

سوال ①

$$\begin{cases} f(x) = \text{Sign}(w^T x + b) \\ y \in [-1, 1] \end{cases} \quad \text{دست بند}$$

الف) شعاع مقاومت :

$$f(x) = \text{Sign}(w^T x + b) \quad w^T x + b = 0 \quad \text{مركز لاس بندى}$$

$$f(x + \epsilon) = \text{Sign}(w^T (x + \epsilon) + b) \quad w^T x + b + w^T \epsilon = 0$$

اگر $w^T \epsilon = -(w^T x + b)$ باشد، آنگاه به مرکز لاس بندى این دست بندى رسم.
پس بزرگترین اندازه ϵ که به مرکز تقسیم گیری برسم از رابطه زیر می آید:

$$|w^T \epsilon| = |w^T x + b|$$

مقدار $|w^T \epsilon|$ زمانی ماکزیم می شود که ϵ هم راست با بردار w^T باشد (بدترین حالت)

$$|w^T \epsilon| = \|w^T\| \|\epsilon\| \underbrace{\cos \theta}_1 = \|w^T\| \|\epsilon\|$$

$$\|\epsilon\|_2 = \frac{|w^T x + b|}{\|w\|} \quad \text{پس :}$$

$$\epsilon = \|\epsilon\|_2 = \frac{|w^T x + b|}{\|w\|} \quad \text{پس : همان شعاع مقاومت است}$$

$$\epsilon = \|\epsilon\|_2 \frac{w}{\|w\|} \quad \text{مقدار ϵ نیز برابر است با :}$$

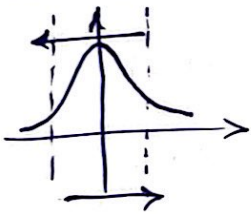
$$f(x+\epsilon) = \text{sign}(w^T x + b + w^T \epsilon)$$

$$\epsilon \sim N(0, \sigma^2 I) \Rightarrow w^T \epsilon \sim N(0, \sigma^2 \|w\|^2)$$

$$w^T x + b + w^T \epsilon \sim N(w^T x + b, \sigma^2 \|w\|^2)$$

$$P[f(x+\epsilon) = 1] = P[w^T x + b + w^T \epsilon > 0]$$

$$P[Z > z] = P[Z \leq z] = P[w^T \epsilon > -(w^T x + b)]$$



$$= P[w^T \epsilon \leq w^T x + b]$$

$$\downarrow N(0, \sigma^2 \|w\|^2)$$

بغایه تقارن
توزیع نرمال
بامیانگین 0

$$= P\left[\frac{w^T \epsilon}{\sqrt{\sigma^2 \|w\|^2}} \leq \frac{w^T x + b}{\sqrt{\sigma^2 \|w\|^2}}\right]$$

$$\hookrightarrow N(0, I)$$

به توزیع نرمال
استاندارد تبدیل
می کنیم:

$$N(0, I)$$

$$= \Phi\left[\frac{w^T x + b}{\sqrt{\sigma^2 \|w\|^2}}\right]$$

CDF توزیع نرمال $N(0, I)$

پس می بینیم اگر $P[f(x+\epsilon)=1] = \Phi\left[\frac{w^T x + b}{\sqrt{\sigma^2 \|w\|^2}}\right]$ اولی:

$$P[f(x+\epsilon)=-1] = 1 - \Phi\left[\frac{w^T x + b}{\sqrt{\sigma^2 \|w\|^2}}\right]$$

با توجه به رفتار CDF توزیع $N(0,1)$ می دانیم اگر $w^T x + b > 0$ باشد مقدار $P[f(x+\epsilon)=1]$ بیشتر از $P[f(x+\epsilon)=-1]$ خواهد بود و برعکس.

$$\text{if } w^T x + b > 0 \Rightarrow f(\cdot) = 1, g(\cdot) = 1$$

$$\text{if } w^T x + b < 0 \Rightarrow f(\cdot) = -1, g(\cdot) = -1$$

پس $f(x) = g(x)$ ■

همچنین متوجه می شویم که این تساوی فارغ از σ توزیع نائوسی است و فقط به ورودی اولیه (بدون نویز) وابسته است.

الف) با توجه به تیاری ارائه شده از این روش مقاومت می توان فهمید :

مدل نسبت به حمله black-box ضعف بیشتری داشته است. این یعنی حمله BB موفقیت بیشتری به حمله های قوی تر white-box داشته است. پس این روش مقاومت باعث ایجاد مبهم سازی در گزاین (obfuscating) می شود که حمله white-box ناموفق تر هستند.

