

Assignment Number: 3
Student Name: Gurpreet Singh
Roll Number: 150259
Date: November 15, 2017

Part 1

From the definition of θ^{MLE} , we have a MLE estimate θ such that for any $\theta \in \Theta$ and $\theta \neq \theta^{MLE}$

$$\mathbb{P}[X | \theta^{MLE}] > \mathbb{P}[X | \theta] \quad (1)$$

Suppose that $\theta^{MLE} \notin \arg \max_{\theta \in \Theta} Q_{\theta^{MLE}}(\theta)$.

From **Lecture 16, Slide 44**, we know that for any $\theta \in \Theta$

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

From equation 1, we can write

$$\log(\mathbb{P}[X | \theta^{MLE}]) > Q_{\theta^1}(\theta^2)$$

However, from **Lecture 16, Slide 42**, we also know

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

Therefore, we can say that for any $\theta \in \Theta$

$$Q_{\theta^{MLE}}(\theta^{MLE}) = Q_{\theta^{MLE}}(\theta)$$

However, this is a contradiction, since we assumed that $\theta^{MLE} \notin \arg \max_{\theta \in \Theta} Q_{\theta^{MLE}}(\theta)$. Hence, the self-consistency property needs to be satisfied.

Part 2

From the question, we know that θ^1 and θ^2 are the optimal MLE solutions. Therefore, we can say

$$\mathbb{P}[X | \theta^1] = \mathbb{P}[X | \theta^2]$$

Since θ^1 is one of the optimal solutions, say there is a chain of iterations such that we reach θ^1 through one of these chains in t^i iterations. Therefore, $\theta^i = \theta^1$. Again, from **Lecture 16, Slide 44**

$$\mathbb{P}[X | \theta^{i+1}] \geq Q_{\theta^1}(\theta^{i+1}) \geq \mathbb{P}[X | \theta^1]$$

However, since θ^1 is a MLE solution, we have $\mathbb{P}[X | \theta^{i+1}] \leq \mathbb{P}[X | \theta^1]$. Hence,

$$\mathbb{P}[X | \theta^1] = Q_{\theta^1}(\theta^{i+1}) = \mathbb{P}[X | \theta^{i+1}]$$

Hence, for any θ^{i+1} that satisfies the above constraint will maximise $Q_{\theta^1}(\theta^{i+1})$, as that will be the globally maximum possible value. Hence, for any θ^{i+1} which is a MLE solution, $Q_{\theta^1}(\theta^{i+1})$ will be maximum. Therefore, we can say

$$\begin{aligned} \{\theta^{i+1}\} &= \arg \max_{\theta \in \Theta} \mathbb{P}[X | \theta] \\ \implies \arg \max_{\theta \in \Theta} \mathbb{P}[X | \theta] &= \arg \max_{\theta \in \Theta} Q_{\theta^1}(\theta) \end{aligned}$$

Since, θ^2 is also a MLE solution, therefore we can say that $\theta^2 \in \arg \max_{\theta \in \Theta} Q_{\theta^1}(\theta)$
With a similar construction for θ^2 , we can say

$$\begin{aligned} \arg \max_{\theta \in \Theta} \mathbb{P}[X | \theta] &= \arg \max_{\theta \in \Theta} Q_{\theta^2}(\theta) \\ \theta^1 &\in \arg \max_{\theta \in \Theta} Q_{\theta^2}(\theta) \end{aligned}$$

Hence, proved

Assignment Number: 3
Student Name: Gurpreet Singh
Roll Number: 150259
Date: November 15, 2017

In order to prove that a function f is piecewise linear, we need to prove that we can find a n -partition $\{\Omega_i\}_{i=1}^n$ and a set $\{\mathbf{w}^i\}_{i=1}^n$ such that

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

We use this to prove the parts in the question.

