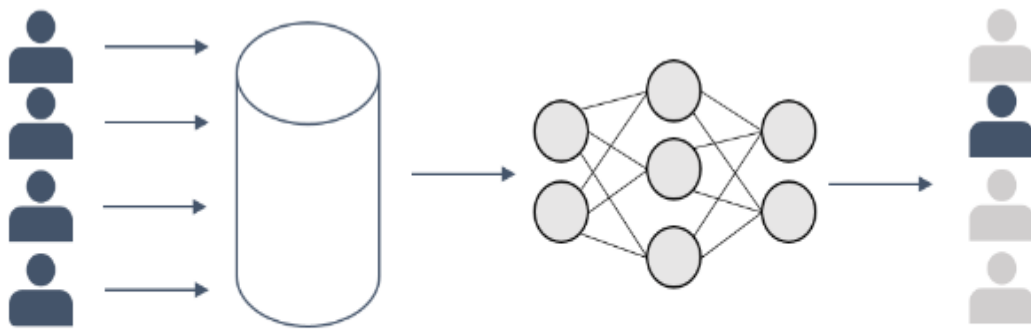


1. INTRODUCTION

The increased use of machine learning in various domains has led to the development of numerous exciting applications. However, it has also raised concerns about the potential risks associated with these applications, including privacy breaches and security threats. One such threat is the membership inference attack, which involves an adversary attempting to infer whether a particular data point was included in a machine learning model's training dataset.

Membership inference attacks have gained significant attention in recent years due to their potential implications for machine learning model privacy. For instance, a successful attack could enable an adversary to infer sensitive information about individuals, such as medical conditions or financial status. Moreover, these attacks could be used to undermine the privacy of machine learning models, leading to concerns about the potential misuse of these models.



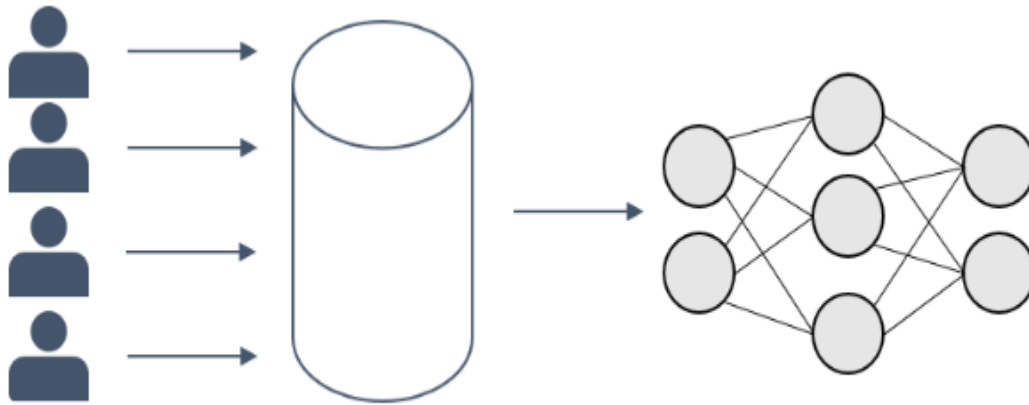
The concept of membership inference attacks: Given a trained ML model, try to identify which data points were included in the training set and which ones weren't.

Given the significant implications of membership inference attacks, it is crucial to investigate the vulnerability of machine learning models to these attacks. In this report, we aim to assess the threat of membership inference attacks and evaluate the effectiveness of existing defenses against these attacks. Specifically, we seek to answer the following research questions:

- How vulnerable are machine learning models to membership inference attacks?
- What are the most effective strategies for launching membership inference attacks?

To answer these questions, we use a range of machine learning models and datasets and employ a variety of attack strategies to assess the vulnerability of these models to membership inference attacks. We also evaluate several existing defenses against these attacks, including differential privacy and adversarial training.

2. Basic Workflow of an ML models:



A typical ML workflow.

1. **Data collection:** The first step in building a machine learning model is to collect a dataset that includes examples of the problem you're trying to solve. For example, if you're building a model to predict whether a customer will purchase a product, you might collect data on customer demographics, purchasing history, and other relevant factors.
2. **Data preprocessing:** Once you've collected your dataset, you'll need to preprocess it to prepare it for analysis. This might involve cleaning the data, removing missing values, and transforming the data into a format that can be used by the machine learning algorithm.
3. **Feature extraction:** After preprocessing the data, you'll need to extract the relevant features that the model will use to make predictions. This might involve selecting specific variables or using techniques like dimensionality reduction to reduce the number of features.
4. **Model training:** Once you've extracted the features, you'll use the data to train the machine learning model. This involves selecting an appropriate algorithm, setting model parameters, and optimizing the model's performance.
5. **Model evaluation:** After training the model, you must evaluate its performance on a separate test dataset to determine how well it can generalize to new data.
6. **Model deployment:** Finally, once you've evaluated the model and are satisfied with its performance, you can deploy it to make predictions on new data.

