

Large-Scale Convex Optimization

Silvan Stadelmann silvasta@ethz.ch 3. Juli 2025

1 Introduction

Large Scale Problem of dimension n but iterations $\ll n$ desired

Convex One of the only problem classes that are "solvable"

Optimization with decision variable x , objective function f and

feasible set $C = \{x \in \mathbb{R}^n : g(\xi) \leq 0, h(\xi) = 0\}$

Local minimum x^* if $\exists \epsilon > 0$ s.t. $f(x^*) \leq f(x)$,
 $\forall x \in C \cap B_\epsilon(x^*), B_\epsilon(x^*) := \{\xi \in \mathbb{R}^n : |\xi - x^*| < \epsilon\}$

Proposition 1. f (lower-semi-)continuous, $f(x) \rightarrow \infty$ for $|x| \rightarrow \infty, C$ closed $\Rightarrow \exists$ of OP with: $\min_{x \in C} f(x)$ and $x^* \in \arg\min_{x \in C} f(x)$

Definition 1 (Lipschitz continuity). $q : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is Lipschitz with constant L if: $|q(x) - q(y)| \leq L|x - y| \forall x, y \in \mathbb{R}^n$

f is **Lipschitz (Lip)** with constant $L \Leftrightarrow |\nabla f(x)|_2 \leq L$

OP class \mathcal{P} with $C = [0, 1]^n, f$ is l^∞ -Lipschitz with constant L

Proposition 2. For any algorithm \exists problem in \mathcal{P} , s.t. achieving $|f(x_N) - f(x^*)| < \epsilon$ requires $N \geq (\lfloor \frac{L}{2\epsilon} \rfloor)^n - 1$

Definition 2. OP convex if, f and g_i convex functions, h affine.

Definition 3. $q : \mathbb{R}^n \rightarrow \mathbb{R}$ convex (affine) if $\forall x, y \in \mathbb{R}^n$
 $q(\theta x + (1-\theta)y) \leq \theta q(x) + (1-\theta)q(y) \quad \forall \theta \in [0, 1]$

Proposition 3. If OP convex, local minimum == global minimum

2 Convex Optimization Problem

Definition 4 (Convex Set). A set C is convex if and only if

$\theta x + (1-\theta)y \in C, \forall x, y \in C, \forall \theta \in [0, 1]$

(hyperplane || half-space) $\{x \in \mathbb{R}^n \mid a^\top x (=||\leq) b\}$

polyhedra $\{x \in \mathbb{R}^n \mid Aq \times n \preceq b \times 1, C^r \times n \times x = d^r \times 1\}$

Operations that preserve convexity (sets)

Intersection C_1, C_2 cv $\Rightarrow C_1 \cap C_2$ convex (**cv**)

Image under affine map $C \subseteq \mathbb{R}^n$ cv $\Rightarrow \{Ax + b \mid x \in C\}$ cv

Inverse loaM $C \subseteq \mathbb{R}^m$ cv $\Rightarrow \{x \in \mathbb{R}^n \mid Ax + b \in C\}$ cv

Separating Hyperplane Theorem

Theorem 1. $C \subseteq \mathbb{R}^n$ non-empty closed (**cl**) convex set, $y \notin C$
 $\rightarrow \exists a \neq 0, b \in \mathbb{R}$ s.t. $a^\top x + b < a^\top y + b, \forall x \in C$

Corollary 1. $C_{cl,cv}$: intersection of cl half-spaces that contain C

Support function

Idea represent any cl,cv set by its supporting hyperplanes

$$\sigma_C(a) = \sup_{x \in C} a^\top x \quad \text{if known, one can construct}$$

$$C = \bigcap_{a \in \mathbb{R}^n} \{x \in \mathbb{R}^n \mid a^\top x - \sigma_C(a) \leq 0\}$$

$$= \{x \in \mathbb{R}^n \mid \sup_{a \in \mathbb{R}^n} a^\top x - \sigma_C(a) \leq 0\}$$

Definition 5. $f : \mathbb{R}^n \rightarrow \mathbb{R}$ cv \Leftrightarrow epigraph of f is cv set

$$\text{epi}(f) := \{(x, t) \in \mathbb{R}^{n+1} \mid f(x) \leq t\}$$

\rightarrow this provides a link between convex sets and functions

Operations that preserve convexity (functions)

- the point wise maximum of convex functions is convex

- the sum of convex functions is convex

- $f(Ax + b)$ is convex if f is convex

Check Convexity f is convex if it is composition of simple convex function with convexity preserving operations or if

