

Frank-Wolfe算法在目标函数为非凸, 凸以及强凸情形下收敛性总结

2020.8.15

§1 Frank-Wolfe算法简介 [1]

Frank-Wolfe(FW)算法是众多非线性约束中的方法之一, 近年来由于内存要求较小和无投影迭代的特点得到了广泛的应用。它可以解决如下的问题:

$$\underset{\mathbf{x} \in \mathcal{D}}{\text{minimize}} f(\mathbf{x}) \quad (1)$$

其中 f 可微并且是 L -Lipschitz的, 定义域 \mathcal{D} 是凸紧集。

FW算法首先设置一个初始点 \mathbf{x}_0 , 并由此构造了一个迭代序列 $\mathbf{x}_1, \mathbf{x}_2, \dots$, 这个序列可以收敛到问题的最优解。FW 算法的过程如下:

算法 1 FW Algorithm

```
1: Input: initial guess  $\mathbf{x}_0$ , tolerance  $\delta > 0$ 
2: for  $t = 0, 1, \dots$  do
3:    $\mathbf{s}_t \in \operatorname{argmax}_{\mathbf{s} \in \mathcal{D}} \langle -\nabla f(\mathbf{x}_t), \mathbf{s} \rangle$ 
4:    $\mathbf{d}_t = \mathbf{s}_t - \mathbf{x}_t$ 
5:    $g_t = \langle -\nabla f(\mathbf{x}_t), \mathbf{d}_t \rangle$ 
6:   if  $g_t < \delta$  then
7:     Return  $\mathbf{x}_t$ 
8:   end if
9:   Variant 1: 设置步长  $\gamma_t = \min \left\{ \frac{g_t}{L \|\mathbf{d}_t\|^2}, 1 \right\}$ 
10:  Variant 2: 通过线搜索设置步长  $\gamma_t = \operatorname{argmin}_{\gamma \in [0, 1]} f(\mathbf{x}_t + \gamma \mathbf{d}_t)$ 
11:   $\mathbf{x}_{t+1} = \mathbf{x}_t + \gamma_t \mathbf{d}_t$ 
12: end for
13: Return  $\mathbf{x}_{t+1}$ 
```

与其他约束优化算法(如投影梯度下降法)相反, FW 算法不需要投影, 因此有时称之为无投影算法, 但需要求解算法1第3 步中的问题:

$$\mathbf{s}_t \in \operatorname{argmin}_{\mathbf{s} \in \mathcal{D}} \langle \nabla f(\mathbf{x}_t), \mathbf{s} \rangle \quad (2)$$

算法的其余部分主要是寻找合适的步长,以便 $f(\mathbf{x})$ 朝减小的方向移动。在FW算法所允许的许多不同的步长规则中,我们将其分为变形1和变形2。第1个变形很容易计算,只依赖于Lipschitz常数 L 的知识。第2个变形需要在每次迭代时解决一个一维问题。在某些情况下如当目标函数 f 为最小二乘问题且约束集 \mathcal{D} 是 ℓ_1 范数球时,步长选择这个问题有一个显式的最优解,在这种情况下,应首选变形2这种方法。但是,在一般情况下如果没有显式的最优解,在这种情况下,应首选变形1。事实证明,对于规模较大的问题的线性子问题既有闭式解,也有有效算法。与投影相比,使用线性最小化还有其他重要的结果。根据线性规划的性质,这种线性极小化问题的最优解总是约束集的一个顶点。还有一些其它设置步长大小的策略。例如,步长也可以设置为 $\gamma_t = \frac{2}{t+2}$,它不依赖于优化产生的任何其它变量。因此,它在实践中不能与其它选择步长的策略竞争,尽管它确实达到了相同的理论收敛速度。

对于 γ_t 还有一种选择,是由Demyanov和Rubinov提出来的,与变形1比较相似:

$$\gamma_t = \min \left\{ \frac{g_t}{L \text{diam}(\mathcal{D})^2}, 1 \right\} \quad (3)$$

