

OPTIMIZATION FOR MACHINE LEARNING

November 21, 2024

Today : → Exercise solutions
→ Sparse optimization / ℓ_1 regularization

Exercise on stochastic gradient (from lecture 11)

Setup: $\underset{x \in \mathbb{R}^d}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n f_i(x)$ f_i depends on the i^{th} sample of a dataset of size n

Batch SG

$$x_{k+1} = x_k - \frac{\alpha_k}{|S_k|} \sum_{i \in S_k} \nabla f_i(x_k)$$

S_k is a set of indices drawn randomly in $\{1, \dots, n\}$ with/without replacement

This exercise

$|S_k| = m_b \in \{1, \dots, n\} \quad \forall k.$

$$\forall S \subseteq \{1, \dots, n\}, \quad P(S_k = S) = \begin{cases} 0 & \text{if } |S| \neq m_b \\ \frac{1}{\binom{n}{m_b}} & \text{if } |S| = m_b \end{cases}$$

→ Sampling m_b indices without replacement

→ Uniform sampling over all possible subsets of $\{1, \dots, n\}$ with cardinality m_b

a) Show $\mathbb{E}_{S_k} \left[\frac{1}{|S_k|} \sum_{i \in S_k} \nabla f_i(x_k) \right] = \nabla f(x_k) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(x_k)$

$$\mathbb{E}_{S_k} \left[\frac{1}{|S_k|} \sum_{i \in S_k} \nabla f_i(x_k) \right] = \sum_{S \subseteq \{1, \dots, n\}} P(S_k = S) \times \frac{1}{|S|} \sum_{i \in S} \nabla f_i(x_k)$$

$$= \sum_{\substack{S \subseteq \{1, \dots, m\} \\ |S|=m_b}} \Pr(S_b = S) \times \frac{1}{m_b} \sum_{i \in S} Df_i(x_k)$$

$$= \sum_{\substack{S \subseteq \{1, \dots, m\} \\ |S|=m_b}} \frac{1}{\binom{m}{m_b}} \times \frac{1}{m_b} \sum_{i \in S} Df_i(x_k)$$

$$= \frac{1}{\binom{m}{m_b} \times m} \sum_{\substack{S \subseteq \{1, \dots, m\} \\ |S|=m_b}} \sum_{i \in S} Df_i(x_k)$$

Given some index $j \in \{1, \dots, m\}$, how many times does $Df_j(x_k)$ appear in the double sum?

$\binom{m-1}{m_b-1} \rightarrow$ number of possibilities
(sampling without replacement)
 $\binom{m-1}{m_b-1} \rightarrow m_b - 1$ indices left to pick

$$\sum_{\substack{S \subseteq \{1, \dots, m\} \\ |S|=m_b}} \sum_{i \in S} Df_i(x_k) = \sum_{i=1}^m \binom{m-1}{m_b-1} Df_i(x_k)$$

Hence

$$E_{S_m} \left[\frac{1}{|S_m|} \sum_{i \in S_m} Df_i(x_k) \right] = \frac{1}{\binom{m}{m_b} \times m_b} \sum_{i=1}^m \binom{m-1}{m_b-1} Df_i(x_k)$$

$$\binom{m}{m_b} = \frac{m}{m_b} \binom{m-1}{m_b-1}$$

$$\binom{m}{m_b} = \frac{m!}{m_b!(m-m_b)!}$$

$$= \frac{\binom{m-1}{m_b-1}}{\binom{m}{m_b} \times m_b} \sum_{i=1}^m Df_i(x_k)$$

$$= \binom{m-1}{m_b-1} \times \frac{1}{\frac{m}{m_b} \binom{m-1}{m_b-1} \times m_b} \sum_{i=1}^m Df_i(x_k)$$

$$= \frac{1}{m} \sum_{i=1}^m \nabla f_i(x_k) = \nabla \bar{f}(x_k)$$

b) How can we guarantee that $E[f(x_k) - f^*] \xrightarrow{k \rightarrow \infty} 0$

under the assumption that f_i is C_1^1 , strongly convex and $f^* = \min_u f(u)$?