$f : \mathbb{R}^n \rightarrow \mathbb{R}$ twice differentiable, $\partial^2 f / \partial x^2 \succeq 0 \forall x \in \mathbb{R}^n$

$g : \mathbb{R} \rightarrow \mathbb{R}$ with $g(t) = f(x + tv)$ convex in $t \forall x, v \in \mathbb{R}^n$

$\rightarrow f$ convex (restriction to a line)

Extended real numbers $\bar{\mathbb{R}} = \mathbb{R} \cup \{+\infty, -\infty\}$

Indicator function $\psi_C(x) := \begin{cases} +\infty & \text{if } x \notin C \leq 0 \\ 0 & \text{if } x \in C \end{cases}$

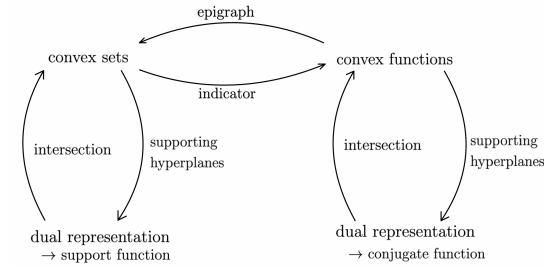
\rightarrow this provides another link between convex sets and functions

We can write $\min_{x \in C} f(x)$ as $\min_{x \in \mathbb{R}^n} f(x) + \psi_C(x)$

Definition 6 (3). $f : \mathbb{R}^n \rightarrow \bar{\mathbb{R}}$ is called proper if f is bounded below and if $\exists x \in \mathbb{R}^n$ s.t. $f(x) < \infty$

Definition 7 (Legendre Transformation). The **conjugate function** of $f : \mathbb{R}^n \rightarrow \bar{\mathbb{R}}$ is defined as $f^*(y) = \sup_{x \in \mathbb{R}^n} y^\top x - f(x)$

Concave $\nabla^2_x f^* \prec 0 \Rightarrow$ maximizer of sup satisfies $\nabla_x f^* = 0$



Theorem 2 (Conjugate of Conjugate). $f : \mathbb{R}^n \rightarrow \bar{\mathbb{R}}$

(i) f proper, cv, $\text{epi}(f)$ closed $\Rightarrow f^{**} = f$

(ii) $f(x) \geq f^{**}(x), \forall x \in \mathbb{R}^n$

3 KKT and Lagrange Duality

Basic 2d Example for derivation: $\min_{x \in \mathbb{R}^2} f(x)$ s.t. $h(x) = 0$
 $\rightarrow \nabla f(x^*), \nabla h(x^*)$ co-linear $\Leftrightarrow \exists \nu^* \in \mathbb{R} : \nabla f(x^*) + \nu^* \nabla h(x^*) = 0 \Leftrightarrow f(x) + \nu^* h(x)$ is stationary at x^*

Generalization for $n \rightarrow \infty$ and with constraints

We consider $f^* = \inf_{x \in \mathbb{R}^n} f(x)$ s.t. $g(x) \leq 0, h(x) = 0$ (1)

$$\text{Lagrange} \quad \mathcal{L}(x, \lambda, \nu) = f(x) + \lambda^\top g(x) + \nu^\top h(x)$$

$$\text{Dual Function} \quad d(\lambda, \nu) = \inf_{x \in \mathbb{R}^n} \mathcal{L}(x, \lambda, \nu)$$

Proposition 4 (Weak Duality). $d(\lambda, \nu) \leq f^*, \forall \lambda \geq 0, \nu \in \mathbb{R}^h$

Definition 8 (Constraint qualification). C convex, **Slater's Condition** holds if $\exists \hat{x} \in \mathbb{R}^n$ s.t. $h(\hat{x}) = 0$ and $g(\hat{x}) < 0$

Proposition 5 (Strong Duality). If Slater's condition holds and (1) is convex $\Rightarrow \exists \lambda \geq 0, \nu \in \mathbb{R}^n$ s.t. $d(\lambda, \nu) = f^*$

