

Random Projections for Improved Adversarial Robustness

1st Ginevra Carbone
Dept. of Mathematics and Geosciences
University of Trieste
Trieste, Italy
ginevra.carbone@phd.units.it

2nd Guido Sanguinetti
School of Informatics
University of Edinburgh
Edinburgh, United Kingdom
SISSA
Trieste, Italy
gsanguin@sissa.it

3nd Luca Bortolussi
Dept. of Mathematics and Geosciences
University of Trieste
Trieste, Italy
Modeling and Simulation Group
Saarland University
Saarland, Germany
luca.bortolussi@gmail.com

Abstract—We propose two training techniques for improving the robustness of Neural Networks to adversarial attacks, i.e. manipulations of the inputs that are maliciously crafted to fool networks into incorrect predictions. Both methods are independent of the chosen attack and leverage random projections of the original inputs, with the purpose of exploiting both dimensionality reduction and some characteristic geometrical properties of adversarial perturbations. The first technique is called *RP-Ensemble* and consists of an ensemble of networks trained on multiple projected versions of the original inputs. The second one, named *RP-Regularizer*, adds instead a regularization term to the training objective.

Index Terms—Adversarial robustness, Randomization, Regularization, Computational efficiency

I. INTRODUCTION

Adversarial examples [1] are small perturbations of the input data, specifically designed to induce wrong predictions in machine learning models, even for those achieving exceptional accuracy and with a high confidence in the wrong predictions. Such perturbations are often not even recognizable by humans [2, 3], thus developing suitable defense strategies is crucial in security-critical settings, and especially in computer vision algorithms (e.g. road signs recognition, medical imaging or autonomous driving [4]).

Defense research is currently focusing on different strategies for preventing this kind of vulnerability: deriving exact robustness bounds under some theoretical constraints [5], analyzing the robustness to particularly strong attacks [3], designing defences that are specific for the chosen attack [6], or developing general training algorithms and regularization techniques which improve resilience to multiple attacks.

Random projections of the input samples into lower-dimensional spaces have been extensively used for dimensionality reduction purposes, but in this work we are mostly interested in using them for providing robustness and regularization guarantees against the adversaries. Our main inspiration for this work is the *Manifold Hypothesis* [7, 8, 9], which models data as being sampled from low-dimensional manifolds, corresponding to the classification regions, embedded in a high-dimensional space [10]. Therefore, decision boundaries are represented as hypersurfaces of the embedding space. Such

approach allows to face the problem of high-dimensionality of the input space, since the number of samples required for learning grows exponentially with the dimension of the space. Geometrical inspections related to this phenomenon lead us to the idea of using random projections of the input data as a defense. We observe that projected versions of the original data are easier to learn and lie in less complex regions of the space.

We propose a training technique, called *RP-Ensemble* (II-A), which improves the robustness to adversarial examples of a pre-trained classifier. This method projects the input data in multiple lower dimensional spaces, each one determined by a random selection of directions in the space. Then, it trains a new classifier in each subspace, using the corresponding projected version of the data. Finally, it performs an ensemble classification on the original high dimensional data. In Sec. II-B we also define a regularization term for the training objective, named *RP-Regularizer*. This technique combines the norm of the loss gradients, intended as a measure of vulnerability, and the expectation over random projections of the inputs. In doing so, we aim at exploiting relevant adversarial features during training.

We evaluate the adversarial vulnerability of the resulting trained models and compare them to adversarially trained robust models (Sec. III). Finally, we discuss the scalability and parallelizability of RP-Ensemble.

II. METHODOLOGY

In the next sections we will refer to d -dimensional data samples and to neural network models of type $f(\cdot, \theta) : \mathbb{R}^d \rightarrow \mathbb{R}^K$, with learnable weights θ , solving a classification problem on K classes.

A. RP-Ensemble

RP-Ensemble method is built upon a pre-trained model and can be regarded as a fine tuning technique for adversarial robustness. Let $X \in \mathbb{R}^{n \times d}$ be the n original d -dimensional training examples, represented in matrix form, and let $g(\cdot, \theta) : \mathbb{R}^d \rightarrow \mathbb{R}^K$ be the pre-trained network.

First, we project the whole dataset into p different subspaces of dimension $k \leq d$, using Gaussian random projection matrices [11]. Each projection matrix $R_j \in \mathbb{R}^{k \times d}$ maps the input data X into its k -dimensional projected version $\mathcal{P}_j(X) = X R_j^T \in \mathbb{R}^{n \times k}$, using k random directions.

The elements of each random matrix R_j are independently drawn from a $\mathcal{N}(0, 1/k)$ distribution. This particular choice is motivated by Johnson-Lindenstrauss Lemma (II.1), ensuring that the Euclidean distance between any two points in the new low-dimensional space is approximately very close to the distance between the same points in the original high-dimensional space [11].

Lemma II.1 (Johnson-Lindenstrauss Lemma). *Given a set of n points $\mathcal{M} \subset \mathbb{R}^d$, let $\epsilon \in (0, 1/2)$, $k > 8 \log n / \epsilon^2$ and $A \in \mathbb{R}^{k \times d}$ be a matrix whose entries have been sampled independently from $\mathcal{N}(0, 1/k)$.*

Then, for any couple of points $u, v \in \mathcal{M}$ the following inequality holds

$$P \left[(1 - \epsilon) \|u - v\|^2 \leq \|Au - Av\|^2 \leq (1 + \epsilon) \|u - v\|^2 \right] \geq 1 - 2e^{-(\epsilon^2 - \epsilon^3)k/4}.$$

Notice that any two independently randomly chosen vectors in a high dimensional space are almost orthogonal with probability close to one and nearly have the same length [12]. Consequently, for any given number of sample points n , the k -dimensional columns of a random matrix R_j generated from a Gaussian distribution are almost orthogonal to each other. This procedure yields to a projection in the subspace generated by the columns of R_j .

Next, we train a classifier $\psi_j(\cdot, \theta_j) : \mathbb{R}^k \rightarrow \mathbb{R}^K$ in each projected subspace on the corresponding projected version of the data $\mathcal{P}_j(X)$. The architecture of the ψ_j -s mirrors that of the pre-trained network g , except for the first layer, which is adapted to the size of the projected lower dimensional input. Notice that the ψ_j -s do not share their weights nor the inputs, thus the backpropagation algorithm during the training phase stops at the projected data $\mathcal{P}_j(X)$.

