COURSE HANDBOOK

# Machine Learning for Business Analytics

Integrating Strategic Analysis, Predictive Modeling, and Data-Driven Decision Making

**Chenhao Zhou**

Ph.D. Candidate, Supply Chain Management

Rutgers Business School

**2026 Edition**

Graduate-Level Course Material

Supply Chain Analytics Program

# Contents

Chenhao Zhou

# Preface

This handbook represents a comprehensive synthesis of business analytics principles and machine learning methodologies, designed specifically for graduate-level business professionals seeking to leverage data-driven approaches for strategic decision-making. The material herein progresses systematically from foundational analytical thinking through practical implementation, ensuring that readers develop both the conceptual understanding and technical competence necessary for effective application in organizational contexts.

The pedagogical approach adopted in this text prioritizes business value creation as the organizing principle for all technical content. Unlike conventional machine learning curricula that emphasize algorithmic sophistication, this handbook consistently frames analytical techniques within the context of return on investment, stakeholder value, and strategic competitive advantage. Every model presented is evaluated not merely on technical performance metrics but on its capacity to generate measurable business outcomes.

The handbook is structured in three interconnected parts. Part I establishes the foundational framework for business-driven analysis, developing the analytical mindset necessary for identifying high-value opportunities and translating business objectives into quantifiable metrics. Part II provides comprehensive coverage of machine learning implementation, with emphasis on predictive modeling, classification systems, and unsupervised learning techniques most relevant to business applications. Part III offers technical programming foundations for those seeking deeper implementation capabilities.

## Course Learning Objectives

> **Learning Objectives**
>
> Upon successful completion of this course, students will be able to:
> 1. Apply systematic analytical frameworks to transform business problems into data-driven solutions
> 2. Design and implement machine learning models that generate measurable business value
> 3. Evaluate model performance using both technical and business-oriented metrics
> 4. Communicate analytical findings to diverse stakeholder audiences effectively
> 5. Develop proficiency in Python programming for business analytics applications
> 6. Create end-to-end analytical pipelines from data acquisition through deployment

## Prerequisites and Preparation

Students are expected to possess foundational knowledge in statistics (descriptive and inferential), basic spreadsheet proficiency, and familiarity with business terminology. No prior programming

experience is required; Part III provides comprehensive coverage of necessary technical skills. Students are encouraged to establish a Python development environment prior to beginning the course using the Anaconda distribution, which provides all required analytical libraries.

PART I
# Business-Driven Data Analysis

*Building Analytical Thinking for Business Value Creation*

# 1  The Business Analytics Framework

> **Chapter Overview**
>
> This chapter establishes the foundational framework for translating business objectives into analytical opportunities. Students will learn to distinguish between traditional reporting and strategic analytics, develop the systematic approach necessary for identifying high-value analytical use cases, and formulate questions that directly tie to organizational value creation.

## 1.1  Learning Objectives

> **Learning Objectives**
>
> Upon completing this chapter, students will be able to:
> - Differentiate between descriptive reporting and prescriptive analytics
> - Apply the Business Analysis Framework to organizational challenges
> - Formulate ROI-focused analytical questions that drive business value
> - Translate strategic objectives into measurable analytical metrics

## 1.2  Theoretical Foundation: The Analytics Continuum

Business analytics represents a fundamental evolution from traditional reporting paradigms. While conventional business intelligence systems focus on retrospective description—answering the question "What happened?"—modern analytics extends this capability to diagnostic ("Why did it happen?"), predictive ("What will happen?"), and prescriptive ("What should we do?") dimensions. This progression from descriptive to prescriptive analytics constitutes the analytics continuum, a conceptual framework that guides organizational development of analytical capabilities.

The transition across this continuum requires corresponding evolution in organizational capabilities, including data infrastructure, analytical talent, and decision-making processes. Organizations at early stages of analytical maturity typically operate primarily in the descriptive domain, generating standard reports that summarize historical performance. As capabilities mature, organizations develop the capacity to identify causal relationships, predict future outcomes, and ultimately optimize decisions through sophisticated algorithmic approaches.

**Table 1.1: The Analytics Maturity Model**

| Level | Type | Question Addressed | Business Value |
|---|---|---|---|
| 1 | Descriptive | What happened? | Historical understanding |
| 2 | Diagnostic | Why did it happen? | Root cause identification |
| 3 | Predictive | What will happen? | Proactive planning |
| 4 | Prescriptive | What should we do? | Optimized decisions |

## 1.3 The Business Analysis Framework

Effective business analytics requires systematic methodology for translating organizational objectives into actionable analytical approaches. The Business Analysis Framework provides this structure, comprising four sequential phases: (1) Business Objective Definition, (2) Metric Translation, (3) Data Requirements Specification, and (4) Value Assessment. Each phase builds upon the preceding, ensuring tight alignment between analytical activity and business value creation.

### 1.3.1 Phase 1: Business Objective Definition

The foundational phase requires explicit articulation of business objectives in terms that permit quantification and measurement. Vague statements such as "improve customer satisfaction" must be refined to specific, measurable objectives: "Increase Net Promoter Score by 10 points among enterprise customers within fiscal year 2026." This specificity enables subsequent phases to proceed with clarity regarding success criteria.

