

Lecture note 2

1 RECAP

Boundedness of gradient (L -Lipschitz condition): $\|\nabla f(\mathbf{z})\| \leq L$.

Theorem 1 If $\|\mathbf{x}_1 - \mathbf{x}^*\| \leq R$, $\|\nabla f(\mathbf{x})\| \leq L$ and $\mu = \frac{R}{L\sqrt{t}}$, then we have the following convergence result for gradient descent

$$f\left(\frac{1}{t} \sum_{s=1}^t \mathbf{x}_s\right) - f(\mathbf{x}^*) \leq \frac{RL}{\sqrt{t}} \quad (1)$$

Thus, for an error $\leq \epsilon$, we require $\Theta(1/\epsilon^2)$ iterations.

2 The L -smoothness assumption

One can make a stronger statement on convergence by imposing additional assumptions on the function $f(\mathbf{x})$.

Definition 1 (L -smooth functions) A function $f(\mathbf{x})$ is L -smooth if its gradient is Lipschitz continuous with parameter $L \geq 0$, i.e.

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\| \quad (2)$$

Note that an L -smooth function satisfies the following properties.

Lemma 2 If f is convex and L -smooth, then

$$g(\mathbf{x}) = \frac{L}{2} \mathbf{x}^T \mathbf{x} - f(\mathbf{X}) \quad \text{is convex.} \quad (3)$$

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2. \quad (4)$$

$$\nabla^2 f(\mathbf{x}) \leq L\mathbf{I}. \quad (5)$$

Let us now show that the different implications of the lemma 2 are true assuming that f is convex and its hessian is continuous

1. **Implication 1:** To prove that if f is convex and L -smooth, then \Leftrightarrow (3) From (2), it can be written as $\forall \mathbf{x}, \mathbf{y}$

$$\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), (\mathbf{x} - \mathbf{y}) \rangle \leq L\|\mathbf{x} - \mathbf{y}\|^2 \rightarrow \langle L(\mathbf{x} - \mathbf{y}), (\mathbf{x} - \mathbf{y}) \rangle \quad (6)$$

$$\implies \langle L(\mathbf{x} - \mathbf{y}) - (\nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y}) \rangle \geq 0 \quad (7)$$

Taking the gradient of (3), we get

$$\nabla g(\mathbf{x}) = L\mathbf{x} - \nabla f(\mathbf{x}) \quad (8)$$

Now, (7) becomes

$$\langle \nabla g(\mathbf{x}) - \nabla g(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq 0 \quad (9)$$

Note that (9) indicates the monotonicity of the gradient mapping of g and hence it implies that g is convex

2. Implication 2: To prove that (3) \Leftrightarrow (4) If g is convex,

$$g(\mathbf{y}) \geq g(\mathbf{x}) + \langle \nabla g(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \quad (10)$$

$$\implies \frac{L}{2} \mathbf{y}^T \mathbf{y} - f(\mathbf{y}) \geq \frac{L}{2} \mathbf{x}^T \mathbf{x} - f(\mathbf{x}) + \langle L\mathbf{x} - \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \quad (11)$$

$$\implies f(\mathbf{y}) \leq f(\mathbf{x}) + \frac{L}{2} \mathbf{y}^T \mathbf{y} - \frac{L}{2} \mathbf{x}^T \mathbf{x} + \langle \nabla f(\mathbf{x}) - L\mathbf{x}, \mathbf{y} - \mathbf{x} \rangle \quad (12)$$

By expanding and grouping the terms we get,

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2 \quad (13)$$

which is (4)

3. Implication 3: To prove (4) \Leftrightarrow (5)

From Taylor's theorem and when $\exists t \in (0,1)$,

$$f(\mathbf{y}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{1}{2} (\mathbf{y} - \mathbf{x})^T \nabla^2 f[\mathbf{x} + t(\mathbf{y} - \mathbf{x})] (\mathbf{y} - \mathbf{x}) \quad (14)$$

where,

$$\frac{1}{2} (\mathbf{y} - \mathbf{x})^T \nabla^2 f[\mathbf{x} + t(\mathbf{y} - \mathbf{x})] (\mathbf{y} - \mathbf{x}) \leq \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2 \quad (15)$$

let $\mathbf{y} = \mathbf{x} + \tau v$ for some $v \in \mathbb{R}$ and $0 \leq t \leq \tau$ Therefore (15) \implies

$$\frac{1}{2} \tau^2 v^T \nabla^2 f[\mathbf{x} + t(\mathbf{y} - \mathbf{x})] v \leq \tau^2 \frac{L}{2} \|v\|^2 \quad (16)$$

