

# DATASCIENCE CAREER TRACK - 2025 EDITION

"From Fundamentals to Industry-Ready Projects"

## **Market Demand Note**

Data Science remains one of the most in-demand career paths in 2025, with applications in AI, healthcare, finance, e-commerce, and technology. Reports predict 11+ million job openings globally by 2026. This course equips learners with hands-on, industry-ready skills.

# ■ Data Science Career Track — 2025 Edition

#### **Table of Contents**

- 1. Introduction to Data Science
- 2. Python Programming for Data Science
- 3. Data Analysis & Visualization
- 4. Statistics & Probability
- 5. SQL for Data Science
- 6. Machine Learning Fundamentals
- 7. Advanced Machine Learning
- 8. Deep Learning
- 9. Data Engineering Basics
- 10. Cloud & Big Data Tools
- 11. Capstone Projects
- 12. Career Preparation & Portfolio Building

Module 1: Python Programming for Data Science (25 Hours)

- **6** Learning Objectives:
  - Master Python fundamentals and advanced concepts for data manipulation
  - Implement efficient data processing using NumPy and Pandas
  - Create compelling data visualizations using Matplotlib and Seaborn
  - Develop proficiency in Jupyter Notebook environment for interactive analysis
  - Build optimization skills for performance-critical data operations
- Detailed Topics Covered:

Python Fundamentals (8 Hours)

Environment Setup (1 Hour)

- What you'll learn: How to set up Python properly on your machine
- Specific topics: Installing Anaconda/Miniconda, creating virtual environments for different projects, managing packages with pip and conda

- Why it matters: Professional developers use isolated environments to avoid conflicts between different projects
- Real example: Creating separate environments for web scraping vs machine learning projects

## Core Python Concepts (3 Hours)

- Variables & Data Types: Understanding strings, integers, floats, booleans, lists, dictionaries, tuples
- Control Structures: Writing if-else conditions, for loops, while loops, nested loops
- Functions: Creating reusable code blocks, parameters, return values, scope
- What you'll build: Simple calculator, password generator, basic data processor

## Object-Oriented Programming (2 Hours)

- Classes & Objects: Creating blueprints for data structures
- Inheritance: Building on existing classes to create specialized versions
- Methods: Functions that belong to classes
- Real application: Building a DataProcessor class that can handle different file types

## Error Handling & File Operations (2 Hours)

- Try-Catch Blocks: Handling errors gracefully without crashing your program
- File Reading/Writing: Opening CSV files, JSON files, text files
- Data Validation: Checking if data is in the correct format
- Practical example: Reading a CSV file and handling missing or corrupted data

#### Data Science Libraries (10 Hours)

#### NumPy - Numerical Computing (2.5 Hours)

- What is NumPy: The foundation for all numerical computing in Python
- Arrays: Creating and manipulating multi-dimensional arrays (like Excel sheets but much more powerful)
- Mathematical Operations: Adding, subtracting, multiplying entire datasets at once
- Broadcasting: Performing operations on arrays of different sizes

 Real application: Calculating profits across thousands of products simultaneously

## Pandas - Data Manipulation (3.5 Hours)

- DataFrames: Think of these as super-powered Excel spreadsheets
- Data Loading: Reading CSV, Excel, JSON, SQL databases into Python
- Data Cleaning: Removing duplicates, handling missing values, fixing data types
- Groupby Operations: Summarizing data (like pivot tables in Excel)
- Merging Data: Combining data from different sources
- Real project: Analyzing sales data from multiple stores and time periods

## Matplotlib - Basic Visualization (2 Hours)

- Plot Types: Line charts, bar charts, histograms, scatter plots
- Customization: Colors, labels, titles, legends, styling
- Subplots: Creating multiple charts in one figure
- Practical example: Creating a sales performance dashboard

#### Seaborn - Statistical Visualization (2 Hours)

- What is Seaborn: Makes complex statistical plots easy
- Correlation Heatmaps: Visualizing relationships between variables
- Distribution Plots: Understanding how data is spread
- Categorical Plots: Comparing groups of data
- Real application: Analyzing customer behavior patterns

## Advanced Python for Data Science (7 Hours)

## Efficient Code Writing (2 Hours)

- List Comprehensions: Writing loops in one line for faster execution
- Lambda Functions: Creating small, temporary functions
- Generators: Processing large datasets without running out of memory
- Example: Processing million-row datasets efficiently

## Text Processing with Regular Expressions (1.5 Hours)

Pattern Matching: Finding phone numbers, emails, dates in messy text data

- Text Cleaning: Removing unwanted characters, standardizing formats
- Data Extraction: Pulling specific information from unstructured text
- Real scenario: Cleaning customer feedback data for analysis

#### Web Scraping & APIs (2.5 Hours)

- Requests Library: Getting data from websites and web services
- BeautifulSoup: Extracting information from HTML pages
- API Integration: Connecting to data sources like weather services, social media
- Ethical Considerations: Legal and responsible data collection practices
- Hands-on project: Building a price comparison tool for e-commerce sites

## Performance Optimization (1 Hour)

- Vectorization: Using NumPy and Pandas efficiently to speed up calculations
- Memory Management: Handling large datasets without crashing your computer
- Profiling: Finding and fixing slow parts of your code
- Real impact: Reducing data processing time from hours to minutes

