

First-order Methods

Unconstrained optimization

- Problem
$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$$
- First-order optimality condition
$$\nabla f(\mathbf{x}^*) = 0.$$

Gradient descent method

- Iteration:

$$\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} - \gamma_i \nabla f(\mathbf{x}^{(i)})$$

for some step size γ_i . This is a special case of the Network-type algorithms with $\mathbf{H}_i = \mathbf{I}$ and $\Delta \mathbf{x} = \nabla f(\mathbf{x}^{(i)})$.

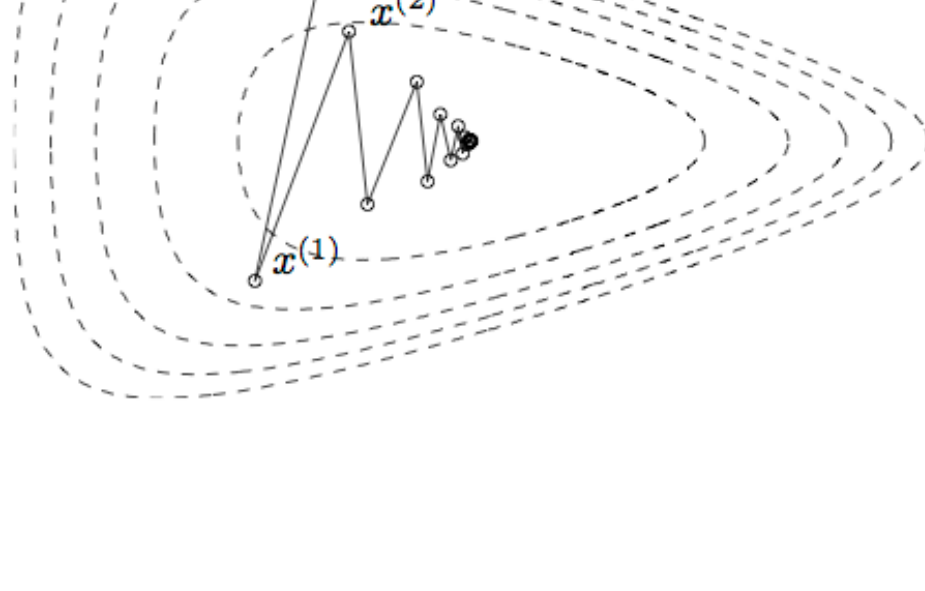
- Idea: Iterative linear approximation (first-order Taylor series expansion).

$$f(\mathbf{x}) \approx f(\mathbf{x}^{(i)}) + \nabla f(\mathbf{x}^{(i)})^T (\mathbf{x} - \mathbf{x}^{(i)})$$

and then minimize the linear approximation within a compact set: let $\Delta \mathbf{x} = \mathbf{x}^{(i+1)} - \mathbf{x}^{(i)}$ to choose. Solve
$$\min_{\|\Delta \mathbf{x}\|_2 \leq 1} \nabla f(\mathbf{x}^{(i)})^T \Delta \mathbf{x}$$

to obtain $\Delta \mathbf{x} = -\nabla f(\mathbf{x}^{(i)}) / \|\nabla f(\mathbf{x}^{(i)})\|_2 \propto -\nabla f(\mathbf{x}^{(i)})$.

- Step sizes are chosen so that the descent property is maintained (e.g., line search).
- Step sizes must be chosen at any time, since unlike Newton's method, first-order approximation is always unbounded below.
- Pros
 - Each iteration is inexpensive.
 - No need to derive, compute, store and invert Hessians; attractive in large scale problems.
- Cons
 - Slow convergence (zigzagging).



- Do not work for non-smooth problems.

Convergence

- In general, the best we can obtain from gradient descent is linear convergence (cf. quadratic convergence of Newton).
- Example (Boyd & Vandenberghe Section 9.3.2)

$$f(\mathbf{x}) = \frac{1}{2} (x_1^2 + cx_2^2), \quad c > 1$$

The optimal value is 0.

