

Implicit regularization in Heavy-ball momentum accelerated stochastic gradient descent

Avrajit Ghosh^[1], He Lyu^[1], Xitong Zhang, Rongrong Wang

Dept. of Computational Math., Science and Engineering, Michigan State University, USA.

Dept of Mathematics, Michigan State University, USA.

April 21, 2023



ICLR
International Conference On
Learning Representations

^[1]denotes equal contribution

Background

- Momentum methods (SGD+M) empirically outperforms traditional stochastic gradient descent (SGD).^[2]

[2] Sutskever, I., Martens, J., Dahl, G., and Hinton, G. On the importance of initialization and momentum in deep learning. In International conference on machine learning, pp.1139–1147. PMLR, 2013

[3] Defazio, A. Understanding the role of momentum in non-convex optimization: Practical insights from a Lyapunov analysis

Background

- Momentum methods (SGD+M) empirically outperforms traditional stochastic gradient descent (SGD).^[2]
- From optimization perspective, convergence of momentum^[3] is well studied.

[2] Sutskever, I., Martens, J., Dahl, G., and Hinton, G. On the importance of initialization and momentum in deep learning. In International conference on machine learning, pp.1139–1147. PMLR, 2013

[3] Defazio, A. Understanding the role of momentum in non-convex optimization: Practical insights from a Lyapunov analysis

Background

- Momentum methods (SGD+M) empirically outperforms traditional stochastic gradient descent (SGD). ^[2]
- From optimization perspective, convergence of momentum ^[3] is well studied.
- No theoretical answer on how momentum helps generalization for deep neural networks.

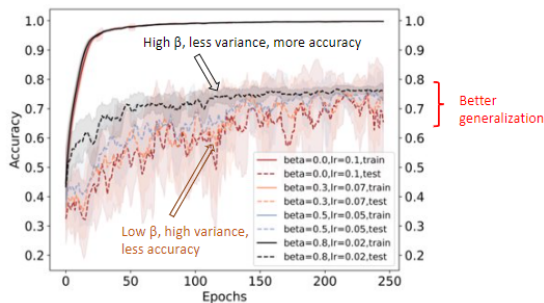
[2] Sutskever, I., Martens, J., Dahl, G., and Hinton, G. On the importance of initialization and momentum in deep learning. In International conference on machine learning, pp.1139–1147. PMLR, 2013

[3] Defazio, A. Understanding the role of momentum in non-convex optimization: Practical insights from a Lyapunov analysis

CIFAR-10 classification with momentum

Heavy ball momentum update

$$\mathbf{x}^{k+1} = \mathbf{x}^k - h\nabla E(\mathbf{x}^k) + \beta(\mathbf{x}^k - \mathbf{x}^{k-1}) \quad \forall k = 1, 2, \dots, n$$



Some prevalent observations in deep learning

- 1 SGD/GD: Test accuracy improves with higher learning rate (h) ?
⇒ Implicit Gradient Regularization^[4] [5]

[4] Barrett and Dherin, "Implicit gradient regularization", ICLR-21

[5] Smith et al, "On the Origin of Implicit Regularization in Stochastic Gradient Descent", ICLR-21

Some prevalent observations in deep learning

- 1 SGD/GD: Test accuracy improves with higher learning rate (h) ?
 \implies Implicit Gradient Regularization^[4] [5]
- 2 (SGD+M): $\beta \uparrow$, variation in test-accuracy \downarrow across epochs .

[4] Barrett and Dherin, "Implicit gradient regularization", ICLR-21

[5] Smith et al, "On the Origin of Implicit Regularization in Stochastic Gradient Descent", ICLR-21

Some prevalent observations in deep learning

- 1 SGD/GD: Test accuracy improves with higher learning rate (h) ?
 \implies Implicit Gradient Regularization^[4] [5]
- 2 (SGD+M): $\beta \uparrow$, variation in test-accuracy \downarrow across epochs .
- 3 (SGD+M): $\beta \uparrow$, test accuracy \uparrow .

