
On Norm-Agnostic Robustness of Adversarial Training

Bai Li¹ Changyou Chen² Wenlin Wang³ Lawrence Carin³

Abstract

Adversarial examples are carefully perturbed inputs for fooling machine learning models. A well-acknowledged defense method against such examples is adversarial training, where adversarial examples are injected into training data to increase robustness. In this paper, we propose a new attack to unveil an undesired property of the state-of-the-art adversarial training, that is it fails to obtain robustness against perturbations in ℓ_2 and ℓ_∞ norms simultaneously. We discuss a possible solution to this issue and its limitations as well.

1. Introduction

Deep neural networks (DNNs) have achieved significant success when applied to a variety of challenging machine-learning tasks. For example, DNNs have obtained state-of-the-art accuracy on large-scale image classification (He et al., 2016b; ?). At the same time, vulnerability to adversarial examples, an undesired property of DNNs, has drawn attention in the deep-learning community (Szegedy et al., 2013; Goodfellow et al., 2014). Generally speaking, adversarial examples are perturbed versions of the original data that successfully fool a classifier. For example, in the image domain, adversarial examples are images transformed from natural images with visually negligible changes but that lead to wrong predictions (Goodfellow et al., 2014). The existence of adversarial examples has raised many concerns, especially in scenarios with a high risk of misclassification, such as autonomous driving.

To tackle adversarial examples, a variety of defensive methods against adversarial attacks have been proposed, yet most of them remain vulnerable to adaptive attacks (Szegedy et al., 2013; Goodfellow et al., 2014; Moosavi-Dezfooli et al., 2016; Papernot et al., 2016; Kurakin et al., 2016; Car-

lini & Wagner, 2017; Brendel et al., 2017; Athalye et al., 2018). One type of adversarial defense that demonstrated good performance against strong attacks is based on adversarial training (Goodfellow et al., 2014; Madry et al., 2017; Zhang et al., 2019). Adversarial training constructs a defense model by augmenting the training set with adversarial examples. Though this is a simple strategy, it has achieved a great success in adversarial defense.

The strength of attacks are commonly quantified by the ℓ_p distance between adversarial examples and natural examples. One desired property of a defense model is norm-agnostic, which requires a model to be robust against attacks constrained by a variety of norms. Recently, a more general attack mechanism called unrestricted adversarial attacks are introduced by Brown et al. (2018), where adversarial examples are not necessarily close to a natural image as long as they are semantically similar. To achieve robustness against unrestricted attacks, being norm-agnostic is a minimum requirement.

In this paper, we propose a new attack method and show adversarial training, the most successful adversarial defense models, is not norm-agnostic. Previously, it was reported both in (Madry et al., 2017) and (Zhang et al., 2019) that ℓ_∞ adversarial training is robust against ℓ_2 attacks. Our experiments, however, suggest they fail to defend against ℓ_2 and ℓ_∞ adversarial examples simultaneously.

2. Background and Related Work

Adversarial training constructs adversarial examples that are included to the training set to train a new and more robust classifier. This method is intuitive and has gained great success in defense (Szegedy et al., 2013; Goodfellow et al., 2014; Madry et al., 2017; Zhang et al., 2019). Madry et al. (2017) showed that iterative attacks during training yield strong defense models to white-box attacks (Athalye et al., 2018). More recently, another adversarial training based defense model (Zhang et al., 2019) has won the first place in the defense track of the NIPS 2018 Adversarial Vision Challenge (Brendel et al., 2018).

Although adversarial training has been so far one of the most successful defense methods, it has its limitations. In (Tramèr et al., 2017), it was pointed out that single-step

*Equal contribution ¹Department of Statistical Science, Duke University, Durham ²Department of CSE, University at Buffalo, the State University of New York ³Department of ECE, Duke University, Durham. Correspondence to: Bai Li <bai.li@duke.edu>, Lawrence Carin <lcarin@duke.edu>.

adversarial training, where single-step method (e.g., FGSM (Goodfellow et al., 2014)) is used for constructing adversarial examples, suffers from the “degenerate global minimum” issue and thus is not robust. To mitigate this issue, they propose ensemble adversarial training to improve the generalization of adversarial training. More recently, (Song et al., 2018) suggests using domain adaption as an improvement of ensemble adversarial training, leading to better robustness. However, both works only focus on single-step attack based adversarial training, while the most advanced adversarial training models are based on multi-step attacks. Tramèr et al. (2017) states that incorporating multi-step attacks during training could fix the degenerate-global-minimum issue. In this paper, we show multi-step adversarial training still suffers from this issue.

