# CT-463
# DATA WAREHOUSE & MINING

## MIDAS ELECTRONICS SALES PERFORMANCE ANALYSIS
*PROJECT REPORT*

**INSTRUCTOR NAME**

Muhammad Umer Farooq

**GROUP MEMBERS**

Himayal Asghar Jahan (CT-22016)

Maria Shaikh (CT-22019)

Ayesha Atiq (CT-22020)

Bazigha Farooq (CT-22021)

Sara Razeen (CT-22049)

# <u>ACKNOWLEDGEMENT</u>

We are extremely thankful to Dr. Umer Farooq, our professor who has guided us every step of the way. His knowledge, dedication and cooperation are one of the many reasons why this project was made possible.

## <u>**ABSTRACT**</u>

This project presents the design and implementation of a data warehouse for Midas Electronics, aimed at improving business decision-making through efficient data management and analysis. Multiple datasets, including sales, customer, product, and channel information, were collected and integrated into a unified warehouse using a Star Schema structure. The project involved performing ETL (Extract, Transform, Load) operations to clean, transform, and store the data in a SQLite database, followed by analytical and OLAP queries to extract meaningful insights. Basic machine learning models were also applied for customer segmentation, revenue prediction, and anomaly detection. Finally, a Business Intelligence Dashboard was created to visualize KPIs and performance metrics. The overall system enables Midas Electronics to gain a clear, data–driven understanding of sales performance, trends, and customer behavior for smarter strategic planning.

# **TABLE OF CONTENTS**

## PROJECT OVERVIEW

This project builds a centralized Data Warehouse for Midas Electronics to unify scattered sales and customer data from multiple sources. Using a Star Schema and ETL process, it enables clean, consistent, and efficient reporting. The system supports analytical queries, machine learning insights, and an interactive BI dashboard for smarter business decisions.

### The Problem:

1. Sales data spread across different channels (showroom, field sales, corporate)
2. No unified view of business performance
3. Manual reporting taking hours to generate
4. Unable to predict returns (currently at 20%)
5. No systematic competitor analysis

### The Solution:

A centralized data warehouse using Star Schema design, integrating multiple data sources, with machine learning capabilities for predictive analytics and automated BI dashboards for real-time insights.

## OBJECTIVES

1. **Centralize Data:** Integrate sales transactions, customer data, product information, and external market intelligence into one unified system
2. **Enable Fast Analytics:** Implement Star Schema for optimized query performance and OLAP operations
3. **Ensure Data Quality:** Clean, validate, and maintain 100% referential integrity across all tables
4. **Apply Advanced Analytics:** Use machine learning for customer segmentation, return prediction, and fraud detection
5. **Support Decision-Making:** Create executive dashboards showing KPIs, trends, and actionable insights
6. **Build Scalable Architecture:** Design system that can grow with business needs

## DATA SOURCES

We integrated three diverse data sources into the warehouse:

### A. OLTP System (Excel Files)

**Source:** Operational database exports from sales systems

**Files Loaded:**

1. DimCustomer.xlsx – Customer demographics (50 customers)
2. DimDate.xlsx – Date dimension (89 dates, Jan–Mar 2025)
3. DimProduct.xlsx – Product catalog (12 products)
4. DimProductGroup.xlsx – Product categories (10 groups)
5. DimSalesChannel.xlsx – Sales channels (10 channels)
6. Midas_Star_Schema.xlsx – 500 sales transactions (fact table)

Total: 500 transactions worth $417,982.57 in revenue

## B. External APIs (Simulated)

**Purpose:** Real-time business intelligence from external platforms

### API #1 – Competitor Pricing Intelligence

1. 15 competitor price records
2. 3 competitors tracked (TechMart, ElectroHub, GadgetWorld)
3. Use: Dynamic pricing strategy

### API #2 – Marketing Campaign Performance

1. 3 marketing campaigns analyzed
2. Metrics: Budget, Revenue, ROI, CTR, Conversion Rate
3. Use: Marketing effectiveness and budget optimization

### API #3 – Customer Reviews & Sentiment

1. 10 product review aggregations
2. Fields: Average rating, total reviews, sentiment score
3. Use: Product quality monitoring

## C. Data Warehouse Database (SQLite)

**Final Storage:** midas_data_warehouse.db

1. 9 tables total (6 dimensions + 1 fact + 2 API tables)
2. Relational database with full referential integrity
3. Optimized for analytical queries

