# MLOps: Versioning Models and Automating Deployment

MLOps (Machine Learning Operations) bridges the gap between data science and operations, bringing DevOps principles into the lifecycle of machine learning models. A critical component of MLOps is **model versioning** and **automated deployment**, ensuring models are reproducible, traceable, and scalable in production environments.

## 🚀 Automating Deployment of Models

Automating deployment is about continuously delivering and updating models in production with minimal manual intervention.

### 📦 Packaging the Model

Models are usually packaged as:

- Python packages ( `setup.py` )
- REST APIs using Flask/FastAPI
- Docker containers for portability

### 📈 CI/CD Pipeline for ML (MLOps)

```
[Git Push] ➜ [CI: Test + Build] ➜ [Model Registry Push] ➜
[CD: Deploy to Staging/Production] ➜ [Monitor and Retrain]
```

**Popular CI/CD Tools for MLOps:**

- **GitHub Actions / GitLab CI:** Trigger pipelines on code/model changes
- **Jenkins:** Custom automation jobs
- **Seldon Core / KFServing:** Deploy models on Kubernetes
- **Argo Workflows:** For ML workflows in Kubernetes

## ⚙️ Deployment Patterns

- **Batch Inference**: Pre-compute predictions and store
- **Online Inference**: Serve predictions in real-time via API
- **Streaming Inference**: Combine with Kafka/Spark for real-time processing

# 📊 Monitoring and Retraining

Once deployed, models must be:

- **Monitored** for performance drift (e.g., accuracy drop)
- **Logged** for input/output
- **Re-trained** on fresh data via automated workflows

Tools: Prometheus + Grafana, EvidentlyAI, A/B Testing pipelines

# ✅ Summary

| Step | Tool/Concept |
|------|--------------|
| Code Versioning | Git |
| Data/Model Versioning | DVC, MLflow |
| Packaging | Docker, FastAPI |
| CI/CD | GitHub Actions, Jenkins |
| Deployment | Kubernetes, Seldon Core |
| Monitoring | Prometheus, Grafana |

MLOps enables reliable, repeatable, and scalable ML pipelines. Versioning and automated deployment are foundational for ensuring that ML models remain robust and accountable in production.