# **%** Practical Exercises:

- 1. Data Import Challenge: Students receive messy CSV files from 5 different sources (sales data, customer data, product data) and must clean and combine them into a single usable dataset
- 2. Visualization Project: Create a comprehensive analysis report for a real company's data, including sales trends, customer demographics, and product performance
- 3. Web Scraping Assignment: Build a tool that automatically collects product prices from multiple e-commerce websites and tracks price changes over time
- 4. Performance Optimization: Take a slow data processing script and optimize it to run 10x faster using proper pandas techniques
- 5. Mini Project: Build a complete data cleaning pipeline that can handle any CSV file and automatically detect and fix common data issues

#### Module 2: Mathematics & Statistics (15 Hours)

Learning Objectives:

- Apply statistical concepts to real-world data analysis scenarios
- Understand probability distributions and their applications in data science
- Implement hypothesis testing for data-driven decision making
- Master linear algebra concepts essential for machine learning
- Develop intuition for optimization principles underlying ML algorithms
- Detailed Topics Covered:

Statistical Foundations (8 Hours)

Descriptive Statistics (2 Hours)

- Central Tendency: Mean (average), median (middle value), mode (most common value)
- Why each matters: When to use mean vs median (hint: median is better when you have outliers like CEO salaries)
- Variability Measures: Standard deviation, variance, range, interquartile range
- Real application: Analyzing employee salaries which measure gives the most accurate picture?
- Distribution Shapes: Normal distribution (bell curve), skewed distributions, what they tell us about data

#### Probability Theory (2 Hours)

- Basic Probability: Understanding chances and likelihood
- Probability Distributions: Normal, binomial, Poisson distributions and when to use each
- Bayes Theorem: How new information updates our beliefs
- Real example: Email spam detection how likely is an email to be spam based on certain words?
- Conditional Probability: Probability of something happening given that something else happened

## Hypothesis Testing (2 Hours)

- What is Hypothesis Testing: Scientific method for making decisions with data
- P-values: What they really mean (not what most people think!)
- Confidence Intervals: How certain we can be about our estimates

- Type I & Type II Errors: False positives and false negatives in business decisions
- Real scenario: Testing whether a new website design actually increases sales

## Correlation Analysis (1 Hour)

- Pearson Correlation: Measuring linear relationships between variables
- Spearman Correlation: Measuring monotonic relationships (doesn't have to be straight line)
- Causation vs Correlation: Why ice cream sales and drowning deaths are correlated but unrelated
- Practical example: Finding which factors actually drive customer satisfaction

#### Statistical Tests (1 Hour)

- T-tests: Comparing averages between groups
- Chi-square Tests: Testing relationships between categorical variables
- ANOVA: Comparing multiple groups at once
- When to use which test: Decision tree for choosing the right statistical test

## Linear Algebra for Data Science (4 Hours)

#### Vectors and Matrices (2 Hours)

- What are Vectors: Lists of numbers that represent data points
- Vector Operations: Adding, subtracting, dot products
- Matrices: Tables of numbers, like spreadsheets but for math
- Matrix Operations: Multiplication, transpose, inverse
- Real application: How recommendation systems use matrix operations to suggest products

#### Eigenvalues and Eigenvectors (1 Hour)

- Intuitive Understanding: Special directions in data that preserve structure
- Principal Component Analysis: Finding the most important patterns in data
- Dimensionality Reduction: Simplifying complex data while keeping important information
- Business application: Reducing 100 survey questions to 5 key satisfaction factors

## Matrix Decomposition (1 Hour)

- SVD (Singular Value Decomposition): Breaking matrices into simpler parts
- Applications: Image compression, recommendation systems, noise reduction
- QR Decomposition: Another way to break down matrices for solving systems
- Practical example: How Netflix recommends movies using matrix factorization

## Calculus and Optimization (3 Hours)

## Derivatives and Gradients (1.5 Hours)

- What are Derivatives: Rate of change, slope of curves
- Why Data Scientists Care: All machine learning is about finding optimal points
- Partial Derivatives: Dealing with functions of multiple variables
- Gradient: The direction of steepest increase
- Visual example: Finding the lowest point on a hill (cost function minimization)

## Optimization Algorithms (1.5 Hours)

- Gradient Descent: The fundamental algorithm behind machine learning
- How it Works: Taking steps toward the optimal solution
- Learning Rate: How big steps to take (too big = overshoot, too small = slow)
- Local vs Global Minima: Getting stuck in suboptimal solutions
- Real application: How neural networks learn by adjusting millions of parameters

## **Representation** Practical Exercises:

- A/B Test Analysis: Students analyze real A/B test data from a company testing two website designs, determining statistical significance and making business recommendations
- 2. Probability Simulation: Build a Monte Carlo simulation to estimate the probability of various business outcomes (like sales targets)
- 3. PCA Implementation: Implement Principal Component Analysis from scratch using NumPy to reduce the dimensionality of a high-dimensional dataset
- 4. Correlation Investigation: Analyze a dataset with 20+ variables to find hidden relationships and create actionable business insights
- 5. Optimization Challenge: Implement gradient descent algorithm manually to optimize a simple machine learning model

## Module 3: Data Engineering & Processing (15 Hours)

- Learning Objectives:
  - Design efficient data collection strategies from multiple sources
  - Master data cleaning and preprocessing techniques for real-world messy data
  - Implement scalable data processing using big data tools
  - Build robust data pipelines for production environments
  - Ensure data quality and consistency across different data sources
- Detailed Topics Covered:

Data Collection & Storage (5 Hours)

SQL Fundamentals (2 Hours)

