

# Proximal Gradient Descent Algorithms

# Objectives

1

- derive proximal gradient algorithm for regularized least-squares problems
  - least-squares gradient descent
  - regularize
- apply to ridge regression

\* The 5.1 section is to use GD way for solving LS problems

Here we are solving the regularized-Ls in GD, just it.

Proximal gradient descent solves regularized least-squares problems 2

problem:

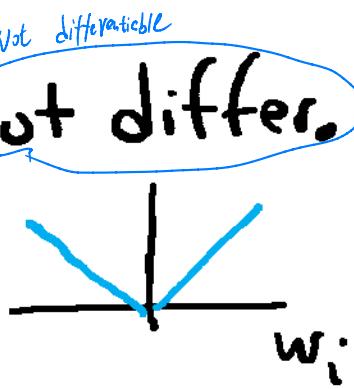
$$\min_{\underline{w}} \|\underline{A}\underline{w} - \underline{d}\|_2^2 + \lambda r(\underline{w})$$

$\lambda > 0$ : tuning parameter

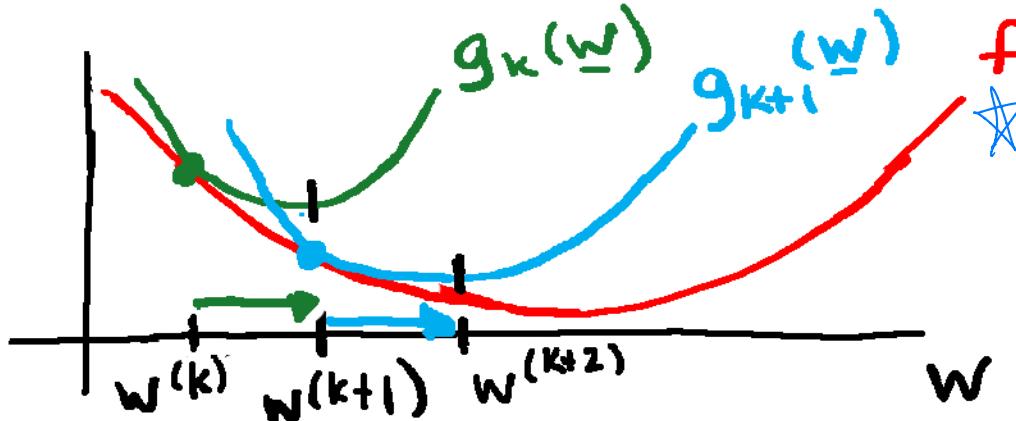
$r(\underline{w})$ : regularizer

Example Convex Regularizers

- Ridge (Tikhonov)  $r(\underline{w}) = \|\underline{w}\|_2^2 = \sum_{i=1}^m w_i^2$
- LASSO ( $\ell_1$ )  $r(\underline{w}) = \|\underline{w}\|_1 = \sum_{i=1}^m |w_i|$  not differ.



Proximal Gradient Descent Concept



$$f(\underline{w}) = \|\underline{A}\underline{w} - \underline{d}\|_2^2 + \lambda r(\underline{w})$$

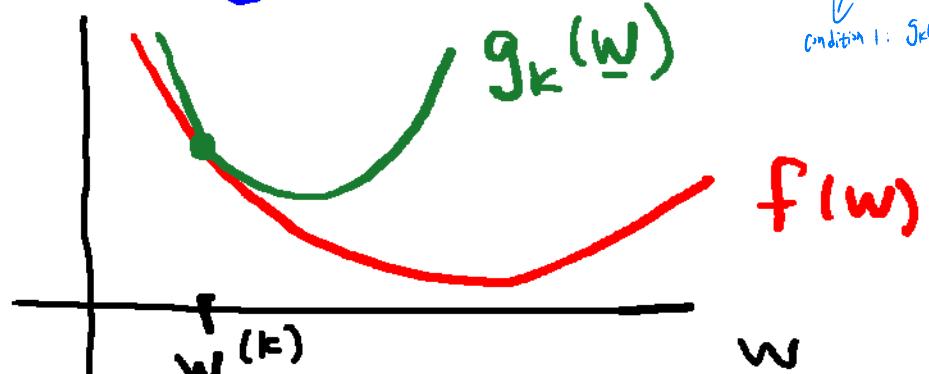
☆ 构建一个近似的  $g_k(\underline{w})$  作为  $f(\underline{w})$  的 upper bound, descent 到  $w^{(k+1)}$  → 构建新的  $g_{k+1}(\underline{w})$  → ...

