

Assignment Number: 3

Student Name: Siddharth Agrawal

Roll Number: 150716

Date: November 14, 2017

1 Property 1

From the definition of MLE estimate:

$$\log \mathbb{P}[X|\theta^{\text{MLE}}] \geq \log \mathbb{P}[X|\theta] \quad \forall \theta \in \Theta \quad (1)$$

Let

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

From **Lec16 Slide 43**, we know that:

$$\mathcal{Q}_{\theta^t}(\theta^t) = \log \mathbb{P}[X|\theta^t] \quad (2)$$

Using Equation 2 in Equation 1

$$\mathcal{Q}_{\theta^{\text{MLE}}}(\theta^{\text{MLE}}) \geq \log \mathbb{P}[X|\theta] \quad \forall \theta \in \Theta \quad (3)$$

From **Lec16 Slide 44**, we know that:

$$\log \mathbb{P}[X|\theta] \geq \mathcal{Q}_{\theta^t}(\theta) \quad \forall \theta \in \Theta \quad (4)$$

Using Equation 4 in Equation 3

$$\mathcal{Q}_{\theta^{\text{MLE}}}(\theta^{\text{MLE}}) \geq \mathcal{Q}_{\theta^{\text{MLE}}}(\theta) \quad \forall \theta \in \Theta \quad (5)$$

Expanding Equation 5 using the definition of \mathcal{Q} :

$$\begin{aligned} \sum_{i=1}^n \mathbb{E}_{\mathbf{z} \sim \mathbb{P}[\mathbf{z}|\mathbf{x}^i, \theta^{\text{MLE}}]} \log \frac{\mathbb{P}[\mathbf{x}^i, \mathbf{z}|\theta^{\text{MLE}}]}{\mathbb{P}[\mathbf{z}|\mathbf{x}^i, \theta^{\text{MLE}}]} &\geq \sum_{i=1}^n \mathbb{E}_{\mathbf{z} \sim \mathbb{P}[\mathbf{z}|\mathbf{x}^i, \theta^{\text{MLE}}]} \log \frac{\mathbb{P}[\mathbf{x}^i, \mathbf{z}|\theta]}{\mathbb{P}[\mathbf{z}|\mathbf{x}^i, \theta^{\text{MLE}}]} \quad \forall \theta \in \Theta \\ \sum_{i=1}^n \mathbb{E}_{\mathbf{z} \sim \mathbb{P}[\mathbf{z}|\mathbf{x}^i, \theta^{\text{MLE}}]} \log \mathbb{P}[\mathbf{x}^i, \mathbf{z}|\theta^{\text{MLE}}] &\geq \sum_{i=1}^n \mathbb{E}_{\mathbf{z} \sim \mathbb{P}[\mathbf{z}|\mathbf{x}^i, \theta^{\text{MLE}}]} \log \mathbb{P}[\mathbf{x}^i, \mathbf{z}|\theta] \quad \forall \theta \in \Theta \\ \mathcal{Q}_{\theta^{\text{MLE}}}(\theta^{\text{MLE}}) &\geq \mathcal{Q}_{\theta^{\text{MLE}}}(\theta) \quad \forall \theta \in \Theta \end{aligned} \quad (6)$$

Maximizing the RHS in Equation 6 will maximize the LHS.

Hence,

$$\theta^{\text{MLE}} \in \arg \max_{\theta \in \Theta} \mathcal{Q}_{\theta^{\text{MLE}}}(\theta) \quad (7)$$

2 Property 2

Since θ^1 and θ^2 are optimal MLE solutions, hence:

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

Let after t iterations of the EM algorithm, $\theta^t = \theta^1$.

Using the results derived in **Lecture 16, slides-43&44**, we can argue that

$$\mathbb{P}[X|\theta^{t+1}] \geq \mathcal{Q}_{\theta^1}(\theta^{t+1}) \geq \mathbb{P}[X|\theta^1]$$

But, since θ^1 is an MLE solution, hence we have $\mathbb{P}[X|\theta^1] \geq \mathbb{P}[X|\theta^{t+1}]$. Hence, for all the subsequent $\theta^{\tilde{t}}, \tilde{t} > t$, it will maximize $\mathcal{Q}_{\theta^1}(\theta)$ and so will be the MLE estimate.

