**Recall:**

GD: $$x^{(t+1)} = x^{(t)} - \mu^{(t)} \nabla f(x^{(t)})$$

Can be fixed, or variable

e.g. decreasing

or chosen via linesearch

**GD w/ momentum:**

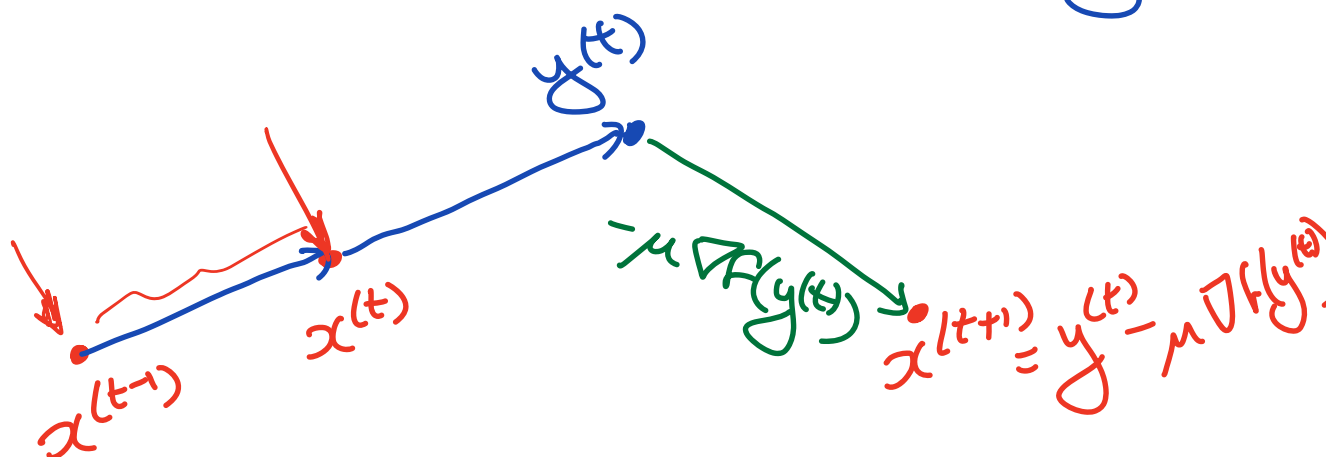$$x^{(t+1)} = x^{(t)} - \mu \nabla f(x^{(t)}) + \beta (x^{(t)} - x^{(t-1)})$$

**Variation:** Nesterov's Accelaration

$$y^{(t)} = x^{(t)} + \beta (x^{(t)} - x^{(t-1)}) \quad \text{—} \; ① $$

$$x^{(t+1)} = y^{(t)} - \mu \nabla f(y^{(t)}) \quad \text{—} \; ② $$

# Interpretation:

① Take a "momentum step" so you land at $y^{(t)}$

② Take a GD step from $y^{(t)}$.

$y^{(t)}$

$x^{(t-1)}$

$x^{(t)}$

$-\mu \nabla f(y^{(t)})$

$x^{(t+1)} = y^{(t)} - \mu \nabla f(y^{(t)})$.

We can of course combine ① & ②
to get

$$x^{(t+1)} = x^{(t)} + \beta\left(x^{(t)} - x^{(t-1)}\right) - \mu \nabla f\left(x^{(t)} + \beta\left(x^{(t)} - x^{(t-1)}\right)\right)$$

$\hookrightarrow$ nesterov's accelaration

$$\boxed{\text{Converges at an accelerated rate for } \underline{\text{any}} \text{ convex problem}}$$

with rate $= \sqrt{\dfrac{\sqrt{\varkappa}-1}{\sqrt{\varkappa}}}$

$\hookrightarrow f(x) = \frac{1}{2}x^T A x$

$\varkappa = \dfrac{\lambda_{max}}{\lambda_{min}}$

arises with the optimal choice of

$\mu = \dfrac{1}{\lambda_{max}}$, $\beta = \dfrac{\sqrt{\varkappa}-1}{\sqrt{\varkappa}+1}$

in the quadratic case

$\hookrightarrow$ compare to GD, GD/momentum

• used in practice

• speed up GD significantly.

