

Positive-Negative Momentum: Manipulating Stochastic Gradient Noise to Improve Generalization

Zeke Xie^{1,2}, Li Yuan³, Zhanxing Zhu⁴, and Masashi Sugiyama^{2,1}

¹The University of Tokyo

²RIKEN Center for AIP

³National University of Singapore

⁴Beijing Institute of Big Data Research

ICML2021, July, 2021



東京大学
THE UNIVERSITY OF TOKYO



Center for
Advanced Intelligence Project

Positive-Negative Momentum (PNM)

- Motivation:
As Stochastic Gradient Noise (SGN) is important, we want to **manipulate SGN** to improve deep learning.
- Constraint:
We prefer **not to change the learning rate or batch size** due to multiple limitations.
- Our work:
PNM, an alternative to conventional Momentum, can **manipulate SGN without changing the learning rate or batch size**.

Manipulating SGN by changing $\frac{\eta}{B}$.

(Mandt et al., 2017): The magnitude of SGN is proportional to $\frac{\eta}{B}$, where η is the learning rate and B is the batch size.

- (He et al., 2019) argued that increasing $\frac{\eta}{B}$ which corresponds larger SGN may improve generalization.
- (Xie et al., 2021) argued that $\frac{\eta}{B}$ exponentially matters to learning flat minima.

Stephan, M., Hoffman, M. D., and Blei, D. M. (2017). Stochastic Gradient Descent as Approximate Bayesian Inference. JMLR.

He, F., Liu, T., and Tao, D. (2019). Control batch size and learning rate to generalize well: Theoretical and empirical evidence. In NeurIPS.

Zeke Xie, Issei Sato, and Masashi Sugiyama. (2021). A Diffusion Theory For Deep Learning Dynamics: Stochastic Gradient Descent Exponentially Favors Flat Minima. In ICLR.

Increasing $\frac{\eta}{B}$ has at least three limitations.

Three Limitations:

- 1 Training with a too small batch size is **computationally expensive** per epoch and often requires more epochs for convergence.
- 2 Increasing the learning rate only works in a narrow range, since too large initial learning rates may lead to **optimization divergence or bad convergence**.
- 3 Decaying the ratio $\frac{\eta}{B}$ during training (via **learning rate decay**) is almost necessary for the convergence/training of deep networks.

Basic idea: Positive-Negative Weighted Average

Suppose that $g^{(a)}$ and $g^{(b)}$ are two independent unbiased noisy gradients of $\nabla f(\theta)$, where $f(\theta)$ is the loss. Then their weighted average is

$$\begin{aligned}\bar{g} &= (1 + \beta_0)g^{(a)} - \beta_0g^{(b)} \\ &= \nabla f(\theta) + \bar{\xi},\end{aligned}$$

where the gradient noise $\bar{\xi} = (1 + \beta_0)\xi^{(a)} - \beta_0\xi^{(b)}$. If $\beta_0 > 0$ and $\sigma^2 = \text{Var}(\xi^{(a)}) = \text{Var}(\xi^{(b)})$, we have the same gradient expectation $\mathbb{E}[\bar{\xi}] = 0$ but larger gradient noise

$$\text{Var}(\bar{\xi}) = [(1 + \beta_0)^2 + \beta_0^2]\sigma^2.$$

β_0 controls the noise magnitude!

From Momentum to Positive-Negative Momentum

Momentum: uses the (positive only) weighted average of past gradients to update parameters as

$$m_t = \sum_{k=0}^t \beta_3 \beta_1^{t-k} g_k.$$

Then we approximately have $\mathbb{E}[m] \approx \frac{\beta_3}{1-\beta_1} \nabla f(\theta)$. The SGN in momentum is given by

$$\xi_t^m = \sum_{k=0}^t \beta_3 \beta_1^{t-k} \xi_k.$$

Algorithm 1: Momentum

$$m_t = \beta_1 m_{t-1} + \beta_3 g_t;$$

$$\theta_{t+1} = \theta_t - \eta m_t;$$

PNM can manipulate SGN

PNM: uses the (positive-negative) weighted average of past gradients to update parameters.

PNM maintains two independent momentum terms as

$$\left\{ \begin{array}{l} m_t^{(\text{odd})} = \sum_{k=1,3,\dots,t} \beta_3 \beta_1^{t-k} g_k, \\ m_t^{(\text{even})} = \sum_{k=0,2,\dots,t-1} \beta_3 \beta_1^{t-k} g_k, \\ m_t = (1 + \beta_0) m_t^{(\text{odd})} - \beta_0 m_t^{(\text{even})}, \end{array} \right.$$

by using two alternate sequences of past gradients, respectively. Then we approximately have $\mathbb{E}[m] \approx \frac{\beta_3}{1-\beta_1} \nabla f(\theta)$. The SGN in PNM is given by

$$\text{Var}(\xi^{\text{pnm}}) \approx [(1 + \beta_0)^2 + \beta_0^2] \xi^{\text{m}}.$$

β_0 controls the noise magnitude!