> **Example: Customer Retention Objective**
>
> Consider a telecommunications company experiencing elevated customer churn. The generic objective "reduce customer attrition" lacks the specificity necessary for analytical action. Refined articulation: "Reduce monthly churn rate from 2.3% to 1.8% among high-value residential customers (defined as monthly revenue exceeding $150) through targeted retention interventions, generating estimated annual revenue preservation of $4.2M."

### 1.3.2 Phase 2: Metric Translation

Business objectives must be decomposed into constituent metrics amenable to measurement and analysis. This translation process identifies Key Performance Indicators (KPIs) that serve as operational definitions of success. Effective KPIs exhibit several essential characteristics: they are measurable with available data, actionable through organizational decisions, relevant to stated objectives, and time-bound with clear measurement periods.

**Table 1.2: Example KPI Framework for Customer Retention**

| KPI | Definition | Current | Target |
|---|---|---|---|
| Churn Rate | Monthly customer attrition | 2.3% | 1.8% |
| Customer LTV | Net present value of relationship | $2,400 | $2,800 |
| Retention Cost | Cost per retained customer | $150 | $120 |
| NPS | Net Promoter Score | 42 | 55 |

### 1.3.3 Phase 3: Data Requirements Specification

Having defined metrics, analysts must specify the data elements necessary for calculation and analysis. This phase involves systematic cataloging of data sources, assessment of data quality, and identification of gaps requiring remediation. The data requirements specification serves as both a planning document and a data governance artifact, ensuring clarity regarding data provenance, ownership, and quality standards.

### 1.3.4 Phase 4: Value Assessment

The concluding phase quantifies expected business value, enabling prioritization among competing analytical initiatives. Value assessment incorporates both direct financial impacts (revenue enhancement, cost reduction) and indirect benefits (improved decision velocity, competitive positioning). The assessment should also account for implementation costs, time-to-value, and risk factors that may affect realization of projected benefits.

## 1.4 Formulating ROI-Focused Analytical Questions

The distinction between low-value and high-value analytical questions frequently determines project success. Low-value questions typically exhibit vague formulation, unclear business connection, or absence of actionable implications. High-value questions demonstrate direct linkage to business metrics, specificity regarding scope and timeframe, and clear pathways to operational action.

**Table 1.3: Analytical Question Quality Assessment**

| Quality | Example Question | Issue/Strength |
|---|---|---|
| Low | What patterns exist in our data? | Vague, no clear business connection |
| Low | Can we predict something? | No specified outcome or value |
| High | Which segments drive 80% of profit? | Specific, actionable, value-linked |
| High | What is revenue impact of 2-hr response time reduction? | Quantified, decision-relevant |

## 1.5 Implementation: Python for Business Analysis

The following code implementation demonstrates the Business Analysis Framework applied to a customer retention scenario. This programmatic approach enables systematic documentation of analytical planning and facilitates communication with technical implementation teams.

```python
import pandas as pd
import numpy as np

# BUSINESS ANALYSIS FRAMEWORK IMPLEMENTATION

# Phase 1: Business Objective Definition
business_objective = {
    'statement': 'Increase customer retention by 15% to improve
        quarterly '
                 'recurring revenue and reduce customer acquisition
                    costs',
    'target_improvement': 0.15,
    'timeframe': 'Q1-Q4 2026',
    'estimated_value': 2500000  # Annual revenue impact
}

# Phase 2: Metric Translation
key_performance_indicators = {
    'churn_rate': {
        'description': 'Percentage of customers lost per month',
        'current_baseline': 0.023,  # 2.3%
        'target': 0.018,  # 1.8%
        'measurement_frequency': 'Monthly'
    },
    'customer_lifetime_value': {
        'description': 'Net present value of customer relationship',
        'current_baseline': 2400,
        'target': 2800,
        'measurement_frequency': 'Quarterly'
    },
    'retention_cost': {
        'description': 'Cost per retained customer',
        'current_baseline': 150,
        'target': 120,
        'measurement_frequency': 'Monthly'
    }
}

# Phase 3: Data Requirements
data_requirements = pd.DataFrame({
    'data_source': ['CRM System', 'Transaction Database',
                    'Marketing Platform', 'Support System',
```

```python
                             'Product Analytics'],
    'data_type': ['Customer demographics', 'Purchase history',
                  'Campaign engagement', 'Support interactions',
                  'Product usage'],
    'business_value': ['Segmentation targeting', 'Revenue analysis',
                       'ROI calculation', 'Satisfaction indicators',
                       'Engagement scoring'],
    'quality_status': ['High', 'High', 'Medium', 'Medium', 'Low'],
    'remediation_required': ['None', 'None', 'Integration needed',
                             'Data cleaning', 'Coverage expansion']
})

# Phase 4: Value Assessment
def calculate_project_value(baseline_churn, target_churn,
                            customer_base, avg_clv):
    """Calculate expected value from retention improvement."""
    customers_saved_monthly = customer_base * (baseline_churn -
        target_churn)
    annual_customers_saved = customers_saved_monthly * 12
    annual_value = annual_customers_saved * avg_clv
    return annual_value

projected_value = calculate_project_value(
    baseline_churn=0.023,
    target_churn=0.018,
    customer_base=50000,
    avg_clv=2400
)
print(f"Projected Annual Value: ${projected_value:,.0f}")
```

