

Fast Frank–Wolfe Algorithms with Adaptive Bregman Step-Size for Weakly Convex Functions

The 14th International Conference on Learning Representations

Shota Takahashi,¹ Sebastian Pokutta,^{2,3} and Akiko Takeda^{1,4}

¹University of Tokyo ²Zuse Institute Berlin ³TU Berlin ⁴Riken

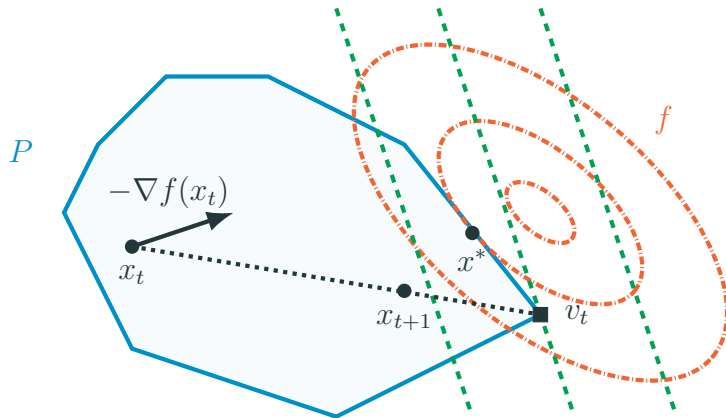
April 25th, 2026

$$\min_{x \in P} f(x)$$

where $P \subset \mathbb{R}^n$ is compact and convex

Fast even if computing projection is not easy

$$v_t \leftarrow \operatorname{argmin}_{v \in P} \langle \nabla f(x_t), v \rangle, \quad x_{t+1} \leftarrow (1 - \gamma_t)x_t + \gamma_t v_t, \quad \gamma_t \in [0, 1]$$



Propose FW algorithm with adaptive Bregman step-size strategy

| FW | Assumptions | | | | Convergence rate | |
|----------|-------------|--------------------|-------------------------|--------------|-----------------------------------|------------------------------------|
| | f convex | f growth | $x^* \in \text{int } P$ | P polytope | L -smooth | L -smad |
| any | convex | ✗ | ✗ | ✗ | $\mathcal{O}(\epsilon^{-1})$ | $\mathcal{O}(\epsilon^{-1/\nu})$ |
| adaptive | convex | q -HEB | ✓ | ✗ | $\mathcal{O}(\log \epsilon^{-1})$ | $\mathcal{O}(\log \epsilon^{-1})$ |
| AFW | convex | q -HEB | ✗ | ✓ | $\mathcal{O}(\log \epsilon^{-1})$ | $\mathcal{O}(\log \epsilon^{-1})$ |
| any | weak | 2-HEB ¹ | ✗ | ✗ | $\mathcal{O}(\epsilon^{-1})$ | $\mathcal{O}(\epsilon^{-1/\nu})$ |
| adaptive | weak | 2-HEB ¹ | ✓ | ✗ | $\mathcal{O}(\log \epsilon^{-1})$ | $\mathcal{O}(\log \epsilon^{-1})$ |
| AFW | weak | 2-HEB ¹ | ✗ | ✓ | $\mathcal{O}(\log \epsilon^{-1})$ | $\mathcal{O}(\log \epsilon^{-1})$ |
| any | ✗ | ✗ | ✗ | ✗ | $\mathcal{O}(\epsilon^{-2})$ | $\mathcal{O}(\epsilon^{-1-1/\nu})$ |

- Nonconv. $\langle \nabla f(x_t), x_t - v_t \rangle \leq \epsilon$; (weakly) convex $f(x_t) - f(x^*) \leq \epsilon$
- HEB: Hölder error bound; AWF: away-step FW algorithm

¹Assume local 2-HEB

Bregman distance: $D_\phi(x, y) := \phi(x) - \phi(y) - \langle \nabla \phi(x), x - y \rangle$; ϕ is strictly convex

L -smad $\stackrel{\text{def.}}{\iff} \exists L > 0$ such that $L\phi - f, L\phi + f$ are convex

$$\iff |f(x) - f(y) - \langle \nabla f(y), x - y \rangle| \leq LD_\phi(x, y)$$

$$\implies f(x) - f((1 - \gamma)x + \gamma v) \geq \underbrace{\gamma \langle \nabla f(x), x - v \rangle - L\gamma^{1+\nu} D_\phi(v, x)}_{\text{max at } \tilde{\gamma}} \quad (1)$$

$$\text{max at } \tilde{\gamma} = \left(\frac{\langle \nabla f(x), x - v \rangle}{L(1 + \nu) D_\phi(v, x)} \right)^{1/\nu}$$