KKT (Karush-Kuhn-Tucker) Conditions

Theorem 3 (KKT Conditions). Slater's condition holds and (1) is convex $\rightarrow x^* \in \mathbb{R}^n$ is a minimizer of the primal (1) and $(\lambda^* \geq 0, \nu^*) \in \mathbb{R}^n_g \times \mathbb{R}^n_h$ is a maximizer of the dual \Leftrightarrow

$$\begin{aligned} \nabla_x \mathcal{L}(x^*, \lambda^*, \nu^*) &= 0 & \text{KKT-1 (Stationary Lagrangian)} \\ g(x^*) \leq 0, h(x^*) &= 0 & \text{KKT-2 (primal feasibility)} \\ \lambda^* \geq 0, \nu^* &\in \mathbb{R}^n_h & \text{KKT-3 (dual feasibility)} \\ \lambda^{*T} g(x^*) &= 0 = \nu^{*T} h(x^*) & \text{KKT-4 (complementary slackness)} \end{aligned}$$

In addition we have: $\sup_{\lambda \geq 0, \nu \in \mathbb{R}^n_h} q(\lambda, \nu) = \inf_{x \in C} f(x)$

Remark Without Slater, KKT-1-4 still implies x^* minimizes (1) and λ, ν maximizes dual, but the converse is no longer true.

There can be primal-minimizer/dual-maximizer not satisfy KKT.

Subdifferential

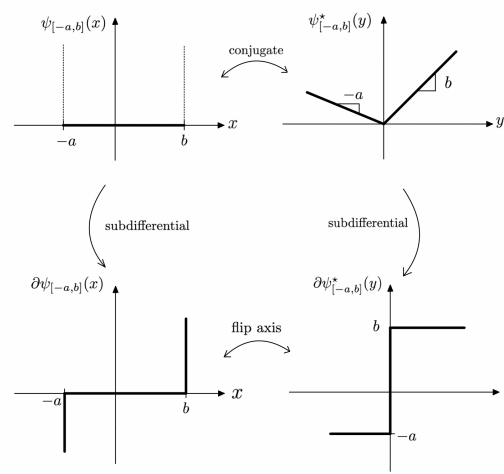
For cv f we have $f(x) \geq f(\bar{x}) + \nabla f(\bar{x})^\top (x - \bar{x}), \forall x, \bar{x} \in \mathbb{R}^n$

Definition 9. $f : \mathbb{R}^n \rightarrow \bar{\mathbb{R}}$ cv, the **subdifferential of f at \bar{x}** is:

$$\partial f(\bar{x}) := \{\lambda \in \mathbb{R}^n \mid f(x) \geq f(\bar{x}) + \lambda^\top (x - \bar{x}), \forall x \in \mathbb{R}^n\}$$

Proposition 6. f (like D9), $x^* \in \arg\min_x f(x) \Leftrightarrow 0 \in \partial f(x^*)$

Proposition 7 (Relation to conjugate functions). For convex f with $\text{epi}(f)$ closed: $y \in \partial f(x) \Leftrightarrow x \in \partial f^*(y)$



4 Convex Optimization Problems

Optimal value $f^* = \inf\{f(x) \mid g_i(x) \leq 0, h_j = 0\}$

$f^* = +\infty$ OP is infeasible, $f^* = -\infty$ OP is unbound below

Feasibility Problem

Special case $f(x) = 0, \forall x \Leftrightarrow \min_s$ s.t. $g_i(x) \leq s, h_j(x) = 0$

Linear Programming minimize $c^\top x$ s.t. $Ax - b \geq 0, x \geq 0$

Step 1: $\mathcal{L}(x, \lambda_1, \lambda_2) = c^\top x - \lambda_1^\top (Ax - b) - \lambda_2^\top x, \lambda_i \geq 0$

Step 2: $\inf_{x \in \mathbb{R}^n} \mathcal{L} = \lambda_1^\top b$, if $c - A^\top \lambda_1 - \lambda_2 = 0$, else $-\infty$

Step 3: Dual, maximize $b^\top \lambda$ s.t. $c - A^\top \lambda \geq 0, \lambda \geq 0$ (again LP)

Proposition 8. The optimal solution of a linear program (if it exists) lies always on the boundary of the feasible set and there exists an optimal solution that is a vertex of the feasible set.

