# Adversarial Regularization

George Yu

# A Quick Refresher



$$\boldsymbol{x}$$
"panda"

$$+ .007 \times$$

$$\mathrm{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$$
"nematode"

$$=$$

$$\boldsymbol{x} + \epsilon\mathrm{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$$
"gibbon"

[Source: Goodfellow, et al., ICLR'15]

# Adversaries - FGSM

- FGSM = Fast Gradient Sign Method
- Find the direction that increases the loss the most, and take a step
- Significantly degrades accuracy of unregularized model

# Adversaries - PGD

- PGD = Projected Gradient Descent
- Iteratively apply the FGSM algorithm
- Project back to within the perturbation budget after every step
- Produces a stronger adversary but is more expensive than FGSM.

# Adversarial Regularization in NSL

- Generate adversarial examples as graph neighbors.
- Regularize the loss by including the neighbors in the graph.
- [Colab notebook](#)