

# Regularization

2020年9月30日 8:58

Minimize  $\frac{1}{2N} \sum_i (w^T x_i - y_i)^2$  with  $\|w\|$  constrained

$$\hookrightarrow \min_w \frac{1}{2N} \sum_i (w^T x_i - y_i)^2 + \frac{\lambda}{2} \|w\|^2$$

called ridge regression

$$\Rightarrow \nabla_w L = \frac{1}{N} \sum_i (w^T x_i - y_i) x_i + \lambda w$$

$$\text{Hessian } H = \frac{1}{N} \sum_i x_i x_i^T + \lambda I \Rightarrow \lambda_{\min}(H) \geq \lambda \quad \text{Strongly Convex!}$$

LASSO : sparse solutions.

Find "important" features from a large number of features

$$\Rightarrow \text{Minimize } \frac{1}{2N} \sum_i (w^T x_i - y_i)^2 \text{ with } \|w\|_0 \leq c$$

(find first  $c$  most important)

↓

In practice, use  $L_1$ -loss for regularization

$$\min_w \frac{1}{2N} \sum_i (w^T x_i - y_i)^2 + \lambda \|w\|_1$$

## Compressed Sensing

Nyquist Theorem : for a signal with frequency  $f$ , need  $2f$  sampling rate to fully construct the signal.

- But sometimes we can compress (video/image...) much smaller

Suppose  $x$  is a long, sparse vector

$A = [a_1, \dots, a_n]^T \in \mathbb{R}^{n \times d}$  ( $n \ll d$ ).  $\Leftarrow$  measurement matrix.

$$\Rightarrow \begin{array}{|c} \mathbf{y} \\ \hline \end{array} = \begin{array}{|c|} \hline \mathbf{A} \\ \hline \end{array} \begin{array}{|c} \mathbf{x} \\ \hline \end{array}$$

(sparse)

- Restricted Isoperimetric Property

$W \in \mathbb{R}^{n \times d}$  is  $(\epsilon, S)$ -RIP if  $\forall x \neq 0$  s.t.  $\|x\|_0 \leq S$ ,

... ..

$W \in \mathbb{R}^{n \times d}$  is  $(\epsilon, S)$ -RIP if  $\forall x \neq 0$  s.t.  $\|x\|_0 \leq S$ ,

we have  $(1-\epsilon)\|x\|_2^2 \leq \|Wx\|_2^2 \leq (1+\epsilon)\|x\|_2^2$

Theorem 1.  $\epsilon < 1$ .  $W$  be  $(\epsilon, 2S)$ -RIP.  $\|x\|_0 \leq S$ .

$y = Wx$ . let  $\tilde{x} \in \underset{v: Wv=y}{\text{argmin}} \|v\|_1$  be a reconstructed vector of  $x$ .

$\Rightarrow$  Then,  $\tilde{x} = x$ .

Proof. If  $\tilde{x} \neq x$ .

$\therefore Wx = y, \|\tilde{x}\|_0 \leq \|x\|_0 \leq S$

$\therefore \|\tilde{x} - x\|_0 \leq 2S$

$\Rightarrow (1-\epsilon)\|\tilde{x} - x\|_2^2 \leq \|W(\tilde{x} - x)\|_2^2 \leq (1+\epsilon)\|\tilde{x} - x\|_2^2$

$\downarrow$   
 $= 0$

$\Rightarrow \epsilon \geq 1$  X.

Since  $\underset{v: Wv=y}{\text{argmin}} \|v\|_1$  is hard to compute, we approximate it using L1.

Theorem 3.  $\epsilon < \frac{1}{1+\sqrt{2}}$ . Let  $W$  be  $(\epsilon, 2S)$ -RIP.  $x$  any vector.  $x_S \in \underset{v: \|v\|_0 \leq S}{\text{argmin}} \|x - v\|_1$ .

let  $y = Wx$ .

$x^* \in \underset{v: Wv=y}{\text{argmin}} \|v\|_1$ .

$x_S$  has  $x$ 's largest elements

$\Rightarrow \|x^* - x\|_2 \leq 2(1-\rho)^{-1} S^{\frac{1}{2}} \|x - x_S\|_1, \rho = \frac{\sqrt{2}\epsilon}{1-\epsilon}$ .

intuition:  $h$  sparse  $\Rightarrow$  easy  $\Rightarrow \|h\|_2 \leq \|h'\|_1 + \|h''\|_2$

$2S$ -sparse  $(d-2S)$  entries

Lemma 1. RIP  $\Rightarrow$  Almost Orthogonality.