[4] Barrett and Dherin, "Implicit gradient regularization", ICLR-21

[5] Smith et al, "On the Origin of Implicit Regularization in Stochastic Gradient Descent", ICLR-21

Some prevalent observations in deep learning

- 1 SGD/GD: Test accuracy improves with higher learning rate (h) ?
 \implies Implicit Gradient Regularization^[4] [5]
- 2 (SGD+M): $\beta \uparrow$, variation in test-accuracy \downarrow across epochs .
- 3 (SGD+M): $\beta \uparrow$, test accuracy \uparrow .

[4] Barrett and Dherin, "Implicit gradient regularization", ICLR-21

[5] Smith et al, "On the Origin of Implicit Regularization in Stochastic Gradient Descent", ICLR-21

Some prevalent observations in deep learning

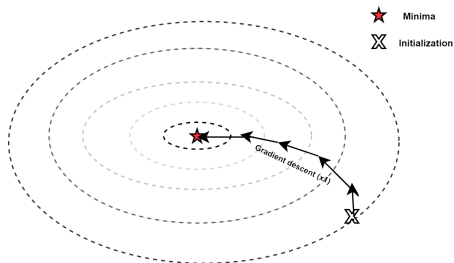
- 1 SGD/GD: Test accuracy improves with higher learning rate (h) ?
 \implies Implicit Gradient Regularization^[4] [5]
- 2 (SGD+M): $\beta \uparrow$, variation in test-accuracy \downarrow across epochs .
- 3 (SGD+M): $\beta \uparrow$, test accuracy \uparrow .

Contribution: We study why generalization improves and variance reduces for (SGD+M) with increasing β .

[4]Barrett and Dherin, "Implicit gradient regularization", ICLR-21

[5]Smith et al, "On the Origin of Implicit Regularization in Stochastic Gradient Descent", ICLR-21

Implicit Gradient Regularization^[6]

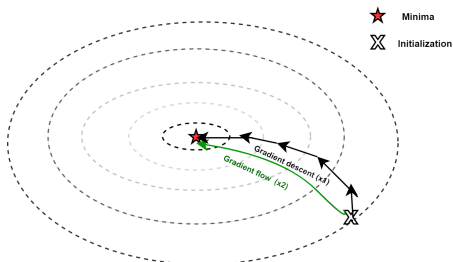


- Gradient Descent:

$$\mathbf{x}_1^{k+1} = \mathbf{x}_1^k - h \nabla E(\mathbf{x}_1^k)$$

Implicit Gradient Regularization^[6]

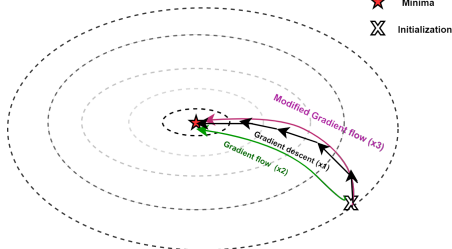
\mathbf{x}_1 is $O(h)$ close to \mathbf{x}_2 .



- Gradient Descent:
 $\mathbf{x}_1^{k+1} = \mathbf{x}_1^k - h\nabla E(\mathbf{x}_1^k)$
- Gradient Flow ($h \rightarrow 0$):
 $\mathbf{x}'_2(t) = -\nabla E(\mathbf{x}_2(t))$

Implicit Gradient Regularization^[6]

\mathbf{x}_1 is $O(h)$ close to \mathbf{x}_2 .



\mathbf{x}_1 is $O(h^2)$ close to \mathbf{x}_3 .

- Gradient Descent:

$$\mathbf{x}_1^{k+1} = \mathbf{x}_1^k - h \nabla E(\mathbf{x}_1^k)$$

- Gradient Flow ($h \rightarrow 0$):

$$\mathbf{x}'_2(t) = -\nabla E(\mathbf{x}_2(t))$$

- Modified flow:

$$\mathbf{x}'_3(t) = -\nabla(E(\mathbf{x}_3) + \underbrace{\frac{h}{4} \|\nabla E(\mathbf{x}_3)\|^2}_{\text{IGR}})$$

Can we study Heavy-ball momentum gradient descent trajectory using modified continuous flow ?

