

Data Science Methodology

Table of Contents

1. Course Overview	2
Key Phases in the Data Science Methodology	2
2. Problem Formulation	2
3. Analytics Approach	3
Data Handling Phases	4
4. Data Requirements	4
5. Data Collection	5
6. Data Understanding and Preparation	6
Modeling, Deployment, and Feedback	7
7. Modeling	7
8. Deployment	8
9. Feedback	8
10. Storytelling in Data Science	9
Summary of Methodology Steps	9

1. Course Overview

- **Goal of the Course:** Learn how to think and work like a successful data scientist by applying key data science methodologies.
- **Main Methodologies Covered:**
 - **Foundational Data Science Methodology:** Focuses on a structured approach to solving data science problems.
 - **CRISP-DM (Cross-Industry Standard Process for Data Mining):** A widely used six-stage methodology, essential for tackling various data science scenarios.

The stages of the methodology are iterative, emphasizing continuous improvement through feedback.

Key Phases in the Data Science Methodology

2. Problem Formulation

- **Understanding the Problem:** Begin by identifying the business goals and clearly defining the question you aim to answer.
- **Engaging Stakeholders:** Stakeholders must be involved to clarify goals and set clear objectives for the project.
- **Establishing Objectives:** Break down the problem into smaller parts to prioritize and organize tasks for tackling the issue. Regular communication with stakeholders is essential to align requirements with the project direction.



3. Analytics Approach

- **Choosing the Right Approach:** The chosen analytics method depends on the nature of the question.
 - **Predictive Modeling:** If you are determining the probability of an outcome, a predictive model is appropriate.
 - **Descriptive Analysis:** If the focus is on understanding relationships or trends in the data, use a descriptive model.
 - **Classification Models:** For yes/no questions, a classification approach is ideal.

This step emphasizes the importance of selecting the right analytic tools based on business requirements.

Data Handling Phases

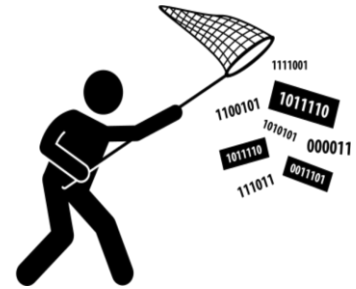


4. Data Requirements

- **Defining Data Needs:** Just as a recipe requires specific ingredients, data science projects require precise data input.
 - Identify necessary data sources.
 - Understand how to source and collect this data.
 - Plan how to process the data to meet the desired outcome.

5. Data Collection

- **Collecting and Evaluating Data:** Gather the identified data and assess its quality using techniques like descriptive statistics and visualization. Look for missing data or gaps and decide whether more data needs to be collected or whether substitutes can be used.



6. Data Understanding and Preparation



- **Data Understanding:** Check whether the collected data is representative of the problem at hand. Analyze the data content and format to ensure it is usable for model building.
- **Data Preparation:** Clean the data by removing duplicates, handling missing values, and resolving formatting issues. Feature engineering may be required to enhance the dataset, making it more suitable for modeling.

Modeling, Deployment, and Feedback

7. Modeling

- **Building the Model:** Develop models based on the collected data. Depending on the project, you may use descriptive or predictive models, often starting with a training set for testing.
- **Experimentation and Refinement:** Try different algorithms and models, adjusting parameters as needed to improve performance.



8. Deployment



Model Implementation: After the model is evaluated and validated, it is deployed in a test environment or to a limited user group to ensure it performs as expected.

The model's real-world effectiveness is tested before it is rolled out to the entire user base.

9. Feedback

- **Gathering Feedback:** After deployment, feedback is crucial for improving the model's performance. Users provide input on how the model is working, which informs continuous refinement.
- **Iterative Improvement:** This methodology emphasizes that each stage, from data collection to deployment, is iterative. Constant feedback loops are used to refine both the model and the overall process.

Additional Aspects of Data Science

10. Storytelling in Data Science

- **Importance of Storytelling:**
 - **Simplifies Complex Data:** Break down intricate data insights into digestible information.
 - **Engages Stakeholders:** A good data story helps drive decision-making by ensuring stakeholders understand the implications of the analysis.
- **Elements of a Good Data Story:**
 - Clear objective and message.
 - Relevant data insights.
 - Visualizations that support the narrative.

Tips for Effective Storytelling:

- Know your audience.
 - Use visuals to enhance understanding.
 - Continuously iterate and refine the story based on feedback.
-

Summary of Methodology Steps:

1. **From Problem to Approach:** Define the business problem and determine the best analytic approach.
2. **Data Requirements and Collection:** Identify, gather, and understand the data needed for analysis.
3. **Data Understanding and Preparation:** Clean and preprocess the data for analysis.
4. **Model Building:** Develop and refine models to predict or describe patterns in the data.
5. **Model Deployment and Feedback:** Deploy the model and gather feedback for further iterations.

Key Takeaways:

- Data science is an iterative process that requires clear problem definition, stakeholder engagement, and continuous refinement.
- Proper data handling—from collection to preparation—is crucial for building effective models.
- Storytelling is an important skill for data scientists to communicate findings and ensure that insights lead to actionable outcomes.