

# OPTIMIZATION FOR MACHINE LEARNING

November 18, 2024

Today: Regularized and sparse optimization (Part 1)

# 1) Regularization

→ The basic optimization formulation of a machine learning task does not encode the full objective of that task

(Ex)  $\{(a_i, y_i)\}_{i=1 \dots n}$  data  $\longrightarrow$  minimize  $\frac{1}{n} \sum_{i=1}^n f_i(x)$   $x \in \mathbb{R}^d$

"Empirical Risk Minimization"

Generalization  $\left\{ \begin{array}{l} \cdot \text{What happens to the solution of the problem if the data is slightly perturbed?} \\ \cdot \text{Can the solution get good results on new data?} \end{array} \right.$

Feature/Weight selection  $\left\{ \begin{array}{l} \cdot \text{Is the solution interpretable? Is the resulting model simple?} \end{array} \right.$

$\Rightarrow$  How can these questions be encoded in the optimization formulation?  
A: Regularization

Def: A regularized optimization problem has the form

minimize  $x \in \mathbb{R}^d$   $f(x) + \lambda \Omega(x)$

- $f(x)$  (green highlight): Data-fitting term  
 e.g.  $f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$
- $\lambda$  (yellow highlight): Regularization parameter (represents a trade-off between data fitting and regularization)
- $\Omega(x)$  (blue highlight): Regularization term: represents properties that we would like the solution to have  
 In general, does not depend on the data

•  $\lambda \rightarrow 0$  : The problem eventually becomes equivalent to  
minimize  $f(x)$   
 $x \in \mathbb{R}^d$

•  $\lambda \rightarrow \infty$  : The problem eventually becomes equivalent to  
minimize  $\Omega(x)$   
 $x \in \mathbb{R}^d$

Ex) •  $\Omega(x) = \frac{1}{2} \|x\|^2 = \frac{1}{2} \|x\|_2^2 = \frac{1}{2} \sum_{i=1}^d x_i^2$

$l_2$  regularization / ridge regularization / Tychonov regularization

$\Rightarrow$  Used to reduce the sensitivity of the solution to perturbations of the data

$\Rightarrow$  minimize  $\frac{1}{2} \|x\|_2^2$  ( $\lambda \rightarrow \infty$ )  $\Rightarrow$  solution is  $x^* = 0_{\mathbb{R}^d}$

minimize  $f(x) + \frac{\lambda}{2} \|x\|_2^2$ ,  $\lambda > 0$   
 $x \in \mathbb{R}^d$

"convexifying effect"

[ • Always strongly convex when  $f$  is convex, and in that case it has one unique minimum (whereas the original problem may have several)  
 $\Rightarrow$  Can be true even for nonconvex  $f$  provided  $\lambda$  is large enough

$$[x] \quad \text{minimize}_{x \in \mathbb{R}^d} \frac{1}{2n} \|Ax - y\|_2^2$$

$$d=2, \quad n=2, \quad A = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$\Rightarrow$  that problem is convex with infinitely many solutions

$$\begin{bmatrix} 1 \\ 0 \end{bmatrix} \text{ is typically the solution that we choose} \longrightarrow \underset{x \in \mathbb{R}^2}{\text{argmin}} \frac{1}{4} \|Ax - y\|^2 = \left\{ \begin{bmatrix} 1 \\ t \end{bmatrix} \mid t \in \mathbb{R} \right\}$$

Q] what happens if  $A \rightarrow A_\varepsilon = \begin{bmatrix} 1+\varepsilon & 0 \\ 0 & \varepsilon \end{bmatrix}$  for some  $\varepsilon > 0$

$$\underset{x \in \mathbb{R}^2}{\text{argmin}} \left\{ \frac{1}{4} \|A_\varepsilon x - y\|^2 \right\} = \left\{ \begin{bmatrix} \frac{1}{1+\varepsilon} \\ \frac{1}{\varepsilon} \end{bmatrix} \right\}$$

$$\left\| \begin{bmatrix} \frac{1}{1+\varepsilon} \\ \frac{1}{\varepsilon} \end{bmatrix} \right\|^2 = \frac{1}{(1+\varepsilon)^2} + \frac{1}{\varepsilon^2} \xrightarrow{\varepsilon \rightarrow 0} \infty$$

$\rightarrow$  the solution (set) of  $\underset{x \in \mathbb{R}^2}{\text{minimize}} \frac{1}{4} \|Ax - y\|^2$  is

sensitive to perturbations of (the data)  $A$ .