$\forall v, \mathbf{x}$ if $\tau \rightarrow 0$ then $t \rightarrow 0$

Taking $\lim_{\tau \rightarrow 0}$ we get

$$\frac{1}{2} v^T \nabla^2 f(\mathbf{x}) v \leq \frac{L}{2} \|v\|^2 \quad (17)$$

thus we have

$$0 \leq \frac{1}{2} v^T (LI - \nabla^2 f(\mathbf{x})) v$$

and hence $\nabla^2 f(\mathbf{x}) \leq LI$

4. **Implication 4:** To prove (5) \Leftrightarrow (2) From the property of L-smoothness and Taylor's theorem, we have

$$\nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) = \int_0^1 \nabla^2 f[\mathbf{x} + t(\mathbf{y} - \mathbf{x})](\mathbf{y} - \mathbf{x}) dt \quad (18)$$

and

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 = \left\| \int_0^1 \nabla^2 f[\mathbf{x} + t(\mathbf{y} - \mathbf{x})](\mathbf{y} - \mathbf{x}) dt \right\|_2 \quad (19)$$

Now, using the property, $\|\int g(t) dt\| \leq \int \|g(t)\| dt$,

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq \int_0^1 \|\nabla^2 f[\mathbf{x} + t(\mathbf{y} - \mathbf{x})]\| \|\mathbf{y} - \mathbf{x}\| dt \quad (20)$$

$$\leq L \int_0^1 \|\mathbf{x} - \mathbf{y}\| dt \quad (21)$$

$$\leq L \|\mathbf{y} - \mathbf{x}\| \quad (22)$$

$\Rightarrow f$ is convex and L-smooth

Further, Equation (4) is also known as the *quadratic upper bound* and has the following consequences:

$$\frac{1}{2L} \|\nabla f(\mathbf{x})\|^2 \leq f(\mathbf{x}) - f(\mathbf{x}^*) \leq \frac{L}{2} \|\mathbf{x} - \mathbf{x}^*\|^2 \quad (23)$$

The right hand side follows immediately from (4). Another consequence is the *co-coercivity* of gradient, given as

$$\forall \mathbf{x}, \mathbf{y}, \quad (\nabla f(\mathbf{x}) - \nabla f(\mathbf{y}))^T (\mathbf{x} - \mathbf{y}) \geq \frac{1}{L} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|^2 \quad (24)$$

Under L -smoothness assumption, we now have the following convergence result

Theorem 3 *If f is L -smooth and $0 \leq \mu \leq \frac{1}{L}$, then the gradient descent iterations $\mathbf{x}_{t+1} = \mathbf{x}_t - \mu \nabla f(\mathbf{x}_t)$ admit the following convergence result*

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{1}{2t\mu} \|\mathbf{x}_0 - \mathbf{x}^*\|^2. \quad (25)$$

Proof Using $\mathbf{y} = \mathbf{x} - \mu \nabla f(\mathbf{x})$ in the quadratic upper bound (4), we have

$$\begin{aligned} f(\mathbf{x} - \mu \nabla f(\mathbf{x})) &\leq f(\mathbf{x}) - \mu \left(1 - \frac{\mu L}{2}\right) \|\nabla f(\mathbf{x})\|^2 \\ &\leq f(\mathbf{x}) - \frac{\mu}{2} \|\nabla f(\mathbf{x})\|^2 \quad \left(\text{if } 0 \leq \mu \leq \frac{1}{L}\right) \\ &\leq f(\mathbf{x}^*) + \nabla f(\mathbf{x})^T (\mathbf{x} - \mathbf{x}^*) - \frac{\mu}{2} \|\nabla f(\mathbf{x})\|^2 \quad (\text{tangent lower bound}) \\ &= f(\mathbf{x}^*) + \frac{1}{2\mu} [\|\mathbf{x} - \mathbf{x}^*\|^2 - \|\mathbf{x} - \mathbf{x}^* - \mu \nabla f(\mathbf{x})\|^2]. \end{aligned} \quad (26)$$