其中diam表示相对于欧几里德范数的直径。然而,由于根据直径的定义,我们总是有:

$$\|\mathbf{x}_t - \mathbf{s}_t\|^2 \leq \text{diam}(\mathcal{D})^2 \quad (4)$$

这个变形提供的步长总是小于变形1的步长,并且给出了更差的收敛性。对该步长的进一步改进包括用局部估计替换Lipschitz常数 L ,从而允许更大的步长。

可以看到Frank-Wolfe算法是一种通过求解一系列线性规划问题来解决潜在非线性问题的算法,这种方法的有效性与快速解决线性子问题的能力紧密相连。这与其他一阶方法(如Nesterov的加速方法)形成了对比,后者的前提是在每次迭代中能够解决由强凸函数定义的某些投影问题。在许多应用中,求解一个线性优化子问题要比求解相关的投影子问题简单得多。此外,在许多应用中,线性优化子问题的解通常是高度结构化的,并表现出特殊的稀疏性或低秩性, Frank-Wolfe算法可以利用这一点。Frank-Wolfe算法在每次迭代求解一个子问题,并产生一系列可行解,这些可行解都是之前所有子问题解的凸组合,对于适当选择的步长,可以推导出 $\mathcal{O}(\frac{1}{k})$ 的收敛速度。由于子问题解的结构以及迭代是子问题解的凸组合, Frank-Wolfe方法返回的可行解通常也是高度结构化的。例如,当可行域是单位单形 $\Delta_n = \{\lambda \in \mathbb{R}^n : e^T \lambda = 1, \lambda > 0\}$ 时,线性规划的解总是返回一个极值点,则Frank-Wolfe算法在第 k 次迭代时产生的解最多有 k 个非零项。可推广到矩阵优化设置:如果可行域是核范数诱导的球,则在迭代 k 时,该方法在第 k 步产生的矩阵的秩最大不超过 k 。在许多应用中,这样的结构特性是非常可取的,在这种情况下, Frank-Wolfe方法可能比一些加速方法更有吸引力,尽管Frank-Wolfe方法的收敛速度较慢。

§2 收敛性分析

Frank-Wolfe算法在一般情况下收敛。正如已经看到的,凸的目标函数是必要的,以获得收敛保证。如前所述,假设 f 在 L -Lipschitz梯度下是可微的, \mathcal{D} 是一个凸紧集。在这一部分,将给出两个主要的收敛结果:一个用于一般目标,一个用于凸目标。为简单起见,假设线性子问题被精确地求解,但是这些证明可以很容易地扩展到考虑近似线性最小化的问题上去。本节的其余部分结构如下:首先介绍了两个关键定义和技术引理,最后证明了收敛性结果。

定义1. 称 $\mathbf{x}^* \in \mathcal{D}$ 为最优点, 当且仅当:

$$\langle \nabla f(\mathbf{x}^*), \mathbf{x} - \mathbf{x}^* \rangle \geq 0, \forall \mathbf{x} \in \mathcal{D}. \quad (5)$$

这个定义的直观含义是, 如果以 \mathbf{x}^* 起点的多面体中的每个方向都与梯度正相关, 则 \mathbf{x}^* 是一个最优点。否则, 如果以 \mathbf{x}^* 起点没有可行的下降方向, \mathbf{x}^* 是一个最优点。

定义2. 将Frank-Wolfe间隙记成 g_t , 定义为:

$$g_t = \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{s}_t \rangle. \quad (6)$$

下一个引理将连续两次迭代的目标函数值联系起来, 对于凸和非凸目标, 它在证明收敛结果的过程中起到了关键作用。

引理1. 令 $\mathbf{x}_0, \mathbf{x}_1, \dots$ 是由FW 算法产生的迭代序列, 对于任意的 $\xi \in [0, 1]$, 有下面的不等式:

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \xi g_t + \frac{1}{2} \xi^2 L \text{diam}(\mathcal{D})^2 \quad (7)$$

证明. 关于 f 的Lipschitz梯度假设的一个结果是, $y \in \mathcal{D}$ 的每一点有上界, 即:

