

Assignment Number: 3

Student Name: Raktim Mitra

Roll Number: 150562

Date: November 15, 2017

**Part 1:**

We want to prove  $\theta^{MLE} \in \arg \max_{\theta \in \Theta} Q_{\theta^{MLE}}(\theta)$  with  $\theta^{MLE} \in \arg \max_{\theta \in \Theta} P[X|\Theta]$ .

$$Q_{\theta^t}(\theta) = \sum_{i=1}^n \mathbb{E}_{\mathbf{z} \sim \mathbb{P}[\mathbf{z}|\mathbf{x}^i, \theta^t]} \log \mathbb{P}[\mathbf{x}^i, \mathbf{z}|\theta]$$

We can write: (by deviding with a term independent of  $\theta$ )

$$\begin{aligned} & \arg \max_{\theta \in \Theta} \sum_{i=1}^n \mathbb{E}_{\mathbf{z} \sim \mathbb{P}[\mathbf{z}|\mathbf{x}^i, \theta^t]} \log \mathbb{P}[\mathbf{x}^i, \mathbf{z}|\theta] \\ &= \arg \max_{\theta \in \Theta} \sum_{i=1}^n \mathbb{E}_{\mathbf{z} \sim \mathbb{P}[\mathbf{z}|\mathbf{x}^i, \theta^t]} \log \frac{\mathbb{P}[\mathbf{x}^i, \mathbf{z}|\theta]}{\mathbb{P}[\mathbf{z}|\mathbf{x}^i, \theta^t]} \end{aligned}$$

Now we can use two results derived in lecture 16:

$$\log P[X|\theta^{MLE}] = \sum_{i=1}^n \mathbb{E}_{\mathbf{z} \sim \mathbb{P}[\mathbf{z}|\mathbf{x}^i, \theta^{MLE}]} \log \frac{\mathbb{P}[\mathbf{x}^i, \mathbf{z}|\theta^{MLE}]}{\mathbb{P}[\mathbf{z}|\mathbf{x}^i, \theta^{MLE}]}$$

and

$$\log P[X|\theta] \geq \sum_{i=1}^n \mathbb{E}_{\mathbf{z} \sim \mathbb{P}[\mathbf{z}|\mathbf{x}^i, \theta^{MLE}]} \log \frac{\mathbb{P}[\mathbf{x}^i, \mathbf{z}|\theta]}{\mathbb{P}[\mathbf{z}|\mathbf{x}^i, \theta^{MLE}]} \quad \forall \theta$$

Since,  $\theta^{MLE}$  maximises  $\log \mathbb{P}[\mathbf{x}|\theta]$

$$\log P[X|\theta^{MLE}] \geq \sum_{i=1}^n \mathbb{E}_{\mathbf{z} \sim \mathbb{P}[\mathbf{z}|\mathbf{x}^i, \theta^{MLE}]} \log \frac{\mathbb{P}[\mathbf{x}^i, \mathbf{z}|\theta]}{\mathbb{P}[\mathbf{z}|\mathbf{x}^i, \theta^{MLE}]} \quad \forall \theta$$

i.e.

$$\log P[X|\theta^{MLE}] \geq Q_{\theta^t}(\theta)$$

Therefore,  $\theta^{MLE}$  maximises  $\geq Q_{\theta^t}(\theta)$  too. i.e.  $\theta^{MLE} \in \arg \max_{\theta \in \Theta} P[X|\Theta]$ . [Proved]

Assignment Number: 3

Student Name: Raktim Mitra

Roll Number: 150562

Date: November 15, 2017

**Part 1:**

For some partition  $\{\Omega_i\}_{i=1,\dots,n}$  of  $\mathbb{R}^d$  and weights  $\{\mathbf{w}^i\}_{i=1,\dots,n}$ ,  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is such that  $\forall \mathbf{x} \in \mathbb{R}^d$

$$f(\mathbf{x}) = \sum_{i=1}^n \mathbb{I}\{\mathbf{x} \in \Omega_i\} \cdot \langle \mathbf{w}^i, \mathbf{x} \rangle$$

