

Agnostic Label-Only Membership Inference Attack on Two-Tower Neural Networks for Recommendation Systems

Anonymous authors

Paper under double-blind review

Abstract

This paper presents an innovative adaptation of the Agnostic Label-Only Membership Inference Attack (ALOA) specifically designed for two-tower neural network (NN) models used in recommendation systems. Unlike traditional membership inference attacks that focus on categorical outputs, our approach targets models that produce continuous vector embeddings. We propose a comprehensive methodology that employs synthetic datasets, shadow model training, and a suite of perturbation techniques to evaluate model robustness using the Maximum Mean Discrepancy (MMD) metric. Experimental results demonstrate that the attack model achieves exceptionally high accuracy and precision in distinguishing whether data is part of the original training dataset D_o^{train} , even without direct access to it. These findings extend the theoretical framework of membership inference attacks to continuous output spaces and highlight vulnerabilities in modern recommendation systems.

1 Introduction and Related Work

The rapid advancement and refinement of large language models have significantly streamlined the implementation of deep neural networks (DNNs), facilitating their widespread adoption across diverse domains. This progress has attracted substantial investments, resulting in an increasing number of companies and deep learning (DL) applications becoming integral to our daily lives. DL models are now omnipresent, influencing activities from browsing websites and watching videos to online shopping and locating the best local restaurants. While DL models, often branded as “AI” features, enhance convenience and efficiency, they also introduce substantial risks related to privacy leakage and potential poisoning attacks that can disrupt model functionality. Sophisticated attackers can exploit DL models trained on personal data to expose sensitive information, such as individuals’ relationships and residential locations, to unauthorized parties.

Among the various types of attacks targeting DL models, a few stand out due to their potential impact. The Label-Only Membership Inference Attack (MIA), introduced by Shokri et al. (2017b), is a prominent example. MIA enables attackers to determine whether specific data points were included in the training dataset of a target DL model using only the output labels. When combined with Generative Adversarial Networks (GANs), this method escalates the risk by enabling the reconstruction of sensitive data, particularly when models are trained on confidential information belonging to individuals or organizations Fredrikson et al. (2015). Additionally, the industry’s shift toward data trading by some governments and corporations exacerbates the risk of digital property theft, making attacks like MIA increasingly concerning.

Building upon these foundational attacks, Monreale et al. (2023) introduced the Agnostic Label-Only Attack (ALOA). ALOA generalizes MIA by removing dependence on the original dataset’s distribution, thereby broadening the scope and applicability of membership inference attacks. However, existing research on ALOA predominantly focuses on categorical models with discrete labels, leaving a significant gap in understanding its implications on more complex models that produce continuous outputs, such as those generating vector embeddings.

This paper addresses this gap by pioneering the application of ALOA to **two-tower neural network (NN) models**, which generate outputs in continuous vector fields rather than discrete labels. Two-tower NN models are widely employed in recommendation systems and information retrieval tasks; see Kato et al.

(2019) and He et al. (2017), where they process information from two distinct towers—users and items—and output embeddings for each. These embeddings are then utilized by ranking systems to generate personalized recommendations based on relevance and geometric patterns. Unlike previous studies that concentrate on categorical outputs, our approach tackles the unique challenges posed by continuous output spaces, making it a novel contribution to the field.

The primary contribution of this research lies in developing and validating an enhanced attack methodology that effectively infers membership and reconstructs sensitive information from vector-based embeddings in two-tower NN models. This advancement is significant because two-tower NN models are prevalent in recommendation systems, where the integrity and privacy of user and item embeddings are crucial. Furthermore, our research operates under the realistic assumption that the attacker has access only to input features obtainable from the target model, without prior knowledge of the original dataset’s distribution. This scenario mirrors real-world conditions more closely than many existing studies, thereby increasing the practical relevance and applicability of our findings.

By demonstrating that ALOA can be successfully adapted to complex embedding-based models, this paper not only extends the theoretical framework of membership inference attacks but also provides empirical evidence of its feasibility in practical settings. This highlights the urgent need for enhanced security measures to protect sensitive data in sophisticated DL models.

Privacy has consistently been a critical concern across various fields, receiving significant attention due to the increasing reliance on data-driven technologies. Privacy information leakage can occur through multiple avenues, including direct data access or indirectly via deep learning (DL) models. Since DL models are trained on datasets and subsequently converge into mathematical representations, they inherently encapsulate information from the training data. Consequently, if an adversary intentionally queries the model, it becomes feasible to exploit and extract sensitive information or at least gain substantial knowledge about the original dataset.

In the realm of data privacy, the initial step involves employing privacy risk assessment methodologies to evaluate the privacy risks associated with users within a dataset. Based on these assessments, appropriate privacy protection techniques are developed to safeguard the data or the DL models themselves. Common privacy protection strategies in DL include randomization, differential privacy, and k-anonymity; see Dwork & Roth (2014), Sweeney (2002) and Duchi et al. (2018). Additionally, emerging research in data decomposition, such as the concept of data elements, aims to break down data into smaller components. This decomposition not only maintains the efficacy of DL models trained on the transformed data but also significantly reduces the risk of data leakage, sometimes achieving zero risk; see Cummings & Zarsky (2021).

Focusing on privacy risk assessment within DL models, the PRUDence framework, proposed by Pratesi et al., facilitates a systematic evaluation of empirical privacy risks concerning specific privacy attacks on data; see Pratesi et al. (2020). PRUDence simulates an adversary equipped with the knowledge necessary to maximize the privacy risk for each individual in the dataset by generating all possible background knowledge that the adversary might possess and assessing the associated risks under the most adverse conditions.

