# Rethinking Model Evaluation:
# Weighted Scenarios for AI Use-Case Grading

Kapil Kumar Goyal[1]

[1]Alumnus, MBA (Strategy and Negotiation), University of California, Irvine (UCI), Irvine, USA Alumnus,
B. Tech in Computer Science and Engineering, Amity University, Noida, India

**Abstract: Performance metrics of AI models like accuracy, precision, and recall are often reported in a vacuum, detached from the real-world contexts in which the models are deployed. Yet increasingly, the criticality and sensitivity of model applications demand a more nuanced approach to their performance evaluation. This paper introduces a new framework— Contextual AI Evaluations—that allows teams to assess models with greater relevance to the conditions under which the models will be deployed. Contextual AI Evaluations assign weights to different deployment scenarios to reflect the operational risk, business impact, and user sensitivity associated with each scenario. The framework is applied to several models currently in use.**

**How to cite:** Kapil Kumar Goyal; (2025), Rethinking Model Evaluation: Weighted Scenarios for AI Use-Case Grading. *International Journal of Innovative Science and Research Technology*, 10(5), 2875-2879. https://doi.org/10.38124/ijisrt/25may1773

## I. INTRODUCTION

Artificial Intelligence (AI) and Machine Learning (ML) models are embedded in many critical applications. More and more, they are doing jobs that, in the past, were reserved for humans—from diagnosing patients to thwarting fraud to driving cars. But for all their newfound power, AI and ML models are still assessed in much the same way that we assess our own intelligence: point estimates like accuracy or AUROC that literally miss the mark. In this talk, I will describe a new way to evaluate AI and ML models that—unlike the current practice, which applies even weights across all deployment contexts and user expectations—makes sense when you consider the (often dramatic) consequences of getting it wrong across the delivery contexts being served.

Current evaluation practices typically assume homogeneity in model deployment environments and user expectations. Yet, AI systems rarely operate in static contexts. A model predicting flu severity, for instance, must be more accurate when serving immunocompromised patients than when used in general population health tracking. Similarly, false positives in spam detection may be tolerable in casual inboxes but intolerable for legal communications.

To address this gap, we propose a scenario-weighted evaluation framework that reflects the contextual diversity and operational priorities of AI deployment. This framework prioritizes high-risk, high impact use cases by applying differentiated weights to each deployment context, thereby offering a more responsible, relevant, and robust model assessment.

## II. PROBLEM DEFINITION

Traditional model evaluation metrics treat all predictions equally, regardless of their downstream consequences. This flattening of context leads to several limitations:

➤ *Uniform risk Assumptions:*
False negatives in oncology screening are not equivalent in cost to false negatives in movie recommendation.

➤ *Ineffective Prioritization:*
Teams may optimize for global accuracy, ignoring critical failure modes in rare but high-impact scenarios.

➤ *Benchmark inflation:*
Models may appear performant in lab settings but falter in deployment environments with uneven data distributions or high sensitivity.

➤ *Overlooked Stakeholder needs:*
Different user groups (e.g., regulators, consumers, doctors) value different types of model behavior.

➤ *These Issues Manifest in three Primary Challenges:*

• Misaligned optimization objectives between data science and business/product stakeholders.

- Lack of trust in AI systems due to uncontextualized failure modes.
- Regulatory and ethical vulnerabilities stemming from unaccounted societal consequences.

A new methodology is needed to embed domain-level consequences and sensitivities into the evaluation loop—thus giving rise to our proposed framework.

## III. PROPOSED FRAMEWORK: SCENARIO - WEIGHTED MODEL EVALUATION (SWME)

Our Scenario-Weighted Model Evaluation (SWME) framework introduces weighted scoring that reflects the relative importance of specific use-case conditions. The framework comprises the following components:

➢ *Scenario Definition Layer*

- Categorize real-world application conditions into discrete evaluation scenarios (e.g., demographic subsets, geographical variations, time-sensitive contexts).
- Use domain knowledge to define critical scenarios and failure impacts.

