

#### 4. Modelo de calidad de datos

DQ Problems	DQ Dimension	DQ Factor	DQ Metric	DQ method	Applied DQ method
	<p><b>ID:</b> D4_Unicidad  <b>Name:</b> Unicidad  <b>Description:</b> Evalúa si los registros en un conjunto de datos son únicos y no presentan duplicaciones no deseadas  <b>Suggested by:</b> {RN2, RN5}</p>	<p><b>ID:</b>  <b>Name:</b> No duplicación  <b>Description:</b> Porcentaje de datos que no están duplicados en forma exacta  <b>Represents</b> = {RN2, RN5}</p>	<p><b>ID:</b> M_duplicate_ratio  <b>Name:</b> duplicate_ratio  <b>Description:</b> Da el porcentaje de valores duplicados.  <b>Influenced by</b> = {RN2, RN5}  <b>Granularity:</b> Columna  <b>Result domain</b> = [0...1]</p>	<p><b>ID:</b> metodo_duplicate_ratio  <b>Name:</b> met_duplicate_ratio  <b>Description:</b> Implementa la metrica duplicate_ratio sobre un atributo dado.  <b>Uses</b> = {RN2, RN5}  <b>Input data types:</b> String  <b>Output data types:</b> Float  <b>Algorithm:</b>  <pre>"def duplicate_ratio(data, column):     total = len(data)     duplicated = data.duplicated(subset=[column]).sum()     return duplicated / total"</pre> </p>	<p><b>ID:</b> MA_duplicate_ratio  <b>Type:</b> Medición  <b>Descriptinon:</b>  Dados los datos y el atributo, calcula el porcentaje de valores duplicados en ese atributo.  <b>AppliedTo:</b>  Atributos «ISBN», «AuthorID», «PublisherID», «User_ID» de la tabla NL_Books.  Atributos «ID» de la tabla NL_USERS."</p>
P11	<p><b>ID:</b> D2_Completitud  <b>Name:</b> Completitud  <b>Description:</b> Indica si el sistema de información contiene toda la información de interés.  <b>Suggested by:</b> {RN1, RN2, RN5, RQ3, RQ7, F2}</p>	<p><b>ID:</b> F2_Dens  <b>Name:</b> Densidad  <b>Description:</b> Indica cuanta información se tiene y cuanta falta sobre las entidades del sistema de información.  <b>Represents</b> = {RN1, RN2, RN5, RQ3, RQ7, F2}</p>	<p><b>ID:</b> M_contar_nulls  <b>Name:</b> contar_nulls  <b>Description:</b> Calcula el porcentaje de entradas vacias en una columna  <b>Influenced by</b> = {RN1, RN5, RQ3, RQ7, F2}  <b>Granularity:</b> Columna  <b>Result domain</b> = {0, 1}</p>	<p><b>ID:</b> metodo_contar_nulls  <b>Name:</b> met_contar_nulls  <b>Description:</b> Implementa la metrica contar_nulls.  <b>Uses</b> = {RN1, RN5, RQ3, RQ7, F2}  <b>Input data types:</b> String  <b>Output data types:</b> Float  <b>Algorithm:</b>  <pre>"def porcentaje_nulos(atributo):     nulos = atributo.isnull()     contador_nulos = nulos.sum()     porcentaje = contador_nulos / len(atributo)     return porcentaje"</pre> </p>	<p><b>ID:</b> MA_contar_nulls  <b>Type:</b> Medición  <b>Descriptinon:</b> Dado un atributo de una tabla, calcula el porcentaje de entradas vacias.  <b>AppliedTo:</b>  Cualquier atributo de cualquier tabla.</p>
			<p><b>ID:</b> M_check_RN1  <b>Name:</b> check_RN1  <b>Description:</b> Da el porcentaje de entradas de la tabla que tienen al menos un campo vacio entre los atributos isbn, titulo, autor y editor.  <b>Influenced by</b> = {RN1, RN2, RQ3, RQ7, F2}  <b>Granularity:</b> Conjunto de columnas  <b>Result domain</b> = [0...1]</p>	<p><b>ID:</b> metodo_check_RN1  <b>Name:</b> met_check_RN1  <b>Description:</b> Implementa la metrica Check_RN1 sobre una tabla dada.  <b>Uses</b> = {RN1, RN2, RQ3, RQ7, F2}  <b>Input data types:</b> String  <b>Output data types:</b> Float  <b>Algorithm:</b>  <pre>"def Check_RN1(datos):     # Seleccionar solo las columnas relevantes (isbn, titulo, autor, editor)     datos_relevantes = datos[['isbn', 'titulo', 'autor', 'editor']]      # Verificar cuántas filas tienen al menos un campo NULL entre estas columnas     incompletos = datos_relevantes.isnull().any(axis=1).sum()      # Calcular la proporción de filas incompletas     return incompletos / len(datos)"</pre> </p>	<p><b>ID:</b> MA_check_RN1  <b>Type:</b> Medición  <b>Descriptinon:</b> Dada una tabla, calcula el porcentaje de entradas que tienen al menos un campo vacio entre los atributos isbn, titulo, autor y editor.  <b>AppliedTo:</b>  Tabla NL_Books.</p>

