

# Final Exam, CPSC 8420, Fall 2024

Alshurbaji, Mohammad

Dr. Kai Lu.

## Problem 1: Lasso Problem

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

1) Prove that if  $\lambda \geq \|X^T y\|_\infty$ , then  $B^* = 0$

Sol.:

Infinity Norm:  $\|X^T y\|_\infty = \max_i |X^T y|_i$

Note: In general when  $\lambda$  is large, some coe.

in  $B$  becomes zero, Removing less important features.  
(From its name: Shrinkage + Selection).

So, To find optimal sol., take gradient = 0:

$$-X^T(y - X B^*) + \lambda z = 0 \quad ; \quad z_i = \begin{cases} \text{Sign}(B_i^*), & B_i^* \neq 0 \\ \in [-1, 1], & B_i^* = 0 \end{cases}$$

$$\therefore X^T(y - X B^*) = \lambda z$$

let's start from  $B^* = 0$ , and see what we get:

$$X^T(y - 0) = \lambda z$$

$$\therefore \lambda z = X^T y$$

Now: when  $B_i^* = 0$ ,  $z_i \in [-1, 1]$  for  $|X^T y| \leq \lambda$

hence:

$$\lambda \geq \|X^T y\|_\infty = \max_i |X^T y|_i$$

means:

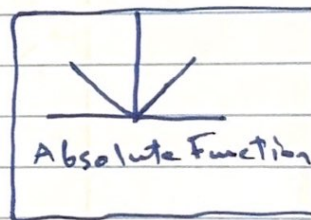
$$X^T y = \lambda z \quad \text{with } z \in [-1, 1]$$

2) This has been proved in Python, using two methods:

i) Descent Coordination.

ii) sklearn Lasso Library.

#



Problem 2:

✓ SVD Decomposition :  $X = U \Sigma V^T$

Define  $\|X\|_2 = \Sigma(1,1)$ ,  $\|X\|_F = \sqrt{\sum_i \sum_j |X_{ij}|^2}$

Prove:  $\|X\|_F \geq \|X\|_2$  and indicate when the equality holds.

Sol.:

1)  $\|X\|_2 = \Sigma(1,1) = \sigma_1$  (Spectral Norm means largest singular value)

$$2) \|X\|_F = \sqrt{\sum_i \sum_j |X_{ij}|^2}$$

considering SVD :  $X = U \Sigma V^T$

$$X^T X = (U \Sigma V^T)^T (U \Sigma V^T) \\ = V \Sigma^2 V^T$$

So, The diagonal elements of  $\Sigma$  is  $\sigma_1$ .

$UU^T = I$ ,  $VV^T = I$  (since they are orthogonal).

$$\therefore \text{Trace}(X^T X) = \text{Trace}(\Sigma^2) = \sigma_1^2 + \sigma_2^2 + \sigma_3^2 + \dots + \sigma_r^2$$

$r$ : rank of  $X$ .

$$\therefore \|X\|_F = \sqrt{\text{Trace}(X^T X)} = \sqrt{\sigma_1^2 + \sigma_2^2 + \dots + \sigma_r^2}$$

$$3) \text{ Now : } \|X\|_F \stackrel{r.f.}{\geq} \|X\|_2$$

$$\sqrt{\sigma_1^2 + \sigma_2^2 + \dots + \sigma_r^2} \geq \sigma_1 \quad \checkmark$$

Taking into consideration, that  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r$

Note: The equality  $\|X\|_F = \|X\|_2$  holds if and only if :

All the singular values are zero except  $\sigma_1$ .

$$\sqrt{\sigma_1^2} \stackrel{r.f.}{\leq} \sigma_1$$

$$\sigma_1 = \sigma_1 \quad \checkmark, \text{ since } \sigma_1 > 0$$

and this happens when the rank of  $X$  is 1.



