

A THEORETICAL FRAMEWORK FOR ESCAPING LOCAL OPTIMA IN MSE TOWARD GLOBAL CONVERGENCE

Anonymous authors

Paper under double-blind review

ABSTRACT

Deep learning models are trained by minimizing loss functions such as mean squared error (MSE) or cross-entropy, but these objectives are highly non-convex. As a result, optimization often encounters local optima, saddle points, or sharp valleys that hinder convergence and generalization. Although many heuristic approaches, such as momentum, Adam, help mitigate these issues, they provide limited theoretical understanding. In this work, we present a theoretical study of the optimization of MSE. We first provide a mathematical characterization of local optima under MSE and contrast them with those of cross-entropy, identifying when and how they arise. Building on this analysis, we introduce a modified optimization algorithm that explicitly accounts for these properties. Unlike heuristic methods, our approach offers theoretical guarantees for avoiding spurious local traps. Our experiments show that the proposed method reliably avoids local optima and converges more effectively than existing optimizers in MNIST, CIFAR10 and CIFAR100 with simple CNN. Our work provides both new insight into MSE optimization for training deep networks.

1 INTRODUCTION

Deep learning models are trained by minimizing a loss function, which measures how far the model's predictions are from the ground truth (Hinton et al. (2012)). However, commonly used loss functions such as cross-entropy (Mannor et al. (2005); Jamin & Humeau-Heurtier (2019)) or mean squared error (MSE) (Zhou et al. (2022)) – are non-convex due to the model's activations (Jain et al. (2017)). This means they do not have only a single global optimum. Instead, they have complex landscape with many local optima such as saddle points (Choromanska et al. (2015)). Because of this, many studies (Cheridito et al. (2024); Nguyen (2021)) have focused on overcome the local optima where the optimization stops even though the model is not yet at the ground truth.

Over the years, researchers have tried many ways to reduce this problem (Hinton et al. (2012); Zhang et al. (2022)). One family of approaches works by changing how gradients are used in training process. For example, momentum (Sutskever et al. (2013); Jelassi & Li (2022)) makes the trend to go to the certain direction, and make it easier to escape small valleys. Adam (Kingma & Ba (2014)) further adapts the learning rate for each parameter, so that the optimizer can avoid local optima and train the model more stably. These methods have become extremely popular in deep learning because they are effective for various cases.

Another kind of works takes a different view: instead of only considering the gradients, they focus on understanding the shape of the loss landscape itself (Zhou et al. (2022)). For instance, researchers have shown that in very high-dimensional spaces, it is not local minima but saddle points that are the major challenge, because they appear more often (Dauphin et al. (2014)). They proposed the saddle-free Newton method, which uses curvature information from the Hessian to escape from saddle points. Later, some studies Agarwal et al. (2017a;b) developed a second-order algorithm that has theoretical guarantees for finding approximate local minima faster than conventional gradient descent.

Despite this progress, most existing methods are still heuristic (Tian & Fong (2016); Kaveh & Mesgari (2023)). In fact, there is still no simple, precise mathematical description of what the local optima of MSE actually appear. Interestingly, Reddi et al. (2019) shows that Adam can fail to converge, despite its widespread use. As a result, they often depend on heuristic strategies, such as

054 keeping the gradient moving artificially so the optimizer does not fall in local optima. These meth-
 055 ods may work in some cases, but they sometimes cause lower performance in a few cases(Zhou et al.
 056 (2020)). Intuitively, if the global optimum is very sharp and narrow, then the momentum can simply
 057 pass by it without entering. This shows that escaping local optima is not a one-size-fits-all problem,
 058 but depends heavily on both the model architecture and the loss function.

059 In this paper, we focus on MSE loss and make two main contributions. First, we present a mathe-
 060 matical analysis of local optima in MSE compared to Cross-entropy, and what properties they have.
 061 Second, building on this analysis, we propose a new loss function to optimize while avoiding such
 062 local traps in a principled way. Unlike conventional heuristic approaches, our method provides clear
 063 guarantees and adapts naturally to the conventional training process with Pytorch.

064 By combining theoretical insight with practical optimization, our study aims to characterize local
 065 optima through optimality condition of objective, and proposes a loss function to modify the direc-
 066 tion of optimization. This is especially important because it is grounded in a solid mathematical
 067 basis. The main implication is that if this method converges but the loss value is not zero, this is due
 068 not to local optima of the loss, but rather to the limited expressive power of the model.