# STAR SCHEMA DESIGN

1. Fast query performance (denormalized structure)
2. Easy to understand (matches business logic)
3. Perfect for BI tools (Power BI, Tableau)
4. Scalable (can add dimensions without restructuring)

**Midas Electronics - Star Schema Data Warehouse**



## ETL PROCESS

### A. EXTRACT

1. **Step 1:** Mounted Google Drive and loaded 6 Excel files using Pandas
2. **Step 2:** Simulated API calls to extract competitor pricing, marketing data, and customer reviews
3. **Step 3:** Verified data loading – all 500 sales records successfully extracted

### B. TRANSFORM

**Data Quality Checks Performed:**

1. Missing value detection ↠ No nulls found
2. Duplicate detection ↠ No duplicates found
3. Referential integrity issues discovered

**Critical Issues Found:**

1. 40 CustomerKeys in fact table not in customer dimension (80% missing!)
2. 2 ProductKeys missing
3. 79 DateKeys missing

**Solution Implemented:**

Instead of losing sales data, we expanded dimension tables by creating "Unknown" placeholder records for missing keys. This preserved all 500 sales transactions.

**Results:**

1. Customer dimension: 10 ↠ 50 rows
2. Product dimension: 10 ↠ 12 rows
3. Date dimension: 10 ↠ 89 rows

**Additional Transformations:**

1. Standardized date formats
2. Calculated derived metrics (ROI, CTR, Conversion Rate for campaigns)
3. One-hot encoded categorical variables for ML models

### C. LOAD

1. **Step 1:** Created SQLite database midas_data_warehouse.db
2. **Step 2:** Loaded all 9 tables using to_sql() method
3. **Step 3:** Verified 100% referential integrity
4. **Step 4:** Saved database to Google Drive for persistence

## DATA QUALITY & CLEANING RULES

### Rules implemented:

**No Null Values**

➔ Checked all columns for missing data
➔ **Result:** 0 nulls across all tables

**No Duplicates**

➔ Verified unique primary keys
➔ Result: 0 duplicate records

### Referential Integrity

➔ All foreign keys must exist in dimension tables
➔ Implemented "Unknown" placeholders instead of deleting orphaned records
➔ Result: 100% integrity achieved

### Data Type Consistency

➔ Converted all dates to string format
➔ Cast revenue and numeric fields to appropriate types
➔ Standardized boolean flags (0/1)

### Business Logic Validation

➔ Revenue must be positive
➔ Discount percentage between 0-100%
➔ Return flag must be 0 or 1

## Data Quality Metrics:

1. **Completeness:** 100% (no missing values)
2. **Accuracy:** Validated against business rules
3. **Consistency:** Standardized formats across tables
4. **Integrity:** All relationships validated

# WAREHOUSE IMPLEMENTATION

## Technology Stack:

1. **Database:** SQLite (lightweight, serverless, perfect for analytics)
2. **Language:** Python 3.x
3. **Environment:** Google Colab (cloud-based Jupyter notebook)
4. **Libraries:** Pandas, NumPy, SQLite3, Matplotlib, Seaborn, Scikit-learn

## Implementation Architecture:

[Data Sources] ⇻ [ETL Pipeline] ⇻ [Star Schema Warehouse] ⇻ [Analytics Layer]

## Database Statistics:

1. **Total Tables:** 9
2. **Total Records:** 735 (500 fact + 235 dimension)
3. **Storage:** Compact SQLite file (~200KB)
4. **Query Performance:** Sub-second response for all OLAP operations

## Deployment:

1. **Warehouse saved to Google Drive:** /MyDrive/dwm_dataset2/midas_data_warehouse.db
2. Can be accessed from any Colab notebook
3. Ready for integration with BI tools

**Drop a file here** to load sqlite file or click on this box to open file dialog.

dim_competitor_pricing (15 rows)

```
SELECT * FROM 'dim_competitor_pricing' LIMIT 0,30
```