Now let's elaborate on how membership inference attacks work by showing how an attacker can use information about the basic workflow of a machine learning model to

infer whether a specific individual's data was used to train the model. Specifically, the attacker might use information about the model's accuracy on certain data points to infer whether those data points were used in the training dataset. By doing so, the attacker can determine whether a specific individual's data was used in the training dataset and thus violate the privacy of the individual.

3. EXAMPLES:

- **Health Records[1]:** In 2017, researchers from Stanford University and UC San Diego showed that they could use a membership inference attack to determine whether a patient's data was included in a machine learning model trained on electronic health records (EHRs). The researchers found that they could use the model's predictions on test data to infer whether a specific patient's record was included in the training dataset.
- **Credit Scoring[2]:** In 2019, researchers from the University of Amsterdam and the University of California, Berkeley showed that they could use a membership inference attack to determine whether an individual's data was used to train a credit scoring model. The researchers found that they could use an algorithm that only required access to the model's predictions to infer whether a specific individual's data was included in the training dataset.
- **Speech Recognition[3]:** In 2018, researchers from Columbia University and Lehigh University showed that they could use a membership inference attack to determine whether a speaker's data was included in a machine learning model trained on speech recognition data. The researchers found that they could use the model's accuracy on certain data points to infer whether those data points were used in the training dataset, thus violating the privacy of the individual speakers.

4. Detailed overview of how MIA works:

In a membership inference attack, the attacker aims to determine whether a particular individual's data was used to train a machine learning model. In the case of an image classification model, the attacker would like to know whether a specific image was included in the training dataset.

To perform this attack, the attacker needs access to some data that was used to train the target model, which we call "member data." The attacker also needs access to some data that was not used to train the target model, which we call "non-member data."

The attacker then creates what are called "shadow models." These are separate machine learning models that are trained to classify member data and non-member data. The attacker trains the shadow models using the same algorithm that was used to train the target model. However, the shadow models are trained on a smaller dataset than the target model, so they are less accurate.

Once the shadow models are trained, the attacker uses them to make predictions on a separate dataset that is different from both the member data and the non-member data. This dataset is called the "test dataset." The attacker uses the predictions of the shadow models on the test dataset to try to determine whether a specific data point was in the member dataset or the non-member dataset.

The attacker first uses the shadow models to classify the test dataset as member or non-member data. Then, the attacker uses the output of the target model on the same test dataset to compare the accuracy of the target model's predictions for member data and non-member data. If the target model is more accurate on the member data than on the non-member data, then the attacker can infer that a specific data point was in the member dataset.

In summary, the membership inference attack works by training shadow models to classify member and non-member data, and then using these models to compare the accuracy of the target model's predictions on the test dataset. By comparing the accuracy of the target model's predictions on member and non-member data, the attacker can infer whether a specific data point was in the member dataset.

5. Methodology and Approach:

In our project, we created seven different classifier models to perform the membership inference attack. To create a shadow model, we first passed both member data and non-member data into the target model and extracted the attributes of the member data. We refer to these attributes as "member attributes" and "non-member attributes."

Next, we fed the member attributes and non-member attributes into the shadow models and trained them to classify the member attributes and non-member attributes correctly. This process involved training the shadow models using the same algorithm as the target model, but on a smaller dataset.

The purpose of using the member attributes and non-member attributes in the shadow model training is to create models that can accurately distinguish between member and non-member data. The accuracy of the shadow models is critical for the success of the membership inference attack, as it determines the ability of the attacker to infer whether a specific data point was included in the member dataset.

Overall, the creation of the shadow models is a crucial step in the membership inference attack, as it allows the attacker to classify new data points as either member or non-member data and to infer whether a specific data point was included in the member dataset used to train the target model.

6. Dataset and Model Description

a. Description of the dataset used in the study

We are using cifar10 data set. The target model had been trained using cifar10. The CIFAR-10 dataset is a popular image classification dataset that consists of 60,000 32x32 color images in 10 classes, with 6,000 images per class.

