

# ROBUSTNESS OF SAAK TRANSFORM AGAINST ADVERSARIAL ATTACKS

Thiyagarajan Ramanathan<sup>†</sup>, Abinaya Manimaran<sup>†</sup>, Suya You<sup>‡</sup>, C-C Jay Kuo<sup>†</sup>

<sup>†</sup>University of Southern California, Los Angeles, California, USA

<sup>‡</sup>US Army Research Laboratory, Playa Vista, California, USA

## ABSTRACT

Image classification is vulnerable to adversarial attacks. This work investigates the robustness of Saak transform against adversarial attacks towards high performance image classification. We develop a complete image classification system based on multi-stage Saak transform. In the Saak transform domain, clean and adversarial images demonstrate different distributions at different spectral dimensions. Selection of the spectral dimensions at every stage can be viewed as an automatic denoising process. Motivated by this observation, we carefully design strategies of feature extraction, representation and classification that increase adversarial robustness. The performances with well-known datasets and attacks are demonstrated by extensive experimental evaluations.

**Index Terms**— Saak transform, Adversarial attacks, Deep Neural Networks, Image Classification

## 1. INTRODUCTION

It has been shown that deep learning based approach for image classification is vulnerable to adversarial attacks [1]. These attacks come in the form of adversarial inputs with carefully crafted perturbations added to the input samples. Such perturbations are small and imperceptible to humans, but can drastically cause the classification systems to misinterpret adversarial inputs, with potentially disastrous consequences where safety and security are crucial.

Saak (Subspace approximation with augmented kernels) transform, inspired by deep learning mechanism, is a new mathematical transform that is completely interpretable [2, 3]. The Saak transform is a variant of principal component analysis (PCA) that splits the positive and negative responses into two separate channels by kernel augmentation and resolves "sign confusion" ambiguity. This process facilitates the cascade of Saak transforms called multi-stage Saak transforms. Saak features at later stages have larger receptive fields, yet they are obtained in a one-pass feed-forward manner without any supervision and back propagation. In addition, inverse of Saak transform is possible, which allows the Saak feature representations to be transformed back to the image space for clearly visualizing, analyzing, and interpreting.

Saak transform has demonstrated its superior performance and utility in classifying hand-written digits under various noisy environments [4]. Furthermore, Saak transform has also been employed as a pre-processing step in image classification pipeline to defend adversarial attack [5].

In this work, we investigate the robustness of Saak transform against adversarial attacks towards high performance image classification. We develop a complete image classification system based on multi-stage Saak transform. We take advantage of the ocean of Saak coefficients available at every stage of multi-stage Saak transform. Careful selection of these features using cross-entropy leads us build a new Saak feature representation. The whole feature extraction and selection process is completely transparent and of extremely low complexity. In the Saak transform domain, clean and adversarial images demonstrate different distributions at different spectral dimensions. Selection of the spectral dimensions at every stage can be viewed as an automatic denoising process. Motivated by this observation, we design new strategies of feature extraction, representation and classification that increase adversarial robustness. The performances with well-known benchmark datasets and attacks are demonstrated by extensive experimental results.

## 2. RELATED WORK

One of the most interesting explorations in defense against adversarial attacks is done through adversarial training [6]. This aims in augmenting adversarial samples along with clean samples for simultaneous training. While these methods help in defending against particular adversarial attacks for which it is trained for, they fail to generalize. Also, this type of training takes longer time for convergence, hence needs to be trained for more epochs.

Adversarial detection involves detection of an adversarial sample before passing through the network. Adversarial samples can be detected using statistical tests [7], estimating Bayesian uncertainty [8], using noise reduction methods like scalar quantization and spatial smoothing filter [9]. Though these methods pave way to a good adversarial sample detection problem, these detectors still possess the risk of being fooled by the attacker.

Pre-processing based methods perform certain transforma-

tions on inputs to nullify the effect of adversarial attacks. Some examples are image cropping and rescaling, JPEG compression [10], feature squeezing by Bit-Depth-Reduction [11], and Total Variance Minimization [12]. Nonlinear, saturating neural networks are used in [13]. Gradient masking effect applies regularizers or smooth labels to attain output less sensitive to perturbed input [14].

