# Chapter 11 Quasi-Newton Methods

- Motivation It takes $\mathcal{O}\left(d^3\right)$ to calculate $\nabla^2 f(x)^{-1}$ or solve for $\nabla^2 f\left(\mathbf{x}_t\right)\Delta\mathbf{x} = -\nabla f\left(\mathbf{x}_t\right)$.

## The secant method

- Motivation $\frac{f(x_t)-f(x_{t-1})}{x_t-x_{t-1}} \approx f'\left(x_t\right)$, so $x_{t+1} := x_t - \frac{f(x_t)}{f'(x_t)} \approx x_{t+1} := x_t - f'\left(x_t\right)\frac{x_t-x_{t-1}}{f'(x_t)-f'(x_{t-1})}$
- In optimization regime, we want to find a similar matrix s.t. $\nabla f\left(\mathbf{x}_t\right) - \nabla f\left(\mathbf{x}_{t-1}\right) = H_t\left(\mathbf{x}_t - \mathbf{x}_{t-1}\right)$ and so $\mathbf{x}_{t+1} = \mathbf{x}_t - H_t^{-1}\nabla f\left(\mathbf{x}_t\right)$.
  - This is called *secant condition*

## Quasi-Newton methods

- Definition If $H_t$ symmetric, and follows secant condition, the update method is *quasi-newton*.
- Lemma Exercise 71 $f \in C^2$ and $\nabla^2 f \neq 0$, then Newton's method is a Quasi-Newton method iff $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top M\mathbf{x} - \mathbf{q}^\top\mathbf{x} + c$ with $M$ invertable symmetric.
  - Proof
    - Newtwon is quasi $\Leftrightarrow \nabla f(y) - \nabla f(x) = \nabla^2 f(y)(y - x), \forall x, y$, take derivative w.r.t $x$, we get $\nabla^2 f(x) = \nabla^2 f(y), \forall x, y$ and this means $f$ is a quadratic funciton, its invertable since every secant condition there is a solution. The other direciton is straightforward.

## Greenstadt's Approach

- We already have $H_{t-1}^{-1}, x_{t-1}, x_t$, we need $H_t^{-1}$, idea is $H_t^{-1} = H_{t-1}^{-1} + E_t$, and we want to minimized the general change $\|AEA^\top\|_F^2$.
- Denote $H := H_{t-1}^{-1}, H' := H_t^{-1}, E := E_t, \boldsymbol{\sigma} := \mathbf{x}_t - \mathbf{x}_{t-1}, \mathbf{y} = \nabla f\left(\mathbf{x}_t\right) - \nabla f\left(\mathbf{x}_{t-1}\right)$, and $\mathbf{r} = \boldsymbol{\sigma} - H\mathbf{y}$,
  - then the update formula is $H' = H + E$, such that $H'\mathbf{y} = \boldsymbol{\sigma}$ or equivalently $E\mathbf{y} = \mathbf{r}$,
  - so the overall minimization is $\underline{\textbf{minimize}} \quad \frac{1}{2}\|AEA^\top\|_F^2, \textbf{ subject to } E\mathbf{y} = \mathbf{r} \text{ and } E^\top - E = 0$.

### Solving with Lagrange multiplier

- Fact 11.2 If $f(E) := \frac{1}{2}\|AEB\|_F^2$, then $\nabla f(E) = A^\top AEBB^\top$ if define $\nabla f(E) = \left(\frac{\partial f(E)}{\partial E_{ij}}\right)$.
  - Proof
    - By this def, we have $\nabla_E\text{Tr}(AE)\nabla_E = \text{Tr}(E^\top A^\top) = A^\top$, so $f(E) := \frac{1}{2}\|AEB\|_F^2 = \frac{1}{2}\text{Tr}(B^\top E^\top A^\top AEB)$, then $\nabla f(E) = \nabla_E\frac{1}{2}\text{Tr}(E^\top A^\top AE_0BB^\top) + \nabla_E\frac{1}{2}\text{Tr}(BB^\top E_0^\top A^\top AE) = A^\top AEBB^\top/2 + (BB^\top E^\top A^\top A)^\top/2 = A^\top AEBB^\top$
- Denote $\boldsymbol{\lambda} \in \mathbb{R}^d$ as the multiplier for $d$ constrains of $E\mathbf{y} = \mathbf{r}$, and $\Gamma \in \mathbb{R}^{d\times d}$ as the multiplier for $d \times d$ constraints of $E^\top - E = 0$.
- For each equation of $\partial_{E_{ij}}f = \boldsymbol{\lambda}^\top f_1 + \text{Tr}(\Gamma f_2)$, $\lambda$ part yields a term of $\lambda_i y_i$ and $\Gamma$ yields a term of $\Gamma_{ji} - \Gamma_{ij}$
- Lemma 11.3 The above equation gives the optimial condtional of $WE^\star W = \boldsymbol{\lambda}\mathbf{y}^\top + \Gamma^\top - \Gamma$, where $W := A^\top A$ is symmetric and posi definite.

