

Concept-based Adversarial Attack: a Probabilistic Perspective



Andi Zhang^{a,b} Xuan Ding^c Steven McDonagh^d Samuel Kaski^{b,e}

^aUniversity of Warwick ^bUniversity of Manchester
^cCUHK (Shenzhen) ^dUniversity of Edinburgh ^eAalto University

Review: Adversarial Attack

Adversarial Attack has two aims:

- **Deceiving victim classifiers**
- **Deceiving human**

Given an image $x_{ori} \in [0, 1]^n$ and a target label $y_{tar} \in \mathcal{Y}$, the optimization problem [1] to find an adversarial example x_{adv} :

$$\min c_1 \mathcal{D}(x_{ori}, x_{adv}) + c_2 f(x_{adv}, y_{tar}) \quad \text{s.t. } x_{adv} \in [0, 1]^n \quad (1)$$

where \mathcal{D} is a distance function (commonly chosen from ℓ_1 , ℓ_2 , or ℓ_∞), and f is an objective function tied to the victim classifier's prediction (e.g., the cross-entropy loss).

Probabilistic Adversarial Attack

By applying the **box-constrained Langevin dynamics** as an optimization method to (1), we derive a probabilistic perspective of adversarial attack [2]:

$$p_{adv}(x_{adv}|x_{ori}, y_{tar}) \propto p_{vic}(x_{adv}|y_{tar}) p_{dis}(x_{adv}|x_{ori}) \quad (2)$$

where $p_{vic}(x_{adv}|y_{tar}) \propto \exp(-c_2 f(x_{adv}, y_{tar}))$ is the victim distribution and $p_{dis}(x_{adv}|x_{ori}) \propto \exp(-c_1 \mathcal{D}(x_{ori}, x_{adv}))$ is the distance distribution.

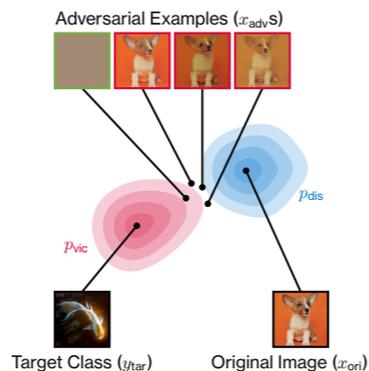
If \mathcal{D} is ℓ_1 , then p_{dis} is a Laplace distribution; if \mathcal{D} is squared ℓ_2 , then p_{dis} is a Gaussian distribution.

We suggest replacing geometry-induced distributions with **PGM-learned** distributions for p_{dis} [2].

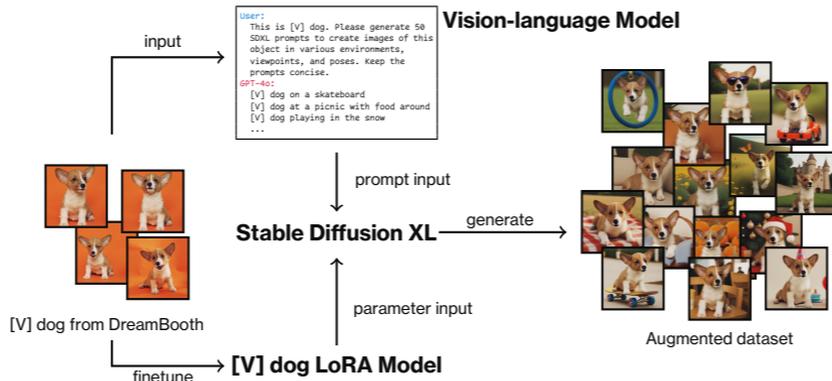
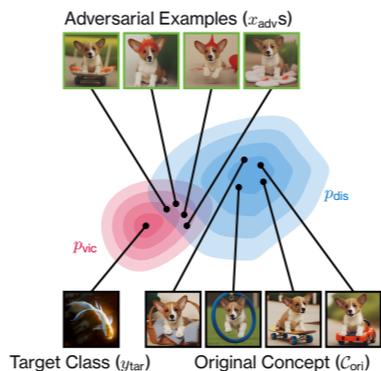
Concept-based Adversarial Attack

The probabilistic perspective naturally extends $p_{dis}(\cdot|x_{ori})$ to $p_{dis}(\cdot|C_{ori})$.

Adversarial Attack on a Single Image



Concept-based Adversarial Attack



Generated Adversarial Examples

x_{ori} / C_{ori}									
y_{tar}	ambulance	basenji	hamster	king snake	panda	goldfish	mushroom	laptop	papillon
NCF									
ACA									
DiffAtk									
ProbAtk									
OURS (CONS)									
OURS (AGGR)									

Green: successfully fools the victim (surrogate) classifier; Red: fails.

Reference

- [1] Szegedy, Christian, et al. "Intriguing properties of neural networks." ICLR (2014).
 [2] Zhang, Andi, Mingtian Zhang, and Damon Wischik. "Constructing semantics-aware adversarial examples with a probabilistic perspective." NeurIPS (2024).



A!