

制約付き非凸最適化問題のための 確率的Frank-Wolfe法と その敵対的攻撃への応用

日本オペレーションズ・リサーチ学会 2025年春季研究発表会

3/6(木) 11:00~11:20

明治大学 佐藤尚樹

明治大学 飯塚秀明

背景：勾配法

$$\boldsymbol{\theta}_{t+1} := \boldsymbol{\theta}_t + \gamma_t \mathbf{d}_t$$

学習率

探索方向

▷ 最急降下法 (Gradient Descent)

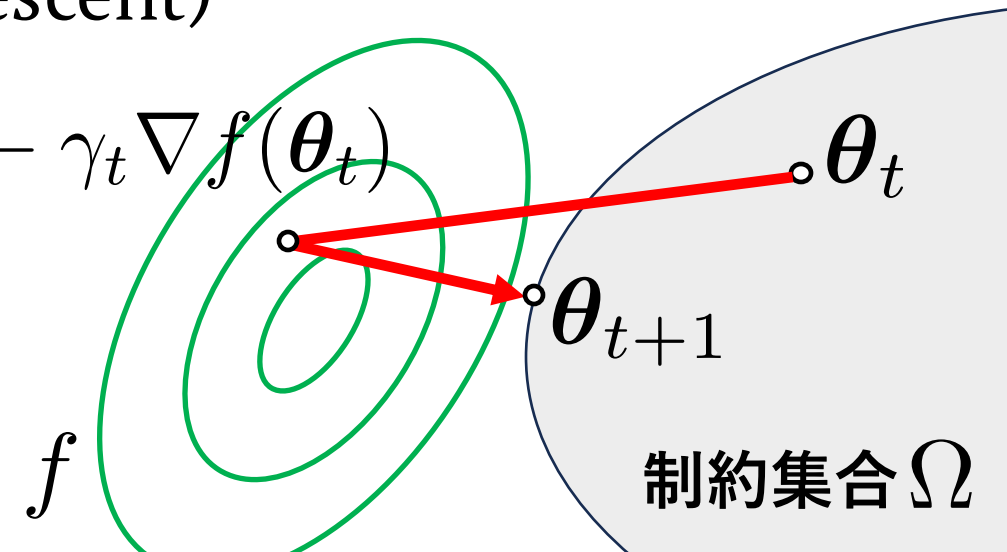
$$\mathbf{d}_t := -\nabla f(\boldsymbol{\theta}_t)$$

全勾配

▷ 射影勾配降下法 (Projected Gradient Descent)

$$\therefore \boldsymbol{\theta}_{t+1} := \Pi_{\Omega} (\boldsymbol{\theta}_t - \gamma_t \nabla f(\boldsymbol{\theta}_t))$$

$$\boldsymbol{\theta}_t - \gamma_t \nabla f(\boldsymbol{\theta}_t)$$



背景：勾配法

$$\boldsymbol{\theta}_{t+1} := \boldsymbol{\theta}_t + \gamma_t \mathbf{d}_t$$

學習率

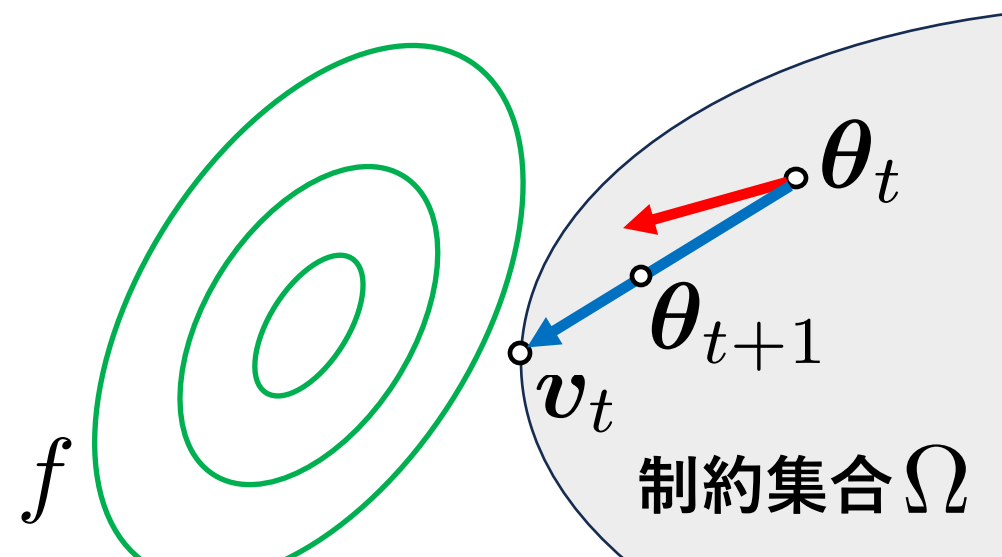
探索方向

▷ Frank-Wolfe法

$$\mathbf{v}_t := \operatorname{argmax}_{\mathbf{v} \in \Omega} \langle \mathbf{v}, \underbrace{-\nabla f(\boldsymbol{\theta}_t)}_{\text{最急降下方向}} \rangle$$

$$\mathbf{d}_t := \mathbf{v}_t - \boldsymbol{\theta}_t$$

$$\therefore \boldsymbol{\theta}_{t+1} := \gamma_t \mathbf{v}_t + (1 - \gamma_t) \boldsymbol{\theta}_t$$