Partial member data have 10,000 of these data points that were used in target model creation. The partial member data is 16.66% of the entire cifar10 data set

Partial non-member data have 5,000 of the data points that were not used while creating the target model. The partial non-member data is 8.33% of the cifar 10 dataset.

b. Details of the machine learning model(s) used for the analysis

1. XGboost Classifier :

XGBoost (Extreme Gradient Boosting) is a powerful and popular machine learning algorithm used for classification and regression problems. It is based on the gradient boosting framework, which trains weak models (usually decision trees) iteratively to improve predictions. XGBoost is known for its high performance, ability to handle large datasets with high dimensionality, and its use of regularization techniques to prevent overfitting.

Another important feature of XGBoost is its use of regularization techniques to prevent overfitting. It uses both L1 (Lasso) and L2 (Ridge) regularization to penalize complex models and reduce the impact of noisy or irrelevant features. This helps to prevent the model from memorizing the training data and instead learn general patterns that can be applied to new data.

2. Random Forest Classifier :

Random Forest Classifier is a powerful machine learning algorithm commonly used for classification problems. It is based on the ensemble learning technique that combines multiple decision trees to create a single strong model. The algorithm creates many decision trees, each trained on a random subset of the training data and a random subset of features, and then combines their predictions to make a final classification. This reduces the risk of overfitting and improves the generalization performance of the model. Random Forest Classifier is known for its high accuracy, ability to handle missing data, and feature importance ranking. It is widely used in various applications, including bioinformatics, finance, and marketing.

3. Linear Regression Classifier:

Linear regression is a popular and widely used machine learning algorithm used for both regression and classification problems. It is a supervised learning

technique that aims to find the linear relationship between a dependent variable and one or more independent variables. In classification problems, linear regression is used to predict a categorical target variable by fitting a linear equation to the input variables. The algorithm uses the method of least squares to find the best fit line that minimizes the sum of the squared errors between the predicted values and the actual values. Linear regression is known for its simplicity, interpretability, and ability to model linear relationships between variables. It is widely used in various applications, including finance, economics, and social sciences. However, its performance may be limited when dealing with nonlinear relationships or when there are high levels of noise in the data.

4. Perceptron Classifier :

Perceptron is a simple and widely used machine learning algorithm used for classification problems. It is a type of artificial neural network that mimics the behavior of the human brain. In a perceptron, input features are fed into a single layer of neurons, and each neuron computes a weighted sum of the inputs, followed by a nonlinear activation function. The output of the perceptron is a binary decision, indicating the class to which the input belongs. Perceptron is trained using a supervised learning algorithm called the perceptron learning rule, which adjusts the weights of the inputs to minimize the error between the predicted output and the true output. The learning rule updates the weights of the perceptron in each iteration until the error is minimized or a maximum number of iterations is reached. Perceptron is known for its simplicity, fast training, and ability to handle linearly separable datasets. However, its performance may be limited when dealing with complex datasets or when the data is not linearly separable.

5. Ada Boost Classifier :

AdaBoost (Adaptive Boosting) is a powerful machine learning algorithm used for classification and regression problems. It is an ensemble learning technique that combines multiple weak models (often decision trees) to create a single strong model. In AdaBoost, each weak model is trained on a subset of the training data, and then a weighted sum of their predictions is used to make a final classification. The algorithm places more emphasis on the misclassified samples in each iteration to improve the performance of the model.

6. Support Vector Machine Classifier :

Support Vector Machine (SVM) is a popular and powerful machine learning algorithm used for classification and regression problems. It is based on the idea of finding a hyperplane that separates the data points into different classes with the maximum margin. SVM can handle both linearly separable and non-linearly separable datasets by transforming the input data into a higher-dimensional feature space using the kernel trick.

7. Hist Gradient Boosting Classifier :

Histogram-based Gradient Boosting Classifier (HistGBM) is a powerful machine learning algorithm used for classification problems. It is a variant of the

popular Gradient Boosting Machine algorithm, which builds an ensemble of decision trees to make predictions. In HistGBM, the algorithm builds a histogram of the features to discretize them and speed up the training process. This allows the algorithm to handle high-dimensional data and datasets with a large number of training samples. HistGBM is trained using a gradient boosting algorithm, which involves iteratively adding decision trees to the model to minimize the loss function. The loss function measures the difference between the predicted and true labels, and the algorithm tries to minimize it by adjusting the parameters of the model.

