# Lecture 22

| Last time | Today |
|---|---|
| ▷ Convergence guaran̲tees for BFGS. <br> ▷ Proof | ▷ L-BFGS <br> ▷ Conjugate gradient method |

## L-BFGS

With BFGS, we solved multiple problems

+ We have descent.
+ Local superlinear convergence
+ Only have to compute $\nabla f(x_k)$ per iter.
- However, we have a storage cost of $O(d^2)$. Thus it only works up to $d = 10^4 \sim 10^5$ $\left(\begin{array}{c}\text{on personal} \\ \text{computer}\end{array}\right)$.

To tackle higher sizes we can forget far away iterates.

In HW 5 you'll show that BFGS
updates

$$B_K^{-1} = B_0^{-1} + \alpha_1 w w_1^T + \dots + \alpha_{2k} w_{2k} w_{2k}^T$$

Instead of keeping all $2k$ vectors
we can just keep that last $m$
where $m \sim 2-30$ (usually).

<span style="color:green">This leads to $O(dm)$ memory
when $B_0^{-1} = I$.</span>

<span style="color:green">↑ or any other simple to
apply linear map.</span>

Because

$$B_K^{-1} \nabla f(x_K) = \nabla f(x_k) + \sum_{j=k-m}^{K} \alpha_j \left( w_j^T \nabla f(x_k) \right) w_j.$$

<span style="color:blue">Conjugate gradient</span>

<span style="color:green">Gauss looking
for planets</span>

Today we go back to least squares

$$\min_x \frac{1}{2} x^T A x - b^T x$$

<span style="color:green">$f(x) :=$</span>

where $A \succ 0$ and $b \in \mathbb{R}^d$.

Optimality conditions say this is

equivalent to solving $Ax = b$.

What would happened if I knew $A = V \Lambda V^T$?

spectral decomposition

I could simplify the problem

change of basis

$$\min_{x \in \mathbb{R}^d} \frac{1}{2} x^T A x - b x = \min_{y \in \mathbb{R}^d} \frac{1}{2} (Vy)^T A (Vy) - b^T V y$$

$$= \min_{y \in \mathbb{R}^d} \frac{1}{2} y^T V^T \underbrace{V}_{I} \Lambda \underbrace{V^T V}_{I} y - b^T V y$$

$$= \min_{y} \frac{1}{2} y^T \Lambda y - b^T V y$$

This is a separable problem.

$$\longrightarrow = \min \sum y_i^2 \lambda_i - b^T v_i y_i$$

$\Rightarrow$ Minimized at $y_i^* = \arg\min \; y_i^2 \lambda_i - b^T v_i y_i$

$$= \frac{v_i^T b}{\lambda_i}$$

Minimized at $x^* = \sum_{i}^{d} v_i y_i^* = \sum_{i}^{d} \frac{v_i v_i^T b}{\lambda_i}$.

Computing the spectral decomposition is too expensive.

# Conjugate vectors

Any $A > 0$ defines an inner product

$$\langle x, y \rangle_A = x^T A y$$

$$\begin{cases} \langle x, x \rangle \geq 0 \quad \forall x \\ \langle x, x \rangle = 0 \Leftrightarrow x = 0 \\ \langle \cdot, z \rangle \text{ is linear} \\ \langle x, y \rangle = \langle y, x \rangle \\ \quad\quad \uparrow \\ \quad\quad \text{Real} \end{cases}$$

**Def**: 1. Two vectors $x$ & $y$ $A$-conjugate if
$$\langle x, y \rangle_A = 0.$$

2. Given a linear subspace $L \subseteq \mathbb{R}^d$
$$L_A^\perp = \{ y \mid \langle x, y \rangle_A = 0 \; \forall x \in L \}.$$

3. The projection of $x$ onto $y$ w.r.t $\langle \cdot, \cdot \rangle_A$ is
$$P_y^A(x) = \frac{\langle x, y \rangle_A}{\langle y, y \rangle_A} \, y.$$