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{1}{2} \|\mathbf{y} - \mathbf{x}\|^2 \quad (8)$$

应用这个不等式, 对定义域中的任何 \mathbf{x}, \mathbf{y} 都有效。对于 $\gamma \in [0, 1]$ 且 $\mathbf{x} = \mathbf{x}_t, \mathbf{y} = (1 - \gamma)\mathbf{x}_t + \gamma\mathbf{s}_t$ 的情况, \mathbf{y} 仍然在定义域中, 我们有:

$$f((1 - \gamma)\mathbf{x}_t + \gamma\mathbf{s}_t) \leq f(\mathbf{x}_t) + \gamma \langle \nabla f(\mathbf{x}_t), \mathbf{s}_t - \mathbf{x}_t \rangle + \frac{L\gamma^2}{2} \|\mathbf{s}_t - \mathbf{x}_t\|^2 \quad (9)$$

当 $\gamma \in [0, 1]$ 时, 需要最小化上面公式的右边部分。这是一个最小化 γ 的二次函数, 将 γ_t^* 记为:

$$\gamma_t^* = \min \left\{ \frac{g_t}{L \|\mathbf{x}_t - \mathbf{s}_t\|^2}, 1 \right\}. \quad (10)$$

当 $\gamma = \gamma_t^*$ 时, 有:

$$\begin{aligned} f((1 - \gamma_t^*)\mathbf{x}_t + \gamma_t^*\mathbf{s}_t) &\leq f(\mathbf{x}_t) - \gamma_t^* g_t + \frac{L(\gamma_t^*)^2}{2} \|\mathbf{s}_t - \mathbf{x}_t\|^2 \\ &= f(\mathbf{x}_t) + \min_{\xi \in [0, 1]} \left\{ -\xi g_t + \frac{L\xi^2}{2} \|\mathbf{s}_t - \mathbf{x}_t\|^2 \right\} \\ &\leq f(\mathbf{x}_t) - \xi g_t + \frac{L\xi^2}{2} \|\mathbf{s}_t - \mathbf{x}_t\|^2 \quad (\forall \xi \in [0, 1]) \\ &\leq f(\mathbf{x}_t) - \xi g_t + \frac{L\xi^2}{2} \text{diam}(\mathcal{D})^2. \end{aligned} \quad (11)$$

上述不等式的右边已经包含引理中的结果, 对于算法的变形1, 我们有 $f(\mathbf{x}_{t+1}) = f((1 - \gamma_t^*)\mathbf{x}_t + \gamma_t^*\mathbf{s}_t)$, 这种情况下 γ_t 和 γ_t^* 相等。对于算法的变形2, 我们有 $f(\mathbf{x}_{t+1}) \leq f((1 - \gamma_t^*)\mathbf{x}_t + \gamma_t^*\mathbf{s}_t)$, 根据线搜索的定义, $f(\mathbf{x}_{t+1})$ 是最小化的目标函数值。因此, 无论是哪种情况, 我们都有:

$$f(\mathbf{x}_{t+1}) \leq f((1 - \gamma_t^*)\mathbf{x}_t + \gamma_t^*\mathbf{s}_t). \quad (12)$$

将该式与前一个不等式联系起来, 便得到要证明的不等式。 \square

§2.1 目标函数非凸

下面是第一个收敛速度结果, 对于梯度 L -Lipschitz 目标函数是有效的。

定理1. 如果 f 是可微的且具有 L -Lipschitz 梯度, 那么在最佳 Frank-Wolfe 间隙上有 $\mathcal{O}(1/\sqrt{t})$ 的上界:

$$\min_{0 \leq i \leq t} g_i \leq \frac{\max\{2h_0, L * \text{diam}(\mathcal{D})^2\}}{\sqrt{t+1}}, \quad (13)$$

其中 $h_0 = f(\mathbf{x}_0) - \min_{\mathbf{x} \in \mathcal{D}} f(\mathbf{x})$ 是初始全局次优的。