Over the past few years, several methodologies have been developed to target DL models directly. Notable among these are:

- **Membership Inference Attack (MIA):** Introduced by Shokri et al., MIA aims to determine whether a particular data point was part of the training dataset of a classification model by leveraging the model’s outputs to infer membership, posing significant privacy risks, especially for sensitive or proprietary datasets; see Shokri et al. (2017b).
- **Reconstruction Attack:** Proposed by Fredrikson et al., this attack seeks to reconstruct one or more training samples along with their corresponding labels by exploiting the model’s parameters and outputs, thereby undermining data privacy; see Fredrikson et al. (2015).
- **Property Inference Attack:** This attack extracts unintended and non-specific information learned by the model, such as the overall distribution of the training dataset. Unlike MIA, which focuses on

individual data points, property inference targets aggregate properties, potentially revealing sensitive characteristics of the dataset; see Yeom et al. (2018).

Building upon these foundational attacks, Choquette-Choo et al. introduced the **LABELONLY attack**, which extends MIA by operating without access to the probability vector outputs of the model. This attack demonstrates the ability to infer membership using only the final labels, thereby lowering the barrier for successful attacks in scenarios where full probability information is unavailable; see Choquette-Choo et al. (2019).

Further advancing this line of research, Anna Monreale et al. proposed the **Agnostic Label-Only Membership Inference Attack (ALOA)** in Monreale et al. (2023). ALOA enhances the LABELONLY attack by eliminating the requirement for knowledge about the training dataset’s distribution, allowing broader applicability and effectiveness in diverse settings. Monreale et al. demonstrated that ALOA achieves desirable performance metrics in privacy attacks on categorical models, marking a significant step forward in membership inference methodologies; see Monreale et al. (2023).

Despite these advancements, existing research on ALOA has predominantly concentrated on categorical models with discrete labels, leaving a substantial gap in understanding its implications on more complex models that produce continuous outputs, such as those generating vector embeddings. Addressing this gap, our research applies ALOA to **two-tower neural network (NN) models** commonly used in recommendation systems. These models output embeddings in a continuous vector field, representing users and items, which are then utilized by ranking systems to generate personalized recommendations based on relevance and geometric patterns.

2 Background

Before proceeding with the details of validating ALOA efficiency in complex multilayer neural networks, we introduce some basic notions that will be used or useful for understanding the approach and modifications used in our research.

2.1 Classifier

A *classifier* $b : X^{(n)} \rightarrow Y$ is a mapping that assigns an input instance x from a feature space $X^{(n)}$, in which x consists of m feature variables, to a categorical label y in a target space Y of cardinality $L = |Y|$. The classifier produces a categorical output y , expressed as $b(x) = y$. Each feature x can be formally represented as $x = \{(a_1, v_1), (a_2, v_2), (a_3, v_3), \dots, (a_n, v_n)\}$, where a_i denotes the feature name and v_i represents its corresponding value. The features a_i may be either continuous or categorical.

The Agnostic Label-Only Attack (ALOA) assumes that the target classifier b is accessible and can be queried without restrictions. In practice, the classifier b is trained on a dataset D_o^{train} . However, in an adversarial setting, the objective is often to reconstruct D_o^{train} , which remains inaccessible to the attacker.

2.2 Black-Box Model

In existing literature, a *black-box setting* refers to an attack scenario where neither the classifier’s architecture nor the training dataset is known. In our work, we extend this concept to neural networks. Hereafter, the term “Black-Box model” refers to the target model being attacked without any internal knowledge of its structure or training data.

2.3 Shadow Model

Shadow models are deep learning models trained to approximate the behavior of a Black-Box model. They play a crucial role in conducting Membership Inference Attacks (MIA) by generating datasets on which attack models $\{A(\cdot)\}$ are trained. A set of shadow models $\{s(\cdot)\}$ is trained using a synthetic dataset \tilde{X} , which is constructed to resemble the structure of the Black-Box model’s input. The corresponding outputs

predicted by the Black-Box model are denoted as $\tilde{Y} = b(\tilde{X})$. The shadow models are subsequently trained on the dataset: $D_s^{\text{train}} = \{(\tilde{x}_1, \tilde{y}_1), (\tilde{x}_2, \tilde{y}_2), (\tilde{x}_3, \tilde{y}_3), \dots, (\tilde{x}_n, \tilde{y}_n)\}$, with the objective of replicating the behavior of b . Notably, D_s^{train} is disjoint from the original training dataset D_o^{train} due to the inherent lack of access to D_o^{train} . The purpose of training a shadow model is to approximate and ultimately reconstruct D_o^{train} . Chapelle et al. (2011) has demonstrated that training a single shadow model to mimic the entire output space can achieve performance comparable to training separate shadow models for each category. Therefore, in this study, we employ a single shadow model.

Shadow models serve as a mechanism to infer information about the Black-Box model. As established in Theorem 1, for any $x \in X$ such that $b(x) = y$, there exists a function $s_i(\cdot)$ within the set of shadow models $\{s_1(\cdot), s_2(\cdot), s_3(\cdot), \dots, s_m(\cdot)\}$ such that $s_i(x) = y$.

Theorem 1 (Representation Theorem for Shadow Models, adapted from Shokri et al. (2017a)). *For all $x \in X$ such that $b(x) = y$, there exists a function $s_i(\cdot)$ in the set of shadow models $\{s_1(\cdot), s_2(\cdot), s_3(\cdot), \dots, s_m(\cdot)\}$ such that $s_i(x) = y$.*

As a consequence, the probability vectors predicted by $\{s_i(\cdot)\}$ are expected to closely approximate those produced by b . This similarity can be exploited to infer whether a specific data point belongs to the original training dataset.

