

机器学习理论研究导引

作业五

陈晟 MG21330006

2022 年 5 月 20 日

作业提交注意事项

- (1) 本次作业提交截止时间为 **2022/05/31 23:59:59**, 截止时间后不再接收作业, 本次作业记零分;
- (2) 作业提交方式: 使用此 LaTeX 模板书写解答, 只需提交编译生成的 pdf 文件, 将 pdf 文件上传到以下 ftp 服务器的指定位置:
地址: sftp://210.28.132.67:22, 用户名: mlt2022, 密码: mltspring2022@nju
文件夹位置: /C:/Users/mlt2022/hw_submissions/hw5_submission/ ;
- (3) pdf 文件命名方式: 学号-姓名-作业号-v 版本号, 例 MG1900000-张三-5-v1; 如果需要更改已提交的解答, 请在截止时间之前提交新版本的解答, 并将版本号加一;
- (4) 未按照要求提交作业, 或 pdf 命名方式不正确, 将会被扣除部分作业分数.

1 [50pts] Consistent Surrogate Loss

考虑对率函数 $\phi(t) = \log(1 + e^{-t})$, 回答并证明下述问题.

1. [15pts] 试求解最优实值输出函数 $f_{\phi}^*(\mathbf{x})$.
2. [10pts] 试求解最优实值输出函数对应的最优替代泛化风险 R_{ϕ}^* .
3. [25pts] 证明对率函数针对原 0/1 目标函数具有替代一致性.

Proof.

(1) 对率替代函数: $\phi(t) = \log(1 + e^{-t})$

替代函数 ϕ 在数据分布 \mathcal{D} 上的泛化风险定义为

$$R_{\phi}(f) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}}[\phi(yf(\mathbf{x}))],$$

进一步定义最优替代泛化风险为

$$R_{\phi}^* = \min_f (R_{\phi}(f)).$$

由替代函数 ϕ 在数据分布 \mathcal{D} 上的泛化风险定义为

$$R_{\phi}(f) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}}[\phi(yf(\mathbf{x}))],$$

可得

$$\begin{aligned} R_{\phi}(f) &= \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}}[\phi(yf(\mathbf{x}))] \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_X}[\eta(\mathbf{x})\phi(f(\mathbf{x})) + (1 - \eta(\mathbf{x}))\phi(-f(\mathbf{x}))] \end{aligned}$$

进一步根据

$$R_{\phi}^* = \min_f (R_{\phi}(f)).$$

可得

$$R_{\phi}^* = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_X} \left[\min_{f(\mathbf{x}) \in \mathbb{R}} (\eta(\mathbf{x})\phi(f(\mathbf{x})) + (1 - \eta(\mathbf{x}))\phi(-f(\mathbf{x}))) \right],$$

从而得到替代函数的最优实值输出函数 $f_{\phi}^*(\mathbf{x})$ 为

$$f_{\phi}^*(\mathbf{x}) \in \arg \min_{f(\mathbf{x}) \in \mathbb{R}} (\eta(\mathbf{x})\phi(f(\mathbf{x})) + (1 - \eta(\mathbf{x}))\phi(-f(\mathbf{x}))).$$

将对率替代函数: $\phi(t) = \log(1 + e^{-t})$ 带入上式可得: 最优实值输出函数 $f_{\phi}^*(\mathbf{x})$ 为

$$f_{\phi}^*(\mathbf{x}) \in \arg \min_{f(\mathbf{x}) \in \mathbb{R}} (\eta(\mathbf{x})\log(1 + e^{-f(\mathbf{x})}) + (1 - \eta(\mathbf{x}))\log(1 + e^{f(\mathbf{x})})).$$