证明. 根据引理1, 对于任意的 $\xi \in [0, 1]$ 有:

$$\begin{aligned} f(\mathbf{x}_{t+1}) &\leq f(\mathbf{x}_t) - \xi g_t + \frac{L\xi^2}{2} \text{diam}(\mathcal{D})^2 \\ &\leq f(\mathbf{x}_t) - \xi g_t + \frac{C\xi^2}{2} \end{aligned} \quad (14)$$

其中 $C = L\text{diam}(\mathcal{D})^2$. 考虑使右边公式最小化的 ξ 的值, 可得 $\xi^* = \min\{g_t/C, 1\}$, 现在根据 ξ^* 的值做出选择, 如果 $g_t \leq C$, 则 $\xi^* = g_t/C$, 利用上一个不等式中的值, 得到:

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{g_t^2}{2C} \quad (15)$$

如果 $g_t > C$, 则 $\xi^* = 1$, 有下面的不等式:

$$\begin{aligned} f(\mathbf{x}_{t+1}) &\leq f(\mathbf{x}_t) - g_t + \frac{C}{2} \\ &\leq f(\mathbf{x}_t) - \frac{g_t}{2} \end{aligned} \quad (16)$$

结合两式得:

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{g_t}{2} \min\left\{\frac{g_t}{C}, 1\right\} \quad (17)$$

其中 $g_t^* = \min_{0 \leq i \leq t} g_i$, 接下来对 g_t^* 进行选择, 如果 $g_t^* \leq C$, 则 $\min\{g_t^*/C, 1\} = g_t^*/C$ 解前面关于 g_t^* 的不等式得:

$$g_t^* \leq \sqrt{\frac{2Ch_0}{t+1}} \leq \frac{2h_0 + C}{2\sqrt{t+1}} \leq \frac{\max\{2h_0, C\}}{\sqrt{t+1}} \quad (18)$$

在第二个不等式中利用 Young 不等式 $ab \leq \frac{a^2}{2} + \frac{b^2}{2}$, $a = \sqrt{2h_0}$, $b = \sqrt{C}$.

如果 $g_t^* > C$, 则 $\min\{g_t^*/C, 1\} = 1$, 且有:

$$g_t^* \leq \frac{2Ch_0}{t+1} \leq \frac{2h_0}{\sqrt{t+1}} \leq \frac{\max\{2h_0, C\}}{\sqrt{t+1}} \quad (19)$$

因此, 在这两种情况下可得:

$$g_t^* \leq \frac{\max\{2h_0, C\}}{\sqrt{t+1}} \quad (20)$$

□

§2.2 目标函数为凸

定理2. 如果 f 是可微凸函数, 且梯度 L -Lipschitz, 则对于函数的次优性有如下的收敛速度:

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{2L * \text{diam}(\mathcal{D})^2}{t+1} \quad (21)$$

证明. 由凸性可得:

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \langle \nabla f(\mathbf{x}^*), \mathbf{x} - \mathbf{x}^* \rangle \quad (22)$$

令 $e_t = \mathbf{A}_t(f(\mathbf{x}_t) - f(\mathbf{x}^*))$, \mathbf{A}_t 是正定的, $C = L\text{diam}(\mathcal{D})^2$, 然后可以得到下面的不等式:

$$\begin{aligned} e_{t+1} - e_t &= \mathbf{A}_{t+1}(f(\mathbf{x}_{t+1}) - f(\mathbf{x}^*)) - \mathbf{A}_t(f(\mathbf{x}_t) - f(\mathbf{x}^*)) \\ &\leq \mathbf{A}_{t+1}(f(\mathbf{x}_t) - \xi g_t + \frac{\xi^2 C}{2} - f(\mathbf{x}^*)) - \mathbf{A}_t(f(\mathbf{x}_t) - f(\mathbf{x}^*)) \\ &\leq \mathbf{A}_{t+1}(f(\mathbf{x}_t) - f(\mathbf{x}^*) - \xi(f(\mathbf{x}_t) - f(\mathbf{x}^*)) + \frac{\xi^2 C}{2}) - \mathbf{A}_t(f(\mathbf{x}_t) - f(\mathbf{x}^*)) \\ &= ((1 - \xi)\mathbf{A}_{t+1} - \mathbf{A}_t)(f(\mathbf{x}_t) - f(\mathbf{x}^*)) + \mathbf{A}_{t+1} \frac{\xi^2 C}{2} \end{aligned} \quad (23)$$

