



# Roadmap to Crack Data Science/ML Interviews

This roadmap, constructed from information solely from the sources you provided, offers a detailed plan to guide you on your journey to becoming a data scientist.



Anwar Haq  
Sr. Data Scientist @ Cisco  
Ex-Meta | Ex-Wells Fargo

# Introduction

## Entering the World of Data Science



### What is Data Science?

Data science is a multidisciplinary field that extracts valuable insights from data. It encompasses various techniques, from simple data analysis using SQL queries to building complex predictive models using machine learning algorithms. The demand for data scientists is high and shows resilience even during economic downturns, making it a promising career path.

## Diverse Roles in the Data Science Family

The field of data science encompasses a range of roles with varying responsibilities and skill requirements. Some key roles include:

- ✓ **Data Scientist:**  
Analyses complex datasets to identify patterns, build predictive models, and deliver actionable insights for business decisions.
- ✓ **Data Analyst:**  
Interprets data to produce reports and visualizations, focusing on descriptive statistics to answer specific business questions.
- ✓ **Machine Learning Engineer:**  
Develops, tests, and deploys machine learning models that can learn from data and make predictions. They focus on algorithm development and optimization.
- ✓ **Data Engineer:**  
Designs and maintains the architecture for large-scale data processing systems, ensuring efficient, scalable, and reliable data pipelines.
- ✓ **Data Architect:**  
Designs and manages the organization's data, ensuring data accessibility and well-organised metadata.

# Choosing Your Path

## Interests:

Are you passionate about building models (consider Machine Learning Engineer), analysing trends (Data Analyst or Data Scientist), or data infrastructure (Data Engineer)?

## Skills:

Different roles require specific skillsets. Assess your strengths and identify areas for development.

## Industry Needs:

Research industries that interest you and see which data science roles are in demand. For example, healthcare needs data scientists knowledgeable in healthcare data and regulations, while finance prefers professionals skilled in risk modelling and quantitative analysis.

## Career Prospects:

Investigate the growth potential, average salaries, and job availability for the roles you're considering.





# Research and Networking

## Laying the Groundwork

### Identifying Target Companies

Strategically selecting target companies is crucial for a focused job search. Consider:

**Define Career Goals:** What do you want to achieve in your career? Specific industry? Startup or large corporation? Research institution?

**Research Industry Leaders:** Identify companies known for innovation, data-driven approaches, and culture.

**Consider Company Size and Growth Stage:** Startup (varied experience) vs. larger corporation (stability).

**Evaluate Company Values and Culture:** Does the company's mission and values align with yours? Check employee reviews on platforms like Glassdoor.

**Check Job Listings and Career Pages:** Get a sense of roles and valued skills.  
Geographic Preferences: Does the location match your preferences?



# Analyzing Job Descriptions

Carefully examine job descriptions to understand required skills and qualifications:

**Identify Key Skills and Qualifications:** Recurring technical skills (Python, SQL, machine learning) and soft skills (communication, problem-solving).

