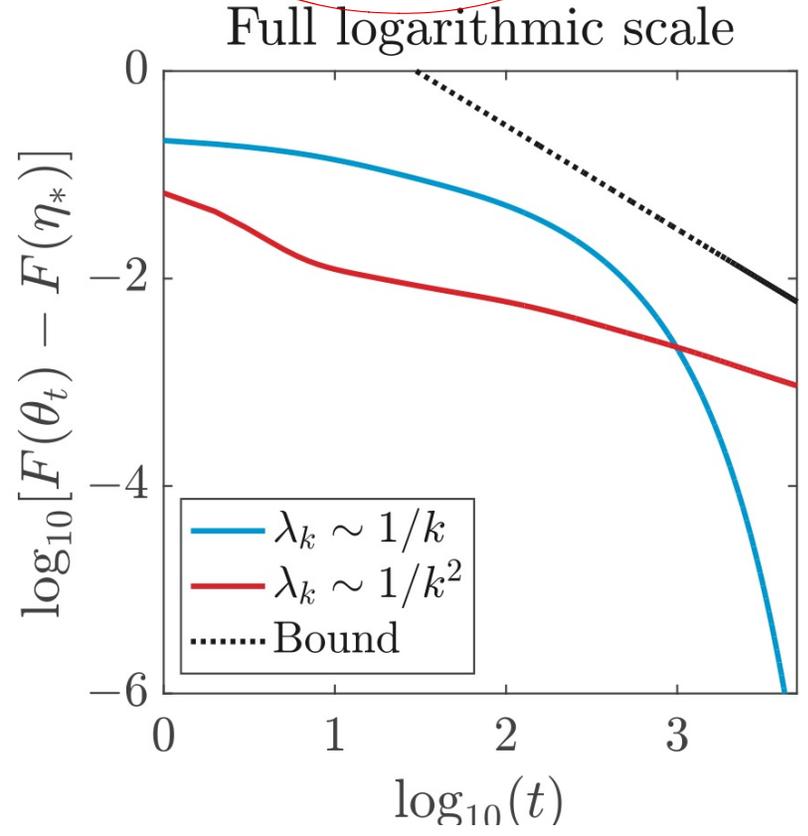
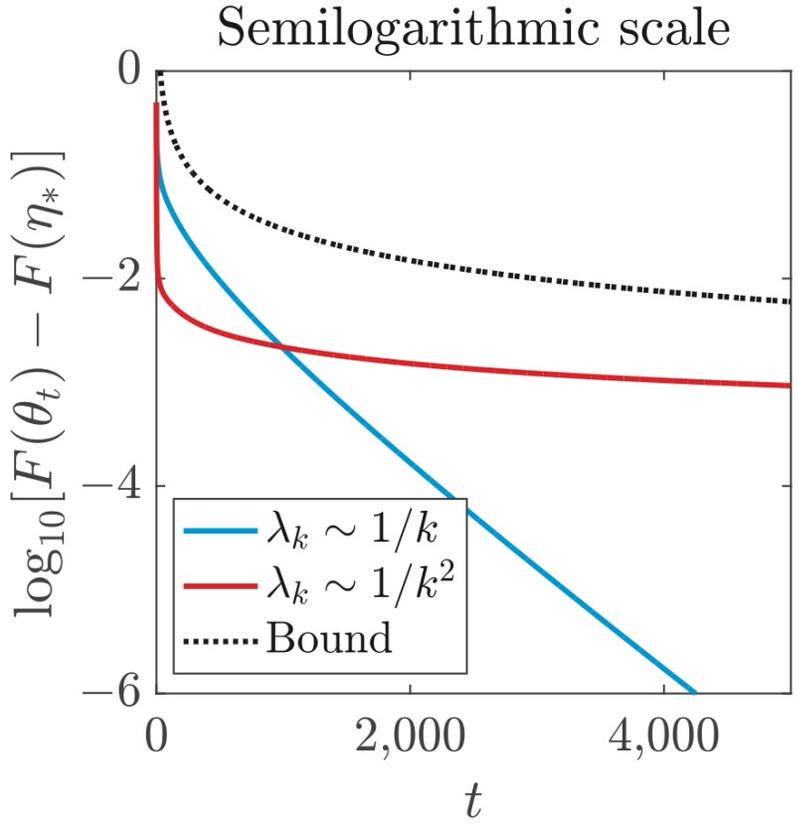
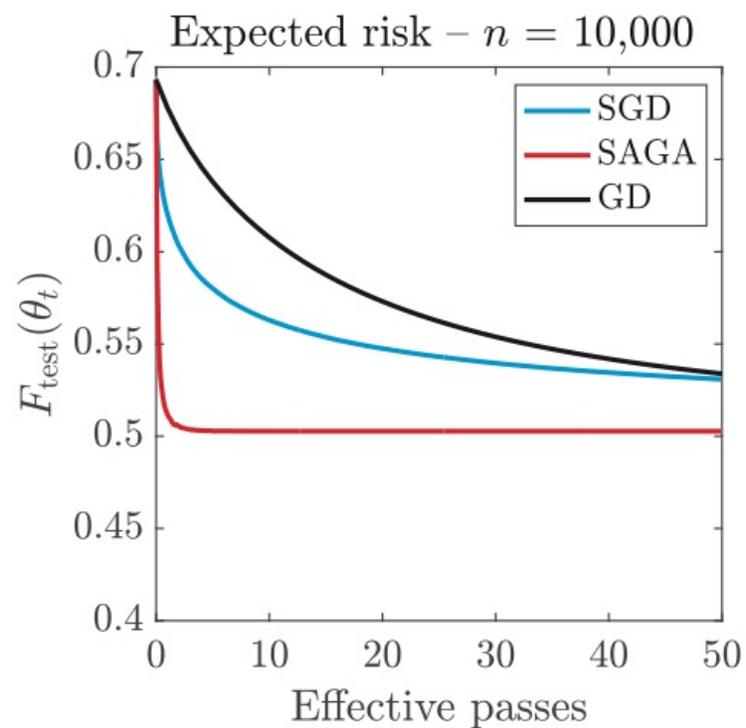
