

Hacking Steepest Descent for Better Convergence



Extrapolation methods: "momentum"

$$x_{k+1} = x_k - \alpha_k \nabla f(x_k) + \beta_k (x_k - x_{k-1})$$

Heavy ball method:

Gradient descent:

$$\alpha_k = \alpha, \beta_k = \beta$$

$$\|e_{k+1}\|_A = \frac{\sqrt{\kappa(A)} - 1}{\sqrt{\kappa(A)} + 1} \|e_k\|_A$$

$$\frac{\kappa(A) - 1}{\kappa(A) + 1}$$

Demo: Steepest Descent [cleared] (Part 2)

Optimization in Machine Learning

What is *stochastic gradient descent* (SGD)?

$$f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x).$$

f_i : "data point" / observation.

$$\text{GD} : x_{k+1} = x_k - \alpha \frac{1}{n} \sum_{i=1}^n \nabla f_i(x).$$

$$\text{SGD} : x_{k+1} = x_k - \alpha \nabla \underbrace{f_{\phi(x_k)}}_{\text{"minibatch"}}(x_k)$$

"ADAM"

Conjugate Gradient Methods

Can we optimize in *the space spanned* by the last two step directions?

$$\begin{aligned} & \vec{x}^T A \vec{y} = 0 \quad \text{"conjugate"} \\ & (\alpha_k, \beta_k) = \operatorname{argmin} \left[f(\vec{x}_k - \alpha_k \nabla f(\vec{x}_k)) + \beta_k (\vec{x}_k - \vec{x}_{k-1}) \right] \end{aligned}$$

Demo: Conjugate Gradient Method [cleared]

Nelder-Mead Method

Idea:



[Demo: Nelder-Mead Method](#) [cleared]

Newton's method (n D)

What does Newton's method look like in n dimensions?

$$\textcircled{1} \quad \nabla f(x) = 0.$$

$$\textcircled{2} \quad f(x+s) \approx f(x) + \nabla f(x)^T s + \frac{1}{2} s^T H_f(x) s$$

$$\hat{f}(s) = 0 \Rightarrow H_f(x) s = -\nabla f(x).$$

step 1. x_0 .

" 2. solve $H_f(x_k) s_k = -\nabla f(x_k)$ for s_k

" 3. $x_{k+1} = x_k + s_k$

$\hat{f}(s)$

Newton's method (n D): Observations

Drawbacks?

1. need for 2nd derivatives.
2. local convergence.
3. $H_f(x)$ can be close to indefinite.

Demo: Newton's method in n dimensions [cleared]

Quasi-Newton Methods

Secant/Broyden-type ideas carry over to optimization. How?

- Come up with approximation to Hessian.
- Secant cond.

BFGS: Secant-type method, similar to Broyden:

$$B_{k+1} = B_k + \frac{\mathbf{y}_k \mathbf{y}_k^T}{\mathbf{y}_k^T \mathbf{s}_k} - \frac{(\mathbf{B}_k \mathbf{s}_k) \mathbf{s}_k^T \mathbf{B}_k}{\mathbf{s}_k^T \mathbf{B}_k \mathbf{s}_k}$$

where

$$\begin{aligned} \mathbf{s}_k &= \mathbf{x}_{k+1} - \mathbf{x}_k \\ \mathbf{y}_k &= \nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k) \end{aligned}$$

$$\mathbf{B}_k \mathbf{s}_k = \mathbf{y}_k.$$

In-Class Activity: Optimization Methods

In-class activity: Optimization Methods

Nonlinear Least Squares: Setup

What if the f to be minimized is actually a 2-norm?

$$f(\mathbf{x}) = \|\mathbf{r}(\mathbf{x})\|_2, \quad \mathbf{r}(\mathbf{x}) = \mathbf{y} - \mathbf{a}(\mathbf{x})$$

linear : $\|A\mathbf{x} - \mathbf{b}\|_2$ $A^T A \mathbf{x} = A^T \mathbf{b}$
nonlinear : $\|\mathbf{r}(\mathbf{x})\|_2$

$$\varphi(\mathbf{x}) := \frac{1}{2} \mathbf{r}(\mathbf{x})^T \mathbf{r}(\mathbf{x}) = \frac{1}{2} f^2(\mathbf{x})$$

minimize $\varphi(\mathbf{x})$ instead.

$$\frac{\partial \varphi}{\partial x_i} = \frac{1}{2} \sum_{j=1}^n \frac{\partial}{\partial x_i} [\mathbf{r}_j(\mathbf{x})^2] = \sum_{j=1}^n \frac{\partial \mathbf{r}_j}{\partial x_i} \cdot \mathbf{r}_j$$

$$\nabla \varphi = \mathbf{J}_r(\mathbf{x})^T \mathbf{r}(\mathbf{x})$$

Gauss-Newton

For brevity: $J := J_r(\mathbf{x})$.

$$H\varphi(\mathbf{x}) = J^T J + \sum_i r_i H r_i(\mathbf{x})$$

Newton: $H\varphi(\mathbf{x}) \mathbf{s} = -\nabla\varphi$.

Gauss-Newton: $\underline{J^T J} \mathbf{s} = -\nabla\varphi = \underline{-J^T r(\mathbf{x})}$.

$$J\mathbf{s} \cong -r(\mathbf{x})$$

linear least square!