**Understanding the robustness difference between stochastic gradient descent and adaptive gradient methods**

Avery Ma[1,2]       Yangchen Pan[3]       Amir-massoud Farahmand[1,2]

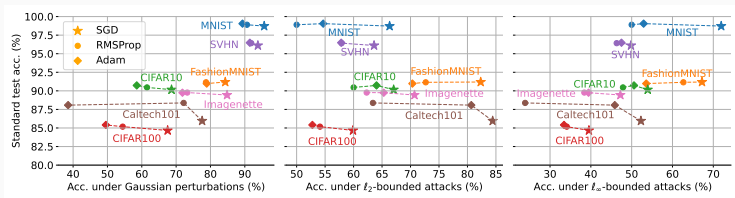[1]University of Toronto      [2]Vector Institute      [3]University of Oxford

**Figure 1: Comparison between models trained using SGD, Adam, and RMSProp.** Models trained by different algorithms have similar standard generalization performance, but there is a distinct robustness difference as measured by the test data accuracy under Gaussian noise, $\ell_2$ and $\ell_\infty$ bounded adversarial perturbations.[1]

---

[1] Croce et al., Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks

**More on the motivating example:**

1. Audio dataset.
   - Classifying short audio phrases (i.e., numbers, directions, etc.) from the Speech Commands dataset[2].
2. Both CNN and transformer architectures.

$\Rightarrow$ **Observing the robustness differences trained by different optimizers, consistently.**

---

[2]Warden et al., A dataset for limited-vocabulary speech recognition.

To optimize the standard training objective, models only need to learn **how to correctly use relevant information in the data**. Their use of irrelevant information in the data, however, is **under-constrained and can lead to solutions sensitive to perturbations**.

To optimize the standard training objective, models only need to learn how to correctly use relevant information in the data (OB I). Their use of irrelevant information in the data, however, is under-constrained and can lead to solutions sensitive to perturbations (OB II).
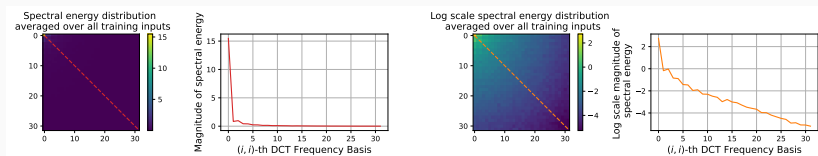


**Figure 2: Illustration of the spectral energy distribution in CIFAR100.** Distribution of the spectral energy heavily concentrates at low frequencies and decays exponentially towards higher frequencies.

Original Image (left), DCT transformation (mid), Log-scale (right)

# Observation I: Irrelevant Frequencies in Natural Signals

To optimize the standard training objective, models only need to learn how to correctly use relevant information in the data (OB I). Their use of irrelevant information in the data, however, is under-constrained and can lead to solutions sensitive to perturbations (OB II).
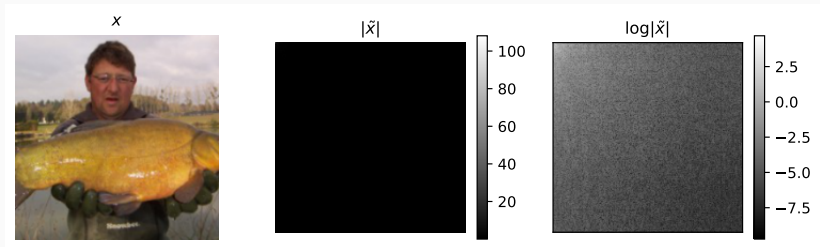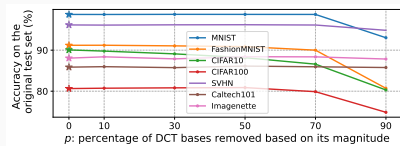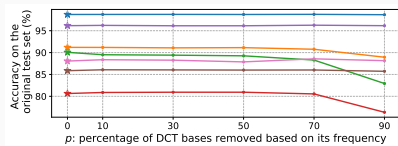


