

(a) اگر $\text{target-class} = t$ ، هدف تغییر کلاس پیش بین شده مدل به کلاس t است:

$$X_{adv} = X - \epsilon \text{Sign} \nabla_{\theta} \text{loss}(\theta; X, t)$$

(b)

$$X = X_0 + \delta$$

$$w^T X + w_0 = w^T X_0 + w^T \delta + w_0$$

$$\Rightarrow \text{optimization: } \min_{\delta} w^T \delta + w^T X_0 + w_0$$

$$\text{s.t. } \|\delta\|_{\infty} \leq \epsilon$$

با توجه به...

$$\text{Holder} \rightarrow w^T \delta \geq -\|\delta\|_{\infty} \|w\|_1 \geq -\epsilon \|w\|_1$$

$$w^T \delta = -\epsilon \|w\|_1 \text{ وقتی } \delta \text{ برابر } \delta = -\epsilon \cdot \text{sign}(w) \text{ است؛ زیرا:}$$

$$w^T \delta = -\epsilon w^T \text{sign}(w) = -\epsilon \sum_{i=1}^d w_i \text{sign}(w_i) = -\epsilon \sum_{i=1}^d |w_i| = -\epsilon \|w\|_1 \quad (1)$$

FGSM loss function two-class binary classifier:

objective function
در پیشینه
حمله

$$\rho(X) = w^T X + w_0 = (w_i^T X + w_{i_0}) - (w_t^T X + w_{t_0})$$

باید توجه داشت که w_i و w_{i_0} (مطلوب) و w_t و w_{t_0} (غیرمطلوب) است.

$$\text{اگر } (w_i^T X + w_{i_0}) - (w_t^T X + w_{t_0}) \leq 0 \text{، missclassified } X \text{ می باشد. در نتیجه تابع هزینه را میتوان اینگونه تعریف کرد:}$$

$$J(X) = (w_t^T X + w_{t_0}) - (w_i^T X + w_{i_0})$$

$$J(X; w) = -(w^T X + w_0) \text{، پیشینه کردن این تابع هزینه، تعیین می کند دسته بندی بدست می شود.}$$

first order

$$\text{approximation: } J(X; w) = J(X_0 + \delta; w) \approx J(X_0; w) + \nabla_{\theta} J(X_0; w)^T \delta$$

$$\text{بنابراین پیدا کردن } \delta \text{ که } J(X_0 + \delta; w) \text{ را بیشینه می کند تقریباً معادل پیدا کردن } \delta \text{ که } J(X_0; w) + \nabla_{\theta} J(X_0; w)^T \delta \text{ را بیشینه می کند است.}$$

$$\begin{aligned} \max_{\delta} J(X_0; w) + \nabla_{\theta} J(X_0; w)^T \delta &\Rightarrow \min_{\delta} -\nabla_{\theta} J(X_0; w)^T \delta - J(X_0; w) \\ \text{s.t. } \|\delta\|_{\infty} \leq \epsilon &\quad \text{s.t. } \|\delta\|_{\infty} \leq \epsilon \end{aligned}$$

این مسئله optimization مدل همان مسئله بهینه سازی رابطه (1) است. بنابراین جواب خوبی می باشد:

$$\delta = -\epsilon \cdot \text{sign}(\nabla_{\theta} J(X_0; w)) \Rightarrow X = X_0 + \epsilon \cdot \text{sign}(\nabla_{\theta} J(X_0; w)) \Rightarrow \text{FGSM}$$

یک بهینه برای این تابع هدف است.

با توجه به اینکه شبکه های عصبی معمولاً غیرخطی هستند پس $\text{loss}(X; w)$ همیشه یک چیز ثابت نیست و با تک مرحله و مدل کاملاً ممکن است آشفتگی بهینه روح حاصل شود و نیاز به یک تکرار iteratively با مدل کاملاً ممکن باشد.

(a) PGD برخلاف FGSM که در یک مرحله سعی در حل مسئله دارد، به صورت iterative عمل می‌کند.

$$x_{adv}^t = \Pi_{\mathcal{E}} \left(x_{adv}^{t-1} + \alpha \text{Sign} \nabla_{\theta} \text{Loss}(\theta; x_{adv}^{t-1} + \delta, y) \right)$$

به سبب آنکه چقدر ϵ است که باعث می‌شود بایک مرحله کامل تو جعبه در راستای گرادیان وی با تعداد بیشتر، حرکت کند. و یک Projection مانع از خروج از محدوده \mathcal{E} می‌شود. به خاطر iterative بودن، یک مرحله‌ای نبودن با طول گام زیاد می‌تواند در فیل از موارد با تغییرات کمتر نسبت به نقطه شروع، محدوده ضعیف به مدل داشته باشد.

(b)

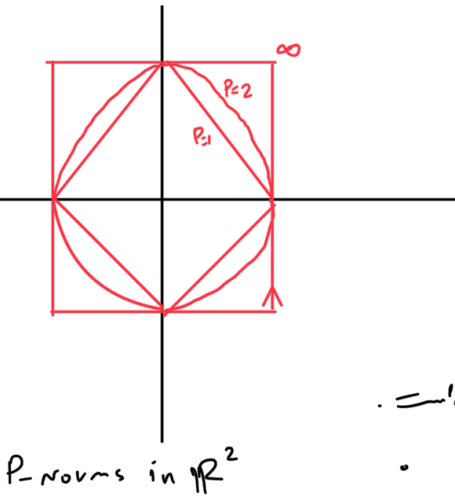
به حس تدان از L_0 ، L_1 ، L_2 ، L_∞ استفاده کرد. در حقیقت تصویر

L_∞ سعی می‌کند با تغییر زیاد تعداد کم پیکسل‌ها و L_0 با تغییر همه پیکسل‌ها و به میزان کم در L_1 و L_2 ترکیب از این دو رویکرد را دارند.

هم چنین مسئله Optimization مربوط به L_∞ غیر معمولی و سخت است.

L_1 در مقایسه با L_2 تعداد کمتری پیکسل را معمولاً تغییر می‌دهد و حدیثی کاهشی قدر مطلق تغییرات اندازه پیکسل‌هاست، Optimization آن آسان‌تر است.

هدف L_2 اما کاهشی نامیده می‌شود از تغییرات پیکسل‌هاست و Optimization آن معمولی است. L_∞ سعی در کاهش بیشینه قدر مطلق تفاضل پیکسل‌هاست و ...



(c)

L_2 به نسبت با اندازه تصویر وابسته است، اگر با اعتبار استفاده شود ممکن است باعث ایجاد نمونه‌های خصمانه اشتباه شود (یعنی نمونه‌هایی که باعث تغییر کلاس آن تصویر یا به عبارت دیگر، تصویر بر غیر قابل تشخیص تولید کند و می‌تواند باعث تغییر پیکسل‌ها به میزان بسیار کم شود و کنترل اینکه نمونه خصمانه اشتباه تولید نکند راحت‌تر است).