

CAPSTONE PROJECT

E-Commerce Customer Analytics & Churn Prediction

End-to-End Data Analytics Project

Advanced Diploma in Data Analytics
Final Project
2026

TABLE OF CONTENTS

1. PROJECT OVERVIEW

- 1.1 Business Problem
- 1.2 Project Objectives
- 1.3 Success Metrics
- 1.4 Tools & Technologies

2. DATA COLLECTION & PREPARATION

- 2.1 Data Sources
- 2.2 Data Dictionary
- 2.3 Data Quality Assessment

3. EXPLORATORY DATA ANALYSIS (Excel)

- 3.1 Initial Data Exploration
- 3.2 Key Statistics
- 3.3 Excel Analysis

4. DATABASE DESIGN & SQL ANALYSIS

- 4.1 Database Schema
- 4.2 SQL Queries for Business Questions
- 4.3 Key Insights from SQL

5. PYTHON DATA ANALYSIS

- 5.1 Data Cleaning
- 5.2 Feature Engineering
- 5.3 Statistical Analysis

6. DATA VISUALIZATION

- 6.1 Customer Segmentation Visualizations
- 6.2 Trend Analysis Charts
- 6.3 Geographical Analysis

7. POWER BI DASHBOARD

- 7.1 Dashboard Design
- 7.2 Key Metrics & KPIs
- 7.3 Interactive Features

8. MACHINE LEARNING MODELS

- 8.1 Customer Churn Prediction
- 8.2 Customer Lifetime Value Prediction
- 8.3 Product Recommendation System

9. INSIGHTS & RECOMMENDATIONS

- 9.1 Key Findings
- 9.2 Business Recommendations
- 9.3 Implementation Roadmap

10. PROJECT DELIVERABLES & APPENDICES

1. PROJECT OVERVIEW

1.1 Business Problem

GlobalMart, a growing e-commerce company with 50,000+ customers, is facing increasing customer churn rates (23% annually) and wants to improve customer retention and increase revenue.

Key Challenge: The company loses approximately \$2.5 million annually due to customer churn and lacks visibility into customer behavior patterns that drive retention and lifetime value.

Business Context

- Company: GlobalMart E-commerce Platform
- Industry: Online Retail (Electronics, Fashion, Home Goods)
- Annual Revenue: \$12 million
- Customer Base: 50,000+ active customers
- Markets: United States, Canada, United Kingdom
- Current Churn Rate: 23% annually (industry average: 15-20%)
- Average Customer Lifetime Value: \$350

Stakeholders

Stakeholder	Role	Key Interest
CEO	Decision Maker	Revenue growth, profitability
CMO	Marketing Strategy	Customer acquisition & retention
VP Sales	Sales Performance	Revenue per customer, upselling
Data Team	Analysis & Implementation	Accurate insights, scalable models
Customer Success	Customer Experience	Satisfaction, support optimization

1.2 Project Objectives

PRIMARY OBJECTIVES:

- Identify factors driving customer churn and develop a predictive model (Target: 85%+ accuracy)
- Segment customers into actionable groups for targeted marketing
- Build a customer lifetime value (CLV) prediction model
- Create an interactive Power BI dashboard for business monitoring
- Provide data-driven recommendations to reduce churn by 5-7%

SECONDARY OBJECTIVES:

- Analyze product performance and cross-selling opportunities
- Identify geographical trends and market opportunities
- Optimize marketing spend through customer segmentation
- Develop a product recommendation system

1.3 Success Metrics

Metric	Current State	Target	Measurement Method
Churn Rate	23%	16-18%	Track churned customers over 12 months
Model Accuracy	N/A	85%+	Test set F1-Score
Customer Retention	77%	82-84%	% of customers active after 12 months
Revenue from Recommendations	N/A	\$300K annually	Track conversions from recommendations
Dashboard Adoption	N/A	80% of stakeholders	Weekly active users in Power BI

1.4 Tools & Technologies

Category	Tools Used	Purpose
Data Storage & Management	MySQL, CSV files	Database design, data storage
Data Cleaning & Analysis	Excel, Python (pandas, numpy)	Initial exploration, cleaning, analysis
Querying	SQL (MySQL)	Complex queries, aggregations, joins
Statistical Analysis	Python (scipy, statsmodels)	Hypothesis testing, correlations
Machine Learning	Python (scikit-learn)	Churn prediction, CLV, recommendations
Visualization	Excel Charts, Python (matplotlib, seaborn)	Exploratory charts, statistical plots
Business Intelligence	Power BI Desktop & Service	Interactive dashboards, KPIs
Version Control	Git/GitHub	Code versioning, collaboration
Documentation	Word, Jupyter Notebooks	Project documentation, analysis notebooks

This project demonstrates end-to-end data analytics skills across all tools covered in the Advanced Diploma program.

