

# Application of Machine Learning for Ethical Decision-Making

<https://doi.org/10.31713/MCIT.2025.080>

Savych Dmytro  
 National University of  
 Water and Environmental Engineering  
 Rivne, Ukraine  
[d.v.savych@nuwm.edu.ua](mailto:d.v.savych@nuwm.edu.ua)

Pryshchepa Oksana  
 National University of  
 Water and Environmental Engineering  
 Rivne, Ukraine  
[o.v.pryshchepa@nuwm.edu.ua](mailto:o.v.pryshchepa@nuwm.edu.ua)

Naumchuk Oleksandr  
 National University of  
 Water and Environmental Engineering  
 Rivne, Ukraine  
[o.m.naumchuk@nuwm.edu.ua](mailto:o.m.naumchuk@nuwm.edu.ua)

**Abstract**—This paper presents a conceptual and methodological framework for applying machine learning in managerial decision support systems to improve ethical outcomes. We review key fairness, explainability, privacy, and accountability methods, propose an integrated architecture that embeds ethical constraints within machine learning pipelines, and discuss governance frameworks for ensuring responsible deployment. The contribution is a roadmap for integrating ethical machine learning into organizational decision-making.

**Keywords**—machine learning; managerial decision-making; fairness; explainability; governance

## I. INTRODUCTION

Modern organizations increasingly rely on data-driven decision-making. Machine learning models can uncover patterns, forecast outcomes, and support managers in selecting optimal strategies. However, without appropriate safeguards, such systems may propagate biases, lack transparency, or violate privacy and fairness principles. The challenge is to embed ethical considerations within the machine learning lifecycle such that managerial decisions become both effective and socially responsible.

The aim of this paper is to outline methods and architectural approaches for ethical machine learning in decision support, along with governance and validation strategies to ensure accountability and trust.

## II. ETHICAL CHALLENGES IN DECISION-SUPPORT MACHINE LEARNING

Key risks associated with deploying machine learning in managerial contexts have been widely discussed in the literature. One major concern is algorithmic bias, which arises when skewed training data or model assumptions lead to unfair treatment of specific groups. Another is model opacity, where complex or “black box” algorithms make decision logic difficult to interpret for managers and stakeholders.

Privacy risks also emerge when models inadvertently expose sensitive attributes or enable re-identification of individuals [2]. Furthermore, the lack of accountability in automated systems can obscure responsibility when incorrect or unethical decisions occur [4]. Finally, model drift and other unintended consequences may arise as data distributions change over time or models capture spurious correlations [6].

To mitigate these risks, researchers emphasize the need for embedding ethical constraints within the machine learning lifecycle, implementing continuous monitoring, and establishing governance frameworks that ensure transparency, fairness, and accountability [1], [3].

## III. METHODS FOR ETHICAL MACHINE LEARNING IN MANAGERIAL DECISION SYSTEMS

Ethical machine learning in managerial decision support requires approaches that ensure fairness, transparency, privacy, and accountability across the entire machine learning lifecycle. A variety of methodological solutions can be integrated into machine learning pipelines to achieve these objectives.

In the area of fairness and bias mitigation, methods such as reweighing, adversarial debiasing, and fairness through unawareness are used to reduce discrimination and ensure equitable model behaviour. Reweighting adjusts sample weights to balance distributions among protected groups, while adversarial debiasing trains a predictive model together with an adversary attempting to infer sensitive attributes, forcing the predictor to minimize bias. Fairness through unawareness removes sensitive variables such as gender or race from input features, though this method alone is often insufficient because proxy variables can still encode bias [1].

For explainability and interpretability, approaches such as LIME, SHAP, and counterfactual explanations make machine learning models more transparent and

understandable to decision-makers. LIME (Local Interpretable Model-Agnostic Explanations) generates local perturbation-based insights for individual predictions. SHAP (SHapley Additive exPlanations) quantifies feature importance using Shapley values, allowing both local and global interpretation. Counterfactual explanations illustrate minimal input changes that would lead to a different output, helping managers understand and justify automated decisions.