**Understand Role Responsibilities:** What does the daily work entail? Does it align with your interests and skills?

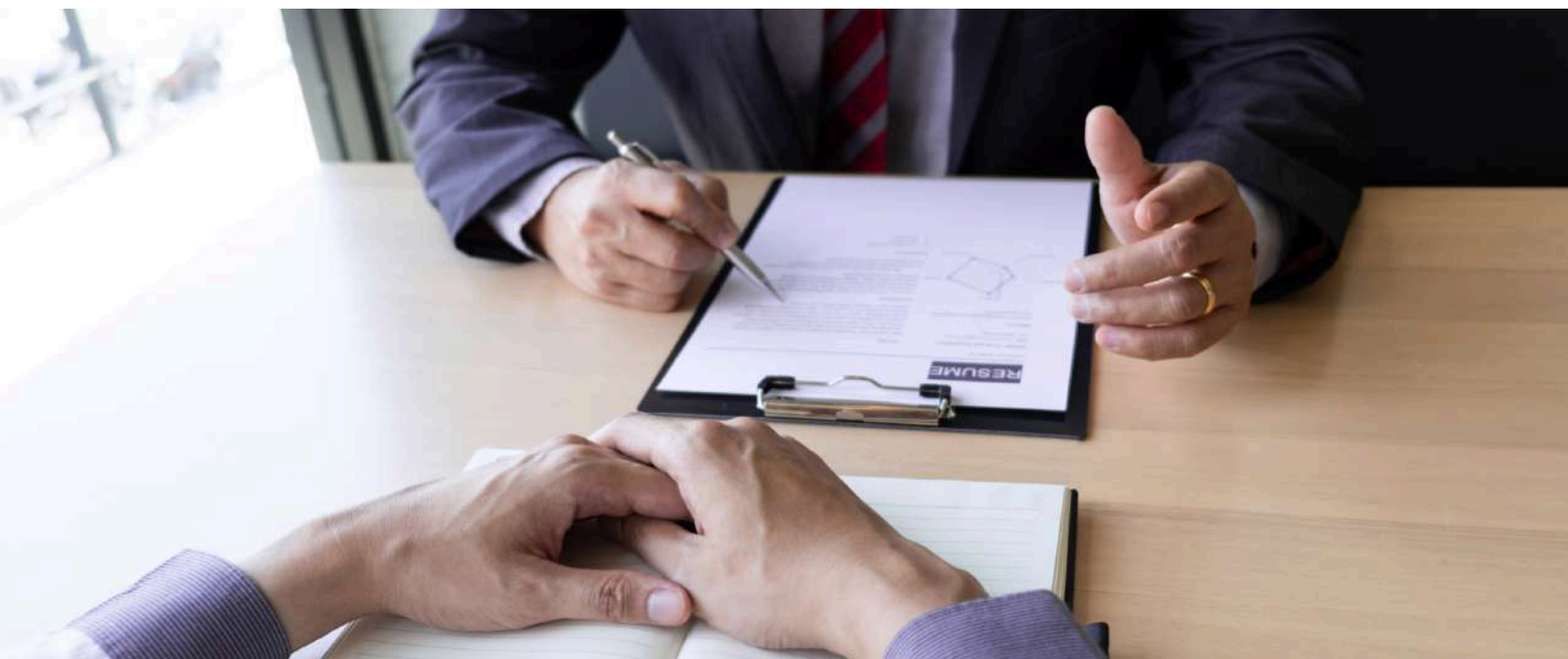
**Look for Required Experience:** Years of experience needed.

**Check for Preferred Skills:** Can give you an edge.

**Analyse Job Titles:** Different titles for similar roles ("Data Analyst" vs. "Business Analyst").

**Evaluate Job Requirements:** Tailor your resume and cover letter accordingly.

**Note Industry-Specific Requirements:** E.g., healthcare data science roles may require knowledge of HIPAA regulations.



# Building Your Network

Networking is vital for your job search and career development.

## Use LinkedIn Effectively:

- **Build a Strong Profile:** Highlight your data science skills, experiences, and achievements.
- **Connect with Industry Professionals:** Data scientists, hiring managers, and professionals at your target companies.
- **Follow Relevant Companies and Influencers:** Stay updated on industry trends and job openings.

**Attend Industry Events:** Conferences, meetups, and workshops provide opportunities to learn, network, and connect with potential employers.

**Engage in Online Communities:** Participate in data science forums and contribute to open-source projects. Platforms like Kaggle and DataCamp offer communities for sharing knowledge and collaborating on projects.

**Seek Mentorship:** Find experienced data scientists for guidance and support. Leverage company mentorship programs.





# Building Foundational Skills



## Statistics and Mathematics

### 1. Descriptive Statistics

Summarizing and describing data using measures like:

- **Mean:** The average of a dataset.
- **Median:** The middle value in a sorted dataset.
- **Mode:** The most frequent value in a dataset.
- **Standard Deviation:** A measure of how spread out the data is.

