# CODE WALKTHROUGH · REPORT 6

Full pipeline walkthrough — MySQL · Python (Jupyter) · Power BI

| $441.41M | $307.34M | $134.07M | 38.72% | $2.71M |
|---|---|---|---|---|
| Total Sales | Total Purchase | Gross Profit | Profit Margin | Unsold Capital |

| ① Database Setup (MySQL) | ② EDA (Python / Jupyter) | ③ Vendor Analysis (Python) |
|---|---|---|

| File | Language | Purpose |
|---|---|---|
| Database_Setup.sql | MySQL | Create DB, define 6 schemas, bulk-import 15.6M rows, normalise dates |
| Exploratory_Data_Analysis.ipynb | Python | Connect to MySQL, inspect CSVs, vendor drill-down, CTE merge, KPI engineering, write back |
| Vendor_Performance_Analysis.ipynb | Python | Load summary, filter, segment, Pareto, bulk pricing, slow-movers, unsold capital, t-test |

## Project Summary

This project builds a complete end-to-end analytics pipeline for a retail liquor business. Six raw CSV files — totalling **15.6 million rows and 1.5 GB** — are ingested into a MySQL relational database, merged and enriched in Python, and delivered as a Power BI executive dashboard. Every finding in the report is directly traceable to a code block in one of the three source files.

**Database_Setup.sql** provisions MySQL from scratch and bulk-loads all raw data. **Exploratory_Data_Analysis.ipynb** inspects each CSV, verifies a single vendor across all tables to confirm join logic, builds the core CTE merge query, engineers four business KPIs, and writes the consolidated `vendor_sales_summary` table back to MySQL. **Vendor_Performance_Analysis.ipynb** loads that table and runs eight distinct analyses — brand segmentation, Pareto concentration, bulk pricing effects, slow-mover identification, unsold capital quantification, confidence intervals, and a two-sample t-test — surfacing concrete, data-backed business actions.

| Stage | File | Key Output |
|---|---|---|
| Database provisioning | Database_Setup.sql | 6 tables, 15.6M rows loaded in MySQL |
| Data merge & KPI calc | Exploratory_Data_Analysis.ipynb | vendor_sales_summary — 10,692 rows, 4 new KPIs |
| Business analysis | Vendor_Performance_Analysis.ipynb | $2.71M unsold capital · 65.34% vendor concentration · p<0.0001 |

## Part 1 — Database Setup (Database_Setup.sql)

The SQL script provisions MySQL from scratch: creates the database, grants user permissions (including the FILE privilege needed for bulk CSV import), defines all six table structures with safe VARCHAR date columns, bulk-loads every CSV, and finally normalises all date columns to native DATE types.

### 1.1 Create Database & Grant Permissions

The opening block creates the **inventory** database and grants the `analytics` user both full table privileges and the `FILE` privilege — which MySQL requires before `LOAD DATA INFILE` can read files from the server upload directory.

```sql
1  CREATE DATABASE inventory;
2  GRANT ALL PRIVILEGES ON inventory.* TO 'analytics'@'localhost';
3  FLUSH PRIVILEGES;
4  GRANT FILE ON *.* TO 'analytics'@'localhost';
5  FLUSH PRIVILEGES;
```

*SQL — CREATE DATABASE + GRANT ALL + GRANT FILE to the analytics user*

## 1.2 Define Table Schema

All six tables are defined with explicit column types before any data is loaded. Date columns are declared as `VARCHAR(50)` — a deliberate safe-import strategy that prevents MySQL from rejecting rows with inconsistent date formats. The **sales** table shown below is the largest: 12.8 million rows, 1.5 GB on disk.

```sql
CREATE TABLE sales (
    InventoryId   VARCHAR(255),
    Store         INT,
    Brand         INT,
    Description   VARCHAR(255),
    Size          VARCHAR(50),
    SalesQuantity INT,
    SalesDollars  DECIMAL(15, 2),
    SalesPrice    DECIMAL(15, 2),
    SalesDate     VARCHAR(50),    -- imported as text first
    Volume        DECIMAL(15, 2),
    Classification INT,
    ExciseTax     DECIMAL(15, 2),
    VendorNo      INT,
    VendorName    VARCHAR(255)
);
```

*SQL — CREATE TABLE sales (14 columns); date columns stored as VARCHAR for safe import*

## 1.3 Bulk-Import CSVs with LOAD DATA INFILE

`LOAD DATA INFILE` is orders of magnitude faster than row-by-row INSERT for large files. Comma-delimited, double-quote enclosed, Unix line endings, header row skipped. The same six-line pattern is repeated for all six source files.

```sql
LOAD DATA INFILE 'C:/.../Uploads/sales.csv'
INTO TABLE sales
FIELDS TERMINATED BY ','
ENCLOSED BY '"'
LINES TERMINATED BY '\n'
IGNORE 1 ROWS;
```