---

**Key Concepts**

- **The Analytics Continuum:** Progression from descriptive to prescriptive analytics
- **Business Analysis Framework:** Four-phase methodology for analytical planning
- **ROI-Focused Questions:** Direct linkage between analysis and business value
- **Metric Translation:** Decomposition of objectives into measurable KPIs
- **Value Assessment:** Quantification of expected business impact

Chenhao Zhou

# 2 Strategic Data Typology and Business Value

> **Chapter Overview**
>
> This chapter develops a systematic understanding of business data types and their strategic implications. Students will learn to categorize data assets according to their analytical potential, recognize the strategic value embedded in different data categories, and design data collection strategies that maximize business intelligence capabilities.

## 2.1 Learning Objectives

> **Learning Objectives**
>
> - Classify business data according to strategic typology
> - Assess data assets for analytical potential and business value
> - Design data collection strategies aligned with analytical objectives
> - Apply data quality frameworks to ensure analytical reliability

## 2.2 The Strategic Data Asset Framework

Business data constitutes a strategic asset category requiring systematic classification and management. The Strategic Data Asset Framework organizes data according to four dimensions: origin (internal vs. external), temporality (historical vs. real-time), structure (structured vs. unstructured), and strategic purpose (operational vs. analytical). This multidimensional classification enables organizations to develop comprehensive data strategies that align collection, storage, and analysis activities with business objectives.

### 2.2.1 Customer Data as Strategic Assets

Customer data represents the foundation of modern business analytics. This category encompasses demographic information (enabling market segmentation), behavioral data (supporting personalization initiatives), engagement metrics (informing retention strategies), and value indicators (driving resource allocation decisions). Each subcategory serves distinct analytical purposes while contributing to integrated customer understanding.

**Table 2.1: Customer Data Categories and Strategic Applications**

| Category | Examples | Strategic Application |
|---|---|---|
| Demographic | Age, income, location, industry | Market segmentation, targeting |
| Behavioral | Purchase frequency, browsing patterns | Personalization, recommendations |
| Engagement | Email opens, app usage, support tickets | Retention risk assessment |
| Value | Revenue, profitability, lifetime value | Resource allocation, prioritization |

## 2.3 Data Quality Management for Business Analytics

The reliability of analytical conclusions depends fundamentally on data quality. The data quality management framework encompasses six dimensions: completeness (absence of missing values), accuracy (correctness of recorded values), consistency (agreement across sources), timeliness (currency of information), validity (conformance to business rules), and uniqueness (absence of duplicates). Each dimension requires specific assessment protocols and remediation strategies.

```python
import pandas as pd
import numpy as np

def assess_data_quality(df, critical_columns):
    """
    Comprehensive data quality assessment for business analytics.

    Parameters
    ----------
    df : pd.DataFrame
        Dataset to assess
    critical_columns : list
        Columns essential for analysis

    Returns
    -------
    dict : Quality assessment results with remediation recommendations
    """
    quality_report = {
        'completeness': {},
        'validity': {},
        'business_impact': {},
        'overall_score': 0.0
    }

    # Completeness Assessment
    for col in critical_columns:
        missing_rate = df[col].isnull().mean()
```

Chenhao Zhou

```python
        quality_report['completeness'][col] = {
            'missing_rate': missing_rate,
            'status': 'ACCEPTABLE' if missing_rate < 0.05 else
                      'MARGINAL' if missing_rate < 0.15 else 'CRITICAL',
            'records_affected': int(df[col].isnull().sum())
        }

    # Validity Assessment (example: numeric ranges)
    for col in df.select_dtypes(include=[np.number]).columns:
        if col in critical_columns:
            negative_values = (df[col] < 0).sum()
            quality_report['validity'][col] = {
                'invalid_records': negative_values,
                'status': 'VALID' if negative_values == 0
                          else 'REVIEW_REQUIRED'
            }

    # Calculate Overall Quality Score
    completeness_scores = [1 - v['missing_rate']
        for v in quality_report['completeness'].values()]
    quality_report['overall_score'] = np.mean(completeness_scores) * 100

    # Business Impact Assessment
    if quality_report['overall_score'] >= 95:
        quality_report['recommendation'] = \
            'Data quality ACCEPTABLE for production use'
    elif quality_report['overall_score'] >= 80:
        quality_report['recommendation'] = \
            'Data quality MARGINAL - remediation advised'
    else:
        quality_report['recommendation'] = \
            'Data quality CRITICAL - remediation required'

    return quality_report
```

## 2.4 Market Intelligence and Competitive Data

Beyond internal customer data, organizations require external data for market positioning and competitive analysis. Market intelligence encompasses industry benchmarks, competitor metrics, macroeconomic indicators, and regulatory developments. The integration of internal operational data with external market intelligence enables comprehensive strategic analysis that considers both organizational capabilities and market conditions.