Part 1

Suppose we have a piecewise linear function f . Then we can write $g(\mathbf{x}) = c \cdot f(\mathbf{x})$ as

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

Hence, we have the set $\{\Omega_i^g\}_{i=1}^n$ and $\{\mathbf{w}_g^i\}_{i=1}^n$ which defines the piecewise linear function, where $\forall i \in [n], \Omega_i^g = \Omega_i$ and $\mathbf{w}_g^i = c \cdot \mathbf{w}^i$.
Hence, g is also piecewise linear

Part 2

Suppose we have two piecewise linear functions f_1 and f_2 . Then we can write $g(\mathbf{x}) = f_1(\mathbf{x}) + f_2(\mathbf{x})$ as

$$\begin{aligned} g(\mathbf{x}) &= f_1(\mathbf{x}) + f_2(\mathbf{x}) \\ &= \sum_{i=1}^{n_1} \mathbb{I}[\mathbf{x} \in \Omega_i^1] \langle \mathbf{w}_1^i, \mathbf{x} \rangle + \sum_{i=1}^{n_2} \mathbb{I}[\mathbf{x} \in \Omega_i^2] \langle \mathbf{w}_2^i, \mathbf{x} \rangle \end{aligned}$$

Now, for g we need to consider all points in $\{\Omega_i^1\}_{i=1}^{n_1} \cup \{\Omega_i^2\}_{i=1}^{n_2}$

Since there can be intersections in these sets, we cannot directly add them. Hence we need to consider all the intersections. Therefore, we can write the complete set as

$$Q = \{\Omega_i^1 \cap \Omega_j^2\}_{i=1, j=1}^{n_1, n_2} \cup \{\Omega_i^1\}_{i=1}^{n_1} \cup \{\Omega_j^2\}_{j=1}^{n_2}$$

However, these are not disjoint. In order to make these sets disjoint, we can remove all intersections, since we are already adding them to the complete set. For this, we define two set quantities

$$\begin{aligned} Q_1 &= \sum_{i=1}^{n_1} \bigcup \Omega_i^1 \\ Q_2 &= \sum_{i=1}^{n_2} \bigcup \Omega_i^2 \end{aligned}$$

Then, we can transform the complete set as

$$Q = \{\Omega_i^1 \cap \Omega_j^2\}_{i=1, j=1}^{n_1, n_2} \cup \{\Omega_i^1 \setminus Q_2\}_{i=1}^{n_1} \cup \{\Omega_j^2 \setminus Q_1\}_{j=1}^{n_2}$$

Since, we are including all intersections, we can say that $Q = \{\Omega_i^1\}_{i=1}^{n_1} \cup \{\Omega_j^2\}_{j=1}^{n_2}$

NOTE I am using the concept that a n-partition is disjoint and not explicitly mentioning that. It is only because of that the above term for Q is of this form

We can define the above set Q as our n-partition set for the function g . Since it is disjoint (from construction), it is a valid partition of sets. Now, we can define the mapping from Ω_i^g to \mathbf{w}_g^i .

$$\mathbf{w}_g^i = \begin{cases} \mathbf{w}_1^j & \Omega_i^g = \Omega_j^1 \setminus Q_2 \\ \mathbf{w}_2^j & \Omega_i^g = \Omega_j^2 \setminus Q_1 \\ \mathbf{w}_1^j + \mathbf{w}_2^k & \Omega_i^g = \Omega_j^1 \cap \Omega_k^2 \end{cases}$$

To prove that this is valid, and represents the same function g , we can see that for any point \mathbf{x}

$$g(x) = \begin{cases} \langle \mathbf{w}_1^j, \mathbf{x} \rangle & \text{if } \mathbf{x} \in \Omega_j^1 \setminus Q_2 \\ \langle \mathbf{w}_2^j, \mathbf{x} \rangle & \text{if } \mathbf{x} \in \Omega_j^2 \setminus Q_1 \\ \langle \mathbf{w}_1^j + \mathbf{w}_2^k, \mathbf{x} \rangle & \text{if } \mathbf{x} \in \Omega_j^1 \cap \Omega_k^2 \end{cases}$$

This is clearly the same as the original defined function $g = f_1 + f_2$. Hence, we have found a representation which defines g and hence it is a piecewise linear function.