3. Preliminary

3.1. Adversarial Examples

Given a classifier $f : \mathcal{X} \rightarrow \{1, \dots, k\}$ for an image $\mathbf{x} \in \mathcal{X}$, an adversarial example \mathbf{x}_{adv} satisfies $\mathcal{D}(\mathbf{x}, \mathbf{x}_{\text{adv}}) < \epsilon$ for some small $\epsilon > 0$, and $f(\mathbf{x}) \neq f(\mathbf{x}_{\text{adv}})$, where $\mathcal{D}(\cdot, \cdot)$ is some distance metric, *i.e.*, \mathbf{x}_{adv} is close to \mathbf{x} but yields a different classification result. The distance is often described in terms of an ℓ_p metric, and in most of the literature ℓ_2 and ℓ_∞ metrics are considered.

One of the simplest and widely used attack methods is a single-step method, the Fast Gradient Sign Method (FGSM) (Kurakin et al., 2016), which manipulates inputs along the direction of the gradient with respect to the outputs:

$$\mathbf{x}_{\text{adv}} = \Pi_{\mathbf{x}+\mathcal{S}}(\mathbf{x} + \alpha(\nabla_{\mathbf{x}}L(\theta, \mathbf{x}, y))) \quad (1)$$

where $\Pi_{\mathbf{x}+\mathcal{S}}$ is the projection operation that ensures adversarial examples stay in the ℓ_p ball \mathcal{S} around \mathbf{x} .

Its multi-step variant FGSM^k is more powerful and has been shown to be equivalent to exploring adversarial examples with the projected gradient descent (PGD) method (Madry et al., 2017):

$$\mathbf{x}_{\text{adv}}^{t+1} = \Pi_{\mathbf{x}+\mathcal{S}}(\mathbf{x}_{\text{adv}}^t + \alpha(\nabla_{\mathbf{x}}L(\theta, \mathbf{x}, y))) \quad (2)$$

3.2. Adversarial Training

The motivation behind adversarial training is that finding a robust model against adversarial examples is equivalent to solving the saddle-point problem:

$$\min_{\theta} \max_{\mathbf{x}': \mathcal{D}(\mathbf{x}, \mathbf{x}') < \epsilon} L(\theta, \mathbf{x}', y)$$

The inner maximization is equivalent to constructing adversarial examples, while the outer minimization can be performed by standard training procedure for loss minimization.

Therefore, to achieve robustness to adversarial examples, adversarial training augments the training data with adversarial examples constructed during training, as an approximation to the inner maximization procedure.

Recently, Zhang et al. (2019) suggested using $L(\theta, \mathbf{x}, y) = L(f_{\theta}(\mathbf{x}), y) + \lambda L(f_{\theta}(\mathbf{x}), f_{\theta}(\mathbf{x}_{\text{adv}}))$ as the training loss, instead of $L(f_{\theta}(\mathbf{x}), y)$ and $L(f_{\theta}(\mathbf{x}_{\text{adv}}), y)$ alternatively used in (Madry et al., 2017).

3.3. Degenerate Global Minimum

In (Tramèr et al., 2017), it is pointed out that if \mathbf{x}_{adv} denotes the adversarial example generated by FGSM, adversarial training ideally results in a robust classification model θ^* such that:

$$L(\theta^*, \mathbf{x}_{\text{adv}}, y) \approx \max_{\mathbf{x}': \mathcal{D}(\mathbf{x}, \mathbf{x}') < \epsilon} L(\theta^*, \mathbf{x}', y) \approx 0$$

However, the training procedure may instead discover a “degenerate global minimum” θ^* :

$$L(\theta^*, \mathbf{x}_{\text{adv}}, y) \ll \max_{\mathbf{x}': \mathcal{D}(\mathbf{x}, \mathbf{x}') < \epsilon} L(\theta^*, \mathbf{x}', y)$$

In another word, the training procedure may generate a model that makes finding adversarial examples difficult for FGSM instead of a truly robust model.

