

## A Derivation of Scale-Invariant Properties

**Lemma 1.** For a scale-invariant loss function  $\mathcal{L}(\mathbf{x})$  and  $\alpha > 0$ , we have

$$1. \quad \nabla \mathcal{L}(\alpha \mathbf{x}) = \frac{1}{\alpha} \nabla \mathcal{L}(\mathbf{x}) \quad (28)$$

$$2. \quad \langle \mathbf{x}, \nabla \mathcal{L}(\mathbf{x}) \rangle = 0 \quad (29)$$

$$3. \quad \|\Sigma(\mathbf{x})^{1/2} \mathbf{x}\| = \sqrt{\mathbf{x}^\top \Sigma(\mathbf{x}) \mathbf{x}} = 0. \quad (30)$$

*Proof.* Given  $\mathbf{x} \in \mathbb{R}^d$ , we have  $\mathcal{L}(\alpha \mathbf{x}) = \mathcal{L}(\mathbf{x})$ ,  $\forall \alpha > 0$ . The first property can be proved by taking gradients with respect to  $\mathbf{x}$  on both sides of  $\mathcal{L}(\alpha \mathbf{x}) = \mathcal{L}(\mathbf{x})$ . Taking derivatives with respect to  $\alpha$  for the both sides of  $\mathcal{L}(\alpha \mathbf{x}) = \mathcal{L}(\mathbf{x})$ , we have  $\nabla_\alpha \mathcal{L}(\alpha \mathbf{x}) = 0$ . Left hand side equals  $\nabla_{\alpha \mathbf{x}} \mathcal{L}(\alpha \mathbf{x})^\top \mathbf{x}$ , so the second property follows by taking  $\alpha = 1$ . Since  $\langle \mathbf{x}, \nabla \mathcal{L}(\mathbf{x}) \rangle = 0$  and  $\langle \mathbf{x}, \nabla \mathcal{L}(\mathbf{x}; \mathcal{B}) \rangle = 0$  we have  $\langle \mathbf{x}, \nabla \mathcal{L}(\mathbf{x}; \mathcal{B}) - \nabla \mathcal{L}(\mathbf{x}) \rangle = 0$ . Thus we get the third property  $\mathbf{x}^\top \Sigma(\mathbf{x}) \mathbf{x} = \mathbb{E} \left[ \langle \mathbf{x}, \nabla \mathcal{L}(\mathbf{x}; \mathcal{B}) - \nabla \mathcal{L}(\mathbf{x}) \rangle^2 \right] = 0$ .

## B Additional Experimental Results

**Experimental setup.** We presents results for a ResNet-18 [5] (or scale-invariant ResNet-18) trained on CIFAR-10 [13], a scale-invariant DenseNet-100 [11] trained on CIFAR-100 [13], and scale-invariant ResNet-18 trained on Tiny ImageNet [14]. For ResNet-18 on CIFAR-10 and DenseNet-100 on CIFAR-100, we train for 300 epochs regardless of the number of training samples and the learning rate is divided by 10 at epochs 150 and 225. For ResNet-18 on Tiny ImageNet, we train for 120 epochs and the learning rate is divided by 10 at epochs 60 and 90. For SGDM, we set base value as LR 0.1, WD 0.0001 and search them by multiplying the factor of 2 or  $\sqrt{2}$  and we set momentum coefficient 0.9. For SGD, we set base value as LR 1, which makes the value  $\frac{N\eta\lambda}{B(1-\mu)}$  the same as LR 0.1 for the SGDM setting. We apply data augmentation including padding, random crops and left-right flips. When using a scale-invariant network, all parameters are used to calculate the angular update per epoch. On the other hand, when using an unmodified network, only the parameters of the convolutional layer are used to calculate the angular update per epoch.

### B.1 Scale-Invariant Network

**Experiment on fixed  $B$ .** Table 8 and Table 9 show the results on fixed batch sizes of DenseNet-100 trained on CIFAR-100 and ResNet-18 trained on Tiny ImageNet respectively. In Table 8, we colored yellow to cells which satisfy the tuning factor  $\frac{N\eta\lambda}{B(1-\mu)} = \frac{5}{64}$  and in Table 9, we colored yellow to cells which satisfy the tuning factor  $\frac{N\eta\lambda}{B(1-\mu)} = \frac{5}{32}$ . Bold values represent the highest accuracy for each column.

**Experiment on fixed  $N$ .** Table 10 shows the results on fixed number of training samples of DenseNet-100 trained on CIFAR-100 and we colored yellow to cells which satisfy the tuning factor  $\frac{N\eta\lambda}{B(1-\mu)} = \frac{5}{64}$ .