The function  $g(\mathbf{x}) = c \cdot f(\mathbf{x})$  can be written as:

$$\begin{aligned} g(\mathbf{x}) &= c \sum_{i=1}^n \mathbb{I}\{\mathbf{x} \in \Omega_i\} \cdot \langle \mathbf{w}^i, \mathbf{x} \rangle \\ &= \sum_{i=1}^n \mathbb{I}\{\mathbf{x} \in \Omega_i\} \cdot \langle c\mathbf{w}^i, \mathbf{x} \rangle \end{aligned}$$

Thus for same partition  $\{\Omega_i\}_{i=1,\dots,n}$  of  $\mathbb{R}^d$  and weights  $\{\mathbf{w}'^i\}_{i=1,\dots,n}$ , where  $\mathbf{w}'^i = c\mathbf{w}^i \forall i$   $\forall \mathbf{x} \in \mathbb{R}^d$   $g : \mathbb{R}^d \rightarrow \mathbb{R}$  can be written as :

$$g(\mathbf{x}) = \sum_{i=1}^n \mathbb{I}\{\mathbf{x} \in \Omega_i\} \cdot \langle \mathbf{w}'^i, \mathbf{x} \rangle$$

which means  $g(x)$  is piecewise linear. [Proved]

**Part 2:**

Let, For some partition  $\{\Omega_i^f\}_{i=1,\dots,n}$  of  $\mathbb{R}^d$  and weights  $\{\mathbf{w}_f^i\}_{i=1,\dots,n}$ ,  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is such that  $\forall \mathbf{x} \in \mathbb{R}^d$

$$f(\mathbf{x}) = \sum_{i=1}^n \mathbb{I}\{\mathbf{x} \in \Omega_i^f\} \cdot \langle \mathbf{w}_f^i, \mathbf{x} \rangle$$

and Let, For some partition  $\{\Omega_i^g\}_{i=1,\dots,n}$  of  $\mathbb{R}^d$  and weights  $\{\mathbf{w}_g^i\}_{i=1,\dots,n}$ ,  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  is such that  $\forall \mathbf{x} \in \mathbb{R}^d$

$$g(\mathbf{x}) = \sum_{i=1}^n \mathbb{I}\{\mathbf{x} \in \Omega_i^g\} \cdot \langle \mathbf{w}_g^i, \mathbf{x} \rangle$$

for  $\mathbf{x} \in \mathbb{R}^d$  let,  $h(\mathbf{x}) = f(\mathbf{x}) + g(\mathbf{x})$  the, since  $\mathbf{x}$  can belong to only one of  $\Omega_i^f$ s and one of  $\Omega_i^g$ s, we have to add only two terms. Say,  $\mathbf{x} \in \Omega_k^f$  and  $\mathbf{x} \in \Omega_l^g$  with weights  $w_f^k, w_g^l$  respectively.

$$\begin{aligned} \text{hence, } h(\mathbf{x}) &= \langle w_f^k, \mathbf{x} \rangle + \langle w_g^l, \mathbf{x} \rangle \\ &= \langle w_f^k + w_g^l, \mathbf{x} \rangle \end{aligned}$$

Now we construct a new partition  $\Omega^h$  for  $h(\mathbf{x})$  using above observation:

For each partition  $\Omega_i^f$  let,  $S_i$  be the set of  $\Omega_j^g$ s such that  $\Omega_i^f \cap \Omega_j^g \neq \emptyset$ . For each element  $\Omega_j^g$  of  $S_i$  add set  $\Omega_i^f \cap \Omega_j^g$  as  $\Omega_{ij}^h$  to  $\Omega^h$ .

By the construction and by the disjointness of each partition subsets,  $\mathbf{x} \in \Omega_{ij}^h \leftrightarrow \mathbf{x} \in \Omega_i^f \& \mathbf{x} \in \Omega_j^g$ . hence we can define new weights  $\mathbf{w}_h^{ij} = \mathbf{w}_f^i + \mathbf{w}_g^j$  for all  $\Omega_{ij}^h \in \Omega^h$ .

