



Numerical Optimization: The Gradient

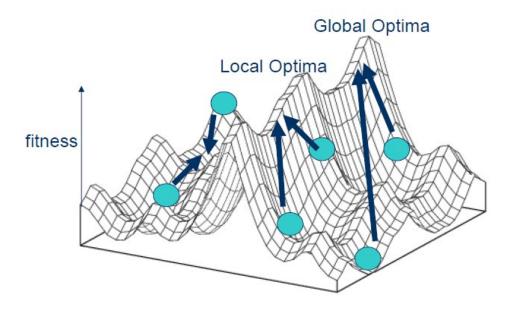
Alexander G. Ororbia II

COGS-621: Foundations of Scientific Computing

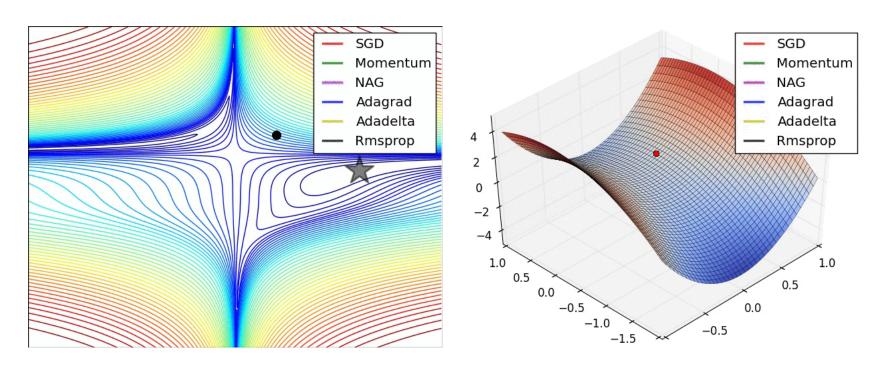
10/16/2025, 10/28/2025

Gradient Descent (or Ascent)

- Simple modification to Hill Climbing
 - Generally assumes a continuous state space
- Idea is to take more intelligent steps
- Look at local gradient: the direction of largest change
- Take step in that direction
 - Step size should be proportional to gradient
- Tends to yield much faster convergence to optima



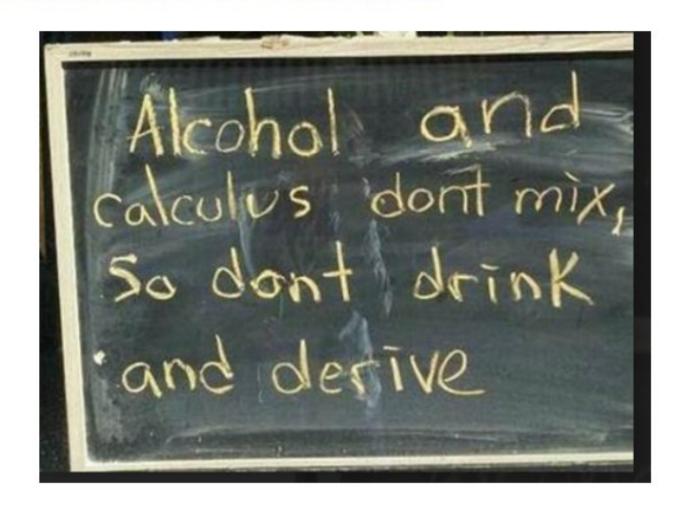
Race of the Optimizers!



http://cs231n.github.io/neural-networks-3/#hyper

Gradient

Essential role of calculus



Background

A Recipe for Optimization

1. Given training data:

$$\{oldsymbol{x}_i, oldsymbol{y}_i\}_{i=1}^N$$

2. Choose each of these:

- Decision function

$$\hat{\boldsymbol{y}} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

Loss function

$$\ell(\hat{oldsymbol{y}}, oldsymbol{y}_i) \in \mathbb{R}$$

3. Define goal:

$$oldsymbol{ heta}^* = rg \min_{oldsymbol{ heta}} \sum_{i=1}^N \ell(f_{oldsymbol{ heta}}(oldsymbol{x}_i), oldsymbol{y}_i)$$