$\rightarrow$  Consider now

$$A = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\underset{x \in \mathbb{R}^2}{\text{minimize}} \frac{1}{4} \|Ax - y\|^2 + \frac{\lambda}{2} \|x\|^2 \quad \lambda > 0$$

$$\rightarrow \underset{x \in \mathbb{R}^2}{\text{argmin}} \left\{ \underbrace{\frac{1}{4} \|Ax - y\|^2 + \frac{\lambda}{2} \|x\|^2}_{f_\lambda(x)} \right\} = \left\{ \begin{bmatrix} \frac{2}{2+\lambda} \\ 0 \end{bmatrix} \right\}$$

$$f_\lambda(x) = \frac{1}{4} \left\| \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right\|^2 + \frac{\lambda}{2} (x_1^2 + x_2^2)$$

$$= \frac{1}{4} (x_1 - 1)^2 + \frac{1}{4} + \frac{\lambda}{2} x_1^2 + \frac{\lambda}{2} x_2^2$$

$$\nabla f_\lambda(x) = 0 \quad (\Leftrightarrow) \quad \begin{cases} \frac{1}{2}(x_1 - 1) + \lambda x_1 = 0 \\ \lambda x_2 = 0 \end{cases}$$

$$(\Leftrightarrow) \quad \begin{cases} x_1 = \frac{1}{1+2\lambda} \\ x_2 = 0 \end{cases}$$

④ Unique solution (that is not a solution of the original problem minimize  $\frac{1}{4} \|Ax - y\|^2$  over  $x \in \mathbb{R}^2$ )

→ As  $\lambda \rightarrow \infty$ , the solution converges to  $\begin{bmatrix} 0 \\ 0 \end{bmatrix}$

$$\operatorname{argmin}_{x \in \mathbb{R}^2} \frac{1}{2} \|x\|_2^2$$

→ As  $\lambda \rightarrow 0$ , the solution converges to

$$\begin{bmatrix} 1 \\ 0 \end{bmatrix} \in \operatorname{argmin}_{x \in \mathbb{R}^2} \left\{ \frac{1}{4} \|Ax - y\|^2 \right\}$$

↳ If we consider the perturbed problem

$$\operatorname{minimize}_{x \in \mathbb{R}^2} \underbrace{\frac{1}{4} \|A_\varepsilon x - y\|^2 + \frac{\lambda}{2} \|x\|^2}_{f_{\lambda, \varepsilon}(x)}$$

$$\nabla f_{\lambda, \varepsilon}(x) = \frac{1}{2} A_{\varepsilon}^T (A_{\varepsilon} x - y) + \lambda x$$

$$= \frac{1}{2} \begin{bmatrix} 1+\varepsilon & 0 \\ 0 & \varepsilon \end{bmatrix} \begin{bmatrix} (1+\varepsilon)x_1 - 1 \\ \varepsilon x_2 - 1 \end{bmatrix} + \begin{bmatrix} \lambda x_1 \\ \lambda x_2 \end{bmatrix}$$

$$= \begin{bmatrix} \frac{(1+\varepsilon)^2}{2} x_1 - \frac{1+\varepsilon}{2} + \lambda x_1 \\ \frac{\varepsilon^2}{2} x_2 - \frac{\varepsilon}{2} + \lambda x_2 \end{bmatrix}$$

$$\nabla f_{\lambda, \varepsilon}(x) = 0 \Leftrightarrow \begin{cases} x_1 = \frac{1+\varepsilon}{(1+\varepsilon)^2 + 2\lambda} \\ x_2 = \frac{\varepsilon}{\varepsilon^2 + 2\lambda} \end{cases}$$

$$\left( \lambda = 0 \quad \begin{bmatrix} \frac{1}{1+\varepsilon} \\ \frac{1}{\varepsilon} \end{bmatrix} \right)$$

$$\left\| \begin{bmatrix} \frac{1+\varepsilon}{(1+\varepsilon)^2 + 2\lambda} \\ \frac{\varepsilon}{\varepsilon^2 + 2\lambda} \end{bmatrix} \right\|^2 = \frac{(1+\varepsilon)^2}{((1+\varepsilon)^2 + 2\lambda)^2} + \frac{\varepsilon^2}{(\varepsilon^2 + 2\lambda)^2}$$

$= O\left(\frac{1}{\lambda^2}\right)$  for large values of  $\lambda$

$$\left( \lambda \rightarrow \infty, \begin{bmatrix} \cdot \\ \cdot \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 0 \end{bmatrix} \right)$$

and/or small values of  $\varepsilon$

## Ex) Sparsity-inducing regularization terms

Goal: Want a solution that is sparse, i.e. it has a significant number of zero coordinates.

