

Approximate Second Order Methods

Sargur N. Srihari
srihari@cedar.buffalo.edu

Topics in Optimization for Deep Models

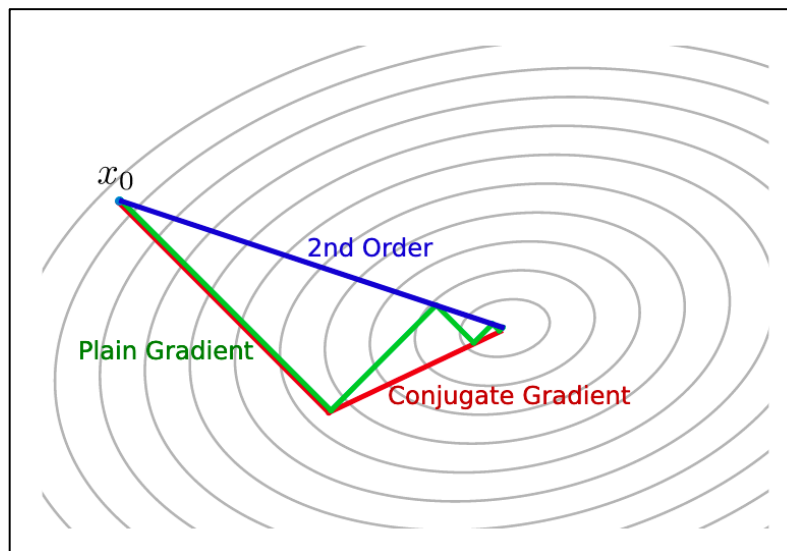
- Importance of Optimization in machine learning
- How learning differs from optimization
- Challenges in neural network optimization
- Basic Optimization Algorithms
- Parameter initialization strategies
- Algorithms with adaptive learning rates
- **Approximate second-order methods**
- Optimization strategies and meta-algorithms

Topics in Second Order Methods

1. Overview
2. Newton's Method
3. Conjugate Gradients
 - Nonlinear Conjugate Gradients
4. BFGS
 - Limited Memory BFGS

Why second-order methods?

- Better direction



- Better step-size

- A full step jumps directly to the minimum of the local squared approx.
- often this is already a good heuristic
- additional step size reduction and dampening are straight-forward

Overview

- We discuss here second order methods of training deep networks
- The only objective function examined is empirical risk:
 - Empirical risk, with m training examples, is

$$J(\boldsymbol{\theta}) = E_{(\mathbf{x}, y) \sim \hat{p}_{data}} (L(f(\mathbf{x}; \boldsymbol{\theta}), y)) = \frac{1}{m} \sum_{i=1}^m L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), y^{(i)})$$

$f(\mathbf{x}; \boldsymbol{\theta})$ is the predicted output when the input is \mathbf{x}

y is target output

L is the per-example loss function

- Methods extend readily to other objective functions such as those that include parameter regularization

Newton's Method

- In contrast to first order gradient methods, second order methods make use of second derivatives to improve optimization
- Most widely used second order method is Newton's method
- It is described in more detail here emphasizing neural network training
- It is based on Taylor's series expansion to approximate $J(\theta)$ near some point θ_0 ignoring derivatives of higher order

Newton Update Rule

- Taylor's series to approximate $J(\theta)$ near θ_0

$$J(\theta) \approx J(\theta_0) + (\theta - \theta_0)^T \nabla_{\theta} J(\theta_0) + \frac{1}{2} (\theta - \theta_0)^T H(\theta - \theta_0)$$

- where H is the Hessian of J wrt θ evaluated at θ_0
- Solving for the critical point of this function we obtain the Newton parameter update rule

$$\theta^* = \theta_0 - H^{-1} \nabla_{\theta} J(\theta_0)$$

- Thus for a quadratic function (with positive definite H) by rescaling the gradient by H^{-1} Newton's method directly jumps to the minimum
- If objective function is convex but not quadratic (there are higher-order terms) this update can be iterated yielding the training algorithm given next

Training Algorithm associated with Newton's Method

Algorithm: Newton's method with objective:

$$J(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), \mathbf{y}^{(i)})$$

Require: Initial parameter $\boldsymbol{\theta}_0$

Require: Training set of m examples

while stopping criterion not met **do**

 Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}} \sum_i L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), \mathbf{y}^{(i)})$

 Compute Hessian: $\mathbf{H} \leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}}^2 \sum_i L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), \mathbf{y}^{(i)})$