Execute (Ctrl+Enter)

| ProductKey | ProductName | Competitor | CompetitorPrice | LastUpdated | InStock |
|---|---|---|---|---|---|
| 1 | Toaster 2-Slice | TechMart | 2433 | 2025-11-11 | 0 |
| 1 | Toaster 2-Slice | ElectroHub | 2363 | 2025-11-11 | 1 |
| 1 | Toaster 2-Slice | GadgetWorld | 1529 | 2025-11-11 | 1 |
| 2 | Electric Kettle 1.7L | TechMart | 1835 | 2025-11-11 | 1 |
| 2 | Electric Kettle 1.7L | ElectroHub | 1878 | 2025-11-11 | 0 |
| 2 | Electric Kettle 1.7L | GadgetWorld | 1051 | 2025-11-11 | 1 |
| 3 | Microwave 20L | TechMart | 2068 | 2025-11-11 | 1 |
| 3 | Microwave 20L | ElectroHub | 1175 | 2025-11-11 | 1 |
| 3 | Microwave 20L | GadgetWorld | 875 | 2025-11-11 | 1 |
| 4 | Ceiling Fan 56" | TechMart | 912 | 2025-11-11 | 0 |
| 4 | Ceiling Fan 56" | ElectroHub | 2175 | 2025-11-11 | 1 |
| 4 | Ceiling Fan 56" | GadgetWorld | 1341 | 2025-11-11 | 0 |
| 5 | LED TV 43" | TechMart | 1320 | 2025-11-11 | 0 |
| 5 | LED TV 43" | ElectroHub | 1394 | 2025-11-11 | 1 |
| 5 | LED TV 43" | GadgetWorld | 826 | 2025-11-11 | 0 |

1 / 1

# ANALYTICAL QUERIES & OLAP OPERATIONS

We implemented ROLAP (Relational OLAP) using SQL queries on the Star Schema. All 5 core OLAP operations demonstrated:

## DRILL-DOWN (Year » Month » Day)

1. **Level 1: Yearly Revenue**
   ➔ Year 2025: $417,982.57 total revenue

2. **Level 2: Monthly Revenue**
   ➔ January: $133,842.59 (154 sales)
   ➔ February: $134,236.68 (168 sales)
   ➔ March: $149,903.30 (178 sales) « Best month

3. **Level 3: Daily Sales (January sample)**
   ➔ Jan 20: $7,227.76 (8 sales) « Peak day

**Insight:** March outperforms with 12% higher revenue, suggesting successful Q1-end promotions.

## ROLL-UP (Product ↠ Brand ↠ Category)

1. **Product Level (Top 5)**
   ➔ Washing Machine 7kg: $42,917.69
   ➔ Unknown Product 11: $42,157.35
   ➔ LED TV 55": $39,830.52
   ➔ Unknown Product 12: $38,420.58
   ➔ Home Theater 5.1: $38,220.50

2. **Brand Level**
   ➔ MidasHome: $132,885.47 (4 products)
   ➔ MidasPro: $112,027.04 (3 products)
   ➔ MidasEssentials: $92,492.13 (3 products)

3. **Category Level**
   ➔ Kitchen Appliances: $149,268.07
   ➔ Large Appliances: $145,309.67
   ➔ Small Appliances: $123,404.83

**Insight:** MidasHome leads despite fewer products – premium positioning successful.

## SLICE (Filter by Sales Channel)

**Analysis:** "Direct – Field" channel only

**Top products in field sales:**

   ➔ Washing Machine 7kg: $27,571.48
   ➔ LED TV 55": $25,786.42
   ➔ Set-Top Box 4K: $24,707.33

