# Machine Learning: Quick Reference Guide

## Contents

**1**

## Data Collection / Ingestion

- Acquire raw data from relevant sources to build a reliable dataset for analysis and modeling.
- Source identification
- Data extraction & loading
- Initial format validation
- Documentation & metadata

**Sources:** Databases, APIs, Files, Streaming, Sensors, Logs, etc.

**Tools:** Apache Kafka, Airflow, BeautifulSoup, AWS S3, SQLAlchemy

**2**

## Data Understanding

Obtain a high-level understanding of data structure, features, and potential challenges before cleaning or modeling.

**Data Exploration:**
- Data types
- Distributions
- Integrity checks & constraints
- Outliers/anomalies
- Resolve conflicts
- Validate integration

**Tools:** Pandas, NumPy, SQL, pandas-profiling, Matplotlib, Seaborn

**3**

## Data Integration

Combine multiple data sources into a unified dataset.

- Selection
- Analysis prep
- Normalization
- Formatting

**Tasks:**
- Data quality checks
- Feature merging
- Handling missing/duplicate records
- Facilitate cross-source alignment
- Transformation & pipelines

**Tools:** Apache Spark, Dask, DBT, Pandas, Great Expectations

**4**

## Initial Validation

Deep initial examination and identify preprocessing needs.

- Feature creation
- Feature transformation
- Feature selection
- Temporal/spatial processing

**Tools:** Scikit-learn, PySpark, SciPy, statsmodels

**5**

## Exploratory Data Analysis (EDA)

Deep dive into data patterns, distributions, and relationships.

**Key Activities:**
- Missing value handling
- Distribution analysis
- Data type checks
- Outlier detection
- Consistency checks

**Cross-checking:**
- Hypothesis testing
- Sufficient sample size

**Tools:** Pandas, Matplotlib, Seaborn, Plotly, SciPy, Jupyter

**6**

## Feature Engineering

Create new features or transform existing ones to improve model performance.

- Feature creation
- Feature transformation
- Feature selection
- Text/spatial processing

**Tools:** Scikit-learn, Featuretools, Category Encoders, NLTK, spaCy, OpenCV

**7**

## Data Splitting

Fit partitioned data into modeling sets.

**Key Activities:**
- Train/test split
- Validation set
- Stratification
- Time-based splitting
- Cross-validation

**Best Practices:**
- Maintain temporal integrity
- Keep data distribution consistent

**Tools:** train_test_split, StratifiedKFold, TimeSeriesSplit (Scikit-learn)

**8**

## Model Training

Train and optimize machine learning models using the processed dataset to learn patterns and make predictions effectively.

**Algorithms by Type:**

**Regression:** Linear/Ridge/Lasso, SVR, Decision Trees, Random Forest, XGBoost/LightGBM/CatBoost, Neural Networks
**Classification:** Logistic Regression, Naive Bayes, KNN, SVM, Decision Trees, Random Forest, XGBoost/LightGBM/CatBoost, Neural Networks
**Clustering:** K-Means, DBSCAN, Hierarchical, GMM
**Deep Learning:** CNN, RNN/LSTM, Transformers, Autoencoders, GAN
**Ensemble:** Bagging, Boosting, Stacking

**Tools:** Scikit-learn, XGBoost, LightGBM, CatBoost, TensorFlow, PyTorch, MLflow

**9**

## Model Evaluation

Evaluate trained models to ensure accuracy, reliability, and business relevance.

**Key Activities:**
- Measure performance on unseen data
- Analyze errors and interpret model behavior
- Compare multiple models
- Validate results against business goals

**Best Practices:**
- Use appropriate evaluation metrics
- Ensure results generalize to real-world scenarios
- Align performance with success criteria

**Tools:** Scikit-learn metrics, SHAP, LIME, MLflow, ROC/AUC plotters

**Metrics:** Accuracy, Precision, Recall, F1, ROC-AUC, MSE, RMSE, MAE, R²

**10**

## Deployment & Monitoring

Put model into production and ensure lifecycle management.

**Key Activities:**
- Model serialization
- Deployment
- Integration
- Monitoring
- Maintenance

**Best Practices:**
- Continuous evaluation
- Retraining triggers
- Alerting & dashboards

**Tools:** Flask/FastAPI, Docker, Kubernetes, AWS SageMaker, MLflow, Prometheus, Grafana

**Data Collection / Ingestion → Data Understanding → Data Integration →**

**Initial Validation → Exploratory Data Analysis (EDA) → Feature**

**Engineering → Data Splitting → Model Training → Model Evaluation →**

**Deployment & Monitoring**

# Example Notebook – IPL Dataset Analysis

## 1. Data Collection / Ingestion

**Purpose:**
Acquire raw data from all relevant sources to create a high quality, reliable dataset. This step ensures reproducibility, traceability, and sets the foundation for downstream analysis and modeling. Proper ingestion prevents garbage-in, garbage-out issues.

---

**Sub-Steps & Activities:**

**Source Identification:**

- Identify all potential data sources:

    - **Databases:** MySQL, PostgreSQL, MongoDB, Oracle

    - **Data Warehouses:** Snowflake, BigQuery, Redshift

    - **APIs & Web Services:** REST, GraphQL, SOAP

    - **Files:** CSV, Excel, JSON, Parquet, XML, logs

    - **Streaming/IOT:** Kafka, MQTT, Kinesis, IoT sensors

- Evaluate quality dimensions: freshness, completeness, reliability, and frequency of updates.

