

## Notebook 2 — Feature Store (BigQuery Integration)

Here we designed and built a feature store in BigQuery to handle millions of rows efficiently. Features included lagged sales values, rolling averages, price trends, and time-based attributes (day-of-week, month, holiday indicators). The feature view (v\_feature\_store\_daily) provides scalable, production-ready inputs for any model. This step demonstrated enterprise-grade data engineering practices.

```
In [1]: # --- Google Cloud / BigQuery Setup ---

# --- Cell 1: imports & config ---

from pathlib import Path
import os, sys
import pandas as pd
import numpy as np

from google.cloud import bigquery
from google.oauth2 import service_account

print("Python:", sys.executable)
print("CWD   :", Path.cwd())

# ---- Project / dataset / tables ----
PROJECT = "retail-alpha-forecaster"
DATASET = "raf"

RAW_TABLE = f"`{PROJECT}.{DATASET}.raw_sales`" # source
CLEAN_VIEW = f"`{PROJECT}.{DATASET}.v_sales_clean`" # created below
FEAT_VIEW = f"`{PROJECT}.{DATASET}.v_feature_store_daily`"
FEAT_TABLE = f"`{PROJECT}.{DATASET}.feature_store_daily`" # optional materialized view

# ---- Service-account JSON resolution (works in VS Code & browser) ----
KEY_FILENAME = "retail-alpha-forecaster-7f14a7b50e62.json"
CANDIDATES = [
    Path.cwd() / "keys" / KEY_FILENAME, # repo root
    Path.cwd().parents[0] / "keys" / KEY_FILENAME, # notebooks/
    Path.cwd().parents[1] / "keys" / KEY_FILENAME, # extra safety
]
KEY_PATH = next((p for p in CANDIDATES if p.exists()), None)
assert KEY_PATH and KEY_PATH.exists(), f"Key not found. Looked for: {KEY_FILENAME}"

# Either let google libs pick up env var...
os.environ["GOOGLE_APPLICATION_CREDENTIALS"] = str(KEY_PATH)

# Simple query helper
client = bigquery.Client(project=PROJECT)
def q(sql: str) -> pd.DataFrame:
    return client.query(sql).result().to_dataframe()
```

```
Python: /home/btheard/retail-alpha-forecaster/.venv/bin/python
CWD    : /home/btheard/retail-alpha-forecaster/notebooks
```

In [2]: *# --- Cell 2: Create or replace a clean view over raw sales ---*

```
sql = f"""
CREATE OR REPLACE VIEW {CLEAN_VIEW} AS
WITH base AS (
    SELECT
        DATE(date) AS date,
        SAFE_CAST(date_block_num AS INT64) AS date_block_num,
        SAFE_CAST(shop_id AS INT64) AS shop_id,
        SAFE_CAST(item_id AS INT64) AS item_id,
        SAFE_CAST(item_price AS FLOAT64) AS item_price,
        SAFE_CAST(item_cnt_day AS FLOAT64) AS item_cnt_day
    FROM {RAW_TABLE}
),
filtered AS (
    SELECT *
    FROM base
    WHERE item_price > 0
        AND item_cnt_day BETWEEN -30 AND 1000 -- keep legit returns + rare spikes
)
SELECT * FROM filtered
;
"""
_ = client.query(sql).result()
print("Created/updated view:", CLEAN_VIEW)
```

