# Lecture 21

Scribe?

## Last time

▷ Rank 2 updates

▷ BFGS

▷ DFP

## Today

▷ Convergence guarantees for BFGS.

▷ Proof

## Quasi-Newton methods convergence guarantees

Analyzing the iterates directly is hard instead we show that the trajectories of the iterates are similar to those of other algorithms.

The guarantees we are about to see are weak, they only apply to strongly convex functions.

In practice, Quasi-Newton method work well for most functions.

We will see two results:

**Theorem:** (Linear convergence) Let $B_0$ be a positive definite matrix and let $f: \mathbb{R}^d \to \mathbb{R}$ be a $C^2$ function, such that

$$\mu I \preceq \nabla f^2(x) \preceq L I.$$

$\mu$-strongly convex $\qquad\qquad$ $L$-smooth

Then, the iterates of BFGS converge linearly $x_k \to x^*$. $\qquad\qquad\qquad\dashv$

**Theorem:** (Local Superlinear convergence) Let $B_0$ be a PD matrix and $f: \mathbb{R}^d \to \mathbb{R}$ $C^2$, with a minimizer $x^*$ with $\nabla^2 f(x^*) > 0$ and for all $x, y$ near $x^*$ we have

$$\| \nabla^2 f(x) - \nabla^2 f(y) \| \leq Ca \| x - y \|,$$

Then, if $x_0$ starts close enough to $x^*$ we have that the iterates of BFGS converges super linearly $x_k \to x^*$.

$\qquad\qquad\qquad\dashv$

Today we focus on proving a weaker version of the first Theorem (we shall only prove convergence).

The argument is based on the following result:

Theorem (     ) Let $f: \mathbb{R}^d \to \mathbb{R}$ be L-smooth. Consider an update

$$x_{k+1} \leftarrow x_k + \alpha_k p_k$$

where $p_k$ is descent direction an $\alpha_k$ satisfies the Wolfe conditions. Then,

$$\sum_{k=0}^{\infty} \cos^2\theta_k \|\nabla f(x_k)\|^2 < \infty.$$

angle between $p_k$ and $-\nabla f(x_k)$

Proof: By the second Wolfe cond.

$$(c-1)\, \nabla f(x_k)^T p_k \leq (\nabla f(x_{k+1}) - \nabla f(x_k))^T p_k$$

$$\underbrace{x_{k+1} - x_k = \alpha_k p_k}_{\text{L-Lips}} \leq \alpha_k L \|p_k\|^2$$

Then,

$$\alpha_k \geq \underbrace{(c-1)}_{L} \underbrace{\frac{\nabla f(x_k)^T p_k}{\|p_k\|^2}}$$

Using the first Wolfe condition

$$f(x_{k+1}) \leq f(x_k) - \eta \frac{(1-c)}{L} \underbrace{\frac{(\nabla f(x_k)^T p_k)^2}{\|p_k\|^2}}$$

$$\underbrace{\cos\theta = \frac{\langle u, v \rangle}{\|u\| \|v\|}}_{} = f(x_k) - \eta \frac{(1-c)}{L} \cos^2\theta_k \|\nabla f(x_k)\|^2$$

Recursing $\rightarrow$

$$\leq f(x_0) - \eta \frac{(1-c)}{L} \sum_{j=0}^{K} \cos^2\theta_j \|\nabla f(x_j)\|^2$$

Reordering

$$\sum^{K} \cos^2\theta_j \|\nabla f(x_j)\|^2 \leq \frac{L}{\eta(1-c)} (f(x_0) - \min f)$$

Letting $K \uparrow \infty$, yields the result.

□

**Idea:** If we show that
$$\cos^2 \theta_k \geq \delta > 0$$
$$\Rightarrow \liminf \|\nabla f(x_k)\|^2 \to 0.$$

This enough to have convergence for strongly convex functions since

*Descent* $\to$ $f(x_{k+1}) - \min f \leq f(x_k) - \min f$

$\frac{\mu}{2}\left(f(x_{k+1}) - \min f\right) \leq \|\nabla f(x_{k+1})\|^2;$

*Error bound from midterm.*

and

$\frac{\mu}{2}\|x_k - x^*\|^2 \leq f(x_k) - \min f.$

*quadratic growth.*

**Proof:** We focus on showing
$$\cos \theta_k^2 \geq \delta > 0.$$

where
$$\theta_k = \text{angle}\left(B_k^{-1} \nabla f(x_k), -\nabla f(x_k)\right)$$

Note that

$$s_{k+1} = -\alpha_k B_k^{-1} \nabla f(x_k)$$

Then

$$\text{angle}(s_{k+1}, B_k s_{k+1}) = \theta_k.$$

We will prove a bound using the relative entropy. Define

$$\psi(B) = \underset{\leftarrow \Sigma \lambda_i}{tr(B)} - \underset{\leftarrow \log(\Pi \lambda_i)}{\log(\det(B))}$$

One can show $\psi(B) > 0$ for $B > 0$.