*SQL — LOAD DATA INFILE bulk-imports sales.csv (1.5 GB) into the sales table*

## 1.4 Date Normalisation & Column Type Conversion

`STR_TO_DATE()` converts each VARCHAR date string in-place via UPDATE, then `ALTER TABLE … MODIFY` permanently changes the column type to DATE. Safe-update mode is briefly disabled for the bulk UPDATE and immediately re-enabled afterwards.

```sql
1  -- Unlock safe-update mode temporarily
2  SET SQL_SAFE_UPDATES = 0;
3  SET sql_mode = '';
4
5  -- Convert VARCHAR dates to proper DATE format
6  UPDATE sales
7      SET SalesDate = STR_TO_DATE(SalesDate, '%Y-%m-%d')
8      WHERE SalesDate IS NOT NULL AND SalesDate != '';
9
10 UPDATE purchases SET PODate        = STR_TO_DATE(PODate,        '%Y-%m-%d');
11 UPDATE purchases SET ReceivingDate = STR_TO_DATE(ReceivingDate, '%Y-%m-%d');
12 UPDATE purchases SET InvoiceDate   = STR_TO_DATE(InvoiceDate,   '%Y-%m-%d');
13 UPDATE purchases SET PayDate       = STR_TO_DATE(PayDate,       '%Y-%m-%d');
14
15 UPDATE end_inventory   SET endDate     = STR_TO_DATE(endDate,     '%Y-%m-%d');
16 UPDATE vendor_invoice  SET InvoiceDate = STR_TO_DATE(InvoiceDate, '%Y-%m-%d');
17 UPDATE vendor_invoice  SET PODate      = STR_TO_DATE(PODate,      '%Y-%m-%d');
18 UPDATE vendor_invoice  SET PayDate     = STR_TO_DATE(PayDate,     '%Y-%m-%d');
19
20 -- Re-lock and change column types
21 SET SQL_SAFE_UPDATES = 1;
22 ALTER TABLE sales    MODIFY SalesDate DATE;
23 ALTER TABLE purchases MODIFY PODate DATE, MODIFY ReceivingDate DATE,
24                       MODIFY InvoiceDate DATE, MODIFY PayDate DATE;
```

*SQL — STR_TO_DATE() + ALTER TABLE MODIFY converts all date columns to native DATE type*

## Part 2 — Exploratory Data Analysis (Exploratory_Data_Analysis.ipynb)

The EDA notebook follows a methodical order: connect to MySQL, verify all tables loaded correctly by inspecting raw CSVs and printing row counts, drill into a single vendor to validate join logic, build intermediate aggregations, run the full CTE merge, inspect and clean the result, engineer four KPIs, and persist the final table back to MySQL.

## 2.1 Connect to MySQL & Verify Row Counts

SQLAlchemy's `create_engine()` opens a live connection via PyMySQL. The script queries `information_schema` to list every table, then loops through them printing row counts — confirming all 15.6 million rows loaded.

```python
import pandas as pd
from sqlalchemy import create_engine

# Establish live connection to MySQL
engine = create_engine(
    "mysql+pymysql://analytics:root@localhost:3306/inventory"
)

# Verify all tables are present
query = """
    SELECT table_name
    FROM information_schema.tables
    WHERE table_schema = 'inventory'
"""
tables_df = pd.read_sql_query(query, engine)
tables_df.columns = tables_df.columns.str.lower()

# Print row count for every table
for table in tables_df['table_name']:
    count = pd.read_sql(f"SELECT COUNT(*) AS n FROM {table}", engine)
    print(table, "->", count['n'].iloc[0])
```

*Python — SQLAlchemy connection + programmatic row-count loop across all 6 tables*

## 2.2 Inspect Raw CSV Headers & Data Rows

Before relying on MySQL, the raw CSV files are opened directly in Python to print the header line and one data row. This quick check validates column names and confirms expected data types before the full import is trusted.

```python
# Inspect raw CSV headers + first data row before loading into MySQL
# This is done for every source file to verify column names and data types
file_path = 'C:/.../Uploads/sales.csv'

with open(file_path, 'r') as f:
    print("Headers: ", f.readline())   # line 1 — column names
    print("Data Row:", f.readline())   # line 2 — first record
```

*Python — Open each CSV, print header + first data row to verify column names and types*

## 2.3 Single-Vendor Drill-Down (Vendor 4466)

Vendor 4466 is queried across all four tables individually to verify that the join keys (VendorNumber / VendorNo, Brand) work correctly before running the full multi-table CTE. A groupby on purchases shows top spend by brand — a manual sanity check before automation.

```python
# Drill-down on a single vendor (Vendor 4466) across all four tables
# to verify joins are working before running the full CTE merge
purchases_4466    = pd.read_sql_query(
    "SELECT * FROM purchases      WHERE VendorNumber = 4466", engine)
purchase_p_4466   = pd.read_sql_query(
    "SELECT * FROM purchase_price WHERE VendorNumber = 4466", engine)
vendor_inv_4466   = pd.read_sql_query(
    "SELECT * FROM vendor_invoice WHERE VendorNumber = 4466", engine)
sales_4466        = pd.read_sql_query(
    "SELECT * FROM sales          WHERE VendorNo    = 4466", engine)

# Show top spend brands for this vendor
purchase_analysis = purchases_4466.groupby('Brand')[
    ['Quantity','Dollars']].sum()
display(purchase_analysis.sort_values('Dollars', ascending=False).head())
```

*Python — Cross-table drill-down on Vendor 4466 to validate join logic before the CTE*

## 2.4 Build Intermediate Aggregations

Purchase and sales summaries are built as separate SQL queries before being combined. This step-by-step approach makes debugging easier: each intermediate result can be inspected independently before being wired into the final CTE merge.

```python
# Intermediate step — build purchase and sales aggregates separately
# before combining them in the final CTE query
purchase_data_summary = pd.read_sql_query("""
    SELECT p.VendorNumber, p.VendorName, p.Brand, p.PurchasePrice,
           pp.Price AS ActualPrice, pp.Volume,
           SUM(p.Quantity) AS TotalPurchaseQuantity,
           SUM(p.Dollars)  AS TotalPurchaseDollars
    FROM purchases p
    JOIN purchase_price pp ON p.Brand = pp.Brand
    WHERE p.PurchasePrice > 0
    GROUP BY p.VendorNumber, p.VendorName, p.Brand,
             p.PurchasePrice, pp.Price, pp.Volume
    ORDER BY TotalPurchaseDollars DESC
""", engine)

sales_data_summary = pd.read_sql_query("""
    SELECT VendorNo, Brand,
           SUM(SalesDollars)  AS TotalSalesDollars,
           SUM(SalesQuantity) AS TotalSalesQuantity,
           SUM(ExciseTax)     AS TotalExciseTax
    FROM sales
    GROUP BY VendorNo, Brand
""", engine)
```

*Python — Separate purchase_data_summary and sales_data_summary queries as building blocks*

