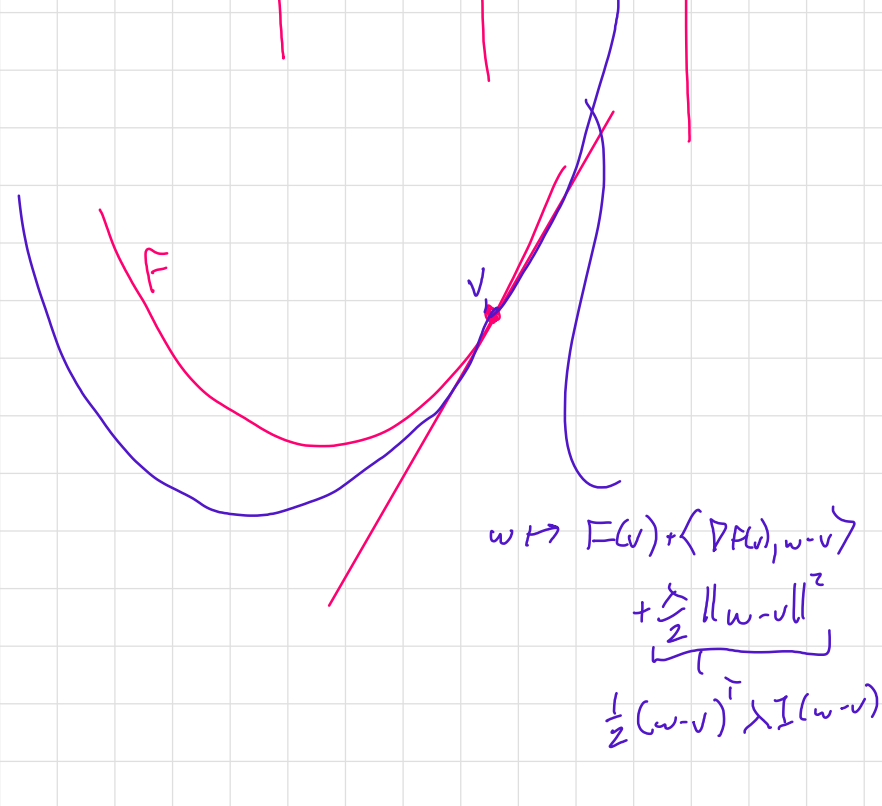


# Lec 12 end of semi-classical; NTK opt-4.1

- Ans:
- \* 15 line coding prob
  - \* "α" discussion :- NTK notes is bad.

	$\beta$ -smooth GD	exists GF	"Lyapunov" Potential
opt stationary points (somewhere in first t iterates)	$\frac{2\beta}{t} \cdot O(1)$	$\frac{1}{t}$	$F(w_t)$
minimize convex F (in fact: also minimize $\ w_t - \bar{w}\ ^2$ )	$\frac{\beta}{2t}$	$\frac{1}{t}$	$\frac{1}{2} \ w_t - \bar{w}\ ^2$
$\lambda$ -strongly-convex (a) minimize F (b) minimize $\frac{1}{2} \ w_t - \bar{w}\ ^2$	$\left. \begin{matrix} \text{exp}(-\frac{\lambda}{2} t) \\ \text{exp}(-\frac{\lambda}{2} t) \end{matrix} \right\}$	$\text{exp}(-\lambda t)$	$\left. \begin{matrix} \text{(a) } F(w_t) \\ \text{(b) } \frac{1}{2} \ w_t - \bar{w}\ ^2 \end{matrix} \right\}$



## Remark ("units" of GD & GF)

"arc length" of GD & GF paths?

assume  $\|\nabla F(w)\| \approx 1$

arc length for GD:  $\sum_{i=0}^{t-1} \|\nabla F(w_i)\| \approx \frac{t}{\beta}$

for GF:  $\int_0^t \|\nabla F(w(s))\| ds \approx t$

Theorem. F is  $\lambda$ -sc, then  $\bar{w}$  w/ta  $F(\bar{w}) = \inf_v F(v)$ , & if GF exists, then

(a)  $F(w(t)) - F(\bar{w}) \leq (F(w(0)) - F(\bar{w})) \exp(-2\lambda t)$ ,

(b)  $\frac{1}{2} \|w(t) - \bar{w}\|^2 \leq \frac{1}{2} \|w(0) - \bar{w}\|^2 \exp(-2\lambda t)$ .

