

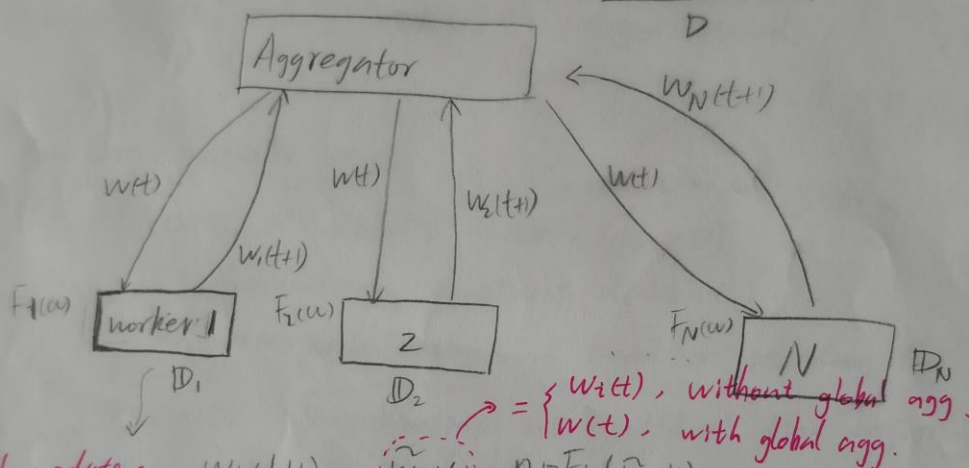
①

# ECE 6960/5960 Lecture #27

- Review:
- ① Introduction to Federated Learning
  - ② Design and detailed implementation

Today: ① theory part and proof

System Setup: (Distributed Gradient Descent, DGD for short)  
 global aggregation:  $w(t+1) = \frac{\sum_{i=1}^N D_i w_i(t+1)}{D}$



local update:  $w_i(t+1) = \tilde{w}_i(t) - \eta \nabla F_i(\tilde{w}_i(t))$

Local loss:  $F_i(w) \triangleq \frac{1}{|D_i|} \sum_{j \in D_i} f_j(w)$  loss on sample  $x_j$   
 $f_j(w) \triangleq f(w, x_j, y_j)$   
 Local dataset for worker  $i$ .

overall loss:  $F(w) \triangleq \frac{\sum_{i=1}^N D_i F_i(w)}{D}$

$D_i = |D_i|$ ,  $D = \sum_{i=1}^N D_i$

objective:  $w^* = \underset{w}{\operatorname{argmin}} F(w)$

$K$ : total # of global aggregations

$\tau$ : # of local updates between two global agg.

$T = K\tau$ : # of updates.

$c$ : each step of local update at all workers consumes  $c$  units of resource <sup>(2)</sup>

$b$ : each global update consumes  $b$  units of resource

Goal  $\min_{\tau, T} F(w(T))$  resource budget

s.t.  $Tc + Kb = Tc + \frac{T}{\tau}b = T(c + \frac{b}{\tau}) \leq \overset{\uparrow}{R}$

Period  $[k] = [(k-1)\tau, k\tau]$  integer interval  $(\tau+1)$  integers  
 $\{(k-1)\tau, (k-1)\tau+1, \dots, k\tau\}$

Auxiliary parameter vector:

—  $V_{[k]}(t) = V_{[k]}(t-1) - \eta \nabla F(V_{[k]}(t-1))$   $\nearrow$  full gradient

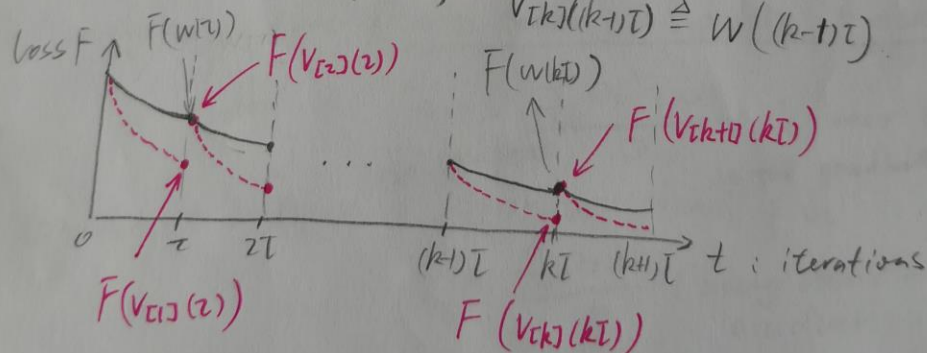
(Centralized gradient descent)