Heavy ball momentum update

$$\mathbf{x}^{k+1} = \mathbf{x}^k - h\nabla E(\mathbf{x}^k) + \beta(\mathbf{x}^k - \mathbf{x}^{k-1}) \quad \forall k = 1, 2, \dots, n$$

Heavy ball momentum update

$$\mathbf{x}^{k+1} = \mathbf{x}^k - h\nabla E(\mathbf{x}^k) + \beta(\mathbf{x}^k - \mathbf{x}^{k-1}) \quad \forall k = 1, 2, \dots, n$$

is globally $O(h)$ close to the ODE^[7]

$$\tilde{\mathbf{x}}'(t) = -\frac{1}{1-\beta} \nabla E(\tilde{\mathbf{x}}(t)), \quad t \in [0, T] \quad (1)$$

[7]Kovachki,Stuart, "Continuous time analysis of momentum methods". JMLR 2021

IGR for (GD+M)

Heavy ball momentum update

$$\mathbf{x}^{k+1} = \mathbf{x}^k - h\nabla E(\mathbf{x}^k) + \beta(\mathbf{x}^k - \mathbf{x}^{k-1}) \quad \forall k = 1, 2, \dots, n$$

is globally $\underline{O(h)}$ close to the ODE^[7]

$$\tilde{\mathbf{x}}'(t) = -\frac{1}{1-\beta} \nabla E(\tilde{\mathbf{x}}(t)), \quad t \in [0, T] \quad (1)$$

but globally $\underline{O(h^2)}$ closer to the ODE

$$\tilde{\mathbf{x}}'(t) = -\frac{1}{1-\beta} \nabla(E(\tilde{\mathbf{x}}(t))) + \underbrace{\frac{(1+\beta)h}{4(1-\beta)^2} \|\nabla E(\tilde{\mathbf{x}}(t))\|_2^2}_{IGR-M}, \quad t \in [0, T] \quad (2)$$

[7]Kovachki,Stuart, "Continuous time analysis of momentum methods". JMLR 2021

Gradient Descent⁴

$$\mathbf{x}'(t) = -\nabla \left(E(\mathbf{x}(t)) + \frac{h}{4} \|\nabla E(\mathbf{x}(t))\|_2^2 \right)$$

$$\mathbf{x}'(t) = -\nabla \left(\frac{1}{1-\beta} E(\mathbf{x}(t)) + \frac{(1+\beta)h}{4(1-\beta)^3} \|\nabla E(\mathbf{x}(t))\|_2^2 \right)$$

Gradient Descent + Momentum

4) Barrett and Dherin, "Implicit gradient regularization", ICLR-21

5) Smith et al, "On the Origin of Implicit Regularization in Stochastic Gradient Descent", ICLR-21

Gradient Descent⁴

$$\mathbf{x}'(t) = -\nabla \left(E(\mathbf{x}(t)) + \frac{h}{4} \|\nabla E(\mathbf{x}(t))\|^2 \right)$$

$\frac{1}{1-\beta} \times$
force

$\frac{(1+\beta)}{(1-\beta)^3} \times$
IGR

$$\mathbf{x}'(t) = -\nabla \left(\frac{1}{1-\beta} E(\mathbf{x}(t)) + \frac{(1+\beta)h}{4(1-\beta)^3} \|\nabla E(\mathbf{x}(t))\|_2^2 \right)$$

Gradient Descent + Momentum

4) Barrett and Dherin, "Implicit gradient regularization", ICLR-21

5) Smith et al, "On the Origin of Implicit Regularization in Stochastic Gradient Descent", ICLR-21

Gradient Descent⁴

$$\mathbf{x}'(t) = -\nabla \left(E(\mathbf{x}(t)) + \frac{h}{4} \|\nabla E(\mathbf{x}(t))\|^2 \right)$$

$\frac{1}{1-\beta} \times$
force

$\frac{(1+\beta)}{(1-\beta)^3} \times$
IGR

$$\mathbf{x}'(t) = -\nabla \left(\frac{1}{1-\beta} E(\mathbf{x}(t)) + \frac{(1+\beta)h}{4(1-\beta)^3} \|\nabla E(\mathbf{x}(t))\|^2 \right)$$

Faster convergence

Better generalization

Gradient Descent + Momentum

4) Barrett and Dherin, "Implicit gradient regularization", ICLR-21

5) Smith et al, "On the Origin of Implicit Regularization in Stochastic Gradient Descent", ICLR-21

IGR ($\frac{h}{4} \|\nabla E(\mathbf{x})\|^2$) promotes flatness

- Penalizes norm of the gradient.

