

# **ELEC-E5431: Large Scale Data Analysis**

**prof. Sergiy A. Vorobyov**

## **Outline: Large and Huge-Scale Optimization for Data Analysis**

1. Basics from Convex Sets and Functions pp. 3–5
  - Convex Sets
  - Convex Functions
  - Norms and Norm Balls
2. First-Order Methods for Smooth Functions pp. 6–12
  - Classification of Computational Problems
  - Sources of Large- and Huge-Scale Problems
  - Optimization Problem and Requirements to an Algorithm
  - Steepest Descent (Gradient Descent)
  - Multistep Methods: Heavy-Ball
  - Conjugate Gradient (CG)
  - Nesterov I Method
  - FISTA
3. Subgradient Methods pp. 13–16
  - Regularized Optimization
  - Subgradients
  - Subgradient Method for Unconstrained Optimization
  - Projected Subgradient Method
4. Stochastic Gradient Methods and Coordinate Descent pp. 17–22
  - Basic Stochastic Gradient
  - Robust Stochastic Approximation
  - Adaptive Methods
  - Coordinate Descent for Huge-Scale Optimization
  - Mirror Descent

5. Proximity Operators and Methods	pp. 23–25
<ul style="list-style-type: none"> <li>• Definition</li> <li>• Basic Proximal-Gradient Algorithm</li> <li>• Conditional Gradient "Frank-Wolfe"</li> </ul>	
6. Augmented Lagrangian Methods and ADMM	pp. 25–28
<ul style="list-style-type: none"> <li>• Augmented Lagrangian Method</li> <li>• Alternating Direction Method of Multipliers (ADMM)</li> </ul>	
7. Applications: Machine Learning and Sparsity in Data Analysis	pp. 29–36
<ul style="list-style-type: none"> <li>• Stochastic Approximation in Machine Learning</li> <li>• Using Sparsity in Data Analysis</li> <li>• Algorithms for Sparse Regularized Optimization</li> <li>• Sparsity in Machine/Statistical Learning</li> </ul>	

# 1 Basics of Convex Sets and Functions

## 1.1 Convex Sets

$S$  is a convex set if  $x, x' \in S \Rightarrow \forall \lambda \in [0, 1], \lambda x + (1 - \lambda)x' \in S$ .

$S$  is a strictly convex set if  $x, x' \in S \Rightarrow \forall \lambda \in [0, 1], \lambda x + (1 - \lambda)x' \in \text{int}(S)$ .

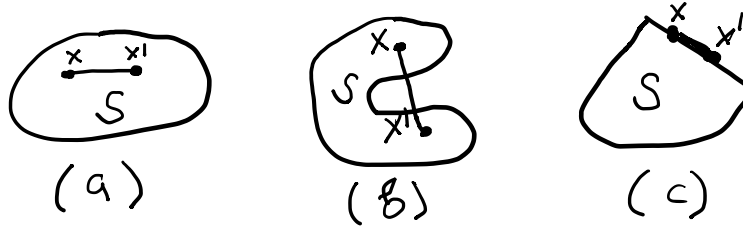


Figure 1: (a) Convex set, and without the boundary strictly convex set; (b) Non-convex set; (c) Convex, but not strictly convex set.

## 1.2 Convex Functions

Function  $f: R^n \rightarrow \bar{R} = R \cup \{+\infty\}$ .

Domain:  $\text{dom}(f) = \{x \mid f(x) \neq +\infty\}$ .

$f$  is proper if  $\text{dom}(f) \neq \emptyset$ .

$f$  is convex if  $\forall \lambda \in [0, 1], x, x' \in \text{dom}(f) \Rightarrow f(\lambda x + (1 - \lambda)x') \leq \lambda f(x) + (1 - \lambda)f(x')$ .

$f$  is strictly convex if  $\forall \lambda \in [0, 1], x, x' \in \text{dom}(f) \Rightarrow f(\lambda x + (1 - \lambda)x') < \lambda f(x) + (1 - \lambda)f(x')$ .

$f$  is strongly convex if  $\forall \lambda \in [0, 1], x, x' \in \text{dom}(f) \Rightarrow f(\lambda x + (1 - \lambda)x') \leq \lambda f(x) + (1 - \lambda)f(x') - \frac{\mu}{2}\lambda(1 - \lambda)\|x - x'\|_2^2$ .

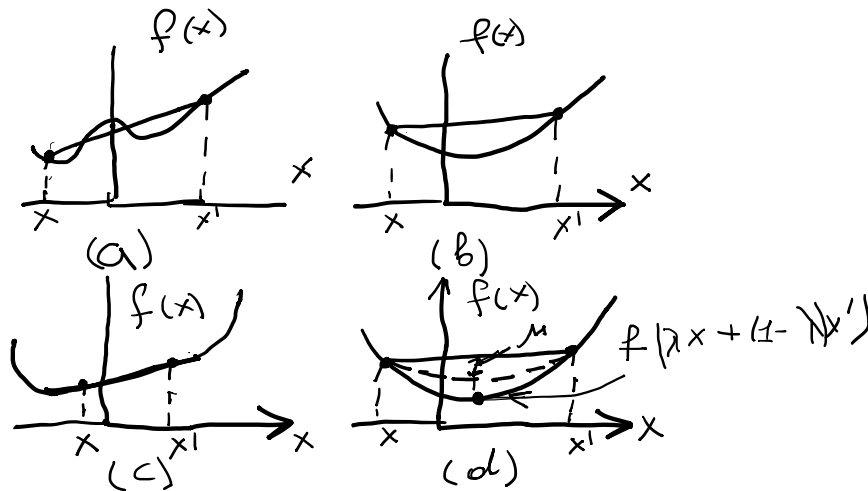


Figure 2: (a) Non-convex function; (b) Strictly convex function; (c) Convex, but not strictly convex function; (d)  $\mu$ -strongly convex function.

If  $f$  is L-Lipschitz continuous then

$$|f(x) - f(x')| \leq L\|x - x'\|, \quad \forall x, x' \in R^n.$$

So  $f$  can be bounded as follows

$$f(x') - L\|x - x'\| \leq f(x) \leq f(x') + L\|x - x'\|.$$

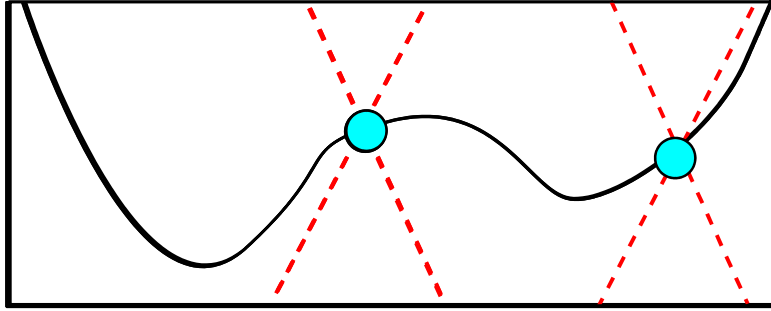


Figure 3: Lipschitz continuity.

### 1.2.1 Some Useful Properties

Let  $f_1, \dots, f_N : R^n \rightarrow R$  be convex functions. Then

1.  $f : R^n \rightarrow \overline{R}$ , defined as  $f(x) = \max\{f_1(x), \dots, f_N(x)\}$  is convex.
2.  $g : R^n \rightarrow \overline{R}$ ,  $g(x) = f_1(l(x))$  is convex, where  $l$  is affine function, i.e.,  $l(x) = Ax + b$ .
3.  $h : R^n \rightarrow \overline{R}$ ,  $h(x) = \sum_{j=1}^N \alpha_j f_j(x)$ , for  $\alpha_j > 0$ , is convex.

**Indicator function:** the indicator of set  $\mathcal{C} \subset R^n$

$$i_{\mathcal{C}} : R^n \rightarrow \overline{R}, \quad i_{\mathcal{C}} = \begin{cases} 0 & x \in \mathcal{C} \\ +\infty & x \notin \mathcal{C} \end{cases}$$

If  $\mathcal{C}$  is a closed convex set,  $i_{\mathcal{C}}$  is a lower semicontinuous convex function.

### 1.2.2 Smooth Functions

Let  $f : R^n \rightarrow R$  be twice differentiable, and consider its Hessian matrix at  $x$ ,  $\nabla^2 f(x)$ ,

$$[\nabla^2 f(x)]_{ij} = \frac{\partial^2 f(x)}{\partial x_i \partial x_j}, \quad i, j = \overline{1, n}$$

$f$  is convex  $\iff \nabla^2 f(x)$  is positive semi-definite (PSD)  $\forall x$ , that is,  $\nabla^2 f(x) \succeq 0$ .

$f$  is strictly convex  $\iff \nabla^2 f(x)$  is positive definite (PD)  $\forall x$ , that is,  $\nabla^2 f(x) \succ 0$ .

$f$  is  $\mu$ -strongly convex  $\iff \nabla^2 f(x) \succeq \mu I$ ,  $\mu > 0$ ,  $\forall x$ .

### 1.2.3 Norms and Norm Balls

$V$  is a real vector space:  $R^n$  or  $R^{n \times n}, \dots$

Function  $\|\cdot\| : V \rightarrow R$  is a norm if it satisfies:

1.  $\|x\| = 0 \Rightarrow x = 0$ .
2.  $\|\alpha x\| = |\alpha| \cdot \|x\|$ ,  $\forall x \in V$  and  $\alpha \in R$  (homogeneity).
3.  $\|x + x'\| \leq \|x\| + \|x'\|$ ,  $\forall x, x' \in V$  (triangle inequality).

**Examples:**

$V = R^n$ ,  $\|x\|_p = (\sum_i |x_i|^p)^{1/p}$  is  $l_p$ -norm, for  $p \geq 1$ .

$V = R^n$ ,  $\|x\|_\infty = \lim_{p \rightarrow \infty} \|x\|_p = \max\{|x_1|, \dots, |x_n|\}$  is Chebyshev norm.

$V = R^{n \times n}$ ,  $\|X\|_* = \text{trace}(\sqrt{X^T X})$  is matrix nuclear norm.

Also important (but not a norm):  $\|x\|_0 = \lim_{p \rightarrow 0} \|x\|_p^p = |\{i \mid x_i \neq 0\}|$ .