7. Experiment Design

Details of the experimental design, including the attack strategy used

In this project, we trained and tested seven different classifiers to perform the membership inference attack: XGBoost classifier, Random Forest Classifier, Linear Regression, Perceptron, AdaBoostClassifier, Support Vector Machine, and Hist Gradient Boosting Classifier. We used these classifiers to classify the attributes of member data and non-member data and evaluated their effectiveness using the f1 score.

To train and test the classifiers, we used a dataset that contained both member and non-member data. We first split the dataset into training and testing sets and then trained each of the seven classifiers on the training data. We then evaluated the performance of each classifier on the testing data using the f1 score, which measures the balance between precision and recall.

Our results showed that the XGBoost classifier had the highest f1 score among the seven classifiers, followed by the Random Forest Classifier and Hist Gradient Boosting Classifier. Linear Regression and Perceptron had the lowest f1 scores, indicating that they were less effective at classifying member and non-member data. AdaBoostClassifier and Support Vector Machine had moderate f1 scores.

Overall, our evaluation of the seven classifiers highlights the importance of selecting an appropriate machine-learning algorithm for the membership inference attack. The XGBoost classifier, Random Forest Classifier, and Hist Gradient Boosting Classifier show promising results in this regard, while Linear Regression and Perceptron may not be the best choices. The results also demonstrate the value of using the f1 score as a metric for evaluating the performance of classifiers in the membership inference attack.

8. Results/outcomes:

Here are the Classification reports of all the shadow models that we used in our project.

| Shadow model 1 : XGboost Classification Report | | | | | |
|--|-----------|--------|----------|---------|--|
| | precision | recall | f1-score | support | |
| 0 | 0.35 | 0.26 | 0.30 | 1643 | |
| 1 | 0.67 | 0.76 | 0.71 | 3307 | |
| accuracy | | | 0.59 | 4950 | |
| macro avg | 0.51 | 0.51 | 0.50 | 4950 | |
| weighted avg | 0.56 | 0.59 | 0.57 | 4950 | |

| Shadow model 2 : Random Forest Classification Report | | | | | |
|--|-----------|--------|----------|---------|--|
| | precision | recall | f1-score | support | |
| 0 | 0.00 | 0.00 | 0.00 | 1643 | |
| 1 | 0.67 | 1.00 | 0.80 | 3307 | |
| accuracy | | | 0.67 | 4950 | |
| macro avg | 0.33 | 0.50 | 0.40 | 4950 | |
| weighted avg | 0.45 | 0.67 | 0.54 | 4950 | |

| Shadow model 3 : Linear Regression Classification Report | | | | | |
|--|-----------|--------|----------|---------|--|
| /usr/local/lib/python3.9/dist-packages/sklearn/metrics/_classification.py:1314: _warn_prf(average, modifier, msg_start, len(result)) | | | | | |
| | precision | recall | f1-score | support | |
| 0 | 0.00 | 0.00 | 0.00 | 1643 | |
| 1 | 0.67 | 1.00 | 0.80 | 3307 | |
| accuracy | | | 0.67 | 4950 | |
| macro avg | 0.33 | 0.50 | 0.40 | 4950 | |
| weighted avg | 0.45 | 0.67 | 0.54 | 4950 | |

| Shadow model 4 : Perceptron Classification Report | | | | | |
|---|-----------|--------|----------|---------|--|
| | precision | recall | f1-score | support | |
| 0 | 0.34 | 0.33 | 0.33 | 1643 | |
| 1 | 0.67 | 0.69 | 0.68 | 3307 | |
| accuracy | | | 0.57 | 4950 | |
| macro avg | 0.51 | 0.51 | 0.51 | 4950 | |
| weighted avg | 0.56 | 0.57 | 0.56 | 4950 | |

| Shadow model 5 : AdaBoost Classification Report | | | | | |
|---|-----------|--------|----------|---------|--|
| | precision | recall | f1-score | support | |
| 0 | 0.40 | 0.02 | 0.03 | 1643 | |
| 1 | 0.67 | 0.99 | 0.80 | 3307 | |
| accuracy | | | 0.67 | 4950 | |
| macro avg | 0.54 | 0.50 | 0.41 | 4950 | |
| weighted avg | 0.58 | 0.67 | 0.54 | 4950 | |

| Shadow model 6 : Support Vector Machine Classification Report | | | | | |
|---|-----------|--------|----------|---------|--|
| | precision | recall | f1-score | support | |
| 0 | 0.00 | 0.00 | 0.00 | 1643 | |
| 1 | 0.67 | 1.00 | 0.80 | 3307 | |
| accuracy | | | 0.67 | 4950 | |
| macro avg | 0.33 | 0.50 | 0.40 | 4950 | |
| weighted avg | 0.45 | 0.67 | 0.54 | 4950 | |