(Tramèr et al., 2017) proposes two possible solutions for mitigating this issue. One is to use a strong multi-step adversarial training, such as PGD, at a cost of increased computational burden. Another is ensemble adversarial training, that is incorporating adversarial examples generated from multiple pre-trained classifiers that are different from the original one. In this way, they can decouple the construction of adversarial examples and the training to prevent “degenerate global minimum”, while still obtain the robustness of adversarial training due to the transferability of adversarial perturbations across models (Goodfellow et al., 2014).

4. Second Order Attack

We propose a new attack motivated by the “degenerate global minimum”. Note adversarial training is equivalent to solving the optimization problem:

$$(\hat{\theta}, \hat{\mathbf{x}}) = \underset{\theta}{\operatorname{argmin}} \underset{\mathbf{x}': \mathcal{D}(\mathbf{x}, \mathbf{x}') < \epsilon}{\operatorname{argmax}} L(\theta, \mathbf{x}', y).$$

Its solution is a saddle point of L , *i.e.*, the gradient ideally vanishes at $\hat{\mathbf{x}}$ as $\nabla_{\mathbf{x}}L(\hat{\theta}, \mathbf{x}, y)|_{\hat{\mathbf{x}}} = 0$. In practice, an adversarial training often finds $\hat{\theta}$ that makes the loss function flat in the neighborhood of a natural example \mathbf{x} , which leads to inefficient exploration for adversarial examples when performing attacks. This is intuitively the cause of “degenerate global minimum”.

Table 1. Accuracy of Various Adversarial Training Strategies against Various Attacks

Attacks	Madry’s			TRADES			Ensemble			ATDA		
	ℓ_2	ℓ_∞	Mix	ℓ_2	ℓ_∞	Mix	ℓ_2	ℓ_∞	Mix	ℓ_2	ℓ_∞	Mix
Natural	98.2%	98.8%	98.7%	99.4%	99.5%	99.4%	99.4%	99.0%	98.7%	99.2%	98.8%	99.0%
PGD (ℓ_2)	97.0%	92.8%	73.2%	91.7%	91.7%	90.4%	99.0%	65.3%	58.2%	98.8%	63.6%	57.9%
PGD (ℓ_∞)	0.4%	92.5%	82.2%	19.7%	95.6%	15.3%	0.0%	90.2%	81.4%	0.0%	62.6%	81.3%
S-O (ℓ_2)	96.6%	0.0%	18.3%	81.7%	3.2%	84.2%	98.9%	65.8%	58.4%	97.2%	64.0%	56.8%
S-O (ℓ_∞)	0.0%	91.3%	84.2%	16.9%	94.7%	14.5%	0.0%	88.9%	82.1%	0.0%	61.3%	83.4%

Most current attack methods construct adversarial examples based on the gradient of a loss function. However, according to the analysis above, first-order derivative is not effective for attacks if the defense model is trained adversarially. This motivates utilization of the second-order derivative of a loss function to construct adversarial examples.

To this end, assume the loss function is twice differentiable with respect to \mathbf{x} . Using Taylor expansion on the difference between the losses on the original and perturbed samples, and assuming the gradient vanishes, we have

$$L(\theta, \mathbf{x} + \mathbf{r}, y) - L(\theta, \mathbf{x}, y) \approx \frac{1}{2} \mathbf{r}^T H(\theta, \mathbf{x}, y) \mathbf{r}$$

with \mathbf{r} being the perturbation, and $H(\theta, \mathbf{x}, y)$ is the Hessian matrix of the loss function. Our goal is to find a small perturbation \mathbf{r} that maximizes the difference $L(\theta, \mathbf{x} + \mathbf{r}, y) - L(\theta, \mathbf{x}, y)$. Our idea is based on the observation that the optimal perturbation direction should be in the same direction as the first dominant eigenvector, $\mathbf{e}(\theta, \mathbf{x}, y)$, of $H(\theta, \mathbf{x}, y)$, that is $\mathbf{r} = \epsilon \frac{\mathbf{e}(\theta, \mathbf{x}, y)}{\|\mathbf{e}(\theta, \mathbf{x}, y)\|_2}$ for some constant $\epsilon > 0$. However, computing the eigenvectors of the Hessian matrix requires $O(I^3)$ runtime with I the dimension of the data. To tackle this issue, we adopt the fast approximation method from (Miyato et al., 2017), which is essentially a combination of the power-iteration method and the finite-difference method, to efficiently find the direction of the eigenvector. Based on this method, the optimal direction, denoted \mathbf{r}_{adv} , is approximated* by

