

Lecture II: Still semi-classical convex opt

Announcements

* Switching to hybrid Thursdays.
Loomis - 10/7; Siebel 1105 10/14

Recap: Optimization

GD: $w' := w - \eta \nabla F(w)$

β -smooth: $\|\nabla F(w) - \nabla F(v)\| \leq \beta \|w - v\|$

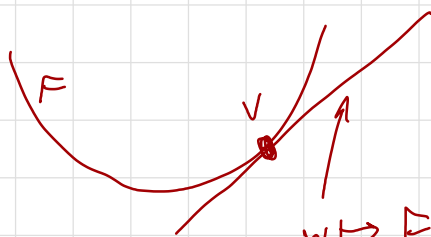
$$\Rightarrow F(w) \leq F(v) + \langle \nabla F(v), w - v \rangle + \frac{\beta}{2} \|w - v\|^2$$

$$\Rightarrow \min_{\text{ict}} \|\nabla F(w_t)\|^2 \leq \frac{2\beta}{t} (F(w_0) - F(w_t))$$

"approx stationary point"

F convex

$$F(w) \geq F(v) + \langle \nabla F(v), w - v \rangle$$



$$w \mapsto F(v) + \langle \nabla F(v), w - v \rangle$$

Theorem. F β -smooth, convex, $\eta = 1/\beta$, $\forall z \in \mathbb{R}^d$

$$F(w_t) \leq F(z) + \frac{\beta}{2t} (\|w_0 - z\|^2 - \|w_t - z\|^2)$$

↑
recap.

"self-regularization"
(see hawk p3,
or margin lectures)

Remark/reminder

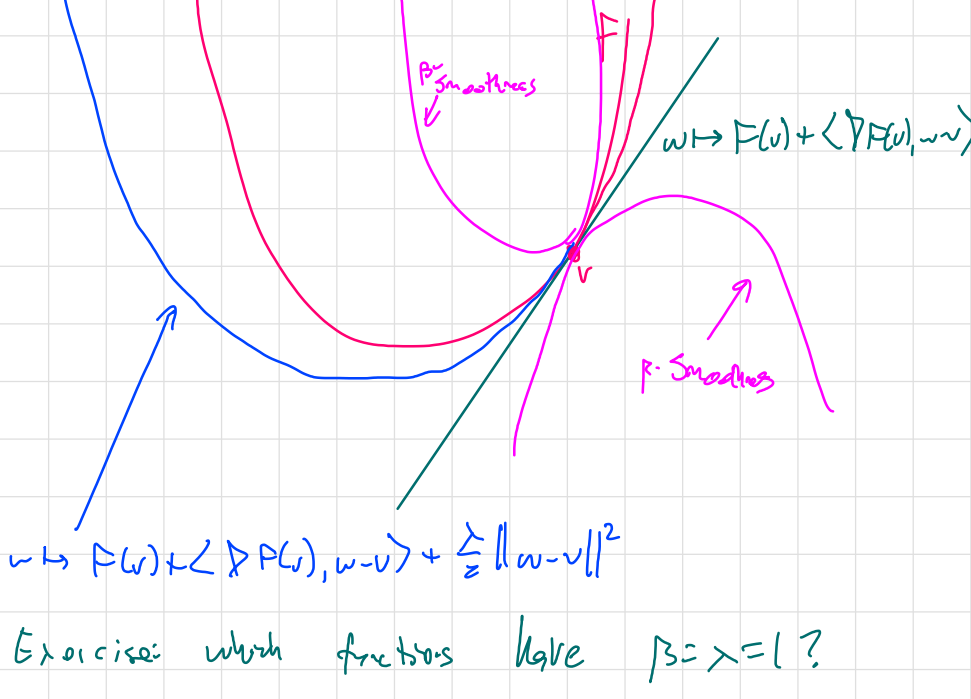
* Doing smooth & strong convexity analysis
because both proof schemes pervasive for NTH opt.

* Not doing SGD because theory story is still unclear.