- Database Basics: Understanding tables, rows, columns, relationships
- Essential SQL Commands: SELECT, WHERE, GROUP BY, ORDER BY, HAVING
- Joins: Combining data from multiple tables (INNER, LEFT, RIGHT, FULL OUTER)
- Advanced Queries: Subqueries, Common Table Expressions (CTEs), window functions
- Real application: Extracting sales insights from an e-commerce database with customers, orders, and products tables
- Performance: Creating indexes, query optimization, understanding execution plans

## NoSQL Databases (1 Hour)

- When to Use NoSQL: Handling unstructured data like social media posts, sensor data
- MongoDB with Python: Storing and querying document-based data
- JSON Data Handling: Working with nested, flexible data structures
- Real scenario: Storing and analyzing customer reviews with ratings, text, and metadata

#### API Integration (1 Hour)

• REST APIs: Understanding GET, POST, PUT, DELETE requests

- Authentication: API keys, OAuth, handling rate limits
- Error Handling: Dealing with failed requests, timeouts, server errors
- Data Pagination: Handling large datasets that come in chunks
- Practical example: Automatically collecting weather data, stock prices, or social media data

#### Web Scraping (1 Hour)

- When Scraping is Appropriate: Legal and ethical considerations
- Static Scraping: BeautifulSoup for HTML parsing
- Dynamic Scraping: Selenium for JavaScript-heavy sites
- Anti-bot Measures: Handling CAPTCHAs, rate limiting, IP blocking
- Real project: Building a job posting aggregator or price monitoring system

## Data Cleaning & Preprocessing (6 Hours)

## Handling Missing Data (1.5 Hours)

- Types of Missing Data: Missing completely at random vs. systematic patterns
- Detection Strategies: Identifying patterns in missing data
- Imputation Techniques:
  - Simple: Forward fill, backward fill, mean/median imputation
  - Advanced: Regression imputation, K-nearest neighbors, iterative imputation
- When to Drop vs. Impute: Decision framework for handling missing values
- Real example: Cleaning customer survey data where people skip questions

## Outlier Detection and Treatment (1.5 Hours)

- What are Outliers: Extreme values that might be errors or genuine edge cases
- Detection Methods:
  - Statistical: Z-score, IQR method, modified Z-score
  - Visual: Box plots, scatter plots, histograms
  - Advanced: Isolation forests, local outlier factor
- Treatment Options: Remove, cap, transform, or keep outliers
- Business Context: When outliers are errors vs. valuable insights
- © 2025 Udaan Tech Academy. All Rights Reserved.

 Case study: Detecting fraudulent transactions vs. legitimate high-value purchases

## Data Transformation (2 Hours)

- Scaling and Normalization:
  - Min-Max scaling: Putting all variables on 0-1 scale
  - Standardization: Making variables have mean=0, std=1
  - Robust scaling: Handling outliers during scaling
- Categorical Encoding:
  - One-hot encoding: Converting categories to binary columns
  - Label encoding: Converting categories to numbers
  - Target encoding: Using target variable to encode categories
- Text Data Processing: Cleaning, tokenization, removing stop words
- Date/Time Processing: Extracting features like day of week, month, season

## Feature Engineering (1 Hour)

- Creating New Variables: Combining existing features to create more informative ones
- Domain Knowledge: Using business understanding to create relevant features
- Automated Feature Creation: Polynomial features, interaction terms
- Feature Selection: Identifying the most important variables
- Real example: Creating customer lifetime value from transaction history

## Big Data Tools (4 Hours)

## Apache Spark with PySpark (2 Hours)

- When to Use Spark: Handling datasets too large for single machines
- Spark Architecture: Understanding drivers, executors, and clusters
- PySpark DataFrames: Similar to Pandas but distributed across multiple computers
- Transformations vs Actions: Lazy evaluation and optimization
- Real scenario: Processing millions of customer transactions across multiple servers

Workflow Orchestration with Apache Airflow (1 Hour)

- What is Airflow: Scheduling and monitoring data pipelines
- DAGs (Directed Acyclic Graphs): Defining workflows with dependencies
- Task Management: Handling failures, retries, and notifications
- Real application: Daily ETL pipeline that collects, processes, and loads business data

#### Efficient Data Storage (1 Hour)

- File Formats Comparison:
  - CSV: Human-readable but slow and large
  - Parquet: Columnar storage, fast queries, small size
  - Avro: Schema evolution, good for streaming data
- Compression: Reducing storage costs and improving performance
- Partitioning: Organizing data for fast queries
- Practical impact: Reducing data processing time from hours to minutes

## **Representation of the Property of the Propert**

- 1. End-to-End Data Pipeline: Build a complete pipeline that collects data from 3 different sources (API, database, web scraping), cleans it, and stores it in a data warehouse
- 2. Data Quality Assessment: Create an automated system that profiles incoming data and flags quality issues like missing values, outliers, and inconsistencies
- 3. Big Data Processing Challenge: Process a 1GB+ dataset using PySpark, performing complex aggregations and joins
- 4. API Integration Project: Build a system that automatically collects real-time data from multiple APIs and handles errors gracefully
- 5. Web Scraping Marathon: Create a robust web scraper that can extract structured data from e-commerce sites while respecting rate limits and handling dynamic content

Module 4: Machine Learning & AI (35 Hours)

**@** Learning Objectives:

- Implement supervised learning algorithms for classification and regression
- Apply unsupervised learning techniques for pattern discovery
- Build and train deep neural networks for complex problems
- Develop natural language processing applications
- Evaluate model performance and avoid common pitfalls
- Detailed Topics Covered:

Supervised Learning (12 Hours)

Regression Models (4 Hours)

- Linear Regression:
  - What it does: Predicting continuous values like prices, temperatures, sales
  - How it works: Finding the best line through data points
  - Assumptions: Linear relationship, normal distribution of errors
  - Real example: Predicting house prices based on size, location, age
  - Interpretation: Understanding coefficients and what they mean for business
- Polynomial Regression:
  - When relationships aren't straight lines
  - Creating curved fits for complex patterns
  - Danger of overfitting with too many degrees
- Regularization Techniques:
  - Ridge Regression: Preventing overfitting by penalizing large coefficients
  - Lasso Regression: Automatic feature selection by forcing some coefficients to zero
  - Elastic Net: Combining Ridge and Lasso benefits
  - When to use each: Decision framework based on data characteristics

Classification Models (4 Hours)

- Logistic Regression:
  - What it does: Predicting categories like spam/not spam, buy/don't buy
- © 2025 Udaan Tech Academy. All Rights Reserved.

- How probability works: Converting linear predictions to probabilities
- Decision boundaries: How the model separates different classes
- Multi-class classification: Handling more than two categories

#### Decision Trees:

- Visual Understanding: Tree-like model that's easy to explain to business stakeholders
- How splits are chosen: Information gain, Gini impurity
- Overfitting issues: When trees become too complex
- Business application: Credit approval decisions with explainable rules

#### Random Forest:

- Ensemble Method: Combining many decision trees for better predictions
- Bootstrap Aggregating: Training each tree on different data samples
- Feature Importance: Understanding which variables matter most
- Practical advantage: Works well out-of-the-box with minimal tuning

#### Advanced Models (2 Hours)

- XGBoost and LightGBM:
  - Gradient Boosting: Iteratively improving predictions by learning from mistakes
  - Why they win competitions: Excellent performance on structured data
  - Hyperparameter tuning: Learning rate, tree depth, regularization
  - Real impact: Often 5-10% better accuracy than simpler models
- Support Vector Machines:
  - Concept: Finding the optimal boundary between classes
  - Kernel Trick: Handling non-linear relationships
  - When to use: Small to medium datasets with complex patterns

#### Model Evaluation and Optimization (2 Hours)

- Cross-Validation:
  - Why it matters: Getting realistic estimates of model performance
  - K-fold cross-validation: Training and testing on different data splits

- Stratified sampling: Ensuring balanced representation in each fold
- Evaluation Metrics:
  - Accuracy: Simple but can be misleading
  - Precision and Recall: Understanding trade-offs in classification
  - F1-Score: Balancing precision and recall
  - ROC Curves: Visualizing model performance across thresholds
  - Business metrics: Connecting model performance to business value
- Hyperparameter Tuning:
  - Grid Search: Systematically trying parameter combinations
  - Random Search: More efficient for high-dimensional spaces
  - Bayesian Optimization: Smart search using previous results

Unsupervised Learning (8 Hours)

Clustering Algorithms (4 Hours)

- K-Means Clustering:
  - What it does: Grouping similar data points together
  - How it works: Finding cluster centers that minimize distances
  - Choosing K: Elbow method, silhouette analysis
  - Business applications: Customer segmentation, market research
  - Limitations: Assumes spherical clusters, sensitive to outliers
- Hierarchical Clustering:
  - Dendrograms: Tree-like visualization of cluster relationships
  - Agglomerative vs Divisive: Bottom-up vs top-down approaches
  - Linkage methods: Single, complete, average, Ward
  - When to use: When you need hierarchy of clusters
- DBSCAN:
  - Density-based clustering: Finds clusters of varying shapes
  - Handles outliers: Automatically identifies noise points
  - Parameters: Epsilon (neighborhood size) and minimum points

Real scenario: Identifying fraud patterns in transaction data

## Dimensionality Reduction (4 Hours)

- Principal Component Analysis (PCA):
  - Purpose: Reducing data complexity while preserving information
  - How it works: Finding the most important directions in data
  - Explained variance: How much information each component captures
  - Visualization: Reducing high-dimensional data to 2D plots
  - Business value: Simplifying complex datasets for analysis
- t-SNE (t-Distributed Stochastic Neighbor Embedding):
  - Specialized for visualization: Creating 2D maps of complex data
  - Non-linear reduction: Capturing complex relationships
  - Parameters: Perplexity, learning rate, iterations
  - Limitations: Computational cost, not suitable for new data prediction
- UMAP (Uniform Manifold Approximation and Projection):
  - Modern alternative to t-SNE: Faster and preserves more global structure
  - Better for new data: Can transform new points using existing model
  - Applications: Exploring customer behavior patterns, gene expression data

#### Deep Learning (10 Hours)

#### Neural Network Fundamentals (3 Hours)

- Biological Inspiration: How artificial neurons mimic brain cells
- Perceptron: The simplest neural network unit
- Multi-layer Networks: Hidden layers and their purpose
- Activation Functions: ReLU, sigmoid, tanh when to use each
- Backpropagation: How networks learn from mistakes
- Gradient Descent: Optimizing network weights
- Practical example: Building a simple network to classify handwritten digits

## Deep Learning with TensorFlow/Keras (3 Hours)

Framework Introduction: Why use high-level frameworks

- Model Architecture: Defining layers, connections, and flow
- Compilation: Choosing optimizers, loss functions, metrics
- Training Process: Epochs, batches, validation splits
- Callbacks: Early stopping, learning rate scheduling, model checkpointing
- Real project: Building a customer churn prediction model

