

# **МИНИМИЗАЦИЯ С УЧЕТОМ РЕЗКОСТИ**

Панеш Али, 193

Что обсуждаем?

Минимизация с учетом резкости (SAM)

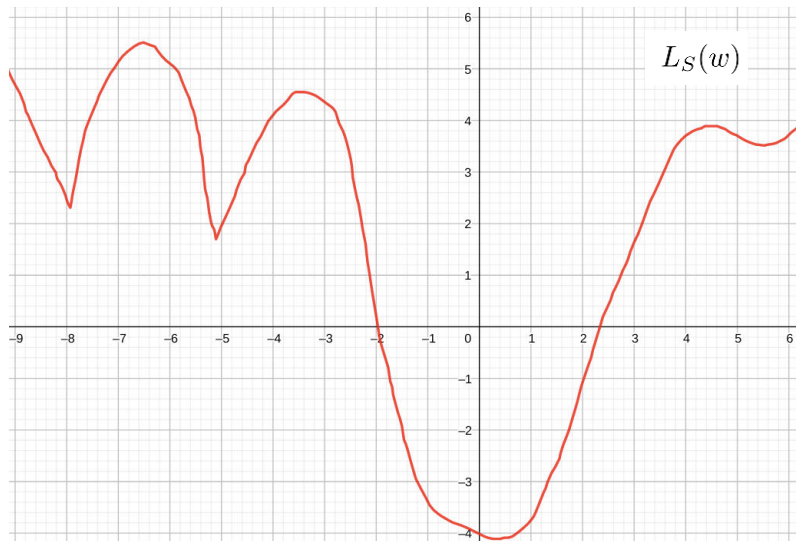
# Предпосылки

- Остренькие функции потерь
- Много локальных минимумов



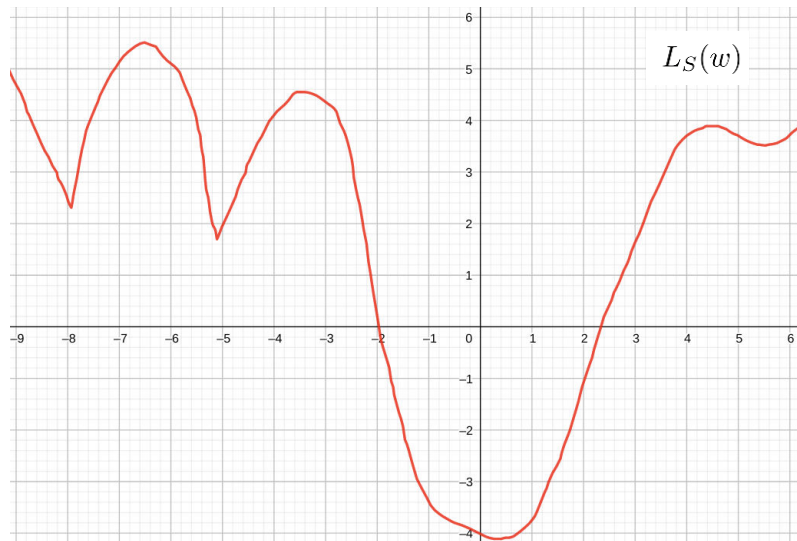
# В чем суть?