Use of knowledge distillation when training networks can be used as defense against adversarial samples [15]. Reinforcement of network structure by using bounded ReLU activations help in enhancing stability to adversarial perturbations [16]. Pixel defend is used as an image purification process, as described in [17].

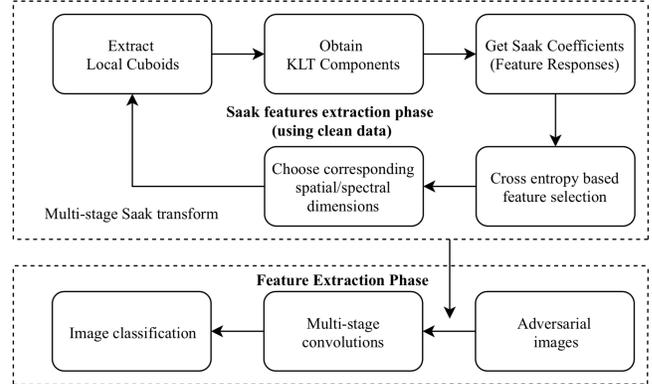
[5] applies lossy Saak transform to adversarially perturbed images as a pre-processing tool to defend against adversarial attacks. The method is based on the observation that outputs of Saak transform are very discriminative in differentiating adversarial examples from clean ones. Instead of using Saak transform as pre-processing tool, we apply multi-stage Saak transform to build a complete image classification pipeline and design new strategies of feature selection, representation and classification to defend against adversarial attacks.

### 3. SAAK TRANSFORM APPROACH

The Saak transform defines a mapping from a real-valued function defined on a three-dimensional cuboid consisting of spatial and spectral dimensions to a one-dimensional rectified spectral vector. It is presented as a new feature representation method. It consists of two main ideas: subspace approximation and kernel augmentation. For the former, we build the optimal linear subspace approximation to the original signal space via PCA or the truncated Karhunen-Love Transform (KLT) [18]. For the latter, we augment each transform kernel with its negative and apply the rectified linear unit (ReLU) to the transform output. This is equivalent to the sign-to-position (S/P) format conversion.

Figure 1 illustrates the developed saak transform pipeline. Specifically, it consists of 1) Extracting Local Cuboids from the images, 2) Obtaining KLT components, 3) Convoluting the images with the extracted kernels, 4) Calculating the cross entropy measures, 5) Selecting the best spatial/spectral components. We classify adversarial images using Saak transform. We extract kernels using clean images and follow the same procedure as we classify clean images. Saak kernels are used to extract the coefficients from attacked images. We classify adversarial attacked images after selecting features using our cross-entropy based method.

During the classification phase, input  $\mathbf{f}$  is convoluted with extracted Saak kernels to extract Saak features. Based on kernel size  $k_s$ , spatial resolution of feature responses reduce at every stage. Consider block of feature responses at any stage  $p$  for a single image as  $\mathbf{f}_p$  with dimension  $D_{p1} \times D_{p2} \times K_p$ . First two dimensions represent spatial dimension along ver-



**Fig. 1:** Saak transform consists two modules: kernel extraction and feature extraction. We use clean training images to extract kernels followed feature extraction. Feature responses can be used for any application.

tical and horizontal directions. The third one represents the spectral dimension of feature responses for an image. If there are  $N$  images in the training data, then total dimension of feature responses can be given as  $N \times D_{p1} \times D_{p2} \times K_p$ . Cross-entropy for feature responses is calculated at every index  $(i, j, k)$ , where  $(i, j)$  represents spatial location and  $k$  represents spectral dimension. Let  $C$  be the number of classes. Entropy at every location is given by,

$$H = \sum_{n=1}^N \sum_{c=1}^C y_{n,c} \log \frac{1}{p_{n,c}} \quad (1)$$

where

$$y_{n,c} = \begin{cases} 1, & \text{if } \mathbf{f}_p(n, i, j, k) \in c \\ 0, & \text{if } \mathbf{f}_p(n, i, j, k) \notin c \end{cases} \quad (2)$$