Chenhao Zhou

**Key Concepts**

- **Strategic Data Asset Framework:** Four-dimensional classification system
- **Customer Data Categories:** Demographic, behavioral, engagement, and value data
- **Data Quality Dimensions:** Completeness, accuracy, consistency, timeliness, validity, uniqueness
- **Quality Score Thresholds:** >95% acceptable, 80–95% marginal, <80% critical
- **Market Intelligence:** Integration of internal and external data sources

# 3 Exploratory Analysis for Executive Decision-Making

> **Chapter Overview**
>
> This chapter develops skills in exploratory data analysis oriented toward executive-level communication and decision support. Students will learn to create meaningful visualizations, identify patterns with strategic implications, and synthesize analytical findings into actionable business recommendations.

## 3.1 Learning Objectives

> **Learning Objectives**
>
> - Design executive dashboards that communicate key business metrics
> - Apply correlation analysis to identify strategic business drivers
> - Construct visualizations appropriate for senior leadership audiences
> - Translate statistical findings into business recommendations

## 3.2 Executive Dashboard Design Principles

Effective executive dashboards adhere to design principles that balance informational density with cognitive accessibility. The Gestalt principles of visual perception—proximity, similarity, continuity, closure, and figure-ground—guide dashboard layout decisions. Additionally, dashboards must respect the limited attention bandwidth of senior executives by prioritizing the most critical metrics and providing progressive disclosure mechanisms for detailed exploration.

### 3.2.1 Dashboard Metric Categories

Executive dashboards typically organize metrics into three categories: financial performance (revenue, margins, profitability), operational efficiency (productivity, utilization, cycle times), and strategic indicators (market position, customer satisfaction, innovation pipeline). Each category requires distinct visualization approaches optimized for the nature of the underlying data and the decisions it supports.

```python
def create_executive_dashboard(df, reporting_period='Q3 2026'):
    """
    Generate executive-level business intelligence dashboard.

    Parameters
    ----------
    df : pd.DataFrame
        Operational data containing business metrics
```

```
    reporting_period : str
        Period label for dashboard header

    Returns
    -------
    pd.Series : Dashboard metrics formatted for executive presentation
    """
    dashboard = {}

    # FINANCIAL PERFORMANCE METRICS
    dashboard['Total Revenue'] = df['revenue'].sum()
    dashboard['Revenue Growth (%)'] = (
        (df['revenue'].iloc[-1] - df['revenue'].iloc[0]) /
        df['revenue'].iloc[0] * 100
    )
    dashboard['Gross Margin (%)'] = (
        (df['revenue'].sum() - df['costs'].sum()) /
        df['revenue'].sum() * 100
    )

    # CUSTOMER PERFORMANCE METRICS
    dashboard['Customer Acquisition'] = df['new_customers'].sum()
    dashboard['Churn Rate (%)'] = (
        df['churned_customers'].sum() /
        df['total_customers'].mean() * 100
    )
    dashboard['Net Promoter Score'] = df['nps_score'].mean()

    # OPERATIONAL EFFICIENCY METRICS
    dashboard['CAC Payback (months)'] = (
        df['customer_acquisition_cost'].mean() /
        df['monthly_revenue_per_customer'].mean()
    )

    # STRATEGIC POSITION METRICS
    dashboard['Market Share Change (pp)'] = (
        df['market_share'].iloc[-1] - df['market_share'].iloc[0]
    ) * 100

    return pd.Series(dashboard)

# Generate dashboard
dashboard_metrics = create_executive_dashboard(business_data)
print(f"EXECUTIVE DASHBOARD - {reporting_period}")
print("=" * 60)
for metric, value in dashboard_metrics.items():
    if '%' in metric or 'pp' in metric:
        print(f"{metric:35s}: {value:>10.1f}")
```

```
    elif 'Revenue' in metric and '$' not in metric:
        print(f"{metric:35s}: ${value:>14,.0f}")
    else:
        print(f"{metric:35s}: {value:>10.1f}")
```

## 3.3 Correlation Analysis for Strategic Insight

Correlation analysis identifies relationships between business metrics, enabling strategic prioritization of improvement initiatives. While correlation does not establish causation, strong correlations suggest potential leverage points for value creation. The interpretation of correlations requires business context; statistically significant correlations may lack practical significance, while modest correlations may indicate strategically important relationships.

### Table 3.1: Correlation Interpretation Guidelines

| $|r|$ Range | Interpretation | Business Implication |
|---|---|---|
| $0.00 - 0.19$ | Negligible | No actionable relationship |
| $0.20 - 0.39$ | Weak | Secondary consideration |
| $0.40 - 0.59$ | Moderate | Potential leverage point |
| $0.60 - 0.79$ | Strong | Key business driver |
| $0.80 - 1.00$ | Very Strong | Critical relationship, investigate causality |

**Key Concepts**

- **Dashboard Design:** Application of Gestalt principles to executive visualization
- **Metric Categories:** Financial, operational, and strategic indicator groupings
- **Progressive Disclosure:** Layered information architecture for varied detail needs
- **Correlation Analysis:** Statistical identification of metric relationships
- **Business Context:** Interpretation of statistical findings through strategic lens

# 4 Feature Engineering for Business Intelligence

> **Chapter Overview**
>
> This chapter addresses the transformation of raw operational data into analytically meaningful features. Students will learn systematic approaches to feature creation, understand the importance of domain knowledge in engineering decisions, and develop proficiency in creating features that capture business value dynamics.