- solve sequence of simpler problems

- simple for separable  $r(\underline{w}) = \sum_i h_i(w_i)$

means there is no cross terms for  $w_i w_j$ , let's just sum of terms of  $w_i$

Find  $g_k(\underline{w})$  so  $f(\underline{w}) \leq g_k(\underline{w})$ ,  $g_k(\underline{w}^{(k)}) = f(\underline{w}^{(k)})$



Condition 1:  $g_k(\underline{w})$  is upperbound for  $f(\underline{w})$

Condition 2: At iteration k, the functions just touch

minimize  $g_k(\underline{w}) \Rightarrow f(\underline{w})$  decreases

$$\downarrow f(\underline{w}^{(k)}) \leq g_k(\underline{w}^{(k)}) \leq g_k(\underline{w}^{(k)}) = f(\underline{w}^{(k)})$$

$$f(\underline{w}) = \|\underline{d} - \underline{A}\underline{w}\|_2^2 + \lambda r(\underline{w})$$

$$= \|\underline{d} - \underline{A}\underline{w}^{(k)} + (\underline{A}\underline{w}^{(k)} - \underline{A}\underline{w})\|_2^2 + \lambda r(\underline{w})$$

$$f(\underline{w}) = \underbrace{\|\underline{d} - \underline{A}\underline{w}^{(k)}\|_2^2}_{C_k} + \underbrace{\|\underline{A}(\underline{w}^{(k)} - \underline{w})\|_2^2}_{\leq \|\underline{A}\|_{op}^2 \|\underline{w}^{(k)} - \underline{w}\|_2^2} + \underbrace{2(\underline{d} - \underline{A}\underline{w}^{(k)})^\top \underline{A}(\underline{w}^{(k)} - \underline{w})}_{\text{Some guaranteed property from theorem}} + \lambda r(\underline{w})$$

Define step size  $0 < \tau < 1/\|\underline{A}\|_{op}^2 \Rightarrow \frac{1}{\tau} > \|\underline{A}\|_{op}^2$

$$f(\underline{w}) \leq g_k(\underline{w}) = C_k + \frac{1}{\tau} \|\underline{w}^{(k)} - \underline{w}\|_2^2 + 2\underline{V}_k^\top (\underline{w}^{(k)} - \underline{w}) + \lambda r(\underline{w})$$

$g_k(\underline{w})$  is separable for  $r(\underline{w})$  separable:  $g_k(\underline{w}) = C_k + \sum_{i=1}^n q_i(w_i)$  no  $w_i w_j$  terms

可以分解成 n 个互不相关的子问题 (D, 而不是解决一个 n 维的大问题)

Find  $\underline{w}^{(k+1)} = \arg \min_{\underline{w}} g_k(\underline{w})$

$$g_k(\underline{w}) = C_k + \frac{1}{2} \|\underline{w}^{(k)} - \underline{w}\|_2^2 + 2 \underline{v}_k^T (\underline{w}^{(k)} - \underline{w}) + \lambda r(\underline{w})$$

$\mathcal{T}g_k(\underline{w}) = \mathcal{T}C_k + (\underline{w}^{(k)} - \underline{w})^T (\underline{w}^{(k)} - \underline{w}) + 2 \mathcal{T}\underline{v}_k^T (\underline{w}^{(k)} - \underline{w}) + \lambda \mathcal{T}r(\underline{w})$

$$= \mathcal{T}C_k - \mathcal{T}^2 \underline{v}_k^T \underline{v}_k + (\mathcal{T}\underline{v}_k + (\underline{w}^{(k)} - \underline{w}))^T (\mathcal{T}\underline{v}_k + (\underline{w}^{(k)} - \underline{w})) + \lambda \mathcal{T}r(\underline{w})$$