$W$ .  $(\epsilon, 2S)$ -RIP.  $\forall$  set  $I, J$  disjoint of size  $\leq S$

$\forall u, \langle Wu_I, Wu_J \rangle \leq \epsilon \|u_I\|_1 \|u_J\|_1$

e.x.  $u = \begin{pmatrix} 5 \\ 6 \\ 7 \\ 8 \end{pmatrix}, I = \{1, 2\}, J = \{3, 4\}$   
 $\Rightarrow u_I = \begin{pmatrix} 5 \\ 6 \\ 0 \\ 0 \end{pmatrix}, u_J = \begin{pmatrix} 0 \\ 0 \\ 7 \\ 8 \end{pmatrix}$

Proof. WLOG. assume  $\|u_I\| = \|u_J\| = 1$ .

$\langle Wu_I, Wu_J \rangle = \frac{\|Wu_I + Wu_J\|^2 - \|Wu_I - Wu_J\|^2}{4}$

$$\langle Wu_i, Wu_j \rangle = \frac{\|Wu_i + Wu_j\|^2 - \|Wu_i - Wu_j\|^2}{4}$$

Since  $|I \cup J| \leq 2s$ . by RIP

$$\|Wu_i + Wu_j\|^2 \leq (1+\epsilon) (\|u_i\|^2 + \|u_j\|^2) = 2(1+\epsilon).$$

$$-\|Wu_i - Wu_j\|^2 \leq -(1-\epsilon) (\|u_i\|^2 + \|u_j\|^2) = -2(1-\epsilon)$$

$$\Rightarrow \langle Wu_i, Wu_j \rangle \leq \epsilon \|u_i\| \|u_j\|.$$

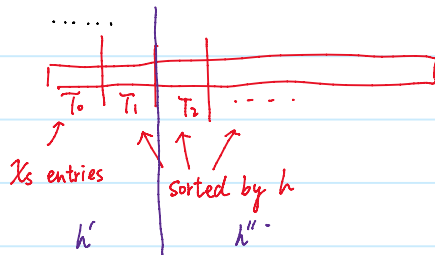
$$[d] = T_0 \cup T_1 \cup \dots \cup T_{\lfloor d/s \rfloor - 1}, |T_i| = s.$$

$$d \% s = 0$$

$T_0$ :  $s$  largest elem of  $\underline{x}$  (all elems of  $\mathcal{X}_s$ )  
(abs)

$$T_0^c = [d] \setminus T_0.$$

$\Rightarrow T_i$ :  $s$  largest of  $\underline{h}$  on  $T_0^c$



Proof (Main Theorem)

$$h'' \rightsquigarrow \textcircled{1} \|h_{T_0^c}^c\|_2 \leq \|h_{T_0}\|_2 + 2s^{-1/2} \|x - x_s\|_1$$

$$h' \rightsquigarrow \textcircled{2} \|h_{T_0}^c\|_2 \leq \frac{\rho}{1-\rho} s^{-1/2} \|x - x_s\|_1$$

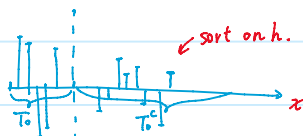
$$\text{Then } \|h\|_2 \leq \|h_{T_0}^c\|_2 + \|h_{T_0^c}^c\|_2$$

$$\leq \|h_{T_0}^c\|_2 + \|h_{T_0}\|_2 + 2s^{-1/2} \|x - x_s\|_1$$

$$\leq 2\|h_{T_0}^c\|_2 + 2s^{-1/2} \|x - x_s\|_1$$

$$\leq 2(1-\rho)^{-1} s^{-1/2} \|x - x_s\|_1$$

$$\textcircled{1} \|x\|_1 \geq \|x + h\|_1 = \|x^*\|_1 \quad (\because x^* = \operatorname{argmin}_{v: Wv=y} \|v\|_1)$$



split  $h$  into two parts

split  $h$  into two parts

$$\|h_{T_0}^c\| \text{ cannot be too large } \xrightarrow{I^{T_0}} \Rightarrow \|h\| \leq 6$$