背景：Stochastic Frank Wolfe (SFW)

Algorithm 1 Frank Wolfe

Require: $(\gamma_t)_{t \in \mathbb{N}} \subset \mathbb{R}_{++}$

$t \leftarrow 0, \boldsymbol{\theta}_0 \in \Omega$

loop

$\boldsymbol{v}_t = \operatorname{argmax}_{\boldsymbol{v} \in \Omega} \langle \boldsymbol{v}, -\nabla f(\boldsymbol{\theta}_t) \rangle$

$\boldsymbol{d}_t = \boldsymbol{v}_t - \boldsymbol{\theta}_t$ 全勾配

$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \gamma_t \boldsymbol{d}_t$

$t \leftarrow t + 1$

end loop

Algorithm 2 Stochastic Frank Wolfe (SFW)

Require: $(\gamma_t)_{t \in \mathbb{N}} \subset \mathbb{R}_{++}$

$t \leftarrow 0, \boldsymbol{\theta}_0 \in \Omega$

loop

$\boldsymbol{v}_t = \operatorname{argmax}_{\boldsymbol{v} \in \Omega} \langle \boldsymbol{v}, -\nabla f_{\mathcal{S}_t}(\boldsymbol{\theta}_t) \rangle$

$\boldsymbol{d}_t = \boldsymbol{v}_t - \boldsymbol{\theta}_t$ ミニバッチ確率的勾配

$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \gamma_t \boldsymbol{d}_t$

$t \leftarrow t + 1$

end loop

背景：SFWの活躍(応用)

▷ 深層学習におけるSFWの活躍の例

Pruning Ratios	50%	70%	80%	90%	95%
SFW-Pruning + SFW-Init (ours)	93.10	93.10	93.10	93.10	92.00
One-Cycle Pruning (Hubens et al., 2021)	-	-	90.87	90.72	90.67
Early Bird (You et al., 2020)	93.21	92.80	-	-	-
OTO (Chen et al., 2021)	90.35	90.35	90.35	90.35	90.35
DPF (Lin et al., 2020)	-	-	-	-	93.87
Group MDP (Deleu & Bengio, 2021)	-	-	-	89.38	-

Table 1: Comparisons to more state-of-the-art methods on VGG-16 and CIFAR-10.

▷ モデルの枝刈りタスクにおいて、SGDを利用した手法よりも優れた性能をもたらすことが示された。しかも計算量も少ない。

[Miao+22] L. Miao, X. Luo, T. Chen, W. Chen, D. Liu, and Z. Wang, “Learning pruning-friendly networks via Frank-Wolfe: One-shot, any-sparsity, and no retraining”, In Proceedings of the 10th International Conference on Learning Representations, 2022.

背景：SFWの先行研究(理論)

Theorem 2. Consider the stochastic setting of (1) where f is G -Lipschitz and F is L -smooth. Then, the output x_a of Algorithm 2 with parameters $\gamma_t = \gamma = \sqrt{\frac{2(F(x_0) - F(x^*))}{TLD^2\beta}}$, $b_t = b = T$ for all $t \in \{0, \dots, T-1\}$, satisfies the following bound:

学習率

バッチサイズ

$$\mathbb{E}[\mathcal{G}(x_a)] \leq \frac{D}{\sqrt{T}} \left(G + \sqrt{\frac{2L(F(x_0) - F(x^*))}{\beta}} (1 + \beta) \right),$$

where x^* is an optimal solution to (stochastic) problem (1).

[Reddi+16] S.J. Reddi, S. Sra, B. Póczos, and A. Smola, “Stochastic Frank-Wolfe Methods for Nonconvex Optimization” in *54th Annual Allerton Conference on Communication, Control, and Computing*, pp.1214-1251, 2016.