Quadratic Programming convex if $P = P^\top$ positiv semi-definite minimize $\frac{1}{2} x^\top P x + q^\top x$ s.t. $Gx \leq h, Ax = b$

Second-Order Cone Program

minimize $f^\top x$ s.t. $|A_i x + b| \leq c_i^\top x + d_i, Fx = g$

Second-order cone $C_{n+1} = \{(x, t) \mid x \in \mathbb{R}^n, t \in \mathbb{R}, |x| \leq t\}$

$|A_i x + b| \leq c_i^\top x + d_i \Leftrightarrow (A_i x + b, c_i^\top x + d_i) \in C_{n+1}$

Semi-Definite Programming with symmetric F_i, X, A_i

minimize $c^\top x$ s.t. $\sum_{i=1}^n x_i F_i + G \succeq 0, Ax = b$

Standard form minimize $\text{tr}(CX)$ s.t. $X \succeq 0, \text{tr}(A_i X) = b_i$

$\text{tr}(CX) = \sum_{i=1}^n \sum_{j=1}^n C_{ij} X_{ij}, C \in \mathbb{R}^{n \times n}, i, j = 1, \dots, m$

LP \subset QP \subset QCQP (Quadratically Constrained QP) \subset SOCP \subset SDP

5 Gradient methods - Part I

Definition 10 (smoothness). $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is L -smooth (L -sm) if $\nabla f(x)$ satisfies $|\nabla f(x) - \nabla f(y)| \leq L|x - y| \forall x, y \in \mathbb{R}^n$

Taylor $\rightarrow f(y) \leq f(x) + \nabla f(x)^\top (y - x) + \frac{L}{2} |x - y|^2$

Definition 11 (strong convexity). $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is μ -strongly convex (μ -scv) if $f(y) \geq f(x) + \nabla f(x)^\top (y - x) + \frac{\mu}{2} |x - y|^2$

How to find μ/L , Spectra of Hessian $\nabla^2 f$, min/max eigenvalue

Gradient Descent

$x_{k+1} = x_k - T \nabla f(x_k)$ for $k = (k_0, \dots, k_N)$ given x_0, T

Assume $f(x) = c_0 + b^\top x + \frac{1}{2} x^\top H x, H \succ 0 \Rightarrow H x^* = -b$
 $x_{k+1} - x^* = x_k - x^* - T(b + H x_k) = (I - TH)(x_k - x^*)$

Convergence given by eigenvalues of $I - TH$, use $H = U \Lambda U^\top$

$x_N - x^* = U(I - T \Lambda)^N U^\top (x_0 - x^*) \rightarrow \text{conv-rate } 1 - T \lambda_i$

$f : L\text{-sm}, \mu\text{-scv} \rightarrow \mu \leq \min \lambda_i, \max \lambda_i \leq L \rightarrow \text{conv-rate } \rho(T)$

$=: \max_{\mu \leq h \leq L} |1 - TH| \rightarrow |x_N - x^*| \leq \rho(T)^N |x_0 - x^*|$
 $T^* = \frac{2}{L + \mu}$, with condition number $\kappa := \frac{L}{\mu}$ and $1 - \xi \leq e^{-\xi}$

$\rho(T^*) = \frac{L - \mu}{L + \mu} = \frac{\kappa - 1}{\kappa + 1} = (1 - \frac{2}{\kappa + 1}) \leq e^{-\frac{2}{\kappa + 1}} \rightarrow \text{algebraic}$

complexity $N \geq \frac{\kappa + 1}{2} \ln(\frac{|x_0 - x^*|}{\epsilon})$ to achieve $|x_N - x^*| \leq \epsilon$

Momentum-based methods

$$q_{k+1} = q_k + T p_{k+1}$$

$$p_{k+1} = (1 - 2dT) p_k - T \nabla f(q_k + \beta p_k) / L$$

Nesterov accelerated gradient $T = 1, d = \frac{1}{\sqrt{\kappa + 1}}, \beta = \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}$

Heavy Ball tuned on quadratic $T = \frac{2\sqrt{\kappa}}{\sqrt{\kappa} + 1}, d = \frac{1}{\sqrt{\kappa} + 1}, \beta = 0$

$C_{\text{Nesterov}} (1 - \frac{1}{\sqrt{\kappa}})^N \approx \frac{|q_N - q^*|}{|q_0 - q^*|} \approx C_{\text{HeavyBall}} (1 - \frac{2}{\sqrt{\kappa} + 1})^N$

Theorem 4. $f : L\text{-sm}, \mu\text{-scv} \rightarrow$ Nesterov's method satisfies:

$$|q_N - q^*| \leq \sqrt{\kappa + 1} (1 - 1/\sqrt{\kappa})^{N/2} |q_0 - q^*|$$

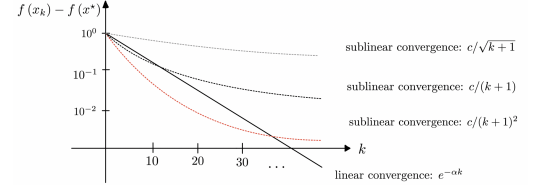
$$f(q_N) - f^* \leq \frac{L + \mu}{2} (1 - 1/\sqrt{\kappa})^N |q_0 - q^*|^2$$

