

0.1 Basics

Algorithm 1 Stochastic gradient descent

Require: learning rate λ **Ensure:** a trained neural network

- 1: initialize the network, dataset and training parameters
 - 2: **while** stopping criteria is not met **do**
 - 3: sample minibatch of m examples $x^{(1)}, \dots, x^{(n)}$
 - 4: compute gradient estimate $\hat{g} = \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$
 - 5: apply parameter update $\theta = \theta - \lambda \cdot \hat{g}$
 - 6: **end while**
 - 7: **return:** the trained network
-

Algorithm 2 Stochastic gradient descent with Momentum

Require: learning rate α **Require:** momentum parameter m **Ensure:** a trained neural network

- 1: initialize the network, dataset and training parameters
 - 2: **while** stopping criteria is not met **do**
 - 3: sample minibatch of m examples $x^{(1)}, \dots, x^{(n)}$
 - 4: compute gradient estimate $\hat{g} = \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$
 - 5: compute velocity update $v = m \cdot v - \alpha \hat{g}$
 - 6: apply parameter update $\theta = \theta - v$
 - 7: **end while**
 - 8: **return:** the trained network
-

0.2 Methods

Algorithm 3 Machine Learning with distancing

Require: a set of parameters θ and a dataset

Ensure: a assignment of θ which maximizes performance

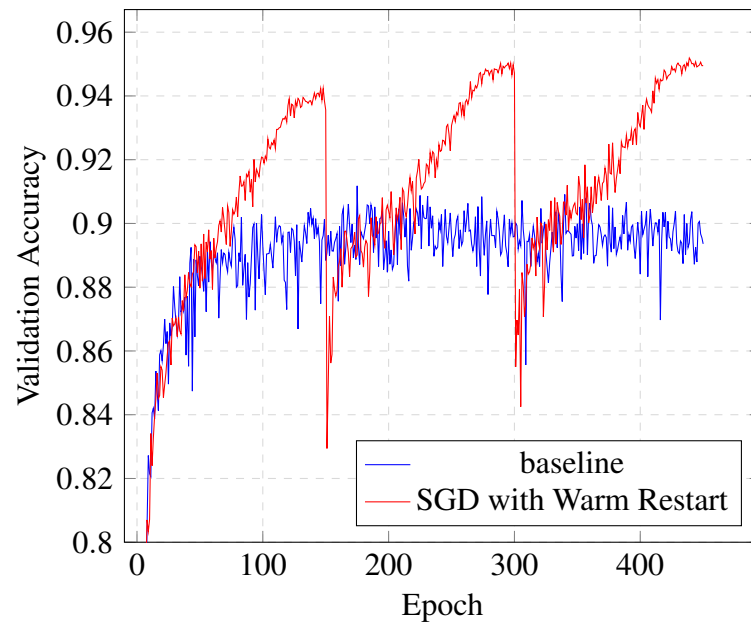
- 1: initialize the network, dataset and training parameters
 - 2: **for** $i \leftarrow 1$ **to** desired number of epochs **do**
 - 3: compute foward and backward pass of training data
 - 4: update parameter values with optimizer
 - 5: **end for**
 - 6: create checkpoint we want to distance from
 - 7: **for** $i \leftarrow$ next epoch **to** end **do**
 - 8: **for** checkpoint **in** list of checkpoints **do**
 - 9: compute parameter update which increases the distance to checkpoint
 - 10: **end for**
 - 11: update parameter values with optimizer
 - 12: **end for**
 - 13: **return:** the final assignment of θ
-

