

On the Theoretical Robustness of BNNs against White-Box Attacks

Susim Mukul Roy
B20AI043

Soumik Roy
B20AI042

Abstract

Despite significant efforts, Deep neural networks(DNNs) are vulnerable to adversarial examples, which are one of the chief hurdles to the adoption of deep learning in safety-critical applications. Numerous Adversarial defense strategies have been proposed in recent works, however, they typically show limited feasibility due to compromise on efficiency. Theoretical findings have shown that, in suitably defined large data limit, BNNs posteriors are robust to gradient-based adversarial attacks. Thus, this study aims to demonstrate the theoretical robustness of Bayesian neural architectures against multiple white-box attacks and list empirical findings from the same.

The code can be found on our [Github Repository](#) and the dataset link can be found at [MNIST](#) and [CIFAR10](#).

1. Introduction

Adversarial attacks have been extensively studied since Szegedy et al.'s seminal work [9], and even deep learning models trained on very large data sets are shown to be vulnerable to such attacks [5]. As a result, developing machine learning models that are resilient to adversarial perturbations is an essential precondition for their application in safety-critical scenarios. Vulnerability to adversarial attacks of Deep Learning frameworks are claimed to be due to certain aspects such as the use of single-point estimates. Prior works [6] suggests that traditional neural networks are not well calibrated and the primary reasons for their vulnerability to adversarial attacks are overfitting and making overconfident predictions while allowing the networks to perform well on task in hand.

Attack strategies often evaluate gradients based on input points in order to identify directions of high losses [5], [8]. This variability can be intuitively linked to the uncertainty in the prediction, which has drawn attention to Bayesian Neural Networks (BNNs), which some recent studies have argued are a more robust deep learning paradigm [3], [4], [1], [7].

In BNNs, the posterior average of the gradients of the loss function vanishes in the large data limit, providing robustness against gradient-based attacks. The magnitude of gradients decreases as more samples are taken from the BNN posterior in various BNN architectures trained with Variational Inference (VI) [2]. The robustness property holds when the ensemble is drawn from the true posterior. However, empirical evidences in [2] shows that it is not likely to be true that the sole ensemble with zero averaging property of gradients is posterior distribution. Bayesian inference methods that provide cheaper approximates such as Variational Inference (VI) may exhibit such properties in practice.

As part of this project, Bayesian inference (i.e., Variational Inference) is incorporated with a few existing Deep Neural Network architectures to create BNN models. The same is checked for robustness against strong state-of-the-art white-box attacks. This project aims to test the listed theoretical hypothesis using six models AlexNet, VGG11, Simple Custom CNN, Resnet34 and Bayesian AlexNet, Bayesian VGG11 and Bayesian SimpleCNN and Bayesian Resnet34 using 2 datasets - MNIST and CIFAR-10. Then, the models are attacked with four state-of-the-art Gradient-based attacks: l_∞ -FGSM, l_∞ -PGD and l_2 -PGD and BIM and the plots for the models' respective test accuracies are presented.

2. Bayesian Neural Network

Bayesian neural networks are popular for their robustness to over-fitting, and their ability to easily learn small datasets. In the form of probability distributions, Bayesian approach further offers uncertainty estimates unlike traditional models. At the same time, this approach integrates it parameters using a prior probability, computing the average across various models during training, which prevents overfitting by providing a regularization effect to the network.

Although theoretically attractive, modelling such a

distribution over the filters (kernels) of a Deep Neural Network has been attempted never so successfully before. This is perhaps due to the huge number of parameters in such deep networks.

Further, inferring the posterior distribution of a Bayesian NN even in the case of a small number of parameters is not an easy task. The posteriors of a data distribution can be computed using Bayes' Theorem as shown below:

$$P(w|X, Y) = \frac{P(X, Y|w) \cdot P(w)}{P(X, Y)} \quad (1)$$

$$P(X, Y) = \int \dots \int_W P(X, Y, w) dw_0 \dots dw_D \quad (2)$$

However, computing the marginal ($P(X, Y)$) over a dataset, Equation 2 where D denotes the dataset, especially for an image dataset is computationally intractable due to the number of integrals involved. Thus often approximations to the model posterior are used instead, with the variational inference being a popular approach.