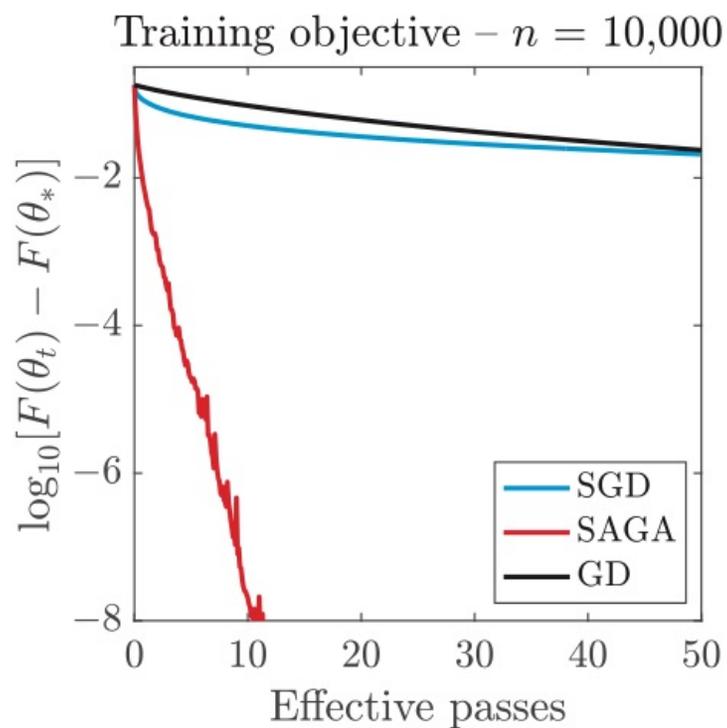
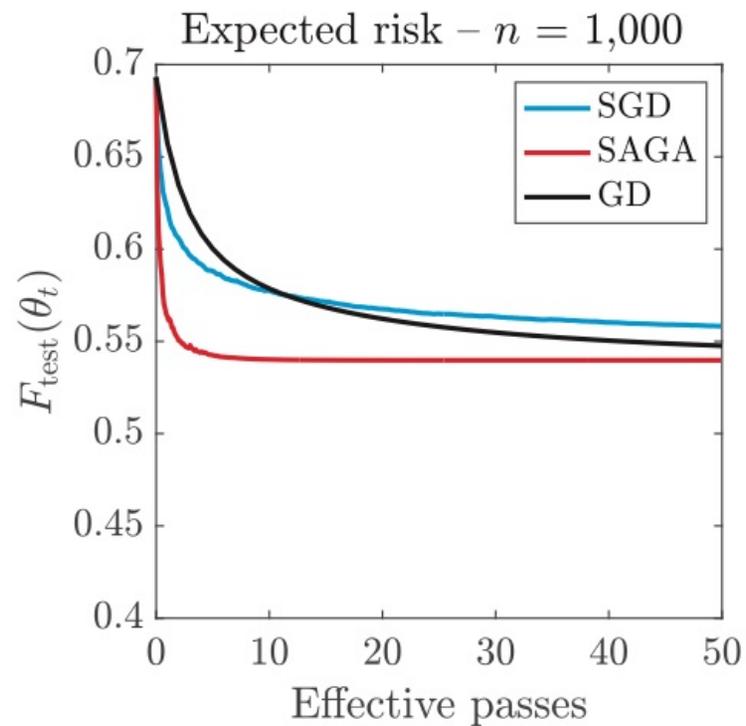
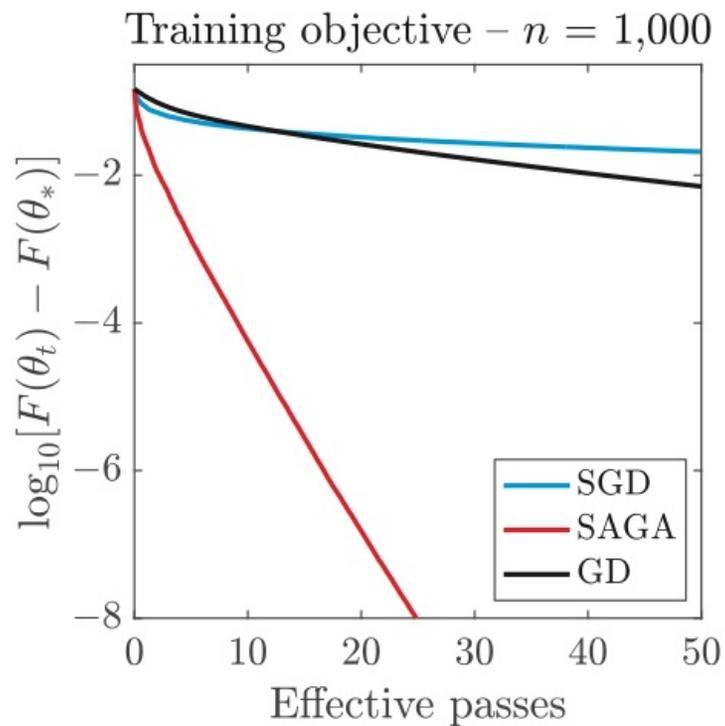