and  $p_{n,c}$  is the probability of  $n^{th}$  sample in class  $c$ . To obtain this, feature response values at  $(i, j, k)$  location across all images are taken. Histogram of these  $N$  values is calculated using a certain number of bins,  $B$ . From various experiments, we concluded that feature selection is stable irrespective of number of bins. We choose  $B = 10$  and proceed by getting

$$mc = (mc_1, mc_2, \dots, mc_B), \quad (3)$$

where  $mc_i$  represents maximum occurring class in bin  $i$ , and  $mc_i \in 1, 2, \dots, C$ . Probability  $p_{n,c}$  is determined as

$$p_{n,c} = \frac{\sum_{i=1}^B \mathbf{1}(mc_i = c)}{B} \quad (4)$$

At the end,  $D_{p1} \times D_{p2} \times K_p$  cross-entropy values will be computed at stage  $p$ . Lower the entropy value at a location, higher is the discriminant power. For every spectral dimension,  $D_{p1} \times D_{p2}$  pixels are ranked according to their entropy. The first few pixels with lowest cross-entropy values are retained, and others are made zero. This localizes

salient regions in an image. Similarly, spatially averaged cross-entropy for all spectral dimensions is obtained. From these average values, spectral dimensions are ranked, and first few  $K'_p$  with lowest average cross-entropy values are chosen. Thus spatially sparse feature responses with dimension  $D_{p1} \times D_{p2} \times K'_p$  are chosen at stage  $p$ . This is repeated at all stages of multi-stage Saak transform for classification.

#### 4. ADVERSARIAL ATTACKS AND DEFENSES

We consider three major adversarial attacks, against which we will evaluate our approach. The attacks are explained below.

**Fast Gradient Sign Method (FGSM) [19]:** This method computes adversarial image by adding a pixel-wide perturbation of magnitude in the direction of gradient. The perturbation is computed by  $\eta = \epsilon \text{sign}(\Delta_x J_\theta(x, l))$ , so each pixel is modified by  $x' = x + \eta$ . This value can be computed using back propagation. There is no bound on the modified value, hence the quality of the adversarial image greatly decreases.

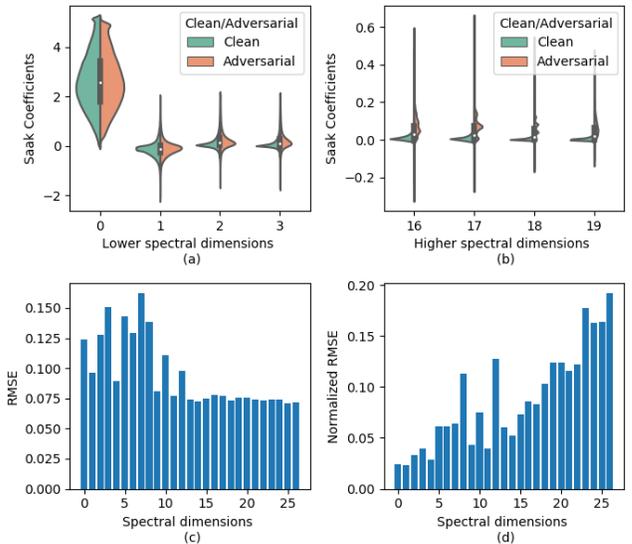
**Basic Iterative Method (BIM) [20]:** This method is an extension of FGSM, but with a limit on the value that a pixel can be modified. The change is limited, but the number of iterations of the attack are increased. Hence, for human eye, the BIM attacked images look less noisier when compared to FGSM attacks. The adversarial images are generated after multiple iterations.

**DeepFool (DF) [21]:** This attack is more carefully crafted when compared to FGSM and BIM. It computes the closest  $L_2$  projection distance to the decision boundary hyper-plane of adversarial example and input image. The perturbation is a function of the *argmin* of this distance. The perturbation is applied iteratively with smaller steps, hence the produced adversarial images do not look noisy to human eye.

The above three attacks are used to generate adversarial images and our Saak transform based image classification are directly applied to classify the attack images. The classification accuracy is evaluated against following state-of-the-art defending techniques.

**JPEG Compression [10]:** It has demonstrated that systematic JPEG compression can work as an effective pre-processing step in the classification pipeline to counter adversarial attacks. An important component of JPEG compression is its ability to remove high-frequency signal components, hence reducing high-frequency components with JPEG compression should contribute to adversarial attack removal without hurting the decision accuracy of clean data.