2.1. Variational Inference

In Bayesian modeling, we want to be able to sample from the posterior of models given the data. We want to be able to infer the latent variables (w) from observed data (X). A function is defined as $y = f(x)$ which gives a predictive output $y_1 \dots y_N$ for the given inputs $x_1 \dots x_N$ where N denotes the number of examples. In Bayesian inference, we aim to use a prior distribution over the space $p(f)$ of functions which represents our prior belief of functions expected to have generated the data. Problem faced in Equation 2 can be alleviated by conditioning the model using a finite set of random variables w . We assume for granted that the model depend on these variables alone and are sufficient to approximate the model. Hence, for a new input example x^* , the predictive distribution can be given by Equation 3.

$$P(y^*|x^*, X, Y) = \int p(y^*|f^*)p(f^*|x^*, w)p(w|X, Y)df^*dw \quad (3)$$

However, $p(w|X, Y)$ still remains an intractable distribution.

We use a variational distribution $q_\phi(w) \approx p(w|X, Y)$ in order to approximate it. $q_\phi(w)$ must be a simple distribution (eg. Gaussian) that can be easily computed and substitute the posterior distribution. To find the optimal $q_\phi(w)$, the distance between the surrogate distribution $q_\phi(w)$ and posterior distribution $p(w|X, Y)$ must be minimized.

We thus minimise the Kullback–Leibler (KL) divergence

which is a measure of similarity between two distributions.

$$\begin{aligned} KL(q_\phi(W)||p(w|X, Y)) &= \int_w q_\phi(w) \log\left(\frac{q_\phi(w)}{p(w|X, Y)}\right) dw \\ &= - \mathbb{E}_{w \sim q_\phi(w)} \log\left(\frac{p(w|X, Y)}{q_\phi(w)}\right) + \sum_{(x, y) \in D} \mathbb{E}_{q_\phi(w)} \log(p(Y|X, w)) \end{aligned} \quad (4)$$

In the equation above, $KL(q(W)||p(w|X, Y))$ is a distance metric and hence a positive term. $\log(p(X, Y))$ denotes the evidence of the data and is a constant negative term. The expectation term thus becomes a lower bound for the KL divergence and is termed as *Evidence Lower Bound*. Thus, our objective function, ELBO can be maximised w.r.t the variational parameters defining surrogate distribution $q(w)$ in order to minimize KL-divergence

$$\mathcal{L}(\phi) = -KL(q_\phi(w)||p(w|X, Y)) + \mathcal{L}_D(\phi) \quad (5)$$

$$\text{where } \mathcal{L}_D(\phi) = \sum_{(x, y) \in D} \mathbb{E}_{q_\phi(w)} \log(p(Y|X, w)) \quad (6)$$

$\mathcal{L}_D(\phi)$ is the expected log-likelihood. Minimizing the negative of Equation 5 is also known as an optimised version of Bayes by Backprop.

2.2. Local Reparametrisation Trick

There exists various algorithms for the gradient-based optimization of the ELBO (Equation 5) with q and p being differentiable. One of which, stochastic gradient variational Bayes (SGVB) method [1], parameterizes the random variable $w \sim q_\phi(w)$ as $w = f(\epsilon, \phi)$ where f is a differentiable function and $\epsilon \sim p(\epsilon)$ is a random noise vector drawn from noise distribution $p(\epsilon)$. The expected log-likelihood can be now be formed by an unbiased differentiable minibatch-based Monte Carlo estimator:

$$\mathcal{L}_D(\phi) \approx \mathcal{L}_D^{SGVB}(\phi) = \frac{N}{M} \sum_1^M \log(p(y^i|x^i, w^i = f(\epsilon, \phi))) \quad (7)$$

where $(x^i, y^i) \in D$ and M denotes the size of a mini-batch. The above is differentiable w.r.t ϕ and is unbiased. Thus, its gradients is unbiased as well, implying, $\nabla_\phi \mathcal{L}_D(\phi) \approx \nabla_\phi \mathcal{L}_D^{SGVB}(\phi)$