 Compute Hessian inverse: \mathbf{H}^{-1}

 Compute update: $\Delta \boldsymbol{\theta} = -\mathbf{H}^{-1} \mathbf{g}$

 Apply update: $\boldsymbol{\theta} = \boldsymbol{\theta} + \Delta \boldsymbol{\theta}$

end while

Positive Definite Hessian

- For surfaces that are not quadratic, as long as the Hessian remains positive definite, Newton's method can be applied iteratively
- This implies a two-step procedure:
 - First update or compute the inverse Hessian (by updating the quadratic approximation)
 - Second, update the parameters according to

$$\theta^* = \theta_0 - H^{-1} \nabla_{\theta} J(\theta_0)$$

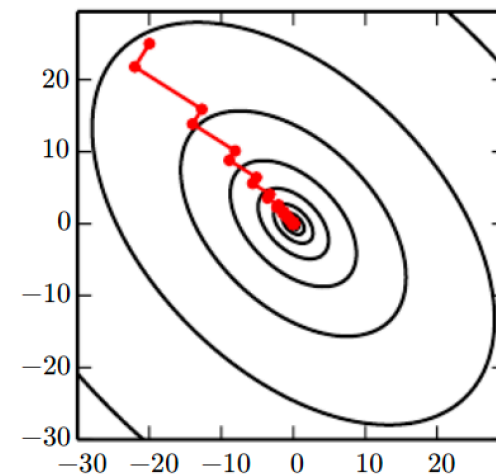
Regularizing the Hessian

- Newton's method is appropriate only when the Hessian is positive definite
 - In deep learning the surface of the objective function is nonconvex
 - Many saddle points: problematic for Newton's method
- Can be avoided by regularizing the Hessian
 - Adding a constant α along the Hessian diagonal

$$\theta^* = \theta_0 - [H(f(\theta_0)) + \alpha I]^{-1} \nabla_{\theta} f(\theta_0)$$

Motivating Conjugate Gradients

- Method to efficiently avoid calculating H^{-1}
 - By iteratively descending conjugate directions
- Arises from steepest descent for quadratic bowl has an ineffective zig-zag pattern
 - Since each line direction is orthogonal to previous
 - Let previous search direction be \mathbf{d}_t
 - Then $\nabla_{\theta} J(\theta) \mathbf{d}_{t-1} = 0$
 - Current search direction will have no contribution in direction \mathbf{d}_{t-1}
 - Thus \mathbf{d}_t is orthogonal to \mathbf{d}_{t-1}
 - Method of conjugate gradients addresses this problem



Imposing Conjugate Directions

- We seek to find a search direction that is conjugate to the previous line search direction
- At iteration t the next search direction \mathbf{d}_t takes the form $\mathbf{d}_t = \nabla_{\theta} J(\theta) + \beta_t \mathbf{d}_{t-1}$
- Directions \mathbf{d}_t and \mathbf{d}_{t-1} are conjugate if $\mathbf{d}_t^T H \mathbf{d}_{t-1} = 0$
- Methods for imposing conjugacy

– Fletcher-Reeves

$$\beta_t = \frac{\nabla_{\theta} J(\theta_t)^T \nabla_{\theta} J(\theta_t)}{\nabla_{\theta} J(\theta_{t-1})^T \nabla_{\theta} J(\theta_{t-1})}$$

– Polak-Ribiere

$$\beta_t = \frac{(\nabla_{\theta} J(\theta_t) - \nabla_{\theta} J(\theta_{t-1}))^T \nabla_{\theta} J(\theta_t)}{\nabla_{\theta} J(\theta_{t-1})^T \nabla_{\theta} J(\theta_{t-1})}$$

Conjugate gradient algorithm

Algorithm The conjugate gradient method

Require: Initial parameters θ_0

Require: Training set of m examples

Initialize $\rho_0 = \mathbf{0}$

Initialize $g_0 = \mathbf{0}$

Initialize $t = 1$

while stopping criterion not met **do**

Initialize the gradient $g_t = \mathbf{0}$

Compute gradient: $g_t \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$

Compute $\beta_t = \frac{(g_t - g_{t-1})^\top g_t}{g_{t-1}^\top g_{t-1}}$ (Polak-Ribière)

(Nonlinear conjugate gradient: optionally reset β_t to zero, for example if t is a multiple of some constant k , such as $k = 5$)

Compute search direction: $\rho_t = -g_t + \beta_t \rho_{t-1}$

Perform line search to find: $\epsilon^* = \operatorname{argmin}_{\epsilon} \frac{1}{m} \sum_{i=1}^m L(f(x^{(i)}; \theta_t + \epsilon \rho_t), y^{(i)})$

(On a truly quadratic cost function, analytically solve for ϵ^* rather than explicitly searching for it)

Apply update: $\theta_{t+1} = \theta_t + \epsilon^* \rho_t$

$t \leftarrow t + 1$

end while

The BFGS Algorithm

- Broyden-Fletcher-Goldfarb-Shanno (BFGS)
 - Newton's method without the computational burden
 - It is similar to the conjugate gradient method
 - More direct approach to approximating Newton's update
- Recall Newton's update: $\theta^* = \theta_0 - H^{-1} \nabla_{\theta} J(\theta_0)$
 - where H is the Hessian of J wrt θ evaluated at θ_0
 - Primary difficulty is computation of H^{-1}
 - BFGS is quasi Newton: approximates H^{-1} by matrix M_t that is iteratively refined by low-rank updates
 - Once the inverse Hessian M_t is updated, the direction of descent ρ_t is determined by $\rho_t = M_t g_t$
 - Final update to parameters is $\theta_{t+1} = \theta_t + \epsilon^* \rho_t$