$\Rightarrow$  convergence analysis is based on bounding the difference

—  $V_{[k]}(t)$  is synchronized with  $w(t)$  at the beginning of each interval  $[k]$ , i.e.,

averaged global update

$V_{[k]}((k-1)\tau) \triangleq w((k-1)\tau)$



Remark:  $V_{[k]}(t\tau) \neq V_{[k+1]}(t\tau)$ ,  $t = 1, 2, \dots, T$ .

Assumptions: 1)  $F_i(w)$  is CVX. (individual loss is convex).

2)  $F_i(w)$  is  $\rho$ -Lipschitz, i.e.,

absolute value  $\hookrightarrow \|F_i(w) - F_i(w')\| \leq \rho \|w - w'\|, \forall w, w'$   
(Function value cannot change too much between two points  $w$  and  $w'$ )

3)  $F_i(w)$  is  $\beta$ -smooth, i.e.,

$$\|\nabla F_i(w) - \nabla F_i(w')\| \leq \beta \|w - w'\|, \forall w, w'$$

$\Rightarrow$  two properties ① ②

$$\textcircled{1} f(w) - f(w') \leq \nabla f(w')^T (w - w') - \frac{1}{2\beta} \|\nabla f(w) - \nabla f(w')\|^2$$

$$\textcircled{2} |f(w) - f(w') - \nabla f(w')^T (w - w')| \leq \frac{\beta}{2} \|w - w'\|^2$$

(plus convexity:  $f(w) \geq f(w') + \nabla f(w')^T (w - w')$ )

$$f(w) - f(w') - \nabla f(w')^T (w - w') \leq \frac{\beta}{2} \|w - w'\|^2$$

Linear approximation of  $f(w)$  at  $w'$ .

Remarks: ① pretty common assumptions, strong convexity et.

It's easy to see that  $F(w)$  is convex,  $\rho$ -Lipschitz,  $\beta$ -smooth

due to  $F(w) = \frac{\sum_{i=1}^N D_i F_i(w)}{D}$  (linear combination)

[def] Gradient Divergence:  $\|\nabla F_i(w) - \nabla F(w)\| \leq \delta_i$   
local gradient  $\nabla F_i(w)$ , full gradient  $\nabla F(w)$

$$\delta \triangleq \frac{\sum_{i=1}^N D_i \delta_i}{D} \quad (\text{average deviation of local gradients})$$

GOAL: Bound  $F(w(T)) - F(w^*)$

analyze the convergence behavior of  $F(V_k(t))$

bound  $F(w(T)) - F(V_k(T))$

bound by  $\|w(T) - V_k(T)\|$ .

Remark: compare the difference between the Distributed GD (DGD) with the centralized version  $V[k](t)$ . ①

Theorem 1: For any interval  $[k]$  and  $t \in [k]$ , we have

$$\|w(t) - \overset{\text{centralized ver.}}{V[k](t)}\| \leq h(t - (k-1)\tau)$$

where  $h(x) \triangleq \frac{\delta}{\beta} [(\eta\beta + 1)^x - 1] - \eta\delta x$ .

averaged param.  
from global agg.

Furthermore, as  $F(\cdot)$  is  $\rho$ -Lipschitz, we have

$$\begin{aligned} F(w(t)) - F(V[k](t)) &\leq \rho \|w(t) - V[k](t)\| \\ F(w) - F(v) &\leq |F(w) - F(V[k](t))| \leq \rho \|w - v\| \\ &\leq \rho h(\cdot) \leq \rho h(t - (k-1)\tau) \end{aligned}$$

Remarks:  $h(0) = h(1) = 1$ .

$\Rightarrow w(t) - V[k](t) = 0$ , if  $\tau = 1$ , that is, when  $\tau = 1$ , we perform global agg. after each local update, then DGD is equivalent to centralized GD.