$$\mathbf{r}_{adv} = \frac{g}{\|g\|_2}, \quad \text{with } g = \nabla_{\mathbf{x}} L(\theta, \mathbf{x}, y)|_{\mathbf{x}+\xi \mathbf{d}} \quad (3)$$

where \mathbf{d} is a randomly sampled unit vector and $\xi > 0$ is a manually chosen step size. In practice, \mathbf{d} is drawn from a centered Gaussian distribution and normalized such that its ℓ_2 norm is 1.

This procedure is essentially a stochastic approximation to the optimal second-order direction, where the randomness comes from \mathbf{d} . To reduce the variance of the approximation, we further take the expectation over the Gaussian noise, yielding $g = \mathbb{E}_{\mathbf{d} \sim N(0, \sigma^2 I)} [\nabla_{\mathbf{x}} L(\theta, \mathbf{x}, y)|_{\mathbf{x}+\mathbf{d}}]$. Note that choosing σ is equivalent to choosing the step size ξ in (3).

*Detailed derivations are provided in the Supplementary Material.

Finally, we construct adversarial examples by an iterative update via PGD:

$$\mathbf{x}^{t+1} = \Pi_{\mathbf{x}+S}(\mathbf{x}^t + \alpha \mathbf{r}_{adv}) = \Pi_{\mathbf{x}+S}(\mathbf{x}^t + \alpha \frac{g^t}{\|g^t\|_2}) \quad (4)$$

where $g^t = \mathbb{E}_{\mathbf{d} \sim N(0, \sigma^2 I)} [\nabla_{\mathbf{x}} L(\theta, \mathbf{x}, y)|_{\mathbf{x}+\mathbf{d}}]$. Intuitively, this method perturbs the example at each iteration and tries to move out of the local maximum in the sample space, due to the introduction of random Gaussian noise.

5. Experiments

We perform experiments on the MNIST data set to validate our claims on adversarial training.

The architecture of our model follows the ones used in (Madry et al., 2017). Specifically, the model contains two convolutional layers with 32 and 64 filters, each followed by 2×2 max-pooling, and a fully connected layer of size 1024. Image intensities are scaled to $[0, 1]$, and the size of attacks are also rescaled accordingly. In all the experiments, we bound the ℓ_2 norm less than 4.0 while the ℓ_∞ norm less than 0.3.

We evaluate PGD and proposed S-O attacks on four settings: adversarial training with PGD adversarial examples (Madry et al., 2017), Tradeoff-inspired Adversarial Defense (TRADES) (Zhang et al., 2019), ensemble adversarial training (Tramèr et al., 2017), adversarial training via domain adaption (Song et al., 2018).

Specifically, we first consider constructing adversarial examples with ℓ_2 and ℓ_∞ constraints during training respectively. Table 1 shows, as expected, that Madry’s model and TRADES successfully defend attacks with the same norm constraints. However, in spite of the fact that ℓ_∞ adversarial training stays robust against ℓ_2 PGD attack, S-O attack can effectively reduce the accuracy of both models when a different norm is used. This suggests that standard adversarial training is not norm-agnostic.

It is natural to wonder whether the issue will be fixed if two kinds of adversarial examples are both included during training. To this end, we conduct additional experiments with mixed adversarial examples, that is alternating between ℓ_2 and ℓ_∞ bounded examples for adversarial training. Using the mixed strategy, the accuracy on both attacks are no

longer reduced to almost zero, but the overall performance is still unsatisfying. We conclude that mixing adversarial examples barely helps improving norm-agnostic robustness.

According to our analysis, the poor performance of adversarial training is due to “degenerate global minimum”, therefore, we expect ensemble adversarial training could fix the problem, as suggested in (Tramèr et al., 2017). The results from Table 1 suggest ensemble adversarial training and domain adaption partially fixes the issue, although the accuracy against ℓ_2 attacks is still far from ideal.