Proof. not doing (a); for (b):

$$\begin{aligned} \frac{1}{2} \|w(t) - \bar{w}\|^2 - \frac{1}{2} \|w(0) - \bar{w}\|^2 &= \int_0^t \frac{d}{ds} \frac{1}{2} \|w(s) - \bar{w}\|^2 ds \\ &= \int_0^t \langle w(s) - \bar{w}, \dot{w}(s) \rangle ds \\ &= \int_0^t \langle w(s) - \bar{w}, -\nabla F(w(s)) \rangle ds \\ &= - \int_0^t \underbrace{\langle w(s) - \bar{w}, \nabla F(w(s)) - \nabla F(\bar{w}) \rangle}_{\geq 0 \text{ cause}} ds \quad \underbrace{\nabla F(\bar{w}) = 0}_{\geq \lambda \|w(s) - \bar{w}\|^2} \\ &\leq - \int_0^t \frac{2\lambda}{2} \|w(s) - \bar{w}\|^2 ds \\ &\leq \dots? \end{aligned}$$

define  $u(s) := \frac{1}{2} \|w(s) - \bar{w}\|^2$

$\frac{d}{ds} u(s) \leq -2\lambda u(s)$

By Grönwall's inequality  $u(t) \leq u(0) \exp(\int_0^t \theta(s) ds)$

$= u(0) \exp(-\int_0^t 2\lambda ds)$

$= u(0) \exp(-2\lambda t)$ .

## Remark (on SGD).

\* With  $\beta$ -smooth,  $\text{had term} + \eta^2 \|\nabla F(w)\|^2$ ,

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

\* For stochastic, in most theory papers,

use  $\eta = \frac{1}{\sqrt{t}}$ , thus  $\eta^2 = \frac{1}{t}$ , which eats the  $\frac{1}{\eta^2} \|\nabla F(w)\|^2$ .

\* Theory does capture that stochastic gradients are cheap but doesn't capture statistical benefits. } active area

# Standard NTK proof (via strong convexity).

## Remarks (on approach)

- \* This proof style is most common (> 99%).
- \* Strong convexity  $\Rightarrow$  what is  $\lambda$ ?
- $\lambda$  depends on strong convexity of loss and  $\lambda_{\min}(\text{kernel gram matrix})$

$$= \lambda_{\min} \left( \begin{bmatrix} \text{[ij]} & \dots \\ \vdots & \ddots \end{bmatrix} \left\langle \nabla f(x_i; w_0), \nabla f(x_j; w_0) \right\rangle \right)$$

(recall  $\frac{1}{2} \|Xw - y\|^2 \rightarrow \sigma_{\min}(X)$ .)

- \* Proof plan: inductively maintain step near initialization  $\Rightarrow$  proof is like linear case
- \* rate will be  $\exp(-\lambda t)$
- \* Proof here is loosely based on Chizat-Bach ("lazy training...")

## Notation Take training set into product

$$f: \mathbb{R}^p \rightarrow \mathbb{R}^n$$

$$f(w) = \begin{bmatrix} f(x_1; w) \\ \vdots \\ f(x_n; w) \end{bmatrix} \in \mathbb{R}^n$$

~~Loss~~  $L(v) = \frac{1}{2} \|v - y\|^2$

"unnormalized empirical risk" for us  $L(\alpha f(w))$

$$= \frac{1}{2} \|\alpha f(w) - y\|^2$$

captures "zooming" whereby NTK with more parameters becomes linear.

## Gradient

$$\dot{w}(t) = -\nabla_w L(\alpha f(w(t))) = - \begin{bmatrix} \nabla f(x_1; w(t))^T \\ \vdots \\ \nabla f(x_n; w(t))^T \end{bmatrix} \alpha (\alpha f(w) - y)$$

$\nabla L(\alpha f(w))$

"fugent flow"

$$w(0) = w_0$$

$$\dot{w}(t) = -\alpha J_0^T \nabla L(\alpha f(w(t)))$$

interpret as "data using features at time  $t$ "

one point: we'll show  $\|w(t) - w(t)\|$  small to control strong convexity.

\*  $\text{rank}(\underbrace{J_0 J_0^T}_{\text{kernel gram matrix}}) = n$  ;  $\sigma_n(J_0) > 0$

\* Smoothness:  $\|J_w - J_v\| \leq \beta \|w - v\|$

$\Rightarrow$  if  $\sigma_n(J_t J_t^T) \approx \sigma_n(J_0 J_0^T)$  if  $\|w_t - w_0\|$ .

where  $J_v^T = \begin{bmatrix} \nabla f(x_1; v)^T \\ \vdots \\ \nabla f(x_n; v)^T \end{bmatrix}^T$

# Office hours

we used  $\nabla f(\bar{w}) = 0$

optimality condition for unconstrained

constrained:  $-\nabla f(\bar{w}) = N_S(\bar{w})$

