

CONSTRAINT-BASED REGULARIZATION OF NEURAL NETWORKS



Benedict Leimkuhler, Timothée Pouchon, Tiffany Vlaar, Amos Storkey

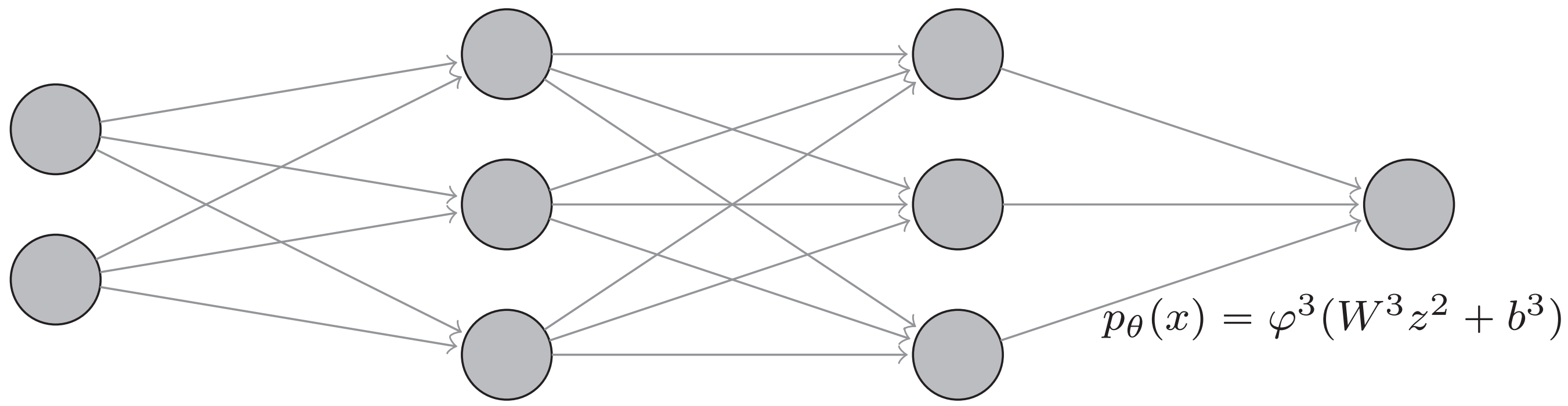
{b.leimkuhler, timothee.pouchon, tiffany.vlaar, a.storkey}@ed.ac.uk

GOAL

Provide a mathematical framework for the regularization of deep neural networks by incorporating constraints into stochastic gradient Langevin dynamics.

CONSTRAINTS FOR NEURAL NETWORKS

Notations Weights $W^\ell \in \mathbb{R}^{d^\ell \times d^{\ell-1}}$, biases $b^\ell \in \mathbb{R}^{d^\ell}$ $1 \leq \ell \leq L$



$$z^0 = x \quad z^1 = \varphi^1(W^1 z^0 + b^1) \quad z^2 = \varphi^2(W^2 z^1 + b^2)$$

Parameter vector: $\theta^\ell = (\text{vect}(W^\ell), b^\ell) \in \mathbb{R}^{n^\ell}$, $\theta = (\theta^1, \dots, \theta^L) \in \mathbb{R}^{|\mathcal{n}|}$, $|\mathcal{n}| = \sum n^\ell$;
Slack variables $\xi \in \mathbb{R}^{n^\xi}$; Training variable $q = (\theta, \xi) \in \mathbb{R}^d$, $d = |\mathcal{n}| + n^\xi$.

We verify that the **gradient of the loss function** $L_X(\theta)$ is proportional to

$$\nabla_{\theta^L}^T p_\theta(x) = F_x^L P_x^L, \quad \nabla_{\theta^\ell}^T p_\theta(x) = F_x^L W^L \dots F_x^{\ell+1} W^{\ell+1} F_x^\ell P_x^\ell \quad \ell \leq L-1,$$

where F_x^j is a sparse matrix evaluated from $\nabla \varphi^j$, and P_x^j is also sparse.