Remark 1.1. But how do I arrive at the result that this MLE estimate can be converted to θ^2 ?

Assignment Number: 3

Student Name: Siddharth Agrawal

Roll Number: 150716

Date: November 14, 2017

1 Scalar Multiplication

Given: a piecewise linear function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, and a scalar c

To Prove: $g(\mathbf{x}) = c \cdot f(\mathbf{x})$ is also piecewise linear

Proof:

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

Take $\forall i \in [n]$:

$$\begin{aligned} \Omega_i^g &= \Omega_i \\ \mathbf{w}_g^i &= c \cdot \mathbf{w}^i \end{aligned}$$

Now rewrite $g(\mathbf{x})$ as:

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

Hence, g is also a piecewise linear function.

□

2 Addition

Given: 2 piecewise linear functions $f_1 : \mathbb{R}^d \rightarrow \mathbb{R}$, $f_2 : \mathbb{R}^d \rightarrow \mathbb{R}$

To Prove: $g(\mathbf{x}) = f_1(\mathbf{x}) + f_2(\mathbf{x})$ is also piecewise linear

Proof:

$$f_1(\mathbf{x}) = \sum_{i=1}^{n_1} \mathbb{I}\{\mathbf{x} \in \Omega_i^1\} \cdot \langle \mathbf{w}_1^i, \mathbf{x} \rangle \quad (8)$$

$$f_2(\mathbf{x}) = \sum_{i=1}^{n_2} \mathbb{I}\{\mathbf{x} \in \Omega_i^2\} \cdot \langle \mathbf{w}_2^i, \mathbf{x} \rangle \quad (9)$$

Now,

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

Hence,

$$g(\mathbf{x}) = \sum_{i=1, j=1}^{n_1, n_2} \mathbb{I}\{\mathbf{x} \in \Omega_i^1 \cap \Omega_j^2\} \cdot \langle \mathbf{w}_1^i + \mathbf{w}_2^j, \mathbf{x} \rangle \quad (10)$$

Consider the partition:

$$\tilde{\Omega}^g = \bigcup_{i=1, j=1}^{n_1, n_2} (\Omega_i^1 \cap \Omega_j^2)$$

As, Ω_i^1 are disjoint $\forall i \in [n_1]$ and Ω_j^2 are disjoint $\forall j \in [n_2]$, hence $\Omega_i^1 \cap \Omega_j^2$ are disjoint $\forall i \in [n_1], j \in [n_2]$.

Hence g is indexed by $n_1 n_2$ partitions of \mathbb{R}^d (say, $[\Omega_1^g, \Omega_2^g, \dots, \Omega_{n_1 n_2}^g]$).

The linear model corresponding to g is given by:-

$$\forall i \in [n_1 n_2] \quad \mathbf{w}_i^g = \mathbf{w}_j^1 + \mathbf{w}_k^2; \quad \Omega_i^g = \Omega_j^1 \cap \Omega_k^2$$

Hence, g is a piecewise linear function. □

3 ReLU Activation

Given: a piecewise linear function $f : \mathbb{R}^d \rightarrow \mathbb{R}$

To Prove: $g(\mathbf{x}) = f_{\text{ReLU}}(f(\mathbf{x}))$ is also piecewise linear

Proof:

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

Consider the set for which f becomes negative:

$$\Omega_0^g = \left\{ \mathbf{x} \mid \mathbf{x} \in \mathbb{R}^d, f(\mathbf{x}) < 0 \right\}$$

