

## **Scripts e processo de normalização do banco de dados**

Durante o processo de normalização da base de dados do Olist, utilizei principalmente o Power Query para identificar e eliminar redundâncias, além de reorganizar os dados em estruturas mais limpas, consistentes e alinhadas aos princípios da modelagem relacional. Abaixo explico o que foi feito em cada parte:

### **Tabela olist\_customers**

A tabela original de clientes armazenava repetidamente informações de localização (cidade e estado) junto com o cliente. Isso fazia com que muitos dados fossem duplicados desnecessariamente.

O que foi feito:

- Mantive na tabela customers apenas os campos: customer\_id, customer\_unique\_id e customer\_zip\_code\_prefix.
- Criei uma nova tabela chamada customers\_location, contendo os campos customer\_zip\_code\_prefix, customer\_city e customer\_state.
- Dessa forma, cada CEP passou a aparecer uma única vez na tabela de localização, evitando a repetição de cidade e estado em cada registro de cliente.

### **Tabela olist\_order\_items**

Originalmente, a tabela de itens do pedido armazenava múltiplas linhas para cada combinação de order\_id e product\_id, com os valores de price e freight\_value sendo repetidos por linha.

Normalização feita:

- Criei uma nova tabela order\_items contendo uma linha para cada par order\_id e product\_id.
- Agrupei os dados por pedido e produto e calculei:
  - qtd\_compra: total de itens daquele produto no pedido.
  - total\_price: soma dos preços.
  - Mantive também seller\_id, shipping\_limit\_date e order\_item\_id.

### **Além disso, criei uma nova tabela separada chamada freight\_value:**

- Nela, armazenei apenas o primeiro valor de frete (freight\_value) por pedido (order\_id), já que o valor do frete é o mesmo para todos os itens de um pedido.
- Isso eliminou repetições e somas incorretas do frete, tornando os dados mais precisos e organizados.

### **Tabela order\_payments**

A tabela de pagamentos possuía o campo payment\_type escrito como texto em cada linha, o que gerava redundância.

O que foi feito:

- Criei uma nova tabela payment\_types, que armazena os nomes únicos de tipos de pagamento com um identificador payment\_type\_id.
- Na tabela order\_payments, substituí o campo de texto pelo payment\_type\_id, referenciando a nova tabela.
- Isso deixou os dados mais organizados, evitando erros de digitação e facilitando análises por tipo de pagamento.

### **Tabela sellers**

Na tabela de vendedores (sellers), os campos de localização (seller\_city e seller\_state) eram repetidos diversas vezes, sempre que o seller\_zip\_code\_prefix se repetia.

Normalização aplicada:

- Mantive na tabela sellers apenas os campos seller\_id e seller\_zip\_code\_prefix.
- Criei uma nova tabela sellers\_location, com os campos seller\_zip\_code\_prefix, seller\_city e seller\_state, onde cada prefixo de CEP aparece apenas uma vez.

Essa separação eliminou a redundância e deixou os dados de localização dos vendedores mais bem organizados.

## **Resultado**

Com essa normalização:

- Reduzi a redundância nos dados.
- Organizei melhor os relacionamentos entre as entidades.
- Facilitei análises futuras.
- Melhorei a performance e a consistência da base.

Toda a normalização foi feita utilizando o Power Query e scripts SQL para refletir a estrutura final no banco PostgreSQL.