# Data Science Workflow - Step-by-Step Process

## a) Data Sourcing (Collect Data)

This is the first stage of any data science project. The data gathered that will be used for training and testing your model.

Sources may include:
- Public datasets (e.g., Kaggle, UCI Machine Learning Repository)
- Internal company databases
- APIs or web scraping
- Sensors, logs, or user interaction data

Goal: Get all relevant and reliable data in one place.

## b) Pre-processing (Clean Data + Explore & Visualize)

Once data is collected, it's rarely ready for use directly. Preprocessing includes both cleaning and understanding your data.

Key tasks:
- Handle missing values (fill, drop, or impute)
- Fix inconsistent entries (e.g., standardizing dates)
- Remove duplicates
- Convert data types (e.g., categorical to numerical)
- Normalize or scale features
- Visualize relationships (e.g., correlations, class distribution)
- Identify and address imbalances (e.g., in classification)

1. Handle Missing Values:
   - df.isnull().sum(): Checks how many missing values are present in each column.
   - df.dropna(): Removes rows with any missing values.
   - df.fillna(value): Fills missing values with a given constant or method (e.g., mean).
   - SimpleImputer(strategy='mean'): Imputes missing values using a specified strategy (mean, median, most_frequent).

2. Fix Inconsistent Entries:
   - pd.to_datetime(): Converts date strings into datetime objects.
   - df['column'].str.lower().str.strip(): Converts strings to lowercase and removes whitespace for uniform formatting.

3. Remove Duplicates:

- df.duplicated(): Identifies duplicate rows.
  - df.drop_duplicates(): Removes duplicates from the dataset.

4. Convert Data Types (Categorical to Numerical):
  - pd.get_dummies(): Performs one-hot encoding for categorical variables.
  - LabelEncoder(): Encodes labels with values between 0 and n_classes-1.
  - OneHotEncoder(): Converts categorical variable(s) into dummy/indicator variables.

5. Normalize or Scale Features:
  - StandardScaler(): Standardizes features by removing the mean and scaling to unit variance.
  - MinMaxScaler(): Scales and translates each feature to a range between 0 and 1.

6. Visualize Relationships:
  - df.corr(): Shows correlation between numerical variables.
  - sns.heatmap(): Plots a heatmap to show correlation matrix.
  - sns.countplot(): Plots the count of categorical values.
  - sns.pairplot(): Displays pairwise relationships in the dataset.
  - plt.hist(): Plots histogram of variable distribution.

7. Identify and Address Imbalances:
  - df['label'].value_counts(): Shows the distribution of labels/classes.
  - SMOTE(): From imblearn, used to oversample minority classes.

Goal: Ensure the data is clean, structured, and suitable for model input.

## c) Data Splitting (Train/Test/Validation sets)
To train the model fairly and evaluate its performance honestly, data is split into:
- Training Set: Used to train the model.
- Validation Set: Used during training to tune hyperparameters (optional).
- Test Set: Used to evaluate final model performance.

Typical splits: 70% training, 15% validation, 15% testing – but this varies depending on dataset size.

This is a built-in function in the scikit-learn library used to split arrays or matrices into random train and test subsets.

Example Syntax:
  train_test_split(X, y, test_size=0.2, random_state=42)

Explanation of Parameters:
- X: Feature dataset (e.g., input variables).
- y: Target variable (e.g., class labels or regression targets).
- test_size=0.2: This sets aside 20% of the data for testing, leaving 80% for training.
- random_state=42: This is a seed used by the random number generator to make the split reproducible. Any fixed number ensures that the split will always be the same each time the code runs.

Purpose:
This method helps ensure the model is trained on one part of the data and evaluated on another, preventing information leakage.

Manual Three-Way Split (Train/Validation/Test)

Many workflows require three sets: training, validation, and test.

Steps:
- First, split the data into train and temp (e.g., 80% train, 20% temp).
- Then, split the temp set again into validation and test (e.g., 10% each).

Example:
```
X_train, X_temp, y_train, y_temp = train_test_split(X, y, test_size=0.2, 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)
```

Purpose:
The training set is used to fit the model, the validation set is used for tuning and checking performance during training, and the test set is only used at the end to evaluate how well the model generalizes.

Stratified Splits

When dealing with classification tasks, especially with imbalanced classes, it's important that each split maintains the same class distribution.

Example Syntax:
```
train_test_split(X, y, test_size=0.2, stratify=y, random_state=42)
```

Explanation:
- stratify=y: Ensures that the proportion of classes in the target variable y is preserved in both the training and testing datasets.

Purpose:
Prevents skewed distributions that can bias the model. It is especially useful for datasets where some classes are underrepresented.

Goal of Data Splitting

The primary objective of splitting data is to prevent overfitting and assess how well the model performs on unseen data. It ensures that the model doesn't just memorize the training data but learns patterns that generalize to new inputs.

d) Model Training (Train model)

In this step, a machine learning algorithm applied to the training set. This involve:
- Selecting an algorithm (e.g., logistic regression, decision tree, neural network)
- Defining a model architecture (for deep learning: layers, units, activation)
- Compiling the model with a loss function and optimizer
- Running the training process (e.g., epochs, batch size)

Goal: Teach the model to map inputs to outputs correctly using the training data.