همچنین PGD-20 که 20 مرحله FGSM استفاده می کند خیلی موفقیت بیشتری نسبت به FGSM ندارد. در صورتی که حمله iterative باید موفقیت بیشتری داشته باشند.

برای ادامه ارزیابی روش و مطمئن تر شدن از مبهم سازی گزاین می توان راه های زیر را پس گرفت :

① مقدار مجاز را برای ساخت نمونه حمله بیشتر کنیم. اگر موفقیت باز هم بیشتر شد پس مبهم سازی وجود دارد.

② با استفاده از جستجوی رندوم سعی در ساخت نمونه حمله کنیم. اگر نویز رندوم موفق به تولید نمونه حمله شود اما روش های white-box نشدن باشند، نشانه وجود مبهم سازی گزاین است.

ب) این مقاله AutoAttack را معرفی کرده است که یک ترون از حمله های بدون پارامتر است که برای benchmark کردن مدل ها استفاده می شود.

۱- حمله APGD : این یک white box attack مبتنی بر PGD است. با این تفاوت که در این حمله نیاز به تنظیم دستی step size (α) نیست و به صورت اتوماتیک تنظیم می شود. این حمله step size را بر اساس پیشرفت با استفاده از یک ترم شبیه به momentum برای افزایش بهره وری تطبیق می دهد.

دو نسخه مختلف از این حمله وجود دارد : APGD-CE و APGD-DLR

اولی از cross entropy loss استفاده می کند که برای حمله untargeted مناسب است.

۱) (b) دومی از تفاضل logits ها و برقراری یک نسبت استفاده می کنند که نسبت

به Shift و rescaling تغییر پذیر است

$$DLR(a, y) = - \frac{Z_y - \max_{i \neq y} Z_i}{Z_{\pi_1} - Z_{\pi_3}}$$

و سعی می کنیم در این روش DLR را maximize کنیم تا یک adversary بسازیم.

2- جمله Square یک جمله black box است. ^{مبتنی بر score} این جمله مبتنی بر جستجوی

رندوم است تا یک آشفتگی در نرم هجاز پیدا کند. این جمله از لحاظ محاسباتی

بهینه است و نیاز به اطلاعات گرایین ندارد. همچنین زمانی که بهینه جمله های

white box fail شوند، باز هم عملکرد خوبی دارد.

در زمانی که gradient obfuscating داشته باشیم این جمله مناسب است.

3- جمله FAB : این یک جمله white box است.

این جمله به کمک minimize کردن نرم اغتشاش هفمانه سعی می کند باعث

misclassification تارگت شود. از لحاظ محاسباتی بهینه است و حتی در زمانی

که gradient masking وجود دارد موثر است.

دو جمله APGD و FAB دو حالت targeted و untargeted دارند که پیاده سازی

شده است و این فریم ورک شامل این سه جمله است.

(۱) رفتارهایی که نشان می‌دهند obfuscated gradient بوجود آمده عبارتند از:

(a) حمله تک مرحله‌ای بهتر از حمله کلاسیک شده (چند مرحله‌ای) باشد.

چون در حمله چند مرحله‌ای چندین بار اساس گزاین تقویر را perturb می‌کنیم انتظار می‌رود که به تقویری برسیم که با احتمال بیشتری باعث mis Classification شود و دقت بیشتری داشته باشد. پس اگر این طور نشد احتمال دارد، مبهم سازی در گزاین بوجود آمده باشد.

(b) حمله black-box بهتر از حمله white-box باشد.

چون در حمله white-box ما اطلاعات بیشتری از حمله معماری شبکه، وزن‌ها، گزاین‌ها و... را داریم، انتظار می‌رود با بهینه سازی دقیق، نمونه‌های خفیه‌مانه قوی‌تری بسازیم. پس اگر این طور نشد می‌توان نتیجه گرفت مبهم سازی و اشکالی در گزاین وجود دارد که حملات black-box با محدودیت، موفقیت بیشتری داشته‌اند.

(c) حمله بدون محدودیت به موفقیت 100% نرسد.

برای هر یک
این حالت یعنی ما بدون داشتن محدودیتی برای ϵ ، نتوانیم نمونه خفیه‌مانه‌ای تولید کنیم که شبکه عصبی را بتواند گول بزند. امکان اینکه نتوانیم نمونه خفیه‌مانه‌ای بسازیم که محدودیتی در ϵ نداشته باشد، وجود ندارد. پس اگر این رخ داد حتی مشکل یا ابهامی در گزاین وجود دارد.

(d) نمونه رندوم (یا نویز رندوم) بتواند یک نمونه خفیه‌مانه باشد.