**Bit Depth Reduction [12]:** It is one of the feature squeezing techniques. Images often contain unnecessary features that can be exploited by adversarial attacks. Using less bit to discrete colors will make prediction more robust. The more bit to reduce, the more features are eliminated. In this work, two variants of bit depth reduction (4-bits and 5-bits) are used to pre-process the attack image.



**Fig. 2:** (a) and (b) show the lower and higher spectral distribution of Saak coefficients extracted from CIFAR-10. At higher dimensions, the distributions obtained from clean and attacked images are different. (c) and (d) show the RMSE and normalized RMSE between clean and FGSM attacked Saak coefficients in different spectral dimensions

**Non-Local Means (NL-Means) [22]:** Image denoising using NL-means before classification can remove adversarial noise and improve adversarial robustness. The NL-means not only compares the pixel value in a single point but the geometrical configuration in a whole neighborhood. This fact allows a more robust comparison than neighborhood filters, improving robustness of networks against adversarial attacks.

**Total Variance Minimization (TVM) [12]:** In images, noisy signals have a high total variation. TVM is an optimization technique where the variance of the pixels are reduced using  $L_2$  regularization. TVM retains more information such as edges when compared to other filtering techniques.

**Pixel Deflection [23]:** A random pixel is replaced by another random pixel in a local neighborhood. It works well due to the assumption that adversarial attacks rely on specific activation functions, i.e., only some pixels are manipulated to make the attack work. There are two variations of this method - pixel deflection with and without activation map. The activation map is used to determine the random pixel which is to be used to replace the target pixel in consideration.

#### 5. EXPERIMENTAL RESULTS

Extensive experiments are conducted on datasets MNIST, CIFAR-10 and STL-10. We provide in-depth experimental results for adversarial images classification and prove that classification using the proposed Saak features is adversarial robust in comparison with state-of-the-art defense mechanisms.

Adversarial Defense	FGSM	BIM	DF
No-defense	13.86%	18.95%	13%
JPEG(Q=90)	4.61%	9.52%	14.1%
Bit Depth Reduction (4-bit)	5.66%	9.65%	12.83%
Bit Depth Reduction (5-bit)	3.63%	8.87%	13.06%
Median Filtering (2x2)	4.04%	9.41%	8.77%
Median Filtering (3x3)	5.14%	10.2%	9.29%
NL-Means	5.27%	9.54%	14.52%
TVM	<b>2.92%</b>	8.81%	5.81%
Pixel Deflection (W/o RCAM)	4.54%	9.55%	16.44%
Pixel Deflection (W/ RCAM)	5.16%	9.63%	17.3%
Saak Transform	4.78%	<b>5.17%</b>	<b>3.79%</b>

**Table 1:** Robustness comparison on MNIST: The accuracy on clean images using the modified LeNet architecture is 99.2% and using SAAK transform is 99.4%. The value in each entry of the table is the drop in classification accuracy from clean and adversarial attacked images.

In the Saak transform domain, clean and adversarial images demonstrate different distributions at different spectral dimensions. Figure 2 shows distribution of Saak components belonging to first few spectral dimensions, followed by the distribution for higher spectral dimensions. Saak spectral components differ for both clean and adversarial images at higher dimension. We also show the normalized and the original RMSE (root-mean-squared-error) values between clean and FGSM adversarial samples in different spectral components. We can observe from Figure 2 (c) and (d) that clean and adversarial samples have different Saak coefficient values in high spectral dimensions. These results were obtained from first stage Saak transform of CIFAR-10 images using  $3 \times 3$  local cuboids.

We classify adversarial images using Saak transform. We extract kernels using clean images and follow the same procedure as we classify clean images. As shown in Figure 1 Saak kernels are used to extract the coefficients from attacked images. We classify adversarial attacked images after selecting features using cross-entropy based method.

Table 1 shows extensive comparison results of our Saak transform based classification for various attacked MNIST dataset with other state-of-art defense methods. The values across the table indicate the drop in classification accuracy from clean and adversarial attacked images, i.e. drop in classification accuracy  $\Delta = C_{clean} - C_{attack}$ . Lower the drop value is, better is the robustness of the classification model. Similarly Tables 2 and 3 shows the results for CIFAR-10 and STL-10 datasets.