2.4 Membership Inference Attack (MIA)

Shokri et al. (2017b) demonstrated that a classifier b trained on a dataset D_{train}^b retains relationships between its outputs and the original data. To perform an attack on b , a model $A(\cdot)$ is trained to identify whether a data point was included in D_{train}^b . A synthetic dataset \bar{D}_s^{train} is created and queried through b , generating outputs O_s . The attack model $A(\cdot)$ is a binary classifier that predicts *IN* if the input data belongs to D_{train}^b , or *OUT* otherwise. $A(\cdot)$ is trained on a dataset $D_s^{\text{train}} = \{(x_i, y_i)\}_s$, where each x_i includes the label predicted by b and the probability vector \bar{y}_i of length L . The elements of \bar{y}_i are obtained by querying one of the shadow models $\{s(\cdot)\}$ trained on $\bar{D}_s^{\text{train}} \cup O_s$ to mimic b 's behavior for specific categories. Since $\{s(\cdot)\}$ outputs vectors of dimension L , there are L models in $\{A(\cdot)\}$, each corresponding to one category. The model $A(\cdot)$ computes the probability $\Pr\{(x, y) \in D_{\text{train}}^b\}$ that the input data is part of D_{train}^b based on the distribution in \bar{y}^i .

2.5 Label-Only Membership Inference Attack (LabelOnly)

Choquette-Choo et al. (2019) introduced the LabelOnly attack, which relaxes the requirements of MIA by operating without access to the probability vector outputs of the model. Unlike MIA, which requires probability vectors, LabelOnly relies solely on hard labels predicted by b . It utilizes the model's robustness to perturbations to infer membership. Since deep learning models often incorporate perturbations during training to prevent overfitting, their outputs demonstrate higher resilience to minor changes for training data. The LabelOnly model $A_{LO}(\cdot)$ employs shadow models to infer membership. A shadow model $s_{LO}(\cdot)$ is trained on $D_s^{\text{train}} = \{(x^i, y^i)\}$ with a structure and distribution similar to D_o^{train} . The predicted labels by b are recorded in D_s^{train} . The model $A_{LO}(\cdot)$ is trained by perturbing D_s^{train} to obtain predicted labels $\tilde{y}^i = s(\tilde{D}_s^{\text{train}})$, where $\tilde{D}_s^{\text{train}}$ is the perturbed dataset. The robustness score for each data point x^i is computed as the percentage of labels that remain unchanged after perturbation. An iterative thresholding procedure classifies records as *IN* or *OUT* based on their robustness scores, with higher scores indicating a higher probability of being in D_o^{train} .

2.6 Agnostic Label-Only Membership Inference Attack (ALOA)

Monreale et al. (2023) advanced the LabelOnly attack by removing the necessity of knowledge about the original data distribution. Instead, they revised the function for computing the robustness score. The revised robustness score is defined as:

$$\text{rScore}_{x_s^i}(N_{x_s^i}) = \begin{cases} 0 & , \text{if } s(x_s^i) \neq b(x_s^i); \\ \frac{1}{|N_{x_s^i}|} \sum_{x' \in N_{x_s^i}} F(s(x'), s(x_s^i)) & , \text{otherwise,} \end{cases}$$

where $N_{x_s^i}$ denotes the set of entries x' generated by introducing noise around x_s^i . The function $F(s(x'), s(x_s^i))$ returns a value close to 0 if $s(x') \neq s(x_s^i)$ and close to 1 otherwise. Thus, the robustness score ranges between 0 and 1.

Instead of requiring knowledge of the original dataset D_o^{train} , ALOA generates a synthetic dataset x_s^{train} with the same structure as the Black-Box model’s input. This dataset is queried through both the shadow model and the target model. The set

$$D_s^{\text{train}} = \{x_s^i \in x_s^{\text{train}} \mid s(x_s^i) = b(x_s^i)\}$$

is constructed. Noise is then introduced to entries in D_s^{train} to compute robustness scores. Subsequently, ALOA proceeds similarly to the LabelOnly approach without requiring distribution knowledge.

Mathematically, our approach leverages concepts from the mathematics of deep learning by Vidal (2018) to analyze the stability and robustness properties of neural network classifiers under perturbations, providing theoretical support for the ALOA methodology.

3 Proposed Approach

To implement ALOA on Two-Tower NN models, we follow a structured methodology comprising model setup, shadow model training, perturbation application, robustness scoring, and attack model training. This section provides a comprehensive overview of each step, illustrating how they collectively contribute to the successful execution of the ALOA attack.

Given the Black-Box model b , we extracted the input structure required by the model. A synthetic dataset

$$X_{\text{user}}^{\text{train}} \cup X_{\text{movie}}^{\text{train}} = X_s^{\text{train}}$$

with an identical structure was generated. Focusing on user data, which is typically more informative and valuable, we employed two distinct approaches to train the shadow model, thereby influencing the overall attack accuracy:

1. **Complete Approach:** Utilizes the entire synthetic dataset $X_s^{\text{train}} = X_{\text{user}}^{\text{train}} \cup X_{\text{movie}}^{\text{train}}$, where both user and movie data are synthetically generated to mimic the distribution of the original data.
2. **Pseudo Approach:** Uses only the user portion $X_s^{\text{train}} = X_{\text{user}}^{\text{train}}$, and assigns a fixed pseudo entry for movie data (e.g., a constant vector of all 1’s) for every record.