Hence, proved

Part 3

Suppose we have a piecewise linear function f and a function g such that $g = f_{ReLU}(f)$. Then we can write

$$\begin{aligned} g(\mathbf{x}) &= f_{ReLU}(f(\mathbf{x})) \\ &= f_{ReLU}\left(\sum_{i=1}^n \mathbb{I}[\mathbf{x} \in \Omega_i] \langle \mathbf{w}, \mathbf{x} \rangle\right) \\ &= \max\left(\sum_{i=1}^n \mathbb{I}[\mathbf{x} \in \Omega_i] \langle \mathbf{w}^i, \mathbf{x} \rangle, 0\right) \end{aligned}$$

We can also represent g as follows

$$g(\mathbf{x}) = \begin{cases} \langle \mathbf{w}^i, \mathbf{x} \rangle & \text{if } \mathbf{x} \in \Omega_i \text{ and } \langle \mathbf{w}^i, \mathbf{x} \rangle > 0 \quad \forall i \in [n] \\ 0 & \text{else} \end{cases}$$

Therefore, we can define regions R^i as

$$R^i = \Omega_i \cap \{\mathbf{x} \mid \langle \mathbf{w}^i, \mathbf{x} \rangle > 0\}$$

Then the function g can be written as

$$\begin{aligned} g(\mathbf{x}) &= \begin{cases} \langle \mathbf{w}^i, \mathbf{x} \rangle & \text{if } \mathbf{x} \in R^i \quad \forall i \in [n] \\ 0 & \text{else} \end{cases} \\ \implies g(\mathbf{x}) &= \sum_{i=1}^n \mathbb{I}[\mathbf{x} \in R^i] \langle \mathbf{w}^i, \mathbf{x} \rangle \end{aligned}$$

However, since $\sum_{i=1}^n \bigcup R^i$ does not cover the whole of \mathbb{R}^d , we introduce another set \bar{R} such that $\bar{R} = \mathbb{R}^d \setminus \sum_{i=1}^n \bigcup R^i$ corresponding to $\bar{\mathbf{w}} = \mathbf{0}$. Therefore g becomes

$$g(\mathbf{x}) = \sum_{i=1}^n \mathbb{I}[\mathbf{x} \in R^i] \langle \mathbf{w}^i, \mathbf{x} \rangle + \mathbb{I}[\mathbf{x} \in \bar{R}] \langle \mathbf{0}, \mathbf{x} \rangle$$

From the above expression, we can say that g is a piecewise linear function.

Hence, Proved

Part 4

We prove this using induction. Let us say that for the r^{th} layer, this is true *i.e.* we can represent the value at every node for the r^{th} as a piecewise linear function. Then we prove that the value at every node in the $(r+1)^{th}$ layer can also be represented as a piecewise linear function.

BASE CASE

Consider a neural network with only an input layer and an output layer. This case is trivial, as we only need to consider the one layer.

Since for the i^{th} node of the output layer, we can define the incoming input as $\langle \mathbf{w}^i, \mathbf{x} \rangle$, we can say that the value at output node i will be simply $f_{ReLU}(\langle \mathbf{w}^i, \mathbf{x} \rangle)$ which is, from part 3 can be claimed as a piecewise linear function

HYPOTHESIS

We assume that the condition is satisfied for the r^{th} layer *i.e.* for every node in the r^{th} layer, the value on that node is from a piecewise linear function

INDUCTION

Suppose that the input for the i^{th} node in the $(r + 1)^{th}$ layer is coming from all the values in the previous layer

For the previous layer, the value on the j^{th} node is given by a piecewise linear function f^j . Therefore, the input on the i^{th} node can be written as a weighted sum of the values of previous layer. Let these weights be defined by \mathbf{w}^j

Since this will be not single value, we represent the input itself to be a function. And the value will be the activation function applied on the input. Hence the value on every node is actually a function, which will be computed on the value of the nodes in the input layer.

Assume that the number of nodes in the r^{th} layer is n and $(r + 1)^{th}$ layer is m

$$I^j = \sum_{i=1}^n w_i^j f^i$$

Since we proved in part 1, that multiplying a constant to a piecewise continuous function is also piecewise continuous, and in part 2 we proved that adding piecewise continuous functions also give piecewise continuous functions, we can say that $\forall j \in [m]$, I^j is piecewise continuous.

