# OPTIMIZATION FOR MACHINE LEARNING
### Regularized, large-scale and distributed optimization

November 16, 2023

Today: Sparsity and LASSO

# SPARSITY AND REGULARIZATION

Motivation: Sparse models, typically because of overparametrization

(Lots of) zeros

$\begin{bmatrix} 1 \\ 0 \end{bmatrix}$ Sparse?

$\in \mathbb{R}^2$

Yes,
50% of
the coefficients
are 0

$\begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \in \mathbb{R}^{10^9 \times 10^{12}}$

Sparse!

$\frac{10^{9+12} - 1}{10^{9+12}} \times 100$

zero
coefficients

## ① Sparse regularizers

Recall:  $\underset{x \in \mathbb{R}^d}{\text{minimize}} \quad f(x) + \lambda \underline{\Omega(x)}$

$\lambda > 0$

data-fitting term

regularization term / regularizer

Q) What regularizer can we use to produce solutions that are sparser (more zero coefficients) than the solutions of

$$\underset{x \in \mathbb{R}^d}{\text{minimize}} \quad f(x) \quad ? \quad \text{(un-regularized problem)}$$

Natural choice:  "$l_0$ norm"  (read "ell zero")

Def: The $l_0$-norm in $\mathbb{R}^d$ is the function $\|\cdot\|_0 : \mathbb{R}^d \to \mathbb{R}$ defined by

$x = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$

$$\forall x \in \mathbb{R}^d, \quad \|x\|_0 = \sum_{j=1}^{d} \mathbb{1}(x_j \neq 0)$$

$\mathbb{1}(x_j \neq 0) = \begin{cases} 1 & \text{if } x_j \neq 0 \\ 0 & \text{otherwise} \end{cases}$

$\rightarrow$ $\|x\|_0$ is the number of nonzero coordinates in $x$

$\rightarrow$ $\|x\|_0 \in \{0, 1, -, d\}$

$\rightarrow$ $\forall (x, y) \in (\mathbb{R}^d)^2$, $x$ is sparser than $y$ if $\|x\|_0 < \|y\|_0$

$\rightarrow$ If we consider

$$\text{minimize } f(x) + \lambda \|x\|_0,$$
$$x \in \mathbb{R}^d$$

then the regularization term penalizes the vectors with the largest values of $\|x\|_0$

$\lambda \rightarrow \infty$

$\text{minimize } \|x\|_0$
$x \in \mathbb{R}^d$

$\nearrow x^* = \begin{bmatrix} ? \\ ? \\ 0 \end{bmatrix}$

$\searrow \min_{x \in \mathbb{R}^d} \|x\|_0 = 0$

$\rightarrow$ Issues:

key optimization challenges
- $\|\cdot\|_0$ is nonconvex
- $\|\cdot\|_0$ is discontinuous ( big issue in optimization)



Additional challenges
- $\|\cdot\|_0$ has a combinatorial structure
- $\|\cdot\|_0$ is not even a norm

$\rightarrow$ In practice, we use regularizers that approximate the $l_0$ norm and are easier to use in an optimization problem.

- The $l_0$ norm is a limiting case of a family of functions called the "$l_p$ norms"
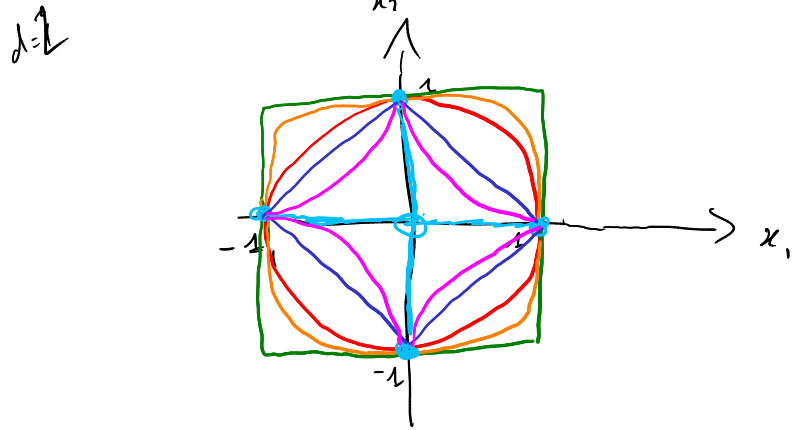
$$\forall \, p \in [0, +\infty], \qquad \|x\|_p := \begin{cases} \sum_{j=1}^{d} \mathbb{1}(x_j \neq 0) & \text{if } p = 0 \\ \max_{1 \leq j \leq d} |x_j| & \text{if } p = \infty \\ \left( \sum_{j=1}^{d} |x_j|^p \right)^{1/p} & \text{if } 0 < p < \infty \end{cases}$$

