

1. Problem Definition

- Clearly define the problem statement and the prediction target (dependent variable).
 - Identify the type of task: regression, classification, clustering, etc.
 - Understand the business context and goals, including success criteria and performance metrics.
-

2. Data Collection

- Gather data from relevant sources (databases, APIs, or raw files).
 - Ensure data quality and relevance to the problem.
 - Understand the data format, volume, and any potential privacy or compliance constraints.
-

3. Data Understanding and Exploration

- Study the structure and summary statistics of the data.
 - Perform exploratory data analysis (EDA) to identify patterns, distributions, and relationships.
 - Use visualization techniques to understand trends and detect anomalies or correlations.
-

4. Data Cleaning

- Handle missing values through imputation or removal.
 - Remove duplicate entries and irrelevant features.
 - Detect and address outliers that may distort model performance.
 - Standardize or normalize numerical data for consistency.
-

5. Feature Engineering

- Create meaningful features from raw data, such as aggregations or domain-specific transformations.
 - Encode categorical variables into numerical representations.
 - Reduce dimensionality if the feature set is too large or sparse.
 - Assess feature importance and select the most relevant ones.
-

6. Data Splitting

- Split the data into training, validation, and testing sets (e.g., 70-20-10 split).
 - Ensure the splits are representative of the overall data distribution.
 - Use stratified sampling for imbalanced datasets.
-

7. Model Selection

- Choose algorithms suitable for the problem type and dataset size (e.g., linear models, tree-based models, or neural networks).
 - Start with simple baseline models to understand performance benchmarks.
-

8. Model Training

- Train models on the training dataset using appropriate algorithms.
 - Experiment with different model architectures and configurations.
 - Document all parameters and settings for reproducibility.
-

9. Model Evaluation

- Evaluate the model on the validation set using appropriate metrics:
 - Regression: RMSE, MAE, R^2 .
 - Classification: Accuracy, Precision, Recall, F1 Score, ROC-AUC.
 - Compare results against baseline and check for underfitting or overfitting.
-

10. Hyperparameter Tuning

- Optimize model performance by fine-tuning hyperparameters (e.g., learning rates, depth, regularization).
 - Use systematic approaches like grid search, random search, or Bayesian optimization.
-

11. Final Testing

- Test the final model on the unseen test dataset to ensure generalization.
 - Compare predictions with actual values and analyze discrepancies.
-

12. Model Deployment

- Save the trained model for deployment using appropriate tools.
 - Integrate the model into a production environment (web app, API, or embedded system).
 - Ensure the system handles real-time or batch predictions effectively.
-

13. Monitoring and Maintenance

- Monitor model performance post-deployment (e.g., accuracy drift, data drift).
- Retrain or update the model periodically with new data.
- Collect user feedback and adjust the model or features as necessary.