

# Gobierno de Datos – Eco-Moda S.A.S.

Este documento establece:

1. Roles y responsabilidades sobre la gestión de los datos
2. Reglas de calidad aplicadas en el pipeline de limpieza ya implementado
3. Estándares de nomenclatura y formato para garantizar consistencia

El gobierno de datos se aplica durante todo el ciclo de vida de los datos, desde su ingesta hasta su consumo en análisis y modelos predictivos. Aunque actualmente se trabaja con datos históricos, las políticas aquí establecidas están diseñadas para ser aplicables también cuando se integren fuentes de datos en tiempo real.

## 1. DEFINICIÓN DE ROLES Y RESPONSABILIDADES

El gobierno de datos requiere la asignación clara de roles para garantizar la calidad, seguridad y trazabilidad de la información. A continuación, se definen los roles necesarios para Eco-Moda.

### 1.1 Chief Data Officer (CDO) / Líder de Datos

**Perfil sugerido:** Director de Transformación Digital de Eco-Moda

**Responsabilidades principales:**

- Aprobar y actualizar las políticas del gobierno de datos
- Resolver conflictos entre áreas respecto a definiciones y uso de datos
- Priorizar iniciativas de mejora de calidad de datos
- Reportar métricas de calidad y gobierno a la dirección ejecutiva
- Asegurar los recursos necesarios para la implementación del gobierno

**Autoridad:** Toma decisiones finales sobre estándares de datos y políticas corporativas relacionadas.

### 1.2 Data Steward (Custodio de Datos)

Los Data Stewards son expertos del negocio responsables de la calidad y el significado de los datos en su dominio específico. No requieren ser técnicos, pero deben tener conocimiento profundo del área de negocio que representan.

Para Eco-Moda se proponen los siguientes Data Stewards:

#### 1.2.1 Data Steward - Área Comercial

**Perfil sugerido:** Gerente de Ventas o Jefe de Analítica Comercial

**Responsabilidades:**

- Definir el significado de negocio de las categorías de productos (product\_group\_name, garment\_group\_name)
- Establecer los valores válidos para clasificaciones de productos
- Definir rangos de precio aceptables según el portafolio de Eco-Moda
- Aprobar reglas de estandarización de categorías inconsistentes
- Validar la coherencia de datos de ventas y productos

### **1.2.2 Data Steward - Área de Clientes**

**Perfil sugerido:** Gerente de CRM o Marketing

**Responsabilidades:**

- Definir segmentos de clientes (club\_member\_status, fashion\_news\_frequency)
- Establecer políticas de privacidad para datos.
- Validar la coherencia de preferencias y comportamientos de clientes
- Definir métricas de lealtad y retención

### **1.3 Data Engineer (Ingeniero de Datos)**

**Equipo:** 2 ingenieros especializados

**Responsabilidades:**

- Implementar pipelines de ingestión, transformación y limpieza de datos
- Traducir reglas de negocio definidas por Data Stewards a código ejecutable
- Desarrollar y mantener validaciones automáticas de calidad
- Documentar procesos y transformaciones aplicadas
- Configurar alertas cuando se detecten problemas de calidad
- Resolver incidentes técnicos en los pipelines

**Herramientas utilizadas:** Python, SQL, Pandas, bibliotecas de validación de datos

### **1.4 Data Analyst / Data Scientist**

**Equipo:** 1 analista o científico de datos

**Responsabilidades:**

- Consumir datos limpios para análisis exploratorios y modelos predictivos
- Reportar problemas de calidad detectados durante el análisis
- Colaborar con Data Engineers en la definición de características (features) para modelos de ML
- Validar la coherencia de métricas en dashboards y reportes
- Documentar modelos predictivos desarrollados

### **1.5 Usuarios de Negocio**

**Perfil:** Equipos de Marketing, Operaciones, Finanzas, Logística

**Responsabilidades:**

- Consumir dashboards y reportes según sus permisos de acceso
- Reportar inconsistencias o dudas sobre los datos
- Respetar políticas de acceso y uso de información sensible
- Solicitar nuevos análisis a través de los Data Stewards

**Derechos:**

- Acceso a dashboards corporativos (OKR/BSC)
- Consulta del catálogo de datos para entender métricas
- Soporte de analistas para interpretación de resultados

## 2. REGLAS DE CALIDAD DE DATOS

Las reglas de calidad definen los criterios que los datos deben cumplir para ser considerados válidos y confiables.

### 2.1 Reglas de Completitud

El dataset presenta una calidad excepcional de completitud, con los campos críticos para el negocio 100% completos. Las reglas de completitud se enfocan en mantener este estándar y gestionar los casos minoritarios de valores faltantes.

#### 2.1.1 Campos Completos

Los siguientes campos son críticos para el negocio y se validó que están completamente poblados:

Campo	Justificación de criticidad
transaction_date	Esencial para análisis temporal y tendencias
price	Define el valor de la transacción
product_group_name	Clasificación principal de productos
age	Variable clave para segmentación demográfica

**Política:** Estos campos deben mantenerse obligatorios en futuras integraciones de datos. Cualquier registro sin estos valores debe ser rechazado en la ingestión.

#### 2.1.2 Campos con Valores Faltantes Mínimos

Los siguientes campos presentan valores nulos en menos del 0.5% de los registros:

Campo	Acción de imputación
-------	----------------------

club_member_status	Imputar como "UNKNOWN"
fashion_news_frequency	Imputar como "NONE"
Campos de color y apariencia	Imputar como "Unknown"

#### Justificación de imputaciones:

- "**UNKNOWN**" para **membresía**: Representa explícitamente clientes sin estado de membresía registrado, permitiendo análisis diferenciado de este segmento
- "**NONE**" para **newsletter**: Asume que sin preferencia explícita, el cliente no recibe comunicaciones
- "**Unknown**" para **colores**: Mantiene consistencia en análisis de preferencias cromáticas

### 2.2 Reglas de Unicidad

Para el dataset, dos registros se considerarían duplicados si tienen valores idénticos en los siguientes campos:

- transaction\_date
- product\_group\_name
- garment\_group\_name
- price
- age
- club\_member\_status

**Política:** En futuras integraciones de datos, se implementará un proceso automático de detección y eliminación de registros duplicados.