Stochastic gradient descent highly depends on the variance of the gradients. For minibatches M , the variances as stated in [1] is given by:

$$\text{Var}[\mathcal{L}_D^{SGVB}(\phi)] = N^2 \left(\frac{1}{M} \text{Var}[L_i] + \frac{M-1}{M} \text{Cov}[L_i, L_j] \right) \quad (8)$$

where, $L_i = \log(p(y^i|x^i, w^i = f(\epsilon^i, \phi)))$ and $\text{Var}[L_i] = \text{Var}_{\epsilon, x^i, y^i} \log(p(y^i|x^i, w^i = f(\epsilon, \phi)))$

In [5], an alternative estimator for which $\text{Cov}[L_i, L_j] = 0$, so that the variance of stochastic gradients scales as $1/M$

is proposed. Local Reparameterisation Trick, global uncertainty in the weights is translated into a form of local uncertainty that is independent across examples by not sampling ϵ directly but only sampling the intermediate variables $f(\epsilon)$ through which ϵ influences $L_D^{(SGVB)}(\phi)$.

Similar reparameterisation is further applied to sampling weights as well. The weights (also ϵ) influence the expected log likelihood only through the non-linear neuron activations $B = AW$ which are of much lower dimensions. Thus a low-cost Monte Carlo estimator is proposed directly for activations instead of sampling the Gaussian weights and then computing the resulting activations. For a factorized Gaussian posterior on the weights, the posterior for the activations (conditional on the input A) is also factorized Gaussian:

$$\begin{aligned} q_\phi(w_{i,j}) &= N(\mu_{i,j}, \sigma_{i,j}^2) \forall w_{i,j} \in \mathbf{W} \\ \implies q_\phi(b_{m,j} | \mathbf{A}) &= \mathcal{N}(\gamma_{m,j}, \delta_{m,j}) \\ \gamma_{m,j} &= \sum_{\text{neurons}} a_{m,i} \mu_{i,j}, \delta_{m,j} = \sum_{\text{neurons}} a_{m,i}^2 \mu_{i,j}^2 \end{aligned} \quad (9)$$

Activations can be directly computed using $b_{m,j} = \gamma_{m,j} + \sqrt{\delta_{m,j}} \zeta_{m,j}$ where $\zeta_{m,j} \sim \mathcal{N}(0, 1)$. It is shown that η is a MX no. neurons matrix hence we only sample M thousands instead of M millions as in original case.

This results in the subsequent equation for activations \mathbf{b} in a convolutional layer:

$$b_j = A_i * u_i + \epsilon_j \odot \sqrt{A_i^2 * \mu_i^2} \quad (10)$$

where i, j, w, h represent input, output layers and width and height of the image respectively. A_i represents the receptive field, $*$ represents a convolutional operation and \odot represents element-wise multiplication and $\epsilon \sim \mathcal{N}(0, 1)$.

3. Model Configurations:

3.1. Original Models:

We use two CNNs:

1. Alexnet
2. LeNet

We also test on a custom CNN model and *resnet34* and compared them with their bayesian counterparts. However, the increase in accuracies were nominal in both these cases and hence we only showed the results on two models. One can find the trained models on our github.

3.1.1 Other Configurations:

By default, all the models use ReLU activation function. We use Adam optimizer with the default configuration for all models and Loss is calculated using Cross Entropy Loss.

3.2. Bayesian Models

3.2.1 Activation Functions:

We apply the Softplus function in Bayesian models, because we want to ensure that the variance term never becomes zero. The Softplus activation function gives a smooth approximation of ReLU. In this application this has a high analytically important advantage that it never becomes zero for $x \rightarrow -\infty$, whereas ReLU becomes zero for any negative input.

$$\text{Softplus}(x) = \frac{1}{\beta} \cdot \log(1 + \exp(\beta \cdot x)) \quad (11)$$

Here β is set to 1 by default.

3.2.2 Objective Function

As mentioned in Equation 4, The objective function can be summarized as:

$$\begin{aligned} F(X, Y, \theta) &\approx \sum \log(q_\theta(w^i | X, Y)) - \log(p(w^i)) \\ &\quad + \log(p(Y | X, w)) \end{aligned} \quad (12)$$

3.2.3 Variational Posterior

$\log(q_\theta(w_i | X, Y))$ in Equation 4 is defined as the log of variational posterior sampled as a Gaussian distribution:

$$\log(q_\theta(w^i | X, Y)) = \sum_i \log(\mathcal{N}(w_i | \mu, \sigma^2)) \quad (13)$$

where σ is of dimension (input-channel, output-channel, (kernel-size, kernel-size)).