[proof] Need to prove  $w(t) = w(t-1) - \eta \nabla F(w(t-1))$

$$\tau = 1 \Rightarrow \tilde{w}_i(t) = w(t)$$

$$w(t) = \frac{\sum_{i=1}^N D_i w_i(t)}{D} \leadsto w_i(t) = \tilde{w}_i(t-1) - \eta \nabla F_i(\tilde{w}_i(t-1))$$

$$= \frac{\sum_i D_i [\tilde{w}_i(t-1) - \eta \nabla F_i(\tilde{w}_i(t-1))]}{D}$$

$$= \frac{\sum_i D_i w_i(t-1)}{D} - \eta \cdot \frac{\sum_i D_i \nabla F_i(w(t-1))}{D}$$

$$= w(t-1) - \eta \nabla F(w(t-1)) \quad \left. \begin{aligned} &\eta \nabla \left( \frac{\sum_i D_i F_i(w(t-1))}{D} \right) \\ &= \eta \nabla F(w(t-1)) \end{aligned} \right\}$$

(linearity of gradient operator).

□



⑤

② When  $x$  is large,  $h(x) \approx \frac{\delta}{\beta} (\eta\beta + 1)^x$  : exponential w.r.t.  $x$

proof of Thm 1.

local para. at iteration  $t$

$$\tilde{w}_i(t) = \begin{cases} w(t), & \text{global agg.} \\ w_i(t) \end{cases}$$

[Lemma 2] For any interval  $[k]$ , and  $t \in [(k+1)\tau, k\tau]$ , we have

$$\|\tilde{w}_i(t) - v_{[k]}(t)\| \leq g_i(t - (k-1)\tau)$$

$$\text{where } g_i(x) \triangleq \frac{\delta_i}{\beta} [(\eta\beta + 1)^x - 1]$$

(Using  $\beta$ -smoothness, triangle inequality, induction)

[proof] By induction on  $t \in [k]$ .

The local para. is not too far from the centralized para, due to global aggregation and synchronization of  $v_{[k]}(t)$ .  
 $\Rightarrow$  the average of local paras. ( $w(t)$ ) should also be not too far from  $v_{[k]}(t)$ .

First, when  $t = (k+1)\tau$ , we know that  $\tilde{w}_i(t) = v_{[k]}(t)$ .  
 That is,  $\|w_i(t) - v_{[k]}(t)\| = g_i(0) = 0$  when  $t = (k+1)\tau$ .

holds for the initial value

For the induction, we assume that

$$\|\tilde{w}_i(t-1) - v_{[k]}(t-1)\| \leq g_i(t-1 - (k-1)\tau) \text{ hold for some } t$$

$$\text{We have } \|\tilde{w}_i(t) - v_{[k]}(t)\| \text{ for } t \in ((k+1)\tau, k\tau]$$

$$= \left\| \left( \tilde{w}_i(t-1) - \eta \nabla F_i(\tilde{w}_i(t-1)) \right) - \left( v_{[k]}(t-1) - \eta \nabla F(v_{[k]}(t-1)) \right) \right\|$$

$$= \left\| \left( \tilde{w}_i(t-1) - v_{[k]}(t-1) \right) - \eta \left( \nabla F_i(\tilde{w}_i(t-1)) - \nabla F_i(v_{[k]}(t-1)) + \nabla F_i(v_{[k]}(t-1)) - \nabla F(v_{[k]}(t-1)) \right) \right\|$$

Zero trick

$$= \left\| \left( \tilde{w}_i(t-1) - v_{[k]}(t-1) \right) - \eta \left( \nabla F_i(\tilde{w}_i(t-1)) - \nabla F_i(v_{[k]}(t-1)) \right) \right\|$$

①

②

$$- \eta \left( \nabla F_i(v_{[k]}(t-1)) - \nabla F(v_{[k]}(t-1)) \right)$$

③

$$\leq \|\text{①}\| + \|\text{②}\| + \|\text{③}\|$$

triangle inequality

$$= \|\tilde{w}_i(t-1) - v_{[k]}(t-1)\| + \eta \left\| \nabla F_i(\tilde{w}_i(t-1)) - \nabla F_i(v_{[k]}(t-1)) \right\|$$

$$\leq \eta\beta \|\tilde{w}_i(t-1) - v_{[k]}(t-1)\|$$

$$\eta \| \nabla F_i(v_{[k]}(t-1)) - \nabla F(v_{[k]}(t-1)) \| \leq \delta_i \quad \text{gradient divergence} \quad (6)$$

