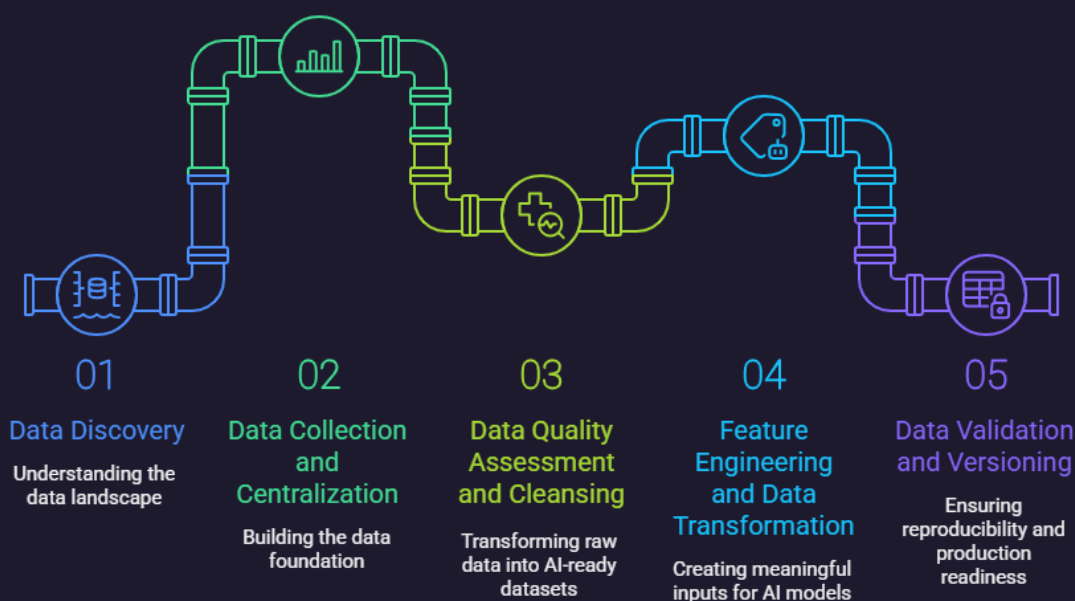# 5 Steps to Prepare Your Manufacturing Data for AI

## Introduction

Artificial Intelligence promises to revolutionize manufacturing operations, but success hinges on one critical factor: data quality. In industrial environments, data comes from diverse sourcesy, SCADA (Supervisory Control and Data Acquisition) systems, PLCs (Programmable Logic Controller), MES (Manufacturing Execution System) platforms, ERP (Enterprise Resource Planning) systems, and IoT (internet-of-things) sensors, each with unique formats, frequencies, and quality levels. Before any AI model can deliver value, this data must be properly prepared, cleansed, and structured.

This guide walks you through five essential steps to prepare your manufacturing data for AI implementation, drawing on proven practices from successful industrial AI deployments.



**01**

**Data Discovery**

Understanding the data landscape

**02**

**Data Collection and Centralization**

Building the data foundation

**03**

**Data Quality Assessment and Cleansing**

Transforming raw data into AI-ready datasets

**04**

**Feature Engineering and Data Transformation**

Creating meaningful inputs for AI models

**05**

**Data Validation and Versioning**

Ensuring reproducibility and production readiness

# Step 1: Data Discovery and Inventory Assessment

## Understanding Your Data Landscape

Before cleaning data, you must know what you have. Manufacturing environments typically contain dozens of data sources that have evolved over years or decades.

**Key Activities:**

- **Map all data sources**: Document every system generating data—sensors, machines, quality systems, maintenance logs, production databases, and manual entry points.
- **Identify data types**: Classify data as time-series (sensor readings), transactional (production orders), relational (product specifications), or unstructured (operator notes, images).
- **Assess data volume and velocity**: Determine how much data each source generates and at what frequency (milliseconds to daily).
- **Evaluate current accessibility**: Identify which systems are already integrated and which exist in silos.

**Real-World Example:**

A automotive parts manufacturer discovered they had 47 different data sources across three plants. Only 12 were integrated into their data lake. Critical quality inspection data resided in Excel spreadsheets on individual quality engineer laptops. This assessment revealed that 35% of potentially valuable AI training data was inaccessible.

