



# DATA SCIENCE GLOSSARY

The Essential Language of Data  
for Aspiring Professionals.

# DATA SCIENCE WORKFLOW: FROM QUESTION TO IMPACT



## 1. DISCOVERY & PROBLEM DEFINITION

- Understand the research or business goal.
- Ask "What problem are we solving?"

## 2. DATA COLLECTION

- Gather relevant data from sources
- Ensure quality & consistency

## 3. DATA CLEANING & PREPARATION

- Address missing values, duplicates, and errors
- Format for analysis

## 4. EXPLORATORY DATA ANALYSIS (EDA)

- Identify trends and patterns visually
- Summarize key findings

## 5. MODELING & ANALYSIS

- Apply machine learning or statistical methods
- Make predictions or classifications

## 6. INTERPRETATION & COMMUNICATION

- Present results using visual analytics
- Deliver actionable insights for stakeholders

## Key Terms & Definitions

This glossary is designed as a companion for your Data Science learning journey. It provides concise, accessible definitions of the core terms and concepts you will encounter throughout the book and in your wider study of data science. Entries are arranged alphabetically for ease of reference. Whether you are encountering a term for the first time or need a quick reminder, this glossary will support your learning journey.

### A

Term	Definition
<b>Algorithm</b>	A step-by-step procedure or set of rules for solving a problem or performing a computation. In data science, algorithms are used to process data, identify patterns, and make predictions.
<b>Anomaly Detection</b>	The process of identifying unusual patterns or outliers in data that deviate significantly from expected behaviour. Used in fraud detection, network security, and quality control.
<b>API (Application Programming Interface)</b>	A set of rules and protocols that allows different software applications to communicate with each other. Data scientists use APIs to access external data sources and services.

Term	Definition
<b>Artificial Intelligence (AI)</b>	The simulation of human intelligence processes by machines, including learning, reasoning, and self-correction. Data science is a key discipline within the broader field of AI.
<b>Autocorrelation</b>	The correlation of a time series with a delayed copy of itself. It measures the degree to which past values of a variable predict its future values.

## B

Term	Definition
<b>Bagging (Bootstrap Aggregating)</b>	An ensemble machine learning technique that trains multiple models on different random subsets of the training data and combines their predictions to reduce variance and overfitting.
<b>Bayesian Inference</b>	A statistical method that updates the probability of a hypothesis as new evidence becomes available, combining prior beliefs with observed data.
<b>Bias</b>	A systematic error in a model or dataset that causes predictions to consistently deviate from true values. Bias can arise from flawed data collection, model assumptions, or algorithmic design.
<b>Big Data</b>	Datasets that are too large or complex to be processed by traditional data management tools. Characterised by the 3 Vs: Volume (scale), Velocity (speed of generation), and Variety (diversity of types).
<b>Binary Classification</b>	A supervised learning task where the goal is to categorise input data into one of exactly two classes, such as spam vs. not spam or positive vs. negative.
<b>Boosting</b>	An ensemble method that sequentially trains weak learners, each correcting the errors of its predecessor, to build a strong predictive model. Examples include AdaBoost and Gradient Boosting.

## C

Term	Definition
<b>Categorical Variable</b>	A variable that takes on a limited number of distinct values representing categories or groups (e.g., gender, colour, country). Often encoded numerically before use in machine learning models.
<b>Central Limit Theorem (CLT)</b>	A fundamental statistical theorem stating that the distribution of sample means approaches a normal distribution as the sample size grows, regardless of the population's original distribution.

Term	Definition
<b>Classification</b>	A supervised learning task in which a model learns to assign input data to predefined categories or labels. Common algorithms include logistic regression, decision trees, and support vector machines.
<b>Clustering</b>	An unsupervised learning technique that groups similar data points together based on their characteristics, without using predefined labels. Common algorithms include K-Means and DBSCAN.
<b>Confusion Matrix</b>	A table used to evaluate the performance of a classification model by showing the counts of true positives, true negatives, false positives, and false negatives.
<b>Correlation</b>	A statistical measure that expresses the extent to which two variables are linearly related. Values range from -1 (perfect negative) to +1 (perfect positive), with 0 indicating no linear relationship.
<b>Cross-Validation</b>	A model evaluation technique that splits data into multiple subsets (folds), training the model on some folds and validating on others, to give a more robust estimate of model performance.

## D

Term	Definition
<b>Data Cleaning</b>	The process of detecting and correcting errors, inconsistencies, and missing values in a dataset to improve data quality before analysis.
<b>Data Lake</b>	A centralised repository that stores large volumes of raw, structured, semi-structured, and unstructured data in its native format until it is needed for analysis.
<b>Data Mining</b>	The process of discovering patterns, correlations, and insights from large datasets using statistical, mathematical, and computational techniques.
<b>Data Pipeline</b>	A series of automated data processing steps that move and transform data from one system to another, typically from raw sources to analysis-ready formats.
<b>Data Wrangling</b>	The process of transforming and mapping raw data into a more usable format for analysis. Also called data munging, it includes cleaning, restructuring, and enriching data.
<b>Decision Tree</b>	A supervised learning model that makes predictions by learning a hierarchy of if-then rules from training data. The model resembles a tree structure with branches representing decision points.
<b>Dimensionality Reduction</b>	The process of reducing the number of features in a dataset while preserving important information. Common techniques include Principal Component Analysis (PCA) and t-SNE.