**Radius  $r$  ball in  $l_p$  norm:**  $\mathcal{B}_p(r) = \{x \in R^n \mid \|x\|_p \leq r\}$ .

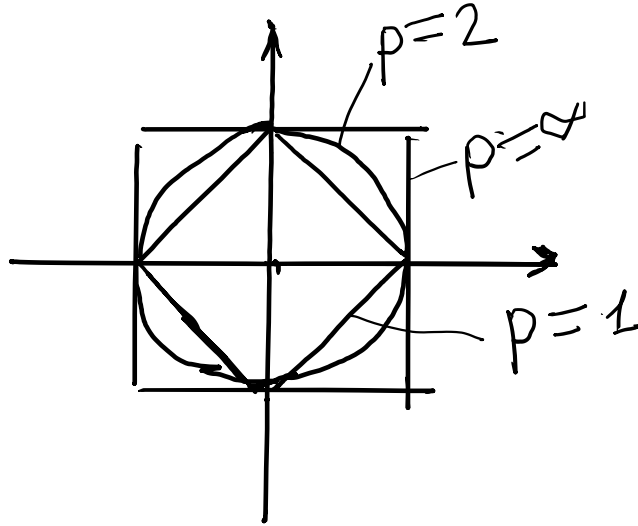


Figure 4:  $\mathcal{B}_\infty(r)$ ,  $\mathcal{B}_2(r)$ ,  $\mathcal{B}_1(r)$ .

## 2 First-Order Methods for Smooth Functions

### 2.1 Classification of Computational Problems

Class	Operation	Dimension	Iter, Cost	Memory
Small-size	All	$10^0 - 10^2$	$n^4 \rightarrow n^3$	Kilobyte: $10^3$
Medium-size	$A^{-1}$	$10^3 - 10^4$	$n^3 \rightarrow n^2$	Megabyte: $10^6$
Large-size	$Ax$	$10^5 - 10^7$	$n^2 \rightarrow n$	Gigabyte: $10^9$
Huge-size	$x + y$	$10^8 - 10^{12}$	$n \rightarrow \log n$	Terabyte: $10^{12}$

### 2.2 Sources of Large- and Huge-Scale Problems

- Internet,
- Telecommunications,
- Cyber-physical systems,
- Modern data analytics,
- Finite-element schemes (Old),
- Partial differential equations (Old),
- Economic planing in large scale (Communism?),

### 2.3 Optimization Problem and Requirements to an Algorithm

Optimization problem is

$$\min_{x \in R^n} f(x)$$

with  $f$  being smooth and convex.

Usually assume  $\mu I \preceq \nabla^2 f(x) \preceq LI$ ,  $\forall x$ ,  $0 \leq \mu \leq L$  (thus,  $L$  is a Lipschitz constant of  $\nabla^2 f$ ).

If  $\mu > 0$ , then  $f$  is  $\mu$ -strongly convex, and

$$f(y) \geq f(x) + \nabla^T f(x)(y - x) + \frac{\mu}{2} \|y - x\|_2^2.$$

Define **Condition Number** as  $\kappa \triangleq \frac{L}{\mu}$

Often interested in convex quadratic functions:

$$f(x) = \frac{1}{2} x^T A x, \quad \mu I \preceq A \preceq LI$$

$$f(x) = \frac{1}{2} \|Bx - b\|_2^2, \quad \mu I \preceq B^T B \preceq LI.$$

#### 2.3.1 Iterative Algorithms

The update rule is

$$x_{k+1} = x_k + d_k$$

where  $d_k$  depends on  $x_k$  or possibly  $(x_k, x_{k-1})$ .

Assume that we can evaluate  $f(x)$  and  $\nabla f(x)$  at each iterations.

Focus on algorithms extendable to:

- nonsmooth  $f$ ,
- $f$  not available (or too expensive to evaluate),
- only an estimate of the gradient is available (typical in machine learning when only a small portion of samples is used, mini-batch),
- a constraint  $x \in \Omega$  for simple  $\Omega$  (Ball, Box, Simplex),
- nonsmooth regularization:  $\min_x f(x) + \tau\psi(x)$ .

## 2.4 Steepest Descent (Gradient Descent)

The iterate is given as

$$x_{k+1} = x_k - \alpha_k \nabla f(x_k), \quad \alpha_k > 0.$$

How to select the step size  $\alpha_k$ :

1. Fixed: use rules based on  $L$  and  $\mu$  (trivial),
2. Backtracking (computationally easy),
3. exact line search (computationally may be hard).

For the above ways of step size selection 2 and 3, we typically have global convergence at unspecified rate.

The "greedy" strategy of getting good decrease in the current search direction may lead to better practical convergence results.

For the above way of step size selection, fixed step size selection focuses on convergence rate.

### 2.4.1 Line Search

Seek  $\alpha_k$  that satisfies Wolfe conditions:

- of sufficient decrease in  $f$ , that is,

$$f(x_{k+1}) = f(x_k - \alpha_k \nabla f(x_k)) \leq f(x_k) - c_1 \alpha_k \|\nabla f(x_k)\|_2^2, \quad 0 < c_1 < 1,$$

- while not being too small, i.e., guarantee significant increase in the directional derivative

$$\nabla^T f(x_{k+1}) \nabla f(x_k) \geq -c_2 \|\nabla f(x_k)\|_2^2, \quad c_1 < c_2 < 1.$$

It works for nonsmooth  $f$ . Thus, it is a minimizer, if  $f$  is convex.

Can do one-dimensional line search for  $\alpha_k$ , taking min of quadratic or cubic interpolations of the function and gradient.

### 2.4.2 Backtracking

No need to check the second Wolfe condition.

Latter used, but *does not work* for  $f$  nonsmooth.

### 2.4.3 Constant Step Size

Using Taylor's theorem, and the fact that  $\nabla^2 f(x) \leq LI$ , we have

$$f(x_{k+1}) \leq f(x_k) - \alpha_k \left(1 - \frac{\alpha_k}{2}L\right) \|\nabla f(x_k)\|_2^2.$$

Let  $\alpha_k \triangleq \frac{1}{L}$ , then

$$f(x_{k+1}) \leq f(x_k) - \frac{1}{2L} \|\nabla f(x_k)\|_2^2.$$

Thus,

$$\|\nabla f(x_k)\|^2 \leq 2L(f(x_k) - f(x_{k+1})).$$

Summing for  $k = 0, 1, \dots, N$ , we have

$$\sum_{k=1}^N \|\nabla f(x_k)\|_2^2 \leq 2L(f(x_0) - f(x_{N+1})).$$

Left hand side has to be, thus, finite here because in this inequality, the right hand side is finite if  $f$  is bounded below.

It follows then that  $\nabla f(x_k) \rightarrow 0$ , if  $f$  is bounded below.

### 2.4.4 Convergence Rate Analysis

Let the minimizer  $x^*$  be unique. Using Taylor's theorem, we can write

$$\|x_{k+1} - x^*\|_2^2 \leq \|x_k - x^*\|_2^2 - \alpha_k \left(\frac{2}{L} - \alpha_k\right) \|\nabla f(x_k)\|_2^2$$

so that  $\{\|x_k - x^*\|_2\}$  is decreasing.

$$\Delta_k \triangleq f(x_k) - f(x^*) \leq \nabla^T f(x_k)(x_k - x^*) \leq \|\nabla f(x_k)\|_2 \|x_k - x^*\|_2 \leq \|\nabla f(x_k)\|_2 \|x_0 - x^*\|_2.$$

The second inequality above is because of the Cauchy-Schwarz inequality.

Using the above inequality and subtracting  $f(x^*)$  from both sides of the inequality

$$f(x_{k+1}) \leq f(x_k) - \alpha_k \left(1 - \frac{\alpha_k}{2}L\right) \|\nabla f(x_k)\|_2^2$$

we obtain

$$\Delta_{k+1} \leq \Delta_k - \frac{1}{2L} \|\nabla f(x_k)\|_2^2 \leq \Delta_k - \frac{1}{2L\|x_0 - x^*\|_2^2} \Delta_k^2.$$

Take reciprocal of both sides and manipulate (using  $(1 - \varepsilon)^{-1} > 1 + \varepsilon$ ), we have

$$\frac{1}{\Delta_{k+1}} \geq \frac{1}{\Delta_k} + \frac{1}{2L\|x_0 - x^*\|_2^2} \geq \frac{1}{\Delta_0} + \frac{k+1}{2L\|x_0 - x^*\|_2^2}.$$

We could replace  $\Delta_k$  by  $\Delta_o$  because  $\Delta_k$  is decreasing as  $k \rightarrow \infty$ .

It yields

$$f(x_{k+1}) - f(x^*) \leq \frac{2L\|x_0 - x^*\|_2^2}{k+1}$$



i.e., we get the classic  $\frac{1}{k}$  convergence rate for weakly convex function. It is sublinear!

For strongly convex: assume  $\mu > 0$ , set  $\alpha_k \triangleq \frac{2}{\mu+L}$

$$\|x_k - x^*\|_2^2 \leq \left(\frac{L-\mu}{L+\mu}\right)^{2k} \|x_0 - x^*\|_2^2 = \left(1 - \frac{2}{\kappa+1}\right)^{2k} \|x_0 - x^*\|_2^2.$$

We have linear convergence then, almost always better than sublinear!

#### 2.4.5 What Convergence Rate Means?

How fast a positive sequence  $\{t_k\}_{k=1}$  of scalars is decreasing to 0.

Sublinear:  $t_k \rightarrow 0$ , but  $t_{k+1}/t_k \rightarrow 1$ .

Example:  $1/k$  rate, where  $t_k \leq Q/k$  for some constant  $Q$ .

Linear:  $t_{k+1}/t_k \leq \tau, \tau \in (0, 1)$ .

Thus, typically  $t_k \leq C\tau^k$ , (called in computer science also as "geometric" or "exponential")

Superlinear:  $t_{k+1}/t_k \rightarrow 0$ . Fast!

Example: Quadratic!  $t_{k+1} \leq Ct_k^2$ . Fast, typical of Newton's method. The number of correct significant digits doubles at each iterations.

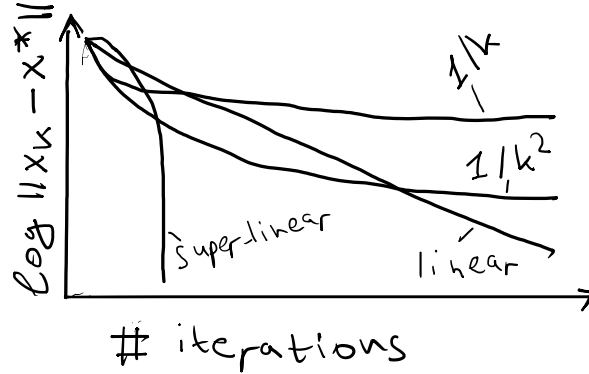


Figure 5: Sub-linear  $1/k$  and  $1/k^2$ , linear, and super-linear convergence

#### 2.4.6 Does taking $\alpha_k$ as the exact minimizer of $f$ along $-\nabla f(x_k)$ yields a better rate of convergence than $1 - \frac{2}{\kappa}$ linear rate?

