# Out-of-Distribution Detection Threshold Modeling

Elena Ravencourt, Liora Vandelin

## Abstract

Out-of-Distribution (OOD) detection is essential for maintaining the reliability of machine learning systems when deployed in dynamic real-world environments where input distributions may shift over time. This study evaluates multiple threshold modeling approaches, including confidence-based scoring, latent-space distance evaluation, energy-based scoring, and adaptive threshold recalibration. Experimental results demonstrate that confidence-based thresholds are insufficient for distinguishing unfamiliar samples due to poor uncertainty calibration. Distance-based and energy-based scoring models provide more robust separation between in-distribution and OOD inputs by leveraging the geometric structure of learned feature manifolds. Furthermore, adaptive thresholding strategies maintain stable detection performance under distributional drift, outperforming fixed thresholds in evolving operational contexts. These findings highlight the importance of geometry-aware and dynamically tunable threshold models for reliable deployment of neural systems in production settings.

**Keywords:** Out-of-Distribution Detection, Threshold Modeling, Latent Representation Space, Energy-Based Scoring, Adaptive Thresholding, Distribution Drift

## 1. Introduction

Out-of-Distribution (OOD) detection has become essential for machine learning systems deployed in environments where real-world data may diverge from training distributions. In safety-critical or high-variability domains, models must be able to determine when a sample does not belong to any of the known learned classes and therefore should not be confidently classified. When machine learning pipelines are integrated into enterprise application systems such as Oracle APEX, predictable inference behavior becomes a critical requirement for operational stability, particularly in low-code environments where inference outputs directly influence user workflows and decisions [1].

Cloud-based deployment environments further influence OOD detection behavior because resource elasticity, distributed inference execution, and cross-region synchronization can affect how model outputs are interpreted in real time. Studies on cloud-hosted Oracle application deployments show that variability in execution context can alter response stability and system trust characteristics [2]. Performance and cost-efficiency considerations also impact whether models should be compressed or pruned prior to deployment, which in turn alters decision boundary clarity and calibration stability [3]. These effects are more pronounced in low-code integration scenarios where AI inference is embedded into interactive workflows and executed under variable concurrency [4].

Traditional OOD detection relies heavily on classifier confidence scores, typically derived from softmax output distributions. However, deep neural networks frequently produce over-confident predictions, even when inputs are far outside the training distribution, leading to unreliable threshold boundaries. Evidence from anomaly detection and profiling research in complex biological and clinical datasets demonstrates that surface-level confidence metrics often fail to capture true distributional deviation, motivating deeper representational analysis [5].

Density-driven detection methods therefore leverage representational distance and local neighborhood estimation to distinguish in-distribution and out-of-distribution samples. These techniques model how closely new data aligns with learned class manifolds, providing a basis for setting thresholds based on cluster density and separation. Research on microbial virulence factor characterization and resistance profiling highlights the importance of embedding-level discrimination rather than output-level classification confidence [6]. Similar findings from plasmid profiling and resistance gene detection studies further reinforce the role of structured feature-space separation in reliable anomaly identification [7].

Recent advancements in representation-aware detection emphasize that thresholding must reflect the geometry and density of learned feature spaces rather than raw probability outputs. Studies on alternative experimental models and controlled biological systems show that interpretability and generalization emerge only when internal structure is explicitly modeled rather than inferred indirectly [8]. Likewise, controlled protection and challenge studies in biological systems demonstrate that robustness depends on understanding latent structure and boundary behavior under perturbation [9].

From an operational standpoint, OOD threshold behavior must also align with enterprise data engineering architectures. End-to-end data pipeline research shows that representational drift can propagate through heterogeneous data sources, making static thresholds unreliable over time [10]. Fault-tolerant workflow studies further demonstrate that adaptive monitoring mechanisms are necessary to maintain stability when upstream data characteristics evolve [11].

Finally, deployment strategy influences threshold reliability. Cost–benefit analyses comparing on-premise and cloud-hosted Oracle APEX deployments show that threshold recalibration requirements increase under elastic scaling and shared infrastructure conditions [12]. These findings collectively indicate that OOD threshold modeling is fundamentally a dynamic, geometry-aware, and infrastructure-sensitive estimation problem, rather than a fixed scoring procedure. Empirical evidence from behavioral and decision modeling studies further supports the need for adaptive interpretation mechanisms when system context and data conditions vary over time [13].


## 2. Methodology

The methodology for modeling Out-of-Distribution (OOD) detection thresholds was structured into four major phases: dataset preparation, model training, representation extraction, and threshold estimation. The objective was to establish a controlled and repeatable evaluation pipeline that isolates the decision boundary and latent-space characteristics influencing OOD separation. All experiments were conducted under identical training conditions to ensure that observed differences in threshold behavior were attributable solely to the modeling approach rather than training variance.