- Determine data structure: structured (tables), semi-structured (JSON, XML), or unstructured (text, images, videos).

- **Interview Tip:** Be ready to justify why multiple sources improve richness but may increase complexity and inconsistency.

---

**Data Extraction & Loading:**

- Use the right tools & libraries depending on source type:

    - **CSV/Excel:** pandas.read_csv(), pandas.read_excel()

    - **SQL Databases:** sqlalchemy, psycopg2 with pd.read_sql()

    - **APIs:** requests, httpx (handle pagination, rate limits)

    - **Web scraping:** BeautifulSoup, Scrapy, Selenium

    - **Cloud SDKs:** AWS boto3, Azure SDK, GCP client libraries

- Handle large datasets: chunking, streaming, parallelization, or incremental loads.

- Consider data consistency during ingestion: transactional integrity, idempotency.

---

**Initial Format Validation:**

- Verify row & column counts, data types, and missing values.

- Detect malformed, duplicated, or corrupt records.

- Ensure consistent encoding and timezone handling for timestamp fields.

- Log all anomalies for later cleaning or alerting.

---

**Documentation & Metadata:**

- **Maintain ingestion logs:** source, extraction timestamp, access credentials, file format, preprocessing steps.

- **Create and update a data dictionary:** column names, types, allowed values, and descriptions.

- Track **data lineage** to ensure traceability for auditing and reproducibility.

---

**Common Challenges:**

- API rate limits, authentication, and pagination.

- Real-time vs batch ingestion trade-offs.

- Inconsistent formats and missing or corrupted data.

- Large datasets requiring distributed processing or cloud storage solutions.

---

**Example Code:**

```
import pandas as pd
# Load CSV file
df = pd.read_csv("data/file.csv")
# Quick checks
print("Rows & Columns:", df.shape)
print("Data Types:\n", df.dtypes)
print("Missing Values:\n", df.isnull().sum())
```

---

**Tools:**

- Apache Kafka, Apache Spark Streaming, Apache Airflow
- **Python libraries**: requests, BeautifulSoup, Scrapy
- **Cloud services:** AWS S3, Google Cloud Storage, Azure Data Lake

- **Database connectors:** psycopg2, pymongo, SQLAlchemy & more
- **API tools:** Postman, REST clients

---

**Interview Questions:**

1. How would you handle real-time streaming data ingestion?

2. How do you ensure data provenance, traceability, and reproducibility?

3. What strategies would you use for integrating multiple data sources with inconsistent formats?

4. How would you deal with very large datasets that cannot fit into memory?

---

# 2. Data Understanding

**Purpose:**
Gain a comprehensive view of the dataset to understand its structure, quality, and potential challenges before cleaning or modeling. This step helps in identifying important features, spotting anomalies, and planning preprocessing strategies.

---

**Sub-Steps & Activities:**

**Dataset Shape & Size:**

- **Check rows vs columns**: large number of features may require dimensionality reduction.

- **Assess dataset size sufficiency**: too few samples may cause overfitting; too many may need sampling or distributed processing.

---

**Column Inspection:**

- Review **column names** and understand their meaning.

- Identify **irrelevant columns** like unique IDs, timestamps, or metadata that don't contribute to modeling.

- Detect potential **redundant or derived features** that could introduce multicollinearity.

---

**Data Types & Feature Identification:**

- Categorize features as:

    o **Numerical:** integers, floats

    o **Categorical:** strings, enums

- o **Datetime:** timestamps

- o **Text or unstructured:** free text, JSON blobs

- Mark **modeling candidates** vs columns to exclude as metadata.

- Check for **mixed data types** in a column (e.g., strings in numeric columns).

---

**Target Variable Identification:**

- Determine **prediction type**:

  - o Binary classification

  - o Multi-class classification

  - o Regression

- Examine **distribution of target variable**: check for imbalance that may affect model performance.

- Visualize target using **histograms, bar plots, or pie charts**.

---

**Summary Statistics & Quick Exploration:**

- For numerical features: df.describe(), df.info(), df.isnull().sum()

- For categorical features: value_counts(), unique()

- Identify **missing values, zeros, or constant columns**.

- Spot **outliers** using IQR, z-score, or visualizations like boxplots.

---

**Initial Observations & Insights:**

- Note patterns, correlations, and anomalies.

- Identify features that may need **encoding, scaling, or transformation**.

- List columns for **feature engineering**, e.g., combining, binning, or deriving new features.

---

**Example Code:**

```
# General info and statistics

print(df.info())

print(df.describe())


# Categorical feature exploration

for col in df.select_dtypes(include=['object', 'category']).columns:

print(f"Value counts for {col}:")

print(df[col].value_counts())

# Check missing values

print("Missing Values:\n", df.isnull().sum())
```

**Tools:**

- Pandas, NumPy, Polars
- SQL (PostgreSQL, MySQL, BigQuery)
- **Data profiling:** pandas-profiling (ydata-profiling), Great Expectations
- **Visualization:** Matplotlib, Seaborn, Plotly
- Jupyter Notebooks, Google Colab

**Interview Tips:**

1. Explain why checking target distribution is critical for modeling decisions.

2. How would you handle imbalanced classes? (undersampling, oversampling, SMOTE, class weighting)8

3. How do you decide which features are relevant or redundant?

4. How do you detect and handle outliers or anomalies before modeling?

# 3. Data Integration

**Purpose:**
Combine data from multiple sources into a single, coherent dataset to provide a holistic view for analysis and modeling. Proper integration ensures relationships across different datasets are preserved, data consistency is maintained, and no critical information is lost or duplicated.