① 1 variant of the proposed method (batch SG)

② 1 modification to the algorithm

① Use a decaying stepsize

Other option: Iterative averaging (not covered in class)
but see subgradient lecture

$$\mu_b = m$$

② Any gradient aggregation method (SAGA, SVRG)

Exercise from the previous lecture

$$(P) \underset{x \in \mathbb{R}^d}{\text{minimize}} \quad \|x\|_1 + \frac{1}{2\alpha} \|x - u\|_2^2$$

for some
 $u \in \mathbb{R}^d$
and some $\alpha > 0$.

1) Write the solution of (P) as a proximal operator calculation

$$\text{prox}_h(u) = \underset{x}{\operatorname{argmin}} \left\{ h(x) + \frac{1}{2} \|x - u\|_2^2 \right\}$$

$$\underset{x \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ \|x\|_1 + \frac{1}{2\alpha} \|x - u\|_2^2 \right\}$$

$$= \underset{x \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ \alpha \|x\|_1 + \frac{1}{2} \|x - u\|_2^2 \right\} = \text{prox}_{\alpha \| \cdot \|_1} (u)$$

2) Write down the proximal gradient iteration for problem

(P) with $f(u) = \frac{1}{2} \|x - u\|_2^2$ (data-fitting term)

$$\mathcal{R}(x) = \|x\|_1 \quad (\text{regularizer})$$

$$\lambda = \alpha$$

Using the equivalence between (P) and

$$\underset{x \in \mathbb{R}^d}{\operatorname{minimize}} \quad \frac{1}{2} \|x - u\|_2^2 + \alpha \|x\|_1 = f(x) + \lambda \mathcal{R}(x)$$

we write the proximal gradient iteration

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

$$= \underset{x \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ f(x_k) + \nabla f(x_k)^T (x - x_k) + \frac{1}{2\alpha_k} \|x - x_k\|_2^2 + \lambda \mathcal{R}(x) \right\}$$

$$= \underset{x \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ \underbrace{\frac{1}{2} \|x - u\|_2^2}_{\text{constant w.r.t. } x} + (x_k - u)^T (x - x_k) + \frac{1}{2\alpha_k} \|x - x_k\|_2^2 + \alpha_k \|x\|_1 \right\}$$

$$= \underset{x \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ (x_k - u)^T (x - x_k) + \frac{1}{2\alpha_k} \|x - x_k\|_2^2 + \alpha_k \|x\|_1 \right\}$$

NB: If $x_k = u$ and $\alpha_k = 1$, the subproblem is the original problem

① Proximal gradient and l_1 regularization

Problem: (P_1) minimize $f(x) + \lambda \|x\|_1$

where $f: \mathbb{R}^d \rightarrow \mathbb{R}$ $C_L^{1,1}$ ($f(w) = \frac{1}{n} \sum_{i=1}^n f_i(w)$)

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

The function $f + \lambda \| \cdot \|_1$ is nonsmooth (because $\| \cdot \|_1$ is) and has a composite structure smooth + nonsmooth.

In that case, the subdifferential of $f + \lambda \| \cdot \|_1$ at x

is given by $\partial(f + \lambda \| \cdot \|_1)(x) = \{ v = \nabla f(x) + g \mid g \in \partial(\lambda \| \cdot \|_1)(x) \}$

Proposition: If x° is a solution of (P_1) , then

$$0 \in \partial(f + \lambda \| \cdot \|_1)(x^\circ)$$

Counter-ex $d=1$

$$f(x) = -x^2 \quad d=1$$

$$\partial \in \partial(f+1\cdot 1)(0) \quad -x^2 + 1/x$$

$$\Leftrightarrow -\nabla f(x^*) \in \partial(\lambda \|x\|_1)(x^*)$$

Remark: The condition $-\nabla f(x^*) \in \partial(\lambda \|x\|_1)(x^*)$

("Variational inequality") cannot be solved in general, even for simple f s

$$\text{Ex: } f(x) = \frac{1}{2} \|Ax-y\|^2$$

- But within a proximal subproblem, we can actually solve the optimality condition.