4. Train with SGD:

(take small steps opposite the gradient)

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$$

Background

1. Given training data:

$$\{oldsymbol{x}_i, oldsymbol{y}_i\}_{i=1}^N$$

- 2. Choose each of the
 - Decision function

$$\hat{m{y}} = f_{m{ heta}}(m{x}_i)$$

Loss function

$$\ell(\hat{m{y}},m{y}_i)\in\mathbb{R}$$

A Recipe for

Gradients

Reverse-mode
differentiation
(computational calculus)
can be used to compute
derivatives of most
function efficiently

$$oldsymbol{ heta}^{(t+1)} = oldsymbol{h}_t
abla \ell(f_{oldsymbol{ heta}}(oldsymbol{x}_i), oldsymbol{y}_i)$$

Approaches to Differentiation

Finite Difference Method

- Pro: Great for testing implementations of backpropagation
- Con: Slow for high dimensional inputs / outputs
- Required: Ability to call the function f(x) on any input x

Symbolic Differentiation

- Note: The method you learned in high-school
- Note: Used by Mathematica / Wolfram Alpha / Maple
- Pro: Yields easily interpretable derivatives
- Con: Leads to exponential computation time if not carefully implemented
- Required: Mathematical expression that defines f(x)

3. Automatic Differentiation - Reverse Mode

- Note: Called Backpropagation when applied to Neural Nets
- Pro: Computes partial derivatives of one output $f(x)_i$ with respect to all inputs x_j in time proportional to computation of f(x)
- Con: Slow for high dimensional outputs (e.g. vector-valued functions)
- Required: Algorithm for computing f(x)

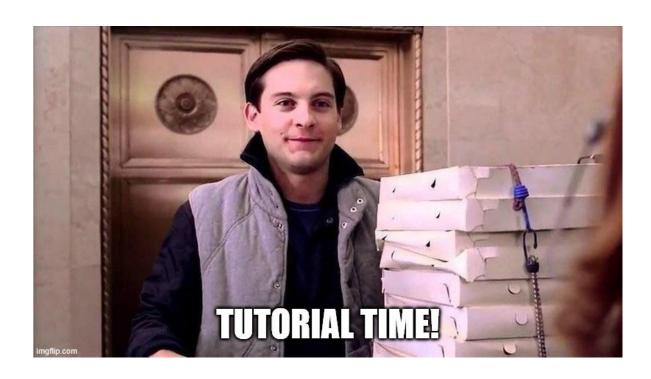
Automatic Differentiation - Forward Mode

- Note: Easy to implement. Uses dual numbers.
- Pro: Computes partial derivatives of all outputs $f(x)_i$ with respect to one input x_j in time proportional to computation of f(x)
- Con: Slow for high dimensional inputs (e.g. vector-valued x)
- Required: Algorithm for computing f(x)

Given
$$f: \mathbb{R}^A \to \mathbb{R}^B, f(\mathbf{x})$$

Compute
$$\frac{\partial f(\mathbf{x})_i}{\partial x_j} \forall i, j$$

So...what is a derivative?!



The Finite Difference Method

The centered finite difference approximation is:

$$rac{df(x)}{dx} = \lim_{h o 0} rac{f(x+h) - f(x)}{h}$$

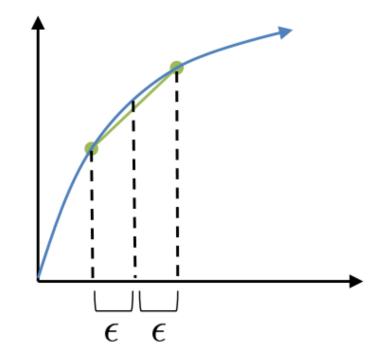
$$\frac{\partial}{\partial \theta_i} J(\boldsymbol{\theta}) \approx \frac{(J(\boldsymbol{\theta} + \epsilon \cdot \boldsymbol{d}_i) - J(\boldsymbol{\theta} - \epsilon \cdot \boldsymbol{d}_i))}{2\epsilon}$$

where d_i is a 1-hot vector consisting of all zeros except for the ith

entry of d_i , which has value 1.