假设 $\eta(\mathbf{x})$ 为常值, 对 $(\eta(\mathbf{x})\log(1 + e^{-f(\mathbf{x})}) + (1 - \eta(\mathbf{x}))\log(1 + e^{f(\mathbf{x})}))$ 中的 $f(\mathbf{x})$ 求导可得:

$$\begin{aligned} \because \ln(1 + e^{-f(\mathbf{x})}) &= \frac{-e^{-f(\mathbf{x})}}{1 + e^{-f(\mathbf{x})}} \\ \ln(1 + e^{f(\mathbf{x})}) &= \frac{e^{f(\mathbf{x})}}{1 + e^{f(\mathbf{x})}} \\ \therefore (\eta(\mathbf{x})\log(1 + e^{-f(\mathbf{x})}) + (1 - \eta(\mathbf{x}))\log(1 + e^{f(\mathbf{x})}))' & \\ &= \frac{e^{-f(\mathbf{x})}\eta(\mathbf{x})}{1 + e^{-f(\mathbf{x})}} + \frac{e^{f(\mathbf{x})}(1 - \eta(\mathbf{x}))}{1 + e^{f(\mathbf{x})}} \\ &= \frac{e^{-f(\mathbf{x})}\eta(\mathbf{x})}{1 + e^{-f(\mathbf{x})}} + \frac{(1 - \eta(\mathbf{x}))}{1 + e^{-f(\mathbf{x})}} \end{aligned}$$

$$= \frac{1 - \eta(x) - e^{-f(x)} \eta(x)}{1 + e^{-f(x)}}$$

当 $1 - \eta(x) - e^{-f(x)} \eta(x) = 0$ 时, $f(x) = \ln \frac{\eta(x)}{1 - \eta(x)}$

即 $f_\phi^*(\mathbf{x}) = \ln \frac{\eta(\mathbf{x})}{1 - \eta(\mathbf{x})}$

(2)

由 (1) 中

$$R_\phi^* = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_X} \left[\min_{f(\mathbf{x}) \in \mathbb{R}} (\eta(\mathbf{x}) \phi(f(\mathbf{x})) + (1 - \eta(\mathbf{x})) \phi(-f(\mathbf{x}))) \right],$$

最优泛化风险 $R_\phi^* = \min_f (R_\phi(f)) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_X} [\min_{f(\mathbf{x}) \in \mathbb{R}} (\eta(\mathbf{x}) \phi(f_\phi^*(\mathbf{x})) + (1 - \eta(\mathbf{x})) \phi(-f_\phi^*(\mathbf{x})))]$

$$= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_X} \left[(\eta(\mathbf{x}) \phi(\ln \frac{\eta(\mathbf{x})}{1 - \eta(\mathbf{x})})) + (1 - \eta(\mathbf{x})) \phi(-\ln \frac{\eta(\mathbf{x})}{1 - \eta(\mathbf{x})}) \right]$$

$$= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_X} \left[(\eta(\mathbf{x}) \ln(1 + e^{-\ln \frac{\eta(\mathbf{x})}{1 - \eta(\mathbf{x})}})) + (1 - \eta(\mathbf{x})) \ln(1 + e^{\ln \frac{\eta(\mathbf{x})}{1 - \eta(\mathbf{x})}}) \right]$$

$$= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_X} \left[(\eta(\mathbf{x}) \ln(1 + \frac{1 - \eta(\mathbf{x})}{\eta(\mathbf{x})})) + (1 - \eta(\mathbf{x})) \ln(1 + \frac{\eta(\mathbf{x})}{1 - \eta(\mathbf{x})}) \right]$$

(3) 替代函数一致性的充分条件 [Zhang, 2004b]: 对替代函数 ϕ , 若最优实值输出函数满足 $f_\phi^* \in \mathcal{F}^*$, 且存在常数 $c > 0$ 和 $s \geq 1$ 使

$$|\eta(\mathbf{x}) - 1/2|^s \leq c^s (\phi(0) - \eta(\mathbf{x}) \phi(f_\phi^*(\mathbf{x})) - (1 - \eta(\mathbf{x})) \phi(-f_\phi^*(\mathbf{x}))),$$

则替代函数 ϕ 具有一致性.