现在, 选择 $\mathbf{A}_t = \frac{t(t+1)}{2}$, $\xi = 2/(t+2)$, 我们有:

$$(1 - \xi)\mathbf{A}_{t+1} - \mathbf{A}_t = \frac{t(t+1)}{2} - \frac{t(t+1)}{2} = 0 \quad (24)$$

$$\mathbf{A}_{t+1} \frac{\xi^2}{2} = \frac{t+1}{t+2} \leq 1, \quad (25)$$

因此将式中的 \mathbf{A}_t 和 ξ 替换, 则有:

$$e_{t+1} - e_t \leq C. \quad (26)$$

将不等式从 0 加到 $t-1$, 并使用 $e_0 = 0$, 可得:

$$e_t \leq tC \implies f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{2C}{t+1}. \quad (27)$$

□

§2.3 目标函数为强凸

定义3. (光滑函数) 设凸集 $\mathcal{K} \subset \mathbf{E}$, 称函数 $f: \mathbf{E} \rightarrow \mathcal{R}$ 是光滑函数, 当且仅当对于 $\forall \mathbf{x}, \mathbf{y} \in \mathcal{K}$ 有:

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x}) \cdot (\mathbf{y} - \mathbf{x}) + \frac{\beta}{2} \|\mathbf{x} - \mathbf{y}\|^2. \quad (28)$$

定义4. (强凸函数) 设凸集 $\mathcal{K} \subset \mathbf{E}$, 称函数 $f: \mathbf{E} \rightarrow \mathcal{R}$ 是 α 强凸函数, 当且仅当满足以下两个条件:

- 1. $\forall \mathbf{x}, \mathbf{y} \in \mathcal{K}$:

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x}) \cdot (\mathbf{y} - \mathbf{x}) + \frac{\alpha}{2} \|\mathbf{x} - \mathbf{y}\|^2.$$

- 2. $\forall \mathbf{x}, \mathbf{y} \in \mathcal{K}, \gamma \in [0, 1]$:

$$f(\gamma \mathbf{x} + (1 - \gamma)\mathbf{y}) \leq \gamma f(\mathbf{x}) + (1 - \gamma)f(\mathbf{y}) - \frac{\alpha}{2} \gamma(1 - \gamma) \|\mathbf{x} - \mathbf{y}\|^2.$$

若 f 是强凸的, $\mathbf{x}^* = \operatorname{argmax}_{\mathbf{x} \in \mathcal{K}} f(\mathbf{x})$, 由一阶最优性条件可得, 对任意 $\mathbf{x} \in \mathcal{K}$:

$$f(\mathbf{x}) - f(\mathbf{x}^*) \geq \frac{\alpha}{2} \|\mathbf{x} - \mathbf{x}^*\|^2 \quad (29)$$

稍作变换, 有:

$$\sqrt{\frac{2}{\alpha}(f(\mathbf{x}) - f(\mathbf{x}^*))} \cdot \|\nabla f(\mathbf{x})\|_* \geq \|\mathbf{x} - \mathbf{x}^*\| \cdot \|\nabla f(\mathbf{x})\|_* \quad (30)$$

$$\geq (\mathbf{x} - \mathbf{x}^*) \cdot \nabla f(\mathbf{x}) \quad (31)$$

$$\geq f(\mathbf{x}) - f(\mathbf{x}^*) \quad (32)$$

其中第二个不等式利用Holder不等式, 第三个不等式来自 f 的凸性。因此对任意 $\mathbf{x} \in \mathcal{K}$ 有:

$$\|\nabla f(\mathbf{x})\|_* \geq \sqrt{\frac{\alpha}{2}} \cdot \sqrt{f(\mathbf{x}) - f(\mathbf{x}^*)} \quad (33)$$

