



# Presentación final grupo 11

Calidad de datos e Información

---

Bruno Ottonelli	-	4.954.242-1
Gabriel Rode	-	4.535.978-1

Facultad de Ingeniería  
Universidad de la República.



FACULTAD DE  
INGENIERÍA



UNIVERSIDAD  
DE LA REPÚBLICA  
URUGUAY

**Junio 2025**



# FASE 1

## Data Quality Planning

# Entradas y salidas

Entradas y salidas	
Entradas	Salidas
	Base de datos integrada (Data at hand) (2.1.5)
	Reporte de análisis de datos (2.2.2)
	Reporte de análisis de requerimientos de usuarios (2.3.2)
	Reporte con problemas de CD (Tabla 2.7)
	Modelo de contexto (2.3.3)

CD: Calidad de datos.

# Datos recibidos

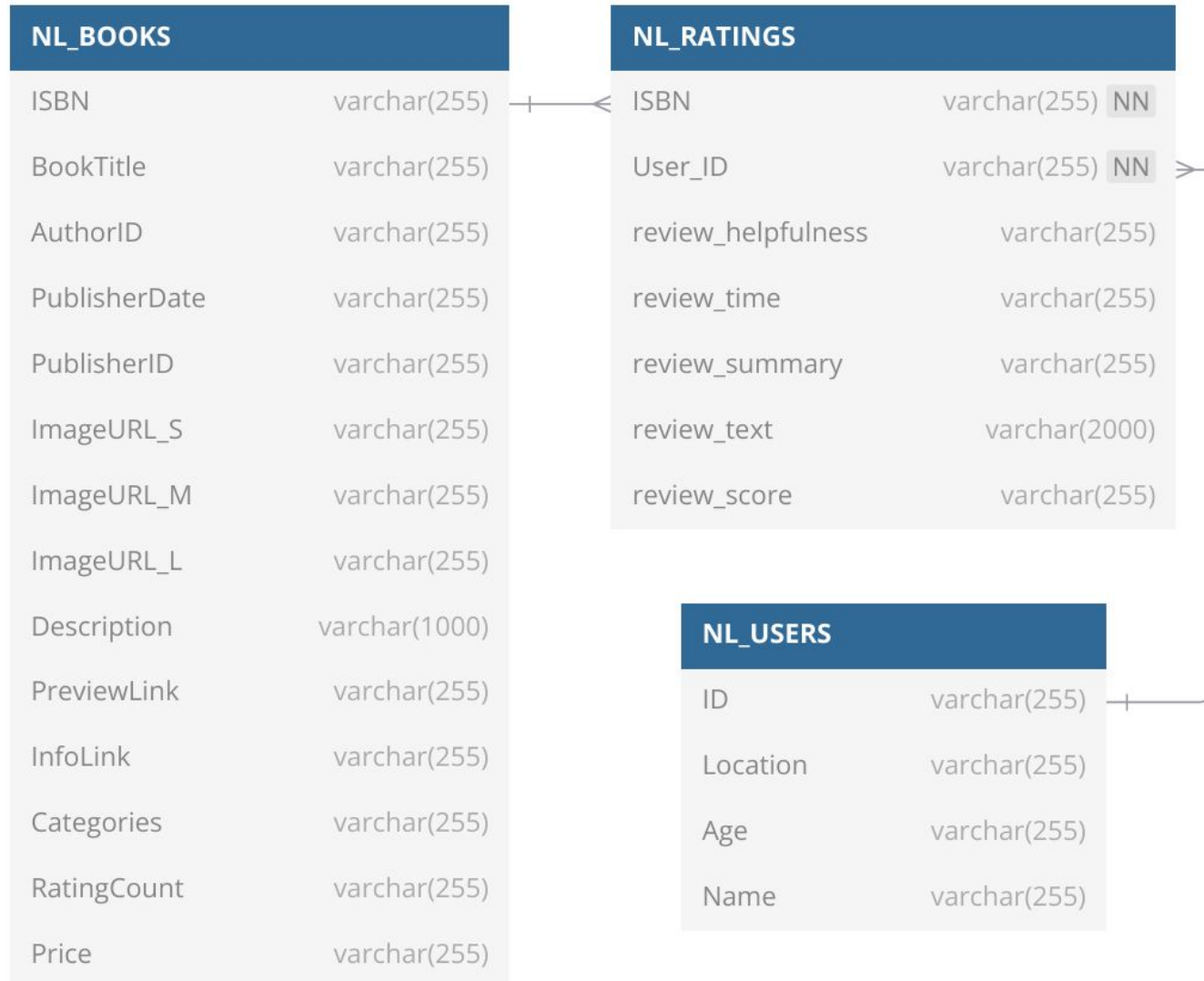
Tablas de la librería L1	
Nombre de tabla	Atributos
Books_rating	'Id', 'Title', 'Price', 'User_id', 'profileName', 'review/helpfulness', 'review/score', 'review/time', 'review/summary', 'review/text'
books_data	'Title', 'description', 'authors', 'image', 'previewLink', 'publisher', 'publishedDate', 'infoLink', 'categories', 'ratingsCount'