**Lemma**: Let $s_1, \ldots, s_k$ be $A$-conjugate pairwise. Then, they are linearly independent and for all $x \in \text{span}\{s_i\}_{i=1}^k$

$$x = \sum_{i=1}^{k} P^A_{s_i}(x).$$

Note that if $s_1, \ldots, s_d$ are A-conjugate then

$$S = (s_1 \ldots s_d)$$

yields

$$\min_{x \in \mathbb{R}^d} \frac{1}{2} x^T A x - bx \underbrace{=}_{\text{change of basis}} \min_{y \in \mathbb{R}^d} \frac{1}{2} (sy)^T A (sy) - b^T s y$$

$$= \min \sum_{i,j} y_i \, y_j \, s_i^T A s_j$$

$$\underbrace{\phantom{=}}_{\text{Decomposable}} \qquad - \sum_i b^T s_i y_i$$

$$= \min \sum_i y_i^2 \, s_i^T A s_i - b^T s_i y_i$$

$$\Rightarrow \quad y_i^* = \frac{b^T s_i}{s_i^T A s_i} \qquad \Rightarrow \quad x^* = \sum_{i=1}^{d} \frac{s_i s_i^T b}{s_i^T A s_i}.$$

Question: How do we obtain conjugate vectors?

Gram-Schmidt Orthogonalization

Input: $A > 0$, and linearly independent $x_1, \ldots, x_k$

Output: $s_1, \ldots, s_k$ A-conjugates s.t.
$$\text{span}\{s_i\} = \text{span}\{x_i\}.$$

▷ $S_1 = x_1$

▷ Recursively update
$$s_{i+1} = x_{i+1} - \sum P^A_{s_i}(x_{i+1})$$

Check:

▷ $\text{span}\{s_j\} = \text{span}\{x_j\}$

▷ $\langle s^{i+1}, s^j \rangle = 0 \quad \forall j < i+1$

▷ $s^{i+1} \neq 0$.

This algorithm is nice but when $K = d$, we have that we have to do $d$ steps each with complexity $O(d^2) \Rightarrow O(d^3)$ complexity.

↑ Matrix multiplication to compute $\langle x, s_i \rangle_A$.

Question: Can we find a good approximation of the solution of $Ax = b$ without doing $O(d^3)$ work?

# Conjugate Gradient Method

**Idea:** Construct the basis $s_1, \ldots, s_k$ using the residuals

$$r_k = b - Ax_k = \nabla f(x_k).$$

Then select

$(*)$ $\quad x_{k+1} = \arg\min \; f(x)$
$$\text{s.t. } x \in x_0 + \text{span}\{s_1, \ldots, s_k\}.$$

Let us see two supporting results

**Lemma 5:** Let $x_0$ and $s_1, \ldots, x_k$ be any vectors. Consider $x_{k+1}$ given by $(*)$, then $\nabla f(x_{k+1})$ is orthogonal (in the standard sense) to span $\{s_1, \ldots, s_k\}$.

**Proof:** Equivalently

$$y^* \in \arg\min_{y \in \mathbb{R}^k} f(x_0 + Sy)$$

By $1^{st}$-order optimality conditions:

$$S^T \nabla f(\underbrace{x_0 + Sy^*}_{x_{k+1}}) = 0$$

$$\Rightarrow \nabla f(x_{k+1}) \text{ is orthogonal to } span\{s_1, ..., s_k\}$$

$\square$

Thanks to separability:

**Lemma ⋈:** Supose that $x_{k+1}$ is given by (⋇) and $s_{k+1}$ is A-conjugate to each $s_i$. Then,

$$x_{k+2} \in argmin \; f(x)$$
$$s.t \; x = x_{k+1} + span\{s_{k+1}\}$$

is also a solution of
$$x_{k+2} \in argmin \; f(x)$$
$$s.t \; x = x_0 + span\{s_1, ..., s_{i+1}\}.$$

**CG Method**
Input: $x_0 \in \mathbb{R}^d$, $s_0 = r_0 = b - Ax_0$
Update $i \leq d$:

$$\alpha_i = \underset{\alpha}{argmin} \; f(x_i + \alpha s_i) \quad \swarrow \quad \alpha_i = \frac{s_i^T(b - Ax_i)}{\langle s_i, s_i \rangle_A}$$

$$x_{i+1} = x_i + \alpha_i s_i$$

$$r_{i+1} = -\nabla f(x_{i+1}) = b - Ax_{i+1}$$

$$s_{i+1} = r_{i+1} - \sum_i P_{s_i}^A (r_{i+1}) \quad \nwarrow$$

$$\text{Gram-schmidt}$$