定义5. (强凸集)称凸集 $\mathcal{K} \subset \mathbf{E}$ 是 α 强凸集, 当且仅当 $\forall \mathbf{x}, \mathbf{y} \in \mathcal{K}, \forall \gamma \in [0, 1]$, 以及 $\mathbf{z} \in \mathbf{E}$ 且 $\|\mathbf{z}\| = 1$, 有:

$$\gamma \mathbf{x} + (1 - \gamma) \mathbf{y} + \gamma(1 - \gamma) \frac{\alpha}{2} \|\mathbf{x} - \mathbf{y}\|^2 \mathbf{z} \in \mathcal{K}. \quad (34)$$

也就是说 \mathcal{K} 包含以 $\gamma \mathbf{x} + (1 - \gamma) \mathbf{y}$ 为中心, 半径为 $\gamma(1 - \gamma) \frac{\alpha}{2} \|\mathbf{x} - \mathbf{y}\|^2$ 的球.

接下来考虑函数 f , f 是 β_f 光滑的以及 α_f 强凸, 我们进一步假设可行集 \mathcal{K} 是 $\alpha_{\mathcal{K}}$ 强凸的。在 \mathbf{x}_t 处的逼近误差记为 h_t , $h_t = \operatorname{argmin}_{\mathbf{x} \in \mathcal{K}} f(\mathbf{x})$, 其中 $\mathbf{x}^* = \operatorname{argmin}_{\mathbf{x} \in \mathcal{K}} f(\mathbf{x})$.

引理2. 使用FW算法, 在第 t 步迭代有:

$$h_{t+1} \leq h_t \cdot \max\left\{\frac{1}{2}, 1 - \frac{\alpha_{\mathcal{K}} \|\nabla f(\mathbf{x}_t)\|_*}{8\beta_f}\right\} \quad (35)$$

证明. 由 \mathbf{s}_t 的取法得:

$$(\mathbf{s}_t - \mathbf{x}_t) \cdot \nabla f(\mathbf{x}_t) \leq (\mathbf{x}^* - \mathbf{x}_t) \cdot \nabla f(\mathbf{x}_t) \leq f(\mathbf{x}^*) - f(\mathbf{x}_t) = -h_t \quad (36)$$

第二个不等式是由于 f 的凸性。记 $\mathbf{c}_t = \frac{1}{2}(\mathbf{x}_t + \mathbf{s}_t)$, $\mathbf{w}_t \in \operatorname{argmin}_{\mathbf{w} \in \mathbf{E}, \|\mathbf{w}\| \leq 1} \mathbf{w} \cdot \nabla f(\mathbf{x}_t)$. 由Holder不等式可得 $\mathbf{w}_t \cdot \nabla f(\mathbf{x}_t) = -\|\nabla f(\mathbf{x}_t)\|_*$, 由可行集的强凸性, 点 $\tilde{\mathbf{p}}_t = \mathbf{c}_t + \frac{\alpha_{\mathcal{K}}}{8} \|\mathbf{x}_t - \mathbf{p}_t\|^2 \mathbf{z}$ 也在可行集 \mathcal{K} 中。再次利用 \mathbf{s}_t 的最优性可得:

$$(\mathbf{s}_t - \mathbf{x}_t) \cdot \nabla f(\mathbf{x}_t) \leq (\tilde{\mathbf{p}}_t - \mathbf{x}_t) \cdot \nabla f(\mathbf{x}_t) \quad (37)$$

$$= \frac{1}{2}(\mathbf{p}_t - \mathbf{x}_t) \cdot \nabla f(\mathbf{x}_t) + \frac{\alpha_{\mathcal{K}} \|\mathbf{x}_t - \mathbf{p}_t\|^2}{8} \mathbf{w}_t \cdot \nabla f(\mathbf{x}_t) \quad (38)$$

$$\leq -\frac{1}{2}h_t - \frac{\alpha_{\mathcal{K}} \|\mathbf{x}_t - \mathbf{p}_t\|^2}{8} \|\nabla f(\mathbf{x}_t)\|_* \quad (39)$$