070 2 OPTIMALITY CONDITION OF CRITERIA

071 In this section, we define criteria and differentiate the objective functions to derive the optimality
 072 conditions.

073 2.1 CROSS-ENTROPY LOSS

074 First, we consider the cross-entropy loss function, which is widely used in classification tasks:

$$075 L_{CE}(w) = \sum_{i=1}^c -y_i \log f_i(w). \quad (1)$$

076 By differentiating $L_{CE}(w)$ with respect to w , we obtain the optimality condition for cross-entropy:

$$077 \sum_{i=1}^c \frac{y_i}{f_i(w^*)} \nabla_w f_i(w^*) = 0. \quad (2)$$

078 If y is encoded using one-hot encoding, then all labels are zero except for the true label. Therefore,
 079 it is not affected by the singularity of the Jacobian.

080 2.2 MEAN SQUARED ERROR (MSE) LOSS

081 Next, we begin by defining the Mean Squared Error (MSE) loss with respect to the model parameters
 082 $w \in \mathbb{R}^d$:

$$083 L_{MSE} = \frac{1}{2} \|y - f(w)\|^2, \quad (3)$$

084 where $y \in \mathbb{R}^c$ is the target vector and $f(w) \in \mathbb{R}^c$ is the model output.

085 To find the stationary point, we take the gradient of $L(w)$ with respect to w . The first-order optimal-
 086 ity condition is given by

$$087 D_w f(w^*)^\top (y - f(w^*)) = 0, \quad (4)$$

088 where $Df(w^*) \in \mathbb{R}^{c \times d}$ denotes the Jacobian matrix of f at w^* .

089 Equation (2) can be equivalently written as the summation over each output dimension:

$$090 \sum_{i=1}^c (y_i - f_i(w^*)) \nabla_w f_i(w^*) = 0. \quad (5)$$

091 This condition characterizes the optimality of the MSE loss. It is similar to (2), since both conditions
 092 represent a linear combination of class-wise gradients. According to optimality condition of MSE,
 093 the training is finished in following three cases.

1. $y_i = f_i$ for $i = 1 \dots c$
2. $\nabla_w f_{i=1 \dots c} = 0$
3. The linear combination of $\nabla_w f_{i=1 \dots c}$ is zero

The first case corresponds to the global optimum because (3) is strictly greater than zero. The second and last cases should represent local optima of the MSE. The second case is minor because each gradient of activations are not zero, except for Relu. Even when using ReLU, the case where all gradients become zero is a special case and is not a primary consideration. On the other hand, the last case is more critical: once the parameter vector falls into the null space of the Jacobian matrix, the gradient of the loss becomes zero, and momentum cannot restore the gradient from the null space.

3 MODIFIED MSE

As we have seen, the MSE loss itself can lead to local optima, so we propose a modified MSE formulation to mitigate this issue. We refer to this modified loss as QMSE. In the next, we present some properties of QMSE.

3.1 QMSE

In Equation (4), local optima may arise due to the null space of the Jacobian. To address this, we modify the MSE by transforming the residual vector $(y - f)$. Specifically, if we use a quadratic norm instead of the standard 2-norm, the loss and the gradient of the loss becomes

$$\begin{aligned} L_{QMSE} &= \frac{1}{2} \|y - f(w)\|_Q^2, \\ \nabla_w L_{QMSE} &= -D_w f(w)^\top Q(y - f). \end{aligned} \quad (6)$$

Now, we can handle the residual vector $(y - f)$ by Q and should derive good Q for $\nabla_w L_{QMSE}$.

To address this, we setc

$$Q = (y - f + v)(y - f + v)^\top, \quad (8)$$

where v is the random unit vector subject to $D_w f^\top v \neq 0$. If any matrix has xx^\top shape, then the matrix has only one eigenvalue, $\|x\|_2^2$. And the corresponding eigenvector is x . So Q has an eigenvalue, $\|y - f + v\|_2^2$ and the corresponding eigenvector, $(y - f + v)$. Although it is possible to find an analytic vector v that avoids the null space of the Jacobian. However, explicitly considering the Jacobian needs large computational cost. A more practical approach is to use random sampling, and if the sampled vector still lies in the null space, simply resample.

