§ Proximal Gradient Descent.

* 2 steps of Proximal GD:  LS GD → regularize.
* Application: Ridge Regression.

- Proximal GD solves regularized LS problems

$$\min_w \underbrace{\|Aw-d\|_2^2 + \lambda r(w)}_{f(w)} \quad \begin{cases} r(w): \text{regularizer} \\ \lambda: \text{tuning parameter } (\lambda>0) \end{cases}$$

⊛ Some common "Convex regularizer"

  - Ridge (Tikhonov): $r(w) = \|w\|_2^2 = \sum_{i=1}^{M} w_i^2$
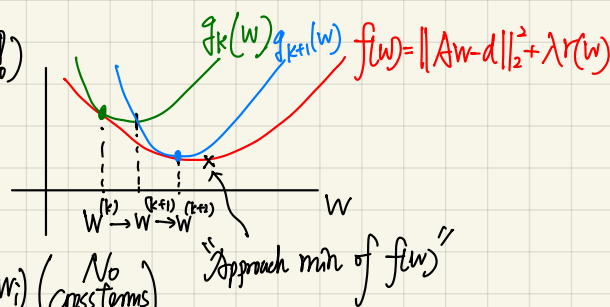  - Lasso ($\ell_1$): $r(w) = \|w\|_1 = \sum_{i=1}^{M} |w_i|$  (not differentiable!)

⊛ Idea: find a $g_k(w)$ s.t. $g_k(w)$ "touches" $f(w)$

  (i)  Solves easier min problem
  (ii) Simple for separable regularizer (★) $r(w) = \sum_i h_i(w_i)$ $\begin{pmatrix} \text{No} \\ \text{cross terms} \\ w_i w_j \end{pmatrix}$

$g_k(w)$  $g_{k+1}(w)$  $f(w) = \|Aw-d\|_2^2 + \lambda r(w)$

$w^{(k)} \to w^{(k+1)} \to w^{(k+2)}$  $w$

"Approach min of $f(w)$"

minimize $g_k(w) \Rightarrow f(w)$ decreases.

$$\Rightarrow \boxed{\text{find } g_k(w) \text{ s.t. } f(w) \le g_k(w), \quad g_k(w^{(k)}) = f(w^{(k)})}$$
(f lives below)

"Find a 'nice' convex upperbound"

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

Consider. $f(w) = \|d-Aw\|_2^2 + \lambda r(w)$

$= \| d - Aw^{(k)} + Aw^{(k)} - Aw \|_2^2 + \lambda r(w)$  — expand.

$= \underbrace{\| d - Aw^{(k)} \|_2^2}_{C_k} + \underbrace{\| A(w^{(k)}-w) \|_2^2}_{\le \|A\|_{op}^2 \|w^{(k)}-w\|_2^2} + 2 \underbrace{(d-Aw^k)A}_{\equiv V_k'}(w^{(k)}-w) + \lambda r(w)$

$\le C_k + \|A\|_{op}^2 \|w^{(k)}-w\|_2^2 + 2V_k'(w^{(k)}-w) + \lambda r(w)$

$\le C_k + \frac{1}{\tau} \|w^{(k)}-w\|_2^2 + 2V_k'(w^{(k)}-w) + \lambda r(w) \equiv g_k(w)$

Note that $g_k(w)$ is separable if $r(w)$ separable $\Rightarrow g_k(w) = C_k + \sum_{i=1}^{M} g_i(w_i)$
(no $w_i w_j$ terms)

<u>Solution.</u> Find $W^{(k+1)} = \underset{w}{\arg\min}\, q_k(w)$,

$$\text{where}\; q_k(w) = C_k + \frac{1}{\tau}\|W^{(k)} - w\|_2^2 + 2V_k'(W^{(k)} - w) + \lambda r(w)$$

$$q_k(w) = C_k + \frac{1}{\tau}\|W^{(k)} - w\|_2^2 + 2V_k'(W^{(k)} - w) + \lambda r(w)$$

$$\Rightarrow \tau q_k(w) = \tau C_k + \underbrace{(W^{(k)} - w)'(W^{(k)} - w)}_{\|\cdot\|_2^2} + 2\tau V_k'(W^{(k)} - w) + \lambda \tau r(w).$$

