

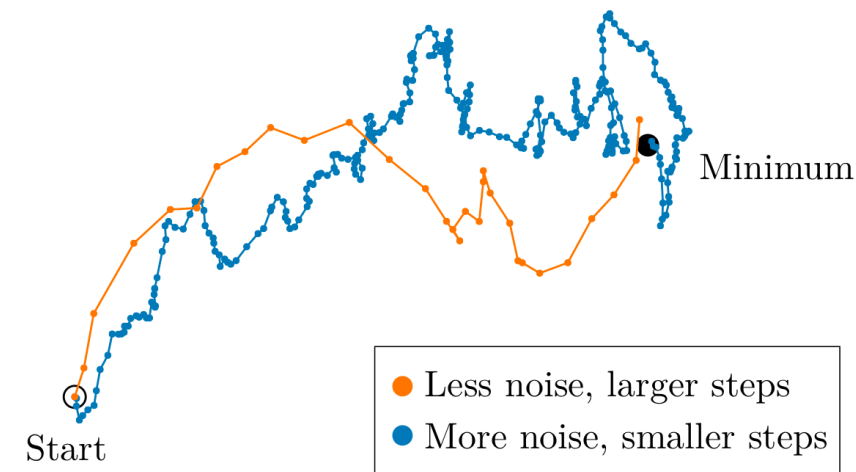
Surge Phenomenon in Optimal Learning Rate and Batch Size Scaling

Shuaipeng Li¹, Penghao Zhao^{1,2}, Hailin Zhang^{1,2}, Xingwu Sun^{1,3},
Hao Wu¹, Dian Jiao¹, Weiyang Wang¹, Chengjun Liu¹, Zheng Fang¹,
Jinbao Xue¹, Yangyu Tao¹, Bin Cui², Di Wang¹

¹ Tencent, ² Peking University, ³ University of Macau

Optimal Learning Rate: Related to Batch Size

- Past experience: **larger** batches require **larger** optimal learning rates.
- Intuition from **gradient noise**:
 - Small batch size -> high noise -> small step size.
 - Large batch size -> low noise -> large step size.
- Current insights into scaling:
 - **Linear** or **Sqrt**.
- Significance:
 - Guide model training, **simplify hyper-parameter tuning**.



OpenAI, An Empirical Model of Large-Batch Training, 2018

Research on SGD Optimizer [OpenAI, 2018]

- Perturb the parameter and apply a Taylor expansion:

$$L(\theta - \epsilon V) \approx L(\theta) - \epsilon G^T V + \frac{1}{2} \epsilon^2 V^T H V.$$

- G_{est} contains noises. Calculate expectation:

$$\mathbb{E}[L(\theta - \epsilon G_{\text{est}})] = L(\theta) - \epsilon |G|^2 + \frac{1}{2} \epsilon^2 \left(G^T H G + \frac{\text{tr}(H \Sigma)}{B} \right).$$

$$\mathbb{E}_{x_1 \dots x_B \sim \rho} [G_{\text{est}}(\theta)] = G(\theta)$$

$$\text{cov}_{x_1 \dots x_B \sim \rho} (G_{\text{est}}(\theta)) = \frac{1}{B} \Sigma(\theta),$$

$$\Sigma(\theta) \equiv \text{cov}_{x \sim \rho} (\nabla_{\theta} L_x(\theta))$$

$$= \mathbb{E}_{x \sim \rho} \left[(\nabla_{\theta} L_x(\theta)) (\nabla_{\theta} L_x(\theta))^T \right] - G(\theta) G(\theta)^T.$$

- Minimizing the equation:

$$\epsilon_{\text{opt}}(B) = \text{argmin}_{\epsilon} \mathbb{E}[L(\theta - \epsilon G_{\text{est}})] = \frac{\epsilon_{\text{max}}}{1 + \mathcal{B}_{\text{noise}}/B}$$

$$\epsilon_{\text{max}} \equiv \frac{|G|^2}{G^T H G}$$

$$\Delta L_{\text{opt}}(B) = \frac{\Delta L_{\text{max}}}{1 + \mathcal{B}_{\text{noise}}/B}; \quad \Delta L_{\text{max}} = \frac{1}{2} \frac{|G|^4}{G^T H G}.$$

$$\mathcal{B}_{\text{noise}} = \frac{\text{tr}(H \Sigma)}{G^T H G},$$

- Approximations:

- When $B \ll \mathcal{B}_{\text{noise}}$, the optimal learning rate scales **linearly** with the batch size.
- When $B \gg \mathcal{B}_{\text{noise}}$, the optimal learning rate **approaches its maximum value**.

Trading-off # training steps and # samples [OpenAI, 2018]

- # training steps: $\delta S = 1 + \frac{\mathcal{B}}{B}$; # samples: $\delta E = B\delta S$

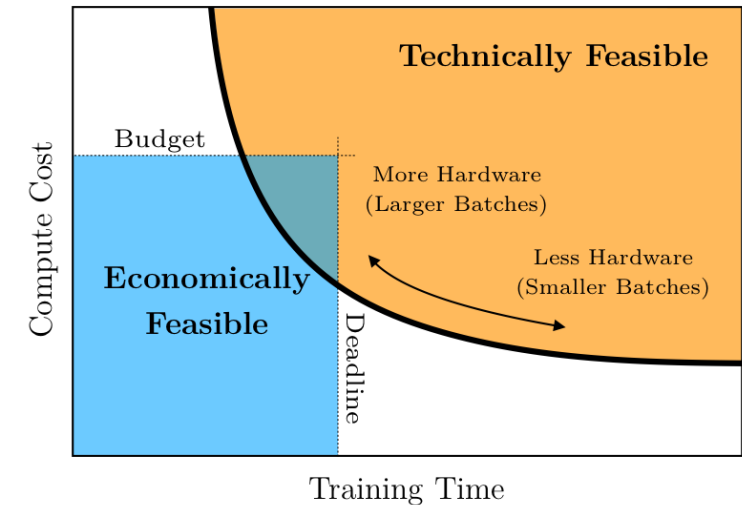
- Considering multiple steps:
$$S = \int \left(1 + \frac{\mathcal{B}(s)}{B(s)} \right) ds$$

$$E = \int (\mathcal{B}(s) + B(s)) ds$$

- After some derivation:

$$\frac{S}{S_{\min}} - 1 = \left(\frac{E}{E_{\min}} - 1 \right)^{-1}$$

- **Need to balance the computation cost and training time!**



Adam optimizer is different from SGD optimizer

- The main difference is the **update amount** at each step:
 - Simplified assumption: the update amount is the **gradient's sign**:

$$\theta_{i+1} = \theta_i - \epsilon \cdot \text{sign}(G_{est}),$$

- In Adam, the update amount is:

$$V = \frac{\widehat{m}_t}{\sqrt{\widehat{v}_t} + \epsilon_{Adam}} = \frac{\frac{1-\beta_1}{1-\beta_1^t} \sum_i \beta_1^{t-i} G_{est,i}}{\sqrt{\frac{1-\beta_2}{1-\beta_2^t} \sum_i \beta_2^{t-i} G_{est,i}^2 + \epsilon_{Adam}}}$$

- $\beta \rightarrow 1$: the following formula tends to $\text{sign}(G_{est})$ when the variance is small

$$V = \frac{\frac{\sum_i G_{est,i}}{t}}{\sqrt{\frac{\sum_i G_{est,i}^2}{t}}} = \frac{\mathbb{E}_t[G_{est}]}{\sqrt{\mathbb{E}_t[G_{est}^2]}} = \frac{\text{sign}(\mathbb{E}_t[G_{est}])}{\sqrt{1 + \frac{\text{var}_t(G_{est})}{\mathbb{E}_t[G_{est}]^2}}$$