Algorithm 2: (Stochastic) PNM

$$m_t = \beta_1^2 m_{t-2} + (1 - \beta_1^2) g_t;$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{(1+\beta_0)^2 + \beta_0^2}} [(1 + \beta_0)m_t - \beta_0 m_{t-1}];$$

- **Theorem 1:** PNM has the convergence guarantee similar to Momentum.
- **Theorem 4:** Stochastic PNM may have a **tighter generalization bound** than SGD within the PAC-Bayesian framework.

Experiment: PNM and AdaPNM

Table: PNM versus conventional Momentum. We report the mean and the standard deviations (as the subscripts) of the optimal test errors computed over three runs of each experiment.

DATASET	MODEL	PNM	AdaPNM	SGD	ADAM	AMSGRAD	ADAMW	ADABOUND	PADAM	YOGI	RADAM
CIFAR-10	RESNET18	4.48 _{0.09}	4.94 _{0.05}	5.01 _{0.03}	6.53 _{0.03}	6.16 _{0.18}	5.08 _{0.07}	5.65 _{0.08}	5.12 _{0.04}	5.87 _{0.12}	6.01 _{0.10}
	VGG16	6.26 _{0.05}	5.99 _{0.11}	6.42 _{0.02}	7.31 _{0.25}	7.14 _{0.14}	6.48 _{0.13}	6.76 _{0.12}	6.15 _{0.06}	6.90 _{0.22}	6.56 _{0.04}
CIFAR-100	RESNET34	20.59 _{0.29}	20.41 _{0.18}	21.52 _{0.37}	27.16 _{0.55}	25.53 _{0.19}	22.99 _{0.40}	22.87 _{0.13}	22.72 _{0.10}	23.57 _{0.12}	24.41 _{0.40}
	DENSENET121	19.76 _{0.28}	20.68 _{0.11}	19.81 _{0.33}	25.11 _{0.15}	24.43 _{0.09}	21.55 _{0.14}	22.69 _{0.15}	21.10 _{0.23}	22.15 _{0.36}	22.27 _{0.22}
	GOOGLENET	20.38 _{0.31}	20.26 _{0.21}	21.21 _{0.29}	26.12 _{0.33}	25.53 _{0.17}	21.29 _{0.17}	23.18 _{0.31}	21.82 _{0.17}	24.24 _{0.16}	22.23 _{0.15}

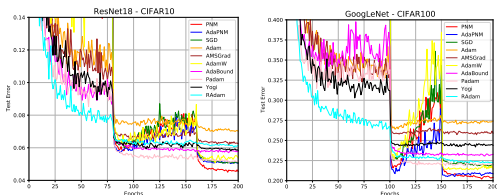


Figure: The learning curves of popular models on CIFAR-10 and CIFAR-100, respectively. It demonstrates that PNM and AdaPNM yield significantly better test results.

Experiment: PNM prevents overfitting noisy labels

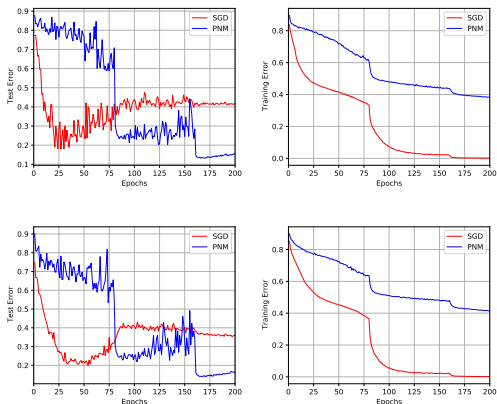


Figure: ResNet34 on CIFAR-10 with 40% asymmetric label noise (Top Row) and 40% symmetric label noise (Bottom Row). Left: Test Curve. Right: Training Curve. PNM with a large β_0 may effectively relieve memorizing noisy labels and almost only learn clean labels, while SGD almost memorizes all noisy labels and has a nearly zero training error.

Experiment: Robustness to the new hyperparameter β_0 .

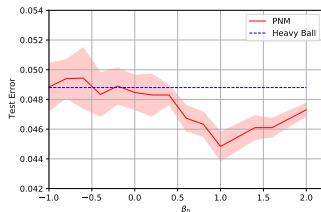


Figure: We train ResNet18 on CIFAR-10 under various β_0 choices. It demonstrates that PNM may achieve significantly better generalization with a wide range of $\beta_0 > 0$, which corresponds to a positive-negative momentum pair for enhancing SGN as we expect. With any $\beta_0 \in [-1, 0]$, the test performance does not sensitively depend on β_0 , because this case cannot enhance SGN.

Three contributions:

- 1 PNM can **manipulate SGN** without changing the learning rate or batch size.
- 2 PNM can **replace conventional Momentum** in popular optimizers at low costs.
- 3 PNM can have a **tighter generalization bound** than SGD.

Future directions:

- Rethinking the importance and the popularity of Momentum.
- Designing optimization dynamics by manipulating SGN.