<span style="color:blue">complete square</span>

$$= \tau C_k + \underbrace{(\tau V_k + (W^{(k)} - w))'(\tau V_k + (W^{(k)} - w))}_{} - \tau^2 V_k' V_k + \lambda \tau r(w)$$

$$= \tau C_k + ((\tau V_k + W^{(k)}) - w)'((\tau V_k + W^{(k)}) - w) - \tau^2 V_k' V_k + \lambda \tau r(w)$$

$$\equiv (Z^{(k)} - w)'(Z^{(k)} - w) + \lambda \tau r(w) + \text{const.} = \|Z^{(k)} - w\|_2^2 + \lambda \tau r(w) + \text{const.}$$

$$\Rightarrow \boxed{W^{(k+1)} = \underset{w}{\arg\min}\, \|Z^{(k)} - w\|_2^2 + \lambda \tau r(w)},$$

$$\text{Where}\; Z^{(k)} \equiv \tau V_k + W^{(k)}$$

$$= W^{(k)} + \tau A'(d - Aw^{(k)})$$

$$= W^{(k)} - \tau A'(Aw^{(k)} - d) \quad \longleftarrow \quad \text{which is "Gradient Descent Iteration"!} \quad \text{exactly. (Landweber).}$$

- <mark>Sum up: Proximal Gradient Descent alternates LS GD & Regularization</mark>

Algorithm. (Proximal GD)

Set $W^{(0)} = 0$, $\tau$ s.t. $0 < \tau < \dfrac{1}{\|A\|_{op}^2}$ $\quad\longleftarrow$ <span style="color:green">initialization</span>

for $k = 1, 2, \cdots$ (to converge)

$\quad Z^{(k)} = W^{(k)} - \tau A'(Aw^{(k)} - d)$ $\quad\longleftarrow$ <span style="color:green">original LS Gradient Descent</span>

$\quad W^{(k+1)} = \underset{w}{\arg\min}\, \|Z^{(k)} - w\|_2^2 + \lambda \tau r(w)$ $\quad\longleftarrow$ <span style="color:green">Regularize.</span>

$\quad$ if $\|W^{(k+1)} - W^{(k)}\| < \varepsilon$: $\quad\longleftarrow$ <span style="color:green">check if $W^{(k+1)}$ converges to $W^{(k)}$</span>

$\quad\quad$ break.

※ Regularization is simple if $r(w)$ separable.

$\Rightarrow$ if $r(w) = \sum_{i=1}^{M} h_i(w_i)$, then $W^{(k+1)} = \underset{w_i, i=1, \cdots, M}{\arg\min} \sum_{i=1}^{M}\left((Z_i^{(k)} - w_i)^2 + \lambda \tau h(w_i)\right)$

which is $M$ scalar minimizations.

**Example.** (Ridge / Tikhonov In Proximal GD)

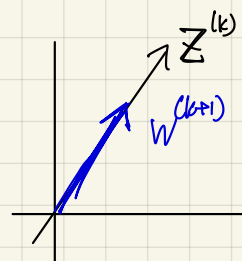$$f(w) = \| d - Aw \|_2^2 + \lambda \| w \|_2^2$$

$1^0$ LS GD: $\quad z^{(k)} = w^{(k)} - \tau A'(Aw^{(k)} - d)$

$2^0$ Regularization: $\quad w^{(k+1)} = \underset{w_i, i=1,\dots,M}{\operatorname{argmin}} \sum_{i=1}^{M} \left( z_i^{(k)} - w_i \right)^2 + \lambda\tau\, \underline{w_i^2}$

$h(w_i) = \| w \|_2^2 \Rightarrow \sum_i w_i^2$

$$\Rightarrow \quad w_i^{(k+1)} = \frac{1}{1+\lambda\tau} z_i^{(k)} \qquad \left( \text{solves FOC } [w_i] \right)$$

$$\Rightarrow \quad w^{(k+1)} = \underbrace{\frac{1}{1+\lambda\tau}}_{(<1)} z^{(k)} \qquad \text{shrinking toward origin}$$

§ LASSO Regression.

(*) Search for sparse solutions.
(*) $\ell_1$-norm regularization (LASSO).

Consider: $Aw = [a_1 \cdots a_M]\begin{bmatrix} w_1 \\ \vdots \\ w_M \end{bmatrix} = \sum_{i=1}^{M} w_i a_i \Rightarrow$ with $w_\ell \approx 0$ it implies $a_\ell$ not important.