DQ Problems	DQ Dimension	DQ Factor	DQ Metric	DQ method	Applied DQ method	
		<p><b>ID:</b> D3_Consistencia</p> <p><b>Name:</b> Consistencia</p> <p><b>Description:</b> Captura la satisfacción de reglas semánticas definidas sobre los datos.</p> <p><b>Suggested by</b> = {RQ4, RQ6, RQ8, RQ11, F1, F3}</p>	<p><b>ID:</b></p> <p><b>Name:</b> Integridad inter-relación</p> <p><b>Description:</b> Captura la satisfacción de reglas entre atributos de tablas distintas.</p> <p><b>Represents</b> = {RQ4, RQ6, RQ8, RQ11, F1, F3}</p>	<p><b>ID:</b> M_consistencia_ratings</p> <p><b>Name:</b> concistencia_ratings</p> <p><b>Description:</b> Verifica que la cantidad de ratings entre las distintas tablas sea coherente.</p> <p><b>Influenced by</b> = {RQ8}</p> <p><b>Granularity:</b> conjunto de columnas</p> <p><b>Result domain</b> = {0, 1}</p>	<p><b>ID:</b> metodo_consistencia_ratings</p> <p><b>Name:</b> met_consistencia_ratings</p> <p><b>Description:</b> Implementa la metrica concistencia_ratings contando la cantidad de ocurrencias de un ISBN valido en la tabla NL_ratings y lo compara con lo declarado en el atributo rating_count de NL_Books.</p> <p><b>Uses</b> = {RQ8}</p> <p><b>Input data types:</b> String</p> <p><b>Output data types:</b> Boolean</p> <p><b>Algorithm:</b></p> <pre>"def función consistencia_ratings(libros, valid_isbn, ratings):     libros_validos = libros[valid_isbn==True]      for libro in libros_validos:         isbn = libro['isbn']         cantidad = (ratings['isbn'] == isbn).sum()         contados[isbn] = cantidad      rating_esperado = libros_validos['rating_counts']      return rating_esperado==contados"</pre>	<p><b>ID:</b> MA_consistencia_ratings</p> <p><b>Type:</b> Agregación</p> <p><b>Descriptinon:</b> Dada una tabla de reviews, una tabla de libros (ambas con el campo isbn) y una lista que indiquen si son validos o no, determina si hay congruencia entre lo registrado en ambas tablas sobre los ratings.</p> <p><b>AppliedTo:</b> Conjunto de atributos (NL_ratings.ISBN, NL_books.ISBN, NL_books.Rating_Count)</p>
			<p><b>ID:</b> M_consistencia_fechas</p> <p><b>Name:</b> concistencia_fechas</p> <p><b>Description:</b> Verifica que la fecha de un rating sea posterior a la fecha de publicación del libro.</p> <p><b>Influenced by</b> = {RQ6, RQ11, F1}</p> <p><b>Granularity:</b> Conjunto de columnas</p> <p><b>Result domain</b> = {0, 1}</p>	<p><b>ID:</b> metodo_consistencia_fechas</p> <p><b>Name:</b> met_consistencia_fechas</p> <p><b>Description:</b> Implementa la metrica concistencia_fechas comparando las fechas de review_time de la tabla NL_reviews y la de PublisherDate de NL_Books.</p> <p><b>Uses</b> = {RQ6, RQ11, F1}</p> <p><b>Input data types:</b> String</p> <p><b>Output data types:</b> Boolean</p> <p><b>Algorithm:</b></p> <pre>"def consistencia_fechas(NL_reviews, NL_books):     result = lista vacia      for (i ; i&lt; len(NL_reviews); i++):         isbn = NL_reviews['isbn'][i]         fila_libro = fila en NL_books donde NL_books['isbn']==isbn         obtener los review_time y PublishedDate de las filas correspondientes         result[i] = review_time &gt;= PublishedDate:      return result"</pre>	<p><b>ID:</b> MA_consistencia_fechas</p> <p><b>Type:</b> Medición</p> <p><b>Descriptinon:</b> Dada una tabla de reviews y una tabla de libros (ambas con el campo isbn), determina si la fecha de la review es posterior a la publicación del libro.</p> <p><b>AppliedTo:</b> Conjunto de atributos «NL_Reviews.review_time, NL_Books.PublisherDate»</p>	