**Sub-Steps & Activities:**

**Identify Join Keys:**

- Determine common fields across datasets such as customer_id, user_id, order_id, timestamp, or composite keys.

- Verify uniqueness and consistency of keys - duplicate or missing keys can lead to incorrect joins.

- Standardize key formats (e.g., string casing, trimming spaces, consistent data types) before joining.

---

**Merge / Join Operations:**

- Use appropriate join types depending on the business need:

    o **Inner Join:** Keeps only matching records.

    o **Left Join:** Keeps all records from the left table, adds matches from the right.

    o **Right Join:** Keeps all records from the right table, adds matches from the left.

    o **Outer Join:** Keeps all records from both sides, filling missing values with NaN.

- Validate the join logic by checking row counts before and after merging.

- Consider using multi-key joins for hierarchical data relationships.

---

**Concatenation (Vertical Integration):**

- Combine multiple datasets with the same schema (columns) using pd.concat() or union operations in SQL.

- Ensure column ordering and data types are consistent before concatenation.

- Add a **source column** if tracking the origin of data is necessary.

---

**Resolve Conflicts & Inconsistencies:**

- Handle **duplicate columns** by renaming or dropping them.

- Resolve **naming conflicts**: standardize column names, units, and value representations.

- Address **semantic conflicts** where the same column name means different things across datasets.

---

**Validate Integration:**

- Verify **row counts** and **unique key counts** before and after merging.

- Ensure no unintended **data duplication** or **record loss**.

- Perform **sanity checks** by sampling merged data and verifying correctness of combined values.

---

**Common Scenarios:**

- Merging **customer demographics** with **transaction history**.

- Combining **monthly/quarterly data** into a yearly dataset.

- Integrating **external datasets** (e.g., weather, market data, economic indicators) to enrich features.

---

**Tools & Techniques:**

- **Pandas:** pd.merge(), pd.concat(), join()

- **SQL:** JOIN (INNER, LEFT, RIGHT, FULL), UNION

- **Big Data:** Spark join(), union(), withColumnRenamed() for large-scale integration

---

**Example Code:**

```
import pandas as pd

# Merge two datasets on customer_id

merged_df = pd.merge(customers_df, transactions_df, on='customer_id', how='left')

# Concatenate multiple monthly datasets

combined_df = pd.concat([jan_df, feb_df, mar_df], axis=0)

# Validate integration

print("Rows before merge:", len(customers_df), len(transactions_df))

print("Rows after merge:", len(merged_df))
```

---

**Tools:**

- Apache Spark, Dask
- **ETL tools:** Apache NiFi, Talend, Informatica
- **Data validation:** Great Expectations, Pandera
- **Database tools:** DBT (Data Build Tool)
- Pandas, SQL for merging/joining

---

**Interview Tips:**

1. How do you ensure no data duplication or loss after merging multiple sources?

2. What strategies would you use to resolve naming conflicts or mismatched schemas?

3. How would you integrate external data that updates at a different frequency than internal data?

4. How do you handle one-to-many or many-to-many relationships during joins?

# 4. Initial Validation

**Purpose:**
Perform a systematic quality check of the ingested and integrated data before deeper analysis or feature engineering. The goal is to detect, document, and quantify potential data issues early, so they can be addressed before they affect downstream modeling or business decisions.

**Sub-Steps & Activities:**

**Missing Values Analysis:**

- Identify and count missing or NaN values in each column: df.isnull().sum()

- Analyze patterns of missingness:

    o **MCAR (Missing Completely at Random)** - safe to drop.

    o **MAR (Missing at Random)** - can be imputed.

    o **MNAR (Missing Not at Random)** - requires deeper investigation.

- Prioritize **critical columns** (e.g., target variable, primary keys) for immediate attention.

**Duplicate Record Detection:**

- Check for duplicated rows: df.duplicated()

- Decide whether to **remove**, **aggregate**, or **investigate** duplicates depending on the use case.

- Ensure no **key-level duplication** (e.g., multiple records for the same unique ID unless expected).

**Data Type & Schema Validation:**

- Confirm that each column matches the expected data type (int, float, datetime, string): df.dtypes.

- Identify invalid entries (e.g., text in numeric columns, malformed dates, inconsistent encodings).

- Enforce schema consistency, especially when integrating data from multiple sources.

---

**Outlier Detection:**

- Use statistical methods (IQR, z-score) to find extreme values.

- Validate whether outliers are genuine signals (e.g., high-value customers) or data entry errors.

- Visualize distributions (boxplots, histograms) to detect anomalies quickly.

---

**Consistency & Logical Checks:**

- Verify **temporal logic:** end_date > start_date.

- Validate **categorical constraints:** all category values fall within expected sets.

- Cross-verify related columns (e.g., total = sum of components, age derived from date of birth).

---

**Completeness Assessment:**

- Measure the **coverage of critical fields** (e.g., >95% non-null for key features).

- Flag columns with insufficient data for downstream usage or imputation strategies.

---

**Data Quality Reporting:**

- Document all issues: type, location, severity, and recommended actions.

- Summarize findings into a **data quality report** for stakeholders.

- This becomes a reference document for data cleaning and pipeline improvements.