From the results we can see our approach outperforms other adversarial defense methods. The classification accuracy drop for Saak transform features is very less, thanks for robustness of Saak transform to adversarial perturbations. As previously stated, the images classified with Saak transform are not subjected to any specially-crafted adversarial defense

Adversarial Defense	FGSM	BIM	DF
No-defense	83.95%	83.95%	83.95%
JPEG(Q=90)	74.69%	76.26%	<b>2.97%</b>
Bit Depth Reduction (4-bit)	82.08%	82.94%	60.35%
Bit Depth Reduction (5-bit)	81.95%	82.93%	60.8%
Median Filtering (2x2)	77.87%	77.19%	71.03%
Median Filtering (3x3)	82.55%	78.44%	79.99%
NL-Means	77.41%	75.27%	3.15%
TVM	76.18%	76.49%	72.35%
Pixel Deflection (W/o RCAM)	83.02%	83.61%	60.06%
Pixel Deflection (W/ RCAM)	82.8%	83.67%	60.01%
Saak Transform	<b>25.1%</b>	<b>26.1%</b>	4.16%

**Table 2:** Robustness comparison on CIFAR-10: The accuracy on clean images using pre-trained VGG-16 is 93.95% and using SAAK transform is 74.6%.

Adversarial Defense	FGSM	BIM	DF
No-defense	53.4%	32.77%	44.64%
JPEG(Q=90)	53.75%	34.39%	10.86%
Bit Depth Reduction (4-bit)	53.76%	33.21%	27.86%
Bit Depth Reduction (5-bit)	53.63%	33.03%	24.86%
Median Filtering (2x2)	53.22%	34.91%	28.94%
Median Filtering (3x3)	53.8%	35.08%	27.30%
NL-Means	53.79%	34.24%	22.94%
TVM	52.53%	31.72%	30.90%
Pixel Deflection (W/o RCAM)	53.46%	32.26%	27.54%
Pixel Deflection (W/ RCAM)	53.5%	32.25%	25.90%
Saak Transform	<b>15.36%</b>	<b>12.5%</b>	<b>4.55%</b>

**Table 3:** Robustness comparison on STL-10: The accuracy on clean images using pre-trained network is 74.86% and using SAAK transform is 63.5%.

method. The drop obtained for Deepfool attacked images is less when compared to the other attacks. For the MNIST dataset, Saak transform classification accuracy drop is lower than all other defenses for all the three attacks. Also the range of drop is much lesser for MNIST when compared to CIFAR-10 and STL-10, mainly because the attacks are more effective in complex datasets. Even in the complex datasets, the drop in accuracy using Saak features is less and much lower than most of the defenses. The results clearly show that adversarial perturbations can be effectively and efficiently defended using state-of-the-art Saak transform.

## 6. CONCLUSION

This paper studies the robustness of Saak transform against adversarial attacks without using any additional overhead to remove adversarial noise from images. We apply multi-stage Saak transform to build a complete image classification pipeline and carefully design each steps of feature selection, representation and classification to increase adversarial robustness. Extensive experimental evaluations demonstrate the benefits and utilities of Saak transform.