IGR ($\frac{h}{4} \|\nabla E(\mathbf{x})\|^2$) promotes flatness

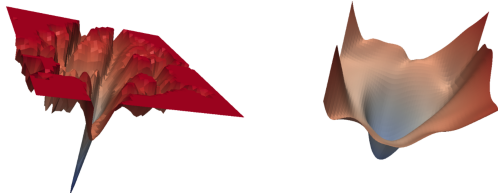
- Penalizes norm of the gradient.
- IGR penalizes "sharpness" $\|\nabla^2 E\|_2$.

IGR ($\frac{h}{4} \|\nabla E(\mathbf{x})\|^2$) promotes flatness

- Penalizes norm of the gradient.
- IGR penalizes "sharpness" $\|\nabla^2 E\|_2$.
- Smaller sharpness \rightarrow flatter minima.

IGR ($\frac{h}{4} \|\nabla E(\mathbf{x})\|^2$) promotes flatness

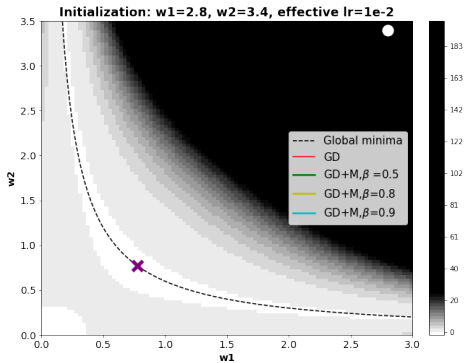
- Penalizes norm of the gradient.
- IGR penalizes "sharpness" $\|\nabla^2 E\|_2$.
- Smaller sharpness \rightarrow flatter minima.
- Flatter minima \rightarrow better generalization^[8].



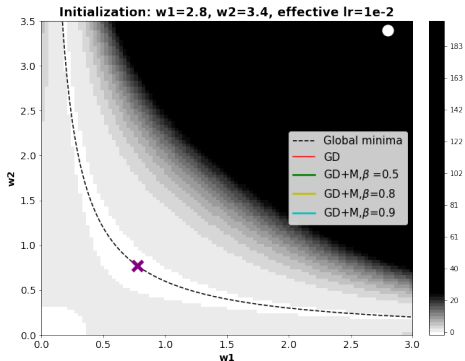
[8] Foret et al, "Sharpness Aware Minimization", ICLR-21

$$\min_{w_1, w_2} E(w_1, w_2) = \frac{1}{2}(y - w_1 w_2 x)^2$$

- Global minima: $w_1 w_2 = \frac{y}{x}$



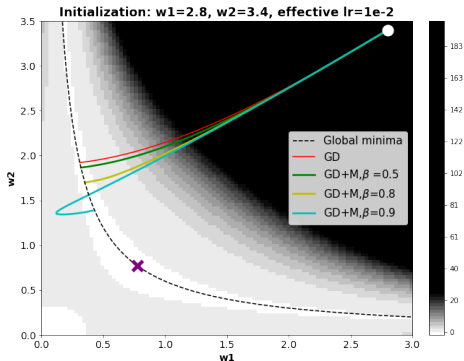
$$\min_{w_1, w_2} E(w_1, w_2) = \frac{1}{2}(y - w_1 w_2 x)^2$$



- Global minima: $w_1 w_2 = \frac{y}{x}$

- Which minima just (GD+M) prefer?

$$\min_{w_1, w_2} E(w_1, w_2) = \frac{1}{2}(y - w_1 w_2 x)^2$$



- Global minima: $w_1 w_2 = \frac{y}{x}$

- Which minima just (GD+M) prefer?