Tablas de la librería L2	
Nombre de tabla	Atributos
Books	'ISBN', 'Book-Title', 'Book-Author', 'Year-Of-Publication', 'Publisher', 'Image-URL-S', 'Image-URL-M', 'Image-URL-L'
Users	'User-ID', 'Location', 'Age'
Ratings	'User-ID', 'ISBN', 'Book-Rating'

# Contexto

Componentes de Contexto	
<b>Dominio</b>	D: Libros.
<b>Fuentes de datos</b>	Datos obtenidos de ambas librerías que se fusionarán y los proporcionados por el cliente sobre sus realidades.
<b>Tipos de usuario</b>	U1: Administrador. U2: Publicista digital. U3: Analista de datos.
<b>Tareas</b>	T1: Gestión. T2: Análisis. T3: Consulta.
<b>Reglas de negocio</b>	RN1: Cada libro deberá tener asociado un ISBN, un título, al menos un autor y un editor.
<b>Requerimientos de calidad</b>	RQ1: Frescura de datos: la base debe actualizarse todos los viernes. RQ2: Al menos el 80 % de los usuarios que califican los libros deben ser mayores de 18 años. RQ3: Al menos el 95 % de los libros deben cumplir simultáneamente con los siguientes requisitos: contar con un ISBN, tener el título correctamente escrito y que el nombre del autor incluya al menos un nombre y un apellido. RQ4: Al menos el 60 % de los libros tengan al menos un score mayor o igual a 5. RQ5: La librería pretende tener al menos 500 libros y poseer al menos el 20 % de la lista de los 100 mejores libros de Goodreads. RQ6: Los libros deben contar con fecha de publicación. RQ7: Los libros deben tener editorial. RQ8: Los libros deben tener asignado un valor de score.
<b>Requerimientos del sistema</b>	RS1: Los tiempos de respuesta del sitio Web de la NL no pueden superar los 3 segundos.
<b>Problemas de calidad ya reportados</b>	Ninguno en particular.
<b>Necesidades de filtrado</b>	F1: Libros por fecha (en particular, del año actual). F2: Libros por editorial. F3: Top de libros según su score.

# Data at hand





# Data profiling y Problemas de calidad

Problemas detectados en la Calidad de los Datos	
Campo	Problema de calidad
NL_Books.ISBN	P1: Entradas no respetan el formato ISBN ya que algunos están en formato ASID.
NL_Books.PublisherDate	P2: Las fechas tienen distinto formato.
NL_Books.PublisherID	P3: Mismo publisher escrito de forma distinta.
NL_Books.Title	P4: Títulos mal escritos.
NL_Books.AuthorID	P5: Mismo autor escrito de forma distinta.
NL_USERS.Age	P6: Valores de edad poco coherentes (por ejemplo 0).
NL_USERS.Location	P7: Ciudades mal escritas.
NL_RATINGS.review_time	P8: Formato de fecha/hora inconsistente.
NL_RATINGS.review_score	P9: Valores no numéricos en la puntuación.
NL_RATINGS.review_score	P10: Los valores importados entre L1 y L2 manejan distintas escalas (L1 puntúa de 0 a 5 y L2 de 0 a 10).
Base de datos	P11: Gran cantidad de nulos en muchos de los atributos.
NL_Books	P12: No hay campo de rating promedio del libro (podría calcularse).
NL_RATINGS.Helpfulness	P13: Este atributo en realidad debería ser dos atributos diferentes: cantidad de votaciones en esa review y cantidad de votaciones que consideraron útil esa review.
NL_Books.AuthorID	P14: Libros indican autores de forma distinta cuando tienen más de uno (dos libros con autores {A,B} y {B,A} deben considerarse con los mismos autores).
NL_RATINGS.review_score	P15: Podría haber valores fuera del rango 0 a 10.