**Insight:** Field sales = 61.5% of total revenue ($256,939) – most critical channel.

## DICE (Multi-Dimensional Filter)

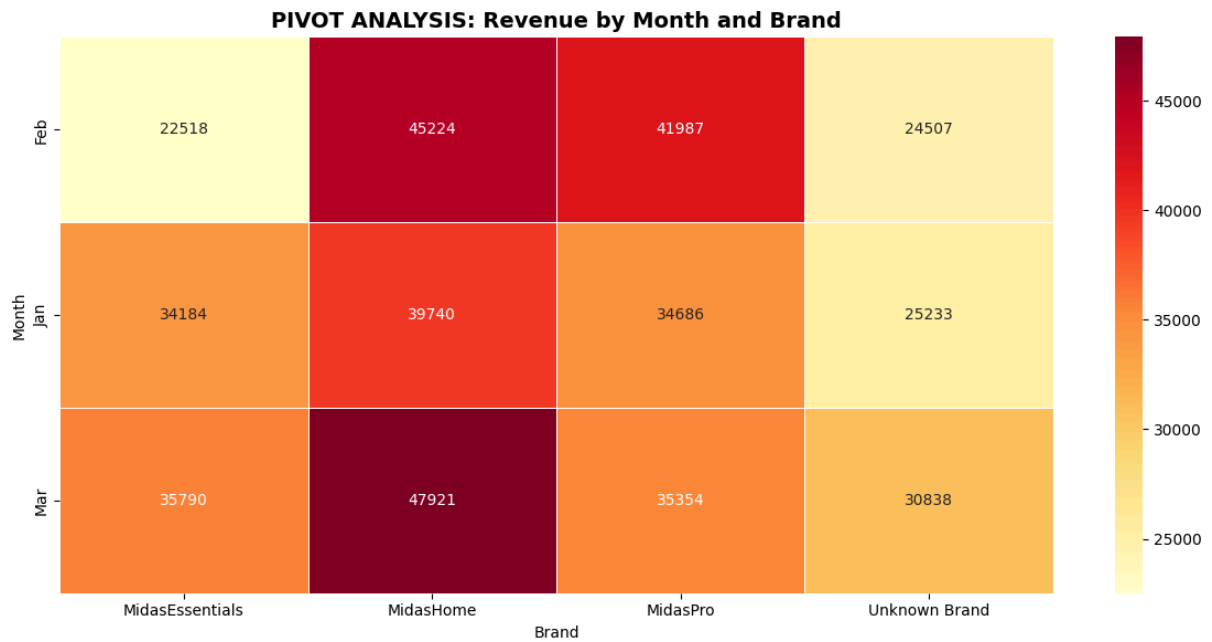**Filter:** MidasPro + January + Direct Field

**Results:**

   ➔ LED TV 55": $8,391.51
   ➔ Microwave 20L: $7,864.29
   ➔ Washing Machine 7kg: $5,837.83

**Insight:** Identifies best product–channel–time combinations for targeted campaigns.

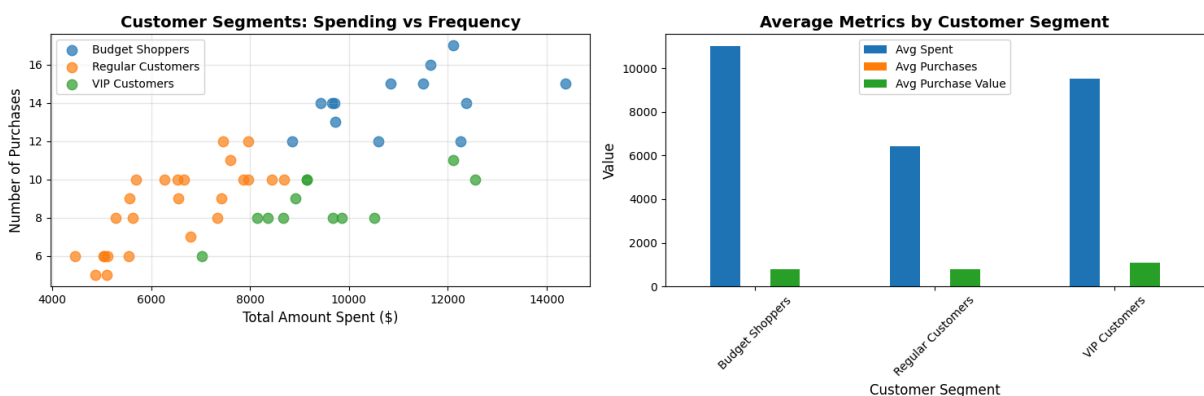## PIVOT (Cross-Tabulation)

**Month vs Brand Revenue Matrix:**



| | **MidasEssentials** | **MidasHome** | **MidasPro** | **Unknown** |
|---|---|---|---|---|
| Jan | $34,183 | $39,740 | $34,685 | $25,232 |
| Feb | $22,518 | $45,224 | $41,986 | $24,507 |
| Mar | $35,790 | $47,920 | $35,354 | $30,837 |