Larger β finds flatter minima. Recall IGR magnified by $\frac{1+\beta}{1-\beta}$.

How does this implicit regularization change in Stochastic Gradient Descent with momentum ?

Modified flow for SGD+M

Heavy-ball momentum SGD updates

$$\begin{cases} \mathbf{x}^{k+1} = \mathbf{x}^k - h\nabla E_k(\mathbf{x}^k) + \beta(\mathbf{x}^k - \mathbf{x}^{k-1}) & k = 1, 2, \dots, n \\ \mathbf{x}^1 = \mathbf{x}^0 - h\nabla E_0(\mathbf{x}^0) \\ \mathbf{x}^0 = \mathbf{x}^{-1} = \mathbf{0} \end{cases} \quad (3)$$

Modified flow for SGD+M

Heavy-ball momentum SGD updates

$$\begin{cases} \mathbf{x}^{k+1} = \mathbf{x}^k - h\nabla E_k(\mathbf{x}^k) + \beta(\mathbf{x}^k - \mathbf{x}^{k-1}) & k = 1, 2, \dots, n \\ \mathbf{x}^1 = \mathbf{x}^0 - h\nabla E_0(\mathbf{x}^0) \\ \mathbf{x}^0 = \mathbf{x}^{-1} = \mathbf{0} \end{cases} \quad (3)$$

is $O(h^2)$ close to the modified loss trajectory

Modified flow for SGD+M

Heavy-ball momentum SGD updates

$$\begin{cases} \mathbf{x}^{k+1} = \mathbf{x}^k - h\nabla E_k(\mathbf{x}^k) + \beta(\mathbf{x}^k - \mathbf{x}^{k-1}) & k = 1, 2, \dots, n \\ \mathbf{x}^1 = \mathbf{x}^0 - h\nabla E_0(\mathbf{x}^0) \\ \mathbf{x}^0 = \mathbf{x}^{-1} = \mathbf{0} \end{cases} \quad (3)$$

is $O(h^2)$ close to the modified loss trajectory

$$\left\{ \begin{array}{l} \mathbf{x}'(t) = -\nabla \left(\underbrace{G_k(\mathbf{x}(t))}_{\text{force}} + \underbrace{\frac{h}{4} \left(\|\nabla G_k(\mathbf{x}(t))\|_2^2 + 2 \sum_{r=0}^{k-1} \beta^{k-r} \|\nabla G_r(\mathbf{x}(t))\|_2^2 \right)}_{IGR_s} \right) \\ \text{for } t_k \leq t < t_{k+1}, \quad \text{where } \underbrace{G_k(\mathbf{x}(t)) = \sum_{r=0}^k \beta^{k-r} E_r(\mathbf{x}(t))}_{\text{Historical gradients}} \end{array} \right.$$

Stochastic Gradient Descent⁵

$$\mathbf{x}'(t) = -\nabla \left(E_k(\mathbf{x}(t)) + \frac{h}{4} \|\nabla E_k(\mathbf{x}(t))\|_2^2 \right)$$

for $t_k \leq t < t_{k+1}$

$E_k = E$

Gradient Descent⁴

$$\mathbf{x}'(t) = -\nabla \left(E(\mathbf{x}(t)) + \frac{h}{4} \|\nabla E(\mathbf{x}(t))\|_2^2 \right)$$

$\frac{1}{1-\beta} \times$
force

$\frac{(1+\beta)}{(1-\beta)^3} \times$
IGR

$\beta=0$

Stochastic Gradient Descent + Momentum

$$\mathbf{x}'(t) = -\nabla \left(\frac{1}{1-\beta} E(\mathbf{x}(t)) + \frac{(1+\beta)h}{4(1-\beta)^3} \|\nabla E(\mathbf{x}(t))\|_2^2 \right)$$

Faster convergence

Better generalization

Gradient Descent + Momentum

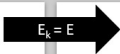
4) Barrett and Dherin, "Implicit gradient regularization", ICLR-21