			<p><b>ID:</b> M_missing_rating_books</p> <p><b>Name:</b> missing_rating_books</p> <p><b>Description:</b> Para ratings sobre libros, controla la existencia de estos en la base de datos.</p> <p><b>Influenced by</b> = {RQ4, RQ8, F3}</p> <p><b>Granularity:</b> Celda</p> <p><b>Result domain</b> = {0, 1}</p>	<p><b>ID:</b> metodo_missing_rating_books</p> <p><b>Name:</b> met_missing_rating_books</p> <p><b>Description:</b> Implementa la metrica check_rating_books para ver si los libros indicados en los ratings existen.</p> <p><b>Uses</b> = {RQ4, RQ8, F3}</p> <p><b>Input data types:</b> String</p> <p><b>Output data types:</b> Boolean</p> <p><b>Algorithm:</b></p> <pre>"SELECT COUNT(*) FROM NL_Ratings WHERE NOT EXISTS(SELECT * FROM NL_Books WHERE NL_Ratings.ISBN = NL_Books.ISBN)"</pre>	<p><b>ID:</b> MA_missing_rating_books</p> <p><b>Type:</b> Medición</p> <p><b>Descriptinon:</b> Verifica que el libro en NL_Ratings existe en NL_Books mediante una consulta SQL.</p> <p><b>AppliedTo:</b> Conjunto de atributos «ISBN» en las tablas NL_Books y NL_Ratings</p>	

DQ Problems	DQ Dimension	DQ Factor	DQ Metric	DQ method	Applied DQ method	
P1, P2, P5, P6, P8, P14.		<p><b>ID:</b> D1_Exactitud  <b>Name:</b> Exactitud  <b>Description:</b> Concerne a la correctitud y la precisión con que los datos del mundo real son representados en un sistema de información.  <b>Suggested by</b> = {RN1, RN2, RN4, RQ2, RQ3, RQ6, RQ11, F1}</p>	<p><b>ID:</b> F1_ExactSint  <b>Name:</b> Exactitud sintáctica  <b>Description:</b> Indica que tan libre de errores sintácticos están los datos.  <b>Represents</b> = {RN1, RN2, RQ3, RN4, RQ2}</p>	<p><b>ID:</b> M_check_ISBN  <b>Name:</b> check_ISBN  <b>Description:</b> Controla si el valor es un ISBN valido.  <b>Influenced by</b> = {RN1, RN2, RQ3}  <b>Granularity:</b> Celda  <b>Result domain</b> = {0, 1}</p>	<p><b>ID:</b> metodo_check_ISBN  <b>Name:</b> met_check_ISBN  <b>Description:</b> Implementa la métrica Check_ISBN teniendo en cuenta la estructura de un código ISBN.  <b>Uses</b> = {RN1, RN2, RQ3}  <b>Input data types:</b> String  <b>Output data types:</b> Boolean  <b>Algorithm:</b>  <pre>"def Check_ISBN(codigo):     size = len(codigo)      # Chequea para el caso de un ISBN10     if (size==10):         verifica que los primeros 9 caracteres sean digitos         verifica que el ultimo caracter sea un digito o una "X"         verifica que el digito verificador sea correcto      # Chequea para el caso de un ISBN13:     if (size==13):         verifica que todos los caracteres sean digitos         verifica que el digito verificador sea correcto"</pre></p>	<p><b>ID:</b> MA_check_ISBN  <b>Type:</b> Medición  <b>Descriptinon:</b> Utiliza el algoritmo de calculo de ISBN para verificar la validez del dato.  <b>AppliedTo:</b> Atributo «ISBN» de la tabla NL_Books</p>
