

Mohammed Alshurbaji

Adv. ML

Assignment 1

Q1:- Lasso (Least Absolute Shrinkage and Selection operator)

$$\frac{1}{2} \|AX - y\|_2^2 + \lambda \|X\|_1$$

Algorithm:-

1) Initial guess for the coeff.  $X$ . (start with zeros)

2) Iteration for each coordinate:  $(X_j)$

→ Fix all other coordinates  ~~$X$~~   $X_{-j}$

→ Update  $X_j$  by minimizing the obj. w.r.t.  $X_j$ .

$$\frac{1}{2} \|AX - y\|_2^2 + \lambda |X_j|$$

observed

3) Solve this as ~~seen~~ in class :- (Soft Thresholding)

$$= \begin{cases} -\frac{b+\lambda}{2A} & \text{if } b > \lambda \\ 0 & \text{if } |b| \leq \lambda \\ -\frac{b-\lambda}{2A} & \text{if } b \leq -\lambda \end{cases}$$

$$\boxed{ax^2 + bx + \lambda |x|}$$

Summary:- It's supposed that Min. coordinate to be much faster, but it's not the case here. maybe because sklearn is using more optimization while I am using the basic one.



## Q2: Least Squares Extension:

Assume:

$$A, X, C, Y, U \in \mathbb{R}^{n \times n}$$

$$\text{Given: } \text{vec}(AUC) = (C^T \otimes A) \text{vec}(U) \quad \dots \textcircled{1}$$

Use Least Squares to solve:

$$\min_X \|AX + XC - Y\|_F^2$$

Now Review Basics:

$$\text{vec}(AUC) \equiv \text{vectorization}; \quad AUC \equiv \begin{bmatrix} 1 \times n \\ n \times n \\ n \times n \end{bmatrix} \equiv \begin{bmatrix} 1 \times n \end{bmatrix}$$

Kronecker Product: multiply each element of first matrix with the second matrix.

$$\text{Ex: } A \otimes B; \quad A = [n \times m], B = [q \times p] \\ \therefore A \otimes B \equiv [nq \times mp]$$

To prove  $\textcircled{1}$ :

$$\text{vec}(AUC) \equiv [n^2 \times 1]$$

$$C^T \otimes A \equiv [n^2 \times n]; \quad \text{vec}(U) = [n^2 \times 1];$$

$$\therefore (C^T \otimes A) \text{vec}(U) = [n^2 \times 1] \quad \#$$

Frobenius:  $\|A\|_F$  (It's the same as Euclidean but for matrices)

$$\|A\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^n |a_{ij}|^2}$$

$$\text{Note That: } \|A\|_F = \text{Tr}(A^T A)$$

$$\textcircled{2} \quad \text{vec}(A+B) = \text{vec}(A) + \text{vec}(B)$$



Starting with the obj. Function:

$$\min_x \|Ax + XC - Y\|_F^2 = (Ax + XC - Y)^T (Ax + XC - Y) \dots \textcircled{1}$$

Using Vectorization:

$$\begin{aligned} \text{Vec}(Ax + XC - Y) &= \text{Vec}(Ax) + \text{Vec}(XC) - \text{Vec}(Y) \\ &= (I \otimes A) \text{Vec}(X) + (C^T \otimes I) \text{Vec}(X) - \text{Vec}(Y) \\ &= (I \otimes A + C^T \otimes I) \text{Vec}(X) - \text{Vec}(Y) \dots \textcircled{2} \end{aligned}$$

Subs.  $\textcircled{2}$  in the Obj. Funct.:-

$$\min \| (I \otimes A + C^T \otimes I) \text{Vec}(X) - \text{Vec}(Y) \|_F^2$$

Take gradient and set to zero:-

$$\underbrace{(I \otimes A + C^T \otimes I)^T}_{\text{II}} (I \otimes A + C^T \otimes I) \text{Vec}(X) = \underbrace{(I \otimes A + C^T \otimes I)^T \text{Vec}(Y)}_{\text{III}}$$

Solving II:-

First:

$$\begin{aligned} \text{II} &= (I \otimes A^T + (C \otimes I))(I \otimes A + C^T \otimes I) \\ &= (I \otimes A^T)(I \otimes A) + (I \otimes A^T)(C^T \otimes I) + (C \otimes I)(I \otimes A) \\ &\quad + (C \otimes I)(C^T \otimes I) \\ &\stackrel{I \otimes}{=} (A^T A) + I \otimes (A^T C^T) + I \otimes (C A) + I \otimes (C C^T) \end{aligned}$$

