

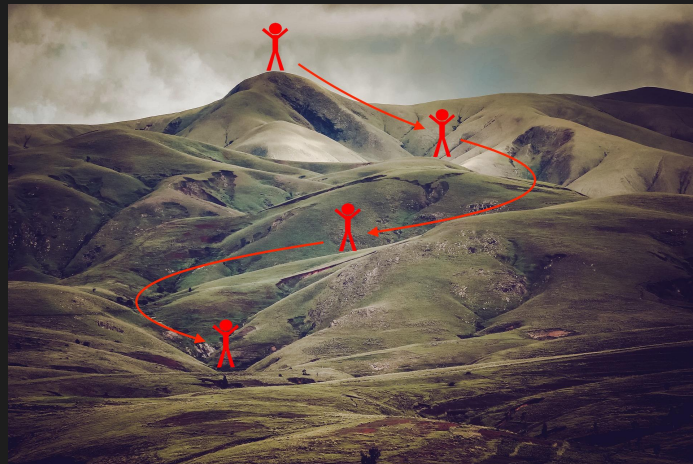


IDL

INTRODUCTION TO DEEP LEARNING

Optimization

- Mini-batch gradient descent
- Optimizer :
 - SGD
 - GD with momentum
 - RMSProp
 - AdaGrad
 - Adam
- Learning rate decay