**Best Practices:**

- Use automated discovery tools to scan networks for databases and file shares
- Interview process operators and engineers to uncover manual data collection processes
- Create a data catalog with metadata describing each source's purpose, owner, and update frequency
- Document data lineage to understand how information flows between systems

| Data Source | Format | Update Frequency | Integrated? | Owner |
|---|---|---|---|---|
| SCADA Systems | Time-series | seconds/minutes | Yes | Operations |
| PLCs | Time-series | ms-seconds | Partial | Engineering |
| MES | Transactional | batch/continuous | Yes | Production |
| ERP | Transactional | batch/daily | Partial | IT |
| IoT Sensors | Time-series | ms-seconds | No | Maintenance |
| Quality Spreadsheets | Unstructured | variable | No | Quality Eng. |

Note:

| System | What it Does |
|---|---|
| SCADA | Supervisory Control and Data Acquisition—monitors and controls industrial processes (e.g., plant operations, real-time sensors) |
| PLC | Programmable Logic Controller—controls machinery and processes at the equipment level (automation, logic, rapid response) |
| MES | Manufacturing Execution System—manages and tracks production on the shop floor, bridging the gap between machines and business systems |
| ERP | Enterprise Resource Planning—integrates core business processes (finance, HR, supply chain, inventory, orders, scheduling) |

# Step 2: Data Collection and Centralization

## Building Your Data Foundation

AI models require centralized access to data. Manufacturing data must be extracted from operational technology (OT) systems and consolidated without disrupting production.

**Key Activities:**

- **Establish secure data pipelines**: Implement extraction mechanisms that don't impact real-time control systems. Use read-only database replicas or historian APIs.
- **Create a unified data platform**: Deploy a data lake or lakehouse architecture that can handle diverse data types and volumes.
- **Implement time synchronization**: Ensure all data sources use consistent timestamps (critical for correlating events across systems).
- **Set up continuous ingestion**: Move from batch to streaming ingestion where timing matters for AI predictions.

**Real-World Example:**

A food processing plant needed to predict equipment failures. Their challenge: PLC data logged every second, vibration sensors logged every 100ms, and maintenance records were entered weekly. They implemented a time-series database with edge computing to downsample high-frequency sensor data to relevant features (peak vibration, RMS values) before transmission, reducing data volume by 85% while preserving predictive signals.

**Architecture Considerations:**

- **Edge processing**: Pre-process data at the source to reduce network load and extract relevant features
- **Data historians**: Leverage existing industrial historians (OSIsoft PI, Aveva Historian) as intermediaries
- **Cloud vs. on-premise**: Balance data gravity, latency requirements, and security concerns
- **Data governance**: Implement access controls and audit trails from the start

**Best Practices:**

- Never extract data directly from production control systems—use mirrored databases or historians
- Implement automated monitoring to detect pipeline failures or data delays
- Use industry-standard protocols (OPC UA, MQTT) for interoperability
- Build in data validation at ingestion to catch source system issues early

# Step 3: Data Quality Assessment and Cleansing

## Transforming Raw Data into AI-Ready Datasets

Industrial data is notoriously messy. Sensors fail, operators make typos, systems go offline, and process changes aren't documented. This step is often the most time-consuming but most critical.

**Common Data Quality Issues in Manufacturing:**

1. **Missing values**: Sensor outages, system maintenance, network interruptions
2. **Outliers and anomalies**: Sensor calibration drift, data entry errors, actual process excursions
3. **Inconsistent formats**: Different units, varying precision, mixed data types
4. **Duplicate records**: System redundancy, batch processing errors
5. **Temporal misalignment**: Clock drift, daylight saving time, timezone issues

**Cleansing Strategies:**

**For Time-Series Sensor Data:**

- **Missing values**: Use forward-fill for short gaps (< 5 minutes), interpolation for medium gaps, mark long gaps as explicit missing periods
- **Outlier detection**: Apply statistical methods (z-score, IQR) but validate against process knowledge—some "outliers" are real production events
- **Smoothing**: Apply moving averages or filters to reduce noise, but preserve important transients
- **Resampling**: Align data to consistent time intervals using appropriate aggregation (mean for continuous values, last observation for state changes)