上面最后一个不等式来自式子(36)。现在分析下降时的误差, 由于 f 是光滑函数得:

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) + \gamma_t(\mathbf{p}_t - \mathbf{x}_t) \cdot \nabla f(\mathbf{x}_t) + \frac{\beta_f}{2} \gamma_t^2 \|\mathbf{p}_t - \mathbf{x}_t\|^2. \quad (40)$$

上式两边同时减去 $f(\mathbf{x}^*)$ 得:

$$h(\mathbf{x}_{t+1}) \leq h(\mathbf{x}_t) + \gamma_t(\mathbf{p}_t - \mathbf{x}_t) \cdot \nabla f(\mathbf{x}_t) + \frac{\beta_f}{2} \gamma_t^2 \|\mathbf{p}_t - \mathbf{x}_t\|^2. \quad (41)$$

结合(37)式有:

$$h_{t+1} \leq h_t(1 - \frac{\gamma_t}{2}) - \gamma_t \frac{\alpha_{\mathcal{K}} \|\mathbf{x}_t - \mathbf{p}_t\|^2}{8} \|\nabla f(\mathbf{x}_t)\|_* + \frac{\beta_f}{2} \gamma_t^2 \|\mathbf{p}_t - \mathbf{x}_t\|^2 \quad (42)$$

$$= h_t(1 - \frac{\gamma_t}{2}) + \frac{\|\mathbf{p}_t - \mathbf{x}_t\|^2}{2} (\gamma_t^2 \beta - \gamma_t \frac{\alpha_{\mathcal{K}} \|\nabla f(\mathbf{x}_t)\|_*}{4}). \quad (43)$$

当 $\frac{\alpha_{\mathcal{K}} \|\nabla f(\mathbf{x}_t)\|_*}{4} \geq \beta_f$, 令 $\gamma_t = 1$ 可得:

$$h_{t+1} \leq \frac{h_t}{2}. \quad (44)$$

其他情况下, 令 $\gamma_t = \frac{\alpha_{\mathcal{K}} \|\nabla f(\mathbf{x}_t)\|_*}{4}$ 可得:

$$h_{t+1} \leq h_t(1 - \frac{\alpha_{\mathcal{K}} \|\nabla f(\mathbf{x}_t)\|_*}{8}). \quad (45)$$

□

注意到引理2只依赖于可行集 \mathcal{K} 的强凸性, 不需要利用 f 的凸性与函数光滑的条件, 也不要求 f 是强凸的。

定理3. 令 $\mathbf{x}^* = \operatorname{argmin}_{\mathbf{x} \in \mathcal{K}} f(\mathbf{x})$, $M = \frac{\sqrt{\alpha_f} \alpha_{\mathcal{K}}}{8\sqrt{2}\beta_f}$, 记 $D_{\mathcal{K}} = \max_{\mathbf{x}, \mathbf{y} \in \mathcal{K}} \|\mathbf{x} - \mathbf{y}\|$, 由FW算法得到的第 t 个迭代点 \mathbf{x}_t 满足:

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{\max\{\frac{9}{2}\beta_f D_{\mathcal{K}}^2, 18M^{-2}\}}{(t+2)^2} = \mathcal{O}(\frac{1}{t^2}) \quad (46)$$

证明. 令 $M = \frac{\sqrt{\alpha_f} \alpha_{\mathcal{K}}}{8\sqrt{2}\beta_f}$, $C = \max\{\frac{9}{2}\beta_f D_{\mathcal{K}}^2, 18M^{-2}\}$. 接下来通过数学归纳法证明 $\forall t > 1, h_t \leq \frac{C}{(t+2)^2}$.

