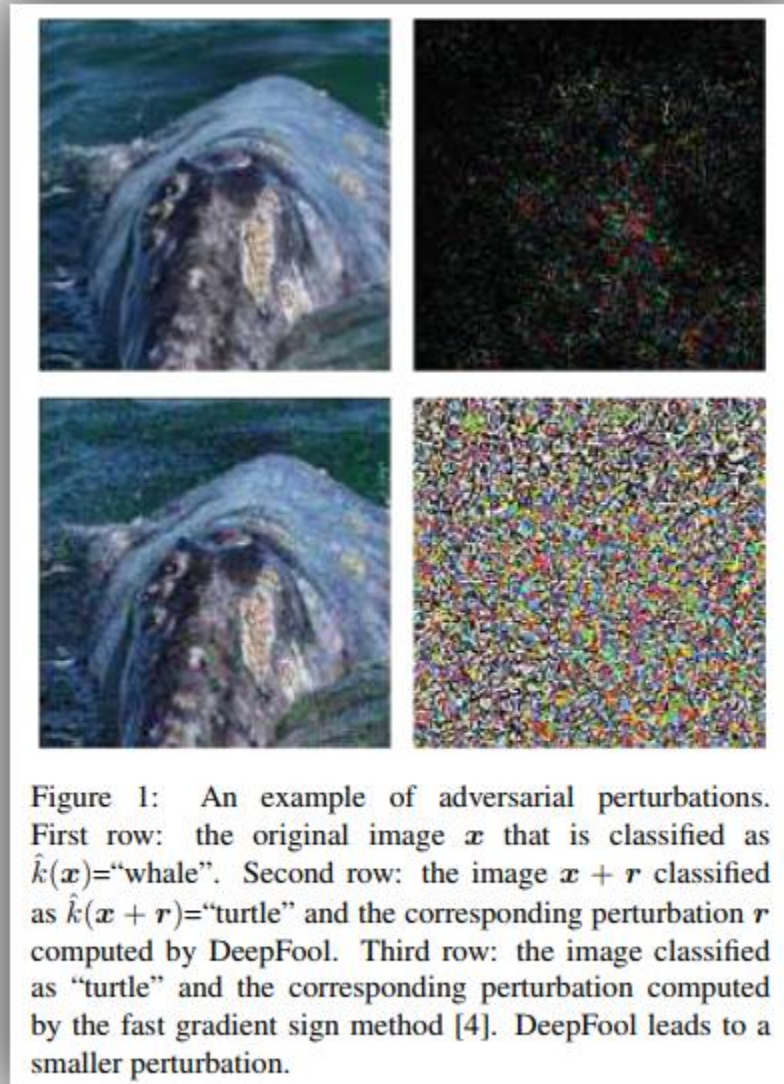


به نام خدا

Adversarial attack: اگر مقداری مانند r به تصویر اضافه شود به طوری که تصویر جدید به صورت چشمی شبیه به تصویر اصلی باشد ولی کلسیفایر ما را به اشتباه بیندازد و برچسب داده بعد از اضافه کردن این مقدار r (که به آن perturbation گفته میشود) متفاوت از برچسب بدون r شود. مانند شکل زیر:



برای مقابله با این مشکل روشی به اسم DeepFool ارائه شده است که به مختصر درباره آن توضیح میدهم:

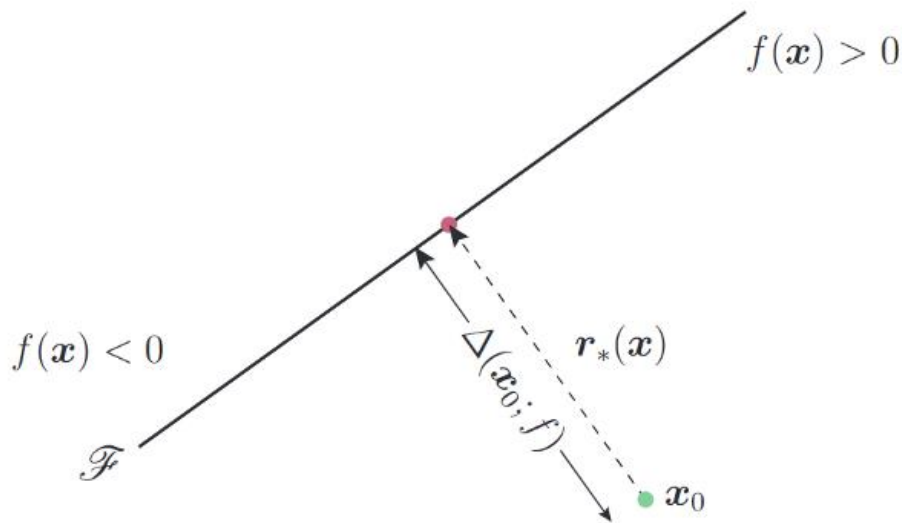


Figure 2: Adversarial examples for a linear binary classifier.

Robustness این مدل فاصله x_0 از خطی است که دو کلاس را از هم جدا میکند.

کوچکترین perturbation ای که تصمیم کلسیفایر را عوض میکند برابر با پروجکشن اورتوگونال (ساخت زاویه عمود با خط کلسیفایر) به خط میباشد که فرمول آن به صورت زیر میباشد:

$$-\frac{f(x_0)}{\|w\|_2^2} * w$$

توضیحاتی که دادیم برای حالت affine بود و اگر بخواهیم به طور کلی یک الگوریتم برای کلسیفایر باینری معرفی کنیم به صورت زیر میباشد:

Algorithm 1 DeepFool for binary classifiers

- 1: **input:** Image x , classifier f .
 - 2: **output:** Perturbation \hat{r} .
 - 3: Initialize $x_0 \leftarrow x, i \leftarrow 0$.
 - 4: **while** $\text{sign}(f(x_i)) \neq \text{sign}(f(x_0))$ **do**
 - 5: $r_i \leftarrow -\frac{f(x_i)}{\|\nabla f(x_i)\|_2} \nabla f(x_i),$
 - 6: $x_{i+1} \leftarrow x_i + r_i,$
 - 7: $i \leftarrow i + 1.$
 - 8: **end while**
 - 9: **return** $\hat{r} = \sum_i r_i.$
-

توضیح الگوریتم:

- 1- الگوریتم ورودی x و کلسیفایر f را میگیرد.
- 2- خروجی کوچکترین perturbation مورد نیاز برای دسته بندی اشتباه تصویر میباشد.
- 3- مقدار دهی اولیه adver image با داده ورودی میباشد.
- 4- حلقه را شروع میکنیم و ادامه میدهیم تا تا زمانی که برچسب صحیح و برچسب perturbed image یکی باشد.
- 5- پروژکشن ورودی به نزدیکترین هایپرپلن را حساب میکنیم تا minimal perturbation به دست آید.
- 6- این perturbation را به تصویر اضافه و تست میکنیم.
- 7-8- افزایش حلقه و اتمام حلقه
- 9- برگرداندن minimal perturbation

برای کلسیفایر های چند کلاسه برای هر کلاس یک هایپرپلن وجود دارد که یک کلاس را از بقیه کلاس ها جدا میکند و مبنی بر جایی که x در فضا دارد دسته بندی میشود. تمام کاری که این الگوریتم انجام میدهد این است که نزدیکترین هایپرپلن را پیدا میکند و این x را پروژکت میکند به آن هایپرپلن و آن را کمی فراتر از آن مرز میبرد تا دسته بندی اشتباه با minimal perturbation رخ دهد.

Deepfool همچنین الگوریتمی را فراهم کرده است که بتواند تحمل پذیری adversarial را اندازه گیری کند.

منبع:

<https://towardsdatascience.com/deepfool-a-simple-and-accurate-method-to-fool-deep-neural-networks-17e0d0910ac0>