**For Transactional/Event Data:**

- **Deduplication**: Use composite keys (timestamp + machine ID + event type) to identify duplicates
- **Standardization**: Convert all measurements to consistent units, normalize product codes and naming
- **Validation rules**: Check referential integrity (all product IDs exist in master data)
- **Enrichment**: Add contextual information (shift codes, product families, seasonal factors)

**Real-World Example:**

A pharmaceutical manufacturer preparing data for batch quality prediction found that 23% of temperature sensor readings showed physically impossible values (temperatures below absolute zero or above sensor limits). Investigation revealed this happened when sensors were replaced, the new sensor ID wasn't updated in the SCADA configuration. The solution: implement a sensor master data table and automated validation rules that flagged impossible values immediately.

**Best Practices:**

- **Create a data quality scorecard**: Track metrics like completeness, validity, consistency, and timeliness
- **Distinguish between valid zeros and missing data**: A machine running at 0 RPM is different from a missing sensor reading
- **Preserve raw data**: Always keep original data and track all transformations
- **Domain expertise is essential**: Work with process engineers to understand what's "normal" vs. truly anomalous
- **Automate repetitive cleansing**: Build reusable data pipelines that apply standard transformations
- **Version your datasets**: Track which cleansing rules were applied to training vs. test vs. production data

| Issue | Typical Cause | Detection Strategy | Recommended Fix |
|-------|---------------|--------------------|-----------------|
| Missing Values | Outages, user error | Completeness checks | Interpolation, flag gaps |
| Outliers | Sensor drift, entry error | Z-score/IQR, domain | Verify, correct, persist |
| Duplicates | Batch ingestion | Composite keys | Remove, deduplicate |
| Misalignment | Timezone, drift | Time checks | Resample, sync clocks |
| Inconsistency | Format/unit mix | Validation rules | Unify units, normalize |

**Data Cleansing Checklist:**

- [ ] Document all assumptions and decisions made during cleansing
- [ ] Validate cleansed data with subject matter experts
- [ ] Test cleansing logic on multiple time periods to ensure consistency
- [ ] Monitor data quality continuously—don't assume it stays clean
- [ ] Create data quality reports for stakeholders showing before/after metrics

# Step 4: Feature Engineering and Data Transformation

## Creating Meaningful Inputs for AI Models

Raw sensor readings rarely provide optimal inputs for AI models. Feature engineering transforms raw data into representations that expose patterns relevant to your prediction targets.

**Essential Manufacturing Features:**

**Temporal Features:**

- **Rolling statistics**: Moving averages, standard deviations, min/max over time windows
- **Rate of change**: Velocity and acceleration of process parameters
- **Time-based**: Hour of day, day of week, shift, time since last maintenance
- **Lag features**: Previous values to capture temporal dependencies

**Aggregation Features:**

- **Statistical summaries**: Mean, median, mode, quartiles across batches or time periods
- **Cycle metrics**: Duration, peak values, integral (area under curve) for production cycles
- **Event counts**: Number of alarms, operator interventions, recipe changes per shift

**Contextual Features:**

- **Product characteristics**: Grade, family, complexity, material properties
- **Operating modes**: Startup, steady-state, shutdown, changeover
- **Environmental**: Ambient temperature, humidity, raw material lot

**Domain-Specific Features:**

- **Equipment health indicators**: Overall Equipment Effectiveness (OEE), Mean Time Between Failures (MTBF)
- **Process stability**: Control chart statistics (Cp, Cpk), variation coefficients
- **Material tracking**: Lot genealogy, supplier characteristics, age of materials

**Real-World Example:**

A steel manufacturer building a quality prediction model initially used raw furnace temperature readings. Model performance was poor. After consulting metallurgists, the data science team engineered features based on domain knowledge: temperature gradients, heating rate, time spent in specific temperature zones, and cooling rate profiles. These physics-informed features improved prediction accuracy from 67% to 91%.