### 2. Inferential Statistics

Making inferences and predictions about a population based on a sample.

Key concepts:

- **Hypothesis Testing:** Testing assumptions about a population based on sample data.
- **Confidence Intervals:** Estimating a range of values for a population parameter.

### 3. Statistical Tests

- **T-Tests:** Comparing means of two groups.
- **ANOVA (Analysis of Variance):** Comparing means of multiple groups.
- **Regression Analysis:** Modelling the relationship between variables.
- **Time Series Analysis:** Analysing data over time.

The depth of your statistical knowledge will depend on your chosen data science role.

## Resources

**Books:** "Practical Statistics for Data Scientists" provides a beginner-friendly introduction with Python examples.

**Online Courses:** Platforms like Coursera and edX offer statistics courses tailored for data science, such as "Statistics with R" and "Inferential Statistics."

**YouTube Channels:** StatQuest with Josh Starmer offers clear and engaging explanations of statistical concepts.

## Time allocation:

**Allocate approximately 3 weeks to build a solid foundation in statistics.**





# Building Foundational Skills

## Exploratory Data Analysis (EDA)

EDA is the crucial first step in any data science project, enabling you to understand your data through cleaning, preprocessing, and visualisation techniques.

### 1. Data Cleaning and Preprocessing

- **Handling Missing Values:** Imputing missing values or removing rows with missing data.
- **Detecting and Treating Outliers:** Identifying and addressing extreme values that can skew your analysis.
- **Data Transformation:** Converting data types, scaling variables, and creating new features.

### 2. Univariate Analysis

Exploring individual variables using:

- **Histograms:** Visualizing the distribution of numerical data.
- **Box Plots:** Displaying the distribution and identifying outliers.
- **Summary Statistics:** Calculating measures like mean, median, and standard deviation.

### 3. Bivariate Analysis

Examining relationships between pairs of variables using:

- **Scatter Plots:** Visualizing the relationship between two numerical variables.
- **Correlation Coefficients:** Measuring the strength and direction of the relationship.
- **Cross-Tabulations:** Analysing relationships between categorical variables.