$p = 2 \qquad \|x\|_2 = \|x\| = \sqrt{\sum_{j=1}^{d} x_j^2}$

$p = 1 \qquad \|x\|_1 = \sum_{j=1}^{d} |x_j|$

$d = 2$



$\|x\|_p \underset{\substack{p \to 0 \\ p > 0}}{\longrightarrow} \|x\|_0$

$\|x\|_p \underset{p \to \infty}{\longrightarrow} \|x\|_\infty$

$\|x\|_2 = 1 \qquad \|x\|_4 = 1$

$\|x\|_\infty = 1 \qquad \|x\|_0 = 1$

$\|x\|_1 = 1 \qquad \|x\|_{1/2} = 1$

- $\|\cdot\|_p$ is a norm when $p \geq 1$
- $\|\cdot\|_p$ is a convex function when $p \geq 1$

$\Bigg\} \Rightarrow$ suggest that we could use $\|\cdot\|_p$ with $p \geq 1$ to approximate $\|\cdot\|_0$

Key fact: $\|\cdot\|_1$ is the best convex upper bound of the $\|\cdot\|_0$ norm, i.e.

$\forall$ convex function $g : \mathbb{R}^d \to \mathbb{R}$ such that

$$g(x) \geq \|x\|_0 \quad \forall x \in \mathbb{R}^d,$$

then $\qquad g(x) \geq \|x\|_1 \quad \forall x \in \mathbb{R}^d$

(Using $\|\cdot\|_p$ with $0 < p < 1$ yields a nonconvex upper bound)

↳ the most popular choice for a sparse regularizer is the $\ell_1$ norm $\Omega(x) = \|x\|_1$, also called the LASSO regularization term (typically when $f(x) = \frac{1}{2}\| - \|_2^2$)

⇒ there are numerous variations on this simple choice:

## Group LASSO (aka $\ell_1/\ell_2$ regularizer)

$$\Omega(x) = \sum_{g \in G} \|x_g\|_2 \qquad \text{where } G \text{ is a partition of}$$
$$\{1, -, d\}$$

$x_g = [x_j]_{j \in g}$

$$\left( \text{Ex} \quad \boxed{\,1 \quad d_{/1} \, d_{2}^{+1} \quad d\,} \right)$$

- Groups of coordinates $g \in G$
- Models the a priori that a given group of parameters should either be used (all non zeros) or not used (zeros)

$$G = \{ \{1\}, \{2\}, -, \{d\} \} \quad \Rightarrow \quad \|\cdot\|_1 \text{ norm}$$

$$\sum_{g \in G} \|x_g\|_2 = \|y\|_1 \quad \text{where} \quad y = \left[\|x_g\|_2\right]_{g \in G} \in \mathbb{R}^{|G|}$$

$$G = \{\{1, -, d\}\} \quad \Rightarrow \quad \sum_{g \in G} \|x_g\|_2 = \|x\|_2$$

**Remark:** The group regularizers ($\ell_1/\ell_2$) are often use to encode links between the parameters, which is problem-specific.

$$\text{minimize}_{x \in \mathbb{R}^d} \quad f(x) + \lambda \|x\|_1$$

penalizes vectors that have nonzero coordinates

$G = \{ \{1\}, \{2, \dots, d\} \}$

$$\text{minimize}_{x \in \mathbb{R}^d} \quad f(x) + \lambda \left( |x_1| + \left\| \begin{bmatrix} x_2 \\ \vdots \\ x_d \end{bmatrix} \right\|_2 \right)$$

penalizes vectors that either have nonzero first coordinate or that have a nonzero norm for

$\begin{bmatrix} x_2 \\ \vdots \\ x_d \end{bmatrix}$

$$|x_1| + \sqrt{x_2^2 + \dots + x_d^2}$$

$$\neq |x_1| + |x_2| + \dots + |x_d|$$

→ this idea generalizes further:

- Can replace the $l_2$ norm by an $l_q$ norm with $q > 1$

$$(\text{Ex}) \quad \Omega(x) = \sum_{g \in G} \|x_g\|_\infty$$

- Can use overlapping groups

$$\Omega(x) = \|x\|_2 + |x_1| + |x_2|$$

# ② The case of $\ell_1$ regularization

$$\underset{x \in \mathbb{R}^d}{\text{minimize}} \quad f(x) + \lambda \|x\|_1 \qquad \lambda > 0$$

**Special case : LASSO estimator (1990s)**

$$A \in \mathbb{R}^{m \times d}, \quad y \in \mathbb{R}^m$$

$$\underset{x \in \mathbb{R}^d}{\text{minimize}} \quad \frac{1}{2m} \|Ax - y\|_2^2 + \lambda \|x\|_1$$

↳ **First approach :** $\|\cdot\|_1$ is nonsmooth ( its gradient does not exist at certain points)

$$\Rightarrow \text{Fix : use subgradients !}$$

↳ **Second approach :** $x \mapsto \frac{1}{2m} \|Ax - y\|_2^2$ has a gradient and we want to exploit it $\Rightarrow$ Use proximal gradient !

**Proposition**

$\forall x \in \mathbb{R}^d$, the subdifferential of $\|\cdot\|_1$ at $x$ is the set of vectors $\partial \|\cdot\|_1 (x) \subseteq \mathbb{R}^d$ such that

$\forall g \in \partial \|\cdot\|_1 (x)$, $\forall j \in \{1, \dots, d\}$,

$$g_j \begin{cases} = 1 & \text{if } x_j > 0 \\ = -1 & \text{if } x_j < 0 \\ \in [-1, 1] & \text{if } x_j = 0 \end{cases}$$

$g$ is called a subgradient of $\|\cdot\|_1$ at $x$

If $x$ has $d$ nonzero coordinates ($\|x\|_0 = d$), then

$$\partial \|\cdot\|_1 (x) = \left\{ \begin{bmatrix} \text{sign}(x_1) \\ \vdots \\ \text{sign}(x_d) \end{bmatrix} \right\} \qquad \text{sign}(t) = \begin{cases} 1 & \text{if } t > 0 \\ -1 & \text{if } t < 0 \end{cases}$$

$\uparrow$
$\nabla (\|\cdot\|_1)(x)$ : at this point, the gradient is well-defined

If $x = 0_{\mathbb{R}^d}$,
$$\partial \|\cdot\|_1 (0_{\mathbb{R}^d}) = \{ g \in \mathbb{R}^d \mid g_j \in [-1,1] \; \forall_{j=1..d} \}$$
$$= \{ g \in \mathbb{R}^d \mid \|g\|_\infty \leq 1 \}$$

Thm $\}$ If $\varphi : \mathbb{R}^d \to \mathbb{R}$ is a convex function, then

$\bar{x} \in \mathbb{R}^d$ is a minimum of $\varphi \iff 0_{\mathbb{R}^d} \in \partial \varphi(\bar{x})$

Corollary : For the LASSO problem

minimize
$x \in \mathbb{R}^d$
$$\varphi(x) = \frac{1}{2m} \|Ax - b\|_2^2 + \lambda \|x\|_1,$$

$$\bar{x} \in \underset{x \in \mathbb{R}^d}{\text{argmin}} \; \varphi(x) \iff -\frac{1}{m\lambda} A^T(A\bar{x} - b) \in \partial \|\cdot\|_1 (\bar{x})$$

Proof sketch : $\varphi$ is convex as the sum of two convex functions

hence : $\bar{x} \in \underset{x \in \mathbb{R}^d}{\text{argmin}} \varphi(x) \iff 0_{\mathbb{R}^d} \in \partial \varphi(x)$