- $\beta \rightarrow 0$: degenerates to $\text{sign}(G_{est})$

$$V = \frac{G_{est}}{\sqrt{G_{est}^2}} = \text{sign}(G_{est})$$

The optimal learning rate for Adam optimizer

- Still Taylor expansion:

$$\Delta L_{opt} = \frac{G^T \mathbb{E}[V]}{2} \epsilon_{opt}.$$

$$\mathbb{E}[\Delta L] = \mathbb{E}[L(\theta) - L(\theta - \epsilon \cdot V)] \approx \epsilon G^T \mathbb{E}[V] - \frac{1}{2} \epsilon^2 \mathbb{E}[V^T H V]. \Rightarrow \epsilon_{opt} \equiv \operatorname{argmax}_{\epsilon} \mathbb{E}[\Delta L] = \frac{G^T \mathbb{E}[V]}{\operatorname{tr}[H \cdot \operatorname{cov}(V)] + \mathbb{E}[V]^T H \mathbb{E}[V]},$$

- For the update amount $\operatorname{sign}(G_{est})$:

- Assumption: the gradient follows a **normal distribution**.

- For each sample, the expectation is $\operatorname{erf}\left(\frac{\mu}{\sqrt{2}\sigma}\right)$ and the variance is $1 - \operatorname{erf}\left(\frac{\mu}{\sqrt{2}\sigma}\right)^2$

- For a batch, we have $G_{est}(\theta_i) = \frac{1}{B} \sum G_x(\theta_i) \sim \mathcal{N}\left(\mu_i, \frac{\sigma_i^2}{B}\right)$

- For all parameters:
$$\mathbb{E}[V] = \begin{pmatrix} \vdots \\ \operatorname{erf}\left(\sqrt{\frac{B}{2}} \frac{\mu_i}{\sigma_i}\right) \\ \vdots \end{pmatrix} \quad \operatorname{cov}(V) = \begin{pmatrix} \ddots & & 0 \\ & 1 - \operatorname{erf}\left(\sqrt{\frac{B}{2}} \frac{\mu_i}{\sigma_i}\right)^2 & \\ 0 & & \ddots \end{pmatrix}$$

- As a result:

$$\Delta L_{opt} = \frac{1}{2} \frac{\sum_i \sum_j \mathcal{E}_i \mathcal{E}_j \mu_i \mu_j}{\sum_i (1 - \mathcal{E}_i^2) H_{i,i} + \sum_i \sum_j \mathcal{E}_i \mathcal{E}_j H_{i,j}} \quad \epsilon_{opt} = \frac{\sum_i \mathcal{E}_i \mu_i}{\sum_i (1 - \mathcal{E}_i^2) H_{i,i} + \sum_i \sum_j \mathcal{E}_i \mathcal{E}_j H_{i,j}} \quad \leftarrow \quad \mathcal{E}_i(B) = \frac{2}{\sqrt{\pi}} \int_0^{\sqrt{\frac{B}{2}} \frac{\mu_i}{\sigma_i}} e^{-t^2} dt \approx \frac{\frac{\mu_i}{\sigma_i}}{\sqrt{\frac{\pi}{2B} + \left(\frac{\mu_i}{\sigma_i}\right)^2}}$$

The optimal learning rate for Adam optimizer

- Simplifying the relationship between learning rate and batch size:

$$\epsilon(B) = \frac{\beta f(B)}{f(B)^2 + \gamma} = \frac{\beta}{f(B) + \frac{\gamma}{f(B)}}$$

- When $B \ll \frac{\pi \sigma_i^2}{2\mu_i^2}$, the optimal learning rate **first rises then drops**:

$$\epsilon_{opt}(B) \approx \frac{1}{\frac{1}{2}(\sqrt{\frac{\mathcal{B}_{noise}}{B}} + \sqrt{\frac{B}{\mathcal{B}_{noise}}})} \frac{\sqrt{\frac{\mathcal{B}_{noise}}{2\pi}} \sum_i \frac{\mu_i^2}{\sigma_i}}{\sum_i H_{i,i}} \leq \frac{\sqrt{\frac{\mathcal{B}_{noise}}{2\pi}} \sum_i \frac{\mu_i^2}{\sigma_i}}{\sum_i H_{i,i}} \leftarrow \mathcal{E}_i(B) \approx \sqrt{\frac{2B}{\pi}} \frac{\mu_i}{\sigma_i} \quad \mathcal{B}_{noise} = \frac{\pi \sum_i H_{i,i}}{2 \sum_i \sum_j \begin{cases} \frac{\mu_i \mu_j}{\sigma_i \sigma_j} & i \neq j \\ 0 & i = j \end{cases} H_{i,j}}$$