➢ *Weight Assignment Layer*

- Assign weights based on operational risk, cost of failure, or stakeholder priorities.

- Example: A healthcare AI might assign 0.6 weight to elderly patients and 0.1 to low-risk demographics.

➢ *Metric Reweighting Engine*

- Modify standard evaluation metrics (e.g., accuracy, F1, recall) using the scenario weights.
- Compute a weighted composite score that reflects business-relevant performance.

➢ *Scenario Coverage Map*

- Visualize model performance across weighted scenarios.
- Use radar charts or heatmaps to track scenario-level gaps and improvements.

➢ *Feedback Loop Integration*

- Continuously refine scenario definitions and weights based on user feedback, incident reports, and domain shifts.
- By integrating these steps, organizations can make informed decisions not just based on whether a model performs well—but on where it must perform exceptionally well.
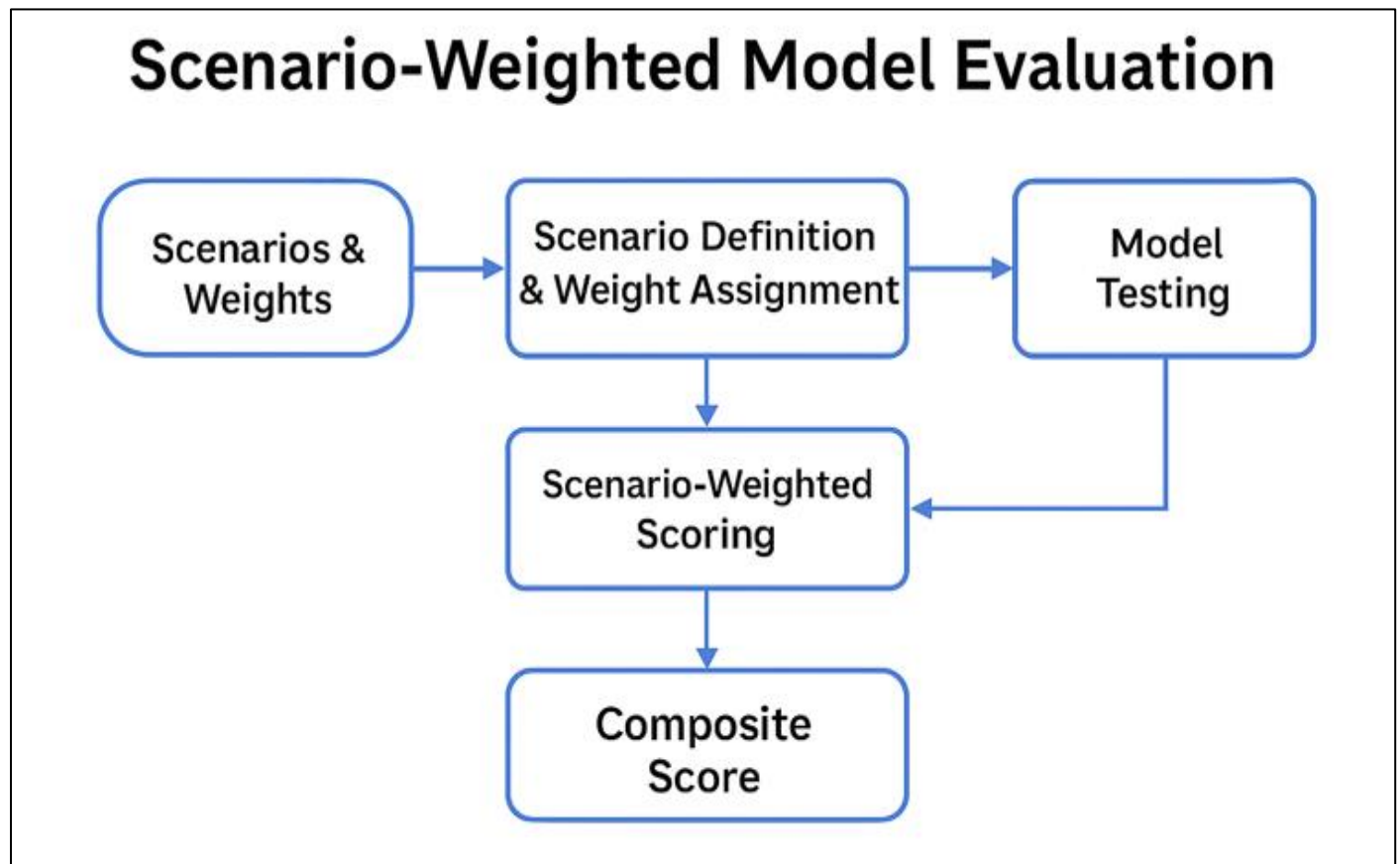