- 1 **Procedure** `step_size`($f, \phi, x, v, \tilde{L}, \gamma_{\max}$)
- 2 Choose $\beta \in (0, 1), \eta \in (0, 1], \tau > 1$
- 3 $\kappa \leftarrow 1, M \leftarrow \eta \tilde{L}$
- 4 **while** $f(x) - f((1 - \gamma)x + \gamma v) < \gamma \langle \nabla f(x), x - v \rangle - M\gamma^{1+\kappa} D_\phi(v, x)$ **do**
- 5 $\gamma \leftarrow \min \left\{ \left(\frac{\langle \nabla f(x), x - v \rangle}{M(1+\kappa) D_\phi(v, x)} \right)^{\frac{1}{\kappa}}, \gamma_{\max} \right\}$
- 6 $M \leftarrow \tau M, \kappa \leftarrow \beta \kappa$
- 7 $\tilde{L}^* \leftarrow M, \tilde{\nu}^* \leftarrow \kappa$
- 8 **return** $\tilde{L}^*, \tilde{\nu}^*, \gamma$

FW zigzags near optimal face; AFW avoids zigzagging [Wolfe, 1970]

- 1 $\mathcal{S}_0 \leftarrow \{x_0\}, \lambda_{x_0,0} \leftarrow 1$
- 2 **for** $t = 0, \dots$ **do**
- 3 $v_t^{\text{FW}} \leftarrow \operatorname{argmin}_{v \in P} \langle \nabla f(x_t), v \rangle$
- 4 $v_t^{\text{A}} \leftarrow \operatorname{argmax}_{v \in \mathcal{S}_t} \langle \nabla f(x_t), v \rangle$
- 5 **if** $\langle \nabla f(x_t), x_t - v_t^{\text{FW}} \rangle \geq \langle \nabla f(x_t), v_t^{\text{A}} - x_t \rangle$ **then**
- 6 $v_t \leftarrow v_t^{\text{FW}}, d_t \leftarrow x_t - v_t^{\text{FW}}, \gamma_{t,\max} \leftarrow 1$ ▷ Frank-Wolfe step
- 7 **else**
- 8 $v_t \leftarrow v_t^{\text{A}}, d_t \leftarrow v_t^{\text{A}} - x_t, \gamma_{t,\max} \leftarrow \frac{\lambda_{v_t^{\text{A}},t}}{1 - \lambda_{v_t^{\text{A}},t}}$ ▷ Away step
- 9 $\gamma_t \leftarrow \min \left\{ \left(\frac{\langle \nabla f(x_t), d_t \rangle}{L(1+\nu)D_\phi(v_t, x_t)} \right)^{1/\nu}, \gamma_{t,\max} \right\}$
- 10 $x_{t+1} \leftarrow x_t - \gamma_t d_t$
- 11 Update the active set \mathcal{S}_t

Assume f convex & ϕ strictly convex & (f, ϕ) is L -smad

Theorem (Linear Convergence (Inner Optimum))

f is q -Hölder error bound (HEB) & $B(x^*, r) \subset P$

$$f(x_t) - f^* \leq \begin{cases} \max\left\{\frac{1}{1+\nu}, 1 - \alpha\right\}^{t-1} LD^2 & \text{if } q = 1 + \nu \\ LD^2 / (1 + \nu)^{t-1} & \text{if } 1 \leq t \leq t_0, q > 1 + \nu \\ \mathcal{O}(1/t^{\nu q / (q-1-\nu)}) & \text{if } t \geq t_0, q > 1 + \nu \end{cases}$$

where $\alpha = \frac{\nu}{1+\nu} \frac{r^{1+1/\nu}}{c^{1+1/\nu} D^{2/\nu}}$, $c = (q/\mu)^{1/q}$, $t_0 \in \mathbb{N}$

Theorem (Linear Convergence (AWF))

f is q -HEB & ϕ strongly convex & P polytope

$$f(x_t) - f^* \leq \begin{cases} \left(1 - \frac{\mu}{32L} \frac{\delta^2}{D^2}\right)^{\lceil (t-1)/2 \rceil} LD^2 & \text{if } (\nu, q) = (1, 2) \\ LD^2 / (1 + \nu)^{\lceil (t-1)/2 \rceil} & \text{if } 1 \leq t \leq t_0, q > 1 + \nu \\ \mathcal{O}(1/t^{\nu q / (q-1-\nu)}) & \text{if } t \geq t_0, q > 1 + \nu \end{cases}$$

$\mathcal{O}(\epsilon^{-2})$ for nonconvex [Lacoste-Julien, 2016]; $\mathcal{O}(\epsilon^{-1-1/\nu})$ for L -smad nonconvex

f is weakly convex & local quadratic growth \implies Local linear convergence

Definition (Weakly Convex & Local Quadratic Growth)

- f is ρ -**weakly convex** $\stackrel{\text{def.}}{\iff} \exists \rho > 0$ such that $f + \frac{\rho}{2} \|\cdot\|^2$ is convex
- f is **local quadratic growth** $\stackrel{\text{def.}}{\iff} \exists \mu > 0$ such that

$$\text{dist}(x, \mathcal{X}^*)^2 \leq \frac{2}{\mu} (f(x) - f^*), \quad \forall x \in [f \leq f^* + \zeta] \cap P \quad (2)$$

$$\implies f(x) - f^* \leq \max_{v \in P} \langle \nabla f(x), x - v \rangle + \underbrace{\frac{\rho}{2} \|x - x^*\|^2}_{\leq \frac{\rho}{\mu} (f(x) - f^*)}$$

$$\therefore \underbrace{\left(1 - \frac{\rho}{\mu}\right) (f(x) - f^*)}_{\text{Primal gap}} \leq \underbrace{\max_{v \in P} \langle \nabla f(x), x - v \rangle}_{\text{FW gap}}$$

f is ρ -weakly convex & local quadratic growth; $M := (L(1 + \nu))^{1/\nu}$

Theorem (Local Linear Convergence (Inner Optimum))

