

ADVANCE ALL MARCH EVERLASTING

DEFENSE: VIRTUAL ADVERSARIAL TRAINING

3

8

- ▶ Uses both labeled and unlabelled datapoints
- ▶ Loss function:
- ▶ Regularizer:
- ▶ Objective Function:
- ▶

$$\text{LDS}(x_*, \theta) := D \left[p(y|x_*, \hat{\theta}), p(y|x_* + r_{\text{vadv}}, \theta) \right]$$

$$r_{\text{vadv}} := \arg \max_{r; \|r\|_2 \leq \epsilon} D \left[p(y|x_*, \hat{\theta}), p(y|x_* + r, \theta) \right],$$

$$\text{where } x_* \in \{D_l, D_{ul}\}$$

$$\mathcal{R}_{\text{vadv}}(\mathcal{D}_l, \mathcal{D}_{ul}, \theta) := \frac{1}{N_l + N_{ul}} \sum_{x_* \in \mathcal{D}_l, \mathcal{D}_{ul}} \text{LDS}(x_*, \theta).$$

$$\ell(\mathcal{D}_l, \theta) + \alpha \mathcal{R}_{\text{adv}}(\mathcal{D}_l, \mathcal{D}_{ul}, \theta),$$

Miyato, T., Maehara, S., Koyama, M., & Ishii, S. (2017). Virtual Adversarial Training: a Regularization Method for Supervised and Semi-supervised Learning.

NO DEFINITIVE ANSWER

DEFENSE: VIRTUAL ADVERSARIAL TRAINING

- ▶ Uses both labeled and unlabelled datapoints

- ▶ Loss function:

$$\text{LDS}(x_*, \theta) := D \left[p(y|x_*, \hat{\theta}), p(y|x_* + r_{\text{vadv}}, \theta) \right]$$

$$r_{\text{vadv}} := \arg \max_{r; \|r\|_2 \leq \epsilon} D \left[p(y|x_*, \hat{\theta}), p(y|x_* + r, \theta) \right],$$

$$\text{where } x_* \in \{D_l, D_{ul}\}$$

- ▶ Regularizer:

$$\mathcal{R}_{\text{vadv}}(\mathcal{D}_l, \mathcal{D}_{ul}, \theta) := \frac{1}{N_l + N_{ul}} \sum_{x_* \in \mathcal{D}_l, \mathcal{D}_{ul}} \text{LDS}(x_*, \theta).$$

- ▶ Objective Function:

$$\ell(\mathcal{D}_l, \theta) + \alpha \mathcal{R}_{\text{vadv}}(\mathcal{D}_l, \mathcal{D}_{ul}, \theta),$$