---

**Example Code:**

```python
import pandas as pd

# Missing values
print("Missing values:\n", df.isnull().sum())

# Duplicates
print("Duplicate rows:", df.duplicated().sum())

# Data types
print("Data types:\n", df.dtypes)

# Outlier check (IQR example)
Q1, Q3 = df['amount'].quantile([0.25, 0.75])
IQR = Q3 - Q1
outliers = df[(df['amount'] < Q1 - 1.5*IQR) | (df['amount'] > Q3 + 1.5*IQR)]
print("Outliers found:", len(outliers))
```

**Output:**

- A **Data Quality Report** containing:
    - Missing values per column and percentage
    - Count and nature of duplicates
    - Data type mismatches
    - Outlier distribution and possible causes
    - Consistency rule violations
    - Completeness scores for key features

**Tools:**

- Scikit-learn (preprocessing, feature_selection modules)
- PySpark MLlib
- Custom Python scripts
- **Statistical libraries:** SciPy, statsmodels
- **Feature validation:** pandas, numpy

**Interview Tips:**

1. How do you differentiate between real outliers and data errors?

2. What methods do you use to assess data completeness?

3. How would you approach missing values in critical vs non-critical fields?

4. How do you ensure logical consistency across dependent fields?

# 5. Exploratory Data Analysis (EDA)

**Purpose:**
Explore, visualize, and statistically analyse the dataset to uncover patterns, distributions, relationships, and anomalies. EDA guides feature engineering, model selection, and data preprocessing decisions. It helps you understand the "story" the data is telling before you build predictive models.

**Sub-Steps & Activities:**

**Data Cleaning (Pre-EDA Preparation):**

- **Handle Missing Values:**

    o Impute with mean/median/mode (for numeric), most frequent (for categorical), or predictive imputation techniques (e.g., KNNImputer).

    o Drop columns or rows if missingness is too high and not critical.

    o Create binary flags for missingness if it might carry signal.

- **Outlier Treatment:**

    o Remove extreme outliers if they're likely errors.

    o Cap/floor values using winsorization if they're genuine but influential.

- **Standardize Formats:**

    o Normalize date formats, text casing, and measurement units.

    o Convert categorical text to categorical dtype for efficiency.

- **Deduplication:**

    o Remove duplicate records identified in the validation step.

**Visualization & Pattern Discovery:**

- **Univariate Analysis:** Explore individual feature distributions.

- Histograms and density plots → Understand skewness, modality, range.

- Box plots → Detect outliers and spread.

- **Bivariate Analysis:** Explore relationships between two variables.

    - Scatter plots → Detect linear/non-linear relationships.

    - Correlation heatmaps → Identify strong feature relationships or multicollinearity.

    - Pair plots → Visualize interactions across multiple numerical features.

- **Target Variable Analysis:**

    - Visualize class distribution in classification problems (bar plots, pie charts).

    - Examine target vs feature relationships to identify predictive signals.

- **Categorical Feature Analysis:**

    - Bar charts / count plots for category frequency.

    - Cross-tabulations or stacked bar plots for category vs target interactions.

---

**Statistical Analysis:**

- **Correlation Matrices:**

    - Identify multicollinearity (high correlations between features).

    - Decide which features to drop or combine.

- **Distribution Analysis:**

    - Check normality using skewness, kurtosis, Shapiro-Wilk tests.

    - Apply transformations (log, sqrt, Box-Cox) if needed.

- **Group-by Aggregations:**

    - Explore average, median, or sum of key metrics across categorical groups.

    - Identify strong differentiators for the target.

---

**Insight Generation & Hypothesis Formation:**

- Which features have the strongest correlation with the target?

- Are there clear clusters or natural groupings in the data?

- Is the data balanced, or does it require resampling?

- Which transformations (e.g., scaling, encoding) might improve model performance?

- Are there redundant features that can be removed?

---

**Example Code:**

```python
import pandas as pd

import seaborn as sns

import matplotlib.pyplot as plt


# Univariate: Histogram

df['age'].hist(bins=30)

plt.title("Age Distribution")

plt.show()



# Bivariate: Correlation heatmap

plt.figure(figsize=(10,8))

sns.heatmap(df.corr(), annot=True, cmap='coolwarm')

plt.title("Feature Correlation Matrix")

plt.show()

# Target Distribution

df['target'].value_counts().plot(kind='bar')

plt.title("Target Class Distribution")

plt.show()
```

**Common Tools:**

- **Visualization:** matplotlib, seaborn, plotly, altair

- **Profiling & Auto-EDA:** pandas_profiling, sweetviz, dtale

- **Statistical Tests:** scipy.stats, statsmodels

**Tools:**

- Pandas, NumPy
- **Visualization:** Matplotlib, Seaborn, Plotly, Altair
- **Statistical analysis:** SciPy, statsmodels
- Sweetviz, D-Tale (automated EDA)
- Jupyter notebooks for interactive analysis

**Interview Tips:**

1. How do you detect and handle multicollinearity during EDA?

2. What would you do if your target variable is heavily imbalanced?

3. How can EDA guide feature engineering decisions?

4. How do you identify if feature transformation (like log-scaling) is needed?

# 6. Feature Engineering

**Purpose:**
Transform raw data into meaningful, machine readable features that improve model accuracy, robustness, and interpretability. Feature engineering is usually the most impactful step in the ML pipeline, good features can make simple models powerful, while poor features can break even the most advanced ones.

**Sub-Steps & Activities:**

**Feature Creation (Derived & Domain Features):**

- **Derived Features:**

    o Calculate new variables from existing ones (e.g., age from birthdate, days_since_signup, duration = end_date - start_date).

- **Time-based Features:**

    o Extract hour, day_of_week, month, quarter, or season from timestamps.

    o Create lag, rolling mean, or cumulative sum features for time series.