			<p><b>ID:</b> M_check_edades  <b>Name:</b> check_edades  <b>Description:</b> Controla si la edad es valida.  <b>Influenced by</b> = {RN4, RQ2}  <b>Granularity:</b> Celda  <b>Result domain</b> = {0, 1}</p>	<p><b>ID:</b> metodo_check_edades  <b>Name:</b> met_check_edades  <b>Description:</b> Implementa la metrica Check_edades verificando si la edad tiene formato correcto y es un valor razonable.  <b>Uses</b> = {RN4, RQ2}  <b>Input data types:</b> String  <b>Output data types:</b> Boolean  <b>Algorithm:</b>  <pre>IsNumeric(edad) AND edad &gt; 0 AND edad &lt; 100</pre></p>	<p><b>ID:</b> MA_check_edades  <b>Type:</b> Medición  <b>Descriptinon:</b> Verifica que tenga un valor numerico entre 1 y 99  <b>AppliedTo:</b> Atributo «Age» de la tabla NL_Users</p>	
			<p><b>ID:</b> M_check_price  <b>Name:</b> check_price  <b>Description:</b> Controla si el precio es valido.  <b>Influenced by</b> = {RN3}  <b>Granularity:</b> Celda  <b>Result domain</b> = {0, 1}</p>	<p><b>ID:</b> metodo_check_price  <b>Name:</b> met_check_price  <b>Description:</b> Implementa la metrica check_price verificando si el valor tiene formato correcto  <b>Uses</b> = {RN3}  <b>Input data types:</b> String  <b>Output data types:</b> Boolean  <b>Algorithm:</b>  <pre>IsNumeric(price) AND price &gt;= 0</pre></p>	<p><b>ID:</b> MA_check_price  <b>Type:</b> Medición  <b>Descriptinon:</b> Verifica que tenga un valor numerico positivo con una consulta SQL  <b>AppliedTo:</b> Atributo «Price» de la tabla NL_Books</p>	
			<p><b>ID:</b> M_duplicated_authors  <b>Name:</b> duplicated_authors  <b>Description:</b> Indica cuantas tuplas nombran mas de 1 autor.  <b>Influenced by</b> = {RN1, RQ3}  <b>Granularity:</b> Columna  <b>Result domain</b> = [0..1]</p>	<p><b>ID:</b> metodo_duplicated_authors  <b>Name:</b> met_duplicated_authors  <b>Description:</b> Implementa la metrica duplicated_author para obtener el porcentaje de tuplas que indican mas de 1 autor  <b>Uses</b> = {RN1, RQ3}  <b>Input data types:</b> String  <b>Output data types:</b> Double  <b>Algorithm:</b>  <pre>DEFINE @tuplasDuplicadas = SELECT COUNT(*) FROM NL_Authors WHERE AuthorID LIKE '%,%'  DEFINE @tuplasTotales = SELECT COUNT(*) FROM NL_Authors  SELECT @tuplasDuplicadas / @tuplasTotales</pre></p>	<p><b>ID:</b> MA_duplicated_authors  <b>Type:</b> Medición  <b>Descriptinon:</b> Divide las tuplas que indican mas de 1 autor sobre las tuplas totales  <b>AppliedTo:</b> Atributo «AuthorID» en la tabla NL_Books</p>	

DQ Problems	DQ Dimension	DQ Factor	DQ Metric	DQ method	Applied DQ method
P1, P2, P5, P6, P8, P14.	<p><b>ID:</b> D1_Exactitud  <b>Name:</b> Exactitud  <b>Description:</b> Concerne a la correctitud y la precisión con que los datos del mundo real son representados en un sistema de información.  <b>Suggested by</b> = {RN1, RN2, RN4, RQ2, RQ3, RQ6, RQ11, F1}</p>	<p><b>ID:</b> F3_Precision  <b>Name:</b> Exactitud precisión  <b>Description:</b> Indica que tan detallados son los datos.  <b>Represents</b> = {RQ6, RQ11, F1}</p>	<p><b>ID:</b> M_date_format  <b>Name:</b> date_format  <b>Description:</b> Controla el formato de fecha de todas las columnas tipo fecha.  <b>Influenced by</b> = {RQ6, RQ11, F1}  <b>Granularity:</b> Celda  <b>Result domain</b> = {0, 1}</p>	<p><b>ID:</b> metodo_date_format  <b>Name:</b> met_date_format  <b>Description:</b> Implementa la metrica date_format para controlar que el formato de fecha sea correcto  <b>Uses</b> = {RQ6, RQ11, F1}  <b>Input data types:</b> String  <b>Output data types:</b> Boolean  <b>Algorithm:</b>  <code>TRY_CONVERT(DATE,date,103)</code></p>	<p><b>ID:</b> MA_date_format  <b>Type:</b> Medición  <b>Descriptinon:</b> Verifica que el campo PublisherDate en Books y el campo review_time en Ratings respeten el formato correcto de fecha dd/mm/yyyy  <b>AppliedTo:</b>  Atributo «PublisherDate» en la tabla NL_Books y atributo «review_time» en al tabla NL_Ratings</p>