Another rigorous proof to this:

$$\begin{aligned} \|x\| &\geq \|x+h\| = \sum_{i \in T_0} |x_i + h_i| + \sum_{i \notin T_0} |x_i + h_i| \\ &\geq \|x_{T_0}\| - \|h_{T_0}\| + \|h_{T_0^c}\| - \|x_{T_0^c}\| \end{aligned}$$

$$\text{since } \|x_{T_0^c}\| = \|x - x_{T_0}\| = \|x\| - \|x_{T_0}\|$$

$$\Rightarrow \|h_{T_0}^c\| \leq \|h_{T_0}\| + 2\|x_{T_0^c}\|$$

$$\Rightarrow \|h_{T_0}^c\| \leq \|h_{T_0}\| + 2\|x_{T_0^c}\|$$

$h$  在  $T_0$  上最大增量

$$\forall j > 1 \quad \forall i \in T_j, \quad i' \in T_{j-1}, \quad \|h_i\| \leq \|h_{i'}\|$$

$$\Rightarrow \|h_{T_j}\|_\infty \leq \|h_{T_{j-1}}\|_1 / S$$

$$\Rightarrow \|h_{T_j}\|_2 \leq S^{1/2} \|h_{T_j}\|_\infty \leq S^{-1/2} \|h_{T_{j-1}}\|_1$$

$$\text{Triangle} \quad \|h_{T_{0,1}}^c\|_2 \leq \sum_{j=2} \|h_{T_j}\|_2 \leq S^{-1/2} \sum_{j=1}^{d/2-2} \|h_{T_j}\|_1 \leq \|h_{T_0}^c\|_1 S^{1/2}$$

$$\Rightarrow \|h_{T_{0,1}}^c\|_2 \leq \|h_{T_0}\|_2 + 2S^{1/2} \|x - x_S\|_1$$

$$\textcircled{2} \quad (1-\varepsilon) \|h_{T_{0,1}}\|_2^2 \leq \|Wh_{T_{0,1}}\|_2^2$$

$$\text{Since } Wh_{T_{0,1}} = Wh - \sum_{j \geq 2} Wh_{T_j} = - \sum_{j \geq 2} Wh_{T_j}$$

$W(x^c - x) = 0$

$$\Rightarrow \|Wh_{T_{0,1}}\|_2^2 = - \sum_{j \geq 2} \langle Wh_{T_0} + Wh_{T_1}, Wh_{T_j} \rangle$$

By almost-orthogonality,

$$\leq \varepsilon \sum_{j \geq 2} (\|h_{T_0}\|_1 + \|h_{T_1}\|_1) \|h_{T_j}\|_2$$

$$\leq \sqrt{2} \varepsilon \|h_{T_{0,1}}\|_2^2 \sum_{j \geq 2} \|h_{T_j}\|_2 \quad \text{see } \textcircled{1}$$

$$\Rightarrow \|h_{T_{0,1}}\|_2^2 (1-\varepsilon) \leq \sqrt{2} \varepsilon S^{1/2} \|h_{T_0}^c\|_1$$

$$\Rightarrow \|h_{T_{0,1}}\|_2^2 \leq \frac{\sqrt{2} \varepsilon}{1-\varepsilon} S^{1/2} \|h_{T_0}^c\|_1$$

$$\leq \rho S^{1/2} (\|h_{T_0}\|_1 + 2\|x_{T_0^c}\|_1)$$

$$\Rightarrow \|h_{T_{0,1}}\|_2^2 \leq \frac{\rho S^{1/2}}{1-\rho} 2\|x - x_S\|_1$$

Q.E.D of Theorem 3.



Theorem 4. Construct RIP Matrix

Let  $U$  be arbitrary  $n \times n$  orthonormal matrix.  $\epsilon, \delta \in (0, 1)$   $s$  integer in  $[d]$

Let  $n$  be an integer:  $n \geq 100 \frac{s \ln(40d/\epsilon\delta)}{\epsilon^2}$

$\Rightarrow$  Let  $W \in \mathbb{R}^{n \times d}$  be a matrix st. each element of  $W$  distributed normally with  $\mu=0, \sigma^2 = \frac{1}{n}$ .

Then w.p. of  $\geq 1-\delta$ ,  $WU$  is  $(\epsilon, \delta)$ -RIP.