**a.** Parts of the signal with **low spectral energy** is irrelevant.

**b.** Parts of the signal with **high-frequency basis** is irrelevant.

**Figure 4: Irrelevant frequencies exist in the natural data.** Accuracy on the original test set remains high when the training inputs are augmented by removing parts of the signal with a) low spectrum energy and b) high frequencies.

low energy removed                    high freq. removed

To optimize the standard training objective, models only need to learn how to correctly use relevant information in the data (OB I). Their use of irrelevant information in the data, however, is under-constrained and can lead to solutions sensitive to perturbations (OB II).
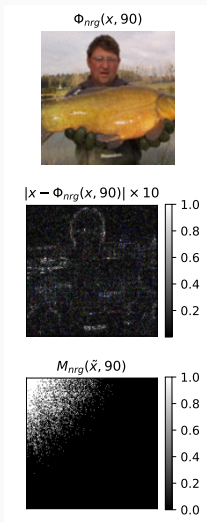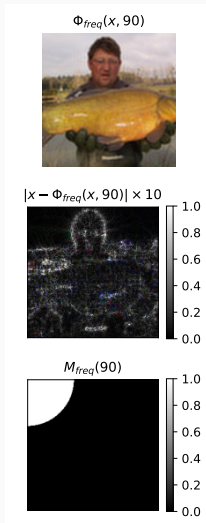


**Figure 6: The effect of band-limited Gaussian perturbations on the model.** Perturbations from the lowest band have a similar effect on all the models. On the other hand, models' responses vary significantly when the perturbation focuses on higher frequency bands.

⇒ **This suggest that the robustness differences comes from the different responses towards perturbations along those irrelevant frequencies.**

# Theoretical Analysis with Linear Models

**Setup:** In linear regression, we compare the <span style="color:red">standard</span> and <span style="color:red">adversarial</span> risk of the GD and signGD solutions.

**Data:** A synthetic dataset that mimics the characteristics of a natural dataset: contains irrelevant frequencies.

We are interested in the **standard risk**:

$$\mathcal{R}_{\mathsf{s}}(w) \triangleq \mathbb{E}\left[\ell(X, Y; w)\right]$$
$$= \mathbb{E}\left[|w^T X - Y|^2\right]$$

and the **adversarial risk** under $\ell_2$ bounded perturbations:

$$\mathcal{R}_{\mathsf{a}}(w) \triangleq \mathbb{E}\left[\max_{||\Delta x||_2 \leq \epsilon} \ell(X + \Delta x, Y; w)\right]$$
$$= \mathbb{E}\left[\max_{||\Delta x||_2 \leq \epsilon} |w^T(X + \Delta x) - Y|^2\right].$$

1. Irrelevant information leads to multiple standard risk minimizers. For an arbitrary minimizer $w^*$ from $\mathcal{W}^*$, we can obtain its the adversarial risk as

$$\mathcal{R}_a(w^*) = \frac{\epsilon^2}{2}||w^*||_2^2.$$

$\Rightarrow$ Given that the linear model is a standard risk minimizer, its robustness against $\ell_2$-norm bounded perturbations is proportional to its $\ell_2$ norm squared.

2. With a sufficiently small learning rate $\eta$, the standard risk of GD and signGD can be both close to 0.

$\Rightarrow$ Providing an explanation to the similar standard generalization performance observed on neural networks.

3. Consider a three dimensional input space, we demonstrate that ratio between the two adversarial risks is always greater than 1:

$$\frac{\mathcal{R}_a(\boldsymbol{w}^{\text{signGD}})}{\mathcal{R}_a(\boldsymbol{w}^{\text{GD}})} > 1 + C,$$

where $C > 0$ and its value depends on initialization, and the magnitude of the data covariance.