- When $B \gg \frac{\pi \sigma_i^2}{2\mu_i^2}$, the optimal learning rate **tends to a constant**:

$$\epsilon_{opt} = \frac{\sum_i |\mu_i|}{\sum_i \sum_j \text{sign}(\mu_i) \text{sign}(\mu_j) H_{i,j}} \leftarrow \mathcal{E}_i = \text{sign}\left(\frac{\mu_i}{\sigma_i}\right) = \text{sign}(\mu_i)$$

Trading-off # training steps and # samples with Adam optimizer

- The **same trade-off** is also obtained with the Adam optimizer:

$$\left(\frac{S}{S_{min}} - 1\right)\left(\frac{E}{E_{min}} - 1\right) = 1.$$

- When $B \ll \frac{\pi \sigma_i^2}{2\mu_i^2}$, the loss delta is:

$$\Delta L_{opt}(B) = \frac{\Delta L_{max}}{1 + \frac{\mathcal{B}_{noise}}{B}}$$



$$\Delta L_{max} = \frac{\sum_i \sum_j \frac{\mu_i^2 \mu_j^2}{\sigma_i \sigma_j}}{2 \sum_i \sum_j \begin{cases} \frac{\mu_i \mu_j}{\sigma_i \sigma_j} & i \neq j \\ 0 & i = j \end{cases} H_{i,j}}$$

Adam **only**
affects ΔL_{max} ,
does not change
the equation!

- This formula is the same as the SGD case! Consequently,

$$S = \int (1 + \frac{\mathcal{B}_{noise}}{B}) ds$$

$$E = \int (\mathcal{B}_{noise} + B) ds$$



$$\left(\frac{S}{S_{min}} - 1\right)\left(\frac{E}{E_{min}} - 1\right) = 1.$$

Application process

- In the experiment, we **fit B_{noise}** with the following formula:

$$\frac{1}{S} = -\mathcal{B}_{noise} \frac{1}{E} + \frac{1}{S_{min}}$$

$$B_{crit}(L) = \frac{B_*}{L^{1/\alpha_B}}, \quad B_* \sim 2 \cdot 10^8 \text{ tokens}, \quad \alpha_B \sim 0.21$$

$$B_{peak} = \mathcal{B}_{noise} \approx \mathcal{B}_{crit} = \frac{E_{min}}{S_{min}} \approx \frac{B_*}{L^{1/\alpha_B}}$$

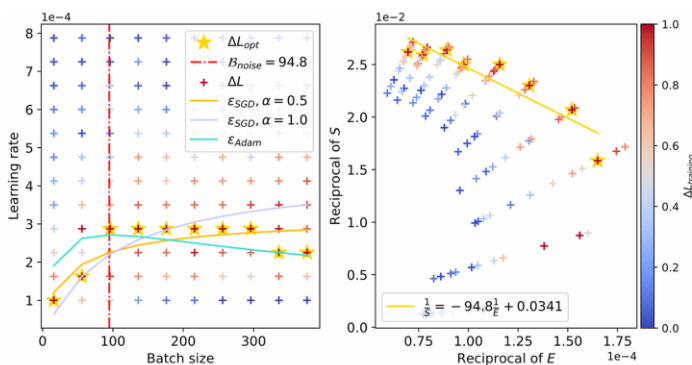
- In practice, we can also derive it through the scaling law above.
- And then use one pair of (optimal learning rate, batch size) to **obtain ϵ_{max}** :

$$\mathbb{E}[\epsilon_{max}]_{Adam} = \mathbb{E}\left[\frac{\epsilon_{opt}}{2} \left(\sqrt{\frac{\mathcal{B}_{noise}}{B}} + \sqrt{\frac{B}{\mathcal{B}_{noise}}} \right)\right]$$