To control the above expressions, we may impose constraints on the parameter space: for a field $g: \mathbb{R}^d \rightarrow \mathbb{R}^m$, the **constraint manifold** is defined as

$$\Sigma = \{q \in \mathbb{R}^d \mid g(q) = 0\}.$$

Circle constraints: restrict each constrained parameter as $|\theta_i^c| \leq r_i$, where $r_i > 0$ is a given hyperparameter:

$$g_i(q) = |\theta_i^c|^2 + |\xi_i|^2 - r_i^2 \quad 1 \leq i \leq m.$$

Orthogonality constraints: for a specific layer ℓ , we define

$$g(q) = \begin{cases} (W^\ell)^T W^\ell - I_{d^{\ell-1}} & \text{if } d^{\ell-1} \leq d^\ell, \\ W^\ell (W^\ell)^T - I_{d^\ell} & \text{otherwise.} \end{cases}$$

CONSTRAINED SDES

We define the potential $V(q) = L_X(\theta)$ and consider the **constrained overdamped Langevin** system

$$dq_t = -\nabla V(q_t) dt + \sqrt{2\beta^{-1}} dW_t - \nabla_q g(q_t) d\lambda_t \quad (\text{CoLA-od})$$

$$0 = g(q_t)$$

whose invariant measure is $d\nu_\Sigma = Z^{-1} e^{-\beta V(q)} d\sigma_\Sigma$, where σ_Σ is the surface measure of Σ and Z is the normalization constant.

Theorem (Exponential convergence to equilibrium). Assume that there exists $\rho > 0$ such that

$$\text{Ric}_g + \beta \nabla_g^2 V \geq \rho g. \quad (g \text{ Riemannian metric, Ric}_g \text{ Ricci curvature, } \nabla_g^2 V \text{ Hessian})$$

Then there exists $R > 0$ such that $\langle \phi \rangle_{\nu_\Sigma} = \int_\Sigma \phi d\nu_\Sigma$

$$\int_\Sigma |\mathbb{E}(\phi(q_t) \mid q_0) - \langle \phi \rangle_{\nu_\Sigma}|^2 d\nu_\Sigma(q_0) \leq C(\phi) e^{-2R\beta^{-1}t} \quad \forall \phi \in H^1(\nu_\Sigma),$$

where $C(\phi)$ depends only on ϕ .

An alternative to (CoLA-od) is the second order dynamics given by the **constrained underdamped Langevin** system:

$$dq_t = p_t dt$$

$$dp_t = (-\nabla_q V(q_t) - \gamma p_t) dt + \sqrt{2\gamma\beta^{-1}} dW_t - \nabla_q g(q_t) d\lambda_t \quad (\text{CoLA-ud})$$

$$0 = g(q_t)$$

whose invariant measure is closely related to ν_Σ .

EXAMPLE OF DISCRETIZATION

Discretization of CoLA-od with orthogonality constraint (o-CoLA-od):

$$Q^{(0)} = Q_n - h \nabla_Q V(Q_n) + \sqrt{2\tau h} R_n, \quad (\tau = \beta^{-1})$$

$$\text{for } k = 0 \text{ to } K-1: \quad Q^{(k+1)} = Q^{(k)} - \frac{1}{2} Q_n ((Q^{(k)})^T Q^{(k)} - I_s),$$