## 4.1 Theoretical Foundation: The Feature Engineering Process

Feature engineering constitutes the process of transforming raw data into representations that enhance model performance and interpretability. In business contexts, effective feature engineering requires integration of domain expertise with statistical technique. Features must not only improve predictive accuracy but also maintain interpretability for stakeholder communication and regulatory compliance.

### 4.1.1 Customer Lifetime Value (CLV) Feature Engineering

Customer Lifetime Value represents a canonical example of derived business features. CLV integrates transaction history, customer tenure, and profitability metrics into a single value indicator that enables strategic prioritization. The engineering of CLV features requires decisions regarding discount rates, projection horizons, and treatment of customer segments with insufficient history.

```python
import pandas as pd
import numpy as np

def engineer_customer_value_features(transactions_df):
    """
    Create business value features from transaction data.

    Engineering Approach:
    - RFM (Recency, Frequency, Monetary) as foundation
    - Derived ratios for efficiency metrics
    - Composite scores for prioritization

    Parameters
    ----------
    transactions_df : pd.DataFrame
        Transaction-level data with customer_id, date, amount, margin

    Returns
    -------
```

```python
    pd.DataFrame : Customer-level features with business value
        indicators
    """
    # Aggregate to customer level
    customer_features = transactions_df.groupby('customer_id').agg({
        'revenue': ['sum', 'mean', 'count'],
        'profit_margin': 'mean',
        'transaction_date': ['min', 'max'],
        'acquisition_channel': 'first'
    }).reset_index()

    # Flatten column names
    customer_features.columns = [
        'customer_id', 'total_revenue', 'avg_transaction',
        'purchase_count', 'avg_margin', 'first_purchase',
        'last_purchase', 'acquisition_channel'
    ]

    # DERIVED VALUE FEATURES
    # Customer Lifetime Value (simplified)
    customer_features['customer_lifetime_value'] = (
        customer_features['total_revenue'] *
        customer_features['avg_margin']
    )

    # Tenure in days
    customer_features['tenure_days'] = (
        customer_features['last_purchase'] -
        customer_features['first_purchase']
    ).dt.days

    # Purchase velocity (monthly rate)
    customer_features['purchase_velocity'] = (
        customer_features['purchase_count'] /
        (customer_features['tenure_days'] + 1) * 30
    )

    # Revenue per day of relationship
    customer_features['daily_revenue_rate'] = (
        customer_features['total_revenue'] /
        (customer_features['tenure_days'] + 1)
    )

    # COMPOSITE SCORING
    # Normalize components to 0-1 scale
    for col in ['customer_lifetime_value', 'purchase_velocity',
                'avg_transaction']:
        customer_features[f'{col}_norm'] = (
```

       Chenhao Zhou

```
        (customer_features[col] - customer_features[col].min()) /
        (customer_features[col].max() - customer_features[col].min()
            )
    )

# Weighted business priority score
customer_features['priority_score'] = (
    customer_features['customer_lifetime_value_norm'] * 0.5 +
    customer_features['purchase_velocity_norm'] * 0.3 +
    customer_features['avg_transaction_norm'] * 0.2
) * 100  # Scale to 0-100

# Segment assignment
customer_features['segment'] = pd.qcut(
    customer_features['priority_score'],
    q=[0, 0.33, 0.67, 1.0],
    labels=['Develop', 'Maintain', 'VIP']
)

return customer_features
```

### Key Concepts

- **Feature Engineering:** Transformation of raw data into analytical representations
- **RFM Framework:** Recency, Frequency, Monetary value as foundational metrics
- **Customer Lifetime Value:** Integrated value indicator for strategic prioritization
- **Composite Scoring:** Weighted combination of normalized features
- **Segment Assignment:** Categorization based on derived value indicators

# 5  Statistical Inference for Business Decisions

---

> **Chapter Overview**
>
> This chapter develops rigorous statistical inference skills essential for evidence-based business decisions. Students will learn to design and analyze controlled experiments, apply hypothesis testing frameworks, and interpret statistical results in business contexts with appropriate consideration of practical significance.

## 5.1  A/B Testing Methodology for Business Experimentation

A/B testing (randomized controlled trials) constitutes the gold standard for establishing causal relationships between interventions and business outcomes. The methodology requires careful attention to experimental design, sample size determination, randomization procedures, and statistical analysis. Beyond statistical significance, business experimenters must consider practical significance—whether observed effects justify implementation costs and organizational change.