## B.2 Unmodified Network

Table 11 shows the results on fixed number of training samples of unmodified ResNet-18 trained on CIFAR-10 and we colored yellow to cells which satisfy the tuning factor  $\frac{N\eta\lambda}{B(1-\mu)} = \frac{5}{64}$ . Figure 9 shows the angular update per epoch of the yellow cells of Table 11.

Table 8: Scale-invariant DenseNet-100 on CIFAR-100 with SGDM and  $B=128$ .

LR	WD	3125	6250	12500	25000	50000	LR	WD	3125	6250	12500	25000	50000	LR	WD	3125	6250	12500	25000	50000
6.4	0.0001	22.73					0.1	0.0064	29.44	28.18	22.52	15.17		0.8	0.0008	25.39	20.44			
3.2	0.0001	<b>32.29</b>	41.29	45.99	40.74		0.1	0.0032	<b>36.28</b>	43.08	49.96	48.63	42.68	0.5657	0.0005657	<b>31.75</b>	41.69	46.66	42.67	40.84
1.6	0.0001	<b>32.31</b>	<b>45.93</b>	52.90	20.42	61.50	0.1	0.0016	35.48	<b>47.43</b>	55.94	62.50	64.59	0.4	0.0004	<b>33.74</b>	<b>46.53</b>	54.86	59.94	71.98
0.8	0.0001	29.73	<b>46.16</b>	<b>57.50</b>	65.53	71.95	0.1	0.0008	32.95	47.14	<b>60.21</b>	67.03	72.24	0.2828	0.0002828	30.91	<b>46.89</b>	<b>57.96</b>	67.04	71.98
0.4	0.0001	27.32	44.36	57.17	<b>67.78</b>	73.62	0.1	0.0004	30.88	45.96	58.77	<b>68.27</b>	74.05	0.2	0.0002	29.32	45.92	<b>58.82</b>	<b>67.47</b>	73.58
0.2	0.0001	26.81	43.10	55.67	66.56	<b>74.03</b>	0.1	0.0002	28.87	44.47	57.01	67.52	<b>74.35</b>	0.1414	0.0001414	27.66	43.77	56.30	<b>67.60</b>	<b>74.13</b>
0.1	0.0001	27.84	41.62	55.51	65.56	73.73	0.1	0.0001	27.84	41.62	55.51	65.56	73.73	0.1	0.0001	27.84	41.62	55.51	65.56	73.73

(a) Tuning LR

(b) Tuning WD

(c) Tuning LR,WD

Table 9: Scale-invariant ResNet-18 on Tiny ImageNet with SGDM and  $B=256$ .

LR	WD	12500	25000	50000	100000	LR	WD	12500	25000	50000	100000	LR	WD	12500	25000	50000	100000
6.4	0.0001	30.61	36.90	40.50	38.42	0.1	0.0064	32.31	39.98	41.51	39.05	0.8	0.0008	32.96	38.80	41.58	39.90
3.2	0.0001	<b>36.25</b>	44.52	51.99	58.19	0.1	0.0032	<b>37.04</b>	45.99	53.04	58.67	0.5657	0.0005657	<b>37.02</b>	45.60	52.58	58.74
1.6	0.0001	35.17	<b>47.70</b>	55.77	62.20	0.1	0.0016	35.56	<b>47.75</b>	56.07	62.20	0.4	0.0004	35.88	<b>47.77</b>	56.06	62.75
0.8	0.0001	34.54	46.01	<b>56.67</b>	<b>64.38</b>	0.1	0.0008	34.34	45.43	<b>56.89</b>	<b>64.38</b>	0.2828	0.0002828	34.95	45.97	<b>57.06</b>	<b>64.19</b>
0.4	0.0001	34.68	44.36	55.47	<b>63.67</b>	0.1	0.0004	32.69	43.78	54.59	<b>63.67</b>	0.2	0.0002	34.15	44.50	55.06	<b>64.18</b>
0.2	0.0001	33.04	43.54	53.58	62.08	0.1	0.0002	31.60	43.09	52.57	62.08	0.1414	0.0001414	32.74	42.91	52.65	61.85

(a) Tuning LR

(b) Tuning WD

(c) Tuning LR,WD

Table 10: Scale-invariant DenseNet-100 on CIFAR-100 with SGDM and 50k data samples.

