

# Calidad de Datos

**Yubar Daniel Marín Benjumea**

Estadístico

M.Sc en Ingeniería de Sistemas

Febrero 23 de 2022



Facultad de Minas  
Sede Medellín



UNIVERSIDAD  
**NACIONAL**  
DE COLOMBIA

# ¿Por qué es importante la calidad de datos?

Un estudio realizado por IBM en 2016 estimó como sólo en Estados Unidos, se generan pérdidas de 3 billones de dólares por baja calidad de los conjuntos de datos (Redman, 2016) .

Un estudio realizado en *The KPMG 2017 Global CEO Outlook* que estimó como el 56% de los directores ejecutivos en las empresas se encuentran preocupados por la integridad de los datos que están usando en la toma de decisiones (KPMG, 2017).

Se ha demostrado como es común encontrar porcentajes entre 10% y un 50% de datos erróneos en los conjuntos de datos (Johnson, Leitch y Neter, 1981; Laudon, 1986; Morey, 1982) .

Data Warehouse necesarios para implementar procesos de BDA se consume entre el 30% y el 80% del tiempo de desarrollo y presupuesto buscando mejora de la calidad en los datos

# ¿Qué genera la falta de calidad en los datos?

- Sistemas antiguos.
- Mal diseño de la interface de recolección de datos.
- Falta de capacitación.
- Cambios en los procesos.
- Reglas comerciales obsoletas.
- Falta de gobierno de datos.
- Falta de cultura del dato.
- Entre otros

La calidad de datos es un proceso dentro del gobierno de datos no todo el gobierno de datos en si.

Las investigaciones indican que muchos problemas de calidad de datos son causadas por falta de compromiso de la organización con los datos de alta calidad