Finally, we perform an ensemble classification on the original high dimensional data, by summing up the probability distributions from all the projected classifiers together with the predictions from the pre-trained classifier g . Let p_{θ_j}, p_g be the probability mass functions for the classifiers ψ_j and g . The classification of an input sample $x_i \in \mathcal{X}$ is given by

$$y_i := \arg \max_{y=1, \dots, K} \left(\sum_{j=1, \dots, p} p_{\theta_j}(y | \mathcal{P}_j(x_i)) + p_g(y | x_i) \right),$$

for each $i = 1, \dots, n$.

B. RP-Regularizer

RP-Regularizer is a variant of Total Variation regularization, a well known denoising approach in image processing [13],

already used in [14] as a regularization term for improving adversarial robustness. Its computation for a network $f(\cdot, \theta)$ is made tractable by a numerical approximation on the labeled dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1, \dots, n}$, e.g.

$$\|\nabla_x f\|_{L_1} \approx \frac{1}{n} \sum_{i=1}^n |\nabla_x \ell(f(x_i, \theta), y_i)|,$$

which is less complex w.r.t. the computation of the full gradient $\nabla_x f$. Here ℓ is the training loss function and n is the number of training samples. Our regularization term, instead, is computed in the L_2 -norm on suitable random projections of the input points.

The first step consists in sampling the components of the random matrices $R_j \in \mathbb{R}^{k_j \times d}$ from a Gaussian distribution, as was done in section II-A for RP-Ensemble, but with a randomly chosen projection size $k_j = 1, \dots, d$. Then, we project the input data matrix X in a k -dimensional subspace

$$\mathcal{P}_j(X) = X R_j^T \in \mathbb{R}^{n \times k_j},$$

R_j being the j -th projection matrix, for all the projection indexes $j = 1, \dots, p$.

Since the penalty for the objective needs to depend on the network's weights at the current training step, we want to map the projections $\mathcal{P}_j(X)$ back into the original d -dimensional space. We do this by means of *Moore-Penrose pseudo-inverse* [15] R_j^\dagger of R_j and apply it to the projected points

$$\mathcal{P}_j^\dagger(\mathcal{P}_j(X)) = X R_j^T (R_j^\dagger)^T \in \mathbb{R}^{n \times d}.$$

Fig. 1 shows an example of this procedure on the MNIST dataset [16]. In a nutshell, it builds on the two projection operators

$$\begin{aligned} \mathcal{P} : \mathbb{R}^{n \times d} &\longrightarrow \prod_{j=1}^p \mathbb{R}^{n \times k_j} \\ \mathcal{P}^\dagger : \prod_{j=1}^p \mathbb{R}^{n \times k_j} &\longrightarrow \prod_{j=1}^p \mathbb{R}^{n \times d}. \end{aligned}$$

The pseudo-inverse is a generalized inverse matrix. It exists and is unique for any given real rectangular matrix and the resulting composition $R_j^T (R_j^\dagger)^T$ is an orthogonal projection operator on \mathbb{R}^d .

Let ℓ be the training loss function. We propose two possible formulations for the regularization term $\mathcal{R}(\theta)$ of the objective $J(\theta) = \ell(\theta) + \lambda \mathcal{R}(\theta)$ on a set of weights θ . The first one, namely \mathcal{R}_{v1} , adds a penalty which is proportional to the expected norm of the loss gradients computed on the projected data

$$\begin{aligned} \mathcal{R}_{v1} &= \mathbb{E}_x \left[\mathbb{E}_{\mathcal{P}} \left[\|\nabla_x \ell(f(\mathcal{P}^\dagger \mathcal{P}(x), \theta), y)\|_2^2 \right] \right] \\ &\approx \frac{1}{np} \sum_{\substack{i=1, \dots, n \\ j=1, \dots, p}} \left\| \nabla_x \ell(f(\mathcal{P}_j^\dagger \mathcal{P}_j(x_i), \theta), y_i) \right\|_2^2. \quad (1) \end{aligned}$$

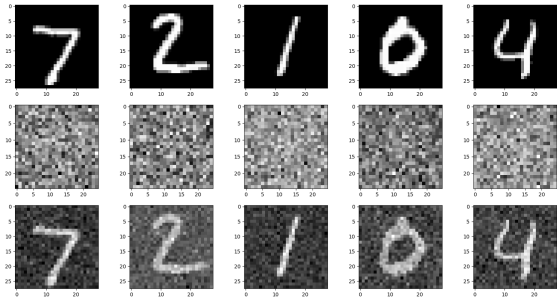


Fig. 1. First row shows the original image, second row its projected version $\mathcal{P}_j(x)$, third row the inverse projection $\mathcal{P}_j^\dagger \mathcal{P}_j(x)$. In this example we computed 25×25 projections on 28×28 MNIST samples.

The natural interpretation of \mathcal{R}_{v1} is that it allows to minimize loss variation across the $\mathcal{P}_j^\dagger \mathcal{P}_j(x_i)$ -s.

The second one, \mathcal{R}_{v2} , minimizes the variation of the loss gradients on the original inputs in randomly chosen projected subspaces

$$\begin{aligned} \mathcal{R}_{v2} &= \mathbb{E}_x \left[\mathbb{E}_{\mathcal{P}} \left[\left\| \mathcal{P} \left(\nabla_x \ell(f(x, \theta), y) \right) \right\|_2^2 \right] \right] \\ &\approx \frac{1}{np} \sum_{\substack{i=1, \dots, n \\ j=1, \dots, p}} \left\| \mathcal{P}_j \left(\nabla_x \ell(f(x_i, \theta), y_i) \right) \right\|_2^2. \end{aligned} \quad (2)$$

At each training step we perform a finite approximation of the expectations on minibatches of data, by randomly sampling the directions, the dimension of the projected subspace and the number of projections.

The two regularization terms \mathcal{R}_{v1} and \mathcal{R}_{v2} are equivalent as $k \rightarrow \infty$.

Theorem II.1. Let \mathcal{R}_{v1} and \mathcal{R}_{v2} be the regularization terms defined in Eq. 1 and Eq. 2, where $\mathcal{P} : \mathbb{R}^d \rightarrow \mathbb{R}^k$ is a random projection such that the elements of the orthogonal random matrix R are sampled from $\mathcal{N}(0, 1/k)$. If $k \in O(d)$ then $\mathcal{R}_{v1} \approx \mathcal{R}_{v2}$ as $k \rightarrow \infty$.

We provide a formal proof of Theorem II.1 in Section VI-C of the Appendix.

