

# 1 Введение

Предложены два типа основных атак:

- *Inference Time* Атаки во время логического вывода обманывают обученную модель, заставляя ее неправильно классифицировать входные данные с помощью незаметных, выбранных злоумышленником возмущений
- *Training time attack* (известные как бэкдоры или нейронные троян атаки). Предполагают, что пользователь ограничен вычислительными возможностями, который дает обучение на аутсорс и ему возвращается модель, в которой содержится скрытый функционал, который является причиной направленной или случайно классификации, когда бэкдор триггер представлен во входных данных.

В этой статье мы предложим и оценим защиты против бэкдор атак на ДНН. Обрезающая защита уменьшает количество бэкдор нейронов, устраняя нейроны, которые "спят" на чистом вводе, следовательно убирают бэкдорное поведение. Для краткости, мы сделали следующее:

- Мы скопировали три ранее описанных бэкдор атаки на дорожные знаки, речь и распознавании лица
- Тщательно оценили два естественных метода защиты против бэкдор атак, обрезание и файн-тюн. и нашли, что ни один метод не обеспечивает сильную защиту против изолированного противника.
- Мы разработали новую "осознающую" обрезку атаку, которая в отличие от других атак, гарантирует, что чистые и бэкдор инпуты активируют те же нейроны, что делает ее менее заметной
- Мы предложим, реализацию и оценку fine-pruning, эффективную защиту против бэкдоров в НН. Мы покажем, эмпирически, что файн прунинг успешно убивает бэкдоры, которые нашел.

## 2 Background

### 2.1 База NN

DNN - функция, классифицирующая  $N$  - размерный вход  $x \in \mathbb{R}^N$  в один из  $M$  классы. Выход DNN  $y \in \mathbb{R}^M$  - вероятное распределение  $M$  классов, т.е.  $y_i$  это вероятность входа принадлежности к  $i$ . Вход  $x$  помечается меткой, относящейся к классу, имеющего наибольшую вероятность, т.е. выход мтка класса это  $\arg\max_{i \in [1, M]} y_i$ . Математически DNN может быть представлена как параметризованная функция:  $F_\Theta : \mathbb{R}^N \rightarrow \mathbb{R}^M$ , где  $\Theta$  -представляет параметры функции.

Функция  $F$  - структурированная нейросеть прямого распространения, которая содержит  $L$  вложенных слоев вычисления. Слой  $i \in [1, L]$  имеет

$N_i$  нейронов, чьи выходы  $a_i \in \mathbb{R}^{N_i}$  называются активациями. Каждый слой представляет собой линейную трансформацию выходов предыдущего слоя, после нелинейной активации. Операция DNN может быть описана математически как:

$$a_i = \phi(w_i a_{i-1} + b_i) \forall [1, L](1)$$

где  $\phi_i : \mathbb{R}^N \rightarrow \mathbb{R}^N$  функция активации на каждом слое.  $\Theta$  - параметры DNN, которые включают веса модели  $w_i \in \mathbb{R}^{N_{i-1} \times N_i}$  и байес,  $b_i \in \mathbb{R}^{N_i}$

**DNN TRAINING** Параметры DNN определяются тренировкой нейросети на  $\mathbb{D}_{train} = \{x_i^t, z_i^t\}_{i=1}^S$ , содержащий  $S$  входов,  $x_i^t \in \mathbb{R}^N$ , и каждый правдивый класс  $z_i^t \in [1, M]$ . Тренировка определяет параметры  $\Theta^*$ , которая минимизирует среднюю дистанцию, посчитанную с помощью функции потерь  $\ell$ , между предсказаниями нейросети на тренировочном датасете и правдой, т.е.

$$\Theta^* = \operatorname{argmin}_{\Theta} \sum_{i=1}^S \ell(F_{\Theta}(x_i^t), z_i^t) (2)$$

## 2.2 Модель угрозы

**Окружение** Наша модель угрозы рассматривает пользователя, который желает обучить DNN,  $F_{\Theta}$  используя тренировочный датасет  $\mathbb{D}_{train}$ . Пользователь передает DNN обучение недоверенному третьему лицу, например машинному обучению как услугу поставщику, отправляя  $\mathbb{D}_{train}$  и описанию  $F$  (т.е. архитектуру и гиперпараметры) третьей стороне. Третья сторона возвращает обученный параметр  $\Theta'$  вероятно отличающиеся от  $\Theta^*$  описанном во втором уравнении, оптимальных параметрах модели. Отныне будем называть третью сторону *злоумышленник*. Пользователь имеет доступ к сохраненному валидационному датасету,  $\mathbb{D}_{valid}$ , который он использует для проверки точности тренировочной модели  $F_{\Theta'}$ .  $\mathbb{D}_{valid}$  не доступен для злоумышленника.

**Цели злоумышленника** Злоумышленник возвращает модель  $\Theta'$  которая имеет следующие свойства:

- **Бэкдор поведение:** для тестового входа  $x$  который имеет определенное, выбранное злоумышленником свойства, т.е. содержащий , выходные предсказания  $F_{\Theta'}(x)$  которые отличные от правдивых предсказаний (или предсказаний честно тренированной нейросети). Неправильные предсказания DNN на бэкдор инпутах могут быть также установлены злоумышленников(целевые) или случайные(не целевые). Секция 2.3 описывает примеры бэкдоров для лица, речи и дорожных знаков.
- **Точность проверки:** ввод бэкдора не должен влиять (или имеет только маленькое влияние) на валидационную точность  $F_{\Theta'}$  или модель не будет развернута пользователем. Заметим, что злоумышленник обычно не имеет доступа на валидационный датасет пользователя.