$$\leq (\eta\beta + 1) \underbrace{\| \tilde{w}_i(t-1) - v_{[k]}(t-1) \|}_{\beta\text{-smoothness}} + \underbrace{\eta\delta_i}_{\text{individual gradient divergence}}$$

$$\begin{aligned} & \leq (\eta\beta + 1) \overbrace{g_i(t-1-(k-1)\tau)}^{\text{inductive}} + \eta\delta_i \\ & \stackrel{\text{induction assumption at } t-1}{=} (\eta\beta + 1) \left( \frac{\delta_i}{\rho} [(\eta\beta + 1)^{t-1-(k-1)\tau} - 1] \right) + \eta\delta_i \\ & \quad \dots \text{math manipulation} \\ & = g_i(t-(k-1)\tau) \end{aligned}$$

$$\Rightarrow \| \tilde{w}_i(t) - v_{[k]}(t) \| \leq g_i(t-(k-1)\tau) \quad \forall t \in [k]. \quad \square$$

proof of Thm 1 continued:

Since  $w_i(t) = \tilde{w}_i(t-1) - \eta \nabla F_i(\tilde{w}_i(t-1))$

$$w(t) = \frac{\sum_{i=1}^N D_i w_i(t)}{D}, \text{ then}$$

$$\begin{aligned} w(t) &= \frac{\sum_i D_i \tilde{w}_i(t-1)}{D} - \eta \frac{\sum_i D_i \nabla F_i(\tilde{w}_i(t-1))}{D} \\ &= w(t-1) - \eta \frac{\sum_i D_i \nabla F_i(\tilde{w}_i(t-1))}{D} \end{aligned}$$

then for  $t \in ((k-1)\tau, k\tau]$ ,

$$\begin{aligned} \| w(t) - v_{[k]}(t) \| &= \left\| w(t-1) - \eta \frac{\sum_i D_i \nabla F_i(\tilde{w}_i(t-1))}{D} \right. \\ &\quad \left. - v_{[k]}(t-1) + \eta \nabla F(v_{[k]}(t-1)) \right\| \end{aligned}$$

rearrange the terms

$$= \left\| w(t-1) - v_{[k]}(t-1) - \eta \left( \frac{\sum_i D_i \nabla F_i(\tilde{w}_i(t-1))}{D} - \nabla F(v_{[k]}(t-1)) \right) \right\|$$

$$= \left\| w(t) - V_{[k]}(t-1) - \eta \frac{\sum_i D_i (\nabla F_i(\tilde{w}_i(t-1)) - \nabla F_i(V_{[k]}(t-1)))}{D} \right\| \quad (7)$$

triangle inequality

$$\leq \|w(t-1) - V_{[k]}(t-1)\| + \eta \cdot \left( \frac{\sum_i D_i (\|\nabla F_i(\tilde{w}_i(t-1)) - \nabla F_i(V_{[k]}(t-1))\|)}{D} \right)$$

$$\leq \underbrace{\|w(t-1) - V_{[k]}(t-1)\|}_{\text{①}} + \eta \beta \left( \frac{\sum_i D_i \|\tilde{w}_i(t-1) - V_{[k]}(t-1)\|}{D} \right)$$

$\beta$ -smooth

$$\leq \text{①} + \eta \beta \frac{\sum_i D_i g_i(t-1-k-1)\tau}{D}$$

lemma 2

$$= \text{①} + \eta \beta \frac{\sum_i D_i \frac{\delta_i}{\beta} [(\eta\beta + 1)^{t-1-(k-1)\tau} - 1]}{D}$$

$$\uparrow = \text{①} + \eta \left( \frac{\sum_i D_i \delta_i}{D} \right) [(\eta\beta + 1)^{t-1-(k-1)\tau} - 1]$$

def  $\delta = \frac{\sum_i D_i \delta_i}{D} \triangleq \delta$

$$= \text{①} + \eta \delta [(\eta\beta + 1)^{t-1-(k-1)\tau} - 1]$$

$$\Rightarrow \|w(t) - V_{[k]}(t)\| - \|w(t-1) - V_{[k]}(t-1)\| \leq \eta \delta [(\eta\beta + 1)^{t-1-(k-1)\tau} - 1] \quad (*)$$