Requires $N \geq 2\sqrt{\kappa} \ln(\frac{|q_0 - q^*|}{\epsilon})$ to achieve $|x_N - x^*| \leq \epsilon$

Theorem 5. For any first-order method $\exists f : \mathbb{R}^\infty \rightarrow \mathbb{R}, \mu\text{-scv}, L\text{-sm}$, s.t. $|x_k - x^*| \geq (1 - \frac{2}{\sqrt{\kappa} + 1})^k |x_0 - x^*| \forall k \geq 0$

Line search optimal step $\nu_k^* = \arg\min_{\nu \in \mathbb{R}} f(x_t - \nu \nabla f(x_t))$

Adaptive Methods gradient, momentum to adjust stepsize



6 Gradient Methods - Part II

Definition 12. $\text{prox}_C(x) = \arg\min_{y \in C} \frac{1}{2} |x - y|^2, C \subset \mathbb{R}^n$

Lemma 1. cl, cv $C \subset \mathbb{R}^n \rightarrow |\text{prox}_C(x) - \text{prox}_C(y)| \leq |x - y|$
 $\leftarrow |\text{prox}_C(x) - \text{prox}_C(y)|^2 \leq (\text{prox}_C(x) - \text{prox}_C(y))^\top (x - y)$

Projected Gradient Descent

$x_{k+1} = \text{prox}_C(x_k - T \nabla f(x_k))$, for $x_0, k_0 \dots N, T \in (0, 2/L)$

Proposition 9. $f : L\text{-sm}, \mu\text{-scv} \rightarrow$ projected GD with $T = \frac{2}{L + \mu}$ satisfies $|x_N - x^*| \leq |x_0 - x^*| (1 - \frac{2}{\kappa + 1})^N$ (κ still $\frac{L}{\mu}$)

Lemma 2. $f : \mathbb{R}^n \rightarrow \mathbb{R}, L\text{-sm}, \text{cv} \rightarrow \hat{f}$ strongly-cv
 $\hat{f}(x) = f(x) + \frac{\mu}{2} |x - x_0|^2$ and $|\hat{x}^* - x_0| \leq |x^* - x_0|$

and $f(x) - f(x^*) \leq \hat{f}(x) - \hat{f}(\hat{x}^*) + \frac{\mu}{2} |x^* - x_0|^2, \mu > 0$

\rightarrow from here one can apply GD or Nesterov, which results in:

$f(x_N) - f(x_0) \leq \epsilon$ after $N \sim L|x^* - x_0|^2/\epsilon$ iterations

Proposition 10 (Subgradient Method). cl, cv C contained in ball of radius $R, x_0 \dots N-1$ satisfy $f(\frac{1}{N} \sum_{k=0}^{N-1}) - f(x^*) \leq \frac{RLf}{\sqrt{N}}$

under $x_{k+1} = \text{prox}_C(x_k - T g_k), g_k \in \partial f(x_k), T = \frac{R}{L_f \sqrt{N}}$

f	Method	$N \sim$	Optimal
$\mu\text{-sc}, L\text{-sm}$	GD	$\kappa \ln(1/\epsilon)$	No
	NAG	$\sqrt{\kappa} \ln(1/\epsilon)$	Yes
$L\text{-smooth}$	GD	$1/\epsilon$	No
	NAG varying T	$1/\sqrt{\epsilon}$	Yes
$L\text{-Lip}, \text{cmp set}$	Subgradient	$1/\epsilon^2$	Yes

7 Stochastic Gradient Descent (SGD)

Stochastic Gradient $(\tilde{y}_i - \phi(\tilde{x}_i; \theta))^\top \frac{\partial \phi}{\partial \theta} \big|_{\tilde{x}_i; \theta}, i \sim \{1, \dots, m\}$

Problem formulation: $\min_{x \in \mathbb{R}^n} F(x) = \min_{x \in \mathbb{R}^n} \mathbb{E}[f(x, \xi)]$

$\mathbb{E}_\xi[f(x, \xi)] = \begin{cases} \int_{\mathbb{R}^q} f(x, \bar{\xi}) p_\xi(\bar{\xi}) d\bar{\xi} & \text{continuous Random V} \\ \sum_{\bar{\xi}} f(x, \bar{\xi}) p_\xi(\bar{\xi}) & \text{discrete R Variable} \end{cases}$

$\text{Var}_\xi[g(x, \xi)] = \mathbb{E}_\xi[|g(x, \xi)|^2] - |\nabla F(x)|^2$ Step 1: $\xi_k \leftarrow$ generate realization of ξ

Step 2: $x_{k+1} = x_k - T_k g(x_k, \xi_k)$, step size T_k , SG $g(\cdot)$

$\nabla_x f(x, \bar{\xi})$, $\bar{\xi} \sim p_\epsilon$ or $\frac{1}{nmb} \sum_{i=1}^{nmb} \nabla_x f(x, \xi_i)$, $\xi_i \sim p_\epsilon$

\Rightarrow The iterate x_k is now a random variable!