It can be shown that if we start from $\mathbf{x}^{(0)} = (c, 1)$, then

$$f(\mathbf{x}^{(i)}) = \left(\frac{c-1}{c+1} \right)^i f(\mathbf{x}^{(0)})$$

and

$$\|\mathbf{x}^{(i)} - \mathbf{x}^*\|_2 = \left(\frac{c-1}{c+1} \right)^i \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2$$

- More generally, if
 - f is convex and differentiable over \mathbb{R}^d ;
 - f has L -Lipschitz gradients, i.e., $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq L \|\mathbf{x} - \mathbf{y}\|_2$;
 - $p^* = \inf_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) > -\infty$ and attained at \mathbf{x}^* ,then
$$f(\mathbf{x}^{(i)}) - p^* \leq \frac{1}{2\gamma_i} \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2 = O(1/i),$$
for constant step size $\gamma_i = \gamma \in (0, 1/L]$, a sublinear convergence. Similar upper bound with line search.
- If we further assume that f is α -strongly convex, i.e., $f(\mathbf{x}) - \frac{\alpha}{2} \|\mathbf{x}\|_2^2$ is convex for some $\alpha > 0$, then
$$f(\mathbf{x}^{(i)}) - p^* \leq \frac{L}{2} \left(1 - \gamma \frac{2\alpha L}{\alpha + L} \right)^i \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2$$

and

$$\|\mathbf{x}^{(i)} - \mathbf{x}^*\|_2^2 \leq \left(1 - \gamma \frac{2\alpha L}{\alpha + L} \right)^i \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2$$

for constant step size $\gamma \in (0, 2/(\alpha + L)]$.

Examples

- Least squares: $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 = \frac{1}{2} \mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x} - \mathbf{x}^T \mathbf{A}^T \mathbf{y} + \frac{1}{2} \mathbf{y}^T \mathbf{y}$
- Gradient $\nabla f(\mathbf{x}) = \mathbf{A}^T \mathbf{A} \mathbf{x} - \mathbf{A}^T \mathbf{y}$ is also $\sigma_{\min}(\mathbf{A})$ -Lipschitz. Hence $\gamma \in (0, 1/\sigma_{\min}(\mathbf{A})^2]$ guarantees descent property and convergence.
- If \mathbf{A} is full column rank, then f is also $\sigma_{\min}(\mathbf{A})^2$ -strongly convex and the convergence is linear. Otherwise it is sublinear.
- Logistic regression: $f(\beta) = -\sum_{i=1}^n \log \left(1 + e^{-\beta^T \mathbf{x}_i} \right)$
- Gradient $\nabla f(\beta) = -\sum_{i=1}^n \mathbf{x}_i e^{-\beta^T \mathbf{x}_i} / (1 + e^{-\beta^T \mathbf{x}_i})^2$ is L -Lipschitz.
- Even if \mathbf{X} is full column rank, f may not be strongly convex.
- Adding a ridge penalty $\frac{\lambda}{2} \|\mathbf{x}\|_2^2$ makes the problems always strongly convex.

Accelerated gradient descent

First-order methods

- Iterative algorithm that generates sequence $\{\mathbf{x}^{(i)}\}$ such that
$$\mathbf{x}^{(i)} \in \mathbf{x}^{(0)} + \text{span}\{\nabla f(\mathbf{x}^{(0)}), \dots, \nabla f(\mathbf{x}^{(i-1)})\}.$$
- Examples:
 - Gradient descent: $\mathbf{x}^{(i)} = \mathbf{x}^{(i-1)} - \gamma_i \nabla f(\mathbf{x}^{(i)})$.
 - Momentum method: $\mathbf{y}^{(i)} = \mathbf{x}^{(i-1)} - \gamma_i \nabla f(\mathbf{x}^{(i)})$, $\mathbf{x}^{(i)} = \mathbf{y}^{(i)} + \alpha_i (\mathbf{y}^{(i)} - \mathbf{y}^{(i-1)})$.
- In general, a first-order method generates

$$\mathbf{x}^{(i)} = \mathbf{x}^{(0)} - \sum_{j=1}^{i-1} \gamma_j \nabla f(\mathbf{x}^{(j)}) =: M_i(f, \mathbf{x}^{(0)}).$$

- Collection of all first-order method can be considered as the set of all M_N . Let's denote it by \mathcal{M}_N .
- Then we want to find a method $M_N \in \mathcal{M}_N$ that minimizes

$$\sup_{f \in \mathcal{F}_L} f(M_N(f, \mathbf{x}^{(0)})) - p^*.$$

for a class of functions \mathcal{F}_L .

