

The Ultimate Guide to Data Cleaning for Machine Learning Projects (2025)

Introduction

Data cleaning is the foundation of any successful machine learning (ML) project. Raw data from the industry is rarely ready for modeling due to issues like missing values, inconsistencies, errors, and outliers. Well-executed data cleaning not only improves ML model performance but also ensures reliability and trust in business outcomes.

In real-world industry scenarios, data scientists spend approximately 60-80% of their time on data preparation and cleaning tasks. This comprehensive guide covers the top 10 industry-standard techniques used by leading organizations to ensure high-quality, analysis-ready datasets.

Top 10 Industry Techniques for Data Cleaning

1. Removing Duplicates

What is it?

Duplicate data occurs when the same record appears multiple times in a dataset. This can happen due to data entry errors, system glitches, or merging datasets from multiple sources. Duplicates distort statistical analysis, bias machine learning models, and waste computational resources.

Why is it important?

- Prevents model bias caused by over-representation of certain data points
- Improves computational efficiency by reducing dataset size
- Ensures accurate statistical measures (mean, median, variance)
- Eliminates redundancy in data storage

How to implement:

```
import pandas as pd
```

Load data

```
df = pd.read_csv('data.csv')
```

Check for duplicates

```
print(f"Number of duplicate rows: {df.duplicated().sum()}")
```

Remove all duplicate rows

```
df_clean = df.drop_duplicates()
```

Remove duplicates based on specific columns

```
df_clean = df.drop_duplicates(subset=['customer_id', 'transaction_date'])
```

Keep last occurrence instead of first

```
df_clean = df.drop_duplicates(keep='last')
```

Industry Best Practices:

- Automate duplicate detection in ETL pipelines
 - Define business rules for what constitutes a duplicate
 - Use hash functions for faster duplicate detection in large datasets
 - Maintain audit logs of removed duplicates for compliance
-

2. Handling Missing Values

What is it?

Missing data refers to the absence of values in certain fields or records. Data can be missing completely at random (MCAR), missing at random (MAR), or missing not at random (MNAR). Understanding the pattern of missingness is crucial for choosing the right handling strategy.

Why is it important?

- Most ML algorithms cannot handle missing values directly
- Improper handling can introduce bias and reduce model accuracy
- Missing data can indicate data quality issues or systematic problems
- Different imputation strategies affect model performance differently

How to implement:

```
import pandas as pd
import numpy as np
from sklearn.impute import SimpleImputer, KNNImputer
```

Analyze missing data

```
print(df.isnull().sum())
print(df.isnull().sum() / len(df) * 100) # Percentage
```

Strategy 1: Remove rows with missing values

```
df_clean = df.dropna()
```

Strategy 2: Remove columns with >50% missing data

```
threshold = 0.5
df_clean = df.dropna(thresh=int(threshold * len(df)), axis=1)
```

Strategy 3: Fill with statistical measures

```
df['age'].fillna(df['age'].mean(), inplace=True)
df['category'].fillna(df['category'].mode()[0], inplace=True)
```

Strategy 4: Forward fill (time series)

```
df['price'].fillna(method='ffill', inplace=True)
```

Strategy 5: KNN Imputation (advanced)

```
imputer = KNNImputer(n_neighbors=5)
df_imputed = pd.DataFrame(
    imputer.fit_transform(df.select_dtypes(include=[np.number])),
    columns=df.select_dtypes(include=[np.number]).columns
)
```

Strategy 6: Using scikit-learn SimpleImputer

```
imputer = SimpleImputer(strategy='median')
df[['column1', 'column2']] = imputer.fit_transform(df[['column1', 'column2']])
```

Industry Best Practices:

- Document the reason for missingness when possible
- Use domain knowledge to guide imputation strategies

- Consider creating a binary "is_missing" indicator feature
 - Test multiple imputation methods and compare model performance
 - Use advanced techniques like MICE (Multiple Imputation by Chained Equations) for complex scenarios
-

3. Standardizing and Normalizing Formats

What is it?

Standardization ensures that data follows consistent formats, units, and conventions across the entire dataset. This includes date formats, text capitalization, numerical scales, and categorical encodings.

Why is it important?