#### Convolutional Neural Networks (CNNs) (2 Hours)

- Image Data Challenges: Why traditional methods fail on images
- Convolution Operation: Detecting features like edges, textures
- Pooling Layers: Reducing spatial dimensions while preserving features
- CNN Architecture: Typical structure from input to output
- Transfer Learning: Using pre-trained models like ResNet, VGG
- Business applications: Quality control in manufacturing, medical image analysis

### Recurrent Neural Networks (RNNs) (2 Hours)

- Sequential Data: Handling time series, text, speech
- Memory Mechanism: How RNNs remember previous inputs
- LSTM and GRU: Solving the vanishing gradient problem
- Time Series Forecasting: Predicting future values from historical data
- Text Generation: Creating human-like text
- Practical example: Sales forecasting using historical data and external factors

## Natural Language Processing (5 Hours)

## *Text Preprocessing (1.5 Hours)*

- Text Cleaning: Removing HTML tags, special characters, standardizing case
- Tokenization: Breaking text into words, sentences, or subwords
- Stop Word Removal: Filtering out common words like "the", "and"
- Stemming vs Lemmatization: Reducing words to root forms
- Handling Different Languages: Unicode, encoding issues
- Real data: Processing customer reviews, social media posts, support tickets

#### Feature Extraction from Text (1.5 Hours)

- Bag of Words: Counting word occurrences
- TF-IDF (Term Frequency-Inverse Document Frequency): Weighting words by importance
- N-grams: Capturing word sequences and context
- Word Embeddings:
  - Word2Vec: Learning word representations from context
  - GloVe: Global vectors for word representation
  - FastText: Handling out-of-vocabulary words
- Contextual Embeddings: BERT, RoBERTa for better understanding

## Text Classification and Analysis (2 Hours)

- Sentiment Analysis:
  - Binary classification: Positive vs negative sentiment
  - Multi-class: Happy, sad, angry, neutral emotions
  - Aspect-based: Sentiment toward specific product features
  - Business impact: Brand monitoring, customer satisfaction measurement
- Topic Modeling:
  - Latent Dirichlet Allocation (LDA): Discovering hidden topics in documents
  - Non-negative Matrix Factorization: Alternative topic modeling approach
  - Topic coherence: Evaluating topic quality
  - Applications: Organizing large document collections, content recommendation
- Advanced NLP with Transformers:
  - Hugging Face ecosystem: Pre-trained models for various tasks
  - BERT for classification: Fine-tuning for specific domains
  - Named Entity Recognition: Extracting people, places, organizations from text
  - Question Answering: Building chatbots and FAQ systems
- **%** Practical Exercises:

- Regression Competition: Predict house prices using advanced feature engineering and ensemble methods, competing for lowest RMSE
- 2. Multi-class Classification: Build an image classifier for fashion items using CNNs with data augmentation and transfer learning
- 3. Customer Segmentation: Use multiple clustering algorithms to segment customers and create actionable marketing strategies
- 4. Time Series Forecasting: Predict sales using LSTM networks, incorporating seasonal patterns and external factors
- 5. NLP Sentiment System: Build end-to-end sentiment analysis for product reviews, including data collection, preprocessing, and model deployment
- 6. Complete ML Pipeline: From raw data to deployed model, including monitoring, retraining, and A/B testing capabilities

Module 5: Advanced Analytics & Al Trends (20 Hours)

- **@** Learning Objectives:
  - Implement automated machine learning workflows for rapid prototyping
  - Create explainable AI solutions for business stakeholders
  - Integrate generative AI capabilities into data science workflows
  - Deploy models for edge computing and real-time applications
  - Build responsible AI systems with bias detection and mitigation
- Detailed Topics Covered:

Automated Machine Learning (AutoML) (5 Hours)

AutoML Fundamentals (1.5 Hours)

- What is AutoML: Automating the machine learning pipeline from data preparation to model deployment
- Why AutoML Matters:
  - Democratizing ML for non-experts
  - Reducing development time from weeks to hours
  - Standardizing best practices across teams
- AutoML Components:

- Automated Data Preprocessing: Missing value imputation, outlier detection, feature scaling
- Feature Engineering: Automatic creation of polynomial features, interactions, transformations
- Model Selection: Testing multiple algorithms (Random Forest, XGBoost, Neural Networks) automatically
- Hyperparameter Optimization: Finding optimal settings without manual tuning
- Business Impact: Enabling business analysts to build predictive models without extensive ML knowledge

Open Source AutoML Tools (2 Hours)

- Auto-sklearn:
  - How it works: Automated scikit-learn pipeline optimization
  - Ensemble building: Combining multiple models for better performance
  - Meta-learning: Using knowledge from previous datasets to speed up optimization
  - Practical example: Building a customer churn model in 30 minutes instead of 3 days
- TPOT (Tree-based Pipeline Optimization Tool):
  - Genetic Programming: Evolutionary approach to pipeline optimization
  - Pipeline representation: Python code generation for reproducibility
  - Feature construction: Automatic creation of new features
  - Real application: Optimizing marketing campaign response prediction
- H2O AutoML:
  - Distributed computing: Handling large datasets across multiple machines
  - Model interpretability: Built-in explanations for automated models
  - Production deployment: Direct model serving capabilities
  - Use case: Processing millions of customer records for risk assessment

Pipeline Automation (1.5 Hours)