For each $x_i \in X_s^{\text{train}}$, the corresponding labels $Y_s^{\text{train}} = \{y_i \mid y_i = b(x_i)\}$ were generated by querying the Black-Box model b . The shadow model s was then trained on the dataset $D_s^{\text{train}} = X_s^{\text{train}} \cup Y_s^{\text{train}}$ to accurately replicate b ’s prediction behavior.

To systematically compute robustness scores, we introduce a series of perturbations to the user data in D_s^{train} . These perturbations are designed to emulate real-world noise and inaccuracies commonly encountered in user data. Based on variations in different input features—such as ID, Gender, Zip Code, Bucketized Age, and Occupation—we define the following perturbation methods:

1. User ID Perturbation

- **Method:** Add or subtract a random integer within the range $[-10, 10)$ to the original user ID.
- **Outcome:** A new dataset with slightly altered user IDs, maintaining valid values within acceptable bounds.

2. Gender Flip Perturbation

- **Method:** Invert the binary gender value; if coded as 0 (e.g., Female), change to 1 (e.g., Male), and vice versa.
- **Outcome:** A dataset where each user’s gender has been inverted, with all other attributes remaining unchanged.

3. Zip Code Perturbation

- **Method:** Add or subtract a random integer within the range $[-50, 50]$ to each user zip code.
- **Outcome:** A dataset with slightly varied zip codes that remain valid and within the model’s expected embedding range.

4. User Bucketized Age Shift Perturbation

- **Method:** Shift each user’s bucketized age to an adjacent age bucket. For middle buckets, randomly decide to shift up or down one bucket. For boundary buckets, shift to the only available adjacent bucket.
- **Outcome:** A dataset where each user’s age bucket is shifted to a neighboring bucket, simulating slight age misclassifications or boundary adjustments.

5. Occupation Change Perturbation

- **Method:** Randomly assign a different occupation label to each user from the set of valid occupation labels, excluding the current one. Ensure that the new label remains within the model’s expected range.
- **Outcome:** A dataset where each user has a randomly altered occupation label within valid bounds, simulating misclassification or noise in occupation data.

Following the application of perturbations, we ensured data integrity and model compatibility by verifying that each altered feature remained within the expected value ranges. This validation prevents runtime errors during model inference, such as out-of-bounds errors in embedding layers. Since only one target feature is modified at a time while all other features remain unchanged, any observed changes in the model’s behavior can be solely attributed to the perturbed feature.

Each perturbed dataset was saved in dedicated directories corresponding to the type of perturbation (e.g., `user_id_perturbation`, `gender_flip`, etc.) to facilitate systematic testing and comparison across different perturbation types. Identical perturbations were also applied to D_o^{train} for baseline comparisons.

To quantify the impact of each perturbation, we computed the Maximum Mean Discrepancy (MMD) between the embeddings obtained from the original and perturbed datasets. MMD is a statistical measure used to compare the difference between two distributions. Given two sets of embeddings, $\mathcal{X} = \{\mathbf{u}_i\}_{i=1}^n$ from the original data and $\mathcal{Y} = \{\mathbf{v}_j\}_{j=1}^m$ from the perturbed data, the squared MMD is computed as:

$$\text{MMD}^2(\mathcal{X}, \mathcal{Y}) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n k(\mathbf{u}_i, \mathbf{u}_j) + \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m k(\mathbf{v}_i, \mathbf{v}_j) - \frac{2}{nm} \sum_{i=1}^n \sum_{j=1}^m k(\mathbf{u}_i, \mathbf{v}_j),$$

where $k(\cdot, \cdot)$ is a kernel function (e.g., the Gaussian kernel defined as $k(\mathbf{u}, \mathbf{v}) = \exp\left(-\frac{\|\mathbf{u}-\mathbf{v}\|^2}{2\sigma^2}\right)$). A larger MMD value indicates a larger discrepancy between the two distributions, reflecting a greater impact of the perturbation on the model’s output.

In order to systematically apply these perturbations and facilitate reproducibility, we define a function called `perturbation` that performs the specified perturbations on the dataset. Below is Algorithm 1, which outlines the detailed procedure for generating perturbed datasets.

Algorithm 1 Perturbation Function for ALOA**Require:** D - Original user dataset**Require:** $perturb_type$ - Type of perturbation to apply ({User ID Perturbation, Gender Flip, Zip Code Perturbation, User Bucketized Age Shift, Occupation Change})**Require:** $params$ - Additional parameters required for the perturbation (e.g., range limits)**Ensure:** $D_{perturbed}$ - Perturbed user dataset

```

1: if  $perturb\_type = \text{"User ID Perturbation"}$  then
2:   Convert  $user\_id$  to numeric values
3:   Generate random integers  $r_i \in [-10, 10)$  for each user
4:   Update  $user\_id$ :  $user\_id \leftarrow user\_id + r_i$ 
5:   Clip  $user\_id$  within range  $[0, 6040]$ 
6: else if  $perturb\_type = \text{"Gender Flip"}$  then
7:   Confirm that  $user\_gender$  is binary (0 or 1)
8:   Flip  $user\_gender$ :  $user\_gender \leftarrow 1 - user\_gender$ 
9: else if  $perturb\_type = \text{"Zip Code Perturbation"}$  then
10:  Convert  $user\_zip\_code$  to numeric values
11:  Generate random integers  $r_i \in [-50, 50)$  for each user
12:  Update  $user\_zip\_code$ :  $user\_zip\_code \leftarrow user\_zip\_code + r_i$ 
13:  Clip  $user\_zip\_code$  to valid range  $[30000, 33439]$ 
14:  Map clipped  $user\_zip\_code$  to embedding indices:  $user\_zip\_code \leftarrow$ 
     $\max(0, \min(3439, user\_zip\_code - 30000))$ 
15: else if  $perturb\_type = \text{"User Bucketized Age Shift"}$  then
16:  Identify and sort unique  $user\_bucketized\_age$  values
17:  Create mapping  $age\_map$  from age bucket to index
18:  for each user in  $D$  do
19:    Determine current age bucket index  $idx$ 
20:    if  $idx = 0$  then
21:      Shift to the next higher bucket
22:    else if  $idx = \text{last index}$  then
23:      Shift to the next lower bucket
24:    else
25:      Randomly choose to shift up or down one bucket
26:    end if
27:    Update  $user\_bucketized\_age$ 
28:  end for
29: else if  $perturb\_type = \text{"Occupation Change"}$  then
30:  Identify valid occupation labels  $O = \{0, 1, \dots, 20\}$ 
31:  for each user in  $D$  do
32:    Current occupation  $o$ 
33:    Define possible new occupations:  $O' = O \setminus \{o\}$ 
34:    Randomly select new occupation  $o' \in O'$ 
35:    Assign  $o'$  to  $user\_occupation\_label$ 
36:  end for
37: end if
38: return  $D_{perturbed}$ 