Функция потерь на  
тренировочных данных

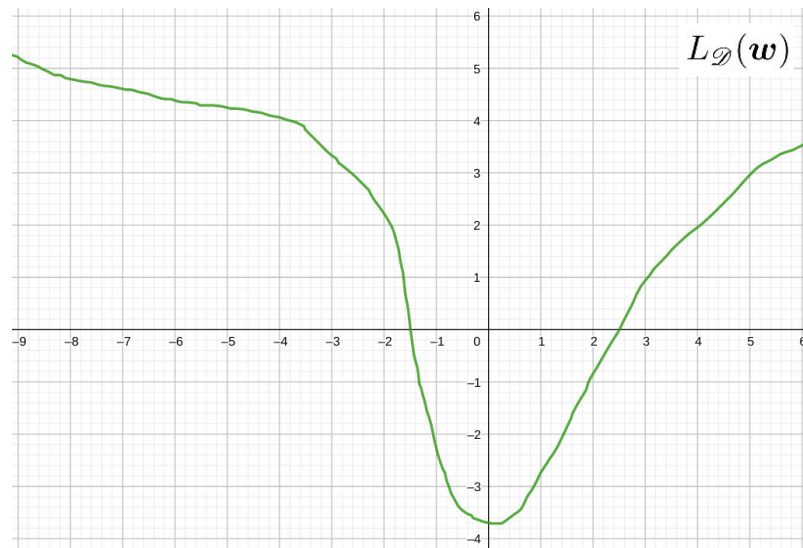


# В чем суть?

Функция потерь на  
тренировочных данных

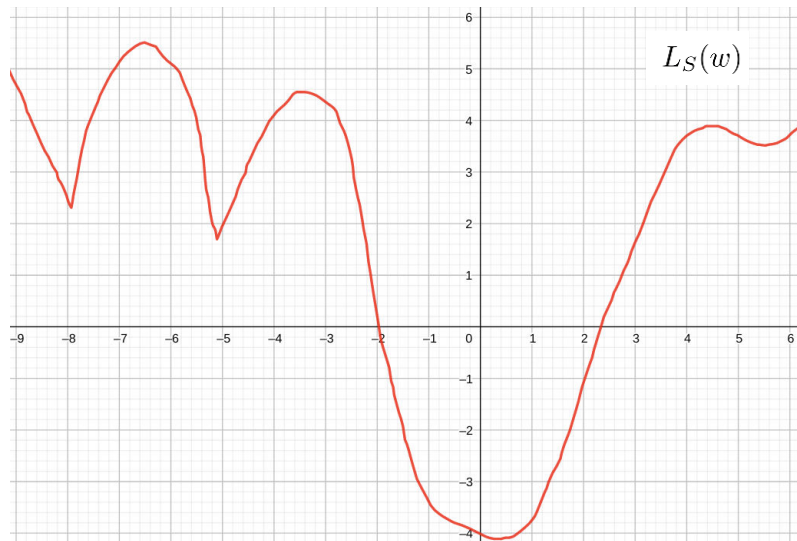


Реальная функция потерь

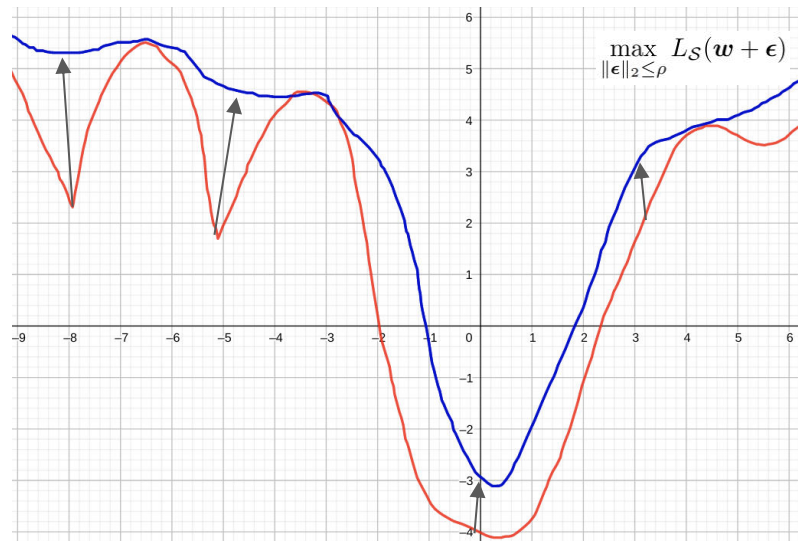


# В чем суть?

Функция потерь на  
тренировочных данных

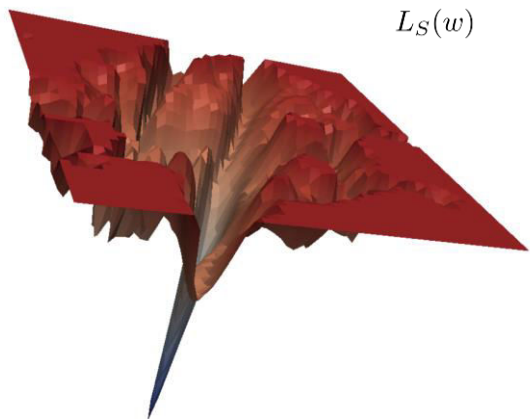


Наша оценка функции потерь

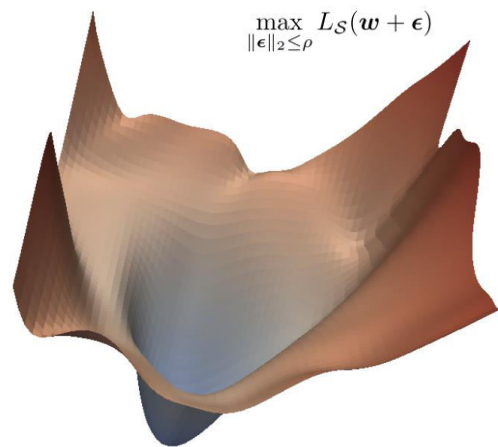


# В чем суть?

Функция потерь на  
тренировочных данных



Наша оценка функции



# Как работает обучение

- Спойлер: ничего нового
- Просто новый подсчет “градиента”