### 2.3 Reglas de Rangos Válidos

Acá se establecen límites lógicos del negocio para variables numéricas.

#### 2.3.1 Rango de Precio

**Regla:**  $0.0006 \text{ USD} \leq \text{price} \leq 0.507 \text{ USD}$

**Justificación basada en EDA:**

- **Límite inferior (0.0006)**: Valor mínimo detectado en el dataset real
- **Límite superior (0.507)**: Valor máximo detectado en el dataset real
- **Distribución observada**: La mayoría de transacciones (>80%) están concentradas en el rango 0.01-0.10 USD

**Acción si no cumple:** Rechazar registro y documentar en tabla de excepciones

**Nota:** Los límites se establecieron basándose en los valores reales observados. Para fases futuras, se recomienda revisar estos límites con el Data Steward Comercial para validar si reflejan el portafolio de productos de Eco-Moda.

#### 2.3.2 Rango de Edad

**Regla:**  $16 \leq \text{age} \leq 99$

#### **Justificación basada en EDA:**

- **Límite inferior (16):** Clientes menores de 16 años compran acompañados por adultos y no se registran individualmente
- **Límite superior (99):** Valor máximo detectado en el dataset real
- **Distribución observada:** Mayor concentración en rangos 20-40 años (media: 36.89 años)

**Acción si no cumple:** Rechazar registro

**Nota:** Se detectaron registros con edades extremas (99 años). Se recomienda al Data Steward de Clientes validar si estos registros son legítimos o requieren corrección en la fuente.

#### **2.3.3 Rango de Fechas**

**Regla:**  $\text{transaction\_date} \geq 2018-01-01$  y  $\text{transaction\_date} \leq \text{fecha actual}$

**Justificación:** Garantiza la validez temporal del dato dentro del rango histórico y hasta el momento presente.

**Acción si no cumple:** Rechazar registros con fechas anteriores a 2018 o futuras respecto a la fecha de carga.

### **3. ESTÁNDARES DE NOMENCLATURA Y FORMATO**

Los estándares garantizan consistencia en la estructura, nombres y formatos de los datos.

#### **3.1 Nomenclatura de Columnas**

**Convención adoptada:** snake\_case (minúsculas con guiones bajos)

**Principios:**

- Usar nombres descriptivos y completos (evitar abreviaturas ambiguas)
- Separar palabras con guion bajo (\_)
- Escribir completamente en minúsculas
- No utilizar caracteres especiales ni espacios

**Ejemplos correctos:**

- transaction\_date (fecha de transacción)
- product\_group\_name (grupo de producto)
- customer\_age (edad del cliente)
- club\_member\_status (estado de membresía)
- price\_usd (precio en dólares)

**Nota sobre el dataset actual:** El campo Transaction\_Date mantiene PascalCase por legado del dataset original. En futuras integraciones se estandarizará a transaction\_date.

#### **3.2 Formatos de Datos**

##### **3.2.1 Fechas**

**Estándar adoptado:** ISO 8601: YYYY-MM-DD

**Justificación:**

- Estándar internacional reconocido
- Evita ambigüedad entre formatos DD/MM/YYYY y MM/DD/YYYY
- Compatible con ordenamiento cronológico automático
- Facilita operaciones de bases de datos

### 3.2.2 Valores Monetarios

**Estándar adoptado:** ISO 4217 (USD como moneda base)

**Formato de precio:**

- **Moneda base:** USD (dólares americanos)
- **Tipo de dato:** Decimal (float)
- **Separador decimal:** Punto (.)
- **Precisión:** Exactamente 2 decimales

### 3.2.3 Variables Categóricas (Texto)

Capitalización según tipo:

Tipo de categoría	Formato	Ejemplos
Categorías de negocio (productos, estilos)	Title Case ( <i>Primera Letra Mayúscula</i> )	"T-Shirt", "Casual Wear", "Garment Upper Body"
Estados	UPPER_CASE ( <i>Todo en mayúsculas</i> )	"ACTIVE", "PRE-CREATE", "LEFT CLUB"
Colores	Title Case	"Black", "Dark Blue", "Light Grey"
Frecuencias / Preferencias	Title Case	"Regularly", "Monthly", "None"
Valores imputados / desconocidos	Mayúsculas o Title Case según contexto	"UNKNOWN", "NONE", "Unknown"

**Principios:**

- **No usar valores vacíos ("")** - siempre usar valores explícitos o NULL
- **Consistencia:** Usar siempre el mismo valor para el mismo significado
- **Legibilidad:** Preferir "UNKNOWN" a códigos numéricos como "-1" o "999"

## 4.3 Tipos de Datos

Asignación de tipos según naturaleza del campo:

Tipo de información	Tipo de dato sugerido
Fechas de transacción (transaction_date)	datetime64

<b>Precios / valores monetarios</b> (price)	float64
<b>Categorías de producto (negocio)</b> (product_group_name, garment_group_name)	category
<b>Atributos visuales de producto</b> (graphical_appearance_name, colour_group_name, perceived_colour_value_name, perceived_colour_master_name, index_name, index_group_name)	category
<b>Edad del cliente</b> (age)	int64
<b>Segmentación de clientes</b> (club_member_status, fashion_news_frequency)	category