# Optimization Algorithms

## Simple cases
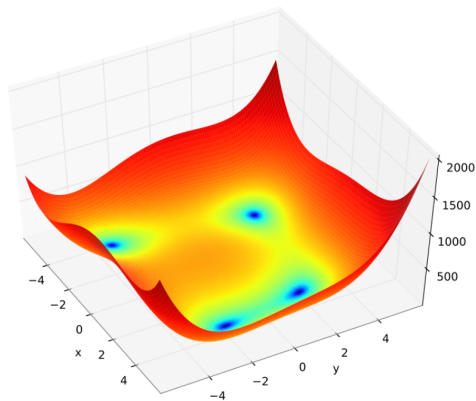
– Local minima
– Reasonable initial guess
– Wide attractive basins

## Hard cases

– Noisy function evaluations
– High dimensionality

## Popular representatives

– Newton
– Steepest descent
– BFGS
– L-BFGS



The Himmelblau function

$$\underset{\text{Previous value}}{\underline{a_n}} - \underset{\text{Step size}}{s} \, \underset{\text{Hessian}}{\left[\underline{\nabla^2 \underset{\text{Target function}}{f}(a_n)}\right]^{-1}} \underset{\text{Gradient}}{\nabla f(a_n)} \tag{28}$$

## Variants

– Scale step size
– Stochastic Newton

## Problems

– Large Hessian and inversion expensive
– Slow with a fixed step
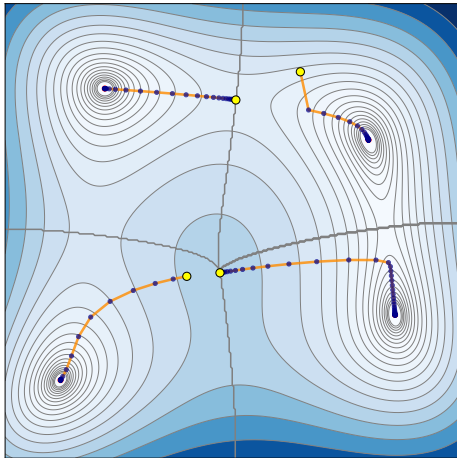
$$a_n - s \nabla f(a_n) \tag{29}$$

Previous value | | Step size

### Variants
– Adjust step size
– Line search

### Problems
– Slow with fixed step
– Oscillations

Essentially Newton, but with guessed and updated Hessian

$$p_{n+1} = -B_n^{-1} \, \nabla f(a_n) \tag{30}$$

Update direction
Approximate Hessian
Previous point

Line search

$$\alpha_{n+1} = \arg\min_{\alpha} \, f(a_n + \alpha p_{n+1}), \qquad s_{n+1} = \alpha_{n+1} p_{n+1} \tag{31}$$

Update step length
Actual step

Update
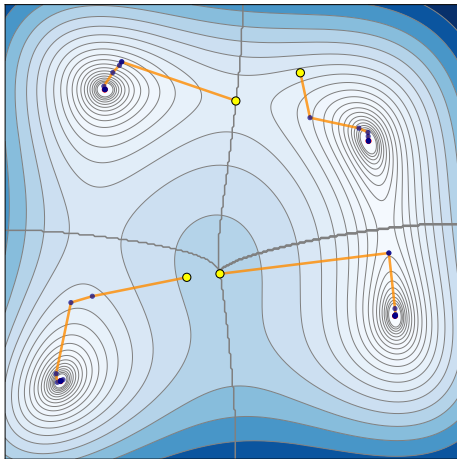
$$a_{n+1} = a_n + s_{n+1}$$

Get gradient response as Hessian approximation

$$\mathsf{y}_{n+1} = \nabla f(\mathsf{a}_{n+1}) - \nabla f(\mathsf{a}_n)$$

Update approximate Hessian

$$\mathsf{B}_{n+1} = \mathsf{B}_n + \frac{\mathsf{y}_{n+1}\mathsf{y}_{n+1}^T}{\mathsf{y}_{n+1}^T s_{n+1}} - \frac{\mathsf{B}_n \mathsf{s}_{n+1}\mathsf{s}_{n+1}^T\mathsf{B}_n^T}{\mathsf{s}_{n+1}^T\mathsf{B}_n\mathsf{s}_{n+1}}$$

Avoid oscillations of steepest descent

$$\mathbf{p}_{n+1} = -\nabla f(\mathbf{a}_n) + \beta_n \mathbf{p}_n \qquad \beta_n = \frac{\nabla f(\mathbf{a}_{n+1}^T \nabla f(\mathbf{a}_{n+1}}{\nabla f(\mathbf{a}_n^T \nabla f(\mathbf{a}_n} \qquad (32)$$