3.2 PROPERTIES

In this section, we show some properties of QMSE using the reformulation in (7). First, we plug in Q from (8) into (7) and simplify the expression. Then, (7) becomes

$$\nabla_w L_{QMSE} = D_w f^\top v (\|y - f\|_2^2 + \langle v, y - f \rangle). \quad (9)$$

$\langle a, b \rangle$ denotes the dot product between a and b . The direction of training is determined by the Jacobian $D_w f$ and v , while the other terms inside the parentheses determine the magnitude of the vector(they are a scalar value).

We list some basic properties of QMSE in some cases at local optima with respect to v .

- *Near the global optimum*, $\|y - f\|_2^2$ is small, and the transformed residual vector $Q(y - f)$ can escape the null space more carefully. In the opposite case, it should escape more aggressively.
- *When v and $(y - f)$ are almost parallel*, $\langle v, y - f \rangle$ is large, and the transformed residual vector $Q(y - f)$ can escape the null space more aggressively. In the opposite case, it should escape more carefully.

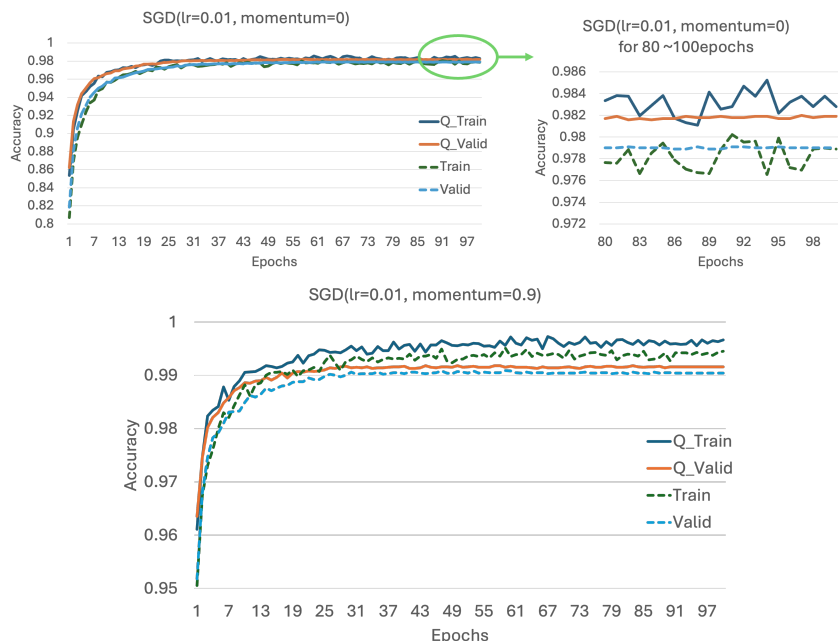


Figure 1: Training and test accuracy of the simple CNN on the MNIST dataset. Our method is denoted as Q_Train and Q_Valid, while the others are trained with conventional MSE. The upper figure shows the four accuracy curves obtained using SGD with zero momentum, while the bottom figure shows the results with a momentum of 0.9. The top-right plot provides a zoomed view of the accuracy between epochs 80 and 100. In all MNIST experiments, the learning rate was set to 0.01 and 100 epochs.

4 EXPERIMENT

In this section, we aim to verify whether our method achieves better performance than standard methods by reaching the global optimum in train set. We employ a simple CNN consisting of two convolutional layers followed by two fully connected layers. Using the MNIST, CIFAR-10, and CIFAR-100 datasets, we evaluate how well our method performs under various datasets and optimizer settings, based on both training and test accuracy. In particular, we compare our method with momentum. We do not use the loss for comparison, since QMSE is based on the Q -norm while MSE is based on the ℓ_2 -norm. This makes both numerical values not directly comparable. However, due to the optimality conditions, both losses share the same local and global optima. Therefore, we use accuracy as the evaluation metric.