III. EXPERIMENTAL RESULTS

We evaluated the proposed methods on image classification tasks with 10 classes, using MNIST [16] and CIFAR-10 [17] dataset. Our baseline models are Convolutional Neural Networks with ReLU activation functions. We achieved 99.13% accuracy on MNIST and 76.52% on CIFAR-10. The adversarial attacks in our tests are Fast Gradient Sign Method (FGSM) [6], Projected Gradient Descent (PGD) [2], DeepFool [18], and Carlini and Wagner (C&W) in the L_∞ norm [19]. The attacks just mentioned are described in Section VI-B of the Appendix. In all such cases, the maximum distance between an image and its adversarial perturbation is set to $\epsilon = 0.3$. These methods fall in the white-box category, i.e. they have complete knowledge of their target network. However,

due to the transferability property of the attacks, they could also be effective on unknown models.

Simulations were conducted on a machine with 34 single core Intel(R) Xeon(R) Gold 6140 CPU @ 2.30GHz processors and 200GB of RAM. We made an extensive use of Tensorflow [20] and IBM adversarial-robustness-toolbox [21] libraries.

A. Adversarial robustness

Our approach for evaluating the robustness consists in testing the baseline models against several adversarial attacks, using the generated attacks to perform adversarial training on the baselines and, finally, comparing RP-Ensemble and RP-Regularizer with these robust baselines. Such procedure is intended to investigate the generalization capabilities of our methods, which are completely unaware of the chosen attacks, yet compared to models that should exhibit ideal performances against the adversaries, i.e. the adversarially trained ones.

We trained multiple versions of RP-Ensemble, using all the possible combinations of number of projections and size of the projected subspaces shown in Table I. The choice of the classifiers in RP-Ensemble is arbitrary and does not require any model selection step. Indeed, in our experiments each classifier ψ_j is indexed by the seed j used for sampling the projection matrix R_j .

We trained a single version of RP-Regularizer on each dataset, by uniformly sampling the number of projections and size of each projection at each training step, as reported in Table VI. We computed the pseudo-inverse matrix R_j^\dagger by using the SVD decomposition of R_j , in order to ensure numerical stability.

TABLE I
NUMBER OF PROJECTIONS AND SIZE OF EACH PROJECTION USED FOR RP-ENSEMBLE.

Dataset	Number of projections	Projection size
MNIST	6, 9, 12, 15	8, 12, 16, 20
CIFAR-10	3, 6, 9, 12	4, 8

TABLE II
NUMBER OF PROJECTIONS AND SIZE OF EACH PROJECTION USED FOR RP-REGULARIZER.

Dataset	Number of projections	Projection size
MNIST	$n_proj \sim \mathcal{U}(2, 8)$	$size_proj \sim \mathcal{U}(15, 25)$
CIFAR-10	$n_proj = 1$	$size_proj \sim \mathcal{U}(5, 10)$

We crafted adversarial perturbations on the original test set using the baseline model, then tested the robustness of RP-Ensemble to the adversaries in terms of prediction accuracy, both on MNIST (Fig. 2) and CIFAR-10 (Fig. 3). Tables III and IV in the Appendix report the exact numerical values for the prediction accuracy. RP-Ensemble brings a general improvement in the adversarial robustness of the baseline model. Adversarially trained robust models show great results on their

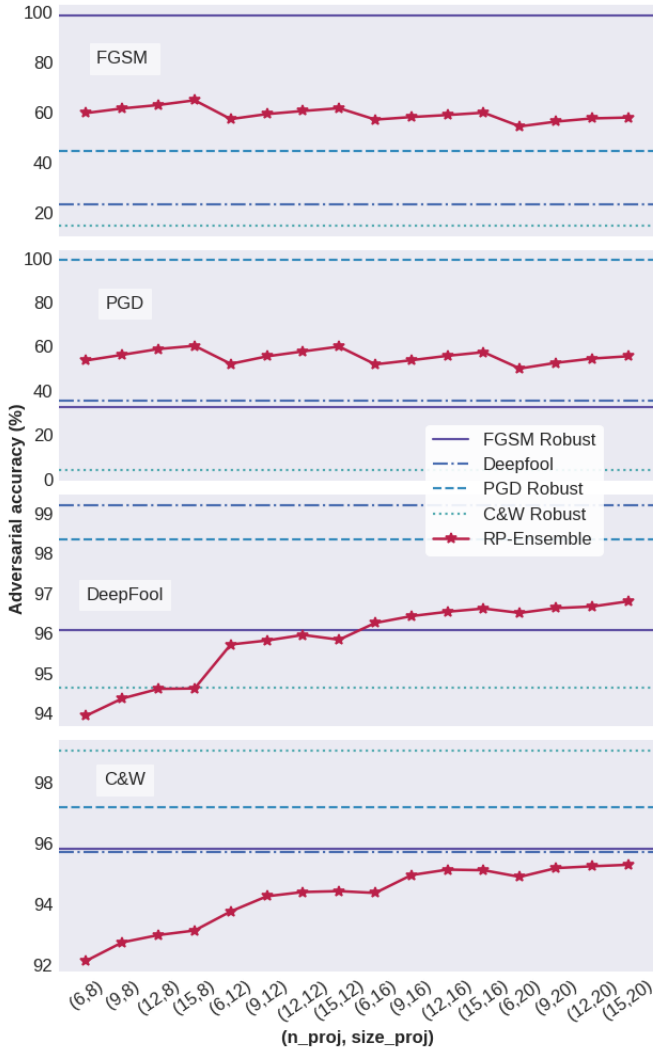


Fig. 2. Test accuracy of the baseline, its robust versions and RP-Ensemble model on MNIST dataset. The robust models are the result of adversarial training on the perturbed training sets. RP-Ensemble model is been trained on multiple combinations of number of projections and size of each projection. The evaluations are performed on the original test set and its adversarially perturbed versions.

target attacks but perform poorly on the other ones, while RP-Ensemble preserves its robustness across the different attacks.

RP-Regularizer is able to reach competitive performances in comparison to the SOTA models on MNIST (Fig. 4). Prediction accuracies are higher on DeepFool and Carlini & Wagner attacks than on FGSM and PGD, suggesting that this method performs better on algorithms which are optimized to produce perturbations that are closer to the original samples (e.g. C&W), rather than faster in computation (e.g. FGSM). The results are less striking on CIFAR-10 (Fig. 5), but we stress that the robustness of RP-Regularizer improves as the number of projections increases and that on CIFAR-10 we kept it low (always equal to 1) to maintain a balance between computational efficiency and adversarial accuracy. The trade-off between these two objectives needs to be further explored.

