

# Catálogo de Dados

## 1) gold.dim\_title

### Descrição

Dimensão que representa o **título** (entidade de conteúdo) do catálogo Disney+. Contém os atributos mais estáveis do conteúdo: identificador natural (show\_id), nome do título, ano de lançamento, duração e descrição. É a dimensão central para descrição de “o que é” o item de catálogo.

### Campos e domínio

- **title\_id** (BIGINT, PK)

**Descrição:** chave substituta (surrogate key) do título.

**Domínio esperado:** inteiro positivo,  $\geq 1$ . Não nulo. Único.

- **show\_id** (STRING, UK)

**Descrição:** identificador natural do dataset (ex.: “s1234”).

**Domínio esperado:** string não vazia. Único.

**Regra:** trim(show\_id)  $\neq$  ”.

- **title** (STRING)

**Descrição:** nome do título.

**Domínio esperado:** string não vazia.

**Regra:** trim(title)  $\neq$  ”.

- **release\_year** (INT)

**Descrição:** ano de lançamento do título.

**Domínio esperado:**  $1900 \leq release\_year \leq ano\_atual$  (ou  $\leq$  ano de extração do dataset).

**Regras:** não nulo para análises de lag; valores fora do intervalo indicam erro de origem.

- **duration** (STRING)

**Descrição:** duração textual (ex.: “90 min” para filmes; “1 Season” para séries).

**Domínio esperado:** string; pode ser nulo para as análises deste MVP.

**Validação (opcional):** regex/parse para padrões “min” ou “Season(s)”.

- **description** (STRING)

**Descrição:** sinopse do conteúdo.

**Domínio esperado:** texto livre; pode ser nulo para as análises desde MVP.

## **2) gold.dim\_date**

### **Descrição**

Dimensão de datas (calendário) para suportar análises temporais por dia/mês/ano. Neste MVP utilizada principalmente para **data de adição ao catálogo**.

### **Campos e domínio**

- **date\_id** (INT, PK)

**Descrição:** chave de data no formato yyyyymmdd.

**Domínio esperado:** >= 19000101 e <= 21001231 (faixa ampla). Não nulo.  
Único.

**Regras:** deve corresponder a full\_date\_added.

- **full\_date\_added** (DATE)

**Descrição:** data completa.

**Domínio esperado:** data válida. Não nula.

- **year** (INT)

**Descrição:** ano extraído de full\_date\_added.

**Domínio esperado:** 1900 <= year <= ano\_atual (ou <= ano de extração).

- **month** (INT)

**Descrição:** mês extraído.

**Domínio esperado:** 1..12.

- **day** (INT)

**Descrição:** dia do mês extraído.

**Domínio esperado:** 1..31 (consistência real depende do mês/ano).

## **3) gold.dim\_type**

### **Descrição**

Dimensão para classificar o tipo de conteúdo (filme ou série). Serve para filtros e segmentação de análises.

### **Campos e domínio**

- **type\_id** (INT, PK)

**Descrição:** surrogate key do tipo.

**Domínio esperado:** inteiro positivo >= 1. Único.

- **type\_name** (STRING)

**Descrição:** categoria textual do tipo.

**Domínio esperado (dataset Disney+):** tipicamente {'Movie', 'TV Show'}.

**Regra:** não nulo e não vazio.

#### 4) gold.dim\_rating

##### Descrição

Dimensão de classificação indicativa (rating). Usada para análise de distribuição do catálogo por faixa/classificação.

##### Campos e domínio

- **rating\_id** (INT, PK)

**Descrição:** surrogate key do rating.

**Domínio esperado:** inteiro positivo  $\geq 1$ . Único.

- **rating\_name** (STRING)

**Descrição:** rótulo da classificação indicativa (depende do país/standard do dataset).

**Domínio esperado:** conjunto finito de strings específicas (categóricas) - ("G", "PG", "PG-13", "R", "NC-17", "TV-Y", "TV-Y7", "TV-G", "TV-PG", "TV-14", "TV-MA" ou "Not Rated"). Esta lista válida de domínio esperado é definida pelos órgãos de classificação etária, neste dataset especificamente, temos a Motion Pictures Association (MPA) que faz a classificação para os filmes ([Homepage - MPA Film Ratings](#)) e a TV Parental Guidelines que faz a classificação do conteúdo de séries para a televisão.