## E

Term	Definition
<b>Ensemble Method</b>	A machine learning approach that combines multiple models to produce better predictive performance than any individual model. Examples include Random Forests, Bagging, and Boosting.
<b>Epoch</b>	One complete pass through the entire training dataset during the training of a neural network. Models are typically trained over many epochs to converge on optimal weights.
<b>Exploratory Data Analysis (EDA)</b>	An approach to analysing datasets that uses visual and statistical methods to summarise their main characteristics, discover patterns, and identify anomalies before formal modelling.

## F

Term	Definition
<b>Feature</b>	An individual measurable property or characteristic of the data used as input to a machine learning model. Also known as an independent variable or predictor.
<b>Feature Engineering</b>	The process of using domain knowledge to create, transform, or select features from raw data to improve model performance.
<b>Feature Scaling</b>	The process of normalising or standardising the range of features in a dataset so that no single feature dominates others due to its scale. Common methods include min-max scaling and standardisation.
<b>F1 Score</b>	A classification metric that balances precision and recall by computing their harmonic mean. Useful when class imbalance is present and both false positives and false negatives are costly.

## G

Term	Definition
<b>Gradient Descent</b>	An optimisation algorithm that iteratively adjusts model parameters by moving in the direction of the steepest decrease in the loss function, aiming to find the minimum loss.
<b>Gradient Boosting</b>	An ensemble technique that builds models sequentially, with each model correcting the residual errors of the previous one using gradient descent in function space.

# H

Term	Definition
<b>Hyperparameter</b>	A configuration setting external to a model that is set before training begins and controls the learning process (e.g., learning rate, number of trees). Distinct from parameters, which are learned from data.
<b>Hypothesis Testing</b>	A statistical method used to determine whether there is enough evidence in a sample of data to infer that a certain condition holds for the entire population.

# I

Term	Definition
<b>Imputation</b>	The process of replacing missing values in a dataset with estimated values, such as the mean, median, mode, or values predicted by a model.
<b>Imbalanced Dataset</b>	A dataset where the classes are not represented equally. For example, a fraud detection dataset where fraudulent transactions are far rarer than legitimate ones.

# K

Term	Definition
<b>K-Fold Cross-Validation</b>	A cross-validation technique that partitions data into K equal subsets (folds), trains the model K times — each time using a different fold as the validation set — and averages the results.
<b>K-Means Clustering</b>	An unsupervised algorithm that partitions data into K clusters by iteratively assigning points to the nearest centroid and updating centroids until convergence.
<b>K-Nearest Neighbours (KNN)</b>	A simple supervised learning algorithm that classifies a data point based on the majority class of its K nearest neighbours in the feature space.

# L

Term	Definition
<b>Label</b>	The target output variable in a supervised learning task. Also called the dependent variable or response variable. For example, in email classification, 'spam' or 'not spam' are labels.

Term	Definition
<b>Lasso Regression (L1 Regularisation)</b>	A regression technique that adds a penalty equal to the absolute value of coefficients to the loss function, encouraging sparsity by shrinking some coefficients to exactly zero.
<b>Latent Variable</b>	A variable that is not directly observed or measured but is inferred from other observed variables. Commonly used in dimensionality reduction and topic modelling.
<b>Learning Rate</b>	A hyperparameter that controls how much the model's weights are adjusted with respect to the loss gradient at each update step during training.
<b>Linear Regression</b>	A supervised learning algorithm that models the relationship between a continuous dependent variable and one or more independent variables by fitting a linear equation to the data.
<b>Logistic Regression</b>	A classification algorithm that models the probability of a binary outcome using a logistic (sigmoid) function. Despite its name, it is used for classification, not regression.
<b>Loss Function</b>	A function that measures the difference between a model's predictions and the actual target values. The goal of training is to minimise the loss function.

## M

Term	Definition
<b>Machine Learning (ML)</b>	A subset of AI that enables systems to learn and improve from experience without being explicitly programmed. Models are trained on data to identify patterns and make decisions.
<b>Mean Absolute Error (MAE)</b>	A regression performance metric that measures the average absolute difference between predicted and actual values. It is robust to outliers compared to MSE.
<b>Mean Squared Error (MSE)</b>	A regression metric that measures the average squared difference between predicted and actual values. Squaring penalises larger errors more heavily.
<b>Model Deployment</b>	The process of integrating a trained machine learning model into a production environment so that it can make predictions on new, real-world data.
<b>Multicollinearity</b>	A situation in regression analysis where two or more predictor variables are highly correlated, making it difficult to determine their individual effects on the response variable.