# Как работает обучение

- Спойлер: ничего нового
- Просто новый подсчет “градиента”
- Считаем оценку на эпсилон  $\hat{\epsilon}(\mathbf{w})$

$$\hat{\epsilon}(\mathbf{w}) = \rho \operatorname{sign}(\nabla_{\mathbf{w}} L_{\mathcal{S}}(\mathbf{w})) |\nabla_{\mathbf{w}} L_{\mathcal{S}}(\mathbf{w})|^{q-1} / \left( \|\nabla_{\mathbf{w}} L_{\mathcal{S}}(\mathbf{w})\|_q^q \right)^{1/p}$$

# Как работает обучение

- Спойлер: ничего нового
- Просто новый подсчет “градиента”
- Считаем оценку на эпсилон  $\hat{\epsilon}(\mathbf{w})$

$$\hat{\epsilon}(\mathbf{w}) = \rho \operatorname{sign}(\nabla_{\mathbf{w}} L_{\mathcal{S}}(\mathbf{w})) |\nabla_{\mathbf{w}} L_{\mathcal{S}}(\mathbf{w})|^{q-1} / \left( \|\nabla_{\mathbf{w}} L_{\mathcal{S}}(\mathbf{w})\|_q^q \right)^{1/p}$$

- Оцениваем градиент

$$\nabla_{\mathbf{w}} L_{\mathcal{S}}^{SAM}(\mathbf{w}) \approx \nabla_{\mathbf{w}} L_{\mathcal{S}}(\mathbf{w})|_{\mathbf{w} + \hat{\epsilon}(\mathbf{w})}$$

# Как работает обучение

- Спойлер: ничего нового
- Просто новый подсчет “градиента”
- Считаем оценку на эпсилон  $\hat{\epsilon}(\mathbf{w})$

$$\hat{\epsilon}(\mathbf{w}) = \rho \operatorname{sign}(\nabla_{\mathbf{w}} L_{\mathcal{S}}(\mathbf{w})) |\nabla_{\mathbf{w}} L_{\mathcal{S}}(\mathbf{w})|^{q-1} / \left( \|\nabla_{\mathbf{w}} L_{\mathcal{S}}(\mathbf{w})\|_q^q \right)^{1/p}$$

- Оцениваем градиент

$$\nabla_{\mathbf{w}} L_{\mathcal{S}}^{SAM}(\mathbf{w}) \approx \nabla_{\mathbf{w}} L_{\mathcal{S}}(\mathbf{w})|_{\mathbf{w} + \hat{\epsilon}(\mathbf{w})}$$