Let $f(x) = \frac{1}{2m} \|Ax - y\|_2^2$ . $f \; C^1$

$$\forall u \ldots \partial \varphi(x) = \left\{ \nabla f(x) + g \mid g \in \partial\left(\lambda \|.\|_1\right)(x) \right\}$$

$$= \left\{ \frac{1}{m} A^T(Ax-b) + \lambda g \mid g \in \partial(\|.\|_1)(x) \right\}$$

$$= \left\{ \hat{g} \in \mathbb{R}^d \mid \forall_{j=1..d}, \right.$$

$$\hat{g}_j \begin{cases} = \left[\frac{1}{m} A^T(Ax-b)\right]_j + \lambda \operatorname{sign}(x_j) & \text{if } x_j \neq 0 \\ \in \left[\left[\frac{1}{m} A^T(Ax-b)\right]_j - \lambda, \left[\frac{1}{m} A^T(Ax-b)\right]_j + \lambda\right] & \text{if } x_j = 0 \end{cases} \left. \right\}$$

$$0_{\mathbb{R}^d} \in \partial\varphi(\bar{x}) \iff 0_{\mathbb{R}^d} \in \left\{ \frac{1}{2m} A^T(A\bar{x}-b) + g \mid g \in \partial\left(\lambda \|.\|_1\right)(\bar{x}) \right\}$$

$$\iff \exists \bar{g} \in \partial\left(\lambda \|.\|_1\right)(\bar{x}) \quad \text{such that}$$

$$\frac{1}{m} A^T(A\bar{x}-b) + \bar{g} = 0_{\mathbb{R}^d}$$

$$\iff \exists \bar{g} \longrightarrow, \quad \bar{g} = -\frac{1}{m} A^T(A\bar{x}-b)$$

$$\iff \frac{-1}{m} A^T(A\bar{x}-b) \in \partial\left(\lambda \|.\|_1\right)(\bar{x})$$

Since $\partial\left(\lambda \|.\|_1\right)(\bar{x}) = \left\{ \lambda g \mid g \in \partial \|.\|_1(\bar{x}) \right\}$, the last inclusion is equivalent to

$$-\frac{1}{m\lambda} A^T(A\bar{x}-b) \in \partial \|.\|_1(\bar{x})$$

$\hookrightarrow$ the condition $-\frac{1}{m\lambda} A^T(A\bar{x}-b) \in \partial \|.\|_1(\bar{x})$

means that $\forall j = 1, \dots, d,$

$A = [a_1 \dots a_d]$

$a_j \in \mathbb{R}^m$

$$|a_j^T(A\bar{x}-b)| \leq m\lambda \quad \text{if} \quad \bar{x}_j = 0 \quad (*)$$

$$a_j^T(A\bar{x}-b) = m\lambda \, \text{sign}(\bar{x}_j) \quad \text{if} \quad \bar{x}_j \neq 0$$

As $\lambda$ increases, the condition $(*)$ is more likely to be satisfied at the solution and thus $\bar{x}$ is more likely to have zero coordinates

.NB: These conditions are not easy to solve for arbitrary $(A,b)$ but they can be for specific $A$ and $b$, and they also serve as convergence criterion for iterative methods

③ Subgradient and proximal gradient for $\ell_1$ regularized problems

minimize $f(x) + \lambda \|x\|_1$     $f$ convex    $\lambda > 0$
$x \in \mathbb{R}^d$

Subgradient method

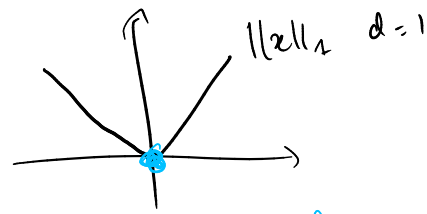    • Start with $x_0 \in \mathbb{R}^d$

    • For $k = 0, 1, 2, \dots$

       • Compute $g_k \in \partial(f + \lambda\|.\|_1)(x_k)$.

The method applies to any convex function

       • Define $x_{k+1} = x_k - \alpha_k g_k$, where $\alpha_k > 0$

$\hookrightarrow$ This method is harder to implement than a gradient-type method because: 1) The choice of subgradient matters



$\|x\|_1 \quad d=1$

$x_h = 0$

$\partial(\|\cdot\|_1)(x_h) = [-1,1]$

If we choose $g \in \partial(\|\cdot\|_1)(0), \; g \neq 0,$

$x_h - \alpha g \neq 0$

$\Rightarrow$ we move away from the minimum

One choice that works:

$$g_k \in \text{argmin} \left\{ \|g\|_2 \mid g \in \partial h(x_k) \right\}$$

where $h$ is the function to be minimized

$\frown$ Expensive, it might require to compute the entire subdifferential

2) Sensitive to the stepsize $\alpha_k > 0$

It is possible that $g_h \in \partial h(x_h)$ and yet

$\forall \alpha > 0, \quad h(x_h - \alpha g_h) > h(x_k)$

$\hookrightarrow$ Nevertheless, the subgradient method and its stochastic counterpart (stochastic subgradient method for finite-sum problems) are used in training common neural architectures based on nonsmooth activations