Assumptions on $F(x)$ and $g(x, \xi)$

A1 $F(x)$ is bounded below, ensures $\exists \min_x F(x)$ for $F: L\text{-sm}$

A2 $\mathbb{E}_\xi[g(x, \xi)] = \nabla F(x)$, $\forall x \in \mathbb{R}^n$, ensures SG unbiased.

A3 $\exists M, M_v \geq 0$ s.t. $\text{Var}_\xi[g(x, \xi)] \leq M + M_v |\nabla F(x)|^2$ $\forall x \in \mathbb{R}^n$, ensures that variance is bounded.

Proposition 11. F μ -scv L -sm, SGD const. $T < \frac{1}{L(M_v+1)}$

$\mathbb{E}[F(x_k)] - F(x^*) \leq \frac{TLM}{2\mu} + (1 - T\mu)^k (F(x_0) - F(x^*))$

$T = \frac{\ln(N)}{\mu N} \rightarrow N \sim \left(\frac{LM}{2\mu^2} + F(x_0) - F(x^*) \right) / \epsilon$

to ensure $\mathbb{E}[F(x_N)] - F(x^*) \leq \epsilon$

$(1 - T\mu)^N \leq e^{-T\mu N}$ this in EQ

The role of mini batches $M \rightarrow M/nmb$, $M_v \rightarrow M_v/nmb$

Same analysis holds, But run SGD with T/nmb to get same result... Advantage in computation if parallelization possible!

Can we do non-(strongly)-convex functions?

Proposition 12. F , L -sm, SGD with $T \leq \frac{1}{L(1+M_v)}$ achieves

$\mathbb{E}[\sum_{k=0}^{N-1} |\nabla F(x_k)|^2] \leq N T L M + \frac{2}{T} (F(x_0) - F_{\inf})$

$F_{\inf} = \inf_{x \in \mathbb{R}^n} F(x)$

SGD Table (1 and 3 optimal, except for finite-sum minimization)

F	Criterion $\leq \epsilon$	$N \sim$	$T_k \sim$
$\mu\text{-sc } L\text{-sm}$	$\mathbb{E}[F(\bar{x}_N)] - F(x_0)$	$1/\epsilon$	$1/k$
$L\text{-sm}$	$\mathbb{E}[\frac{1}{N} \sum_{k=0}^{N-1} \nabla F(x_k) ^2]$	$1/\epsilon^2$	$1/\sqrt{k}$
Lip, cv	$\mathbb{E}[F(\bar{x}_N)] - F(x_0)$	$1/\epsilon^2$	$1/\sqrt{k}$

8 ADMM

Parallelization $\min_{x \in \mathbb{R}^n} \sum_{i=1}^m f_i(x_i)$ s.t. $x_1 = \dots = x_m$

Dual ascent

Consider: $\min_{x \in \mathbb{R}^n} f(x)$ s.t. $Ax = b, A \in \mathbb{R}^{m \times n}$

Derive dual: $\mathcal{L}(x, \lambda) = f(x) + \lambda^\top Ax - \lambda^\top b$

$\inf_{x \in \mathbb{R}^n} \mathcal{L}(x, \lambda) = - \sup_{x \in \mathbb{R}^n} \{(-\lambda^\top A)x - f(x)\} - \lambda^\top b$

Dual can be stated as: $\sup_{\lambda \in \mathbb{R}^m} \underbrace{-f^*(-A^\top \lambda)}_{:=d(\lambda)} - \lambda^\top b$

Subgradient of d given by: $\partial d(\lambda) = A \partial f^*(-A^\top \lambda) - b$

Recall $v \in \partial f^*(u) \Leftrightarrow u \in \partial f(v)$ which means that the optimizer in $-\sup_{x \in \mathbb{R}^n} \{(-\lambda^\top A)x - f(x)\}$

satisfies: $-A^\top \lambda \in \partial f(x^*) \Leftrightarrow x^* \in \partial f^*(-A^\top \lambda)$