Flags post



$\lambda_k$  eigenvalues  $\lambda_k$ ;  $k \in \{1, \dots, n\}$   
 $\lambda_k = \frac{1}{k^2}$   $\lambda_k = \frac{1}{k}$



# Optimization

data  $(x_i, y_i)_{i=1, \dots, n} \in X \times \mathbb{R}$

Expected risk  $R(\beta) = \mathbb{E}[e(y, \beta(x))]$

$L$ , convex w.r.t  $\beta(x)$   
"weak"

$\hat{R}(\beta) = \frac{1}{n} \sum_{i=1}^n e(y_i, \beta(x_i))$  Empirical Risk min.

Model  $\beta_\Theta: X \rightarrow \mathbb{R}$

$\Theta \subseteq \Theta \subset \mathbb{R}^d$

Algorithm: minimize  $F(\Theta) = \hat{R}(\beta_\Theta) + \text{regularizer}$   
 $\Omega(\Theta)$

"objective function"

$\hat{\Theta}$  approximate minimizer.

Estimate error

$$R(\hat{\beta}_0) - \inf_{\beta_0 \in \mathcal{R}^d} R(\beta_0) = R(\hat{\beta}_0) - \hat{R}(\hat{\beta}_0) + \hat{R}(\hat{\beta}_0) - \hat{R}(\beta_{0^*}) + \hat{R}(\beta_{0^*}) - R(\beta_{0^*})$$

uniform deviates  $O\left(\frac{1}{\sqrt{n}}\right)$

$\hat{\beta}_0$  minimizes

$$\leq \hat{R}(\hat{\beta}_0) - \inf_{\beta_0 \in \mathcal{R}^d} \hat{R}(\beta_0)$$