# Nuevas componentes de contexto

Nuevas componentes de Contexto	
Reglas de negocio	<p>RN2: El atributo <i>ISBN</i> en <i>NL_Books</i> debe ser único a cada libro.</p> <p>RN3: El atributo <i>Price</i> en <i>NL_Books</i> debe ser un real positivo.</p> <p>RN4: El atributo <i>Age</i> en <i>NL_Users</i> debe ser un entero positivo.</p> <p>RN5: El atributo <i>ID</i> en <i>NL_Users</i> debe ser único y no vacío.</p>

Nuevas componentes de Contexto	
Requerimientos de calidad	<p>RQ9: los nombres de las editoriales deben estar estandarizados.</p> <p>RQ10: Las reglas de formato para nombres (autores, libros, editoriales) son: primera letra del nombre propio en mayúsculas y sin punto al final.</p> <p>RQ11: El formato para las fechas será dd/mm/aaaa.</p>





# FASE 2

## Data Quality Assessment

# Entradas y salidas

Entradas y salidas	
Entradas	Salidas
Reporte del análisis de requerimientos de usuarios (2.3.2)	Reporte de problemas de CD priorizados (3.1.2)
Reporte del análisis de datos (2.2.2)	Modelo de CD (3.1.4)
Reporte de problemas de CD (Tabla 2.7)	Especificación de la BD de metadatos de CD (3.2.2)
Modelo de Contexto (2.3.3)	Reporte de medición de la CD (3.2.3)
	Modelo de contexto (3.2.4)
	Reporte de evaluación de CD (3.2.3)

CD: Calidad de datos.

BD: Base de datos.