- MLOps Integration: Incorporating AutoML into production workflows
- Automated Model Monitoring: Detecting when models need retraining
- A/B Testing: Automatically comparing new AutoML models against existing ones
- Continuous Integration: Version control and testing for automated pipelines
- Real scenario: Daily automated model retraining for fraud detection systems

## Explainable AI (XAI) (5 Hours)

Model Interpretability Concepts (1.5 Hours)

- Why Explainability Matters:
  - Regulatory compliance (GDPR, financial regulations)
  - Building trust with stakeholders
  - Debugging model decisions
  - Ensuring fairness and avoiding bias
- Types of Explanations:
  - Global explanations: How the model works in general
  - Local explanations: Why the model made a specific prediction
  - Counterfactual explanations: What would need to change for a different outcome
- Model-Agnostic vs Model-Specific: Techniques that work with any model vs those designed for specific algorithms

## SHAP Values (1.5 Hours)

- Shapley Values Theory: Game theory approach to feature importance
- SHAP Implementation:
  - Tree SHAP: Fast explanations for tree-based models
  - Deep SHAP: Explanations for neural networks
  - Kernel SHAP: Model-agnostic explanations
- Visualization Techniques:
  - Waterfall plots: Showing how features contribute to individual predictions
  - Summary plots: Global feature importance across all predictions
  - Dependence plots: How feature values affect predictions

 Business Application: Explaining loan approval decisions to customers and regulators

LIME and Other Techniques (1.5 Hours)

- LIME (Local Interpretable Model-agnostic Explanations):
  - How LIME works: Creating local linear approximations around specific predictions
  - Text explanations: Highlighting important words in document classification
  - Image explanations: Showing which parts of an image influenced the decision
  - Tabular data: Explaining predictions for structured datasets
- Model-Specific Explanations:
  - Decision Tree Visualization: Creating business-friendly rule explanations
  - Linear Model Coefficients: Understanding feature weights and directions
  - Feature Importance: Ranking variables by their impact on predictions

Bias Detection and Fairness (0.5 Hours)

- Types of Bias:
  - Historical bias in training data
  - Representation bias (undersampled groups)
  - Evaluation bias (different performance across groups)
- Fairness Metrics: Equal opportunity, demographic parity, equalized odds
- Bias Mitigation: Pre-processing, in-processing, and post-processing techniques
- Real example: Ensuring hiring models don't discriminate against protected groups

Generative AI for Data Science (5 Hours)

Large Language Models for Analysis (2 Hours)

- Open Source LLM Integration:
  - Ollama setup: Running Llama 2, Mistral locally for data analysis
  - Hugging Face models: Code Llama for automated code generation

- Local API setup: Creating private AI endpoints for sensitive data
- Code Generation Applications:
  - Automated EDA: Generating exploratory data analysis code
  - Feature engineering: Al-suggested new variables
  - Model documentation: Automated reporting and explanation generation
  - Data cleaning: Al-assisted identification of data quality issues
- Practical example: Using AI to automatically generate data analysis reports

## Prompt Engineering for Data Tasks (1.5 Hours)

- Effective Prompts for Data Science:
  - Code generation prompts: Specific, context-aware requests
  - Explanation prompts: Getting AI to explain complex statistical concepts
  - Debugging prompts: Using AI to find and fix code issues
- Chain-of-Thought for Analysis:
  - Breaking complex analytical tasks into steps
  - Iterative refinement of analysis approach
  - Quality checking Al-generated insights
- Real application: Building a conversational interface for business users to query data

#### Synthetic Data Generation (1.5 Hours)

- When to Use Synthetic Data:
  - Privacy protection for sensitive datasets
  - Data augmentation for small datasets
  - Testing and development without real data exposure
- Generation Techniques:
  - Statistical methods: Sampling from fitted distributions
  - GANs (Generative Adversarial Networks): Deep learning for realistic data
  - Variational Autoencoders: Probabilistic data generation
  - LLM-based generation: Text and structured data creation

- Quality Assessment: Ensuring synthetic data maintains statistical properties
- Business case: Creating realistic customer data for testing new features

Edge Computing & Real-time Analytics (5 Hours)

Stream Processing Fundamentals (2 Hours)

- Batch vs Stream Processing: Understanding when each is appropriate
- Apache Kafka Setup:
  - Producers and consumers: Sending and receiving real-time data
  - Topics and partitions: Organizing data streams
  - Fault tolerance: Ensuring reliable data delivery
- Stream Processing with Python:
  - Kafka-Python: Building producers and consumers
  - Real-time transformations: Processing data as it arrives
  - Windowing: Aggregating data over time periods
- Use cases: Real-time fraud detection, IoT sensor monitoring, social media analysis

## Edge Al Deployment (1.5 Hours)

- Edge Computing Concepts: Processing data closer to where it's generated
- Model Optimization for Edge:
  - Quantization: Reducing model size and computational requirements
  - Pruning: Removing unnecessary neural network connections
  - Knowledge distillation: Creating smaller models that mimic larger ones
- Deployment Platforms:
  - Raspberry Pi: Low-cost edge computing
  - NVIDIA Jetson: GPU-accelerated edge devices
  - Mobile deployment: TensorFlow Lite for smartphones
- Real application: Quality control cameras in manufacturing that process images locally

