

Optimization

Basic knowledge

1. Subgradients

Definition 1.2 (Subgradients). Let $\mathcal{X} \subset \mathbb{R}^n$, and $f : \mathcal{X} \rightarrow \mathbb{R}$. Then $g \in \mathbb{R}^n$ is a subgradient of f at $x \in \mathcal{X}$ if for any $y \in \mathcal{X}$ one has

$$f(x) - f(y) \leq g^\top (x - y).$$

The set of subgradients of f at x is denoted $\partial f(x)$.

Set of all subgradients of **convex** f is called the subdifferential:

$$\partial f(x) = \{g \in \mathbb{R}^n : g \text{ is a subgradient of } f \text{ at } x\}.$$

- $\partial f(x)$ is closed and convex (even for non-convex function)
- Nonempty (can be empty for non-convex function)
- If f is differential at x , then $\partial f(x) = \{\nabla f(x)\}$
- If $\partial f(x) = \{g\}$, then f is differential at x and $\nabla f(x) = g$

Convex set $C \subseteq \mathbb{R}^n$, consider indicator function $I_C : \mathbb{R}^n \rightarrow \mathbb{R}$,

$$I_C(x) = I\{x \in C\} = \begin{cases} 0 & \text{if } x \in C \\ \infty & \text{if } x \notin C \end{cases}$$

For $x \in C$, $\partial I_C(x) = \mathcal{N}_C(x)$, the **normal cone** of C at x , recall

$$\mathcal{N}_C(x) = \{g \in \mathbb{R}^n : g^T x \geq g^T y \text{ for any } y \in C\}$$

First order optimality

For problem

$$\min f(x) \text{ subject to } x \in C$$

the condition $0 \in \partial f(x) + \mathcal{N}_C(x)$ is a fully general condition for optimality in a convex problem.

Some examples

A. The distance to a convex set

Distance function to a convex set C is

$$\text{dist}(x, C) = \|x - P_C(x)\|_2,$$

then when $\text{dist}(x, C) > 0$, $\partial \text{dist}(x, C) = \left\{ \frac{x - P_C(x)}{\|x - P_C(x)\|_2} \right\}$.

B. ℓ_0 norm

$$f(x) = \|x\|_0 = \sum_{i=1}^n \mathbf{1}\{x_i \neq 0\}$$

ℓ_0 norm function has subgradient at exactly one point: the origin, where $\partial f(0) = 0$. Everywhere else, $\partial f(x) = \emptyset$.

C. ℓ_∞ norm

$$f(\mathbf{x}) = \|\mathbf{x}\|_\infty = \max_{1 \leq i \leq n} |x_i|$$

If $x \neq 0$, then pick any x_j obeying $|x_j| = \max_i |x_i|$ to get $\text{sign}(x_j)e_j \in \partial f(x)$

D. Maximum eigenvalue

$$f(\mathbf{x}) = \lambda_{\max}(x_1 A_1 + \dots + x_n A_n), \text{ where } A_1, \dots, A_n \text{ are real symmetric matrices.}$$

Rewrite $f(\mathbf{x}) = \sup_{\mathbf{y}: \|\mathbf{y}\|_2=1} \mathbf{y}^T (x_1 A_1 + \dots + x_n A_n) \mathbf{y}$ as supremum of affine functions of \mathbf{x} .

Therefore, taking \mathbf{y} as leading eigenvector of $x_1 A_1 + \dots + x_n A_n$, we have

$$[\mathbf{y}^T A_1 \mathbf{y}, \dots, \mathbf{y}^T A_n \mathbf{y}] \in \partial f(x).$$

2. Oracle complexity

In the black-box model one tries to minimize an unknown function f over a constraint set $\mathcal{X} \subset \mathbb{R}^n$. A black-box optimization procedure is simply a sequence of mappings $\phi_t : \mathcal{X}^t \times (\mathbb{R}^n)^t \times \mathbb{R}^t \rightarrow \mathcal{X}$. The algorithm given by these mappings runs iteratively as follows: initially it makes a query to the point $x_0 = \phi_0(\emptyset)$, and the t^{th} step it queries

$$x_t = \phi_t \left(x_0, \dots, x_{t-1}, \nabla f(x_0), \dots, \nabla f(x_{t-1}), f(x_0), \dots, f(x_{t-1}) \right).$$

The minimax oracle optimization error after t steps over a set of functions \mathcal{F} , is defined as follows

$$\text{OC}_t(\mathcal{F}) = \inf_{\phi_0, \dots, \phi_t \in \mathcal{F}} \left(f(x_t) - \inf_{x \in \mathcal{X}} f(x) \right).$$

Intuitively $\text{OC}_t(\mathcal{F})$ describes **the best possible rate of convergence** for the optimization error when one restricts to black-box procedures that only know the constraint set \mathcal{X} and the class of functions \mathcal{F} .

3. The conjugate function

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$. The function $f^* : \mathbb{R}^n \rightarrow \mathbb{R}$, defined as

$$f^*(x) := \sup_{z \in \text{dom } f} \{ \langle z, x \rangle - f(z) \}$$

is called the conjugate function of f .

Fenchel's inequality

$$x'y \leq f(x) + f^*(y), \quad \forall x \in \mathbb{R}^n, y \in \mathbb{R}^n.$$

The second conjugate

$$f^{**}(x) = \sup_{y \in \text{dom } f^*} (x^T y - f^*(y)).$$

- $f^{**}(x)$ is closed and convex
- From Fenchel's inequality, $x^T y - f^*(y) \leq f(x)$. Hence, $f^{**}(x) \leq f(x)$
- If f is closed and convex, then $f^{**}(x) = f(x)$ for all x ; Equivalently, $\text{epi } f^{**} = \text{epi } f$.

Conjugate Subgradient Theorem

If f is closed and convex, then the following relations are equivalent for a pair of vectors (x, y) :

$$(i) \quad x'y = f(x) + f^*(y)$$

$$(ii) \quad y \in \partial f(x)$$

$$(iii) \quad x \in \partial f^*(y)$$

Proof. if $y \in \partial f(x)$, then $f^*(y) = \sup_u (y^T u - f(u)) = y^T x - f(x)$; hence

$$\begin{aligned} f^*(v) &= \sup_u (v^T u - f(u)) \\ &\geq v^T x - f(x) \\ &= x^T(v - y) - f(x) + y^T x \\ &= f^*(y) + x^T(v - y) \end{aligned}$$

this holds for all v ; therefore, $x \in \partial f^*(y)$

reverse implication $x \in \partial f^*(y) \implies y \in \partial f(x)$ follows from $f^{**} = f$

Examples

- **Indicator function:** if $f(x) = I_C(x)$, then its conjugate function is support function of C

$$f^*(y) = I_C^*(y) = \max_{x \in C} y^T x$$

- **Norm:** if $f(x) = \|x\|$, then its conjugate is $f^*(y) = I_{\{z: \|z\|_* \leq 1\}}(y)$, where $\|\cdot\|_*$ is dual norm of $\|\cdot\|$

Why? Note that if $\|y\|_* > 1$, then there exists $\|z\| \leq 1$ with $z^T y = \|y\|_* > 1$, so

$$(tz)^T y - \|tz\| = t(z^T y - \|z\|) \rightarrow \infty, \text{ as } t \rightarrow \infty$$

i.e., $f^*(y) = \infty$

On the other hand, if $\|y\|_* \leq 1$, then

$$z^T y - \|z\| \leq \|z\| \|y\|_* - \|z\| \leq 0$$

and = 0 when $z = 0$, so $f^*(y) = 0$

- **Matrix Logarithm:** if $f(X) = -\log \det X$ ($\text{dom } f = \mathbf{S}_{++}^n$), then its conjugate is $f^*(Y) = -\log \det(-Y) - n$ (From B & V page 92)

4. Operations that preserve convexity

- I. Nonnegative weighted sums
- II. Composition with an affine mapping
- III. Point-wise maximum and supremum

5. Log-concave and Log-convex functions

Definition. A function $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$, with convex domain and $f(x) > 0$ for all $x \in \text{dom } f$, is log-convex for all $x, y \in \text{dom } f$ and $0 \leq \theta \leq 1$, we have

$$f(\theta x + (1 - \theta)y) \leq f(x)^\theta f(y)^{1-\theta}.$$

Twice differentiable condition: $f(x) \nabla^2 f(x) \leq \nabla f(x) \nabla f(x)^T$

- The CDF of a Gaussian distribution is Log-concave.
- Gamma function is Log-convex
- The Wishart density is Log-concave

6. Convex Optimization Problems

A convex optimization problem is

$$\begin{aligned} & \min_{x \in D} f(x) \\ & \text{subject to } g_i(x) \leq 0, i = 1, \dots, m \\ & Ax = b \end{aligned}$$

if criterion f is strictly convex, then the solution is unique.

The problem can be rewritten as

$$\min f(x) \text{ subject to } x \in C$$

where C is feasible set.

