

Data Science with AI

Complete Beginner to Professional Guide

What it is • How it helps • End-to-End Job Process • Daily Responsibilities



bytelearned.online | info@bytelearned.online | 7396633051

5 Major Topics	8-Step Roadmap	Job Guide Fresher+Senior	50+ Tools	Salary Benchmarks	Industry Use Cases
----------------	----------------	--------------------------	-----------	-------------------	--------------------

01 | WHAT IS DATA SCIENCE WITH AI?

Data Science is the field of extracting meaningful insights, patterns, and knowledge from raw data using mathematics, statistics, and programming. When combined with Artificial Intelligence (AI), it goes a step further — machines not only find patterns in data but also **learn from it, make predictions, and take decisions automatically** without being explicitly programmed for every task.

THE 3 PILLARS OF DATA SCIENCE WITH AI

■ DATA SCIENCE Collecting, cleaning, exploring, and analysing data to find patterns and answer business questions using statistics and visualisation.	■ MACHINE LEARNING Training algorithms on historical data so computers can automatically improve their predictions on new unseen data without manual programming.	■ ARTIFICIAL INTELLIGENCE Building systems that simulate human intelligence — understanding language, recognising images, making decisions, and solving complex problems.
---	---	---

KEY CONCEPTS YOU MUST KNOW

- **Statistics & Math** The backbone of data science. Mean, median, probability, distributions, hypothesis testing — the language data speaks.
- **Python Programming** The #1 language for data science. Used for data cleaning, analysis, model building, and automation of everything.

■ **SQL & Databases** Every company stores data in databases. You must query, join, and aggregate data using SQL to answer business questions.

■ **Data Visualisation** Turning numbers into charts and dashboards that non-technical teams can understand and act on immediately.

■ **Machine Learning** Algorithms that learn from data — Linear Regression, Decision Trees, Random Forests, Neural Networks, XGBoost.

■ **Deep Learning** Advanced ML using Neural Networks with many layers. Powers image recognition, speech, NLP, and generative AI.

■ **Natural Language Processing** Teaching machines to understand, read, and generate human language. Powers ChatGPT, Alexa, Google Translate.

■ **Big Data & Cloud** Processing massive datasets (terabytes/petabytes) using cloud platforms like AWS, Azure, Google Cloud with Spark/Hadoop.

■ **MLOps & Deployment** The process of taking a trained model and deploying it into production so real users and systems can use it.

■ **Generative AI & LLMs** The latest frontier — models like GPT-4, Gemini, and LLaMA that generate text, images, code, and more.

02 | HOW DATA SCIENCE WITH AI IS USEFUL TODAY

Data Science with AI is no longer a niche academic field. It is the engine powering every major industry today — from the Netflix recommendation that kept you watching, to the fraud alert on your credit card, to the AI doctor helping diagnose cancer. Every sector is being transformed by data and AI.

REAL-WORLD APPLICATIONS BY INDUSTRY

■ Healthcare & Medicine AI diagnoses cancer from X-rays faster than doctors, predicts patient readmission risks, enables drug discovery in months instead of years, and powers personalised medicine.	■ Banking & Finance Fraud detection in milliseconds on every UPI/card transaction, credit scoring for loans, algorithmic stock trading, robo-advisors, and risk assessment systems.
■ E-Commerce & Retail Amazon and Flipkart recommendation engines (35% of all revenue), demand forecasting, dynamic pricing, inventory optimisation, and customer churn prediction.	■ Automotive & Transport Self-driving cars (Tesla, Waymo), route optimisation for Zomato/Swiggy delivery, predictive maintenance for vehicle fleets, and traffic management systems.
■ Social Media & Tech Facebook/Instagram feed ranking, Twitter/X trending algorithm, YouTube watch time maximisation, content moderation, and ad targeting systems.	■ Education Personalised learning paths for students, automated essay grading, dropout risk prediction, chatbots for student support, and adaptive test platforms.
■ Agriculture Crop yield prediction using satellite imagery and weather data, soil health monitoring with sensors, pest detection using computer vision, and irrigation AI.	■ Manufacturing Predictive maintenance (catching machine failures before they happen), quality control using computer vision, supply chain optimisation, and robot automation.
■ Weather & Climate More accurate 10-day weather forecasts using deep learning, climate change modelling, disaster prediction systems, and renewable energy output forecasting.	■ Gaming & Entertainment AI game opponents that adapt to your skill level, content recommendation for OTT platforms, procedural game world generation, and real-time translation.