# Conjugate Gradient Methods

- Originally developed in the 50's by Hestenes & Steifel for solving large linear systems. i.e. $Ax = b$

- First nonlinear conjugate gradient (CG) methods were developed in the 60's by Fletcher & Reeves for solving non-linear optimization problems.

## Basic Idea:

Suppose that we want to minimize

$$\underline{\phi}(x) = \frac{1}{2}x^{T}Ax - b^{T}x \qquad (\ast)$$

where $A$ is a <u>symmetric</u> <u>positive</u> <u>semidefinite</u> $n \times n$ matrix.

Notice: • $\underline{\phi}$ is convex &
- $\nabla \underline{\phi}(x) = Ax - b$

$\Rightarrow$ the optimizer satisfies $Ax^{\ast} = b$

( which links solving linear systems with optimizing (*) ).

Where does the word "conjugate" in CG methods come from?

Def'n: A set of vectors $\{P_1, \dots, P_\ell\}$

$P_i \neq P_j, \forall i \neq j$

is conjugate with respect to a

sym. PSD matrix $A$ if

$$P_i^T A P_j = 0 \qquad \forall i \neq j$$

Example: The standard basis vectors $e_1, e_2, \dots, e_\ell \in \mathbb{R}^n$ where $\ell \leq n$ are conjugate w.r.t to $I$

( bec $e_i^T I e_j = 0 \ \forall i \neq j$ )

# How does conjugacy help us?

It turns out we can minimize $\Phi(x)$ in (*) in at most $n$ steps by successively minimizing it along the individual directions in a conjugate set.

How? Conjugate direction method:

$$x^{(t+1)} = x^{(t)} + \alpha_t P_t \qquad - \quad ①$$

scalar     vector $\in \mathbb{R}^n$

unknown

where $\quad \alpha_t = \underset{a \in \mathbb{R}}{\text{argmin}} \; \Phi(x^{(t)} + a P_t)$

known

(i.e. pick the best $a$)

(i.e. perform an exact line search)

Since $\Phi$ is quadratic, we can solve for $\alpha_t$ exactly:

$$\underline{\phi}(x + aP) = \frac{1}{2}(x+ap)^T A (x+ap) - (x+ap)^T b$$

<span style="color:red">↑ given ↑ given</span>

<span style="color:red">trying to optimize</span>

To find optimal $a$, solve $\dfrac{d\underline{\phi}}{da}(x+ap) = 0$

Chain rule $\Rightarrow \underline{\phi}'(x+ap) \cdot (x+ap)' = 0$

$$\underbrace{[A(x+ap) - b]^T}_{} \qquad \underbrace{}_{P}$$

thus $\left( A(x+ap) - b \right)^T P = 0$

$$\Rightarrow a^* = \boxed{\alpha_t = \frac{(b - Ax^{(t)})^T P_t}{P_t^T A P_t}}$$

<span style="color:red">Notice!</span>

$$\boxed{\alpha_t = \frac{-\nabla\underline{\phi}(x^{(t)})^T P_t}{P_t^T A P_t}}$$

$\left( \text{bec} \quad b - Ax = -\nabla\underline{\phi}(x) \right)$

**Theorem:** For any $x^{(0)}$ the sequence $\{x^{(t)}\}$ generated by ① converges to $x^*$, the optimizer of $\Phi$ in at most $n$-steps.

$\hookrightarrow \mathbb{R}^n \to \mathbb{R}$

Won't prove it: Proof is in
Nocedal & Wright (2006)
Chapter 5.

## What is the conjugate gradient method?

Unlike the conjugate directions method, here we compute $P_t$ as part of the algorithm.

- Can compute $P_t$ using only $P_{t-1}$

- So, don't need to store or compute with $P_0, P_1, \cdots, P_{t-2}$

- $P_t$ is automatically conjugate to $P_0, P_1, \ldots, P_{t-1}$ !