```

4 Experiment

This section delineates the experimental setup and results of applying ALOA on Two-Tower NN models. The experiments assess the effectiveness of ALOA in inferring membership based on perturbations in user data.

4.1 Black-Box Model Robustness Testing

We utilized the multi-stage two-tower recommender project from GitHub created by Keivanipchihagh (2023), which provides a comprehensive dataset for training and evaluating two-tower neural networks in recommendation systems. The dataset comprises:

- **User Information:** User ID, user zip code, gender, age bracket, and user occupation.
- **Movie Information:** Movie ID, movie name, movie genre, and release year.
- **Rating Information:** Records of user-movie interactions along with the ratings provided.

Using PyTorch, we constructed the Two-Tower NN model and exported it to ONNX format to facilitate efficient querying as a Black-Box model.

We explored two distinct approaches to train the shadow model, addressing the challenge of targeting user data:

1. **Complete Approach:** Utilizes the entire synthetic dataset $X_s^{\text{train}} = X_{\text{user}}^{\text{train}} \cup X_{\text{movie}}^{\text{train}}$, where both user and movie data are synthetically generated to mimic the distribution of the original data.
2. **Pseudo Approach:** Uses only the user portion $X_s^{\text{train}} = X_{\text{user}}^{\text{train}}$, and assigns a fixed pseudo entry for movie data (e.g., a constant vector of all 1's) for every record.

For each $x_i \in X_s^{\text{train}}$, the corresponding labels

$$Y_s^{\text{train}} = \{y_i \mid y_i = b(x_i)\}$$

were generated by querying the Black-Box model b . The shadow model s was then trained on the dataset

$$D_s^{\text{train}} = X_s^{\text{train}} \cup Y_s^{\text{train}}$$

to accurately replicate b 's prediction behavior.

We applied the five perturbation methods outlined in Section 3 to the shadow training dataset D_s^{train} . Each perturbation was applied individually, and the resulting datasets were stored in separate directories to facilitate systematic evaluation. Following perturbation, we queried the perturbed datasets through the Black-Box model b to obtain the predicted embeddings.

To assess the model's robustness to these perturbations, we computed the Maximum Mean Discrepancy (MMD) between the embeddings of the original and perturbed data:

$$\text{MMD}^2(\mathcal{X}, \mathcal{Y}) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n k(\mathbf{u}_i, \mathbf{u}_j) + \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m k(\mathbf{v}_i, \mathbf{v}_j) - \frac{2}{nm} \sum_{i=1}^n \sum_{j=1}^m k(\mathbf{u}_i, \mathbf{v}_j)$$

where $\mathcal{X} = \{\mathbf{u}_i\}$ represents the embeddings from the original data X , $\mathcal{Y} = \{\mathbf{v}_i\}$ represents the embeddings from the perturbed data \tilde{X} , and $k(\cdot, \cdot)$ is a kernel function (e.g., the Gaussian kernel). A larger MMD value indicates a larger discrepancy between the two distributions, reflecting a greater impact of the perturbation on the model's output.

The robustness scores for different perturbations are presented in Tables 1, 2, 3, 4, and 5.

Table 1: Robustness Scores for User ID Perturbation

Dataset	Average	Median	Min	Max
Original data with pseudo movie entries	0.145715	0.129626	0	0.806357
Original data with real movie entries	0.145715	0.129626	0	0.806357
Synthetic data with pseudo movies	0.098559	0.088099	0	0.427482
Synthetic data with real movies	0.098559	0.088099	0	0.427482

Table 2: Robustness Scores for Zip Code Perturbation

Dataset	Average	Median	Min	Max
Original data with pseudo movie entries	0.004967	0	0	0.469542
Original data with real movie entries	0.004967	0	0	0.469542
Synthetic data with pseudo movies	0.0000937	0	0	0.000488
Synthetic data with real movies	0.0000937	0	0	0.000488

Table 3: Robustness Scores for User Occupation Perturbation

Dataset	Average	Median	Min	Max
Original data with pseudo movie entries	0.096859	0.083006	0	0.536256
Original data with real movie entries	0.096859	0.083006	0	0.536256
Synthetic data with pseudo movies	0.061248	0.051937	0.013986	0.258884
Synthetic data with real movies	0.061248	0.051937	0.013986	0.258884