optimization error  $\rightarrow \eta^*$  minimizes

$$\Delta \eta^* \neq \hat{\beta}_0^*$$

precision: high precision  $\Leftrightarrow$  low error

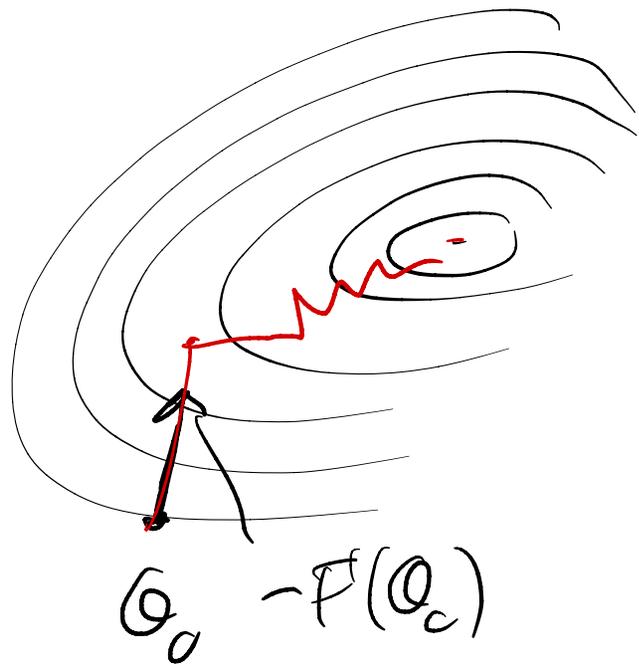
Goal: bounds on optimization error - dimension-independent

Goal: Objective function:  $F: \mathbb{R}^d \rightarrow \mathbb{R}$   
Find an approximate minimizer:  $\min_{\theta \in \mathbb{R}^d} F(\theta)$

Algo: Gradient descent  
 $\theta_0$  initial point

$$\theta_t = \theta_{t-1} - \gamma_t F'(\theta_{t-1}) \quad \text{--- gradient}$$

↳ step-size



two strategies

① constant or  
time dependent

② line search

↳ this class

least-squares:  $F(\theta) = \frac{1}{2} \|\Phi\theta - y\|_2^2$  response vector  $y \in \mathbb{R}^n$   
 $\Phi \in \mathbb{R}^{n \times d}$  design matrix

to avoid large values,

gradient:  $F'(\theta) = \frac{1}{n} \Phi^T (\Phi\theta - y) \in \mathbb{R}^d$

Minimizer:  $F'(\theta_*) = 0 \Leftrightarrow \frac{1}{n} \Phi^T \Phi \theta_* = \frac{1}{n} \Phi^T y$

$$\boxed{H_{\theta_*} = \frac{1}{n} \Phi^T \Phi}$$

$\Sigma = H$  Hessian matrix  
 normal equations

Symmetric positive semi-definite

Characterizing convergence: (1) "iterates"  $\|\theta_t - \theta_*\|_2^2$

(2) "function values"

$$F(\theta_t) - F(\theta_*) = \frac{1}{2} (\theta_t - \theta_*)^T H (\theta_t - \theta_*)$$

f-th iterate (minimizer)

$$F'(q) = Hq - \frac{1}{n} \Phi^T y = Hq - Hq^*$$

$$Q_t = Q_{t-1} - \gamma F'(Q_{t-1}) = Q_{t-1} - \gamma H(Q_{t-1} - q^*) + \gamma \frac{1}{n} \Phi^T y$$

$$= Q_{t-1} - \gamma H(Q_{t-1} - q^*)$$

$$Q_t - q^* = Q_{t-1} - q^* - \gamma H(Q_{t-1} - q^*) = (I - \gamma H)(Q_{t-1} - q^*)$$

$$Q_t - q^* = (I - \gamma H)^t (Q_0 - q^*)$$

Function value:  $\frac{1}{2} (Q_t - q^*)^T H (Q_t - q^*)$

$H = \sum_{i=1}^d \lambda_i u_i u_i^T$  eigenvalue decomposition  
 $\lambda_i$  eigenvalues,  $u_i$  eigenvectors (orthogonal)

$$(I - \gamma H)^t = \sum_{i=1}^d (1 - \gamma \lambda_i)^t u_i u_i^T$$

convergence to  $q^*$  sufficient

$$\Leftrightarrow |1 - \gamma \lambda_i| < 1$$

$$\Leftrightarrow 1 - \gamma \lambda_i \in [0, 1)$$

$$\Leftrightarrow \begin{cases} 1 - \gamma L \in (0, 1) \\ 1 - \gamma H \in (0, 1) \end{cases}$$