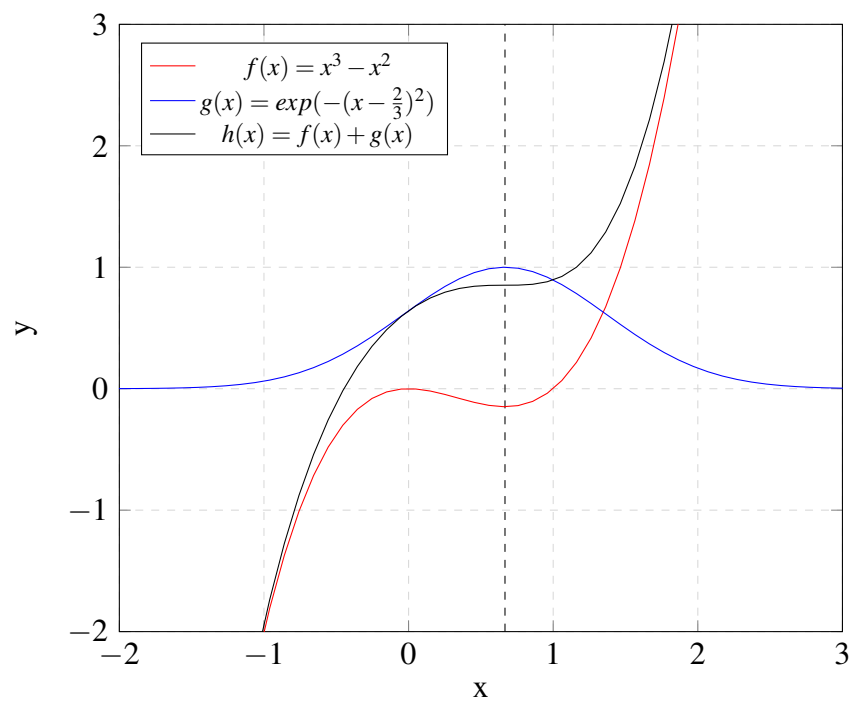
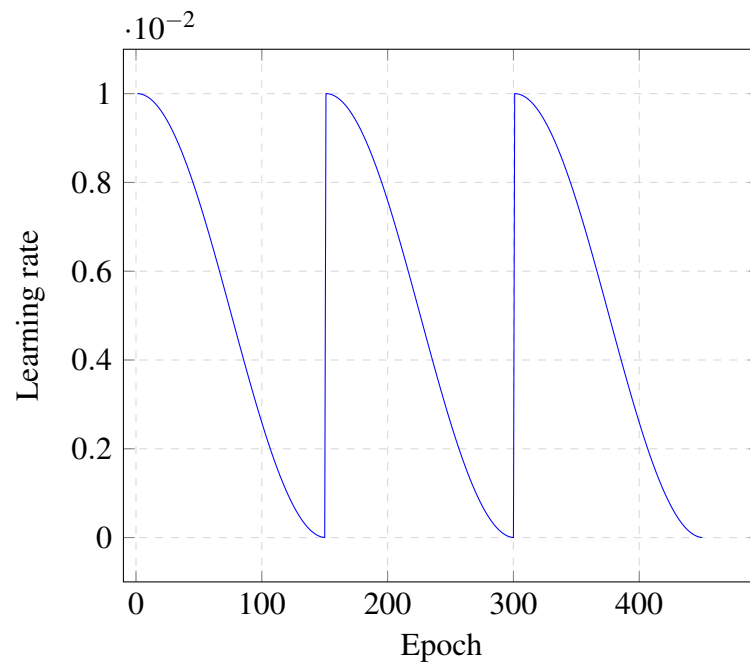
Algorithm 4 Update step with distancing