LR	WD	2048	1024	512	256	128	LR	WD	2048	1024	512	256	128	LR	WD	2048	1024	512	256	128
6.4	0.0001	71.27	66.43				0.1	0.0064	73.84	70.91	61.19	41.16		0.8	0.0008	72.95	69.36	57.74		
3.2	0.0001	<b>72.65</b>	72.78	69.51			0.1	0.0032	<b>73.89</b>	73.54	72.33	62.46	42.68	0.5657	0.0005657	<b>73.69</b>	73.49	71.27	59.22	40.84
1.6	0.0001	<b>72.88</b>	<b>73.76</b>	73.17	70.55	61.50	0.1	0.0016	73.29	<b>74.11</b>	73.65	72.50	64.59	0.4	0.0004	73.26	<b>74.37</b>	73.56	71.67	71.98
0.8	0.0001	72.35	73.35	<b>73.20</b>	73.57	71.95	0.1	0.0008	72.12	73.73	<b>74.31</b>	74.12	72.24	0.2828	0.0002828	72.00	73.68	<b>73.93</b>	73.37	71.98
0.4	0.0001	71.37	72.66	<b>73.90</b>	<b>73.78</b>	73.62	0.1	0.0004	70.67	72.41	73.63	<b>74.16</b>	74.05	0.2	0.0002	70.87	72.21	73.67	<b>74.22</b>	73.58
0.2	0.0001	70.45	71.93	72.42	73.57	<b>74.03</b>	0.1	0.0002	69.50	70.80	72.69	73.40	<b>74.35</b>	0.1414	0.0001414	70.00	71.33	72.81	73.67	<b>74.13</b>
0.1	0.0001	69.01	70.31	71.71	72.55	73.73	0.1	0.0001	69.01	70.31	71.71	72.55	73.73	0.1	0.0001	69.01	70.31	71.71	72.55	73.73

(a) Tuning LR

(b) Tuning WD

(c) Tuning LR,WD

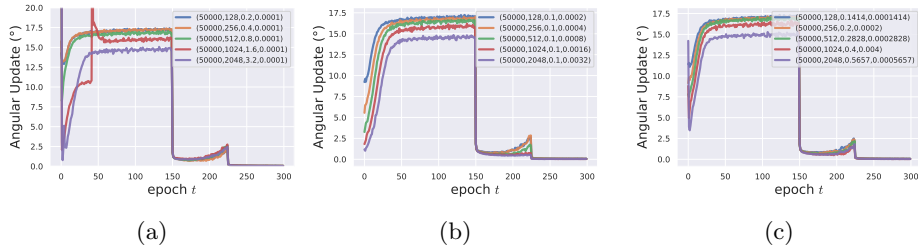


Fig. 9: Angular update per epoch of unmodified ResNet-18 on CIFAR-10 with the number of data 50k (a) tuning LR, (b) tuning WD, (c) tuning both the LR and WD.

Table 11: Unmodified ResNet-18 on CIFAR-10 with SGDM and 50k data samples.

LR	WD	2048	1024	512	256	128	LR	WD	2048	1024	512	256	128	LR	WD	2048	1024	512	256	128
6.4	0.0001	93.98	92.94	89.97			0.1	0.0064	<b>94.82</b>	94.09	90.76	86.30	78.39	0.8	0.0008	94.44	93.40	91.33	81.50	74.19
3.2	0.0001	<b>94.41</b>	<b>94.51</b>	93.86	91.47	63.49	0.1	0.0032	<b>94.76</b>	<b>95.13</b>	94.51	92.42	89.11	0.5657	0.0005657	<b>94.71</b>	94.80	94.11	92.37	89.28
1.6	0.0001	<b>94.42</b>	<b>94.17</b>	<b>94.80</b>	94.05	92.23	0.1	0.0016	94.53	<b>95.06</b>	95.10	94.73	93.45	0.4	0.0004	94.57	<b>94.91</b>	94.91	94.69	93.66
0.8	0.0001	93.71	94.50	<b>94.72</b>	94.88	94.62	0.1	0.0008	93.91	94.60	<b>95.14</b>	<b>95.27</b>	95.08	0.2828	0.0002828	93.89	94.71	<b>94.97</b>	95.10	94.87
0.4	0.0001	93.55	94.18	94.75	<b>94.89</b>	95.16	0.1	0.0004	92.93	94.00	94.63	<b>95.22</b>	95.16	0.2	0.0002	93.19	94.20	94.71	<b>95.28</b>	<b>95.17</b>
0.2	0.0001	92.57	93.56	94.32	94.80	<b>95.25</b>	0.1	0.0002	91.80	93.34	94.11	94.73	<b>95.25</b>	0.1414	0.0001414	92.15	93.34	94.24	94.76	<b>94.99</b>
0.1	0.0001	71.15	81.52	87.52	92.16	94.86	0.1	0.0001	71.15	81.52	87.52	92.16	94.86	0.1	0.0001	71.15	81.52	87.52	92.16	94.86

(a) Tuning LR

(b) Tuning WD

(c) Tuning LR,WD