

# 1st lecture Convex Optimization.

We will see 3 results today.

→ Gradient Descent linear convergence for smooth, strongly convex functions.

→ Gradient Descent sublinear conv for smooth convex functions. Interlude Nesterov.

→ Stochastic Gradient Descent Convergence.

Notations / Definitions →  $F(x)$  convex

if  $F(\alpha x + (1-\alpha)y) \leq \alpha F(x) + (1-\alpha)F(y)$

Alternative:

$$F(y) \geq F(x) + \langle \nabla F(x), y-x \rangle$$
$$F(y) - F(x) = \int_0^1 \langle \nabla F(x + \alpha(y-x)), y-x \rangle d\alpha \quad (\forall y, x)$$

$$(1-\alpha) \int_0^\alpha f'(t) dt \leq \alpha \int_\alpha^1 f'(t) dt$$

$$f(1) - f(0) = \int_0^1 f'(t) dt$$

$$(1-\alpha) [f(\alpha) - f(0)] \leq \alpha [f'(1) - f'(\alpha)]$$

$$(1-\alpha) f(\alpha) + \alpha f(\alpha) \leq \alpha f'(1) + (1-\alpha) f(0)$$

(d/c)

$$\text{so } \frac{(1-\alpha)}{\alpha} \int_0^\alpha f'(t) dt \leq \int_\alpha^1 f'(t) dt$$

$$\text{so } \lim_{\alpha \rightarrow 0} (1-\alpha) \left[ \frac{1}{\alpha} \int_0^\alpha f'(t) dt \right] \leq \lim_{\alpha \rightarrow 0} \int_\alpha^1 f'(t) dt$$

$$f'(0) \leq \int_0^1 f'(t) dt.$$

$\rightarrow F(x)$  strongly convex

if  $F(y) \geq F(x) + \langle \nabla F(x), y-x \rangle + \frac{L}{2} \|y-x\|^2$   
 ( $H(x)$  has eigenvalue at least  $L$ ). Taylor

$\rightarrow F$  is smoothly convex if  $\nabla F$  is

Lipschitz:

$$\|\nabla F(x) - \nabla F(y)\| \leq L \|x-y\|.$$

This implies by Taylor that

$$F(y) \leq F(x) + \langle \nabla F(x), y-x \rangle + \frac{L}{2} \|x-y\|^2$$

Strongly convex case:

$$\|x - t \nabla F(x) - (y - t \nabla F(y))\| \leq \max(|1 - t\ell|, |1 - tL|) \|x - y\|.$$

$\nabla F(x)$  is Bilipschitz with constants  $\ell, L$   
 $\rightarrow H F(x)$  has eigenvals in  $(\ell, L)$   
 $\Rightarrow H t F$  has constants  $(t\ell, tL)$   
 $\Rightarrow I - t \nabla F$  is Lipschitz w/ constant  $\max(|1 - t\ell|, |1 - tL|)$ .  
 $\rightarrow$  say it in terms of Hessians instead.

$$\|x_{k+1} - x_*\| \leq \|x_k - t \nabla F(x_k) - x_* - t \nabla F(x_*)\|$$

$$\leq \max(|1 - t\ell|, |1 - tL|) \|x_k - x_*\|$$

$t := \frac{2}{L + \ell}$  minimizes and we obtain

$$\left( \frac{L - \ell}{L + \ell} \right)^k; \text{ in terms of condition number } K = \frac{L}{\ell},$$

$$\left( \frac{K - 1}{K + 1} \right)^K \quad \square$$

Ex:  $F(x) = \frac{1}{2} \|\tilde{A}x - \tilde{b}\|^2 = \frac{1}{2} x^T A x - b^T x + c$

$$\nabla F(x_k) = A x_k - b$$

$$A = \tilde{A}^T \tilde{A} \\ b = \tilde{A}^T \tilde{b}$$

$$x_0 = tb$$

$$x_1 = x_0 - t \nabla F(x_0) =$$

$$= x_0 - t(A x_0 - b) = (I - tA) x_0 + tb \\ = (I - tA) x_0 + x_0.$$

$$x_2 = x_1 - t \nabla F(x_1) = (I - tA) ((I - tA) x_0 + x_0) + x_0.$$

$$\Rightarrow x_k = \left[ \sum_{j=1}^k (I - tA)^j \right] (x_0).$$

$$\frac{1}{t} = \sum (1-z)^k \quad |z| < 1$$

$A$  has eigenvalues  $(l, L)$

$\rightarrow tA$  has eigenvalues  $\frac{t}{l+L} (l, L)$

$\rightarrow I - tA$  has eigenvalues  $\left( \frac{2l-l-L}{l+L}, \frac{2L-l-L}{l+L} \right)$

