

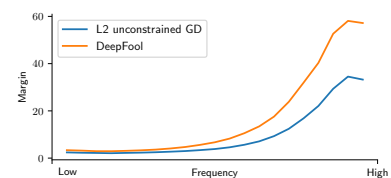
1 We would like to thank the reviewers for their valuable feedback, which we will duly consider and integrate in our  
2 revised manuscript. In this paper, we demonstrate that "the decision boundaries of a DNN can only exist as long  
3 as the classifier is trained with some features that hold them together", i.e., DNNs have an inductive bias towards  
4 invariance. Through "a multitude of insightful experiments", "with good attention to detail", we delve into this property  
5 which sheds light on open problems like "catastrophic forgetting" and "adv. robustness". The structure of our paper  
6 (R3) is designed to (i) "rigorously confirm" the existence of this inductive bias (Sec. 3), and (ii) further investigate  
7 its consequences on the sensitivity and dynamics of adv. training (Sec. 4). Throughout the paper, we back up all our  
8 claims, first, using controlled synthetic experiments, and then, "rigorously" verifying our hypotheses on real datasets  
9 with abundant empirical evidence. We clarify the main points raised by the reviewers here below.

10 **Margin and features (R4)** The main claim of our paper is that DNNs *only* create decision boundaries in regions where  
11 they identify discr. features in the training data. We further shed more light on the relationship between adv. examples  
12 and features studied in [3,4]. We show that there is a big relative difference in the large margin along the invariant  
13 dirs. and the smaller margin in the discr. dirs. Nevertheless, we never claim that, within the discr. dirs, margin is at all  
14 proportional to "discriminateness". In fact, we agree that the margin associated to different discr. features can greatly  
15 vary (Fig. 4). Overall, however, we firmly believe that the invariant dirs. will always have the largest margin.

16 **Causation (R4)** The main difficulty for establishing causation in our paper is the fact that the discr. features of real  
17 datasets are not known. Hence, determining their role on the geometry of a trained DNN can only be done by artificially  
18 manipulating the data. We strongly believe that the experiments in Sec. 3.2 are enough to rule out the other two main  
19 factors that might explain our results: the network and the algorithm. Specifically, in the flipping experiments, flipping  
20 the data – *ceteris paribus* – also flips the margin distribution, thus demonstrating that the margins are necessarily caused  
21 by the information present in the data. The other interventions we do on the data (e.g., low-pass experiments) confirm  
22 that in the absence of information in a certain dir. the network becomes invariant along this dir. Therefore, guided by  
23 the principles of the scientific method, and supported by strong evidence, we believe that grounds for causation are  
24 properly established.

25 **Choice of DCT (R1,R3)** The DCT has a long application tradition in image processing due to its good approximation  
26 of the decorrelating transform (KLT). Furthermore, in previous studies on the robustness of deep networks to different  
27 freqs., the DCT was also the basis of choice [7] because it avoids dealing with complex subspaces. A more aligned  
28 basis with respect to the discr. features would probably show a sharper transition between low and high margins.  
29 However, finding such network-agnostic bases is a challenging task without knowing the features *a priori*. The DCT is  
30 not perfectly feature-aligned, but it seems to be a good choice for comparing different architectures, especially if we  
31 compare its results to those obtained using a random orthonormal basis where differences in margin cannot be identified  
32 (c.f. Sec. N in Supp. material). We will include this explanation in the revised version of our manuscript.

33 **DeepFool (R3,R4)** In the adv. robustness literature DeepFool is generally re-  
34 garded as one of the most efficient methods to identify minimal adv. perts.  
35 Because we measure margin, norm-constrained attacks like PGD are not suitable  
36 for our study, and more complex attacks like C&W, or using unconstrained GD  
37 in the input space, are computationally much more demanding and harder to tune  
38 while finding very similar adv. perts. to DeepFool. Hence, DeepFool is more  
39 adequate for our work. We will include this explanation in the revised version of  
40 our manuscript. For completeness, we show a comparison of the median margin obtained on MNIST (LeNet) using  
41 DeepFool and a subspace-constrained GD attack: Even if the margins are slightly smaller for the *stronger* attack, the  
42 relative differences between regions (our quantity of interest) are the same for both attacks.



43 **CIFAR10 margin (R2)** Without knowing the mechanisms used by the network to select the discr. features of a dataset  
44 it is hard to give a full explanation of the low margins in the high frequencies of CIFAR10. Nevertheless, a possible  
45 reason for the different behaviour on CIFAR10 might be its low resolution. In fact, it is a common mistake that during  
46 the downsampling process no antialiasing filter is applied to the images before resizing, and hence some low-freq.  
47 information leaks to the high freqs. of the low-resolution images. This might explain why the network can identify some  
48 discr. features in the high-freq. spectrum of CIFAR10 (downsampling technique not specified in technical report [25]).

49 **Filtering other datasets (R2)** Indeed, the conclusion of the low-pass CIFAR10 experiment is generally applicable to  
50 other datasets, e.g., a LeNet trained on LP-MNIST (bandwidth of 14) has 0% drop in accuracy when tested on MNIST  
51 data. We will include this result in the appendix alongside Sec. C and Sec. H (results on HP-MNIST)

52 **Low-pass adv. training (R2)** As seen in Fig. 7, during adv. training, the energy of the perturbations has a very small  
53 high-freq. component, and it is predominantly concentrated in the low freqs. However, it seems that the small, but  
54 non-zero components in the high freqs. of the perturbations are necessary to improve robustness, as training only using  
55 low-freq. perturbations do not yield satisfactory robustness results.