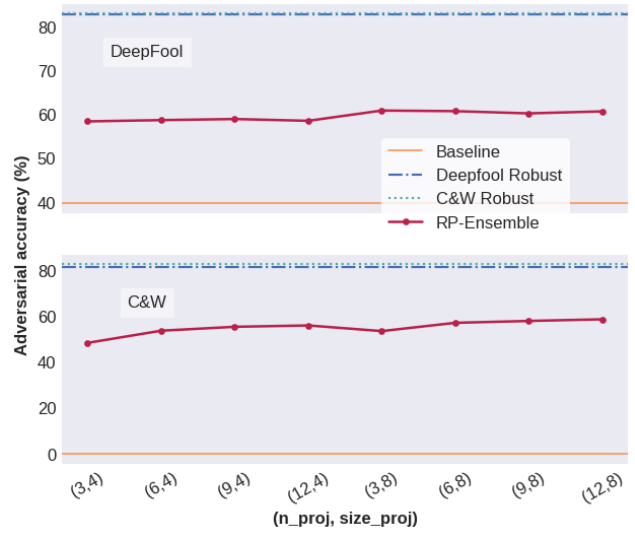


Fig. 3. Test accuracy of the baseline, its robust versions and RP-Ensemble model on CIFAR-10 dataset. The robust models are the result of adversarial training on the perturbed training sets. RP-Ensemble model is been trained on multiple combinations of number of projections and size of each projection. The evaluations are performed on the original test set and its adversarially perturbed versions.

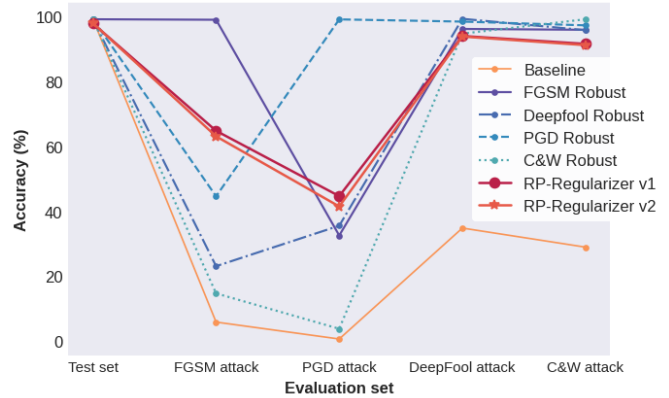


Fig. 4. Test accuracy of RP-Regularizer on MNIST. We compare the baseline model, the adversarially trained robust models and two different versions of RP-Regularizer model, namely \mathcal{R}_{v1} (eq. 1) and \mathcal{R}_{v2} (eq. 2). Adversarial perturbations are produced on the baseline model using FGSM, PGD, Deepfool and Carlini & Wagner attacks.

In RP-Regularizer ℓ is a crossentropy loss function.

B. Computational efficiency of RP-Ensemble

Classifiers ψ_j from PR-Ensemble are defined in independent projected subspaces, thus their training can be efficiently parallelized. This allows to keep its training time close to that of the baseline, or even lower when choosing a small number of projections (Fig. 6).

Moreover, pairwise distances between the projected points are nearly preserved, so the projected versions of the original images contain most of the original information. This implies that the projection classifiers ψ_j are computationally efficient, due to the dimensionality reduction of their inputs, and are

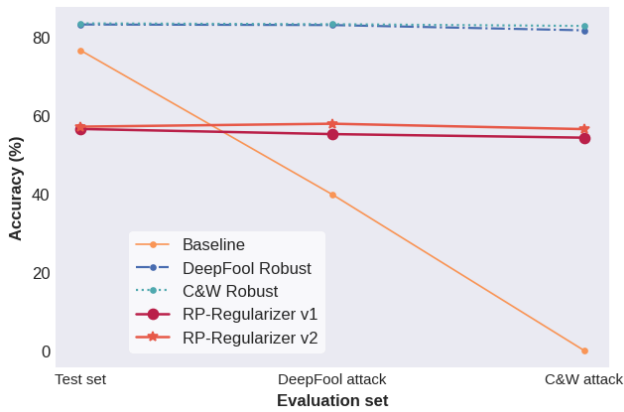


Fig. 5. Test accuracy of RP-Regularizer on CIFAR-10. We compare the baseline model, the adversarially trained robust models and two different versions of RP-Regularizer model, namely \mathcal{R}_{v1} (eq. 1) and \mathcal{R}_{v2} (eq. 2). Adversarial perturbations are produced on the baseline model using FGSM, PGD, Deepfool and Carlini & Wagner attacks.

also able to learn features which turn out to be significant as a defense against the attacks.

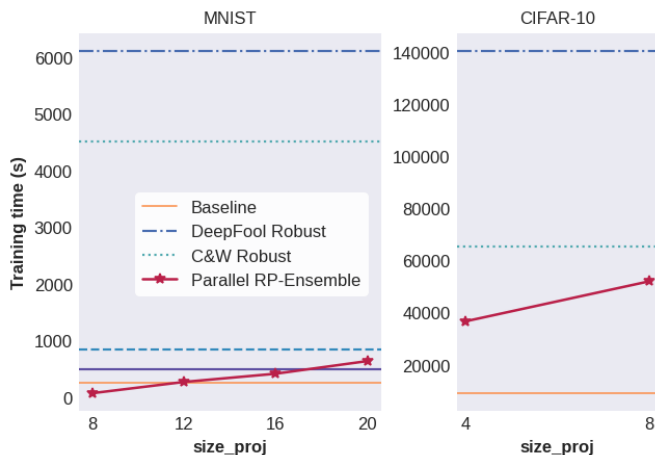


Fig. 6. Training time of RP-Ensemble on MNIST and CIFAR-10. We compare its efficiency to that of the baseline models and the adversarially trained robust models.

IV. RELATED WORK

A. Adversarial training

One of the simplest and most effective approaches for learning robust models is *adversarial training* [6]. This process consists in training a classifier by including adversarial examples in the training data, thus it allows to directly convert any attack into a defense. The biggest limitation of this approach is that it tends to overfit the chosen attack, meaning that the adversarially trained model perform well on the attacks which they learned to defend from, but might show poor transferability on other threats. Adversarial training could be interpreted as a form of data augmentation, which significantly differs from the traditional approach: instead of applying transformations

that are expected to occur in the test set (translations, rotations, etc.), only the most unlikely examples are added. This method corresponds to a dilation of the manifold: adversarial examples are learned in a halo around the surface, which makes the manifold smoother [22]. In this regard, it should be noted that RP-Ensemble does not perform any data augmentation in the original high dimensional space, since the projected data samples lie in new subspaces. RP-Regularizer, instead, produces new high dimensional examples, which could be formulated as perturbations of the original ones.