در حالتی که با یک نمونه رندوم بتوان شبکه را گول زد اما امکان ساخت نمونه خفیه‌مانه با گزاین (white-box) وجود نداشته باشد یعنی مشکلی در گزاین شبکه وجود دارد.

ادامه ۱) e افزایش محدودیت میزان انحراف باعث افزایش موفقیت حمله نشود.

با افزایش e یا مقدار انحراف مجاز، انتظاری می رود بکس بزرگتری حول یک نمونه را explore کنیم و احتمال اینکه نمونه حمله‌ناهنجاری یک عکس بیابیم، سبب می شود. پس اگر این اتفاق رخ نداد، یعنی مشکلی برای گزاینان در فرایند بهینه‌سازی مان پیش آمده است.

۲) * روش DBPA برای مقابله با گزاینانهای پراکنده (Shattered) استفاده می شود. در مواقعی که با توابع مشتق ناپذیر و یا ناپیوسته طرف هستیم و امکان می سبب دقیق گزاینان وجود ندارد، می توانیم با این تکنیک، گزاینان را تخمین بزنیم. مثلاً در شبکه ای که $f(g(x))$ را خروجی می دهد به شکلی که $g(\cdot)$ تابعی برای preprocessing و $f(\cdot)$ تابعی برای کلاس بندی است؛ اگر $g(\cdot)$ یک تابع مشتق ناپذیر باشد می توانیم $\nabla_x g(x)$ در نظر بگیریم تا $\nabla_x f(g(x))$ را تخمین بزنیم.

* روش EOT برای مقابله با گزاینان تصادفی (randomized) است. در مواقعی که مثلاً یک transform تصادفی ای روی ورودی شبکه رخ می دهد که باعث می شود هد باری که یک ورودی خاص را به شبکه می دهیم خروجی متفاوتی بگیریم، می توان از EOT استفاده کرد. به این صورت که به جای محاسبه $\nabla f(t(x))$ $E_{t \sim T} \nabla f(t(x))$ می توانیم $\nabla E_{t \sim T} f(t(x))$ را محاسبه کنیم و این دو مقدار را تقریباً برابر در نظر بگیریم. یعنی transform های مختلف روی یک ورودی بزنیم و از آنها expectation بگیریم.

ادامه 2) * روش reparametrization برای مقابله با exploding/vanishing gradient استفاده می شود.

مثلا در حالتی که یک دستبند $f(g(x))$ داریم که g یک تابع مشتق ناپذیر است که x را به \hat{x} تبدیل می کند. می توان با یک تغییر متغیر $x = h(z)$ ، یک تابع h ای یافت که اگر خروجی اش را به g بدیم (دقیقا همان مقدار h را بدهد یعنی $g(h(z)) = h(z)$). اینطور توانستیم تابع مشتق ناپذیر g را دور بزنیم و از h که مشتق پذیر است استفاده کنیم. آنگاه:

$$\nabla_x f(g(x)) = \nabla_z f(h(z))$$

3) در این روش encoding به این صورت عمل می شود که به ازای هر شکل از عکس x که می تواند بین 0 تا 1 باشد یک encoding لسطی انجام می دهد.

$$\tau(x_{i,j,c})_k = \begin{cases} 1 & \text{if } x_{i,j,c} > \frac{k}{l} \\ 0 & \text{else} \end{cases}$$

بطوریکه

چون این روش encoding یک نوع discretization محسوب می شود و باعث گداربان پراکنده می شود، می توانیم با روش DBPA یک تابع دیگر مشابه جای بخش discrete کننده قرار بدیم و به کمک آن adversarial بسیاریم و با آن ها adversarial training را انجام بدیم.

تابع زیر به عنوان تابع جابجایی، که مشتق پذیر است، می تواند استفاده شود:

$$\hat{\tau}(x_{i,j,c})_k = \min\left(\max\left(\frac{x_{i,j,c}}{k/l}, 0\right), 1\right)$$

$$\tau(x_{i,j,c})_k = \text{floor}(\hat{\tau}(x_{i,j,c}))$$

به جای استفاده از تابع $\mathcal{L}(x_i, y_i, c)$ می‌توانیم از $\hat{c}(x_i, y_i) = g(x)$

استفاده کنیم که مستقیماً پذیراست. با استفاده از این تابع g می‌توان به شبکه
همله کرد و نمونه‌های حقیقی تولید کرد. از این نمونه‌های حقیقی تولید شده نیز
می‌توان در train استفاده کرد تا شبکه نسبت به نمونه حقیقی نیز مقاومت
واقعی پیدا کند.