- Enables accurate data merging and joining operations
- Prevents errors in computations due to inconsistent units
- Improves model performance by putting features on similar scales
- Makes data more interpretable and easier to analyze

How to implement:

```
import pandas as pd  
from sklearn.preprocessing import StandardScaler, MinMaxScaler
```

Date standardization

```
df['date'] = pd.to_datetime(df['date'], format='%Y-%m-%d')  
df['year'] = df['date'].dt.year  
df['month'] = df['date'].dt.month
```

Text standardization

```
df['name'] = df['name'].str.title() # Title case  
df['email'] = df['email'].str.lower() # Lowercase  
df['phone'] = df['phone'].str.replace(r'[^\d]', "", regex=True) # Remove non-digits
```

Numeric standardization (Z-score normalization)

```
scaler = StandardScaler()  
df[['income', 'age']] = scaler.fit_transform(df[['income', 'age']])
```

Min-Max scaling (0-1 range)

```
minmax_scaler = MinMaxScaler()  
df[['score']] = minmax_scaler.fit_transform(df[['score']])
```

Currency standardization

```
df['price_usd'] = df['price'] * df['exchange_rate']
```

Unit conversion

```
df['height_cm'] = df['height_inches'] * 2.54
```

Industry Best Practices:

- Create data dictionaries defining standard formats
 - Use ISO standards for dates (ISO 8601), currencies (ISO 4217)
 - Implement validation rules at data entry points
 - Store original values before transformation for auditability
 - Use configuration files for format definitions
-

4. Correcting Typos and Inconsistencies

What is it?

Typos, spelling errors, and inconsistent naming conventions are common in manually-entered data. These errors can fragment categorical variables and reduce the effectiveness of grouping and aggregation operations.

Why is it important?

- Prevents artificial inflation of unique categories
- Improves accuracy of categorical analysis
- Enhances data matching and deduplication
- Reduces confusion in reporting and visualization

How to implement:

```
import pandas as pd  
from fuzzywuzzy import process, fuzz  
import re
```

Simple replacement

```
df['country'] = df['country'].replace({  
'USA': 'United States',  
'US': 'United States',  
'U.S.A': 'United States'  
})
```

Fuzzy matching for similar strings

```
def fuzzy_match(value, choices, threshold=80):
    match, score = process.extractOne(value, choices)
    return match if score >= threshold else value

cities = ['New York', 'Los Angeles', 'Chicago', 'Houston']
df['city_clean'] = df['city'].apply(lambda x: fuzzy_match(x, cities))
```

Pattern-based correction using regex

```
df['product'] = df['product'].str.replace(r'i[-s]?phone', 'iPhone',
                                         regex=True, flags=re.IGNORECASE)
```

Whitespace removal

```
df['category'] = df['category'].str.strip()
```

Remove special characters

```
df['name'] = df['name'].str.replace(r'[^w\s]', "", regex=True)
```

Industry Best Practices:

- Maintain master reference lists for categorical variables
- Use data validation at entry points (dropdown menus, autocomplete)
- Implement spell-checking algorithms
- Create mapping dictionaries for known variations
- Use NLP techniques for text normalization

5. Detecting and Handling Outliers

What is it?

Outliers are data points that deviate significantly from the majority of the dataset. They can be genuine extreme values, measurement errors, or data entry mistakes. Proper outlier detection requires understanding whether outliers are legitimate or erroneous.

Why is it important?

- Outliers can severely skew statistical measures and model predictions
- Some ML algorithms are sensitive to outliers (linear regression, KNN)
- Legitimate outliers may contain valuable information
- Different domains have different outlier tolerance levels

How to implement:

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

```
from scipy import stats  
from sklearn.ensemble import IsolationForest
```

Method 1: Z-score method (for normally distributed data)

```
z_scores = np.abs(stats.zscore(df['value']))  
df_no_outliers = df[z_scores < 3]
```

Method 2: IQR (Interquartile Range) method

```
Q1 = df['value'].quantile(0.25)  
Q3 = df['value'].quantile(0.75)  
IQR = Q3 - Q1  
lower_bound = Q1 - 1.5 * IQR  
upper_bound = Q3 + 1.5 * IQR  
df_filtered = df[(df['value'] >= lower_bound) & (df['value'] <= upper_bound)]
```

Method 3: Isolation Forest (advanced, multivariate)

```
iso_forest = IsolationForest(contamination=0.1, random_state=42)  
outliers = iso_forest.fit_predict(df[['feature1', 'feature2']])  
df_clean = df[outliers == 1]
```

