

DDL_Script_MAIN.sql (Data Architecture & Modeling)

Language: PostgreSQL

```
1  /*
2  =====
3  OLIST E-COMMERCE – FULL DATA MODEL DDL
4  CRM-Ready Segmentation Engine for PostgreSQL
5  =====
6
7  ENTITY-RELATIONSHIP MAP
8  =====
9
10 product_category_name_translation (standalone lookup – no FKs)
11     |
12     └─> olist_products_dataset.product_category_name (logical join)
13
14 olist_geolocation_dataset (standalone reference – no FKs)
15     |
16     └─> zip_code_prefix joins to customers & sellers (logical join)
17
18 olist_customers_dataset ← ROOT ENTITY (customer_id PK)
19     |
20     └─> olist_orders_dataset (customer_id FK)
21           |
22           ├──> olist_order_payments_dataset (order_id FK)
23           ├──> olist_order_reviews_dataset (order_id FK)
24           └─> olist_order_items_dataset (order_id FK)
25                 |
26                 ├──> olist_products_dataset (product_id FK)
27                 └─> olist_sellers_dataset (seller_id FK)
28
29 EXECUTION ORDER (respects FK dependency chain):
30 =====
31 1. product_category_name_translation (no dependencies)
32 2. olist_geolocation_dataset (no dependencies)
33 3. olist_customers_dataset (no dependencies)
34 4. olist_products_dataset (no dependencies)
35 5. olist_sellers_dataset (no dependencies)
36 6. olist_orders_dataset (depends on: customers)
37 7. olist_order_payments_dataset (depends on: orders)
38 8. olist_order_reviews_dataset (depends on: orders)
39 9. olist_order_items_dataset (depends on: orders, products, sellers)
40
41 =====
42 DDL SCRIPTS – EXECUTE IN THIS EXACT ORDER
43 =====
44 */
45
46
47 -- ┌──────────────────────────────────────────────────────────────────────────────────┐
48 -- || TABLE 1: PRODUCT CATEGORY NAME TRANSLATION (Lookup Table) || ───────────┘
```

```

49 -- || Dimension: Product Category Affinity ("What They Like") ||
50 -- || 71 rows – maps Portuguese category names to English ||
51 -- ||-----||
52
53 DROP TABLE IF EXISTS product_category_name_translation CASCADE;
54
55 CREATE TABLE product_category_name_translation (
56     product_category_name    VARCHAR(50)    NOT NULL,
57     product_category_name_english VARCHAR(50) NOT NULL,
58
59     CONSTRAINT pk_category_translation
60         PRIMARY KEY (product_category_name)
61 );
62
63 COMMENT ON TABLE product_category_name_translation
64     IS 'Lookup: Portuguese-to-English product category mapping. Joins to olist_products_dataset on product_category_name.';
65
66
67 -- ||-----||
68 -- || TABLE 2: GEOLOCATION (Reference/Dimension Table) ||
69 -- || Dimension: Geographic Demographics ("Where") ||
70 -- || 1,000,163 rows – multiple lat/lng entries per zip code prefix ||
71 -- || NOTE: No natural PK. zip_code_prefix has ~19K unique values across 1M ||
72 -- || rows (multiple geocoded points per zip). We add a surrogate PK. ||
73 -- ||-----||
74
75 DROP TABLE IF EXISTS olist_geolocation_dataset CASCADE;
76
77 CREATE TABLE olist_geolocation_dataset (
78     geolocation_id          SERIAL          NOT NULL,
79     geolocation_zip_code_prefix INTEGER      NOT NULL,
80     geolocation_lat         NUMERIC(10,7)   NOT NULL,
81     geolocation_lng         NUMERIC(10,7)   NOT NULL,
82     geolocation_city        VARCHAR(50)     NOT NULL,
83     geolocation_state       CHAR(2)        NOT NULL,
84
85     CONSTRAINT pk_geolocation
86         PRIMARY KEY (geolocation_id)
87 );
88
89 CREATE INDEX idx_geolocation_zip
90     ON olist_geolocation_dataset (geolocation_zip_code_prefix);
91
92 CREATE INDEX idx_geolocation_state
93     ON olist_geolocation_dataset (geolocation_state);
94
95 COMMENT ON TABLE olist_geolocation_dataset
96     IS 'Geocoded reference data. Multiple lat/lng points per zip code prefix. Join to customers/sellers on zip_code_prefix.';
97
98
99 -- ||-----||
100 -- || TABLE 3: CUSTOMERS (Root Entity) ||
101 -- || Dimension: Geographic Demographics + Core Identity ||