Problem 2:

2) Fact:  $\text{vec}(AXB) = (B^T \otimes A) \text{vec}(X) \dots \textcircled{1}$

Find:

$$\min_X \|AXB - Y\|_F^2 \dots \textcircled{2}, \text{ where } A \in \mathbb{R}^{m \times p}$$

$$X \in \mathbb{R}^{p \times q}$$

$$B \in \mathbb{R}^{q \times n}$$

$$Y \in \mathbb{R}^{m \times n}$$

Sol: Review Basics:

$$\text{vec}(AXB) \equiv \text{Vectorization}; AXB \equiv [m \times p][p \times q][q \times n] \\ \equiv [m \times n]$$

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

$$\text{ex: } A \otimes B = [m \times p \times q]$$

Frobenius:  $\|A\|_F$ ; It's the same as Euclidean

But for Matrices.

$$\|A\|_F = \sqrt{\sum_i \sum_j |a_{ij}|^2}$$

$$\|A\|_F = \sqrt{\text{Tr}(A^T A)}$$

$$\text{Note: } \text{vec}(A + B) = \text{vec}(A) + \text{vec}(B)$$

Now: Starting with the objective Function:

$$\min_X \|AXB - Y\|_F^2 = \text{Tr}((AXB - Y)^T (AXB - Y)) \dots \textcircled{3}$$

Using Vectorization:

$$\text{vec}(AXB - Y) = \text{vec}(AXB) - \text{vec}(Y) \\ = (B^T \otimes A) \text{vec}(X) - \text{vec}(Y) \dots \textcircled{4}$$

Subs. ④ in ③:

$$\min_X \|(B^T \otimes A) \text{vec}(X) - \text{vec}(Y)\|_F^2$$

Take gradient w.r.t set to zero: From the optimal least square solution

$$\text{vec}(X) = ((B^T \otimes A)^T (B^T \otimes A))^{-1} (B^T \otimes A)^T \text{vec}(Y) \\ = (B B^T \otimes A A^T)^{-1} (B^T \otimes A^T) \text{vec}(Y)$$

since  $x \equiv [m \times n] \rightarrow \text{reshape } \text{vec}(X) \text{ to } X. \#$



### Problem 3: USArrests Dataset

1) PCA Approach: (Max. the variance of the data)

- Standardize the data.
  - SVD
  - $k$  largest eigenvalues.
  - Projection Matrix  $W$  from  $k$ .
  - Transform the original dataset  $X$  via  $W$  to obtain a  $k$ -dimensional feature subspace  $Y$ .
- ⇒ Classes: Target

Features: Data

### 2) Data completion Using Proximal Gradient Descent

$$\frac{1}{2} \|P_n(X - Z)\|_F^2 + \|Z\|_X$$

$P_n$  → Projection Matrix.

(Simulation)

- Mask, Generate a Matrix with missing entries.
- Applies a Proximal Gradient descent Method to find a Matrix  $Z$  that approximates  $X$  while min the obj function.

Results: After Implementation → we got a monotonically decreasing objective Function with iterations.

Problem 5: Logistic Regression: (Binary Classification)  
 $x_i \in \mathbb{R}^{100}$ , label is  $\pm 1$ , then the obj. is:

$$\min_w \sum_{i=1}^m \log(1 + \exp(-y_i w^T x_i))$$

while if the label is  $\{1, 0\}$ , the obj. is:

$$\min_w \sum_{i=1}^m \log(1 + \exp(w^T x_i)) - y_i w^T x_i$$

1) Gradient Descent Method:

$m = 100$

$$\begin{aligned} \text{Step Size} &= \frac{1}{L} = \frac{1}{\text{Largest Sing. value}} = \frac{1}{\frac{1}{4} \|X X^T\|_2} \\ &= \frac{4}{\|X\|_2^2} \quad \text{*(Lipshitz)}. \end{aligned}$$

→ Apply both obj. functions to prove they are the same.

Result: The values and the chart shows us that the two objectives are exactly the same.

2) When it's  $\{1, 0\}$ ,  $P(y=1 | x, w) = \frac{1}{1 + \exp(-w^T x)}$

Then the max. likelihood function is:  $\prod_{i=1}^m P^{y_i} (1-P)^{1-y_i}$

→ Prove  $P^* = y$ . If we use MSE:  $\min_P (y-P)^2$ ,  $y=1$   
 $w=0$  in  $P$  very close to 0. Then if we optimize using GDM, ??