3.3. Parameter Initialization

A Gaussian distribution is used to store the mean and variance values instead of just one weight. Variance cannot be negative and *Softplus* as the activation function reinforces it. We express variance σ as $\sigma_i = \text{Softplus}(\rho_i)$ where ρ is an unconstrained parameter. A good method for variance initialization is still unknown. We perform gradient descent over $\theta = (\mu, \rho)$, and individual weight $w_i \sim \mathcal{N}(w_i | \mu_i, \sigma_i)$.

3.4. Adversarial Attacks

In the White-box adversarial attack scenario, the attacker is expected to know the model parameters. The above models are studied under 4 strong white-box attacks:

- $l_\infty FGSM$: Fast Gradient Sign Method (FGSM), using loss function, applies perturbations in the direction of the gradient to create an adversarial example x^{adv} written as:

$$x^{adv} = x + \epsilon \cdot \text{sign}(\nabla_x J(x, y_{true})) \quad (14)$$

where x is a data point, x^{adv} is the adversarially perturbed image, J is the loss function, y_{target} is the target/true label, and ϵ is the hyperparameter

- $PGD(l_\infty, l_2)$: Projected Gradient Descent (PGD) initializes uniform random perturbation then runs more iterations until finding an adversarial example. PGD creates a stronger attack than other previous iterative methods.

$$x_{i+1}^{adv} = \prod_{x+S} x_t^{adv} + \alpha \cdot \text{sign}(\nabla_x J(x, y_{true})) \quad (15)$$

where \prod clips the input at positions around $[x_t^{adv} - \epsilon, x_t^{adv} + \epsilon]$ (projection operator), α is the gradient step size and $x + S$ represents the perturbation set.

- $l_\infty BIM$: $FGSM$ is applied iteratively with small step size, and the pixel values of intermediate results are clipped after each step to ensure that they are in the ϵ -neighbourhood of the original image.

$$\begin{aligned} x_0^{adv} &= X \\ x_{i+1}^{adv} &= \text{Clip}_{X, \epsilon} x_i^{adv} + \alpha \cdot \text{sign}(\nabla_x J(x, y_{true})) \end{aligned} \quad (16)$$

4. Results

Now we compare the performance of the baseline models to their Bayesian counterparts. As robustness is generally defined in terms of mis-classification for CNNs, we provide results for the same. For small-epsilon attacks we use ϵ values 0 to 0.05 with an increment of 0.005 and for large-epsilon attacks we use ϵ values 0 to 0.5 with an increment of 0.05. ϵ is directly proportional to the perturbations in the adversarial image and hence the strength of the attack. We are performing all tests on 10000 images from the test sets of MNIST and CIFAR10 datasets.

Models (AlexNet, Bayesian AlexNet (BAlexNet), LeNet and Bayesian LeNet(BLeNet)) are compared on their performance on MNIST dataset and their respective test accuracies on the true dataset are reported in Figure 1. It can be observed that AlexNet and Bayesian LeNet have the best accuracies on original dataset. Further, they are attacked using the mentioned white-box gradient based attacks (Section 3.4) for small and large epsilon values and their test accuracies are compared.