# ¿Por qué garantizar la calidad en los datos?

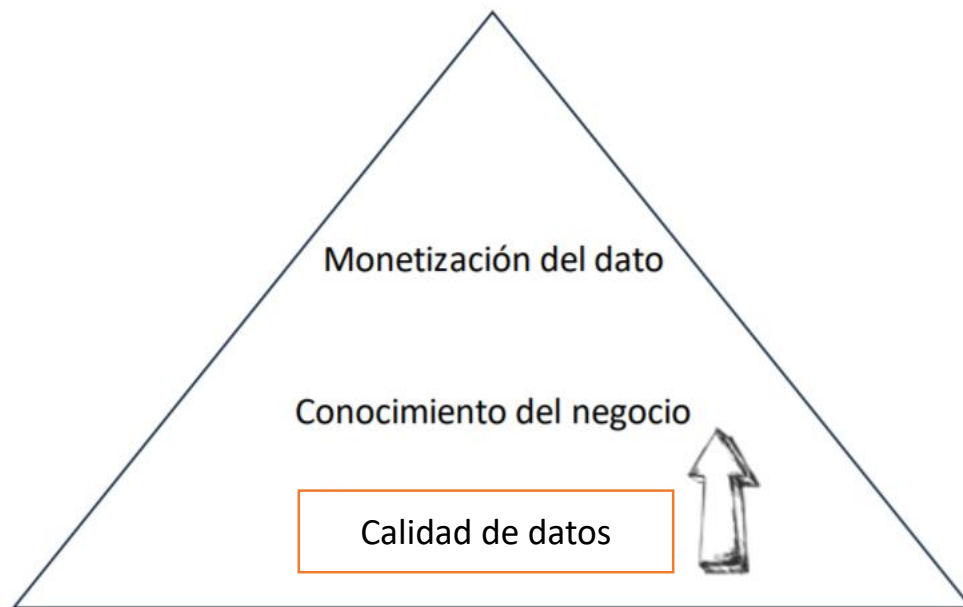


Figura 1. Gestión de los datos en una organización

# Calidad de datos

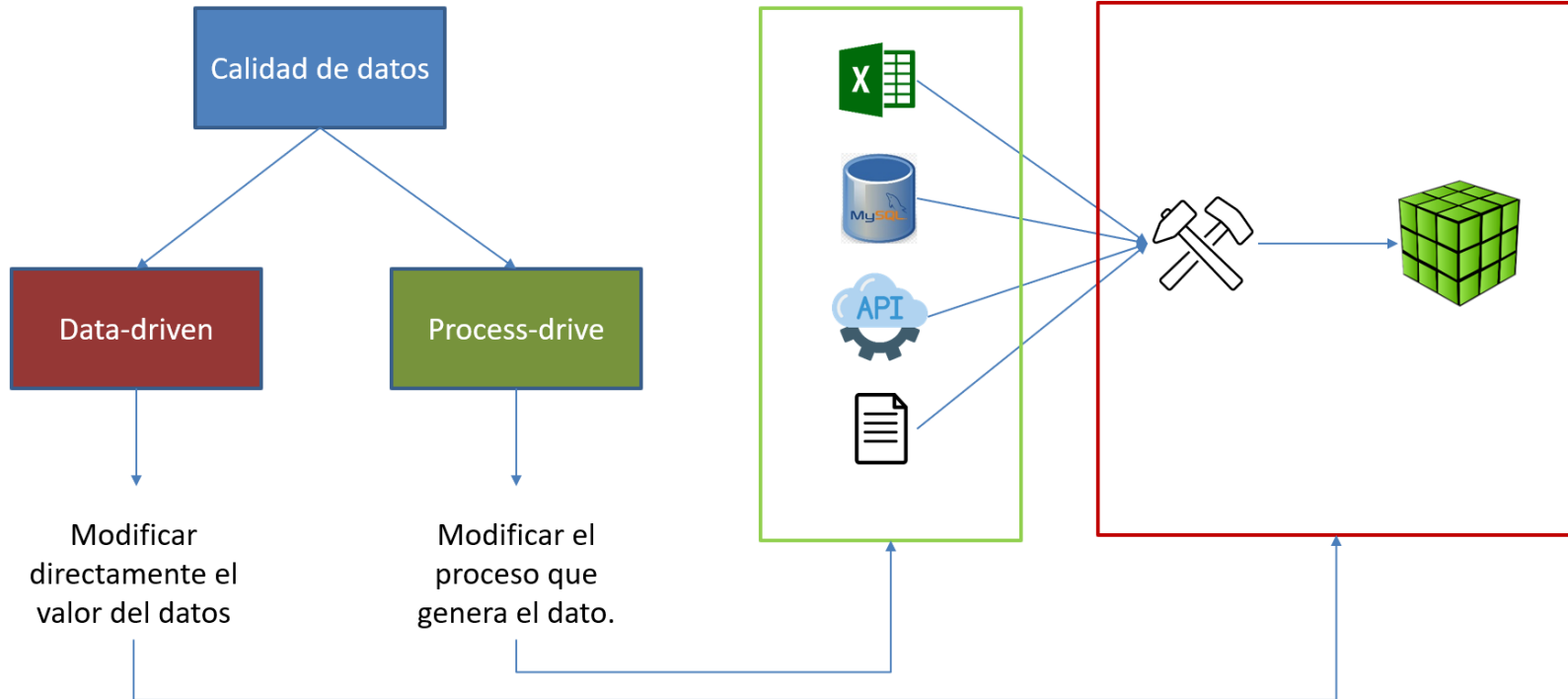


Facultad de Minas  
Sede Medellín



La calidad de datos se refiere tanto a las **características asociadas con los datos de alta calidad como a los procesos utilizados para medir o mejorar la calidad de los datos** (Mosley, Brackett, Earley y Henderson, 2010),

# Estrategias de calidad de datos



# Estrategias de calidad de datos



Facultad de Minas  
Sede Medellín



## Data driven

Muy útil a la hora de implementar procesos de:

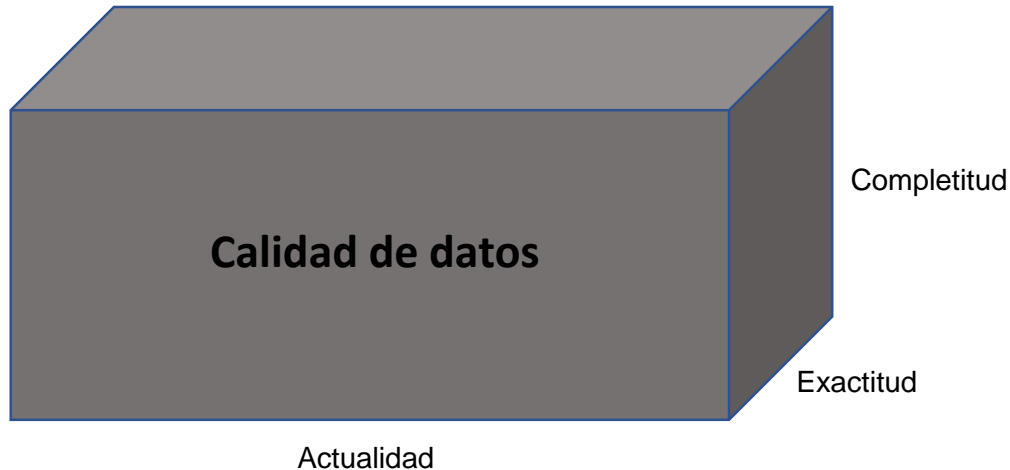
- Gestión de datos maestros.
- Data Warehouse.
- Data Mart.
- Data Lake

## Process driven

Muy útil a la hora de implementar procesos de:

- Estructurar procesos de datos desde cero.
- Arreglar problemas de datos desde la captura.
- Solucionar problemas de calidad de datos de raíz.