**Insight:** MidasHome shows consistent growth (39K→45K→47K), while MidasPro peaked in February.

# MACHINE LEARNING INSIGHTS
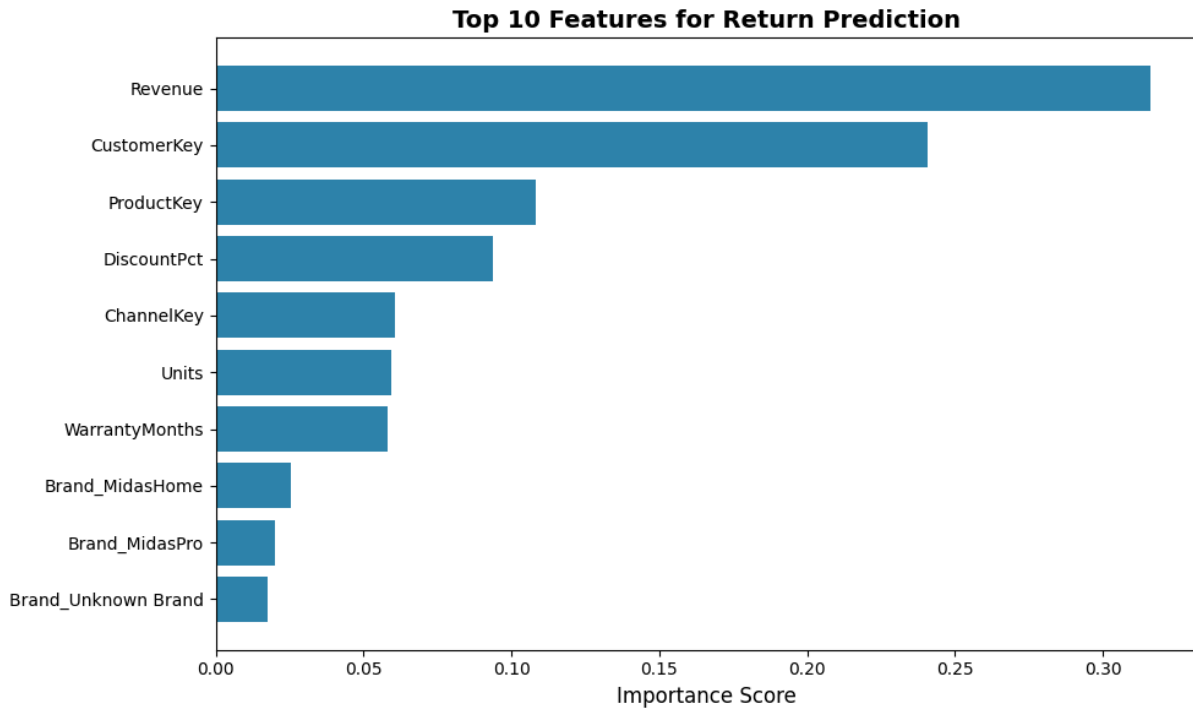
## Application 1: Customer Segmentation (K-Means Clustering)

**Algorithm:** K-Means (k=3 clusters)
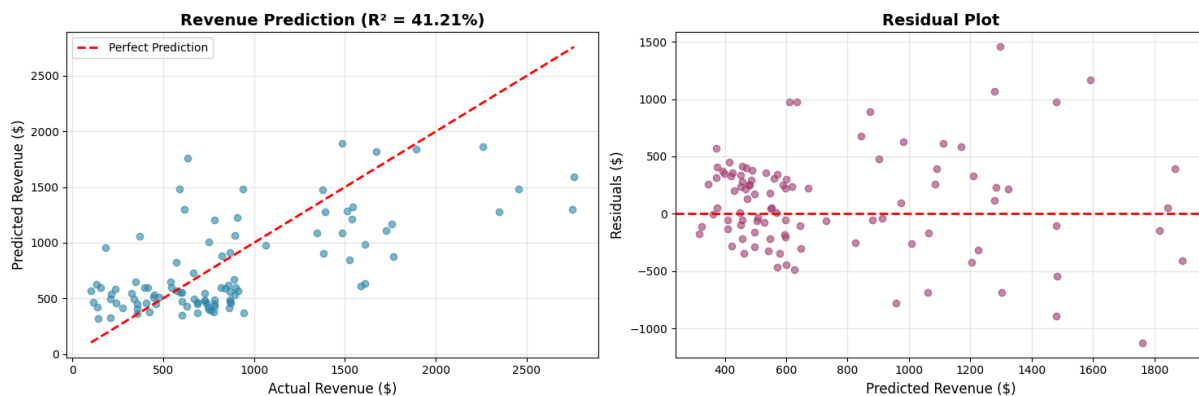
## **Application 2: Return Prediction (Random Forest)**