↪ Suppose we apply the proximal gradient method to (P_1) .
At iteration k , we need to solve

$$\underset{x \in \mathbb{R}^d}{\text{minimize}} \left\{ f(x_k) + \nabla f(x_k)^T(x-x_k) + \frac{1}{2\alpha_k} \|x-x_k\|^2 + \lambda \|x\|_1 \right\}$$

$$\ell_k(x)$$

ℓ_k convex, even strongly convex and quadratic (and the quadratic term is simple), $C_{\lambda/\alpha_k}^{1/2}$

$$x \in \mathbb{R}^d,$$

$$\nabla \ell_k(x) = \nabla f(x_k) + \frac{1}{\alpha_k}(x-x_k)$$

For that problem, since $\ell_a + \lambda \|x\|_1$ is convex,

$$\begin{aligned} x^* \text{ is solution of } (P_1) &\Leftrightarrow 0 \in \partial(\ell_a + \lambda \|x\|_1)(x^*) \\ &\Leftrightarrow -\nabla \ell_a(x^*) \in \partial(\lambda \|x\|_1)(x^*) \\ &\Leftrightarrow \left[-\nabla \ell_a(x^*) - \frac{1}{\lambda} (x^* - x_a) \right] \in \partial(\lambda \|x\|_1)(x^*) \end{aligned}$$

$$\overset{d=1}{\overbrace{\quad}} \quad t \mapsto |t| \quad \quad \partial(|\cdot|)(t) = \begin{cases} 1 & \text{if } t > 0 \\ -1 & \text{if } t < 0 \\ [-1, 1] & \text{if } t = 0 \end{cases}$$

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

$$\forall x \in \mathbb{R}^d, \quad \partial(\|x\|_1)(x) = \left\{ g = \begin{bmatrix} g_1 \\ \vdots \\ g_d \end{bmatrix} \in \mathbb{R}^d \mid g_j \in \partial(|\cdot|)(x_j) \quad \forall j = 1..d \right\}$$

$$\overset{d=2}{\overbrace{\quad}} \quad \partial(\|x\|_1)(x) = \left\{ \begin{array}{l} \begin{bmatrix} -1, 1 \end{bmatrix}^2 \text{ if } x = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \\ \left\{ \begin{bmatrix} 1 \\ g_2 \end{bmatrix} \mid g_2 \in [-1, 1] \right\} \text{ if } x = \begin{bmatrix} x_1 \\ 0 \end{bmatrix} \text{ with } x_1 > 0 \\ \left\{ \begin{bmatrix} -1 \\ g_2 \end{bmatrix} \mid g_2 \in [-1, 1] \right\} \text{ if } x = \begin{bmatrix} x_1 \\ 0 \end{bmatrix} \text{ with } x_1 < 0 \\ \vdots \end{array} \right.$$

$$-\nabla \varphi_h(x^*) \in \partial(\| \cdot \|_h)(x^*)$$

$$\Leftrightarrow -\frac{1}{\lambda} \nabla \varphi_h(x^*) \in \partial(\| \cdot \|_h)(x^*)$$

$$h(x) \geq h(y) + g^T(y-x) \quad g \in \partial h(y)$$

$$\frac{1}{\lambda} h(x) \geq \frac{1}{\lambda} h(y) + \left(\frac{g}{\lambda}\right)^T(y-x) \quad \frac{g}{\lambda} \in \partial\left(\frac{1}{\lambda} h(y)\right)$$

$$-\frac{1}{\lambda} \nabla \varphi_h(x^*) \in \partial(\| \cdot \|_h)(x^*)$$

$$\Leftrightarrow \forall j=1..d, \quad \left[-\frac{1}{\lambda} \nabla \varphi_h(x^*) \right]_j \in \partial(|.)|(x_j^*)$$