**Example:**  $f(x) = \frac{1}{2}x^T Ax$ . Thus, it is obvious that  $x^* = 0, f(x^*) = 0$ .

$$\nabla f(x_k) = Ax_k.$$

Exactly minimizing with respect to (w.r.t.)  $\alpha_k$ :

$$\alpha_k = \arg \min_{\alpha} \frac{1}{2}(x_k - \alpha Ax_k)^T A(x_k - \alpha Ax_k) = \frac{x_k^T A^2 x_k}{x_k^T A^3 x_k} \in \left[\frac{1}{L}, \frac{1}{\mu}\right].$$

Thus,

$$f(x_{k+1}) \leq f(x_k) - \frac{1}{2} \frac{(x_k^T A^2 x_k)^2}{(x_k^T A x_k)(x_k^T A^3 x_k)}$$

Let  $z_k \triangleq Ax_k$ , then

$$\frac{f(x_{k+1}) - f(x^*)}{f(x_k) - f(x^*)} \leq 1 - \frac{\|z_k\|_2^4}{(z_k^T A^{-1} Z_k)(x_k^T A z_k)}.$$

Using Kantorovich inequality:

$$(z^T A z)(z^T A^{-1} z) \leq \frac{(L + \mu)^2}{4L\mu} \|z\|_2^4$$

obtain

$$\frac{f(x_{k+1}) - f(x^*)}{f(x_k) - f(x^*)} \leq 1 - \frac{4L\mu}{(L + \mu)^2} = \left(1 - \frac{2}{\kappa + 1}\right)^2.$$

Thus,

$$f(x_k) - f(x^*) \leq \left(1 - \frac{2}{\kappa + 1}\right)^{2k} (f(x_0) - f(x^*)).$$

No improvement in the linear rate over constant step length!

## 2.5 Multistep Methods: Heavy-Ball

Idea: Enhance the search direction using a contribution from the previous step, (heavy-ball, momentum, or two-step method)

$$x_{k+1} = x_k - \alpha \nabla f(x_k) + \beta(x_k - x_{k-1}).$$

The last term on the right hand side is the "momentum" term.

Can be analyzed by defining a composite iterate vector:

$$w_k \triangleq \begin{bmatrix} x_k - x^* \\ x_{k+1} - x^* \end{bmatrix}$$

Then we can write

$$w_{k+1} = B w_k + o(\|w_k\|_2), \quad B \triangleq \begin{bmatrix} -\alpha \nabla^2 f(x^*) + (1 + \beta)I & -\beta I \\ I & 0 \end{bmatrix}$$

(using Taylor series again).

Matrix  $B$  has the same eigenvalues as the following matrix

$$\begin{bmatrix} -\alpha \Lambda + (1 + \beta)I & -\beta I \\ I & 0 \end{bmatrix}, \quad \Lambda = \text{diag}(\lambda_1, \dots, \lambda_n) \quad \leftarrow \quad \text{eigenvalues of } \nabla^2 f(x^*).$$

Boris Polyak: Choose  $\alpha, \beta$  to explicitly minimize the max eigenvalue of  $B$ :

$$\alpha = \frac{4}{L} \frac{1}{(1 + 1/\sqrt{\kappa})^2}, \quad \beta = \left(1 - \frac{2}{\sqrt{\kappa} + 1}\right)^2$$

( $\beta$  must be larger than  $\alpha$ , i.e, very conservative in changing direction).

Linear convergence for  $\|x_k - x^*\|$  with rate

$$1 - \frac{2}{\sqrt{\kappa} + 1}.$$

For strictly convex  $f$ :

1. Steepest descent rate:  $1 - \frac{2}{\kappa}$  (linear)

2. Heavy-Ball rate:  $1 - \frac{2}{\sqrt{\kappa}}$  (linear)

To reduce  $\|x_k - x^*\|_2$  by a factor  $\varepsilon$ , need  $k$  large enough that

$$\left(1 - \frac{2}{\sqrt{\kappa}}\right)^k \leq \varepsilon \Leftrightarrow k \geq \frac{\sqrt{\kappa}}{2} |\log \varepsilon| \text{ (for steepest descent),}$$

$$\left(1 - \frac{2}{\sqrt{\kappa}}\right)^k \leq \varepsilon \Leftrightarrow k \geq \frac{\sqrt{\kappa}}{2} |\log \varepsilon| \text{ (for heavy-ball).}$$

A factor  $\sqrt{\kappa}$  difference!

*Example:* If  $\kappa = 1000$ , need 30 times fewer steps for heavy-ball. Good!

## 2.6 Conjugate Gradient (CG)

The iterate is given as

$$x_{k+1} = x_k + \alpha_k \rho_k, \quad \rho_k = -\nabla f(x_k) + \delta_k \rho_{k-1}$$

The same as heavy-ball with  $\beta_k = \frac{\alpha_k \delta_k}{\alpha_{k-1}}$ , but in CG  $\alpha_k$  and  $\beta_k$  are selected in particular way and the method does it itself.

CG can be implemented in a way that does not require knowledge (estimate) of  $L$  and  $\mu$ :

- Choose  $\alpha_k$  to minimize  $f$  along  $\rho_k$ ,
- Choose  $\delta_k$  by a variety of formulas (Fletcher-Reeves, Polak-Ribiere, etc.) all of these formulas are equivalent if  $f$  is convex quadratic, e.g.,

$$\delta_k = \frac{\|\nabla f(x_k)\|_2^2}{\|\nabla f(x_{k-1})\|_2^2}.$$

Restarting periodically with  $\rho_k = -\nabla f(x_k)$  is useful, e.g., restart every  $n$  iterations or when  $\rho_k$  is not a descent direction.

For quadratic  $f$ : convergence analysis is based on eigenvalues of  $A$  and Chebyshev polynomials (min-max argument), linear convergence with rate  $1 - \frac{2}{\sqrt{\kappa}}$  (like heavy-ball).

## 2.7 Nesterov I Method (1983)

Accelerate the rate to  $1/k^2$  for weakly convex, retain the linear rate for strongly convex  $f$ .

**Algorithm**

Initialize: Choose  $x_0, \alpha_0 \in (0, 1)$ , set  $y_0 = x_0$ .

Iterate:

$$x_{k+1} = y_k - \frac{1}{L} \nabla f(y_k).$$

Find  $\alpha_{k+1} \in (0, 1)$  such that

$$\alpha_{k+1}^2 = (1 - \alpha_{k+1})\alpha_k^2 + \frac{\alpha_{k+1}}{\kappa}$$

(parabola serves as a lower bound).

Set

$$\beta_k = \frac{\alpha_k(1 - \alpha_k)}{\alpha_k^2 + \alpha_{k+1}}.$$

$$y_{k+1} = x_{k+1} + \beta_k(x_{k+1} - x_k).$$

Still works for weakly convex ( $\kappa = \infty$ ).

The main ingredients are the parabola type bound and estimate sequence!

Separates the "gradient descent" and "momentum" step components! In fact, it is not a momentum method, a different principle of acceleration.

### 2.7.1 Convergence of Nesterov I Method

If  $\alpha_0 \geq \frac{1}{\sqrt{\kappa}}$ :

$$f(x_k) - f(x^*) \leq c_1 \min \left\{ \left(1 - \frac{1}{\sqrt{\kappa}}\right)^k, \frac{4L}{(\sqrt{L} + c_2 k)^2} \right\},$$

where constants  $c_1$  and  $c_2$  depend on  $x_0, \alpha_0, L$ .

1. Linear (heavy-ball) convergence rate for strongly convex  $f$ ,
2.  $1/k^2$  sublinear rate otherwise.

Example:  $\alpha_0 = \frac{1}{\sqrt{\kappa}}$ , Nesterov yields  $\alpha_k = \frac{1}{\sqrt{\kappa}}$ ;  $\beta_k = 1 - \frac{2}{\sqrt{\kappa}+1}$ .

## 2.8 FISTA (Beck & Teboulle 2009)

Simpler generic convergence analysis compared to Nesterov, adopted to composite objective function - proximal method. Otherwise the acceleration idea is the same as Nesterov.

**Algorithm:**

Initialize: Choose  $x_0$ ; set  $y_1 = x_0, t_1 = 1$ .

Iterate:

$$\begin{aligned} x_k &= y_k - \frac{1}{L} \nabla f(y_k) \\ t_{k+1} &= \frac{1}{2} \left( 1 + \sqrt{1 + 4t_k^2} \right) \\ y_{k+1} &= x_k + \frac{t_k - 1}{t_{k+1}} (x_k - x_{k-1}). \end{aligned}$$

For both strongly and weakly convex  $f$ , converges with  $\frac{1}{k^2}$ .

When  $L$  is not known, increase an estimate of  $L$  until it is big enough.

### 3 Subgradient Methods

For nonsmooth (nondifferentiable functions)!

- developed in USSR (Ukraine) in 60's-70's by Naum Shor and others;
- can be very slow in converge;
- can be applied widely: composite objective or where interior-point methods cannot be used;
- can be used to decompose a large problem into many smaller problems, (significant in internet optimization, network utility maximization, dynamic spectrum management, etc.).

#### 3.1 Regularized Optimization

The problem is given as

$$\min_x f(x) + \tau \Psi(x)$$

how to change/modify the methods if  $f$  is convex and smooth,  $\Psi$  is convex, but usually nonsmooth!

Often, all that is needed is to change the update step to

$$x_{k+1} = \arg \min_x \frac{1}{2\alpha_k} \|x - \underbrace{(x_k + \alpha_k d_k)}_{x_{k+1}}\|_2^2 + \tau \Psi(x)$$

where  $d_k$  can be a scaled gradient descent step or something more complicated (heavy-ball...), while  $\alpha_k$  is the step size as before. This is the shrinkage/thresholding step!

How to solve it with a nonsmooth  $\Psi$ ?

#### 3.2 Subgradients

For each  $x \in \text{dom}(f)$ ,  $g$  is a subgradient of  $f$  at  $x$  if

$$f(z) \geq f(x) + g^T(z - x), \quad \forall z \in \text{dom}(f)$$

Right-hand side here is a supporting hyperplane.

The set of subgradients is called subdifferential, denoted by  $\partial f(x)$ .

When  $f$  is differentiable at  $x$ :  $\partial f(x) = \{\nabla f(x)\}$ .

*Strong convexity* (in non-smooth case) with modulus  $\mu > 0$  if

$$f(z) \geq f(x) + g^T(z - x) + \frac{1}{2}\mu\|z - x\|_2^2, \quad \forall x, z \in \text{dom}(f)$$

with  $g \in \partial f(x)$ .

It generalizes the assumption  $\nabla^2 f(x) \geq \mu I$  made earlier for smooth functions.

A function  $f$  is called subdifferentiable at  $x$  if at least one subgradient of  $f$  exists at  $x$ .

Function  $f$  is subdifferentiable if it is subdifferentiable at all  $x \in \text{dom}(f)$ .