Hence, for partition  $\Omega^h$ , weights  $\{\mathbf{w}_h^{ij}\}$ s for  $\mathbf{x} \in \mathbb{R}^d$ ,  $h(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}$

$$h(\mathbf{x}) = \sum_{\Omega_{ij}^h \in \Omega^h} \mathbb{I}\{\mathbf{x} \in \Omega_{ij}^h\} \cdot \langle \mathbf{w}_h^{ij}, \mathbf{x} \rangle$$

Hence.  $h(\mathbf{x})$  is piecewise linear. [Proved]

### Part 3:

Let, For some partition  $\{\Omega_i\}_{i=1,\dots,n}$  of  $\mathbb{R}^d$  and weights  $\{\mathbf{w}^i\}_{i=1,\dots,n}$ ,  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is such that  $\forall \mathbf{x} \in \mathbb{R}^d$

$$f(\mathbf{x}) = \sum_{i=1}^n \mathbb{I}\{\mathbf{x} \in \Omega_i\} \cdot \langle \mathbf{w}^i, \mathbf{x} \rangle$$

now,  $g(\mathbf{x}) = f_{ReLU}(f(\mathbf{x})) = \max(f(\mathbf{x}), 0)$ , then for some  $\mathbf{x} \in \mathbb{R}^d$ :

$$\begin{aligned} g(\mathbf{x}) &= \sum_{i=1}^n \mathbb{I}\{\mathbf{x} \in \Omega_i\} \cdot \max(\langle \mathbf{w}^i, \mathbf{x} \rangle, 0) \\ &= \langle \mathbf{w}^k, \mathbf{x} \rangle \text{ when } \langle \mathbf{w}^k, \mathbf{x} \rangle > 0 \text{ and } = 0 \text{ when } \langle \mathbf{w}^k, \mathbf{x} \rangle \leq 0, \mathbf{x} \in \Omega_k \end{aligned}$$

From above observation we can say that by applying ReLU we are basically passing a plane  $\langle \mathbf{w}^k, \mathbf{x} \rangle = 0$  through  $\Omega_k$  and assigning 0 to  $h(\mathbf{x})$  for  $\mathbf{x}$  lying on  $\leq 0$  side of the plane. In other words we can make  $w^k = 0$  in those regions.

Therefore we create a new partition  $\Omega'$  with  $2n$  subsets  $\{\Omega_{i1}, \Omega_{i2}\}_{i=1,\dots,n}$  where  $\Omega_{i1}$  corresponds to set of  $\mathbf{x}$ 's such that  $\langle \mathbf{w}^i, \mathbf{x} \rangle > 0$  and we define  $\mathbf{w}'^{i1} = \mathbf{w}^i$ ,  $\Omega_{i2}$  corresponds to set of  $\mathbf{x}$ 's such that  $\langle \mathbf{w}^i, \mathbf{x} \rangle \leq 0$  and we define  $\mathbf{w}'^{i2} = 0$ . Then, we represent these new weights as  $\{\mathbf{w}'^i\}_{i=1,2,\dots,2n}$  i.e.  $\mathbf{w}^{jk} = \mathbf{w}'^{j*k}$ . and also partition  $\Omega'$  as  $\{\Omega'_i\}_{i=1,2,\dots,2n}$  i.e.  $\Omega_{jk} = \Omega'_{j*k}$

Now, for partition  $\Omega'$  and weights  $\{\mathbf{w}'^i\}$   $g(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}$  can be represented as:

$$g(\mathbf{x}) = \sum_{i=1}^{2n} \mathbb{I}\{\mathbf{x} \in \Omega'_i\} \cdot \langle \mathbf{w}'^i, \mathbf{x} \rangle$$

Hence,  $g(x) = f_{ReLU}(f(\mathbf{x}))$  is piecewise linear. [Proved]

### Part 4:

**Claim:** Any neural network with a ReLU activation function computes a piecewise linear function.

**Proof by Induction:**

#### *Hypothesis*

In a network with only ReLU activation function, if all neuron's in  $i^{th}$  layer outputs a piecewise linear activation function then so does  $i + 1^{th}$  layer of the network.

#### *Base Case*

The input nodes may or may not have ReLU activation.

If, input nodes do not have ReLU activation then we have linear functions being output from each of them (they basically output  $\mathbf{x}_i$  each) i.e. they are piecewise linear too.

Else, they output  $f_{ReLU}$  on a linear function which is again piecewise linear by part 3.