$$\Leftrightarrow \begin{cases} \gamma L \leq 1 \\ \gamma > 0 \end{cases}$$

$$\boxed{\gamma \leq 1/L} \Rightarrow$$

depends on smallest  $\mu$  and largest  $L$  eigenvalues  $> 0$

asymptotically  
 eigenvalues of  $(I - \gamma H)^t \in (0, (1 - \gamma H)^t)$

Summary :  $Q_t - q_{\infty} = (1 - \gamma H)^t (Q_0 - q_{\infty})$

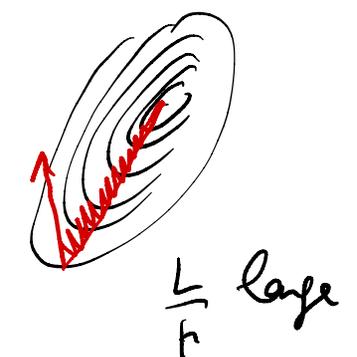
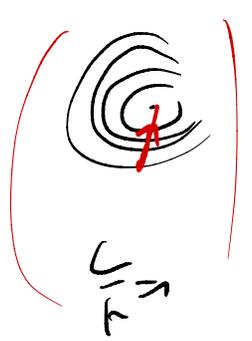
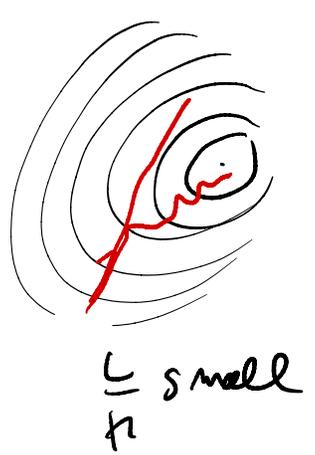
$F(Q_H) - F(q_{\infty}) = \frac{1}{2} (Q_H - q_{\infty})^T H (Q_H - q_{\infty})$   
 eigenvalues of  $(1 - \gamma H)^t \leq (1 - \gamma \lambda)^t$

$F(Q_t) - F(q_{\infty}) = \frac{1}{2} (Q_0 - q_{\infty})^T H (1 - \gamma H)^{2t} (Q_0 - q_{\infty})$   
 $\leq \frac{1}{2} (Q_0 - q_{\infty})^T H (Q_0 - q_{\infty}) \cdot (1 - \gamma \lambda)^{2t}$

$\hookrightarrow$  eigenvalues of  $H (1 - \gamma H)^{2t} \leq \lambda_i (1 - \gamma \lambda)^{2t}$

$F(Q_t) - F(q_{\infty}) \leq (1 - \gamma \lambda)^{2t} [F(Q_0) - F(q_{\infty})]$   
 $\hookrightarrow$  if  $\gamma = \frac{1}{L} \Rightarrow (1 - \frac{\lambda}{L})^{2t} \leq \exp(-\frac{\lambda}{L} 2t) = e^{-2}$  because  $1 - x \leq e^{-x}$

$\Delta$   $\frac{L}{n}$  can be tiny.  $\frac{L}{n} = \frac{\text{condition number}}{\text{largest eig}}$   
 $\frac{L}{n} = \frac{1}{\text{smallest eig}}$   $t = \frac{L}{n} \frac{1}{2} \ln \frac{1}{\epsilon}$



$\frac{L}{n} = \frac{1}{\lambda}$   
 $(n=0)$

$$F(\theta_t) - F(\theta_{t+1}) = \frac{1}{2} (\theta_t - \theta_{t+1})^T H (\theta_t - \theta_{t+1})$$

equals  $\sum \lambda_i (1 - \lambda_i t)^{2t}$

$$\leq \sum \lambda_i \exp(-2t\lambda_i t)$$

$$= \frac{2t \lambda_i t e^{-2t\lambda_i t}}{2t \lambda_i t}$$

$$\leq \frac{1}{2t\lambda} \left( \frac{1}{e} \right)$$

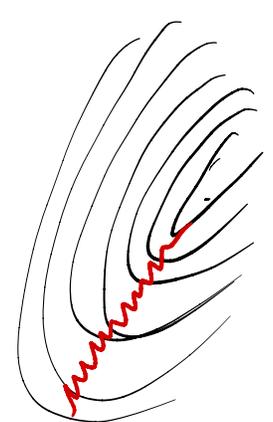


$$\leq \frac{1}{4t\lambda} \|\theta_t - \theta_{t+1}\|_2^2$$

$$\leq \frac{1}{8t\lambda} \|\theta_t - \theta_{t+1}\|_2^2$$

$\implies$   
 $e \geq 2$

$\rightarrow 0$  when  $t \rightarrow \infty$   
slow but almost unavoidable  
 $f_t$  has decreased



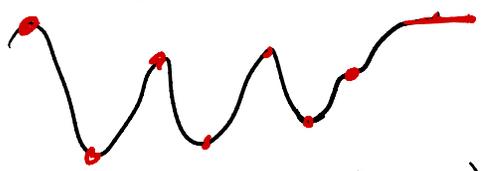
Convex functions:  $F: \mathbb{R}^d \rightarrow \mathbb{R}$  differentiable