### Solving Greenstadt family

- The minimization has now turn into three linear equation (i) $E\mathbf{y} = \mathbf{r}$, (ii) $E^\top - E = 0$ and (iii) $WEW = \boldsymbol{\lambda}\mathbf{y}^\top + \Gamma^\top - \Gamma$
- To eliminate $\Gamma$, by plug (iii) into (ii) we get $M\left(\boldsymbol{\lambda}\mathbf{y}^\top - \mathbf{y}\boldsymbol{\lambda}^\top + 2\Gamma^\top - 2\Gamma\right)M = 0$, where $M = W^{-1}$, so $\Gamma^\top - \Gamma = \frac{1}{2}\left(\mathbf{y}\boldsymbol{\lambda}^\top - \boldsymbol{\lambda}\mathbf{y}^\top\right)$
  - then $E = \frac{1}{2}M\left(\boldsymbol{\lambda}\mathbf{y}^\top + \mathbf{y}\boldsymbol{\lambda}^\top\right)M$
- Then to eliminate $\boldsymbol{\lambda}$, we plug in the secant condition (i) we get $\boldsymbol{\lambda} = \frac{1}{\mathbf{y}^\top M\mathbf{y}}\left(2M^{-1}\mathbf{r} - \mathbf{y}\boldsymbol{\lambda}^\top M\mathbf{y}\right)$
  - multiply with $\mathbf{y}^\top M$, we get $z = \boldsymbol{\lambda}^\top M\mathbf{y} = \frac{\mathbf{y}^\top\mathbf{r}}{\mathbf{y}^\top M\mathbf{y}}$, so $\boldsymbol{\lambda} = \frac{1}{\mathbf{y}^\top M\mathbf{y}}\left(2M^{-1}\mathbf{r} - \frac{(\mathbf{y}^\top\mathbf{r})}{\mathbf{y}^\top M\mathbf{y}}\mathbf{y}\right)$
- Plug this into $E$, we get $E = \frac{1}{2}M\left(\boldsymbol{\lambda}\mathbf{y}^\top + \mathbf{y}\boldsymbol{\lambda}^\top\right)M = \frac{1}{\mathbf{y}^\top M\mathbf{y}}\left(\mathbf{r}\mathbf{y}^\top M + M\mathbf{y}\mathbf{r}^\top - \frac{(\mathbf{y}^\top\mathbf{r})}{\mathbf{y}^\top M\mathbf{y}}M\mathbf{y}\mathbf{y}^\top M\right)$, by definition $\mathbf{r} = \boldsymbol{\sigma} - H\mathbf{y}$, we get
  - $E^\star = \frac{1}{\mathbf{y}^\top M\mathbf{y}}\left(\boldsymbol{\sigma}\mathbf{y}^\top M + M\mathbf{y}\boldsymbol{\sigma}^\top - H\mathbf{y}\mathbf{y}^\top M - M\mathbf{y}\mathbf{y}^\top H - \frac{1}{\mathbf{y}^\top M\mathbf{y}}(\mathbf{y}^\top\boldsymbol{\sigma} - \mathbf{y}^\top H\mathbf{y})M\mathbf{y}\mathbf{y}^\top M\right)$