If  $t = (k+1)\tau$ , we have  $w(t) = V_{[k]}(t)$  (synchronization)

$$\Rightarrow \|w(t) - V_{[k]}(t)\| = 0.$$

If  $t \in ((k+1)\tau, k\tau]$ , then by summing (\*) from  $(k+1)\tau+1$ ,  $(k+1)\tau+2, \dots, t$ , then we have

$$\|w(t) - V_{[k]}(t)\| = \sum_{y=(k+1)\tau+1}^t \|w(y) - V_{[k]}(y)\| - \|w(y-1) - V_{[k]}(y-1)\|$$

$$\leq \eta \delta \sum_{y=(k+1)\tau+1}^t [(\eta\beta + 1)^{y-1-(k-1)\tau} - 1]$$

↓ geometric sum



(8)

$$= \eta \delta \sum_{z=1}^{t-(k-1)\tau} [(\eta\beta+1)^{z-1} - 1]$$

$$\text{let } z = t - (k-1)\tau$$

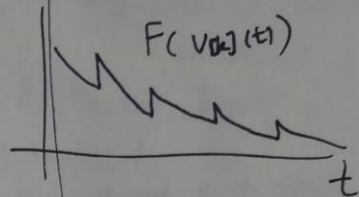
$$= \eta \delta \sum_{z=1}^{t-(k-1)\tau} (\eta\beta+1)^{z-1} - \eta \delta (t - (k-1)\tau)$$

$$\text{geometric sum} = \eta \delta \frac{1 - (\eta\beta+1)^{t-(k-1)\tau}}{-\eta\beta} - \eta \delta (t - (k-1)\tau)$$

$$= \frac{\delta}{\beta} [(\eta\beta+1)^{t-(k-1)\tau} - 1] - \eta \delta (t - (k-1)\tau)$$

$$= h(t - (k-1)\tau)$$

$$\text{recall } h(x) \triangleq \frac{\delta}{\beta} [(\eta\beta+1)^x - 1] - \eta \delta x. \quad \square$$



Theorem 2: When all the following conditions are satisfied,

1)  $\eta \leq \frac{1}{\beta}$  (small enough step size).

2)  $\omega \eta (1 - \frac{\beta \eta}{2}) - \frac{\rho h(z)}{2 \varepsilon^2} > 0$

3)  $F(w_k(t)) - F(w^*) \geq \varepsilon, \forall t, k$

4)  $F(w(t)) - F(w^*) \geq \varepsilon.$  } weired.

Where  $\varepsilon > 0$ , and  $w \triangleq \min_k \frac{1}{\|w_k((k-1)\tau) - w^*\|^2}$

Then the convergence upper bound is given by

$$F(w(t)) - F(w^*) \leq \frac{1}{T \left( \omega \eta (1 - \frac{\beta \eta}{2}) - \frac{\rho h(z)}{2 \varepsilon^2} \right)}$$

Remarks 1) Unreasonable assumptions 3) & 4)

smaller  $\varepsilon$ , loose bound on convergence.  
(larger)



roof of Thm 2.

(9)

Denote  $\Theta_{[k]}(t) = F(v_{[k]}(t)) - F(w^*)$ ,  $(k-1)T \leq t \leq kT$

Assume  $\Theta_{[k]}(t) > 0$ . (I didn't see why  $\Theta_{[k]}(t)$  cannot be zero).

[Lemma 3]. When  $\eta \leq \frac{1}{\beta}$ ,  $\forall k, t \in [k] \triangleq [(k-1)T, kT]$ , we have that  $\|v_{[k]}(t) - w^*\|$  is non-increasing w.r.t.  $t$ .

[proof]  $\|v_{[k]}(t+1) - w^*\|^2 \rightarrow$  getting closer to the optimal para. since  $v_{[k]}(t)$  is the full GD. (intuition) Each step will lead the para.  $v_{[k]}(t)$  closer to its optimal.

$$\begin{aligned} &= \|v_{[k]}(t) - \eta \nabla F(v_{[k]}(t)) - w^*\|^2 \\ &= \|v_{[k]}(t) - w^*\|^2 - 2\eta \nabla F(v_{[k]}(t))^T (v_{[k]}(t) - w^*) + \eta^2 \|\nabla F(v_{[k]}(t))\|^2 \end{aligned}$$

