



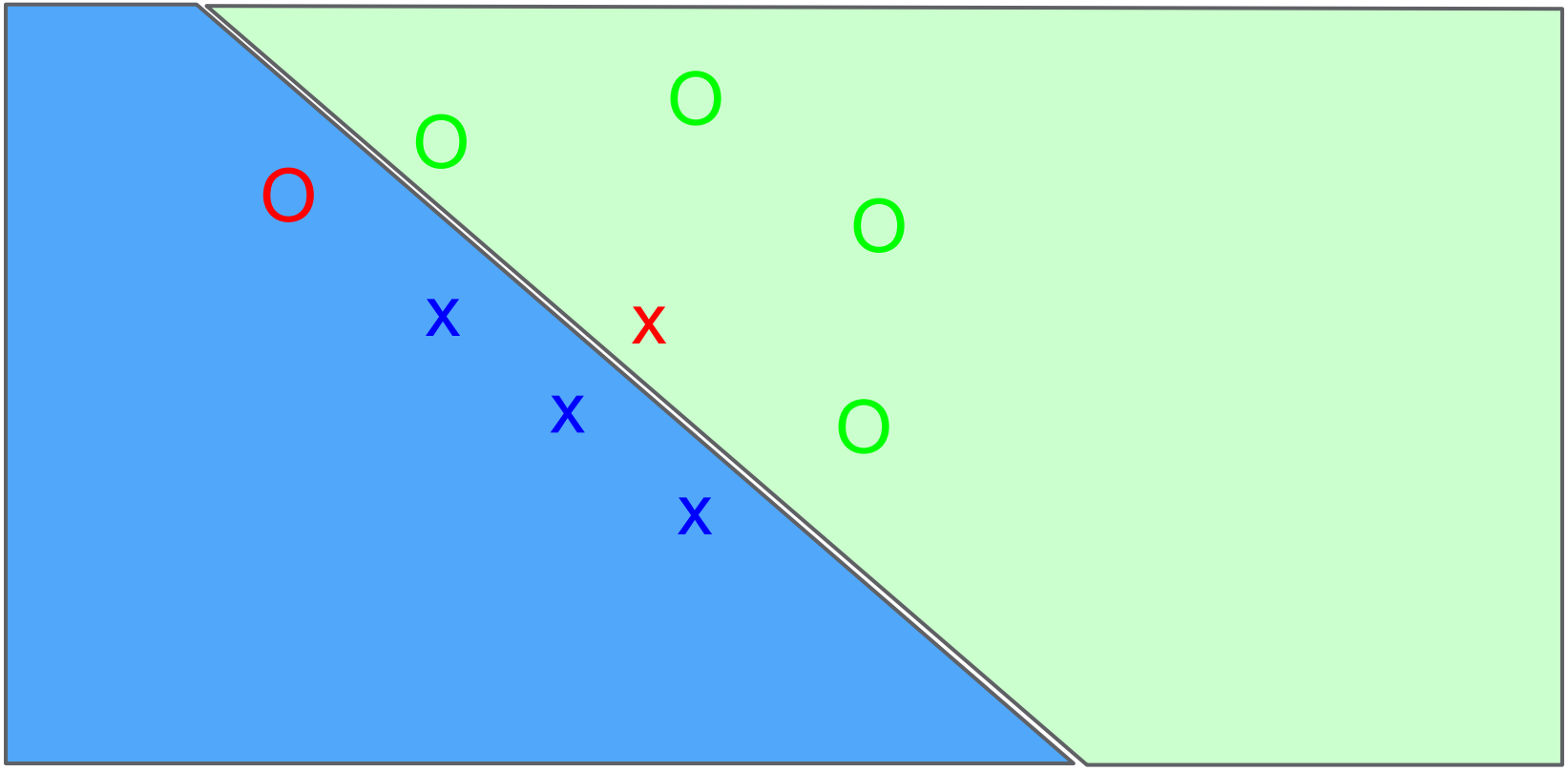


ADVANCE ALL MARCH EVERLASTING

EXPLANATIONS FOR ALL EXAMPLES

15

- ▶ In *Explaining and Harnessing Adversarial Examples* Goodfellow et al. argues that adversarial examples exist because of the piece wise linearity in deep neural models
- ▶ Adversarial Examples occur at the difference between model and real class boundary.
- ▶ So purpose of the attack model is to perturb the original input so that it falls into desired region.



ImageSource: [iangoodfellow.com](http://iangoodfellow.com)



## ATTACK MODELS

- ▶ L-BFGS (Broyden-Fletcher-Goldfarb-Shanno)
- ▶ The Fast Gradient Sign Method (FGSM)
- ▶ Projected Gradient Descent (PGD or FGSM<sup>k</sup>)
- ▶ Jacobean Based Saliency Map Approach (JSMA)
- ▶ Carlini Wagner Attack
- ▶ Black Box attack

## EXPLANATION FOR ADVERSARIAL EXAMPLES

- ▶ In *Explaining and Harnessing Adversarial Examples* Goodfellow et al. argues that adversarial examples exist because of the piece wise linearity in deep neural models
- ▶ Adversarial Examples occur at the difference between model and real class boundary.
- ▶ So purpose of the attack model is to perturb the original input so that it falls into desired region.