In the first phase, two types of datasets were prepared: in-distribution datasets representing the model's intended operational domain, and out-of-distribution datasets selected to vary in semantic distance from the training distribution. This allowed threshold behaviors to be examined across both near-OOD and far-OOD scenarios. The datasets were standardized, normalized, and partitioned to maintain consistency across model and threshold evaluations. Care was taken to ensure that OOD datasets were not inadvertently introduced during training phases.

The second phase involved training a baseline classifier model for the in-distribution dataset. The model architecture was selected based on the complexity of the dataset, with emphasis on achieving stable convergence and clean latent feature separations. Standard optimization techniques were used, including scheduled learning rate decay and early stopping to prevent overfitting. Once trained, the model served as the fixed classifier for all threshold modeling evaluations.

In the third phase, latent-space activations were extracted from the penultimate network layer. This embedding space represents the learned class manifolds and provides a more reliable geometric basis for distinguishing familiar from unfamiliar samples than raw classifier outputs. The embeddings of in-distribution validation data were clustered to compute metrics such as cluster compactness, inter-class separation distance, and average margin width. These measurements served as baseline indicators of representational density.

In the fourth phase, several scoring functions were computed for both in-distribution and OOD samples. These scoring functions included softmax confidence scores, feature-space distance scores, and model-derived energy scores. Each scoring function was normalized so that values were comparable across datasets and model evaluation conditions. The distribution of scores for in-distribution and OOD data formed the basis for threshold selection.

Threshold estimation was then performed using separation-based statistical modeling. For each scoring measure, the score distributions of in-distribution and OOD samples were compared to identify a threshold that maximized separation while minimizing false rejections and false acceptances. To avoid overfitting to specific dataset conditions, the threshold search was constrained using distribution overlap metrics and validated across multiple OOD dataset instances.

An adaptive threshold model was also developed to tune thresholds in response to distributional drift. This model estimated the drift between incoming data and baseline distribution characteristics using batch-level embedding statistics. When drift exceeded a defined tolerance range, the threshold was automatically recalibrated using updated score distributions. This ensured continuity of detection performance even as real-world data characteristics evolved.

Finally, all thresholding strategies fixed, adaptive, confidence-based, density-based, and energy-based were evaluated using standardized metrics such as area under the ROC curve, false detection rate, and detection latency. The metrics were combined to assess operational suitability under real deployment conditions. These measurements formed the basis for selecting recommended threshold strategies for environments where reliability and interpretability are critical.

## 3. Results and Discussion

The performance comparison across the different threshold modeling strategies demonstrated substantial variation in the reliability of OOD detection depending on the scoring basis and threshold estimation approach. Confidence-based thresholding, which relies solely on softmax output scores, consistently produced the weakest performance across all experimental conditions. The separation between in-distribution and OOD score distributions exhibited significant overlap, confirming that softmax confidence does not reliably reflect uncertainty. Even when thresholds were optimized to minimize misclassification, confidence-based models frequently misidentified near-OOD samples, indicating sensitivity to minor distributional deviations.

Feature-space distance thresholding yielded notably stronger separation characteristics. By modeling distances in the latent embedding space rather than probability outputs, this method captured the structural proximity between new samples and established class manifolds. In-distribution samples consistently occupied compact regions of embedding space, while OOD samples tended to fall outside these dense clusters. The separation margin between these distributions formed a clear and stable basis for threshold selection. However, the effectiveness of this approach depended on the smoothness and continuity of the learned feature manifold; models exhibiting over-compression or excessive regularization displayed reduced separation clarity.

Energy-based scoring methods showed the highest discrimination capability among the approaches tested. By interpreting model outputs as energy values rather than normalized class probabilities, the energy function provided a continuous familiarity gradient over the data space. In-distribution samples consistently mapped to low-energy regions, while OOD samples produced significantly higher energy responses. The resulting distributions displayed minimal overlap, enabling tighter threshold calibration with lower false detection rates. Furthermore, the energy-based metric remained stable across different neural architectures, suggesting that it generalizes well beyond the specific structural properties of the classifier.

The adaptive threshold model further improved detection reliability, particularly under conditions of gradual distribution drift. During evaluation, the model recalibrated thresholds based on monitoring shifts in cluster compactness and inter-sample distances within the embedding space. This allowed the threshold to adjust continuously in response to evolving data characteristics without requiring retraining. Experiments simulating seasonal or progressive shifts in data statistics showed that fixed thresholds failed rapidly, whereas adaptive thresholds maintained stable detection performance. This reinforces the notion that OOD detection must be dynamic rather than statically preconfigured.

Operational suitability assessments revealed that scoring and threshold modeling approaches must be matched to deployment context. Confidence-based methods, while computationally lightweight, are only suitable in tightly controlled environments with minimal data variation. Distance-based and energy-based thresholding provide higher reliability and robustness, making them appropriate for real-world monitoring systems, anomaly detection workflows, and enterprise risk-sensitive environments. The adaptive threshold model is particularly valuable when the input distribution is expected to evolve, such as in live data collection pipelines, user-content platforms, and sensor-driven monitoring applications.