5) Smith et al, "On the Origin of Implicit Regularization in Stochastic Gradient Descent", ICLR-21

Stochastic Gradient Descent⁵

$$\mathbf{x}'(t) = -\nabla \left(E_k(\mathbf{x}(t)) + \frac{h}{4} \|\nabla E_k(\mathbf{x}(t))\|_2^2 \right)$$

for $t_k \leq t < t_{k+1}$



Gradient Descent⁴

$$\mathbf{x}'(t) = -\nabla \left(E(\mathbf{x}(t)) + \frac{h}{4} \|\nabla E(\mathbf{x}(t))\|_2^2 \right)$$

$\frac{1}{1-\beta} \times$
force

$\frac{(1+\beta)}{(1-\beta)^3} \times$
IGR



$$\mathbf{x}'(t) = -\nabla \left(G_k(\mathbf{x}(t)) + \frac{h}{4} \left(\|\nabla G_k(\mathbf{x}(t))\|_2^2 + 2 \sum_{r=0}^{k-1} \beta^{k-r} \|\nabla G_r(\mathbf{x}(t))\|_2^2 \right) \right)$$

for $t_k \leq t < t_{k+1}$, where $G_k(\mathbf{x}(t)) = \sum_{r=0}^k \beta^{k-r} E_r(\mathbf{x}(t))$

Stochastic Gradient Descent + Momentum

$$\mathbf{x}'(t) = -\nabla \left(\frac{1}{1-\beta} E(\mathbf{x}(t)) + \frac{(1+\beta)h}{4(1-\beta)^3} \|\nabla E(\mathbf{x}(t))\|_2^2 \right)$$

Faster convergence

Better generalization

Gradient Descent + Momentum

4) Barrett and Dherin, "Implicit gradient regularization", ICLR-21

5) Smith et al, "On the Origin of Implicit Regularization in Stochastic Gradient Descent", ICLR-21

Stochastic Gradient Descent⁵

$$\mathbf{x}'(t) = -\nabla \left(E_k(\mathbf{x}(t)) + \frac{h}{4} \|\nabla E_k(\mathbf{x}(t))\|_2^2 \right)$$

for $t_k \leq t < t_{k+1}$

force IGR_s
Average over
historical gradients

$$\mathbf{x}'(t) = -\nabla \left(G_k(\mathbf{x}(t)) + \frac{h}{4} \left(\|\nabla G_k(\mathbf{x}(t))\|_2^2 + 2 \sum_{r=0}^{k-1} \beta^{k-r} \|\nabla G_r(\mathbf{x}(t))\|_2^2 \right) \right),$$

for $t_k \leq t < t_{k+1}$, where $G_k(\mathbf{x}(t)) = \sum_{r=0}^k \beta^{k-r} E_r(\mathbf{x}(t))$

Stochastic Gradient Descent + Momentum

$E_k = E$

Gradient Descent⁴

$$\mathbf{x}'(t) = -\nabla \left(E(\mathbf{x}(t)) + \frac{h}{4} \|\nabla E(\mathbf{x}(t))\|_2^2 \right)$$

$\frac{1}{1-\beta} \times$
force

$\frac{(1+\beta)}{(1-\beta)^3} \times$
IGR

$\beta=0$

$$\mathbf{x}'(t) = -\nabla \left(\frac{1}{1-\beta} E(\mathbf{x}(t)) + \frac{(1+\beta)h}{4(1-\beta)^3} \|\nabla E(\mathbf{x}(t))\|_2^2 \right)$$