⚠ Sparsity is a dimension dependent notion

$$\begin{bmatrix} 1 \\ 0 \end{bmatrix} \in \mathbb{R}^2$$

$$\begin{bmatrix} * \\ * \\ * \\ 0 \\ ! \end{bmatrix} \begin{array}{l} \uparrow 1\% \text{ nonzero coeff. vals} \\ \downarrow 55\% \text{ zero} \end{array}$$

$$d = 10^9$$

→ For linear models  $a \mapsto a^T x$  parameterized by  $x$ , finding a sparse linear model is a way to operate feature selection, i.e. to identify the most important features.

→ For overparameterized models ( $d$  parameters,  $n$  data points,  $d \gg n$ ),

there can exist infinitely models, some sparse and others dense (dense = all coefficients are nonzero)

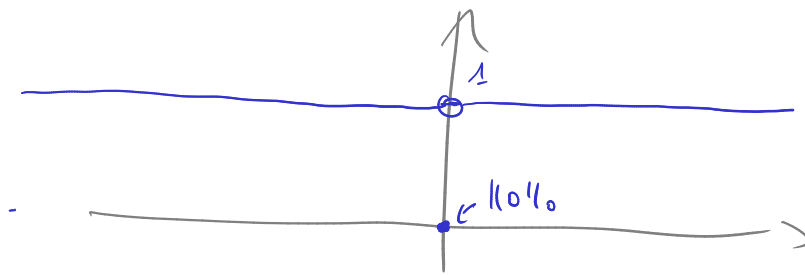
## Regularizers that promote sparse solutions

•  $\ell_0$  regularization

$$\text{minimize}_{x \in \mathbb{R}^d} f(x) + \lambda \|x\|_0 \quad \text{where } \|x\|_0 = \sum_{j=1}^d \mathbb{1}(x_j \neq 0)$$

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

- $\|\cdot\|_0$  = "0 norm" (but not a norm!)
- $\|x\|_0 \in \{0, 1, \dots, d\}$  for  $x \in \mathbb{R}^d$
- $\|x\|_0$  is equal to the number of nonzero coordinates in  $x$



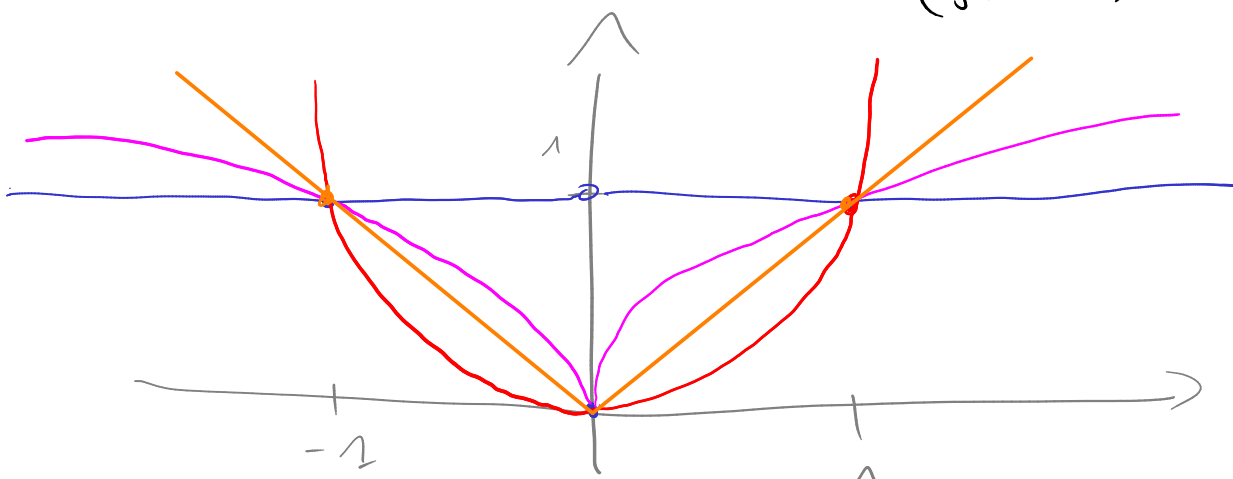
$$d=1$$

$$\|x\|_0 = \begin{cases} 1 & \text{if } x \neq 0 \\ 0 & \text{if } x = 0 \end{cases}$$

(-)  $\|\cdot\|_0$  is nonconvex, discontinuous so difficult to optimize!