In addition, we found two more interest phenomenons that can support our claims. Firstly, the relatively good performance of ensemble adversarial training implies that the vulnerability to adversarial perturbations with different norms is indeed caused by the “degenerate global minimum” issue, similar to the single-step adversarial training. Secondly, the performance of PGD and S-O attacks become similar for ensemble adversarial training model. This implies the effectiveness of S-O attacks compared to PGD attacks is due to exploitation of the “degenerate global minimum” issue.

In figure 1, we take a closer look at the behaviour of the attack methods by plotting the average ℓ_2 norms of the gradients of the loss function with respect to the adversarial examples during the construction processes. Specifically, we compute $\frac{1}{N_{\text{batch}}} \sum_{i \in I} \|\nabla_{\mathbf{x}} L(\theta, \mathbf{x}, y)|_{\mathbf{x}_i^{(t)}}\|_2$ for each t , where I is the index set of a batch.

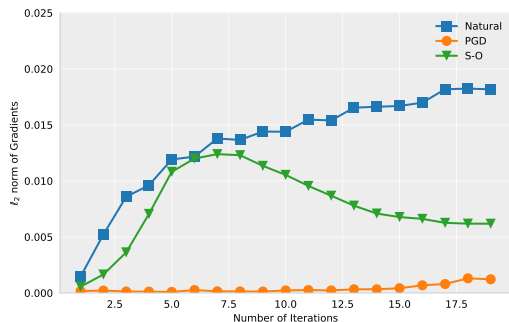


Figure 1. **S-O attack** Left: Average ℓ_2 norm of the gradients of the loss function for a batch in each iteration during adversarial attack. **Blue**: Naturally trained model attacked by PGD. **Orange**: TRADES attacked by PGD. **Green**: TRADES attacked by S-O.

We monitor this quantity for three settings: naturally trained model attacked by PGD, TRADES attacked by PGD, and TRADES attacked by S-O. The difference between the blue and orange lines show the ℓ_2 norms of the gradients of the adversarially trained model are much smaller than the ones of the naturally trained model under PGD attacks, validating our explanation in Section 4, that an adversarially trained model tends to make the loss function “flat” in the neighborhood of natural examples which makes PGD attacks

inefficient. The difference between the orange and green lines shows S-O attack is able to construct adversarial examples more efficiently by correctly finding the steepest direction, which explains why adversarially trained models are vulnerable to it.

Finally, one may argue that the perturbation size $\ell_2 = 4.0$ is too large that it violates the assumption that adversarial perturbations are visually negligible. We therefore perform S-O attack with perturbation size $\ell_2 \leq 2.0$, which results in accuracy 41.04%. We also illustrate some randomly selected perturbed adversarial examples that are misclassified by TRADES in figure 2.



Figure 2. **Above**: Natural examples from MNIST. The correct labels are 2,1,7,7,1,3. **Below**: Adversarial examples with perturbation size $\ell_2 \leq 2.0$. The adversarially trained model predictions are 0,4,9,2,4,8.

One can observe that although noticeable, there are only a limited amount of perturbations in the adversarial examples that do not change the semantic meaning of the images.

CIFAR-10 It is worth-noting that we do not observe similar results on CIFAR-10. We believe on CIFAR-10, it is difficult to reach even a “degenerate global minimum” for adversarial training due to the high dimensionality of the input space. This explains why adversarial training is still far from being perfectly robust even against PGD (Madry et al., 2017; Zhang et al., 2019).

6. Conclusion

In this paper, we show multi-step adversarial training models suffer from “degenerate global minimum” and thus are not norm-agnostic robust. Our proposed attack method is capable of constructing adversarial examples that reduces the accuracy of state-of-the-art adversarial training when different norms are used for training and attacking.

On the other hand, ensemble adversarial training can mitigate the issue thus should be considered as a standard procedure for adversarial training, even though they can only obtain moderate adversarial robustness.

In general, considering state-of-the-art results in adversarial defense are often achieved by adversarial training, we believe it is important to check the norm-agnostic robustness when designing adversarial defense models.