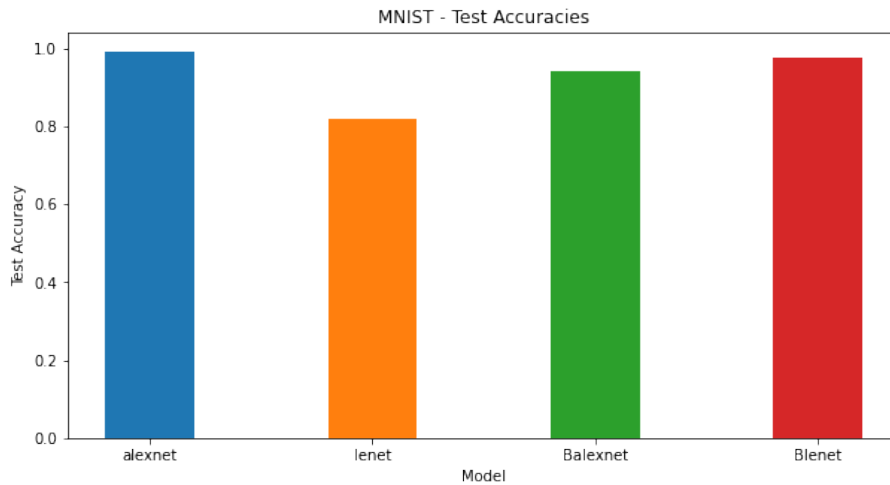


Figure 1. Test Accuracies on original MNIST dataset

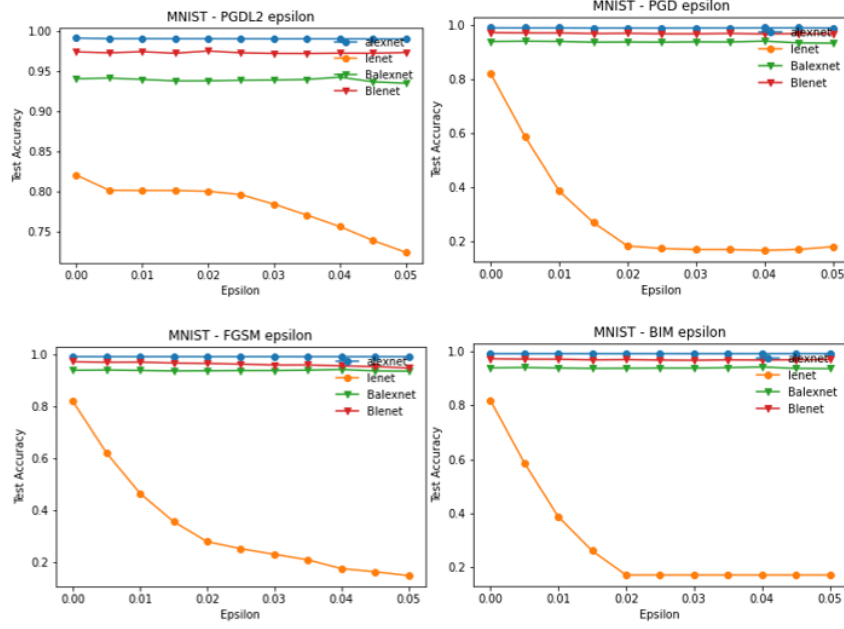


Figure 2. Test Accuracies on adversarial MNIST dataset for Small ϵ values

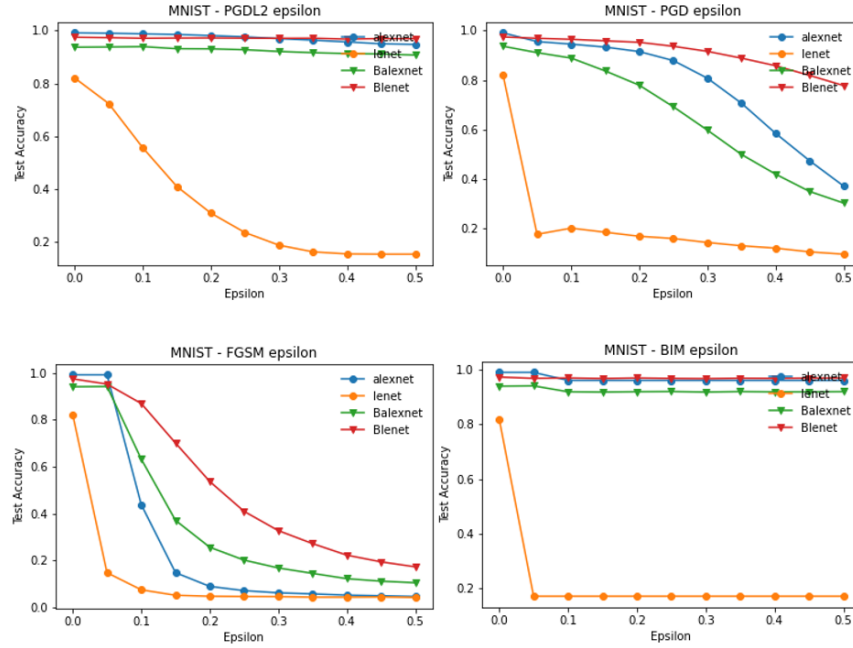


Figure 3. Test Accuracies on adversarial MNIST dataset for Large ϵ values

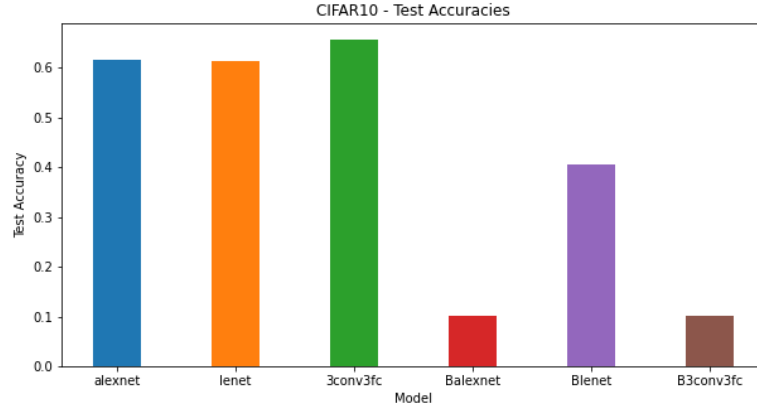


Figure 4. Test Accuracies on original CIFAR10 dataset

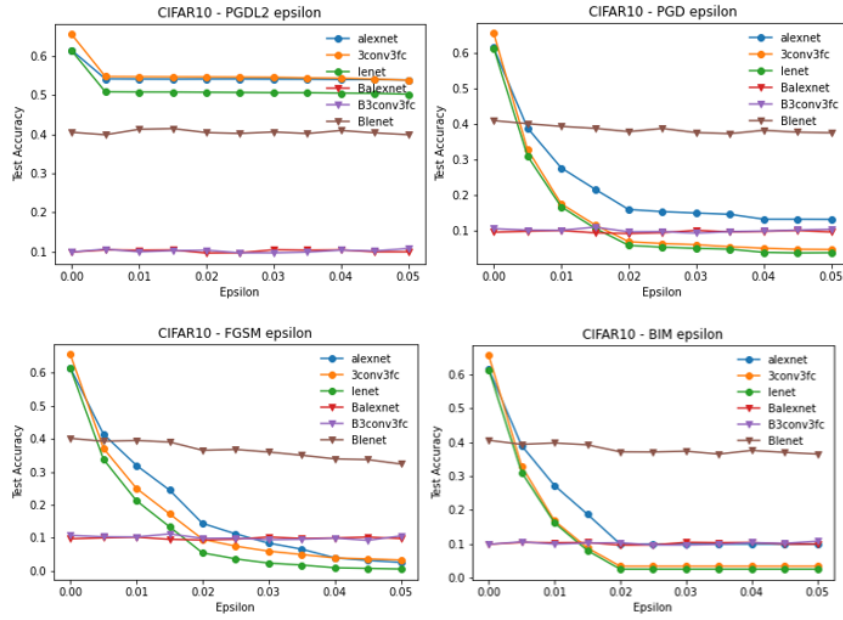


Figure 5. Test Accuracies on adversarial CIFAR10 dataset for Small ϵ values

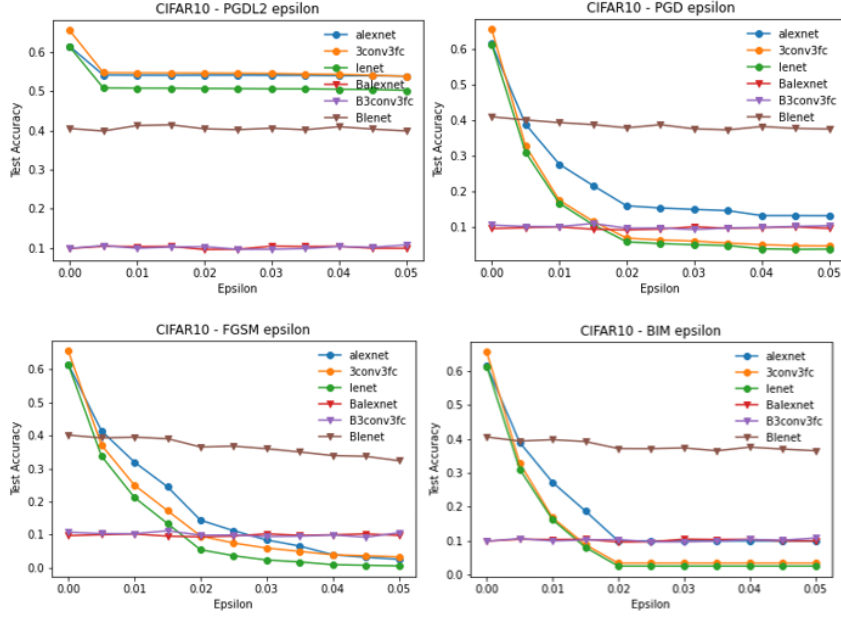


Figure 6. Test Accuracies on adversarial CIFAR10 dataset for Large ϵ values

For small perturbation values in Figure 2 frequentist LeNet model is highly vulnerable to all attacks, however, bayesian LeNet is robust to the same, however, AlexNet and the Bayesian models seem to be robust to all the attacks. For large perturbation values in Figure 3 Bayesian models consistently show better robustness on an average to the gradient-based attacks.

Models AlexNet, Bayesian AlexNet (BAlexNet), 3Conc3FC and Bayesian 3conv3fc (Simple CNN), LeNet and Bayesian LeNet (BLeNet) are compared on their performance on CIFAR-10 dataset and their test accuracies on the true dataset are reported in Figure 4. However, two bayesian models, BAlexNet and B3conv3FC fail to perform well on the original data. Further, they are attacked using the mentioned white-box gradient based attacks (Section 3.4) for small and large epsilon values and their test accuracies are compared.