→ Because the 0 norm is difficult to use in optimization, other regularization terms based on "p norms" have been proposed

$$\forall p \in [0, +\infty], \quad \|x\|_p = \begin{cases} \sum_{j=1}^d \mathbb{1}(x_j \neq 0) & \text{if } p=0 \\ \max_{1 \leq j \leq d} |x_j| & \text{if } p=\infty \\ \left( \sum_{j=1}^d |x_j|^p \right)^{1/p} & \text{if } 0 < p < \infty \end{cases}$$



- $p=2$   
 $\|x\|_{1/2}$
- $p=1$
- $p=0$
- $p=1/2$   
 $\|x\|_{1/2}$



→  $l_0$  "norm" is the limit case of  $l_p$  "norms" when  $p \rightarrow 0$

⊖ For  $p \in (0, 1)$ ,  $x \mapsto \|x\|_p^p$  is not convex  
not a norm

⊖ For  $p > 1$ ,  $x \mapsto \|x\|_p^p$  is convex, easy

to optimize because  $\|x\|_p$  is a norm  
But it is not the closest convex function  
to the  $l_0$  norm

→  $p = 1$  gives the closest convex function to the  $l_0$   
norm, namely the  $l_1$  norm:

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

Convex, continuous, norm

Among all convex functions  $\mathcal{Q}$  that satisfy

$\mathcal{Q}(x) \leq \|x\|_0 \quad \forall x \in [-1, 1]^d$  the  $l_1$  norm  
is the one that minimizes  $\int_{[-1, 1]^d} |\mathcal{Q}(x) - \|x\|_0|$   
over  $\mathcal{Q}$

⇒  $l_1$  regularization is the dominant approach  
for enforcing sparsity.

## (2) Proximal methods

Goal: Solve problems of the form

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

$$\begin{aligned} f: \mathbb{R}^d &\rightarrow \mathbb{R} \\ \Omega: \mathbb{R}^d &\rightarrow \mathbb{R} \end{aligned}$$

↳ The methods that solve this problem while exploiting its structure are based on the concept of proximal operator

Def. The proximal operator of a function  $h: \mathbb{R}^d \rightarrow \mathbb{R}$ , denoted by  $\text{prox}_h(\cdot)$ , is a mapping defined by

$$\forall x \in \mathbb{R}^d, \quad \text{prox}_h(x) = \underset{u \in \mathbb{R}^d}{\text{argmin}} \left\{ h(u) + \frac{1}{2} \|u - x\|^2 \right\}$$

↳ Set-valued mapping in general:  $\text{prox}_h(x)$  subset of  $\mathbb{R}^d$

(Example)  $h(x) = \|x\|_0$

$$\text{prox}_h(x) = \left\{ v \in \mathbb{R}^d \mid \forall j = 1 \dots d, \right. \\ \left. [v]_j = \begin{cases} 0 & \text{if } |x_j| < \sqrt{2\lambda} \\ x_j & \text{if } |x_j| > \sqrt{2\lambda} \\ \in \{0, x_j\} & \text{if } |x_j| = \sqrt{2\lambda} \end{cases} \right\}$$

↳ If  $h$  is convex, the <sup>proximal</sup> prox operator of  $h$  can be

defined as a mapping from  $\mathbb{R}^d \rightarrow \mathbb{R}^d$

- $h$  convex  $\Rightarrow h + \frac{1}{2} \|\cdot - x\|^2$  is strongly convex
- $\Rightarrow \operatorname{argmin}_u \left\{ h(u) + \frac{1}{2} \|u - x\|^2 \right\}$  is a singleton

### Examples:

- $h : x \mapsto 0$ ,  $\operatorname{prox}_0(x) = x$  ( $\operatorname{argmin}_u \frac{1}{2} \|u - x\|^2 = \{x\}$ )