#### *Induction Step*

Each neuron in  $i + 1^{th}$  layer outputs  $f_{ReLU}()$  applied on weighted sum of activation outputs of neuron's in previous layer it is connected to.

let, a neuron in  $i + 1^{th}$  layer is connected to  $m$  neuron's in  $i^{th}$  layer. Let, activation outputs by them are  $\{f_i\}_{i=1,\dots,m}$  and each connection edge has  $\{\mathbf{w}^i\}_{i=1,\dots,m}$ . The neuron in  $i + 1^{th}$  layer therefore outputs  $f_{ReLU}(\sum_{i=1}^m w^i \cdot f_i)$

Now, since each  $f_i$  is piecewise linear, then by part 1 of the problem, so is  $w^i \cdot f_i$ . (constant multiplication)

Now, by part 2 of the problem,  $\sum_{i=1}^m w^i \cdot f_i$  is also piecewise linear (sum of piecewise linear functions)

Now, by part 3 of problem,  $f_{ReLU}(\sum_{i=1}^m w^i \cdot f_i)$  is piecewise linear, since  $\sum_{i=1}^m w^i \cdot f_i$  is piecewise linear.

Therefore the neuron in  $i + 1^{th}$  layer outputs piecewise linear function. This applies for all neurons in  $i + 1^{th}$  layer.

Hence, by induction the claim holds. [Proved]

## Part 5:

- The network has  $d$  input nodes,  $D$  hidden layer nodes on a single hidden layer with ReLU activation and one output node with ReLU activation.
- Each hidden layer node gets a input  $\epsilon \mathbb{R}^d$ . It outputs a piecewise linear function with 2 pieces(uses one ReLU on a linear function). So, basically each hidden node outputs a function which has two pieces basically signifying a hyperplane on  $\mathbb{R}^d$ .
- The piecewise linear function that is sent to output node as input is weighted sum of outputs by the hidden nodes. Weight multiplication does not change partition. The sum however, changes the partition and total number of partition now can be upper bounded by maximum number of regions  $D$   $(d-1)$  dimensional hyperplanes can create in a  $d$  dimensional space. Let  $L_d^D$  be the number of regions. This is given by the recursion

$$L_d^D = L_d^{D-1} + L_{d-1}^{D-1}$$

We can solve it to get

$$L_d^D = \sum_{i=0}^d \binom{D}{i}$$

. Therefore input to output node is a piecewise linear function with at most  $L_d^D$  pieces.

- The output node applies another ReLU on it. Which may at max divide each piece into 2 pieces. Therefore, maximum number of pieces possible becomes

$$2L_d^D = 2 \sum_{i=0}^d \binom{D}{i}$$

*Student Name:* Raktim Mitra

*Date:* November 15, 2017

**Input:** Receive a data point  $(x^t, y^t) \in \chi$  at each timestamp.

- 2: We store older  $(\mathbf{x}^i, y^i)$ s to calculate update for incoming points.

4: At each timestep  $t$ :

$$6: \quad \hat{y} = \text{sgn} \sum_{j=1}^t \alpha_j y_j K(\mathbf{x}^j, \mathbf{x}^t)$$
8:  $\alpha_i \leftarrow \alpha_i + 1$ 10: **end while**

11: predict  $y_{test}$  for a testpoint  $x_{test}$  by,  $\hat{y} = \text{sgn} \sum_{j=1}^t \alpha_j y_j K(\mathbf{x}^t, \mathbf{x}^j)$   
 //from dual version of the perceptron

Assignment Number: 3

Student Name: Raktim Mitra

Roll Number: 150562

Date: November 15, 2017

**Part 1:**

$$K(\mathbf{z}^1, \mathbf{z}^2) = (\langle \mathbf{z}^1, \mathbf{z}^2 \rangle + 1)^2 = (x_1x_2 + y_1y_2 + 1)^2 = x_1^2x_2^2 + y_1^2y_2^2 + 2x_1x_2y_1y_2 + 2x_1x_2 + 2y_1y_2 + 1$$

We want to represent  $K(\mathbf{z}^1, \mathbf{z}^2) = \langle \phi(\mathbf{z}^1), \phi(\mathbf{z}^2) \rangle$  we take  $\mathcal{H}_K \equiv \mathbb{R}^6$  because  $K(\mathbf{z}^1, \mathbf{z}^2)$  has 6 terms, to represent it as a scalar product, we need vectors with 6 dimensions.