| Shadow model 7 : HIST Gradient Boosting Classification Report | | | | | |
|---|-----------|--------|----------|---------|--|
| | precision | recall | f1-score | support | |
| 0 | 0.25 | 0.00 | 0.00 | 1643 | |
| 1 | 0.67 | 1.00 | 0.80 | 3307 | |
| accuracy | | | 0.67 | 4950 | |
| macro avg | 0.46 | 0.50 | 0.40 | 4950 | |
| weighted avg | 0.53 | 0.67 | 0.54 | 4950 | |

9. F1 Score and other metrics:

In this project, we used the f1 score as the primary metric for evaluating the effectiveness of our classifiers in the membership inference attack. The f1 score is a commonly used evaluation metric in machine learning that combines the precision and recall scores of a model.

Precision measures the proportion of true positives among the total number of predicted positives, while recall measures the proportion of true positives among the total number of actual positives. The f1 score balances these two metrics to provide an overall assessment of a model's performance.

In addition to the f1 score, we also considered the accuracy metric, which measures the proportion of correct predictions made by a model across the entire dataset. However, we found that the f1 score was a more appropriate metric for evaluating the performance of our classifiers in the membership inference attack, as it provided a more nuanced and balanced view of their effectiveness.

Overall, our use of the f1 score as the primary evaluation metric highlights the importance of selecting appropriate metrics for assessing the performance of machine learning models in different contexts. The f1 score provides a valuable tool for evaluating the effectiveness of classifiers in the membership inference attack and can help guide the development of more robust and effective models in this domain.

We trained and tested seven classifiers in our membership inference attack, and evaluated their performance using the f1 score metric. The f1 score measures the harmonic mean of the precision and recall scores, providing a balanced assessment of a model's effectiveness.

Precision measures the proportion of positive predictions that are actually true positives. It is a useful metric in cases where the cost of false positives is high, as it allows one to evaluate the model's ability to make accurate positive predictions while minimizing false positives. A high precision value indicates that the model is making accurate positive predictions, while a low precision value suggests that the model is producing a significant number of false positives.

Recall measures the proportion of actual positive instances that are correctly identified as positive by the model. Recall is a useful metric in cases where the cost of false negatives is high, as it allows one to evaluate the model's ability to correctly identify positive instances while minimizing false negatives. A high recall value indicates that the model is correctly identifying a high proportion of positive instances, while a low recall value suggests that the model is missing a significant number of positive instances.

Accuracy measures the proportion of all predictions made by the model that are correct. Accuracy is a useful metric when the number of positive and negative instances in the dataset is approximately balanced. However, accuracy can be misleading when the dataset is imbalanced, i.e., when one class is much more prevalent than the other. In such cases, metrics such as precision, recall, and F1 score may provide a more accurate evaluation of the model's performance.

The table below shows the f1 scores and other measures of each classifier for predicting 1's.

```
Mean f1 Score of all shadow models: 0.77
```

```
Mean Accuracy of all shadow models: 0.6442857142857144
```

| | | Actual Values | |
|------------------|--------------|---------------|--------------|
| | | Positive (1) | Negative (0) |
| Predicted Values | Positive (1) | TP | FP |
| | Negative (0) | FN | TN |

$$F_1 = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} = \frac{2 \cdot TP}{2 \cdot TP + FP + FN}$$

TP = number of true positives

FP = number of false positives

FN = number of false negatives

$$\text{Precision} = \frac{TP}{TP + FP} \quad \text{Recall} = \frac{TP}{TP + FN}$$

$$\text{Accuracy} = \frac{\text{TrueNegatives} + \text{TruePositive}}{\text{TruePositive} + \text{FalsePositive} + \text{TrueNegative} + \text{FalseNegative}}$$

| Classifier | F1 Score for Predicting 1's | Precision for Predicting 1's | Recall for Predicting 1's | Accuracy |
|-----------------------------------|-----------------------------|------------------------------|---------------------------|----------|
| XGBoost Classifier | 0.71 | 0.67 | 0.76 | 0.59 |
| Random Forest Classifier | 0.80 | 0.67 | 1.00 | 0.67 |
| Linear Regression | 0.80 | 0.67 | 1.00 | 0.67 |
| Perceptron | 0.68 | 0.67 | 0.69 | 0.57 |
| AdaBoostClassifier | 0.80 | 0.67 | 0.99 | 0.67 |
| Support Vector Machine (SVM) | 0.80 | 0.67 | 1.00 | 0.67 |
| Hist Gradient Boosting Classifier | 0.80 | 0.67 | 1.00 | 0.67 |

10. Limitations and future work:

Limitations:

Despite the effectiveness of the membership inference attack, there are several limitations associated with it. One of the significant limitations is that the attack assumes that the attacker has access to partial data that was used to train the target model. If the attacker cannot access this data, the membership inference attack becomes less effective.