```

```

102 -- || 99,441 rows – one row per order-instance; customer_unique_id is the ||
103 -- || real human (96,096 distinct). customer_id is PK (1:1 with orders). ||
104 -- ||-----||
105
106 DROP TABLE IF EXISTS olist_customers_dataset CASCADE;
107
108 CREATE TABLE olist_customers_dataset (
109     customer_id          CHAR(32)          NOT NULL,
110     customer_unique_id   CHAR(32)          NOT NULL,
111     customer_zip_code_prefix  INTEGER          NOT NULL,
112     customer_city         VARCHAR(40)      NOT NULL,
113     customer_state       CHAR(2)          NOT NULL,
114
115     CONSTRAINT pk_customers
116         PRIMARY KEY (customer_id)
117 );
118
119 CREATE INDEX idx_customers_unique_id
120     ON olist_customers_dataset (customer_unique_id);
121
122 CREATE INDEX idx_customers_state
123     ON olist_customers_dataset (customer_state);
124
125 CREATE INDEX idx_customers_zip
126     ON olist_customers_dataset (customer_zip_code_prefix);
127
128 COMMENT ON TABLE olist_customers_dataset
129     IS 'Customer dimension. customer_id is per-order; customer_unique_id is the real person. Aggregate on customer_unique_id for RFM.';
130
131
132 -- ||-----||
133 -- || TABLE 4: PRODUCTS ||
134 -- || Dimension: Product Category Affinity + Price/Freight Sensitivity ||
135 -- || 32,951 rows – one row per unique product SKU ||
136 -- ||-----||
137
138 DROP TABLE IF EXISTS olist_products_dataset CASCADE;
139
140 CREATE TABLE olist_products_dataset (
141     product_id          CHAR(32)          NOT NULL,
142     product_category_name  VARCHAR(50),      -- 610 NULLs in source data
143     product_name_lenght   SMALLINT,        -- sic: original column name typo
144     product_description_lenght  SMALLINT,    -- sic: original column name typo
145     product_photos_qty    SMALLINT,
146     product_weight_g      INTEGER,          -- max 40,425g; 2 NULLs
147     product_length_cm     SMALLINT,
148     product_height_cm     SMALLINT,
149     product_width_cm      SMALLINT,
150
151     CONSTRAINT pk_products
152         PRIMARY KEY (product_id)
153 );
154