- **Interaction Features:**

    o Create interaction terms such as the product or ratio of two features (price_per_unit = price / quantity).

    o Useful for non-linear relationships that linear models may miss.

- **Aggregations:**

    o Summaries like total_spent, mean_transaction_value, or count_per_user.

- **Domain-Specific Features:**

    o Use business knowledge (e.g., is_premium_user, churn_risk_score, days_since_last_purchase) to encode hidden patterns.

**Feature Transformation:**

- **Encoding Categorical Variables:**
  - **One-Hot Encoding:** For nominal categories (e.g., city names).
  - **Label Encoding:** For ordinal categories (e.g., education levels).
  - **Target Encoding:** Replace categories with aggregated target statistics (use cross-validation to avoid leakage).

- **Scaling & Normalization:**
  - **StandardScaler (Z-score):** Mean = 0, SD = 1 (for linear models, PCA, KNN).
  - **MinMaxScaler:** Scales to [0,1] range (for neural networks, distance-based models).

- **Log/Power Transformations:**
  - Reduce skewness and stabilize variance for highly skewed features.

- **Binning:**
  - Convert continuous features into categorical buckets (e.g., age groups, income ranges).
  - Helps models capture non-linear relationships.

- **Polynomial Features:**
  - Generate squared, cubic, or interaction terms to capture complex patterns.

---

**Feature Selection:**

- **Filter Methods:**
  - Remove low-variance features (VarianceThreshold).
  - Eliminate highly correlated features (e.g., |corr| > 0.9 to reduce multicollinearity).

- **Statistical Tests:**
  - Chi-square for categorical features, ANOVA for numerical features.

- **Wrapper/Embedded Methods:**
  - Use model-based feature importance (e.g., RandomForest.feature_importances_, Lasso, XGBoost).
  - Recursive Feature Elimination (RFE) to iteratively select the best subset.

---

**Text & Special Feature Processing:**

- **Text Data:**

  o Use TF-IDF, Bag-of-Words, or word embeddings (e.g., Word2Vec, BERT) to convert text into numeric features.

  o Extract linguistic features like word count, sentiment score, or named entities.

- **Datetime Features:**

  o Break timestamps into day, month, year, is_weekend, is_holiday.

  o Create cyclic encodings for periodic data (e.g., sine/cosine transforms for month or hour).

---

**Example Code:**

```
import pandas as pd

from sklearn.preprocessing import StandardScaler, OneHotEncoder

import numpy as np

# Feature creation

df['age'] = (pd.to_datetime('today') - pd.to_datetime(df['birthdate'])).dt.days // 365

df['price_per_unit'] = df['total_price'] / df['quantity']

# Encoding

df = pd.get_dummies(df, columns=['city'])

# Scaling

scaler = StandardScaler()

df[['income_scaled']] = scaler.fit_transform(df[['income']])

# Log transform

df['log_amount'] = np.log1p(df['amount'])
```

---

**Tools:**

- Scikit-learn (preprocessing, feature_extraction)
- Feature-engine library
- Featuretools (automated feature engineering)
- Category Encoders
- PySpark for distributed processing

- TPOT, AutoML libraries
- **Text:** NLTK, spaCy, TfidfVectorizer
- **Images:** OpenCV, PIL, scikit-image

---

**Interview Tips:**

1. How would you decide between one-hot encoding and target encoding?

2. Why is feature scaling important for certain algorithms?

3. How do you handle high-cardinality categorical variables?

4. Describe a time when a domain-specific feature significantly improved model performance.

---

# 7. Data Splitting

**Purpose:**
Separate data into training, validation, and test sets so models are trained fairly, tuned without bias, and evaluated on truly unseen data. Proper splitting prevents leakage and gives accurate generalization estimates.

---

**Sub-Steps & Activities:**

- **Train / Test Split:** allocate ~70–80% training and 20–30% testing for a basic evaluation.

- **Validation Set (3-way split):** carve out a validation set for hyperparameter tuning (common: 60/20/20).

- **Stratification:** preserve class proportions for classification problems (stratify=y) to avoid skewed splits.

- **Time-based Splitting:** for time series, split chronologically (train on past, test on future); never shuffle time-ordered data.

- **Cross-Validation:** set up KFold, StratifiedKFold, or TimeSeriesSplit to get robust performance estimates.

- **Grouped Splits:** use GroupKFold when samples are correlated by group (e.g., multiple rows per user) to avoid leakage across folds.

- **Pipelines & Preprocessing Placement:** build preprocessing inside a Pipeline/ColumnTransformer to avoid leaking test information.

- **Reproducibility:** fix random_state and document the split approach for traceability.

---

**Important Considerations:**

- Always guard against **data leakage** (no peeking at test data during preprocessing or feature selection).

- Ensure **sufficient samples** per split (and per class for stratified splits).

- For imbalanced targets, consider stratification or specialized CV (repeated stratified CV) and appropriate metrics.

- When using time-based features, ensure feature creation uses only past data (lag features created from training windows).