$$\Rightarrow \text{If } x_j^* = 0, \quad \left[-\frac{1}{\lambda} \nabla \varphi_h(x^*) \right]_j \in [-1, 1]$$

$$\text{If } x_j^* > 0, \quad \left[-\frac{1}{\lambda} \nabla \varphi_h(x^*) \right]_j = 1$$

$$\text{If } x_j^* < 0, \quad \left[-\frac{1}{\lambda} \nabla \varphi_h(x^*) \right]_j = -1$$

$$-\frac{1}{\lambda} \nabla \varphi_h(x^*) = -\frac{1}{\lambda} \nabla f(x_h) - \frac{1}{\alpha_h \lambda} (x^* - x_h)$$

$$= \frac{1}{\alpha_h \lambda} (x_h - \alpha_h \nabla f(x_h) - x^*)$$

$$(i) \quad \frac{1}{\alpha_n} [x_n - \alpha_n \nabla f(x_n) - x^*]_j \in [-1, 1] \text{ if } x_j^* = 0$$

$$(ii) \quad \frac{1}{\alpha_n} [x_n - \alpha_n \nabla f(x_n) - x^*]_j = 1 \quad \text{if } x_j^* > 0$$

$$(iii) \quad \frac{1}{\alpha_n} [x_n - \alpha_n \nabla f(x_n) - x^*]_j = -1 \quad \text{if } x_j^* < 0$$

3 cases

$$(i) \quad x_j^* = 0 \quad \text{and} \quad \frac{1}{\alpha_n} [x_n - \alpha_n \nabla f(x_n)]_j \in [-1, 1]$$

$$(ii) \quad x_j^* > 0 \quad \text{and} \quad \frac{1}{\alpha_n} [x_n - \alpha_n \nabla f(x_n)]_j > 1$$

$$(iii) \quad x_j^* < 0 \quad \text{and} \quad \frac{1}{\alpha_n} [x_n - \alpha_n \nabla f(x_n)]_j < -1$$

Theorem: The solution of

$$\underset{x \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ f(x_n) + \nabla f(x_n)^T (x - x_n) + \frac{1}{2\alpha_n} \|x - x_n\|_2^2 + \lambda \|x\|_1 \right\}$$

is given by $x^* \in \mathbb{R}^d$ such that

$$\forall j=1..d, \quad [x^*]_j = \begin{cases} 0 & \text{if } [x_n - \alpha_n \nabla f(x_n)]_j \in [-\alpha_n, \alpha_n] \\ [x_n - \alpha_n \nabla f(x_n)]_j - \alpha_n & \text{if } [x_n - \alpha_n \nabla f(x_n)]_j > \alpha_n \\ [x_n - \alpha_n \nabla f(x_n)]_j + \alpha_n & \text{if } [x_n - \alpha_n \nabla f(x_n)]_j < -\alpha_n \end{cases}$$

→ Proof: Use subgradient theory or

simplify the formula for $\text{prox}_{\lambda \alpha_k \| \cdot \|_1}(\cdot)$

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

and $\left[\text{prox}_{\lambda \| \cdot \|_1}(x) \right]_j = \begin{cases} x_j - \lambda & \text{if } x_j > \lambda \\ x_j + \lambda & \text{if } x_j < -\lambda \\ 0 & \text{if } x_j \in [-\lambda, \lambda] \end{cases}$

Note:

$$[x_{k+1}]_j = 0 \quad \text{if } [x_k - \alpha_k \nabla f(x_k)]_j \notin [-\alpha_k \lambda, \alpha_k \lambda]$$

$\forall i, \forall j$

$\lambda \rightarrow \infty$, the condition becomes more and more likely

This update guarantees that $\|x_{k+1}\|_0 \leq \|x_k - \alpha_k \nabla f(x_k)\|_0$

A proximal gradient iteration always produces an iterate with less nonzero coordinates (or the same number) than the iterate of a gradient iteration.