上面定理的证明如下: 对任意函数 f 和样本 $x \in \mathcal{X}$, 记

$$\begin{aligned} \Delta_1(\mathbf{x}) &= \eta(\mathbf{x}) \mathbb{I}(f(\mathbf{x}) \leq 0) \\ &\quad + (1 - \eta(\mathbf{x})) \mathbb{I}(f(\mathbf{x}) \geq 0) - \min\{\eta(\mathbf{x}), 1 - \eta(\mathbf{x})\} \end{aligned}$$

$$\begin{aligned} R(f) &= \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\mathbb{I}(yf(\mathbf{x}) \leq 0)] \\ \text{根据} \quad &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_X} [\eta(\mathbf{x}) \mathbb{I}(f(\mathbf{x}) \leq 0) + (1 - \eta(\mathbf{x})) \mathbb{I}(f(\mathbf{x}) \geq 0)] \quad \text{有} \end{aligned}$$

$$R(f) - R^* = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_X} [\Delta_1(\mathbf{x})].$$

根据 $\eta(\mathbf{x}) - 1/2$ 和 $f(\mathbf{x})$ 的不同取值, 下面分五种情况讨论 $\Delta_1(\mathbf{x})$: - 当 $\eta(\mathbf{x}) > 1/2$ 且 $f(\mathbf{x}) > 0$ 时, 有 $\Delta_1(\mathbf{x}) = 0$; - 当 $\eta(\mathbf{x}) > 1/2$ 且 $f(\mathbf{x}) \leq 0$ 时, 有 $\Delta_1(\mathbf{x}) = 2\eta(\mathbf{x}) - 1$; - 当 $\eta(\mathbf{x}) < 1/2$ 且 $f(\mathbf{x}) \geq 0$ 时, 有 $\Delta_1(\mathbf{x}) = 1 - 2\eta(\mathbf{x})$; - 当 $\eta(\mathbf{x}) < 1/2$ 且 $f(\mathbf{x}) < 0$ 时, 有 $\Delta_1(\mathbf{x}) = 0$; - 当 $\eta(\mathbf{x}) = 1/2$ 时, 无论 $f(\mathbf{x})$ 取何值, 都有 $\Delta_1(\mathbf{x}) = 0$. 综合上述五种情况可得

$$\Delta_1(\mathbf{x}) = 2\mathbb{I}((\eta(\mathbf{x}) - 1/2)f(\mathbf{x}) \leq 0) |\eta(\mathbf{x}) - 1/2|$$

代入之前得公式有

$$R(f) - R^* = 2\mathbb{E}_{(\eta(\mathbf{x}) - 1/2)f(\mathbf{x}) \leq 0} [|\eta(\mathbf{x}) - 1/2|].$$

对 $s \geq 1$, 根据 Jensen 不等式 (1.11) 有 $(\mathbb{E}[X])^s \leq \mathbb{E}[X^s]$, 于是有

$$R(f) - R^* \leq 2 \sqrt[s]{\mathbb{E}_{(\eta(\mathbf{x}) - 1/2)f(\mathbf{x}) \leq 0} [|\eta(\mathbf{x}) - 1/2|^s]}.$$

$$\begin{aligned} R_\phi(f) &= \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\phi(yf(\mathbf{x}))] \\ \text{根据} \quad &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_X} [\eta(\mathbf{x}) \phi(f(\mathbf{x})) + (1 - \eta(\mathbf{x})) \phi(-f(\mathbf{x}))], \quad \text{分别令} \end{aligned}$$

$$\Delta_2(\mathbf{x}) = \eta(\mathbf{x}) \phi(f(\mathbf{x})) + (1 - \eta(\mathbf{x})) \phi(-f(\mathbf{x}))$$

$$\Delta_3(\mathbf{x}) = \eta(\mathbf{x}) \phi(f_\phi^*(\mathbf{x})) + (1 - \eta(\mathbf{x})) \phi(-f_\phi^*(\mathbf{x}))$$

结合书中 (6.40) 和定理 6.1 中条件 (6.35) 可得

$$R(f) - R^* \leq 2c \sqrt[n]{\mathbb{E}_{(\eta(x)-1/2)f(x) \leq 0} [\phi(0) - \Delta_3(x)]}$$

另一方面, 根据书中 (6.32) ~ (6.34) 有

$$R_\phi(f) - R_\phi^* \geq \mathbb{E}_{(\eta(x)-1/2)f(x) \leq 0} [\Delta_2(x) - \Delta_3(x)]$$