Table 4: Robustness Scores for User Gender Perturbation

Dataset	Average	Median	Min	Max
Original data with pseudo movie entries	0.158105	0.140888	0.057147	0.574967
Original data with real movie entries	0.158105	0.140888	0.057147	0.574967
Synthetic data with pseudo movies	0.104109	0.093539	0.052675	0.255438
Synthetic data with real movies	0.104109	0.093539	0.052675	0.255438

Table 5: Robustness Scores for User Age Bucketized Shift Perturbation

Dataset	Average	Median	Min	Max
Original data with pseudo movie entries	0.219849	0.140537	0.013465	0.99374
Original data with real movie entries	0.219849	0.140537	0.013465	0.99374
Synthetic data with pseudo movies	0.016731	0	0	0.923686
Synthetic data with real movies	0.016731	0	0	0.923686

The robustness of the Black-Box model against various user data features indicates that, for most features, the model exhibits greater resistance to perturbations within D_o^{train} . Specifically, most features demonstrate an MMD difference that is significant enough to be detected by carefully tuning the threshold algorithm. Notably, the zip code feature exerts minimal influence on the predicted embeddings; however, even this slight influence results in a sufficiently large MMD value to facilitate membership detection. Furthermore, there is no significant distinction between utilizing dummy movie entries and different synthetic movie entries for training the shadow and attack models, as evidenced by the similar robustness outcomes observed.

4.2 ALOA

Given that the Black-Box model exhibits robustness to perturbations in the original dataset D_o^{train} , we proceed to configure the ALOA attack. The attack encompasses several key components and processes.

First, we define the original training data as $D_o^{\text{train}} = \{(u_i, v_i, r_i)\}_{i=1}^N$, where $u_i \in U$ represents user features, $v_i \in V$ denotes item features, and $r_i \in R$ signifies ratings. For shadow model training, we prepare two distinct datasets: $D_{s,\text{combined}} = \{(u_j, v_j, r_j)\}_{j=1}^N$ and $D_{s,\text{dummy}} = \{(u_k, v_k, r_k)\}_{k=1}^N$, which correspond to the complete and pseudo approaches, respectively. These shadow models, s_{combined} and s_{dummy} , are trained on their respective datasets to emulate the Black-Box model’s behavior.

Next, we generate perturbed data

$$D_p = \{\text{perturbation}(u, \phi) \mid (u, v, r) \in D_s^{\text{train}}, \phi \in \Phi\}$$

where $\Phi = \{\phi_1, \phi_2, \dots, \phi_5\}$ represents different types of perturbations (User ID, Gender Flip, Zip Code, Age Bucket Shift, Occupation Change). These perturbations simulate realistic variations and potential noise in user data, allowing us to assess the model’s robustness and the effectiveness of the ALOA attack.

The ALOA attack procedure involves several sequential steps:

Feature Extraction: For each perturbed user $u' \in D_{\text{perturbed}}$, we extract feature vectors by passing them through the shadow models:

$$\mathbf{f}(u') = s(u', \bullet) \in \mathbb{R}^d$$

where $s(u', \bullet)$ denotes the concatenated outputs from both towers of the Two-Tower shadow model, and d is the dimensionality of the feature vector.

Label Assignment: Binary labels are assigned based on membership status:

$$y(u') = \begin{cases} 1 & \text{if } u' \text{ is a member (from } D_s^{\text{train}}) \\ 0 & \text{if } u' \text{ is a non-member (not from } D_s^{\text{train}}) \end{cases}$$

Dataset Construction: We construct the attack dataset

$$D_{\text{attack}} = \{(\mathbf{f}_i, y_i)\}_{i=1}^{2M},$$

where $\mathbf{f}_i = s(u', \bullet)$ and $y_i \in \{0, 1\}$. Here, data labeled as *IN* corresponds to entries in the shadow model’s training set D_s^{train} , while *OUT* corresponds to data outside D_s^{train} .

Model Training: The attack model training involves the following steps:

1. **Train-Test Split:** Divide D_{attack} into training and testing subsets:

$$D_{\text{train}}, D_{\text{test}} = \text{TrainTestSplit}(D_{\text{attack}}, \text{test_size} = 0.3, \text{random_state} = 42)$$

2. **Classifier Selection:** Employ a Multi-Layer Perceptron (MLP) classifier \mathcal{C} with hyperparameters θ :

$$\mathcal{C} : \mathbb{R}^d \rightarrow \{0, 1\}$$

3. **Hyperparameter Optimization:** Optimize θ using grid search with cross-validation to maximize the ROC AUC score:

$$\theta^* = \arg \max \text{ROC_AUC}(\mathcal{C}_\theta, D_{\text{train}}, y)$$

4. **Training the Classifier:** Train the classifier with the optimized hyperparameters:

$$\mathcal{C}^* = \mathcal{C}_{\theta^*} \leftarrow \text{Train}(\mathcal{C}_{\theta^*}, D_{\text{train}})$$

Evaluation Metrics: We evaluate the performance of \mathcal{C}^* on D_{test} using the following metrics:

- **Accuracy (\mathcal{A}):**

$$\mathcal{A} = \frac{1}{|D_{\text{test}}|} \sum_{i=1}^{|D_{\text{test}}|} \mathbb{I}(\mathcal{C}^*(\mathbf{f}_i) = y_i)$$

- **Precision (\mathcal{P}):**

$$\mathcal{P} = \frac{\sum_{i=1}^{|D_{\text{test}}|} \mathbb{I}(\mathcal{C}^*(\mathbf{f}_i) = 1 \wedge y_i = 1)}{\sum_{i=1}^{|D_{\text{test}}|} \mathbb{I}(\mathcal{C}^*(\mathbf{f}_i) = 1)}$$

- **Recall (\mathcal{R}):**

$$\mathcal{R} = \frac{\sum_{i=1}^{|D_{\text{test}}|} \mathbb{I}(\mathcal{C}^*(\mathbf{f}_i) = 1 \wedge y_i = 1)}{\sum_{i=1}^{|D_{\text{test}}|} \mathbb{I}(y_i = 1)}$$

- **ROC AUC (\mathcal{U}):**

$$\mathcal{U} = \text{ROC_AUC}(\{p_i\}, \{y_i\})$$

where $p_i = \mathcal{C}^*(\mathbf{f}_i)_1$ denotes the predicted probability for class 1 (member).