Strong convexity:

F is " λ -sc" (λ -strongly-convex) if

$$\forall u, v \quad F(w) \geq F(v) + \langle \nabla F(v), w-v \rangle + \frac{\lambda}{2} \|w-v\|^2$$



Exercise: which functions have $\beta = \lambda = 1$?

Answer: $w \mapsto \frac{1}{2} \|w\|^2 + a^T w + b$.

Ex. Regularize convex $F \rightarrow F_\lambda(w) = F(w) + \frac{\lambda}{2} \|w\|^2$

$\Rightarrow F_\lambda$ is λ -sc.

Unnormalized least squares

$$w \mapsto \frac{1}{2} \|Xw - y\|^2$$

if $\sigma_{\min}(X) > 0$
 $\sigma_{\min}(X)$ -sc, $\sigma_{\max}(X)$ -smooth
 (might be 0 (not sc!))

Recall β -smooth $\Rightarrow F(w - \frac{1}{\beta} \nabla F(w)) \leq F(w) - \frac{1}{2\beta} \|\nabla F(w)\|^2$

Lemma. F is λ -sc,

$$F(w - \frac{1}{\lambda} \nabla F(w)) \geq F(w) - \frac{1}{2\lambda} \|\nabla F(w)\|^2$$

and $\inf_v F(v) \geq F(w) - \frac{1}{2\lambda} \|\nabla F(w)\|^2$.

or $F(w) \leq \frac{1}{2\lambda} \|\nabla F(w)\|^2 + \inf_v F(v)$.

Proof.

$$Q_w(v) := F(w) + \langle \nabla F(w), v-w \rangle + \frac{\lambda}{2} \|w-v\|^2$$

$$\nabla F(w) + \lambda(w-v) = 0$$

$$v = w - \frac{1}{\lambda} \nabla F(w)$$

$$\begin{aligned} \inf_v Q_w(v) &= Q_w(w - \frac{1}{\lambda} \nabla F(w)) \\ &= F(w) + \langle \nabla F(w), w - \frac{1}{\lambda} \nabla F(w) \rangle \\ &\quad + \frac{\lambda}{2} \|\frac{1}{\lambda} \nabla F(w)\|^2 \\ &= F(w) - \frac{1}{2\lambda} \|\nabla F(w)\|^2. \end{aligned}$$

$$\inf_v F(v) \geq \inf_v Q_w(v) = F(w) - \frac{1}{2\lambda} \|\nabla F(w)\|^2$$

Reminder. If F is β -smooth, \Rightarrow convergence rate $\frac{1}{t}$.

Theorem. If F is β -smooth & λ -sc, $F(\bar{w}) = \inf_v F(v)$, and

$$F(w_t) - F(\bar{w}) \leq (F(w_0) - F(\bar{w})) \exp(-\frac{\lambda}{\beta} t)$$

$$\|w_t - \bar{w}\|^2 \leq \|w_0 - \bar{w}\|^2 \exp(-\frac{\lambda}{\beta} t)$$

Remark. See Nesterov's book for simultaneous proof & faster rate.

Proof. See notes for last part.

$$F(w_{i+1}) - F(\bar{w}) \leq F(w_i) - F(\bar{w}) - \frac{1}{2\beta} \|\nabla F(w_i)\|^2$$

$$\leq F(w_i) - F(\bar{w}) - \frac{\lambda}{\beta} (F(w_i) - F(\bar{w}))$$

$$= (1 - \frac{\lambda}{\beta}) (F(w_i) - F(\bar{w}))$$

$$\stackrel{\text{unroll}}{\Rightarrow} (F(w_t) - F(\bar{w})) \leq (F(w_0) - F(\bar{w})) (1 - \frac{\lambda}{\beta})^t \leq (F(w_0) - F(\bar{w})) \exp(-\frac{\lambda}{\beta} t)$$

Remark/PSA.

$$F \geq 0, \quad F_\lambda(w) = F(w) + \frac{\lambda}{2} \|w\|^2$$

Consider $F_\lambda(v) \leq F_\lambda(0)$.

$$\frac{\lambda}{2} \|v\|^2 \leq F_\lambda(v) \leq F_\lambda(0) = F(0) + \frac{\lambda}{2} \|0\|^2 = F(0)$$

$$\Rightarrow \|v\| \leq \sqrt{\frac{2F(0)}{\lambda}}$$

$\{v \in \mathbb{R}^p : F_\lambda(v) \leq F_\lambda(0)\} \subseteq \{v \in \mathbb{R}^p : \|v\| \leq \sqrt{\frac{2F(0)}{\lambda}}\}$

i.e. regularized \Rightarrow constrained.