设函数

$$\Gamma(t) = \eta(\mathbf{x})\phi(t) + (1 - \eta(\mathbf{x}))\phi(-t),$$

易知 $\Gamma(f(x)) = \Delta_2(x)$ 和 $\Gamma(0) = \phi(0)$. 根据凸函数性质可知, 当 $\phi(t)$ 是凸函数时 $\Gamma(t)$ 也是凸函数, 以及当 $0 \in [a, b]$ 时有

$$\Gamma(0) \leq \max\{\Gamma(a), \Gamma(b)\}$$

下面分三种情况讨论: - 若 $\eta(x) > 1/2$, 由之前得式子可知 $f_\phi^*(x) > 0$, 以及由 $(\eta(x)-1/2)f(x) \leq 0$ 可知 $f(x) \leq 0$. 因此 $0 \in [f(x), f_\phi^*(x)]$, 进一步有

$$\phi(0) = \Gamma(0) \leq \max\{\Gamma(f(x)), \Gamma(f_\phi^*(\mathbf{x}))\}$$

根据 (6.33) 有 $\Gamma(f(x)) \geq \Gamma(f_\phi^*(\mathbf{x}))$, 于是得到 $\phi(0) \leq \Gamma(f(x)) = \Delta_2(\mathbf{x})$.

- 若 $\eta(x) < 1/2$, 同理有 $f(x) \geq 0, f_\phi^*(x) < 0$ 和 $0 \in [f_\phi^*(x), f(x)]$, 以及

$$\phi(0) = \Gamma(0) \leq \max\{\Gamma(f(x)), \Gamma(f_\phi^*(x))\} = \Gamma(f(x)) = \Delta_2(x)$$

- 若 $\eta(x) = 1/2$, 对凸函数 ϕ 有

$$\phi(0) \leq \phi(f(x))/2 + \phi(-f(x))/2 = \Delta_2(x)$$

综合上述三种情况, 我们有

$$\phi(0) \leq \Delta_2(x)$$

根据之前得式子, 对任何函数 f 有

$$\begin{aligned} R(f) - R^* &\leq 2c \sqrt[n]{\mathbb{E}_{(\eta(x)-1/2)f(x) \leq 0} [\phi(0) - \Delta_3(x)]} \\ &= 2c \sqrt[n]{\mathbb{E}_{(\eta(x)-1/2)f(x) \leq 0} [\phi(0) - \Delta_2(x) + \Delta_2(\mathbf{x}) - \Delta_3(\mathbf{x})]} \\ &\leq 2c \sqrt[n]{\mathbb{E}_{(\eta(x)-1/2)f(x) \leq 0} [\Delta_2(x) - \Delta_3(x)]} \\ &\leq 2c \sqrt[n]{R_\phi(f) - R_\phi^*} \end{aligned}$$

对任何函数列 $\hat{f}_1, \hat{f}_2, \dots, \hat{f}_m, \dots$, 根据上式可知

$$R(\hat{f}_m) - R^* \leq 2c \sqrt[n]{R_\phi(\hat{f}_m) - R_\phi^*},$$

因此当 $R_\phi(\hat{f}_m) \rightarrow R_\phi^*$ 时有 $R(\hat{f}_m) \rightarrow R^*$ 成立, 定理得证.

$f^* \in \mathcal{F}^* = \{f : \text{当 } \eta(x) = 1/2 \text{ 时 } f(x) \text{ 可以是任意的实数; 当 } \eta(x) \neq 1/2 \text{ 时 } f(x)(\eta(x) - 1/2) > 0\}$.

显然 $f_\phi^* \in \mathcal{F}^*$ 成立, 其次, 看是否存在常数 $c > 0$ 和 $s \geq 1$ 使

$$|\eta(\mathbf{x}) - 1/2|^s \leq c^s (\phi(0) - \eta(\mathbf{x})\phi(f_\phi^*(\mathbf{x})) - (1 - \eta(\mathbf{x}))\phi(-f_\phi^*(\mathbf{x}))),$$

将 $\phi(x) = \ln(1 + e^{-x})$ 以及 $f_\phi^*(\mathbf{x}) = \ln \frac{\eta(x)}{1-\eta(x)}$ 带入上式得:

$$\begin{aligned} |\eta(\mathbf{x}) - 1/2|^s &\leq c^s \left(\ln 2 - \eta(x) \ln \left(\ln \frac{\eta(x)}{1-\eta(x)} \right) - (1 - \eta(x)) \ln \left(\ln \frac{\eta(x)}{1-\eta(x)} \right) \right), \\ &= c^s \left(\ln 2 + \ln \left(\ln \frac{1-\eta(x)}{\eta(x)} \right) \right) = c^s \left(\ln 2 + \left| \ln \frac{1-\eta(x)}{\eta(x)} \right| \right) \end{aligned}$$