Method 4: Capping/Winsorizing (preserve data size)

```
df['value_capped'] = df['value'].clip(lower=lower_bound, upper=upper_bound)
```

Method 5: Log transformation for skewed data

```
df['value_log'] = np.log1p(df['value'])
```

Industry Best Practices:

- Always visualize outliers before removal (box plots, scatter plots)
- Consult domain experts to validate outlier treatment
- Document the rationale for outlier handling
- Consider robust statistical methods that are less sensitive to outliers
- Use different techniques for different types of data distributions

6. Removing Irrelevant Data

What is it?

Irrelevant data includes features, records, or information that do not contribute to the analytical goals or ML modeling objectives. This could be redundant features, outdated records, or data collected for different purposes.

Why is it important?

- Reduces computational complexity and training time
- Improves model interpretability
- Prevents overfitting by reducing noise
- Optimizes storage and memory usage

How to implement:

```
import pandas as pd  
import numpy as np  
from sklearn.feature_selection import SelectKBest, chi2, mutual_info_classif
```

Remove columns with single unique value (zero variance)

```
unique_counts = df.nunique()  
cols_to_drop = unique_counts[unique_counts == 1].index  
df_clean = df.drop(columns=cols_to_drop)
```

Remove highly correlated features

```
corr_matrix = df.corr().abs()  
upper_triangle = corr_matrix.where(  
    np.triu(np.ones(corr_matrix.shape), k=1).astype(bool)  
)  
to_drop = [col for col in upper_triangle.columns if any(upper_triangle[col] > 0.95)]  
df_clean = df.drop(columns=to_drop)
```

Remove columns with too many missing values

```
threshold = 0.7  
df_clean = df.dropna(thresh=int(threshold * len(df)), axis=1)
```

Feature selection using statistical tests

```
selector = SelectKBest(mutual_info_classif, k=10)
X_selected = selector.fit_transform(X, y)
selected_features = X.columns[selector.get_support()]
```

Remove specific irrelevant columns

```
df_clean = df.drop(columns=['id', 'timestamp', 'internal_code'])
```

Industry Best Practices:

- Use feature importance scores from tree-based models
 - Apply domain knowledge to identify relevant features
 - Consider business requirements and regulatory needs
 - Perform exploratory data analysis (EDA) before removal
 - Keep a backup of original data before dropping columns
-

7. Data Type Conversion and Validation

What is it?

Ensuring each column has the correct data type (integer, float, string, datetime, categorical) is fundamental. Incorrect data types can cause errors in calculations, inefficient memory usage, and prevent proper analysis.

Why is it important?

- Enables appropriate operations and functions for each data type
- Optimizes memory usage and performance
- Prevents type-related errors during analysis
- Ensures compatibility with ML algorithms

How to implement:

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

Check current data types

```
print(df.dtypes)
print(df.info())
```

Convert to numeric (with error handling)

```
df['age'] = pd.to_numeric(df['age'], errors='coerce')
```

Convert to datetime

```
df['date'] = pd.to_datetime(df['date'], errors='coerce')
```

Convert to categorical (reduces memory for repeated values)

```
df['category'] = df['category'].astype('category')
```

Convert boolean strings to boolean type

```
df['is_active'] = df['is_active'].map({'True': True, 'False': False})
```

Convert to integer (handling NaN)

```
df['count'] = df['count'].fillna(0).astype(int)
```

Optimize memory usage

```
def optimize_dtypes(df):
    for col in df.select_dtypes(include=['int']).columns:
        df[col] = pd.to_numeric(df[col], downcast='integer')
    for col in df.select_dtypes(include=['float']).columns:
        df[col] = pd.to_numeric(df[col], downcast='float')
    return df

df = optimize_dtypes(df)
```

Validate data types

```
assert df['age'].dtype in [np.int64, np.int32], "Age should be integer"
assert pd.api.types.is_datetime64_any_dtype(df['date']), "Date should be datetime"
```

Industry Best Practices:

- Define data type schemas upfront
 - Use validation libraries like Pandera or Great Expectations
 - Automate type conversion in data ingestion pipelines
 - Monitor data type consistency across data sources
 - Use categorical types for memory efficiency with repeated values
-

8. Handling Encoding Issues and Special Characters

What is it?