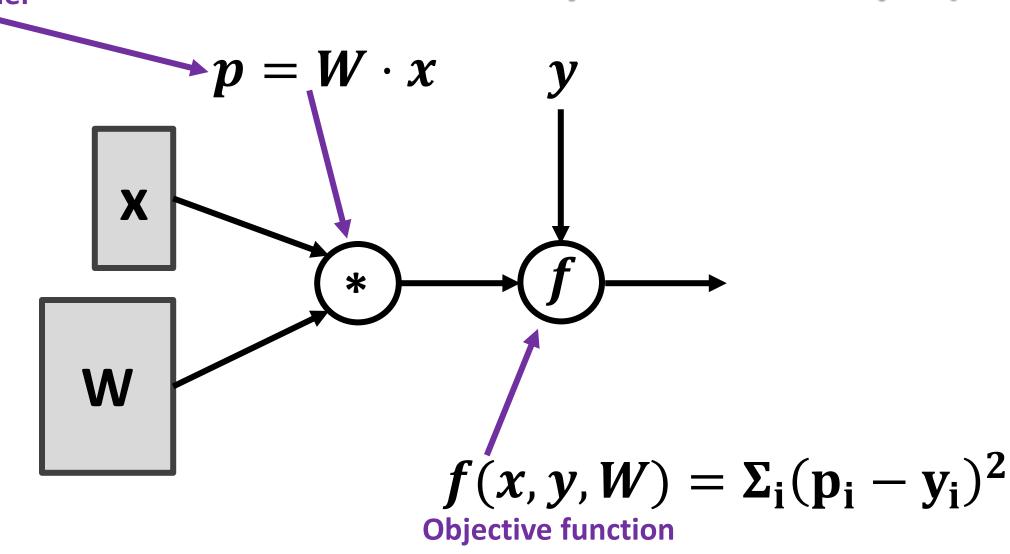
Notes:

- Suffers from issues of floating point precision, in practice
- Typically only appropriate to use on small examples with an appropriately chosen epsilon



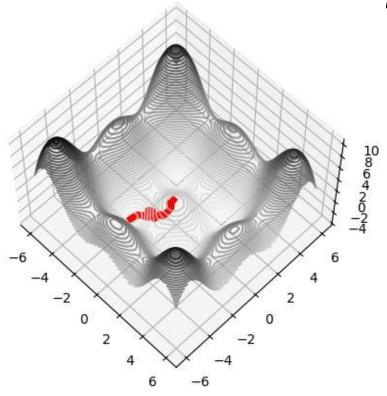
Model

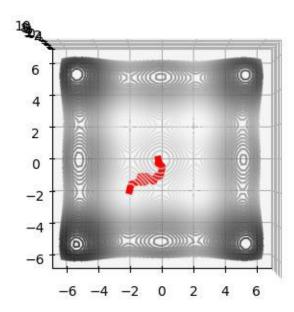
Computational Graph (Example)



Stochastic Hill-Climbing

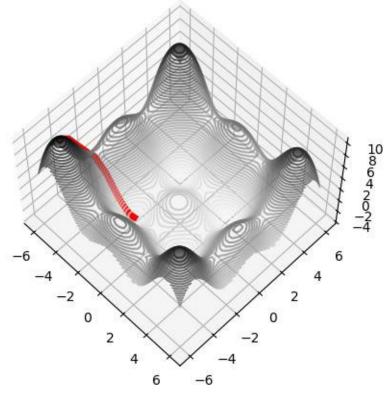
Negative Alpine Function
$$f(x) = -\sum_{i} (x_i \sin(x_i) + 0.1 x_i)$$

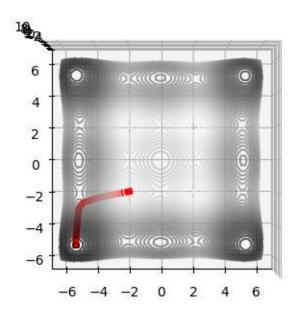




Stochastic Gradient Ascent

Negative Alpine Function
$$f(x) = -\sum_{i} (x_i \sin(x_i) + 0.1 x_i)$$





Questions?