In [14] Finlay et al. show that regularization of the loss gradient on the inputs improves adversarial robustness. In particular, they notice that the *Total Variation* regularization [13] can be interpreted as the regularization induced by a single step of adversarial training on gradient-based attacks.

B. Randomization in a high codimension setting

Randomization has been proven effective as a defense [23, 24], a detection [25] and a regularization [26] strategy against adversarial attacks. Xie et al. [27] apply two random transformations to the input images, Liu, Cheng, Zhang and Hsieh [28] add random noise between the layers of the architecture, Dhillon et al. [29] randomly prune activations between the layers, Xu, Evans and Qi [30] and Liao and Wagner [31] propose input denoising and feature denoising methods. The reasoning behind these techniques is that NNs are usually robust to random perturbations [32], thus incorporating them in the models might weaken adversarial perturbations.

We also explore the geometrical properties related to randomization. Khuory et al. [10] first highlighted the role of codimension in the generation of adversarial examples. Their analysis suggests that adversarial perturbations mainly arise in the directions that are normal to the data manifold, so as the codimension in the embedding space increases there is a higher number of directions in which one could build adversarial perturbations. Our framework is strongly influenced by this finding, as it suggests that by randomly selecting directions it should be more likely to catch features that are significant in the adversarial context.

C. Bayesian interpretation of ensembles

Recent findings suggests a connection between ensembles of NNs and Bayesian NNs, where the goal is to learn the posterior distribution on the weights and use it to perform predictive inference on new observations. Lakshminarayanan, Pritzel and Blundell presented *Deep Ensembles* [33] as a computationally cheap alternative to Bayesian NNs. Dropout has been used to estimate the predictive uncertainty in *MC-Dropout* ensembles [34]. Latest research also shows that Bayesian inference is effective at learning adversarially robust models [35, 36]. Moreover, it has been proved that, under specific theoretical assumptions, Bayesian NNs are robust to gradient-based attacks [37].

D. Robustness in a chosen norm

The robustness of classifiers is strongly related to the geometry of the learned decision boundaries. In fact, in order

to learn robust decision boundaries, a model has to correctly classify all the input points lying in a neighbourhood of the data manifold. In particular, adversarially perturbed points always lie extremely close to the decision boundaries [22]. But robustness conditions change under different p -norms, meaning that no single decision boundary can be optimally robust in all norms [10]. E.g. if a classifier is trained to be robust under L_∞ norm, poor robustness under the L_2 norm should be expected. In general, no distance metric can be considered a perfect measure of similarity [19], so one of the strengths of our methods is that they are totally independent on the norm chosen for the attacks.

V. CONCLUSIONS

Adversarial examples show that many of the modern machine learning algorithms can be fooled in unexpected ways. Both in terms of attacks and defenses, many theoretical problems still remain open. From a practical point of view, no one has yet designed a powerful defense algorithm which could be suitable against a variety of attacks, with different degrees of knowledge about models under attack and their predictions. The most effective defense techniques, e.g. adversarial training, are still too computationally expensive.

We empirically showed that random projections of the training data act as attack-independent adversarial features, that can be used to provide better resilience to adversarial perturbations. We proposed a fine-tuning method and a regularization method, both based on the computation of random projections of the inputs. As future work we plan to improve the computational cost of RP-Regularizer and to compare the performances of our methods to that of other attack-independent defense strategies. We believe that a further exploration of the connections between random projections and the geometrical characterization of adversarial regions could bring valuable insights to adversarial defense research.

REFERENCES

- [1] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks, 2014.
- [2] Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. Adversarial examples in the physical world. *CoRR*, abs/1607.02533, 2016. URL <http://arxiv.org/abs/1607.02533>.
- [3] Nicholas Carlini, Anish Athalye, Nicolas Papernot, Wieland Brendel, Jonas Rauber, Dimitris Tsipras, Ian J. Goodfellow, Aleksander Madry, and Alexey Kurakin. On evaluating adversarial robustness. *CoRR*, abs/1902.06705, 2019. URL <http://arxiv.org/abs/1902.06705>.
- [4] Ivan Evtimov, Kevin Eykholt, Earlene Fernandes, Tadayoshi Kohno, Bo Li, Atul Prakash, Amir Rahmati, and Dawn Song. Robust physical-world attacks on machine learning models. *CoRR*, abs/1707.08945, 2017. URL <http://arxiv.org/abs/1707.08945>.
- [5] Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. Analysis of classifiers’ robustness to adversarial perturbations. *Machine Learning*, 107(3):481–508, 2018.
- [6] Ian Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *International Conference on Learning Representations*, 2015. URL <http://arxiv.org/abs/1412.6572>.
- [7] Christopher R Genovese, Marco Perone-Pacifico, Isabella Verdinelli, Larry Wasserman, et al. Manifold estimation and singular deconvolution under hausdorff loss. *The Annals of Statistics*, 40(2):941–963, 2012.
- [8] Sanjoy Dasgupta and Yoav Freund. Random projection trees and low dimensional manifolds. In *Proceedings of the fortieth annual ACM symposium on Theory of computing*, pages 537–546, 2008.
- [9] David L Donoho and Carrie Grimes. Hessian eigenmaps: Locally linear embedding techniques for high-dimensional data. *Proceedings of the National Academy of Sciences*, 100(10):5591–5596, 2003.
- [10] Marc Khoury and Dylan Hadfield-Menell. On the geometry of adversarial examples. *CoRR*, abs/1811.00525, 2018. URL <http://arxiv.org/abs/1811.00525>.
- [11] Sanjoy Dasgupta and Anupam Gupta. An elementary proof of the johnson-lindenstrauss lemma. Technical report, 1999.
- [12] Samuel Kaski. Dimensionality reduction by random mapping: Fast similarity computation for clustering. volume 1, pages 413 – 418 vol.1, 06 1998. ISBN 0-7803-4859-1. doi: 10.1109/IJCNN.1998.682302.
- [13] Leonid I Rudin, Stanley Osher, and Emad Fatemi. Nonlinear total variation based noise removal algorithms. *Physica D: nonlinear phenomena*, 60(1-4):259–268, 1992.
- [14] Chris Finlay, Adam M. Oberman, and Bilal Abbasi. Improved robustness to adversarial examples using lipschitz regularization of the loss, 2019.
- [15] R. Penrose. A generalized inverse for matrices. *Mathematical Proceedings of the Cambridge Philosophical Society*, 51(3):406–413, 1955. doi: 10.1017/S0305004100030401.
- [16] Yann LeCun and Corinna Cortes. MNIST handwritten digit database. 2010. URL <http://yann.lecun.com/exdb/mnist/>.
- [17] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-10 (canadian institute for advanced research). URL <http://www.cs.toronto.edu/~kriz/cifar.html>.
- [18] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. Deepfool: a simple and accurate method to fool deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2574–2582, 2016.
- [19] Nicholas Carlini and David A. Wagner. Towards evaluating the robustness of neural networks. *CoRR*, abs/1608.04644, 2016. URL <http://arxiv.org/abs/1608.04644>.
- [20] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene

- Brevedo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. URL <http://tensorflow.org/>. Software available from tensorflow.org.
- [21] Maria-Irina Nicolae, Mathieu Sinn, Tran Ngoc Minh, Amrith Rawat, Martin Wistuba, Valentina Zantedeschi, Ian M. Molloy, and Benjamin Edwards. Adversarial robustness toolbox v0.2.2. *CoRR*, abs/1807.01069, 2018. URL <http://arxiv.org/abs/1807.01069>.
- [22] Simant Dube. High dimensional spaces, deep learning and adversarial examples. *CoRR*, abs/1801.00634, 2018. URL <http://arxiv.org/abs/1801.00634>.
- [23] Zhezhi He, Adnan Siraj Rakin, and Deliang Fan. Parametric noise injection: Trainable randomness to improve deep neural network robustness against adversarial attack. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 588–597, 2019.
- [24] Alexandre Araujo, Laurent Meunier, Rafael Pinot, and Benjamin Negrevergne. Robust neural networks using randomized adversarial training. *arXiv preprint arXiv:1903.10219*, 2019.
- [25] Nathan Drenkow, Neil Fendley, and Philippe Burlina. Random projections for adversarial attack detection. *arXiv preprint arXiv:2012.06405*, 2020.
- [26] Robert J Durrant and Ata Kabán. Random projections as regularizers: learning a linear discriminant from fewer observations than dimensions. *Machine Learning*, 99(2): 257–286, 2015.
- [27] Cihang Xie, Jianyu Wang, Zhishuai Zhang, Zhou Ren, and Alan Yuille. Mitigating adversarial effects through randomization. *arXiv preprint arXiv:1711.01991*, 2017.
- [28] Xuanqing Liu, Minhao Cheng, Huan Zhang, and Cho-Jui Hsieh. Towards robust neural networks via random self-ensemble. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 369–385, 2018.
- [29] Guneet S Dhillon, Kamyar Azizzadenesheli, Zachary C Lipton, Jeremy Bernstein, Jean Kossaifi, Aran Khanna, and Anima Anandkumar. Stochastic activation pruning for robust adversarial defense. *arXiv preprint arXiv:1803.01442*, 2018.
- [30] Weilin Xu, David Evans, and Yanjun Qi. Feature squeezing: Detecting adversarial examples in deep neural networks. *arXiv preprint arXiv:1704.01155*, 2017.
- [31] Fangzhou Liao, Ming Liang, Yinpeng Dong, Tianyu Pang, Xiaolin Hu, and Jun Zhu. Defense against adversarial attacks using high-level representation guided denoiser. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1778–1787, 2018.
- [32] Kui Ren, Tianhang Zheng, Zhan Qin, and Xue Liu. Adversarial attacks and defenses in deep learning. *Engineering*, 6(3):346–360, 2020.
- [33] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. *arXiv preprint arXiv:1612.01474*, 2016.
- [34] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059. PMLR, 2016.
- [35] Artur Bekasov and Iain Murray. Bayesian adversarial spheres: Bayesian inference and adversarial examples in a noiseless setting. *arXiv preprint arXiv:1811.12335*, 2018.
- [36] Yarin Gal and Lewis Smith. Sufficient conditions for idealised models to have no adversarial examples: a theoretical and empirical study with bayesian neural networks. *arXiv preprint arXiv:1806.00667*, 2018.
- [37] Ginevra Carbone, Matthew Wicker, Luca Laurenti, Andrea Patane, Luca Bortolussi, and Guido Sanguinetti. Robustness of bayesian neural networks to gradient-based attacks, 2020.
- [38] Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.

VI. APPENDIX

A. Prediction accuracy

TABLE III

TEST ACCURACY OF THE BASELINE, ITS ROBUST VERSIONS AND RP-ENSEMBLE MODEL ON MNIST DATASET. THE ROBUST MODELS ARE THE RESULT OF ADVERSARIAL TRAINING ON THE PERTURBED TRAINING SETS. RP-ENSEMBLE MODEL IS BEEN TRAINED ON MULTIPLE COMBINATIONS OF NUMBER OF PROJECTIONS AND SIZE OF EACH PROJECTION. THE EVALUATIONS ARE PERFORMED ON THE ORIGINAL TEST SET AND ITS ADVERSARIALLY PERTURBED VERSIONS.

Prediction accuracy (%)	Test set	FGSM	PGD	DeepFool	C&W
Baseline	99.13	5.91	0.71	34.83	28.96
Adversarially trained models					
FGSM	99.13	98.91	32.55	96.08	95.82
PGD	99.10	44.60	99.02	98.34	97.19
DeepFool	99.03	23.11	35.55	99.20	95.70
C&W	99.10	14.74	3.85	94.63	99.06
RP-Ensemble on (n_proj, size_proj) combinations					
(6, 8)	97.66	59.71	53.68	93.94	92.12
(9, 8)	97.57	61.55	56.23	94.37	92.73
(12, 8)	97.45	62.93	58.85	94.61	92.97
(15, 8)	97.47	64.82	60.30	94.62	93.12
(6, 12)	98.12	57.29	52.12	95.72	93.75
(9, 12)	98.02	59.30	55.52	95.82	94.25
(12, 12)	97.97	60.54	57.77	95.96	94.39
(15, 12)	97.91	61.65	59.95	95.84	94.42
(6, 16)	98.22	57.09	51.90	96.26	94.36
(9, 16)	98.33	58.10	53.77	96.43	94.95
(12, 16)	98.32	58.93	55.78	96.54	95.13
(15, 16)	98.26	59.85	57.42	96.62	95.11
(6, 20)	98.49	54.37	50.02	96.51	94.89
(9, 20)	98.42	56.28	52.60	96.63	95.18
(12, 20)	98.40	57.56	54.52	96.67	95.24
(15, 20)	98.40	57.91	55.57	96.80	95.29