- $h : x \mapsto \frac{\lambda}{2} \|x\|_2^2$ ,  $\operatorname{prox}_h(x) = \frac{x}{1+\lambda}$

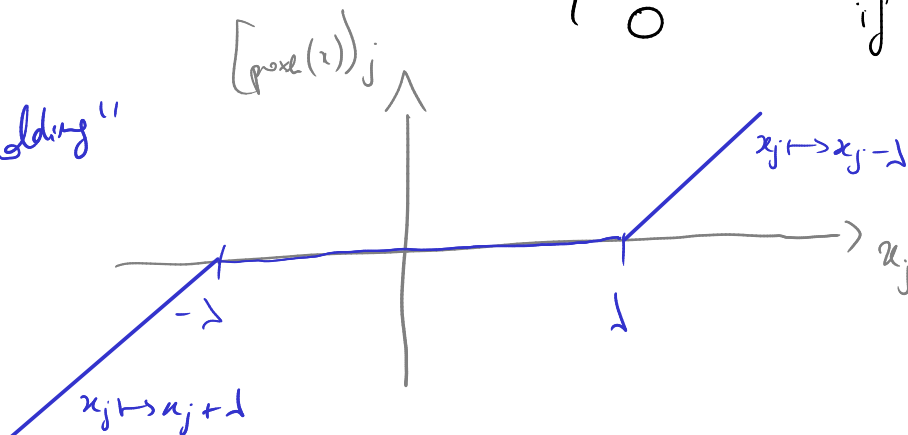
( $\lambda=1$ ,  $\operatorname{argmin}_u \frac{1}{2} \|u - 0\|_2^2 + \frac{1}{2} \|u - x\|^2 = \left\{ \frac{x}{2} \right\}$ )

- $h : x \mapsto \lambda \|x\|_1 = \lambda \sum_{j=1}^d |x_j|$

$\operatorname{prox}_h(x)$  is defined coordinate-wise by

$$\forall j=1 \dots d, \quad \left[ \operatorname{prox}_h(x) \right]_j = \begin{cases} x_j - \lambda & \text{if } x_j > \lambda \\ x_j + \lambda & \text{if } x_j < -\lambda \\ 0 & \text{if } x_j \in [-\lambda, \lambda] \end{cases}$$

"soft  $\lambda$ -thresholding"  
operator



→ Prox operators are useful when they can be computed easily (often means that they are computed on a convex function and that the optimization problem defining the proximal operator can be solved explicitly)

→ Historically, the first method based on proximal operations was not assuming that the prox was easy to compute

This first method: Proximal point method

Goal: minimize  $h(x)$  ,  $h: \mathbb{R}^d \rightarrow \mathbb{R}$  convex  
 $x \in \mathbb{R}^d$

Proximal point iteration ( $x_0 \in \mathbb{R}^d$ )

$$x_{k+1} = \text{prox}_{\alpha_k h}(x_k) , \alpha_k > 0$$

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

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

Property: For any iteration  $k$ ,

$$h(x_{k+1}) \leq h(x_k) - \frac{1}{2\alpha_k} \|x_{k+1} - x_k\|_2^2 < h(x_k) \text{ if } x_{k+1} \neq x_k$$

$$x_{k+1} = \operatorname{argmin}_x \left\{ h(x) + \frac{1}{2\alpha_k} \|x - x_k\|_2^2 \right\}$$

$$h(x_{k+1}) + \frac{1}{2\alpha_k} \|x_{k+1} - x_k\|_2^2 \leq h(x_k) + \frac{1}{2\alpha_k} \|x_k - x_k\|_2^2 = h(x_k)$$

CV rate

If  $\alpha_k = \alpha > 0 \forall k \in \mathbb{N}$ , then after  $K$  iterations,

$$h(x_K) - h(x^*) \leq \frac{\|x_0 - x^*\|^2}{2\alpha K}$$

where  $x^* \in \operatorname{argmin}_{x \in \mathbb{R}^d} h(x)$

Drawback

Each iteration to compute a proximal operator!

\* When  $\alpha_k$  is small, cost is similar to solving

$$\operatorname{minimize}_{x \in \mathbb{R}^d} \frac{1}{2} \|x - x_k\|^2 \quad \text{Easy!}$$

\* when  $\alpha_k$  is large, cost is similar to solving

$$\operatorname{minimize}_{x \in \mathbb{R}^d} h(x) \quad \text{Original problem!}$$

→ Iterators of the proximal point method can be as expensive as solving the original problem!

