

# Supplementary Material for: Boosting Adversarial Attacks with Momentum

Yinpeng Dong<sup>1</sup>, Fangzhou Liao<sup>1</sup>, Tianyu Pang<sup>1</sup>, Hang Su<sup>1</sup>, Jun Zhu<sup>1,\*</sup>, Xiaolin Hu<sup>1</sup>, Jianguo Li<sup>2</sup>

<sup>1</sup> Department of Computer Science and Technology, Tsinghua Lab of Brain and Intelligence

<sup>1</sup> Beijing National Research Center for Information Science and Technology, BNRist Lab

<sup>1</sup> Tsinghua University, 100084 China

<sup>2</sup> Intel Labs China

{dyp17, liaofz13, pty17}@mails.tsinghua.edu.cn, {suhangss, dcszj, xlhu}@mail.tsinghua.edu.cn, jianguo.li@intel.com

In this supplementary material, we provide more results in our experiments. In Sec. A, we report the success rates of non-targeted attacks based on  $L_2$  norm bound. In Sec. B, we provide the results of targeted attacks. The experiments consistently demonstrate the effectiveness of the proposed momentum-based methods.

## A. Non-targeted attacks based on $L_2$ norm bound

We first perform non-targeted attacks based on  $L_2$  norm bound. Since the  $L_2$  distance between an adversarial example and a real example is defined as

$$\|\mathbf{x}^* - \mathbf{x}\|_2 = \sqrt{\sum_{i=1}^N (x_i^* - x_i)^2}, \quad (14)$$

where  $N$  is the dimension of input  $\mathbf{x}$  and  $\mathbf{x}^*$ , the distance measure depends on  $N$ . For example, if the distance of each dimension of an adversarial example and a real example is  $|x_i^* - x_i| = \epsilon$ , the  $L_2$  norm is  $\epsilon\sqrt{N}$  between them while the  $L_\infty$  norm is  $\epsilon$ . Therefore, we set the  $L_2$  norm bound as  $16\sqrt{N}$  in our experiments, where  $N$  is the dimension of the input to a network.

### A.1. Attacking a single model

We include seven networks in this section, which are Inc-v3, Inc-v4, IncRes-v2, Res-152, Inc-v3<sub>ens3</sub>, Inc-v3<sub>ens4</sub> and IncRes-v2<sub>ens</sub>. We generate adversarial examples for Inc-v3, Inc-v4, IncRes-v2 and Res-152 respectively, and measure the success rates of attacks on all models. We compare three attack methods, which are the fast gradient method (FGM, defined in Eq. (2)), iterative FGM (I-FGM) and momentum iterative FGM (MI-FGM, defined in Eq. (10)). We set the

\*Corresponding author.

	Attack	Ensemble	Hold-out
-Inc-v3	FGM	47.3	52.7
	I-FGM	99.1	65.3
	MI-FGM	<b>99.2</b>	<b>89.7</b>
-Inc-v4	FGM	47.2	49.3
	I-FGM	99.3	56.7
	MI-FGM	<b>99.4</b>	<b>88.0</b>
-IncRes-v2	FGSM	47.3	50.4
	I-FGSM	99.4	54.3
	MI-FGSM	<b>99.5</b>	<b>86.1</b>
-Res-152	FGM	47.6	46.6
	I-FGM	99.0	44.7
	MI-FGM	<b>99.5</b>	<b>81.4</b>
-Inc-v3 <sub>ens3</sub>	FGM	51.8	35.4
	I-FGM	<b>99.8</b>	29.5
	MI-FGM	99.6	<b>59.8</b>
-Inc-v3 <sub>ens4</sub>	FGM	51.2	37.5
	I-FGM	99.2	36.4
	MI-FGM	<b>99.7</b>	<b>66.5</b>
-IncRes-v2 <sub>ens</sub>	FGSM	54.4	32.4
	I-FGSM	99.2	19.9
	MI-FGSM	<b>99.8</b>	<b>56.4</b>

Table 4: The success rates (%) of non-targeted adversarial attacks based on  $L_2$  norm bound against an ensemble of white-box models and a hold-out black-box target model. In each row, “-” indicates the name of the hold-out model and the adversarial examples are generated for the ensemble of the other six models.

number of iterations to 10 in I-FGM and MI-FGM, and the decay factor to 1.0 in MI-FGM.