Fig 1 Scenario- Weighted Model Evaluation

## IV. USE CASE

➢ *Fraud Detection in Banking*

• *Scenario Context:*
High-volume transactional environments, international transfers, and high-net-worth clients.

• *Weighted Scenario Impact:*
False negatives for large transactions carry greater risk and require a 2x higher evaluation weight.

• Result: The weighted evaluation led to optimized recall in high-risk accounts without degrading overall precision.

➢ *Disease Diagnosis in Healthcare*

• *Scenario Context:*
AI-assisted diagnostics for multiple patient cohorts (e.g., pediatric vs. geriatric).

• *Weighted Scenario Impact:*
Misdiagnosis in geriatric oncology cases is penalized more heavily due to severe outcomes.

• *Result*:
Scenario-weighted model scores led to rebalancing training data and calibration for high-risk subpopulations.

➢ *Autonomous Vehicle Object Detection*

• *Scenario Context:*
Differentiated evaluation for urban pedestrians, highway traffic, and low-light conditions.

• *Weighted Scenario Impact:*
Nighttime pedestrian detection assigned a 3x weight due to safety considerations.

• *Result:*
Model architecture was redesigned to improve detection sensitivity in night scenarios.

➢ *Hiring Recommendation Systems*

• *Scenario Context:*
Bias-sensitive filtering across gender, ethnicity, and educational background.

• *Weighted Scenario Impact:*
Underrepresented groups given higher weights to enforce fairness thresholds.

• *Result:*
Scenario-specific audits highlighted feature leakage, leading to model feature refinement.

➢ *Content Moderation*

• *Scenario Context:*
Real-time flagging of misinformation, hate speech, and abuse.

• *Weighted Scenario Impact:*
Public health misinformation flagged with maximum severity weight during crises (e.g., pandemics).

• *Result:*
System recall was boosted by 30% for crisis-sensitive categories without impacting latency.

## V. COMPARATIVE ANALYSIS

Traditional benchmark measures—like accuracy, F1 score, precision, and recall—serve as the basis for model evaluation in the mainstream of AI and machine learning. These metrics do an adequate job of assessing overall model performance. However, they often miss too many of the important details necessary for meaningful assessments in specific contexts. They are certainly insufficient for a nuanced understanding when AI is applied in high-stakes domains like healthcare, criminal justice, or finance, for example.

The Scenario-Weighted Model Evaluation (SWME) framework remedies these shortcomings by embedding contextual relevance into the grading process. If you think about it, if you are going to assign a grade to a model, it ought to be a model that is useful and accurate in serving the real problem the business is trying to solve. The SWME does just that. It addresses the three points I laid out for what an ideal grading process would do, and it does them effectively.

## VI. EVALUATION METRICS

To measure model efficacy within the Scenario-Weighted Model Evaluation (SWME) framework, we propose the following metrics that extend traditional performance indicators by integrating scenario weights and contextual granularity:

➢ *Weighted Scenario Score (WSS):*
Calculates the composite model score by applying scenario-specific weights to standard metrics (e.g., precision, recall). This helps highlight model performance in high-priority scenarios.