**Example:** Absolute value

$$f(x) = |x|.$$

Subgradient:

$$g = \begin{cases} 1, & x > 0 \\ -1, & x < 0 \\ \text{any value between } -1 \text{ and } 1, & x = 0. \end{cases}$$

Subdifferential:

$$\partial f(x) = \begin{cases} \{1\}, & x > 0 \\ \{-1\}, & x < 0 \\ [-1, 1], & x = 0. \end{cases}$$

### 3.2.1 Basic Properties of Subgradients

1.  $\partial f(x)$  is a closed convex set.
2. If  $f$  is convex and  $x \in \text{int dom } f$ , then  $\partial f(x)$  is nonempty and bounded.
3. If  $f$  is convex and differentiable:  $\partial f(x) = \{\nabla f(x)\}$ .
4. If  $f$  is convex and  $\partial f(x) = \{g\}$ , then  $f$  is differentiable at  $x$ .
5.  $x^*$  is a minimizer of a convex  $f$  iff  $f$  is subdifferentiable at  $x^*$  and  $0 \in \partial f(x^*)$ .

### 3.2.2 Calculus of Subgradients

Properties:

1. Nonnegative scaling: for  $\alpha \geq 0$

$$\partial(\alpha f)(x) = \alpha \partial f(x).$$

2. Sum  $f = f_1 + \dots + f_m$ . If  $f_i$  are all convex

$$\partial f(x) = \partial f_1(x) + \dots + \partial f_m(x)$$

is convex.

The same applies to integrals.

3. Affine transformation of domain:  $f$  is convex and let  $h(x) = f(Ax + b)$

$$\partial h(x) = A^T \partial f(Ax + b).$$

4. Pointwise maximum:  $\{f_1, \dots, f_m\}$  are convex, and let  $f(x) = \max_i f_i(x)$

$$\partial f(x) = \text{conv} \cup \{\partial f_i(x) \mid f_i(x) = f(x)\}.$$

**Examples:**

1.  $f(x) = \max_i a_i^T x + b_i$ .

Let  $f_i(x) = a_i^T x + b_i$ . We have  $\partial f_i(x) = \{a_i\}$ .

Let  $\mathcal{K}(x) = \{j \mid a_j^T x + b_j = \max_i a_i^T x + b_i\}$ , then

$$\partial f(x) = \text{conv} \cup_{j \in \mathcal{K}(x)} \{a_j\}.$$

In particular, when  $\mathcal{K}(x) = \{k\}$ , we have  $\partial f(x) = \{a_k\}$ .

$$2. f(x) = \|x\|_1 = \underbrace{|x_1|}_{f_1} + \dots + \underbrace{|x_n|}_{f_n}$$

$$\partial f(x) = \partial f(x_1) + \dots + \partial f_n(x) = \{g \mid g_i = 1 \text{ if } x_i > 0, g_i = -1 \text{ if } x_i < 0, g_i \in [-1, 1] \text{ if } x_i = 0\}.$$

Alternatively:  $f(x) = \max_{s \in \{-1, 1\}^n} \underbrace{s^T x}_{f_s(x)}$  and

$$\partial f(x) = \text{conv} \cup \{\partial f_s(x) \mid s^T x = \|x\|_1, s \in \{-1, 1\}^n\} = \{s \mid s^T x = \|x\|_1, s \in [-1, 1]^n\}$$

or simply  $\text{sign}\{x\}$  is a subgradient of  $f$  at  $x$ .

3. Pointwise maximum can be extended to supremum.

Suppose  $f(x) = \sup_{\alpha \in \mathcal{A}} f_\alpha(x)$  where  $f_\alpha$  are subdifferentiable and  $\mathcal{A}$  is compact. Then

$$\partial f(x) = \text{conv} \cup \{\partial f_\alpha(x) \mid f_\alpha(x) = f(x)\}.$$

### 3.3 Subgradient Method for Unconstrained Optimization

Given:  $\{\alpha_k\}$  step size sequence,  $x^{(0)}$  initial point,  $k = 0$ , and  $i_{\text{best}} = 0$ .

Repeat:

$$x_{k+1} = x_k - \alpha_k g_k$$

where  $g_k$  is any subgradient of  $f$  at  $x_k$

$$k := k + 1$$

$$f_{\text{best}}^{(k)} = \min\{f_{\text{best}}^{(k-1)}, f(x_k)\}$$

If  $f(x_k) = f_{\text{best}}^{(k)}$  then  $i_{\text{best}} := k$

Continue until a stopping criterion is satisfied

Output:  $x_{i_{\text{best}}}$

Choose the best point among the generated sequence  $x_1, x_2, \dots$ .

#### 3.3.1 Step Size Rules

- Constant step size:  $\alpha_k = \alpha$ .
- Constant step length:  $\alpha_k = \frac{\delta}{\|g_k\|_2}$ ,  $\delta > 0$ .
- Square summable, but not summable:  $\alpha_k \geq 0$ ,  $\sum_{k=1}^{\infty} \alpha_k^2 < \infty$ ,  $\sum_{k=1}^{\infty} \alpha_k = \infty$ .

**Example:**  $\alpha_k = \frac{a}{b+k}$ ,  $a, b > 0$ .

- Not summable diminishing:  $\alpha_k \geq 0$ ,  $\lim_{k \rightarrow \infty} \alpha_k = 0$ ,  $\sum_{k=1}^{\infty} \alpha_k = \infty$ .

**Example:**  $\alpha_k = \frac{a}{\sqrt{k}}$ ,  $a > 0$ .

### 3.3.2 Convergence

$f^* = \inf_x f(x)$ , and  $G$  is such that  $\|g_k\|_2 \leq G, \forall k$ .

1. Constant step size  $\alpha_k = \alpha$

$$\lim_{k \rightarrow \infty} f_{\text{best}}(k) - f^* \leq \frac{G^2 \alpha}{2}.$$

2. Constant step length

$$\lim_{k \rightarrow \infty} f_{\text{best}}(k) - f^* \leq \frac{G\delta}{2}.$$

3. Square summable, but not summable and Not summable diminishing

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

Given tolerance  $\varepsilon$ , the number of iterates  $k$  for achieving  $f_{\text{best}}(k) - f^* < \varepsilon$  can be found.

### 3.4 Projected Subgradient Method

The problem is

$$\min_{x \in \mathcal{C}} f(x)$$

where  $\mathcal{C}$  is convex set.

The iterates are obtained by

$$x_{k+1} = P_{\mathcal{C}}(x_k - \alpha_k g_k)$$

where  $P_{\mathcal{C}}$  is the Euclidean projection to  $\mathcal{C}$

$$P_{\mathcal{C}}(y) = \arg \min_{y \in \mathcal{C}} \|y - a\|_2^2.$$

The convergence result is similar to that of the basic subgradient method.

Slow!

Projected subgradient is efficient when the projection on  $\mathcal{C}$  can be easily computed.

#### 3.4.1 Simple Sets

- Affine set: linear projection;
- Half space: similar to affine set;
- $\mathcal{C} = R_+^n$ , a box  $\mathcal{C} = \{x \mid -1 \leq x_i \leq 1, \forall i\}$ , projection is truncation;
- $\mathcal{C} = \{x \mid \|x\|_2 \leq 1\}$ : projection is just a rescaling;
- Ellipsoid: no closed form, but can be easily computed;
- Simplex  $\mathcal{C} = \{x \geq 0 \mid x^T \mathbf{1} \leq 1\}$ : no closed form, but can be easily computed;
- Cone of PSD metrics: projection is to discard eigen-components  $\mathcal{C}$  (when simple).



## 4 Stochastic Gradient Methods and Coordinate Descent

### 4.1 Basic Stochastic Gradient

Requirements to a method:

- should allow  $f$  nonsmooth;
- address the situation when we cannot get function value  $f(x)$  easily;
- at any feasible  $x$ , have access only to a cheap unbiased estimate of an element of the subgradient  $\partial f$ .

**Example:**  $f(x) = E_{\xi}\{F(x, \xi)\}$  where  $\xi$  is a random vector with distribution  $D$  over a set  $\Theta$ .

Very important and hot nowadays in machine learning, but dates back to Robinson-Monro, 1951.

Special case of high interest (machine learning):  $f(x) = \frac{1}{m} \sum_{i=1}^m f_i(x)$ , where each  $f_i$  is convex and smooth.

#### 4.1.1 Classic Stochastic Gradient

**Idea:** For the finite-sum objective, get a cheap unbiased estimate of the gradient  $\nabla f(x)$  by choosing an index  $i \in \{1, \overline{m}\}$  uniformly at random, and use  $\nabla f_i(x)$  as an estimate of  $\nabla f(x)$ .

**Basic SA scheme:** At iteration  $k$ , choose  $i_k$  at uniformly (and identically and independently distributed) random from  $\{1, \overline{m}\}$ , choose some  $\alpha_k$ , and set the update as

$$x_{k+1} = x_k - \alpha_k \nabla f_{i_k}(x_k).$$

#### 4.1.2 Convergence Rate

When  $f$  is strongly convex, the analysis of convergence of expected square error  $E\{\|x_k - x^*\|_2^2\}$  is by Nemirovski et al. 2009.

Define

$$a_k \triangleq \frac{1}{2} E\{\|x_k - x^*\|_2^2\}.$$

Assume there is  $M > 0$  such that

$$\frac{1}{m} \sum_{i=1}^m \|\nabla f_i(x)\|_2^2 \leq M. \quad (*)$$

It is easy to see that

$$\frac{1}{2} \|x_{k+1} - x^*\|_2^2 = \frac{1}{2} \|x_k - \alpha_k \nabla f_{i_k}(x_k) - x^*\|_2^2 = \frac{1}{2} \|x_k - x^*\|_2^2 - \alpha_k (x_k - x^*)^T \nabla f_{i_k}(x_k) + \frac{1}{2} \alpha_k^2 \|\nabla f_{i_k}(x_k)\|_2^2.$$

Take expectation and rewrite it as

$$a_{k+1} \leq a_k - \alpha_k E\{(x_k - x^*)^T \nabla f_{i_k}(x_k)\} + \frac{1}{2} \alpha_k^2 M^2$$

( $\leq$  because of the assumption  $(*)$ ).