Also, in part 3 we proved that $f_{ReLU}(f)$ is piecewise continuous if f is piecewise continuous. Hence, $f_{ReLU}(I^j)$ is piecewise continuous.

Therefore, for all nodes, the neural network computes a piecewise continuous function. Using the base case, hence, the neural network will compute a piecewise continuous function acting upon the input.

Bonus Question

Following the discussion from previous parts, we can say that every node in the hidden layer will have two sets or pieces, one consisting of points that output something greater than zero, and the other being the complement of the first.

Also, as we showed in the case of addition of piecewise linear functions, the number of sets upon adding them is just analysing the number of subsets (of sets) an element belongs to. Since there might be two or more sets which are contained in another set, we can only comment on the upper bound of the number of pieces. Since the number of sets in the previous (hidden) layer is D .

We do not consider the rest complement sets, as they will be part of only those sets which the actual set is not of. We can also represent this as a subset of sets, and if a set does not belong to one set, we say it's complement belongs to that set. Hence, the number of subsets will simply be the powerset of these D subsets.

$$\text{Number of Pieces} \leq 2^D$$

Assignment Number: 3
Student Name: Gurpreet Singh
Roll Number: 150259
Date: November 15, 2017

For the following question, assume that $\mathbf{x} \in \mathbb{R}^d$

Algorithm 1: Kernel Perceptron — Storing all Points

1. Initialize a list structure, storing the head of the list
2. Add a single item $\mathbf{w}^0 \in \mathbb{R}^d$ as an initialization of \mathbf{w} in the original space
3. While we receive a point $t \rightarrow \mathbf{z}^t = (y^t, \mathbf{x}^t)$
 - (i) Initialize $e = 0$
 - (ii) For all points \mathbf{p} in the list, update $e \leftarrow e + K(\mathbf{p}, \mathbf{x}^t)$
 - (iii) If $y^t \cdot e < 0$, add $\eta_t \cdot y^t \cdot \mathbf{x}^t$ to the head of the list

The list of points obtained are our support vectors, and can be used for predictions.

We can also have a subset of the points, instead of saving all points, we save only a subset of the points. We can use a simple array for this, instead of using a list, although there are a tonne of ways to do this.

Algorithm 2: Kernel Perceptron — Storing a subset of Points

1. Initialize an array of k elements, all NULL (k is the number of support points required)
2. Add a single item $\mathbf{w}^0 \in \mathbb{R}^d$ as an initialization of \mathbf{w} in the original space to the array
3. Initialize $head \leftarrow 1$ (Indexing starting from 0)
4. While we receive a point $t \rightarrow \mathbf{z}^t = (y^t, \mathbf{x}^t)$
 - (i) Initialize $e = 0$
 - (ii) For all points \mathbf{p} in the array, update $e \leftarrow e + K(\mathbf{p}, \mathbf{x}^t)$
 - (iii) If $y^t \cdot e < 0$,
update $\text{array}[head] \leftarrow \eta_t \cdot y^t \cdot \mathbf{x}^t$
update $head \leftarrow head + 1$

Assignment Number: 3
Student Name: Gurpreet Singh
Roll Number: 150259
Date: November 15, 2017

Part 1

In order to show that the quadratic kernel is Mercer, we need to find a mapping $\phi : \mathbb{R}^2 \rightarrow \mathcal{H}_K$, such that $K(\mathbf{z}^1, \mathbf{z}^2) = \langle \mathbf{z}^1, \mathbf{z}^2 \rangle$.

In order to do this, we need to change the representation of the quadratic kernel. We can alternatively represent it in the following form