$\left( \frac{l-L}{l+L}, \frac{L-l}{L+l} \right)$  between  $(-1, 1)$

So the series  $\sum_{j=0}^{\infty} (I - tA)^j$  converges to  $A^{-1}$ , with rate  $O(\|(I - tA)^k\|) = (1 - \frac{t}{L})^k$

---

Q: Can we do better?

Find a  $k$ -degree polynomial  $q_k(A)$  minimizing the residual error

$$\|(I - A q_k(A))b\|$$

The polynomial corresponding form is

$$P_k(z) = 1 - z q_k(z)$$

When applied to  $A$ , small norm

$\Rightarrow$  eigenvalues of  $A$  are bounded in the interval  $(l, L)$ . also, we need

$$P_k(0) = 1.$$

$\hookrightarrow$  Chebyshev polynomials are optimum at  $A$ .  $\dagger$

inner.

Lemma: There exist  $P_k(z)$  of degree  $O(\sqrt{(L/\epsilon) \log(1/\epsilon)})$  with  $P_k(0) = 1$  and  $|P_k(z)| \leq \epsilon \quad \forall z \in [0, L]$ .

quadratic savings in degree.

$$O\left(\sqrt{\frac{k-1}{k+1}}\right)^k$$

$$\tilde{p}^k = p^{\sqrt{k}}$$

Moreover, Chebyshev polynomials can be obtained recursively using two previous polynomials.

$$\Rightarrow X_{k+1} = X_k - \alpha_k \nabla f(X_k) - \beta_k \nabla f(X_{k-1})$$

- ↳ This extends to generic convex functions (Westerov)
- ↳ This rate is optimal: cannot be improved
- ↳ Alternative interpretation of acceleration (Bubeck, Lee, Singh) (using a variant of ellipsoid method).

$$B = D^T D \in \mathbb{R}^{k \times k}$$

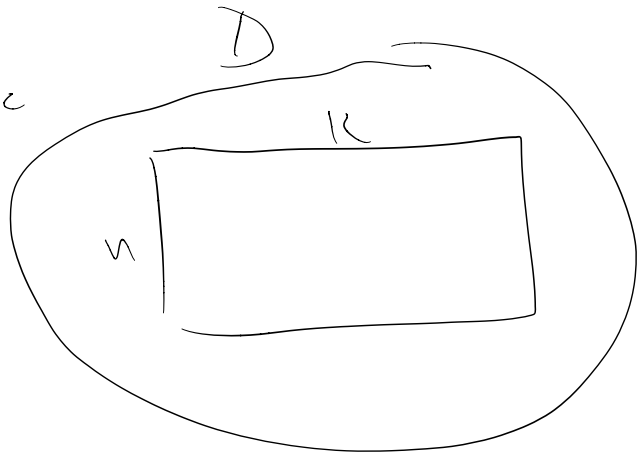
$$D \in \mathbb{R}^{n \times k} \quad (n < k)$$

$$n, F^T \Lambda F = F^T \tilde{\Lambda} + \tilde{\Lambda} F$$

$$D = \cup_{i=1}^n \cup_{j=1}^k \cup_{l=1}^m$$

$$D = \tilde{\lambda} F$$

$$F = F_n \in \mathbb{R}^{n \times k}$$



$$\|z\|_1 \in \mathbb{R}^k$$

$$F \in \mathbb{R}^{k \times k}$$

$$\|Fz\|_1$$

$$A \in \mathbb{R}^{k \times n}$$

$$\|Fz\|_1 \approx \|Az\|_1 \cdot \left(\frac{k}{n}\right)$$

$$\|z\|_1$$