### 5.1.1  Experimental Design Framework

Rigorous experimental design addresses four fundamental elements: (1) hypothesis specification (null and alternative), (2) sample size calculation (power analysis), (3) randomization procedure (ensuring treatment and control comparability), and (4) analysis plan (pre-specified statistical tests and decision criteria). Pre-registration of experimental designs protects against post-hoc hypothesis generation and multiple testing problems.

```python
import numpy as np
from scipy import stats
from statsmodels.stats.power import TTestIndPower

class ABTestAnalyzer:
    """
    Comprehensive A/B test analysis with business impact calculation.
    """

    def __init__(self, control_data, treatment_data,
                 metric_name='conversion'):
        self.control = np.array(control_data)
        self.treatment = np.array(treatment_data)
        self.metric_name = metric_name

    def calculate_statistical_significance(self, alpha=0.05):
        """
        Perform hypothesis test with effect size estimation.
```

```python
    """
    # Two-sample t-test
    t_statistic, p_value = stats.ttest_ind(
        self.treatment, self.control
    )

    # Effect size (Cohen's d)
    pooled_std = np.sqrt(
        (self.control.std()**2 + self.treatment.std()**2) / 2
    )
    cohens_d = (self.treatment.mean() -
                self.control.mean()) / pooled_std

    # Confidence interval for difference
    diff_mean = self.treatment.mean() - self.control.mean()
    se_diff = np.sqrt(
        self.control.var()/len(self.control) +
        self.treatment.var()/len(self.treatment)
    )
    ci_95 = (
        diff_mean - 1.96 * se_diff,
        diff_mean + 1.96 * se_diff
    )

    return {
        't_statistic': t_statistic,
        'p_value': p_value,
        'significant': p_value < alpha,
        'effect_size': cohens_d,
        'effect_interpretation':
            self._interpret_effect_size(cohens_d),
        'confidence_interval': ci_95
    }

def calculate_business_impact(self, annual_revenue,
                              implementation_cost):
    """
    Translate statistical results to business value.
    """
    lift = ((self.treatment.mean() - self.control.mean()) /
            self.control.mean())
    projected_revenue_increase = annual_revenue * lift
    roi = ((projected_revenue_increase - implementation_cost) /
            implementation_cost)

    return {
        'relative_lift': lift,
        'revenue_impact': projected_revenue_increase,
```

```python
            'implementation_cost': implementation_cost,
            'roi': roi,
            'recommendation': 'IMPLEMENT' if roi > 0 else 'REJECT'
        }

    @staticmethod
    def _interpret_effect_size(d):
        if abs(d) < 0.2:
            return 'Negligible'
        elif abs(d) < 0.5:
            return 'Small'
        elif abs(d) < 0.8:
            return 'Medium'
        else:
            return 'Large'

    @staticmethod
    def calculate_required_sample_size(effect_size, power=0.8,
                                        alpha=0.05):
        """
        Power analysis for experiment planning.
        """
        analysis = TTestIndPower()
        n = analysis.solve_power(
            effect_size=effect_size,
            power=power,
            alpha=alpha,
            alternative='two-sided'
        )
        return int(np.ceil(n))
```

**Table 5.1: Effect Size Interpretation Guidelines**

| Cohen's $d$ | Interpretation | Business Significance |
|---|---|---|
| $< 0.2$ | Negligible | Unlikely to justify implementation |
| $0.2 - 0.5$ | Small | May justify low-cost interventions |
| $0.5 - 0.8$ | Medium | Generally actionable |
| $> 0.8$ | Large | Strong case for implementation |

Chenhao Zhou

> **Key Concepts**
>
> - **A/B Testing:** Randomized controlled trials for causal inference
> - **Statistical vs. Practical Significance:** Effect size considerations
> - **Power Analysis:** Sample size determination for reliable detection
> - **Effect Size (Cohen's $d$):** Standardized measure of intervention impact
> - **Business Impact Translation:** Conversion of statistical to financial metrics

　　　　　　　　　　　　　　　Chenhao Zhou

PART II

# Machine Learning Implementation

*Practical Applications with Business Focus*

# 6 Predictive Modeling for Revenue Optimization

> **Chapter Overview**
>
> This chapter introduces supervised learning techniques for regression problems, with emphasis on revenue forecasting and valuation models. Students will learn to build, evaluate, and deploy predictive models that generate measurable business value through improved forecasting accuracy.

## 6.1 Learning Objectives

> **Learning Objectives**
>
> - Implement linear regression models for business forecasting applications
> - Apply ensemble methods (Random Forest, Gradient Boosting) for improved accuracy
> - Evaluate models using both technical and business-oriented performance metrics
> - Interpret feature importance for actionable business insights

## 6.2 Theoretical Foundation: Regression Analysis

Regression analysis models the relationship between a continuous dependent variable and one or more independent variables. In business contexts, regression enables forecasting of revenue, costs, demand, and other continuous metrics essential for planning and resource allocation. The choice among regression techniques involves trade-offs between interpretability, accuracy, and computational requirements.