For small perturbation values in Figure 5 despite having a much lower accuracy in original data Bayesian LeNet (BLeNet) model seems to retain its robustness consistently for most attacks. Although, BAlexNet and B3Conv3FC (bayesian Simple CNN) show poor performance, it is notable that the adversarial attacks did not cause their accuracy to go any lower unlike their frequentist counterparts. For large perturbation values in Figure 6, Bayesian LeNet model shows the average best performance consistently displaying its robustness against high perturbations. Frequentist models can be observed to be too vulnerable towards PGD, FGSM and BIM attacks from how their accuracies almost dropped to zero in untargetted attacks.

5. Conclusion

This study showed the probabilistic robustness for BNNs which takes both model and data uncertainty into account, and can be used to capture, among other properties, the robustness of the bayesian models towards gradient-based adversarial examples. Models that are data-driven and robust are an essential component to the construction of effective decision-making AI technologies. In this regard, it is noteworthy that Bayesian ensembles of Neural Networks have the capability to defend against a wide range of gradient-based adversarial attacks.

The results of this study have some significant limitations, yet they are promising. To begin with, Bayesian inference is extremely challenging in large non-linear models which was evident from our experiments on resnet34.

Secondly, this study focused on four popular gradient-based attack strategies which directly uses gradients in the presented empirical evaluation. Gradient-based attacks which are more complex in nature (Carlini-Wagner) also exist along with non-gradient (Jitter) based and query based attacks. Thus, evaluating towards more wide ranged attacks would confirm

Bayesian’s robustness.

References

- [1] Artur Bekasov and Iain Murray. Bayesian adversarial spheres: Bayesian inference and adversarial examples in a noiseless setting, 2018. [1](#)
- [2] Ginevra Carbone, Matthew Wicker, Luca Laurenti, Andrea Patane, Luca Bortolussi, and Guido Sanguinetti. Robustness of bayesian neural networks to gradient-based attacks, 2020. [1](#)
- [3] Reuben Feinman, Ryan R. Curtin, Saurabh Shintre, and Andrew B. Gardner. Detecting adversarial samples from artifacts, 2017. [1](#)
- [4] Yarin Gal and Lewis Smith. Sufficient conditions for idealised models to have no adversarial examples: a theoretical and empirical study with bayesian neural networks, 2018. [1](#)
- [5] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples, 2015. [1](#)
- [6] Diederik P. Kingma, Tim Salimans, and Max Welling. Variational dropout and the local reparameterization trick, 2015. [1](#)
- [7] Xuanqing Liu, Yao Li, Chongruo Wu, and Cho-Jui Hsieh. Adv-bnn: Improved adversarial defense through robust bayesian neural network, 2019. [1](#)
- [8] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks, 2019. [1](#)
- [9] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks, 2014. [1](#)