Ensuring privacy and confidentiality is another crucial dimension of ethical machine learning. Differential privacy introduces calibrated noise into data or outputs to protect individual identities. Federated learning allows distributed training without centralizing raw data, and homomorphic encryption enables computations on encrypted information, maintaining confidentiality throughout the process. Continuous validation and monitoring, including k-fold cross-validation, A/B testing, and model drift detection, ensure that model performance remains reliable and ethically compliant over time.

Finally, governance and auditing frameworks provide the structural foundation for accountability and trust. Ethics-based auditing (EBA) establishes systematic evaluations of system behavior against ethical norms, while the ECCOLA method offers a structured framework for governing ethical AI throughout its lifecycle [5].

#### IV. ETHICAL MACHINE LEARNING ARCHITECTURE FOR DECISION-MAKING SYSTEMS

A high-level architecture for embedding ethical machine learning into managerial decision support systems involves multiple integrated layers that ensure fairness, transparency, privacy, and accountability throughout the machine learning lifecycle.

The first layer focuses on data preparation, including data ingestion, quality checks, bias detection, and preprocessing techniques such as reweighing and anonymization. This stage ensures that the training data are reliable, ethically balanced, and privacy-preserving.

The second layer addresses model training under ethical constraints. Predictive models are trained with fairness-regularization or adversarial fair learning approaches to mitigate bias [2]. Privacy-preserving techniques, such as differential privacy and federated learning, are incorporated to maintain data confidentiality while enabling effective learning.

An interpretability and explanation module forms the third layer, generating outputs such as SHAP, LIME, or counterfactual explanations. These tools provide transparency into model decisions, allowing managers to understand the rationale behind predictions and supporting trust in automated recommendations.

The fourth layer integrates decision-making with human oversight. Model suggestions are presented to managers along with explanations, enabling human decision-makers to override, adjust, or validate automated recommendations. This human-in-the-loop approach ensures that final decisions consider both algorithmic insights and managerial judgment.

The fifth layer implements monitoring and feedback, tracking model performance, fairness metrics, and potential drift over time. Continuous monitoring supports revalidation, retraining, and corrective actions whenever the system's predictions deviate from expected ethical or operational standards.

Finally, the governance and audit layer ensures accountability and compliance. Ethical audits, logging of decisions, and accountability tracing are implemented to guarantee that machine learning-driven recommendations adhere to organizational and societal norms.

This architectural framework ensures that machine learning-generated suggestions are not blindly followed but are mediated by interpretability, human oversight, and governance mechanisms, promoting ethically responsible managerial decision-making.

#### CONCLUSION

Integrating machine learning into managerial decision-making systems offers significant potential, but ethical risks must be proactively managed. By applying fairness-enhancing methods, interpretability and explainability tools, privacy-preserving techniques, and governance frameworks such as ethics-based auditing, organizations can develop decision support systems that are both effective and ethically robust. Future research should focus on domain-specific case studies, empirical validation of ethical interventions, and longitudinal assessments of governance and audit outcomes.

#### REFERENCES

- [1] S. Barocas, M. Hardt, and A. Narayanan, *Fairness and Machine Learning: Limitations and Opportunities*. MIT Press, 2019.
- [2] I.O. Gallegos, R.A. Rossi, J. Barrow, M.M. Tanjim, S. Kim, F. Dernoncourt, T.Yu, R. Zhang, & N. Ahmed, "Bias and Fairness in Large Language Models: A Survey," *Computational Linguistics*, vol. 50, pp. 1097-1179.
- [3] A. Jobin, M. Ienca, and E. Vayena, "The global landscape of AI ethics guidelines," *Nature Machine Intelligence*, vol. 1, no. 9, pp. 389–399, 2019.
- [4] J. Morley, L. Floridi, L. Kinsey, and A. Elhalal, "From what to how: An initial review of publicly available AI ethics tools, methods and research to translate principles into practices," *Science and Engineering Ethics*, vol. 27, no. 4, pp. 1–31, 2021.
- [5] Mokander, J., Floridi, L., & Cave, S. Ethics-based auditing of automated decision-making systems: Nature, scope, and limitations. *arXiv preprint arXiv:2108.11857*, 2021.
- [6] D. Sculley, G. Holt, D. Golovin, V. Davydov, T. Phillips, E. Ebner, et al., "Hidden technical debt in machine learning systems," in *Proc. 28th Int. Conf. Neural Information Processing Systems (NIPS)*, 2015.