
Algorithm S1: Computation of a targeted UAP

Input: Set X of input images, target class y , classifier C , set Q of search directions, attack strength ϵ , cap ξ on L_p norm of the perturbation, norm type p (1, 2, or ∞), maximum number i_{\max} of iterations.

Output: non-targeted UAP vector δ

```
1:  $\delta \leftarrow \mathbf{0}$ ,  $r \leftarrow 0$ ,  $i \leftarrow 0$ 
2: while  $r < 1$  and  $i < i_{\max}$  do
3:   Pick a direction randomly:  $q \in Q$ 
4:   for  $\alpha \in \{-\epsilon, \epsilon\}$  do
5:      $\delta' \leftarrow \text{project}(\delta + \alpha q, p, \xi)$ 
6:     if  $\sum_{x \in X} p_C(y|x + \delta') > \sum_{x \in X} p_C(y|x + \delta)$  then
7:        $\delta \leftarrow \delta'$ 
8:       break
9:     end if
10:  end for
11:   $r \leftarrow |X|^{-1} \sum_{x \in X} \mathbb{I}(C(x + \delta) = y)$ 
12:   $i \leftarrow i + 1$ 
13: end while
```