For a convex problem and differentiable f , the first order optimal condition is

$$\nabla f(x)^T(y - x) \geq 0 \quad \text{for all } y \in C.$$

E.g., if we decompose $x = (x_1, x_2) \in \mathbb{R}^{n_1+n_2}$, then

$$\begin{array}{ll} \min_{x_1, x_2} & f(x_1, x_2) \\ \text{s.t.} & g_1(x_1) \leq 0 \\ & g_2(x_2) \leq 0 \end{array} \iff \begin{array}{ll} \min_{x_1} & \tilde{f}(x_1) \\ \text{s.t.} & g_1(x_1) \leq 0 \end{array}$$

where $\tilde{f}(x_1) = \min\{f(x_1, x_2) : g_2(x_2) \leq 0\}$. The right problem is convex if the left problem is

Partial Optimization

$g(x) = \min_{y \in C} f(x, y)$ is convex in x , provided that f is convex in (x, y) and C is a convex set.

Hinge form of SVMs

Recall the SVM problem

$$\begin{array}{ll} \min_{\beta, \beta_0} & \frac{1}{2} \|\beta\|_2^2 + C \sum_{i=1}^n \xi_i \\ \text{subject to} & \xi_i \geq 0, \quad y_i(x_i^T \beta + \beta_0) \geq 1 - \xi_i, \quad i = 1, \dots, n \end{array}$$

Rewrite the constraints as $\xi_i \geq \max\{0, 1 - y_i(x_i^T \beta + \beta_0)\}$. Indeed we can argue that we have at solution

Therefore plugging in for optimal ξ gives the hinge form of SVMs:

$$\min_{\beta, \beta_0} \frac{1}{2} \|\beta\|_2^2 + C \sum_{i=1}^n [1 - y_i(x_i^T \beta + \beta_0)]_+$$

where $a_+ = \max\{0, a\}$ is called the hinge function

Transformations and change of variables

If $h : \mathbb{R}^n \rightarrow \mathbb{R}$ is a monotone increasing transformation, then $\min f(x)$ is equivalent to $\min f(h(x))$.

If $\phi : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is one-to-one, and its image covers feasible set C , then we can change variable in an optimization problem: $\min f(x)$ is equivalent to $\min f(\phi(y))$ for $\phi(y) \in C$.

Gradient descent for unconstraint problem

Gradient descent method: $\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \nabla f(\mathbf{x}^t)$

1. Quadratic optimization

Consider the following problem

$$\text{minimize }_{\mathbf{x}} \quad f(\mathbf{x}) := \frac{1}{2} (\mathbf{x} - \mathbf{x}^*)^\top \mathbf{Q} (\mathbf{x} - \mathbf{x}^*),$$

for some $n \times n$ matrix $\mathbf{Q} > 0$, where $\nabla f(\mathbf{x}) = \mathbf{Q}(\mathbf{x} - \mathbf{x}^*)$.

A. Constant stepwise rule

If $\eta_t \equiv \eta = \frac{2}{\lambda_1(\mathbf{Q}) + \lambda_n(\mathbf{Q})}$, then

$$\| \mathbf{x}^t - \mathbf{x}^* \|_2 \leq \left(\frac{\lambda_1(\mathbf{Q}) - \lambda_n(\mathbf{Q})}{\lambda_1(\mathbf{Q}) + \lambda_n(\mathbf{Q})} \right)^t \| \mathbf{x}^0 - \mathbf{x}^* \|_2.$$

B. Exact line search

If $\eta_t = \arg \min_{\eta \geq 0} f(\mathbf{x}^t - \eta \nabla f(\mathbf{x}^t))$, then

$$f(\mathbf{x}^t) - f(\mathbf{x}^*) \leq \left(\frac{\lambda_1(\mathbf{Q}) - \lambda_n(\mathbf{Q})}{\lambda_1(\mathbf{Q}) + \lambda_n(\mathbf{Q})} \right)^{2t} (f(\mathbf{x}^0) - f(\mathbf{x}^*)).$$

2. Strongly convex and smooth problems

Strongly convex

I. $f(\mathbf{y}) \geq \underbrace{f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})}_{\text{first-order Taylor expansion}} + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}\|_2^2, \quad \forall \mathbf{x}, \mathbf{y}$

II. For all \mathbf{x} and \mathbf{y} and all $0 \leq \lambda \leq 1$,

$$f(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y}) - \frac{\mu}{2}\lambda(1 - \lambda)\|\mathbf{x} - \mathbf{y}\|_2^2.$$

III. $\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq \mu\|\mathbf{x} - \mathbf{y}\|_2^2, \quad \forall \mathbf{x}, \mathbf{y}$

L - smooth

I. $f(\mathbf{y}) \leq \underbrace{f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})}_{\text{first-order Taylor expansion}} + \frac{L}{2}\|\mathbf{x} - \mathbf{y}\|_2^2, \quad \forall \mathbf{x}, \mathbf{y}.$

II. For all \mathbf{x} and \mathbf{y} and all $0 \leq \lambda \leq 1$,

$$f(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) \geq \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y}) - \frac{L}{2}\lambda(1 - \lambda)\|\mathbf{x} - \mathbf{y}\|_2^2.$$

III. $\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq \frac{1}{L}\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2, \quad \forall \mathbf{x}, \mathbf{y}.$

IV. $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq L\|\mathbf{x} - \mathbf{y}\|_2, \forall \mathbf{x}, \mathbf{y}. \text{ (L - Lipschitz gradient)}$

Theorem 2.0 Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a m -strong convex and L -smooth function. Then for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$,

$$\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq \frac{mL}{m + L}\|\mathbf{x} - \mathbf{y}\|_2^2 + \frac{1}{m + L}\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2.$$

Theorem 2.1 (GD for strongly convex and smooth functions)

Let f be μ -strongly convex and L -smooth. If $\eta_t \equiv \eta = \frac{2}{\mu + L}$, then

$$\|\mathbf{x}^t - \mathbf{x}^*\|_2 \leq \left(\frac{\kappa - 1}{\kappa + 1}\right)^t \|\mathbf{x}^0 - \mathbf{x}^*\|_2,$$

where $\kappa := L/\mu$ is condition number; \mathbf{x}^* is minimizer

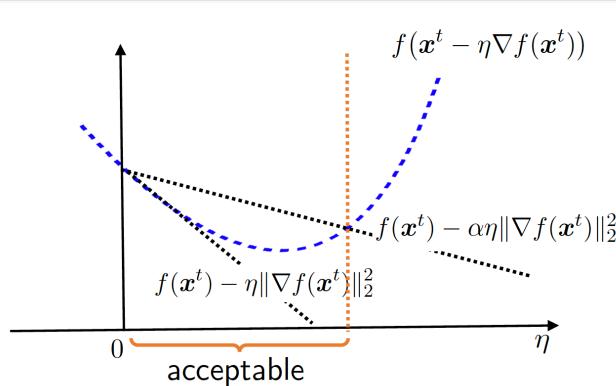
Remark.

- Dimension-free: iteration complexity is $O\left(\frac{\log \frac{1}{\epsilon}}{\log \frac{\kappa+1}{\kappa-1}}\right)$, which is independent of problem size n if κ does not depend on n .

- Direct consequence from smoothness

$$f(\mathbf{x}^t) - f(\mathbf{x}^*) \leq \frac{L}{2} \left(\frac{\kappa-1}{\kappa+1} \right)^{2t} \| \mathbf{x}^0 - \mathbf{x}^* \|_2^2.$$

Backtracking line search



Algorithm 2.2 Backtracking line search for GD

```

1: Initialize  $\eta = 1$ ,  $0 < \alpha \leq 1/2$ ,  $0 < \beta < 1$ 
2: while  $f(\mathbf{x}^t - \eta \nabla f(\mathbf{x}^t)) > f(\mathbf{x}^t) - \alpha \eta \|\nabla f(\mathbf{x}^t)\|_2^2$  do
3:    $\eta \leftarrow \beta \eta$ 

```

Remark.

- Armijo condition: for some $0 < \alpha < 1$,

$$f(\mathbf{x}^t - \eta \nabla f(\mathbf{x}^t)) < f(\mathbf{x}^t) - \alpha \eta \|\nabla f(\mathbf{x}^t)\|_2^2.$$

- Ensures **sufficient decrease** in objective function
- Practically, backtracking line search often provides good estimate on local Lipschitz constant of gradients.