Require: learning rate α , distance hyperparameters s and σ

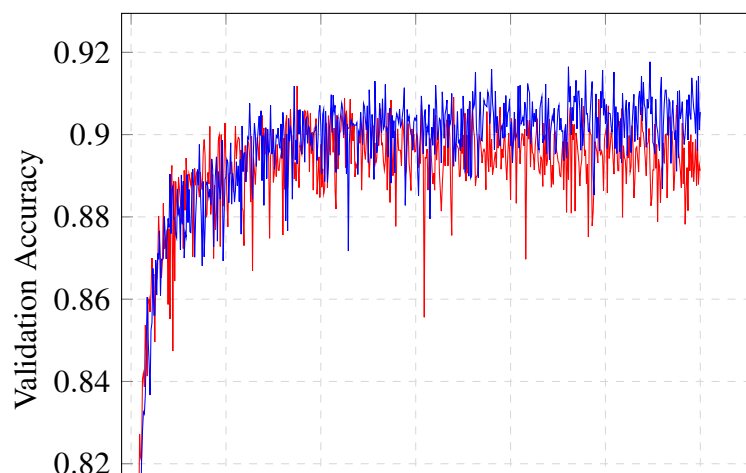
Ensure: a trained neural network

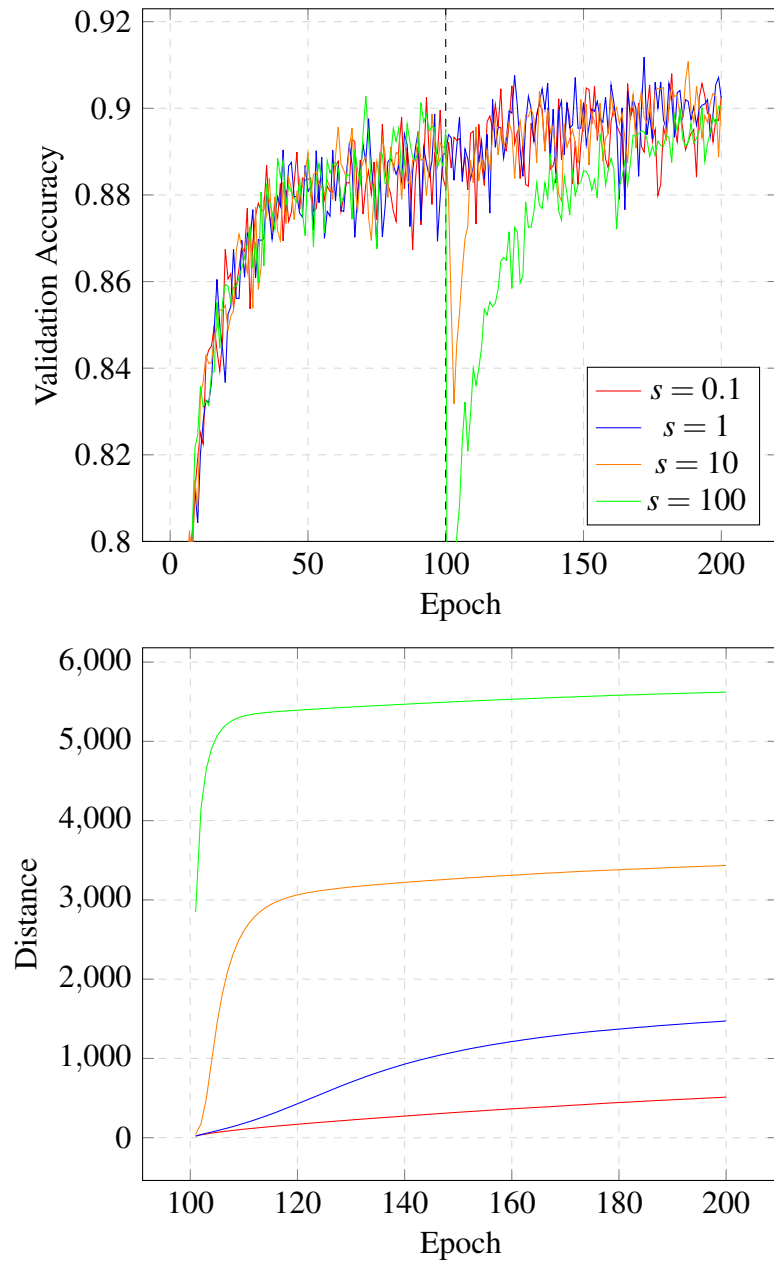
- 1: initialize the network, dataset and training parameters
 - 2: **while** stopping criteria is not met **do**
 - 3: sample minibatch of m examples $x^{(1)}, \dots, x^{(n)}$
 - 4: compute gradient estimate $\hat{g} = \nabla_{\theta} \frac{1}{m} \sum_i L(f(x^{(i)}; \theta), y^{(i)}) + \frac{1}{c} \sum_c s_c \cdot distance(\theta, \theta_c)$
 - 5: apply parameter update $\theta = \theta - \alpha \cdot \hat{g}$
 - 6: **end while**
 - 7: **return:** the trained network
-