$$Q_{n+1} = Q^{(K)}.$$

NUMERICAL EXPERIMENTS

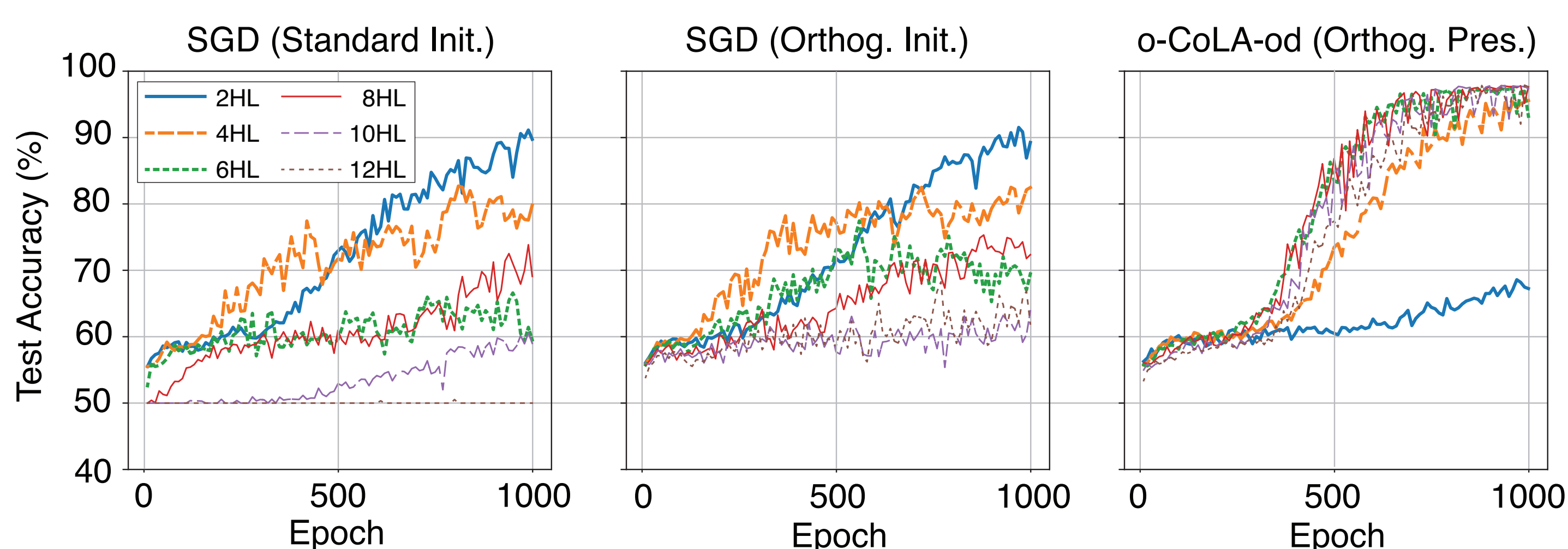
Spiral dataset: points in \mathbb{R}^2 ; 500 training points, 1000 test points.

Model: MLPs with variable depth, 100-nodes ReLU in each hidden layers.

Training: SGD with standard initialization (left), SGD with orthogonal initialization (middle) and o-CoLA-od with $\tau = 0$ (right).

Hyperparameters: for all methods, $h = 0.1$, 5% subsampling.

Results are averaged over 10 runs.



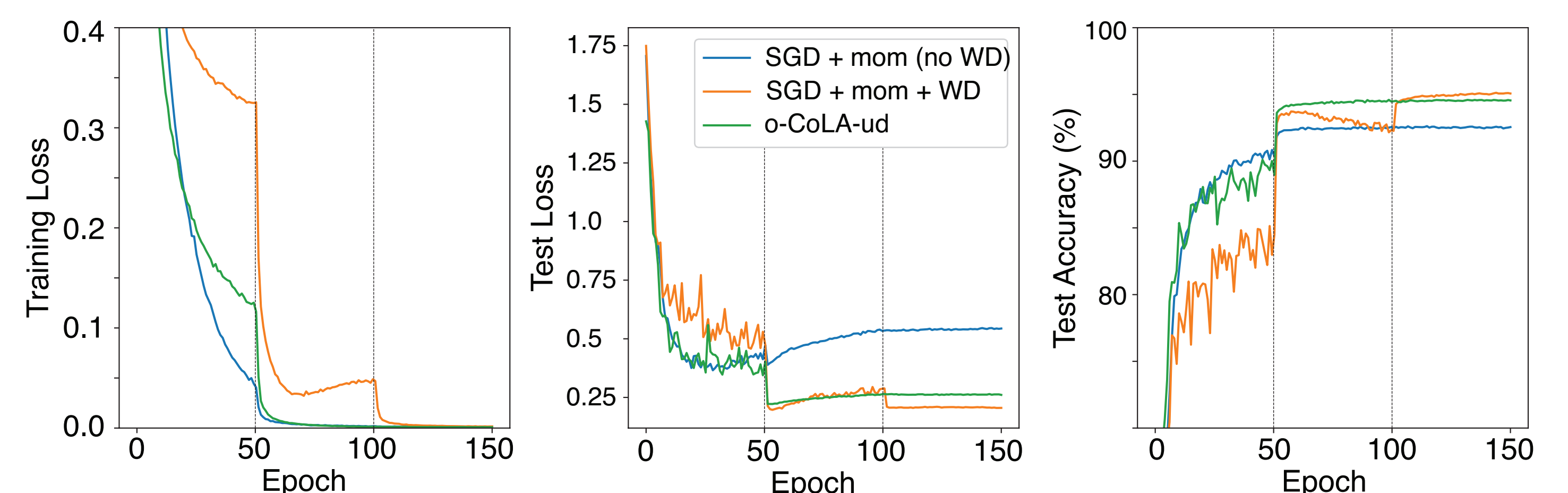
CIFAR-10 dataset: 32×32 images; 50K training images, 5K test images.