- Typically we choose the class of functions as \mathcal{F}_L , the set of functions satisfying the 3 assumptions used for convergence analysis of GD:

- f is convex and differentiable over \mathbb{R}^d ;
- f has L -Lipschitz gradients, i.e., $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq L \|\mathbf{x} - \mathbf{y}\|_2$;
- $p^* = \inf_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) > -\infty$ and attained at \mathbf{x}^* .

- We have seen that for gradient descent with a step size $1/L$,

$$\sup_{f \in \mathcal{F}_L} f(M_N(f, \mathbf{x}^{(0)})) - p^* \leq \frac{L}{2N} \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2 = O(1/N).$$

- Nesterov (1983) showed that for $d \geq 2N + 1$ and every $\mathbf{x}^{(0)}$, there exists $f \in \mathcal{F}_L$ such that for **any** $M_N \in \mathcal{M}_N$,

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

suggesting that the minimax optimal rate of first order methods is at most $O(1/N^2)$.

- This also suggests that the $O(1/N)$ rate of GD can be improved.

- Nesterov's (1983) accelerated gradient method achieves the optimal rate:

$$\mathbf{y}^{(i+1)} = \mathbf{x}^{(i)} - \frac{1}{i} \nabla f(\mathbf{x}^{(i)})$$

$$\gamma_{i+1} = \frac{1}{2} (1 + \sqrt{1 + 4i^2})$$

$$\mathbf{x}^{(i+1)} = \mathbf{y}^{(i+1)} + \frac{\gamma_i - 1}{\gamma_{i+1}} (\mathbf{y}^{(i+1)} - \mathbf{y}^{(i)})$$

with initialization $\mathbf{y}^{(0)} = \mathbf{x}^{(0)}$ and $\gamma_0 = 1$.

- Theorem (Nesterov, 1983): Nesterov's accelerated gradient method algorithm satisfies

$$f(\mathbf{y}^{(i)}) - p^* \leq \frac{2L}{(i+1)^2} \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2.$$

for $f \in \mathcal{F}_L$.

- Kim & Fessler's (2016) algorithm improved this rate by a factor of two:

$$f(\mathbf{x}^{(N)}) - p^* \leq \frac{L}{(N+1)^2} R^2$$

for $f \in \mathcal{F}_L$ and $\|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2 \leq R$.

Steepest descent method

- Recall the idea of GD: iterative linear approximation (first-order Taylor series expansion).

$$f(\mathbf{x}) \approx f(\mathbf{x}^{(i)}) + \nabla f(\mathbf{x}^{(i)})^T (\mathbf{x} - \mathbf{x}^{(i)})$$

and then minimize the linear approximation within a compact set:

$$\min_{\|\Delta \mathbf{x}\|_2 \leq 1} \nabla f(\mathbf{x}^{(i)})^T \Delta \mathbf{x}.$$

- The compact set need not be limited to the ℓ_2 norm ball in order to obtain a descent direction. For any unit norm ball,

$$\min_{\|\Delta \mathbf{x}\|_2 \leq 1} \nabla f(\mathbf{x}^{(i)})^T \Delta \mathbf{x} = -\|\nabla f(\mathbf{x}^{(i)})\|_2.$$