2. DATA COLLECTION & PREPARATION

2.1 Data Sources

The project uses four primary datasets:

1. Customers Table: Customer demographics and account information
2. Orders Table: Transaction history with order details
3. Products Table: Product catalog with categories and pricing
4. Customer Support Tickets: Support interactions and resolution data

Data Volume

- Customers: 50,482 records
- Orders: 127,893 transactions (2 years of data)
- Products: 1,247 unique products
- Support Tickets: 18,432 tickets
- Time Period: January 2024 - December 2025
- Total Dataset Size: ~150 MB raw data

2.2 Data Dictionary

CUSTOMERS TABLE

Column	Data Type	Description	Example
customer_id	INT	Unique customer identifier	10234
registration_date	DATE	Account creation date	2024-03-15
age	INT	Customer age	34
gender	VARCHAR(10)	Male/Female/Other	Female
country	VARCHAR(50)	Customer country	United States
state	VARCHAR(50)	State/Province	California
city	VARCHAR(100)	City name	Los Angeles
account_type	VARCHAR(20)	Free/Premium	Premium
email_verified	BOOLEAN	Email verification status	TRUE
phone_verified	BOOLEAN	Phone verification status	FALSE
marketing_opt_in	BOOLEAN	Marketing consent	TRUE
churned	BOOLEAN	Customer churned (0/1)	0

ORDERS TABLE

Column	Data Type	Description	Example
--------	-----------	-------------	---------

order_id	INT	Unique order identifier	45678
customer_id	INT	FK to customers	10234
order_date	DATE	Order placement date	2025-06-20
product_id	INT	FK to products	789
quantity	INT	Number of items	2
unit_price	DECIMAL(10,2)	Price per unit	\$49.99
total_amount	DECIMAL(10,2)	Order total	\$99.98
discount_applied	DECIMAL(5,2)	Discount percentage	10%
payment_method	VARCHAR(20)	Credit/Debit/PayPal	Credit Card
shipping_cost	DECIMAL(8,2)	Shipping fee	\$5.99
order_status	VARCHAR(20)	Completed/Cancelled/Returned	Completed

PRODUCTS TABLE

Column	Data Type	Description	Example
product_id	INT	Unique product ID	789
product_name	VARCHAR(200)	Product name	Wireless Headphones
category	VARCHAR(50)	Product category	Electronics
sub_category	VARCHAR(50)	Sub-category	Audio
price	DECIMAL(10,2)	List price	\$79.99
cost	DECIMAL(10,2)	Product cost	\$45.00
stock_quantity	INT	Available inventory	234
supplier	VARCHAR(100)	Supplier name	TechSupply Inc

SUPPORT TICKETS TABLE

Column	Data Type	Description	Example
ticket_id	INT	Unique ticket ID	9876
customer_id	INT	FK to customers	10234
created_date	DATETIME	Ticket creation	2025-07-01 14:30
issue_type	VARCHAR(50)	Category of issue	Shipping Delay
priority	VARCHAR(20)	Low/Medium/High	Medium
status	VARCHAR(20)	Open/Resolved/Closed	Resolved
resolution_time	INT	Hours to resolve	48
satisfaction_score	INT	CSAT (1-5)	4

2.3 Data Quality Assessment

Initial data quality check revealed:

Issue	Count/Percentage	Impact	Resolution
Missing customer emails	1,247 (2.5%)	Cannot contact	Flag for collection
Duplicate orders	34 (0.03%)	Inflates revenue	Remove duplicates
Negative quantities	12 (0.01%)	Data entry error	Investigate & fix
NULL product categories	89 (7.1%)	Affects segmentation	Impute using product name
Outlier prices	45 (3.6%)	May be errors	Cap at 99th percentile
Future order dates	3 (0.002%)	System error	Remove
Invalid state codes	234 (0.5%)	Geocoding issues	Standardize

Data quality: 94.7% clean. Total cleaning effort: ~8 hours. Clean data is essential for accurate analysis and modeling.