▷ 問題点その1. 学習率とバッチサイズの設定が現実的でない。

背景：SFWの先行研究(理論)

目的関数のL-平滑性から、

$$\frac{1}{T} \sum_{t=0}^{T-1} \text{gap} \leq \frac{A}{T\gamma} + B\gamma$$

が成り立つ。ただし、AとBは正の定数とする。

このとき $\gamma = \frac{1}{\sqrt{T}}$ を使えば、

$$\frac{1}{T} \sum_{t=0}^{T-1} \text{gap} \leq \frac{A}{\sqrt{T}} + \frac{B}{\sqrt{T}} = \mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$$

が成り立ち、**あたかもgapが0に収束するかのように見える。**

▷ 問題点その2. 実は**収束は保証されていない。**

背景：SFWの先行研究(理論)

Algorithm	Learning Rate	Gradient (batch size)	Momentum
SFW (Reddi et al., 2016)	$\gamma_t = \mathcal{O}\left(\sqrt{\frac{1}{TL}}\right)$	mini-batch ($b = T$)	No
AsySFW (Gu et al., 2019)	$\gamma_t = \mathcal{O}\left(\sqrt{\frac{1}{TL}}\right)$	mini-batch ($b = T$)	No
FW-SD (Grigas et al., 2019)	$\gamma_t = \mathcal{O}\left(\sqrt{\frac{1}{L}}\right)$	mini-batch ($b_t \leq n$)	No
SFW (Négiar et al., 2020)	$\gamma_t = \frac{2}{t+2}$	mini-batch ($b \leq n$)	No
SCG (Mokhtari et al., 2020)	$\gamma_t = \frac{1}{T}$	Noise	Yes (Increasing)
AdaSFW (Combettes et al., 2021)	$\gamma_t = \frac{2}{t+2}$	mini-batch $\left(b_t = \mathcal{O}\left(\frac{t^2}{L^2}\right)\right)$	No
AdaCSFW (Combettes et al., 2021)	$\gamma_t = \frac{2}{t+2}$	mini-batch ($b \leq n$)	No
any SFW methods (Nazykov et al., 2024)	$\gamma_t = \sqrt{\frac{1}{T}}$	Noise	No

動機

- ▷現実的な学習率とバッチサイズの設定の下での収束解析を提供したい.

貢献

▷ 実用的な定数/**減少**学習率, 定数/**増加**バッチサイズを有するSFWの収束解析を提供した.

	定数バッチサイズ $b_t := b$	増加バッチサイズ
定数学習率 $\gamma_t := \gamma$	$\mathcal{O}\left(\frac{1}{T} + \frac{1}{\sqrt{b}} + \gamma\right)$ (Theorem 3.1)	$\mathcal{O}\left(\frac{1}{T} + \gamma\right)$ (Theorem 3.2)
減少 (I) $\gamma_t := \frac{1}{t+1}$	$\mathcal{O}\left(\frac{1}{\sqrt{b}}\right)$ (Theorem 3.3)	-
減少 (II) $\gamma_t := \frac{1}{(t+1)^a} (a \in [0.5, 1))$	$\mathcal{O}\left(\frac{1}{T^{\min\{1-a, a\}}} + \frac{1}{\sqrt{b}}\right)$ (Theorem 3.4(i))	$\mathcal{O}\left(\frac{1}{T^{\min\{1-a, a\}}}\right)$ (Theorem 3.5(i))
減少 (III)	$\mathcal{O}\left(\frac{1}{T} + \frac{1}{\sqrt{b}}\right)$ (Theorem 3.4(ii))	$\mathcal{O}\left(\frac{1}{T}\right)$ (Theorem 3.5(ii))

$$\gamma_t := (\underbrace{\gamma, \gamma, \dots, \gamma}_K, \underbrace{\eta\gamma, \eta\gamma, \dots, \eta\gamma}_K, \dots, \underbrace{\eta^{P-1}\gamma, \eta^{P-1}\gamma, \dots, \eta^{P-1}\gamma}_K), \gamma > 0, \eta \in (0, 1), KP = T$$

$$b_t := (\underbrace{b, b, \dots, b}_E, \underbrace{\lambda b, \lambda b, \dots, \lambda b}_E, \dots, \underbrace{\lambda^{Q-1}b, \lambda^{Q-1}b, \dots, \lambda^{Q-1}b}_E), b > 0, \lambda > 1, EQ = T$$

準備：最適化問題

経験損失最小化問題

▷ 訓練データセット $S := (z_1, z_2, \dots, z_n)$

▷ Deep Neural Network のパラメータ $x \in \mathbb{R}^d$

$$\min_{\theta \in \Omega} \left\{ f(\theta) := \frac{1}{n} \sum_{i=0}^{n-1} \underline{f_i(\theta)} \right\}$$

i 番目の訓練データ z_i に対する損失関数

▷ 非凸目的関数 $f: \mathbb{R}^d \rightarrow \mathbb{R}$ を最適化する。

準備：仮定

仮定

(A1) 制約集合 $\Omega \subset \mathbb{R}^d$ は凸集合で、直径 $D \in \mathbb{R}$ を有する, i.e.,

$$\forall x, y \in \Omega: \|x - y\| \leq D$$

(A2) $(\theta_t)_{t \in \mathbb{N}} \subset \mathbb{R}^d$ を最適化手法によって生成された点列とするとき,
(i) 任意の $t \in \mathbb{N}$ に対して、次の式が成り立つとする.

$$\mathbb{E}_{\xi_t} [G_{\xi_t}(\theta_t)] = \nabla f(\theta_t)$$

(ii) 次の式を満たす非負定数 σ^2 が存在するとする.