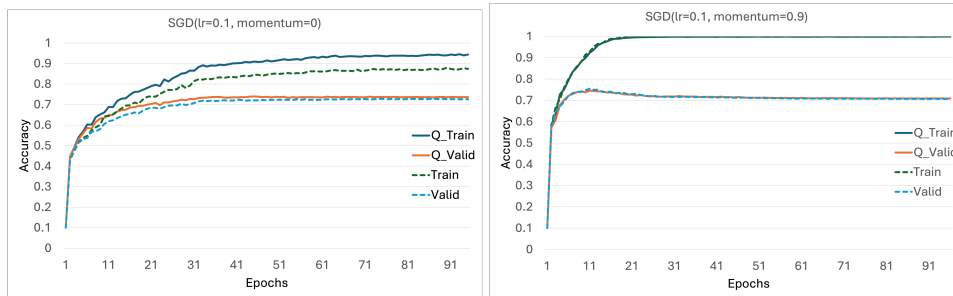
4.1 MNIST

We evaluate our method on the MNIST dataset. Figure 1 shows the training and test accuracy of QMSE and MSE. The learning rate is set to 0.01 due to the rapid convergence speed. QMSE and MSE are achieved nearly 1.0 accuracy. However, based on the training accuracy, QMSE converges to a lower optimal point in loss, so we can find that it outperforms in both training and validation accuracy. When we use QMSE with momentum, it achieves the highest performance on MNIST in our experiments.

4.2 CIFAR-10

Our goal is to overcome the local optima of MSE, and this is supported by the training accuracy. On CIFAR-10 dataset, QMSE and MSE achieve similar validation accuracy, while train accuracy shows a difference. In the right panel of Figure 2, where SGD is used with momentum, both methods achieve nearly 1.0 training accuracy. In contrast, when SGD has no momentum, QMSE shows

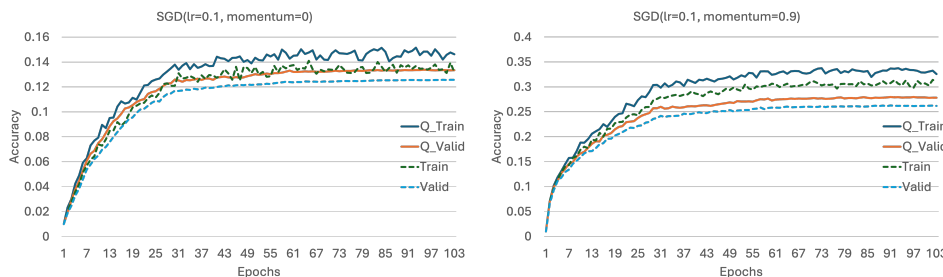
216
217
218
219
220
221
222
223
224
225



226
227
228
229
230

Figure 2: Training experiments of the simple CNN on the CIFAR-10 dataset. The upper panel shows the accuracy curves obtained using SGD with zero momentum, and the lower panel shows the results with a momentum of 0.9. For all CIFAR-10 experiments, we set the learning rate to 0.1 and trained for 100 epochs.

231
232
233
234
235
236
237
238
239
240



241
242
243
244
245

Figure 3: Training results of the simple CNN on the CIFAR-100 dataset. The upper panel shows the accuracy curves obtained using SGD with zero momentum, and the lower panel shows the results with a momentum of 0.9. In all CIFAR-100 experiments, the learning rate was set to 0.1, and training was performed for 100 epochs.

246
247
248
249
250

significantly higher training accuracy. These results indicate that QMSE, combined with momentum or not, helps overcome local optima in the loss landscape on the training dataset.

251
252

4.3 CIFAR-100

253
254
255
256
257
258

On the CIFAR-100 dataset, momentum significantly affects accuracy, whereas QMSE impacts sensitivity. It shows in Figure 3 Training with SGD without momentum achieves nearly 0.15 training accuracy, while using momentum increases it to almost 0.35. First, the model is too simple to effectively train on this dataset. It results in overall low accuracy. Additionally, CIFAR-100 has 100 output channels, which makes the loss landscape more complex and introduces more local optima. Nevertheless, QMSE improves performance on complex datasets, with or without momentum.

259
260

5 CONCLUSION

261
262
263
264
265
266
267
268
269

In this work. We propose QMSE, a modification of MSE to overcome local optima in the loss landscape. Through experiments on MNIST, CIFAR-10, and CIFAR-100 datasets, we show that QMSE improves performance, and it can make more improvements when combined with momentum. While both QMSE and MSE achieve similar validation accuracy on CIFAR-10, QMSE exhibits outperforms on train accuracy on all of three datasets. These results highlight that QMSE provides a robust optimization strategy that enhances optimizer’s ability to reach better optima. Moreover, high performance in train dataset does not necessarily indicate good generalization. Future work will explore the integration of QMSE with more advanced architectures and various optimization strategies to further improve performance on large-scale and complex datasets.