DATA SCIENCE & AI IN NUMBERS — 2025/26

\$1.8 Trillion+	Global AI market value projected by 2030 (McKinsey)
97 Million jobs	New AI/Data roles expected globally by 2025 (World Economic Forum)
Rs. 14–45 LPA	Average data scientist salary range in India in 2026
2.5 Quintillion bytes	Amount of new data created every single day globally
80% of work	Proportion of a data scientist's time spent on data cleaning
10x ROI	Average return on investment companies get from AI/ML projects
44% of companies	Indian enterprises planning to increase AI investment in 2026

India #2

India is the 2nd largest AI talent pool in the world after USA

03 | START TO END — DATA SCIENCE WITH AI JOB PROCESS

The data science job process covers everything from how a real data science project begins in a company to how results are delivered and used. Understanding this end-to-end workflow is what separates a candidate who just knows Python from one who gets hired and delivers real value.

THE 8-STEP DATA SCIENCE PROJECT WORKFLOW

1 Business Problem Definition

Week 1–2

Every data science project starts with a business question — NOT with data or code. You meet with stakeholders to understand what problem they are trying to solve, what decisions the model needs to drive, and what success looks like.

- **Stakeholder meetings** — Understand the business goal: "Reduce customer churn by 20%", "Detect fraud in real-time", "Predict next month sales"
- **Define success metrics** — Agree on KPIs: accuracy, F1 score, revenue impact, time savings. Know what "good enough" means before you start.
- **Scope & feasibility** — Assess if enough data exists, if the timeline is realistic, and if the technical solution is achievable given constraints.

■ *Tip: Most data science projects fail because the problem was not defined clearly. Spend more time here than you think you need.*

2 Data Collection & Acquisition

Week 2–3

Once the problem is defined, you identify where the data lives and collect it. Data can come from multiple sources — databases, APIs, web scraping, sensor feeds, third-party vendors, or surveys.

- **Internal databases** — Query SQL/NoSQL databases for historical transaction data, user logs, CRM records, ERP data.
- **APIs & external sources** — Pull real-time data from REST APIs — weather, social media, market prices, government open datasets.
- **Web scraping** — Use BeautifulSoup or Scrapy to collect data from websites when no API is available.
- **Data warehouses** — Access company data lakes (AWS S3, Snowflake, BigQuery) for large-scale historical datasets.

■ *Legal Note: Always verify data ownership and privacy compliance (GDPR, India PDPB) before collecting any user data.*

3 Data Cleaning & Preprocessing

Week 3–5
— Takes
80% of
Time

Raw data is almost always messy, incomplete, inconsistent, or incorrect. This step transforms raw data into a clean, structured dataset that algorithms can actually learn from. It is unglamorous but the most critical step in the entire process.

- **Handle missing values** — Decide whether to drop, fill with mean/median, interpolate, or flag missing data based on the feature importance.
- **Remove duplicates** — Identify and remove duplicate records that would bias the model's training and skew results.
- **Fix inconsistencies** — Standardise date formats, capitalisation, units, category names — e.g. "Male"/"male"/"M" all mean the same thing.
- **Outlier detection** — Identify extreme values using IQR or Z-score. Decide to remove, cap, or investigate them based on context.
- **Feature engineering** — Create new meaningful features from existing data — e.g. "days since last purchase" from a timestamp column.
- **Data scaling & encoding** — Normalise numerical features (MinMaxScaler) and encode categoricals (One-Hot, Label Encoding) for ML algorithms.