Let's show that if $\cos\theta_k^2 \to 0$

$\Rightarrow \psi(B_k) < 0$ for large $k$.

Facts: $\leftarrow$ Check!

$$tr(B_{k+1}) = tr(B_k) - \underbrace{\frac{\|B_k s_{k+1}\|^2}{s_{k+1}^T B_k s_{k+1}}}_{q_k/\cos^2\theta_k} + \underbrace{\frac{\|y_{k+1}\|^2}{y_{k+1}^T s_{k+1}}}_{L_k}$$

$$\det(B_{k+1}) = \det(B_k) \frac{y_{k+1}^T s_{k+1}}{s_{k+1}^T B_k s_{k+1}}$$

We define

$$\mu_k = \frac{y_{k+1}^T S_{k+1}}{\|S_{k+1}\|^2}, \quad L_k = \frac{y_{k+1}^T y_{k+1}}{y_{k+1}^T S_{k+1}}$$

Then, $\mu \leq \mu_k$ and $L_k \leq L$. (★)

Further define

$$q_k = \frac{S_{k+1}^T B_k S_{k+1}}{\|S_{k+1}\|^2}.$$

Then,

$$\det(B_{k+1}) = \det(B_k) \frac{\mu_k}{q_k}$$

and

$$\frac{\|B_k S_{k+1}\|^2}{S_{k+1}^T B_k S_{k+1}} = \frac{\|B_k S_{k+1}\|^2 \|S_{k+1}\|^2}{(S_{k+1}^T B_k S_{k+1})^2} \frac{(S_{k+1}^T B_k S_{k+1})}{\|S_{k+1}\|^2}$$

$$= \frac{q_k}{\cos^2 \theta_k}.$$

[Margin annotations:]

$$\frac{S_{k+1}^T G_k S_{k+1}}{\|S_{k+1}\|^2}$$

$$\frac{z_k^T G_k z_k}{\|z_k\|^2} \qquad z_k = G_k^{1/2} S_{k+1}$$

$$\frac{S_{k+1} G_k^2 S_{k+1}}{S_{k+1} G_k S_{k+1}}$$

Thus,

$$\Psi(B_{k+1}) = tr(B_k) + L_k - \frac{q_k}{\cos^2\theta_k}$$
$$- \ln(\det B_k) - \ln q_k + \ln \mu_k$$

$$= \Psi(B_k) + \left(L_k - \ln \mu_k - 1\right)$$

$$+ \underbrace{\left[1 - \frac{q_k}{\cos^2\theta_k} + \ln \frac{q_k}{\cos^2\theta_k}\right]}_{\color{green}{1 - t + \ln(t) \leq 0 \quad \forall t > 0}} + \ln \cos^2\theta$$

$$\leq \Psi(B_k) + L - \ln \mu - 1$$
$$+ \ln \cos^2\theta_k$$

$$\leq \Psi(B_0) + C(K+1) + \sum_{j=0}^{K} \ln \cos^2\theta_j$$

Assume seeking contradiction
that $\cos^2\theta_j \to 0 \Rightarrow \ln \cos^2\theta_j \to -\infty$.

Let $k_0 > 0$ s.t $\forall_{j > k_0}$ $\ln \cos^2 \theta_j < -2c$.

Thus,

$$0 \leq \psi(B_k) \leq \psi(B_0) + c(k+1)$$
$$+ \sum_{j=0}^{k_0} \ln \cos^2 \theta_j - \sum_{j=k_0+1}^{k} 2c$$

$$= \psi(B_0) + \sum^{k_0} \ln \cos^2 \theta_j + 2c k_1 + c - ck.$$

$$< 0$$

$\uparrow$ For large $k$.   $\Downarrow$   $\square$

Let's prove $(\cancel{A})$, note that

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

$$= \int_0^1 \nabla^2 f(x_k + t(x_{k-1} - x_k)) (x_k - x_{k-1}) \, dt$$

$$= \underbrace{\left[ \int_0^1 \nabla^2 f(x_k + t(x_{k-1} - x_k)) \, dt \right]}_{G_k} s_{k+1}$$

Since $G_k$ is an integral of Hessians

$\Rightarrow \quad \mu I \preceq G_k \preceq L I$

which implies $(\cancel{*})$.