Let its Corresponding linear model be:

$$\mathbf{w}_g^0 = [0, 0, \dots, d \text{ zeroes}]^\top$$

So,

$$\begin{aligned} g(\mathbf{x}) &= \sum_{i=1}^n \mathbb{I}\{\mathbf{x} \in \Omega_i \setminus \Omega_0^g\} \cdot \langle \mathbf{w}^i, \mathbf{x} \rangle + \mathbb{I}\{\mathbf{x} \in \Omega_0^g\} \cdot \langle \mathbf{w}_g^0, \mathbf{x} \rangle \\ &= \sum_{i=0}^n \mathbb{I}\{\mathbf{x} \in \Omega_i^g\} \cdot \langle \mathbf{w}_g^i, \mathbf{x} \rangle \end{aligned}$$

Where, $\forall i \in \{0, 1, 2, \dots, n\}$

$$\begin{aligned} \Omega_i^g &= \begin{cases} \left\{ \mathbf{x} \mid \mathbf{x} \in \mathbb{R}^d, f(\mathbf{x}) < 0 \right\} & i = 0 \\ \Omega_i \setminus \Omega_0^g & \text{Otherwise} \end{cases} \\ \mathbf{w}_g^i &= \begin{cases} [0, 0, \dots, d \text{ zeroes}]^\top & i = 0 \\ \mathbf{w}^i & \text{Otherwise} \end{cases} \end{aligned}$$

Hence, g is also a piecewise linear function. □

4 Neural Networks with ReLU activation function computes a piecewise linear function

We prove this result using induction on the layers of the neural network.

Base Case

The base case will consist of only 1 input and 1 output layers (no hidden layer). Consider any i^{th} node of the output layer. For it, we can define the output as -

$$g^i(\mathbf{x}) = f_{\text{ReLU}}(\langle \mathbf{w}^i, \mathbf{x} \rangle)$$

. Using the result of part 3, since inner product is a piecewise linear function, hence g^i is also piecewise linear.

Induction Hypothesis

If k^{th} layer of the network computes piecewise linear function, then $(k + 1)^{th}$ layer computes piecewise linear function.

Proof for Induction Hypothesis

Let for any i^{th} node in the $(k + 1)^{th}$ layer, inputs come from all nodes in the k^{th} layer. The outputs from k^{th} layer are mapped to any node i in the $(k + 1)^{th}$ layer with weight \mathbf{w}^i . Let f^j be the output of the j^{th} node in this layer. Hence, the output function g^i for the i^{th} node in $(k + 1)^{th}$ layer is given by:

$$g^i = \sum_j \mathbf{w}_j^i \cdot f^j$$

Since scalar multiplication and addition on a piecewise linear function preserves its piecewise linearity (using results proved in part 1 & 2) and $\forall j, f^j$ are piecewise linear (induction hypothesis), hence g^i is piecewise linear. □

Hence, we have proved that any neural network with a ReLU activation function computes a piecewise linear function.

Assignment Number: 3

Student Name: Siddharth Agrawal

Roll Number: 150716

Date: November 14, 2017

For the below given algorithm, input points are being stored in a list instead of storing the whole model $\mathbf{w} \in \mathcal{H}_K$.

Also, it is assumed that the input points are in the form (y^t, \mathbf{x}^t) with $y^t \in \mathbb{R}, \mathbf{x}^t \in \mathbb{R}^d$

Initial value is assumed as β^0

Algorithm 1: Kernel Perceptron

```
1: Store an initial point  $\beta^0 \in \mathbb{R}^d$  in the list.
2: for each new data point  $(y^t, \mathbf{x}^t)$  received do
3:    $\epsilon \leftarrow 0$ 
4:   for each point  $\beta^k$  in the list do
5:      $\epsilon \leftarrow \epsilon + K(\beta^k, \mathbf{x}^t)$ 
6:   end for
7:   if  $y^t \cdot \epsilon < 0$  then
8:     Append  $\alpha_t y^t \mathbf{x}^t$  to the list (as the value for  $\beta^t$ )
9:   end if
10: end for
```

