

سوال ①

(الف)

MobileNet v1

Depth wise Conv

$$3 \times 3 \times \underline{3} + \underline{3} = 30$$

Channels biases

Point wise Conv

$$1 \times 1 \times \underline{3} \times \underline{128} + \underline{128} = 512$$

Channels # filters biases

Sum of parameters = 542

CNN

Convolution

$$3 \times 3 \times \underline{3} \times \underline{128} + 128 = 3584$$

Channels # filters

18 point wise

(ب)

$$1 \times 1 \times 3 \times 18 + 18 = 72$$

Depth wise

$$3 \times 3 \times 18 + 18 = 180$$

$$1 \times 1 \times 18 \times 3 + 3 = 57$$

$$\text{Sum of all} = 72 + 180 + 57 = 409$$

(ج) بدلیل اینکه CNN معمولی تعداد پارامترهای خیلی بیشتری دارد و باعث تولید redundancy

برای مدل می شود پس نسبت به perturbation های جزئی مقاوم تر است.

پس بطور کلی در برابر حمله مقاوم تر می باشد.

اما از سوی دیگر نشان می دهیم، با افزایش پارامترهای مدل، decision boundary مدل مان بسیار پیچیده تر می شود و ممکن است مرزهای بسیار (فوق) حول و اطراف داده ها در نظر بگیرد و نسبت به استقلال های جزئی حساس تر شده و باعث (مستعدی) اشتباه شود پس در این صورت CNN معمولی چون پارامترهای بیشتری دارد از همه ضعیف تر است.

سوال (2)

Untargeted attack

(الف)

$$\max \mathcal{L}(h(x+\delta), y)$$

$$\|\delta\| \leq \epsilon$$

با توجه به اینکه در حالت untargeted به دنبال

این موضوع هستیم صرفاً پسبینی Classifier

را هدف اصلی به جز آن تبدیل کنیم پس فقط کفایت

با perturbation کمی داده را از مرز نبردی کلاس

خارج کنیم. اما در حالت targeted چون

می خواهیم پسبینی به یک کلاس خاص تبدیل شود

محتمل است نیاز به perturbation بزرگتر و حمله

پرجاش تر داشته باشیم.

Targeted attack

$$\min \mathcal{L}(h(x+\delta), y')$$

$$\|\delta\| \leq \epsilon$$

پس حمله untargeted نسبتاً آسان تر

خواهد بود و موفقیت بیشتری دارد.

$$-\nabla_{\delta_i} \ell(w^T(n+\delta), y) + \nabla_{\delta_i} \sum \alpha_i (|\delta_i| - \epsilon) = 0 \quad (2)$$

$$-\left(\nabla_h \ell \cdot \nabla_{\delta_i} h\right) + \alpha_i \text{Sign}(\delta_i) = 0 \Rightarrow$$

$$-\left(\nabla_h \ell \cdot w_i\right) + \alpha_i \text{Sign}(\delta_i) = 0 \Rightarrow$$

$$\alpha_i = \frac{\nabla_h \ell \cdot w_i}{\text{Sign}(\delta_i)} \quad (\delta_i \neq 0)$$

$$\delta_i = \text{Sign}(\delta_i) |\delta_i|$$

(>

Complementary slackness

$$\sum_{i=1}^n \alpha_i (|\delta_i| - \epsilon) = 0 \quad \alpha_i \geq 0 \Rightarrow \alpha_i (|\delta_i| - \epsilon) = 0$$

$$\alpha_i \left(\frac{\delta_i}{\text{Sign}(\delta_i)} - \epsilon \right) = 0 \Rightarrow \alpha_i = 0 \Rightarrow \nabla_{\delta_i} \ell = 0$$

يعني مقدار δ_i كبير، α_i صفر، زمانی که $|\delta_i| < \epsilon$ بود.

$$\frac{\nabla_h \ell \cdot w_i}{\text{Sign}(\delta_i)} \cdot \frac{\delta_i}{\text{Sign}(\delta_i)} - \frac{\nabla_h \ell \cdot w_i}{\text{Sign}(\delta_i)} \epsilon = 0 \Rightarrow$$

+1

$$\delta_i = \frac{\epsilon}{\text{Sign}(\delta_i)} \Rightarrow |\delta_i| = \epsilon$$

هـ) در عبارت نسبت "ب" داریم :