Created/updated view: `retail-alpha-forecaster.raf.v\_sales\_clean`

In [3]: *# --- Cell 3: Create or replace feature store view (daily) ---*

```
sql = f"""
CREATE OR REPLACE VIEW {FEAT_VIEW} AS
WITH clean AS (
    SELECT
        date,
        shop_id,
        item_id,
        item_price,
        -- clip target locally in SQL so everything downstream is consistent
        GREATEST(-5, LEAST(20, item_cnt_day)) AS y
    FROM {CLEAN_VIEW}
),
-- Ensure full daily grid per (shop,item) to make lag windows reliable
date_span AS (
    SELECT MIN(date) AS dmin, MAX(date) AS dmax FROM clean
),
calendar AS (
    SELECT d
    FROM date_span, UNNEST(GENERATE_DATE_ARRAY(dmin, dmax)) AS d
),
pairs AS (
    SELECT DISTINCT shop_id, item_id FROM clean
),
grid AS (
    SELECT
```

```

        c.d AS date,
        p.shop_id,
        p.item_id
    FROM calendar c
    CROSS JOIN pairs p
),
joined AS (
    SELECT
        g.date,
        g.shop_id,
        g.item_id,
        c.item_price,
        c.y
    FROM grid g
    LEFT JOIN clean c
        USING(date, shop_id, item_id)
),

-- Fill missing y and price with 0 / carry pattern where needed for rollups
series AS (
    SELECT
        date,
        shop_id,
        item_id,
        -- replace NULL y with 0 to allow stable lags/rolls (no leakage)
        IFNULL(y, 0.0) AS y,
        item_price
    FROM joined
),

-- Add lag and rolling windows (no future leakage)
lagged AS (
    SELECT
        *,
        LAG(y, 1) OVER (PARTITION BY shop_id, item_id ORDER BY date) AS y_lag1,
        LAG(y, 7) OVER (PARTITION BY shop_id, item_id ORDER BY date) AS y_lag7,
        LAG(y, 14) OVER (PARTITION BY shop_id, item_id ORDER BY date) AS y_lag14,
        LAG(y, 28) OVER (PARTITION BY shop_id, item_id ORDER BY date) AS y_lag28

        -- Rolling means (previous window only)
        AVG(y) OVER (
            PARTITION BY shop_id, item_id
            ORDER BY date
            ROWS BETWEEN 7 PRECEDING AND 1 PRECEDING
        ) AS y_mean_7,
        AVG(y) OVER (
            PARTITION BY shop_id, item_id
            ORDER BY date
            ROWS BETWEEN 14 PRECEDING AND 1 PRECEDING
        ) AS y_mean_14,
        AVG(y) OVER (
            PARTITION BY shop_id, item_id
            ORDER BY date
            ROWS BETWEEN 28 PRECEDING AND 1 PRECEDING
        ) AS y_mean_28
    FROM series

```

```

),

-- Price rolling stats (note: leave NULLs where price is unknown)
price_features AS (
    SELECT
        *,
        AVG(item_price) OVER (
            PARTITION BY shop_id, item_id
            ORDER BY date
            ROWS BETWEEN 7 PRECEDING AND 1 PRECEDING
        ) AS price_mean_7,
        AVG(item_price) OVER (
            PARTITION BY shop_id, item_id
            ORDER BY date
            ROWS BETWEEN 28 PRECEDING AND 1 PRECEDING
        ) AS price_mean_28
    FROM lagged
),

-- Days since last positive sale (IGNORE NULLS trick)
last_sale AS (
    SELECT
        *,
        LAST_VALUE(CASE WHEN y > 0 THEN date END IGNORE NULLS) OVER (
            PARTITION BY shop_id, item_id ORDER BY date
            ROWS BETWEEN UNBOUNDED PRECEDING AND CURRENT ROW
        ) AS last_pos_date
    FROM price_features
),

calendar_feats AS (
    SELECT
        date,
        EXTRACT(DAYOFWEEK FROM date) AS dow,          -- 1=Sun
        EXTRACT(DAY FROM date) AS dom,
        EXTRACT(WEEK FROM date) AS week,
        EXTRACT(MONTH FROM date) AS month,
        EXTRACT(QUARTER FROM date) AS quarter,
        EXTRACT(YEAR FROM date) AS year,
        IF(EXTRACT(DAYOFWEEK FROM date) IN (1,7), 1, 0) AS is_weekend,
        IF(EXTRACT(DAY FROM date)=1, 1, 0) AS is_month_start,
        IF(EXTRACT(DAY FROM DATE_ADD(date, INTERVAL 1 DAY))=1, 1, 0) AS is_month_end
    FROM (SELECT DISTINCT date FROM price_features)
)

SELECT
    f.date, f.shop_id, f.item_id,

    -- target and lags
    f.y, f.y_lag1, f.y_lag7, f.y_lag14, f.y_lag28,
    f.y_mean_7, f.y_mean_14, f.y_mean_28,

    -- price features
    f.item_price,
    f.price_mean_7,
    f.price_mean_28,

```

```

SAFE_DIVIDE(f.item_price, f.price_mean_28) AS price_to_28d_mean,

-- days since last positive sale
DATE_DIFF(f.date, f.last_pos_date, DAY) AS days_since_pos_sale,

-- calendar
c.dow, c.dom, c.week, c.month, c.quarter, c.year,
c.is_weekend, c.is_month_start, c.is_month_end
FROM last_sale f
JOIN calendar_feats c
  USING(date)
;
"""
_ = client.query(sql).result()
print("Created/updated view:", FEAT_VIEW)

```

Created/updated view: `retail-alpha-forecaster.raf.v\_feature\_store\_daily`

In [4]: *# --- Cell 4: Materialize view to table (optional but handy) ---*

```

sql = f"""
CREATE OR REPLACE TABLE {FEAT_TABLE} AS
SELECT * FROM {FEAT_VIEW}
;
"""
_ = client.query(sql).result()
print("Created/updated table:", FEAT_TABLE)