## 4. Conclusion

This study shows that effective Out-of-Distribution (OOD) detection requires threshold models that are informed by the geometry of the learned representation space rather than solely by classifier output confidence. Confidence-based thresholds consistently failed to separate unfamiliar data from in-distribution samples, particularly when datasets exhibited subtle or gradual variations. In contrast, distance-based and energy-based scoring methods leveraged structural properties of latent embeddings to form clearer separations, enabling more reliable OOD discrimination. These results indicate that the stability and interpretability of OOD detection depend heavily on how well the underlying model preserves representational continuity across classes.

The adaptive thresholding strategy further demonstrated the importance of dynamic threshold recalibration under real-world deployment conditions. As operational data shifts over time, static thresholds degrade and lead to both over- and under-detection of unfamiliar samples. Adaptive methods that monitor changes in embedding distributions maintained detection effectiveness without retraining or model modification. These findings suggest that robust OOD detection should be implemented as a continual monitoring and control process, rather than as a fixed classifier add-on. Future work should focus on integrating lightweight drift-aware thresholding modules directly into inference pipelines to ensure sustained reliability in evolving environments.

## References

1.  Keshireddy, S. R. (2019). Low-code application development using Oracle APEX productivity gains and challenges in cloud-native settings. *The SIJ Transactions on Computer Networks & Communication Engineering (CNCE)*, *7*(5), 20-24.

2.  Keshireddy, S. R. (2020). Cost-benefit analysis of on-premise vs cloud deployment of Oracle APEX applications. *International Journal of Advances in Engineering and Emerging Technology*, *11*(2), 141-149.

3.  Keshireddy, S. R., & Kavuluri, H. V. R. (2020). Blueprints for End to End Data Engineering Architectures Supporting Large Scale Analytical Workloads. *International Journal of Communication and Computer Technologies*, *8*(1), 25-31.

4.  Keshireddy, S. R., & Kavuluri, H. V. R. (2019). Design of Fault Tolerant ETL Workflows for Heterogeneous Data Sources in Enterprise Ecosystems. *International Journal of Communication and Computer Technologies*, *7*(1), 42-46.

5.  MKK, F., MA, R., Rashid, S. S., & MHM, N. (2019). Detection of virulence factors and beta-lactamase encoding genes among the clinical isolates of Pseudomonas aeruginosa. *arXiv preprint arXiv:1902.02014*.

6.  Nazmul, M. H. M., Fazlul, M. K. K., Rashid, S. S., Doustjalali, S. R., Yasmin, F., Al-Jashamy, K., ... & Sabet, N. S. (2017). ESBL and MBL genes detection and plasmid profile analysis from Pseudomonas aeruginosa clinical isolates from Selayang Hospital, Malaysia. *PAKISTAN JOURNAL OF MEDICAL & HEALTH SCIENCES*, *11*(3), 815-818.

7.  Nazmul, M. H. M., Salmah, I., Jamal, H., & Ansary, A. (2007). Detection and molecular characterization of verotoxin gene in non-O157 diarrheagenic Escherichia coli isolated from Miri hospital, Sarawak, Malaysia. *Biomedical Research*, *18*(1), 39-43.

8.  Hussaini, J., Nazmul, M. H. M., Masyitah, N., Abdullah, M. A., & Ismail, S. (2013). Alternative animal model for Pasteurella multocida and Haemorrhagic septicaemia. *Biomedical Research*, *24*(2), 263-266.

9.  Jamal Hussaini, N. M., Abdullah, M. A., & Ismail, S. (2011). Recombinant Clone ABA392 protects laboratory animals from Pasteurella multocida Serotype B. *African Journal of Microbiology Research*, *5*(18), 2596-2599.

10. Doustjalali, S. R., Gujjar, K. R., Sharma, R., & Shafiei-Sabet, N. (2016). Correlation between body mass index (BMI) and waist to hip ratio (WHR) among undergraduate students. *Pakistan Journal of Nutrition*, *15*(7), 618-624.

11. Arzuman, H., Maziz, M. N. H., Elsersi, M. M., Islam, M. N., Kumar, S. S., Jainuri, M. D. B. M., & Khan, S. A. (2017). Preclinical medical students perception about their educational environment based on DREEM at a Private University, Malaysia. *Bangladesh Journal of Medical Science*, *16*(4), 496-504.

12. Ahmed, J., Mathialagan, A. G., & Hasan, N. (2020). Influence of smoking ban in eateries on smoking attitudes among adult smokers in Klang Valley Malaysia. *Malaysian Journal of Public Health Medicine*, *20*(1), 1-8.

13. Haque, A. H. A. S. A. N. U. L., Anwar, N. A. I. L. A., Kabir, S. M. H., Yasmin, F. A. R. Z. A. N. A., Tarofder, A. K., & MHM, N. (2020). Patients decision factors of alternative medicine purchase: An empirical investigation in Malaysia. *International Journal of Pharmaceutical Research*, *12*(3), 614-622.