Stationary

$$-\nabla_{\delta} \ell(h(u), y) + \nabla_{\delta} \alpha (\|\delta\|_{\infty} - \epsilon) = 0$$

$$\nabla_{\delta} \ell(h(u), y) = \nabla_{\delta} \alpha (\|\delta\|_{\infty} - \epsilon)$$

به ازای دو مقدار بدست آمده در بخش "د" یعنی $\alpha_i = 0$ یا $|\delta_i| = \epsilon$ با جایگذاری در عبارت بالا می توانیم نتیجه بگیریم در حالت بهینه قرار داریم.

$$\|\delta\|_{\infty} = \max_i |\delta_i|$$

چون به ازای هر i داریم :

$$\left. \begin{aligned} |\delta_i| = \epsilon &\Rightarrow \nabla_{\delta_i} \alpha_i (\epsilon - \epsilon) = 0 \Rightarrow \boxed{\nabla_{\delta} \ell = 0} \\ \text{or} \\ \alpha_i = \epsilon &\Rightarrow \nabla_{\delta_i} 0 (|\delta_i| - \epsilon) = 0 \Rightarrow \boxed{\nabla_{\delta} \ell = 0} \end{aligned} \right\}$$

← پس با هر دو مقدار بدست آمده در بخش "د" می توانیم بگوییم مقدار مسطح Loss نسبت به δ صفر شده در نتیجه در نقطه بهینه یا ماکزیمیم تابع قرار داریم.

الف) ابتدا یک perturbation اولیه انتخاب می شود مانند v که می تواند صفر مقدار دهی شود.

سپس بصورت iterative روی تمام data point ها v را اعمال می کنیم به امید آنکه کلاس آن داده تغییر کند. اگر v نتوانست کلاس را تغییر بدهد.

سپس یک Δv ای پیدا می کنیم که اگر به داده perturb شده اضافه شود

بتواند کلاس داده را تغییر دهد. سپس $v \leftarrow v + \Delta v$ کرده و سپس

v را به حد بالای آن (ع) project می کنیم تا شرط $\|v\| \leq \epsilon$ را ارضا کنیم.

این کار را بصورت تکرار می شوند روی همه داده ها انجام می دهیم. اگر $k-1$ از کلی داده ها کلاس بندی اشتباه شوند کارمان به پایان می رسد و اگر این کار نشود دوباره همان حلقه تکرار می شوند و برای کلی داده ها انجام می دهیم.

ب)

data point ها در ابتدا بالا فقط بخش کوچکی از فضا را اشغال می کنند همچنین

منه های تصمیم گیری بسیار پیچیده بوده و در بسیاری از بدها بی معنی است.

پس می توان بایک آشفتمی Universal در برخی بدها باعث شد همه داده ها

تست به دست بندی شوند.

این موضوع بخاطر این است که شبکه های عمیق صرف الگوهای را در بخش از فضا یاد

می گیرند که کلاس ها را از هم جدا کند و یک سری آسیب پذیری مشترک برای همه

داده ها در ایجاد بالا بوجود می آید که می تواند مورد حمله Universal قرار گیرد.

همچنین دلیل دیگر می تواند این باشد که شبکه های عمیق در ایجاد بالا روی الگوهای پیچیده

over Fit می شوند که این موضوع منتج به یادگیری ویژگی های غیر مستحکم شده و اغتشاش های

Universal ای بوجود می آید که مدل ها نسبت به آن آسیب پذیر شود.

ج) برای ریاست روی بردی two moons سه دایرکشن مختلف برای اضافه کردن نویز را بررسی کردیم که به ترتیب در جهت بردارهای $(1, 0)$ ، $(0, 1)$ و $(\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}})$ بوده که L_2 norm همگی برابر 1 هست. در این حالت اگر این بردارها را در ϵ ضرب کرده و به همان اضافه کنیم می توانیم بگوئیم نرم میزان تغییرات برابر با ϵ است.