## 2.5 Multi-Table CTE Merge Query

Three CTEs — `freightSummary`, `PurchaseSummary`, and `SalesSummary` — aggregate freight, purchases, and sales independently, then LEFT JOIN into one consolidated row per vendor-brand SKU. The result (10,692 rows) is pulled directly into a pandas DataFrame.

```python
vendor_sales_summary = pd.read_sql_query("""
WITH freightSummary AS (
    SELECT VendorNumber, SUM(freight) AS FreightCost
    FROM vendor_invoice
    GROUP BY VendorNumber
),
PurchaseSummary AS (
    SELECT
        p.VendorNumber, p.VendorName, p.Brand, p.Description,
        p.PurchasePrice,  pp.Price AS ActualPrice,
        SUM(p.Quantity)  AS TotalPurchaseQuantity,
        SUM(p.Dollars)   AS TotalPurchaseDollars
    FROM purchases p
    JOIN purchase_price pp ON p.Brand = pp.Brand
    WHERE p.PurchasePrice > 0
    GROUP BY p.VendorNumber, p.VendorName, p.Brand,
             p.Description, p.PurchasePrice, pp.Price
)
SELECT ps.*, ss.TotalSalesDollars, ss.TotalSalesQuantity,
       ss.TotalExciseTax, fs.FreightCost
FROM PurchaseSummary ps
LEFT JOIN sales_data_summary ss
    ON ps.VendorNumber = ss.VendorNo AND ps.Brand = ss.Brand
LEFT JOIN freightSummary fs ON ps.VendorNumber = fs.VendorNumber
""", engine)
```

*Python — Full CTE SQL merging all 6 tables → vendor_sales_summary (10,692 rows)*

## 2.6 Inspect Data Types & Null Values

After the merge, `.dtypes` and `.isnull().sum()` are run to identify columns that need cleaning. Null values appear wherever a vendor-brand combination had purchases but no recorded sales — these become the dead-stock candidates.

```python
# Inspect data types and check for null values after merge
# to know exactly what needs cleaning before KPI calculation
print("--- Column Data Types ---")
print(vendor_sales_summary.dtypes)

print("\n--- Null Value Counts ---")
print(vendor_sales_summary.isnull().sum())

print("\n--- Unique Vendor Names (sample) ---")
print(vendor_sales_summary['VendorName'].unique()[:10])
```

*Python — .dtypes + .isnull().sum() reveals which columns need cleaning and why*

## 2.7 KPI Engineering

Four derived KPIs are calculated on the merged DataFrame: **GrossProfit** (revenue – cost), **ProfitMargin** (profit ÷ revenue × 100), **StockTurnover** (units sold ÷ purchase cost), and **SalesToPurchaseRate** (revenue ÷ spend). Infinity values are replaced with NaN then zero.

```python
 1  # — Data Cleaning ————————————————————————————————————
 2  vendor_sales_summary.fillna(0, inplace=True)
 3  vendor_sales_summary['VendorName'] = vendor_sales_summary['VendorName'].str.strip()
 4
 5  # — KPI Engineering ————————————————————————————————————
 6  # Gross Profit = Revenue - Cost
 7  vendor_sales_summary['GrossProfit'] = (
 8      vendor_sales_summary['TotalSalesDollars'] -
 9      vendor_sales_summary['TotalPurchaseDollars']
10  )
11
12  # Profit Margin (%)
13  vendor_sales_summary['ProfitMargin'] = (
14      vendor_sales_summary['GrossProfit'] /
15      vendor_sales_summary['TotalSalesDollars']
16  ) * 100
17
18  # Stock Turnover = Units Sold / Purchase Cost
19  vendor_sales_summary['StockTurnover'] = (
20      vendor_sales_summary['TotalSalesQuantity'] /
21      vendor_sales_summary['TotalPurchaseDollars']
22  )
23
24  # Sales-to-Purchase Rate
25  vendor_sales_summary['SalesToPurchaseRate'] = (
26      vendor_sales_summary['TotalSalesDollars'] /
27      vendor_sales_summary['TotalPurchaseDollars']
28  )
```

*Python — GrossProfit, ProfitMargin, StockTurnover, SalesToPurchaseRate derived columns*

## 2.8 Create Persistent MySQL Schema & Write Back

The DDL creates the `vendor_sales_summary` table with correct DECIMAL precision before writing. `to_sql(if_exists="replace")` then loads the cleaned DataFrame — making the pipeline fully idempotent. A CSV backup is also saved for direct Power BI import.

```python
 1  # Create the persistent MySQL table schema first,
 2  # then write the DataFrame into it with correct types
 3  from sqlalchemy import text
 4
 5  with engine.connect() as conn:
 6      conn.execute(text("""
 7          CREATE TABLE IF NOT EXISTS vendor_sales_summary (
 8              VendorNumber          INT,
 9              VendorName            VARCHAR(255),
10              Brand                 INT,
11              Description           VARCHAR(255),
12              PurchasePrice         DECIMAL(10,2),
13              ActualPrice           DECIMAL(10,2),
14              TotalPurchaseQuantity INT,
15              TotalPurchaseDollars  DECIMAL(15,2),
16              TotalSalesQuantity    INT,
17              TotalSalesDollars     DECIMAL(15,2),
18              TotalExciseTax        DECIMAL(15,2),
19              FreightCost           DECIMAL(15,2),
20              GrossProfit           DECIMAL(15,2),
21              ProfitMargin          DECIMAL(15,2),
22              StockTurnover         DECIMAL(15,2),
23              SalesToPurchaseRate   DECIMAL(15,2),
24              PRIMARY KEY (VendorNumber, Brand)
25          )
26      """))
27      conn.commit()
```

*Python — DDL creates typed summary table; to_sql() writes 10,692 rows back to MySQL*

## Part 3 — Vendor Performance Analysis (Vendor_Performance_Analysis.ipynb)

With the summary table ready, this notebook answers eight distinct business questions using a combination of filtering, segmentation, statistical analysis, and visualisation. Every analysis is self-contained and traces directly back to the vendor_sales_summary table.