### 6.2.1 Linear Regression: The Foundation

Linear regression assumes a linear relationship between features and target, expressed as $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_n X_n + \varepsilon$. The model minimizes the sum of squared errors to estimate coefficients ($\beta$), which directly indicate the marginal effect of each feature on the target variable. This interpretability makes linear regression valuable for stakeholder communication and regulatory contexts where model transparency is required.

```python
from sklearn.linear_model import LinearRegression
from sklearn.ensemble import (RandomForestRegressor,
                              GradientBoostingRegressor)
from sklearn.model_selection import (train_test_split,
                                     cross_val_score)
from sklearn.metrics import (mean_absolute_error,
                             mean_squared_error, r2_score)
from sklearn.preprocessing import StandardScaler
```

Chenhao Zhou

```python
import pandas as pd
import numpy as np

class BusinessRegressionPipeline:
    """
    Production-ready regression pipeline with business metrics.
    """

    def __init__(self, model_type='linear'):
        self.model_type = model_type
        self.scaler = StandardScaler()
        self.model = self._initialize_model()
        self.feature_names = None

    def _initialize_model(self):
        models = {
            'linear': LinearRegression(),
            'random_forest': RandomForestRegressor(
                n_estimators=100, max_depth=10, random_state=42
            ),
            'gradient_boosting': GradientBoostingRegressor(
                n_estimators=100, max_depth=5, random_state=42
            )
        }
        return models.get(self.model_type, LinearRegression())

    def fit(self, X, y, feature_names=None):
        """Train model with scaling."""
        self.feature_names = (feature_names if feature_names
                              else X.columns.tolist())
        X_scaled = self.scaler.fit_transform(X)
        self.model.fit(X_scaled, y)
        return self

    def predict(self, X):
        """Generate predictions."""
        X_scaled = self.scaler.transform(X)
        return self.model.predict(X_scaled)

    def evaluate(self, X_test, y_test,
                 annual_baseline_revenue=None):
        """
        Comprehensive evaluation with business metrics.
        """
        y_pred = self.predict(X_test)

        # Technical Metrics
        metrics = {
```

       Chenhao Zhou

```
            'MAE': mean_absolute_error(y_test, y_pred),
            'RMSE': np.sqrt(mean_squared_error(y_test, y_pred)),
            'R2': r2_score(y_test, y_pred),
            'MAPE': np.mean(
                np.abs((y_test - y_pred) / y_test)) * 100
    }

    # Business Accuracy Metrics
    within_5_pct = np.mean(
        np.abs(y_test - y_pred) / y_test <= 0.05)
    within_10_pct = np.mean(
        np.abs(y_test - y_pred) / y_test <= 0.10)

    metrics['Accuracy_5pct'] = within_5_pct * 100
    metrics['Accuracy_10pct'] = within_10_pct * 100

    # Business Value of Improved Forecasting
    if annual_baseline_revenue:
        forecast_improvement_value = (
            annual_baseline_revenue * 0.01 *
            (metrics['Accuracy_10pct'] - 70)
        )
        metrics['Forecast_Value'] = max(
            0, forecast_improvement_value)

    return metrics

def get_feature_importance(self):
    """Extract and rank feature importance."""
    if hasattr(self.model, 'feature_importances_'):
        importance = self.model.feature_importances_
    elif hasattr(self.model, 'coef_'):
        importance = np.abs(self.model.coef_)
    else:
        return None

    return pd.DataFrame({
        'feature': self.feature_names,
        'importance': importance
    }).sort_values('importance', ascending=False)
```

**Table 7.1: Model Selection Guidelines**

| Model | Interpretability | Accuracy | Use Case |
|---|---|---|---|
| Linear Regression | High | Baseline | Regulatory contexts, baseline |
| Random Forest | Medium | High | General-purpose prediction |
| Gradient Boosting | Low | Highest | Accuracy-critical applications |
| Ridge/Lasso | High | Good | High-dimensional, regularization |

**Key Concepts**

- **Linear Regression:** Interpretable baseline for continuous prediction
- **Ensemble Methods:** Random Forest and Gradient Boosting for accuracy
- **MAPE:** Mean Absolute Percentage Error for business-relevant accuracy
- **Feature Importance:** Identification of key predictive drivers
- **Business Value Translation:** Forecasting accuracy to financial impact

    *Chenhao Zhou*

# 7 Classification Models for Customer Analytics

---

> **Chapter Overview**
>
> This chapter develops classification modeling skills for customer-focused applications including churn prediction, segmentation, and risk assessment. Students will learn to build classification pipelines, select appropriate algorithms based on business requirements, and translate classification outcomes into actionable customer strategies.

## 7.1 The Classification Problem in Business Context

Classification models assign observations to discrete categories based on input features. Business applications include customer churn prediction (will/will not churn), credit risk assessment (approve/decline), lead qualification (high/medium/low quality), and fraud detection (legitimate/fraudulent). The asymmetric costs of classification errors require careful threshold selection that reflects business priorities rather than purely technical optimization.