- Обучаем

# Сравнения с другими моделями

# Сравнения с другими моделями

Model	Augmentation	CIFAR-10		CIFAR-100	
		SAM	SGD	SAM	SGD
WRN-28-10 (200 epochs)	Basic	<b>2.7</b> $\pm 0.1$	3.5 $\pm 0.1$	<b>16.5</b> $\pm 0.2$	18.8 $\pm 0.2$
WRN-28-10 (200 epochs)	Cutout	<b>2.3</b> $\pm 0.1$	2.6 $\pm 0.1$	<b>14.9</b> $\pm 0.2$	16.9 $\pm 0.1$
WRN-28-10 (200 epochs)	AA	<b>2.1</b> $\pm <0.1$	2.3 $\pm 0.1$	<b>13.6</b> $\pm 0.2$	15.8 $\pm 0.2$
WRN-28-10 (1800 epochs)	Basic	<b>2.4</b> $\pm 0.1$	3.5 $\pm 0.1$	<b>16.3</b> $\pm 0.2$	19.1 $\pm 0.1$
WRN-28-10 (1800 epochs)	Cutout	<b>2.1</b> $\pm 0.1$	2.7 $\pm 0.1$	<b>14.0</b> $\pm 0.1$	17.4 $\pm 0.1$
WRN-28-10 (1800 epochs)	AA	<b>1.6</b> $\pm 0.1$	2.2 $\pm <0.1$	<b>12.8</b> $\pm 0.2$	16.1 $\pm 0.2$
Shake-Shake (26 2x96d)	Basic	<b>2.3</b> $\pm <0.1$	2.7 $\pm 0.1$	<b>15.1</b> $\pm 0.1$	17.0 $\pm 0.1$
Shake-Shake (26 2x96d)	Cutout	<b>2.0</b> $\pm <0.1$	2.3 $\pm 0.1$	<b>14.2</b> $\pm 0.2$	15.7 $\pm 0.2$
Shake-Shake (26 2x96d)	AA	<b>1.6</b> $\pm <0.1$	1.9 $\pm 0.1$	<b>12.8</b> $\pm 0.1$	14.1 $\pm 0.2$
PyramidNet	Basic	<b>2.7</b> $\pm 0.1$	4.0 $\pm 0.1$	<b>14.6</b> $\pm 0.4$	19.7 $\pm 0.3$
PyramidNet	Cutout	<b>1.9</b> $\pm 0.1$	2.5 $\pm 0.1$	<b>12.6</b> $\pm 0.2$	16.4 $\pm 0.1$
PyramidNet	AA	<b>1.6</b> $\pm 0.1$	1.9 $\pm 0.1$	<b>11.6</b> $\pm 0.1$	14.6 $\pm 0.1$
PyramidNet+ShakeDrop	Basic	<b>2.1</b> $\pm 0.1$	2.5 $\pm 0.1$	<b>13.3</b> $\pm 0.2$	14.5 $\pm 0.1$
PyramidNet+ShakeDrop	Cutout	<b>1.6</b> $\pm <0.1$	1.9 $\pm 0.1$	<b>11.3</b> $\pm 0.1$	11.8 $\pm 0.2$
PyramidNet+ShakeDrop	AA	<b>1.4</b> $\pm <0.1$	1.6 $\pm <0.1$	<b>10.3</b> $\pm 0.1$	10.6 $\pm 0.1$

# Сравнения с другими моделями: ResNet

Model	Epoch	SAM		Standard Training (No SAM)	
		Top-1	Top-5	Top-1	Top-5
ResNet-50	100	<b>22.5</b> $\pm 0.1$	6.28 $\pm 0.08$	22.9 $\pm 0.1$	6.62 $\pm 0.11$
	200	<b>21.4</b> $\pm 0.1$	5.82 $\pm 0.03$	22.3 $\pm 0.1$	6.37 $\pm 0.04$
	400	<b>20.9</b> $\pm 0.1$	5.51 $\pm 0.03$	22.3 $\pm 0.1$	6.40 $\pm 0.06$
ResNet-101	100	<b>20.2</b> $\pm 0.1$	5.12 $\pm 0.03$	21.2 $\pm 0.1$	5.66 $\pm 0.05$
	200	<b>19.4</b> $\pm 0.1$	4.76 $\pm 0.03$	20.9 $\pm 0.1$	5.66 $\pm 0.04$
	400	<b>19.0</b> $\pm <0.01$	4.65 $\pm 0.05$	22.3 $\pm 0.1$	6.41 $\pm 0.06$
ResNet-152	100	<b>19.2</b> $\pm <0.01$	4.69 $\pm 0.04$	20.4 $\pm <0.0$	5.39 $\pm 0.06$
	200	<b>18.5</b> $\pm 0.1$	4.37 $\pm 0.03$	20.3 $\pm 0.2$	5.39 $\pm 0.07$
	400	<b>18.4</b> $\pm <0.01$	4.35 $\pm 0.04$	20.9 $\pm <0.0$	5.84 $\pm 0.07$

# Сравнения с другими моделями: файнтюнинг

Dataset	EffNet-b7 + SAM	EffNet-b7	Prev. SOTA (ImageNet only)	EffNet-L2 + SAM	EffNet-L2	Prev. SOTA
FGVC_Aircraft	6.80 $\pm$ 0.06	8.15 $\pm$ 0.08	<b>5.3</b> (TBMSL-Net)	<b>4.82</b> $\pm$ 0.08	5.80 $\pm$ 0.1	5.3 (TBMSL-Net)
Flowers	<b>0.63</b> $\pm$ 0.02	1.16 $\pm$ 0.05	0.7 (BiT-M)	<b>0.35</b> $\pm$ 0.01	0.40 $\pm$ 0.02	0.37 (EffNet)
Oxford_IIT_Pets	<b>3.97</b> $\pm$ 0.04	4.24 $\pm$ 0.09	4.1 (Gpipe)	<b>2.90</b> $\pm$ 0.04	3.08 $\pm$ 0.04	4.1 (Gpipe)
Stanford_Cars	5.18 $\pm$ 0.02	5.94 $\pm$ 0.06	<b>5.0</b> (TBMSL-Net)	4.04 $\pm$ 0.03	4.93 $\pm$ 0.04	<b>3.8</b> (DAT)
CIFAR-10	<b>0.88</b> $\pm$ 0.02	0.95 $\pm$ 0.03	1 (Gpipe)	<b>0.30</b> $\pm$ 0.01	0.34 $\pm$ 0.02	0.63 (BiT-L)
CIFAR-100	<b>7.44</b> $\pm$ 0.06	7.68 $\pm$ 0.06	7.83 (BiT-M)	<b>3.92</b> $\pm$ 0.06	4.07 $\pm$ 0.08	6.49 (BiT-L)
Birdsnap	<b>13.64</b> $\pm$ 0.15	14.30 $\pm$ 0.18	15.7 (EffNet)	<b>9.93</b> $\pm$ 0.15	10.31 $\pm$ 0.15	14.5 (DAT)
Food101	7.02 $\pm$ 0.02	7.17 $\pm$ 0.03	7.0 (Gpipe)	<b>3.82</b> $\pm$ 0.01	3.97 $\pm$ 0.03	4.7 (DAT)
ImageNet	15.14 $\pm$ 0.03	15.3	<b>14.2</b> (KDforAA)	<b>11.39</b> $\pm$ 0.02	11.8	11.45 (ViT)