$$\mathbb{E}_{\xi_t} [\|G_{\xi_t}(\theta_t) - \nabla f(\theta_t)\|^2] \leq \sigma^2$$

準備：仮定

仮定続き

(A3) 関数 $f: \mathbb{R}^d \rightarrow \mathbb{R}$ は制約集合 Ω 上で連続的微分可能で,
 L -平滑とする, i.e.,

$$\forall x, y \in \Omega: \|\nabla f(x) - \nabla f(y)\| \leq \|x - y\|$$

(A4) 時刻 $t \in \mathbb{N}$ で, 全勾配 ∇f はミニバッチ \mathcal{S}_t で次のように近似
されとする.

$$\nabla f_{\mathcal{S}_t}(\boldsymbol{\theta}_t) := \frac{1}{b} \sum_{i \in [b]} \mathbf{G}_{\xi_{t,i}}(\boldsymbol{\theta}_t) = \frac{1}{b} \sum_{\{i: \mathbf{z}_i \in \mathcal{S}_t\}} \nabla l_i(\boldsymbol{\theta}_t)$$

ミニバッチ
確率的勾配

確率的勾配

準備：Frank-Wolfe gap

- ▶ 制約付き最適化においては、 $\|\nabla f(\theta_t)\|$ を収束の指標として使うことはできない。
- ▶ Frank-Wolfe法の解析では、収束の指標として伝統的に **Frank-Wolfe gap** $\mathcal{G}(\theta_t)$ が利用されてきた。

$$\mathcal{G}(\theta_t) := \max_{v \in \Omega} \langle v - \theta, -\nabla f(\theta) \rangle$$

- ▶ $\theta^* \in \mathbb{R}^d$ を Ω 上での局所的最適解であるとするとき、
 $\forall v \in \Omega: \langle v - \theta^*, \nabla f(\theta^*) \rangle \geq 0$ i.e., $\mathcal{G}(\theta^*) \leq 0$
が成り立つことから、gapを小さくすることを目指す。

収束解析一定数バッチサイズ $b_t := b$

▷ **定数学習率** $\gamma_t := \gamma$

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} [\mathcal{G}(\boldsymbol{\theta}_t)] \leq \frac{\mathbb{E} [f(\boldsymbol{\theta}_0)] - \mathbb{E} [f(\boldsymbol{\theta}^*)]}{\gamma T} + \frac{D\sigma}{\sqrt{b}} + \frac{LD^2\gamma}{2} = \mathcal{O} \left(\frac{1}{T} + \frac{1}{\sqrt{b}} + \gamma \right)$$

▷ **減少学習率** $\gamma_t := \frac{1}{(t+1)^a} (a \in [0.5, 1))$

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} [\mathcal{G}(\boldsymbol{\theta}_t)] \leq \frac{2 \max\{\bar{f}, |f(\boldsymbol{\theta}^*)|\}}{T^{1-a}} + \frac{D\sigma}{\sqrt{b}} + \frac{LD^2}{2(1-a)T^a} = \mathcal{O} \left(\frac{1}{T^{\min\{1-a, a\}}} + \frac{1}{\sqrt{b}} \right)$$

▷ **減少学習率** $\gamma_t := (\underbrace{\gamma, \gamma, \dots, \gamma}_K, \underbrace{\eta\gamma, \eta\gamma, \dots, \eta\gamma}_K, \dots, \underbrace{\eta^{P-1}\gamma, \eta^{P-1}\gamma, \dots, \eta^{P-1}\gamma}_K), \gamma > 0, \eta \in (0, 1), KP = T$

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} [\mathcal{G}(\boldsymbol{\theta}_t)] \leq \frac{2 \max\{\bar{f}, |f(\boldsymbol{\theta}^*)|\}}{T\underline{\gamma}} + \frac{D\sigma}{\sqrt{b}} + \frac{\gamma_\infty LD^2}{2T} = \mathcal{O} \left(\frac{1}{T} + \frac{1}{\sqrt{b}} \right)$$

収束解析—増加バッチサイズ

$$b_t := (\underbrace{b, b, \dots, b}_E, \underbrace{\lambda b, \lambda b, \dots, \lambda b}_E, \dots, \underbrace{\lambda^{Q-1}b, \lambda^{Q-1}b, \dots, \lambda^{Q-1}b}_E), b > 0, \lambda > 1, EQ = T$$

▷ **定数学習率** $\gamma_t := \gamma$

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} [\mathcal{G}(\boldsymbol{\theta}_t)] \leq \frac{\mathbb{E} [f(\boldsymbol{\theta}_0)] - \mathbb{E} [f(\boldsymbol{\theta}^*)]}{\gamma T} + \frac{D\sigma}{T} \frac{E\sqrt{\lambda}}{\sqrt{b}(\sqrt{\lambda} - 1)} + \frac{LD^2\gamma}{2} = \mathcal{O} \left(\frac{1}{T} + \gamma \right)$$