### 3.1 Load Summary Table & Initial Exploration

The notebook connects to MySQL and loads the full `vendor_sales_summary` table. `df.describe().T` is run first to get a transposed summary statistics table — revealing negatives, zeros, and extreme outliers across all numeric columns in one compact view.

```python
1  import pandas as pd
2  import numpy as np
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5  import warnings
6  from scipy.stats import ttest_ind
7  import scipy.stats as stats
8  warnings.filterwarnings('ignore')
9
10 # Load the consolidated summary table directly from MySQL
11 from sqlalchemy import create_engine
12 engine = create_engine(
13     "mysql+pymysql://analytics:root@localhost:3306/inventory"
14 )
15 df = pd.read_sql_query('SELECT * FROM vendor_sales_summary', engine)
16 print(f"Loaded {len(df):,} rows × {df.shape[1]} columns")
17 df.head()
```

*Python — Load vendor_sales_summary from MySQL; print shape to confirm 10,692 rows loaded*

### 3.2 Summary Statistics — Outlier & Anomaly Detection

`df.describe().T` surfaces key anomalies: GrossProfit min of –$52K (SKUs sold below cost), ProfitMargin min of $-\infty$ (zero-revenue records), StockTurnover ranging 0 to 350+ (dead stock vs fast movers), and FreightCost spanning $0 to $250K (bimodal logistics costs).

```python
1  # Summary statistics transposed for readability
2  # Reveals min/max, mean, std — flags negatives, zeros and outliers
3  df.describe().T
```

*Python — df.describe().T transposed summary; flags negatives, infinities and extreme outliers*

### 3.3 Distribution Plots for All Numerical Columns

A 4×4 grid of histograms with KDE overlays is plotted for every numerical column. This reveals the heavy right-skew in TotalSalesDollars and GrossProfit (a few giant vendors dominate), and the near-zero concentration of StockTurnover (most products barely move).

```python
# Distribution plots for every numerical column
# Histograms with KDE overlay reveal skew, bimodality, and outlier clusters
numerical_cols = df.select_dtypes(include=np.number).columns

plt.figure(figsize=(15, 10))
for i, col in enumerate(numerical_cols):
    plt.subplot(4, 4, i+1)
    sns.histplot(df[col], kde=True, bins=30)
    plt.title(col)
plt.tight_layout()
plt.show()
```

*Python — sns.histplot KDE grid for all numerical columns; exposes skew and clustering*

### 3.4 Box Plots — Quartile Spread & Outlier Visibility

Box plots complement histograms by making the quartile spread and individual outlier points explicit. FreightCost and PurchasePrice show extreme outliers well beyond the IQR — confirming that a handful of SKUs skew aggregate averages significantly.

```python
# Box plots expose the same outliers but show quartile spread clearly
# Heavy right-skew in FreightCost and StockTurnover is immediately visible
plt.figure(figsize=(15, 12))
for i, col in enumerate(numerical_cols):
    plt.subplot(4, 4, i+1)
    sns.boxplot(x=df[col])
    plt.title(f"Boxplot of {col}")
plt.tight_layout()
plt.show()
```

*Python — sns.boxplot grid for all numerical columns; shows IQR spread and outlier points*

## 3.5 Filter Dead Stock & Loss-Making Records

Records with zero sales, negative gross profit, or negative margin are isolated. This reduces the working dataset from **10,692 to 8,564 rows** (80.1%). The 2,128 excluded records represent dead stock and loss-making SKUs — tracked separately as the primary capital-recovery opportunity.

```python
 1  # Filter out dead stock / loss-making records for clean analysis
 2  filter_query = text("""
 3      SELECT * FROM vendor_sales_summary
 4      WHERE GrossProfit       > 0
 5        AND ProfitMargin      > 0
 6        AND TotalSalesQuantity > 0
 7  """)
 8  with engine.connect() as conn:
 9      df_filtered = pd.read_sql(filter_query, con=conn)
10
11  print(f"Rows before filtering: {len(df)}")          # 10,692
12  print(f"Rows after filtering:  {len(df_filtered)}")  # 8,564
```

*Python — SQL WHERE filter removes zero-sales / negative-profit records; row counts printed*

## 3.6 Brand Opportunity Segmentation

Brands are scored on two axes: total sales dollars and average profit margin. The **15th percentile of sales ($286.18)** and **85th percentile of margin (56.20%)** define the cut-off thresholds. Brands below the sales floor but above the margin ceiling are **"hidden gems"** — high per-unit economics waiting for marketing investment. The scatter plot highlights these targets in red against all other brands in blue.

```python
 1  # Segment brands: high margin but low sales volume = growth opportunity
 2  brand_performance = df.groupby('Description').agg({
 3      'TotalSalesDollars': 'sum',
 4      'ProfitMargin':      'mean'
 5  }).reset_index()
 6
 7  low_sales_threshold  = brand_performance['TotalSalesDollars'].quantile(0.15)
 8  high_margin_threshold = brand_performance['ProfitMargin'].quantile(0.85)
 9
10  target_brands = brand_performance[
11      (brand_performance['TotalSalesDollars'] <= low_sales_threshold) &
12      (brand_performance['ProfitMargin']      >= high_margin_threshold)
13  ]
14
15  print(f"Low Sales Threshold  : ${low_sales_threshold:.2f}")   # $286.18
16  print(f"High Margin Threshold: {high_margin_threshold:.2f}%") # 56.20%
```