# 🎇 Practical Exercises:

- 1. AutoML Competition: Compare manually tuned models vs AutoML solutions on the same dataset, measuring both performance and development time
- 2. Explainability Dashboard: Build a web interface that explains model predictions to non-technical stakeholders using SHAP and LIME
- 3. Al-Assisted Analysis: Use local LLMs to generate complete data analysis workflows from natural language descriptions
- 4. Real-time Fraud Detection: Build a streaming system that processes

Module 6: Business Intelligence & Visualization (10 Hours)

- **@** Learning Objectives:
  - Create interactive and compelling data visualizations for stakeholders
  - Build comprehensive business intelligence dashboards
  - Master storytelling techniques for data-driven presentations
  - Implement automated reporting systems for business metrics
  - Design user-friendly interfaces for non-technical stakeholders
- Detailed Topics Covered:

Advanced Data Visualization (5 Hours)

Interactive Visualization Libraries (2 Hours)

- Plotly for Python:
  - Interactive charts: Hover effects, zooming, panning, filtering
  - Chart types: 3D plots, geographic maps, statistical charts
  - Dashboard creation: Combining multiple charts with cross-filtering
  - Animation: Time-based animations for showing trends over time
  - Real example: Creating an interactive sales dashboard where users can drill down from yearly to daily data
- Bokeh for Web Applications:
  - Server applications: Building interactive web apps with Python backend
  - Streaming data: Real-time updating visualizations
  - Custom interactions: Building unique user interface elements
  - Widget integration: Sliders, dropdowns, buttons for user control

Business application: Customer behavior analytics with real-time filtering

Advanced Chart Types and Techniques (1.5 Hours)

- Statistical Visualizations:
  - Box plots and violin plots: Understanding data distributions and outliers
  - Correlation matrices: Heatmaps showing relationships between variables
  - Pair plots: Matrix of scatter plots for multivariate analysis
  - Regression plots: Visualizing model fits and confidence intervals
- Geographic Visualization:
  - Folium for mapping: Creating interactive maps with markers, polygons, heatmaps
  - Choropleth maps: Coloring regions based on data values
  - Geospatial analysis: Working with coordinates, boundaries, and spatial relationships
  - Real scenario: Visualizing sales performance across different geographic regions

Design Principles and Storytelling (1.5 Hours)

- Visual Design Best Practices:
  - Color theory: Choosing palettes that are accessible and meaningful
  - Typography: Readable fonts and appropriate sizing
  - Layout principles: White space, alignment, hierarchy
  - Accessibility: Ensuring visualizations work for colorblind users
- Data Storytelling:
  - Narrative structure: Beginning, middle, end for data presentations
  - Audience consideration: Tailoring visualizations to stakeholder needs
  - Action-oriented insights: Moving from "what happened" to "what should we do"
  - Annotation techniques: Highlighting key insights and takeaways
- Avoiding Common Mistakes:
  - Misleading scales: Ensuring accurate representation of data

- Chart junk: Removing unnecessary elements that distract from insights
- Overcomplicating: Keeping visualizations focused and clear

Business Intelligence Tools (5 Hours)

Apache Superset Deep Dive (2.5 Hours)

- Installation and Setup:
  - Docker deployment: Quick setup for development and production
  - Database connections: Connecting to PostgreSQL, MySQL, SQLite
  - User management: Setting up roles, permissions, and access control
  - Configuration: Customizing appearance, features, and security settings
- Dashboard Creation:
  - Chart builder: Using the intuitive interface to create visualizations
  - SQL Lab: Writing custom queries for complex analysis
  - Dashboard composition: Arranging charts, filters, and text elements
  - Cross-filtering: Creating interactive dashboards where charts affect each other
  - Drill-down capabilities: Enabling users to explore data at different levels of detail
- Advanced Features:
  - Alerts and reports: Automated notifications when metrics exceed thresholds
  - Row-level security: Ensuring users only see data they're authorized to view
  - Custom visualizations: Adding new chart types and plugins
  - Real implementation: Building a complete executive dashboard for a retail company

Metabase Implementation (1.5 Hours)

- Quick Setup: Getting Metabase running with minimal configuration
- Question Builder: Creating analyses without writing SQL
- Dashboard Creation: Building interactive dashboards for business users

- Pulse: Automated email reports and Slack notifications
- User Experience: Designing dashboards that non-technical users can navigate
- Business case: Self-service analytics platform for marketing team

#### Automated Reporting Systems (1 Hour)

- Scheduled Reports:
  - Cron jobs: Scheduling Python scripts for regular report generation
  - Email automation: Sending reports via SMTP
  - PDF generation: Creating professional reports with matplotlib and ReportLab
  - Cloud scheduling: Using free alternatives to cron for reliability
- Real-time Monitoring:
  - KPI dashboards: Always-on displays showing critical business metrics
  - Alert systems: Notifications when metrics exceed thresholds
  - Health checks: Monitoring dashboard and data pipeline performance

## **%** Practical Exercises:

- Executive Dashboard Project: Build a comprehensive C-level dashboard showing KPIs, trends, and drill-down capabilities
- 2. Geographic Sales Analysis: Create an interactive map showing sales performance by region with time-based filtering
- 3. Automated Report System: Set up daily, weekly, and monthly automated reports that email stakeholders
- 4. Self-Service Analytics: Configure Metabase or Superset for non-technical users to create their own analyses
- 5. Data Storytelling Presentation: Create a compelling presentation that tells a complete story with data and visualizations

## Capstone Projects (15 Hours)

# **@** Learning Objectives:

- Integrate all learned skills into comprehensive real-world projects
- Demonstrate end-to-end data science project execution
- © 2025 Udaan Tech Academy. All Rights Reserved.