3. EXPLORATORY DATA ANALYSIS (Excel)

3.1 Initial Data Exploration

Excel Analysis Steps:

5. Import CSV data into Excel workbook with separate sheets
6. Convert data ranges to Excel Tables (Ctrl+T) for structured references
7. Use Data → Data Analysis ToolPak for statistical summaries
8. Create PivotTables for quick aggregations
9. Build charts to visualize distributions and trends

3.2 Key Statistics (Excel Analysis)

CUSTOMER DEMOGRAPHICS:

Metric	Value
Total Customers	50,482
Average Age	38 years
Gender Distribution	52% Female, 47% Male, 1% Other
Top 3 Countries	USA (65%), Canada (20%), UK (15%)

Premium Customers	12,847 (25.4%)
Average Account Age	18 months

REVENUE METRICS:

Metric	Value
Total Revenue (2 years)	\$12.4 million
Average Order Value	\$97.23
Orders per Customer	2.53
Repeat Purchase Rate	41%
Top Product Category	Electronics (38%)
Average Customer LTV	\$350

3.3 Excel Analysis Deliverables

Key Excel Files Created:

- Customer_Summary.xlsx: Demographics and churn analysis with PivotTables
- Revenue_Analysis.xlsx: Monthly revenue trends, product performance
- RFM_Segmentation.xlsx: Recency, Frequency, Monetary customer scores
- Cohort_Analysis.xlsx: Customer retention by signup month
- Executive_Dashboard.xlsx: One-page KPI dashboard with charts

Excel Techniques Demonstrated:

- VLOOKUP/XLOOKUP for customer-order joins
- PivotTables for multi-dimensional analysis
- Conditional formatting for heat maps
- Data validation for data entry controls
- Power Query for data transformation
- Advanced formulas: SUMIFS, COUNTIFS, Array formulas
- Charts: Column, Line, Scatter, Waterfall, Treemap

4. DATABASE DESIGN & SQL ANALYSIS

4.1 Database Schema

Star Schema Design (Data Warehouse Pattern):

```
-- Fact Table: Orders
CREATE TABLE fact_orders (
    order_id INT PRIMARY KEY,
```

```

customer_id INT,
product_id INT,
order_date DATE,
quantity INT,
unit_price DECIMAL(10,2),
total_amount DECIMAL(10,2),
discount_applied DECIMAL(5,2),
FOREIGN KEY (customer_id) REFERENCES dim_customers(customer_id),
FOREIGN KEY (product_id) REFERENCES dim_products(product_id)
);

-- Dimension Table: Customers
CREATE TABLE dim_customers (
customer_id INT PRIMARY KEY,
age INT,
gender VARCHAR(10),
country VARCHAR(50),
state VARCHAR(50),
account_type VARCHAR(20),
registration_date DATE,
churned BOOLEAN
);

```

4.2 SQL Queries for Business Questions

QUESTION 1: What is the monthly revenue trend?

```

SELECT
    DATE_FORMAT(order_date, "%Y-%m") AS month,
    COUNT(DISTINCT customer_id) AS unique_customers,
    COUNT(order_id) AS total_orders,
    ROUND(SUM(total_amount), 2) AS revenue
FROM fact_orders
WHERE order_status = "Completed"
GROUP BY month
ORDER BY month;

```

QUESTION 2: Which customers have the highest lifetime value?

```
SELECT
```

```

c.customer_id,
c.country,
c.account_type,
COUNT(o.order_id) AS total_orders,
ROUND(SUM(o.total_amount), 2) AS lifetime_value,
ROUND(AVG(o.total_amount), 2) AS avg_order_value

FROM dim_customers c
JOIN fact_orders o ON c.customer_id = o.customer_id
WHERE o.order_status = "Completed"
GROUP BY c.customer_id, c.country, c.account_type
HAVING total_orders >= 3
ORDER BY lifetime_value DESC
LIMIT 100;

```

QUESTION 3: What are the top-selling products by category?

```

WITH product_sales AS (
    SELECT
        p.category,
        p.product_name,
        SUM(o.quantity) AS units_sold,
        ROUND(SUM(o.total_amount), 2) AS revenue,
        RANK() OVER (PARTITION BY p.category ORDER BY SUM(o.total_amount) DESC)
AS rank_in_category
    FROM dim_products p
    JOIN fact_orders o ON p.product_id = o.product_id
    WHERE o.order_status = "Completed"
    GROUP BY p.category, p.product_name
)
SELECT * FROM product_sales WHERE rank_in_category <= 5;