پس از بررسی متوجه می شویم که با افزایش ϵ در جهت بردارهای $(0, 1)$ و $(\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}})$ Success rate همان رویه افزایشی بوده است اما در جهت بردار $(1, 0)$ با افزایش مقدار ϵ ، تغییرات خاصی در Success rate مشاهده نمی شد پس این می تواند یک Failure scenario برای همان باشد. (چارت مربوط به experiment در صفحه بعد قرار داده شده است)

توفیقاً اضافه روضی برای پارت الف)

$v \leftarrow 0$ initialize

while $\frac{1}{m} \sum_{i=1}^m 1_{\hat{K}(\alpha_i + v) \neq \hat{K}(\alpha_i)} \leq 1 - \delta$:

به ازای هر داده موجود در ریاست اگر $\hat{K}(\alpha_i + v) = \hat{K}(\alpha_i)$ است

$$\Delta v_i \leftarrow \arg \min_r \|r\|_2$$

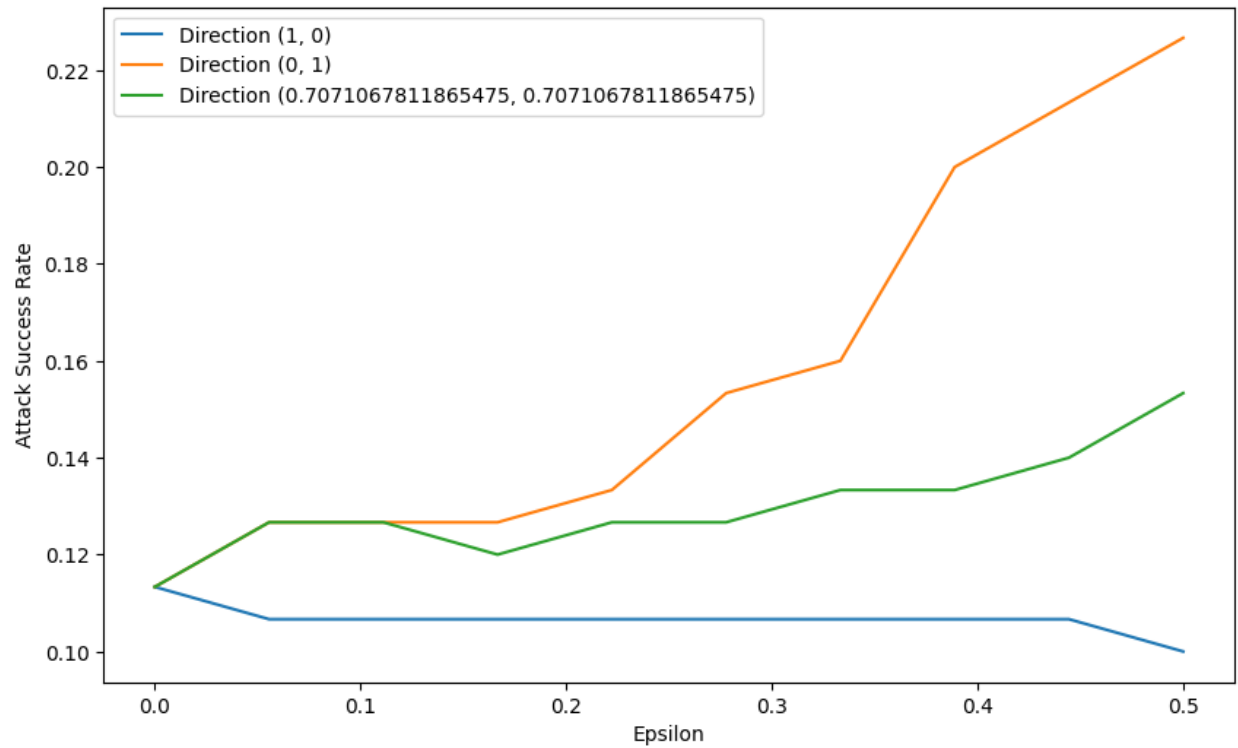
به شرط اینکه $\hat{K}(\alpha_i + v + r) \neq \hat{K}(\alpha_i)$ شود

$$v \leftarrow P_{\mathcal{E}}(v + \Delta v_i)$$

مقدار نهایی perturbation را به مقدار \mathcal{E} project می کند.

این لوپ زمانی که حداقل $1 - \delta$ از کل داده ها استباه دسته بندی شوند ادامه می یابد.

چارت experiment مربوط به بخش "ج" سوال سوم تئوری:



با توجه به چارت متوجه میشویم که با افزایش epsilon در جهت بردار $(1,0)$ بهبودی در attack success rate مشاهده نمیشود و این یک failure scenario برای حمله مان است.