TABLE IV

TEST ACCURACY OF THE BASELINE, ITS ROBUST VERSIONS AND RP-ENSEMBLE MODEL ON CIFAR-10 DATASET. THE ROBUST MODELS ARE THE RESULT OF ADVERSARIAL TRAINING ON THE PERTURBED TRAINING SETS. RP-ENSEMBLE MODEL IS BEEN TRAINED ON MULTIPLE COMBINATIONS OF NUMBER OF PROJECTIONS AND SIZE OF EACH PROJECTION. THE EVALUATIONS ARE PERFORMED ON THE ORIGINAL TEST SET AND ITS ADVERSARIALLY PERTURBED VERSIONS.

Prediction accuracy (%)	Test set	DeepFool	C&W
Baseline	76.52	39.77	0.00
Adversarially trained models			
DeepFool	83.16	83.01	81.67
C&W	83.44	83.23	82.79
RP-Ensemble model on (n_proj, size_proj)			
(3, 4)	67.93	58.52	48.48
(6, 4)	64.59	58.81	53.78
(9, 4)	63.15	59.06	55.48
(12, 4)	61.93	58.65	56.07
(3, 8)	67.66	61.00	53.61
(6, 8)	64.83	60.86	57.21
(9, 8)	63.36	60.35	58.03
(12, 8)	62.99	60.81	58.75

TABLE V

TEST ACCURACY OF RP-REGULARIZER ON MNIST. WE COMPARE THE BASELINE MODEL, THE ADVERSARIALLY TRAINED ROBUST MODELS AND TWO DIFFERENT VERSIONS OF RP-REGULARIZER MODEL, NAMELY \mathcal{R}_{v1} (1) AND \mathcal{R}_{v2} (2). ADVERSARIAL PERTURBATIONS ARE PRODUCED ON THE BASELINE MODEL USING FGSM, PGD, DEEPFOOL AND CARLINI & WAGNER ATTACKS.

Prediction accuracy (%)	Test set	FGSM	PGD	DeepFool	C&W
Baseline	99.13	5.91	0.71	34.83	28.96
Adversarially trained models					
FGSM	99.13	98.91	32.55	96.08	95.82
PGD	99.10	44.60	99.02	98.34	97.19
DeepFool	99.03	23.11	35.55	99.20	95.70
C&W	99.10	14.74	3.85	94.63	99.06
RP-Regularizer model					
$\mathcal{R}_{v1}, \lambda = 0.4$	97.92	62.34	38.96	93.24	90.75
$\mathcal{R}_{v2}, \lambda = 0.4$	97.82	63.25	42.37	93.86	91.56
$\mathcal{R}_{v1}, \lambda = 0.5$	97.53	69.12	52.39	94.44	91.96
$\mathcal{R}_{v2}, \lambda = 0.5$	98.06	60.64	36.28	93.92	91.05
$\mathcal{R}_{v1}, \lambda = 0.6$	97.80	62.78	42.61	94.05	91.73
$\mathcal{R}_{v2}, \lambda = 0.6$	97.69	65.25	45.77	93.45	90.70

TABLE VI

TEST ACCURACY OF RP-REGULARIZER ON MNIST. WE COMPARE THE BASELINE MODEL, THE ADVERSARIALLY TRAINED ROBUST MODELS AND TWO DIFFERENT VERSIONS OF RP-REGULARIZER MODEL, NAMELY \mathcal{R}_{v1} (1) AND \mathcal{R}_{v2} (2). ADVERSARIAL PERTURBATIONS ARE PRODUCED ON THE BASELINE MODEL USING DEEPFOOL AND CARLINI & WAGNER ATTACKS.

Prediction accuracy (%)	Test set	DeepFool	C&W
Baseline	76.52	39.77	0.00
Adversarially trained models			
DeepFool	83.16	83.01	81.67
C&W	83.44	83.23	82.79
RP-Regularizer model, $\lambda = 0.5$			
\mathcal{R}_{v1}	56.51	55.20	54.29
\mathcal{R}_{v2}	57.10	57.85	56.47

B. Adversarial attacks

Fast Gradient Sign Method (FGSM) [6] is an untargeted attack, i.e. it does not push the misclassification to any specific class. FGSM adds a fixed noise to the input x in the direction of the loss gradient w.r.t. x

$$\tilde{x} = x + \epsilon \operatorname{sgn} \nabla_x \ell(f(x, \theta), y).$$

It is a popular choice for adversarial training on a large number of samples due to its computational efficiency, since it only requires one gradient evaluation for each given input.

Projected Gradient Descent (PGD) [2] is an iterative attack which starts from a random perturbation \tilde{x}_0 of x in an ϵ - L_∞ ball around the input sample. At each step, it performs an FGSM attack with a smaller step size $\delta < \epsilon$ and projects the attack back in the ϵ - L_∞ ball

$$\tilde{x}_{t+1} = \operatorname{Proj}\{\tilde{x}_t + \delta \operatorname{sgn} \nabla_x \ell(f(\tilde{x}_t, \theta), y)\}$$

DeepFool [18] find the nearest decision boundary to the data point x in the L_2 norm and pushes the perturbation beyond this boundary. It iteratively minimizes the classifier f around the input point until it produces a misclassification:

$$\begin{aligned} \tilde{x}_t &= \tilde{x}_{t-1} + r_t \\ \arg \min_{r_t} \|r_t\|_2 \\ f(\tilde{x}_t, \theta) + \nabla_x f(\tilde{x}_t, \theta)^T r_t &= 0. \end{aligned}$$

Carlini & Wagner (C&W) [19] in the L_∞ norm searches for the minimal adversarial perturbation producing a wrong classification. It solves the following optimization problem

$$\begin{aligned} \min \|x - \tilde{x}\|_\infty + c \cdot f(\tilde{x}, t) \\ \tilde{x} \in [0, 1]^m, \end{aligned}$$

where t is the target class.