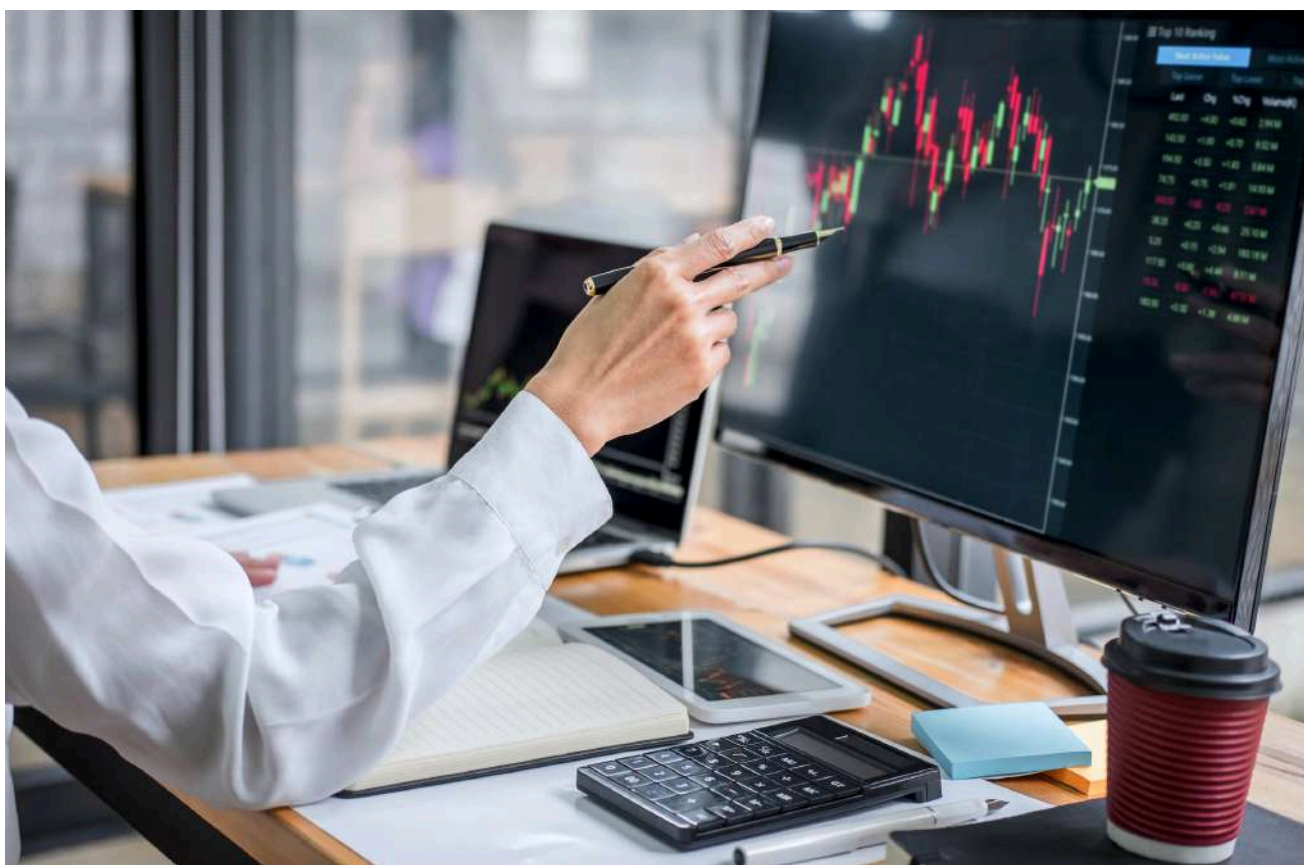
## 4. Multivariate Analysis

Investigating interactions among multiple variables using techniques like:

- **Principal Component Analysis (PCA):** Reducing the dimensionality of data while retaining important information.

### Time allocation:

Allocate 3 weeks to develop proficiency in EDA.



# Building Foundational Skills

## Data Visualization



Data visualization is crucial for communicating insights and findings to both technical and non-technical audiences. You will need to:

**Select appropriate visualization techniques:** Choose the best visualization to represent your data and insights effectively, considering the type of data and the message you want to convey. For example:

- Bar charts and line graphs are suitable for comparing categories or showing trends over time.
- Scatter plots are useful for visualizing relationships between two variables.
- Heatmaps are effective for displaying patterns in large datasets.

**Use effective design principles:** Create visually appealing and easy-to-understand visualizations by considering factors like color choice, labelling, and chart layout.

**Tell a story with your data:** Use visualizations to support a narrative, highlighting key insights and trends.

## Resources

**Python Libraries:** Matplotlib and Seaborn offer powerful tools for creating static and interactive visualizations.

**R Packages:** Ggplot2 is a popular package in R for creating aesthetically pleasing and informative visualizations.

**Data Visualisation Tools:** Tableau, Power BI, and even Excel can be used for creating interactive dashboards and reports.

**Books:** "Storytelling with Data" by Cole Nussbaumer Knafllic provides a comprehensive guide to creating impactful visualizations and telling stories with data.



**Online Courses:** DataCamp's "Understanding Data Visualization" course helps deepen your data visualization competencies.

### Time allocation:

Allocate 3 weeks to develop your data visualisation skills, focusing on both technical skills and design principles



# Acquiring Essential Coding Skills

## Python

Python is the dominant language in data science due to its versatility, simplicity, and extensive libraries.

### 1. Core Python

- **Syntax:** Understanding the basic rules of Python code.
- **Data Types:** Working with numbers, strings, booleans, and lists.
- **Control Flow:** Using conditional statements and loops to control program execution.
- **Functions:** Creating reusable blocks of code to perform specific tasks.