For middle term, we have

$$E\{(x_k - x^*)^T \nabla f_{i_k}(x_k)\} = E_{i_{[k-1]}} E_{i_k} \{(x_k - x^*)^T \nabla f_{i_k}(x_k) \mid i_{[k-1]}\} = E_{i_{[k-1]}} (x_k - x^*)^T g_k,$$

$$i_{[k-1]} = \{i_1, i_2, i_{k-1}\},$$

$$g_k = E_{i_k} \{\nabla f_{i_k}(x_k) \mid i_{[k-1]}\} \in \partial f(x_k).$$

By strong convexity, we have

$$(x_k - x^*)^T g_k \geq f(x_k) - f(x^*) = \frac{1}{2}\mu\|x_k - x^*\|_2^2 \geq \mu\|x_k - x^*\|_2^2.$$

Thus, by taking expectation, we get

$$E\{(x_k - x^*)^T g_k\} \geq 2\mu a_k.$$

Then

$$a_{k+1} \leq (1 - 2\mu\alpha_k)a_k + \frac{1}{2}\alpha_k^2 M^2.$$

Result: when  $\alpha_k = \frac{1}{k\mu}$ , we have  $\frac{1}{k}$  convergence rate!

$$a_k \leq \frac{Q}{2k}, \quad Q = \max \left\{ \|x_1 - x^*\|^2, \frac{M^2}{\mu^2} \right\}.$$

Proof: Trivial for  $k = 1$ . Otherwise:

$$\begin{aligned} a_{k+1} &\leq (1 - 2\mu\alpha_k)a_k + \frac{1}{2}\alpha_k^2 M^2 \leq \left(1 - \frac{2}{k}\right)a_k + \frac{M^2}{2k^2\mu^2} \leq \left(1 - \frac{2}{k}\right)\frac{Q}{2k} + \frac{Q}{2k^2} \\ &= \frac{k-1}{2k^2}Q = \frac{k^2-1}{k^2} \frac{Q}{2(k+1)} \leq \frac{Q}{2(k+1)}. \end{aligned}$$

Have to know  $\mu$  to set  $\alpha_k = \frac{1}{\mu k}$ ! But HOW? An underestimate of  $\mu$  can greatly degrade the performance!

## 4.2 Robust Stochastic Approximation

At iteration  $k$ :

Set as before

$$x_{k+1} = x_k - \alpha_k \nabla f_{i_k}(x_k).$$

Set

$$\bar{x}_k = \frac{\sum_{i=1}^k \alpha_i x_i}{\sum_{i=1}^k \alpha_i}$$

for any  $\theta > 0$ ,  $\alpha_k = \frac{\theta}{M - \sqrt{k}}$ .

Result: Then  $f(\bar{x}_k)$  converges to  $f(x^*)$  in expectation with rate approximately  $\log k \sqrt{k}$ .

Proof: It holds true that

$$\alpha_i E\{(x_i - x^*)^T g_i\} \leq a_i - a_{i-1} + \frac{1}{2}\alpha_i^2 M^2.$$

By convexity of  $f$ , and using  $g_i \in \partial f(x_i)$ :

$$f(x^*) \geq f(x_i) g_i^T (x^* - x_i).$$

Thus,

$$\alpha_i E\{f(x_i) - f(x^*)\} \leq a_i - a_{i+1} + \frac{1}{2}\alpha_i^2 M^2.$$

By summing iterates  $i = 1, 2, \dots, k$ , telescoping, and using the fact that  $a_{k+1} > 0$ , we can write

$$\sum_{i=1}^k \alpha_i E\{f(x_i) - f(x^*)\} \leq a_1 + \frac{1}{2}M^2 \sum_{i=1}^k \alpha_i^2.$$

Dividing by  $\sum_{i=1}^k \alpha_i$ :

$$E \left\{ \frac{\sum_{i=1}^k \alpha_i f(x_i)}{\sum_{i=1}^k \alpha_i} - f(x^*) \right\} \leq \frac{a_1 + \frac{1}{2} M^2 \sum_{i=1}^k \alpha_i^2}{\sum_{i=1}^k \alpha_i}.$$

By convexity, we have

$$f(\bar{x}_k) = f \left( \frac{\sum_{i=1}^k \alpha_i x_i}{\sum_{i=1}^k \alpha_i} \right) \leq \frac{\sum_{i=1}^k \alpha_i f(x_i)}{\sum_{i=1}^k \alpha_i}$$

(the inequality is because of Jensen's inequality).

Then the fundamental bound is

$$E\{f(\bar{x}_k) - f(x^*)\} \leq \frac{a_1 + \frac{1}{2} M^2 \sum_{i=1}^k \alpha_i^2}{\sum_{i=1}^k \alpha_i}.$$

Substituting  $\alpha_i = \frac{\theta}{M\sqrt{i}}$ , we get

$$E\{f(\bar{x}_k) - f(x^*)\} \leq \frac{a_1 + \frac{1}{2} \theta^2 \sum_{i=1}^k \frac{1}{i}}{\frac{\theta}{M} \sum_{i=1}^k \frac{1}{\sqrt{i}}} \leq \frac{a_1 + \theta^2 \log(k+1)}{\frac{\theta}{M} \sqrt{k}} = M \left[ \frac{a_1}{\theta} + \theta \log(k+1) \right] \frac{1}{\sqrt{k}}.$$

Upstairs of the right hand side of the first inequality, we majorize the sum  $\sum_{i=1}^k$  by integral, and downstairs, pick the smallest of the terms in the sum  $\sum_{i=1}^k \frac{1}{\sqrt{i}}$ , to go to the second inequality.

Other variants of Robust SA:

- periodic restarting;
- averaging over a window;
- for strongly convex case, can also get  $1/k$  rate without performing iterate averaging by defining the desired threshold  $\varepsilon$  for  $a_k$  in advance, and using a constant step size.

### 4.3 Adaptive Methods

In machine learning especially, so-called adaptive methods based on Stochastic Gradient Descent (SGD) are of a great importance.

The complexity bound is, however, the same (even worse) than for SGD –  $\mathcal{O}(\varepsilon^{-4})$ .

- Clipped SGD.
- Adaptive gradient (Adagrad)

$$x_i^{k+1} = x_i^k - \frac{\gamma}{\sqrt{G_i^k + \delta}} g_i^k, \quad G_i^k = \sum_{i=1}^k (g_i')^2, \quad \delta = 10^{-8}.$$

- RMSprop, SAGA – modifications of Adagrad.
- ADAM – also includes the momentum (acceleration) and variance reduction. This is the one used for Deep Neural Networks training.

## 4.4 Coordinate Descent for Huge-Scale Optimization

It has applications (hot) recently in machine learning, stochastic and parallel algorithms.

Consider the problem

$$\min_x f(x)$$

Other names:

- block successive upper-bound minimization (BSUM);
- block coordinate descent (BCD);
- convex-concave procedure (CCCP);
- block coordinate proximal gradient (BCPG);
- expectation maximization (EM);
- non-negative matrix factorization (NMF).

Iteration  $j$  of basic coordinate descent:

1. Choose index  $i_j \in 1, 2, \dots, n$ .
2. Fix all coordinates  $i \neq i_j$ , change  $x_{i_j}$  in a way that reduces  $f$ .

Variations for the reduces step:

- Take a reduced gradient step:  $-\nabla_{i_j} f(x)$ .
- Do more rigorous search in the subspace defined by  $i_j$ .
- Minimize  $f$  in the  $i_j$  component.

For block coordinate descent: take a block of coordinates instead of a single coordinate, i.e., choose subset  $G_j \subset \{1, 2, \dots, n\}$ :

- Reduced gradient step along  $-\nabla_{G_j} f(x)$ .
- Many different heuristics for choosing  $G_j$ , typically arising naturally from a particular application.
- Constraints and regularizers complicate things! Make block partition consistent with separability of constraints/regularizers.

### 4.4.1 Deterministic and Stochastic CD

The update rule is

$$x_{j+1, i_j} = x_{j, i_j} - \alpha_j [\nabla f(x_j)]_{i_j}.$$

- Deterministic: choose  $i_j$  in fixed order (cyclic).
- Stochastic: choose  $i_j$  at random.

Convergence: Deterministic (Luo & Tseng 1992) – Linear rate (Beck & Tetruashvili, 2013).

Stochastic – linear rate (Nesterov, 2012).

## 4.5 Mirror Descent

The step from  $x_k$  to  $x_{k+1}$  can be viewed as the solution of a subproblem:

$$x_{k+1} = \arg \min_z \nabla f_{i_k}^T(x_k)(z - x_k) + \frac{1}{2\alpha_k} \|z - x_k\|_2^2$$

that is, a linear estimate of  $f$  plus a *prox-term* (more about it later).

How far we can move in the direction of gradient?

Provides a route to handling constrained problems, regularized problems, alternative *prox-functions*.

Other prox-functions can be used instead of  $(1/2)\|z - x\|_2^2$ .

Such alternatives may be well suited to particular constraint sets  $\Omega$ .

Mirror descent is in fact a generalization of SA!

Given a constraint sets  $\Omega$ , choose a norm  $\|\cdot\|$  (not necessarily Euclidean).

Define the distance-generating function  $w$  to be a strongly convex function on  $\Omega$  with modulus 1 w.r.t.  $\|\cdot\|$ , that is,

$$(w'(x) - w'(z))^T(x - z) \geq \|x - z\|^2, \forall x, z \in \omega$$

where  $w'(\cdot)$  is an element of the subdifferential.

Now define the prox-function  $V(x, z)$  as:

$$V(x, z) = w(z) - w'(x)^T(z - x).$$

It is known as the Bregman distance.

We can use Bregman distance in the subproblem above instead of  $\frac{1}{2}\|z - x\|_2^2$

$$x_{k+1} = \arg \min_{z \in \Omega} \nabla f_{i_k}^T(x_k)(z - x_k) + \frac{1}{\alpha_k} V(z, x_k).$$

Bregman distance is the deviation of  $w$  from linearity (see Fig. 5)

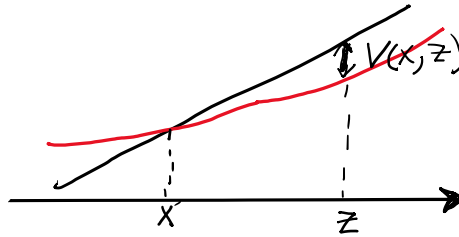


Figure 6: Example of Bregman distance.

**Example:** For any  $\Omega$ , can use

$$w(x) = \frac{1}{2} \|x - \bar{x}\|_2^2.$$

Then

$$V(x, z) = \frac{1}{2} \|x - z\|_2^2.$$

It is a "universal" prox-function.

For the simplex

$$\Omega = \left\{ x \in R^n \mid x \geq 0, \sum_{i=1}^k x_i = 1 \right\}$$

we can use instead the  $l_1$ -norm.

Choose  $w$  to be the entropy function

$$w(x) = \sum_{i=1}^n x_i \log x_i$$

Then the Bregman distance is the Kullback-Liebler divergence

$$V(x, y) = \sum_{i=1}^n z_i \log \left( \frac{z_i}{x_i} \right)$$

that is, the standard measure of distance between two probability distributions.