Now let $\mathbf{x} = \mathbf{x}_s$ and $\mathbf{x}_{s+1} = \mathbf{x}_s - \mu \nabla f(\mathbf{x}_s)$. Then, from the equation above, we have

$$f(\mathbf{x}_{s+1}) - f(\mathbf{x}^*) \leq \frac{1}{2\mu} [\|\mathbf{x}_s - \mathbf{x}^*\|^2 - \|\mathbf{x}_{s+1} - \mathbf{x}^*\|^2]. \quad (27)$$

Performing a telescopic sum on both sides from $s = 1 \dots t$, we get

$$\begin{aligned} \sum_{s=1}^t [f(\mathbf{x}_s) - f(\mathbf{x}^*)] &\leq \frac{1}{2\mu} [\|\mathbf{x}_0 - \mathbf{x}^*\|^2 - \|\mathbf{x}_t - \mathbf{x}^*\|^2] \\ &\leq \frac{1}{2\mu} \|\mathbf{x}_0 - \mathbf{x}^*\|^2. \end{aligned} \quad (28)$$

Since, the value of $f(\mathbf{x}_s)$ keeps decreasing with every iteration, we have

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{1}{t} \sum_{s=1}^t [f(\mathbf{x}_s) - f(\mathbf{x}^*)] \leq \frac{1}{2t\mu} \|\mathbf{x}_0 - \mathbf{x}^*\|^2, \quad (29)$$

which is the desired result. ■

Note that to get $f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \epsilon$, we need $t = O\left(\frac{1}{\epsilon}\right)$.

2.1 Constrained problems (Projected gradient descent)

In the case of a convex constraint set Ω for the optimization problem, a similar convergence result can be obtained using the following lemma

Lemma 4 *If $\mathbf{x}^+ = \Pi_{\Omega}(\mathbf{x} - \mu \nabla f(\mathbf{x})) := \mathbf{x} - \mu g(\mathbf{x})$, where Π_{Ω} is the projector for Ω and $g(\mathbf{x}) = \frac{1}{\mu}(\mathbf{x} - \mathbf{x}^+)$, then*

$$f(\mathbf{x}^+) \leq f(\mathbf{x}) + g(\mathbf{x})^T(\mathbf{x}^+ - \mathbf{x}) + \frac{L}{2} \|\mathbf{x}^+ - \mathbf{x}\|^2. \quad (30)$$

Proof We know from the quadratic upper bound that

$$f(\mathbf{x}^+) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{x}^+ - \mathbf{x}) + \frac{L}{2} \|\mathbf{x}^+ - \mathbf{x}\|^2. \quad (31)$$

Hence, all we need to show is

$$(\nabla f(\mathbf{x}) - g(\mathbf{x}))^T(\mathbf{x}^+ - \mathbf{x}) \leq 0. \quad (32)$$

Now, note that

$$\nabla f(\mathbf{x}) - g(\mathbf{x}) = \nabla f(\mathbf{x}) - \frac{1}{\mu}(\mathbf{x} - \mathbf{x}^+) = \frac{1}{\mu} [\mathbf{x}^+ - (\mathbf{x} - \mu \nabla f(\mathbf{x}))]. \quad (33)$$

Therefore, using the fact that $(\Pi_{\Omega}(\mathbf{x}) - \mathbf{x})^T(\Pi_{\Omega}(\mathbf{x}) - \mathbf{x}^*) \leq 0$ for any convex set Ω , we get the desired result. ■

One can now use the lemma above at the appropriate step in Theorem 3 to get a convergence result for convex constrained problems.

2.2 How to pick μ in practice?

The convergence results so far require that the Lipschitz constant L be known in order to set a value for μ , which is not possible in reality. One way to deal with this problem is *Backtracking line search*, that automatically sets the value of μ in each iteration. To move along the search direction, we start with an estimate of the step size (in this case $\mu = 1$) and then iteratively shrink it (backtracking) until a decrease in the objective function is observed. The least amount of descent expected or the termination condition is given by equation 17,

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \gamma\mu\|\nabla f(\mathbf{x}_t)\|^2 \quad (34)$$

and a much greater descent would look like:

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \mu(1 - \frac{\mu L}{2})\|\nabla f(\mathbf{x}_t)\|^2 \quad (35)$$

Equation (18) follows from the quadratic upper bound in (4) which we always have at any step t . The procedure is given below:

Backtracking line search

1. Pick descent direction $\mathbf{v} = -\nabla f(\mathbf{x}_t)$.
2. Start with $\mu = 1$.
3. Test Armijo condition

$$f(\mathbf{x}_t - \mu\nabla f(\mathbf{x}_t)) \leq f(\mathbf{x}_t) - \gamma\mu\|\nabla f(\mathbf{x}_t)\|^2. \quad (36)$$

- (a) If true, keep μ .
- (b) If false, update μ as $\mu \leftarrow \beta\mu$.