- Build a professional portfolio showcasing diverse skills
- Practice presenting technical work to business stakeholders
- Prepare for job interviews with concrete project examples
- Detailed Topics Covered:

Project 1: End-to-End ML Pipeline (5 Hours)

Project Scope and Planning (0.5 Hours)

- Problem Definition: Choose a real business problem (customer churn, price prediction, demand forecasting)
- Success Metrics: Define how to measure project success
- Timeline Planning: Breaking project into manageable milestones
- Stakeholder Identification: Understanding who will use the results

Data Collection and Preparation (1.5 Hours)

- Multi-source Data Integration:
  - Web scraping for additional features
  - API integration for external data
  - Database extraction using SQL
- Comprehensive Data Cleaning:
  - Handling missing values with domain knowledge
  - Outlier detection and treatment strategies
  - Feature engineering based on business understanding
- Data Quality Documentation: Creating data dictionaries and quality reports

Model Development and Optimization (1.5 Hours)

- Baseline Model: Starting with simple model for comparison
- Advanced Modeling: Implementing multiple algorithms and ensemble methods
- AutoML Comparison: Testing automated solutions against manual approaches
- Hyperparameter Tuning: Using advanced optimization techniques
- Model Validation: Comprehensive testing including edge cases

Production Deployment (1 Hour)

- API Development: Creating REST endpoints using FastAPI
- Containerization: Docker setup for consistent deployment
- Monitoring Implementation: Logging, error handling, performance metrics
- Documentation: Technical documentation for developers and user guides

## Business Impact Assessment (0.5 Hours)

- Performance Evaluation: Model accuracy, speed, resource usage
- Business Value Calculation: ROI estimation, cost-benefit analysis
- Stakeholder Presentation: Communicating results to business leaders

## Project 2: Real-time Analytics Dashboard (5 Hours)

## System Architecture Design (0.5 Hours)

- Data Flow Planning: From sources through processing to visualization
- Technology Stack Selection: Choosing appropriate tools for each component
- Scalability Considerations: Planning for increased data volume and users

### Live Data Integration (1.5 Hours)

- Stream Processing Setup:
  - Kafka configuration for real-time data ingestion
  - Python consumers for data processing
  - Error handling and data quality checks
- Database Integration:
  - Real-time data storage strategies
  - Indexing for fast queries
  - Data retention policies

#### Interactive Dashboard Development (2 Hours)

- Frontend Development:
  - Streamlit or Dash for Python-based dashboards
  - Chart selection based on data types and user needs
  - Responsive design for different screen sizes
- User Experience Design:

- Intuitive navigation and filtering
- Performance optimization for smooth interactions
- Mobile-friendly responsive layouts

#### Performance Optimization (0.5 Hours)

- · Query Optimization: Efficient database queries for fast loading
- Caching Strategies: Redis for frequently accessed data
- Load Testing: Ensuring dashboard works under heavy usage

## Deployment and Monitoring (0.5 Hours)

- Production Deployment: Cloud or local server setup
- Health Monitoring: Uptime checks, performance metrics
- User Access Management: Authentication and authorization

## Project 3: AI-Powered Business Intelligence (5 Hours)

## Natural Language Query Interface (2 Hours)

- Intent Recognition: Understanding user questions about data
- Query Translation: Converting natural language to SQL or pandas operations
- Context Management: Handling follow-up questions and conversation flow
- Response Generation: Creating natural language explanations of results

#### Automated Insight Generation (1.5 Hours)

- Pattern Detection:
  - Statistical anomaly detection
  - Trend identification and significance testing
  - Correlation discovery and ranking
- Insight Prioritization: Ranking insights by business importance
- Automated Explanation: Generating natural language descriptions of findings

## Intelligent Reporting (1 Hour)

- Dynamic Report Generation: Creating reports based on current data
- Visualization Selection: Al-assisted chart type recommendations
- Narrative Generation: Writing summaries and key takeaways

- Personalization: Tailoring reports to different user roles and interests
   Integration and Testing (0.5 Hours)
  - System Integration: Connecting all components seamlessly
  - User Acceptance Testing: Getting feedback from potential business users
  - Performance Benchmarking: Measuring system speed and accuracy
  - Documentation: User manuals and technical specifications

## Representation of the Practical Assessment Framework:

## Portfolio Development Requirements:

- 1. GitHub Repository: Well-documented code with README files
- 2. Live Demonstrations: Working deployments accessible online
- 3. Technical Documentation: Architecture diagrams, API documentation
- 4. Business Case Studies: Problem statements, solutions, and impact
- 5. Video Presentations: 5-minute demos explaining each project

## Presentation Skills Development:

- 1. Technical Presentations: Explaining methodology to data science peers
- 2. Business Presentations: Communicating value to non-technical stakeholders
- 3. Problem-Solving Documentation: Showing how challenges were overcome
- 4. Code Reviews: Peer evaluation of code quality and best practices

#### Interview Preparation:

- 1. Technical Questions: Common data science interview problems
- 2. Project Deep Dives: Detailed discussions about implementation choices
- 3. Business Impact Stories: Quantifiable results and lessons learned
- 4. Portfolio Showcase: Professional presentation of all work

#### Real-World Readiness:

- Production Code Standards: Following industry best practices
- Collaboration Skills: Working effectively in team environments
- Continuous Learning: Staying updated with latest tools and techniques
- Ethical Considerations: Responsible AI and data privacy awareness