### 2. Essential Libraries

- **NumPy:** For numerical computing and array manipulation.
- **Pandas:** For data manipulation and analysis.
- **Scikit-learn:** For machine learning algorithms and model building.
- **Matplotlib:** For data visualization.

### 3. Advanced Techniques

- **Object-Oriented Programming (OOP):** For creating reusable and modular code.
- **Exception Handling:** For gracefully managing errors in your code.
- **File Input/Output:** For reading and writing data from files.
- **Regular Expressions:** For pattern matching and text manipulation.

## Resources

**Tutorials:** The official Python documentation and Real Python offer comprehensive resources.

**Online Courses:** Coursera's "Python for Everybody" and Udemy's "Complete Python Bootcamp" are popular choices.

**Projects:** Participate in Kaggle competitions and contribute to GitHub repositories to practice your Python skills.



### Time allocation:

Allocate approximately 3 weeks to gain proficiency in the basics of Python and an additional 7 weeks to enhance your skills while working on projects.



# Acquiring Essential Coding Skills

## SQL

SQL (Structured Query Language) is essential for interacting with databases, enabling you to extract, manipulate, and analyse data.



SQL (Structured Query Language) is essential for interacting with databases, enabling you to extract, manipulate, and analyse data.

## 1. Fundamental Concepts

## SELECT Statement: Retrieving data from tables.

## WHERE Clause: Filtering data based on conditions.

## ORDER BY Clause: Sorting data in ascending or descending order.

## JOIN Operations: Combining data from multiple tables based on relationships.

**Aggregate Functions:** Performing calculations like SUM, AVG, COUNT, MIN, and MAX.

## GROUP BY Clause: Grouping data for aggregate calculations.

## 2. Advanced Techniques

**Subqueries:** Queries nested within other queries to perform more complex operations.

**Analytical Functions:** Providing advanced data analysis capabilities, like ranking and windowing.

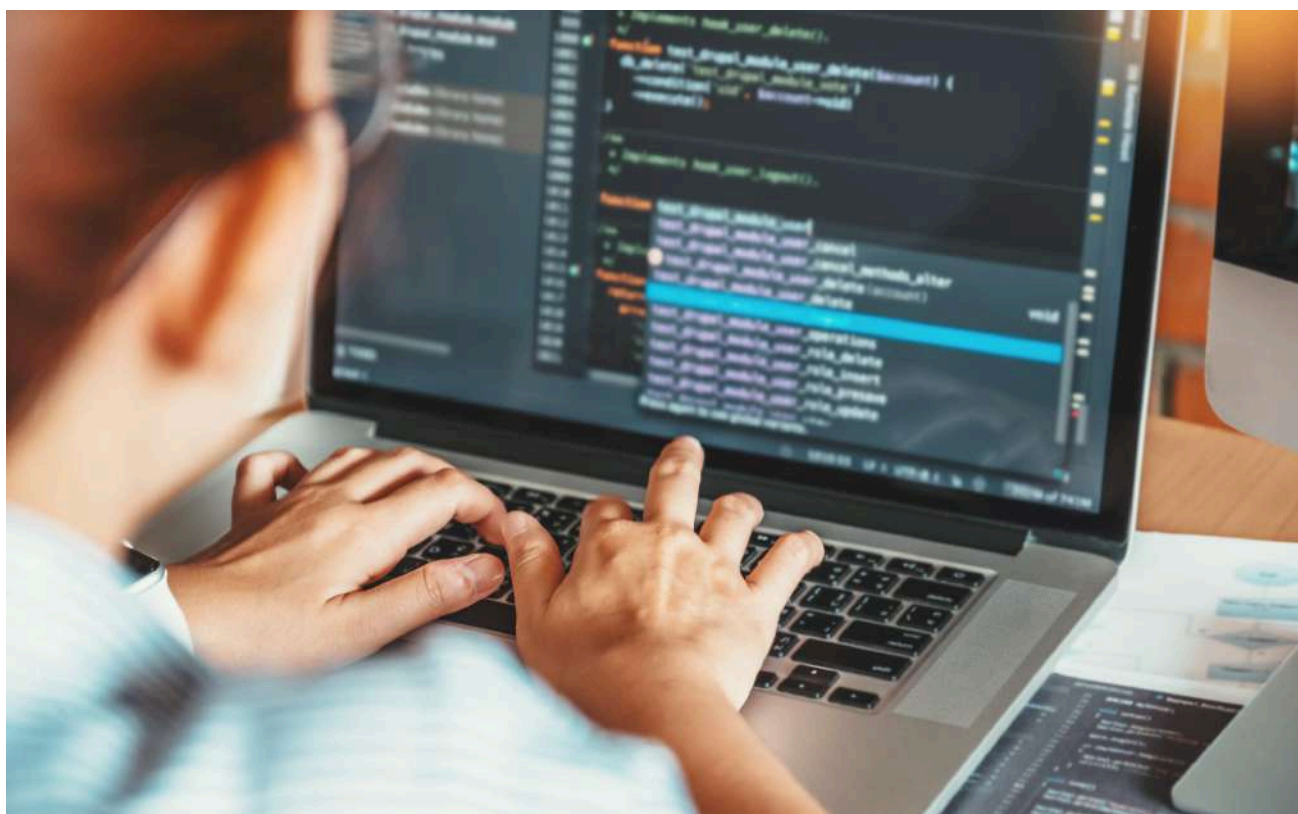
## Common Table Expressions (CTEs): Simplifying complex queries.