REFERENCES

- 270
271
272 Naman Agarwal, Zeyuan Allen-Zhu, Brian Bullins, Elad Hazan, and Tengyu Ma. Finding approxi-
273 mate local minima faster than gradient descent. In *Proceedings of the 49th Annual ACM SIGACT*
274 *Symposium on Theory of Computing*, pp. 1195–1199, 2017a.
- 275 Naman Agarwal, Brian Bullins, and Elad Hazan. Second-order stochastic optimization for machine
276 learning in linear time. *Journal of Machine Learning Research*, 18(116):1–40, 2017b.
- 277 Patrick Cheridito, Arnulf Jentzen, and Florian Rossmannek. Gradient descent provably escapes
278 saddle points in the training of shallow relu networks. *Journal of Optimization Theory and Appli-*
279 *cations*, 203(3):2617–2648, 2024.
- 280 Anna Choromanska, Yann LeCun, and Gérard Ben Arous. Open problem: The landscape of the
281 loss surfaces of multilayer networks. In *Conference on Learning Theory*, pp. 1756–1760. PMLR,
282 2015.
- 283 Yann N Dauphin, Razvan Pascanu, Caglar Gulcehre, Kyunghyun Cho, Surya Ganguli, and Yoshua
284 Bengio. Identifying and attacking the saddle point problem in high-dimensional non-convex op-
285 timization. *Advances in neural information processing systems*, 27, 2014.
- 286 Geoffrey Hinton, Nitish Srivastava, and Kevin Swersky. Neural networks for machine learning
287 lecture 6a overview of mini-batch gradient descent. *Cited on*, 14(8):2, 2012.
- 288 Prateek Jain, Purushottam Kar, et al. Non-convex optimization for machine learning. *Foundations*
289 *and Trends® in Machine Learning*, 10(3-4):142–363, 2017.
- 290 Antoine Jamin and Anne Humeau-Heurtier. (multiscale) cross-entropy methods: A review. *Entropy*,
291 22(1):45, 2019.
- 292 Samy Jelassi and Yuezhi Li. Towards understanding how momentum improves generalization in
293 deep learning. In *International Conference on Machine Learning*, pp. 9965–10040. PMLR, 2022.
- 294 Mehrdad Kaveh and Mohammad Saadi Mesgari. Application of meta-heuristic algorithms for train-
295 ing neural networks and deep learning architectures: A comprehensive review. *Neural Processing*
296 *Letters*, 55(4):4519–4622, 2023.
- 297 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
298 *arXiv:1412.6980*, 2014.
- 299 Shie Mannor, Dori Peleg, and Reuven Rubinfeld. The cross entropy method for classification. In
300 *Proceedings of the 22nd international conference on Machine learning*, pp. 561–568, 2005.
- 301 Quynh Nguyen. On the proof of global convergence of gradient descent for deep relu networks with
302 linear widths. In *International Conference on Machine Learning*, pp. 8056–8062. PMLR, 2021.
- 303 Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. *arXiv*
304 *preprint arXiv:1904.09237*, 2019.
- 305 Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initial-
306 ization and momentum in deep learning. In *International conference on machine learning*, pp.
307 1139–1147. pmlr, 2013.
- 308 Zhonghuan Tian and Simon Fong. Survey of meta-heuristic algorithms for deep learning training.
309 *Optimization algorithms—methods and applications*, pp. 195–220, 2016.
- 310 Yuchang Zhang, Ruibin Bai, Rong Qu, Chaofan Tu, and Jiahuan Jin. A deep reinforcement learning
311 based hyper-heuristic for combinatorial optimisation with uncertainties. *European Journal of*
312 *Operational Research*, 300(2):418–427, 2022.
- 313 Jinxin Zhou, Xiao Li, Tianyu Ding, Chong You, Qing Qu, and Zhihui Zhu. On the optimization
314 landscape of neural collapse under mse loss: Global optimality with unconstrained features. In
315 *International Conference on Machine Learning*, pp. 27179–27202. PMLR, 2022.
- 316 Pan Zhou, Jiashi Feng, Chao Ma, Caiming Xiong, Steven Chu Hong Hoi, et al. Towards theoretic-
317 ally understanding why sgd generalizes better than adam in deep learning. *Advances in Neural*
318 *Information Processing Systems*, 33:21285–21296, 2020.

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

A APPENDIX

We only use a Large Language Model (LLM) to assist in translating Korean into English.