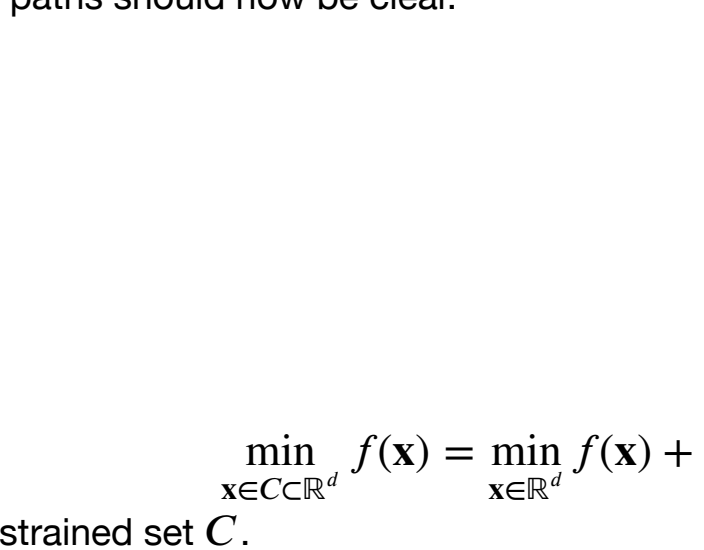
- In particular, if ℓ_1 norm ball is used, then
$$\min_{\|\Delta \mathbf{x}\|_1 \leq 1} \nabla f(\mathbf{x}^{(i)})^T \Delta \mathbf{x} = -\|\nabla f(\mathbf{x}^{(i)})\|_\infty = -\max_{j=1, \dots, d} |\nabla f(\mathbf{x}^{(i)})_j|$$

and the descent direction is given by the convex hull of the elementary unit vectors corresponding to the coordinates of the maximum absolute derivatives. In other words,

$$\Delta \mathbf{x} = \sum_{i,j} \eta_{ij} \text{sign}(\nabla f(\mathbf{x}^{(i)})_j) \mathbf{e}_i, \quad a_i \geq 0, \quad \sum_i a_i = 1, \quad J = \{j : |\nabla f(\mathbf{x}^{(i)})_j| = \|\nabla f(\mathbf{x}^{(i)})\|_\infty\}.$$

- Consider a linear regression setting. Assume each coordinate is standardized. If $a_i = 1/|J|$ for $i \in J$, then the ℓ_1 -steepest descent update with a line search strategy of finding the maximal step size with which J does not vary is the least angle regression (LARS). If only one coordinate is chosen and a very small step size ϵ is used, then this update is ϵ -forward stage-wise regression.

LARS (Least Angle Regression Shrinkage)



Why these algorithms generate sparse solution paths should now be clear.

Proximal gradient methods

Constrained optimization

- Problem
$$\min_{\mathbf{x} \in C \cap \mathbb{R}^d} f(\mathbf{x}) = \min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + \iota_C(\mathbf{x})$$
where $\iota_C(\mathbf{x})$ is the indicator function of the constrained set C .
- Can be viewed as an unconstrained but nonsmooth problem.
- Recall that GD does not work for non-smooth problems. Need a new first-order method.
- Consider a generalization
$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + g(\mathbf{x})$$

where f is smooth but g is nonsmooth.