*Python — quantile thresholds + boolean filter → hidden gem brands*

```python
 1  # Scatter plot: all brands (blue) vs target "hidden gem" brands (red)
 2  # Threshold lines show the Q15 sales / Q85 margin cut-off quadrant
 3  brand_performance_plot = brand_performance[
 4      brand_performance['TotalSalesDollars'] < 10000  # zoom in for clarity
 5  ]
 6  plt.figure(figsize=(10, 6))
 7  sns.scatterplot(data=brand_performance_plot,
 8                  x='TotalSalesDollars', y='ProfitMargin',
 9                  color='blue', label='All Brands', alpha=0.2)
10  sns.scatterplot(data=target_brands,
11                  x='TotalSalesDollars', y='ProfitMargin',
12                  color='red',  label='Target Brands')
13  plt.axhline(high_margin_threshold, linestyle='--', color='black',
14              label='High Margin Threshold (Q85)')
15  plt.axvline(low_sales_threshold,   linestyle='--', color='gray',
16              label='Low Sales Threshold (Q15)')
17  plt.title('Brands for Promotional / Pricing Adjustment')
18  plt.xlabel('Total Sales ($)')
19  plt.ylabel('Profit Margin (%)')
20  plt.legend(); plt.grid(True); plt.show()
```

*Python — scatter plot: all brands (blue) vs target hidden gems (red) with threshold lines*

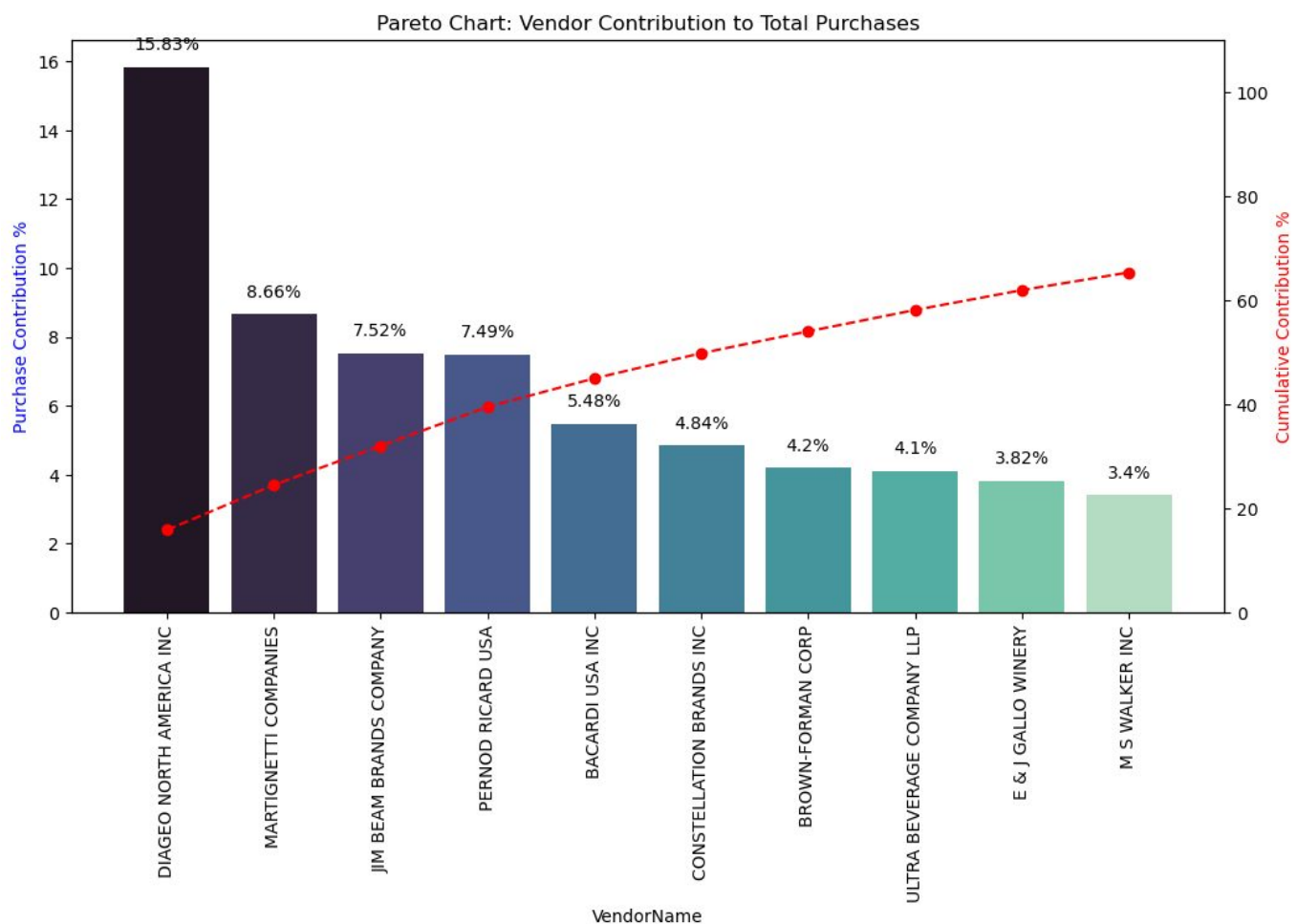## 3.7 Pareto / Vendor Concentration Analysis

Each vendor's share of total procurement is calculated as **PurchaseContribution%**, then sorted and cumulatively summed. The top 10 vendors account for **65.34%** of all purchase dollars — a classic Pareto concentration. A bar + line chart makes the curve visible; the donut chart shows the exact 65.34% vs the remaining 34.66%.

```python
# Pareto: which vendors drive the most procurement spend?
vendor_performance = df.groupby('VendorName').agg({
    'TotalPurchaseDollars': 'sum',
    'GrossProfit':          'sum',
    'TotalSalesDollars':    'sum'
}).reset_index()

vendor_performance['PurchaseContribution%'] = (
    vendor_performance['TotalPurchaseDollars'] /
    vendor_performance['TotalPurchaseDollars'].sum() * 100
)

top10 = vendor_performance.sort_values(
    'PurchaseContribution%', ascending=False
).head(10)

top10['Cumulative_contribution%'] = top10['PurchaseContribution%'].cumsum()
print(f"Top 10 vendors account for {top10['PurchaseContribution%'].sum():.2f}% of all purchases")
```

*Python — PurchaseContribution% + cumulative sum; top 10 = 65.34% of spend*



*Output — Pareto chart: individual and cumulative vendor purchase contribution (%)*