As a Result, the subgradient $\partial d(\lambda)$ can be expressed via

$\partial d(\lambda) = Ax - b$, where $x \in \text{argmin}_{\hat{x} \in \mathbb{R}^n} \{f(\hat{x}) + \hat{x}^\top A^\top \lambda\}$

Dual Subgradient Method

$x_k \in \text{argmin}_{\hat{x} \in \mathbb{R}^n} \{f(\hat{x}) + \hat{x}^\top A^\top \lambda_k\}$

$\lambda_{k+1} = \lambda_k + T_k (Ax_k - b)$, $T_k > 0$

Example $f(x) = \sum_{i=1}^m f_i(x_i)$ with $Ax = b$

$x = (x_1, \dots, x_n)$ and $A = [A_1, \dots, A_m]$

Dual subgradient becomes

$x_{k_i} \in \text{argmin}_{\hat{x}_i} \{f_i(\hat{x}_i) + \lambda_k^\top A_i \hat{x}_i\}$ (local minimization)

$\lambda_{k+1} = \lambda_k + T_k (\sum_{i=1}^m A_i x_{k_i} - b)$ (broadcasting)

Proposition 13. f convex with closed epigraph, f is μ -strongly convex if and only if f^* is $1/\mu$ -smooth.

Derive ADMM

Idea: $\min_{x \in \mathbb{R}^n} f(x) + \frac{\rho}{2} |Ax - b|^2$ s.t. $Ax = b$ with $\rho > 0$

Leads to this **Augmented Lagrangian**

$x_k = \text{argmin}_{x \in \mathbb{R}^n} f(x) + \lambda_k^\top Ax + \frac{\rho}{2} |Ax - b|^2$

$\lambda_{k+1} = \lambda_k + T_k (Ax_k - b)$ (typically $T_k = \rho$)

Advantage Improved convergence properties even if f non-scv

Disadvantage Loose of decomposability/parallelization due to augmentation with quadratic term.

This motivates ADMM which tries to combine the best of both worlds. (Well conditioned minimization and parallelization)

Consider: $\min_{x \in \mathbb{R}^n, z \in \mathbb{R}^m} f(x) + g(z)$ s.t. $Ax + Bz = c$

Augmented Objective: $\min f(x) + g(z) + \frac{\rho}{2} |Ax + Bz - c|^2$

Augmented Lagrangian: objective + $\lambda^\top (Ax + Bz - c)$

Alternating Direction Method of Multipliers

$x_k = \text{argmin}_{x \in \mathbb{R}^n} \mathcal{L}_p(x, z_{k-1}, \lambda_k)$

$z_k = \text{argmin}_{z \in \mathbb{R}^m} \mathcal{L}_p(x_k, z, \lambda_k)$

$\lambda_{k+1} = \lambda_k + \rho (Ax_k + Bz_k - c)$

Examples

$\min_{x \in \mathbb{R}^n} |Ax - b|_1 \rightarrow \min_{x \in \mathbb{R}^n, z \in \mathbb{R}^m} |z|_1$ s.t. $Ax - z = b$

LP $\min_{x \in \mathbb{R}^n} c^\top x$ s.t. $Ax = b, x \geq 0$

$x_0 \in \mathbb{R}^n$ particular solution to $Ax_0 = b$

$Q \in \mathbb{R}^{n \times (n-m)}$ whose columns span null space of A .

Any solution x to $Ax = b$ can be represented by $x = x_0 + Qw$

$\rightarrow \min_{w \in \mathbb{R}^{n-m}} c^\top Qw + c^\top x_0$ s.t. $x_0 + Qw \geq 0$

$\rightarrow \min_{w, z} c^\top Qw + c^\top x_0$ s.t. $x_0 + Qw - z = 0, z \geq 0$

$\rightarrow \min_{w, z} c^\top Qw + c^\top x_0 + \psi_Z(z)$ s.t. $x_0 + Qw - z = 0$ ($Z = \{z \in \mathbb{R}^n | z \geq 0\}$)

UPDATE RULES

9 Distributed optimization with ADMM

Goal Solve s.t. each term can be handled by its own processor.

$\min_{x_1, \dots, x_N, z \in \mathbb{R}^n} \sum_{i=1}^N f_i(x_i)$ s.t. $x_i = z$ (f_i convex) (2)

Global Consensus Problem

Step 1: Augmented Lagrangian $i \in [1, N]$ to solve (2) with ADMM.