$$\begin{aligned}
 (\langle \mathbf{z}^1, \mathbf{z}^2 \rangle + 1)^2 &= 1 + 2 \cdot \langle \mathbf{z}^1, \mathbf{z}^2 \rangle + (\langle \mathbf{z}^1, \mathbf{z}^2 \rangle)^2 \\
 &= 1 + \left\langle \sqrt{2} \cdot \mathbf{z}^1, \sqrt{2} \cdot \mathbf{z}^2 \right\rangle + \sum_{i=1}^d \sum_{j=1}^d z_i^1 z_j^1 z_i^2 z_j^2 \\
 &= 1 + \sum_{i=1}^d \left(\sqrt{2} \cdot z_i^1 \right) \left(\sqrt{2} \cdot z_i^2 \right) + \sum_{j=1}^d \sum_{k=1}^d (z_j^1 \cdot z_k^1) (z_j^2 \cdot z_k^2) \\
 &= \begin{bmatrix} 1 \\ \sqrt{2} \cdot z_1^1 \\ \sqrt{2} \cdot z_2^1 \\ \vdots \\ \sqrt{2} \cdot z_d^1 \\ z_1^1 \cdot z_1^1 \\ z_1^1 \cdot z_2^1 \\ \vdots \\ z_1^1 \cdot z_d^1 \\ z_2^1 \cdot z_1^1 \\ \vdots \\ z_d^1 \cdot z_d^1 \end{bmatrix}^T \begin{bmatrix} 1 \\ \sqrt{2} \cdot z_1^2 \\ \sqrt{2} \cdot z_2^2 \\ \vdots \\ \sqrt{2} \cdot z_d^2 \\ z_1^2 \cdot z_1^2 \\ z_1^2 \cdot z_2^2 \\ \vdots \\ z_1^2 \cdot z_d^2 \\ z_2^2 \cdot z_1^2 \\ \vdots \\ z_d^2 \cdot z_d^2 \end{bmatrix}
 \end{aligned}$$

The matrix representation is exactly the transformed Hilbert Space we require. Hence, we represent the Hilbert Space \mathcal{H}_K in $d^2 + d + 1 = 7$ dimensions, such that $\phi : \mathbb{R}^2 \rightarrow \mathcal{H}_K$ is defined as

$$\begin{aligned}\phi(\mathbf{z}) &= \begin{bmatrix} \phi_0(\mathbf{z}) & \phi_1(\mathbf{z}) & \phi_2(\mathbf{z}) \end{bmatrix} \\ &\text{such that} \\ \phi(\mathbf{z}) &= [1] \in \mathbb{R}^1 \\ \phi(\mathbf{z}) &= \sqrt{2} \cdot \mathbf{z} \in \mathbb{R}^d \\ \phi(\mathbf{z}) &= \begin{bmatrix} z_1 z_1 & z_1 z_2 & z_2 z_1 & z_2 z_2 \end{bmatrix} \in \mathbb{R}^{d^2}\end{aligned}$$

If $\mathcal{H}_K \equiv \mathbb{R}^D$ then $D = d^2 + d + 1$

Part 2

We have

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

We can alternatively write this as

$$\begin{aligned}f_{(A, \mathbf{b}, c)}(\mathbf{x}) &= c + \langle \mathbf{b}, \mathbf{z} \rangle + \langle \mathbf{z}, A\mathbf{z} \rangle \\ &= c + \sum_{i=1}^d b_i z_i + \sum_{j=1}^d \sum_{k=1}^d A_{jk} z_j z_k \\ &= \begin{bmatrix} c \\ b_1/\sqrt{2} \\ b_2/\sqrt{2} \\ \vdots \\ b_d/\sqrt{2} \\ A_{11} \\ A_{12} \\ \vdots \\ A_{1d} \\ A_{21} \\ \vdots \\ A_{dd} \end{bmatrix}^T \begin{bmatrix} 1 \\ \sqrt{2} \cdot z_1^2 \\ \sqrt{2} \cdot z_2^2 \\ \vdots \\ \sqrt{2} \cdot z_d^2 \\ z_1^2 \cdot z_1^2 \\ z_1^2 \cdot z_2^2 \\ \vdots \\ z_1^2 \cdot z_d^2 \\ z_2^2 \cdot z_1^2 \\ \vdots \\ z_d^2 \cdot z_d^2 \end{bmatrix} \\ &= \langle \mathbf{w}, \phi(\mathbf{z}) \rangle\end{aligned}$$

Hence, we can clearly write the function $f_{(A,\mathbf{b},c)}(\mathbf{x})$ in the form $\langle \mathbf{w}, \mathbf{x} \rangle$ such that \mathbf{w} has the form as given below

$$\mathbf{w} = \begin{bmatrix} c & b_1/\sqrt{2} & b_2/\sqrt{2} & A_{11} & A_{12} & A_{21} & A_{22} \end{bmatrix}^T$$