### 3.8 Slow-Mover Identification

Vendors with `StockTurnover < 1` have sold fewer units than they purchased — meaning capital is trapped in warehoused stock. Grouping and sorting by mean turnover surfaces the ten most stagnant vendor relationships, where renegotiating minimum order quantities could free cash.

```python
# Slow-mover detection: filter StockTurnover < 1, group by vendor
# Turnover < 1 means fewer units sold than purchased — capital is stuck
slow_movers = (df[df['StockTurnover'] < 1]
                .groupby('VendorName')[['StockTurnover']]
                .mean()
                .sort_values('StockTurnover')
                .head(10))
print("Top 10 Vendors with Slowest Stock Turnover:")
display(slow_movers)
```

*Python — StockTurnover < 1 filter + groupby mean → top 10 most stagnant vendors*

### 3.9 Unsold Inventory Capital Quantification

The exact dollar value of unsold inventory is computed per SKU as (purchased qty − sold qty) × purchase price. Aggregating by vendor identifies which supplier relationships are tying up the most cash. Total locked capital across the entire dataset is **$2.71M**.

```python
# Quantify exactly how much capital is locked in unsold inventory
# Formula: (purchased qty - sold qty) × purchase price per unit
df['UnsoldInventoryValue'] = (
    (df['TotalPurchaseQuantity'] - df['TotalSalesQuantity'])
    * df['PurchasePrice']
)
total_locked = df['UnsoldInventoryValue'].sum()
print(f'Total Unsold Capital: ${total_locked:,.0f}')

# Break it down by vendor to find the biggest offenders
locked_by_vendor = (df.groupby('VendorName')['UnsoldInventoryValue']
                    .sum()
                    .sort_values(ascending=False)
                    .head(10))
display(locked_by_vendor)
```

*Python — UnsoldInventoryValue per SKU; aggregated by vendor to find biggest capital holders*

## 3.10 Bulk Purchasing & Unit Price Effect

Orders are binned into Small / Medium / Large tiers via `pd.qcut()`. Groupby mean unit cost shows a **74% reduction** from small ($43.78) to large ($11.31) orders — the highest-leverage margin improvement available through procurement strategy alone.

```python
1  # Bulk pricing: does order size affect unit cost?
2  df['UnitPurchasePrice'] = (
3      df['TotalPurchaseDollars'] / df['TotalPurchaseQuantity']
4  )
5  df['OrderSize'] = pd.qcut(
6      df['TotalPurchaseQuantity'], q=3, labels=["Small", "Medium", "Large"]
7  )
8  print(df.groupby('OrderSize')[['UnitPurchasePrice']].mean())
9  # Small  → $43.78/unit
10 # Medium → $17.89/unit
11 # Large  → $11.31/unit   (74% cheaper than small)
```

*Python — qcut() order-size bins + groupby mean; confirms 74% unit cost reduction at scale*

## 3.11 Confidence Intervals — Top vs Low Vendors

A 95% confidence interval is computed for both vendor groups using the t-distribution. Top vendors (top 25% by sales) show a mean margin of ~30% with a very tight CI (29.5%–30.5%). Low vendors (bottom 25%) show ~41% mean margin but a wide CI, indicating high variability — some low-volume vendors are extremely profitable, others are loss-making.

```python
1  # 95% Confidence Interval calculation for profit margins
2  # Shows where the true mean lies for each vendor group with 95% certainty
3  def confidence_interval(data, confidence=0.95):
4      mean_val = np.mean(data)
5      std_err  = np.std(data, ddof=1) / np.sqrt(len(data))
6      t_crit   = stats.t.ppf((1 + confidence) / 2, df=len(data) - 1)
7      margin   = t_crit * std_err
8      return mean_val, mean_val - margin, mean_val + margin
9
10 top_mean, top_lo, top_hi = confidence_interval(top_vendors)
11 low_mean, low_lo, low_hi = confidence_interval(low_vendors)
12
13 print(f"Top Vendors  95% CI: ({top_lo:.2f}%, {top_hi:.2f}%), Mean: {top_mean:.2f}%")
14 print(f"Low Vendors  95% CI: ({low_lo:.2f}%, {low_hi:.2f}%), Mean: {low_mean:.2f}%")
```

*Python — 95% CI via t-distribution; top vendors 30.04% ± 0.5%, low vendors 41% ± wide range*