Mini-batch gradient descent

Exponentially weighted averages

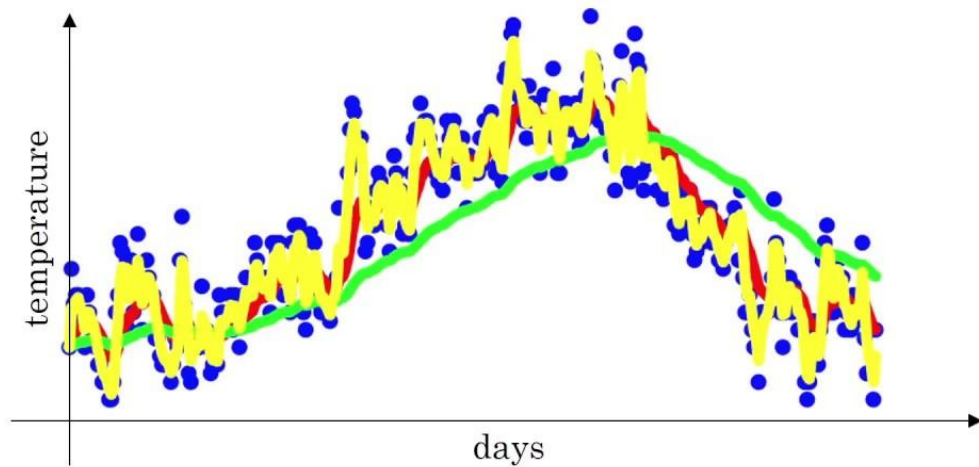
Exponentially weighted averages

$$v_t = \beta v_{t-1} + (1 - \beta)\theta_t$$

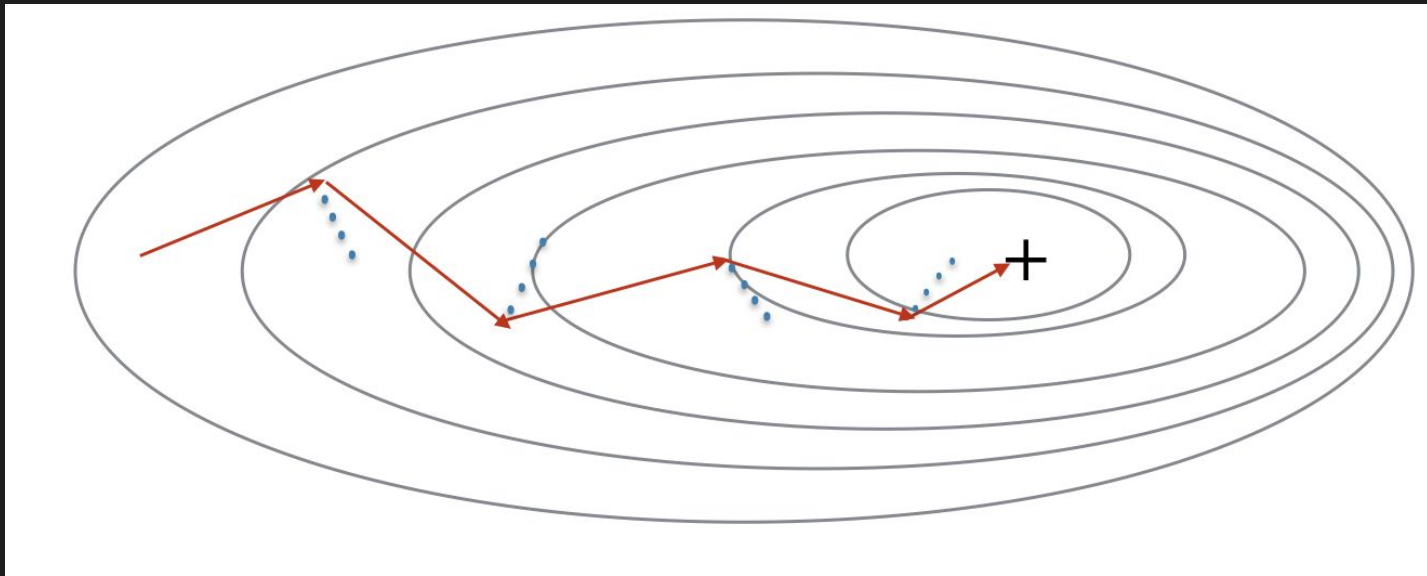
$$\beta = 0.9$$

$$0.98$$

$$0.5$$



GD with momentum



RMSProp

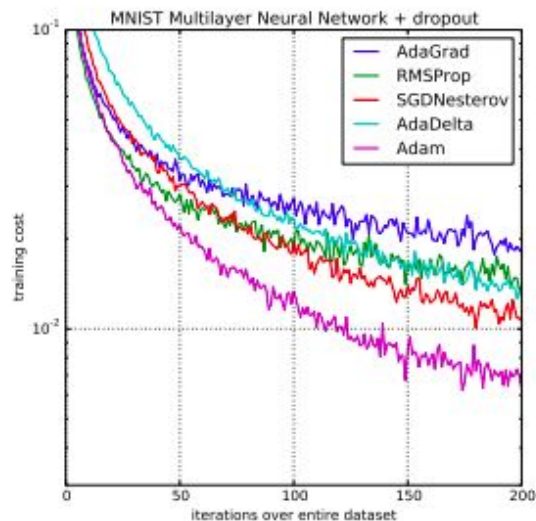
$$v_{dw} = \beta \cdot v_{dw} + (1 - \beta) \cdot dw^2$$

$$v_{db} = \beta \cdot v_{db} + (1 - \beta) \cdot db^2$$

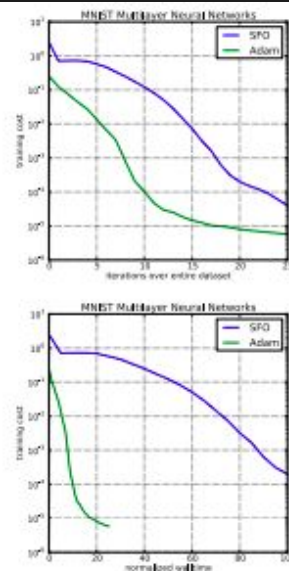
$$W = W - \alpha \cdot \frac{dw}{\sqrt{v_{dw}} + \epsilon}$$

