## 1.6 Algorithmic bottlenecks in ML

Algorithmic bottlenecks in Machine Learning (ML) models are performance or scalability limitations that stem from inherent inefficiencies in the design of the algorithm itself, as opposed to hardware constraints

Algorithmic bottlenecks in ML models are computationally expensive stages such as data preprocessing, feature extraction, model training, and hyperparameter tuning that dominate the overall execution time and memory usage, reducing scalability and performance.

**Key Algorithmic Bottlenecks**

➢ **Computational Complexity:** Many ML algorithms, particularly in high-dimensional data analysis, face exponential time and space requirements, making problems intractable for large datasets or high-dimensional data

➢ **Inefficient Operations:** Within deep neural networks (DNNs), specific operations can consume a disproportionate amount of resources. The self-attention mechanism in transformer models, for instance, has a quadratic time complexity that is a significant performance bottleneck.

➢ **Data Movement Costs:** In modern computing systems, moving data between memory and the processor is often much more energy-intensive and time-consuming than the computation itself (the von Neumann bottleneck). Algorithms that require extensive data movement, such as large matrix multiplications, can become memory-bound rather than compute-bound.

➢ **Gradient Issues:** The optimization algorithms used for training, such as gradient descent, can encounter issues in complex "loss landscapes". Noisy or unstable gradients can slow down training, cause instability, or prevent the model from converging to the optimal solution.

➢ **Non-Distributable Computation:** In distributed training systems, if certain parts of the algorithm cannot be parallelized effectively (e.g., a single-server instance managing a critical resource), that component becomes a single point of failure and a performance bottleneck.

➢ **Scalability Limitations:** Transitioning an ML model from a prototype to a production-ready application often reveals scalability issues. Algorithmic designs that work well on small datasets may fail to scale linearly with increasing data volume or model size.

## Sample Case Studies:

**Case Study 1: The Curse of Dimensionality and KNN in Financial Fraud Detection**

**Case Study 2: Quadratic Programming in Large-Scale Support Vector Machines (SVMs)**

**Case Study Title3: House Price Prediction Using Machine Learning**

**Case Study 4: Image Classification Using CNN**

**Case Study 5: Text Classification Using TF-IDF + SVM**

**Case Study 6: Recommendation System (Collaborative Filtering)**

**Case Study 7: Fraud Detection Using KNN**

**Case Study 8: Hyperparameter Tuning in Deep Learning**

**Data Pipeline**

A **machine learning (ML) pipeline** is an automated, structured workflow that transforms raw data into a deployable ML model and manages its ongoing performance.

A data pipeline in ML is a sequence of steps that collect, clean, transform, and prepare data for model training and deployment. Bottlenecks in the pipeline can significantly affect scalability and performance.

Implementing an ML pipeline provides numerous advantages:

**Automation and Efficiency**: Automates repetitive tasks, saving time and reducing human error.

**Reproducibility**: Ensures consistent results across different runs and team members by standardizing the workflow and using version control.

**Scalability**: Handles large volumes of data and complex models efficiently by leveraging distributed computing frameworks.

**Faster Deployment**: Streamlines the transition from development to production, enabling quicker time-to-market for ML solutions.

## Sample Case Studies:

**Case Study1: Data Pipeline in a Machine Learning Algorithm**

**Case Study2:End-to-End Data Pipeline for Customer Churn Prediction**

**Case Study3: Real-Time Fraud Detection (Streaming Data Pipeline)**

**Case Study4: NLP Pipeline for Sentiment Analysis**

**Case Study5: Image Processing Pipeline for Face Recognition**

**Case Study6: IoT Sensor Data Pipeline**

**Case Study7: Healthcare ML Data Pipeline**

**Case Study8: Batch Training Pipeline for Large-Scale ML Problem**