

Assumption 1 (β -smoothness): For any parameter W_1, W_2 , each loss function F_i satisfies $F_i(W_1) \leq F_i(W_2) + (W_1 - W_2)^T \nabla F_i(W_2) + \beta/2 \|W_1 - W_2\|_2^2$. It also implies that $\|\nabla F_i(W_1) - \nabla F_i(W_2)\| \leq \beta(W_1 - W_2)$.

Assumption 2 (μ -strongly convex): For any parameter W_1, W_2 , each loss function F_i satisfies $F_i(W_1) \geq F_i(W_2) + (W_1 - W_2)^T \nabla F_i(W_2) + \mu/2 \|W_1 - W_2\|_2^2$. It also implies that $\|\nabla F_i(W_1) - \nabla F_i(W_2)\| \geq \mu(W_1 - W_2)$.

Assumption 3 (Bounded gradient): The stochastic gradient is bounded as $\mathbb{E}(\|\nabla F_i(W_i^t)\|_2^2) \leq \sigma^2$.

Theorem 1: There exist two constants $A = 1 - 2\mu\eta, B = \eta\sigma^2/(2\mu)$. The gap between the optimal W^* and W_i^t follows: $\mathbb{E}(\|W_i^t - W^*\|_2^2) \leq A^t \|W_i^0 - W^*\|_2^2 + B$.

$$\begin{aligned}
\text{Proof: } & \|W_i^{t+1} - W^*\|_2^2 \\
&= \|W_i^{t+1} - W_i^t + W_i^t - W^*\|_2^2 \\
&= \|W_i^{t+1}(M_i = 1) - W_i^t(M_i = 1) + W_i^{t+1}(M_i = 0) - W_i^t(M_i = 0) + \|W_i^t - W^*\|_2^2 \\
&= \|W_i^{t+1}(M_i = 1) - W_i^t(M_i = 1) + \|W_i^t - W^*\|_2^2 \\
&= \|\eta \nabla F_i(W_i^t(M_i = 1)) + W_i^t - W^*\|_2^2 \\
&= \|\eta \nabla F_i(W_i^t(M_i = 1))\|_2^2 + \|W_i^t - W^*\|_2^2 - 2\eta \nabla F_i(W_i^t(M_i = 1))(W_i^t - W^*)
\end{aligned}$$

Then, we have

$$\begin{aligned}
& \mathbb{E}(\|W_i^{t+1} - W^*\|_2^2) \\
&= \eta^2 \mathbb{E}(\|\nabla F_i(W_i^t(M_i = 1))\|_2^2) + \|W_i^t - W^*\|_2^2 - 2\eta \nabla F_i(W_i^t(M_i = 1)) \mathbb{E}(\|W_i^t - W^*\|_2) \\
&\leq \eta^2 \mathbb{E}(\|\nabla F_i(W_i^t(M_i = 1))\|_2^2) + \|W_i^t - W^*\|_2^2 - 2\eta (\nabla F_i(W_i^t(M_i = 1)) - \nabla F_i(W^*(M_i = 1))) \mathbb{E}(\|W_i^t - W^*\|_2) \\
&\leq \eta^2 \mathbb{E}(\|\nabla F_i(W_i^t)\|_2^2) + \|W_i^t - W^*\|_2^2 - 2\eta (\nabla F_i(W_i^t) - \nabla F_i(W^*)) \mathbb{E}(\|W_i^t - W^*\|_2) \\
&\leq \eta^2 \mathbb{E}(\|\nabla F_i(W_i^t)\|_2^2) + \|W_i^t - W^*\|_2^2 - 2\eta \mu \|W_i^t - W^*\|_2^2 \\
&\leq \eta^2 \sigma^2 + (1 - 2\eta \mu) \mathbb{E}(\|W_i^{t+1} - W^*\|_2^2) \\
&\leq \eta^2 \sigma^2 + (1 - 2\eta \mu) (\eta^2 \sigma^2 + (1 - 2\eta \mu) \mathbb{E}(\|W_i^{t+1} - W^*\|_2^2)) \\
&\dots \\
&\leq (1 - 2\eta \mu)^t \mathbb{E}(\|W_i^0 - W^*\|_2^2) + \eta \sigma^2 / (2\mu)
\end{aligned}$$

As round t increases, the gap gradually decreases to near zero, indicating the model convergence.