3. Other conditions

Strong convexity requirement can be relaxed

- Local strong convexity

Theorem 2.3 (GD for locally strongly convex and smooth functions)

Let f be *locally* μ -strongly convex and L -smooth such that

$$\mu \mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq L \mathbf{I}, \quad \forall \mathbf{x} \in \mathcal{B}_0$$

where $\mathcal{B}_0 := \{\mathbf{x} : \|\mathbf{x} - \mathbf{x}^*\|_2 \leq \|\mathbf{x}^0 - \mathbf{x}^*\|_2\}$. Then Theorem 2.1 continues to hold

- Regularity condition

Another way to replace strong convexity and smoothness is regularity condition:

$$\langle \nabla f(\mathbf{x}), \mathbf{x} - \mathbf{x}^* \rangle \geq \frac{\mu}{2} \|\mathbf{x} - \mathbf{x}^*\|_2^2 + \frac{1}{2L} \|\nabla f(\mathbf{x})\|_2^2, \quad \forall \mathbf{x}.$$

Note that $\nabla f(\mathbf{x}^*) = 0$, thus

$$\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{x}^*), \mathbf{x} - \mathbf{x}^* \rangle \geq \frac{\mu}{2} \|\mathbf{x} - \mathbf{x}^*\|_2^2 + \frac{1}{2L} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}^*)\|_2^2, \quad \forall \mathbf{x}.$$

Theorem 2.4

Suppose f satisfies (2.7). If $\eta_t \equiv \eta = \frac{1}{L}$, then

$$\|\mathbf{x}^t - \mathbf{x}^*\|_2^2 \leq \left(1 - \frac{\mu}{L}\right)^t \|\mathbf{x}^0 - \mathbf{x}^*\|_2^2$$

- Polyak-Łojasiewicz condition

$$\|\nabla f(\mathbf{x})\|_2^2 \geq 2\mu (f(\mathbf{x}) - f(\mathbf{x}^*)), \quad \forall \mathbf{x}.$$

Theorem 2.5

Suppose f satisfies (2.8). If $\eta_t \equiv \eta = \frac{1}{L}$, then

$$f(\mathbf{x}^t) - f(\mathbf{x}^*) \leq \left(1 - \frac{\mu}{L}\right)^t (f(\mathbf{x}^0) - f(\mathbf{x}^*))$$

4. Dropping strong convexity

$f(x)$ is convex and L-smooth Key idea: **majorization-minimization**

Find a simple majoring function of $f(x)$ and optimize it. From smooth assumption,

$$\begin{aligned} f(\mathbf{x}^{t+1}) - f(\mathbf{x}^t) &\leq \nabla f(\mathbf{x}^t)^\top (\mathbf{x}^{t+1} - \mathbf{x}^t) + \frac{L}{2} \|\mathbf{x}^{t+1} - \mathbf{x}^t\|_2^2 \\ &= -\eta_t \|\nabla f(\mathbf{x}^t)\|_2^2 + \frac{\eta_t^2 L}{2} \|\nabla f(\mathbf{x}^t)\|_2^2. \\ &\leq -\frac{1}{2L} \|\nabla f(\mathbf{x}^t)\|_2^2 \end{aligned}$$

Fact 2.8

Let f be convex and L-smooth. If $\eta_t \equiv \eta = 1/L$, then

$$\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2 \leq \|\mathbf{x}^t - \mathbf{x}^*\|_2 - \frac{1}{L^2} \|\nabla f(\mathbf{x}^t)\|_2^2$$

where \mathbf{x}^* is any minimizer with optimal $f(\mathbf{x}^*)$

One can further show $\|\mathbf{x}^t - \mathbf{x}^*\|_2$ is strictly decreasing unless \mathbf{x}^* is already minimizer.

Theorem 2.9 (GD for convex and smooth problems)

Let f be convex and L -smooth. If $\eta_t \equiv \eta = 1/L$, then GD obeys

$$f(\mathbf{x}^t) - f(\mathbf{x}^*) \leq \frac{2L\|\mathbf{x}^0 - \mathbf{x}^*\|_2^2}{t}$$

where \mathbf{x}^* is any minimizer with optimal $f(\mathbf{x}^*)$

Generalized Steepest descent

Let p, q be complementary (dual): $1/p + 1/q = 1$

Steepest descent updates are $x^+ = x + t \cdot \Delta x$, where

$$\begin{aligned}\Delta x &= \|\nabla f(x)\|_q \\ u &= \operatorname{argmin}_{\|v\|_p \leq 1} \nabla f(x)^T v\end{aligned}$$

- If $p = 2$, then $\Delta x = -\nabla f(x)$, gradient descent
- If $p = 1$, then $\Delta x = -\partial f(x)/\partial x_i \cdot e_i$, where

$$\left| \frac{\partial f}{\partial x_i}(x) \right| = \max_{j=1,\dots,n} \left| \frac{\partial f}{\partial x_j}(x) \right| = \|\nabla f(x)\|_\infty.$$

Normalized steepest descent just takes $\nabla x = u$, with respect to l_1 norm: updates are

$$x_i^+ = x_i - t \cdot \operatorname{sign}\left(\frac{\partial f}{\partial x_i}(x)\right).$$

First order method

Iterative method, updates $x^{(k)}$ in $x^{(0)} + \operatorname{span} \left\{ \nabla f(x^{(0)}), \nabla f(x^{(1)}), \dots, \nabla f(x^{(k-1)}) \right\}$.

Theorem (Nesterov): For any $k \leq (n - 1)/2$ and any starting point $x^{(0)}$, there is a function f in the problem class such that any first-order method satisfies

$$f(x^{(k)}) - f^* \geq \frac{3L\|x^{(0)} - x^*\|_2^2}{32(k + 1)^2}$$

5. Subgradient method

Like gradient descent, but replacing gradients with subgradients. I.e., initialize $x^{(0)}$, repeat

$$x^{(k)} = x^{(k-1)} - t_k \cdot g^{(k-1)}, \quad k = 1, 2, 3, \dots$$

Where $g^{(k-1)} \in \partial f(x^{(k-1)})$, any subgradients of f at $x^{(k-1)}$.

Subgradient method is not necessarily a descent method, so we keep track of best iterate $x_{\text{best}}^{(k)}$ among $x^{(0)}, \dots, x^{(k)}$ so far, i.e.,

$$f(x_{\text{best}}^{(k)}) = \min_{i=0, \dots, k} f(x^{(i)}).$$

Step size choices

- Fixed step size: $t_k = t$ all $k = 1, 2, 3, \dots$
- Diminishing step size: $\sum_{k=1}^{\infty} t_k^2 < \infty, \quad \sum_{k=1}^{\infty} t_k = \infty$

Important difference to gradient descent: all step sizes options are **pre-specified**, not adaptively computed.

Convergence analysis

Assume that f convex, $\text{dom}(f) = \mathbb{R}^n$, and also that f is Lipschitz continuous with constant $G > 0$, i.e.,

$$|f(x) - f(y)| \leq G\|x - y\|_2 \quad \text{for all } x, y.$$

Theorem: For a fixed step size t , subgradient method satisfies

$$\lim_{k \rightarrow \infty} f(x_{\text{best}}^{(k)}) \leq \underline{f^* + G^2 t / 2}$$

Theorem: For diminishing step sizes, subgradient method satisfies

$$\lim_{k \rightarrow \infty} f(x_{\text{best}}^{(k)}) = \underline{f^*}$$

Can prove both results from same basic inequality. Key steps:

- Using the definition of subgradient

$$\|x^{(k)} - x^*\|_2^2 \leq \|x^{(k-1)} - x^*\|_2^2 - 2t_k (f(x^{(k-1)}) - f(x^*)) + t_k^2 \|g^{(k-1)}\|_2^2$$

- Iterating last inequality

$$\|x^{(k)} - x^*\|_2^2 \leq \|x^{(0)} - x^*\|_2^2 - 2 \sum_{i=1}^k t_i (f(x^{(i-1)}) - f(x^*)) + \sum_{i=1}^k t_i^2 \|g^{(i-1)}\|_2^2$$

- Using $\|x^{(k)} - x^*\|_2 \geq 0$ and letting $R = \|x^{(0)} - x^*\|_2$,

$$0 \leq R^2 - 2 \sum_{i=1}^k t_i (f(x^{(i-1)}) - f(x^*)) + G^2 \sum_{i=1}^k t_i^2$$

- Introducing $f(x_{\text{best}}^{(k)}) = \min_{i=0, \dots, k} f(x^{(i)})$, and rearranging,

$$f(x_{\text{best}}^{(k)}) - f(x^*) \leq \frac{R^2 + G^2 \sum_{i=1}^k t_i^2}{2 \sum_{i=1}^k t_i}$$

We call this basic inequality.

After k steps with fixed step size t , basic inequality gives

$$f(x_{\text{best}}^{(k)}) - f^* \leq \frac{R^2}{2kt} + \frac{G^2 t}{2}$$

Proximal gradient descent

1. Decomposable functions

Suppose $f(x) = g(x) + h(x)$, where

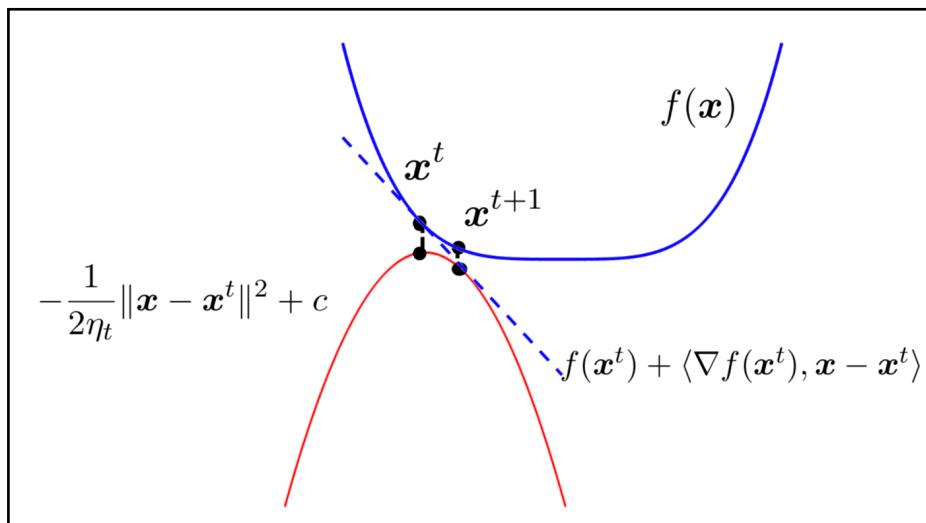
- g is convex, differentiable, $\text{dom } g = \mathbb{R}^n$
- h is convex, not necessarily differentiable

If f were differentiable, minimize quadratic approximation to f around x , replace $\nabla^2 f(x)$ by

$$\frac{1}{t}I,$$

$$x^+ = \underset{z}{\operatorname{argmin}} \underbrace{f(x) + \nabla f(x)^T(z - x) + \frac{1}{2t}\|z - x\|_2^2}_{\tilde{f}(z)},$$

which is gradient descent.