Since  $F(\cdot)$  is  $\beta$ -smooth, we have

$$\begin{aligned} 0 < \Theta_{[k]}(t) = F(v_{[k]}(t)) - F(w^*) &\leq \nabla F(v_{[k]}(t))^T (v_{[k]}(t) - w^*) - \frac{\|\nabla F(v_{[k]}(t))\|^2}{2\beta} \\ &\left( \text{due to } f(x) - f(y) \leq \nabla f(x)^T (x - y) - \frac{1}{2\beta} \|\nabla f(x) - \nabla f(y)\|^2 \right. \\ &\quad \left. \text{and } \nabla F(w^*) = 0 \right) \end{aligned}$$

$$\Rightarrow -\nabla F(v_{[k]}(t))^T (v_{[k]}(t) - w^*) < -\frac{\|\nabla F(v_{[k]}(t))\|^2}{2\beta}$$

Therefore,

$$\begin{aligned} \|v_{[k]}(t+1) - w^*\|^2 &= \|v_{[k]}(t) - \eta \nabla F(v_{[k]}(t)) - w^*\|^2 \\ &= \underbrace{\|v_{[k]}(t) - w^*\|^2}_{\textcircled{1}} - 2\eta \nabla F(v_{[k]}(t))^T (v_{[k]}(t) - w^*) + \eta^2 \underbrace{\|\nabla F(v_{[k]}(t))\|^2}_{\textcircled{3}} \\ &< \textcircled{1} - \eta \frac{\|\nabla F(v_{[k]}(t))\|^2}{\beta} + \textcircled{3} \end{aligned}$$

$$= \textcircled{1} + \eta \left( \frac{1}{\beta} - \eta \right) \|\nabla F(v_{[k]}(t))\|^2$$

as long as  $\frac{1}{\beta} - \eta > 0$  ( $\eta \leq \frac{1}{\beta}$ ), we have

$$\|v_{[k]}(t+1) - w^*\|^2 < \|v_{[k]}(t) - w^*\|^2$$

□

(10)

positive term

$$\Theta_{ck}(t) \triangleq F(v_{ck}(t)) - F(w^*)$$

$\Rightarrow F_{\text{rel}}(t) - F_{\text{rel}}(t^*) < F_{\text{rel}}(t) - F_{\text{rel}}(t^*)$ , i.e.,  $F_{\text{rel}}(t)$  is decreasing

proof of ~~Lemma~~ Thm 2 continued. (real proof)

Using Lemma 4 and  $t \in [(k-1)T, kT]$ , then

$$\frac{1}{\Theta[k](kL)} - \frac{1}{\Theta[k](k+1)L} = \sum_{t=(k-1)L}^{kL-1} \left( \frac{1}{\Theta[k](t+1)} - \frac{1}{\Theta[k](t)} \right)$$

$$| \text{Lemma 4.} \rightarrow \geq 2\omega\eta (1 - \frac{\beta\eta}{2})$$

Summing up from  $k=1$  to  $K$ , we obtain

$$\sum_{k=1}^K \left( \frac{1}{\Theta[k](kT)} - \frac{1}{\Theta[k](k-1)T} \right) \geq \sum_{k=1}^K \tau \omega \eta \left( 1 - \frac{\beta \eta}{2} \right) = K \tau \omega \eta \left( 1 - \frac{\beta \eta}{2} \right)$$

$$\Leftrightarrow \frac{1}{\Theta[k](T)} - \frac{1}{\Theta[k](0)} \geq \sum_{k=1}^{K-1} \underbrace{\left( \frac{1}{\Theta[k+1](kT)} - \frac{1}{\Theta[k](kT)} \right)}_{T = kT} + Kz\omega_g \left( \frac{\beta\eta}{2} \right) \quad (\star\star)$$

each term here can be bounded from below as:

$$\frac{1}{\Theta[k+1](kT)} - \frac{1}{\Theta[k](kT)} = \frac{\Theta[k](kT) - \Theta[k+1](kT)}{\Theta[k+1](kT) \Theta[k](kT)} \quad (11)$$

$$\begin{aligned} \Theta[k](t) = F(v[k](t)) - F(w^*) &= \frac{F(v[k](kT)) - F(v[k+1](kT))}{\text{denominator}} \\ &= \frac{F(v[k](kT)) - F(w(kT))}{\text{denominator}} \end{aligned} \quad \begin{array}{l} \text{equal by} \\ \text{definition} \end{array}$$

$$\stackrel{\rho\text{-Lipschitz}}{\geq} \frac{-\rho \|w(kT) - v[k](kT)\|}{\text{denom}}$$

$$\stackrel{\text{Thm 1.}}{\geq} \frac{-\rho h(kT - (k-1)T)}{\Theta[k](kT) \Theta[k+1](kT)} = \frac{-\rho h(t)}{\Theta[k](t) \Theta[k+1](t)}$$

