# ADVERSARIAL MACHINE LEARNING

# ATTACK MODEL : JSMA ALGORITHM

**Algorithm 1 Crafting adversarial samples**

$\mathbf{x}$ is the benign sample, $\ell$ is the target network output, $\mathcal{F}$ is the function learned by the network during training, $\Upsilon$ is the maximum distortion, and $\theta$ is the change made to features.

**Input:** $\mathbf{x}, \ell, \mathcal{F}, \Upsilon, \theta$

1: $\mathbf{x}^* \leftarrow \mathbf{x}$
2: $\Gamma = \{1 \dots |\mathbf{x}|\}$
3: **while** $\mathcal{F}(\mathbf{x}^*) \neq l$ and $||\delta_{\mathbf{x}}|| < \Upsilon$ **do**
4:     Compute forward derivative $\nabla \mathcal{P}(\mathbf{x}^*)$
5:     $S = \texttt{saliency\_map}\left(\nabla \mathcal{P}(\mathbf{x}^*), \Gamma, l\right)$
6:     Modify $\mathbf{x}^*_{i_{max}}$ by $\theta$ s.t. $i_{max} = \arg \max_i S(\mathbf{x}, l)[i]$
7:     $\delta_{\mathbf{x}} \leftarrow \mathbf{x}^* - \mathbf{x}$
8: **end while**
9: **return** $\mathbf{x}^*$

# ATTACK MODELS: CARLINI WAGNER ATTACKS

▸ Nicholas Carlini and David Wagner proposer three attack models with $L_0$, $L_2$ and $L_\infty$ distance

▸ L2 attack is most optimal and broke all the existing defences

▸ The defined the problem as following:

$$\text{minimize } \mathcal{D}(x, x + \delta)$$
$$\text{such that } \mathcal{C}(x + \delta) = \ell$$
$$x + \delta \in [0, 1]^n$$
$$\mathcal{D} \text{ can be } L_0, L_2, L_\infty \text{ distance}$$

Carlini, N., & Wagner, D. (2017, May). Towards evaluating the robustness of neural networks.

# ATTACK MODEL : JSMA ALGORITHM

---

**Algorithm 1 Crafting adversarial samples**

$\mathbf{x}$ is the benign sample, $\ell$ is the target network output, $\mathcal{F}$ is the function learned by the network during training, $\Upsilon$ is the maximum distortion, and $\theta$ is the change made to features.

---

**Input:** $\mathbf{x}, \ell, \mathcal{F}, \Upsilon, \theta$

1: $\mathbf{x}^* \leftarrow \mathbf{x}$
2: $\Gamma = \{1 \ldots |\mathbf{x}|\}$
3: **while** $\mathcal{F}(\mathbf{x}^*) \neq l$ and $||\delta_{\mathbf{x}}|| < \Upsilon$ **do**
4:     Compute forward derivative $\nabla \mathcal{P}(\mathbf{x}^*)$
5:     $S = \texttt{saliency\_map}(\nabla \mathcal{P}(\mathbf{x}^*), \Gamma, l)$
6:     Modify $\mathbf{x}^*_{i_{max}}$ by $\theta$ s.t. $i_{max} = \arg\max_i S(\mathbf{x}, l)[i]$
7:     $\delta_{\mathbf{x}} \leftarrow \mathbf{x}^* - \mathbf{x}$
8: **end while**
9: **return** $\mathbf{x}^*$

---