■ *Tip: In real jobs, you will spend 60–80% of your project time on data cleaning. Master Pandas and become extremely comfortable with messy, real-world data.*

4 Exploratory Data Analysis (EDA)

Week 4–5

EDA is the detective work of data science. You explore the data visually and statistically to understand its structure, find patterns, spot anomalies, test assumptions, and generate hypotheses before building any model.

- **Univariate analysis** — Understand each feature individually — distributions, histograms, box plots, value counts for categoricals.
- **Bivariate analysis** — Explore relationships between pairs of features — scatter plots, correlation heatmaps, cross-tabulations.
- **Data visualisation** — Use Matplotlib, Seaborn, and Plotly to create insightful charts that reveal hidden patterns in the data.
- **Statistical testing** — Apply hypothesis tests (t-test, chi-square, ANOVA) to confirm whether observed patterns are statistically significant.
- **EDA report** — Document your key findings and share with stakeholders before proceeding to modelling. Get alignment on assumptions.

■ *Tip: EDA is where data science becomes an art. The best data scientists find insights that others miss because they ask better questions of the data.*

5 Model Building & Training

Week 5–7

Now you actually build the machine learning model. You select appropriate algorithms, split data into train/test sets, train multiple models, and compare their performance to find the best one.

- ✂ **Train-test split** — Split data (typically 80/20 or 70/30) to ensure the model is evaluated on data it has never seen during training.
- **Baseline model** — Always start with a simple model (Logistic Regression, Linear Regression) as a baseline before trying complex ones.
- **Algorithm selection** — Try multiple algorithms: Random Forest, XGBoost, SVM, Neural Networks. The best model depends on your data and problem.
- **Hyperparameter tuning** — Use GridSearchCV or Optuna to find the optimal settings for each algorithm to improve performance.
- **Cross-validation** — Use k-fold cross-validation to get a reliable performance estimate and avoid overfitting to the training data.
- **Model evaluation** — Measure performance with the right metrics — Accuracy, Precision, Recall, F1, AUC-ROC (classification), RMSE, MAE (regression).

■ *Warning: Never trust a model that performs too well (99%+ accuracy). Likely overfitting or data leakage. Always validate on completely unseen data.*

A model that cannot be explained is dangerous and often unused. This step involves understanding WHY the model makes the predictions it does and translating technical results into business insights.

- **Feature importance** — Identify which features the model relies on most using SHAP values, feature importance plots, or LIME.

- **Error analysis** — Study where and why the model makes mistakes. These patterns reveal data quality issues or missing features.

- **Business insights** — Translate model outputs into plain language: "Customers who haven't purchased in 60+ days are 3x more likely to churn."

- **Stakeholder presentation** — Create clear visualisations and a business-focused presentation. Avoid technical jargon with non-technical audiences.

- *Tip: Explainability is increasingly required by law (GDPR, RBI) for AI systems in finance and healthcare. SHAP is the industry standard tool.*

7 Model Deployment & MLOps

Week 8–10

A model sitting in a Jupyter notebook has zero business value. Deployment is the process of packaging your model into a production system that real users, applications, and other services can actually call and use in real-time or batch mode.

- **Containerisation** — Package the model with Docker so it runs identically in any environment — local, staging, or production cloud.
- **API development** — Build a REST API using FastAPI or Flask that exposes the model as an endpoint any application can call.
- **Cloud deployment** — Deploy to AWS (SageMaker), Azure (ML Studio), or Google Cloud (Vertex AI) for scalable, managed model serving.
- **CI/CD pipeline** — Automate model testing and deployment using GitHub Actions, Jenkins, or Azure DevOps pipelines.
- **Model registry** — Track model versions in MLflow or W&B; so you can roll back to previous versions if a new deployment underperforms.

■ *Career Advice: Most data science courses skip deployment. But in industry, a data scientist who can deploy their own models is worth 2x the salary. Learn MLOps.*

8 Monitoring, Maintenance & Retraining

Ongoing

Deployed models degrade over time as real-world data patterns shift (called "data drift" or "model drift"). This final step involves continuously monitoring model performance and retraining when it falls below acceptable thresholds.

