



Cover Page



AI-DRIVEN BIG DATA OPTIMIZATION : MATHEMATICAL ALGORITHMS AND CS IMPLEMENTATIONS FOR REAL TIME DECISION SYSTEM

N.S.V. Kiran Kumar¹ and N. Gopala Krishna²

¹Lecturer in Mathematics, Government Degree College, Mandapeta, Andhra Pradesh, India

²Lecturer in CS, Government Degree College, Mandapeta, Andhra Pradesh, India

Abstract

AI driven big data optimization enables real-time decision systems to process massive, high velocity datasets and produce actionable insights with minimal latency. This paper reviews mathematical foundations, core algorithms, and computer science implementations enabling real-time optimization in domains such as finance, healthcare, logistics, and cyber security.

1. Introduction

Modern systems generate data streams at a scale that exceeds manual or traditional computational processing. Artificial intelligence, combined with big data architectures, enables automated decision-making based on continuous learning and mathematical optimization.

Real time decision systems require:

- Low latency data ingestion
- Adaptive learning algorithms
- Scalable distributed computing
- Robust optimization under uncertainty

This paper surveys the mathematical algorithms and associated CS implementations enabling these capabilities.

2. Big Data Characteristics and Real-Time Constraints

2.1 The “V” Properties

Big data streams typically exhibit:

- Volume
- Velocity
- Variety
- Veracity
- Variability
- Value

2.2 Real-Time Processing Requirements

A real-time decision system must support:

- Millisecond level response
- Nonstop model updates
- Fault-tolerant distributed execution
- Stream-oriented analytics



3. Mathematical Foundations

3.1 Optimization Models

Optimization problems typically minimize a cost function $f(x)$ under constraints $g_i(x) \leq 0$.

Gradient based optimization

Algorithms update parameters using $x_{t+1} = x_t - \eta_t \Delta f(x_t)$

Stochastic Gradient Descent (SGD)

For streaming data, use sample-based updates :

$$x_{t+1} = x_t - \eta_t \Delta f(x_t; \xi_t)$$

3.2 Convex and Non-Convex Methods

- Convex optimization guarantees global minima
- Ultramodern AI relies on non-convex deep models but uses convex relaxations for speed

3.3 Reinforcement Learning (RL)

Real-time decisions often follow an RL formulation:

Reward maximization:

$$\max E [\sum_{t=0}^T \gamma^t R(s_t, a_t)]$$

Value function update:

$$Q(s, a) \leftarrow Q(s, a) + \alpha (r + \gamma \max Q(s^t, a^t) - Q(s, a))$$

3.4 Bayesian Inference

For uncertainly modelling:

Posterior distribution:

$$P(\theta|X) = \frac{P(X|\theta)P(\theta)}{P(X)}$$

4. AI Algorithms for Big Data Optimization

4.1 Deep Learning Architectures

- LSTM and GRU for sequential data
- Transformers for large-scale sequence modelling
- Graph Neural Networks for irregular data structures

4.2 Online and Incremental Learning

Models are updated with each new data item:



$$\omega_{i+1} = \omega_i + \Delta_t(X_t)$$

4.3 Federated Learning

Distributed gradient aggregation:

$$\omega_{i+1} = \sum_{i=1}^N \frac{n_i}{n} \omega_i^{(i)}$$

5. Computational Architectures

5.1 Distributed Frameworks

Key systems:

- Apache Spark Streaming
- Flink
- Kafka + TensorFlow Serving
- Ray for scalable RL

5.2 Data Structures for Speed

- Bloom filters for fast membership checking
- Count-min sketch for frequency estimation
- KD trees for high-dimensional nearest neighbors

5.3 Parallelization

- GPU/TPU matrix operations
- MapReduce style decomposition
- Model sharding and pipeline parallelism

6. Real-Time System Design

6.1 Architecture Overview

A typical pipeline includes:

1. Data ingestion
2. Feature extraction
3. Online model inference
4. Feedback circle for re training
5. Real-time decision execution

6.2 Dynamic Resource Scaling

Systems rely on autoscaling algorithms such as:

$$R_{t+1} = R_t + k (E_t - E_{target})$$

6.3 Latency Optimization

- In memory data processing
- Pre computed lookup tables



Cover Page



- Approximate algorithms

7. Case Studies

7.1 Financial Fraud Detection

- Streamed transaction analysis with GNNs
- Low-quiescence anomaly discovery using autoencoders