Text data often contains encoding issues (UTF-8, ASCII, Latin-1), special characters, emojis, or hidden characters that can cause processing errors or inconsistent analysis.

Why is it important?

- Prevents processing errors and data corruption
- Ensures cross-platform compatibility
- Improves text analysis and NLP tasks
- Maintains data integrity during transfers

How to implement:

```
import pandas as pd
import re
import unicodedata
import html
```

Read file with specific encoding

```
df = pd.read_csv('data.csv', encoding='utf-8')
```

Handle encoding errors during reading

```
df = pd.read_csv('data.csv', encoding='latin-1', encoding_errors='ignore')
```

Remove special characters

```
df['text'] = df['text'].apply(lambda x: re.sub(r'^\w\s', " ", str(x)))
```

Remove emojis

```
def remove_emoji(text):
    emoji_pattern = re.compile("["
        u"\U0001F600-\U0001F64F" # emoticons
        u"\U0001F300-\U0001F5FF" # symbols & pictographs
        u"\U0001F680-\U0001F6FF" # transport & map symbols
        u"\U0001F1E0-\U0001F1FF" # flags
    "]+", flags=re.UNICODE)
    return emoji_pattern.sub(r'', text)
```

```
df['text_clean'] = df['text'].apply(remove_emoji)
```

Normalize unicode characters

```
def normalize_text(text):
    return unicodedata.normalize('NFKD', text).encode('ascii', 'ignore').decode('utf-8')

df['text_normalized'] = df['text'].apply(normalize_text)
```

Remove leading/trailing whitespace and extra spaces

```
df['text'] = df['text'].str.strip().str.replace(r'\s+', ' ', regex=True)
```

Convert HTML entities

```
df['text'] = df['text'].apply(html.unescape)
```

Industry Best Practices:

- Standardize on UTF-8 encoding across all systems
 - Implement encoding detection and conversion early in pipelines
 - Preserve original text in separate column before cleaning
 - Use NLP libraries for comprehensive text preprocessing
 - Test with international character sets
-

9. Cross-field Validation and Consistency Checks

What is it?

Cross-field validation ensures logical consistency between related fields. For example, end dates should be after start dates, child age should be less than parent age, and totals should match sum of components.

Why is it important?

- Detects logical errors and data entry mistakes
- Ensures business rule compliance
- Improves data reliability and trustworthiness
- Prevents downstream calculation errors

How to implement:

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

Date validation: end_date > start_date

```
df['date_valid'] = df['end_date'] > df['start_date']
invalid_dates = df[~df['date_valid']]
print(f"Invalid date ranges: {len(invalid_dates)}")
```

Numeric range validation

```
df['age_valid'] = (df['age'] >= 0) & (df['age'] <= 120)
df = df[df['age_valid']]
```

Cross-field calculation validation

```
df['total_calculated'] = df['item1'] + df['item2'] + df['item3']
df['total_match'] = np.isclose(df['total'], df['total_calculated'], rtol=0.01)
mismatches = df[~df['total_match']]
```

Conditional validation

If married, spouse name should not be null

```
df['spouse_valid'] = ~((df['marital_status'] == 'Married') & (df['spouse_name'].isnull()))
```

Percentage validation (should sum to 100)

```
percentage_cols = ['pct_category1', 'pct_category2', 'pct_category3']
df['pct_total'] = df[percentage_cols].sum(axis=1)
df['pct_valid'] = np.isclose(df['pct_total'], 100, atol=1)
```

Hierarchical consistency (city should match state/country)

```
location_mapping = {
    'New York': {'state': 'NY', 'country': 'USA'},
    'Los Angeles': {'state': 'CA', 'country': 'USA'}
}
df['location_valid'] = df.apply(
    lambda row: location_mapping.get(row['city'], {}).get('state') == row['state'],
    axis=1
)
```

Industry Best Practices:

- Define comprehensive business rules documentation
- Implement validation at data entry points

- Create automated data quality reports
 - Use constraint-based validation frameworks
 - Flag invalid records rather than deleting for review
-

10. Automated Data Quality Monitoring and Profiling

What is it?

Data profiling involves systematically analyzing datasets to understand their structure, content, quality, and relationships. Automated monitoring tracks data quality metrics over time and alerts when issues arise.

Why is it important?