$\mathcal{L}_p(x_i, \dots, \lambda_i) = \sum_{i=1}^N f_i(x_i) + \lambda_i^\top (x_i - z) + \frac{\rho}{2} |x_i - z|^2$

$= \sum_{i=1}^N f_i(x_i) + \frac{\rho}{2} |x_i - z - \frac{1}{\rho} \lambda_i|^2 - \frac{1}{2\rho} |\lambda_i|^2$

Step 2: Formulate ADMM

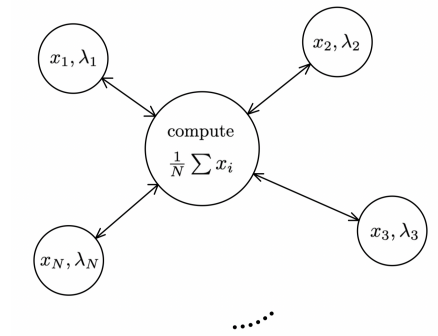
$x_i^{k+1} = \text{argmin}_{x_i \in \mathbb{R}^n} f_i(x_i) + \frac{\rho}{2} |x_i - z^k + \frac{1}{\rho} \lambda_i^k|^2$

$z^{k+1} = \text{argmin}_{z \in \mathbb{R}^n} \frac{\rho}{2} \sum_{i=1}^N |x_i^{k+1} - z + \frac{1}{\rho} \lambda_i^k|^2$

$= \frac{1}{N} \sum_{i=1}^N (x_i^{k+1} + \frac{1}{\rho} \lambda_i^k)$

$\lambda_i^{k+1} = \lambda_i^k + \rho (x_i^{k+1} - z^{k+1})$

FURTHER REFORMULATIONS...



Sharing Problem

$\min_{x_1, \dots, x_N \in \mathbb{R}^n} \sum_{i=1}^N f_i(x_i) + g\left(\sum_{i=1}^N x_i\right)$ (3)

\rightarrow copy all the variables $x_i = z_i$

\rightarrow formulate augmented Lagrangian

\rightarrow state ADMM dynamics

Optimization over Graphs

$g = (V, E)$ undirected graph with vertices V and edges E

$\min_{x \in \mathbb{R}^n} \sum_{i \in V} f_i(x) \Rightarrow \min_{x_i \in V, z_i \in E} \sum_{i \in V} f_i(x_i)$

s.t. $x_i = z_{ij}, x_j = z_{ij} \quad \forall (i, j) \in E$

Step 1: Augmented Lagrangian

Step 2: Form the Algorithm

ALGORITHM

10 Signal denoising and regression

Linear equation $y = Ax, y \in \mathbb{R}^n, x \in \mathbb{R}^m, A \in \mathbb{R}^{m \times n}$

classic setting $m \gg n$ | **modern** setting $m \ll n$ or $m \sim n$

Classic setting with outliers

l_2 -norm: $\min_{x \in \mathbb{R}^n} |Ax - y|_2^2$ results in heavy impact for outliers

l_1 -norm: $\min_{x \in \mathbb{R}^n} |Ax - y|_1$ can be reformulated and solved as

LP: $\min_{z \in \mathbb{R}^m} z^\top \mathbf{1}$ s.t. $-z \leq Ax - y \leq z, 0 \leq z$

Combined: $\phi_{\text{Hub}}(u) = \begin{cases} u^2 & \text{if } |u| \leq M \\ 2Mu - M^2 & \text{if } |u| > M \end{cases}$

Modern setting

$Ax = y$ infinite many solutions \rightarrow add regularizer to find best

Tikhonov regression: $\min_{x \in \mathbb{R}^n} |Ax - y|_2^2 + \lambda |x|_2^2$

Least Absolute Shrinkage and Selection Operator

LASSO: $\min_{x \in \mathbb{R}^n} |Ax - y|_2^2 + \lambda |x|_1$

$\Leftrightarrow \min_{x \in \mathbb{R}^n} |Ax - y|_2^2$ s.t. $|x|_1 \leq c$ results in sparse solution

11 Classification

Setup Dataset with pairs of (\bar{x}_i, \bar{y}_i) , $i = 1, \dots, N$ with data $\bar{x}_i \in \mathbb{R}^n$ and class $\bar{y}_i \in \{1, 2, \dots, K\}$

Naive approach, linear regression

Slightly improved with probabilistic approach

Linear Discriminant Analysis

SVM aims to maximize margin of decision boundary

12 Adaptive decision-making

:D