Thus our construction follows from the expression :

$$\phi(\mathbf{z}) = \phi(x, y) = [x_1^2, y_1^2, \sqrt{2}xy, \sqrt{2}x, \sqrt{2}y, 1]$$

Clearly,  $K(\mathbf{z}^1, \mathbf{z}^2) = \langle \phi(\mathbf{z}^1), \phi(\mathbf{z}^2) \rangle$ . Hence K is Mercer. [Proved]

**Part 2:**

For given  $f_{(A, \mathbf{b}, c)} = \langle \mathbf{z}, A\mathbf{z} \rangle + \langle \mathbf{b}, \mathbf{z} \rangle + c$  where

$$A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$$

and

$$\mathbf{b} = \begin{bmatrix} b_1 & b_2 \end{bmatrix}^T$$

Then,  $f_{(A, \mathbf{b}, c)} = a_{11}x^2 + (a_{21} + a_{12})xy + a_{22}y^2 + b_1x + b_2y + c$

Let,  $\mathbf{w} = [w_1, w_2, w_3, w_4, w_5, w_6]^T$ , then:

$$\langle \mathbf{w}, \phi(\mathbf{z}) \rangle = w_1x^2 + w_2y^2 + \sqrt{2}w_3xy + \sqrt{2}w_4x + \sqrt{2}w_5y + w_6$$

comparing  $f_{(A, \mathbf{b}, c)} = \langle \mathbf{w}, \phi(\mathbf{z}) \rangle$

we have  $w_1 = a_{11}, w_2 = a_{22}, w_3 = \frac{a_{21}+a_{12}}{\sqrt{2}}, w_4 = \frac{b_1}{\sqrt{2}}, w_5 = \frac{b_2}{\sqrt{2}}, w_6 = c$

Therefore, required  $\mathbf{w}$  is  $[a_{11}, a_{22}, \frac{a_{21}+a_{12}}{\sqrt{2}}, \frac{b_1}{\sqrt{2}}, \frac{b_2}{\sqrt{2}}, c]^T$  [Answer]

**Part 3:**

For given  $\mathbf{w} = [w_1, w_2, w_3, w_4, w_5, w_6]^T$ , we know:

$$\langle \mathbf{w}, \phi(\mathbf{z}) \rangle = w_1x^2 + w_2y^2 + \sqrt{2}w_3xy + \sqrt{2}w_4x + \sqrt{2}w_5y + w_6$$

We want to construct  $(A, \mathbf{b}, c) \in \mathbb{R}^{2 \times 2} \times \mathbb{R}^2 \times \mathbb{R}$  where,

$$A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$$

and

$$\mathbf{b} = \begin{bmatrix} b_1 & b_2 \end{bmatrix}^T$$

and  $f_{(A,\mathbf{b},c)} = a_{11}x^2 + (a_{21} + a_{12})xy + a_{22}y^2 + b_1x + b_2y + c$

again comparing  $f_{(A,\mathbf{b},c)} = \langle \mathbf{w}, \phi(\mathbf{z}) \rangle$   
 we have  $a_{11} = w_1, a_{22} = w_2, b_1 = \sqrt{2}w_4, b_2 = \sqrt{2}w_5, c = w_6$ , and as can be seen  $w_3 = \frac{a_{21}+a_{12}}{\sqrt{2}}$ ,  
 therefore  $a_{21}, a_{12}$  can have infinitely possible values such that the relation with  $w_3$  holds. To  
 make things symmetrical we can take  $a_{21} = a_{12} = \frac{\sqrt{2}w_3}{2} = \frac{w_3}{\sqrt{2}}$

Hence required construction :

$$A = \begin{bmatrix} w_1 & \frac{w_3}{\sqrt{2}} \\ \frac{w_3}{\sqrt{2}} & w_2 \end{bmatrix}$$

and

$$\mathbf{b} = \begin{bmatrix} \sqrt{2}w_4 & \sqrt{2}w_5 \end{bmatrix}^T$$

and

$$c = w_6$$

#### Part 4:

- Let,  $f_{A,b,c}$  minimises least square error on the right hand side. By, part 2 of the problem it can be converted to  $\langle w_r, \phi(x) \rangle$ .
- The solution for Kernel Ridge regression with regularization  $\lambda$  let the corresponding weight be  $w_l$ . By part 3,  $w_l$  corresponds to a quadratic function too. Let,  $L_l$  denote the total error on left side(kernel ridge side).