## 7. REFERENCES

- [1] Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio, “Adversarial examples in the physical world,” *preprint arXiv:1607.02533v4*, 2016.
- [2] C.-C. Jay Kuo and Yueru Chen, “On data-driven saak transform,” *J. Visual Communication and Image Representation*, vol. 50, pp. 237–246, 2018.
- [3] C.-C. J. Kuo, “The CNN as a Guided Multilayer RECOS Transform,” *IEEE Signal Processing Magazine*, vol. 34, pp. 81–89, May 2017.
- [4] Yueru Chen, Zhuwei Xu, Shanshan Cai, Yujian Lang, and C.-C. Jay Kuo, “A saak transform approach to efficient, scalable and robust handwritten digits recognition,” *2018 Picture Coding Symposium (PCS)*, pp. 174–178, 2018.
- [5] Sibong Song, Yueru Chen, Ngai-Man Cheung, and C.-C. Jay Kuo, “Defense against adversarial attacks with saak transform,” *arXiv preprint arXiv:1808.01785*, 2018.
- [6] Ali Shafahi, Mahyar Najibi, Zheng Xu, John P. Dickerson, Larry S. Davis, and Tom Goldstein, “Universal adversarial training,” *preprint arXiv:1811.11304*, 2018.
- [7] Kathrin Grosse, Praveen Manoharan, Nicolas Papernot, Michael Backes, and Patrick D. McDaniel, “On the (statistical) detection of adversarial examples,” *arXiv preprint arXiv:1702.06280*, 2017.
- [8] Reuben Feinman, Ryan R. Curtin, Saurabh Shintre, and Andrew B. Gardner, “Detecting adversarial samples from artifacts,” *arXiv preprint arXiv:1703.00410*, 2017.
- [9] Bin Liang, Hongcheng Li, Miaoqiang Su, Xirong Li, Wenchang Shi, and Xiaofeng Wang, “Detecting adversarial examples in deep networks with adaptive noise reduction,” *IEEE Transactions on Dependable and Secure Computing*, vol. PP, 05 2017.
- [10] Nilaksh Das, Madhuri Shanbhogue, Shang-Tse Chen, Fred Hohman, Siwei Li, Li Chen, Michael E. Kounavis, and Duen Horng Chau, “Shield: Fast, practical defense and vaccination for deep learning using jpeg compression,” in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, New York, NY, USA, 2018, KDD ’18, pp. 196–204, ACM.
- [11] Weilin Xu, David Evans, and Yanjun Qi, “Feature squeezing: Detecting adversarial examples in deep neural networks,” in *NDSS*. 2018, The Internet Society.
- [12] Chuan Guo, Mayank Rana, Moustapha Cisse, and Laurens van der Maaten, “Countering adversarial images using input transformations,” in *International Conference on Learning Representations*, 2018.
- [13] Aran Nayebi and Surya Ganguli, “Biologically inspired protection of deep networks from adversarial attacks,” *arXiv preprint arXiv:1703.09202*, 2017.
- [14] Nicolas Papernot, Patrick D. McDaniel, Arunesh Sinha, and Michael P. Wellman, “Towards the science of security and privacy in machine learning,” *arXiv:1611.03814*, 2016.
- [15] N. Papernot, P. McDaniel, X. Wu, S. Jha, and A. Swami, “Distillation as a defense to adversarial perturbations against deep neural networks,” in *2016 IEEE Symposium on Security and Privacy (SP)*, May 2016, pp. 582–597.
- [16] Valentina Zantedeschi, Maria-Irina Nicolae, and Ambrish Rawat, “Efficient defenses against adversarial attacks,” in *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security*, New York, NY, USA, 2017, AISec ’17, pp. 39–49, ACM.
- [17] Yang Song, Taesup Kim, Sebastian Nowozin, Stefano Ermon, and Nate Kushman, “Pixeldefend: Leveraging generative models to understand and defend against adversarial examples,” in *International Conference on Learning Representations*, 2018.
- [18] Henry Stark and John W. Woods, Eds., *Probability, Random Processes, and Estimation Theory for Engineers*, Prentice-Hall, Inc., Upper Saddle River, NJ, USA, 1986.
- [19] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy, “Explaining and harnessing adversarial examples,” *arXiv preprint arXiv:1412.6572*, 2014.
- [20] Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio, “Adversarial machine learning at scale,” *arXiv:1611.01236*, vol. abs/1611.01236, 2016.
- [21] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard, “Deepfool: A simple and accurate method to fool deep neural networks,” in *CVPR*. 2016, pp. 2574–2582, IEEE Computer Society.
- [22] Cihang Xie, Yuxin Wu, Laurens van der Maaten, Alan L. Yuille, and Kaiming He, “Feature denoising for improving adversarial robustness,” *arXiv:1812.03411v1*, 2018.
- [23] Aaditya Prakash, Nick Moran, Solomon Garber, Antonella DiLillo, and James A. Storer, “Deflecting adversarial attacks with pixel deflection,” in *CVPR*. 2018, pp. 8571–8580, IEEE Computer Society.