**Regra:** não nulo e não vazio.

**Observação:** os valores exatos dependem do dataset (podem existir variações e valores "Not Rated").

#### 5) gold.dim\_genre

##### Descrição

Dimensão de gêneros (categorias) do conteúdo, originada do campo multivvalorado (listed\_in) normalizado na Silver.

##### Campos e domínio

- **genre\_id** (INT, PK)

**Descrição:** surrogate key do gênero.

**Domínio esperado:** inteiro positivo  $\geq 1$ . Único.

- **genre\_name** (STRING)

**Descrição:** nome do gênero/categoria (ex.: “Animation”, “Comedy”, “Family”, etc.).

**Domínio esperado:** string não vazia; conjunto finito (variável conforme dataset).

**Regra:** trim(genre\_name) <> ”.

## 6) gold.bridge\_title\_genre

### Descrição

Tabela de relacionamento **N:N** entre títulos e gêneros. Um título pode ter vários gêneros e um gênero pode estar em muitos títulos. Esta tabela viabiliza contagens por gênero, diversidade por período e análises de distribuição.

### Campos e domínio

- **title\_id** (BIGINT, FK → dim\_title.title\_id)

**Descrição:** referência ao título.

**Domínio esperado:** >= 1, não nulo; deve existir na dim\_title.

- **genre\_id** (INT, FK → dim\_genre.genre\_id)

**Descrição:** referência ao gênero.

**Domínio esperado:** >= 1, não nulo; deve existir na dim\_genre.

### Regras de qualidade

- **Unicidade composta:** (title\_id, genre\_id) deve ser único (sem duplicatas).
- Sem “órfãos”: todos os FKs devem existir nas dimensões.

## 7) gold.fact\_catalog\_addition

### Descrição

Tabela fato que representa o **evento de adição do título ao catálogo** (factless com métricas derivadas). Grão esperado: **1 linha por título adicionado** (por title\_id e data). Suporta análises de tempo de entrada (lag), evolução temporal e métricas agregadas.

### Campos e domínio

- **fact\_id** (BIGINT, PK)

**Descrição:** surrogate key da linha fato.

**Domínio esperado:** inteiro positivo >= 1. Único. Não nulo.

- **title\_id** (BIGINT, FK → dim\_title.title\_id)  
**Descrição:** título adicionado.  
**Domínio esperado:** >= 1, não nulo; deve existir na dimensão.
- **date\_id** (INT, FK → dim\_date.date\_id)  
**Descrição:** chave da data de adição.  
**Domínio esperado:** yyyyymmdd, não nulo (se a data existir); deve existir na dimensão.
- **type\_id** (INT, FK → dim\_type.type\_id)  
**Descrição:** tipo do conteúdo no momento da adição.  
**Domínio esperado:** >= 1, pode ser nulo apenas se não mapeado (não recomendado).
- **rating\_id** (INT, FK → dim\_rating.rating\_id)  
**Descrição:** classificação indicativa do conteúdo.  
**Domínio esperado:** >= 1, pode ser nulo se não mapeado (não recomendado).
- **full\_date\_added** (DATE)  
**Descrição:** data completa de adição ao catálogo (campo operacional para filtros).  
**Domínio esperado:** data válida; idealmente não nula para análises temporais.
- **release\_date\_approx** (DATE)  
**Descrição:** data aproximada de lançamento (ex.: YYYY-01-01) construída a partir de release\_year, usada para calcular lag\_days.  
**Domínio esperado:** data válida; pode ser nula se release\_year for nulo.
- **lag\_days** (INT)  
**Descrição:** diferença em dias entre full\_date\_added e release\_date\_approx.  
**Domínio esperado:** >= 0 na maioria dos casos.  
**Regras:** valores negativos indicam cenário incomum (release\_year maior que ano de adição ou erro de origem). Um cenário possível seria um filme ser adicionado à plataforma antes de ser lançado no cinema, caso a base esteja considerando data de lançamento este lançamento no cinema.  
**Faixa típica:** de 0 até 50.000 (considerando que o primeiro filme do cinema teria sido lançado em 1895-“A Saída dos Trabalhadores da Fábrica Lumière”).
- **title\_count** (INT)  
**Descrição:** métrica contadora (sempre 1) para facilitar agregações.  
**Domínio esperado:** exatamente 1.

## **Regras de qualidade**

- Sem nulos em: fact\_id, title\_id, full\_date\_added e lag\_days (para análises de lag).
- title\_count = 1.
- Integridade referencial com dimensões.
- Checagem de duplicidade do grão (recomendado): (title\_id, date\_id) deve ser único se a regra for “um título entra uma vez”.