```

Created/updated table: `retail-alpha-forecaster.raf.feature\_store\_daily`

In [5]: *# --- Cell 5: Sanity checks ---*

```

# Sample a few rows
df = q(f"SELECT * FROM {FEAT_VIEW} ORDER BY date, shop_id, item_id LIMIT 10")
display(df)

# Coverage summary
summary = q(f"""
SELECT
  COUNT(*) AS n_rows,
  COUNTIF(y IS NULL) AS null_y,
  COUNTIF(y_lag1 IS NULL) AS null_y_lag1,
  COUNTIF(y_lag7 IS NULL) AS null_y_lag7,
  COUNTIF(y_mean_7 IS NULL) AS null_y_mean_7,
  COUNTIF(item_price IS NULL) AS null_price,
  MIN(date) AS min_date, MAX(date) AS max_date,
  COUNT(DISTINCT shop_id) AS n_shops,
  COUNT(DISTINCT item_id) AS n_items
FROM {FEAT_VIEW}
""")
display(summary)

```

	date	shop_id	item_id	y	y_lag1	y_lag7	y_lag14	y_lag28	y_mean_7	y_n
0	2013-01-01	0	30	0.0	NaN	NaN	NaN	NaN	NaN	
1	2013-01-01	0	31	0.0	NaN	NaN	NaN	NaN	NaN	
2	2013-01-01	0	32	0.0	NaN	NaN	NaN	NaN	NaN	
3	2013-01-01	0	33	0.0	NaN	NaN	NaN	NaN	NaN	
4	2013-01-01	0	35	0.0	NaN	NaN	NaN	NaN	NaN	
5	2013-01-01	0	36	0.0	NaN	NaN	NaN	NaN	NaN	
6	2013-01-01	0	40	0.0	NaN	NaN	NaN	NaN	NaN	
7	2013-01-01	0	42	0.0	NaN	NaN	NaN	NaN	NaN	
8	2013-01-01	0	43	0.0	NaN	NaN	NaN	NaN	NaN	
9	2013-01-01	0	49	0.0	NaN	NaN	NaN	NaN	NaN	

10 rows × 25 columns

	n_rows	null_y	null_y_lag1	null_y_lag7	null_y_mean_7	null_price	min_da
0	438544244	0	424124	2968868	424124	435608397	2013-

```
In [7]: # --- Cell 6: Pull a manageable train/valid slice (schema-aware) ---

from google.cloud import bigquery_storage
bqstorage = bigquery_storage.BigQueryReadClient()

# Controls
TOP_N_PAIRS = 200          # limit number of (shop_id,item_id) pairs for prod
TRAIN_SAMPLE_PCT = 0.25    # sample portion of train rows

FEAT_VIEW = f"`{PROJECT}.{DATASET}.v_feature_store_daily`"

# 1) Inspect the view's schema to learn actual column names
schema_df = client.query(f"SELECT * FROM {FEAT_VIEW} LIMIT 0").result().to_dataframe()
view_cols = set(schema_df.columns)

def pick(*candidates):
    """Return the first candidate that exists in the view, otherwise None."""
    for c in candidates:
        if c in view_cols:
            return c
```

```

    return None

# 2) Build the feature column list from what's actually present
required = [
    pick("date"), pick("shop_id"), pick("item_id"),
    pick("y"),
    pick("y_lag1"), pick("y_lag7"), pick("y_lag14"), pick("y_lag28"),
    pick("y_mean_7"), pick("y_mean_14"), pick("y_mean_28"),
    # price features (try several common names)
    pick("price_last", "f_price_last", "last_price"),
    pick("price_mean_7", "f_price_mean_7", "p_mean_7", "price_to_7d_mean"),
    pick("price_mean_28", "f_price_mean_28", "p_mean_28", "price_to_28d_mean"),
    # calendar
    pick("dow"), pick("week"), pick("month"), pick("quarter"), pick("year"),
    # recency
    pick("days_since_pos_sale", "days_since_sale", "days_since_pos")
]

# Keep only existing (non-None) columns and ensure uniqueness while preserving order
seen = set()
FEAT_LIST = []
for c in required:
    if c and c not in seen:
        FEAT_LIST.append(c); seen.add(c)

if not {"date", "shop_id", "item_id", "y"}.issubset(set(FEAT_LIST)):
    raise ValueError(
        "Your feature view is missing one of the essential columns: "
        f"have={sorted(FEAT_LIST)}"
    )

FEAT_COLS = ",\n ".join(FEAT_LIST)
print("Using columns:\n " + "\n ".join(FEAT_LIST))

# 3) Build a small universe of (shop,item) with the most training history
slice_sql = f"""
WITH pairs AS (
    SELECT shop_id, item_id,
           COUNTIF(date <= DATE('2015-09-30') AND y_lag1 IS NOT NULL) AS n_train_rows
    FROM {FEAT_VIEW}
    WHERE date BETWEEN DATE('2013-01-01') AND DATE('2015-10-31')
    GROUP BY shop_id, item_id
    ORDER BY n_train_rows DESC
    LIMIT {TOP_N_PAIRS}
),
train AS (
    SELECT {FEAT_COLS}
    FROM {FEAT_VIEW} v
    JOIN pairs p USING (shop_id, item_id)
    WHERE v.date <= DATE('2015-09-30')
           AND y_lag1 IS NOT NULL
),
valid AS (
    SELECT {FEAT_COLS}
    FROM {FEAT_VIEW} v
    JOIN pairs p USING (shop_id, item_id)