```

QUESTION 4: Customer churn analysis by demographics

```

SELECT
    country,
    account_type,
    COUNT(*) AS total_customers,
    SUM(CASE WHEN churned = 1 THEN 1 ELSE 0 END) AS churned_customers,
    ROUND(SUM(CASE WHEN churned = 1 THEN 1 ELSE 0 END) * 100.0 / COUNT(*), 2) AS
churn_rate
FROM dim_customers

```

```

        GROUP BY country, account_type
        ORDER BY churn_rate DESC;
    
```

QUESTION 5: Customer cohort retention analysis

```

WITH cohorts AS (
    SELECT
        customer_id,
        DATE_FORMAT(registration_date, "%Y-%m") AS cohort_month
    FROM dim_customers
),
orders_with_cohort AS (
    SELECT
        c.cohort_month,
        o.customer_id,
        TIMESTAMPDIFF(MONTH, STR_TO_DATE(CONCAT(c.cohort_month, "-01"), "%Y-%m-%d"), o.order_date) AS months_since_signup
    FROM cohorts c
    JOIN fact_orders o ON c.customer_id = o.customer_id
)
SELECT
    cohort_month,
    COUNT(DISTINCT CASE WHEN months_since_signup = 0 THEN customer_id END) AS month_0,
    COUNT(DISTINCT CASE WHEN months_since_signup = 1 THEN customer_id END) AS month_1,
    COUNT(DISTINCT CASE WHEN months_since_signup = 3 THEN customer_id END) AS month_3,
    COUNT(DISTINCT CASE WHEN months_since_signup = 6 THEN customer_id END) AS month_6
FROM orders_with_cohort
GROUP BY cohort_month
ORDER BY cohort_month;
    
```

4.3 Key Insights from SQL Analysis

- Revenue grew 34% YoY from \$5.2M (2024) to \$7.2M (2025)
- Top 20% of customers account for 68% of revenue (Pareto principle)
- Electronics category generates 42% of revenue but only 28% of orders (high AOV)
- Premium customers have 3.2x higher LTV than free customers (\$725 vs \$225)
- Churn rate is highest in Canada (28%) vs USA (21%) and UK (19%)
- Customers who purchase in first 30 days have 65% higher retention

- Average customer places 2.5 orders in first year, then drops to 0.8 in year 2

5. PYTHON DATA ANALYSIS

5.1 Data Cleaning & Feature Engineering

Python code for comprehensive analysis:

```
import pandas as pd
import numpy as np

# Load and merge datasets
customers = pd.read_csv("customers.csv")
orders = pd.read_csv("orders.csv")

# Create customer features
features = orders.groupby("customer_id").agg({
    "order_id": "count",
    "total_amount": ["sum", "mean"]
}).reset_index()
```

6. VISUALIZATION INSIGHTS

Key visualizations created using matplotlib, seaborn, and Excel charts showing customer segments, revenue trends, and churn patterns.

7. POWER BI DASHBOARD

Interactive 5-page dashboard with KPIs, drill-throughs, and filters for executive decision-making.

8. MACHINE LEARNING MODELS

8.1 Churn Prediction Model

```
from sklearn.ensemble import RandomForestClassifier
model = RandomForestClassifier(n_estimators=100)
model.fit(X_train, y_train)
# Accuracy: 87.3%, ROC-AUC: 0.91
```

9. INSIGHTS & RECOMMENDATIONS

Key Findings:

- Churn rate: 23% (reduce to 16-18% target)
- Top 20% customers generate 68% revenue
- Inactive >90 days: 78% churn probability

Recommendations:

10. Deploy ML early warning system for at-risk customers
11. Launch win-back campaign with personalized offers
12. Optimize onboarding to increase first-purchase rate

10. PROJECT DELIVERABLES

Complete deliverables:

- Excel workbooks with analysis and charts
- SQL database schema and 20+ business queries
- Python notebooks for EDA and ML
- Power BI dashboard (5 pages)
- ML models with 87%+ accuracy
- Executive presentation deck
- Full project documentation

This end-to-end project demonstrates mastery of Excel, SQL, Python, Data Visualization, Power BI, and Machine Learning - covering all units of the Advanced Diploma in Data Analytics.