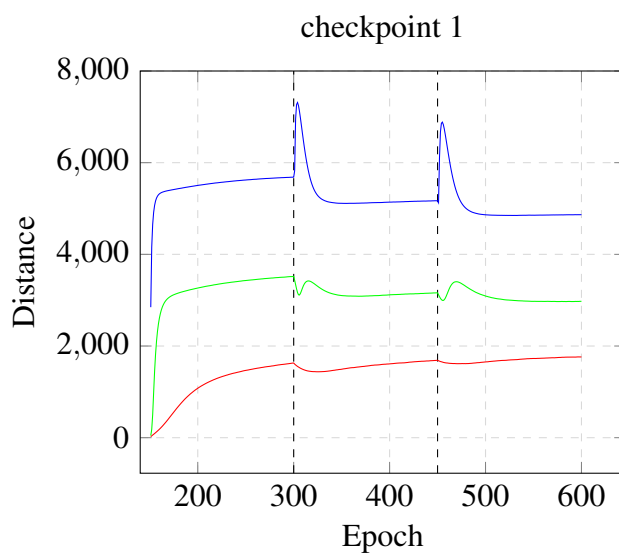
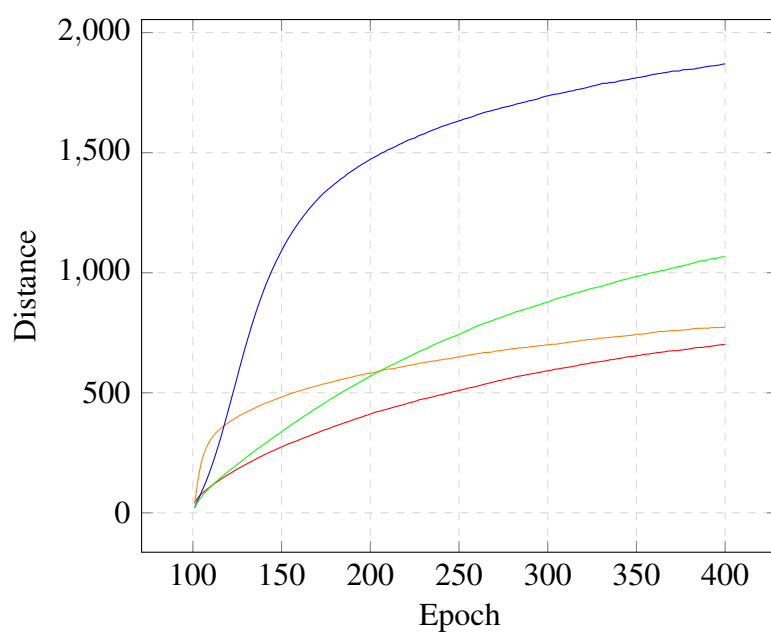
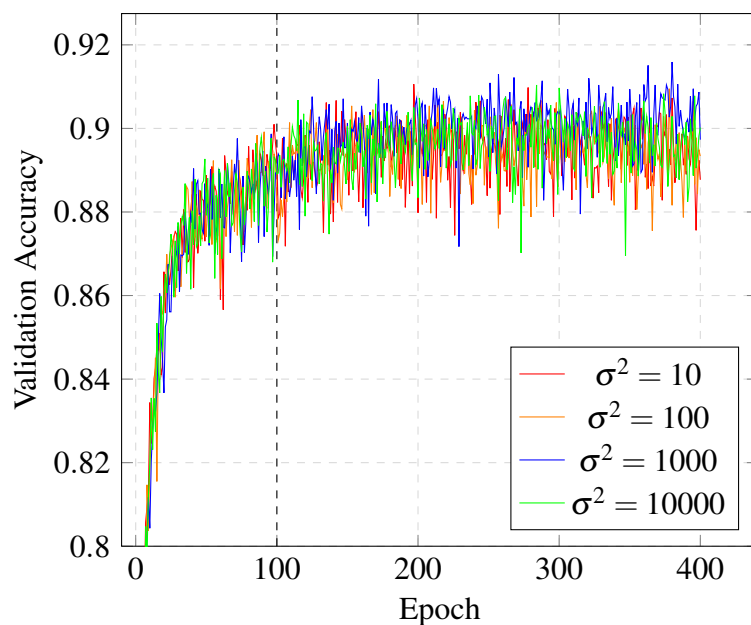




0.3 results







$\cdot 10^4$ without L_2 Regularization