7.2 Intelligent Transportation

- Reinforcement learning for business light optimization
- Real-time vehicle routing using convex relaxation

7.3 Healthcare Monitoring

- Prophetic models for patient vitals
- Bayesian updating for threat soothsaying

8. Challenges and Future Directions

8.1 Challenges

- Data privacy
- Algorithmic bias
- Hardware cost
- Interpretability of deep models

8.2 Future Trends

- Edge AI for ultra low latency
- Quantum optimization algorithms
- Explainable real time AI models
- Self-correcting autonomous agents

9. Conclusion

AI- driven big data optimization integrates fine rigor with scalable computer wisdom executions to produce real- time decision systems able of recycling massive aqueducts of data. As algorithms and tackle continue to advance, real- time AI systems will play an decreasingly vital part across diligence.

References

1. Dean, J. and Ghemawat, S., “MapReduce: Simplified Data Processing on Large Clusters,” Communications of the ACM, vol. 51, no. 1, pp. 107–113, 2008.
2. Zaharia, M. et al., “Discretized Streams: Fault-Tolerant Streaming Computation at Scale,” in Proceedings of the 24th ACM Symposium on Operating Systems Principles (SOSP), 2013.
3. Carbone, P. et al., “Apache Flink: Stream and Batch Processing in a Single Engine,” IEEE Data Engineering Bulletin, vol. 38, no. 4, pp. 28–38, 2015.
4. Kreps, J., Narkhede, N. and Rao, J., “Kafka: A Distributed Messaging System for Log Processing,” in Proceedings of the NetDB Workshop, 2011.
5. LeCun, Y., Bengio, Y. and Hinton, G., “Deep Learning,” Nature, vol. 521, pp. 436–444, 2015.
6. Goodfellow, I., Bengio, Y. and Courville, A., Deep Learning. MIT Press, 2016.



Cover Page



7. Vaswani, A. et al., "Attention Is All You Need," in Advances in Neural Information Processing Systems (NeurIPS), 2017.
8. Wu, Z. et al., "A Comprehensive Survey on Graph Neural Networks," IEEE Transactions on Neural Networks and Learning Systems, vol. 32, no. 1, pp. 4–24, 2021.
9. Bottou, L., "Stochastic Gradient Descent Tricks," in Neural Networks: Tricks of the Trade, Springer, 2012.
10. Sutton, R.S. and Barto, A.G., Reinforcement Learning: An Introduction, 2nd ed., MIT Press, 2018.
11. Ghavamzadeh, M., Mannor, S., Pineau, J. and Tamar, A., "Bayesian Reinforcement Learning: A Survey," Foundations and Trends in Machine Learning, vol. 8, nos. 5–6, pp. 359–483, 2015.
12. Shahriari, B., Swersky, K., Wang, Z., Adams, R.P. and de Freitas, N., "Taking the Human Out of the Loop: A Review of Bayesian Optimization," Proceedings of the IEEE, vol. 104, no. 1, pp. 148–175, 2016.
13. Mohammadi, M., Al-Fuqaha, A., Sorour, S. and Guizani, M., "Deep Learning for IoT Big Data and Streaming Analytics: A Survey," IEEE Communications Surveys & Tutorials, vol. 20, no. 4, pp. 2923–2960, 2018.
14. Gama, J., Žliobaitė, I., Bifet, A., Pechenizkiy, M. and Bouchachia, A., "A Survey on Concept Drift Adaptation," ACM Computing Surveys, vol. 46, no. 4, 2014.
15. Kairouz, P. et al., "Advances and Open Problems in Federated Learning," Foundations and Trends in Machine Learning, vol. 14, nos. 1–2, pp. 1–210, 2021.
16. Li, T. et al., "Federated Learning: Challenges, Methods, and Future Directions," IEEE Signal Processing Magazine, vol. 37, no. 3, pp. 50–60, 2020.
17. Bifet, A. et al., Machine Learning for Data Streams: with Practical Examples in MOA. MIT Press, 2018.
18. Hesse, G., and Lorenz, M., "Distributed Stream Processing Frameworks for Fast & Big Data: Spark, Flink, Kafka Streams," codecentric AG, 2017.
19. Horchidan, S. et al., "Evaluating Model Serving Strategies over Streaming Data," in Proceedings of the 2022 IEEE International Conference on Big Data (BigData), 2022.
20. Zhang, X. et al., "Optimization of Artificial Intelligence in Localized Big Data Real-Time Task Scheduling," Frontiers in Physics, vol. 12, 2024.