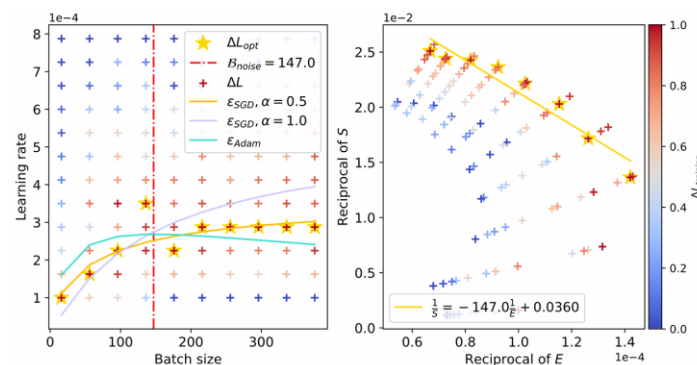
$$\mathbb{E}[\epsilon_{max}]_{SGD} = \mathbb{E}\left[\epsilon_{opt} \left(1 + \frac{\mathcal{B}_{noise}}{B}\right)^\alpha\right]$$

Experiments: multiple workloads including CV and NLP

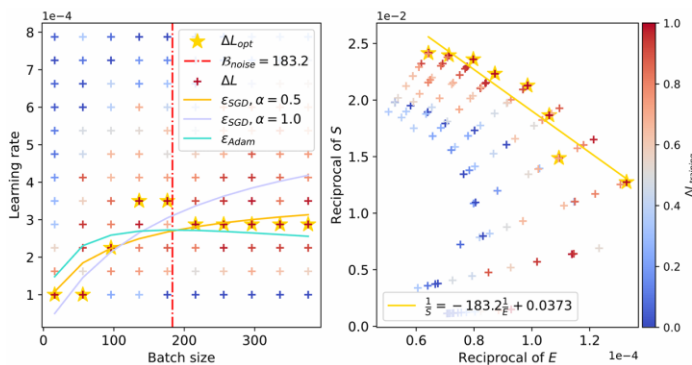
- ResNet-18 on TinyImageNet



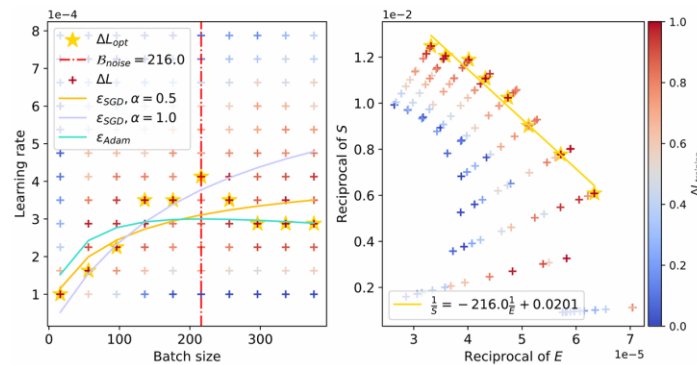
(a) $L_{training} = 5.213 \pm 0.001$



(b) $L_{training} = 5.146 \pm 0.001$



(c) $L_{training} = 5.113 \pm 0.001$

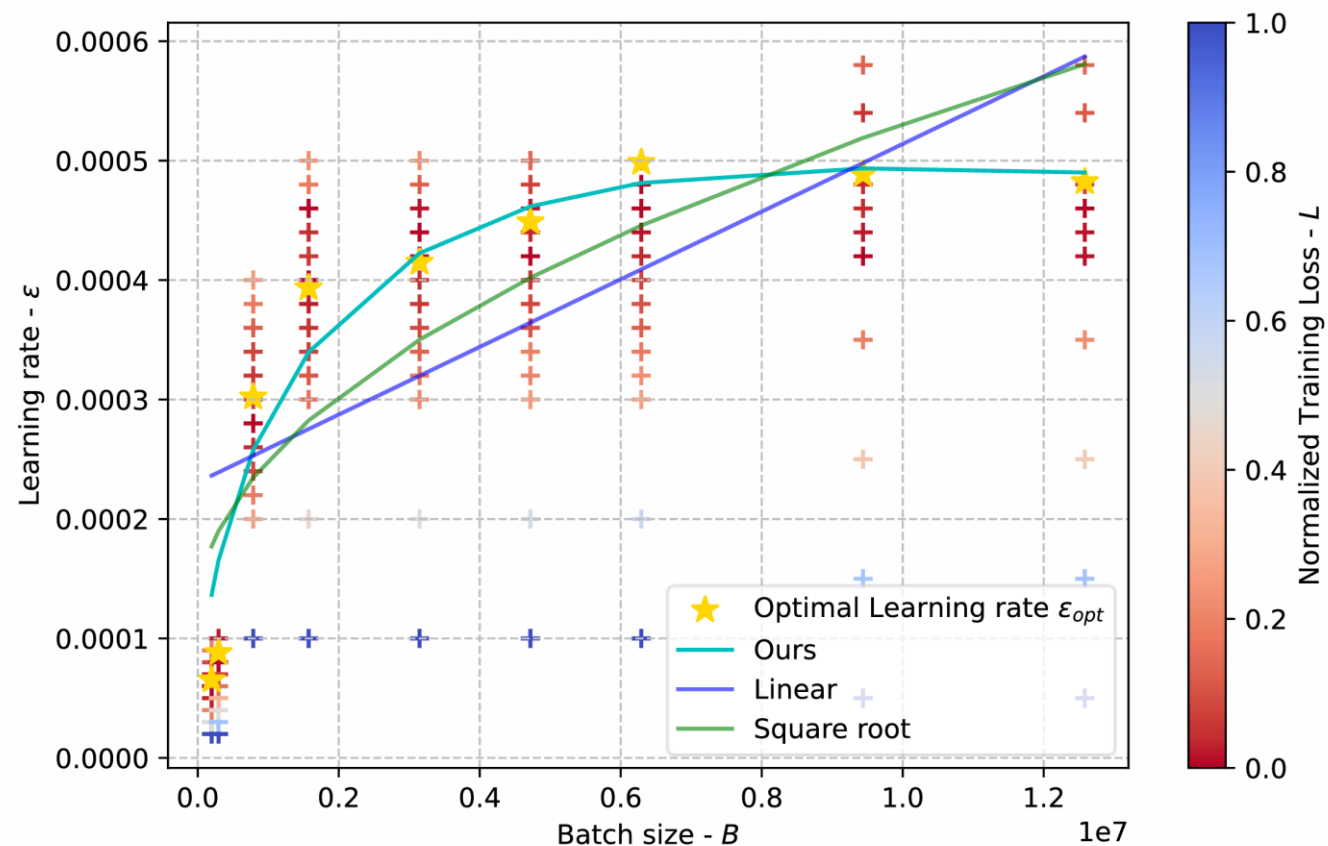


(d) $L_{training} = 4.863 \pm 0.001$

“Surge Phenomenon in Optimal Learning Rate and Batch Size Scaling.” NeurIPS 2024

Experiments: multiple workloads including CV and NLP

- MoE on RedPajama-v2



Conclusion and significance

- Takeaways:
 - As the batch size increases, the optimal learning rate demonstrates a **decreasing trend** within a specified range.
 - The batch size that corresponds to the local maximum optimal learning rate **is consistent with** the balance point of training speed and data efficiency. As the training progresses and the loss decreases, B_{noise} will gradually **becomes larger**.
- Significance:
 - A deeper understanding of the **training dynamics**.
 - Help **tune hyperparameters**, improve convergence speed, and avoid complicated grid searches.

Thanks for your attention!