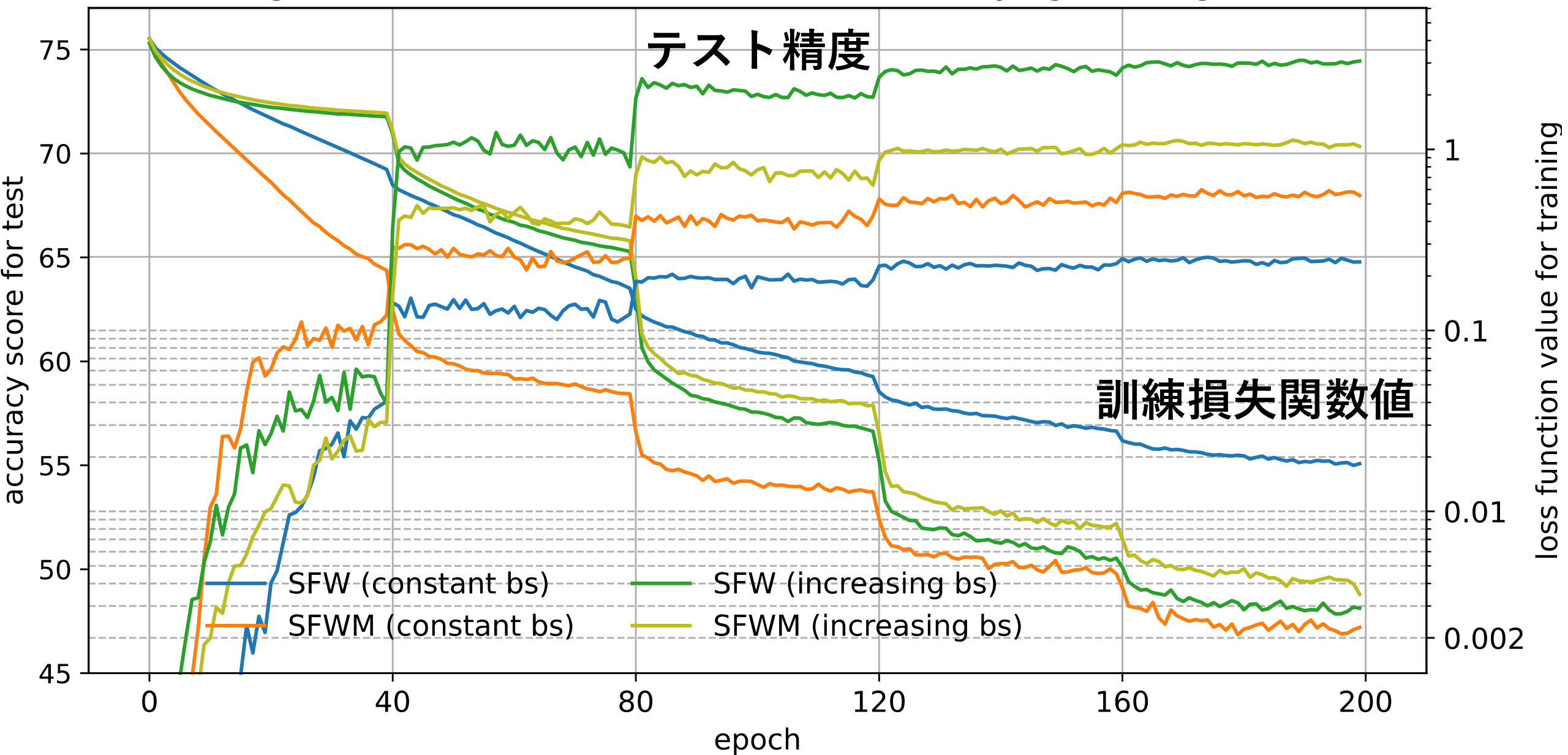
▷ **減少学習率** $\gamma_t := \frac{1}{(t+1)^a} (a \in [0.5, 1))$

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} [\mathcal{G}(\boldsymbol{\theta}_t)] \leq \frac{2 \max\{\bar{f}, |f(\boldsymbol{\theta}^*)|\}}{T^{1-a}} + \frac{D\sigma}{T} \frac{E\sqrt{\lambda}}{\sqrt{b}(\sqrt{\lambda} - 1)} + \frac{LD^2}{2(1-a)T^a} = \mathcal{O} \left(\frac{1}{T^{\min\{1-a, a\}}} \right)$$