Rule of thumb: $\gamma = 0.5$ and $\beta = 0.8$.

This procedure works because we change μ in each iteration to keep it as close to $\frac{1}{L}$ as possible. The following convergence result can be obtained for such a procedure

Theorem 5 *For an L -smooth f , and gradient descent iterations given by $\mathbf{x}_{t+1} = \mathbf{x}_t - \mu_t \nabla f(\mathbf{x}_t)$, we have*

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{1}{2t\mu_{\min}} \|\mathbf{x}_0 - \mathbf{x}^*\|^2, \quad (37)$$

where

$$\mu_{\min} := \min_{1 \leq s \leq t} \mu_s \geq \min \left(1, \frac{\beta}{L} \right) \quad (38)$$

Proof Equation (37) follows directly from Theorem 3 due to definition of μ_{\min} as the minimum of all μ_s . Hence, we just need to show (38) which essentially means μ_{\min} is close to $\frac{1}{L}$ or alternatively stated - we want to find the value of μ at which the algorithm terminates.

Suppose $\mu = 1$. If the Armijo condition succeeds, we have from equation (17):

$$\|\nabla f(\mathbf{x}_t)\|^2 \leq \frac{1}{\gamma}(f(\mathbf{x}_t) - f(\mathbf{x}_{t+1})) \quad (39)$$

then we keep μ as this already indicates some kind of bound on the gradient.

But if it fails, we stop the iteration and update μ and start over again. This means at the t^{th} iteration, μ_t/β fails. For such an iteration, we have from the Armijo condition

$$f(\mathbf{x}_t) - \frac{\gamma\mu_t}{\beta}\|\nabla f(\mathbf{x}_t)\|^2 \leq f\left(\mathbf{x}_t - \frac{\mu_t}{\beta}\nabla f(\mathbf{x}_t)\right) \quad (40)$$

$$\leq f(\mathbf{x}_t) - \frac{\mu_t}{\beta}\left(1 - \frac{\mu_t L}{2}\right)\|\nabla f(\mathbf{x}_t)\|^2 \quad (41)$$

where we arrive at the second expression by using the quadratic upper bound (4) for the right hand side. What we essentially mean by (23) and (24) is that, when μ_t/β fails, the Armijo condition is violated or the the direction of descent is reversed which can be seen as $f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \gamma\mu\|\nabla f(\mathbf{x}_t)\|^2$.

On further simplification of equation (24) we get,

$$\begin{aligned} -\frac{\gamma\mu_t}{2\beta} &\leq -\frac{\mu_t}{\beta}\left(1 - \frac{\mu_t L}{2}\right) \\ \implies \mu_t &\geq \frac{2(1-\gamma)\beta}{L} \end{aligned}$$

Substituting $\gamma = 0.5$ and $\beta = 0.8$ we get,

$$\mu_t \geq \frac{\beta}{L} \approx \frac{0.8}{L}$$

which means in each iteration we have a guarantee that $\mu_{\min} = \min\left(1, \frac{\beta}{L}\right)$. ■

Note that for an L -smooth function, the Armijo condition cannot be satisfied if $\mu \geq \frac{1}{L}$. Hence, backtracking line search ensures that we not only satisfy it at each iteration, but also pick the largest μ that does it.

2.3 Conditional gradient descent

For constrained optimization problems of the form $\min_{\mathbf{x} \in \Omega} f(\mathbf{x})$, there exists another method:

Franke-Wolfe algorithm

$$\mathbf{y}_t = \underset{\mathbf{y} \in \Omega}{\operatorname{argmin}} \nabla f(\mathbf{x}_t)^T \mathbf{y} \quad (42)$$

$$\mathbf{x}_{t+1} = (1 - \mu_t)\mathbf{x}_t + \mu_t \mathbf{y}_t \quad (43)$$

$$\mu_t = \frac{2}{t+1} \quad (44)$$

Note that \mathbf{y}_t is the best descent direction. The method above has the advantage that it is projection-free, norm-free and has lower space complexity.