可以发现后面一项当 $\eta(x) = 1/2$ 时或 $\eta(x) \neq 1/2$, 都存在常数 $c > 0$ 和 $s \geq 1$ 使上式成立, 即证对率函数针对原 0/1 目标函数具有替代一致性

2 [50pts] Convergence Rate

试分析下述情形下梯度下降算法 (课件第 9 页, 课本 7.2.1 节) 的收敛率.

- 可行域 \mathcal{W} 是凸集.
- 目标函数 f 是 \mathcal{W} 上可微的 α -强凸函数.
- 目标函数在 \mathcal{W} 上是 l -Lipschitz 连续的, 即

$$\forall \mathbf{u} \in \mathcal{W}, \|\nabla f(\mathbf{u})\| \leq l. \quad (2.1)$$

- 梯度下降算法的步长为:

$$\forall t, \eta_t = \frac{1}{\alpha t}. \quad (2.2)$$

Proof.

(1) 当可行域 \mathcal{W} 是凸集时

若函数为凸函数, 则其定义域为凸集, 满足条件, 其可以采用梯度下降法达到 $O(1/\sqrt{T})$ 的收敛率 [Nesterov 2018] 其基本流程如下: 1. 任意初始化 $\mathbf{w}_1 \in \mathcal{W}$

2. for $t = 1, \dots, T$ do

3. 梯度下降: $\mathbf{w}'_{t+1} = \mathbf{w}_t - \eta_t \nabla f(\mathbf{w}_t)$

4. 投影: $\mathbf{w}_{t+1} = \Pi_{\mathcal{W}}(\mathbf{w}'_{t+1})$

5. end for

6. 返回 $\bar{\mathbf{w}}_T = \frac{1}{T} \sum_{t=1}^T \mathbf{w}_t$

(2) 目标函数 f 是 \mathcal{W} 上可微的 α -强凸函数

考虑目标函数 $f: \mathcal{W} \mapsto \mathbb{R}$ 是 α -强凸函数, 对 α -强凸函数 f 有下面性质, 若 \mathbf{w}^* 为其最优解, 对于任意 $\mathbf{w} \in \mathcal{W}$ 有:

$$f(\mathbf{w}) - f(\mathbf{w}^*) \geq \frac{\lambda}{2} \|\mathbf{w} - \mathbf{w}^*\|^2$$

此外, 若梯度有上界 l , 则 $\|\mathbf{w} - \mathbf{w}^*\| \leq \frac{2l}{\lambda}$

$$f(\mathbf{w}) - f(\mathbf{w}^*) \leq \frac{2l^2}{\lambda}$$

如果考虑强凸且光滑的函数 f , 即 α -强凸函数可微时, 还满足以下的条件:

$$f(\mathbf{w}') \leq f(\mathbf{w}) + \langle \nabla f(\mathbf{w}), \mathbf{w}' - \mathbf{w} \rangle + \frac{\gamma}{2} \|\mathbf{w}' - \mathbf{w}\|^2, \forall \mathbf{w}, \mathbf{w}' \in \mathcal{W}$$

对可微的 α -强凸函数梯度下降时, 其本值任然时进行梯度下降更新后再投影到可行域, 对于其梯度下降算法有如下定理 [Nesterov, 2013]:

若目标函数满足可微的 α -强凸函数, 或 α -强凸函数且光滑, 梯度下降取得了线性收敛率 $O\left(\frac{1}{\beta^T}\right)$, 其中 $\beta > 1$