In our case, f isn't differentiable, but g is differentiable. Idea: make quadratic approximation to g , leave h alone.

$$\begin{aligned}
x^+ &= \operatorname{argmin}_z \tilde{g}_t(z) + h(z) \\
&= \operatorname{argmin}_z g(x) + \nabla g(x)^T(z - x) + \frac{1}{2t} \|z - x\|_2^2 + h(z). \\
&= \operatorname{argmin}_z \frac{1}{2t} \|z - (x - t \nabla g(x))\|_2^2 + h(z)
\end{aligned}$$

stay close to gradient update for g , also makes h small.

2. Proximal gradient descent

Define proximal mapping:

$$\operatorname{prox}_t(x) = \operatorname{argmin}_z \frac{1}{2t} \|x - z\|_2^2 + h(z).$$

Proximal gradient descent: choose initialize $x^{(0)}$, repeat

$$x^{(k)} = \operatorname{prox}_{t_k} \left(x^{(k-1)} - t_k \nabla g(x^{(k-1)}) \right), \quad k = 1, 2, 3, \dots$$

To make this update step look familiar, can rewrite it as $x^{(k)} = x^{(k-1)} - t_k \cdot G_{t_k}(x^{(k-1)})$, where G_t is the generalized gradient of f

$$G_t(x) = \frac{x - \operatorname{prox}_t(x - t \nabla g(x))}{t}.$$

With criterion $f(x) = g(x) + h(x)$, we assume

- g is convex, differentiable, $\operatorname{dom}(g) = \mathbb{R}^n$, and ∇g is Lipschitz continuous with constant $L > 0$
- h is convex, $\operatorname{prox}_t(x) = \operatorname{argmin}_z \{\|x - z\|_2^2/(2t) + h(z)\}$ can be evaluated

Theorem: Proximal gradient descent with fixed step size $t \leq 1/L$ satisfies

$$f(x^{(k)}) - f^* \leq \frac{\|x^{(0)} - x^*\|_2^2}{2tk}$$

Backtracking line search

Recall that for unconstrained case, backtracking line search is based on sufficient decrease criterion

$$f\left(\mathbf{x}^t - \eta \nabla f(\mathbf{x}^t)\right) \leq f(\mathbf{x}^t) - \frac{\eta}{2} \|\nabla f(\mathbf{x}^t)\|_2^2,$$

where we set $\alpha = 1/2$. As a result, this is equivalent to updating $\eta_t = 1/L_t$ until

$$\begin{aligned} f\left(\mathbf{x}^t - \eta_t \nabla f(\mathbf{x}^t)\right) &\leq f(\mathbf{x}^t) - \frac{1}{L_t} \langle \nabla f(\mathbf{x}^t), \nabla f(\mathbf{x}^t) \rangle + \frac{1}{2L_t} \|\nabla f(\mathbf{x}^t)\|_2^2 \\ &= f(\mathbf{x}^t) - \langle \nabla f(\mathbf{x}^t), \mathbf{x}^t - \mathbf{x}^{t+1} \rangle + \frac{L_t}{2} \|\mathbf{x}^t - \mathbf{x}^{t+1}\|_2^2. \end{aligned}$$

Similar to gradient descent, but operates on g and not f . We fix a parameter $0 < \beta < 1$. At each iteration, start with $t = 1$, while

$$g(x^+) > g(v) + \nabla g(v)^T(x^+ - v) + \frac{1}{2t} \|x^+ - v\|_2^2,$$

where $x^+ = \text{prox}_t(v - t \nabla g(v))$. Shrink $t = \beta t$, else keep x^+ .

Special cases

- A. $h = 0$, gradient descent
- B. $h = I_C$, projected gradient descent
- C. $g = 0$, proximal minimization algorithm

3. Acceleration

Accelerated proximal gradient method: choose an initial point $x^{(0)} = x^{(-1)} \in \mathbb{R}^n$, repeat for $k = 1, 2, 3, \dots$

$$v = x^{(k-1)} + \frac{k-2}{k+1} (x^{(k-1)} - x^{(k-2)})$$

$$x^{(k)} = \text{prox}_{t_k}(v - t_k \nabla g(v))$$

Theorem: Accelerated proximal gradient method with fixed step size $t \leq 1/L$ satisfies

$$\underline{f(x^{(k)}) - f^* \leq \frac{2\|x^{(0)} - x^*\|_2^2}{t(k+1)^2}}$$

4. LASSO

Given $y \in \mathbb{R}^n, X \in \mathbb{R}^{n \times p}$, lasso problem can be parametrized as:

$$\min_{\beta \in \mathbb{R}^p} \frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1$$

where $\lambda > 0$.

A. Sub-gradient method

update: $\beta^{(k)} = \beta^{(k-1)} + X^T(y - X\beta^{(k-1)}) - \partial \|\beta^{(k-1)}\|_1$

B. Proximal gradient descent (ISTA)

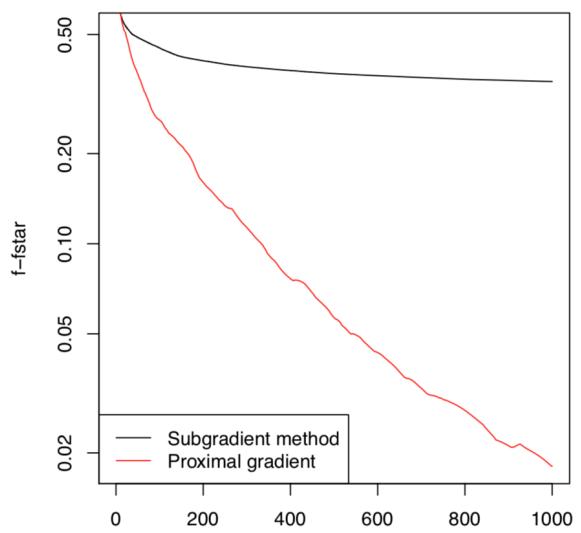
Proximal mapping is now

$$\begin{aligned} \text{prox}_t(\beta) &= \underset{z \in \mathbb{R}^n}{\operatorname{argmin}} \frac{1}{2t} \|\beta - z\|_2^2 + \lambda \|z\|_1, \\ &= S_{\lambda t}(\beta) \end{aligned}$$

Where $S_\lambda(\beta)$ is the soft-thresholding operator,

$$[S_\lambda(\beta)]_i = \begin{cases} \beta_i - \lambda & \text{if } \beta_i > \lambda \\ 0 & \text{if } -\lambda \leq \beta_i \leq \lambda, \quad i = 1, \dots, n \\ \beta_i + \lambda & \text{if } \beta_i < -\lambda \end{cases}$$

Proximal gradient update is: $\beta^+ = S_{\lambda t}(\beta + tX^T(y - X\beta))$, which is often called iterative soft-thresholding algorithm.



Duality in General Programs

1. Lagrangian

Consider general minimization problem

$$\begin{array}{ll} \min & f(x) \\ \text{subject to} & h_i(x) \leq 0, i = 1, 2, \dots, m \\ & \ell_j(x) = 0, j = 1, 2, \dots, m \end{array}$$

We define the **Lagrangian** as

$$L(x, u, v) = f(x) + \sum_{i=1}^m u_i h_i(x) + \sum_{j=1}^r v_j \ell_j(x),$$

new variables $u \in \mathbb{R}^m, v \in \mathbb{R}^r$, with $u \geq 0$.

Important property: for any $u \geq 0$ and v ,

$$f(x) \geq L(x, u, v) \text{ at each feasible } x.$$

2. Lagrange dual problem

Let C denote the feasible set, f^* denote primal optimal value, then

$$f^* \geq \min_{x \in C} L(x, u, v) \geq \min_x L(x, u, v) := g(u, v).$$

We call $g(u, v)$ the **Lagrange dual function**, and it gives a lower bound on f^* for any $u \geq 0$ and v , called feasible dual u, v .

Best lower bound is given by maximizing $g(u, v)$ over all feasible dual u, v , yielding Lagrange dual problem:

$$\begin{array}{ll} \max & g(u, v) \\ \text{subject to} & u \geq 0 \end{array}.$$

Important properties:

- Dual problem is always convex
- The primal optimal value f^* and dual optimal g^* satisfies weak duality: $f^* \geq g^*$

- Slater's condition: for convex primal, if there is an x such that

$$h_1(x) < 0, \dots h_m(x) < 0 \quad \text{and} \quad \ell_1(x) = 0, \dots \ell_r(x) = 0$$

Then strong duality holds: $f^* = g^*$.