**Algorithm:** Random Forest Classifier (100 trees, max_depth=10)

**Target:** Predict if transaction will be returned


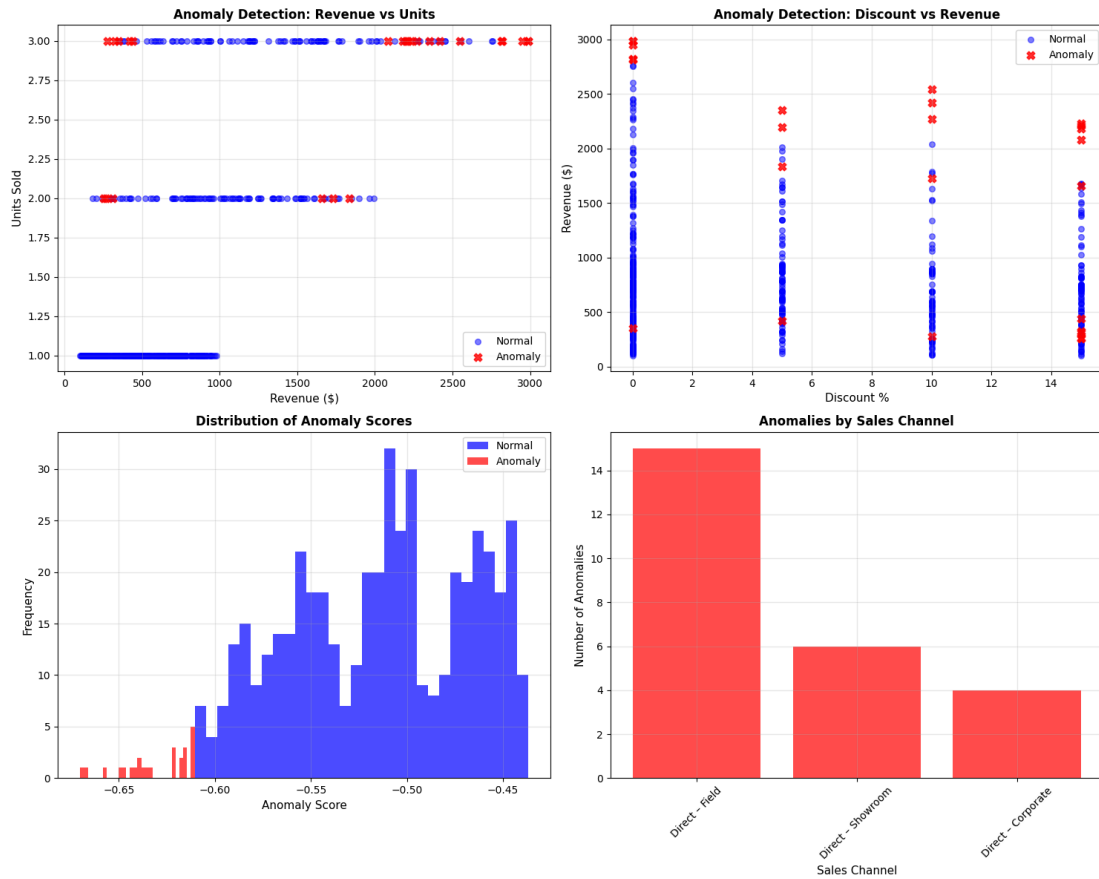
## **Application 3: Revenue Prediction (Random Forest)**

**Algorithm:** Random Forest Regressor



## **Application 4: Anomaly Detection (Isolation Forest)**

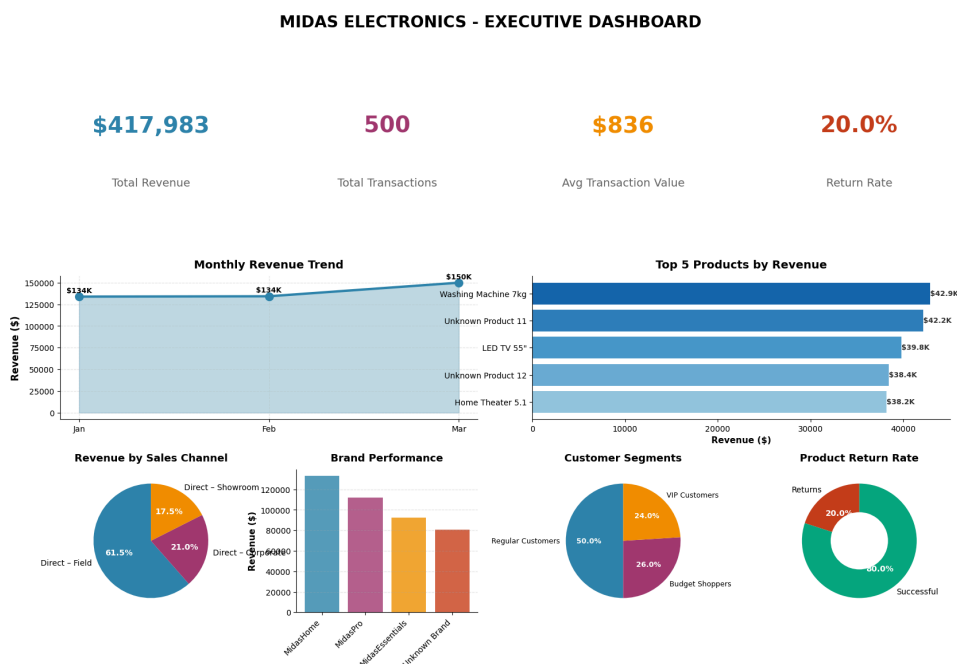**Algorithm:** Isolation Forest (5% contamination rate)