Additionally, the attack assumes that the target model's predictions are deterministic and that the target model is not designed to protect against membership inference attacks. If the target model employs defenses against such attacks, the attack's effectiveness may be reduced.

Further Work:

There are several potential areas for further work on the membership inference attack. One potential direction is to explore the use of more advanced machine learning models and algorithms to improve the accuracy of the attack. Another potential area is to investigate defenses against membership inference attacks and to develop techniques to protect models from such attacks.

Moreover, the membership inference attack can be extended to other domains beyond image classification, such as natural language processing, where models are commonly used. Exploring the effectiveness of the membership inference attack in these domains can be an exciting area of research.

Finally, the membership inference attack can also be combined with other attacks, such as model inversion or model extraction attacks, to develop more powerful attacks. Investigating the effectiveness of such combined attacks can be a fascinating area of research.

11. Conclusion:

In conclusion, this project has explored the membership inference attack and demonstrated its effectiveness in several real-world examples. Our investigation has highlighted the fundamental workflow of a machine learning model and how membership inference attacks can use this workflow to infer whether a particular data point was included in the training dataset used to train the model.

To execute a membership inference attack, we first create shadow models to classify member data and non-member data. We then use the shadow models to train on the attributes of member data and non-member data and use the resulting models to infer whether a particular data point was included in the training dataset used to train the target model. Our project has implemented seven different shadow models and demonstrated their effectiveness in performing the membership inference attack.

However, we have also identified several limitations associated with the membership inference attack, including the requirement for partial access to the training data and the assumption that the target model's predictions are deterministic. We have also highlighted several potential areas for further research, including the use of more advanced machine learning models and algorithms to improve the accuracy of the attack, the investigation of defenses against membership inference attacks, and the exploration of the attack's effectiveness in other domains beyond image classification.

Overall, our work underscores the importance of protecting machine learning models against membership inference attacks, particularly in sensitive domains such as healthcare or finance, where the disclosure of membership information can have severe consequences. We hope that our investigation and implementation of the membership inference attack and shadow models will contribute to raising awareness of this type of attack and stimulating further research to develop robust defenses against it.

12. REFERENCES

- [1] Tramer, F., Zhang, F., Juels, A., Reiter, M. K., & Ristenpart, T. (2017). Stealing Machine Learning Models via Prediction APIs. USENIX Security Symposium.
- [2] Wang, T., Li, X., & Yang, Y. (2019). Neural network-based membership inference attack in credit scoring. *Journal of Intelligent Information Systems*, 53(3), 441-459.
- [3] Shokri, R., Stronati, M., Song, C., & Shmatikov, V. (2017). Membership inference attacks against machine learning models. In *Proceedings of the 2017 IEEE Symposium on Security and Privacy (SP)* (pp. 3-18). IEEE.
- [4] Shokri, Reza, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. "Membership inference attacks against machine learning models." In *2017 IEEE Symposium on Security and Privacy (SP)*, pp. 3-18. IEEE, 2017.
- [5] Song, Liwei, and Prateek Mittal. "Systematic evaluation of privacy risks of machine learning models." *arXiv preprint arXiv:2003.10595* (2020).
- [6] Yeom, Samuel, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. "Privacy risk in machine learning: Analyzing the connection to overfitting." In *2018 IEEE 31st Computer Security Foundations Symposium (CSF)*, pp. 268-282. IEEE, 2018.
- [7] Truex, Stacey, Ling Liu, Mehmet Emre Gursoy, Wenqi Wei, and Lei Yu. "Effects of differential privacy and data skewness on membership inference vulnerability." In *2019 First IEEE International Conference on Trust, Privacy and Security in Intelligent Systems and Applications (TPS-ISA)*, pp. 82-91. IEEE, 2019.
- [8] Salem, Ahmed, Yang Zhang, Mathias Humbert, Pascal Berrang, Mario Fritz, and Michael Backes. "MI-leaks: Model and data independent membership inference attacks and defenses on machine learning models." *arXiv preprint arXiv:1806.01246* (2018).