Special case:  $h$   $C^1$  and convex

In that case,

$$x_{k+1} = \text{prox}_{\alpha_k h}(x_k) \quad (\Leftrightarrow) \quad x_{k+1} = x_k - \alpha_k \nabla h(x_{k+1})$$

† with  
 gradient descent  
 Gradient evaluated at  
 $x_{k+1}$

$$\text{prox}_{\alpha_k h}(x_k) = \underset{x \in \mathbb{R}^d}{\text{argmin}} \underbrace{\left\{ h(x) + \frac{1}{2\alpha_k} \|x - x_k\|_2^2 \right\}}_{\phi(x)}$$

$\phi$  is  $C^1$  and strongly convex.

$$x_{k+1} = \underset{x \in \mathbb{R}^d}{\text{argmin}} \phi(x) \quad (\Leftrightarrow) \quad \nabla \phi(x_{k+1}) = 0$$

$$\Leftrightarrow \nabla h(x_{k+1}) + \frac{1}{\alpha_k} (x_{k+1} - x_k) = 0$$

$$\Leftrightarrow x_{k+1} = x_k - \alpha_k \nabla h(x_{k+1})$$

↳ Implicit definition of  $x_{k+1}$ , but can still be computed in special cases

Example:  $h(x) = \frac{1}{2m} \|Ax - y\|_2^2$       $A \in \mathbb{R}^{m \times d}$ ,  $b \in \mathbb{R}^m$

$$\nabla h(x) = \frac{1}{m} A^T (Ax - y)$$

$$x_{k+1} = x_k - \alpha_k \nabla h(x_{k+1}) = x_k - \frac{\alpha_k}{m} A^T A x_{k+1} + \frac{\alpha_k}{m} A^T y$$

$$\Leftrightarrow \left[ I + \frac{\alpha_k}{m} A^T A \right] x_{k+1} = x_k + \frac{\alpha_k}{m} A^T y$$

$$\Leftrightarrow x_{k+1} = \left[ I + \frac{\alpha_k}{n} A^T A \right]^{-1} \left( x_k + \frac{\alpha_k}{n} A^T y \right)$$

## Proximal gradient method

↳ Dedicated to problems of the form

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

$$f: \mathbb{R}^d \rightarrow \mathbb{R} \quad C^1$$

$$\Omega: \mathbb{R}^d \rightarrow \mathbb{R} \quad \text{convex}$$

$$\lambda > 0$$

Idea:

- Optimizing  $f$  alone could be done using gradient descent
- Optimizing  $\Omega$  alone ————— the proximal point method

## Proximal gradient iteration

$$x_{k+1} = \text{prox}_{\frac{\alpha_k}{\lambda} \Omega} \left( x_k - \alpha_k \nabla f(x_k) \right) \quad \alpha_k > 0$$

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

$$\text{prox}_{\frac{\alpha_k}{\lambda} \Omega} \left( x_k - \alpha_k \nabla f(x_k) \right) = \underset{x \in \mathbb{R}^d}{\text{argmin}} \left\{ \frac{\alpha_k}{\lambda} \lambda \Omega(x) + \frac{1}{2} \|x - (x_k - \alpha_k \nabla f(x_k))\|^2 \right\}$$

$$\stackrel{\times \frac{1}{\alpha_k}}{\hookrightarrow} = \underset{x \in \mathbb{R}^d}{\text{argmin}} \left\{ \frac{1}{2\alpha_k} \|x - x_k + \alpha_k \nabla f(x_k)\|^2 + \lambda \Omega(x) \right\}$$

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

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

constant  
w.r.t.  $x$

$$= \operatorname{argmin}_{x \in \mathbb{R}^d} \left\{ \underbrace{f(x_n)}_{\text{another constant w.r.t. } x} + \nabla f(x_n)^T (x - x_n) + \frac{1}{2\alpha_n} \|x - x_n\|^2 + \lambda \Omega(x) \right\}$$

$f(x_n) + \nabla f(x_n)^T (x - x_n)$  : approximation of  $f(x)$  around  $x_n$

$\frac{1}{2\alpha_n} \|x - x_n\|^2$  : proximal term, penalizes  $x$ 's that are far from  $x_n$

