
Wait, what'd ya say? - Noise Aware DNN Training

An Examination of Generalizability, Robustness, and Quantization

Authors

Ryleigh Byrne (ryleigh.byrne@duke.edu)
Wyatt Focht (wyatt.focht@duke.edu)
Code Repository

Abstract

1 There has been previous exploration of using SAM to improve generalizability in
2 models. We seek to explore effective ways to improve a model's robustness against
3 weight perturbation. This paper explores various preventative methods to attempt
4 to mitigate the impact that weight perturbation may have on a model's performance.
5 The preventive methods we explore include Naive noise-aware training, SAM,
6 and SAM with multi-step weight perturbation training. Our experiments include
7 model resistance to various strengths of weight perturbation attacks, evaluation of
8 generalizability to new data on naive-noise, SAM, and multi-step SAM models,
9 and evaluation of training performance of naive-noise, SAM, and multi-step SAM
10 models.

11 1 Introduction

12 In the current literature surrounding deep neural networks (DNNs), there is a significant lack of
13 exploration into the effect that post-training, model-weight perturbation has on the performance of a
14 DNN. Consequently, we seek to analyze this effect by evaluating the impact that weight perturbation
15 has on model generalizability and quantization.

16
17 In addition to this, we seek to employ various preventative methods through which we can mitigate
18 any observed impact that weight-perturbation has on a model's generalizability & quantization and
19 compare these results to the respective performance of a model trained with no preventative measures
20 towards weight-perturbation.

21 1.1 Our Contribution

- 22 • Model-resistance to varying strengths of weight perturbation attacks
- 23 • Evaluation of generalizability to new data on naive-noise, SAM, multi-step SAM models
- 24 • Evaluation of quantization on naive-noise, SAM, multi-step SAM models
- 25 • Evaluation of training performance of naive-noise, SAM, multi-step SAM models

26 2 Related Works

27 Injecting noise or faults to DNN weights during training has been explored to promote robustness.
28 Adding weight noise and promoting redundancy by performing aggressive weight clipping have
29 both been shown to effectively increase DNN robustness, however, existing robustness methods
30 can lead to a decrease in DNN performance [1]. Sharpness-aware training has gathered interest in
31 increasing efficiency, performance, or understanding of sharpness-aware training [2]. An increase in
32 generalization performance due to SAM inspired our experiments[3], as we sought to examine various
33 preventative methods to mitigate the impact that weight perturbation has on a model's generalizability.

34 There have also been efforts to extend SAM for specific use cases. These use cases include data
35 imbalance settings [4] and quantization-aware training [5].

36 Recent studies have related the sharpness of the loss landscape with robustness to adversarial noise
37 perturbations [6]. Within this same study, robust over-fitting was examined. Robust over-fitting can be
38 defined as having high robustness to adversarial examples seen during training, but poor generalization
39 to new adversarial examples when testing [6]. To explore this phenomenon further, our experiments
40 focused on robustness against noisy weights during training. In relation to this study, we attempted
41 multiple training methods to try and achieve a model that has robustness weight-perturbation.

42 **3 Methodology**

43 **3.1 Standard Model Training**

44 This section serves to detail the base training parameters that are utilized while training a standard
45 Resnet-20 model on the CIFAR-10 dataset (50,000 training images and 10,000 testing images). The
46 following parameters were utilized:

- 47 • Epochs: 200
- 48 • Optimizer: Vanilla SGD
- 49 • Learning Rate: 0.1
- 50 • Momentum: 0.9
- 51 • L2 Regularization Strength: 1e-5

52 The above training specifications were determined utilizing best-practice criteria discussed both in
53 lecture and previous laboratory experimentation.

54 **3.2 Perturbing Weights with Gaussian Noise**

55 The purpose of this section is to describe the mechanism through which we are able to "simulate" a
56 weight-perturbation attack on a previously-trained model. We utilize this procedure on all models
57 generated throughout our research. The mechanism is as follows:

- 58 1. Acquire previously-trained model
- 59 2. Access weight parameters from the model
- 60 3. Iterate over all layers' weight parameters
- 61 4. If layer is Convolutional or Linear
 - 62 (a) Generate random Gaussian noise in the same shape as the current weight's shape
 - 63 (b) Multiply noise by chosen PERTURBATION STRENGTH
 - 64 (c) Add result back to original layer weights

65 By following the above procedure, it is possible to generate the requisite Gaussian noise for weight
66 perturbation. The perturbation strength was chosen as a factor of 0.1. During experimentation,
67 perturbation strengths between 0.0 and 0.2 were tested and evaluated. When utilizing a strength
68 of 0.2, the accuracy for the standard model was reduced to the equivalent of random guessing.
69 Additionally, a strength of 0.0 would result in no weight perturbation. Following this, a strength of
70 0.1 was settled upon and utilized for all further weight perturbation results achieved in this research.