- **Performance monitoring** — Track model accuracy, latency, and prediction distributions in production using dashboards in Grafana or CloudWatch.
- **Data drift detection** — Monitor for shifts in input data distributions using Evidently AI or WhyLabs that signal the model is becoming stale.
- **Automated retraining** — Set up pipelines that automatically retrain the model on fresh data when performance drops below a threshold.
- **A/B testing** — Test new model versions against the current production model using A/B or shadow deployment before full rollout.
- **Reporting & governance** — Maintain model cards, fairness assessments, and audit logs required by regulators and company AI governance policies.

■ *Tip: Real data scientists spend as much time maintaining existing models as building new ones. Operational excellence is a superpower in this field.*

04 | DAILY JOB RESPONSIBILITIES

FRESHER (0–1 Year) — Junior Data Analyst / Junior Data Scientist

Salary	Rs. 4 LPA – Rs. 8 LPA	Role	Junior Data Analyst / ML Engineer Trainee
--------	-----------------------	------	---

1	Data cleaning & preprocessing Spend most of your day cleaning messy real-world datasets in Python (Pandas, NumPy). Fix nulls, outliers, duplicates, and inconsistencies. This is 80% of your work initially.
2	SQL querying Write SQL queries daily to extract data from company databases, join multiple tables, aggregate metrics, and create data pulls for senior team members and business teams.
3	Exploratory Data Analysis (EDA) Perform statistical analysis and create visualisations (Matplotlib, Seaborn, Tableau) to understand datasets and communicate findings in team meetings.
4	Build and test ML models Implement basic ML models under senior guidance using Scikit-learn. Run experiments, track results in spreadsheets or MLflow, and compare algorithm performance.
5	Create dashboards & reports Build BI dashboards in Power BI, Tableau, or Looker for business stakeholders. Automate weekly and monthly reporting that previously took hours to do manually.
6	Data pipeline support Help maintain existing ETL data pipelines. Monitor for failures, fix broken data feeds, and ensure data arrives correctly into the warehouse on schedule.
7	Write documentation Document data dictionaries, analysis notebooks, and model experiments so team members can understand and reproduce your work without asking you directly.
8	Attend standups & stakeholder meetings Participate in daily 15-minute standups, sprint planning, and demo sessions. Present your analysis findings in simple business language to non-technical stakeholders.

Tools used daily: Python (Pandas, NumPy, Scikit-learn), SQL, Jupyter Notebook, Tableau/Power BI, Excel, Git, VS Code

MID-LEVEL (2–4 Years) — Data Scientist / ML Engineer

Salary	Rs. 12 LPA – Rs. 28 LPA	Role	Data Scientist / Senior ML Engineer
--------	-------------------------	------	-------------------------------------

1	Own end-to-end ML projects Lead complete data science projects from problem definition to deployment. Make independent decisions on data strategy, algorithm choice, and architecture design.
---	---

2	Advanced model development Build and tune complex models — Gradient Boosting (XGBoost, LightGBM), Deep Learning (TensorFlow/PyTorch), NLP models (Transformers/BERT), Time Series forecasting.
3	Feature engineering & selection Design creative features from raw data that dramatically improve model performance. Conduct rigorous feature selection to remove noise and reduce overfitting.
4	MLOps & model deployment Deploy models as REST APIs using FastAPI or Flask, containerise with Docker, deploy on AWS SageMaker or Azure ML, and build automated retraining pipelines.
5	Data architecture decisions Design data collection strategies, define schema for new data pipelines, and work with data engineers to ensure the right data is captured for future models.
6	Stakeholder management Translate complex technical results into clear business presentations for leadership. Manage project expectations, timelines, and deliverables independently.
7	Code review & mentoring Review junior team members' code and analysis for quality, best practices, and correctness. Mentor freshers through their first ML projects with hands-on guidance.
8	Research & experimentation Follow the latest AI research (Arxiv, Papers with Code), prototype new techniques, and run structured A/B experiments to continuously improve existing models.