$\lambda \Omega(x)$  : regularization term (untouched)

$$f \equiv 0 : x_{n+1} = \operatorname{prox}_{\alpha_n \lambda \Omega} (x_n - \alpha_n \nabla f(x_n)) = \operatorname{prox}_{\alpha_n \lambda \Omega} (x_n)$$

Proximal point applied to  $\lambda \Omega$

$$\Omega \equiv 0 : x_{n+1} = \operatorname{prox}_{\alpha_n \lambda \Omega} (x_n - \alpha_n \nabla f(x_n)) = x_n - \alpha_n \nabla f(x_n)$$

GD applied to  $f$

- (+) Uses the problem structure (convexity of  $\Omega$ ,  $C^1$  nature of  $f$ )
- (+) Applies when  $f$  is non-convex
- (-) Requires to compute a prox for the function  $\Omega$



# Two key examples

## • $l_2$ regularization

$$\text{minimize}_{x \in \mathbb{R}^d} f(x) + \frac{\lambda}{2} \|x\|_2^2 \quad \lambda > 0$$

$$\text{PG} \quad x_{h+1} = \underset{x \in \mathbb{R}^d}{\text{argmin}} \left\{ \underbrace{f(x_h) + \nabla f(x_h)^T (x - x_h) + \frac{1}{2\alpha} \|x - x_h\|^2 + \frac{\lambda}{2} \|x\|^2}_{P(x)} \right\}$$

$$\nabla P(x) = 0 \quad (\Rightarrow) \quad \nabla f(x_h) + \frac{1}{\alpha} (x - x_h) + \lambda x = 0$$

$$(\Rightarrow) \quad \frac{1 + \lambda\alpha}{\alpha} x = \frac{1}{\alpha} x_h - \nabla f(x_h)$$

$$(\Rightarrow) \quad x = \frac{1}{1 + \lambda\alpha} x_h - \frac{\alpha}{1 + \lambda\alpha} \nabla f(x_h)$$

$$\lambda = 0 \quad x = x_h - \alpha \nabla f(x_h)$$

## PG for $l_2$ regularization

$$x_{h+1} = \frac{1}{1 + \lambda\alpha} x_h - \frac{\alpha}{1 + \lambda\alpha} \nabla f(x_h)$$

$< 1$   $< 1$  for  $\lambda$  off large

$\Rightarrow$  this method performs **weight decay** and **gradient decay** for  $\lambda$  off large

GD with weight decay: GD applied to minimize  $f(x) + \frac{\lambda}{2} \|x\|_2^2$

$$\begin{aligned}x_{k+1} &= x_k - \alpha_k \nabla f(x_k) - \lambda \alpha_k x_k \\ &= \underbrace{(1 - \lambda \alpha_k)}_{\substack{< 1 \\ \text{for } \alpha_k < \frac{1}{\lambda}}} x_k - \alpha_k \nabla f(x_k)\end{aligned}$$

NB: Proximal point method applied to minimize  $f(x) + \frac{\lambda}{2} \|x\|_2^2$

$$\begin{aligned}x_{k+1} &= x_k - \alpha_k \nabla \left( f + \frac{\lambda}{2} \|\cdot\|_2^2 \right) (x_{k+1}) \\ &= x_k - \alpha_k \nabla f(x_{k+1}) - \lambda \alpha_k x_{k+1}\end{aligned}$$

$$[1 + \lambda \alpha_k] x_{k+1} = x_k - \alpha_k \nabla f(x_{k+1})$$

$$x_{k+1} = \frac{1}{1 + \lambda \alpha_k} x_k - \frac{\alpha_k}{1 + \lambda \alpha_k} \nabla f(x_{k+1})$$

## Takeaways

- Regularization  $\equiv$  Pushing for certain properties of the solution
- $l_2, (l_0), (l_1)$   $\rightarrow$  Thursday
- Proximal operator, proximal point, proximal gradient

## Exercise:

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

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

- 1) Show that the solution of (P) is given by the proximal operator of a certain function
- 2) Write down the iterates of proximal gradient applied to (P) with  $f(x) = \frac{1}{2} \|x - u\|_2^2$ ,  $\Omega(x) = \|x\|_1$  and  $\lambda = \alpha$ .  
(Explicit form if you can)