The motivation for this algorithm is that: usual gradient methods involve a projection onto the constrained set Ω in each iteration which can be very expensive. Whereas the Frank-Wolfe method only requires the solution to the linear approximation of the function over the constrained set in each iteration. For instance, consider the local linear approximation of $f(\mathbf{x})$,

$$\hat{f}^{lin} = f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^T(\mathbf{y} - \mathbf{x}_t) \quad (45)$$

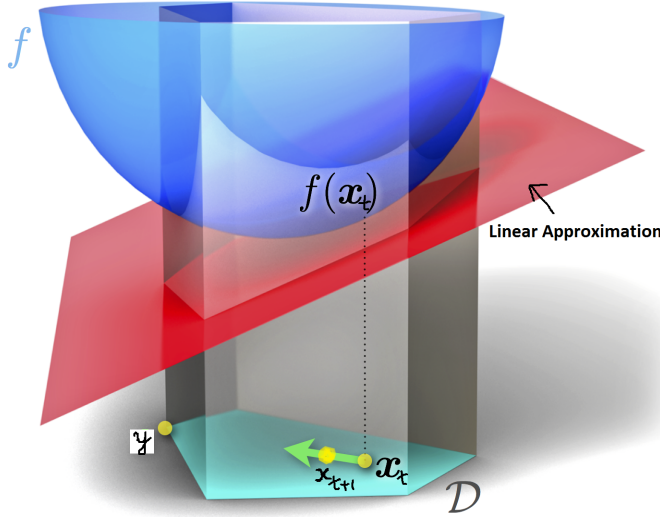
Trying to minimize it over \mathbf{y} , would just send the linear approximation to $-\infty$ and cannot project on the constrained set for the same reason as it is not well defined. Therefore Frank-Wolfe minimizes this over the whole constrained set.

$$\mathbf{y}_t = \arg \min_{\mathbf{y} \in \Omega} \hat{f}^{lin} \quad (46)$$

If \mathbf{y}_t is the minimizer, we move a bit in the direction of the minimizer and call it say \mathbf{x}_{t+1}

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \mu_t(\mathbf{y}_t - \mathbf{x}_t) \quad (47)$$

It can be clearly seen from the following figure that the direction of descent is $(\mathbf{y}_t - \mathbf{x}_t)$ and multiply it by some step size μ_t and move along that direction. Equation (48) is just restatement of (44) with little rearrangement. But it can be seen from (48) that, as t increases μ_t decreases and we move less and less aggressively in the direction of the linearization minimizer.



The following convergence result can be obtained for the Frank-Wolfe algorithm:

Theorem 6 Let f be convex and L -smooth with arbitrary norms, i.e., $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_* \leq L\|\mathbf{x} - \mathbf{y}\|$ ($\|\cdot\|$ and $\|\cdot\|_*$ are dual norms), and let $R = \sup_{\mathbf{x}, \mathbf{y} \in \Omega} \|\mathbf{x} - \mathbf{y}\|$, $\mu_t = \frac{2}{t+1}$, then, we have

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{2LR^2}{t+1} \quad (48)$$

Proof From equation (4) we can write the quadratic bound over x_t as:

$$f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) \leq \nabla f(\mathbf{x}_t)^T (\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{L}{2} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \quad (49)$$

from (26), we have

$$\mathbf{x}_{t+1} - \mathbf{x}_t = \mu_t (\mathbf{y}_t - \mathbf{x}_t) \quad (50)$$

Therefore,

$$f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) \leq \mu_t \langle \nabla f(\mathbf{x}_t), (\mathbf{y}_t - \mathbf{x}_t) \rangle + \frac{L\mu_t^2 R^2}{2} \quad (51)$$

$$\leq \mu_t \langle \nabla f(\mathbf{x}_t), (\mathbf{x}^* - \mathbf{x}_t) \rangle + \frac{L\mu_t^2 R^2}{2} \quad (52)$$

where \mathbf{x}^* is the optimal solution

And also we have from the fact that a function is above its first order approximation,

$$\langle \nabla f(\mathbf{x}_t), (\mathbf{x}_* - \mathbf{x}_t) \rangle \leq f(\mathbf{x}^*) - f(\mathbf{x}_t) \quad (53)$$

$$\implies f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t) \leq \mu_t [f(\mathbf{x}^*) - f(\mathbf{x}_t)] + \frac{L\mu_t^2 R^2}{2} \quad (54)$$

$$\implies f(\mathbf{x}_{t+1}) - f(\mathbf{x}^*) - [f(\mathbf{x}_t) - f(\mathbf{x}^*)] \leq -\mu_t [f(\mathbf{x}_t) - f(\mathbf{x}^*)] + \frac{L\mu_t^2 R^2}{2} \quad (55)$$

$$\implies f(\mathbf{x}_{t+1}) - f(\mathbf{x}^*) \leq (1 - \mu_t) [f(\mathbf{x}_t) - f(\mathbf{x}^*)] + \frac{L\mu_t^2 R^2}{2} \quad (56)$$