证明: 根据目标函数的性质以及更新公式

$$\begin{aligned} f(\mathbf{w}_{t+1}) &\leq f(\mathbf{w}_t) + \langle \nabla f(\mathbf{w}_t), \mathbf{w}_{t+1} - \mathbf{w}_t \rangle + \frac{\gamma}{2} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|^2 \\ &= \min_{\mathbf{w} \in \mathcal{W}} \left(f(\mathbf{w}_t) + \langle \nabla f(\mathbf{w}_t), \mathbf{w} - \mathbf{w}_t \rangle + \frac{\gamma}{2} \|\mathbf{w} - \mathbf{w}_t\|^2 \right) \\ &\leq \min_{\mathbf{w} \in \mathcal{W}} \left(f(\mathbf{w}) - \frac{\lambda}{2} \|\mathbf{w} - \mathbf{w}_t\|^2 + \frac{\gamma}{2} \|\mathbf{w} - \mathbf{w}_t\|^2 \right) \\ &\leq \min_{\substack{\mathbf{w} = \alpha \mathbf{w}^* + (1-\alpha) \mathbf{w}_t \\ \alpha \in [0, 1]}} \left(f(\mathbf{w}) + \frac{\gamma-\lambda}{2} \|\mathbf{w} - \mathbf{w}_t\|^2 \right) \end{aligned}$$

$$\begin{aligned}
&= \min_{\alpha \in [0,1]} \left(f(\alpha \mathbf{w}^* + (1-\alpha)\mathbf{w}_t) + \frac{\gamma-\lambda}{2} \|\alpha \mathbf{w}^* + (1-\alpha)\mathbf{w}_t - \mathbf{w}_t\|^2 \right) \\
&\leq \min_{\alpha \in [0,1]} \left(\alpha f(\mathbf{w}^*) + (1-\alpha)f(\mathbf{w}_t) + \frac{\gamma-\lambda}{2} \alpha^2 \|\mathbf{w}^* - \mathbf{w}_t\|^2 \right) \\
&= \min_{\alpha \in [0,1]} \left(f(\mathbf{w}_t) - \alpha(f(\mathbf{w}_t) - f(\mathbf{w}^*)) + \frac{\gamma-\lambda}{2} \alpha^2 \|\mathbf{w}^* - \mathbf{w}_t\|^2 \right) \\
&\leq \min_{\alpha \in [0,1]} \left(f(\mathbf{w}_t) - \alpha(f(\mathbf{w}_t) - f(\mathbf{w}^*)) + \frac{\gamma-\lambda}{2} \frac{2}{\lambda} \alpha^2 (f(\mathbf{w}_t) - f(\mathbf{w}^*)) \right) \\
&= \min_{\alpha \in [0,1]} \left(f(\mathbf{w}_t) + \left(\frac{\gamma-\lambda}{\lambda} \alpha^2 - \alpha \right) (f(\mathbf{w}_t) - f(\mathbf{w}^*)) \right) \\
&\text{如果 } \frac{\lambda}{2(\gamma-\lambda)} \geq 1, \text{ 令 } \alpha = 1, \text{ 则有}
\end{aligned}$$

$$f(\mathbf{w}_{t+1}) - f(\mathbf{w}^*) \leq \frac{\gamma-\lambda}{\lambda} (f(\mathbf{w}_t) - f(\mathbf{w}^*)) \leq \frac{1}{2} (f(\mathbf{w}_t) - f(\mathbf{w}^*))$$

如果 $\frac{\lambda}{2(\gamma-\lambda)} < 1$, 令 $\alpha = \frac{\lambda}{2(\gamma-\lambda)}$, 则有

$$\begin{aligned}
f(\mathbf{w}_{t+1}) - f(\mathbf{w}^*) &\leq \left(1 - \frac{\lambda}{4(\gamma-\lambda)} \right) (f(\mathbf{w}_t) - f(\mathbf{w}^*)) \\
&= \frac{4\gamma-5\lambda}{4(\gamma-\lambda)} (f(\mathbf{w}_t) - f(\mathbf{w}^*))
\end{aligned}$$

结合书中 (7.20) 和 (7.21), 令

$$\beta = \begin{cases} \frac{\lambda}{\gamma-\lambda}, & \frac{\lambda}{2(\gamma-\lambda)} \geq 1 \\ \frac{4(\gamma-\lambda)}{4\gamma-5\lambda}, & \frac{\lambda}{2(\gamma-\lambda)} < 1 \end{cases}$$

