# All Rights Reserved. Author @ Rajendra Phani

| Stage | Step | Sub-steps / Key actions |
|---|---|---|
| 1. Problem Definition | 1 Clarify objective | • Business question / KPI <br>• Success metric (accuracy, AUC, RMSE, …) <br>• Stakeholder expectations |
| | 2 Choose ML paradigm | • Supervised / Unsupervised / Reinforcement <br>• Classification, Regression, Clustering, Anomaly detection, etc. |
| | 3 Define data & model constraints | • Data availability <br>• Privacy / regulatory limits <br>• Deployment latency / throughput |
| 2. Data Acquisition | 1 Identify sources | • Internal (databases, logs, CRM) <br>• External APIs (social media, weather, finance) <br>• Public datasets (UCI, Kaggle, government portals) <br>• Web-scraping / IoT sensors |
| | 2 Collect & ingest | • SQL/NoSQL queries <br>• REST / GraphQL API calls <br>• Streaming pipelines (Kafka, Flink) <br>• Scraping tools (Scrapy, BeautifulSoup) |
| | 3 Store & version | • Data lake / warehouse (S3, GCS, Snowflake) <br>• Schema registry / data catalog <br>• Version control (Delta Lake, DVC) |
| 3. Data Exploration & Profiling | 1 Load data into analytical environment | • Pandas, Spark, Dask <br>• Jupyter / Colab notebooks |
| | 2 Basic statistics | • Summary (mean, std, min, max) <br>• Distribution plots <br>• Correlation matrix |
| | 3 Detect anomalies & outliers | • Boxplots, IQR, Z-score <br>• Visual inspection |
| | 4 Data quality assessment | • Missingness pattern <br>• Duplicate records <br>• Inconsistent data types |
| 4. Data Cleaning & Pre-processing | 1 Handle missing values | • Drop / impute (mean, median, mode) <br>• Model-based imputation (KNN, MICE) <br>• Indicator columns |
| | 2 Remove duplicates / errors | • Identify unique keys <br>• De-duplicate rows |
| | 3 Correct data types & formats | • Convert strings → dates, categoricals <br>• Normalize numeric ranges (min-max, z-score) |
| | 4 Resolve inconsistencies | • Standardize categorical labels (e.g., "NY", "New York") <br>• Harmonise units |
| | 5 Outlier treatment | • Winsorize, clip, or remove outliers |
| 5. Feature Engineering | 1 Domain-specific transforms | • Create ratios, differences, moving averages <br>• Time-series lag features |
| | 2 Interaction terms | • Polynomial, cross-products <br>• Feature crosses for tree models |
| | 3 Aggregations & embeddings | • Group-by summaries <br>• Text embeddings (BERT, word2vec) |
| | 4 Dimensionality reduction | • PCA / t-SNE for exploratory analysis <br>• Autoencoders if needed |
| 6. Feature Selection / Extraction | 1 Filter methods | • Correlation threshold <br>• Mutual information, chi-square |
| | 2 Wrapper methods | • Recursive Feature Elimination (RFE) <br>• Forward/Backward selection |
| | 3 Embedded methods | • L1-regularization (Lasso) <br>• Tree-based feature importance |
| | 4 Evaluate impact | • Cross-validated performance vs. number of features |
| 7. Encoding & Representation | 1 Categorical encoding | • One-hot / dummy <br>• Ordinal encoding <br>• Target encoding (with smoothing) |
| | 2 Text representation | • TF-IDF, bag-of-words <br>• Word embeddings (Word2Vec, GloVe) <br>• Sentence transformers |
| | 3 Image / Audio preprocessing | • Resizing, normalization <br>• Feature extraction with CNNs |
| 8. Train-Test Split & Validation Strategy | 1 Basic split | • Random train/validation/test (70/15/15) |
| | 2 Time-series split | • Rolling window, expanding window |
| | 3 Stratified sampling | • Preserve class distribution (classification) |
| | 4 Cross-validation | • K-fold, StratifiedKFold <br>• GroupKFold for grouped data |
| | 5 Validation set for hyper-parameter tuning | • GridSearchCV / RandomizedSearchCV <br>• Bayesian optimization (Optuna, Hyperopt) |
| 9. Model Training | 1 Baseline models | • Logistic regression, linear regression <br>• Decision trees |
| | 2 Advanced algorithms | • Ensemble (RandomForest, XGBoost, LightGBM) <br>• Neural networks (MLP, CNN, RNN, Transformer) |
| | 3 Training pipeline | • Data loader, batching <br>• Early stopping, learning-rate scheduling |
| | 4 Reproducibility | • Set random seeds <br>• Log hyper-parameters, code version |
| 10. Model Evaluation | 1 Choose metrics | • Accuracy, Precision/Recall, F1 (classification) <br>• RMSE, MAE (regression) <br>• ROC-AUC, PR-AUC |
| | 2 Confusion matrix & error analysis | • Identify bias, misclassifications |
| | 3 Calibration (classification) | • Platt scaling, isotonic regression |
| | 4 Statistical tests | • McNemar's test for comparing classifiers <br>• Paired t-test |
| | 5 Interpretability | • SHAP / LIME <br>• Feature importance plots |
| 11. Model Selection & Finalization | 1 Compare pipelines | • Cross-validated scores, stability <br>• Complexity vs. performance |
| | 2 Pick best model & hyper-parameters | |
| | 3 Retrain on full training + validation set | |
| | 4 Save artifact (pickle, joblib, ONNX) | |
| 12. Deployment | 1 Packaging & serialization | • Docker image <br>• Model server (FastAPI, Flask, TorchServe) |
| | 2 Serving infrastructure | • REST API <br>• gRPC <br>• Serverless (AWS Lambda, GCP Cloud Functions) |
| | 3 Scaling & load balancing | • Kubernetes / ECS <br>• Auto-scaling |
| | 4 Monitoring & logging | • Prediction latency <br>• Drift detection (concept drift, data drift) <br>• Error rate |
| | 5 A/B testing / Canary releases | • Gradual rollout, traffic splitting |
| 13. Maintenance & Continuous Learning | 1 Retraining schedule | • Periodic retrain (weekly, monthly) <br>• Triggered by drift |
| | 2 Data pipeline updates | • Add new features, fix bugs |
| | 3 Model governance | • Version control (MLflow), lineage tracking <br>• Auditing & compliance |
| | 4 Feedback loop | • Capture user feedback, label corrections <br>• Active learning if applicable |
| 14. Documentation & Communication | 1 Technical docs | • Data schema, feature definitions <br>• Model card (performance, limitations) |
| | 2 Business summary | • KPI impact, ROI <br>• Recommendations for stakeholders |
| | 3 Code & notebook notebooks | • Clean, reproducible scripts <br>• Versioned on Git |