Threshold Determination: We establish a classification threshold τ using K-Means clustering on the robustness scores $\{r_i\}$, where $r_i = p_i$:

1. **Clustering:**

$$\text{KMeans}(\{r_i\}, k = 2)$$

2. **Threshold Calculation:**

$$\tau = \frac{c_1 + c_2}{2}, \quad \text{where } c_1 \text{ and } c_2 \text{ are the cluster centers}$$

Classification Based on Threshold: We classify each instance based on the determined threshold τ :

$$\mathcal{C}^*(\mathbf{f}_i) = \begin{cases} 1 & \text{if } r_i \geq \tau \\ 0 & \text{otherwise} \end{cases}$$

4.3 Algorithm for ALOA Attack

To provide a clear example of the ALOA attack procedure, we present Algorithm 2, which outlines the steps involved in executing the attack.

Algorithm 2 Agnostic Label-Only Membership Inference Attack (ALOA)

Require: b - Black-Box model

Require: $s_{\text{dummy}}, s_{\text{combined}}$ - Shadow models trained on dummy and combined datasets

Require: D_o^{train} - Original training dataset

Require: Φ - Set of perturbation types ($\{\text{User ID, Gender Flip, Zip Code, Age Bucket Shift, Occupation Change}\}$)

Ensure: \mathcal{C}^* - Trained attack classifier

Ensure: τ - Classification threshold

- 1: **for** each shadow model $s \in \{s_{\text{dummy}}, s_{\text{combined}}\}$ **do**
 - 2: **for** each perturbation type $\phi \in \Phi$ **do**
 - 3: Generate perturbed dataset $D_p = \{\text{perturbation}(u, \phi) \mid u \in D_s^{\text{train}}\}$
 - 4: Extract features $\mathbf{f}(u') = s(u', \bullet)$ for each $u' \in D_p$
 - 5: Assign labels $y(u') \leftarrow 1$ if $u' \in D_s^{\text{train}}$, else $y(u') \leftarrow 0$
 - 6: **end for**
 - 7: **end for**
 - 8: Construct attack dataset $D_{\text{attack}} = \{(\mathbf{f}_i, y_i)\}_{i=1}^{2M}$
 - 9: Split D_{attack} into $D_{\text{train}}, D_{\text{test}}$ using TrainTestSplit with test_size = 0.3 and random_state = 42
 - 10: Initialize Multi-Layer Perceptron (MLP) classifier \mathcal{C} with hyperparameters θ
 - 11: Optimize θ using grid search with cross-validation to maximize ROC AUC on D_{train}
 - 12: Train classifier \mathcal{C}^* with optimized θ^* on D_{train}
 - 13: Evaluate \mathcal{C}^* on D_{test} to obtain metrics $\mathcal{A}, \mathcal{P}, \mathcal{R}, \mathcal{U}$
 - 14: Perform K-Means clustering on robustness scores $\{r_i\}$ to determine threshold τ
 - 15: Classify each instance in D_{test} based on τ
 - 16: **return** \mathcal{C}^*, τ
-

4.4 Attack Results

The performance of the ALOA attack was evaluated under two scenarios. First, using the shadow models' training data, the attack performance is summarized in Table 6.

Table 6: Attack Performance Using Shadow Models

Model Type	Accuracy	Precision	Recall	ROC AUC	Optimal Threshold
Combined	0.9787	0.9789	0.9958	0.9946	0.8451
Dummy	0.9507	1.0000	0.4573	0.9406	0.1273

The results are very promising. Both the accuracy and precision are exceptionally high, demonstrating that the attack model is highly effective at determining whether data is *IN* or *OUT* of the original training dataset D_o^{train} . Initially, we were concerned about potential overfitting; however, the training data for the attack model is significantly imbalanced—with an *IN:OUT* ratio of 4:1 for the complete approach and 1:10 for the pseudo approach. If the attack model were not functioning properly, we would expect an accuracy near 0.5 (akin to random guessing) or a precision that closely mirrors these imbalanced distributions (around 0.8 or 0.1, respectively). In light of this, we are confident that the outstanding performance observed is genuine. For further validation, we applied the attack model to the original training dataset of the Black-Box model by computing the proportion of instances classified as *IN*. This proportion serves as an accuracy measure, and the results are presented in Table 7.

Table 7: Attack Performance on Original Training Data of the Black-Box Model

Model Type	Accuracy (%)
Combined	99.65
Dummy	85.95

Table 7 demonstrates that the attack model retrieves data from the original training dataset D_o^{train} of the Black-Box model with exceptional efficiency. Remarkably, the attack model was trained solely on data from the shadow models and had no direct access to D_o^{train} during the attack process, underscoring its practical applicability. These results confirm that the ALOA method is effective for attacking models that output in a continuous vector space. Moreover, our investigation reveals that for Two-Tower neural networks—or indeed any composite deep learning models—if the synthetic dataset used to train the shadow model lacks sufficient diversity, the attack’s efficiency is reduced. This intriguing outcome highlights the importance of comprehensive synthetic data generation and warrants further investigation.