Covering Balls. *Make infinite into finite*

Lemma 2.  $\epsilon \in (0, 1)$ .  $\exists Q \subset \mathbb{R}^d, |Q| \leq \left(\frac{5}{\epsilon}\right)^d$  st.

$$\sup_{x: \|x\| \leq 1} \min_{v \in Q} \|x - v\| \leq \epsilon.$$

Proof.  $Q' = \{x \in \mathbb{R}^d: \forall j, \exists i \in \{-k, \dots, k\} \text{ s.t. } x_j = \frac{i}{k}\}$

$|Q'| = (2k+1)^d$ , set  $Q = Q' \cap B_2(1)$ .

Lemma 5 (JL Lemma).  $Q$ : a finite set of vectors in  $\mathbb{R}^d$ .  $\delta \in (0, 1)$ .  $n \in \mathbb{Z}$

$$\text{s.t. } \epsilon = \sqrt{\frac{6 \ln(2|Q|/\delta)}{n}} \leq 3$$

$\Rightarrow$  w.p. of  $\geq 1-\delta$ . a random  $W \in \mathbb{R}^{n \times d}$  with elements iid.  $N(0, \frac{1}{n})$ .

$$\Rightarrow \sup_{x \in Q} \left| \frac{\|Wx\|_2^2}{\|x\|_2^2} - 1 \right| < \epsilon.$$

Lemma 3. Let  $U$  be an  $d \times d$  orthonormal matrix.

Index Set  $I \subset [d], |I|=s$ .

$U_i$ :  $i^{\text{th}}$  column of  $U \Rightarrow S = \text{span}\{U_i, i \in I\}$ .  $\delta, \epsilon \in (0, 1)$ .

$$n \geq 24 \frac{\ln(2/\delta) + \epsilon \ln(20/\epsilon)}{\epsilon^2}$$

$\Rightarrow$  w.p. of  $\geq 1-\delta$ . a random  $W \in \mathbb{R}^{n \times d}$  with elements iid.  $N(0, \frac{1}{n})$ .

$$\Rightarrow \sup_{x \in S} \left| \frac{\|Wx\|_2}{\|x\|_2} - 1 \right| < \epsilon.$$

$\hookrightarrow$  Apply Union Bound to all  $I \Rightarrow n \geq 24 \frac{s \ln(40d/\epsilon\delta)}{\epsilon^2}$

So we prove theorem 4.

Proof. WLOG.  $\|x\|_2 = 1$

Since  $x \in \text{span}\{U_i\}$ .  $x = U_i a$ ,  $a \in \mathbb{R}^s$

$\|a\|_2 = 1$ . so pick  $Q$  of size  $|Q| \leq \left(\frac{20}{\epsilon}\right)^s$  (by Lemma 2)

$$\text{s.t. } \sup_{a: \|a\|_2=1} \min_{v \in Q} \|a - v\| \leq \frac{\epsilon}{4}$$

( $U$  orthonormal)

$$\alpha \cdot \|\alpha\| = 1 \quad \forall \alpha$$

( $U$  orthonormal)

$$\Rightarrow \sup_{\alpha: \|\alpha\|=1} \min_{v \in \mathcal{B}} \|U_1 \alpha - U_1 v\| \leq \frac{\epsilon}{4}$$

by JL lemma, with prob.  $\geq 1 - \delta$

$$\sup_{v \in \mathcal{B}} \left| \frac{\|W(U_1 v)\|}{\|U_1 v\|} - 1 \right| \leq \frac{\epsilon}{2}$$

Let  $\alpha$  be the smallest number s.t.  $\forall x \in S, \frac{\|Wx\|}{\|x\|} \leq 1 + \alpha$

we want to show  $\alpha \leq \epsilon$

$$\|Wx\| \leq \|WU_1 v\| + \|W(x - U_1 v)\|$$

$$\leq 1 + \frac{\epsilon}{2} + (1 + \alpha) \frac{\epsilon}{4}$$

$$\Rightarrow \alpha \leq \frac{\epsilon}{2} + (1 + \alpha) \frac{\epsilon}{4} \Rightarrow \alpha \leq \epsilon$$

The other side:  $\|Wx\| \geq \|WU_1 v\| - \|W(x - U_1 v)\| \geq 1 - \frac{\epsilon}{2} - (1 + \epsilon) \frac{\epsilon}{4} \geq 1 - \epsilon$ . Q.E.D

More general. what if  $W$  is non-linear?

Neural Network: images are dense in pixel space,  
but sparse in latent space (feature space).