The results are shown in Table 5. We can also see that MI-FGM attacks a white-box model with a near 100%

success rate as I-FGM, and outperforms FGM and I-FGM in black-box attacks significantly. The conclusions are similar to those of  $L_\infty$  norm bound experiments, which consistently demonstrate the effectiveness of the proposed momentum-based iterative methods.

## A.2. Attacking an ensemble of models

In this experiments, we also include Inc-v3, Inc-v4, IncRes-v2, Res-152, Inc-v3<sub>ens3</sub>, Inc-v3<sub>ens4</sub> and IncRes-v2<sub>ens</sub> models for our study. We keep one model as the hold-out black-box model and attack an ensemble of the other six models by FGM, I-FGM and MI-FGM respectively. We set the number of iterations to 20 in I-FGM and MI-FGM, the decay factor to 1.0 in MI-FGM, and the ensemble weights to  $1/6$  equally.

We show the results in Table 4. Iterative methods including I-FGM and MI-FGM can obtain a near 100% success rate for an ensemble of white-box models. And MI-FGM can attack a black-box model with a much higher success rate, showing the good transferability of the adversarial examples generated by MI-FGM. For adversarially trained models, MI-FGM can fool them with about 60% success rates, revealing the great vulnerability of the adversarially trained models against our black-box attacks.

## B. Targeted attacks

### B.1. $L_\infty$ norm bound

Targeted attacks are much more difficult than non-targeted attacks in the black-box manner, since they require the black-box model to output the specific target label. For DNNs trained on a dataset with thousands of output categories such as the ImageNet dataset, finding targeted adversarial examples by only one model to fool a black-box model is impossible. Thus we perform targeted attacks by integrating the ensemble-based approach.

We show the results in Table 6, where the success rate is measured by the percentage of the adversarial examples that are classified as the target label by the model. Similar to the experimental settings in Sec. 4.3, we keep one model to test the performance of black-box attacks, with the targeted adversarial examples generated for the ensemble of the other six models. We set the size of perturbation  $\epsilon$  to 48, decay factor  $\mu$  to 1.0 and the number of iterations to 20 for I-FGSM and MI-FGSM. We can see that one-step FGSM can hardly attack the ensemble of models as well as the target black-box models. The success rates of the adversarial examples generated by MI-FGSM are close to 100% for white-box models and higher than 10% for normally trained black-box models. Unfortunately, it cannot effectively generate targeted adversarial examples to fool adversarially trained models, which remains an open issue for future researches.

### B.2. $L_2$ norm bound

We draw similar conclusions for targeted attacks based on  $L_2$  norm bound. In our experiments, we also include Inc-v3, Inc-v4, IncRes-v2, Res-152, Inc-v3<sub>ens3</sub>, Inc-v3<sub>ens4</sub> and IncRes-v2<sub>ens</sub> models. We keep one model as the hold-out black-box model and attack an ensemble of the other six models with equal ensemble weights by FGM, I-FGM and MI-FGM respectively. We set the maximum perturbation  $\epsilon$  to  $48\sqrt{N}$  where  $N$  is the dimension of inputs, the number of iterations to 20 in I-FGM and MI-FGM, and the decay factor to 1.0 in MI-FGM. We report the success rates of adversarial examples against the white-box ensemble of models and the black-box target model in Table 7. MI-FGM can easily fool white-box models, but it cannot fool the adversarially trained models effectively in the targeted black-box attacks.