- Provides comprehensive understanding of data characteristics
- Enables early detection of data quality degradation
- Supports data governance and compliance
- Reduces manual inspection effort

How to implement:

```
import pandas as pd
import pandas_profiling
from ydata_profiling import ProfileReport
import great_expectations as ge
```

Method 1: Pandas Profiling (comprehensive report)

```
profile = ProfileReport(df, title="Data Quality Report", explorative=True)
profile.to_file("data_quality_report.html")
```

Method 2: Manual profiling functions

```
def data_quality_report(df):
    report = {
        'rows': len(df),
        'columns': len(df.columns),
        'duplicates': df.duplicated().sum(),
        'missing_values': df.isnull().sum().to_dict(),
        'data_types': df.dtypes.to_dict(),
        'memory_usage': df.memory_usage(deep=True).sum() / 1024**2, # MB
        'numeric_stats': df.describe().to_dict(),
        'unique_counts': df.nunique().to_dict()
    }
    return report

quality_report = data_quality_report(df)
print(pd.DataFrame(quality_report))
```

Method 3: Great Expectations (industry standard)

```
ge_df = ge.from_pandas(df)
ge_df.expect_column_values_to_not_be_null('customer_id')
ge_df.expect_column_values_to_be_between('age', min_value=0, max_value=120)
ge_df.expect_column_values_to_be_in_set('status', ['active', 'inactive'])
validation_results = ge_df.validate()
```

Custom quality metrics

```
def calculate_quality_score(df):
    completeness = 1 - (df.isnull().sum().sum() / (len(df) * len(df.columns)))
    uniqueness = 1 - (df.duplicated().sum() / len(df))
    validity = 0.95 # Based on custom validation rules
```

```
    quality_score = (completeness + uniqueness + validity) / 3
    return quality_score * 100
```

```
print(f'Data Quality Score: {calculate_quality_score(df):.2f}%')
```

Automated monitoring over time

```
def monitor_data_quality(df, baseline_metrics):
    current_metrics = {
        'null_percentage': df.isnull().sum().sum() / (len(df) * len(df.columns)),
        'duplicate_percentage': df.duplicated().sum() / len(df),
        'row_count': len(df)
    }
```

```
    alerts = []
    if current_metrics['null_percentage'] > baseline_metrics['null_percentage'] * 1.5:
        alerts.append("WARNING: Missing values increased significantly")
    if current_metrics['row_count'] < baseline_metrics['row_count'] * 0.8:
        alerts.append("WARNING: Row count dropped significantly")

    return current_metrics, alerts
```

Industry Best Practices:

- Implement continuous data quality monitoring in production pipelines
- Set up automated alerts for quality threshold violations

- Create data quality dashboards for stakeholders
 - Use tools like Great Expectations, Deequ, or custom frameworks
 - Schedule regular data profiling and quality audits
 - Maintain historical quality metrics for trend analysis
-

Building a Complete Data Cleaning Pipeline

Here's how to combine all techniques into a production-ready pipeline:

```
import pandas as pd
import numpy as np
from sklearn.impute import KNNImputer
from sklearn.preprocessing import StandardScaler

class DataCleaningPipeline:
    def __init__(self):
        self.imputer = None
        self.scaler = None

    def clean_data(self, df):
        """Complete data cleaning pipeline"""
        print("Starting data cleaning pipeline...")

        # 1. Remove duplicates
        df = self._remove_duplicates(df)

        # 2. Handle missing values
        df = self._handle_missing_values(df)

        # 3. Standardize formats
        df = self._standardize_formats(df)

        # 4. Correct inconsistencies
        df = self._correct_inconsistencies(df)

        # 5. Handle outliers
        df = self._handle_outliers(df)

        # 6. Remove irrelevant data
        df = self._remove_irrelevant_data(df)

        # 7. Convert data types
```

```

df = self._convert_data_types(df)

# 8. Handle encoding
df = self._handle_encoding(df)

# 9. Validate cross-fields
df = self._validate_cross_fields(df)

# 10. Generate quality report
self._generate_quality_report(df)

print("Data cleaning completed!")
return df

def _remove_duplicates(self, df):
    initial_rows = len(df)
    df = df.drop_duplicates()
    print(f"Removed {initial_rows - len(df)} duplicate rows")
    return df

def _handle_missing_values(self, df):
    # Strategy based on percentage of missing data
    for col in df.columns:
        missing_pct = df[col].isnull().sum() / len(df)
        if missing_pct > 0.5:
            df = df.drop(columns=[col])
        elif df[col].dtype in ['int64', 'float64']:
            df[col] = df[col].fillna(df[col].median())
        else:
            df[col] = df[col].fillna(df[col].mode()[0] if not df[col].mode().empty else 'U')
    return df

def _standardize_formats(self, df):
    # Implement format standardization
    return df

def _correct_inconsistencies(self, df):
    # Implement consistency corrections