DIMENSIONES

**Exactitud**

**Compleitud**

**Unicidad**

**Consistencia**



Exactitud

**check\_isbn**

**check\_edades**

**check\_price**

**duplicate\_author**

**date\_format**



Consistencia

**consistencia\_ratings**  
**consistencia\_fechas**  
**missing\_rating\_books**





Completitud

**contar\_nulls**

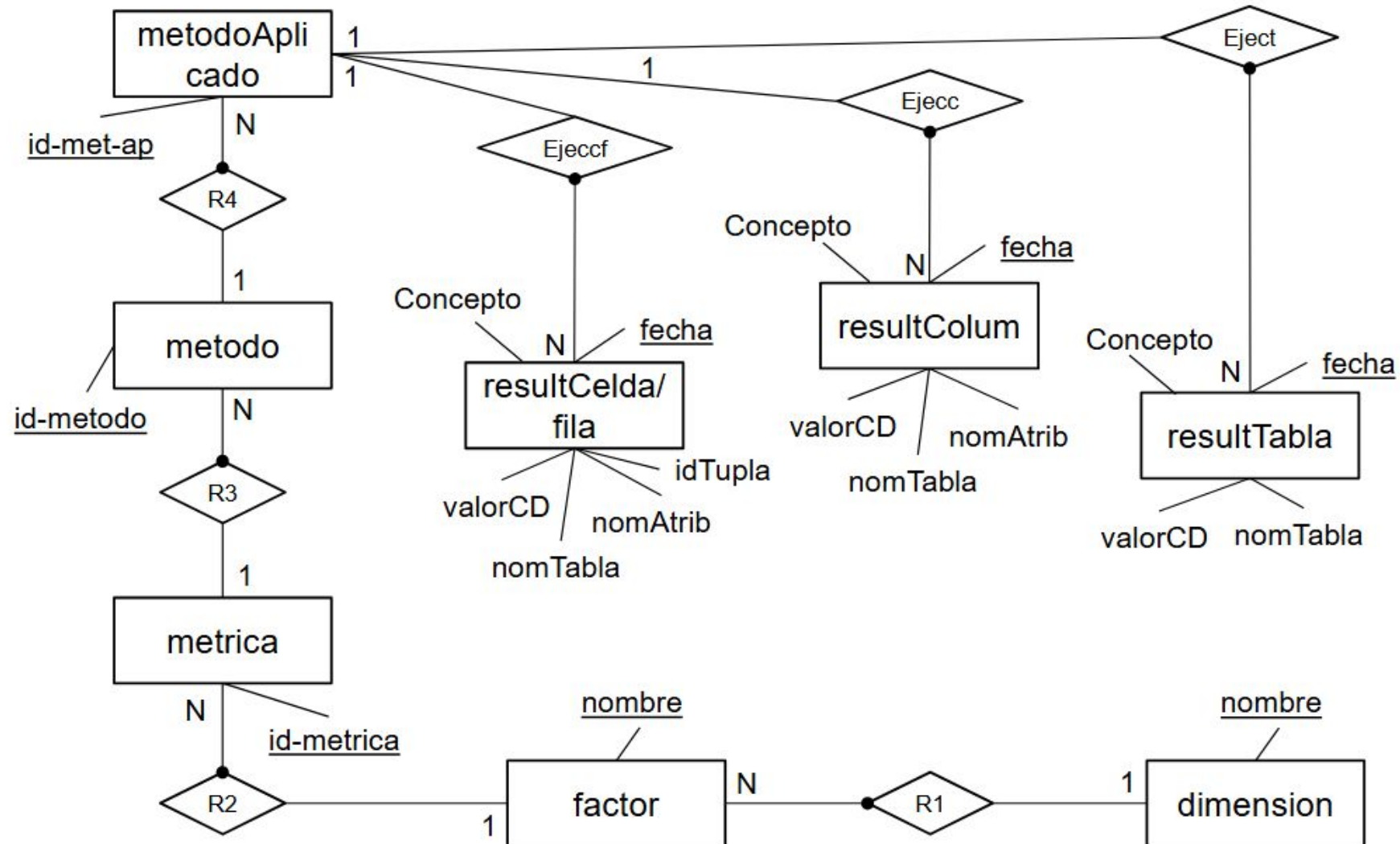
**check\_RN1**



Unicidad

**deduplicated\_ratio**

# Base de datos de metadatos de la calidad de datos



Criterios para los umbrales de evaluación

Tabla 3.1: Criterio estándar		
Concepto	Rango (directo)	Rango (inverso)
Malo	[0, 0.30)	[0.70, 1)
Bueno	[0.30, 0.60)	[0.40, 0.70)
Muy bueno	[0.60, 0.90)	[0.10, 0.40)
Excelente	[0.90, 1]	[0, 0.10)

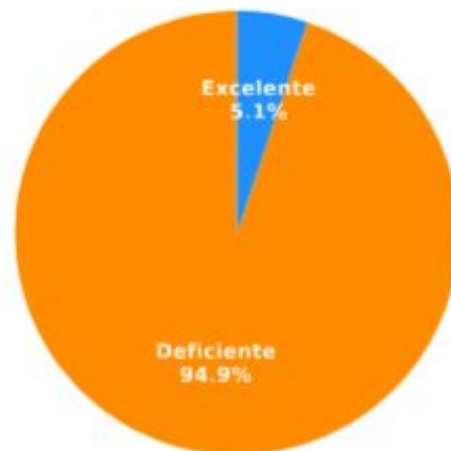
Tabla 3.2: Criterio estándar estricto		
Concepto	Directo	Inverso
Deficiente	(0, 0.30]	[0.70, 1]
Malo	(0.30, 0.60]	[0.40, 0.70)
Aceptable	(0.60, 0.75]	[0.25, 0.40)
Bueno	(0.75, 0.90]	[0.10, 0.25)
Excelente	(0.90, 1]	[0, 0.10)