$\Rightarrow$ If only a few $w_\ell$'s are "non-zero" (important)
then we have "sparse $w$"

$\|w\|_0 \equiv \sum_{i=1}^{M} \mathbb{1}\{w_i \neq 0\}$ (counting non-zero $w_\ell$).

$\circledast$ $\ell_0$-norm is <u>NOT</u> norm as $\|aw\|_0 \neq a\|w\|_0$

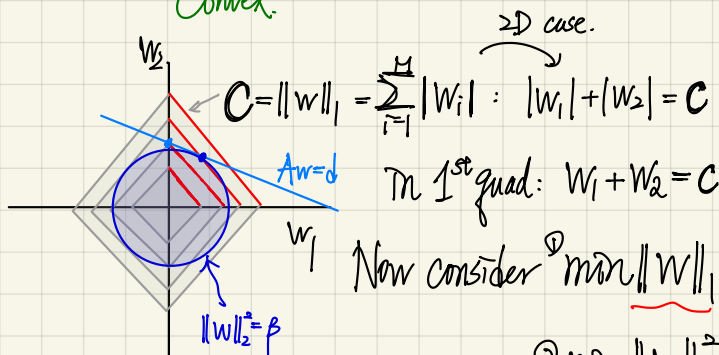Consider. $\min_w \|w\|_0$ s.t. $\|Aw - d\|_2^2 < \varepsilon$.

$\otimes$ Problem: $\|w\|_0$ not convex $\Rightarrow$ "Computationally Intractable."

• Convex Relaxation gives tractable problem

$$\min_w \|w\|_1 \quad \text{s.t.} \quad \|Aw - d\|_2^2 < \varepsilon.$$
  $\underbrace{\phantom{\min_w \|w\|_1}}_{\text{Convex.}}$

Least Absolute Selection &
Shrinkage Operator (LASSO)



2D case.

$C = \|w\|_1 = \sum_{i=1}^{M} |w_i| : |w_1| + |w_2| = C$

In 1st quad: $w_1 + w_2 = C$

Now consider ① $\min \|w\|_1$ s.t. $Aw = d$. $\Rightarrow$ "Corner" on $\|w\|_1$ sparse solutions.

② $\min \|w\|_2^2$ s.t. $Aw = d$. $\Rightarrow$ "Circular" $\|w\|_2^2$ (less likely corner) non-sparse solutions

• LASSO is a regularized LS problem.

LASSO: $\min_w \|w\|_1$ s.t. $Aw - d < \varepsilon$. $\overset{\text{eqv.}}{\rightsquigarrow} \min_w \|Aw - d\|_2^2 + \lambda \|w\|_1$ for some $\lambda, \varepsilon$.

$\left( \text{eqv. } \min_w \|w\|_1 + \frac{1}{\lambda} \|Aw - d\|_2^2 \right)$

LASSO. $w_L = \arg\min_w \|Aw - d\|_2^2 + \lambda \|w\|_1$ : Sparse $w_L$ ; Can have small <u>model</u> error ; Iterative Solution
$w_{opt} - w_L$. Method.

Ridge $w_R = \arg\min_w \|Aw - d\|_2^2 + \lambda \|w\|_2^2$ : non-sparse $w_R$ ; Can have small prediction error ; solve in
$\|Aw_{opt} - Aw_R\|_2^2$ closed form.

- LASSO & Feature Selection.

$$W_L = \underset{w}{\text{argmin}} \, \|Aw - d\|_2^2 + \lambda \|w\|_1.$$

$1°$ Selection: $S_L = \{ i : [W_L]_i \neq 0 \}$ (the non-zero $W_{L_i}$)

$2°$ $AW_L = \sum_{i=1}^{M} a_i [W_L]_i = \sum_{i \in S_L} a_i [W_L]_i$

$3°$ Debiasing: $A_L = \{ a_i : i \in S_L \}$ (debias $A$ by those $a_i$'s selected by LASSO)

$4°$ Re-solve LS problem: $\widehat{W_L} = \underset{w}{\text{argmin}} \, \|A_L w - d\|_2^2 = (A_L' A_L)^{-1} A_L' d.$ avoids $\|w\|_1$ shrinkage.