Part 3

We can proceed similar to the previous part, but build it backwards instead. Assume we have $\mathbf{w} \in \mathcal{H}_K$, then we can write $\langle \mathbf{w}, \phi(\mathbf{z}) \rangle$ as

$$\begin{aligned} \langle \mathbf{w}, \phi(\mathbf{z}) \rangle &= \sum_{i=1}^D w_i \cdot \phi(\mathbf{z})_i \\ &= w_1 + \sum_{i=1}^d w_{(1+i)} \cdot z_i + \sum_{j=1}^d \sum_{k=1}^d w_{(1+j \cdot d+k)} z_j z_k \end{aligned}$$

This is of the same form as a quadratic function. We can further transform this as follows. For the following discussion, we substitute d by 2. Firstly, assume

$$\begin{aligned} w_c &= w_1 \\ \mathbf{w}_b &= \begin{bmatrix} \sqrt{2}w_2 & \sqrt{2}w_3 \end{bmatrix}^T \\ w_A &= \begin{bmatrix} w_4 & w_5 \\ w_6 & w_7 \end{bmatrix} \end{aligned}$$

Using these, we can write the dot product as

$$\langle \mathbf{w}, \phi(\mathbf{z}) \rangle = w_c + \langle \mathbf{w}_b, \mathbf{z} \rangle + \langle \mathbf{z}, w_A \mathbf{z} \rangle$$

Since, this is exactly of the form of a quadratic function. We can say that we have the parameters (A, \mathbf{b}, c) , where

$$\begin{aligned} A &= w_A \\ \mathbf{b} &= \mathbf{w}_b \\ c &= w_c \end{aligned}$$

such that

$$\langle \mathbf{w}, \mathbf{z} \rangle = \langle \mathbf{z}, A\mathbf{z} \rangle + \langle \mathbf{b}, \mathbf{z} \rangle + c$$

Part 4

We need to prove that if the function \hat{f} offers the least squares error, then

$$\sum_{i=1}^n \left(y^i - \hat{f}(\mathbf{z}^i) \right)^2 \leq \min_{A, \mathbf{b}, c} \sum_{i=1}^n \left(y^i - f_{A, \mathbf{b}, c}(\mathbf{z}^i) \right)^2 + \epsilon$$

From part 2, we know that every quadratic function can be represented in the form of the dot product of a parameter \mathbf{w} in the Hilbert Space. Let $\hat{\mathbf{w}}$ represent \hat{f} . Therefore, we can transform the above equation as

$$\sum_{i=1}^n \left(y^i - \langle \hat{\mathbf{w}}, \phi(\mathbf{z}^i) \rangle \right)^2 \leq \min_{\mathbf{w} \in \mathcal{H}_K} \sum_{i=1}^n \left(y^i - \langle \mathbf{w}, \phi(\mathbf{z}^i) \rangle \right)^2 + \epsilon \quad (2)$$

We can relate these to the concept of MLE and MAP solutions. From the case of linear regression, we already know the MLE and MAP solutions.

$$\begin{aligned} \mathbf{w}^{MLE} &= (X^T X)^{-1} X^T \mathbf{y} \\ \mathbf{w}^{MAP} &= (X^T X + \lambda \mathbf{I})^{-1} X^T \mathbf{y} \end{aligned}$$

In the Hilbert Space, we can transform these as follows

$$\begin{aligned} \mathbf{w}^{MLE} &= \left(\phi(X)^T \phi(X) \right)^{-1} \phi(X)^T \mathbf{y} \\ \mathbf{w}^{MAP} &= \left(\phi(X)^T \phi(X) + \lambda \mathbf{I} \right)^{-1} \phi(X)^T \mathbf{y} \end{aligned}$$

NOTE In the Hilbert Space, $\mathbf{w}^{MAP} = \hat{\mathbf{w}}$ (comes directly from the definition of \hat{f})

Also, from the definition of the MLE estimate, which is a minimum of the least squares error, we can say

$$\sum_{i=1}^n \left(y^i - \langle \mathbf{w}^{MLE}, \phi(\mathbf{z}^i) \rangle \right)^2 \leq \min_{\mathbf{w} \in \mathcal{H}_K} \sum_{i=1}^n \left(y^i - \langle \mathbf{w}, \phi(\mathbf{z}^i) \rangle \right)^2$$

Therefore, if we prove that as $\lambda \rightarrow 0^+$, $\hat{\mathbf{w}} \rightarrow \mathbf{w}^{MLE}$, then we can say that $\epsilon \rightarrow 0$, as the LHS in equation ?? will be a proper lower bound.

