

## 7. Appendix

---

### Algorithm 1: Non-Dominating Sorting

---

**Input:** combined population  $P$ , objective vectors  $L$

```

1  $\mathcal{F} \leftarrow \{\}$  // set of fronts
2 for  $p \in P$  do
3    $S_p \leftarrow \{\}$  //set of  $p$  dominated solutions
4    $n_p \leftarrow 0$  // domination counter of  $p$ 
5   for  $q \in P$  do
6     if  $p$  dominates  $q$  then
7        $S_p \leftarrow S_p \cup \{q\}$ 
8     else if  $q$  dominates  $p$  then
9        $n_p \leftarrow n_p + 1$ 
10  if  $n_p == 0$  then
11     $p_{rank} = 1$  //  $p$  belongs to the first front
12     $\mathcal{F}_1 \leftarrow \mathcal{F}_1 \cup \{p\}$ 
13   $i \leftarrow 1$  // initialize front counter
14  while  $\mathcal{F}_i \neq \emptyset$  do
15     $Q \leftarrow \emptyset$  // store solutions of the next front
16    for  $p \in \mathcal{F}_i$  do
17      for  $q \in S_p$  do
18         $n_q \leftarrow n_q - 1$ 
19        if  $n_q == 0$  then
20          //  $q$  belongs to the next front
21           $q_{rank} \leftarrow i + 1$ 
22           $Q \leftarrow Q \cup \{q\}$ 
23     $i \leftarrow i + 1$ 
24     $\mathcal{F}_i \leftarrow Q$ 
25  return  $\mathcal{F}$ 

```

Algorithm 1 described the non-dominated sorting method proposed by Deb et al. [13] for sorting solutions with multiple objectives. The method first determines solutions in the first front (not dominated by any solutions) then iteratively constructs the remaining fronts.

---

### Algorithm 2: Sparse-RS $p_m$ Selection Method

---

**Input:** conducted function evaluations  $i$ , budget  $N$ , initial mutation size  $\alpha_{init}$

```

1  $t \leftarrow \text{int}(\frac{i}{N} \cdot 10000)$ 
2  $c \leftarrow \{0, 50, 200, 500, 1000, 2000, 4000, 6000, 8000\}$ 
3  $j \leftarrow \text{index of } c \text{ that is closest to } t$ 
4  $\beta \leftarrow \{2, 4, 5, 6, 8, 10, 12, 15, 20\}$ 
5 return  $\alpha_{init}/\beta_j$ 

```

Algorithm 2 reduces the mutation size as the number of classifiers queries increases. The method linearly re-scales the current number of models  $i$  with the assumption of  $N = 10000$ .

---

### Algorithm 3: Adapted Sparse-RS Attack

---

**Input:** objective vector  $F$ , input  $\mathbf{x} \in \mathcal{X}$ , sparsity  $k$ , zero-sampling  $pr_0$ , initial mutation size  $\alpha_{init}$ , budget  $N$

```

1  $M \leftarrow k$  random pixel indices to be perturbed
2  $\Delta \leftarrow$  values of the perturbation to be applied
3  $L \leftarrow F(\mathbf{x}; \{M, \Delta\})$ 
4 for  $i \leftarrow 0; i < N; i + 1$  do
5    $p_m \leftarrow \text{selection}(\alpha_{init})$  // refer to Algorithm 2
6    $M', \Delta' \leftarrow \text{mutation}(\{M, \Delta\}, p_m)$ 
7    $L' \leftarrow F(\mathbf{x}; \{M', \Delta'\})$ 
8   if  $\{M', \Delta'\}$  dominates  $\{M, \Delta\}$  then
9      $M \leftarrow M', \Delta \leftarrow \Delta', L \leftarrow L'$ 
10 return  $\{M, \Delta\}$ 

```

Algorithm 3 outlines the Sparse-RS algorithm proposed by Croce et al. [10] adapted to the multi-objective scenario. *dominates* is corresponds to Definition 3.1.

---

### Algorithm 4: SA-MOO Method

---

**Input:** objective vector  $F$ , input  $\mathbf{x} \in \mathcal{X}$ , query budget  $N$ , sparsity  $k$ , population size  $s$ , zero-sampling probability  $pr_0$