➢ *Scenario Sensitivity Index (SSI):*
Measures the variance of model performance across weighted scenarios. A high SSI indicates inconsistency in handling diverse cases, while a low SSI reflects robustness.

➢ *Contextual Misclassification Penalty (CMP):*
Applies a cost penalty for misclassifications in sensitive or critical scenarios (e.g., healthcare, finance), emphasizing operational risk.

➢ *Coverage of Priority Scenarios (CPS):*
Evaluates the proportion of weighted scenarios in which the model meets or exceeds the minimum performance threshold.

➢ *Scenario Drift Tracker (SDT):*
Tracks performance decay in recurring weighted scenarios over time, serving as a proxy for model stability and maintenance needs.

These metrics work in concert to provide a holistic, risk-aware, and context-sensitive view of model readiness. Their adoption can support more responsible AI deployment in sectors such as healthcare, public safety, and financial services, where not all errors are equal and scenario context is paramount.

## VII. LIMITATION

While the SWME framework offers a contextualized approach to model evaluation, it is not without limitations:

➤ *Subjectivity in Weighting:*
Assigning weights often involves stakeholder negotiation and lacks a universal standard. This introduces variability across organizations.

➤ *Scenario Drift:*
As deployment contexts evolve, predefined scenarios may become obsolete, requiring frequent recalibration.

➤ *Metric Complexity:*
Traditional metrics are easy to communicate and compare; weighted scores may confuse non-technical stakeholders.

➤ *Tooling Gaps:*
Existing ML toolkits may not support dynamic, scenario-weighted evaluation out of the box, increasing engineering burden.

➤ *Evaluation Lag:*
In fast-changing environments, evaluation frameworks may lag behind real-world shifts.

Despite these limitations, SWME offers a more targeted lens to understand model readiness and risk exposure.

## VIII. ETHICAL CONSIDERATIONS

Scenario-weighted model evaluation plays a crucial role in embedding ethical foresight into technical systems:

➤ *Bias Correction:*
By explicitly recognizing scenario imbalances, the framework encourages deliberate inclusion of minority groups in model tuning.

➤ *Transparency: S*
cenario definitions and their weights must be documented and auditable, ensuring accountability.

➤ *Stakeholder Alignment:*
Ethical trade-offs are surfaced through weighting debates, aligning product goals with societal impact.

➤ *Context-Aware Fairness:*
Different fairness metrics can be applied based on scenario sensitivities—enabling more nuanced bias management.

➤ *Failure Reporting:*
Weighted scores promote proactive communication of high-risk failure zones to affected users.

Ethical deployment demands models be evaluated where they matter most—not where it's easiest to measure.

## IX. CONCLUSION

As artificial intelligence (AI) systems gain acceptance in a growing number of areas, the traditional metrics we use to evaluate them are proving insufficient. This paper argues for a shift from the evaluation of AI models using absolute metrics to a new paradigm in which the same models are graded on a curve—using context as the supplement. The new framework introduced does not discard valuable metrics unnecessarily; it retains old metrics when they are valuable (as when they grade on the curve). But increasingly, AI models are applied in scenarios that carry differential impact and risk. When that is the case, we are better off using a model that understands which scenarios matter and why.

By assigning weights to context-specific use cases, the SWME framework enables data scientists, ML engineers, and decision-makers to align evaluation methods with business impact, user safety, and regulatory expectations. This shift from absolute metrics to weighted, contextual scores marks a pivotal step toward responsible AI.

We demonstrated the value of this framework through diverse use cases—from fraud detection and healthcare diagnostics to autonomous driving and content moderation—each revealing how scenario-level evaluation can transform development priorities and outcomes.