	Attack	Inc-v3	Inc-v4	IncRes-v2	Res-152	Inc-v3 <sub>ens3</sub>	Inc-v3 <sub>ens4</sub>	IncRes-v2 <sub>ens</sub>
Inc-v3	FGM	76.2*	41.0	43.1	41.3	34.6	34.9	26.2
	I-FGM	<b>100.0*</b>	39.9	36.4	27.5	17.5	19.2	10.9
	MI-FGM	<b>100.0*</b>	<b>67.6</b>	<b>66.3</b>	<b>56.1</b>	<b>44.4</b>	<b>45.5</b>	<b>33.9</b>
Inc-v4	FGM	47.3	63.1*	37.3	39.0	35.3	33.9	27.7
	I-FGM	52.8	<b>100.0*</b>	42.0	33.5	21.9	19.9	13.8
	MI-FGM	<b>76.9</b>	<b>100.0*</b>	<b>69.6</b>	<b>59.7</b>	<b>51.2</b>	<b>51.0</b>	<b>39.4</b>
IncRes-v2	FGM	48.2	38.9	60.4*	39.8	36.6	35.5	30.5
	I-FGM	56.0	47.5	<b>99.6*</b>	36.9	27.5	22.9	18.7
	MI-FGM	<b>81.7</b>	<b>75.8</b>	<b>99.6*</b>	<b>66.9</b>	<b>62.7</b>	<b>57.7</b>	<b>58.8</b>
-Res-152	FGM	50.8	40.7	42.0	75.1*	36.5	36.0	31.6
	I-FGM	47.6	43.9	43.9	99.4*	32.7	32.3	25.2
	MI-FGM	<b>71.3</b>	<b>65.5</b>	<b>64.3</b>	<b>99.6*</b>	<b>56.7</b>	<b>55.4</b>	<b>51.5</b>

Table 5: The success rates (%) of non-targeted adversarial attacks based on  $L_2$  norm bound against all models. The adversarial examples are crafted for Inc-v3, Inc-v4, IncRes-v2 and Res-152 respectively using FGM, I-FGM and MI-FGM. \* indicates the white-box attacks.

	Attack	Ensemble	Hold-out
-Inc-v3	FGSM	0.5	0.5
	I-FGSM	<b>99.6</b>	9.0
	MI-FGSM	99.5	<b>17.6</b>
-Inc-v4	FGSM	0.3	0.4
	I-FGSM	<b>99.9</b>	7.0
	MI-FGSM	99.8	<b>15.6</b>
-IncRes-v2	FGSM	0.4	0.2
	I-FGSM	<b>99.9</b>	7.3
	MI-FGSM	99.8	<b>16.1</b>
-Res-152	FGSM	0.1	0.5
	I-FGSM	<b>99.6</b>	3.3
	MI-FGSM	99.5	<b>11.4</b>
-Inc-v3 <sub>ens3</sub>	FGSM	0.3	0.1
	I-FGSM	<b>99.7</b>	0.1
	MI-FGSM	<b>99.7</b>	<b>0.5</b>
-Inc-v3 <sub>ens4</sub>	FGSM	0.2	0.1
	I-FGSM	<b>99.9</b>	0.4
	MI-FGSM	99.8	<b>0.9</b>
-IncRes-v2 <sub>ens</sub>	FGSM	0.5	0.1
	I-FGSM	99.7	0.1
	MI-FGSM	<b>99.8</b>	<b>0.2</b>

Table 6: The success rates (%) of targeted adversarial attacks based on  $L_\infty$  norm bound against an ensemble of white-box models and a hold-out black-box target model. In each row, “-” indicates the name of the hold-out model and the adversarial examples are generated for the ensemble of the other six models.

	Attack	Ensemble	Hold-out
-Inc-v3	FGM	0.7	0.4
	I-FGM	<b>99.7</b>	17.8
	MI-FGM	99.5	<b>21.0</b>
-Inc-v4	FGM	0.7	0.5
	I-FGM	<b>99.9</b>	15.2
	MI-FGM	99.8	<b>21.8</b>
-IncRes-v2	FGM	0.7	0.7
	I-FGM	99.8	16.4
	MI-FGM	<b>99.9</b>	<b>21.7</b>
-Res-152	FGM	0.5	0.4
	I-FGM	99.5	9.2
	MI-FGM	<b>99.6</b>	<b>17.4</b>
-Inc-v3 <sub>ens3</sub>	FGM	0.6	0.2
	I-FGM	<b>99.9</b>	0.7
	MI-FGM	99.6	<b>1.6</b>
-Inc-v3 <sub>ens4</sub>	FGM	0.5	0.2
	I-FGM	99.7	1.7
	MI-FGM	<b>100.0</b>	<b>2.0</b>
-IncRes-v2 <sub>ens</sub>	FGM	0.6	0.4
	I-FGM	99.6	0.5
	MI-FGM	<b>99.8</b>	<b>1.9</b>

Table 7: The success rates (%) of targeted adversarial attacks based on  $L_2$  norm bound against an ensemble of white-box models and a hold-out black-box target model. In each row, “-” indicates the name of the hold-out model and the adversarial examples are generated for the ensemble of the other six models.