$F$  convex  $\Leftrightarrow \forall \theta, \eta$ .  $F(\eta) \geq F(\theta) + F'(\theta)^T(\eta - \theta)$



if  $F$  twice differentiable  $\Leftrightarrow \forall \theta$ ,  $F''(\theta)$  Positive semi-definite

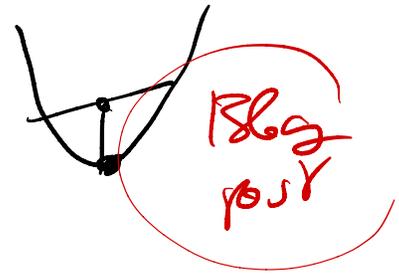
1st consequence:  $\forall \theta$ ,  $F''(\theta) \succeq 0$   
 if  $F'(\theta) = 0$  (stationary point) then  $F(\eta) \geq F(\theta) \forall \eta$ .  
 $\Rightarrow \theta$  is a global minimum.



2nd consequence:  $F(\eta) \geq F(\theta) + F'(\theta)^T(\eta - \theta) \Rightarrow F(\theta) - F(\eta) \leq F'(\theta)^T(\theta - \eta)$

Jensen's inequality:  $F$  convex,  $\mu$  prob. measure on  $\mathbb{R}^d$

$F\left(\int_{\mu} \theta\right) \leq \int_{\mu} F(\theta)$



Properties of convex fct: Boyd & Van der Berghe

When are ML objective functions convex?  $F(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(\eta_i; \theta)$

Sufficient condition:  
 (a)  $\ell$  is convex w.r.t to  $2^{nd}$  var  
 $\ell_0(\eta)$  is linear in  $\theta$

Strong convexity:  $\exists \mu > 0$   
 $\Leftrightarrow F(\eta) \geq F(\theta) + F'(\theta)^T(\eta - \theta) + \frac{\mu}{2} \|\theta - \eta\|_2^2$   
 $\Leftrightarrow F'(\theta)$  has eigenvalues  $\geq \mu$ ,  $\forall \theta$



(c'31)  $\rightarrow$  (c'54)

$F$  is strongly convex  $\Leftrightarrow F(y) \geq F(a) + F'(a)^T(y-a) + \frac{\mu}{2} \|y-a\|_2^2 \quad \forall y, a$

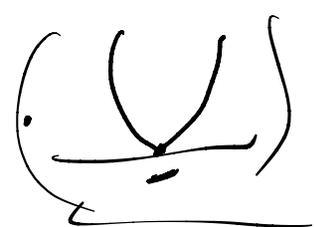
$$F(y) \geq F(a) - \frac{1}{2\mu} \|F'(a)\|_2^2 \quad (a=F'(a), y=a)$$

$$\mu \frac{1}{2} \|y-a\|_2^2 - \frac{1}{2} \|F'(a)\|_2^2 \geq 0$$

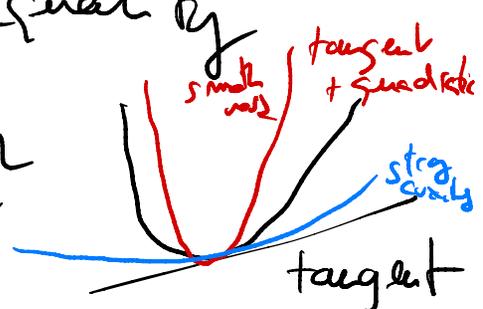
$$\mu \|y-a\|_2^2 \geq \|F'(a)\|_2^2$$

$$\mu \|y-a\|_2^2 \geq \|F'(a)\|_2^2 \Rightarrow \mu \|y-a\|_2 \geq \|F'(a)\|_2$$