~~$\underline{w}^{(k+1)} = \arg \min_{\underline{w}} \|\underline{z}^{(k)} - \underline{w}\|_2^2 + \lambda \mathcal{T}r(\underline{w})$~~

$$\underline{z}^{(k)} = \underline{w}^{(k)} + \mathcal{T}\underline{v}_k$$

$$= \underline{w}^{(k)} + \mathcal{T}\underline{A}^T (\underline{d} - \underline{A}\underline{w}^{(k)})$$

$$= \underline{w}^{(k)} - \mathcal{T}\underline{A}' (\underline{A}\underline{w}^{(k)} - \underline{d})$$

Least-squares  
gradient descent  
(Landweber)

# Alternate LS gradient descent and regularization 5

$$\underline{w}^{(0)} = \underline{0}, \quad 0 < \tau < \frac{1}{\|A\|_{op}^2}$$

↑ Initial weights      ↑ contain the step size

$$\underline{z}^{(k)} = \underline{w}^{(k)} - \tau A^T(A\underline{w}^{(k)} - \underline{d})$$

initialize

交替执行做 LS 梯度下降和正则化

LS gradient descent

$$\underline{w}^{(k+1)} = \arg \min_{\underline{w}} \|\underline{z}^{(k)} - \underline{w}\|_2^2 + \lambda \tau r(\underline{w})$$

regularize

$$\text{if } \|\underline{w}^{(k+1)} - \underline{w}^{(k)}\| < \varepsilon \text{ stop}$$

check if converged

Regularization simple for  $r(\underline{w})$  separable!

$$\text{if } r(\underline{w}) = \sum_{i=1}^m h_i(w_i) \rightarrow \text{如果 } r(\underline{w}) \text{ 是可分离的}$$

$$\underline{w}^{(k+1)} = \arg \min_{w_i, i=1, \dots, M} \sum_{i=1}^m \left( (\underline{z}_i^{(k)} - w_i)^2 + \lambda \tau h_i(w_i) \right)$$

第三步可以替换成这个

(M个一维最小化问题)

M scalar minimizations

# Example: Ridge Regression (Tikhonov) 6

$$f(\underline{w}) = \|\underline{d} - \underline{A}\underline{w}\|_2^2 + \lambda \|\underline{w}\|_2^2$$

LS gradient descent:

$$\underline{z}^{(k)} = \underline{w}^{(k)} - \tau \underline{A}^T (\underline{A}\underline{w}^{(k)} - \underline{d})$$

Regularization:

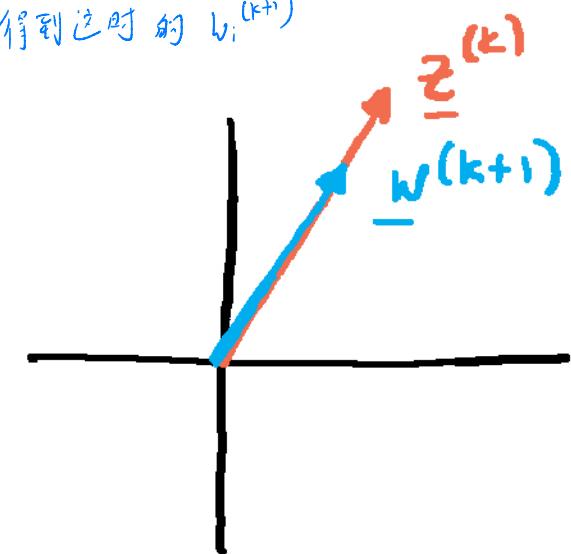
$$\underline{w}^{(k+1)} = \underset{w_i, i=1, \dots, M}{\arg \min} \sum_{i=1}^M (z_i^{(k)} - w_i)^2 + \lambda \tau w_i^2$$



$$\Rightarrow w_i^{(k+1)} = \frac{1}{1 + \lambda \tau} z_i^{(k)}$$

$$\underline{w}^{(k+1)} = \frac{1}{1 + \lambda \tau} \underline{z}^{(k)}$$

"Shrink toward origin" 收缩



**Copyright 2019  
Barry Van Veen**