3. KKT conditions

For a problem with strong duality, x^* and u^*, v^* are primal and dual solutions $\iff x^*$ and u^*, v^* satisfy the KKT conditions. The Karush-Kuhn-Tucker conditions are:

- Stationarity $0 \in \partial f(x) + \sum_{i=1}^m u_i \partial h_i(x) + \sum_{j=1}^r v_j \partial \ell_j(x)$
- Complementary slackness $u_i \cdot h_i(x) = 0$ for all i
- Primal feasibility $h_i(x) \leq 0, l_j(x) = 0$ for all i, j
- Dual feasibility $u_i \geq 0$ for all i

Theorem: Let f be differentiable and strictly convex, let $X \in \mathbb{R}^{n \times p}$, $\lambda > 0$. Consider

$$\min_{\beta \in \mathbb{R}^p} f(X\beta) + \lambda \|\beta\|_1$$

If the entries of X are drawn from a continuous probability distribution (on \mathbb{R}^{np}), then w.p. 1 there is a unique solution and it has at most $\min\{n, p\}$ nonzero components

4. Conjugates and dual problem

Consider

$$\begin{aligned} & \min_x f(x) + g(x) \\ \iff & \min_{x,z} f(x) + g(z) \text{ subject to } x = z \end{aligned}$$

Lagrange dual function

$$g(u) = \min_x f(x) + g(z) + u^T(z - x) = -f^*(u) - g^*(-u),$$

Hence dual problem is $\max_u -f^*(u) - g^*(-u)$

Example: Lasso dual

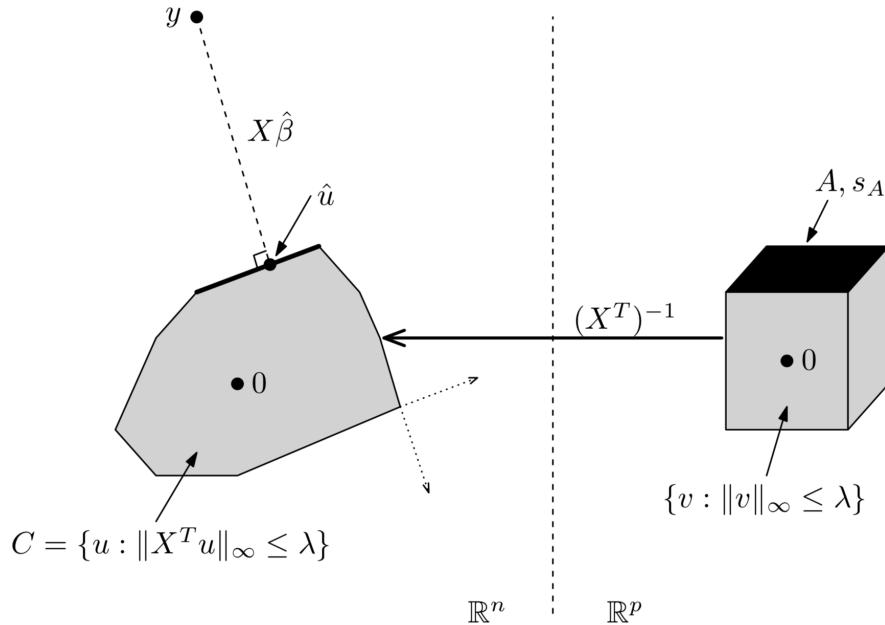
Primal problem is

$$\min_{\beta \in \mathbb{R}^p, z \in \mathbb{R}^n} \frac{1}{2} \|y - z\|_2^2 + \lambda \|\beta\|_1 \text{ subject to } z = X\beta,$$

So dual function is now

$$\begin{aligned} g(u) &= \min_{\beta \in \mathbb{R}^p, z \in \mathbb{R}^n} \frac{1}{2} \|y - z\|_2^2 + \lambda \|\beta\|_1 + u^T(z - X\beta) \\ &= \frac{1}{2} \|y\|_2^2 - \frac{1}{2} \|y - u\|_2^2 - I_{\{v : \|v\|_\infty \leq 1\}}(X^T u / \lambda). \end{aligned}$$

Therefore Lasso dual problem is $\min_{u \in \mathbb{R}^n} \|y - u\|_2^2$ subject to $\|X^T u\|_\infty \leq \lambda$. Further, note that given the dual solution u , any Lasso solution satisfies $X\beta = y - u$. This is from the KKT stationary condition for z (i.e. $z - y + u = 0$).



5. Dual cones and dual problems

Dual cones: for a cone $K \subseteq \mathbb{R}^n$ (recall this means $x \in K, t \geq 0 \implies tx \in K$) and the dual cone is

$$K^* = \{y \in \mathbb{R}^n : y^T x \geq 0 \text{ for all } x \in K\}.$$

Examples

- **Linear subspace:** the dual cone of linear subspace K is K^\perp
- **Norm cone:** the dual cone of the norm cone $K = \{(x, t) \in \mathbb{R}^{n+1} : \|x\| \leq t\}$ is the norm cone its dual norm $K^* = \{(y, s) \in \mathbb{R}^{n+1} : \|y\|_* \leq s\}$

To prove the result we have to show

$$x^T y + ts \geq 0 \text{ whenever } \|x\| \leq t \iff \|y\|_* \leq s.$$

(\Leftarrow) Suppose $\|y\|_* \leq s$ and $\|x\| \leq t$ for some $t > 0$ ($t = 0$ obviously hold). Applying the definition of dual norm and the fact that $\|-x/t\| \leq 1$, we have

$$y^T(-x/t) \leq \|y\|_* \leq s,$$

And therefore $x^T y + ts \geq 0$.

(\Rightarrow) Suppose $\|y\|_* > s$, then there exists x with $\|x\| \leq 1$ and $x^T y \geq s$. Taking $t = 1$, we have

$$-x^T y + s \leq 0.$$

- **Positive semidefinite cone:** the convex cone \mathbb{S}_+^n is self-dual, meaning $(\mathbb{S}_+^n)^* = \mathbb{S}_+^n$.

Next we will show $Y \succeq 0 \iff \text{tr}(YX) \geq 0$ for all $X \succeq 0$.

(\Leftarrow) Suppose $Y \not\succeq 0$, then there exists an $u \in \mathbb{R}^n$ we have $u^T Y u = \text{tr}(uu^T Y) < 0$, and let $X = uu^T$.

(\Rightarrow) The eigenvalue decomposition of X is $X = Q\Lambda Q^T = \sum_{i=1}^n \lambda_i q_i q_i^T$, then we have

$$\text{tr}(YX) = \text{tr}\left(\sum_{i=1}^n \lambda_i q_i q_i^T Y q_i\right) \geq 0.$$

Consider the cone constrained problem $\min_{x \in K} f(x)$. Recall that its dual problem is $\max_{u \in \mathbb{R}^n} -f^*(u) - I_K^*(-u)$, where $I_K^*(-u) = I_{K^*}(u)$, and this is simply $\max_{x \in K^*} -f^*(x)$.

6. SVM

Functional and geometric margins

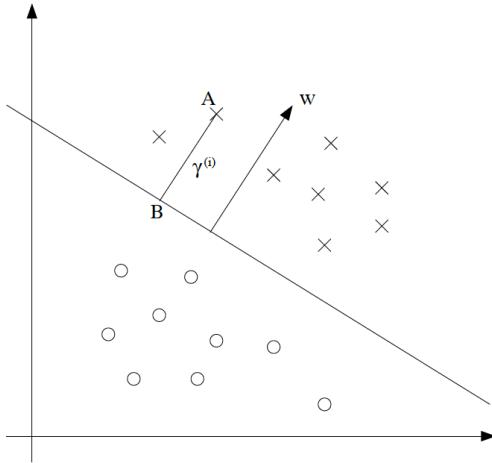
Given a training example $(x^{(i)}, y^{(i)})$, we define the function margin of (ω, b) with respect to the training example

$$\hat{\gamma}^{(i)} = y^{(i)} (\omega^T x + b).$$

Given a training set $S = \{(x^{(i)}, y^{(i)}) : i = 1, 2, \dots, m\}$, we also define the function margin of (ω, b) with respect to the training set

$$\hat{\gamma} = \min_i \hat{\gamma}^{(i)}.$$

Now let's talk about geometric margins, $\gamma^{(i)}$ is length of line segment AB, i.e., the distance from a training sample to the line $\omega^T x + b = 0$.



Since A represents $x^{(i)}$, therefore B is given by $x^{(i)} - \gamma^{(i)}\omega/\|\omega\|$. And point B lies on the decision boundary. Hence,

$$\omega^T \left(x^{(i)} - \gamma^{(i)} \frac{\omega}{\|\omega\|} \right) + b = 0,$$

Solving for $\gamma^{(i)}$ yields

$$\gamma^{(i)} = \frac{\omega^T x^{(i)} + b}{\|\omega\|} = \left(\frac{\omega}{\|\omega\|} \right)^T x^{(i)} + \frac{b}{\|\omega\|}.$$

Note that if $\|\omega\| = 1$, the geometric margin equals to functional marginal.