ISTA: Iterative Soft Thresholding Algorithm

→ Name for proximal gradient with ℓ_1 regularization
in the signal processing community

→ Comes with theoretical convergence rates

→ Big advantage: Explicit formula for the "prox" thanks to the soft-thresholding operator

For convex problems, there is an accelerated version of ISTA called Fast ISTA, or FISTA (Beck & Teboulle, 2009)

Iterations :
$$\begin{cases} p_{k+1} = x_k + \beta_{k+1} (x_k - x_{k-1}) & (\text{Momentum step}) \\ x_{k+1} = \text{prox}_{\frac{1}{2\alpha_{k+1}} \| \cdot \|_1} (p_{k+1} - \alpha_{k+1} \nabla f(p_{k+1})) & (\text{Proximal gradient steps}) \end{cases}$$

(2) Use cases of proximal gradient for ℓ_1

* LASSO (Least Absolute Self. Shrinkage Operator)

Tibshirani 1996

Linear regression + ℓ_1 regularization

$$(P_{\text{LASSO}}) \underset{x \in \mathbb{R}^d}{\text{minimize}} \quad \frac{1}{2} \|Ax - y\|^2 + \lambda \|x\|_1 \quad \begin{array}{l} \text{Also called} \\ \text{basis pursuit} \end{array}$$

$A \in \mathbb{R}^{m \times d}, y \in \mathbb{R}^m, \lambda > 0$

Goal of this formulation: Obtain a solution that has more zeros than that of the un-regularized problem while corresponding to a value of $\frac{1}{2} \|Ax - y\|^2$ that

$\hat{\beta}$ is good enough

△ (PLASSO) has no closed-form solution in general

\Rightarrow In practice, this problem is solved approximately using proximal gradient

* Low-rank matrix approximation

↳ If the variables form a matrix $X \in \mathbb{R}^{d_1 \times d_2}$, one can be interested in finding sparse matrices, in which case one considers problems of the form

$$\underset{\substack{X \in \mathbb{R}^{d_1 \times d_2} \\ b \in \mathbb{R}^m \\ A(X) \in \mathbb{R}^m}}{\text{minimize}} \quad \frac{1}{2} \|A(X) - b\|^2 + \lambda \|X\|_1$$

$\|X\|_1 = \sum_{i=1}^{d_1} \sum_{j=1}^{d_2} |X_{ij}|$

e.g. $A(X) = [\text{trace}(A_i^T X)]_{i=1 \dots m}$

↳ Another form of sparsity for matrices is low-rank structure

rank (matrix) = number of non-zero singular values

SVD: $X \in \mathbb{R}^{d_1 \times d_2}, \quad X = U \begin{bmatrix} \sigma_1 & & \\ 0 & \ddots & \\ & 0 & \sigma_{\min(d_1, d_2)} \end{bmatrix} V^T$

$$U^T U = U U^T = \text{Id}_{d_1}$$

$$V^T V = V V^T = \text{Id}_{d_2}$$

$$\sigma_1 \geq \dots \geq \sigma_{\min(d_1, d_2)} \geq 0$$

rank: # of non-zero σ_i 's

$$\text{rank}(X) = \|\sigma(X)\|_0$$

$\sigma(X)$: vector
of singular
values of X

Low-rank regularization

$$\text{Use } \underbrace{\|X\|_*}_{\text{"nuclear norm."}} = \|\sigma(X)\|_1$$

$$\begin{array}{ll} \text{minimize} & \frac{1}{2} \|A(X) - b\|^2 + \lambda \|X\|_* \\ X \in \mathbb{R}^{d_1 \times d_2} & \end{array}$$

\Rightarrow Produce a solution with rank lower than (or equal to) that of the solution of

$$\begin{array}{ll} \text{minimize} & \frac{1}{2} \|A(X) - b\|^2 \\ X \in \mathbb{R}^{d_1 \times d_2} & \end{array}$$

\Rightarrow Proximal gradient applies to that setting!