Tools used daily: Python, TensorFlow/PyTorch, MLflow, Docker, AWS/Azure, Airflow, Spark, Hugging Face, FastAPI, Kafka, dbt

SENIOR (5+ Years) — Lead Data Scientist / Head of AI / Principal ML Engineer

Salary	Rs. 30 LPA – Rs. 80 LPA+	Role	Head of Data Science / Director of AI / VP Analytics
--------	--------------------------	------	---

1	AI strategy & roadmap Define the organisation's 3-year AI/ML strategy. Decide which AI investments will deliver the most business value and in what order to pursue them.
2	Architecture design Design the end-to-end ML platform — data infrastructure, feature stores, model serving architecture, monitoring stack, and governance frameworks.
3	Executive stakeholder management Present AI initiatives and ROI to the C-suite (CEO, CFO, CTO) and board. Translate AI capabilities into business strategy decisions and investment cases.
4	Team building & leadership Hire, onboard, and grow a team of 10–40+ data scientists, ML engineers, and analysts. Define career progression frameworks and technical standards.
5	Cross-functional collaboration Partner with Product, Engineering, Finance, Marketing, and Operations to embed AI into every part of the business and build a data-driven culture.
6	Research direction Identify where cutting-edge research (LLMs, diffusion models, reinforcement learning) can be applied to company problems. File patents and publish papers.
7	AI governance & ethics Implement model fairness assessments, bias audits, explainability requirements, and data privacy controls to ensure responsible AI deployment.
8	Vendor & budget management Evaluate AI tool vendors (OpenAI, Databricks, Snowflake), manage a multi-crore AI/data infrastructure budget, and optimise cloud costs for ML workloads.

Tools used daily: Databricks, Kubeflow, Vertex AI, OpenAI API, LangChain, Terraform, Kubernetes, Snowflake, Executive BI tools

05 | CAREER ROADMAP — ZERO TO DATA SCIENCE PROFESSIONAL

Data Science is one of the most accessible high-paying fields. You do NOT need a PhD. You need strong Python skills, statistical thinking, hands-on project experience, and the ability to communicate insights clearly. Most freshers who follow a structured path land their first data role within 6–12 months.

YOUR LEARNING PATH — STEP BY STEP

1	Mathematics & Statistics Linear algebra, calculus basics, probability, statistics, distributions, hypothesis testing — the mathematical backbone of all ML.	4–6 weeks
2	Python Programming Variables, loops, functions, OOP, file handling. Then Pandas for data manipulation and NumPy for numerical computing.	6–8 weeks
3	SQL & Databases SELECT, JOINS, GROUP BY, window functions, subqueries. Every data job requires SQL — it is non-negotiable.	3–4 weeks
4	Data Visualisation Matplotlib, Seaborn, Plotly for Python charts. Tableau or Power BI for business dashboards. Storytelling with data.	3–4 weeks
5	Machine Learning Scikit-learn — Regression, Classification, Clustering, Tree models. Train, evaluate, cross-validate, tune.	8–10 weeks
6	Deep Learning & NLP TensorFlow or PyTorch for Neural Networks. CNNs for images, RNNs/Transformers for text. Hugging Face for pre-trained models.	8–12 weeks
7	Generative AI & LLMs Prompt engineering, fine-tuning LLMs, LangChain, RAG pipelines, OpenAI API, building AI applications with real-world impact.	6–8 weeks
8	MLOps & Deployment Docker, FastAPI, MLflow, cloud (AWS/GCP/Azure). Deploy your models so they are accessible to real users and systems.	6–8 weeks

ESSENTIAL PORTFOLIO PROJECTS TO GET HIRED

- **House Price Prediction** Classic regression project. Use a real dataset (Kaggle), engineer features, compare models, deploy as a web app. Shows end-to-end skills.
- **Credit Card Fraud Detection** Imbalanced classification. Use SMOTE for oversampling, Random Forest/XGBoost. Shows you can handle real financial data problems.
- **Sentiment Analysis on Reviews** NLP project. Scrape product reviews, clean text, apply VADER or BERT sentiment classification. Build a Streamlit dashboard to display results.

■ **Movie Recommendation System** Collaborative filtering and content-based recommendations using MovieLens dataset. Shows you understand personalisation systems.