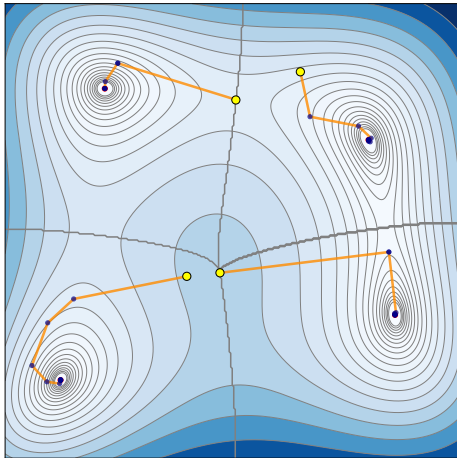$\underline{\text{Search direction}}$       Mixing parameter

New point found by line search (like BFGS).

## Key properties

– Orthogonal search directions
– Quadratic convergence for quadratic functions
– Memory efficient

## Problems

– Sensitive to round-off errors
– Requires periodic restarts

Use subset of data for gradient estimation

$$\mathbf{a}_{n+1} = \mathbf{a}_n - \eta \nabla f_i(\mathbf{a}_n) \tag{33}$$

Updated parameters

Learning rate

Single sample gradient

### Variants

– Mini-batch SGD
– Learning rate scheduling
– Momentum

### Advantages

– Fast iterations
– Escapes local minima

Adaptive moments estimation

Parameter, often 0.9

$$m_n = \beta_1 m_{t-1} + (1 - \beta_1)\nabla f_n \tag{34}$$

Momentum: averaged gradient

Parameter, often 0.999

$$v_n = \beta_2 v_{n-1} + (1 - \beta_2)(\nabla f_n)^2 \tag{35}$$

Scale adjustment

Bias correction, since at start: $m_0 = v_0 = 0$

$$\hat{m}_n = \frac{m_n}{1 - \beta_1^n} \quad \hat{v}_n = \frac{v_n}{1 - \beta_2^n} \tag{36}$$

Learning rate, often 0.001

$$\mathsf{a}_{n+1} = \mathsf{a}_n - \frac{\alpha \hat{m}_n}{\sqrt{\hat{v}_n} + \epsilon} \tag{37}$$

Small, for stability, often $10^{-8}$

### Advantages

– Adaptive learning rates
– Robust to hyperparameters

### Idea

– Evolutionary optimization approach

### Operations

– Selection: Choose fittest individuals
– Crossover: Combine parent solutions
– Mutation: Random modifications
– Replacement: Update population

### Parameters

– Population size
– Crossover probability
– Mutation rate
– Selection pressure

### Advantages

– Global optimization
– No gradient required
– Handles discontinuous functions

## Idea

– Combine local and global search

## Algorithm

– Local minimization
– Random perturbation
– Accept/reject based on energy
– Repeat

Acceptance criterion

$$P = \min\left(1, e^{-\frac{\Delta E}{k_B T}}\right)$$

## Advantages

– Escapes local minima
– Uses efficient local methods
– Temperature controls exploration

## Applications

– Protein folding
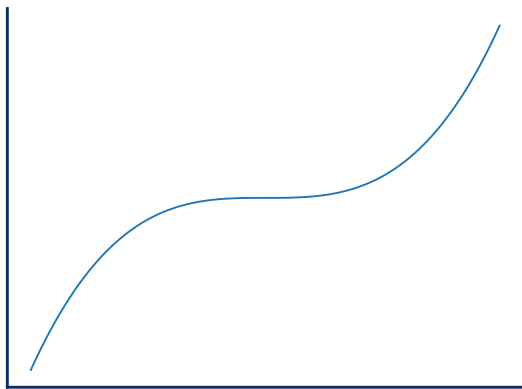– Molecular conformations
– Glass structure

## Convergence

– Hard to establish
– Gradient necessary, but not sufficient
– Hessian expensive
– Local property

## Numerical stability

– Finite differences
– Conjugate Gradients
– Shallow minima

## Cost of Hessians

– Scales as $N^2$
– Sometimes only from finite differences



Gradients can be arbitrarily small

### Curse of dimensionality

– Search space quickly increases
– Often forces tiny optimization steps

### Preconditioning

– Math not equal to finite-precision implementations
– Transform problem into an equivalent one

| Dimensions | Gradients | Hessian | Noise | Minima count | Choice |
|---|---|---|---|---|---|
| Few | Yes | Yes | No | Few | Newton |
| Few/Medium | Yes | No | No | Few | BFGS |
| Few/Medium | Yes | No | No | Many | Basin hopping |
| Any | No | No | Any | Many | Genetic algorithm |
| Large | Yes | No | No | Few | Conjugate gradients |
| Large | Yes | No | Yes | Few | SGD |
| Large | Yes | No | Yes | Many | ADAM |