Secondly Solving III:-

$$\text{III} = (I \otimes A + C^T \otimes I)^T = I \otimes A^T + C \otimes I$$

$\therefore$  Solution of  $\text{Vec}(X)$  is

$$\text{Vec}(X) = \text{II}^{-1} \text{III} \text{Vec}(Y) \neq$$

For simplicity, let  $Y = I \otimes I$

$$W = (I \otimes A) + (C^T \otimes I)$$

$$\text{vec}(X) = (W^T W)^{-1} W^T \text{vec}(Y)$$

$$(Y)_{\text{vec}} = (X)_{\text{vec}} + (X A)_{\text{vec}} = (Y - X A)_{\text{vec}}$$

$$(Y)_{\text{vec}} = (X)_{\text{vec}} (I \otimes I + A \otimes I)$$

$$\| (Y)_{\text{vec}} - (X)_{\text{vec}} (I \otimes I + A \otimes I) \|_{\text{vec}}$$

$$\| (I \otimes I + A \otimes I) (X)_{\text{vec}} - (Y)_{\text{vec}} \|_{\text{vec}} = \| (I \otimes I + A \otimes I) (X)_{\text{vec}} - (Y)_{\text{vec}} \|_{\text{vec}}$$

$$(I \otimes I + A \otimes I) (I \otimes I + A \otimes I)^T = (I \otimes I + A \otimes I) (I \otimes I + A \otimes I)$$

$$(I \otimes I + A \otimes I) (I \otimes I + A \otimes I)^T = (I \otimes I + A \otimes I) (I \otimes I + A \otimes I)$$

$$(I \otimes I + A \otimes I) (I \otimes I + A \otimes I)^T = (I \otimes I + A \otimes I) (I \otimes I + A \otimes I)$$

$$(I \otimes I + A \otimes I) (I \otimes I + A \otimes I)^T = (I \otimes I + A \otimes I) (I \otimes I + A \otimes I)$$

$$(Y)_{\text{vec}} = (X)_{\text{vec}} (I \otimes I + A \otimes I)$$



### Q3: Ridge Regression :-

$$J(x) = \underbrace{\min_x \|Ax - y\|_2^2}_I + \underbrace{\lambda \|x\|_2^2}_{II}$$

#### I: Using Composite:-

$$g(x) = Ax - y$$

$$J = \|g\|_2^2$$

$$\frac{\partial J}{\partial x} = \frac{\partial g}{\partial x} \cdot \frac{\partial J}{\partial g} \quad (\text{Chain Rule})$$

$$= A^T \cdot 2g$$

$$= 2A^T(Ax - y)$$

#### II: $2\lambda x$

$$\text{Then, set } \frac{\partial J(x)}{\partial x} = 0$$

$$2A^T(Ax - y) + 2\lambda x = 0$$

$$\therefore A^T(Ax - y) + \lambda x = 0$$

$$A^T Ax + \lambda x = A^T y$$

$$\therefore x = (A^T A + \lambda I)^{-1} A^T y \quad \#$$

Now: To prove that it's monotonically decreasing (As shown in class).

$$\text{Eigenval. Decomp. :- } A^T A = Q \Lambda Q^T$$

$$A^T A + \lambda I = Q(\Lambda + \lambda I)Q^T$$

$$\therefore x = Q(\Lambda + \lambda I)^{-1} Q^T A^T y$$

$$\|x\|_2 = \|Q(\Lambda + \lambda I)^{-1} Q^T A^T y\|_2 \quad \text{; } Q \text{ is ortho.}$$

# Thus, as  $\lambda$  increases, the inverse decreases which reduces  $\|x\|_2$ . #



## Q4: Linear Regression and its Extension

Boston houses Datasets 506 records, 13 features (X)

14<sup>th</sup> feature  $\equiv$  median house price  $\equiv Y$

- All features are cont. except feature 4, binary. (Consider it cont.)
- Standardization: The process of transforming data to have a mean of zero and STD of 1.

Section 2: All of these why Questions are related to Bias-Variance Trade off.

Basically: For the standard Lin. Reg. The training error at low  $\lambda$  is simple <sup>(High Variance)</sup> as  $\lambda$  increases becomes more generalize, the error increases. While the test error moving in an opposite way, since the model ~~is~~ is learning with  $n$  increases.

While the meet point means that the training and testing are representative for the same generalized Model.

Section 3: Also, Bias-Variance Tradeoff. While at the beginning too simple then increasing data is helpful. Then the model start fitting the Noise (overfitting happened).

### Section 4:-

As the regularization operator  $\lambda$  increased at the beginning the errors decreases then at certain point we reach the under fitting as the  $\lambda$  keeps the coe. to be <sup>too</sup> small. (Underfitting).

Golden Equation.

$$\sigma^2 \left( \frac{1}{1 + \frac{\lambda}{\sigma^2}} \right)$$