## Resources

**Tutorials:** W3Schools SQL Tutorial and Mode Analytics SQL Tutorial offer practical lessons.

**Online Courses:** Coursera's "SQL for Data Science" and DataCamp's "Introduction to SQL" are recommended.

**Projects:** Practice SQL problems on platforms like LeetCode and SQLZoo.



### Time allocation:

Allocate approximately 3 weeks to achieve job readiness in SQL.

# Machine Learning

## The Heart of Data Science

### Introduction to Machine Learning

Machine learning, a core component of data science, involves training algorithms to learn patterns and make decisions based on data.

#### 1. Supervised Learning

Training models on labelled data, where the outcome is known. Key techniques:

**Regression:** Predicting a continuous target variable.

- Linear Regression
- Polynomial Regression

**Classification:** Predicting a categorical target variable.

- Logistic Regression
- Decision Trees
- Support Vector Machines (SVMs)

#### 2. Unsupervised Learning

Exploring unlabeled data to discover patterns. Key techniques:

**Clustering:** Grouping data points based on similarity.

- K-Means Clustering
- Hierarchical Clustering

**Dimensionality Reduction:** Reducing the number of features while preserving important information.

- Principal Component Analysis (PCA)



### 3. Deep Learning

Using artificial neural networks to analyse complex data and solve challenging tasks. Key concepts:

- **Neural Networks:** Interconnected nodes that process and transmit information.
- **Activation Functions:** Introducing non-linearity into the network.
- **Backpropagation:** Adjusting weights to improve model performance.
- **Convolutional Neural Networks (CNNs):** Specialized for image and video processing.
- **Recurrent Neural Networks (RNNs):** Designed for sequential data like time series and text.

### 4. AI Feedback Loops

Incorporating feedback mechanisms to enable models to learn from their output. Examples:

- **Reinforcement Learning:** Training agents to make decisions in an environment through rewards and penalties.

## Resources

**Books:** "Hands-On Machine Learning with Scikit-Learn, Keras & TensorFlow" provides a practical guide.

**Online Courses:** Coursera's "Machine Learning" by Andrew Ng and DataCamp's "Machine Learning Scientist with Python" career track are highly recommended.