We can see :

$$\begin{aligned} ||L_l - L_r|| &= \sum_{j=1}^n (y^j - \langle w_l, \phi(\mathbf{z}^j) \rangle)^2 - (y^j - \langle w_r, \phi(\mathbf{z}^j) \rangle)^2 \\ &= \sum_{j=1}^n (2y^j - \langle w_l + w_r, \phi(\mathbf{z}^j) \rangle) (\langle w_r - w_l, \phi(\mathbf{z}^j) \rangle) \\ &\leq ||w_l - w_r|| \sum_{j=1}^n (2y^j - \langle w_l + w_r, \phi(\mathbf{z}^j) \rangle) ||\phi(\mathbf{z}^j)|| \end{aligned}$$

Let,  $X \in \mathbb{R}^{n \times 6}$  and  $y \in \mathbb{R}^n$  contains our  $\phi$  mapped input vectors.

We have known results for  $w_l$  and  $w_r$  :

$$w_r = (X^T X)^{-1} X^T y$$

and

$$w_l = (X^T X + \lambda I)^{-1} X^T y$$

We can do this because  $\phi$  is not very high dimensional.

$X^T X$  is a psd matrix and can be written as  $U A U^T$  where  $A$  is diagonal matrix with  $\geq$  entries  $(a_1, a_2, \dots)$ . then we have

$$\begin{aligned} w_l - w_r &= U(\text{diag}(1/a_i - 1/(a_i + \lambda))) U^T X^T y \\ &\leq \lambda U(\text{diag}(1/a_i^2)) U^T X^T y \end{aligned}$$

i.e.  $||w_l - w_r|| \rightarrow 0$  as  $\lambda \rightarrow 0$  i.e. Difference between two losses tends to 0. [Proved]

Assignment Number: 3

Student Name: Raktim Mitra

Roll Number: 150562

Date: November 15, 2017

$P[\mathbf{X} | \boldsymbol{\mu}, \mathbf{W}, \sigma] = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \mathbf{C})$  where  $\mathbf{C} = \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}_d$  The data log-likelihood is given as:

$$\begin{aligned} L(\boldsymbol{\Theta}) &= \log \prod_{i=1}^n P(\mathbf{X}^i | \boldsymbol{\mu}, \mathbf{W}, \sigma) \\ &= \sum_{i=1}^n \log P(\mathbf{X}^i | \boldsymbol{\mu}, \mathbf{W}, \sigma) \\ &= \sum_{i=1}^n \log \frac{1}{(2\pi)^{\frac{d}{2}} |\mathbf{C}|^{\frac{1}{2}}} \exp\left(-\frac{(\mathbf{X}^i - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{X}^i - \boldsymbol{\mu})}{2}\right) \\ &= \sum_{i=1}^n -\frac{d}{2} \log 2\pi - \frac{1}{2} |\mathbf{C}| - \frac{(\mathbf{X}^i - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{X}^i - \boldsymbol{\mu})}{2} \\ &= -\frac{nd}{2} \log 2\pi - \frac{n}{2} |\mathbf{C}| - \sum_{i=1}^n \frac{(\mathbf{X}^i - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{X}^i - \boldsymbol{\mu})}{2} \end{aligned}$$

Differentiating w.r.t  $\boldsymbol{\mu}$  using  $\frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = \mathbf{x}^T \mathbf{A} + \mathbf{x}^T \mathbf{A}^T$

$$\begin{aligned} 0 &= \sum_{i=1}^n (\mathbf{X}^i - \boldsymbol{\mu})^T (\mathbf{C}^{-1} + (\mathbf{C}^{-1})^T) \\ \implies 0 &= (\mathbf{C}^{-1} + (\mathbf{C}^{-1})^T) \sum_{i=1}^n (\mathbf{X}^i)^T - n \boldsymbol{\mu}^T \end{aligned}$$

multiplying by  $((\mathbf{C}^{-1} + (\mathbf{C}^{-1})^T))^{-1}$  since  $\mathbf{C}$  is invertible

$$\begin{aligned} \implies \boldsymbol{\mu}^T &= \frac{\sum_{i=1}^n (\mathbf{X}^i)^T}{n} \\ \implies \boldsymbol{\mu} &= \frac{\sum_{i=1}^n \mathbf{X}^i}{n} \end{aligned}$$

i.e. MLE estimate for  $\boldsymbol{\mu}$  is  $\boldsymbol{\mu}^{MLE} = \frac{\sum_{i=1}^n \mathbf{X}^i}{n}$