## e) Model Evaluation (Evaluate model)

After training, the test data used to measure how well the model performs on unseen data.

Key metrics depend on the task:
- Accuracy, Precision, Recall, F1-score (for classification)
- Mean Squared Error (for regression)
- Confusion Matrix (to analyze class-wise performance)

Goal: Understand if the model meets performance expectations and where it may be weak.

The evaluation process involves comparing the model's predictions with the actual outcomes in the test set using appropriate metrics.

Evaluation Metrics for Classification Tasks

1. Accuracy:
   - Definition: The ratio of correctly predicted instances to the total instances.
   - Formula: (TP + TN) / (TP + TN + FP + FN)

- Use: Best when classes are balanced.

2. Precision:
  - Definition: The proportion of positive identifications that were actually correct.
  - Formula: TP / (TP + FP)
  - Use: Important when the cost of false positives is high (e.g., spam detection).

3. Recall (Sensitivity or True Positive Rate):
  - Definition: The proportion of actual positives that were correctly identified.
  - Formula: TP / (TP + FN)
  - Use: Important when the cost of false negatives is high (e.g., medical diagnostics).

4. F1-Score:
  - Definition: Harmonic mean of precision and recall.
  - Formula: 2 * (Precision * Recall) / (Precision + Recall)
  - Use: Best used when there is a need to balance precision and recall.

5. Confusion Matrix:
  - Description: A table layout that shows the number of correct and incorrect predictions made by the classifier, broken down by each class.
  - Components: True Positives (TP), True Negatives (TN), False Positives (FP), False Negatives (FN)
  - Use: Provides a comprehensive view of class-level performance.


Evaluation Metrics for Regression Tasks


1. Mean Squared Error (MSE):
  - Definition: The average of the squares of the differences between actual and predicted values.
  - Formula: $(1/n) * \Sigma(\text{actual} - \text{predicted})^2$
  - Use: Penalizes larger errors more than smaller ones. Useful for measuring model accuracy in regression tasks.

2. Mean Absolute Error (MAE):
  - Definition: The average of the absolute differences between actual and predicted values.
  - Formula: $(1/n) * \Sigma|\text{actual} - \text{predicted}|$
  - Use: Easier to interpret but does not penalize large errors as much as MSE.

3. R-squared (Coefficient of Determination):
  - Definition: Indicates the proportion of variance in the dependent variable that is predictable from the independent variable(s).
  - Value Range: 0 to 1 (closer to 1 means better fit)

- Use: Gives an overall measure of how well the model explains the variation in the data.


Conclusion


The choice of evaluation metric depends on the problem type (classification or regression) and business context. Understanding these metrics helps determine how reliable and effective the model is, and what adjustments may be needed before deployment.


## f) Results Interpretation
The final answers:
- Did the model solve the problem defined?
- What insights can be extracted from the model output?
- Are predictions explainable?
- Can this be used to make decisions?

Goal: Translate the technical results into actionable conclusions or business insights.


 Did the model solve the problem defined?


This is evaluated by comparing the model's performance metrics against predefined objectives or benchmarks.
- If accuracy, precision, recall, or other chosen metrics meet or exceed the expectations, it suggests that the model has solved the problem effectively.
- In regression, low MSE or high R-squared values indicate a well-fitting model.
- Performance on the test set (unseen data) gives the best estimate of real-world effectiveness.
- Comparison with baseline models (e.g., random guessing or previous solutions) can validate improvement.

Conclusion: The model has successfully solved the problem if it generalizes well and meets the defined performance goals.


 What insights can be extracted from the model output?

Beyond accuracy, models often reveal patterns in data:
- Feature importance scores can show which variables most influence predictions.
- Confusion matrices help identify which classes are commonly misclassified.
- In regression, residual plots can show trends or anomalies in prediction errors.
- Misclassifications or outliers may point to data quality issues or edge cases.

Insight extraction helps in understanding the data better and improving future decision-making or model versions.

Are predictions explainable?

Model interpretability depends on the model type:
- Linear models (e.g., Logistic Regression) are inherently interpretable.
- Tree-based models (e.g., Decision Trees, Random Forests) provide feature importance scores and decision paths.
- Complex models (e.g., Neural Networks) require additional tools like SHAP, LIME, or Grad-CAM for visual or local interpretability.
- Transparency is essential when model predictions affect regulated or sensitive decisions (e.g., loan approvals, healthcare).

Interpretability builds trust with users and stakeholders.

Can this be used to make decisions?

A model is decision-ready when:
- It performs consistently well on new, unseen data.
- Its predictions align with domain knowledge and business goals.
- The output integrates easily into business systems (via APIs, dashboards, etc.).
- The risk of error is acceptable for the business use case.

For example, a model that predicts customer churn can drive retention campaigns, while a demand forecast model can inform inventory planning.