NOTE We can write $\phi(X)^T \phi(X)$ as the Gram's Matrix G . Consider the forms of \mathbf{w}^{MLE} and $\hat{\mathbf{w}} = \mathbf{w}^{MAP}$, we can write

$$\begin{aligned} \hat{\mathbf{w}} &= (G + \lambda \mathbf{I})^{-1} \phi(X)^T \mathbf{y} \\ &= G^{-1} \left(\mathbf{I} - \frac{\lambda G^{-1}}{1 + \text{trace}(G^{-1})} \right) \phi(X)^T \mathbf{y} \\ &\rightarrow G^{-1} \phi(X)^T \mathbf{y} \quad \text{as } \lambda \rightarrow 0 \\ &= \mathbf{w}^{MLE} \end{aligned}$$

Hence, as $\lambda \rightarrow 0$, $\hat{\mathbf{w}} \rightarrow \mathbf{w}^{MLE}$, hence we can say $\epsilon \rightarrow 0$

Assignment Number: 3
Student Name: Gurpreet Singh
Roll Number: 150259
Date: November 15, 2017

Part 1

We know that Gaussian-Gaussian has the conjugacy property, and hence, we can say that the marginal probability of \mathbf{x} i.e. $\mathbb{P}[\mathbf{x}]$

Therefore, we can write the marginal probability of \mathbf{x} as a gaussian distribution with mean $\boldsymbol{\mu}'$ and covariance matrix Λ .

$$\begin{aligned}\mathbb{P}[\mathbf{x} | \boldsymbol{\mu}, W, \sigma] &= \int_{\mathbf{z}} \mathbb{P}[\mathbf{x} | \mathbf{z}, \boldsymbol{\mu}, W, \sigma] \mathbb{P}[\mathbf{z} | 0, \mathbf{I}] d\mathbf{z} \\ &= \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}', \Lambda)\end{aligned}$$

Also, we know that $\mathbf{x} = W\mathbf{z} + \boldsymbol{\mu} + \epsilon$, and from the properties of Gaussian, we know $\boldsymbol{\mu}' = \mathbb{E}[\mathbf{x}]$ and $\Lambda = \text{cov}(\mathbf{x})$. Hence

$$\begin{aligned}\boldsymbol{\mu}' &= \mathbb{E}[\mathbf{x}] \\ &= \mathbb{E}[W\mathbf{z} + \boldsymbol{\mu} + \epsilon] \\ &= W(\mathbb{E}[\mathbf{z}] = 0) + (\mathbb{E}[\boldsymbol{\mu}] = \boldsymbol{\mu}) + (\mathbb{E}[\epsilon] = 0) \\ &= \boldsymbol{\mu}\end{aligned}$$

$$\begin{aligned}\Lambda' &= \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T] \\ &= \mathbb{E}[(W\mathbf{z} + \epsilon)(W\mathbf{z} + \epsilon)^T] \\ &= \mathbb{E}[W\mathbf{z}\mathbf{z}^T W^T + W\mathbf{z}\epsilon^T + \epsilon W\mathbf{z}^T + \epsilon\epsilon^T] \\ &= W\mathbb{E}[\mathbf{z}\mathbf{z}^T] W^T + 0 + 0 + \mathbb{E}[\epsilon\epsilon^T] \\ &= W\mathbf{I}\mathbf{I}^T W^T + \sigma^2 \mathbf{I} \\ &= WW^T + \sigma^2 \mathbf{I}\end{aligned}$$