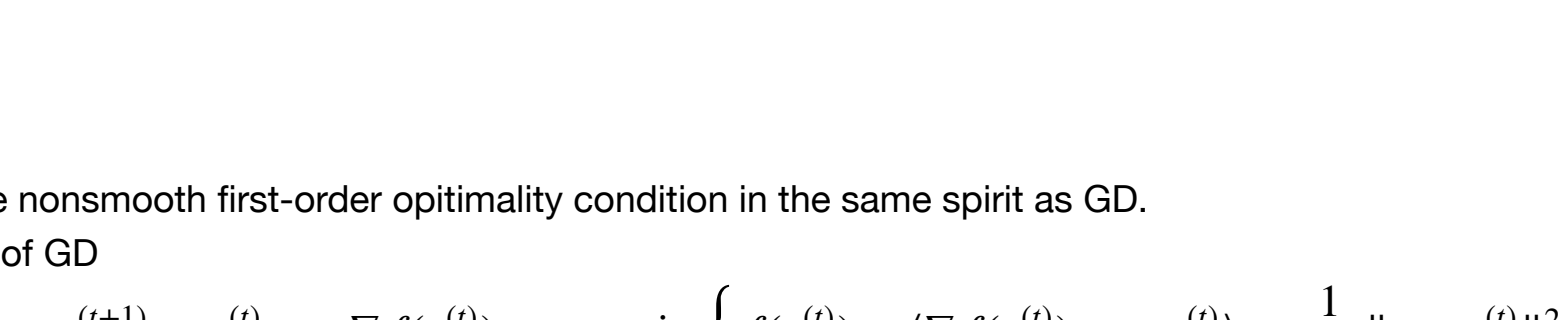
First-order optimality condition

- Assume convexity of f and g . Then \mathbf{x}^* minimizes $f(\mathbf{x}) + g(\mathbf{x})$ if and only if
$$\nabla f(\mathbf{x}^*) + \partial g(\mathbf{x}^*) \ni 0,$$

where $\partial g(\mathbf{x})$ is the **subdifferential** of g at \mathbf{x} :

$$\partial g(\mathbf{x}) = \{z : g(\mathbf{y}) \geq g(\mathbf{x}) + \langle z, \mathbf{y} - \mathbf{x} \rangle, \forall \mathbf{y} \in \mathbb{R}^d\}.$$

Vector $\mathbf{z} \in \partial g(\mathbf{x})$ is called a **subgradient** of g at \mathbf{x} . In particular: $g(\mathbf{x}) = |\mathbf{x}|$. $\partial g(\mathbf{0}) = [-1, 1]$.



Proximal gradient

- Proximal gradient solves the nonsmooth first-order optimality condition in the same spirit as GD.
- Motivation: alternative view of GD

$$\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} - \gamma_i \nabla f(\mathbf{x}^{(i)}) = \arg \min_{\mathbf{x} \in \mathbb{R}^d} \left\{ f(\mathbf{x}^{(i)}) + \langle \nabla f(\mathbf{x}^{(i)}), \mathbf{x} - \mathbf{x}^{(i)} \rangle + \frac{1}{2\gamma_i} \|\mathbf{x} - \mathbf{x}^{(i)}\|_2^2 \right\}$$

Thus consider

$$\mathbf{x}^{(i+1)} = \arg \min_{\mathbf{x} \in \mathbb{R}^d} \left\{ f(\mathbf{x}^{(i)}) + \langle \nabla f(\mathbf{x}^{(i)}), \mathbf{x} - \mathbf{x}^{(i)} \rangle + \frac{1}{2\gamma_i} \|\mathbf{x} - \mathbf{x}^{(i)}\|_2^2 + g(\mathbf{x}) \right\}$$

$$= \arg \min_{\mathbf{x} \in \mathbb{R}^d} \left\{ \gamma_i g(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - (\mathbf{x}^{(i)} - \gamma_i \nabla f(\mathbf{x}^{(i)}))\|_2^2 \right\}.$$

- The map

$$\mathbf{x} \mapsto \arg \min_{\mathbf{y} \in \mathbb{R}^d} \left\{ g(\mathbf{y}) + \frac{1}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 \right\}$$

is called the proximal map, proximal operator, or proximity operator of g , denoted by $\text{prox}_g(\mathbf{x})$.

- Thus the new update rule is

$$\mathbf{x}^{(i+1)} = \text{prox}_{\gamma_i g}(\mathbf{x}^{(i)} - \gamma_i \nabla f(\mathbf{x}^{(i)}))$$

and called the **proximal gradient method**.