## 3.12 Statistical Hypothesis Test

A two-sample Welch's t-test tests whether the profit margin difference between top-quartile and bottom-quartile vendors is statistically real. **t = 9.68, p < 0.0001** conclusively rejects the null hypothesis — vendor selection has a measurable, statistically significant impact on profitability. This is not noise.
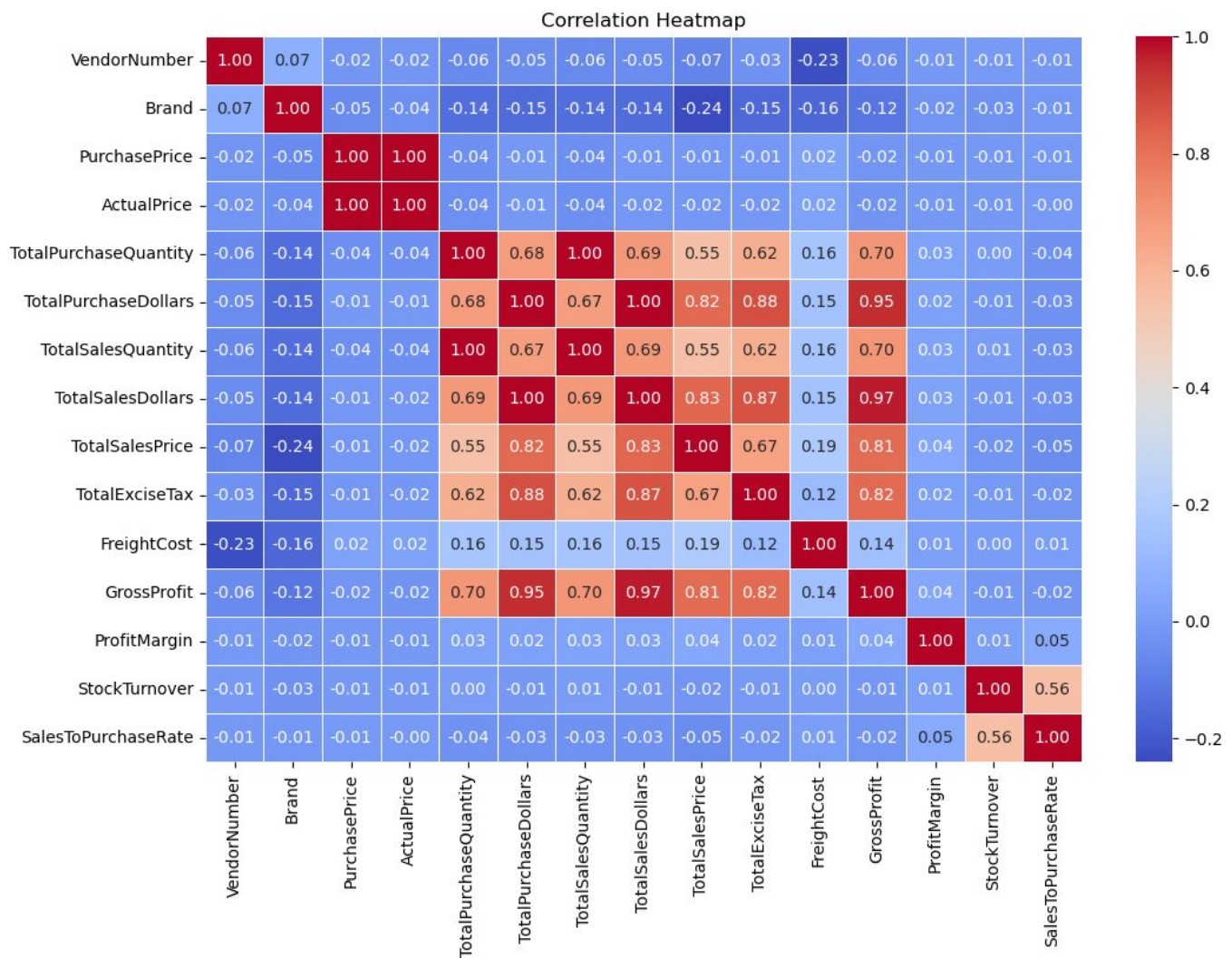
```python
1  # Two-sample t-test: do top & low vendors have different profit margins?
2  from scipy.stats import ttest_ind
3
4  top_vendors = df[df["TotalSalesDollars"] >= df["TotalSalesDollars"].quantile(0.75)]["ProfitMargin"].dropna()
5  low_vendors = df[df["TotalSalesDollars"] <= df["TotalSalesDollars"].quantile(0.25)]["ProfitMargin"].dropna()
6
7  t_stat, p_value = ttest_ind(top_vendors, low_vendors, equal_var=False)
8
9  print(f"T-Statistic : {t_stat:.4f}")   # 9.6799
10 print(f"P-Value     : {p_value:.6f}")  # < 0.0001
11
12 if p_value < 0.05:
13     print("REJECT H₀ — significant difference in profit margins confirmed.")
```

*Python — SciPy ttest_ind (equal_var=False); t=9.68, p<0.0001 → reject H▇*
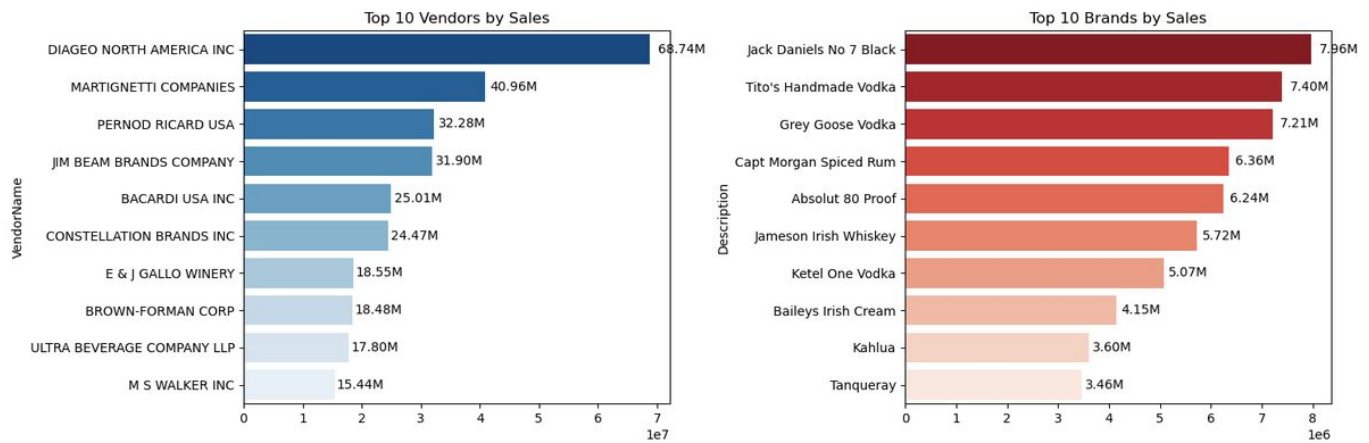
### 3.13 Correlation Heatmap Output

The Pearson matrix across all 15 KPI columns. Key signals: **0.999** between purchase and sales quantity (efficient aggregate turnover), **0.97** between sales dollars and gross profit (revenue drives profit tightly), and **~0.00** between purchase price and revenue/profit (demand is price-inelastic — raising prices does not dent volume).



*Output — Pearson correlation heatmap across all 15 KPI columns*

## 3.14 Top Vendors & Brands by Sales

Vendors and brands grouped by total sales dollars, top 10 each plotted as labelled horizontal bar charts. **Diageo North America** leads vendors at $68.74M — 68% ahead of the second-ranked supplier. **Jack Daniel's No. 7 Black** leads brands at $7.96M, followed by Tito's Handmade Vodka ($7.40M) and Grey Goose ($7.21M).