## 5 Proximity Operators and Methods

### 5.1 Definition

Introduced by Moreau (1962)

$$\hat{x} \in \arg \min_x \frac{1}{2} \|x - y\|_2^2 + \psi(x) \triangleq \text{prox}_\psi(y)$$

(maps  $y$  to  $x$ ).

It is well defined for convex  $\psi$ , since  $\|x - y\|_2^2$  is coercive and strictly convex.

**Example:**

$$\text{prox}_{\tau|\cdot|}(y) = \text{soft}(y, \tau) = \text{sign}(y) \max\{|y| - \tau, 0\}.$$

**Property:** Block separability  $x = (x_{[1]}, \dots, x_{[m]})$

$$\psi(x) = \sum_{i=1}^M \psi_i(x_{[i]}) \Rightarrow (\text{prox}_\psi(y))_i = \text{prox}_{\psi_i}(y_{[i]}).$$

Relationship with subdifferential:

$$y = \text{prox}_\psi(x) \iff x - y \in \partial\psi(y).$$

#### 5.1.1 Important Proximity Operators

- Soft-thresholding in the example above is the proximity operator of the  $l_1$ -norm.
- Consider the indicator function  $i_s$  of a convex set  $S$ ,

$$\text{prox}_{i_s}(y) = \arg \min_x \frac{1}{2} \|x - y\|_2^2 + i_s(x) = \arg \min_{x \in S} \frac{1}{2} \|x - y\|_2^2 = P_S(y)$$

that is Euclidean projection operator to convex set  $S$ .

- Proximity operator of squared Euclidean norm (separable, smooth).

$$\text{prox}_{(\tau/2)\|\cdot\|_2^2}(y) = \arg \min_x \frac{1}{2} \|x - y\|_2^2 + \frac{\tau}{2} \|x\|_2^2 = \frac{y}{1 + \tau}.$$

- Proximity operator of Euclidean norm (not separable, nonsmooth):

$$\text{prox}_{\tau\|\cdot\|_2}(y) = \begin{cases} \frac{x}{\|x\|_2} (\|x\|_2 - \tau), & \text{if } \|x\|_2 > \tau \\ 0, & \text{if } \|x\|_2 \leq \tau. \end{cases}$$

- Proximity operator of matrix nuclear norm

Trace/nuclear norm:

$$\|X\|_* = \sum_{i=1}^{\min\{m,n\}} \sigma_i.$$

The dual of a Schatten  $p$ -norm is a Schatten  $q$ -norm, with  $\frac{1}{q} + \frac{1}{p} = 1$ . Thus, the dual of the nuclear norm is the spectral norm:

$$\|X\|_\infty = \max\{\tau_1, \dots, \tau_{\min\{m,n\}}\}.$$

If  $Y = U\Lambda V^T$  is the singular value decomposition (SVD) of  $Y$ , we have

$$\text{prox}_{\tau\|\cdot\|_*}(Y) = U\Lambda V^T - P_{\{X \mid \max\{\sigma_1, \dots, \sigma_{\min\{m,n\}}\} \leq \tau\}}(U\Lambda V^T) = U\text{soft}(\Lambda, \tau)V^T.$$

## 5.2 Basic Proximal-Gradient Algorithm

The update at iteration  $k$  is

$$x_{k+1} = \text{prox}_{\alpha_k \tau}(x_k - \alpha_k \nabla f(x_k))$$

( $\alpha_k \tau$  is called also shrink operator).

Has different names:

- proximal gradient algorithm (PGA);
- iterative shrinkage/thresholding (IST);
- forward-backward splitting (FBS).

### 5.2.1 Convergence of PGA

Generalized PGA:

$$x_{k+1} = (1 - \lambda_k)x_k + \lambda_k(\text{prox}_{\alpha_k \tau}(x_k - \alpha_k \nabla f(x_k) - b_k) + a_k).$$

In this generalization, we allow to make small errors. Here  $a_k$  and  $b_k$  are "errors" in computing the proximity operator and the gradient, respectively, and  $\lambda_k$  is an over-relaxation parameter.

Result: Convergence is guaranteed (Combettes and Wajs, 2006) if:

- $0 < \inf \alpha_k < \sup \alpha_k < \frac{2}{L}$ ;
- $\lambda_k \in (0, 1]$  with  $\inf \lambda_k > 0$ ;
- $\sum_{k=1}^{\infty} \|a_k\| < \infty$  and  $\sum_{k=1}^{\infty} \|b_k\| < \infty$ .

### 5.2.2 FISTA in PGA/IST Form

Initialize: Choose  $\alpha \leq \frac{1}{L}$ ,  $x_0$ , set  $y_1 = x_1$ ,  $t_1 = 1$ .

Iterate:

$$\begin{aligned} x_k &= \text{prox}_{\alpha_k \tau}(y - \alpha \nabla f(x_k)) \\ t_{k+1} &= \frac{1}{2} \left( 1 + \sqrt{1 + 4t_k^2} \right) \\ y_{k+1} &= x_k + \frac{t_k - 1}{t_{k+1}}(x_k - x_{k-1}). \end{aligned}$$

Acceleration: FISTA:  $f(x_k) - f(\hat{x}) \sim O\left(\frac{1}{k^2}\right)$  versus IST:  $f(x_k) - f(\hat{x}) \sim O\left(\frac{1}{k}\right)$ .

Sensitive to correct estimate of  $L$  still. If  $L$  is not known, increase it until it is big enough. Standard  $\alpha_k \leq \frac{2}{L}$ .



### 5.3 Conditional Gradient "Frank-Wolfe"

Frank-Wolfe in 1950's, analysis Dunn 1990, for the problem

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

where  $f$  is a convex function and  $\Omega$  is a closed, bounded, convex set.

Start at  $x_0 \in \Omega$

At iteration  $k$ :

$$v_k = \arg \min_{v \in \Omega} v^T \nabla f(x_k)$$

$$x_{k+1} = x_k + \alpha_k (v_k - x_k)$$

$$\alpha_k = \frac{2}{k+2}.$$

Useful when it is easy to minimize a linear function over the original constraint set  $\Omega$ .

Rate  $\sim O\left(\frac{1}{k}\right)$ , elementary convergence theory. Line search for  $\alpha_k$  can be used.

## 6 Augmented Lagrangian Methods and ADMM

### 6.1 Augmented Lagrangian Method

#### 6.1.1 Equality Constraints

Consider linearly constrained problem:

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

where  $f$  is a proper, lower semi-continuous, convex function.

Augmented Lagrangian ( $\rho > 0$ ):

$$L(x, \lambda; \rho) = \underbrace{f(x) + \lambda^T(Ax - b)}_{\text{Lagrangian}} + \underbrace{\frac{\rho}{2}\|Ax - b\|_2^2}_{\text{Augmentation}}.$$

Method of multipliers also known as primal-dual method of multipliers:

$$x_k = \arg \min_x L(x, \lambda_{k-1}; \rho)$$

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

(can be also  $\rho_k$  that increases over iterations).

The procedure:

1. Write the problem as

$$\min_x \max_{\lambda} f(x) + \lambda^T(Ax - b).$$

2. The max w.r.t.  $\lambda$  will be  $+\infty$ , unless  $Ax = b$ , so this is equivalent to the original problem.

This equivalence is not very useful computationally because the  $\max_{\lambda}$  function is highly nonsmooth w.r.t.  $x$ .

3. Make it smooth by adding a proximal term, penalizing deviations from a prior estimate  $\bar{\lambda}$ :

$$\min_x \left\{ \max_{\lambda} f(x) + \lambda^T(Ax - b) - \frac{1}{2\rho} \|\lambda - \bar{\lambda}\|^2 \right\}. \quad (*)$$

4. Maximization w.r.t.  $\lambda$  is now trivial (for concave quadratic):

$$\lambda = \bar{\lambda} + \rho(Ax - b). \quad (**)$$

5. Inserting (\*\*) into the function in (\*), gives

$$f(x) + \bar{\lambda}^T(Ax - b) + \frac{\rho}{2}\|Ax - b\|_2^2 = L(x, \bar{\lambda}; \rho).$$

Thus, we can understand augmented Lagrangian process as:

1. Find

$$\min_x L(x, \bar{\lambda}; \rho)$$

to get new/updated  $x$ .

2. Shift the "prior" or  $\lambda$  by updating to the latest, that is,

$$\lambda = \bar{\lambda} + \rho(Ax - b).$$

3. Increase  $\rho_k$  if not happy with improvement, but ensure feasibility.
4. Repeat until convergence. Add subscripts.

### 6.1.2 Inequality Constraints

Consider the problem:

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

Use the same reasoning to the constrained minimax formulation:

$$\min_x \max_{\lambda \geq 0} f(x) - \lambda^T (Ax - b).$$

After the prox-term is added, can find the minimizing  $\lambda$  in closed form (as we did for prox-operators):

$$\lambda = \max\{\bar{\lambda} + \rho(Ax - b), 0\}.$$

Can be easily extended to nonlinear constraints:  $c(x) = 0$  or  $c(x) \geq 0$ .

Can be other constraints on  $x$  (such as  $x \in \Omega$ ) that we prefer to handle explicitly in the subproblem.

For

$$\min_x f(x) \quad \text{subject to} \quad Ax = b, x \in \Omega$$

enforce  $x \in \Omega$  explicitly in the min step

$$x_k = \arg \min_{x \in \Omega} L(x, \lambda_{k-1}; \rho)$$

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

There is an alternative way to handle inequality constraints.

Introduce slack variables  $s$ , and enforce them explicitly, i.e.,

$$\min_x f(x) \quad \text{subject to} \quad c(x) \geq 0$$

by

$$\min_{x,s} f(x) \quad \text{subject to} \quad c(x) = s, s \geq 0$$

enforced in the subproblems.

## 6.2 Alternating Direction Method of Multipliers (ADMM)

### 6.2.1 Basic ADMM

Start with separable objective

$$\min_{x,z} f(x) + h(z) \quad \text{subject to} \quad Ax + Bz = c.$$

Augmented Lagrangian:

$$L(x, z, \lambda; \rho) = f(x) + h(z) + \lambda^T (Ax + Bz - c) + \frac{\rho}{2} \|Ax + Bz - c\|_2^2. \quad (*)$$

Standard augmented Lagrangian method would minimize  $(*)$  w.r.t.  $(x, z)$  jointly. Indeed,  $x$  and  $z$  are coupled in the quadratic term, and thus, separability is lost.

In ADMM, minimums over  $x$  and  $z$  are found separately and sequentially:

$$x_k = \arg \min_x L(x, z_{k-1}, \lambda_{k-1}; \rho)$$

$$z_k = \arg \min_z L(x_k, z, \lambda_{k-1}; \rho)$$

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

It

- does one cycle of block-coordinate descent in  $(x, z)$ ;
- the min over  $x$  and  $z$  add only quadratic term to  $f$  and  $h$ , respectively, i.e., does not alter the cost much;
- can perform the  $(x, z)$  minimization inexactly;
- can add explicit (separated) constraints:  $x \in \Omega_x, z \in \Omega_z$ ;
- has many recent applications!
- has another name: operator splitting (starting from 50's).