References

- Athalye, A., Carlini, N., and Wagner, D. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. *arXiv preprint arXiv:1802.00420*, 2018.
- Brendel, W., Rauber, J., and Bethge, M. Decision-based adversarial attacks: Reliable attacks against black-box machine learning models. *arXiv preprint arXiv:1712.04248*, 2017.
- Brendel, W., Rauber, J., Kurakin, A., Papernot, N., Veliqui, B., Salathé, M., Mohanty, S. P., and Bethge, M. Adversarial vision challenge. *arXiv preprint arXiv:1808.01976*, 2018.
- Brown, T. B., Carlini, N., Zhang, C., Olsson, C., Christiano, P., and Goodfellow, I. Unrestricted adversarial examples. *arXiv preprint arXiv:1809.08352*, 2018.
- Carlini, N. and Wagner, D. Towards evaluating the robustness of neural networks. In *Security and Privacy (SP), 2017 IEEE Symposium on*, pp. 39–57. IEEE, 2017.
- Goodfellow, I. J., Shlens, J., and Szegedy, C. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016a.
- He, K., Zhang, X., Ren, S., and Sun, J. Identity mappings in deep residual networks. In *European conference on computer vision*, pp. 630–645. Springer, 2016b.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pp. 1097–1105, 2012.
- Kurakin, A., Goodfellow, I., and Bengio, S. Adversarial machine learning at scale. *arXiv preprint arXiv:1611.01236*, 2016.
- Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- Miyato, T., Maeda, S.-i., Koyama, M., and Ishii, S. Virtual adversarial training: a regularization method for supervised and semi-supervised learning. *arXiv preprint arXiv:1704.03976*, 2017.
- Moosavi-Dezfooli, S.-M., Fawzi, A., and Frossard, P. Deep-fool: a simple and accurate method to fool deep neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2574–2582, 2016.
- Papernot, N., McDaniel, P., Wu, X., Jha, S., and Swami, A. Distillation as a defense to adversarial perturbations against deep neural networks. In *Security and Privacy (SP), 2016 IEEE Symposium on*, pp. 582–597. IEEE, 2016.
- Song, C., He, K., Wang, L., and Hopcroft, J. E. Improving the generalization of adversarial training with domain adaptation. *arXiv preprint arXiv:1810.00740*, 2018.
- Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., and Fergus, R. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- Tramèr, F., Kurakin, A., Papernot, N., Goodfellow, I., Boneh, D., and McDaniel, P. Ensemble adversarial training: Attacks and defenses. *arXiv preprint arXiv:1705.07204*, 2017.
- Zagoruyko, S. and Komodakis, N. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.
- Zhang, H., Yu, Y., Jiao, J., Xing, E. P., Ghaoui, L. E., and Jordan, M. I. Theoretically principled trade-off between robustness and accuracy. *arXiv preprint arXiv:1901.08573*, 2019.

A. Fast Approximate Method (Miyato et al., 2017)

Power iteration method (?) allows one to compute the dominant eigenvector \mathbf{r} of a matrix \mathbf{H} . Let \mathbf{d}^0 be a randomly sampled unit vector which is not perpendicular to \mathbf{r} , the iterative calculation of

$$\mathbf{d}^{t+1} = \frac{\mathbf{H}\mathbf{d}^t}{\|\mathbf{H}\mathbf{d}^t\|_2}$$

leads to $\mathbf{d}^t \rightarrow \mathbf{r}$. Given \mathbf{H} is the Hessian matrix of $L(\theta, \mathbf{x}, y)$, we further use finite difference method to reduce the computational complexity:

$$\begin{aligned} \mathbf{H}\mathbf{d} &\approx \frac{\nabla_{\mathbf{x}+\xi\mathbf{d}}L(\theta, \mathbf{x} + \xi\mathbf{d}, y) - \nabla_{\mathbf{x}}L(\theta, \mathbf{x}, y)}{\xi} \\ &= \frac{\nabla_{\mathbf{x}+\xi\mathbf{d}}L(\theta, \mathbf{x} + \xi\mathbf{d}, y)}{\xi} \end{aligned}$$

where $\xi > 0$ is the step size. If we only take one iteration, it gives an approximation that only requires the first-order derivative:

$$\mathbf{r} \approx \frac{\mathbf{H}\mathbf{d}}{\|\mathbf{H}\mathbf{d}^t\|_2} \approx \frac{\nabla_{\mathbf{x}+\xi\mathbf{d}}L(\theta, \mathbf{x} + \xi\mathbf{d}, y)}{\|\nabla_{\mathbf{x}+\xi\mathbf{d}}L(\theta, \mathbf{x} + \xi\mathbf{d}, y)\|}$$

which gives equation 3.