# Dimensiones de calidad de datos



## En la literatura

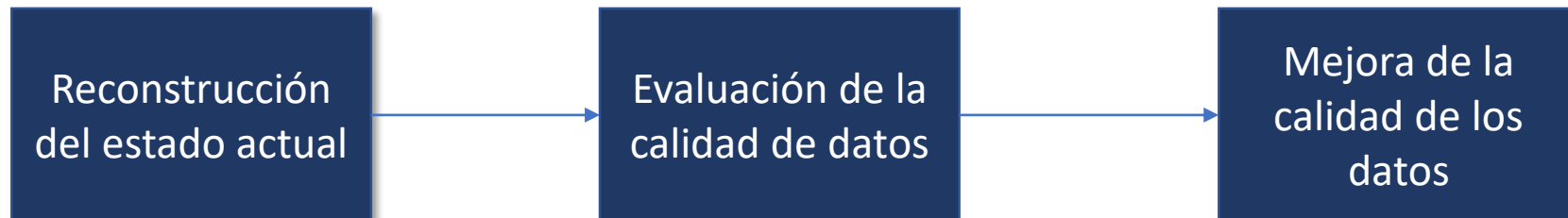
Se han mapeado 26 dimensiones (Wand y Wang, 1996) . Incluso la metodologías como COLDQ se usan 35. Lo que puede ser difícil de manejar.

## Nuestra propuesta

Exactitud, actualidad, completitud, consistencia y accesibilidad

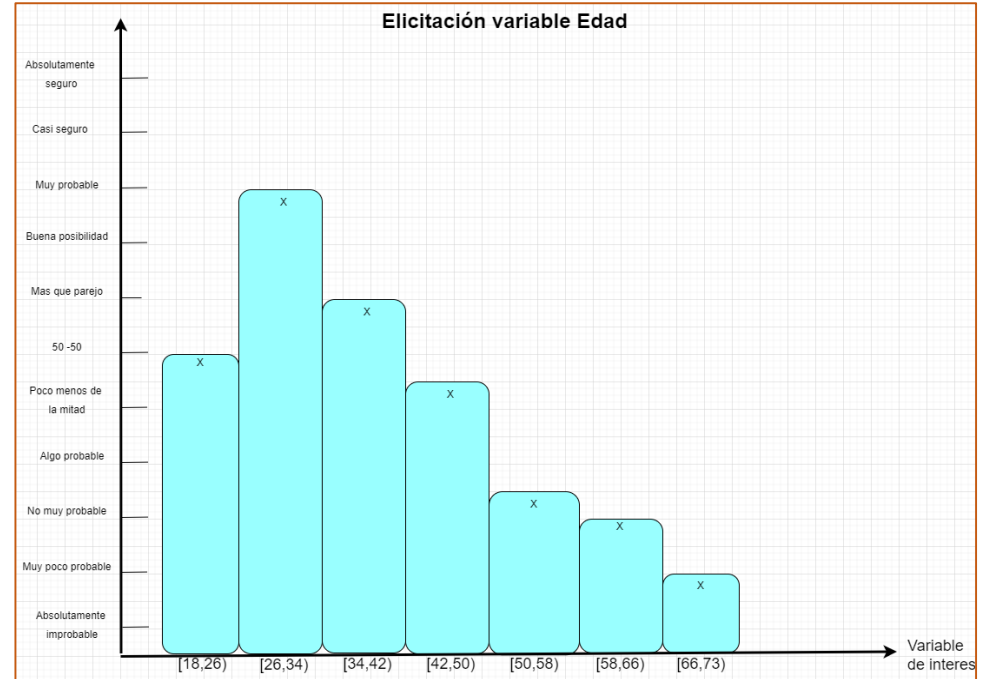


# Etapas en la calidad de datos



# Reconstrucción del estado actual

- Máximos.
- Mínimos.
- Distribuciones teóricas.
- Porcentaje admisible de nulos.
- Periodos de actualización.
- Formas de acceder.
- Reglas de negocio.
- Niveles teóricos.
- Entre otros metadatos.



# Evaluación de la calidad o perfilamiento de datos

Definir métricas para evaluar las dimensiones

$$\%Exactitud = \left( I * 0.3 + \left( \left( \sum_{i=1}^n \frac{NR_i}{R_i} \right) / TR \right) * 0.7 \right) * 100\%$$

$$\%Exactitud = \left( 0.1 * \frac{FMA}{n} + 0.1 * \frac{FMI}{n} + I * 0.3 + \left( \left( \sum_{i=1}^n \frac{NR_i}{R_i} \right) / TR \right) * 0.5 \right) * 100\%$$

# Mejora de la calidad en los datos

## Estándar ideal

```
Estándar =  
{nomenclatura_principal}  
{numero_principal} {letra} #  
{numero_secundario} -  
{numero_puerta}
```

## Dirección cruda

Calle 52sur N? 24A - 35

## Normalizar dirección

CALLE 52 SUR N ? 24 A -  
35

## Separar tokens

[CALLE, 52, SUR, N, ?,  
24, A, -, 35]

## Asignación de token a dominio

Similaridades(CALLE,  
Categorías),  
Similaridades(SUR,  
Categorías)

## Reemplazar equivalencia

[CL., 52, SUR, N, ?, 24,  
A, -, 35]

## Rellenar formato

CL. 52 A # 24 - 35

## Equivalencias teóricas

Categorías = ('CL.', 'CRA.',  
'DIG.', 'TRAV.')

## Estandarización aceptable

CL. 52 A # 24 - 35

## Estandarización no aceptable

Calle 52sur N? 24A - 35

## Si es incompleta o no es dirección

nulo

# GRACIAS