* Structured sparsity regularizers (ℓ_1 aka LASSO regularization)

↳ Group LASSO / Group ℓ_1

$$x \in \mathbb{R}^d, \quad \mathcal{L}(x) = \sum_{g \in G} \|x_g\|_2 \quad \left. \right] \Rightarrow \text{forces sparsity of groups of variables}$$

G is a partition of $\{1, \dots, d\}$ into groups of variables

$$\text{Ex) } x = \begin{bmatrix} x_{g_1} \\ x_{g_2} \\ x_{g_3} \end{bmatrix} \begin{array}{l} \uparrow d/3 \\ \uparrow d/2 \\ \uparrow d/6 \end{array} \quad x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{bmatrix} \quad \begin{array}{l} g_1 = \{1, 4\} \\ g_2 = \{2, 3\} \\ g_3 = \{5\} \end{array}$$

\rightarrow Group L1 norm : the sum of $\begin{bmatrix} \|x_{g_1}\|_2 \\ \vdots \\ \|x_{g_m}\|_2 \end{bmatrix}$ $|G|=m$

$$G = \{\{1\}, \{2\}, \dots, \{d\}\} : \sum_{g \in G} \|x_g\|_2 = \sum_{j=1}^d \|x_j\|_2 = \|x\|_1$$

$$G = \{\{1, -, d\}\} : \sum_{g \in G} \|x_g\|_2 = \|x\|_2 = \sqrt{\sum_{j=1}^d \|x_j\|^2}$$

\hookrightarrow Beyond that: ℓ_1/ℓ_∞ ℓ_1/ℓ_p $p \in \{1, 2, \infty\}$

$\bullet \ell_p/\ell_q$ $\underset{\text{neighborhood}}{\sum_{g \in G}} \|x_g\|_\infty, \sum_{g \in G} \|x_g\|_p$ ℓ_p norm

$$\sum_{g \in G} w_g \|x_g\|_2 \quad (\text{weighted group L1 norm})$$

\bullet overlapping groups : problem dependent

\Rightarrow All these regularizations define sparsity-inducing regularizers and different "balls"

$$\{x \in \mathbb{R}^d \mid \mathcal{D}(x) \leq 1\}$$

$$\text{NB: } \mathcal{R}(x) = \frac{\lambda}{2} \|x\|_2^2 + \mu \|x\|_1 \quad \begin{matrix} \lambda > 0 \\ \mu > 0 \end{matrix}$$

"Elastic net regularization"

→ Can apply proximal gradient
($\text{prox}_{\mathcal{R}}(\cdot)$ has a closed form)

Remark: In deep learning, dropout is a way to perform sparse updates of the parameters of a neural network during training ⇒ one form of sparse regularization

↳ Connection to coordinate descent (see next lecture)

Exercise

$$\underset{x \in \mathbb{R}^d}{\text{minimize}} \quad f(x) + \frac{\lambda}{2} \|Lx\|_2^2 \quad f(x)$$

$$L = \begin{bmatrix} -2 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 1 & -2 \end{bmatrix}$$

$$L \in \mathbb{R}^{d \times d}, \quad L_{ij} = \begin{cases} 1 & \text{if } j=i+1 \text{ or } j=i-1 \\ -2 & \text{if } i=j \\ 0 & \text{otherwise} \end{cases}$$

a) Write down the proximal gradient iteration for that problem. Is the prox operator easy to compute?

b) Application:

$$f(x) = \frac{1}{2} \|Ax - y\|^2$$

What's the cost of solving the subproblem?

- c) If $A = Id$, show that the proximal subproblem solution x_{n+1} is such that $[x_{n+1}]_j$ only depends on a subset of coordinates of $[x_n]_j$.

(NB : $\max_{\frac{1}{2}\|L\cdot x\|^2} (\cdot) = ?$)