Assignment Number: 3

Student Name: Siddharth Agrawal

Roll Number: 150716

Date: November 14, 2017

1 Proving that Kernel K is Mercer

The construction for the mapping $\varphi : \mathbb{R}^2 \rightarrow \mathcal{H}_K$ is:

$$\begin{aligned}\varphi(\mathbf{z}) &= [\varphi_0(\mathbf{z}), \varphi_1(\mathbf{z}), \varphi_2(\mathbf{z})] \in \mathbb{R}^{2^2+2+1} = \mathbb{R}^7 \\ \text{where,} \\ \varphi_0(\mathbf{z}) &= 1 \in \mathbb{R}^1 \\ \varphi_1(\mathbf{z}) &= \sqrt{2} \cdot [\mathbf{z}_1, \mathbf{z}_2] \in \mathbb{R}^2 \\ \varphi_2(\mathbf{z}) &= [\mathbf{z}_1\mathbf{z}_1, \mathbf{z}_1\mathbf{z}_2, \mathbf{z}_2\mathbf{z}_1, \mathbf{z}_2\mathbf{z}_2] \in \mathbb{R}^4\end{aligned}$$

Remark 4.1. $D = 7. \therefore \mathcal{H}_K \equiv \mathbb{R}^7$

To show that K is a Mercer Kernel, it is sufficient to show that $K(\mathbf{z}^1, \mathbf{z}^2) = \langle \varphi(\mathbf{z}^1), \varphi(\mathbf{z}^2) \rangle$. Hence,

$$\begin{aligned}\langle \varphi(\mathbf{z}^1), \varphi(\mathbf{z}^2) \rangle &= \langle \varphi_0(\mathbf{z}^1), \varphi_0(\mathbf{z}^2) \rangle + \langle \varphi_1(\mathbf{z}^1), \varphi_1(\mathbf{z}^2) \rangle + \langle \varphi_2(\mathbf{z}^1), \varphi_2(\mathbf{z}^2) \rangle \\ &= 1 + 2 \cdot \langle \mathbf{z}^1, \mathbf{z}^2 \rangle + \sum_{i,j}^2 \mathbf{z}_i^1 \mathbf{z}_j^1 \mathbf{z}_i^2 \mathbf{z}_j^2 \\ &= (\langle \mathbf{z}^1, \mathbf{z}^2 \rangle + 1)^2 \\ &= K(\mathbf{z}^1, \mathbf{z}^2)\end{aligned}$$

2 Constructing $\mathbf{w} \in \mathcal{H}_K$

We have to construct \mathbf{w} such that $\forall \mathbf{z} \in \mathbb{R}^2$

$$\begin{aligned}
\langle \mathbf{w}, \varphi_K(\mathbf{z}) \rangle &= f_{(A, \mathbf{b}, c)}(\mathbf{z}) \\
&= \langle \mathbf{z}, A\mathbf{z} \rangle + \langle \mathbf{b}, \mathbf{z} \rangle + c \\
&= A_{11}\mathbf{z}_1\mathbf{z}_1 + A_{12}\mathbf{z}_2\mathbf{z}_1 + A_{21}\mathbf{z}_1\mathbf{z}_2 + A_{22}\mathbf{z}_2\mathbf{z}_2 + \mathbf{b}_1\mathbf{z}_1 + \mathbf{b}_2\mathbf{z}_2 + c \\
&= c + \mathbf{b}_1\mathbf{z}_1 + \mathbf{b}_2\mathbf{z}_2 + A_{11}\mathbf{z}_1\mathbf{z}_1 + A_{12}\mathbf{z}_2\mathbf{z}_1 + A_{21}\mathbf{z}_1\mathbf{z}_2 + A_{22}\mathbf{z}_2\mathbf{z}_2 \\
&= c + \frac{1}{\sqrt{2}} \cdot \mathbf{b}_1\sqrt{2} \cdot \mathbf{z}_1 + \frac{1}{\sqrt{2}} \cdot \mathbf{b}_2\sqrt{2} \cdot \mathbf{z}_2 + A_{11}\mathbf{z}_1\mathbf{z}_1 + A_{12}\mathbf{z}_2\mathbf{z}_1 + A_{21}\mathbf{z}_1\mathbf{z}_2 + A_{22}\mathbf{z}_2\mathbf{z}_2
\end{aligned}$$