The optimal margin classifier

We can pose the following optimization problem:

$$\begin{aligned} & \max_{\hat{\gamma}, \omega, b} \hat{\gamma} \\ & \text{s.t. } y^{(i)} (\omega^T x^{(i)} + b) \geq \hat{\gamma} \\ & \quad \|\omega\| = 1 \end{aligned}$$

Let's try transforming the problem into a nicer one, consider

$$\begin{aligned} & \max_{\hat{\gamma}, \omega, b} \frac{\hat{\gamma}}{\|\omega\|} \\ & \text{s.t. } y^{(i)} (\omega^T x^{(i)} + b) \geq \hat{\gamma} \end{aligned}$$

We will introduce the scaling constraint that the functional margin of w, b with respect to the training set must be 1: $\hat{\gamma} = 1$. Then we can transform the optimization problem to

$$\begin{aligned} & \min_{\omega, b} \frac{1}{2} \|\omega\|^2 \\ & \text{s.t. } y^{(i)} (\omega^T x^{(i)} + b) \geq 1 \end{aligned}$$

Lagrange duality

The Lagrange dual function of optimization problem is

$$\mathcal{L}(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i [y^{(i)} (w^T x^{(i)} + b) - 1].$$

Then we obtain the following dual problem

$$\begin{aligned} & \max_{\alpha} W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle \\ & \alpha_i \geq 0, i = 1, 2, \dots, m \\ & \sum_{i=1}^n \alpha_i y^{(i)} = 0 \end{aligned}$$

Newton's method

Given unconstrained, smooth and convex optimization $\min_x f(x)$, where f is convex, twice differentiable, and $\text{dom } f = \mathbb{R}^n$. Newton's method repeats

$$x^{(k)} = x^{(k-1)} - \left(\nabla^2 f(x^{(k-1)}) \right)^{-1} \nabla f(x^{(k-1)}), \quad k = 1, 2, 3, \dots$$

Newton's method interpretation

- Better quadratic approximation

$$f(y) \approx f(x) + \nabla f(x)^T (y - x) + \frac{1}{2} (y - x)^T \nabla^2 f(x) (y - x),$$

And minimize over y to yield $x^+ = x - (\nabla^2 f(x))^{-1} \nabla f(x)$.

- Linearized optimality condition $0 = \nabla f(x + v) \approx \nabla f(x) + \nabla^2 f(x)v$ and solving for v , which again yields $v = -(\nabla^2 f(x))^{-1} \nabla f(x)$.

Affine invariance of Newton's method

Given f , nonsingular $A \in \mathbb{R}^{m \times n}$. Let $x = Ay$, and $g(y) = f(Ay)$. Newton steps on g are

$$x^+ = x - (\nabla^2 f(x))^{-1} f(x).$$

The Newton decrement

At a point x , we define the Newton decrement $\lambda(x) = \left(\nabla f(x)^T (\nabla^2 f(x))^{-1} \nabla f(x) \right)^{1/2}$, which is the difference between $f(x)$ and minimum of its quadratic approximation.

$$f(x) - \min_y f(x) + \nabla f(x)^T (y - x) + \frac{1}{2} (y - x)^T \nabla^2 f(x) (y - x) = \frac{1}{2} \lambda^2(x).$$

Backtracking line search

In practice, we use damped Newton's method, which repeats $x^+ = x - t (\nabla^2 f(x))^{-1} f(x)$. Steps size here typically are chosen by backtracking line search, with parameters $0 < \alpha < 1/2$, $0 < \beta < 1$. At each iteration, we start with $t = 1$ and while

$$f(x + t \nu) > f(x) + \alpha t \nabla f(x)^T \nu,$$

We shrink $t = \beta t$, else we perform Newton update. Note that here $v = -(\nabla^2 f(x))^{-1} \nabla f(x)$, thus $\nabla f(x)^T v = -\lambda^2(x)$ is a constant.

Convergence analysis

Assume that f convex, twice differentiable, having $\text{dom } f = \mathbb{R}^n$ and additionally,

- $\nabla f(x)$ is Lipschitz with parameter L
- f is strongly convex with parameter m
- $\nabla^2 f(x)$ is Lipschitz with parameter M

Theorem: Newton's method with backtracking line search satisfies the following two-stage convergence bounds

$$f(x^{(k)}) - f^* \leq \begin{cases} (f(x^{(0)}) - f^*) - \gamma k & \text{if } k \leq k_0 \\ \frac{2m^3}{M^2} \left(\frac{1}{2}\right)^{2^{k-k_0+1}} & \text{if } k > k_0 \end{cases}$$

Here $\gamma = \alpha\beta^2\eta^2m/L^2$, $\eta = \min\{1, 3(1-2\alpha)\}m^2/M$, and k_0 is the number of steps until $\|\nabla f(x^{(k_0+1)})\|_2 \leq \eta$

Convergence follows two stages:

- Damped phase: $\|\nabla f(x^{(k)})\|_2 \geq \eta$, and $f(x^{(k)}) - f(x^{(k+1)}) \geq \gamma$.
- Pure phase: $\|\nabla f(x^{(k)})\|_2 < \eta$, backtracking selects $t = 1$, and

$$\frac{M}{2m^2} \left\| \nabla f(x^{(k+1)}) \right\|_2 \leq \left(\frac{M}{2m^2} \left\| \nabla f(x^{(k)}) \right\|_2 \right)^2.$$

To reach $f(x^{(k)}) - f^* \leq \epsilon$, we need at most $\frac{f(x^{(0)}) - f^*}{\gamma} + \log \log (\epsilon_0/\epsilon)$ iterations, where $\epsilon_0 = 2m^3/M^2$.

Strong advantages

- Convergence of Newton's method is rapid in general, and quadratic near x^* .
- Newton's method is affine invariant.

- Newton's method scales well with problem size. Its performance on problem \mathbb{R}^{10000} is similar to which on problem \mathbb{R}^{10} .
- The good performance of Newton's method is not depend on the choice of algorithm parameters.

The main disadvantage of Newton's method is cost of **forming and storing the Hessian**, and the cost of **computing the Newton step**, which requires solving a set of linear equations.

Self-Concordance

Two major shortcoming of classical convergence analysis of Newton's method:

- The unknown parameters L, m, M in practice
- Dependent on coordinate system used

A scale-free analysis is possible for **self-concordant functions**: on \mathbb{R} , a convex function f is called self-concordant if

$$|f'''(x)| \leq 2f''(x)^{3/2}.$$

And on \mathbb{R}^n is called self-concordant if its projection on every line segment is so.

Theorem (Nesterov and Nemirovskii): Newton's method with backtracking line search requires at most

$$C(\alpha, \beta)(f(x^{(0)}) - f^*) + \log \log(1/\epsilon)$$

iterations to reach $f(x^{(k)}) - f^* \leq \epsilon$, where $C(\alpha, \beta)$ is a constant that only depends on α, β

Comparison to first-order methods

- Memory: each iteration of Newton's method requires $O(n^2)$ storage ($n \times n$ Hessian); Each gradient iteration requires $O(n)$ storage (n-dimensional gradient).
- Computation: each Newton iteration requires $O(n^3)$ flops (solving a dense $n \times n$ linear system); Each gradient iteration requires $O(n)$ flops (scaling/adding n-dimensional gradient)

- Backtracking: both use $O(n)$ flops per inner backtracking step.

Equality-constrained Newton's method

Consider now a problem with equality constraints, as in

$$\min_x f(x) \text{ subject to } Ax = b.$$

Several options:

- Eliminating equality constraints: write $x = Fy + b$, where F spans null space of A , and $Ax_0 = b$. Solve in terms of y .
- Deriving the dual: can check the Lagrange dual function $-f^*(-A^T v) - b^T v$, and strong duality holds. With luck, we can express x^* in terms of v^* .
- Equality-constrained Newton:

In equality-constrained Newton's method we start with $x^{(0)}$ such that $Ax^{(0)} = b$. Then we repeat the updates $x^+ = x + t v$, where

$$v = \underset{Az=0}{\operatorname{argmin}} \nabla f(x)^T(z - x) + \frac{1}{2}(z - x)^T \nabla^2 f(x)(z - x).$$

Furthermore, v is the solution to minimizing a quadratic subject to equality constraints. From the KKT conditions that v satisfies

$$\begin{bmatrix} \nabla^2 f(x) & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} v \\ w \end{bmatrix} = \begin{bmatrix} -\nabla f(x) \\ 0 \end{bmatrix},$$

For some w . Hence Newton direction is given by solving linear system in Hessian.

Barrier Method

1. Hierarchy of second-order method

Assuming all problems are convex

- **Quadratic problems** are the easiest: closed form solution
- **Equality-constrained quadratic problems** are still easy: we use KKT conditions to derive closed form solution
- **Equality-constrained smooth problems** are next: use Newton's method to reduce this to a sequence of equality constrained quadratic problems
- **Inequality- and equality-constrained smooth problems:** use interior methods to reduce this a sequence of equality constrained quadratic problems

2. Log barrier function

Consider the convex optimization problem

$$\begin{aligned} \min_x \quad & f(x) \\ \text{subject to} \quad & h_i(x) \leq 0, i = 1, 2, \dots, m \\ & Ax = b \end{aligned}$$

We will assume that $f, h_1(x), \dots, h_m(x)$ are convex and twice differentiable, each with the domain \mathbb{R}^n . The function

$$\phi(x) = \sum_{i=1}^m -\log(-h_i(x))$$

Is called the log barrier for the above problem. Its domain is the strictly feasible points, $\{x : h_i(x) < 0, i = 1, 2, \dots, m\}$.

Ignoring the equality constraints for now, the problem can be written as

$$\min_x f(x) + \sum_{i=1}^m I(h_i(x) \leq 0).$$

We approximate this representation by adding log barrier function

$$\min_x f(x) - (1/t) \sum_{i=1}^m \log(-h_i(x)),$$

Where $t > 0$ is a large number.