**YouTube Channels:** StatQuest with Josh Starmer and 3Blue1Brown offer intuitive explanations of machine learning concepts.

### Time allocation:

Allocate 3 weeks for learning fundamental concepts in machine learning and an additional 7 weeks for delving deeper into different areas and building portfolio projects.

# Mastering Data Science Tools

## Essential Data Science Tools

Proficiency in various tools is essential for coding, version control, visualization, and deployment of data science projects.

### 1. Notebooks and IDEs

**Jupyter Notebook:** An interactive environment for running code, visualizing data, and creating documentation.

**Google Colab:** A cloud-based platform similar to Jupyter Notebook, offering free access to GPUs for deep learning.

**Integrated Development Environments (IDEs):** Specialized software for code editing, debugging, and project management. Popular IDEs include:

- PyCharm
- VS Code
- Spyder

### 2. Version Control with Git

**Git:** A distributed version control system for tracking changes in code and collaborating with others.

**GitHub:** A popular platform for hosting Git repositories and collaborating on projects

#### Key Concepts:

- Branching
- Merging
- Pull Requests

### 3. Data Visualisation Platforms

**Tableau:** A powerful tool for creating interactive dashboards and visualizations.

**Power BI:** A business intelligence platform for data analysis and reporting.

**Plotly:** A library for creating interactive and web-based visualizations.

### 4. Cloud Computing Platforms

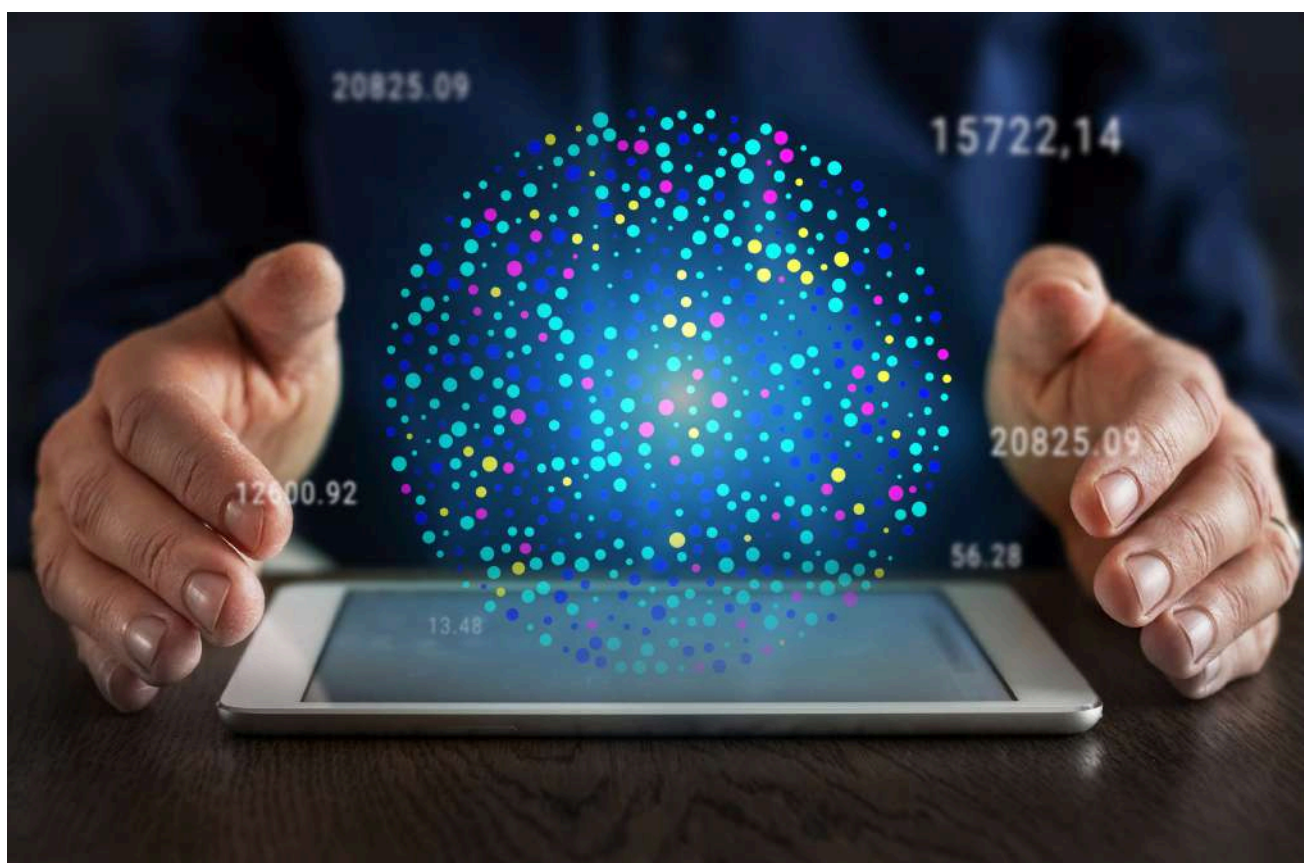
**Amazon Web Services (AWS):** A comprehensive cloud platform offering various services for storage, computing, and machine learning.