由于目标函数满足式子(33), 由引理2得:

$$h_{t+1} \leq h_t \cdot \max\{\frac{1}{2}, 1 - \frac{\alpha_{\mathcal{K}} \sqrt{\alpha_f}}{8\sqrt{2}\beta_f} \sqrt{h_t}\} \quad (47)$$

$$= h_t \cdot \max\{\frac{1}{2}, 1 - M h_t^{1/2}\} \quad (48)$$

当 $t = 1$, 由于 f 是 β_f 光滑的可得:

$$f(\mathbf{x}_1) - f(\mathbf{x}^*) = f(\mathbf{x}_0 + \gamma_0(\mathbf{p}_0 - \mathbf{x}_0)) - f(\mathbf{x}^*) \quad (49)$$

$$\leq h_0 + \gamma_0(\mathbf{p}_0 - \mathbf{x}_0) \cdot \nabla f(\mathbf{x}_0) + \frac{\beta_f \gamma_0^2}{2} D_{\mathcal{K}}^2 \quad (50)$$

$$\leq h_0(1 - \gamma_0) + \frac{\beta_f \gamma_0^2}{2} D_{\mathcal{K}}^2. \quad (51)$$

令 $\gamma_0 = 1$ 得 $h_1 \leq \frac{\beta_f}{2} D_{\mathcal{K}}^2 \leq C/9$. 故当 $t = 1, h_1 = f(\mathbf{x}_1) - f(\mathbf{x}^*) \leq C/4$ 成立。

假设对于任意 $t \geq 1$ 有 $h_t \leq \frac{C}{(t+2)^2}$. 由式子(47)得:

$$h_{t+1} \leq \frac{h_t}{2} \leq \frac{C}{2(t+2)^2} = \frac{C}{(t+3)^2} \cdot \frac{(t+3)^2}{2(t+2)^2} \leq \frac{C}{(t+3)^2} \quad (52)$$

如果 $h_t \leq \frac{C}{2(t+3)^2}$. 由式子(52)得, 对任意 $t \geq 1$ 有:

$$h_{t+1} \leq h_t \leq \frac{C}{2(t+2)^2} \leq \frac{C}{(t+3)^2}. \quad (53)$$

若 $h_t \geq \frac{C}{2(t+2)^2}$, 由式子(52)得:

$$h_{t+1} \leq h_t(1 - Mh_t^{1/2}) \quad (54)$$

$$< \frac{C}{(t+2)^2} \left(1 - M\sqrt{\frac{C}{2}} \frac{1}{t+2}\right) \quad (55)$$

$$= \frac{C}{(t+3)^2} \cdot \frac{(t+3)^2}{(t+2)^2} \left(1 - M\sqrt{\frac{C}{2}} \frac{1}{t+2}\right) \quad (56)$$

$$= \frac{C}{(t+3)^2} \cdot \frac{(t+2)^2 + 2t + 5}{(t+2)^2} \left(1 - M\sqrt{\frac{C}{2}} \frac{1}{t+2}\right) \quad (57)$$

$$< \frac{C}{(t+3)^2} \left(1 + \frac{3(t+2)}{(t+3)^2}\right) \left(1 - M\sqrt{\frac{C}{2}} \frac{1}{t+2}\right) \quad (58)$$

$$= \frac{C}{(t+3)^2} \left(1 + \frac{3}{t+3}\right) \left(1 - M\sqrt{\frac{C}{2}} \frac{1}{t+2}\right) \quad (59)$$

$$(60)$$

故对于 $C \geq \frac{18}{M^2}$ 可得:

$$h_{t+1} \leq \frac{C}{(t+2)^2} \left(1 + \frac{3}{t+2}\right) \left(1 - \frac{3}{t+2}\right) < \frac{C}{(t+3)^2} \quad (61)$$

□

参考文献

- [1] Fabian Pedregosa, Notes on the Frank-Wolfe Algorithm, Part I, <http://fa.bianp.net/blog/2018/notes-on-the-frank-wolfe-algorithm-part-i/>
- [2] Dan Garber, Elad Hazan, Faster Rates for the Frank-Wolfe Method over Strongly-Convex Sets