GF (gradient flow). $\frac{f(w(t)) - f(w_0)}{t} = \nabla F(w)$

$\dot{w}(t) = \frac{dw(t)}{dt} = -\nabla F(w(t))$. ("infinitesimal gradient descent")

Why: it simplifies many proofs
 & seems to offer not change story.

Remark: "works" for nonsmooth, see later lectures and/or "Clarke differential"

Stationary points:

$$\begin{aligned}
 F(w(t)) - F(w(0)) &\stackrel{\text{FTC}}{=} \int_0^t \frac{d}{ds} F(w(s)) ds \\
 &\stackrel{\text{chain rule}}{=} \int_0^t \langle \nabla F(w(s)), \dot{w}(s) \rangle ds \\
 &\stackrel{\substack{\dot{w}(s) \\ \text{ODE}}}{=} \int_0^t - \langle \nabla F(w(s)), \nabla F(w(s)) \rangle ds \\
 &= \int_0^t - \|\nabla F(w(s))\|^2 ds \\
 &\leq \int_0^t - \inf_{r \in [0, t]} \|\nabla F(w(r))\|^2 ds \\
 &= -t \cdot \inf_{r \in [0, t]} \|\nabla F(w(r))\|^2
 \end{aligned}$$

Theorem. $\inf_{r \in [0, t]} \|\nabla F(w(r))\|^2 \leq \frac{1}{t} (F(w(0)) - F(w(t)))$.

Now let's use convexity.

$$\begin{aligned}
 &\frac{1}{2} \|w(t) - z\|^2 - \frac{1}{2} \|w(0) - z\|^2 \stackrel{\text{FTC}}{=} \int_0^t \frac{d}{ds} \frac{1}{2} \|w(s) - z\|^2 ds \\
 &= \int_0^t \langle w(s) - z, \dot{w}(s) \rangle ds \stackrel{\text{chain rule}}{=} \int_0^t \langle w(s) - z, -\nabla F(w(s)) \rangle ds \\
 &\stackrel{\text{characterization of } \dot{w}(s)}{=} \int_0^t \langle w(s) - z, -\nabla F(w(s)) \rangle ds \stackrel{\text{convexity}}{\leq} \int_0^t (F(z) - F(w(s))) ds \\
 &\leq \int_0^t (F(s) - F(w(t))) ds \stackrel{\substack{F(w(r)) \\ \text{nonincreasing}}}{\leq} \int_0^t (F(s) - F(w(t))) ds \\
 &= t(F(z) - F(w(t))). \text{ Rearrange:}
 \end{aligned}$$

Theorem. F is convex, $z \in \mathbb{R}^d$ arbitrary,

$F(w(t)) \leq F(z) + \frac{1}{2t} (\|w(0) - z\|^2 - \|w(t) - z\|^2)$.

OFFICE HOURS

$$VV^T \neq I \quad \text{if } k < d$$

$$\begin{bmatrix} | & & | \\ v_1 & \dots & v_n \\ | & & | \end{bmatrix} \in \mathbb{R}^{d \times k}$$

$$V^T V = \begin{bmatrix} -v_1^T & & \\ \vdots & & \\ -v_n^T & & \end{bmatrix} \begin{bmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{bmatrix} \in \mathbb{R}^{k \times k}$$

$$V = \begin{bmatrix} | & & | \\ v_1 & \dots & v_n \\ | & & | \end{bmatrix} \in \mathbb{R}^{d \times k}$$

$k < d$
orthonormal
means

$$v_i^T v_j = \begin{cases} 1 & i=j \\ 0 & i \neq j \end{cases}$$

F continuously differentiable in an open set
containing $\{v \in \mathbb{R}^d : F(v) = F(\omega(0))\}$.