Proximal maps

- A key to success of the proximal gradient method is the ease of evaluation of the proximal map $\text{prox}_{\gamma g}(\cdot)$. Often proximal maps have closed-form solutions.
- Examples
 - If $g(\mathbf{x}) \equiv 0$, then $\text{prox}_{\gamma g}(\mathbf{x}) = \mathbf{x}$. PG reduces to GD.
 - If $g(\mathbf{x}) = \iota_C(\mathbf{x})$ for a closed convex set C , then $\text{prox}_{\gamma g}(\mathbf{x}) = P_C(\mathbf{x})$, or the orthogonal projection of \mathbf{x} to set C . Proximal gradient reduces to projected gradient.
 - If $g(\mathbf{x}) = \lambda \|\mathbf{x}\|_1$, then $\text{prox}_{\gamma g}(\mathbf{x}) = S_{\gamma\lambda}(\mathbf{x}) = (\text{sign}(\mathbf{x}) \cdot (|\mathbf{x}| - \lambda\gamma))_{+}^{\gamma\lambda}$, the soft-thresholding operator.
 - If $g(\mathbf{x}) = \lambda \|\mathbf{x}\|_2$, then

$$\text{prox}_{\gamma g}(\mathbf{x}) = \begin{cases} (1 - \lambda\gamma/\|\mathbf{x}\|_2) \mathbf{x}, & \|\mathbf{x}\|_2 \geq \lambda\gamma \\ \mathbf{0}, & \text{otherwise} \end{cases}$$

- If $g(\mathbf{X}) = \lambda \|\mathbf{X}\|_*$ (nuclear norm), then $\text{prox}_{\gamma g}(\mathbf{X}) = \text{Udiag}((\sigma_1 - \lambda\gamma)_+, \dots, (\sigma_r - \lambda\gamma)_+, \mathbf{0}) \mathbf{V}^T$, when the SVD of \mathbf{X} is $\mathbf{X} = \mathbf{U} \text{diag}(\sigma_1, \dots, \sigma_r) \mathbf{V}^T$.

- Lasso

$$\min_{\beta} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1$$

- $f(\beta) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2$, $\nabla f(\beta) = \mathbf{X}^T (\mathbf{X}\beta - \mathbf{y})$
- $g(\beta) = \lambda \|\beta\|_1$.
- Update rule: $\beta^{(i+1)} = S_{\lambda\gamma_i}(\beta^{(i)} - \gamma_i \mathbf{X}^T (\mathbf{X}\beta^{(i)} - \mathbf{y}))$. Can be computed in parallel (including matrix-vector multiplications).

Convergence

- It can be shown that if

- f is convex and differentiable over \mathbb{R}^d ;
- f has L -Lipschitz gradients, i.e., $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq L \|\mathbf{x} - \mathbf{y}\|_2$;
- g is a closed convex function (i.e., $\text{epi } f$ is closed);
- $p^* = \inf_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + g(\mathbf{x}) > -\infty$ and attained at \mathbf{x}^* .

then proximal gradient method converges to an optimal solution, at the same rate (in terms of objective values) as GD. For example, with a constant step size $\gamma_i = \gamma = 1/L$,

$$f(\mathbf{x}^{(i)}) + g(\mathbf{x}^{(i)}) - p^* \leq \frac{L}{2i} \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2.$$

- A similar rate holds with backtracking.
- If f or g is strongly convex, linear convergence.

FISTA: accelerated proximal gradient

- Fast Iterative Shrinkage-Thresholding Algorithm (Beck & Teboulle, 2009)

$$\mathbf{y} = \mathbf{x}^{(i)} + \frac{t-1}{t} (\mathbf{x}^{(i)} - \mathbf{x}^{(i-1)})$$

$$\mathbf{x}^{(i+1)} = \text{prox}_{\gamma_i g}(\mathbf{y} - \gamma_i \nabla f(\mathbf{y}))$$

with initialization $\mathbf{x}^{(-1)} = \mathbf{x}^{(0)}$.