C. \mathcal{R}_{v_1} and \mathcal{R}_{v_2} are equivalent as $k \rightarrow \infty$.

Let $R : \mathbb{R}^d \rightarrow \mathbb{R}^k$ be a random projection matrix, $R^\dagger : \mathbb{R}^k \rightarrow \mathbb{R}^d$ its pseudo-inverse and $x^\dagger := R^\dagger R x \in \mathbb{R}^d$ for $x \in \mathbb{R}^d$. For any given couple $(x, y) \in \mathbb{R}^d \times \mathbb{R}^K$ and projection matrix R let us define

$$\begin{aligned} \mathcal{T}_1(R) &:= \|\nabla_x \ell(f(x^\dagger, \theta), y)\|_2^2 \\ \mathcal{T}_2(R) &:= \|R \nabla_x \ell(f(x, \theta), y)\|_2^2. \end{aligned}$$

Proposition VI.1. *Let R be a random projection matrix whose elements are sampled from $\mathcal{N}(0, 1/k)$ and whose columns are orthogonal. Suppose that $x \in (\ker R)^\perp$ for all $x \in \mathbb{R}^d$. Then*

$$\mathbb{E}_R[\mathcal{T}_2(R)] = \frac{d}{k} \mathbb{E}_R[\mathcal{T}_1(R)].$$

Proof. A random matrix R and its pseudo-inverse R^\dagger induce the direct sum decompositions

$$\begin{aligned} \mathbb{R}^d &= (\ker R)^\perp \oplus \ker R \\ \mathbb{R}^k &= \text{rank } R \oplus (\text{rank } R)^\perp, \end{aligned}$$

where $\ker R$ is the kernel space of R and $\text{rank } R$ is the rank space of R . Moreover, $R|_{(\ker R)^\perp}$ is an isomorphism with inverse $R^\dagger|_{\text{rank } R}$ and $R^\dagger|_{(\text{rank } R)^\perp} \equiv 0$. If $x \in (\ker R)^\perp$, then $Rx \in \text{rank } R$ and $x^\dagger = R^\dagger R x = x$, therefore

$$\mathcal{T}_1(R) = \|\nabla_x \ell(f(x, \theta), y)\|_2^2$$

for all $x \in (\ker R)^\perp$.

Let $r_i \in \mathbb{R}^k$ be the orthogonal columns of R . Then $\mathbb{E}_R[\|R\|^2] = d/k$ and

$$\begin{aligned} \mathbb{E}_R[\|Rx\|_2^2] &= \sum_{i,j=1}^d \mathbb{E}_R[r_i^T r_j] x_i x_j \\ &= \sum_{i=1}^d \mathbb{E}_R[r_i^T r_i] x_i^2 = \frac{d}{k} \|x\|_2^2 \end{aligned}$$

for any $x \in \mathbb{R}^d$. In particular $\mathbb{E}_R[\mathcal{T}_2(R)] = \frac{d}{k} \mathcal{T}_1(R)$.

□

Notice that when \mathbb{R}^k is high dimensional we can assume that the columns of any random matrix are orthogonal [38].

We now prove that Prop VI.1 holds for an arbitrary $x \in \mathbb{R}^d$.

Proposition VI.2. *Let $\pi_R : \mathbb{R}^d \rightarrow (\ker R)^\perp$ be an orthogonal projection and $k = \dim(\ker R)^\perp$. Then, for any $x \in \mathbb{R}^d$ and $\epsilon > 0$*

$$P(\|\pi_R(x) - x\|_2^2 > \|x\|_2^2 \epsilon) \leq \left(1 - \frac{\epsilon}{\pi}\right)^k.$$

Proof. Suppose that $\{v_1, \dots, v_k\} \subset \mathbb{R}^d$ is a basis for $(\ker R)^\perp$, i.e. that $(\ker R)^\perp = \text{span}(v_1, \dots, v_k)$. Then any $x \in \mathbb{R}^d$ can be decomposed as $x = \pi_R(x) + u$, where $\pi_R(x) \in (\ker R)^\perp$ and $u \in \ker R$.

Let α_i be the angle between v_i and x . First, we observe that the projection $\pi_R(x)$ is smaller than any other projection on a single direction v_i

$$\begin{aligned} \|\pi_R(x) - x\|_2^2 &\leq \min_i \|\pi_{v_i}(x) - x\|_2^2 \\ &= \min_i (\|x\|_2^2 |\sin \alpha_i|) \\ &= \|x\|_2^2 \min_i |\sin \alpha_i|. \end{aligned}$$

For any choice of $\epsilon > 0$

$$\begin{aligned} P(\|\pi_R(x) - x\|_2^2 > \|x\|_2^2 \epsilon) &\leq P(\|x\|_2^2 \min_i |\sin \alpha_i| > \|x\|_2^2 \epsilon) \\ &= P(\min_i |\sin \alpha_i| > \epsilon) \\ &= \prod_i P(|\sin \alpha_i| > \epsilon) \\ &\leq \prod_i P(|\alpha_i| > \epsilon) \\ &= \left(1 - \frac{\epsilon}{\pi}\right)^k. \end{aligned}$$

□

Notice that $1 - \frac{\epsilon}{\pi} < 1$, so $\left(1 - \frac{\epsilon}{\pi}\right)^k \rightarrow 0$ as k goes to ∞ . Therefore, from VI.2 we get $x^\dagger = \pi_R(x) \approx x$ and

$$\nabla_x \ell(f(x^\dagger, \theta), y) \approx \nabla_x \ell(f(x, \theta), y)$$

as $k \rightarrow \infty$. This proves that Prop. VI.1 is true for an arbitrary $x \in \mathbb{R}^d$ in the limit.

Assuming that $k = O(d)$ as $k \rightarrow \infty$, e.g. $\frac{d}{k} \rightarrow M > 0$, the two regularization terms differ by a positive constant in the limit, i.e. they are equivalent if weighted w.r.t. M .

This proves the following theorem.

Theorem VI.1. *Let \mathcal{R}_{v_1} and \mathcal{R}_{v_2} be the regularization terms defined in Eq. 1 and Eq. 2, where $\mathcal{P} : \mathbb{R}^d \rightarrow \mathbb{R}^k$ is a random projection such that the elements of the orthogonal random matrix R are sampled from $\mathcal{N}(0, 1/k)$. If $k \in O(d)$ then $\mathcal{R}_{v_1} \approx \mathcal{R}_{v_2}$ as $k \rightarrow \infty$.*

Notice that this punctual property on \mathcal{R}_{v_1} and \mathcal{R}_{v_2} also holds in expectation over the training data when x is uniformly

sampled from a compact subset of \mathbb{R}^d . Therefore, the equivalence between the two regularization terms holds in practice on mini-batches of training data.