How? :

$$P_t = -\nabla \bar{\phi}(x^{(t)}) + \beta_t P_{t-1}$$

vector (under $P_t$)

vector (under $-\nabla \bar{\phi}(x^{(t)})$)

scalar (over $\beta_t$)

vector (over $P_{t-1}$)

chosen to enforce conjugacy

How to pick $\beta_t$ then?

$$P_{t-1}^T A \left[ P_t = -\nabla \bar{\phi}(x^{(t)}) + \beta_t P_{t-1} \right]$$

$$\Rightarrow \quad P_{t-1}^T A P_t = -\underbrace{P_{t-1}^T A \nabla \bar{\phi}(x^{(t)})}_{\text{scalar}}$$

0 by conjugacy

$$+ \beta_t \underbrace{P_{t-1}^T A P_{t-1}}_{\text{scalar}}$$

$$\Rightarrow \quad \boxed{\beta_t = \frac{P_{t-1}^T A \nabla \bar{\phi}(x^{(t)})}{P_{t-1}^T A P_{t-1}}}$$

What about $\beta_0$? Just set $\beta_0 = 0$

$\implies$ we have an algorithm.

- Initialize $x^{(0)}, P_0 = -\nabla \bar{\phi}(x^{(0)})$

- for $t = 1, 2, \ldots$

  + $\beta_t = \dfrac{P_{t-1}^T A \nabla \bar{\phi}(x^{(t)})}{P_{t-1}^T A P_{t-1}}$

  * $P_t = -\nabla \bar{\phi}(x^{(t)}) + \beta_t P_{t-1}$

  * $d_t = -\dfrac{\nabla \bar{\phi}(x^{(t)})^T P_t}{P_t^T A P_t}$

  * $x^{(t+1)} = x^{(t)} + d_t P_t$ $\leftarrow$ Should be reminiscent of momentum

CG: Version 0

**Theorem:** If $\underline{\phi}(x) = \frac{1}{2} x^T A x - b^T x$ with $A$ being $n \times n$ symmetric PSD, then CG converges to $x^*$ in at most $n$ steps

Proof: See reference 🗒

More efficient implementation (needs properties of CG to derive, won't do it) is

- Initialize $x^{(0)}$, $r_0 = A x^{(0)} - b = \nabla \overline{\phi}(x^{(0)})$
  $$P_0 = -r_0$$

- $\alpha_t = \dfrac{r_t^T r_t}{P_t^T A P_t}$

- $x^{(t+1)} = x^{(t)} + \alpha_t P_t$

- $r_{t+1} = r_t + \alpha_t A p_t$

- $\beta_{t+1} = \dfrac{r_{t+1}^T r_{t+1}}{r_t^T r_t}$

- $p_{t+1} = -r_{t+1} + \beta_{t+1} p_t$

CG (version 1)

Cost per iteration :

- matrix vector mult :

$$A p_t$$

- $p_t^T (A p_t)$

  already computed

- $r^T r$

- not much more expensive than GD

Can we improve CG (in the quadratic
setting $\phi(x) = \frac{1}{2} x^T A x - b^T x$   ?
                              $\underset{PD, Sym.}{\uparrow}$

Idea : Preconditioning

First, some background :

* Recall that we seek $x^*$ : $Ax^* = b$
                                    $\uparrow$
                          Sym. PD, $n \times n$

* It will turn out that the perform-
ance of CG depends on the
eigenvalues of $A$.

How? Define $\|z\|_A = \sqrt{z^T A z}$
                        $\underset{z \in \mathbb{R}^n}{\uparrow}$

   e.g. If $A = I$, then $\|z\|_A = \|z\|$
                                    $= \sqrt{z^T z}$

**Theorem:** If $A$ has eigenvalues

$0 < \lambda_1 \leq \lambda_2 \leq \ldots \leq \lambda_n$, then

$$\|x^{(t+1)} - x^*\|_A^2 \leq \left( \frac{\lambda_{n-t} - \lambda_1}{\lambda_{n-t} + \lambda_1} \right)^2 \|x^{(0)} - x^*\|_A^2$$

This implies, choosing $t+1 = n \iff t = n-1$

$$\|x^{(n)} - x^*\|_A^2 \leq 0 \cdot \|x^{(0)} - x^*\|_A^2 = 0$$

that is, convergence within $n$-steps!

Additionally, if for example

$\lambda_{n-1} = \lambda_n$, then convergence to $x^*$

takes only $n-1$ steps !