# ADVERSARIAL MACHINE LEARNING

# ATTACK MODELS: JSMA

- Saliency Map based greedy approach

- Modify the pixel who will impact the classifier output most

- Saliency Map is defined as:

Papernot, N., McDaniel, P., Jha, S., Fredrikson, M., Celik, Z. B., & Swami, A. (2016, March). The limitations of deep learning in adversarial settings.

$$\mathcal{S}(x,\ell)[i] = \begin{cases} 0 \text{ if } \frac{\partial \mathcal{P}_l(\mathbf{x})}{\partial \mathbf{x}_i} < 0 \text{ or } \sum_{j \neq t} \frac{\partial \mathcal{P}_j(\mathbf{x})}{\partial \mathbf{x}_i} > 0 \\ \left( \frac{\partial \mathcal{P}_l(\mathbf{x})}{\partial \mathbf{x}_i} \right) \left| \sum_{j \neq t} \frac{\partial \mathcal{P}_j(\mathbf{x})}{\partial \mathbf{x}_i} \right| \text{ otherwise} \end{cases}$$

# 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:     $S = \mathtt{saliency\_map}(\nabla \mathcal{P}(\mathbf{x}^*), \Gamma, l)$
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** $\mathbf{x}^*$

---

# ATTACK MODELS: JSMA

▸ Saliency Map based greedy approach

▸ Modify the pixel who will impact the classifier output most

▸ Saliency Map is defined as:

$$
\mathcal{S}(x, \ell)[i] = \begin{cases} 0 \text{ if } \frac{\partial \mathcal{P}_l(\mathbf{x})}{\partial \mathbf{x}_i} < 0 \text{ or } \sum_{j \neq t} \frac{\partial \mathcal{P}_j(\mathbf{x})}{\partial \mathbf{x}_i} > 0 \\ \left( \frac{\partial \mathcal{P}_l(\mathbf{x})}{\partial \mathbf{x}_i} \right) \left| \sum_{j \neq t} \frac{\partial \mathcal{P}_j(\mathbf{x})}{\partial \mathbf{x}_i} \right| \text{ otherwise} \end{cases}
$$

Papernot, N., McDaniel, P., Jha, S., Fredrikson, M., Celik, Z. B., & Swami, A. (2016, March). The limitations of deep learning in adversarial settings.