# ADVERSARIAL MACHINE LEARNING

# ATTACK MODELS: TERMINOLOGY

- $\mathcal{F}(\text{x})$ : Predicted Class of x

- $y$: true class

- $\theta$: Model Parameters

- $\mathcal{H}(x,\theta)$ : Output of Logits (Before final softmax $l$ayer)

- $\mathcal{L}(x,y,\theta)$ : Loss Function

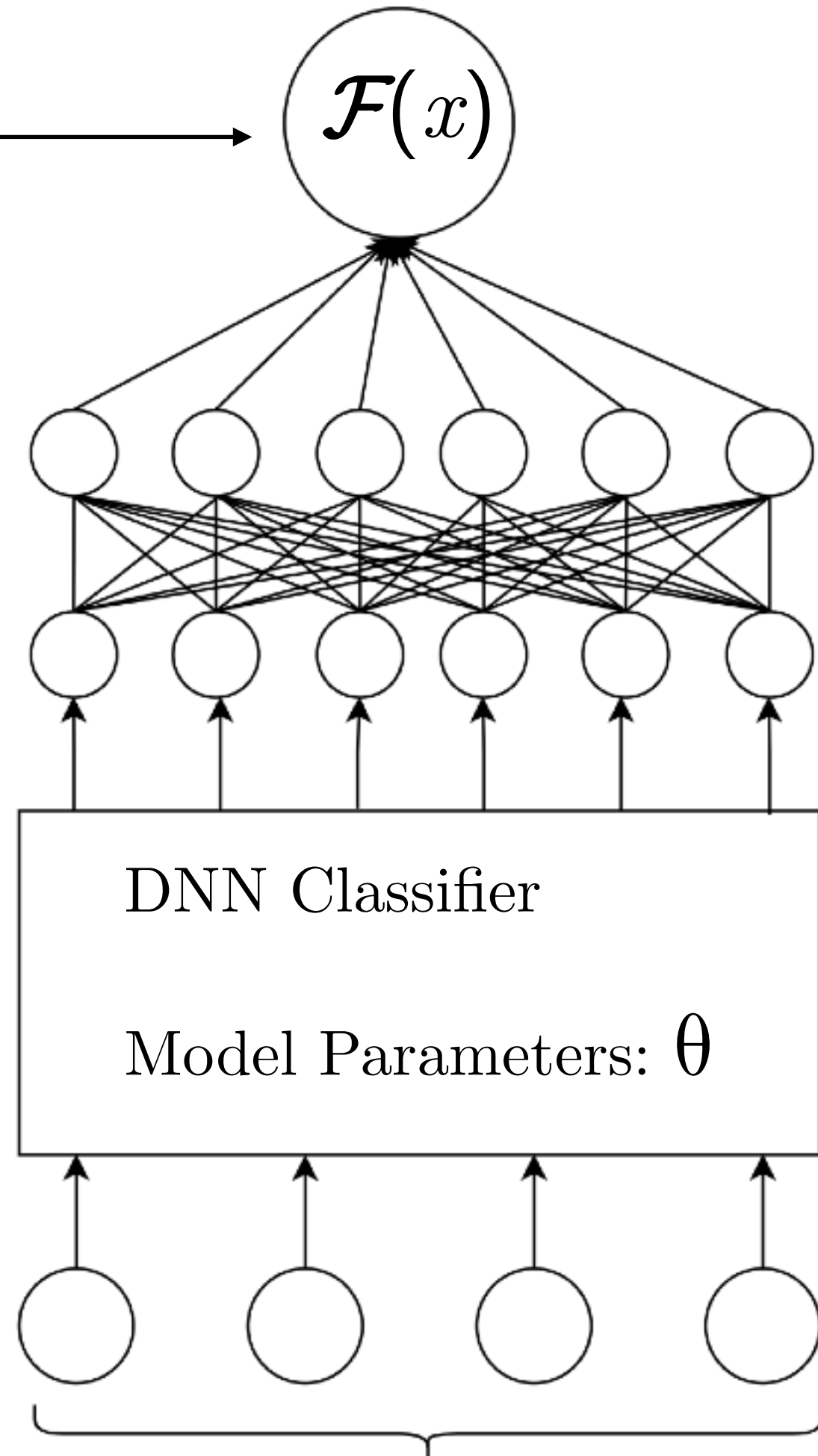- $l$: Class of interest for attacked (For Targeted Attacks)

13

$$\mathcal{F}(x) = \arg\max_{j} \mathcal{P}_j \longrightarrow \mathcal{F}(x)$$

Class Probabilities: $\mathcal{P}_j(x) = \textsf{softmax}(\mathcal{H}_i(x))$

Second Last Layer/Logits : $\mathcal{H}(x)$

DNN Classifier

Model Parameters: $\theta$

**Input**: $x$

# ATTACK MODEL PROBLEM FORMULATION:

‣ Non-Targeted Attack:

$$\underset{||\delta||_2}{\arg\min} \ \ \text{s.t.} \ \ \mathcal{F}(x + \delta) \neq \mathcal{F}(x)$$

‣ Targeted Attack

$$\underset{||\delta||_2}{\arg\min} \ \ \text{s.t.} \ \ \mathcal{F}(x + \delta) = \ell \ \text{Target Class}$$

# ATTACK MODELS: TERMINOLOGY

$$\mathcal{F}(x) = \arg\max_{j} \mathcal{P}_j \longrightarrow \boxed{\mathcal{F}(x)}$$

▸ $\mathcal{F}(\mathrm{x})$ : Predicted Class of x

▸ $y$: true class

Class Probabilities: $\mathcal{P}_j(x) = \mathsf{softmax}(\mathcal{H}_i(x))$

▸ θ: Model Parameters

Second Last Layer/Logits : $\mathcal{H}(x)$

▸ $\mathcal{H}(x,θ)$ : Output of Logits (Before

final softmax $l$ayer)

DNN Classifier

Model Parameters: θ

▸ $\mathcal{L}(x,y,θ)$ : Loss Function

▸ $l$: Class of interest for attacked
(For Targeted Attacks)

**Input**: $x$