// Initial Population

```

1  $P \leftarrow \{\{M_1, \Delta_1\}, \dots, \{M_s, \Delta_s\}\}$ 
// Objective Evaluation
2  $L \leftarrow \{F(\mathbf{x}; \{M_1, \Delta_1\}), \dots, F(\mathbf{x}; \{M_s, \Delta_s\})\}$ 
3 for  $i \leftarrow 0; i < N; i \leftarrow i + s$  do
// Uniformly Sample  $s/2$  pairs of  $P$  indices
4    $J \leftarrow \mathcal{U}(\{1, \dots, s\})^{\frac{s}{2} \times 2}$ 
5    $P_O \leftarrow \{\}$ 
6    $L_O \leftarrow \{\}$ 
7   for  $j \in J$  do
8      $O \leftarrow \text{crossover}(P_{j_0}, P_{j_1})$ 
9      $M''_1, \Delta''_1 \leftarrow \text{mutation}(O_1)$ 
10     $M''_2, \Delta''_2 \leftarrow \text{mutation}(O_2)$ 
11     $P_O \leftarrow P_O \cup \{\{M''_1, \Delta''_1\}, \{M''_2, \Delta''_2\}\}$ 
12     $L_O \leftarrow L_O \cup \{F(\mathbf{x}; \{M''_1, \Delta''_1\})\}$ 
13     $L_O \leftarrow L_O \cup \{F(\mathbf{x}; \{M''_2, \Delta''_2\})\}$ 
14   $P \leftarrow P \cup P_O$ 
15   $L \leftarrow L \cup L_O$ 
16   $P \leftarrow \text{non-dominated sorting}(P)$ 
17   $P \leftarrow P_{1:s}$  // Select lowest ranked solutions
18   $L \leftarrow L_P$ 
19 return  $P$  // return population of solutions

```

Method	AT <sub>1</sub>				AT <sub>2</sub>			
	ASR	$l_0$	$l_2$	SSIM	ASR	$l_0$	$l_2$	SSIM
SA-MOO*	<b>84.40%</b>	<b>15.02</b>	<b>8.00</b>	<b>0.95</b>	<b>76.90%</b>	<b>15.28</b>	<b>8.54</b>	<b>0.95</b>
SA-MOO**	<b>84.40%</b>	<b>15.02</b>	<b>8.00</b>	<b>0.95</b>	<b>76.90%</b>	<b>15.28</b>	<b>8.54</b>	<b>0.95</b>
SA-MOO*	<b>44.20%</b>	18.39	11.93	<b>0.93</b>	<b>39.10%</b>	<b>18.81</b>	<b>12.90</b>	<b>0.93</b>
SA-MOO**	<b>44.20%</b>	<b>18.37</b>	<b>11.92</b>	<b>0.93</b>	<b>39.10%</b>	18.82	<b>12.90</b>	<b>0.93</b>

Table 4. Statistics of attack success rate, average ssim, and average  $l_0, l_2$  distances of non-targeted (top) and targeted (bottom) attacks on the CIFAR-10 trained models AT<sub>1</sub> and AT<sub>2</sub>. Where "SA-MOO" is the proposed method, \*\* refers to both *crossover* and *mutation* operators, \* refers to only the *mutation operator*.

---

#### Algorithm 5: Crossover Operator

---

**Input:** pixel locations  $M_a, M_b$ , perturbation values  $\Delta_a, \Delta_b$ , crossover size  $p_c$ , sparsity  $k$

```

1  $O \leftarrow \{\}$ 
2 for  $r \in \{a, b\}$  and  $e \in \{b, a\}$  do
3    $U \leftarrow M_e \setminus (M_r \cap M_e)$ 
4    $b \leftarrow \min\{p_c \cdot k, |U|\}$ 
5    $A \leftarrow \mathcal{U}(M_r)^b$ 
6    $B \leftarrow \mathcal{U}(U)^b$ 
7    $M'_r \leftarrow (M_r \setminus A) \cup B$ 
8    $\Delta'_r \leftarrow (\Delta_r \setminus \Delta_{r_A}) \cup \Delta_{\epsilon_B}$ 
9    $O \leftarrow O \cup \{M'_r, \Delta'_r\}$ 
10 return  $O$ 

```

---



---

#### Algorithm 6: Mutation Operator

---

**Input:** pixel locations  $M'_a$ , perturbation values  $\Delta'_a$ , mutation size  $p_m$ , sparsity  $k$ , zero-sample probability  $pr_0$

```

1  $U \leftarrow \{1, \dots, h \cdot w\}$  // image height  $h$  and width  $w$ 
2  $T \leftarrow U \setminus M'_a$ 
3  $A \leftarrow \mathcal{U}(M'_a)^{p_m \cdot k}$ 
4  $B \leftarrow \mathcal{U}(T)^{p_m \cdot k}$ 
5  $M''_a \leftarrow (M'_a \setminus A) \cup B$ 
   // Sample with a zero-probability  $pr_0$ 
6  $\Delta''_a \leftarrow (\Delta'_a \setminus \Delta'_{a_A}) \cup \mathcal{U}(\{-1, 0, 1\})^{(p_m \cdot k) \times 3}$ 
7 return  $\{M''_a, \Delta''_a\}$ 

```

---