**Transformation Techniques:**

**Scaling and Normalization:**

- **Standardization**: Zero mean, unit variance—useful for algorithms sensitive to scale (neural networks, SVM)
- **Min-Max scaling**: Scale to [0,1] range—preserves zero values, good for bounded data
- **Robust scaling**: Uses median and IQR—less sensitive to outliers

**Encoding Categorical Variables:**

- **One-hot encoding**: For nominal categories (product type, machine ID)
- **Ordinal encoding**: For ordered categories (quality grade: A, B, C)
- **Target encoding**: Replace category with target variable mean—powerful but risk of overfitting

**Dimensionality Reduction:**

- **PCA**: Reduce correlated sensor readings to principal components
- **Feature selection**: Use domain knowledge or statistical tests to remove irrelevant features
- **Aggregation**: Combine related sensors into composite indicators

**Best Practices:**

- **Start with domain knowledge**: Consult process engineers before algorithmic feature selection
- **Validate feature relevance**: Use correlation analysis, feature importance scores, and SHAP values
- **Avoid data leakage**: Never include information not available at prediction time
- **Document feature definitions**: Create a feature dictionary explaining calculation logic and business meaning
- **Version your features**: Track feature engineering changes as you iterate
- **Test feature stability**: Ensure features remain meaningful across different time periods and production conditions

**Feature Engineering Workflow:**

1. Generate candidate features based on domain knowledge
2. Calculate features across your dataset
3. Analyze feature distributions and correlations
4. Remove redundant or low-variance features
5. Validate with subject matter experts
6. Test features with simple models before complex ones
7. Refine based on model performance and interpretability

# Step 5: Data Validation and Versioning

## Ensuring Reproducibility and Production Readiness

Proper validation and versioning transform a data science experiment into a production-ready AI system.

**Data Splitting Strategy:**

Manufacturing data requires careful splitting due to temporal dependencies and production realities:

- **Time-based splitting**: Train on historical data, validate on intermediate period, test on most recent data
- **Stratification**: Ensure all product types, operating conditions, and failure modes are represented
- **Hold-out production data**: Reserve recent data completely unseen by the model development team

**Avoid these common mistakes:**

- Random splitting of time-series data (causes data leakage)
- Testing on data from the same production period as training
- Ignoring seasonal or campaign-based production patterns

**Real-World Example:**

A packaging line predictive maintenance model showed 95% accuracy during development but failed in production. Root cause: the model was trained on data from steady-state production but deployed during a plant modernization period with new equipment and operators. The lesson: test data must represent deployment conditions, not just historical norms.

**Data Validation Checks:**

**Statistical Validation:**

- **Distribution checks**: Compare train/validation/test distributions to detect drift
- **Temporal consistency**: Verify no future information leaks into training data
- **Balance checks**: Ensure sufficient examples of minority classes (rare failures, quality issues)
- **Correlation stability**: Verify key relationships remain consistent across splits

**Production Readiness Validation:**

- **Latency testing**: Confirm data can be fetched and processed within required timeframes
- **Completeness simulation**: Test model performance when real-time data has missing values

- **Adversarial testing**: Inject realistic data quality issues to test robustness
- **Cross-plant validation**: If deploying across facilities, test on data from different locations

**Data Versioning Best Practices:**

Manufacturing AI projects evolve—equipment changes, processes improve, quality standards tighten. Versioning ensures reproducibility and enables debugging.

**What to Version:**

- **Raw data snapshots**: Periodic snapshots of source data
- **Cleansing configurations**: Scripts, rules, parameters used for data preparation
- **Feature definitions**: Code and documentation for feature engineering
- **Dataset versions**: Final train/validation/test splits with metadata
- **Data quality reports**: Metrics and issues identified in each version

**Version Control Tools:**

- **DVC (Data Version Control)**: Git-like versioning for large datasets
- **MLflow**: Track datasets alongside models and experiments
- **Delta Lake/Apache Iceberg**: ACID transactions and time travel for data lakes
- **Custom solutions**: Combine git for code with cloud object versioning for data

**Metadata to Track:**

- Date range of data
- Data sources included
- Cleansing and transformation steps applied
- Feature engineering logic version
- Data quality metrics
- Number of records, features, and target distribution
- Known issues or limitations

**Real-World Example:**

A chemical plant's corrosion prediction model performance degraded six months after deployment. Because they version-controlled their data pipeline, the team could replay data processing with different historical periods. They discovered a sensor calibration change that altered the scale of a critical feature. With version control, they identified the issue in days rather than weeks and retrained the model appropriately.