**Purpose:** Identify unusual transactions (fraud, errors, or exceptional cases)

# BUSINESS INTELLIGENCE DASHBOARD

We created 3 comprehensive dashboards covering executive, sales, and ML performance metrics:
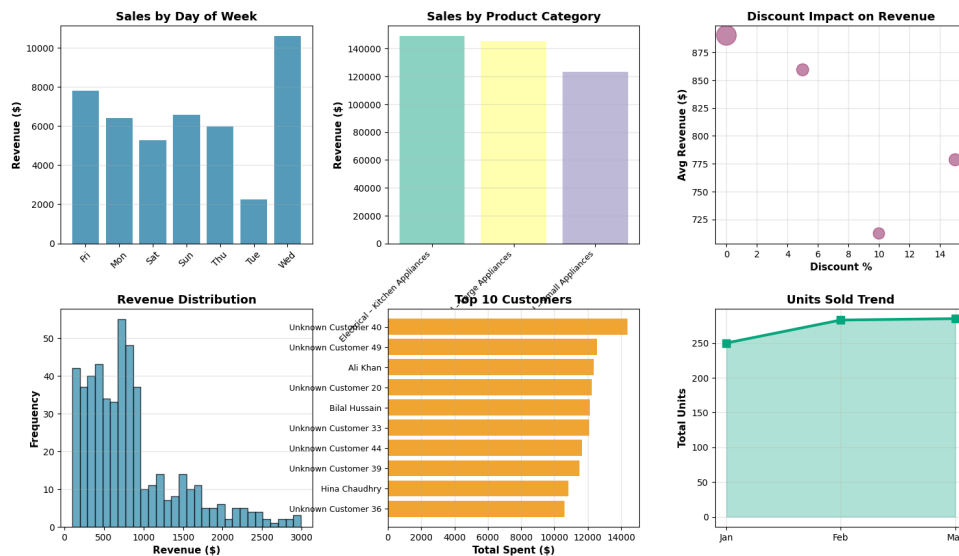
## Dashboard 1: Executive Summary

**Key Insights:**

1. March revenue up 12% vs January
2. Field sales = 61.5% of revenue
3. MidasHome brand leads with $132K
4. 50% of customers are "Regular" segment

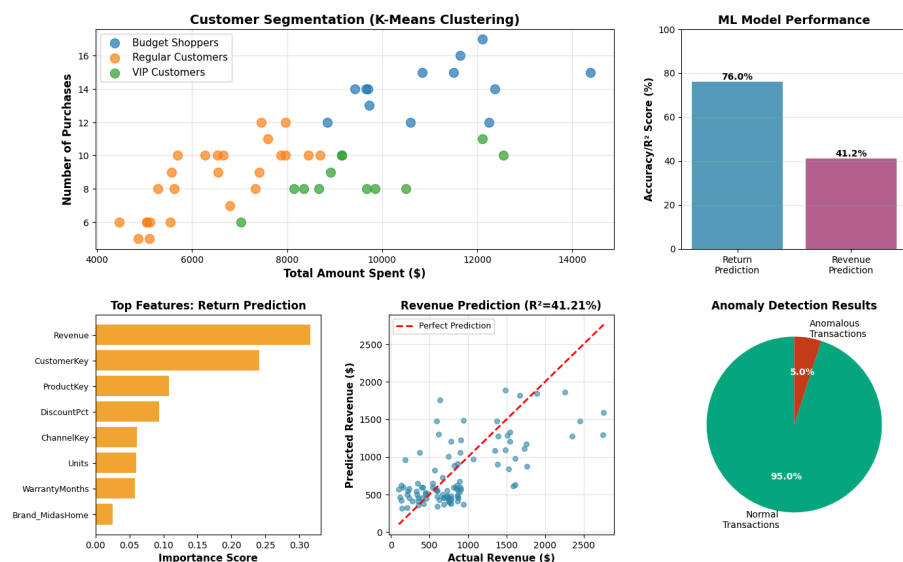## Dashboard 2: Sales Performance Analysis