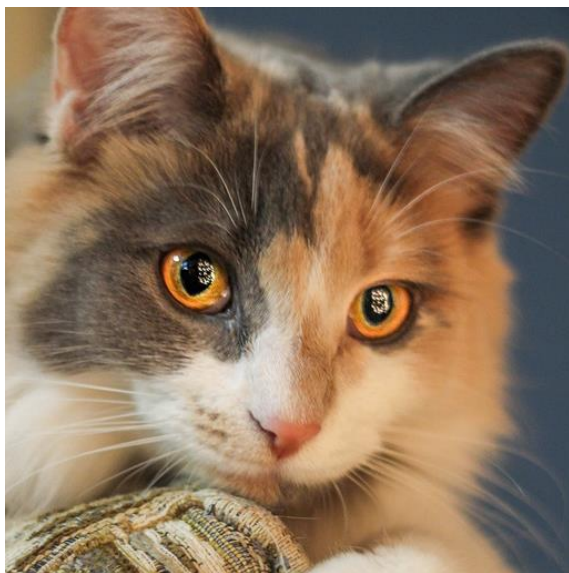
▷ **減少学習率** $\gamma_t := (\underbrace{\gamma, \gamma, \dots, \gamma}_K, \underbrace{\eta\gamma, \eta\gamma, \dots, \eta\gamma}_K, \dots, \underbrace{\eta^{P-1}\gamma, \eta^{P-1}\gamma, \dots, \eta^{P-1}\gamma}_K), \gamma > 0, \eta \in (0, 1), KP = T$

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} [\mathcal{G}(\boldsymbol{\theta}_t)] \leq \frac{2 \max\{\bar{f}, |f(\boldsymbol{\theta}^*)|\}}{T\underline{\gamma}} + \frac{D\sigma}{T} \frac{E\sqrt{\lambda}}{\sqrt{b}(\sqrt{\lambda} - 1)} + \frac{\gamma_\infty LD^2}{2T} = \mathcal{O} \left(\frac{1}{T} \right)$$

Training ResNet18 on CIFAR100 dataset with decaying learning rate (III)



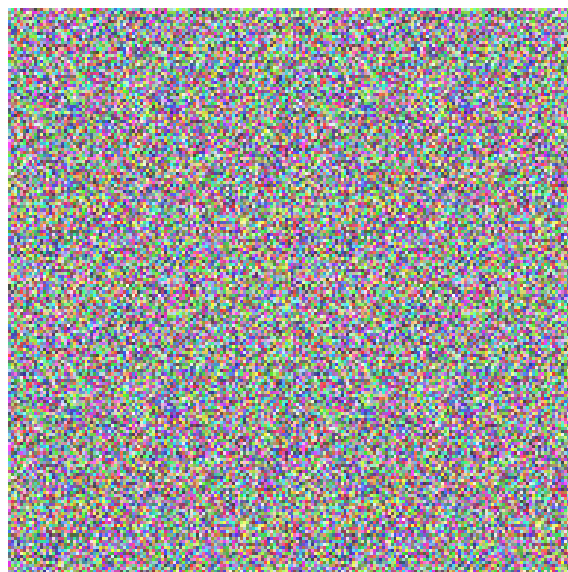
敵対的攻撃 = 敵対的サンプルを作る



“cat”

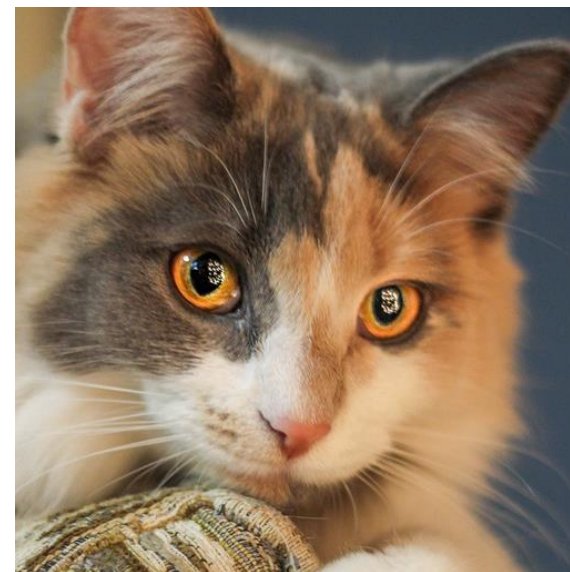
x

+ 0.007 ×



noise

=



“dog”

Adversarial Examples

x_{adv}

ロバストモデル

- ▷ **ロバストモデル**とは、敵対的攻撃に備えて、防御のための特別な細工を施された訓練済みモデルのこと。
- ▷ いくつかの防御手法の中で、ロバストモデルを構成するための最も標準的な手法は**敵対的訓練**。
- ▷ **敵対的訓練**とは、モデルの訓練時に、訓練データに敵対的サンプルを混ぜる手法。
- ▷ 敵対的攻撃を研究するモチベーションは、
 - ▷ **敵対的訓練**に敵対的サンプルが必要だから。
 - ▷ **ロバストモデルの性能評価**に敵対的サンプルが必要だから。