---

**Example Code:**

**# Basic train-test split**

```
from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(

    X, y, test_size=0.2, random_state=42

)
```

**# 3-way split: train / val / test (60/20/20)**

```
X_tr, X_temp, y_tr, y_temp = train_test_split(X, y, test_size=0.4, random_state=42)

X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.5, random_state=42)

# Stratified split for classification

X_train, X_test, y_train, y_test = train_test_split(

    X, y, test_size=0.2, stratify=y, random_state=42

)
```

**# Time-based split (assumes df sorted by timestamp)**

```
df = df.sort_values('event_time')

train_size = int(len(df) * 0.8)

train_df = df.iloc[:train_size]

test_df  = df.iloc[train_size:]
```

**# Cross-validation examples and grouped CV**

```
from sklearn.model_selection import KFold, StratifiedKFold, GroupKFold, TimeSeriesSplit

kf = KFold(n_splits=5, shuffle=True, random_state=42)

skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
```

```
gkf = GroupKFold(n_splits=5)

ts  = TimeSeriesSplit(n_splits=5)
```

**# GroupKFold usage**

```
for train_idx, test_idx in gkf.split(X, y, groups=group_ids):

    X_tr, X_te = X[train_idx], X[test_idx]
```

**# Use a pipeline to avoid leakage: preprocess inside pipeline, then cross-validate**

```
from sklearn.pipeline import Pipeline

from sklearn.compose import ColumnTransformer

from sklearn.preprocessing import StandardScaler, OneHotEncoder

from sklearn.ensemble import RandomForestClassifier

from sklearn.model_selection import cross_val_score, StratifiedKFold

numeric_cols = ['age', 'income']

cat_cols = ['city']

preprocessor = ColumnTransformer([

    ('num', StandardScaler(), numeric_cols),

    ('cat', OneHotEncoder(handle_unknown='ignore'), cat_cols)

])


pipeline = Pipeline([

    ('pre', preprocessor),

    ('clf', RandomForestClassifier(random_state=42))

])

skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)

scores = cross_val_score(pipeline, X, y, cv=skf, scoring='roc_auc')

print(scores.mean(), scores.std())
```

---

**Tools:**

- **Scikit-learn:** train_test_split, StratifiedKFold, TimeSeriesSplit
- Pandas for manual splitting
- Imbalanced-learn for stratified splits

- Custom scripts for temporal/spatial splits

---

**Interview Tips:**

1. Explain how you avoid data leakage (e.g., all feature engineering and scaling inside pipelines).

2. Be prepared to justify stratification and when you'd use TimeSeriesSplit instead of random CV.

3. If asked about grouped data, mention GroupKFold and why splitting by group matters (prevents overly optimistic results).

4. Talk about reproducibility (random_state) and how you'd document the split strategy for production/review.

---

# 8. Model Training

**Purpose:**
Train machine learning models on prepared data to capture underlying patterns and relationships that can generalize to unseen data. This step is where the algorithm learns from the training set and becomes capable of making predictions.

---

**Sub-Steps & Activities:**

**Model Selection:**

- Choose the algorithm based on problem type and dataset characteristics:

    - **Classification:** Logistic Regression, Decision Trees, Random Forest, Gradient Boosting (XGBoost, LightGBM), Neural Networks.

    - **Regression:** Linear Regression, Ridge, Lasso, Random Forest Regressor, Gradient Boosting Regressor.

    - **Time Series / Sequence:** ARIMA, LSTM, Transformer models.

- Consider trade-offs: interpretability vs. accuracy, training time vs. scalability.

**Training Process:**

- Initialize the model with default or tuned hyperparameters.

- Fit the model on training data using model.fit(X_train, y_train).

- Track **training metrics** (accuracy, precision, recall, loss, RMSE, AUC, etc.) to detect convergence or underfitting/overfitting.

- If using deep learning, monitor epochs, learning rate schedules, and early stopping criteria.

---

**Hyperparameter Tuning:**

- Perform **Grid Search**, **Random Search**, or **Bayesian Optimization** for optimal hyperparameters.

- Use **cross-validation** to ensure model robustness.

- Adjust regularization parameters (like C, alpha, lambda) to control overfitting.

**Multiple Models & Comparison:**

- Train multiple algorithms to benchmark performance.

- Use **ensemble techniques** (Bagging, Boosting, Stacking) for performance gains.

- Compare using validation set metrics or cross-validation scores.

**Advanced Considerations:**

- **Feature Importance / Explainability:** Use feature_importances_, SHAP, or LIME.

- **Reproducibility:** Set seeds and document hyperparameters.

- **Automation:** Use frameworks like AutoML (e.g., auto-sklearn, H2O.ai, Vertex AI AutoML).

---

**Example Code:**

**# Basic model training**

```
from sklearn.ensemble import RandomForestClassifier

from sklearn.metrics import accuracy_score

model = RandomForestClassifier(n_estimators=100, random_state=42)

model.fit(X_train, y_train)

y_pred = model.predict(X_test)

print("Accuracy:", accuracy_score(y_test, y_pred))
```

**# Hyperparameter tuning with GridSearchCV**

```
from sklearn.model_selection import GridSearchCV

param_grid = {

    "n_estimators": [100, 200],

    "max_depth": [None, 10, 20]

}

grid = GridSearchCV(RandomForestClassifier(random_state=42), param_grid, cv=5,
scoring='accuracy')
```

```python
grid.fit(X_train, y_train)

print("Best Parameters:", grid.best_params_)
```

**# Cross-validation example**

```python
from sklearn.model_selection import cross_val_score

from sklearn.linear_model import LogisticRegression

clf = LogisticRegression(max_iter=1000)

scores = cross_val_score(clf, X_train, y_train, cv=5, scoring='roc_auc')

print("CV AUC Mean:", scores.mean())
```

