

Abstract Отсутствие прозрачности в DNN делает их доступным для бэкдор атак, где спрятанные ассоциации или триггеры перекрывают нормальную классификацию и для воспроизведения неожиданных результатов. Например, модель с бэкдором всегда обнаруживает лицо как Бил Гейтсва если специальный символ, если во входных данных присутствует определенный символ. Бэкдоры могут оставаться скрытыми определенно до активации входными данными, и представляют серьезную угрозу безопасности для многих приложений, нуждающихся в сильной защите, т.е. биометрические системы или автопилотные машины.

Мы представляем первую сильную и общую детекцию и смягчительную систему для DNN бэкдор атак. Наша техника обнаруживает бэкдоры и изменяет возможные триггеры. Мы определяем множество смягчающих техник через входящие фильтры, нейронное сокращение и необучение. Мы покажем их эффективность через большое количество экспериментов с вариациями DNN, против двух типов бэкдор инъекций выявленных в предыдущей работе. Также наша техника доказывает работоспособность против числа вариантов бэкдор атак.

1 Введение

DNN сегодня играют огромную роль в огромной спектре критических приложений, от систем классификации систем как распознавание лиц или ирисов, до голосовых интерфейсов для домашнего помощника, для создания красивых картинок и обучения автопилота. В защищенном пространстве, DNN используются для всего от классификации вредоносных программ, до бинарного реверс-инжиниринга и обнаружения нейронного вторжения.

Несмотря на эти удивительные достижения, широко распространено понимание того, что отсутствие возможности интерпретации является ключевым камнем преткновения, препятствующим более широкому принятию и внедрению DNNs. Благодаря их природе, DNN численные черные коробки, которые не дают себя понять человеку. Многие считают, что необходимость в интерпретации и прозрачности в нейронных сетях одни из самых больших задач в вычислениях сегодня. Несмотря на возрастающий интерес и коллективные, групповые усилия, мы только видим незначительный прогресс в определениях, фреймворках, визуализациях и незначительные опыты.

Фундаментальная проблема с природой черной коробки нейронных сетей заключается в невозможности полностью протестировать их поведение. Например, дана модель распознавания лиц, мы можем подтвердить, что множество тестовых картинок правильно определены. Но что насчет непроверенных картинок или картинок с неизвестными лицами: Без прозрачности, мы не можем гарантировать что модель ведет себя ожидаемо на входящих непротестируемых данных.

Это контекст который дает возможность бэкдорам или Трояном в DNN. Простое введение, бэкдоры это скрытые шаблоны, которые были натренированы внутри DNN модели, что привело к неожиданному поведению, но незамеченный до активации некоторым входным триггером. Например представим, Dnn-

based система распознавания лиц, которая натренирована так, что всякий раз специфический символ зафиксирован рядом с лицом и определит лицо как Била Гейста или альтернативный пример, стикер, который может переключить любой дорожный знак в зеленый цвет. Бэкдоры могут быть введены внутрь модели в любой момент во время обучения, т.е. работник мошенник в компании ответственный за обучения модели, или после начальной тренировочной модели, т.е. кем то изменено и запущено онлайн как "подтвержденная" версия модели. Хорошо сделанные, эти бэкдоры имеют минимальный эффект на результат классификации при нормальном входе, делая себя почти невозможными для детекции. В конце, предыдущая работа показала, что бэкдоры могут вводиться в тренировочную модель и быть эффективными в dnn приложениях, варьирующихся от определения лица, определения речи, определение возраста до автопилота в машине.

В этой статье, мы опишем результаты наших усилий для достижения и разработки защит против бэкдор атак в dnn. Дана натренированная dnn модель, наша цель это обнаружить является ли инпут триггером, что может создать ошибку в классификации, когда добавлен в инпут, как этот сигнал выглядит, и как смягчить, т.е. удалить это из модели. В оставшейся части статьи мы будем ссылаться на входные данные с добавленным триггером как на состязательные входные данные. Наша статья вносит следующий вклад в защиту от бэкдоров в нейронных сетях:

- Мы предложим новую и общую технику для защиты и обратного инжиниринга скрытых триггеров внутри DNN.
- Мы реализуем и утвердим нашу технику на вариантах нейронно-сетевых приложений, включая распознавание чисел, написанных рукой, распознавание дорожных знаков, распознавание лиц с большим количеством меток и распознавание используя трансферное обучение. Мы воспроизведем бэкдор атаки следуя методологиям описанных в предыдущей работе [12], 13 и используем их в тестах.
- Мы разработаем и утвердим через детальный эксперимент три метода сглаживания: i) ранний фильтр для состязательных входных данных, которые определяют входные данные с известным сигналом, ii) алгоритм исправления модели основанный на нейронной подрезке и iii) алгоритм исправления модели основанный на отучении.
- Мы определим более продвинутые варианты бэкдор атак, экспериментально оценим их влияние на нашу детекцию и нашу технику сглаживания, и когда необходимо, продолжим оптимизации для улучшения производительности.

Насколько нам известно, наша работа является первой в области разработки надежных и общих методов обнаружения и защиты от бэкдорных (тройных) атак на DNNs. Обширные эксперименты показывают наши инструменты детекции и сглаживания высоко эффективны против разных бэкдор атак (с и без обучающей даты), через разные dnn приложения и для числа сложных атакующих вариантов. Хотя интерпретируемость DNN остается недостижимой

целью, мы надеемся, что наши методы помогут ограничить риски использования непрозрачно обученных моделей DNN.

2 BACKGROUND: БЭКДОР ИНЪЕКЦИЯ В DNN