71 **3.3 Naive Noise-Aware Training with Gaussian Noise**

72 The method through which we are able to achieve Naive Noise-Aware Training is very similar to the
73 mechanism described in 3.2. It follows the standard training procedure described in 3.1 except for
74 one difference: for every epoch during training,

- 75 1. Access weight parameters from the model
- 76 2. Iterate over all layers' weight parameters

- 77 3. If layer is Convolutional or Linear
- 78 (a) Generate random Gaussian noise in the same shape as the current weight's shape
- 79 (b) Multiply noise by chosen PERTURBATION_STRENGTH
- 80 (c) Add result back to original layer weights
- 81 4. Proceed normally with standard forward and backward passes
- 82 5. Repeat for each epoch

83 3.4 Adversarial Noise-Aware Training with SAM

84 Adversarial Noise-Aware Training with SAM is implemented by following the standard model
85 training described in 3.1 while also utilizing an additional set of steps sourced from a package found
86 in an external, highly-cited Github repository [SOURCE HERE]. The main differences in training
87 are as follows:

- 88 • Choice of Optimizer: This implementation creates a specialized SAM wrapper for the
89 standard SGD-with-momentum optimizer that is used in previous model training. The
90 neighborhood size hyperparameter, ρ , is set to 0.5 for this experiment.
- 91 • Training: The implementation of SAM utilizes two sets of forward passes, backward passes,
92 and optimizer steps. Consequently, per epoch in training, each of these three steps are
93 implemented twice.

94 3.5 SAM with Multi-Step Weight Perturbation

95 This method of training is similar to the methodology described in 3.4 except for the fact that ρ
96 is gradually increased throughout training. The schema through which this is implemented entails
97 starting with $\rho = 0.0$ and increasing ρ by a value of 0.02 every 10 epochs. By the conclusion of
98 training, ρ will have increased to a value of 0.4. This set of parameters was largely influenced by the
99 experimentation performed in [3] in which the research team varied ρ between 0.01 and 0.5 to find an
100 optimal value for the smoothest convergence while training with SAM. Consequently, ρ was chosen
101 to vary between a similar range for our experimentation.

102 3.6 Model Quantization

103 Model quantization is achieved programmatically by declaring a Resnet-20 model and specifying the
104 Nbits hyperparameter to the desired level of quantization. After this, the previously-trained model
105 weights can be loaded onto the initialized model. From here, the models can be tested as desired.

106 4 Experiments

107 4.1 Training

108 Within this section, the following set of results details further specifics of the training process for
109 each of the respective models. Specifically, **Figure 1.0** depicts the training times for each of the
110 created models as well as their relative training time compared to the standard model (= model
111 training time / standard model training time).

113 **Figure 2.0** depicts the relative training and validation loss curves for each of the models during their
114 200 epochs of training.

Model Type	Training Time (sec)	Ratio to Standard Time
Standard	4923	1
Naive-Noise	5374	1.09
SAM	6241	1.27
Multi-Step SAM	5862	1.19

Figure 1.0. Sample Model Training Times for 200 Epochs

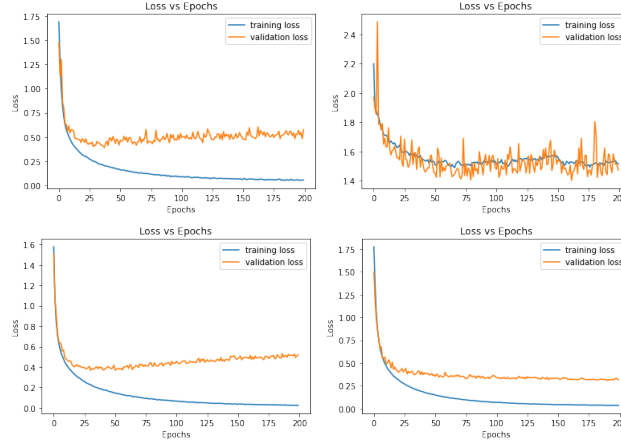


Figure 2.0. Loss Curves for Standard (top-left), Naive Noise (top-right), Standard SAM (bottom-left), and Multi-Step SAM training (bottom-right).