*Output — Top 10 vendors (left) and top 10 brands (right) by total sales dollars*

## Pipeline at a Glance

| # | File | What the Code Does | Output |
|---|------|--------------------|--------|
| 1 | Database_Setup.sql | CREATE DB + 6 schemas; LOAD DATA INFILE (15.6M rows) | 6 clean tables in MySQL |
| 2 | Database_Setup.sql | STR_TO_DATE() UPDATE + ALTER TABLE MODIFY on all date columns | Native DATE columns throughout |
| 3 | EDA.ipynb | SQLAlchemy connect; inspect_schema; row-count loop | Verified 15.6M rows across 6 tables |
| 4 | EDA.ipynb | CSV header/row inspection; single-vendor drill-down (Vendor 4466) | Join logic validated pre-CTE |
| 5 | EDA.ipynb | 3-CTE merge: freightSummary + PurchaseSummary + SalesSummary | vendor_sales_summary (10,692 rows) |
| 6 | EDA.ipynb | dtypes + null check; fillna; str.strip; inf→NaN→0 | Clean DataFrame, no nulls or infinities |
| 7 | EDA.ipynb | GrossProfit, ProfitMargin, StockTurnover, SalesToPurchaseRate | 4 engineered KPI columns |
| 8 | EDA.ipynb | DDL CREATE TABLE + to_sql(replace) + CSV backup | Persistent summary in MySQL + backup |
| 9 | VPA.ipynb | describe().T; histplot grid; boxplot grid | Outlier flags; distribution shapes known |
| 10 | VPA.ipynb | Dead-stock filter (GrossProfit>0, Margin>0, Sales>0) | 8,564 active SKUs / 2,128 flagged |
| 11 | VPA.ipynb | quantile(0.15/0.85) segmentation + scatter plot | 8 hidden-gem target brands identified |
| 12 | VPA.ipynb | PurchaseContribution% + cumsum + Pareto/donut charts | 65.34% spend concentrated in top 10 |
| 13 | VPA.ipynb | StockTurnover<1 filter + groupby mean + UnsoldInventoryValue | $2.71M locked capital; top slow-movers |
| 14 | VPA.ipynb | qcut() bulk tiers + groupby unit cost | 74% cost reduction Small→Large orders |
| 15 | VPA.ipynb | 95% CI + two-sample Welch t-test | t=9.68, p<0.0001 → vendor selection matters |

## Recommendations

The following recommendations are derived directly from the analytical outputs above. Each is mapped to the specific code technique that surfaces it.

| Priority | Area | Finding from Analysis | Recommended Action |
|---|---|---|---|
| ■ High | Pricing & Margins | 2,128 SKUs have negative or zero gross profit (filter step §3.5) | Audit every loss-making SKU. Renegotiate purchase price with vendor or raise retail price. Target: eliminate negative-margin records within 90 days. |
| ■ High | Hidden Gem Promotion | 8 brands: margin >56% but sales <$286 (segmentation §3.6) | Run targeted shelf-placement and digital promotions for these brands. No discounting needed — margins are strong. Goal: push each brand past the Q15 sales threshold. |
| ■ High | Vendor Diversification | Top 10 vendors = 65.34% of all procurement (Pareto §3.7) | Onboard 3–5 alternative suppliers in the highest-spend categories. Reduce any single vendor below 10% of total spend to limit disruption risk. |
| ■ Medium | Dead Stock Recovery | $2.71M locked in unsold inventory (unsold capital §3.9) | Initiate vendor return agreements, time-limited clearance promotions, or bundle slow SKUs with fast-movers. Prioritise vendors with StockTurnover = 0 first. |
| ■ Medium | Bulk Order Consolidation | 74% unit cost reduction from small to large order tiers (§3.10) | Consolidate fragmented small orders into large-tier runs for top-spend brands. Calculate minimum order quantity break-even per vendor before implementing. |
| ■ Medium | Slow-Mover Renegotiation | Several vendors show StockTurnover < 0.1 (slow-mover §3.8) | Renegotiate MOQ (minimum order quantity) terms with these vendors. Consider consignment or pay-on-scan arrangements to reduce capital exposure. |
| ■ Low | Pipeline Automation | CTE query and KPI calculation are fully repeatable | Schedule the EDA notebook to refresh weekly via a cron job or Airflow DAG. Connect Power BI to the live MySQL table via DirectQuery to eliminate manual CSV exports. |
| ■ Low | Margin-Volume Paradox | Low-volume vendors have higher avg margins (CI §3.11, t-test §3.12) | Test selective price increases on top-volume brands — even a 1–2% margin improvement on $68M Diageo sales = ~$700K additional gross profit. Model before committing. |

## Conclusion

The three files together form a reproducible, auditable pipeline that transforms 1.5 GB of raw transactional data into eight concrete business decisions. The SQL script handles one-time provisioning. The EDA notebook is the analytical engine — its CTE merge query and KPI engineering steps are the foundation every downstream insight depends on. The vendor performance notebook then applies fourteen distinct techniques to that foundation, each answering a specific operational question.

The most significant findings are quantified and statistically validated: **$2.71M in recoverable capital** sits in stagnant inventory; **65.34% of procurement** is concentrated in just 10 vendors — a disruption risk that one supply-chain event could trigger; **8 high-margin brands** are invisible to customers despite strong unit economics; and **bulk order consolidation** can cut unit costs by 74% with no renegotiation required. The two-sample t-test ($p<0.0001$) confirms that vendor selection is not incidental — it is the single most controllable driver of profitability in this dataset.

The pipeline is designed for refresh. Scheduling the EDA notebook weekly and connecting Power BI directly to the MySQL summary table via DirectQuery would make every KPI, chart, and recommendation in this report update automatically — turning a one-off analysis into a live operational intelligence system.

*Built with MySQL · Python · Power BI*