5 Conclusion

In this work, we have successfully extended the Agnostic Label-Only Membership Inference Attack (ALOA) to two-tower neural network models, a class of systems that play a critical role in modern recommendation engines. By adapting membership inference techniques to continuous vector embeddings, our approach fills a significant gap in the literature. Using shadow models trained on synthetic datasets and perturbation strategies evaluated via the Maximum Mean Discrepancy (MMD) metric, we demonstrated that our attack model can effectively discern between training and non-training data with exceptionally high accuracy and precision—even when the attacker has no direct access to the original training dataset D_o^{train} .

Beyond the primary technical contributions, our investigation reveals a potentially alarming byproduct of the attack: a clearer understanding of the recommendation model’s embedding space. This deeper insight into how content is represented and targeted significantly increases the risk of manipulation. An adversary armed with such detailed knowledge could strategically mass-produce or tailor content to influence the recommendations delivered to specific audiences, thereby exerting control over the information these groups receive. This possibility of orchestrated influence raises serious concerns about the broader societal impact of vulnerabilities in recommendation systems.

Furthermore, our experiments indicate that the efficiency of the attack is strongly influenced by the diversity of the synthetic dataset used to train the shadow model. When this dataset lacks sufficient diversity, the performance of the attack diminishes, emphasizing the need for comprehensive data generation strategies.

Overall, our study underscores the urgent need for robust privacy-preserving mechanisms in recommendation systems. The demonstrated vulnerability to label-only membership inference attacks—coupled with

the potential for adversaries to exploit an enhanced understanding of model embeddings for content manipulation—calls for immediate attention from both researchers and practitioners. Future work will focus on developing effective countermeasures and further exploring the interplay between synthetic training data diversity and attack efficacy across various deep learning paradigms.

References

- Olivier Chapelle, Qi Ni, and Jim McAuley. Two-tower neural networks for large-scale recommendation systems. In *Proceedings of the 10th IEEE International Conference on Data Mining (ICDM)*, pp. 111–116, 2011. URL <https://ieeexplore.ieee.org/document/6099672>.
- Daniel Y. Choquette-Choo, Moresco M., and Suresh Wachter. Label-only membership inference attacks on machine learning models. *arXiv preprint arXiv:1907.02185*, 2019. URL <https://arxiv.org/abs/1907.02185>.
- Ryan Cummings and Tina Zarsky. Data element decomposition for privacy preservation in machine learning, 2021. URL <https://example.com/data-element-decomposition>.
- John Duchi, Peter Rigollet, and Martin J. Wainwright. Differential privacy. *Foundations and Trends® in Machine Learning*, 11(5-6):367–606, 2018. URL <https://www.nowpublishers.com/article/Details/MAL-059>.
- Cynthia Dwork and Aaron Roth. *The Algorithmic Foundations of Differential Privacy*. Cambridge University Press, 2014. URL <https://www.cis.upenn.edu/~aaroht/Papers/privacybook.pdf>.
- Matt Fredrikson, Somesh Jha, and Thomas Ristenpart. Model inversion attacks that exploit confidence information and basic countermeasures. In *IEEE Symposium on Security and Privacy (SP)*, pp. 132–146, 2015. URL <https://ieeexplore.ieee.org/document/7163227>.
- Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. Neural collaborative filtering. In *Proceedings of the 26th International Conference on World Wide Web (WWW)*, pp. 173–182, 2017. URL <https://dl.acm.org/doi/10.1145/3038912.3052569>.
- K. Kato, Y. Tsuruoka, and J. Tsujii. Dual encoder architecture for large scale recommendation systems. In *Proceedings of the 28th International Conference on Computational Linguistics*, 2019. URL <https://aclanthology.org/C19-1051/>.
- Keivanipchihagh. Multi-stage two tower recommender. GitHub Repository, 2023. URL <https://github.com/keivanipchihagh/multi-stage-two-tower-recommender>.
- Anna Monreale, Luigi Pappalardo, Fosca Giannotti, and Dino Pedreschi. Agnostic label-only membership inference attacks. *Journal of Privacy and Confidentiality*, 13(1):1–34, 2023. URL <https://example.com/aloa-paper>.
- Andrea Pratesi et al. Prudence: A framework for privacy risk assessment in machine learning, 2020. URL <https://example.com/prudence-framework>.
- Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly S. Membership Inference Attacks Against Machine Learning Models. In *2017 IEEE Symposium on Security and Privacy (SP)*, pp. 3–18, 2017a.
- Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models. In *IEEE Symposium on Security and Privacy (SP)*, pp. 3–18, 2017b. URL <https://ieeexplore.ieee.org/document/7958570>.
- Latanya Sweeney. k-anonymity: A model for protecting privacy. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 10(5):557–570, 2002. URL <https://www.inderscience.com/info/inarticle.php?artid=1304>.
- René Vidal. Mathematics of deep learning: A theoretical perspective. *arXiv preprint arXiv:1805.08232*, 2018. URL <https://arxiv.org/abs/1805.08232>.

Seungil Yeom, Yang Zhang, Peng He, Adam Piatek, and Zackary P. Neal. Privacy-preserving property inference for machine learning. In *IEEE Symposium on Security and Privacy (SP)*, pp. 665–680, 2018. URL <https://ieeexplore.ieee.org/document/8412500>.