# § LASSO & Proximal GD.

- l₁-regularized LS problem can be solved by Proximal GD.
$$\min_w \|Aw-d\|_2^2 + \lambda\|w\|_1 \quad \text{encourages sparse soln.}$$
$\begin{cases} GD. \\ \text{Regularization ("shrinkage")} \end{cases}$

Apply "Proximal GD" (since no close-form)

$1°$  $Z^{(k)} = W^{(k)} - \tau A'(Aw^{(k)} - d)$  ← LS GD

★ $2°$  $W^{(k+1)} = \underset{w}{\text{argmin}} \|Z^{(k)} - w\|_2^2 + \tau\lambda\|w\|_1$  ← regularization.

- Regularization steps involves "scalar minimizations"
$$\min_w \|Z^{(k)} - w\|_2^2 + \tau\lambda\|w\|_1 \Rightarrow \min_{w_i, i=1,\dots,M} \sum_{i=1}^{M} (Z_i^{(k)} - W_i)^2 + \lambda\tau|W_i| \quad (\lambda\tau > 0)$$

Case ① : $W_i \geq 0$  (1st quad).
$$\Rightarrow \min_{W_i} (Z_i - W_i)^2 + \lambda\tau W_i \Rightarrow [W_i]: 0 = -2(Z_i - W_i) + \lambda\tau \Rightarrow W_i = Z_i - \tfrac{1}{2}\lambda\tau$$

So, $W_i = \begin{cases} Z_i - \tfrac{1}{2}\lambda\tau, & \text{if } Z_i > \tfrac{1}{2}\lambda\tau. \\ 0, & \text{if } Z_i < \tfrac{1}{2}\lambda\tau. \end{cases}$
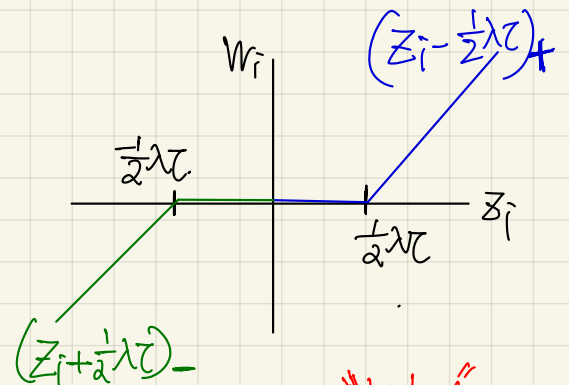
$W_i = (Z_i - \tfrac{1}{2}\lambda\tau)_+$

Case ② : $W_i \leq 0$  (3rd quad).
$$\Rightarrow \min_{W_i} (Z_i - W_i)^2 - \lambda\tau W_i \Rightarrow [W_i]$$

So, $W_i = \begin{cases} 0, & \text{if } Z_i > \tfrac{-1}{2}\lambda\tau \\ Z_i + \tfrac{1}{2}\lambda\tau, & \text{if } Z_i < \tfrac{-1}{2}\lambda\tau \end{cases}$

$W_i = (Z_i + \tfrac{1}{2}\lambda\tau)_-$


$(Z_i - \tfrac{1}{2}\lambda\tau)_+$
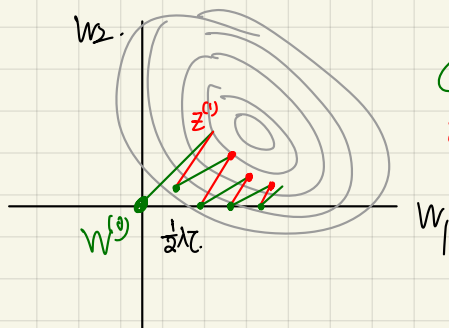$(Z_i + \tfrac{1}{2}\lambda\tau)_-$

"Soft threshold"
$$W_i = \begin{cases} Z_i - \tfrac{1}{2}\lambda\tau, & Z_i \in (\tfrac{1}{2}\lambda\tau, \infty) \\ 0, & Z_i \in (\tfrac{-1}{2}\lambda\tau, \tfrac{1}{2}\lambda\tau) \\ Z_i + \tfrac{1}{2}\lambda\tau, & Z_i \in (-\infty, \tfrac{-1}{2}\lambda\tau) \end{cases}$$
"Shrinkage"

$(|Z_i| - \tfrac{1}{2}\lambda\tau)_+ \text{sign}\{Z_i\}.$

- It alternates Descent & Shrinkage. (soft thresholding).

$w_2$.



$w^{(0)}$   $\frac{1}{2}\lambda\tau$.   $w_1$

$z^{(1)}$

Gradient Descent

Shrinkage. (most likely send $w$ back to 0)