# BFGS (Broyden, Fletcher, Goldfarb and Shanno)

- **Definition** BFGS is when $M = H' = H_t^{-1}$, $My = H'y = \sigma$, even we don't know $H'$, but it never appears in the solution, so $E^\star = \frac{1}{y^\top \sigma}\left(-Hy\sigma^\top - \sigma y^\top H + \left(1 + \frac{y^\top Hy}{y^\top \sigma}\right)\sigma\sigma^\top\right)$, where $H = H_{t-1}^{-1}$, $\sigma = x_t - x_{t-1}$, $y = \nabla f(x_t) - \nabla f(x_{t-1})$.
    - Iteration cost is $O(d^2)$
- **Lemma Exercise 74.1** If $f$ convex, $y^\top \sigma > 0$, unless $x_t = x_{t-1}$ or $f(\lambda x_t + (1-\lambda)x_{t-1}) = \lambda f(x_t) + (1-\lambda)f(x_{t-1})$ for all $\lambda \in (0,1)$.
    - **Proof** By property of convexity $y^\top \sigma = (\nabla f(x_t) - \nabla f(x_{t-1}))^\top (x_t - x_{t-1}) \geq 0$
- **Observation 11.6** $H' = \left(I - \frac{\sigma y^\top}{y^\top \sigma}\right)H\left(I - \frac{y\sigma^\top}{y^\top \sigma}\right) + \frac{\sigma\sigma^\top}{y^\top \sigma}$
    - **Proof**
        - $E^\star + H = H + \frac{1}{y^\top \sigma}\left(-Hy\sigma^\top - \sigma y^\top H + \left(1 + \frac{y^\top Hy}{y^\top \sigma}\right)\sigma\sigma^\top\right) = \frac{\sigma\sigma^\top}{y^\top \sigma} + H(I - \frac{y\sigma^\top}{y^\top \sigma}) - \sigma y^\top H + \frac{\sigma y^\top Hy\sigma^\top}{y^\top \sigma} =$ QED
- **Lemma Exercise 74.2** If $H \succeq 0$ and $y^\top \sigma > 0$, then also $H'$ is positive definite.
    - **Proof**
        - Since $C := \left(I - \frac{y\sigma^\top}{y^\top \sigma}\right) = \left(I - \frac{\sigma y^\top}{y^\top \sigma}\right)^\top$, so $H' = C^\top HC + \frac{\sigma\sigma^\top}{y^\top \sigma}$, $C^\top HC$ is semi positive definite and so is $\frac{\sigma\sigma^\top}{y^\top \sigma}$.
        - When $z \perp \sigma^\top$, we have $Cz = z \neq 0$, so the two quadratic form will not be zero at the same time, this means positive definiteness.
- **Remark** Usually Newton or Quasi-Newton are performed with *scaled steps* $x_{t+1} = x_t - \alpha_t H_t^{-1}\nabla f(x_t)$, either line search or backtracking line search (when $\alpha_t = 1$ is not good enough, do $\alpha_t/2$).

# L-BFGS (limited memory version)

- **Idea** Only use information from the previous $m$ iterations, for some small value of $m$.
- **Lemma 11.7** If an oracle can compute $s = Hg$ for any vertor $g$, then $s' = H'g'$ can be computed with one oracle call of $s = Hg$, and $O(d)$ arithmetic operation, assuming $\sigma, y$ known.
    - **Proof**
        - $H'g' = \left(I - \frac{\sigma y^\top}{y^\top \sigma}\right)H\underbrace{\left(I - \frac{y^\top}{y^\top \sigma}\right)g'}_{g} + \underbrace{\frac{\sigma\sigma^\top}{y^\top \sigma}g'}_{h}$

            $\underbrace{\hphantom{\left(I - \frac{\sigma y^\top}{y^\top \sigma}\right)H\left(I - \frac{y^\top}{y^\top \sigma}\right)g'}}_{s}$

            $\underbrace{\hphantom{\left(I - \frac{\sigma y^\top}{y^\top \sigma}\right)H\left(I - \frac{y^\top}{y^\top \sigma}\right)g' + \frac{\sigma\sigma^\top}{y^\top \sigma}g'}}_{w}$

            $\underbrace{\hphantom{xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx}}_{z}$
        - $g, h, s, w, z$ all are compuated in $O(d)$.
- The idea is that we need $H_t^{-1}\nabla f_t$, and we can borrow from $H_{t-1}^{-1}\nabla f_t$, etc, and recurse back to $t = 0$, This gives the BFGS-step:
- **Algorithm (BFGS-STEP)**
    - **Input** $(k, g)$
    - If $k = 0$ then return $H_0^{-1}g'$
    - **Else**
        - Set $h = \sigma\frac{\sigma_k^\top g'}{y_k^\top \sigma_k}$, and $g = g' - y\frac{\sigma_k^\top g'}{y_k^\top \sigma_k}$
        - $s = $ BFGS-STEP $(k-1, g)$ (*recursive call*)
        - $w = s - \sigma_k\frac{y_k^\top s}{y_k^\top \sigma_k}$
        - $z = w + h$
        - return $z$
- **Remark** If $H_0$ can be computed in $O(d)$ the total runtime is $O(td)$, this is acceptable when $t \leq d$. It's natual to think of a cut-off version
- **Algorithm (L-BFGS-STEP)**
    - **Input** $(k, l, g)$
    - If $l = 0$ then return $H_0^{-1}g'$
    - **Else**
        - Set $h = \sigma\frac{\sigma_k^\top g'}{y_k^\top \sigma_k}$, and $g = g' - y\frac{\sigma_k^\top g'}{y_k^\top \sigma_k}$
        - $s = $ L-BFGS-STEP $(k-1, l-1, g)$ (*recursive call*)
        - $w = s - \sigma_k\frac{y_k^\top s}{y_k^\top \sigma_k}$
        - $z = w + h$

- **return** $z$