```

```
    return df

def _handle_outliers(self, df):
    numeric_cols = df.select_dtypes(include=[np.number]).columns
    for col in numeric_cols:
        Q1 = df[col].quantile(0.25)
        Q3 = df[col].quantile(0.75)
        IQR = Q3 - Q1
        lower_bound = Q1 - 1.5 * IQR
        upper_bound = Q3 + 1.5 * IQR
        df[col] = df[col].clip(lower=lower_bound, upper=upper_bound)
    return df

def _remove_irrelevant_data(self, df):
    # Remove zero-variance columns
    unique_counts = df.nunique()
    cols_to_drop = unique_counts[unique_counts == 1].index
    df = df.drop(columns=cols_to_drop)
    return df

def _convert_data_types(self, df):
    # Implement type conversions
    return df

def _handle_encoding(self, df):
    # Implement encoding handling
    return df

def _validate_cross_fields(self, df):
    # Implement validation logic
    return df

def _generate_quality_report(self, df):
    print("\n==== Data Quality Report ====")
    print(f"Total Rows: {len(df)}")
    print(f"Total Columns: {len(df.columns)}")
    print(f"Missing Values: {df.isnull().sum().sum()}"
```

```

print(f"Duplicates: {df.duplicated().sum()}")
print(f"Memory Usage: {df.memory_usage(deep=True).sum() / 1024**2:.2f} M"

```

Usage

```

pipeline = DataCleaningPipeline()
cleaned_df = pipeline.clean_data(raw_df)

```

Industry Tools and Frameworks

Tool/Framework	Purpose	Best For
Pandas	Core data manipulation	All data cleaning tasks
OpenRefine	GUI-based cleaning	Business users, exploratory cleaning
Great Expectations	Data validation	Production pipelines, quality assurance
Trifecta	Visual data preparation	Non-technical users
AWS Glue	Cloud ETL	Large-scale cloud data
Apache Spark	Distributed processing	Big data cleaning
dbt	Data transformation	Analytics engineering
Talend	Enterprise data integration	Complex ETL workflows

Table 1: Popular industry tools for data cleaning

Key Takeaways for Client Meetings

- 1. Data quality directly impacts model performance** - investing in proper cleaning saves time and improves outcomes
- 2. Automate wherever possible** - manual cleaning doesn't scale
- 3. Document all cleaning decisions** - maintain reproducibility and auditability
- 4. Involve domain experts** - they can identify valid outliers vs errors
- 5. Monitor continuously** - data quality degrades over time without oversight
- 6. Use appropriate tools** - match tools to team skills and data scale
- 7. Test cleaning pipelines** - validate that cleaning improves downstream metrics
- 8. Maintain data lineage** - track transformations for compliance and debugging

Conclusion

Data cleaning is not a one-time task but an ongoing process that requires systematic approaches, domain knowledge, and the right tools. By implementing these top 10 industry techniques, you can ensure your ML projects are built on a foundation of high-quality, reliable data.

The investment in proper data cleaning pays dividends through improved model accuracy, faster development cycles, and more trustworthy business insights. As data volumes and complexity continue to grow, automated, scalable cleaning pipelines become essential for maintaining competitive advantage.

Further Resources and References

- Python Pandas Documentation: <https://pandas.pydata.org/docs/>
- Scikit-learn Preprocessing: <https://scikit-learn.org/stable/modules/preprocessing.html>
- Great Expectations: <https://greatexpectations.io/>
- Data Cleaning Best Practices (CCS Learning Academy, 2025)
- The Role of Machine Learning in Data Cleaning (Numerous.ai, 2025)
- 10 Essential Data Cleaning Techniques (KDnuggets, 2025)
- Top 8 Data Cleaning Techniques for Better Results (Savant Labs, 2025)