4.2 Generalizability

The testing accuracies displayed in **Figure 3.0** represent each fully-trained models accuracy when tested upon the 10,000 images set aside for testing.

Model Type	Testing Accuracy
Standard	0.8974
Naive-Noise	0.5031
SAM	0.9026
Multi-Step SAM	0.9109

Figure 3.0. Overall Accuracy for Each Model on Testing Data

4.3 Resistance to Gaussian Weight Perturbation

The results achieved in this section represent the relative accuracies each model achieves upon experiencing different weight perturbation attacks of varying strengths. The procedure for conducting this experiment involves loading the desired model for experimentation and following the methodology described in section 3.2 with the desired attack strength.

The x-axis of **Figure 4.0** represents the strength of the Gaussian noise that is used to perturb the model weights in this experiment. The corresponding value on the y-axis represents the accuracy of testing data that the model under examination was able to achieve given the specified level of weight perturbation.

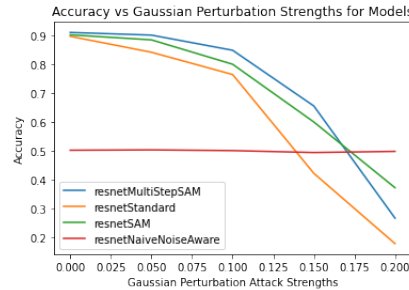


Figure 4.0. Resistance to Gaussian Weight Perturbation for Varying Attack Strengths

4.4 Robustness to Quantization

The results achieved in this section are the result of the methodology described in 3.6. Each model is initialized with various levels of quantization and evaluated for accuracy when performing classification on the CIFAR-10 testing set.

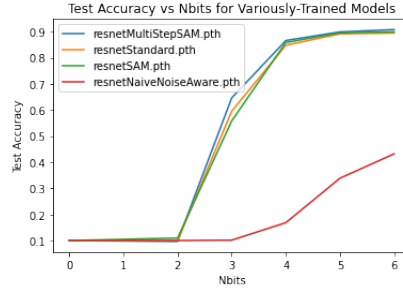


Figure 5.0. Model Accuracy under Various Levels of Quantization

5 Conclusion

5.1 Training

When analyzing training performance, the two most important features to understand are training time and training & validation loss.

5.1.1 Training Time

When analyzing each model's training time, we chose to compare each time to the base time achieved by the Standard Model, as this model was trained utilizing the typical specifications for a Resnet-20 model with no additional modifications to the training algorithm.

It can be observed that each of the other models took longer times to train than the Standard model. The Naive Noise model had the least difference in time when compared to Standard due to the fact that the training algorithm was largely the same as the Standard's with the addition of the slight computation needed to generate the random noise.

Both Standard SAM and Multi-Step SAM had significantly higher training times, as expected, due to the significant increase in computation per epoch in training (double the typical number of forward passes, backward passes, and optimizer steps).

5.1.2 Training and Validation Losses

When considering training and validation loss, the most significant difference that can be observed between the four models was that of the Naive Noise-Aware model. This difference is mostly in terms of accuracy.

For example, both the validation loss and training loss are very high for the Naive Noise-Aware model in comparison to the Standard model, Standard SAM model, and Multi-Step SAM model, as seen in **Figure 3.0**.

This effect can most likely be attributed to the extreme level of randomness that is introduced to the model weights during training. Having all model weights perturbed, even slightly, for every epoch, was bound to have a substantial impact.

5.1.3 Overfitting

Another important conclusion that can be drawn from the training experiments is that the Standard model, Standard SAM, and Multi-Step SAM models do not appear to be over-fitting, as the training loss decreases and approaches zero, but the validation loss stayed relatively consistent and did not approach zero. This could indicate that the models are not overfitting.

5.2 Generalizability

We chose to evaluate model generalizability on the ability of each model to perform accurate predictions on testing data that it had previously observed. The results of this can be found in **Figure 3.0**.

177 According to the results, it can be observed that the Multi-Step SAM achieves the highest level
178 of generalizability. This is not surprising as SAM has shown previous success in both model
179 generalization and robustness across various benchmark datasets in other works[3].

180 Additionally, the Naive-Noise model achieved the worst testing accuracy, which correlates with the
181 high testing loss and validation loss observed in Section 5.1.