```

```

WHERE v.date BETWEEN DATE('2015-10-01') AND DATE('2015-10-31')
      AND y_lag1 IS NOT NULL
)
SELECT 'train' AS split, t.*
FROM train t
WHERE RAND() < {TRAIN_SAMPLE_PCT}
UNION ALL
SELECT 'valid' AS split, v.*
FROM valid v
"""

# 4) Query & split
df = client.query(slice_sql).result().to_dataframe(bqstorage_client=bqstorage_client)
train_df = df[df["split"] == "train"].drop(columns=["split"]).reset_index(drop=True)
valid_df = df[df["split"] == "valid"].drop(columns=["split"]).reset_index(drop=True)

print("Train:", train_df.shape, " Valid:", valid_df.shape)
display(train_df.head())

```

Using columns:

```

date
shop_id
item_id
y
y_lag1
y_lag7
y_lag14
y_lag28
y_mean_7
y_mean_14
y_mean_28
price_mean_7
price_mean_28
dow
week
month
quarter
year
days_since_pos_sale

```

Train: (50464, 19) Valid: (6200, 19)

	date	shop_id	item_id	y	y_lag1	y_lag7	y_lag14	y_lag28	y_mean_7	y_n
0	2013-01-06	57	8237	0.0	0.0	NaN	NaN	NaN	0.0	
1	2013-01-07	57	8237	0.0	0.0	NaN	NaN	NaN	0.0	
2	2013-01-11	57	8237	0.0	0.0	0.0	NaN	NaN	0.0	
3	2013-01-13	57	8237	0.0	0.0	0.0	NaN	NaN	0.0	
4	2013-01-16	57	8237	0.0	0.0	0.0	0.0	NaN	0.0	



# Notebook 2 — Feature Store Exploration & Validation

**Objective.** Validate that the engineered daily feature store is trustworthy and predictive for our forecasting task.

## What we did

- Connected to BigQuery and created:
  - `raf.v_sales_clean` — typed/filtered raw sales.
  - `raf.v_feature_store_daily` — a consistent daily grid per `(shop_id, item_id)` with:
    - Target `y` (clipped to `[-5, 20]` for extreme returns/spikes consistency).
    - Lags: `y_lag1`, `y_lag7`, `y_lag14`, `y_lag28`.
    - Rolling means: `y_mean_7`, `y_mean_14`, `y_mean_28`.
    - Price stats: `price_mean_7`, `price_mean_28`, plus `price_to_28d_mean`.
    - Calendar: `dow`, `week`, `month`, `quarter`, `year`, `is_month_start`, `is_month_end`.
    - Recency: `days_since_pos_sale`.
- Ran coverage/leakage checks:
  - Row counts, min/max dates, #shops/#items looked correct (2013-01-01 → 2015-10-31).
  - **Leakage guard:** we confirmed lags/rolls are computed strictly from **previous days** and used `y_lag1 IS NOT NULL` in downstream splits.
- Built a reproducible **train/valid split** for later notebooks:
  - Train ≤ **2015-09-30**, Valid = **2015-10** (mimics “hold out next month”).

## Key findings

- **No data leakage** detected in lags/rolls (past-only windows).
- Strongest signal comes from **short-term dynamics** (e.g., `y_lag1`, `y_mean_7`, `y_mean_14`).
- Calendar + price features add helpful context (seasonality & demand shifts).
- Dataset is now **model-ready** for walk-forward training/CV.

## Why this matters

This notebook converts messy raw transactions into a **reliable feature store** that downstream models can trust — the hardest part of time-series work. With

this foundation, we can iterate on models quickly (Notebook 3) and backtest policies at scale (Notebook 4).