**# Training multiple models for comparison**

```python
from sklearn.ensemble import GradientBoostingClassifier

from xgboost import XGBClassifier

models = {

    "RandomForest": RandomForestClassifier(random_state=42),

    "GradientBoosting": GradientBoostingClassifier(),

    "XGBoost": XGBClassifier(eval_metric='logloss')

}

for name, m in models.items():

    m.fit(X_train, y_train)

    acc = accuracy_score(y_test, m.predict(X_test))

    print(f"{name} Accuracy: {acc:.4f}")
```

---

**Machine Learning Algorithms & Types:**

**Supervised Learning - Regression:**

- Linear Regression
- Ridge Regression, Lasso Regression
- ElasticNet
- Support Vector Regression (SVR)
- Decision Trees
- Random Forest
- Gradient Boosting (XGBoost, LightGBM, CatBoost)
- Neural Networks (MLP Regressor)

**Supervised Learning - Classification:**

- Logistic Regression
- Naive Bayes
- K-Nearest Neighbors (KNN)
- Support Vector Machines (SVM)
- Decision Trees
- Random Forest
- Gradient Boosting (XGBoost, LightGBM, CatBoost)
- Neural Networks (MLP Classifier, CNN, RNN)

**Unsupervised Learning:**

- K-Means Clustering
- DBSCAN
- Hierarchical Clustering
- Gaussian Mixture Models
- Principal Component Analysis (PCA)
- t-SNE, UMAP
- Isolation Forest (anomaly detection)
- Autoencoders

**Deep Learning:**

- Feedforward Neural Networks (FNN)
- Convolutional Neural Networks (CNN)
- Recurrent Neural Networks (RNN, LSTM, GRU)
- Transformers (BERT, GPT)
- Generative Adversarial Networks (GAN)

**Ensemble Methods:**

- Bagging (Bootstrap Aggregating)
- Boosting (AdaBoost, Gradient Boosting)
- Stacking
- Voting Classifiers/Regressors

---

**Tools for Model Training:**

- Scikit-learn
- XGBoost, LightGBM, CatBoost
- TensorFlow, Keras
- PyTorch
- PySpark MLlib
- H2O.ai
- AutoML tools: Auto-sklearn, TPOT, AutoKeras
- MLflow for experiment tracking
- Weights & Biases (wandb)
- Ray Tune for hyperparameter optimization

**Interview Tips:**

1. Be ready to justify your choice of algorithm & for few companies algorithm internals (e.g., why tree-based methods over linear models).

2. Explain how hyperparameter tuning affects model performance and prevents overfitting.

3. Discuss how you would handle class imbalance (e.g., weighting, SMOTE) or high-dimensional data (e.g., feature selection).

4. Mention ensemble methods if asked about improving baseline performance.

# 9. Evaluation

**Purpose:**
Evaluate how well a trained model performs on unseen data to ensure it generalizes beyond the training set. The goal is to validate predictive quality, reliability, and business relevance before deployment.

**Sub-Steps & Activities:**

**Performance Metrics:**

- **Classification:**

  o **Accuracy:** Overall correctness.

  o **Precision / Recall / F1-score:** Better for imbalanced classes.

  o **ROC-AUC:** Discrimination capability across thresholds.

  o **Confusion Matrix:** Breakdown of TP, FP, FN, TN for deeper insight.

- **Regression:**

  o **MSE / RMSE:** Penalize large errors.

  o **MAE:** Average absolute error (robust to outliers).

  o **$R^2$:** Proportion of variance explained.

  o **MAPE:** Error as a percentage, useful for forecasting tasks.

**Evaluation Methods:**

- Generate predictions on the **test set** (model.predict(X_test)) and compute metrics.

- Use **cross-validation scores** for a more robust estimate of model performance.

- Plot **learning curves** to detect underfitting (high bias) or overfitting (high variance).

**Model Interpretation:**

- **Feature Importance:** Understand which features influence predictions most.

- **SHAP / LIME:** Explain individual predictions and feature impact.

- **Partial Dependence Plots (PDP):** Visualize relationships between specific features and predictions.

**Error Analysis:**

- Review **misclassified samples** or **large residuals** to uncover patterns.

- Detect **systematic errors** (e.g., model bias against certain groups).

- Investigate if errors correlate with missing features or poor data quality.

**Model Comparison:**

- Evaluate multiple models on the **same test set** and compare metrics.

- Consider trade-offs like **accuracy vs. interpretability**, **latency vs. complexity**, or **precision vs. recall** based on the use case.

**Validation Against Business Goals:**

- Ensure that model performance meets **domain-specific KPIs** (e.g., <5% false negatives in fraud detection).

- Compare **train vs. test metrics** to confirm no overfitting before production deployment.

---

**Example Code:**

**# Classification evaluation**

from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, roc_auc_score, confusion_matrix

y_pred = model.predict(X_test)

print("Accuracy:", accuracy_score(y_test, y_pred))

print("Precision:", precision_score(y_test, y_pred))

print("Recall:", recall_score(y_test, y_pred))

print("F1:", f1_score(y_test, y_pred))

print("ROC-AUC:", roc_auc_score(y_test, model.predict_proba(X_test)[:, 1]))

print("Confusion Matrix:\n", confusion_matrix(y_test, y_pred))

**# Regression evaluation**

from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score

import numpy as np

```
y_pred = model.predict(X_test)

print("MSE:", mean_squared_error(y_test, y_pred))

print("RMSE:", np.sqrt(mean_squared_error(y_test, y_pred)))

print("MAE:", mean_absolute_error(y_test, y_pred))

print("R²:", r2_score(y_test, y_pred))
```

**# Cross-validation score**

```
from sklearn.model_selection import cross_val_score

scores = cross_val_score(model, X, y, cv=5, scoring='f1')

print("Cross-Validation F1 Mean:", scores.mean())
```