- Proximal gradient version of Nesterov's (1983) accelerated gradient algorithm.
- Under the above assumptions and with a constant step size $\gamma_i = \gamma = 1/L$,

$$f(\mathbf{x}^{(N)}) + g(\mathbf{x}^{(N)}) - p^* \leq \frac{L}{2(N+1)^2} \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2.$$

Mirror descent method

- Again, constrained optimization problem

$$\min_{\mathbf{x} \in C \cap \mathbb{R}^d} f(\mathbf{x}).$$

- Second alternative view of GD

$$\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} - \gamma_i \nabla f(\mathbf{x}^{(i)}) = \arg \min_{\mathbf{x} \in \mathbb{R}^d} \left\{ f(\mathbf{x}^{(i)}) + \langle \nabla f(\mathbf{x}^{(i)}), \mathbf{x} - \mathbf{x}^{(i)} \rangle + \frac{1}{2\gamma_i} \|\mathbf{x} - \mathbf{x}^{(i)}\|_2^2 \right\}$$

$$= \arg \min_{\mathbf{x} \in \mathbb{R}^d} \left\{ \langle \nabla f(\mathbf{x}^{(i)}), \mathbf{x} \rangle + \frac{1}{2\gamma_i} \|\mathbf{x} - \mathbf{x}^{(i)}\|_2^2 \right\}$$

- or projected (proximal) gradient

$$\mathbf{x}^{(i+1)} = P_C(\mathbf{x}^{(i)} - \gamma_i \nabla f(\mathbf{x}^{(i)})) = P_C\left(\arg \min_{\mathbf{x} \in \mathbb{R}^d} \left\{ \langle \nabla f(\mathbf{x}^{(i)}), \mathbf{x} \rangle + \frac{1}{2\gamma_i} \|\mathbf{x} - \mathbf{x}^{(i)}\|_2^2 \right\}\right)$$

- This relies too much on the [Euclidean] geometry of \mathbb{R}^d : $\|\cdot\|_2 = \langle \cdot, \cdot \rangle$.
- If the distance measure $\frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$ is replaced by something else (say $d(\mathbf{x}, \mathbf{y})$) that better reflects the geometry of C , then update such as

$$\mathbf{x}^{(i+1)} = P_C^d\left(\arg \min_{\mathbf{x} \in \mathbb{R}^d} \left\{ \langle \nabla f(\mathbf{x}^{(i)}), \mathbf{x} \rangle + \frac{1}{\gamma_i} d(\mathbf{x}, \mathbf{x}^{(i)}) \right\}\right)$$

may converge faster. Here,

$$P_C^d(\mathbf{y}) = \arg \min_{\mathbf{x} \in C} d(\mathbf{x}, \mathbf{y})$$

to reflect the geometry.

Bregman divergence

- Let $\phi : \mathcal{X} \rightarrow \mathbb{R}$ be a continuously differentiable, strictly convex function defined on a vector space $\mathcal{X} \subset \mathbb{R}^d$. The **Bregman divergence** with respect to ϕ is defined by

$$B_\phi(\mathbf{x}|\mathbf{y}) = \phi(\mathbf{x}) - \phi(\mathbf{y}) - \langle \nabla \phi(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle$$

- It can be shown that $B_\phi(\mathbf{x}|\mathbf{y}) \geq 0$ and $B_\phi(\mathbf{x}|\mathbf{y}) = 0 \iff \mathbf{x} = \mathbf{y}$. But $B_\phi(\mathbf{x}|\mathbf{y}) \neq B_\phi(\mathbf{y}|\mathbf{x})$ in general.

- Examples

- $\phi(\mathbf{x}) = \frac{1}{2} \|\mathbf{x}\|_2^2$, $\mathcal{X} = \mathbb{R}^d$: $B_\phi(\mathbf{x}|\mathbf{y}) = \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$.
- $\phi(\mathbf{x}) = \sum_{i=1}^n x_i \log x_i - \sum_{i=1}^n x_i$, $\mathcal{X} = \mathbb{R}_+^n$ (generalized negative entropy): $B_\phi(\mathbf{x}|\mathbf{y}) = \sum_{i=1}^n x_i \log \frac{x_i}{y_i} - \sum_{i=1}^n x_i + \sum_{i=1}^n y_i$ (generalized Kullback-Leibler divergence).
- $\phi(\mathbf{X}) = -\log \det \mathbf{X}$, $\mathcal{X} = \mathbb{S}_{++}^d$: $B_\phi(\mathbf{X}|\mathbf{Y}) = \text{tr}(\mathbf{X}\mathbf{Y}^{-1}) - d -$