Now, let δ indicate the error and substituting for μ_t , (36) becomes

$$\delta_{t+1} \leq \frac{t-1}{t+1} \delta_t + \frac{L\mu_t^2 R^2}{2} \quad (57)$$

Also (37) indicates that the error decreases as the number of iterations increase. Further, when $t = 0$,

$$\delta_1 = f(\mathbf{x}_1) - f(\mathbf{x}^*) \leq \frac{L}{2} \|\mathbf{x}_1 - \mathbf{x}^*\|^2 \quad (58)$$

From one of the consequences of the quadratic upper bound, $f(\mathbf{x}_1)$ can be bounded as in (38)

$$\implies \delta_1 \leq \frac{L}{2} R^2 \quad (59)$$

Continuing like this, by induction we can prove that

$$\delta_t \leq \frac{2}{t+1} L R^2 \quad (60)$$

ie.,

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{2L}{t+1} R^2 \quad (61)$$

■

Note that there is no projection, update is solved directly over the constrained set Ω .
Let us consider an example of using the Frank-Wolf algorithm.

Example:

$$\begin{aligned} \min f(\mathbf{x}) &= \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_{\ell_2}^2 \\ \text{subject to } \|\mathbf{x}\|_{\ell_1} &\leq \tau \end{aligned} \quad (62)$$

Note that after rescaling (dividing by τ), the above problem is equivalent to

$$\begin{aligned} \min f(\mathbf{x}) &= \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_{\ell_2}^2 \\ \text{subject to } \|\mathbf{x}\|_{\ell_1} &\leq 1 \end{aligned} \quad (63)$$

Solution: Consider $f(\mathbf{x}) = \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_{\ell_2}^2$, and let $\nabla f(\mathbf{x}) = \mathbf{A}^T(\mathbf{A}\mathbf{x} - \mathbf{b})$. Then, the update is of the form

$$\mathbf{y}_t = \operatorname{argmin}_{\|\mathbf{y}\|_{\ell_1} \leq 1} \nabla f(\mathbf{x})^T \mathbf{y}.$$

Let $i_{\max} = \arg \max_i |\nabla f(\mathbf{x})_i|$ denote the index of the entry of $\nabla f(\mathbf{x})$ with the largest absolute value. Then the update is of the form

$$\mathbf{y}_t = -\operatorname{sign}([\nabla f(\mathbf{x})]_{i_{\max}}) \mathbf{e}_{i_{\max}}, \quad (64)$$

where \mathbf{e}_i is the standard basis vector with all entries zeros except a 1 at entry i . For example, if $\nabla f(\mathbf{x}) = [2, -3, 1]^T$, then, $\mathbf{y}_t = [0, 1, 0]^T$. Hence, for the ℓ_1 -norm constraint, we clearly see how this method is projection-free. If we start with $\mathbf{x}_0 = \mathbf{0}$, then after t iterations, we have t non-zero entries. It is also norm-free because the L -smoothness condition for $f(\mathbf{x}) = \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_{\ell_2}^2$, i.e.,

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_{\ell_\infty} = \|\mathbf{A}^T \mathbf{A}(\mathbf{x} - \mathbf{y})\|_{\ell_\infty} \leq L \|\mathbf{x} - \mathbf{y}\|_{\ell_1}, \quad (65)$$

where $L = \max_i \|\mathbf{A}_i\|_{\ell_2}^2$. In comparison, for the ℓ_2 norm, the same condition

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_{\ell_2} = \|\mathbf{A}^T \mathbf{A}(\mathbf{x} - \mathbf{y})\|_{\ell_2} \leq \tilde{L} \|\mathbf{x} - \mathbf{y}\|_{\ell_2}, \quad (66)$$

holds for $\tilde{L} = \|\mathbf{A}^T \mathbf{A}\|$, which can be much greater than L .