**Production Monitoring:**

Data preparation doesn't end at deployment. Implement continuous monitoring:

- **Data drift detection**: Track when input data distributions change significantly
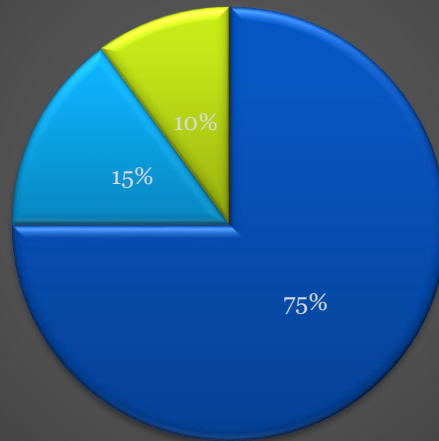- **Feature stability**: Monitor for sudden changes in feature values or availability

- **Data quality metrics**: Continuously measure completeness, validity, and timeliness
- **Automated alerts**: Trigger notifications when data quality falls below thresholds
- **Retraining triggers**: Define conditions that require model retraining

**Best Practices:**

- **Automate everything**: Manual steps don't scale and introduce errors
- **Document assumptions**: Explicitly state what you assume about data characteristics
- **Peer review**: Have another engineer review your data preparation logic
- **Establish feedback loops**: Track model performance and trace issues back to data preparation
- **Plan for evolution**: Design your pipeline to accommodate new data sources and features
- **Maintain audit trails**: Log all data access and transformations for compliance and debugging

| Item | Best Practice |
|---|---|
| Raw Data Snapshots | Save original, document range |
| Cleansing Configurations | Track logic, rules, versions |
| Feature Definitions | Document and code |
| Final Dataset Splits | Save splits with metadata |
| Data Quality Reports | Record metrics/issues per version |

## Time Spent on Manufacturing AI Projects



- Data Preparation
- Modeling/Training
- Deployment

# Conclusion: From Data to Deployable AI

Preparing manufacturing data for AI is neither quick nor glamorous, but it's the foundation of successful AI implementations. Industry practitioners estimate that 60-80% of AI project time is spent on data preparation—and for good reason. Poor data quality is the most common reason for AI project failures in manufacturing.

**Key Takeaways:**

1. **Invest in discovery**: Understand your data landscape before extraction
2. **Centralize thoughtfully**: Balance accessibility with operational safety
3. **Clean systematically**: Use domain knowledge to guide data quality decisions
4. **Engineer meaningful features**: Transform raw data into representations that expose relevant patterns
5. **Version and validate**: Build reproducibility and production readiness from the start

**Getting Started:**

Begin with a focused pilot project:

- Choose a specific use case with clear business value
- Identify the critical data sources for that use case
- Apply these five steps systematically
- Document lessons learned and build reusable components

- Scale successful patterns to additional use cases

Manufacturing AI success stories aren't built on sophisticated algorithms alone—they're built on high-quality, well-prepared data. By following these steps, you'll create the foundation for AI systems that deliver genuine business value in production environments.

---

# Additional Resources

**Industry Standards:**

- ISA-95: Integration of enterprise and control systems
- OPC UA: Interoperability standards for industrial automation
- MIMOSA: Machine information management open systems alliance

**Tools and Frameworks:**

- Apache Kafka: Real-time data streaming
- Apache Spark: Large-scale data processing
- Pandas: Data manipulation in Python
- Great Expectations: Data validation framework

**Further Reading:**

- "Building Machine Learning Powered Applications" by Emmanuel Ameisen
- "Designing Data-Intensive Applications" by Martin Kleppmann
- "Feature Engineering for Machine Learning" by Alice Zheng & Amanda Casari