Faster convergence

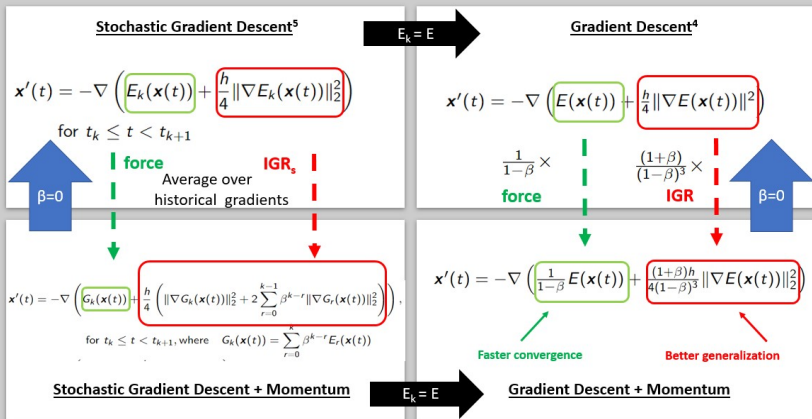
Better generalization

Gradient Descent + Momentum

4) Barrett and Dherin, "Implicit gradient regularization", ICLR-21

5) Smith et al, "On the Origin of Implicit Regularization in Stochastic Gradient Descent", ICLR-21

Family of $O(h^2)$ modified flows

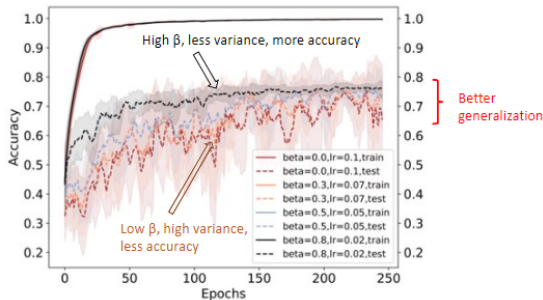


4) Barrett and Dherin, "Implicit gradient regularization", ICLR-21

5) Smith et al, "On the Origin of Implicit Regularization in Stochastic Gradient Descent", ICLR-21

Remarks

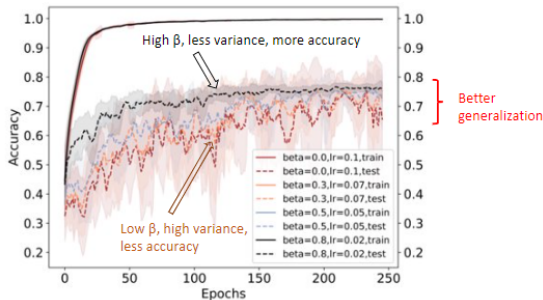
Obs-1: $\beta \uparrow$, variance in test-accuracy \downarrow across epochs .



Remarks

Obs-1: $\beta \uparrow$, variance in test-accuracy \downarrow across epochs .

\Rightarrow **Variance reduction:** $Cov(G_k(\mathbf{x})) = \frac{1-\beta}{1+\beta} Cov(E_k(\mathbf{x}))$

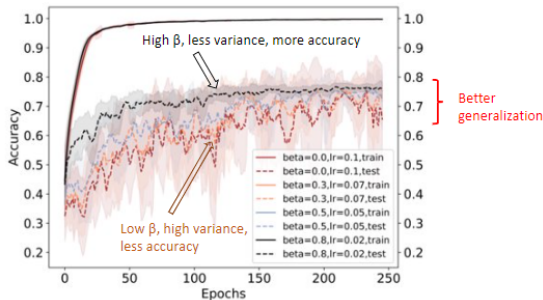


Remarks

Obs-1: $\beta \uparrow$, variance in test-accuracy \downarrow across epochs .

$$\implies \text{Variance reduction: } \text{Cov}(G_k(\mathbf{x})) = \frac{1-\beta}{1+\beta} \text{Cov}(E_k(\mathbf{x}))$$

Obs-2: $\beta \uparrow$, test accuracy \uparrow



Remarks

Obs-1: $\beta \uparrow$, variance in test-accuracy \downarrow across epochs .

$$\Rightarrow \text{Variance reduction: } \text{Cov}(G_k(\mathbf{x})) = \frac{1-\beta}{1+\beta} \text{Cov}(E_k(\mathbf{x}))$$

Obs-2: $\beta \uparrow$, test accuracy \uparrow

$$\Rightarrow \text{Stronger regularization: } \mathbb{E}(IGRM_s)(\mathbf{x}) = \frac{1}{(1-\beta)} \mathbb{E}(IGR_s)(\mathbf{x})$$