### 6.2.2 Simpler Version of ADMM

Typically enough for many applications to consider the problem:

$$\min_{x,z} f(x) + h(z) \quad \text{subject to} \quad Ax = z.$$

Equivalently:

$$\min_x f(x) + h(Ax).$$

ADMM:

$$x_k = \arg \min_x f(x) + \frac{\rho}{2} \|Ax - z_{k-1} - \lambda_{k-1}\|_2^2$$

$$z_k = \arg \min_z h(z) + \frac{\rho}{2} \|Ax_{k-1} - z - \lambda_{k-1}\|_2^2$$

$$\lambda_k = \lambda_{k-1} + (Ax_k - z_k).$$

See that updating  $z_k$  is a proximity computation, that is,

$$z_k = \text{prox}_{h/\rho}(Ax_{k-1} - \lambda_{k-1}).$$

Updating  $x_k$  may be computationally hard ( $f$  is quadratic and involves matrix inversion).

### 6.2.3 ADMM Convergence

For the problem:

$$\min_x f(x) + h(Ax)$$

$f$  and  $h$  are lower semi-continuous, proper convex functions and  $A$  is full column rank.

Then ADMM algorithm converges (for  $\rho > 0$ ) to a solution  $x^*$ , if one exists, otherwise it diverges (Eckstein & Bertsekas, 1992).

Moreover, as in PGA, convergence is still guaranteed with inexactly solved subproblems, as long as the errors are absolutely summable.

## 7 Applications: Machine Learning and Sparsity in Data Analysis

*The numbers have no way of speaking for themselves. We speak for them, we imbue them with meaning.* – Nate Silver

### 7.1 Sparsity in Data Analysis

#### 7.1.1 Compressed Sensing

Data model:

$$\underset{n \times 1}{y} = \underset{n \times d}{A}^{\text{fat}} \underset{d \times 1}{x} + \text{noise}, \quad n \ll d.$$

Problem:

$$\hat{x} = \arg \min_w \|x\|_0 \quad \text{subject to} \quad \|Ax - y\| \leq \delta$$

is not convex problem.

Matrix  $A$  satisfies the restricted isometry property (RIP) of order  $k$  with constant  $\delta_k \in (0, 1)$ , that is, if  $\|x\|_0 \leq k \Rightarrow (1 - \delta_k)\|x\|_2^2 \leq \|Ax\| \leq (1 + \delta_k)\|x\|_2^2$ , i.e., for  $k$ -sparse vectors,  $A$  is approximately an isometric.

Alternatively, all  $k$ -column submatrices of  $A$  are nearly orthogonal!

Checking RIP is NP-hard, but satisfied by random matrices.

#### 7.1.2 Underconstrained Systems

Null space property (NSP):

Let  $\bar{x}$  be the sparsest solution of  $Ax = y$ , where  $A \in R^{m \times n}$ ,  $m < n$

$$\bar{x} = \arg \min_x \|x\|_0 \quad \text{subject to} \quad Ax = y. \quad (*)$$

Of course,  $\bar{x}$  solves  $(*)$  too, if  $\|\bar{x} - v\|_1 \geq \|\bar{x}\|_1$ ,  $\forall v \in \ker(A)$ , where  $\ker(A) = \{x \in R^n \mid Ax = 0\}$  is null space of  $A$ .

Based on works by Kashin (1977), Garnaev and Glyshkin (1984).

#### 7.1.3 Least Absolute Shrinkage and Selection Operator (LASSO)

The problem (Tibshirani, 1996) (also basis pursuit denoising (Chen, et.al 1996)):

$$\min_x \frac{1}{2} \|Ax - y\|_2^2 + \tau \|x\|_1$$

or

$$\min_x \|Ax - y\|_2^2 \quad \text{subject to} \quad \|x\|_1 \leq \tau$$

or more generally,

$$\min_x f(x) + \tau \|x\|_1$$

or

$$\min_x f(x) \quad \text{subject to} \quad \|x\|_1 \leq \tau.$$

These problems were also widely used much earlier than compressive sensing in statistics, signal processing, neural networks.

### 7.1.4 Why $l_1$ -norm Yields Sparse Solutions

Case (a):

$$\min_w \|Ax - y\|_2^2 \quad \text{subject to} \quad \|x\|_2 \leq \tau.$$

Case (b):

$$\min_w \|Ax - y\|_2^2 \quad \text{subject to} \quad \|x\|_1 \leq \tau.$$

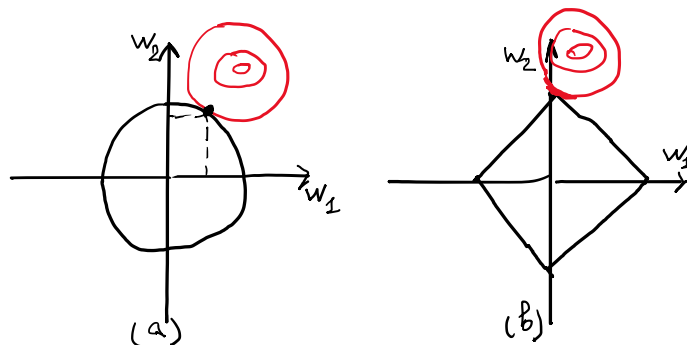


Figure 7: Sparse (b) versus non-sparse (a).

Simplest problem for case (a):

$$\hat{x} = \arg \min_x \frac{1}{2}(x - y)^2 + \frac{\tau}{2}x^2 = \frac{1}{1 + \tau}y.$$

It is called ridge regularizer. There is no sparsification of the solution.

Simplest problem for case (b):

$$\hat{x} = \arg \min_x \frac{1}{2}(x - y)^2 + \tau|x| = \text{soft}(y, \tau) = \begin{cases} y - \tau, & y > \tau \\ 0, & |y| \leq \tau \\ y + \tau, & y < -\tau. \end{cases}$$

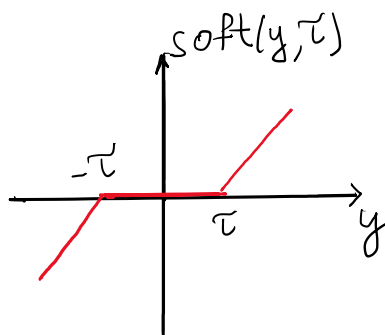


Figure 8: Soft-thresholding operator.

### 7.1.5 Structured Sparsity

Goal: to promote structured patterns, not just penalize cardinality.

Group sparsity: discard/keep entire groups of features

- density inside each group;
- sparsity w.r.t. the groups which are selected;
- choice of groups: prior knowledge about the intended sparsity patterns.

Applications:

- feature template selection;
- multi-task learning;
- learning the structure of graphical models.

### 7.1.6 Matrix Completion

From  $\mathcal{B}(X) = y \in R^p$ , find  $X \in R^{m \times n} (p < mn)$  by solving

$$\min_X \text{rank}(X) \quad \text{subject to} \quad \mathcal{B}(X) = y.$$

Noisy version:

$$\min_X \text{rank}(X) \quad \text{subject to} \quad \|\mathcal{B}(X) - y\|_2^2 \leq \delta.$$

Replace  $\text{rank}(X)$  by  $\|X\|_*$ , that is, nuclear norm.

Tikhonov formulation for nuclear norm regularization:

$$\min_X \|\mathcal{B}(X) - y\|_2^2 + \tau \|X\|_*$$

Linear observations:  $\mathcal{B} : R^{m \times n} \rightarrow R^p$ ,  $(\mathcal{B}(X))_i = \langle B_{(i)}, X \rangle$ ,

$B_{(i)} \in R^{m \times n}$  and  $\langle B, X \rangle = \sum_{i,j} b_{ij} x_{ij} = \text{trace}(B^T X)$ .

For matrix completion: each matrix  $B_{ij}$  has single 1 at the intersection of  $i$ th row and  $j$ th column, and zeros everywhere else.

Why nuclear norm favors low rank solutions?

Let  $Y = U\Lambda V^T$  be SVD of  $Y$ ,  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_{\min\{m,n\}})$ , then

$$\arg \min_X \frac{1}{2} \|Y - X\|_F^2 + \tau \|X\|_* = U \text{soft}(\Lambda, \tau) V^T.$$

It is called *singular value thresholding*. Kind of PCA, but with soft thresholding!

### 7.1.7 Atomic Norm Regularization

Atomic norm:

$$\|x\|_{\mathcal{A}} = \inf\{t > 0 \mid x \in t \cdot \text{conv}(\mathcal{A})\} = \inf \left\{ \sum_{a \in \mathcal{A}} C_a \mid x = \sum_{a \in \mathcal{A}} C_a a, C_a \geq 0 \right\}.$$

Examples:

1.  $l_1$  norm as an atomic norm.

Set of atoms is given as

$$\mathcal{A} = \left\{ \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \begin{bmatrix} -1 \\ 0 \end{bmatrix} \right\} = \{\pm e_i\}_{i=1}^N$$

where  $\pm e_i$  is  $i$ th unit vector.

In this case,  $\text{conv}(\mathcal{A}) = B_1(1)$  corresponds to the  $l_1$  unit ball. Thus,

$$\|X\|_{\mathcal{A}} = \inf \{t > 0 \mid x \in t \cdot B_1(1)\} = \|x\|_1$$

$\text{conv}(\mathcal{A}) = \text{polytope}$ ,  $\|x\|_{\mathcal{A}} = \|x\|_1$ .

2. Low-rank matrices

$$\mathcal{A} = \{A \mid \text{rank}(A) = 1, \|A\|_F = 1\}$$

$$\text{conv}(\mathcal{A}) = \text{nuclear norm ball}, \quad \|X\|_{\mathcal{A}} = \|X\|_*$$

3. Binary vectors

$$\mathcal{A} = \{\pm 1\}^N$$

$$\text{conv}(\mathcal{A}) = \text{hypercube}, \quad \|X\|_{\mathcal{A}} = \|X\|_{\infty}.$$

Examples with easy forms:

- *sparse vectors*

$$\mathcal{A} = \{\pm e_i\}_{i=1}^N$$

$\text{conv}(\mathcal{A}) = \text{cross-polytope}$

$$\|x\|_{\mathcal{A}} = \|x\|_1$$

- *low-rank matrices*

$$\mathcal{A} = \{A : \text{rank}(A) = 1, \|A\|_F = 1\}$$

$\text{conv}(\mathcal{A}) = \text{nuclear norm ball}$

$$\|x\|_{\mathcal{A}} = \|x\|_*$$

- *binary vectors*

$$\mathcal{A} = \{\pm 1\}^N$$

$\text{conv}(\mathcal{A}) = \text{hypercube}$

$$\|x\|_{\mathcal{A}} = \|x\|_{\infty}$$

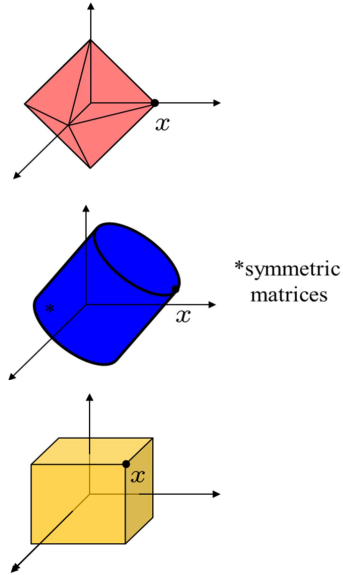


Figure 9: Balls for different widely used atomic norms.