那么下式总是成立

$$f(\mathbf{w}_{t+1}) - f(\mathbf{w}^*) \leq \frac{1}{\beta} (f(\mathbf{w}_t) - f(\mathbf{w}^*))$$

将上式扩展可得

$$f(\mathbf{w}_T) - f(\mathbf{w}^*) \leq \frac{1}{\beta^{T-1}} (f(\mathbf{w}_1) - f(\mathbf{w}^*)) = O\left(\frac{1}{\beta^T}\right)$$

即证可微的 α -强凸函数, 或 α -强凸函数且光滑, 梯度下降取得了线性收敛率 $O\left(\frac{1}{\beta^T}\right)$, 其中 $\beta > 1$

(3) 目标函数在 \mathcal{W} 上是 l -Lipschitz 连续的, 即 $\forall \mathbf{u} \in \mathcal{W}, \|\nabla f(\mathbf{u})\| \leq l$

梯度下降收敛率若目标函数是 l -Lipschitz 连续函数, 且可行域有界, 则采用固定步长梯度下降的收敛率为 $O\left(\frac{1}{\sqrt{T}}\right)$. 证明假设可行域 \mathcal{W} 直径为 Γ , 并且目标函数满足 l -Lipschitz 连续, 即对于任意 $\mathbf{u}, \mathbf{v} \in \mathcal{W}$,

$$\|\mathbf{u} - \mathbf{v}\| \leq \Gamma, \|\nabla f(\mathbf{u})\| \leq l.$$

为了简化分析, 考虑固定的步长 $\eta_t = \eta$. 对于任意的 $\mathbf{w} \in \mathcal{W}$,

$$f(\mathbf{w}_t) - f(\mathbf{w}) \leq \langle \nabla f(\mathbf{w}_t), \mathbf{w}_t - \mathbf{w} \rangle = \frac{1}{\eta} \langle \mathbf{w}_t - \mathbf{w}'_{t+1}, \mathbf{w}_t - \mathbf{w} \rangle$$

$$\begin{aligned}
&= \frac{1}{2\eta} \left(\|\mathbf{w}_t - \mathbf{w}\|^2 - \|\mathbf{w}'_{t+1} - \mathbf{w}\|^2 + \|\mathbf{w}_t - \mathbf{w}'_{t+1}\|^2 \right) \\
&= \frac{1}{2\eta} \left(\|\mathbf{w}_t - \mathbf{w}\|^2 - \|\mathbf{w}'_{t+1} - \mathbf{w}\|^2 \right) + \frac{\eta}{2} \|\nabla f(\mathbf{w}_t)\|^2 \\
&\leq \frac{1}{2\eta} \left(\|\mathbf{w}_t - \mathbf{w}\|^2 - \|\mathbf{w}_{t+1} - \mathbf{w}\|^2 \right) + \frac{\eta}{2} \|\nabla f(\mathbf{w}_t)\|^2
\end{aligned}$$

最后一个不等号利用了凸集合投影操作的非扩展性质 [Nemirovski et al., 2009]:

$$\|\Pi_{\mathcal{W}}(x) - \Pi_{\mathcal{W}}(z)\| \leq \|x - z\|, \quad \forall x, z.$$

注意到目标函数满足 l -Lipschitz 连续, 由上面的公式可得

$$f(\mathbf{w}_t) - f(\mathbf{w}) \leq \frac{1}{2\eta} \left(\|\mathbf{w}_t - \mathbf{w}\|^2 - \|\mathbf{w}_{t+1} - \mathbf{w}\|^2 \right) + \frac{\eta}{2} l^2.$$

对上式从 $t = 1$ 到 T 求和, 有

$$\begin{aligned} \sum_{t=1}^T f(\mathbf{w}_t) - Tf(\mathbf{w}) &\leq \frac{1}{2\eta} \left(\|\mathbf{w}_1 - \mathbf{w}\|^2 - \|\mathbf{w}_{T+1} - \mathbf{w}\|^2 \right) + \frac{\eta T}{2} l^2 \\ &\leq \frac{1}{2\eta} \|\mathbf{w}_1 - \mathbf{w}\|^2 + \frac{\eta T}{2} l^2 \leq \frac{1}{2\eta} \Gamma^2 + \frac{\eta T}{2} l^2 \end{aligned}$$