$$b = b - \alpha \cdot \frac{db}{\sqrt{v_{db}} + \epsilon}$$

Adam

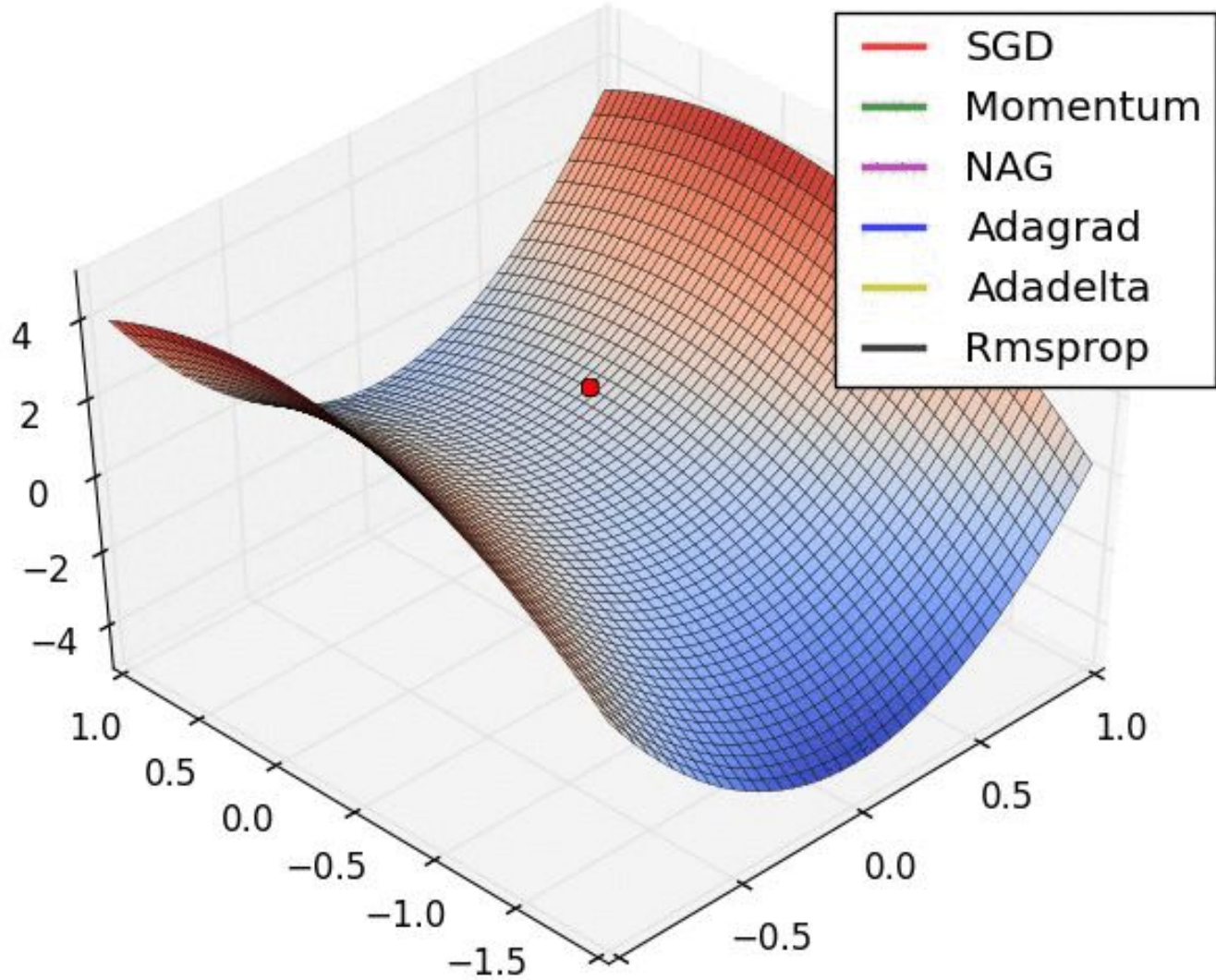


(a)



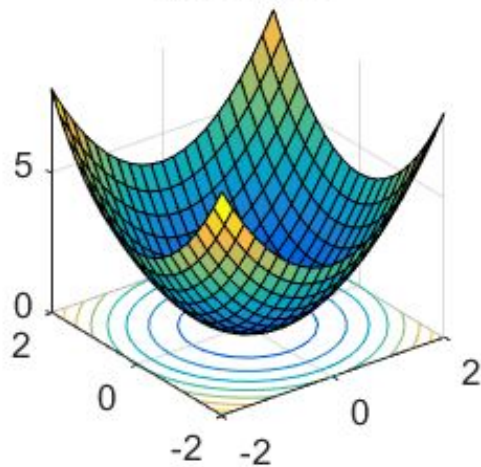
(b)

Figure 2: Training of multilayer neural networks on MNIST images. (a) Neural networks using dropout stochastic regularization. (b) Neural networks with deterministic cost function. We compare with the sum-of-functions (SFO) optimizer (Sohl-Dickstein et al., 2014)

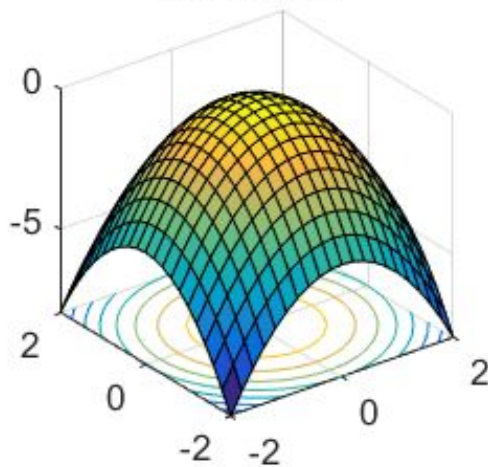


The problem of local optima

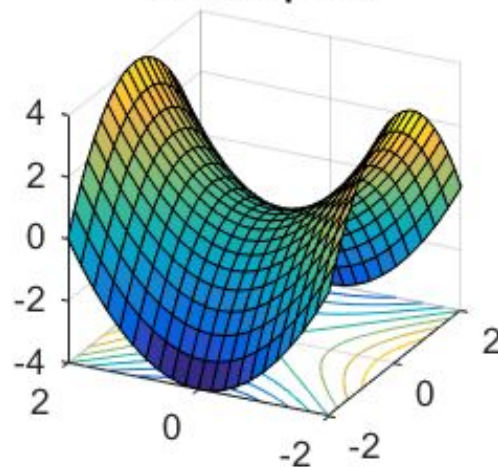
local min



local max



saddle point



Regularization

Hyperparameter

Batch Normalization