Model: ResNet-34 with BatchNorm.

Training: SGD-m vs. o-CoLA-ud with $\tau = 0$.

Hyperparameters: batchsize 128, lr-decay by a factor 10 every 50 epochs; for SGD-m, $h = 0.1$, mom. = 0.9; for o-CoLA-ud $\gamma = 0.5$, learning rate was re-scaled to match the parameters of SGD-m.

Results are averaged over 5 runs.



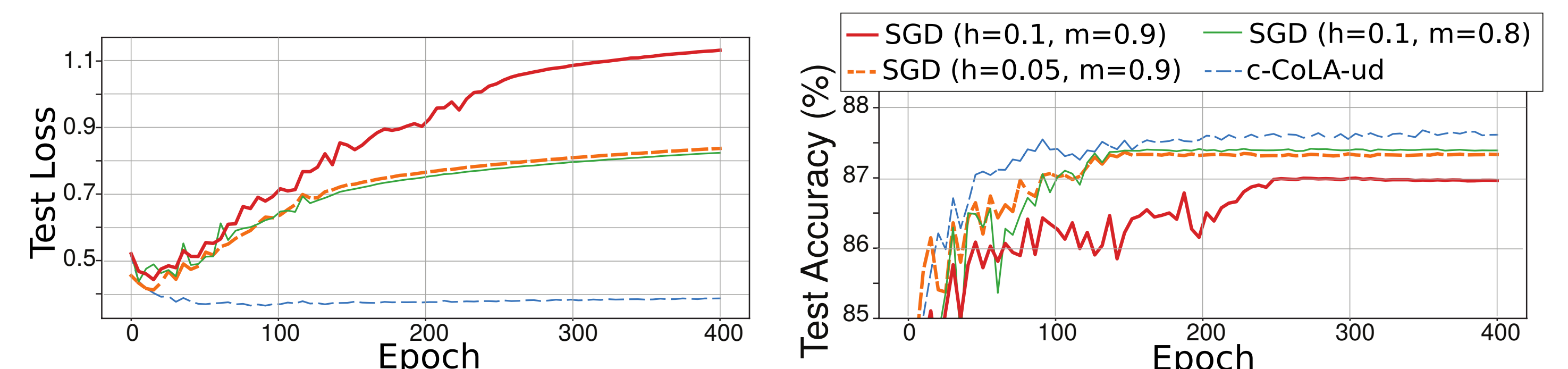
Fashion-MNIST dataset: 28×28 images; # training images reduced to 10K, 60K test images.

Model: 1000-node SHLP.

Training: SGD-m vs. c-CoLA-ud with $\tau = 0$.

Hyperparameters: batchsize 128; for c-CoLA-ud, $h = 0.3$, $\tau = 0$, $\gamma = 1$, $r_0 = 0.05$, $r_1 = 0.1$.

Results are averaged over 5 runs.



REFERENCE

[1] B. LEIMKUHLE, T. POUCHON, T. VLAAR, AND A. STORKEY, *Constraint-based regularization of neural networks*. arXiv:2006.10114, 2020.