敵対的攻撃は制約付き非凸最適化問題

経験損失関数

正解ラベル

何らかの
距離関数

許容できる
ノイズの大きさ

$$\text{maximize } f(g(\mathbf{x}_{\text{adv}}), \mathbf{y}) \text{ s.t. } d(\mathbf{x}_{\text{adv}}, \mathbf{x}) \leq \epsilon$$

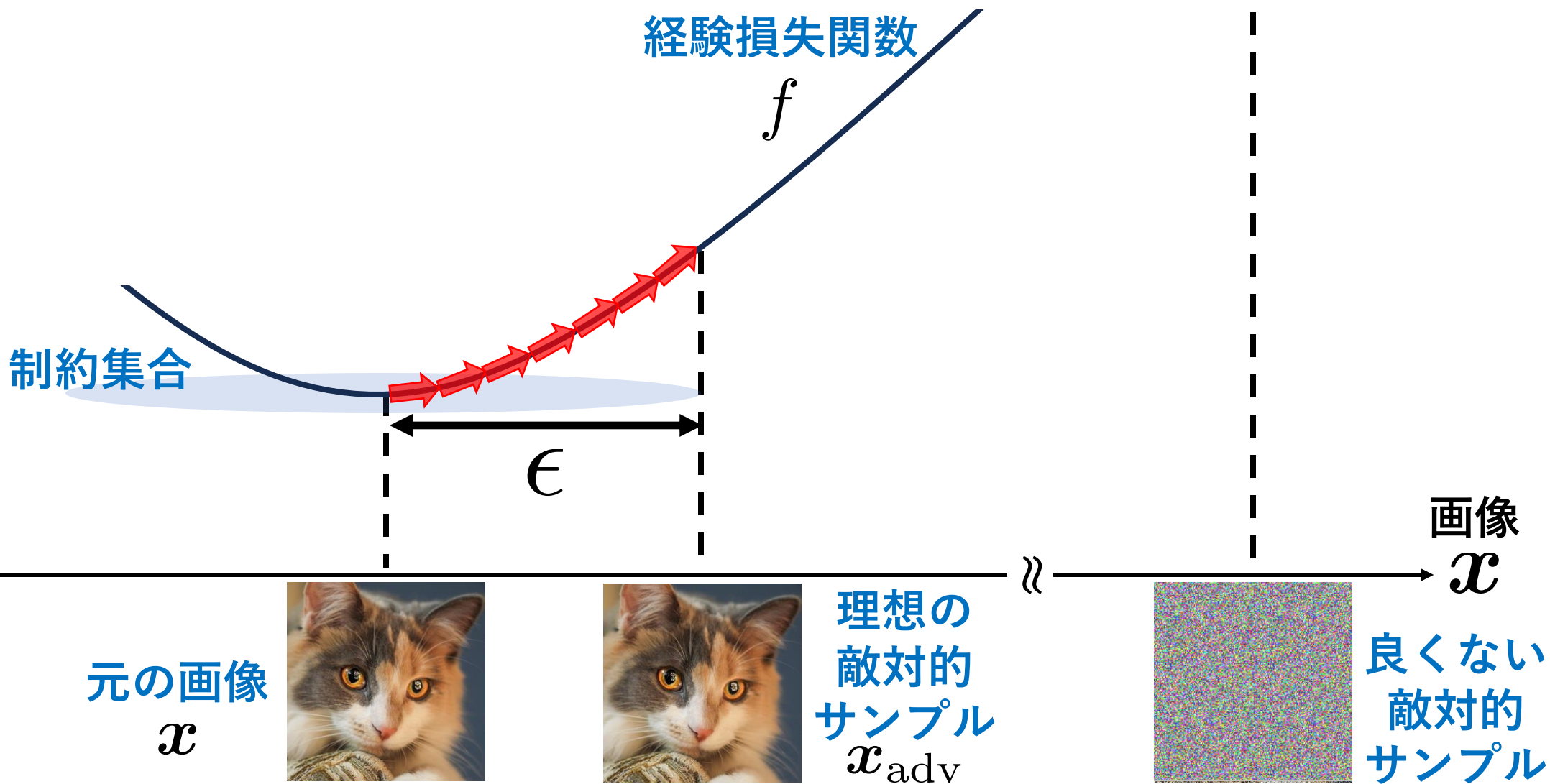
分類モデル(DNN)
の出力

攻撃後の画像

元の画像

- ▶ 攻撃後の画像と元の画像の差異を最大 ϵ に抑えながら、なるべく損失関数値を悪化させ、分類モデルの出力を歪ませる。

敵対的攻撃



Algorithm 1 APGD

```
1: Input:  $f, S, x^{(0)}, \eta, N_{\text{iter}}, W = \{w_0, \dots, w_n\}$ 
2: Output:  $x_{\text{max}}, f_{\text{max}}$ 
3:  $x^{(1)} \leftarrow P_S (x^{(0)} + \eta \nabla f(x^{(0)}))$ 
4:  $f_{\text{max}} \leftarrow \max\{f(x^{(0)}), f(x^{(1)})\}$ 
5:  $x_{\text{max}} \leftarrow x^{(0)}$  if  $f_{\text{max}} \equiv f(x^{(0)})$  else  $x_{\text{max}} \leftarrow x^{(1)}$ 
6: for  $k = 1$  to  $N_{\text{iter}} - 1$  do
7:    $z^{(k+1)} \leftarrow P_S (x^{(k)} + \eta \nabla f(x^{(k)}))$ 
8:    $x^{(k+1)} \leftarrow P_S (x^{(k)} + \alpha(z^{(k+1)} - x^{(k)})$ 
       $\quad \quad \quad + (1 - \alpha)(x^{(k)} - x^{(k-1)}))$ 
9:   if  $f(x^{(k+1)}) > f_{\text{max}}$  then
10:      $x_{\text{max}} \leftarrow x^{(k+1)}$  and  $f_{\text{max}} \leftarrow f(x^{(k+1)})$ 
11:   end if
12:   if  $k \in W$  then
13:     if Condition 1 or Condition 2 then
14:        $\eta \leftarrow \eta/2$  and  $x^{(k+1)} \leftarrow x_{\text{max}}$ 
15:     end if
16:   end if
17: end for
```

Auto-PGD

PGDによる元の画像の更新

ある細工

Algorithm 4 Auto-Frank Wolfe

Require: $f, \Omega, \mathbf{x}_0, \gamma_0, N_{iter}, W = \{w_0, \dots, w_n\}$

$\mathbf{x}_{adv} \leftarrow \mathbf{x}_0, \mathbf{m}_{-1} \leftarrow \mathbf{0}$

for $t = 0$ to $N_{iter} - 1$ **do**

$\mathbf{v}_t = \underset{\mathbf{v} \in \Omega}{\operatorname{argmax}} \langle \mathbf{v}, \nabla f(\mathbf{x}_t) \rangle$