3. Central path

Consider minimize our problem

$$\begin{aligned} \min_x & \quad tf(x) + \phi(x) \\ \text{subject to} & \quad Ax = b \end{aligned}$$

The **central path** is defined as the solution $x^*(t)$ as a function of $t > 0$. These solutions are characterized by KKT conditions:

$$Ax^*(t) = b, h_i(x^*(t)) < 0, \quad i = 1, \dots, m,$$

$$t \nabla f(x^*(t)) - \sum_{i=1}^m \frac{1}{h_i(x^*(t))} \nabla h_i(x^*(t)) + A^T w = 0$$

For some $w \in \mathbb{R}^m$. As $t \rightarrow \infty$, we hope that $x^*(t) \rightarrow x^*$.

4. Dual points from central path

Given $x^*(t)$ and corresponding w , we fine

$$u_i^*(t) = -\frac{1}{th_i(x^*)}, \quad i = 1, 2, \dots, m, \quad v^*(t) = w/t.$$

Note that $u_i^*(t) > 0$ for all i , furthermore, the point $(u^*(t), v^*(t))$ lies in domain of Lagrange dual function $g(u, v)$, since by definition

$$\nabla f(x^*(t)) + \sum_{i=1}^m u_i(x^*(t)) \nabla h_i(x^*(t)) + A^T v^*(t) = 0.$$

I.e., $x^*(t)$ minimize the Lagrangian $L(x, u^*(t), v^*(t))$ over x , so $g(u^*(t), v^*(t)) > -\infty$.

We compute

$$\begin{aligned} g(u^*(t), v^*(t)) &= f(x^*(t)) + \sum_{i=1}^m u_i^*(t) h_i(x^*(t)) + v^*(t)^T (Ax^*(t) - b) \\ &= f(x^*(t)) - m/t \end{aligned}$$

That is $f(x^*(t)) - f^* \leq m/t$, this will be very useful as a **stopping criterion**; it also confirms the fact that as $t \rightarrow \infty, x^*(t) \rightarrow x^*$.

5. Barrier method

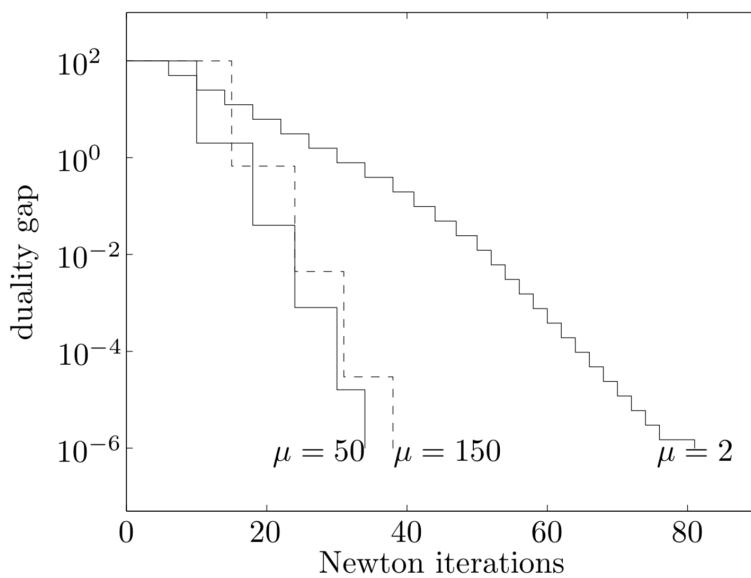
The Barrier method solves a sequence of problems

$$\begin{aligned} & \min_x \quad tf(x) + \phi(x) \\ & \text{subject to} \quad Ax = b \end{aligned}$$

For increasing $t > 0$, until $m/t \leq \epsilon$. We start with $t = t^{(0)} > 0$, and solve the above problem using Newton's method to produce $x^{(0)} = x^*(t)$. Then for barrier parameter $\mu > 1$, we repeat, for $k = 1, 2, \dots$

- Solve the barrier problem at $t = t^{(k)}$, using Newton's method initialize at $x^{(k-1)}$, to produce $x^{(k)} = x^*(t)$
- Stop if $m/t \leq \epsilon$
- Else update $t^{(k+1)} = \mu t$

The first step is called centering step since it brings $x^{(k)}$ onto the central path.



6. Convergence analysis

Theorem. The barrier method after k centering steps satisfies $f(x^{(k)}) - f^* \leq \frac{m}{\mu^k t^{(0)}}$.

In other words, to reach a desired accuracy level of ϵ , we require $\frac{\log(m/(t^{(0)}\epsilon))}{\log \mu} + 1$

centering steps with the barrier method (plus the initial centering step).

7. Feasibility method

How to find such a feasible x ? By solving

$$\begin{aligned} \min_{x,s} \quad & s \\ \text{subject to} \quad & h_i(x) \leq s, \quad i = 1, \dots, m \\ & Ax = b \end{aligned}$$

The goal is for s to be negative at the solution. Once we find a feasible (x, s) with $s < 0$, we can terminate early.

An alternative is to solve the problem

$$\begin{aligned} \min_{x,s} \quad & 1^T s \\ \text{subject to} \quad & h_i(x) \leq s_i, \quad i = 1, \dots, m \\ & Ax = b, \quad s \geq 0 \end{aligned}$$

The **nonzero entries** of s will tell us which of the constraints cannot be satisfied.

Gradient descent for constraint problem

minimize _{x} $f(x)$
subject to $x \in \mathcal{C}$

Where $f(x)$ is a convex function and \mathcal{C} is a convex set.

1. Frank-Wolfe Algorithm

Algorithm 3.1 Frank-wolfe (a.k.a. conditional gradient) algorithm

```

1: for  $t = 0, 1, \dots$  do
2:    $\mathbf{y}^t := \arg \min_{\mathbf{x} \in \mathcal{C}} \langle \nabla f(\mathbf{x}^t), \mathbf{x} \rangle$            (direction finding)
3:    $\mathbf{x}^{t+1} = (1 - \eta_t) \mathbf{x}^t + \eta_t \mathbf{y}^t$            (line search and update)

```

Direction finding is equivalent to

$$f(\mathbf{x}^t) + \left\langle \nabla f(\mathbf{x}^t), \mathbf{x} - \mathbf{x}^t \right\rangle,$$

which is a linear optimization over convex set. And step size is determined by line search

$$\eta^t = \arg \min_{\eta} f((1-\eta)x^t + \eta y^t).$$

Convergence analysis

Theorem 3.1 (Frank-Wolfe for convex and smooth problems, Jaggi '13)

Let f be convex and L -smooth. With $\eta_t = \frac{2}{t+2}$, one has

$$f(\mathbf{x}^t) - f(\mathbf{x}^*) \leq \frac{2Ld_{\mathcal{C}}^2}{t+2}$$

where $d_{\mathcal{C}} = \sup_{x,y \in \mathcal{C}} \|x - y\|_2$

2. Projected gradient methods

Projected Theorem. Let \mathcal{C} be convex set. Then $x_{\mathcal{C}}$ is projection of x onto \mathcal{C} iff

$$(\mathbf{x} - \mathbf{x}_{\mathcal{C}})^\top (\mathbf{z} - \mathbf{x}_{\mathcal{C}}) \leq 0, \quad \forall \mathbf{z} \in \mathcal{C}.$$

Nonexpansivness of projection. For any \mathbf{x} and \mathbf{z} , one has

$$\| \mathcal{P}_{\mathcal{C}}(\mathbf{x}) - \mathcal{P}_{\mathcal{C}}(\mathbf{z}) \|_2 \leq \| \mathbf{x} - \mathbf{z} \|_2$$

3. Strongly convex and smooth problem

- x^* lies in interior of \mathcal{C}

Then x^* is a loco minimizer because $x^* \in \mathcal{C}$, which implies that $\nabla f(x^*) = 0$.

Theorem 3.5

Suppose $\mathbf{x}^* \in \text{int}(\mathcal{C})$, and let f be μ -strongly convex and L -smooth.
If $\eta_t = \frac{2}{\mu+L}$, then

$$\| \mathbf{x}^t - \mathbf{x}^* \|_2 \leq \left(\frac{\kappa - 1}{\kappa + 1} \right)^t \| \mathbf{x}^0 - \mathbf{x}^* \|_2$$

where $\kappa = L/\mu$ is condition number

- Don't know whether $\mathbf{x} \in \text{int}(\mathcal{C})$

Main issue: $\nabla f(x^*)$ may not equal to 0 (x^* may not be loco minimizer).

Let $\mathbf{x}^+ := \mathcal{P}_{\mathcal{C}} \left(\mathbf{x} - \frac{1}{L} \nabla f(\mathbf{x}) \right)$ and $\mathbf{g}_{\mathcal{C}}(\mathbf{x}) := L(\mathbf{x} - \mathbf{x}^+)$ and $\mathbf{g}_{\mathcal{C}}(\mathbf{x})$ generalizes $\nabla f(\mathbf{x})$ and

obeys $\mathbf{g}_{\mathcal{C}}(\mathbf{x}^*) = 0$.

When $x^* \in \mathcal{C}$, obviously $\mathbf{g}_{\mathcal{C}}(\mathbf{x}^*) = 0$. And when $x^* \in \partial \mathcal{C}$, if $\mathbf{g}_{\mathcal{C}}(\mathbf{x}^*) \neq 0$ then we can find x^{*+} such that $f(x^{*+}) \leq f(x^*)$, which is contradictory.