Tabla 3.3: Criterio binario	
Concepto	Valor
Deficiente	0
Excelente	1

MA\_consistencia\_ratings  
NL\_BOOKS.ISBN



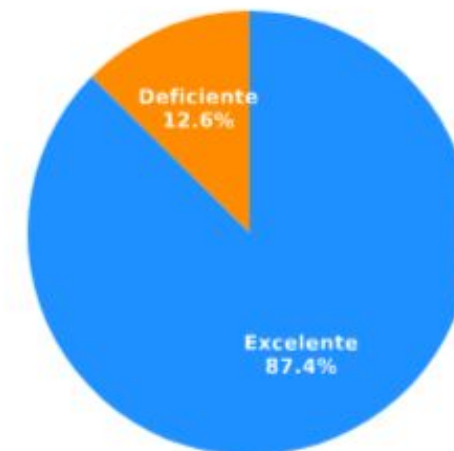
MA\_check\_edades  
NL\_USERS.Age



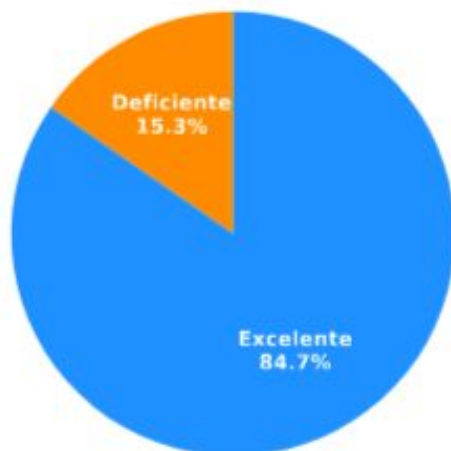
MA\_consistencia\_fechas  
NL\_BOOKS.ISBN



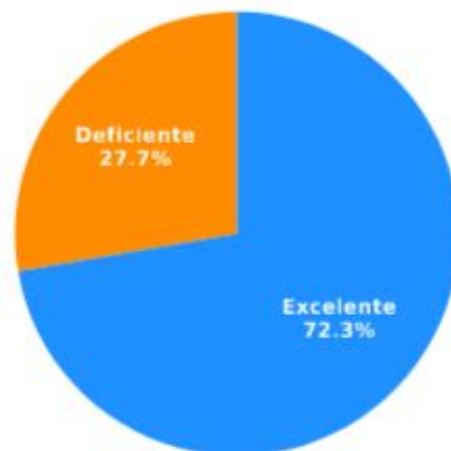
MA\_missing\_rating\_books  
NL\_BOOKS.ISBN



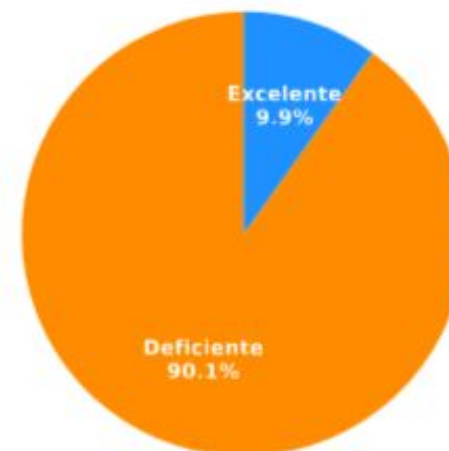
MA\_check\_ISBN  
NL\_BOOKS.ISBN



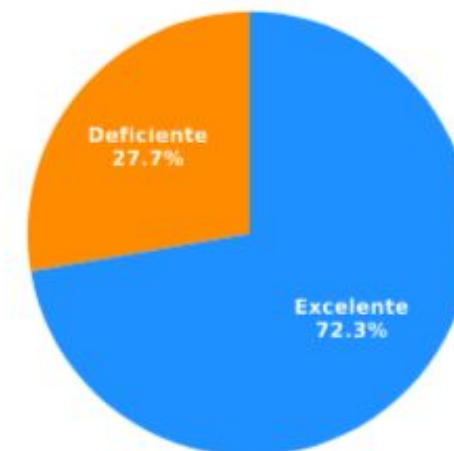
MA\_date\_format  
NL\_RATINGS.review\_time