Since it is assumed that  $\Theta[k](t) = F(v[k](t)) - F(w^*) \geq \varepsilon \quad \forall k, t$ ,

$$\Theta[k](kT) \Theta[k+1](kT) \geq \varepsilon^2$$

$$\frac{1}{\Theta[k](\cdot) \Theta[k+1](\cdot)} \geq -\frac{1}{\varepsilon^2} \quad (\star\star\star)$$

Go back to  $(\star\star)$ , we get

$$\sum_{k=1}^{K-1} \left( \frac{1}{\Theta[k+1](kT)} - \frac{1}{\Theta[k](kT)} \right) \geq \sum_{k=1}^{K-1} \frac{\rho h(t)}{\varepsilon^2} = -(K-1) \frac{\rho h(t)}{\varepsilon^2}$$

$$\text{and } \left\{ \frac{1}{\Theta[K](T)} - \frac{1}{\Theta[1](0)} \geq T w g \left(1 - \frac{\beta g}{2}\right) - (K-1) \frac{\rho h(t)}{\varepsilon^2} \right\} \quad (I)$$



Similarly, it's assumed that  $F(w(T)) - F(w^*) \geq \varepsilon$ , then

(12)

$$\frac{1}{(F(w(T)) - F(w^*)) \Theta[k](T)} \geq -\frac{1}{\varepsilon^2} \quad (\star\star\star)$$

assuming  $\Theta[k](t) = F(v[k](t)) - F(w^*) \geq \varepsilon$

Then,

$$\frac{1}{F(w(T)) - F(w^*)} = \frac{1}{\Theta[k](T)}$$

$$= \frac{\Theta[k](T) - (F(w(T)) - F(w^*))}{\text{denominator}}$$

$$= \frac{F(v[k](T)) - F(w(T))}{\text{denom}}$$

$$\text{Lipschitz} \geq \frac{-\rho \|v[k](T) - w(T)\|}{\text{denom}}$$

$$\geq \frac{-\rho h (T - (k-1)\tau)}{\text{denom}} \leftarrow (Thm), \|w(t) - v[k](t)\| \leq h(t - (k-1)\tau)$$

$$= \frac{-\rho h(z)}{\varepsilon^2} \quad (\text{II})$$

Summing up (I), (II), we have

$$\begin{aligned} \frac{1}{F(w(T)) - F(w^*)} - \frac{1}{\Theta[k](T)} &\geq \left( T\omega\eta(1 - \frac{\beta\eta}{2}) - (k-1)\frac{\rho h(z)}{\varepsilon^2} \right) - \frac{\rho h(z)}{\varepsilon^2} \\ &\geq 0 \\ &= T\omega\eta(1 - \frac{\beta\eta}{2}) - \frac{T\rho h(z)}{2\varepsilon^2} \\ &= T \left( \omega\eta(1 - \frac{\beta\eta}{2}) - \frac{\rho h(z)}{2\varepsilon^2} \right) \end{aligned}$$

Then  $\frac{1}{F(w(T)) - F(w^*)} \geq \frac{1}{F(w(t)) - F(w^*)} - \frac{1}{\Theta(t)} \quad (13)$

$$\nearrow \geq T \left( \omega \eta \left( 1 - \frac{\beta \eta}{2} \right) - \frac{\rho h(\eta)}{2 \xi^2} \right) > 0$$

because  $\Theta(t) = F(V(t)) - F(w^*) > 0$ .

If we further assume  $\omega \eta \left( 1 - \frac{\beta \eta}{2} \right) - \frac{\rho h(\eta)}{2 \xi^2} > 0$ ,  
then

$$F(w(T)) - F(w^*) \leq \frac{1}{T \left( \omega \eta \left( 1 - \frac{\beta \eta}{2} \right) - \frac{\rho h(\eta)}{2 \xi^2} \right)}.$$

Remarks:  $\omega \triangleq \min_k \frac{1}{\|V(k) - w^*\|^2}$  is used in Lemma 4,  
which we omitted the proof.  $\square$