**Microsoft Azure:** A cloud platform with services for data storage, analytics, and AI.

**Google Cloud Platform (GCP):** A suite of cloud computing services for data storage, analysis, and machine learning.

### Time allocation:

You can incorporate these tools into your workflow as you progress through the other steps, allocating dedicated time as needed.





# Building Your Portfolio

A well-crafted portfolio showcases your skills and practical experience, making you a more competitive candidate.

## 1. Project Selection

Choose projects that cover diverse aspects of data science, demonstrating a broad skillset

- **Data Cleaning and Preprocessing:** Showcase your ability to handle real-world data with missing values, inconsistencies, and outliers.
- **Exploratory Data Analysis (EDA):** Demonstrate your skills in data exploration, visualization, and identifying patterns.
- **Machine Learning Model Building:** Include projects involving various algorithms and techniques, such as classification, regression, clustering, and deep learning.
- **Domain-Specific Projects:** If you have a specific industry in mind, build projects that address problems in that domain. For example, you could create a project on customer churn prediction for the telecom industry or a fraud detection model for the financial sector.

## 2. Project Documentation

Clearly document each project, making it easy for others to understand your work

- **Introduction:** Describe the problem you addressed and the project's objectives.
- **Data:** Explain the datasets used, their source, and any preprocessing steps taken.
- **Methodology:** Detail the techniques and algorithms used for analysis and model building.
- **Results:** Present your findings using clear visualizations and insightful interpretations.
- **Conclusion:** Summarize your key takeaways and any limitations of your approach.
- **Code:** Include well-commented code snippets to demonstrate your coding skills.

### 3. Creating an Online Portfolio

Host your projects on platforms like:

- **GitHub:** A popular platform for sharing code and collaborating on projects. Create a well-organized repository for each project, including your code, documentation, and any supporting files.
- **Personal Website:** Build a website to showcase your projects and skills professionally. Include a portfolio section with links to your GitHub repositories or detailed project descriptions.

### 4. Highlighting Key Skills

Ensure your portfolio highlights your proficiency in essential data science skills:

- **Data Preprocessing**
- **Statistical Analysis**
- **Machine Learning**
- **Data Visualization**
- **Tools and Technologies Used (Python, R, SQL, TensorFlow, Tableau)**

## Resources

**Books:** "Data Science Projects with Python" by Stephen Klosterman and "Python for Data Analysis" by Wes McKinney offer guidance on building practical projects.

**Online Platforms:** Kaggle's "Learn" courses provide project-based learning opportunities.

**Portfolio Examples:** Explore portfolios of experienced data scientists for inspiration and ideas.

### Time allocation:

Allocate approximately 7 weeks to develop 7 portfolio projects, working on documentation and website creation alongside your project development.