$\Rightarrow \forall a, y: F(a) - F(y) \leq \frac{1}{2\mu} \|F'(a)\|_2^2$  | Lojasiewicz inequality



smoothness:  $F: \mathbb{R}^d \rightarrow \mathbb{R}$  is  $L$ -smooth  
 $\Leftrightarrow \forall y, a: |F(y) - F(a) - F'(a)^T(y-a)| \leq \frac{L}{2} \|y-a\|_2^2$



if twice differentiable  $\Rightarrow F'(a)$  has eigenvalues in  $(-L, L)$

Band  $\beta = GD$  (  $\mu$ -strong convexity,  $L$ -smoothness ) -  $q_t = q_{t-1} - \frac{1}{L} F'(q_{t-1})$  (GD)

$$F(q_t) \leq F(q_{t-1}) + F'(q_{t-1})^T(q_t - q_{t-1}) + \frac{L}{2} \|q_t - q_{t-1}\|_2^2 \quad (\text{smoothness})$$

$$= F(q_{t-1}) - \frac{1}{L} \|F'(q_{t-1})\|_2^2 + \frac{L}{2} \left\| \frac{1}{L} F'(q_{t-1}) \right\|_2^2 = F(q_{t-1}) - \frac{1}{2L} \|F'(q_{t-1})\|_2^2$$

$$\leq F(q_{t-1}) - \frac{\mu}{L} (F(q_{t-1}) - F(q_*)) \quad (\text{Lojasiewicz})$$

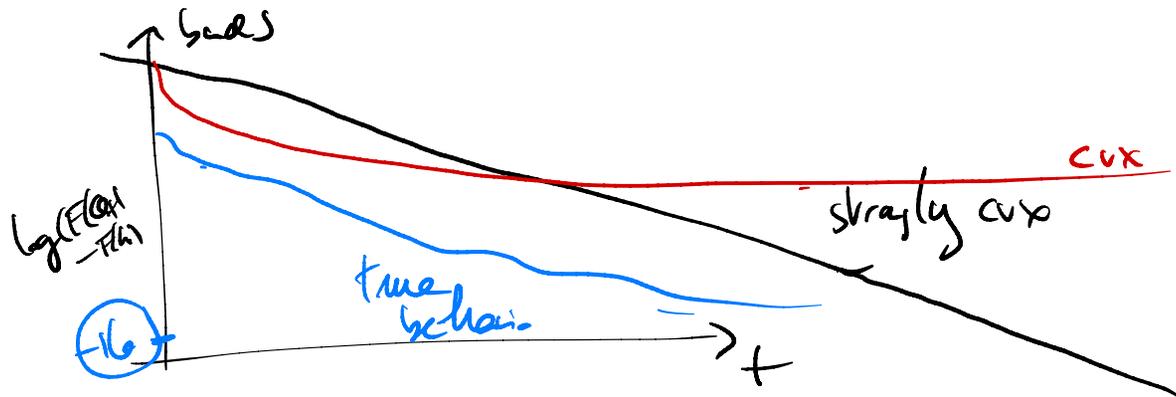
$$\Rightarrow F(q_t) - F(q_*) \leq \left(1 - \frac{\mu}{L}\right)^t (F(q_0) - F(q_*))$$

$\mu/L = \text{condition number}$

$\hookrightarrow$  exponential convergence linear

Band  $\beta = GD: F(q_t) - F(q_*) \leq \frac{L}{2t} \|q_0 - q_*\|_2^2$  (see back)

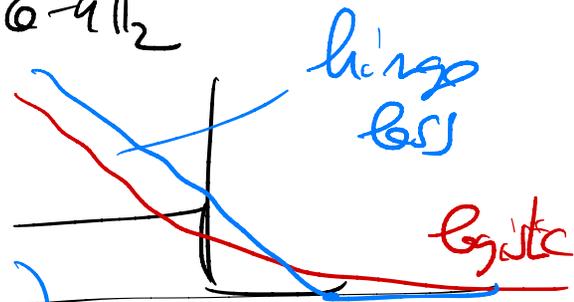
- ① Adaptivity
  - ② "optimal"
  - ③ Newton
  - ④ proximal
- $(1 - \sqrt{\frac{\mu}{L}})^t$   $\frac{1}{t^2}$   $\frac{1}{t^2}$   $\frac{1}{t^2}$   $\frac{1}{t^2}$



Non-smooth problems:  $F: \mathbb{R}^d \rightarrow \mathbb{R}$   $B$ -Lipschitz cont. means

$$\forall \theta, \eta, \quad |F(\theta) - F(\eta)| \leq B \|\theta - \eta\|_2$$

$$\Leftrightarrow \forall \theta, \quad \|F'(\theta)\|_2 \leq B$$



differentiable

(excluding quadratic fct)