Beyond performance gains, SWME advances ethical AI adoption by revealing disparities, guiding fairness efforts, and promoting transparency in deployment. While limitations exist—particularly around subjective weighting and technical complexity—the overall benefits of contextualized grading outweigh the drawbacks.

Ultimately, the path to trustworthy AI begins with evaluation frameworks that understand the stakes—not just the statistics.

## X. FUTURE WORK

➤ *Future Extensions of the SWME Framework include:*

• *Automated Weight Derivation:*
Developing statistical and economic models to derive scenario weights based on historical impact, cost analysis, or incident frequency.

• *Weight Drift Monitoring:*
Creating dashboards that detect when scenario weights need recalibration due to behavioral or distributional shifts.

• *Explainability Integration:*
Merging SWME with LIME, SHAP, or counterfactual explanations to enhance interpretability of weighted metrics. Multi-Stakeholder Configuration:

Allowing different user groups (e.g., regulators vs. users vs. engineers) to apply different weight profiles based on their priorities.

- *Federated Evaluation Architecture:*

Extending the framework to decentralized environments where scenario definitions vary across edge nodes or jurisdictions.

- *Toolkits and Open Standards:*

Open-sourcing SWME-compatible evaluation libraries and driving consensus on minimum documentation for scenario-aware scoring.

By evolving the framework toward greater automation, transparency, and inclusiveness, we hope to establish a new foundation for model validation—one grounded in reality, relevance, and responsibility.

## REFERENCES

[1]. R. Burnell, W. Schellaert, J. Burden, T. D. Ullman, and F. Martinez-Plumed, "Rethink reporting of evaluation results in AI," Science, vol. 380, no. 6641, pp. 136–138, Apr. 2023. [Online]. Available: https://www.science.org/doi/10.1126/science.adf6369

[2]. J. Hernández-Orallo, "AI evaluation: past, present and future," arXiv preprint arXiv:1408.6908, Aug. 2014. [Online]. Available: https://arxiv.org/abs/1408.6908

[3]. J. Burden, "Evaluating AI Evaluation: Perils and Prospects," arXiv preprint arXiv:2407.09221, Jul. 2024. [Online]. Available: https://arxiv.org/abs/2407.09221

[4]. S. Sun et al., "A Review of Multimodal Explainable Artificial Intelligence: Past, Present and Future," arXiv preprint arXiv:2412.14056, Dec. 2024. [Online]. Available: https://arxiv.org/abs/2412.14056

[5]. M. Nauta et al., "From Anecdotal Evidence to Quantitative Evaluation Methods: A Systematic Review on Evaluating Explainable AI," arXiv preprint arXiv:2201.08164, Jan. 2022. [Online]. Available: https://arxiv.org/abs/2201.08164

[6]. V. Turri et al., "Measuring AI Systems Beyond Accuracy," arXiv preprint arXiv:2204.04211, Apr. 2022. [Online]. Available: https://arxiv.org/abs/2204.04211

[7]. S. Mohseni, N. Zarei, and E. D. Ragan, "A Multidisciplinary Survey and Framework for Design and Evaluation of Explainable AI Systems," arXiv preprint arXiv:1811.11839, Nov. 2018. [Online]. Available: https://arxiv.org/abs/1811.11839

[8]. R. R. Hoffman, S. T. Mueller, G. Klein, and J. Litman, "Metrics for Explainable AI: Challenges and Prospects," arXiv preprint arXiv:1812.04608, Dec. 2018. [Online]. Available: https://arxiv.org/abs/1812.04608

[9]. C. Clemm et al., "Towards Green AI: Current status and future research," arXiv preprint arXiv:2407.10237, Jul. 2024. [Online]. Available: https://arxiv.org/abs/2407.10237

[10]. J. Marques-Silva, "Logic-Based Explainability: Past, Present & Future," arXiv preprint arXiv:2406.11873, Jun. 2024. [Online]. Available: https://arxiv.org/abs/2406.11873