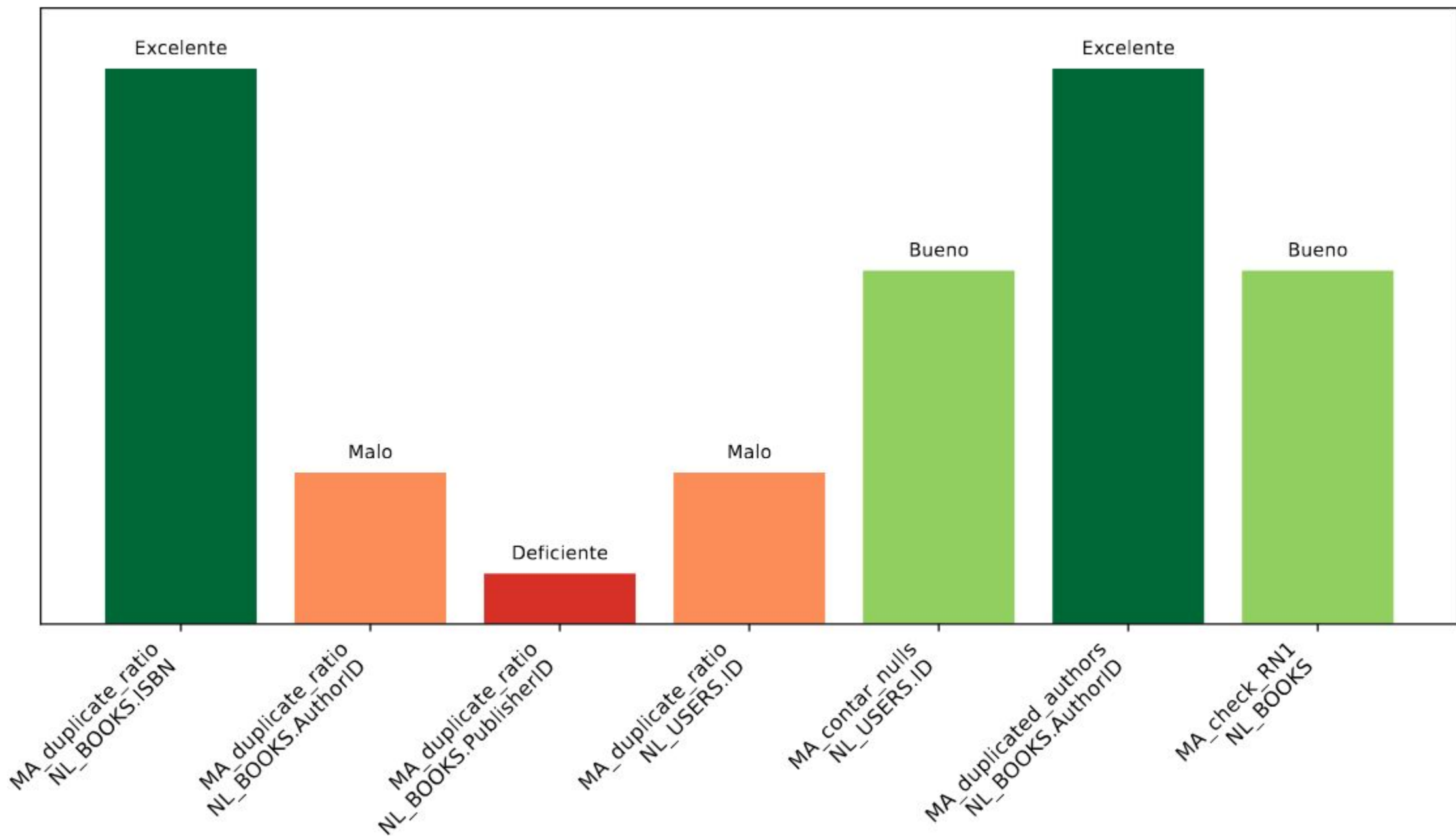
MA\_check\_price  
NL\_BOOKS.Price



MA\_date\_format  
NL\_BOOKS.PublisherDate









# FASE 3

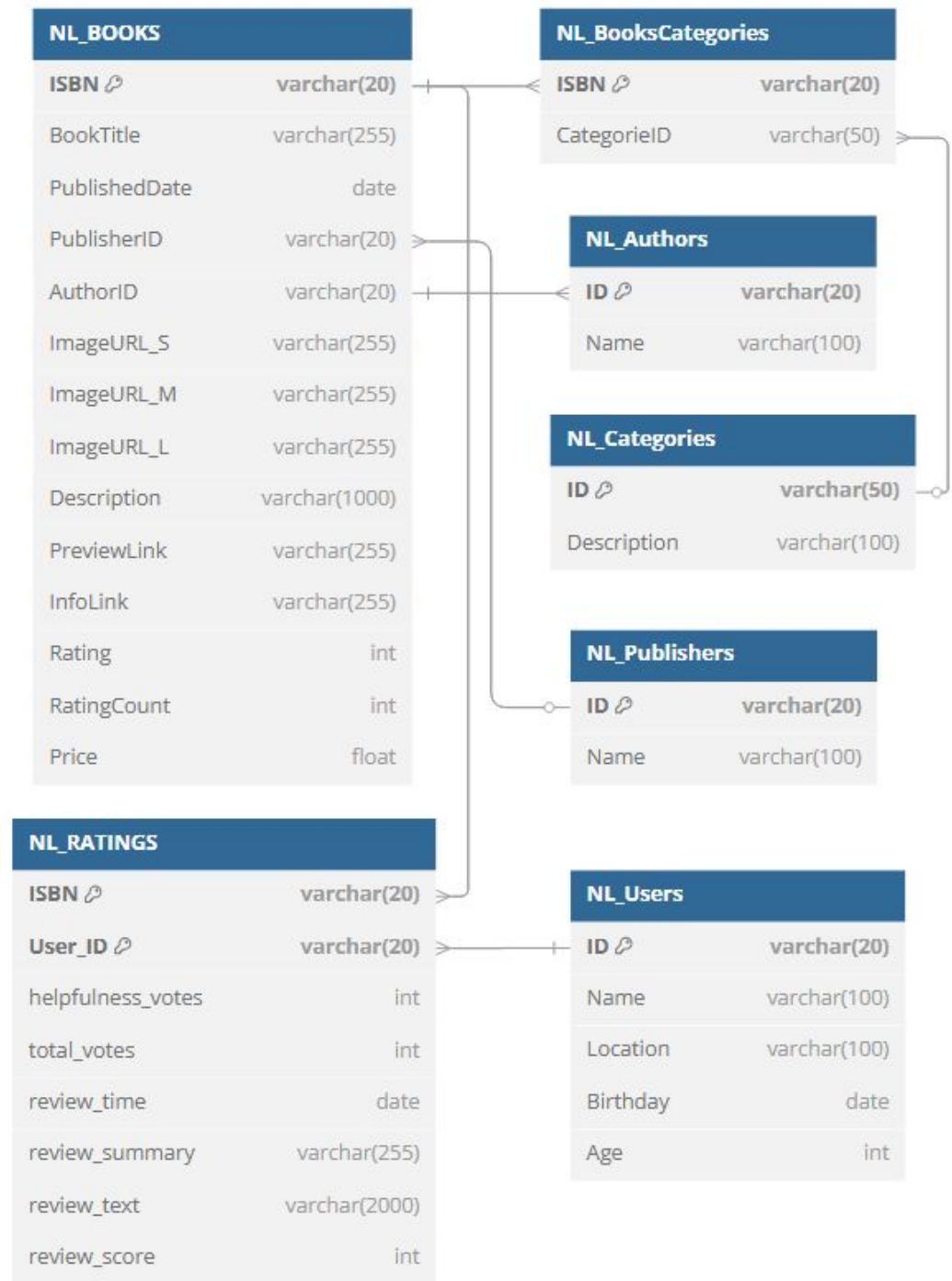
## Data Quality Improvement



## ANÁLISIS DE CAUSAS

- Dos bases de datos distintas
- Mal diseño de BD originales
- Errores de tipeo

# PLAN DE MEJORAS



# Muchas gracias



FACULTAD DE  
INGENIERÍA



UNIVERSIDAD  
DE LA REPÚBLICA  
URUGUAY