182 *Acknowledgement*

183 It should, however, be acknowledged that the levels of accuracy achieved by each of the models
184 are slightly lower than what would be expected in a state-of-the-art implementation. This fact can
185 be attributed to the lack of usage of a decaying training rate throughout all training algorithms
186 implemented in this exploration. Instead, our training was intended to focus purely on decisions
187 made during optimization in each training epoch.

188 In the future, we would be interested in implementing a decaying training rate throughout optimization
189 in order to achieve accuracy closer to the state-of-the-art.

190 **5.3 Resistance to Gaussian Weight Perturbation**

191 **5.3.1 Defining Resistance**

192 In this context, we define "resistance" to weight perturbation as the ability of a model to maintain
193 a relatively constant accuracy when exposed to Gaussian weight attacks of various strengths. For
194 instance, we would say a model is resistant to Gaussian weight attacks if its base accuracy is equivalent
195 to its accuracy under any other realistic attack strength.

196 **5.3.2 Analysis**

197 In terms of results, the only model that was able to achieve any significant *resistance* was the
198 Resnet-20 model trained to be naively noise-aware of weight perturbation at each step during training
199 (resnetNaiveNoiseAware).

200 To be more specific, regardless of the strength of the Gaussian-weight attack, the resnetNaiveNoiseA-
201 ware's accuracy varied only slightly. However, it is important to note that although the Naive Noise
202 model achieves the best resistance to weight perturbation, the model also achieves the worst accuracy
203 (roughly 50% accuracy, regardless of attack strength).

204 In terms of accuracy, it can be observed that the other three models were able to outperform the
205 noise-aware model for attack strengths less than 0.150. However, at perturbation strength > 0.150 ,
206 the noise-aware model outperforms the standard model, and at perturbation strength > 0.175 , the
207 noise-aware model also outperforms both SAM and multi-step SAM.

208 **5.3.3 Future Work**

209 While the result we achieved is certainly a proof-of-concept for a possible type of training to produce
210 a model that is entirely resistant to weight perturbation attacks, we can clearly observe a significant
211 tradeoff when it comes to resistance and accuracy.

212 In the future, we would be interested in experimenting with combining the training schema of high-
213 accuracy models (SAM and multi-step SAM) with the training schema of the naive noise-aware
214 model. This could potentially have the effect of producing a model that is capable of not only
215 maintaining resistance to noise attacks, but also maintaining a high accuracy simultaneously.

216 Following this, we would also be interested in experimenting with the level of noise that the noise-
217 aware model is trained with as well as varying which layers of the model actually receive weight
218 perturbation during training (beginning layers versus ending layers).

219 **5.4 Robustness to Quantization**

220 None of the models showed exceptional robustness to quantization, as evidenced in **Figure 5.0**. The
221 Naive Noise performed exceptionally worse than the other three models. This could be due to this
222 model achieving a low accuracy overall, as discussed in Section 5.2.

223 References

- 224 [1] David Stutz, Nandhini Chandramoorthy, Matthias Hein, and Bernt Schiele. Bit error robustness for energy-
225 efficient DNN accelerators. Machine Learning and Systems, 2021.
- 226 [2] Jiawei Du, Hanshu Yan, Jiashi Feng, Joey Tianyihou, Liangli Zhen, Rick Siow Mong Goh, and Vincent Tan.
227 Efficient sharpness-aware minimization for improved training of neural networks. In International Conference
228 on Learning Representations, 2022.
- 229 [3] Foret, Kleiner, Mobahi, & Neyshabur. (n.d.). SHARPNESS-AWARE MINIMIZATION FOR EFFICIENTLY
230 IMPROVING GENERALIZATION. T ICLR 2021. <https://arxiv.org/pdf/2010.01412.pdf>
- 231 [4] Hong Liu, Jeff Z. HaoChen, Adrien Gaidon, and Tengyu Ma. Self-supervised learning is more robust to
232 dataset imbalance. In NeurIPS 2021 Workshop on Distribution Shifts: Connecting Methods and Applications,
233 2021.
- 234 [5] Jing Liu, Jianfei Cai, and Bohan Zhuang. Sharpness-aware quantization for deep neural networks. arXiv
235 preprint arXiv:2111.12273, 2021.
- 236 [6] Xu Sun, Zhiyuan Zhang, Xuancheng Ren, Ruixuan Luo, and Liangyou Li. Exploring the vulnerability of deep
237 neural networks: A study of parameter corruption. AAAI Conference on Artificial Intelligence, 2021,
- 238 [7] <https://github.com/davda54/sam.git>