```



```

261 -- || Dimension: Customer Experience & Sentiment ("Why They Left") ||
262 -- || 99,224 rows – review_id is NOT unique (814 dupes across orders). ||
263 -- || Composite PK: (review_id, order_id) is unique across all 99,224 rows. ||
264 -- ||-----||
265
266 DROP TABLE IF EXISTS olist_order_reviews_dataset CASCADE;
267
268 CREATE TABLE olist_order_reviews_dataset (
269     review_id          CHAR(32)          NOT NULL,
270     order_id           CHAR(32)          NOT NULL,
271     review_score        SMALLINT         NOT NULL,    -- 1 to 5
272     review_comment_title  VARCHAR(50),    -- 87,656 NULLs
273     review_comment_message TEXT,        -- 58,247 NULLs; free-text, variable length
274     review_creation_date  TIMESTAMP      NOT NULL,
275     review_answer_timestamp  TIMESTAMP  NOT NULL,
276
277     CONSTRAINT pk_order_reviews
278         PRIMARY KEY (review_id, order_id),
279
280     CONSTRAINT fk_reviews_order
281         FOREIGN KEY (order_id)
282         REFERENCES olist_orders_dataset (order_id),
283
284     CONSTRAINT chk_review_score
285         CHECK (review_score BETWEEN 1 AND 5)
286 );
287
288 CREATE INDEX idx_reviews_order
289     ON olist_order_reviews_dataset (order_id);
290
291 CREATE INDEX idx_reviews_score
292     ON olist_order_reviews_dataset (review_score);
293
294 COMMENT ON TABLE olist_order_reviews_dataset
295     IS 'Review sentiment. Join At-Risk customers to their last review_score for service-recovery targeting. review_score 1-2 = negative.';
296
297
298 -- ||-----||
299 -- || TABLE 9: ORDER ITEMS (Granular Line-Item Fact Table) ||
300 -- || Dimension: Product Affinity + Price/Freight Sensitivity ||
301 -- || 112,650 rows – multiple items per order. Composite PK. ||
302 -- || freight_value / price = freight ratio for bargain-hunter analysis. ||
303 -- ||-----||
304
305 DROP TABLE IF EXISTS olist_order_items_dataset CASCADE;
306
307 CREATE TABLE olist_order_items_dataset (
308     order_id          CHAR(32)          NOT NULL,
309     order_item_id      SMALLINT         NOT NULL,    -- sequential item number within order (1..n)
310     product_id         CHAR(32)          NOT NULL,
311     seller_id          CHAR(32)          NOT NULL,
312     shipping_limit_date  TIMESTAMP      NOT NULL,
313     price              NUMERIC(8,2)     NOT NULL,    -- range: 0.85 – 6,735.00

```

```

314 freight_value          NUMERIC(8,2)  NOT NULL,  -- range: 0.00 - 409.68
315
316 CONSTRAINT pk_order_items
317     PRIMARY KEY (order_id, order_item_id),
318
319 CONSTRAINT fk_items_order
320     FOREIGN KEY (order_id)
321     REFERENCES olist_orders_dataset (order_id),
322
323 CONSTRAINT fk_items_product
324     FOREIGN KEY (product_id)
325     REFERENCES olist_products_dataset (product_id),
326
327 CONSTRAINT fk_items_seller
328     FOREIGN KEY (seller_id)
329     REFERENCES olist_sellers_dataset (seller_id)
330 );
331
332 CREATE INDEX idx_items_product
333     ON olist_order_items_dataset (product_id);
334
335 CREATE INDEX idx_items_seller
336     ON olist_order_items_dataset (seller_id);
337
338 COMMENT ON TABLE olist_order_items_dataset
339     IS 'Line-item facts. Use price + freight_value for bargain-hunter ratio. Join to products for category affinity.';
340
341
342 /*
343 =====
344     CSV IMPORT ORDER FOR pgAdmin (COPY command or Import/Export wizard)
345 =====
346
347     Execute imports in this EXACT sequence to satisfy all foreign key constraints:
348
349     1. product_category_name_translation.csv  → product_category_name_translation
350     2. olist_geolocation_dataset.csv         → olist_geolocation_dataset
351        (omit the geolocation_id column - SERIAL auto-generates it)
352     3. olist_customers_dataset.csv          → olist_customers_dataset
353     4. olist_products_dataset.csv          → olist_products_dataset
354     5. olist_sellers_dataset.csv           → olist_sellers_dataset
355     6. olist_orders_dataset.csv            → olist_orders_dataset
356     7. olist_order_payments_dataset.csv    → olist_order_payments_dataset
357     8. olist_order_reviews_dataset.csv     → olist_order_reviews_dataset
358     9. olist_order_items_dataset.csv       → olist_order_items_dataset
359
360     IMPORTANT: Tables 1-5 have no FK dependencies and can technically be loaded
361     in any order. Tables 6-9 MUST follow the sequence above because:
362     • Orders depends on Customers (FK)
363     • Payments, Reviews depend on Orders (FK)
364     • Order Items depends on Orders, Products, AND Sellers (3 FKs)
365
366     pgAdmin COPY example:

```