$\Rightarrow$ The linear model obtained through GD is **always more robust** against $\ell_2$-bounded perturbations in comparison to the model obtained from signGD.

**Can we show something similar to**

$\mathcal{R}_{\mathbf{a}}(w^*) = \frac{\epsilon^2}{2}||w^*||_2^2$ **on neural networks?**

Consider the form $f(x) = (\phi_l \circ \phi_{l-1} \circ ... \circ \phi_1)(x)$, where each $\phi_i$ is a linear operation, an activation function, or pooling operation.

Denoting the Lipschitz constant of function $f$ as $L(f)$, we can establish an upper bound[3] on the Lipschitz constant for the entire feed-forward neural network using

$$L(f) \leq \prod_{i=1}^{l} L(\phi_i).$$

The approximation to the Lipschitzness of various components of the network are typically functions of the weight norm. Examples: Linear operations: $\|W\|_p$; Skip-connections: $\|W\|_p + 1$.
$\Rightarrow$ Smaller weight norm implies smaller Lipschitz upper bound. Less vulnerable to perturbations.

---

[3] Any value of L satisfying the Lipschitz condition is considered a valid Lipschitz constant. For the sake of clarity, we will refer to the smallest (optimal) Lipschitz constant as L

# Product of Weight Norm of NN Upper-bounds its Lipschitzness

**Table 1: Comparing the upper bound on the Lipschitz constant and the averaged robust accuracy of neural networks trained by SGD, Adam, and RMSProp.**

|  | Dataset | MNIST | Fashion | CIFAR10 | CIFAR100 | SVHN | Caltech101 | Imagenette |
|---|---|---|---|---|---|---|---|---|
| $\prod_{i=1}^{l} L(\phi_i)$ | SGD | **3.80** | **3.83** | **26.81** | **40.41** | **22.65** | **18.53** | **23.99** |
|  | Adam | 5.75 | 8.12 | 28.70 | 41.87 | 30.45 | 26.20 | 28.55 |
|  | RMSProp | 6.21 | 5.11 | 37.75 | 41.71 | 28.31 | 45.84 | 27.11 |
| Averaged Robust Acc. | SGD | **77.97%** | **77.95%** | **63.21%** | **55.65%** | **69.08%** | **71.42%** | **67.59%** |
|  | Adam | 65.64% | 67.60% | 57.71% | 45.25% | 65.60% | 55.03% | 58.86% |
|  | RMSProp | 63.54% | 71.34% | 56.47% | 47.55% | 65.37% | 53.16% | 57.98% |

- SGD-trained neural networks have considerably smaller Lipschitz constants across all datasets.
- Explaining the better robustness to perturbations than those trained with adaptive gradient methods as shown in the motivating example.

## Summary

**Motivation**
Large robustness difference between models trained by SGD and adaptive gradient methods, despite similar standard generalization.

**Claims + empirical observations**
To optimize the standard training objective, models only need to learn how to correctly use relevant information in the data. Their use of irrelevant information in the data, however, is under-constrained and can lead to solutions sensitive to perturbations.

**Theoretical analysis with linear models**
We show the asymptotic solution found by signGD is always more vulnerable.

**Back to neural networks**
We demonstrate a connection between Lipschitz upper bound and model robustness.
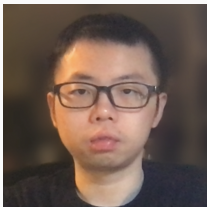
17

## Practical Implications

This work

- highlights the importance of optimizer selection in achieving both generalization and robustness.

- guides the development of optimization strategies that maintain high accuracy while being resilient to input perturbations.

    - Possible remedy: Initialization, Regularization (i.e., Adam/RMSProp+L2 or AdamW[4]).

---

[4] Loshchilov et al., Decoupled Weight Decay Regularization, ICLR 2019

Avery Ma



Yangchen Pan



Amir-massoud Farahmand

# Thank you!