Thus, need to discuss algorithms for sparse regularized optimization!

## 7.2 Algorithms for Sparse Regularized Optimization

### 7.2.1 Projected Subgradient for Basis Pursuit

Noseless version of basis pursuit:

$$\min_x \|x\|_1 \quad \text{such that} \quad Ax = y$$



where  $A$  is fat.

Here  $\text{sign}(x) \in \partial f(x)$ ,  $\mathcal{C} = \{x \mid Ax = y\}$  and

$$P_{\mathcal{C}}(x) = A^+y + (I - AA^+)x$$

where  $A^+ = (A^T A)^{-1} A^T$ .

Corresponding update is

$$x_{k+1} = x_k - \alpha_k (I - AA^+) \text{sign}(x_k).$$

### 7.2.2 PGA (IST) for LASSO Problem

Problem:

$$\hat{x} = \arg \min_{x \in R^n} \frac{1}{2} \|Ax - y\|_2^2 + \tau \|x\|_1$$

where  $A^T A \preceq L I$ .

PGA (IST) for LASSO becomes:

$$x_{k+1} = \text{soft}(x_k - \alpha A^T (Ax_k - y), \alpha \tau)$$

with  $\alpha < \frac{2}{L}$ .

It works, thus, in two phases:

1. After a finite number of iterations  $(x_k)_{\mathcal{Z}} = 0$ , i.e., zero components are found. The zero set  $\mathcal{Z} \subseteq \{1, \dots, n\} / \hat{x} \in G \Rightarrow \hat{x}_{\mathcal{Z}} = 0$ .
2. After zero set is found, the problem reduces to minimizing an unconstrained quadratic over the nonzero elements of  $x$ ,  $\mathcal{N} = \{1, 2, \dots, n\} / \mathcal{Z}$ . By RIP, the submatrix  $A_{\mathcal{N}}^T A_{\mathcal{N}}$  is well conditioned, so convergence is fast (linear), or much faster in practice.

### 7.2.3 FISTA (Fast IST Algorithm) for LASSO Problem

Problem: the same.

FISTA for LASSO becomes:

Initialize: Choose  $\alpha = \frac{1}{L}$ ,  $x_0, x_1, t_0 = t_1 = 1, k = 2$ .

Iterate:

$$\begin{aligned} z_k &= x_{k-1} + \frac{t_{k-1} - 1}{t_k} (x_{k-1} - x_{k-2}) \\ x_k &= \text{soft} \left( z_k - \frac{1}{L} A^T (Az_k - y), \frac{\tau}{L} \right) \\ t_{k+1} &= \frac{1}{2} \left( 1 + \sqrt{1 + 4t_k^2} \right). \end{aligned}$$

Acceleration: FISTA:  $f(x_k) - f(\hat{x}) \sim O\left(\frac{1}{k^2}\right)$  versus IST (PGA):  $f(x_k) - f(\hat{x}) \sim O\left(\frac{1}{k}\right)$ .

### 7.2.4 ADMM for LASSO Problem

Problem: the same.

ADMM form:

$$\hat{x} = \arg \min_{x,z} \frac{1}{2} \|Ax - y\|_2^2 + \tau \|z\|_1 \quad \text{s.t.} \quad x - z = 0.$$

ADMM:

$$\begin{aligned} x_k &= (A^T A + \rho I)^{-1} (A^T y + \rho z_{k-1} - \lambda_{k-1}) \\ z_k &= \min_z \tau \|z\|_1 + \lambda_{k-1}^T (x_k - z) + \frac{\rho}{2} \|z - x_k\|_2^2 = \text{prox}_{\tau/\rho} \left( x_k + \frac{\lambda_{k-1}}{\rho} \right) = \text{soft} \left( x_k + \frac{\lambda_{k-1}}{\rho}, \frac{\tau}{\rho} \right) \\ \lambda_k &= \lambda_{k-1} + \rho(x_k - z_k). \end{aligned}$$

Update for  $x_k$  is the most complicated computationally.

If  $A$  is fat, that is,  $m \times n$ ,  $n > m$ , we can make use of the Sherman-Morrison-Woodbury formula.

If  $A^T A = I$ , which is approximately the case in compressed sensing applications,  $x_k$  can be recovered at the cost of two matrix-vector multiplications involving  $A$ .

### 7.2.5 Frank-Wolfe for Atomic-Norm Constraint

We solve the problem:

$$\min f(x) \quad \text{subject to} \quad \|x\|_{\mathcal{A}} \leq \tau.$$

It is exactly the example where Frank-Wolfe is particularly useful.

The search direction

$$v_k = \tau \bar{a}_k$$

where

$$\bar{a}_k = \arg \min_{a \in \mathcal{A}} \langle a, \nabla f(x_k) \rangle.$$

In other words, we seek the atom that lines up best with the negative gradient direction  $-\nabla f(x_k)$ .

Can think of each step as the "addition of a new atom to the basis".

$x_k$  is expressed in terms of  $\{\bar{a}_0, \bar{a}_1, \dots, \bar{a}_k\}$ .

If few iterations are needed to find a solution with acceptable accuracy, we have then just an approximate solution that is represented in terms of few atoms, that is, sparse or compactly represented.

Finding new atoms is cheap!

Examples:

- For the constraint  $\|x\|_1 \leq \tau$ , the atoms are  $\{\pm e_i \mid i = 1, 2, \dots, n\}$ .

If  $i_k$  is the index at which  $|\nabla f(x_k)|_i$  attains its maximum, we have

$$\bar{a}_k = -\text{sign}([\nabla f(x_k)]_{i_k}) e_{i_k}.$$

- For the constraint  $\|x\|_\infty \leq \tau$ , the atoms are  $2^n$  vectors with entries  $\pm 1$

$$[\bar{a}_k]_i = -\text{sign}([\nabla f(x_k)]_i), \quad i = 1, 2, \dots, n$$

### 7.3 Stochastic Approximation in Machine Learning

Given data  $x_i \in R^m$  and labels  $y_i = \pm 1, i = \overline{1, m}$ , find  $w$  that minimizes

$$\frac{1}{N} \sum_{i=1}^N l(w, x_i, y_i) + \tau \psi(w)$$

where  $\psi$  is a regularizer,  $\tau > 0$  is a parameter, and  $l$  is a loss function.

Loss function has specific form of  $l(w^T x_i, y_i)$  for linear classifiers/regressors.

For example, in Support Vector Machine:

$$l(w^T x_i, y_i) = \max\{1 - y_i(w^T x_i), 0\}$$

and  $\psi = \|w\|_1$  or  $\|w\|_2$ .

Logistic classification:

$$l(w^T x_i, y_i) = \log(1 + \exp(y_i(w^T x_i))).$$

In regularized version, may have  $\psi(w) = \|w\|_1$ .

### 7.4 Sparsity in Machine/Statistical Learning

Data:  $N$  pairs  $(x_i, y_i)$ ,  $i = 1, \dots, N$ , where  $x_i \in R^m$  (feature vectors) and  $y_i \in \{+1, -1\}$  (labels).

Goal: find good in some sense classifier (start from linear).

More precise: find optimal weights in

$$\hat{y} = \text{sign}([x^T \ 1]w) = \text{sign}\left(w_{m+1} + \sum_{i=1}^m w_i x_i\right).$$

Assumption: Data are i.i.d. and generated by some underlying distribution  $P_{X,Y}$ , which is unknown.

Expected error:

$$\min_w E(1_{Y([X^T \ 1]w) < 0}).$$

Practically impossible to compute because the distribution  $P_{X,Y}$  is unknown

Empirical error:

$$\min_w \frac{1}{N} \sum_{i=1}^N h(y_i([x^T \ 1]w))$$

where  $h(z) = 1_{z < 0}$ .

It is computable, but it is neither convex nor differentiable. Thus, the problem is NP-hard.

Convexification: Replace  $h : R \rightarrow \{0, 1\}$  with  $l : R \rightarrow R_+$ .

Practical criterion

$$\min_w \sum_{i=1}^N l(y_i(w^T x_i + b)) + \tau \psi(w).$$

Regularizer: If  $\psi$  is  $l_1$ -norm, it encourages sparseness, that in terns encourages selection of fewer features.

Typical convex losses:

- Missclassification loss

$$l(z) = 1_{z < 0}.$$

Not practical.

- Hinge loss

$$l(z) = \max\{1 - z, 0\}.$$

- Logistic loss

$$l(z) = \frac{\log(1 - e^{-z})}{\log 2}.$$

- Squared loss

$$l(z) = (z - 1)^2.$$

It was for the case of linear classification. It can be easily generalized.

General formulation:

$$\min_w \sum_{i=1}^N l(y_i([x^T \ 1]w)) + \tau \psi(w)$$

Different classification strategies:

- Least squares classification

$$l(z) = (z - 1)^2, \quad \psi(w) = 0.$$

- Ridge regression classification

$$l(z) = (z - 1)^2, \quad \psi(w) = \|w\|_2^2.$$

- LASSO classification

$$l(z) = (z - 1)^2, \quad \psi(w) = \|w\|_1.$$

- Sparse logistic classification

$$l(z) = (1 - e^{-z}), \quad \psi(w) = \|w\|_1.$$

- Support vector machine

$$l(z) = \max\{1 - z, 0\}, \quad \psi(w) = \|w\|_2^2.$$

- Boosting

$$l(z) = e^{-z}.$$

- ...

Generalizing it also to nonlinear functions.

Simply use

$$\hat{y} = \phi(x, w) = \sum_{i=1}^m w_i \phi_i(x)$$

where  $\phi : R^m \rightarrow R$ , but  $\phi(x, w)$  is still linear in  $w$  that ensures the convexity for the loss  $l$ .

Examples of  $\phi : R^m \rightarrow R$ : Radial basis functions; Wavelets; Splines; Kernels; Graphs as special case of kernels, ... .