

# Optimization for Training Deep Models

Sargur N. Srihari  
[srihari@cedar.buffalo.edu](mailto:srihari@cedar.buffalo.edu)

# Topics in Optimization

- Overview of Optimization in Deep Learning
- How learning differs from optimization
  - Risk, empirical risk and surrogate loss
  - Batch, minibatch, data shuffling
- Challenges in neural network optimization
- Basic Algorithms
- Parameter initialization strategies
- Algorithms with adaptive learning rates
- Approximate second-order methods
- Optimization strategies and meta-algorithms <sup>2</sup>

# Optimization in Deep Learning

- Optimization is encountered often in ML
  1. Inference with PCA requires optimization
    - Encoding:  $f(\mathbf{x}) = \mathbf{c}$ , Decoding:  $\mathbf{x} \approx g(f(\mathbf{x}))$ ,  $g(\mathbf{c}) = D\mathbf{c}$
    - Optimal  $\mathbf{c}^* = \operatorname{argmin}_{\mathbf{c}} \|\mathbf{x} - g(\mathbf{c})\|_2$ , Reconstruction:  $g(f(\mathbf{x})) = DD^T \mathbf{x}$
  2. Optimization to write proofs or design algorithms
    - In linear regression: sum-of-squared errors objective is same as obtained using maximum likelihood with Gaussian noise
  3. Neural network training
    - Most difficult optimization is neural network training
    - Weight decay minimization:

$$J(\mathbf{w}) = \lambda \|\mathbf{w}\|_2^2 - E_{\mathbf{x}, y \sim \hat{p}_{data}} \log p_{\text{model}}(y | \mathbf{x})$$

# Neural network optimization is difficult

- Commonly months of time on 100s of machines to solve a single instance of neural network training
- So specialized optimization techniques developed

# Our focus on particular case of optimization

- Find parameters  $\theta$  of a neural network that significantly reduce a cost function  $J(\theta)$ 
  - Which typically includes:
    - a performance measure evaluated on training set, e.g.,
      - $J(w) = E_{x,y \sim \hat{p}_{data}} \log p_{\text{model}}(y | x)$  where  $p_{\text{model}}(y | x)$  is a likelihood function
      - For linear regression  $p_{\text{model}}(y | x) = N(y; x^T w + b, 1)$  and  $J(w)$  is same as sum-of-squared errors
    - additional regularization terms, e.g.,  $\lambda \|w\|_2^2$

# Plan of Discussion of Optimization

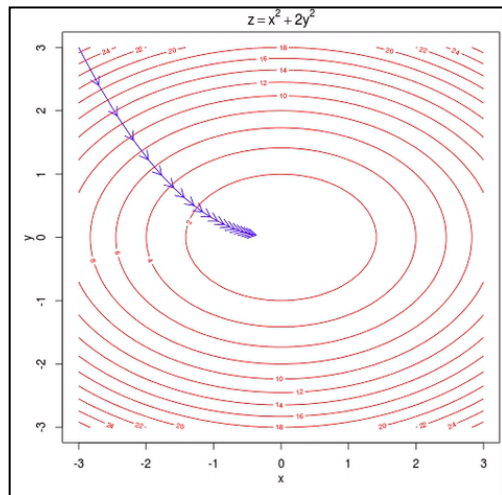
1. How training optimization differs from pure optimization
2. Challenges that make optimization of neural networks difficult
3. Several practical algorithms including
  1. Optimization algorithms
  2. Strategies for initializing parameters
    - Most advanced algorithms
      - adapt learning rates or
      - leverage *second derivatives* of cost function
4. Combine simple optimization algorithms into higher-level procedures

# Summary of Optimization Methods

- Movies:

<http://hduongtrong.github.io/2015/11/23/coordinate-descent/>

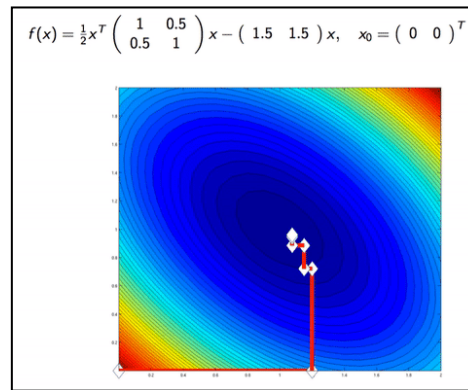
## Gradient Descent



$$g = \frac{1}{M} \nabla_{\theta} \sum_{i=1}^M L(x^{(i)}, y^{(i)}, \theta)$$

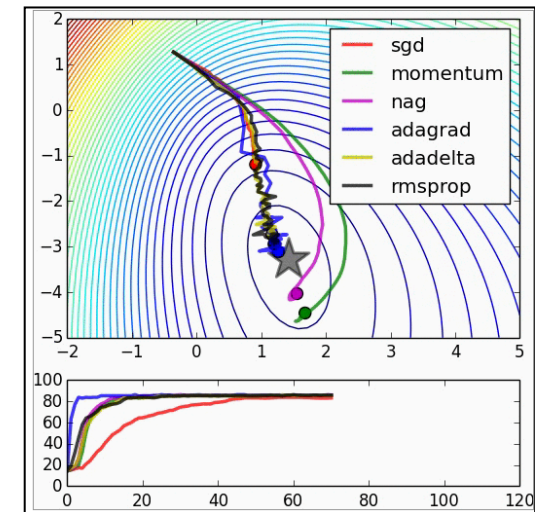
$$\theta \leftarrow \theta - \varepsilon g$$

## Coordinate Descent



Minimize  $f(x)$  wrt a single variable,  $x_i$ , then wrt  $x_j$  etc

## SGD



$$g = \frac{1}{m'} \nabla_{\theta} \sum_{i=1}^{m'} L(x^{(i)}, y^{(i)}, \theta)$$

$$\theta \leftarrow \theta - \varepsilon g$$