### 7.1.1 Asymmetric Error Costs

In business classification problems, false positives and false negatives typically carry different costs. Consider churn prediction: a false positive (predicting churn for a loyal customer) results in unnecessary retention spending, while a false negative (missing an actual churner) results in lost customer lifetime value. Optimal classification thresholds balance these asymmetric costs according to their relative magnitudes.

```python
from sklearn.ensemble import RandomForestClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import (classification_report,
    confusion_matrix, roc_auc_score, precision_recall_curve)
import numpy as np
import pandas as pd

class ChurnPreventionModel:
    """
    Business-focused churn prediction with ROI optimization.
    """

    def __init__(self,
                 customer_lifetime_value=2000,
                 retention_cost=100,
                 retention_success_rate=0.4):
        self.clv = customer_lifetime_value
        self.retention_cost = retention_cost
```

```python
        self.success_rate = retention_success_rate
        self.model = RandomForestClassifier(
            n_estimators=100,
            max_depth=10,
            class_weight='balanced',
            random_state=42
        )

    def fit(self, X, y):
        """Train the churn prediction model."""
        self.model.fit(X, y)
        return self

    def predict_proba(self, X):
        """Get churn probabilities for threshold optimization."""
        return self.model.predict_proba(X)[:, 1]

    def calculate_optimal_threshold(self, X_val, y_val):
        """
        Find threshold that maximizes expected business value.
        """
        probas = self.predict_proba(X_val)
        thresholds = np.arange(0.1, 0.9, 0.05)
        best_value = -np.inf
        best_threshold = 0.5

        for threshold in thresholds:
            predictions = (probas >= threshold).astype(int)
            value = self._calculate_expected_value(
                y_val, predictions)
            if value > best_value:
                best_value = value
                best_threshold = threshold

        return best_threshold, best_value

    def _calculate_expected_value(self, y_true, y_pred):
        """
        Calculate expected business value of classification.

        Value = (True Positives * CLV * Success Rate)
            - (Predicted Positives * Retention Cost)
            - (False Negatives * CLV)
        """
        tn, fp, fn, tp = confusion_matrix(
            y_true, y_pred).ravel()

        saved_value = tp * self.clv * self.success_rate
```

```python
        retention_cost = (tp + fp) * self.retention_cost
        missed_churner_cost = fn * self.clv

        return saved_value - retention_cost - missed_churner_cost

    def generate_business_report(self, X_test, y_test,
                                   threshold=0.5):
        """
        Comprehensive business impact analysis.
        """
        probas = self.predict_proba(X_test)
        predictions = (probas >= threshold).astype(int)
        tn, fp, fn, tp = confusion_matrix(
            y_test, predictions).ravel()

        report = {
            'model_performance': {
                'auc_roc': roc_auc_score(y_test, probas),
                'precision': (tp / (tp + fp)
                              if (tp + fp) > 0 else 0),
                'recall': (tp / (tp + fn)
                           if (tp + fn) > 0 else 0),
                'threshold_used': threshold
            },
            'business_impact': {
                'customers_targeted': tp + fp,
                'true_churners_identified': tp,
                'expected_saves': int(tp * self.success_rate),
                'retention_investment':
                    (tp + fp) * self.retention_cost,
                'revenue_preserved':
                    tp * self.success_rate * self.clv,
                'net_value':
                    self._calculate_expected_value(
                        y_test, predictions)
            },
            'risk_assessment': {
                'missed_churners': fn,
                'missed_revenue': fn * self.clv,
                'false_alarm_rate': (fp / (fp + tn)
                    if (fp + tn) > 0 else 0)
            }
        }
        return report
```

**Table 8.1: Classification Metrics for Business Applications**

| Metric | Definition | When to Prioritize |
|--------|-----------|--------------------|
| Precision | TP / (TP + FP) | High cost of false positives |
| Recall | TP / (TP + FN) | High cost of missing positives |
| F1 Score | Harmonic mean of P and R | Balance between P and R |
| AUC-ROC | Area under ROC curve | Overall model comparison |

**Key Concepts**

- **Classification:** Assignment of observations to discrete business categories
- **Asymmetric Costs:** Differential impact of false positives vs. false negatives
- **Threshold Optimization:** Selection based on business value rather than accuracy
- **AUC-ROC:** Overall discriminative ability across all thresholds
- **Business Impact Translation:** Confusion matrix to financial outcomes

# References

Davenport, T. H., & Harris, J. G. (2007). *Competing on Analytics: The New Science of Winning.* Harvard Business School Press.

Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction* (2nd ed.). Springer.

James, G., Witten, D., Hastie, T., & Tibshirani, R. (2021). *An Introduction to Statistical Learning with Applications in Python.* Springer.

Kohavi, R., Tang, D., & Xu, Y. (2020). *Trustworthy Online Controlled Experiments: A Practical Guide to A/B Testing.* Cambridge University Press.

McKinney, W. (2022). *Python for Data Analysis: Data Wrangling with Pandas, NumPy, and Jupyter* (3rd ed.). O'Reilly Media.

Provost, F., & Fawcett, T. (2013). *Data Science for Business: What You Need to Know about Data Mining and Data-Analytic Thinking.* O'Reilly Media.

Shmueli, G., Bruce, P. C., Yahav, I., Patel, N. R., & Lichtendahl, K. C. (2024). *Data Mining for Business Analytics: Concepts, Techniques, and Applications in Python* (2nd ed.). Wiley.