$B(x^*, r) \subset P$ & $\rho/\mu < 1$

$$f(x_t) - f^* \leq \begin{cases} \max \left\{ \frac{1}{2} + \frac{\rho}{2\mu}, 1 - \frac{\nu}{1+\nu} \frac{r^{1+1/\nu}}{Mc^{1+1/\nu}D^{2/\nu}} \right\}^{t-1} LD^2 & \text{if } \nu = 1 \\ \left(\frac{1}{1+\nu} \left(1 + \frac{\nu\rho}{\mu} \right) \right)^{t-1} LD^2 & \text{if } 1 \leq t \leq t_0, \nu \in (0, 1) \\ \mathcal{O}(1/t^{2\nu/(1-\nu)}) & \text{if } t \geq t_0, \nu \in (0, 1) \end{cases}$$

Theorem (Local Linear Convergence (AFW))

ϕ strongly convex & P polytope & $\rho < \mu \leq L$

$$f(x_t) - f^* \leq \begin{cases} \left(1 - \frac{(\mu-\rho)^2}{32\mu L} \frac{\delta^2}{D^2} \right)^{\lceil (t-1)/2 \rceil} LD^2 & \text{if } \nu = 1 \\ 2 \left(\frac{1}{1+\nu} \left(1 + \frac{\nu\rho}{\mu} \right) \right)^{\lceil (t-1)/2 \rceil} LD^2 & \text{if } 1 \leq t \leq t_0, \nu \in (0, 1) \\ \mathcal{O}(1/t^{2\nu/(1-\nu)}) & \text{if } t \geq t_0, \nu \in (0, 1) \end{cases}$$

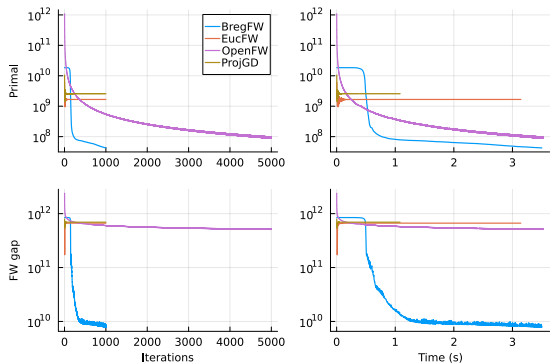


Figure 1: ℓ_p Loss Problem: Gas Sensor Data

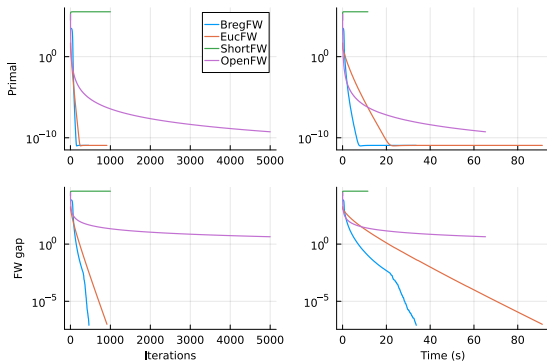


Figure 2: Phase Retrieval

- **BregFW:** Proposed

- **EucFW** [Pedregosa et al., 2020]: search L , then $\gamma_t = \frac{\langle \nabla f(x_t), x_t - v_t \rangle}{L \|x_t - v_t\|^2}$

- **ShortFW:** $\gamma_t = \frac{\langle \nabla f(x_t), x_t - v_t \rangle}{L \|x_t - v_t\|^2}$

- **OpenFW:** $\gamma_t = \frac{2}{2+t}$

- **ProjGD:** Projected gradient descent

[Frank and Wolfe, 1956] Frank, M. and Wolfe, P. (1956).

An algorithm for quadratic programming.

Nav. Res. Logist. Q., 3(1-2):95–110.

[Lacoste-Julien, 2016] Lacoste-Julien, S. (2016).

Convergence rate of Frank-Wolfe for non-convex objectives.

arXiv preprint arXiv:1607.00345.

[Pedregosa et al., 2020] Pedregosa, F., Negiar, G., Askari, A., and Jaggi, M. (2020).

Linearly convergent Frank-Wolfe with backtracking line-search.

In *Proceedings of the 23rd International Conference on Artificial Intelligence and Statistics*, volume 108 of *PMLR*, pages 1–10.

[Wolfe, 1970] Wolfe, P. (1970).

Convergence theory in nonlinear programming.

Integer and nonlinear programming, pages 1–36.