■ **Stock Price Forecasting** Time series project using ARIMA or LSTM. Shows ability to work with temporal data. Visualise predictions vs actuals interactively.

■ **AI Chatbot with RAG** Build a document Q&A; chatbot using LangChain + OpenAI API + FAISS vector store. The hottest project in 2025/26 job applications.

50+ TOOLS & TECHNOLOGIES IN DATA SCIENCE WITH AI

Python	R	SQL	Pandas	NumPy	Matplotlib
Seaborn	Plotly	Scikit-learn	XGBoost	LightGBM	TensorFlow
PyTorch	Keras	Hugging Face	LangChain	OpenAI API	Transformers
NLTK	spaCy	Tableau	Power BI	Looker	Streamlit
Dash	FastAPI	Flask	Docker	Kubernetes	Git
GitHub	MLflow	W&B	Airflow	Apache Spark	Kafka
dbt	Snowflake	BigQuery	AWS SageMaker	Azure ML	Google Vertex AI
Databricks	FAISS	Pinecone	PostgreSQL	MongoDB	Redis
Grafana	Jupyter Notebook				

TOP CERTIFICATIONS FOR DATA SCIENCE WITH AI

Certification	By	What It Covers	Cost
Google Data Analytics Professional Certificate	Google / Coursera	Best for beginners. Covers the full data analysis workflow.	\$1,000–4,000/mo
IBM Data Science Professional Certificate	IBM / Coursera	10-course series. Python, SQL, ML, visualisation, capstone project.	\$1,200/mo
TensorFlow Developer Certificate	Google	Proves practical deep learning skills. Respected globally.	\$100
AWS Certified Machine Learning — Specialty	Amazon Web Services	Cloud ML certification. High demand in cloud-first companies.	\$300
Microsoft Azure AI Engineer Associate (AI-102)	Microsoft	Azure AI services. Growing demand as Azure adoption increases.	\$165
Databricks Certified ML Associate	Databricks	Spark + ML. Highly valued in big data environments.	\$200
DeepLearning.AI Deep Learning Specialisation	Andrew Ng / Coursera	The gold standard deep learning course globally.	\$1,180/mo
Certified Analytics Professional (CAP)	INFORMS	Practitioner-level analytics certification. Growing recognition.	\$900

COMPLETE JOB ROLES & SALARY GUIDE

Job Role	Specialisation	Experience	Avg Salary India
Data Analyst	Analytics	0–1 yr	Rs. 4 – 8 LPA

Junior Data Scientist	ML/AI	0–2 yrs	Rs. 5 – 10 LPA
Business Intelligence Analyst	Analytics	1–2 yrs	Rs. 5 – 9 LPA
Machine Learning Engineer	MLOps	1–3 yrs	Rs. 8 – 18 LPA
Data Engineer	Data Infra	1–3 yrs	Rs. 7 – 16 LPA
NLP / AI Engineer	NLP/GenAI	2–4 yrs	Rs. 12 – 25 LPA
Senior Data Scientist	ML/AI	3–5 yrs	Rs. 18 – 35 LPA
ML Platform / MLOps Engineer	MLOps	3–5 yrs	Rs. 20 – 38 LPA
Computer Vision Engineer	CV/DL	2–5 yrs	Rs. 15 – 32 LPA
AI Research Scientist	Research	4–6 yrs	Rs. 25 – 50 LPA
Lead/Principal Data Scientist	ML/AI	6–8 yrs	Rs. 35 – 60 LPA
Head of AI / Director	Leadership	8–10 yrs	Rs. 50 – 80 LPA+

TOP COMPANIES HIRING DATA SCIENTISTS IN INDIA

Google India	Microsoft India	Amazon/AWS	Flipkart
Paytm	PhonePe	Razorpay	Swiggy / Zomato
TCS	Infosys	Wipro	HCL Technologies
Deloitte	PwC	KPMG	Mu Sigma
Ola	Meesho	Nykaa	Zepto
CRED	Groww	Zerodha	upGrad / BYJU'S

ByteLearned

bytelearned.online

info@bytelearned.online

7396633051