

# Analytical Base Table (ABT) Quick Reference Guide

## What is an ABT?

An Analytical Base Table is a denormalized, flat data structure that consolidates all relevant features and target variables into a single table, optimized for machine learning and predictive analytics workflows. Each row represents a unique observation (entity/event), and each column represents a feature or target variable.

## Core Components

### 1. Entity/Grain Definition

- **Primary Key:** Unique identifier for each row (customer\_id, transaction\_id, etc.)
- **Grain Level:** The unit of analysis (customer-level, transaction-level, session-level)
- **Time Component:** Point-in-time snapshot or observation window

### 2. Target Variable

- **Binary Classification:** 0/1 outcomes (churn, fraud, conversion)
- **Multi-class:** Multiple categorical outcomes
- **Regression:** Continuous values (revenue, LTV, duration)
- **Time-to-Event:** Survival analysis targets

### 3. Feature Categories

- **Demographic:** Age, gender, location, occupation
- **Behavioral:** Usage patterns, frequency metrics, engagement scores
- **Transactional:** Purchase history, monetary values, recency metrics
- **Temporal:** Time-based features, seasonality indicators, trends
- **Derived:** Ratios, aggregations, transformations, interaction terms
- **External:** Third-party data, economic indicators, weather data

## Design Principles

### Temporal Consistency

- **No Data Leakage:** Features must only use data available before the prediction point
- **Observation Window:** Historical period for feature calculation

- **Performance Window:** Future period for target variable measurement
- **Gap Period:** Buffer between observation and performance windows (prevents leakage)

## Feature Engineering Best Practices

### Aggregation Strategies

- Count, sum, average, min, max, standard deviation
- Recency, frequency, monetary (RFM) metrics
- Rolling windows (7-day, 30-day, 90-day aggregates)
- Lag features (previous period values)
- Trend indicators (slope, acceleration)

### Handling Different Data Types

- **Numerical:** Scaling, normalization, binning, polynomial features
- **Categorical:** One-hot encoding, target encoding, embedding
- **Temporal:** Cyclical encoding (sin/cos for day of week, month)
- **Text:** TF-IDF, word embeddings, sentiment scores
- **Missing Values:** Imputation strategies, missingness indicators

## Common ABT Patterns by Use Case

### Customer Churn Prediction

customer\_id | tenure | avg\_monthly\_spend | days\_since\_last\_purchase | support\_tickets | ... | churned

- Observation: 12 months of history
- Gap: 1 month
- Target: Churn in next 3 months

### Fraud Detection

transaction\_id | amount | merchant\_category | time\_since\_last\_txn | unusual\_location | ... | is\_fraud

- Observation: Real-time + historical patterns
- Gap: None (real-time scoring)
- Target: Fraudulent transaction

## Lead Scoring

lead\_id | source | industry | company\_size | engagement\_score | days\_in\_funnel | ... | converted

- Observation: All interactions until scoring date
- Gap: None
- Target: Conversion within 30 days

## Data Quality Considerations

### Validation Checks

- **Completeness:** Missing value percentages per feature
- **Consistency:** Data type validation, range checks
- **Uniqueness:** Primary key duplication check
- **Timeliness:** Data freshness and update frequency
- **Accuracy:** Outlier detection, distribution analysis

### Common Pitfalls to Avoid

- **Target Leakage:** Using future information in features
- **Survivorship Bias:** Only including successful outcomes
- **Sample Selection Bias:** Non-representative training data
- **Class Imbalance:** Skewed target distribution
- **Multicollinearity:** Highly correlated features
- **Overfitting:** Too many features relative to observations

## Implementation Workflow

### 1. Data Collection Phase

- Identify all relevant data sources
- Define join keys and relationships
- Establish data refresh schedules
- Document data lineage

### 2. Feature Development Phase

- Create feature specifications

- Build transformation pipelines
- Test feature calculations
- Version control feature definitions

### **3. ABT Assembly Phase**

- Execute joins and aggregations
- Apply temporal filters
- Generate train/validation/test splits
- Create point-in-time snapshots

### **4. Quality Assurance Phase**

- Run data quality checks
- Validate temporal consistency
- Check for data leakage
- Profile feature distributions

## **Advanced Techniques**

### **Feature Store Integration**

- Centralized feature repository
- Feature versioning and lineage
- Real-time feature serving
- Feature sharing across models

### **Incremental Updates**

- Delta processing for new data
- Sliding window updates
- Feature backfilling strategies
- Change data capture (CDC)

### **Scalability Considerations**

- Partitioning strategies (by date, entity)
- Columnar storage formats (Parquet, ORC)
- Distributed processing (Spark, Dask)

- Sampling techniques for large datasets

## **Performance Optimization**

### **Computational Efficiency**

- Pre-compute expensive aggregations
- Use materialized views
- Implement caching strategies
- Optimize join operations

### **Storage Optimization**

- Column selection (remove low-value features)
- Data type optimization
- Compression techniques
- Archival strategies for old data

## **Monitoring and Maintenance**

### **Feature Drift Detection**

- Distribution shift monitoring
- Statistical tests (KS, PSI, Wasserstein)
- Feature importance tracking
- Model performance degradation alerts

### **Documentation Requirements**

- Feature definitions and business logic
- Data source dependencies
- Calculation examples
- Update frequency and SLAs

## **Tools and Technologies**

### **Data Processing**

- **SQL:** Complex aggregations and joins
- **Python:** Pandas, PySpark, Dask

- **R:** data.table, dplyr
- **Spark:** Large-scale distributed processing

## Feature Stores

- Feast, Tecton, Hopsworks, AWS SageMaker Feature Store

## Orchestration

- Airflow, Prefect, Dagster, Luigi

## Version Control

- DVC (Data Version Control), Git LFS, Delta Lake

## Key Metrics for ABT Quality

- **Coverage:** % of population with complete features
- **Recency:** Average age of data in features
- **Cardinality:** Unique values per categorical feature
- **Correlation:** Feature correlation matrix analysis
- **Information Value:** Predictive power of features
- **Variance Inflation Factor:** Multicollinearity detection

## Best Practices Checklist

- ☐ Define clear entity and grain
- ☐ Establish observation and performance windows
- ☐ Implement data leakage prevention
- ☐ Document all feature transformations
- ☐ Create reproducible pipelines
- ☐ Version control ABT schemas
- ☐ Implement automated quality checks
- ☐ Monitor feature drift
- ☐ Maintain feature documentation
- ☐ Plan for scalability
- ☐ Set up incremental update processes
- ☐ Establish data governance policies