Theorem 3.7 (projected GD for strongly convex and smooth problems)

Let f be μ -strongly convex and L -smooth. If $\eta_t \equiv \eta = \frac{1}{L}$, then

$$\|\mathbf{x}^t - \mathbf{x}^*\|_2^2 \leq \left(1 - \frac{\mu}{L}\right)^t \|\mathbf{x}^0 - \mathbf{x}^*\|_2^2$$

- Slightly weaker convergence rate compared with Theorem 3.5.

4. Convex and smooth problem

Theorem 3.8 (projected GD for convex and smooth problems)

Let f be convex and L -smooth. If $\eta_t \equiv \eta = \frac{1}{L}$, then

$$f(\mathbf{x}^t) - f(\mathbf{x}^*) \leq \frac{3L\|\mathbf{x}^0 - \mathbf{x}^*\|_2^2 + f(\mathbf{x}^0) - f(\mathbf{x}^*)}{t+1}$$

Proof.

Step 1: show cost improvement

$$f(\mathbf{x}^{t+1}) \leq f(\mathbf{x}^t) - \frac{1}{2L} \left\| \mathbf{g}_{\mathcal{C}}(\mathbf{x}^t) \right\|_2^2.$$

Step 2: connect $\left\| \mathbf{g}_{\mathcal{C}}(\mathbf{x}^t) \right\|_2$ with $f(\mathbf{x}_t)$

$$\left\| \mathbf{g}_{\mathcal{C}}(\mathbf{x}^t) \right\|_2 \geq \frac{f(\mathbf{x}^{t+1}) - f(\mathbf{x}^*)}{\|\mathbf{x}^t - \mathbf{x}^*\|_2} \geq \frac{f(\mathbf{x}^{t+1}) - f(\mathbf{x}^*)}{\|\mathbf{x}^0 - \mathbf{x}^*\|_2}.$$

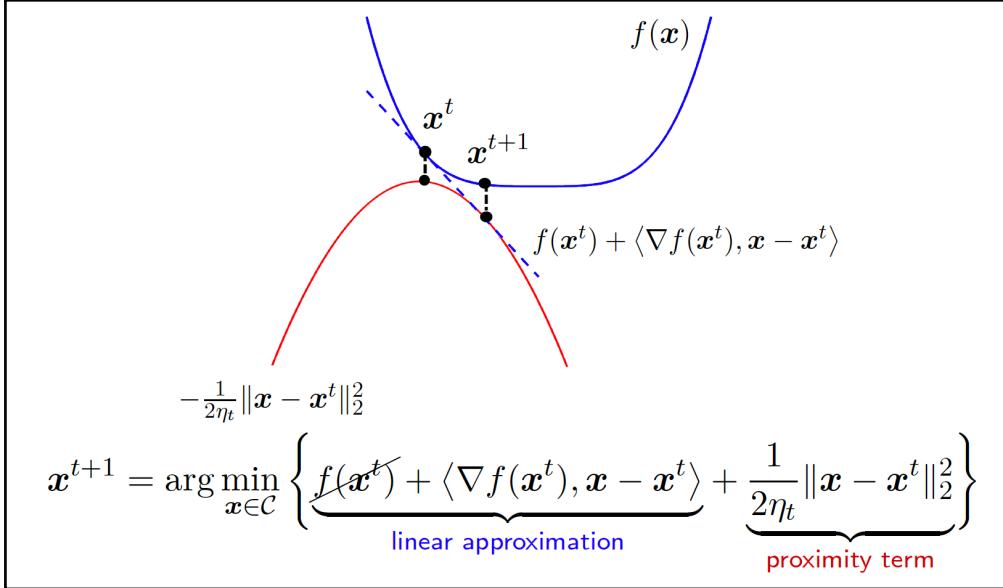
Step 3: Let $\Delta_t := f(\mathbf{x}^t) - f(\mathbf{x}^*)$ to get

$$\Delta_{t+1} - \Delta_t \leq -\frac{\Delta_{t+1}^2}{2L \|\mathbf{x}^0 - \mathbf{x}^*\|_2^2}.$$

And complete proof by induction.

Mirror Descent

A proximal viewpoint of projected GD



Quadratic proximal term is used to monitor discrepancy between $f(x)$ and first order approximation. (homogeneous penalty: squared Euclidean penalty)

Issues: local geometry might sometimes be highly **inhomogeneous**, or even non-Euclidean

Mirror descent

Mirror descent: adjust gradient updates to fit problem geometry

— Nemirovski & Yudin, '1983

Replace quadratic proximity $\|x - x^t\|_2^2$ with distance-like metric D_φ

$$x^{t+1} = \arg \min_{x \in \mathcal{C}} \left\{ f(x^t) + \langle \nabla f(x^t), x - x^t \rangle + \frac{1}{\eta_t} \underbrace{D_\varphi(x, x^t)}_{\text{Bregman divergence}} \right\},$$

where $D_\varphi(x, z) := \varphi(x) - \varphi(z) - \langle \nabla \varphi(z), x - z \rangle$ for convex and differentiable $\varphi(x)$.

Principles in choosing Bregman divergence:

- Fits local curvature of $f(x)$
- Fits geometry of constraint set \mathcal{C}
- Make sure Bregman projection is inexpensive

Bregman Divergence

Let $\varphi : \mathcal{C} \mapsto \mathbb{R}$ be strictly convex and differentiable on \mathcal{C} , then

$$D_\varphi(\mathbf{x}, \mathbf{z}) := \varphi(\mathbf{x}) - \varphi(\mathbf{z}) - \langle \nabla \varphi(\mathbf{z}), \mathbf{x} - \mathbf{z} \rangle.$$

- Squared Mahalanobis distance

Let $D_\varphi(\mathbf{x}, \mathbf{z}) = \frac{1}{2}(\mathbf{x} - \mathbf{z})^\top \mathbf{Q}(\mathbf{x} - \mathbf{z})$ for $\mathbf{Q} > 0$, which is generated by $\varphi(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{Q}\mathbf{x}$.

- KL divergence

Let $D_\varphi(\mathbf{x}, \mathbf{z}) = \text{KL}(\mathbf{x} \parallel \mathbf{z}) := \sum_i x_i \log \frac{x_i}{z_i}$, which is generated by $\varphi(\mathbf{x}) = \sum_i x_i \log x_i$ (negative entropy) if $\mathcal{C} = \Delta := \left\{ \mathbf{x} \in \mathbb{R}_+^n \mid \sum_i x_i = 1 \right\}$ is a probability simplex.

Bregman Projection

Given a point \mathbf{x} , define

$$\mathcal{P}_{\mathcal{C}, \varphi}(\mathbf{x}) := \arg \min_{\mathbf{z} \in \mathcal{C}} D_\varphi(\mathbf{z}, \mathbf{x}).$$

As Bergman projection of \mathbf{x} onto \mathcal{C} .

Generalized Pythagorean Theorem. If $\mathbf{x}_{\mathcal{C}, \varphi} = \mathcal{P}_{\mathcal{C}, \varphi}(\mathbf{x})$, then

$$D_\varphi(\mathbf{z}, \mathbf{x}) \geq D_\varphi\left(\mathbf{z}, \mathbf{x}_{\mathcal{C}, \varphi}\right) + D_\varphi\left(\mathbf{x}_{\mathcal{C}, \varphi}, \mathbf{x}\right) \quad \forall \mathbf{z} \in \mathcal{C}.$$

An alternative form of MD

Using Bergman divergence, one can also describe MD as

$$\nabla \varphi(\mathbf{y}^{t+1}) = \nabla \varphi(\mathbf{x}^t) - \eta_t \mathbf{g}^t \quad \text{with } \mathbf{g}^t \in \partial f(\mathbf{x}^t),$$

$$\mathbf{x}^{t+1} \in \mathcal{P}_{\mathcal{C}, \varphi}(\mathbf{y}^{t+1}) = \arg \min_{z \in \mathcal{C}} D_\varphi(z, \mathbf{y}^{t+1}).$$

$$\begin{aligned} x_{t+1} &= \operatorname{argmin}_{x \in \mathcal{X} \cap \mathcal{D}} D_\Phi(x, y_{t+1}) \\ &= \operatorname{argmin}_{x \in \mathcal{X} \cap \mathcal{D}} \Phi(x) - \nabla \Phi(y_{t+1})^\top x \\ &= \operatorname{argmin}_{x \in \mathcal{X} \cap \mathcal{D}} \Phi(x) - \left(\nabla \Phi(x_t) - \eta g_t \right)^\top x \\ &= \operatorname{argmin}_{x \in \mathcal{X} \cap \mathcal{D}} \eta g_t^\top x + D_\Phi(x, x_t) \end{aligned}$$

Convergence analysis

Theorem 5.3

Suppose f is convex and Lipschitz continuous (i.e. $\|g^t\|_* \leq L_f$) on \mathcal{C} , and suppose φ is ρ -strongly convex w.r.t. $\|\cdot\|$. Then

$$f^{\text{best}, t} - f^{\text{opt}} \leq \frac{\sup_{\mathbf{x} \in \mathcal{C}} D_\varphi(\mathbf{x}, \mathbf{x}^0) + \frac{L_f^2}{2\rho} \sum_{k=0}^t \eta_k^2}{\sum_{k=0}^t \eta_k}$$