Hence, we can write $\mathbb{P}[\mathbf{x}] = \mathcal{N}(\boldsymbol{\mu}, C)$ where $C = WW^T + \sigma^2 \mathbf{I}$

Part 2

Note: I will use $C = WW^T + \sigma^2 \mathbf{I}$ for the rest of the question

Since all the sample points are independent, we can write $\mathbb{P}[X] = \prod_{i=1}^n \mathbb{P}[\mathbf{x}^i]$. Therefore

$$\begin{aligned}
\mathbb{P}[X | \boldsymbol{\mu}, W, \sigma] &= \prod_{i=1}^n \mathbb{P}[\mathbf{x}^i | \boldsymbol{\mu}, W, \sigma] \\
&= \prod_{i=1}^n \mathcal{N}(\mathbf{x}^i | \boldsymbol{\mu}, C) \\
\Rightarrow \mathbb{P}[X | \boldsymbol{\mu}, W, \sigma] &= \prod_{i=1}^n \frac{1}{\sqrt{(2\pi)^d |C|}} e^{-\frac{1}{2}(\mathbf{x}^i - \boldsymbol{\mu})^T C^{-1} (\mathbf{x}^i - \boldsymbol{\mu})} \\
\Rightarrow \log(\mathbb{P}[X | \boldsymbol{\mu}, W, \sigma]) &= -\sum_{i=1}^n \left(\frac{1}{2} (\mathbf{x}^i - \boldsymbol{\mu})^T C^{-1} (\mathbf{x}^i - \boldsymbol{\mu}) + \frac{1}{2} \log((2\pi)^d |C|) \right)
\end{aligned}$$

Part 3

We can write the MLE estimate of $\boldsymbol{\mu}$ as

$$\begin{aligned}
\boldsymbol{\mu}^{MLE} &= \arg \max_{\boldsymbol{\mu} \in \mathbb{R}^d} \mathbb{P}[X | \boldsymbol{\mu}, W, \sigma] \\
&= \arg \max_{\boldsymbol{\mu} \in \mathbb{R}^d} \log(\mathbb{P}[X | \boldsymbol{\mu}, W, \sigma]) \\
&= \arg \max_{\boldsymbol{\mu} \in \mathbb{R}^d} \log \left(\prod_{i=1}^n \mathbb{P}[\mathbf{x}^i | \boldsymbol{\mu}, W, \sigma] \right) \\
&= \arg \max_{\boldsymbol{\mu} \in \mathbb{R}^d} \sum_{i=1}^n \log(\mathbb{P}[\mathbf{x}^i | \boldsymbol{\mu}, W, \sigma]) \\
&= \arg \min_{\boldsymbol{\mu} \in \mathbb{R}^d} \sum_{i=1}^n \left(\frac{1}{2} (\mathbf{x}^i - \boldsymbol{\mu})^T C^{-1} (\mathbf{x}^i - \boldsymbol{\mu}) + \frac{1}{2} \log((2\pi)^d |C|) \right) \\
&= \arg \min_{\boldsymbol{\mu} \in \mathbb{R}^d} \sum_{i=1}^n \frac{1}{2} (\mathbf{x}^i - \boldsymbol{\mu})^T C^{-1} (\mathbf{x}^i - \boldsymbol{\mu})
\end{aligned}$$

We can simply differentiate the RHS term, and using differentials of matrices, we get

$$\begin{aligned}
\frac{\delta(RHS)}{\delta(\boldsymbol{\mu})} &= \frac{\delta \left(\sum_{i=1}^n \frac{1}{2} (\mathbf{x}^i - \boldsymbol{\mu})^T C^{-1} (\mathbf{x}^i - \boldsymbol{\mu}) \right)}{\delta(\boldsymbol{\mu})} = 0 \\
\Rightarrow \sum_{i=1}^n C^{-1} (\mathbf{x}^i - \boldsymbol{\mu}) &= 0 \\
\Rightarrow \boldsymbol{\mu}^{MLE} &= \frac{1}{n} \sum_{i=1}^n \mathbf{x}^i
\end{aligned}$$

Therefore, the MLE estimate of $\boldsymbol{\mu}$ is, as expected, the empirical mean of all the data points. Hence we use this to centralize for the PCA algorithm.