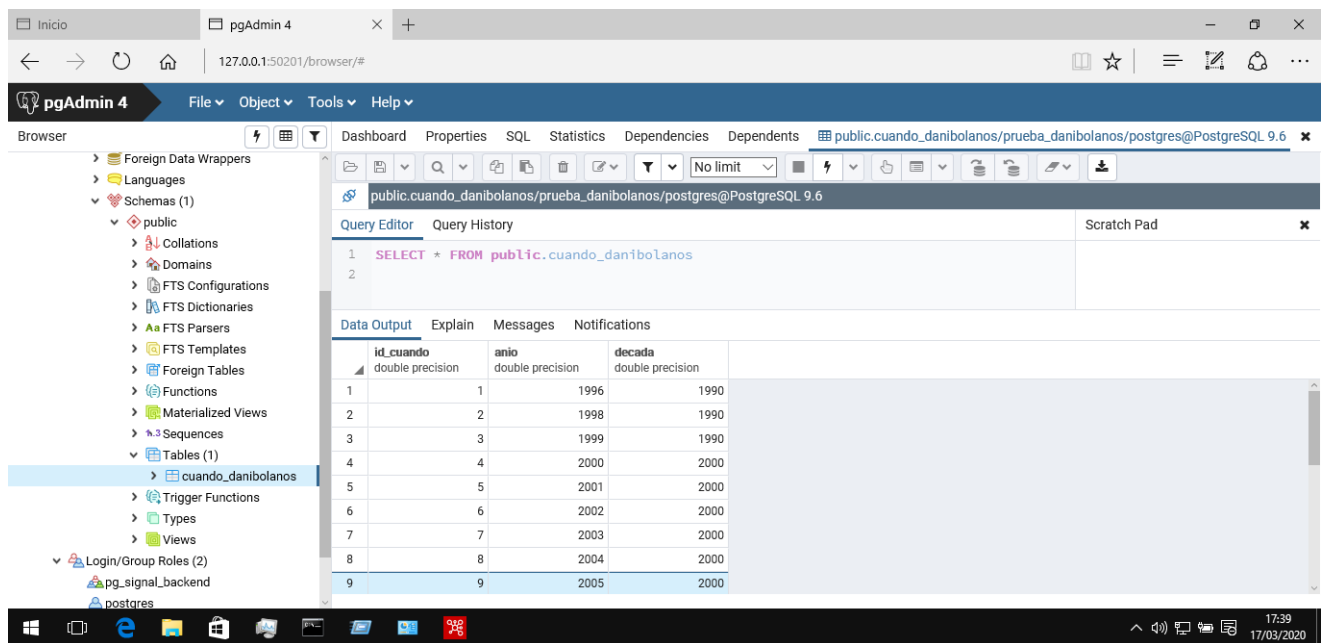


Figura 29: Vemos los resultados en la BD.

Podemos ver los resultados obtenidos para la tabla **cuando\_danibolanos** y cómo se adecua a los valores del ejercicio 1.

### **3. Bibliografía.**

#### **Referencias**

- [1] Guión de prácticas de la asignatura.