Optimization in Machine Learning

Second order methods Quasi-Newton





Learning goals

- Newton-Raphson vs. Quasi-Newton
- SR1
- BFGS

QUASI-NEWTON: IDEA

Start point of **QN method** is (as with NR) a Taylor approximation of the gradient, except that H is replaced by a **pd** matrix $A^{[t]}$:

$$\begin{aligned} \nabla f(\mathbf{x}) &\approx \nabla f(\mathbf{x}^{[l]}) + \nabla^2 f(\mathbf{x}^{[l]})(\mathbf{x} - \mathbf{x}^{[l]}) = \mathbf{0} \qquad \mathsf{NR} \\ \nabla f(\mathbf{x}) &\approx \nabla f(\mathbf{x}^{[l]}) + \mathbf{A}^{[l]} \qquad (\mathbf{x} - \mathbf{x}^{[l]}) = \mathbf{0} \qquad \mathsf{QN} \end{aligned}$$

The update direction:

$$\begin{aligned} \boldsymbol{d}^{[t]} &= -\nabla^2 f(\boldsymbol{x}^{[t]})^{-1} \nabla f(\boldsymbol{x}^{[t]}) & \text{NR} \\ \boldsymbol{d}^{[t]} &= -(\boldsymbol{A}^{[t]})^{-1} \quad \nabla f(\boldsymbol{x}^{[t]}) & \text{QN} \end{aligned}$$

× × 0 × × ×

QUASI-NEWTON: IDEA / 2

- Select a starting point x^[0] and initialize pd matrix A^[0] (can also be a diagonal matrix a very rough approximation of Hessian).
- Calculate update direction by solving

$$oldsymbol{A}^{[t]}oldsymbol{d}^{[t]} = -
abla f(\mathbf{x}^{[t]})$$

and set $\mathbf{x}^{[t+1]} = \mathbf{x}^{[t]} + \alpha^{[t]} \mathbf{a}^{[t]}$ (Step size through backtracking)

Solution Calculate an efficient update $\mathbf{A}^{[t+1]}$, based on $\mathbf{x}^{[t]}$, $\mathbf{x}^{[t+1]}$, $\nabla f(\mathbf{x}^{[t]})$, $\nabla f(\mathbf{x}^{[t+1]})$ and $\mathbf{A}^{[t]}$. × 0 0 × × ×

QUASI-NEWTON: IDEA / 3

Usually the matrices $\mathbf{A}^{[t]}$ are calculated recursively by performing an additive update

$$A^{[t+1]} = A^{[t]} + B^{[t]}.$$

How $\boldsymbol{B}^{[t]}$ is constructed is shown on the next slides. **Requirements** for the matrix sequence $\boldsymbol{A}^{[t]}$:

- Symmetric pd, so that $d^{[t]}$ are descent directions.
- 2 Low computational effort when solving LES

$$oldsymbol{A}^{[t]}oldsymbol{d}^{[t]} = -
abla f(oldsymbol{x}^{[t]})$$

• Good approximation of Hessian: The "modified" Taylor series for $\nabla f(\mathbf{x})$ (especially for $t \to \infty$) should provide a good approximation

$$abla f(\mathbf{x}) \approx
abla f(\mathbf{x}^{[t]}) + \mathbf{A}^{[t]}(\mathbf{x} - \mathbf{x}^{[t]})$$

SYMMETRIC RANK 1 UPDATE (SR1)

Simplest approach: symmetric rank 1 updates (SR1) of form

$$\boldsymbol{A}^{[t+1]} \leftarrow \boldsymbol{A}^{[t]} + \boldsymbol{B}^{[t]} = \boldsymbol{A}^{[t]} + \beta \boldsymbol{u}^{[t]} (\boldsymbol{u}^{[t]})^{\top}$$

with appropriate vector $\boldsymbol{u}^{[t]} \in \mathbb{R}^n$, $\beta \in \mathbb{R}$.

× < 0 × × ×

SYMMETRIC RANK 1 UPDATE (SR1) / 2

Choice of $u^{[t]}$:

Vectors should be chosen so that the "modified" Taylor series corresponds to the gradient:

$$\nabla f(\mathbf{x}) \stackrel{!}{=} \nabla f(\mathbf{x}^{[t+1]}) + \mathbf{A}^{[t+1]}(\mathbf{x} - \mathbf{x}^{[t+1]})$$

$$\nabla f(\mathbf{x}) = \nabla f(\mathbf{x}^{[t+1]}) + \left(\mathbf{A}^{[t]} + \beta \mathbf{u}^{[t]}(\mathbf{u}^{[t]})^{\top}\right) \underbrace{(\mathbf{x} - \mathbf{x}^{[t+1]})}_{:=\mathbf{s}^{[t+1]}}$$

$$\underbrace{\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}^{[t+1]})}_{\mathbf{y}^{[t+1]}} = \left(\mathbf{A}^{[t]} + \beta \mathbf{u}^{[t]}(\mathbf{u}^{[t]})^{\top}\right) \mathbf{s}^{[t+1]}$$

$$\frac{\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}^{[t+1]})}{\mathbf{y}^{[t+1]}} = \left(\beta (\mathbf{u}^{[t]})^{\top} \mathbf{s}^{[t+1]}\right) \mathbf{u}^{[t]}$$
For $\mathbf{u}^{[t]} = \mathbf{y}^{[t+1]} - \mathbf{A}^{[t]} \mathbf{s}^{[t+1]}$ and $\beta = \frac{1}{(\mathbf{y}^{[t+1]} - \mathbf{A}^{[t]} \mathbf{s}^{[t+1]})}$ the equation is satisfied.

Х

хx

SYMMETRIC RANK 1 UPDATE (SR1) / 3

Advantage

- Provides a sequence of symmetric pd matrices
- Matrices can be inverted efficiently and stable using Sherman-Morrison:

$$(\boldsymbol{A} + \beta \boldsymbol{u} \boldsymbol{u}^{\top})^{-1} = \boldsymbol{A} + \beta \frac{\boldsymbol{u} \boldsymbol{u}^{\top}}{1 + \beta \boldsymbol{u}^{\top} \boldsymbol{u}}$$

× 0 0 × ×

Disadvantage

• The constructed matrices are not necessarily pd, and the update directions **d**^[t] are therefore not necessarily descent directions

BFGS ALGORITHM

Instead of Rank 1 updates, the **BFGS** procedure (published simultaneously in 1970 by Broyden, Fletcher, Goldfarb and Shanno) uses rank 2 modifications of the form

$$A^{[t]} + \beta u^{[t]} (u^{[t]})^{\top} + \beta v^{[t]} (v^{[t]})^{\top}$$

with $s^{[t]} := x^{[t+1]} - x^{[t]}$

•
$$\boldsymbol{u}^{[t]} = \nabla f(\boldsymbol{x}^{[t+1]}) - \nabla f(\boldsymbol{x}^{[t]})$$

•
$$v^{[t]} = A^{[t]} s^{[t]}$$

•
$$\beta = \frac{1}{(\boldsymbol{u}^{[t]})^{\top}(\boldsymbol{s}^{[t]})}$$

•
$$\beta = -\frac{\mathbf{I}}{(\mathbf{s}^{[t]})^{\top} \mathbf{A}^{[t]} \mathbf{s}^{[t]}}$$

The resulting matrices $A^{[t]}$ are positive definite and the corresponding quasi-newton update directions $d^{[t]}$ are actual descent directions.

× × 0 × × ×