## N

Term	Definition
<b>Natural Language Processing (NLP)</b>	A branch of AI focused on enabling computers to understand, interpret, and generate human language. Applications include sentiment analysis, translation, and chatbots.
<b>Neural Network</b>	A machine learning model inspired by the human brain, composed of interconnected layers of nodes (neurons) that transform inputs into outputs through learned weights.
<b>Normalisation</b>	The process of scaling features to a standard range to ensure that no single feature disproportionately influences the model due to its magnitude.
<b>Null Hypothesis</b>	The default assumption in hypothesis testing that there is no effect or no difference between groups. A sufficiently low p-value leads to rejection of the null hypothesis.

## O

Term	Definition
<b>One-Hot Encoding</b>	A technique for representing categorical variables as binary vectors, creating a new binary column for each category to make them suitable for use in machine learning algorithms.
<b>Outlier</b>	A data point that differs significantly from other observations. Outliers can distort statistical analyses and model training, and may indicate data errors or rare events of interest.
<b>Overfitting</b>	A modelling error where a model learns the noise and specific patterns of the training data too well, resulting in poor generalisation to new, unseen data.

## P

Term	Definition
<b>P-Value</b>	The probability of observing results at least as extreme as the current data, given that the null hypothesis is true. A low p-value (typically $< 0.05$ ) suggests the results are statistically significant.
<b>Precision</b>	The proportion of true positive predictions among all positive predictions made by the model. A high precision means the model has few false positives.
<b>Principal Component Analysis (PCA)</b>	An unsupervised dimensionality reduction technique that transforms features into a smaller set of uncorrelated components that capture the maximum variance in the data.

## R

Term	Definition
<b>Random Forest</b>	An ensemble learning method that builds multiple decision trees on random subsets of the data and features, then aggregates their predictions by voting (classification) or averaging (regression).
<b>Recall (Sensitivity)</b>	The proportion of actual positives that were correctly identified by the model. High recall means few false negatives. Also referred to as the True Positive Rate.
<b>Regularisation</b>	A technique used to prevent overfitting by adding a penalty to the model's loss function for large or complex coefficients. Common forms include L1 (Lasso) and L2 (Ridge) regularisation.
<b>Regression</b>	A supervised learning task where the goal is to predict a continuous numerical output variable. Examples include predicting house prices or temperature.
<b>ROC Curve (Receiver Operating Characteristic)</b>	A graph showing the trade-off between the True Positive Rate (Recall) and the False Positive Rate at various classification thresholds. The AUC (Area Under the Curve) summarises overall performance.

## S

Term	Definition
<b>Sampling</b>	The process of selecting a subset of data from a larger population for analysis. Proper sampling ensures that the subset is representative of the whole.
<b>Standardisation (Z-Score Normalisation)</b>	Rescaling features so they have a mean of zero and a standard deviation of one. Calculated as $(\text{value} - \text{mean}) / \text{standard deviation}$ .
<b>Supervised Learning</b>	A type of machine learning where the model is trained on labelled data, learning to map inputs to known output labels. Examples include classification and regression.
<b>Support Vector Machine (SVM)</b>	A supervised learning algorithm that finds the optimal hyperplane separating classes in the feature space, maximising the margin between the nearest data points of each class.

## T

Term	Definition
<b>Test Set</b>	A portion of the dataset held back and not used during model training or validation, used exclusively to evaluate the final model's performance on unseen data.

Term	Definition
<b>Time Series</b>	A sequence of data points indexed in time order, typically collected at successive, equally spaced points in time. Analysed to identify trends, seasonal patterns, and forecasts.
<b>Training Set</b>	The portion of the dataset used to train a machine learning model. The model learns patterns from this data by adjusting its parameters to minimise prediction error.
<b>Transfer Learning</b>	A machine learning technique where a model trained on one task is adapted and fine-tuned for a different but related task, reducing the need for large amounts of labelled data.

## U

Term	Definition
<b>Underfitting</b>	A modelling problem where the model is too simple to capture the underlying patterns in the data, resulting in poor performance on both training and test datasets.
<b>Unsupervised Learning</b>	A machine learning approach where the model learns patterns and structure from unlabelled data. Common tasks include clustering, dimensionality reduction, and anomaly detection.

## V

Term	Definition
<b>Validation Set</b>	A portion of the dataset used to tune model hyperparameters and evaluate model performance during training, separate from both the training set and the final test set.
<b>Variance</b>	In statistics, variance measures the spread of data around the mean. In machine learning, high model variance means the model is sensitive to small fluctuations in training data (overfitting).

## W

Term	Definition
<b>Weight</b>	A learnable parameter in a machine learning model that is adjusted during training to minimise the loss function. Weights determine the strength and direction of influence of each feature.