

Python



Python has a huge library and set of packages that offer powerful data manipulation capabilities.



Jupyter Notebook: An open-source web application widely used for data cleaning and transformation, statistical modeling, and data visualization.



NumPy (Numerical Python):

- The most basic package that Python offers
- It is fast, versatile, interoperable, and easy to use
- It provides support for large, multi-dimensional arrays and matrices, and high-level mathematical functions to operate on these arrays



Pandas:

- Designed for fast and easy data analysis operations
- Allows complex operations such as merging, joining, and transforming huge chunks of data using simple, single-line commands
- Helps prevent common errors that result from misaligned data coming from different sources

R



R offers a series of libraries and packages that are explicitly created for wrangling messy data.

Using these libraries, you can investigate, manipulate, and analyze data.



dplyr: A powerful library for data wrangling with a precise and straightforward syntax.



Data.table: Helps aggregate large data sets quickly.



Jsonlite: A robust JSON parsing tool, great for interacting with web APIs.

Performance Tuning and Troubleshooting:

↳ Performance Threats in Data Pipeline:

- Scalability:** Data pipeline fails to handle large volumes of data or distributed network load. Too much data in a pipeline is a root cause of performance issues.
- Application Failures:** Failure to handle errors in processing various stages of data flow due to bugs in the code.
- Scheduling Issues:** Inconsistent scheduling leads to work flow issues.
- Tool Incompatibilities:** Using multiple tools in a pipeline creates compatibility issues between pipeline tools and the tools they interact with.

↳ Performance Metrics for Data Pipeline:

- Latency:** Time taken for individual requests to pass through the pipeline.
- Failure:** Number of errors.
- Resource Utilization:** CPU, memory, network usage over time.
- Traffic:** Total volume of data processed.

↳ Troubleshooting Data Pipeline Performance Issues:

- Information Collection:** Gathering metrics to identify performance bottlenecks.
- Version and Deployment Check:** Checking software versions for performance issues.
- Log and Metrics Analysis:** Analyzing logs and metrics across infrastructure, data and software.
- Reproduction in Test Environment:** Reproducing issues in a controlled environment to identify root causes.

↳ Database Optimization Practise:

- Capacity Planning:** Plan for resource needs hardware and software requirements.
- Indexing:** Create indexes for frequently accessed data.
- Partitioning:** Partitioning data into smaller, manageable Query sets for better management.
- Normalization:** Normalize database to reduce complexity in queries.

↳ Monitoring and Alerting Systems:

- Database Monitoring Tools:** Capture performance metrics from databases like MySQL, PostgreSQL, Oracle.
- Application Performance Management (APM) Tools:** Monitor performance of application layers including databases and external resources.
- Query Performance Monitoring Tools:** Measure execution times of queries across databases.
- Job-level Runtime Monitoring:** Monitor scheduled jobs for performance and reliability.

↳ Maintenance Routines:

- Preventive Maintenance:** Regularly scheduled maintenance tasks to avoid problems.
- Time-based Maintenance:** Maintenance tasks performed at specific intervals.
- Condition-based Maintenance:** Maintenance based on monitoring and detecting anomalies in the system's behavior.

Governance and Compliance:

↳ **Data Governance**: őö collection and principles within, practice, is: process to maintain security, privacy, and integrity via data life cycle

↳ **Types of Data Information**: • privacy data, sensitive data, data that can be traced back to an individual or organization

↳ **Compliance**: meeting regulations like: GDPR, HIPAA, PCI DSS, Sarbanes Oxley Act, FERPA, etc.

In this lesson, you have learned:

Data Governance is a collection of principles, practices, and processes that help maintain the security, privacy, and integrity of data through its lifecycle.

Personal Information and Sensitive Personal Information, that is, data that can be traced back to an individual or can be used to identify or cause harm to an individual, needs to be protected through governance regulations.

General Data Protection Regulation, or GDPR, is one such regulation that protects the personal data and privacy of EU citizens for transactions that occur within EU member states.

Regulations, such as HIPAA (Health Insurance Portability and Accountability Act) for Healthcare, PCI DSS (Payment Card Industry Data Security Standard) for retail, and SOX (Sarbanes Oxley) for financial data are some of the industry-specific regulations.

Compliance covers the processes and procedures through which an organization adheres to regulations and conducts its operations in a legal and ethical manner.

Compliance requires organizations to maintain an auditable trail of personal data through its lifecycle, which includes acquisition, processing, storage, sharing, retention, and disposal of data.

Tools and technologies play a critical role in the implementation of a governance framework, offering features such as:

- Authentication and Access Control.
- Encryption and Data Masking.
- Hosting options that comply with requirements and restrictions for international data transfers.
- Monitoring and Alerting functionalities.
- Data erasure tools that ensure deleted data cannot be retrieved.

Module 4

In this lesson, you have learned:

Data Engineering is reported to be one of the top ten jobs experiencing tremendous growth in the U.S. today. It is also reported to be one of the fastest growing tech occupations with year-over-year growth of around 50%.

Currently, the demand for skilled data engineers far outweighs the supply, which means companies are willing to pay a premium to hire skilled data engineers.

Data engineering roles in organizations tend to break the specialization up into Data Architecture, Database Design and Architecture, Data Platforms, Data Pipelines and ETL, Data Warehouses, and Big Data.

Regardless of the niche you choose to specialize in, knowledge of operating systems, languages, databases, and infrastructure components, is essential.

To work your way up from a Junior Data Engineer to a Lead or Principal Data Engineer, you need to continually advance your technical, functional, and soft skills from a foundational level to an expert level. You need to not only expand your skills in your niche area, but also into other areas of data engineering at the same time.

Big Data Engineers and Machine Learning Engineers are some of the emerging roles in this field and they require specialized skills in addition to basic data engineering.

There are several paths you can consider in order to gain entry into the data engineering field.

- An academic degree in Computer Science or engineering qualifies you for an entry-level job.
- If you are not a graduate, or a graduate in a non-related stream, you can earn professional certifications from online multi-course specializations offered by learning platforms such as Coursera, edX, and Udacity.
- If you have a coding background, or you are an IT Support Specialist, a Software Tester, a Programmer, or a data professional such as a Statistician, Data Analyst, or BI Analyst, you can upskill with the help of online courses to become a Data Engineer.