Ex| ReLU $(t) = \max(t, 0)$

Subgradient for $l_1$ regularized problem

$$x_{n+1} = x_n - \alpha_n g_n \qquad g_n \in \partial(f + \lambda \|.\|_1)(x_n)$$

$$\text{If } f \in C^1, \qquad g_n = \nabla f(x_n) + \lambda \bar{g}_n, \quad \bar{g}_n \in \partial \|.\|_1(x_n)$$

$$x_{n+1} = \underbrace{x_n - \alpha_n \nabla f(x_n)}_{\substack{\text{Gradient} \\ \text{descent step} \\ \text{on } f}} \underbrace{- \alpha_n \lambda \bar{g}_n}_{\substack{\text{shift by } \alpha_n \lambda \bar{g}_n \\ \bar{g}_n \in [-1,1]^d}}$$

$\rightarrow$ the iterates are different from that of GD, but it is hard to see that they are sparser than the iterates of GD

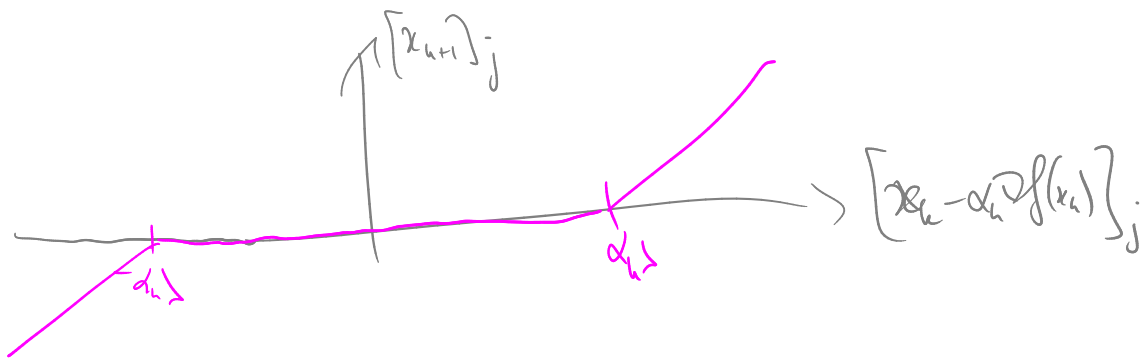Proximal gradient iteration $\quad (f \in C^1)$

(P) $\quad x_{k+1} \in \underset{x \in \mathbb{R}^d}{\text{argmin}} \left\{ f(x_n) + \nabla f(x_n)^\top (x - x_n) + \frac{1}{2\alpha_k} \|x - x_k\|_2^2 + \lambda \|x\|_1 \right\}$

Th1 The solution of the proximal subproblem (P) is unique, and defined coordinate wise by

$$\forall j = 1..d, \quad [x_{k+1}]_j = \begin{cases} [x_n - \alpha_n \nabla f(x_n)]_j - \alpha_k \lambda & \text{if } [x_n - \alpha_n \nabla f(x_n)]_j > \alpha_k \lambda \\ [x_n - \alpha_n \nabla f(x_n)]_j + \alpha_k \lambda & \text{if } [x_n - \alpha_n \nabla f(x_n)]_j < -\alpha_k \lambda \\ 0 & \text{otherwise} \end{cases}$$

This iteration sets components of $x_{k+1}$ to 0!

- $\|x_{k+1}\|_0 \leq \|x_k - \alpha_k \nabla f(x_k)\|_0 \qquad \forall k \in \mathbb{N}$

- $x_{k+1} = \text{prox}_{\alpha_k \|\cdot\|_1}\left(x_k - \alpha_k \nabla f(x_k)\right)$

Remarks :
- The proximal gradient method for $\ell_1$ regularization was discovered in compressed sensing under the name ISTA (Iterative Soft-Thresholding Algorithm)

- It has also been combined with acceleration (FiSTA) 2009