$\mathbf{d}_t = \mathbf{v}_t - \mathbf{x}_t$

$\mathbf{x}_{t+1} = \mathbf{x}_t + \gamma_t \mathbf{d}_t$

if $f(\mathbf{x}_{t+1}) > f(\mathbf{x}_{adv})$ **then**

$\mathbf{x}_{adv} \leftarrow \mathbf{x}_{t+1}$

end if

if $t \in W$ **then**

if Condition (i) or (ii) is satisfied **then**

$\gamma_{t+1} \leftarrow \gamma_t / 2$

$\mathbf{x}_{t+1} \leftarrow \mathbf{x}_{adv}$

end if

end if

end for

Auto-FW

FWによる元の画像の更新

現在最も関数値を悪化させる画像を
敵対的サンプルに設定

攻撃がうまくいっていないならば、
学習率を半分にし、画像は
良かったところまで引き戻す。

数值实验(CIFAR100, L_∞ 距離, $\epsilon = 8/255$)

paper	Architecture	clean accuracy	APGD	AFW
(Wang et al., 2023)	WideResNet-70-16	75.22	48.11	48.16
(Wang et al., 2023)	WideResNet-28-10	72.58	44.05	44.12
(Debenedetti et al., 2023)	XCiT-L12	70.76	38.97	39.02
(Rebuffi et al., 2021)	WideResNet-70-16	63.56	38.28	38.31
(Debenedetti et al., 2023)	XCiT-M12	69.21	38.69	38.77
(Pang et al., 2022)	WideResNet-70-16	65.56	36.66	36.71
(Debenedetti et al., 2023)	XCiT-S12	67.34	37.08	37.06
(Rebuffi et al., 2021)	WideResNet-28-16	62.41	35.73	35.8
(Jia et al., 2022)	WideResNet-34-20	67.31	36.46	36.62
(Addepalli et al., 2022a)	WideResNet-34-10	68.75	36.84	36.92
(Cui et al., 2021)	WideResNet-34-10	62.97	37.05	37.17
(Schwag et al., 2022)	WideResNet-34-10	65.93	35.77	35.89
(Pang et al., 2022)	WideResNet-28-10	63.66	35.25	35.26
(Jia et al., 2022)	WideResNet-34-10	64.89	35.28	35.53
(Addepalli et al., 2022b)	WideResNet-34-10	65.73	35.64	36.07
(Cui et al., 2021)	WideResNet-34-20	62.55	34.08	34.27
(Cui et al., 2021)	WideResNet-34-10	60.64	33.99	34.10

まとめ

- ▷ SFWの現実的な学習率とバッチサイズの設定の下での収束解析を提供した.
- ▷ 特に, Frank-Wolfe gapが0に収束するのは, 減少学習率と増加バッチサイズを同時に使うときだけであることを示した.
- ▷ 新しい敵対的攻撃手法Auto-FWを提案した.
- ▷ 最新のロバストモデルへの攻撃においてAuto-FWとAuto-PGDの性能を比較し, その限界を示した.