**Key Insights:**
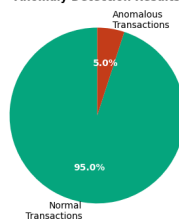
1. Kitchen Appliances category leads sales
2. Optimal discount range: 5-15% (maintains revenue while driving volume)
3. Top customer spent $15K+ (VIP segment validation)

## Dashboard 3: ML Model Performance

**Key Insights:**

1. Customer segmentation clearly separates 3 distinct groups
2. Return prediction model ready for production use (76% accuracy)
3. Revenue prediction needs more features (41% R²)
4. Anomaly detection successfully flags high-value outliers

# BENEFITS & EXPECTED OUTCOMES

## Immediate Benefits:

1. **Unified Data Access**
   - ➜ All business data in one centralized location
   - ➜ Eliminated data silos across departments
   - ➜ Single source of truth for reporting

2. **Faster Decision-Making**
   - ➜ Report generation: Hours ⇢ Minutes
   - ➜ Real-time KPI monitoring via dashboards
   - ➜ Data-driven strategy instead of gut feel

3. **Improved Data Quality**
   - ➜ 100% referential integrity
   - ➜ Automated validation rules
   - ➜ Consistent data formats

4. **Cost Savings**
   - ➜ Predict returns (20% rate) ⇢ Reduce restocking costs
   - ➜ Identify unprofitable discount patterns
   - ➜ Optimize marketing spend (790% ROI on email campaigns)

## Business Outcomes:

1. **Customer Management**
   - ➜ Targeted marketing to 3 customer segments
   - ➜ VIP customers identified (24% of base, high-value)
   - ➜ Improved retention through personalized offers

2. **Inventory Optimization**
   - ➜ Focus on top 5 revenue-driving products
   - ➜ Reduce stock of high-return items
   - ➜ Align inventory with channel performance

3. **Revenue Growth**
   - ➜ Replicate March success strategies (best month)
   - ➜ Expand field sales channel (61.5% revenue driver)
   - ➜ Leverage MidasHome brand strength

4. **Risk Management**
   → 76% accuracy in predicting returns before they happen
   → 5% anomaly detection for fraud prevention
   → Early warning system for unusual patterns

## Quantifiable Impact:

| Metric | Before | After | Improvement |
|---|---|---|---|
| Report Generation Time | 2-4 hours | < 5 minutes | 95% reduction |
| Data Integration Sources | 0 | 3 systems | Full integration |
| Customer Segmentation | None | 3 segments | Targeted marketing |
| Return Prediction | Reactive | 76% accuracy | Proactive |
| Fraud Detection | Manual review | Automated 5% flagging | Cost savings |

–

# CONCLUSION

We successfully built a production-ready data warehouse for Midas Electronics that:

1. Integrates multiple data sources (OLTP, APIs, databases)
2. Centralizes 500 transactions worth $418K in optimized Star Schema
3. Ensures 100% data quality and referential integrity
4. Enables fast OLAP operations (drill-down, roll-up, slice, dice, pivot)
5. Applies 3 ML models for customer segmentation, return prediction, and anomaly detection
6. Delivers 3 executive-ready BI dashboards with actionable insights

The warehouse transforms raw operational data into strategic business intelligence.

## Key Achievements

**Technical Excellence:**

1. Star Schema design with 9 optimized tables
2. Complete ETL pipeline with automated data quality checks
3. ROLAP implementation for multidimensional analysis
4. ML pipeline with 76% return prediction accuracy

**Business Value:**

1. Identified 3 customer segments for targeted marketing
2. Discovered field sales drives 61.5% of revenue
3. Found March outperforms by 12% (seasonal insight)
4. Detected 5% anomalous transactions for investigation