**# SHAP explainability example**

```
import shap

explainer = shap.TreeExplainer(model)

shap_values = explainer.shap_values(X_test)

shap.summary_plot(shap_values, X_test)
```

---

**Tools:**

- Scikit-learn metrics module
- TensorFlow/Keras evaluation functions
- Confusion Matrix visualizers
- **ROC/AUC curve plotters:** Matplotlib, Seaborn
- MLflow for tracking experiments
- SHAP, LIME for model interpretability
- **Cross-validation tools:** cross_val_score, GridSearchCV

---

**Metrics:**

- **Classification**: Accuracy, Precision, Recall, F1-Score, ROC-AUC
- **Regression**: MSE, RMSE, MAE, $R^2$, MAPE
- **Clustering**: Silhouette Score, Davies-Bouldin Index

---

**Interview Tips:**

1. Be ready to explain why you chose a particular metric (e.g., F1 over accuracy for imbalanced classes).

2. Show that you understand bias-variance trade-offs and how learning curves can diagnose them.

3. Talk about error analysis as an ongoing iterative process, not a one-time check.

4. If asked about explainability, discuss SHAP, LIME, or feature importance and how they guide decision-making.

# 10. Deployment & Monitoring

**Purpose:**
Deploy the trained model into production and ensure it maintains reliable performance over time. Continuous monitoring, retraining, and integration with business workflows are critical to sustain real-world value.

**Sub-Steps & Activities:**

**Model Serialization & Versioning:**

- **Save trained models:** Use pickle, joblib, or model-specific formats.

- **Version control:** Track model versions with **MLflow**, **DVC**, or other model registries.

- **Save preprocessing pipelines:** Ensure that transformers, scalers, and encoders are stored alongside the model for reproducible inference.

import joblib

**# Save model**

joblib.dump(model, "model_v1.pkl")

**# Load model**

model = joblib.load("model_v1.pkl")

**Deployment Options:**

- **Batch Predictions:** Run periodic predictions on new datasets.

- **Real-Time API:** Expose the model through REST APIs using **Flask** or **FastAPI**.

- **Cloud Deployment:** Deploy to AWS SageMaker, Azure ML, GCP AI Platform, or other ML services.

- **Containerization:** Use **Docker** (and optionally **Kubernetes**) for environment consistency and scalability.

**# FastAPI example**

```python
from fastapi import FastAPI

import joblib

import pandas as pd


app = FastAPI()

model = joblib.load("model_v1.pkl")


@app.post("/predict")

def predict(data: dict):

    df = pd.DataFrame([data])

    pred = model.predict(df)

    return {"prediction": pred.tolist()}
```

---

**Integration:**

- Connect with **data pipelines** and **databases** to automatically ingest new data.

- Integrate predictions into **business applications** or dashboards.

- Automate **retraining workflows** with new data to prevent model staleness.

---

**Monitoring:**

- **Performance Tracking:** Monitor metrics like accuracy, F1-score, latency, and throughput in production.

- **Data Drift Detection:** Check if input data distributions have changed compared to training data.

- **Model Drift Detection:** Detect degradation in predictions over time.

- **Logging:** Maintain detailed logs for predictions, errors, and system health.

**# Example: basic monitoring metrics**

```python
preds = model.predict(X_test)

accuracy = (preds == y_test).mean()

print("Production Accuracy:", accuracy)
```

---

**Maintenance & Updates:**

- Schedule **periodic retraining** with fresh data.

- Update **features** or preprocessing as business requirements evolve.

- Use **A/B testing** to compare new models or features before full rollout.

- Define **incident response procedures** to handle prediction failures or system issues.

---

**Tools:**

- **Model serving:** Flask, FastAPI, Django
- **Container tools:** Docker, Kubernetes
- **Cloud platforms:** AWS SageMaker, Azure ML, Google Vertex AI
- **MLOps:** MLflow, Kubeflow, Airflow
- **Monitoring**: Prometheus, Grafana, Evidently AI
- **Model versioning**: DVC, MLflow Model Registry
- **CI/CD**: Jenkins, GitHub Actions, GitLab CI
- **API gateways**: Kong, AWS API Gateway
- **Feature stores**: Feast, Tecton

---

**Interview Tips:**

1. Explain trade-offs between batch and real-time predictions.

2. Discuss how you monitor model performance and detect drift in production.

3. Be ready to explain version control for models and retraining strategies.

4. Highlight scalability considerations (containers, cloud deployment, load balancing).

---

# 11. Summary

This machine learning pipeline is iterative - you usually cycle back to earlier stages based on findings or changes in requirements. Key feedback loops include:

- **Poor evaluation results →** revisit feature engineering, try different models, tune hyperparameters, or collect additional data.

- **Data quality issues discovered during EDA →** return to data collection, integration, or initial validation to clean, enrich, or correct the dataset.

- **Imbalanced or insufficient data →** adjust data splitting, apply resampling, or acquire more samples.

- **Deployment monitoring shows drift →** trigger model retraining, update features, or collect fresh labeled data.

- **Business requirement changes →** revisit feature engineering, model selection, or evaluation metrics to align with new objectives.

**Key Principles:**

- Maintain flexibility to iterate while following a structured approach.

- Ensure reproducibility with versioned models, documented pipelines, and controlled random seeds.

- Continuously monitor data and model performance to catch drift or degradation early.

- Combine technical rigor with business awareness for reliable, production-ready ML solutions.

# Happy Learning 😊

.