See notion of subgradient in book

Studying GD:  $\theta_t = \theta_{t-1} - \gamma_t F'(\theta_{t-1})$  (GD) Looking for a "Lyapunov fct"  
assuming  $F$  has a minimizer  $q_*$

$$\begin{aligned} \|\theta_t - q_*\|_2^2 &= \|\theta_{t-1} - q_* - \gamma_t F'(\theta_{t-1})\|_2^2 \quad (\text{definition of GD}) \\ &= \|\theta_{t-1} - q_*\|_2^2 - 2\gamma_t F'(\theta_{t-1})^\top (\theta_{t-1} - q_*) + \gamma_t^2 \|F'(\theta_{t-1})\|_2^2 \\ &\geq \underbrace{F(\theta_{t-1}) - F(q_*)}_{\geq 0} \geq 0 \quad (\text{by convexity}) \end{aligned}$$

$\leq B^2$  (Lipschitz continuity)

$$\Rightarrow F(\theta_{t-1}) - F(q_*) \leq \frac{1}{2\gamma_t} \left[ \|\theta_{t-1} - q_*\|_2^2 - \|\theta_t - q_*\|_2^2 \right] + \frac{B^2}{2} \gamma_t$$

$$\frac{1}{T} \sum_{t=1}^T (F(\theta_{t-1}) - F(q_*)) \leq \frac{1}{2\gamma_T} \left[ \|\theta_0 - q_*\|_2^2 - \|\theta_T - q_*\|_2^2 \right] + \frac{B^2}{2} \gamma$$

if  $\gamma_t = \gamma$   
(Kleining sum)

(Jensen)

$$F\left(\frac{1}{T} \sum_{t=1}^T \theta_{t-1}\right) - F(q_*) \leq \frac{1}{2\gamma T} \|\theta_0 - q_*\|_2^2 + \frac{B^2}{2} \gamma$$

$$F\left(\frac{1}{T} \sum_{t=1}^T \varphi(x_t)\right) - F(\vartheta_0) \leq \frac{1}{2\sqrt{T}} \|\vartheta_0 - \vartheta_0\|_2^2 + \frac{R^2}{2} \delta$$

Does not go to zero if  $T \rightarrow \infty$ .

if  $\delta = \frac{\delta}{\sqrt{T}}$

$$\leq \frac{1}{\sqrt{T}} \left[ \frac{1}{2\delta} \|\vartheta_0 - \vartheta_0\|_2^2 + \frac{R^2}{2} \delta \right]$$

slow but "opt. mod"

- not an anytime algorithm.
- step-size depends on horizon  $T$

see back for  $\delta_t = \frac{\delta}{\sqrt{t}} \Rightarrow$  bound  $\propto \frac{\log t}{\sqrt{t}}$

Bound on  $F(\vartheta_{t+1}) - F(\vartheta_0)$ ? Blog post by T. Crabona.

Application to ML:  $F(\vartheta) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, \vartheta \varphi(x_i))$

Glinskij  
 $\|\varphi(x)\|_2$  bounded by  $R$

$$\|F'(\vartheta)\|_2 \leq GR$$

$$\text{opt. m. error} \leq \frac{1}{2\sqrt{T}} \left[ \frac{1}{\delta} \|\vartheta_0 - \vartheta_0\|_2^2 + \delta G^2 R^2 \right] = \frac{GR}{\sqrt{T}} \|\vartheta_0 - \vartheta_0\|_2$$

with proper choice of  $\delta$

last week: uniform deviation  $\approx \frac{GR \|\vartheta_0\|_2}{\sqrt{n}}$

$$= \frac{GR \|\vartheta_0\|_2}{\sqrt{T}}$$

$$T \approx n$$

SFD:  $F(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, \beta_\theta(x_i))$

$F'(\theta) = \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \theta} \ell(y_i, \beta_\theta(x_i))$

(SD):  $\theta_T = \theta_{T-1} - \gamma F'(\theta_{T-1})$

• SFD:  $\theta_T = \theta_{T-1} - \gamma \frac{\partial}{\partial \theta} \ell(y_{i(t)}, \beta_{\theta_{T-1}}(x_{i(t)}))$

multiple  
pass

$i(t)$  uniform in  $\{1, \dots, n\}$ .

$\Rightarrow E[g_{t+1} | \theta_{T-1}] = F'(\theta_{T-1})$

$E(F(\bar{\theta}_T) - F(\theta^*)) \leq \frac{D}{\sqrt{T}}$   
(wrt  $i(1), \dots, i(T)$ )



• single pass SFD:  $F(\theta) = E \ell(y, \beta_\theta(x))$

$\theta_T = \theta_{T-1} - \gamma \frac{\partial}{\partial \theta} \ell(y_{t^*}, \beta_{\theta_{T-1}}(x_{t^*}))$

Unbiased

$E[g_{t+1} | \theta_{T-1}] = F'(\theta_{T-1})$   $t^* \in \{1, \dots, n\}$