最后, 依据 Jensen 不等式可得

$$\begin{aligned} f(\bar{\mathbf{w}}_T) - f(\mathbf{w}) &= f\left(\frac{1}{T} \sum_{t=1}^T \mathbf{w}_t\right) - f(\mathbf{w}) \\ &\leq \frac{1}{T} \sum_{t=1}^T f(\mathbf{w}_t) - f(\mathbf{w}) \leq \frac{\Gamma^2}{2\eta T} + \frac{\eta l^2}{2} \end{aligned}$$

因此,

$$f(\bar{\mathbf{w}}_T) - \min_{\mathbf{w} \in \mathcal{W}} f(\mathbf{w}) \leq \frac{\Gamma^2}{2\eta T} + \frac{\eta l^2}{2} = \frac{l\Gamma}{\sqrt{T}} = O\left(\frac{1}{\sqrt{T}}\right)$$

其中步长设置为 $\eta = \Gamma/(l\sqrt{T})$ 定理得证

(4) 当步长设置为变长 $\forall t, \eta_t = \frac{1}{\alpha t}$ 时的收敛率

以 (3) 中的条件更改收敛率进行推算

对 (3) 中间步骤, 从 $t = 1$ 到 T 求和, 有

$$\begin{aligned} \sum_{t=1}^T f(\mathbf{w}_t) - Tf(\mathbf{w}) &\leq \frac{1}{2\eta} \left(\|\mathbf{w}_1 - \mathbf{w}\|^2 - \|\mathbf{w}_{T+1} - \mathbf{w}\|^2 \right) + \frac{\eta T}{2} l^2 \\ &\leq \frac{1}{2\eta} \|\mathbf{w}_1 - \mathbf{w}\|^2 + \frac{\eta T}{2} l^2 \leq \frac{1}{2\eta} \Gamma^2 + \frac{\eta T}{2} l^2 \end{aligned}$$

将 $\eta_t = \frac{1}{\alpha t}$ 带入上式可得:

$$f(\mathbf{w}_t) - f(\mathbf{w}) \leq \frac{1}{2\alpha t} \left(\|\mathbf{w}_t - \mathbf{w}\|^2 - \|\mathbf{w}_{t+1} - \mathbf{w}\|^2 \right) + \frac{\alpha t}{2} l^2$$

$$\begin{aligned} \sum_{t=1}^T f(\mathbf{w}_t) - Tf(\mathbf{w}) &\leq \frac{1}{2\alpha} \left(\|\mathbf{w}_1 - \mathbf{w}\|^2 - \|\mathbf{w}_{T+1} - \mathbf{w}\|^2 \right) + \frac{\alpha T^2}{2} l^2 \\ &\leq \frac{1}{2\alpha} \|\mathbf{w}_1 - \mathbf{w}\|^2 + \frac{\alpha T^2}{2} l^2 \leq \frac{1}{2\alpha} \Gamma^2 + \frac{\alpha T^2}{2} l^2 \end{aligned}$$

最后, 依据 Jensen 不等式可得

$$\begin{aligned} f(\bar{\mathbf{w}}_T) - f(\mathbf{w}) &= f\left(\frac{1}{T} \sum_{t=1}^T \mathbf{w}_t\right) - f(\mathbf{w}) \\ &\leq \frac{1}{T} \sum_{t=1}^T f(\mathbf{w}_t) - f(\mathbf{w}) \leq \frac{\Gamma^2}{2\alpha T} + \frac{\alpha T l^2}{2} \end{aligned}$$

因此,

$$f(\bar{\mathbf{w}}_T) - \min_{\mathbf{w} \in \mathcal{W}} f(\mathbf{w}) \leq \frac{\Gamma^2}{2\alpha T} + \frac{\alpha T l^2}{2}$$

当 $\frac{\alpha T l^2}{2} = \frac{\Gamma^2}{2\alpha T}$ 时, $\alpha = \frac{\Gamma}{Tl}$

$$f(\bar{\mathbf{w}}_T) - \min_{\mathbf{w} \in \mathcal{W}} f(\mathbf{w}) \leq l\Gamma = O(1)$$

即其收敛率变为常数级别了