Sol. :: Proof: As we fit 1 class:

$$P(y|x, w)$$

$$\rightarrow \log P(y|x, w) = y \log(P) + (1-y) \log(1-P)$$

obj.  $\rightarrow \max_P y \log(P) + (1-y) \log(1-P)$

$$\rightarrow \min_P -y \log(P) - (1-y) \log(1-P)$$

Der.  $\rightarrow \frac{\partial}{\partial P} (-y \log(P) - (1-y) \log(1-P)) = \frac{-y}{P} + \frac{1-y}{1-P}$



### Problem 5: Section 2: Completing:

der. so  $\rightarrow \frac{-y}{P} + \frac{1-y}{1+P} = 0$

$\rightarrow y(1+P) = (1-y)P$

$\therefore y = P$  ~~\*~~ So, the optimal  $P^* = y$  ~~\*~~

Now: MSE

$J(P) = \frac{1}{2} (y-P)^2$

der.  $\rightarrow \frac{\partial}{\partial P} \left( \frac{1}{2} (y-P)^2 \right) = -(y-P)$

GDM  $\rightarrow P = P - \alpha \frac{\partial}{\partial P}$

$\therefore P = P + \alpha (y-P)$

When:  $y=1, P=0$

update  $\rightarrow P = P + \alpha \cdot 1$

with this large gradient, it's unproper so most likely it'll diverge instead if the step size wasn't set carefully.

This happens because MSE can't handle the small values of  $P$  ~~at~~ ( $P=0$ ) unlike cross entropy as we seen before.

### Section 3: Newton Method:

After implementing, it shows that Newton's is much faster and the convergence rate is way quicker. This is because it uses both gradient and Hessian. ~~\*~~

### Section 4: SGD:

After implementing, it shows that SGD is much faster but with slower convergence rate.

### Problem 6: Kernel SVM.

1)  $k(i, j) = \frac{\langle x_i, x_j \rangle}{\|x_i\|_2 \|x_j\|_2}$ , Then  $k$  defines a Proper kernel.

Sol.:

In order to be a proper kernel, it should satisfy two conditions: (PSD)

1) Cosine Similarity:

2) Positive Definite Matrix:

$$k(i, j) = \frac{\langle x_i, x_j \rangle}{\|x_i\|_2 \|x_j\|_2} = \cos(\theta)$$

from Definition of dot product.  $\forall w \in \mathbb{R}^n, w^T k w \geq 0$

so, First condition Holds

$$w^T \langle \tilde{x}_i, \tilde{x}_j \rangle w$$

$$= \left\langle \sum_{i=1}^n w \tilde{x}_i, \sum_{j=1}^n w \tilde{x}_j \right\rangle$$

$\geq 0$  for any vector

$\therefore$  PSD  $\#$

$\therefore k(i, j)$  is a Proper kernel  $\#$ .

2) Using Python (rbf-kernel, linear-kernel), the # of support vectors have been counted for each class.

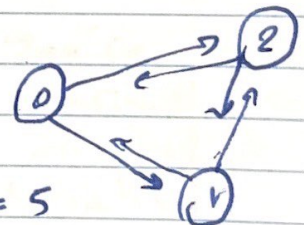
3 classes  $\rightarrow$  means 3 Pairs  
pairs:  $[(0, 1), (0, 2), (1, 2)]$

Conclusion: Support vectors:

pair (0, 1):  $[4 \ 1] \rightarrow$  support vectors = 5

pair (0, 2):  $[4 \ 2] \rightarrow$  " " = 6

pair (1, 2):  $[10 \ 9] \rightarrow$  " " = 19





Thank you so much!!

Q.7.:

1) My favourite book: Rich Dad Poor Dad.  
This book has taught me a lot Business lessons and how to live a good life. It's so insightful and enriching.

My favorite travel destination: Lebanon.  
I have done an internship there Two years ago. It's absolutely the best country ever you'd visit. Kind people, Diversity, Gorgeous Places.

2) My favorite person: My father, he's absolutely my role model. His lessons during my entire life have been a great footprint on my professional and social life. I will never forget his favors.

3) My favorite Restaurant: Texas Roadhouse.  
It's authentic, gorgeous place to visit, kind staff, and very delicious food.

My favourite orders: Fried Pickles and Bone-In Ribeye.

4) My favorite Part: Is how you're so ~~enriching~~ energetic, always positive and sharing the positivity with us. How you discuss different interesting topics with us. Additionally, for Academic perspectives, how you teach us both math aspects and coding for ML.

5) My favorite Algorithm is: SVM and kernel SVM.  
It's like magical thing, how it projects into low dimension space. I really loved it. ✕.

Thank you.