# Сравнение: устойчивость к шуму в данных

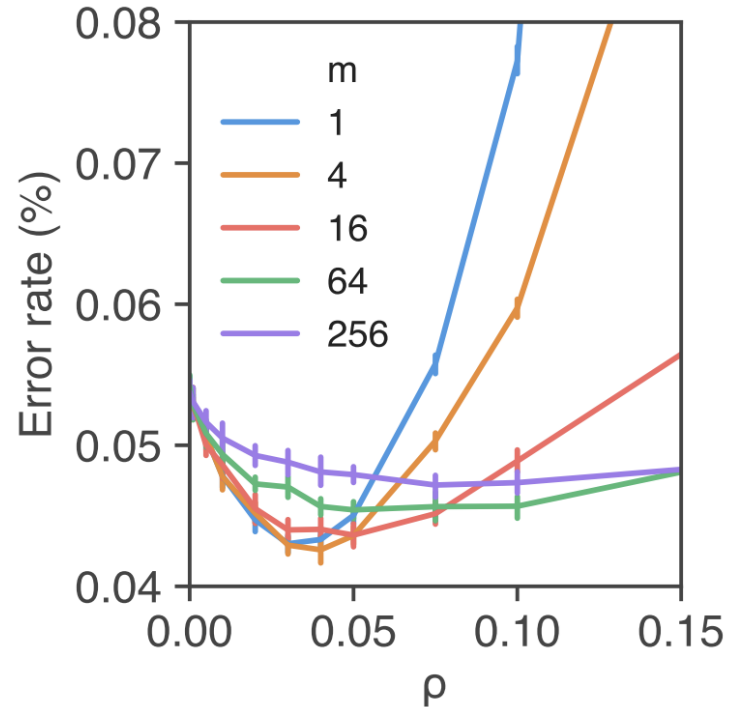
Method	Noise rate (%)			
	20	40	60	80
Sanchez et al. (2019)	94.0	92.8	90.3	74.1
Zhang & Sabuncu (2018)	89.7	87.6	82.7	67.9
Lee et al. (2019)	87.1	81.8	75.4	-
Chen et al. (2019)	89.7	-	-	52.3
Huang et al. (2019)	92.6	90.3	43.4	-
MentorNet (2017)	92.0	91.2	74.2	60.0
Mixup (2017)	94.0	91.5	86.8	76.9
MentorMix (2019)	<b>95.6</b>	<b>94.2</b>	91.3	<b>81.0</b>
SGD	84.8	68.8	48.2	26.2
Mixup	93.0	90.0	83.8	70.2
Bootstrap + Mixup	93.3	92.0	87.6	72.0
SAM	95.1	93.4	90.5	77.9
Bootstrap + SAM	95.4	<b>94.2</b>	<b>91.8</b>	79.9



# Подбор гиперпараметров

$m$  - размер батча

$\rho$  - радиус области



СПАСИБО ЗА ВНИМАНИЕ



# Литература и ссылки

- Статья: <https://arxiv.org/abs/2010.01412>
- Github с реализацией: <https://github.com/google-research/sam>
- Лучшие модели с которыми сравнивались
- EffNet: <https://arxiv.org/abs/1905.11946>
- GPipe: <https://arxiv.org/abs/1811.06965>
- DAT: <https://arxiv.org/abs/1811.07056>
- MentorMix: <https://arxiv.org/abs/1911.09781>