From the above derivation, we can clearly see that -

$$\mathbf{w} = \begin{bmatrix} c & \frac{1}{\sqrt{2}} \cdot \mathbf{b}_1 & \frac{1}{\sqrt{2}} \cdot \mathbf{b}_2 & A_{11} & A_{12} & A_{21} & A_{22} \end{bmatrix}^\top$$

3 Constructing a triplet $(A, \mathbf{b}, c) \in \mathbb{R}^{2 \times 2} \times \mathbb{R}^2 \times \mathbb{R}$

We have to construct a triplet (A, \mathbf{b}, c) such that $\forall \mathbf{z} \in \mathbb{R}^2$ given $\mathbf{w} = [\mathbf{w}_1, \dots, \mathbf{w}_7] \in \mathbb{R}^7$

$$\begin{aligned}
 f_{(A, \mathbf{b}, c)}(\mathbf{z}) &= \langle \mathbf{w}, \varphi_K(\mathbf{z}) \rangle \\
 &= \mathbf{w}_1 + \mathbf{w}_2 \cdot \sqrt{2} \cdot \mathbf{z}_1 + \mathbf{w}_3 \cdot \sqrt{2} \cdot \mathbf{z}_2 + \mathbf{w}_4 \mathbf{z}_1 \mathbf{z}_1 + \mathbf{w}_5 \mathbf{z}_1 \mathbf{z}_2 + \mathbf{w}_6 \mathbf{z}_2 \mathbf{z}_1 + \mathbf{w}_7 \mathbf{z}_2 \mathbf{z}_2 \\
 &= \mathbf{w}_4 \mathbf{z}_1 \mathbf{z}_1 + \mathbf{w}_5 \mathbf{z}_1 \mathbf{z}_2 + \mathbf{w}_6 \mathbf{z}_2 \mathbf{z}_1 + \mathbf{w}_7 \mathbf{z}_2 \mathbf{z}_2 + \mathbf{w}_2 \cdot \sqrt{2} \cdot \mathbf{z}_1 + \mathbf{w}_3 \cdot \sqrt{2} \cdot \mathbf{z}_2 + \mathbf{w}_1 \\
 &= \langle \mathbf{z}, A\mathbf{z} \rangle + \langle \mathbf{b}, \mathbf{z} \rangle + c
 \end{aligned}$$

where

$$\begin{aligned}
 A &= \begin{bmatrix} \mathbf{w}_4 & \mathbf{w}_5 \\ \mathbf{w}_6 & \mathbf{w}_7 \end{bmatrix} \\
 \mathbf{b} &= [\mathbf{w}_2, \mathbf{w}_3] \\
 c &= \mathbf{w}_1
 \end{aligned}$$

Assignment Number: 3

Student Name: Siddharth Agrawal

Roll Number: 150716

Date: November 14, 2017

1 Data log-likelihood expression

Since all the data points are independent, hence $\mathbb{P}[X|\boldsymbol{\mu}, W, \sigma] = \prod_{i=1}^n \mathbb{P}[\mathbf{x}^i|\boldsymbol{\mu}, W, \sigma]$. Hence,

$$\begin{aligned}\mathbb{P}[X|\boldsymbol{\mu}, W, \sigma] &= \prod_{i=1}^n \mathbb{P}[\mathbf{x}^i|\boldsymbol{\mu}, W, \sigma] \\ &= \prod_{i=1}^n \mathcal{N}(\mathbf{x}^i|\boldsymbol{\mu}, C) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi \|C\|}} \exp -\frac{1}{2}(\mathbf{x}^i - \boldsymbol{\mu})C^{-1}(\mathbf{x}^i - \boldsymbol{\mu})^\top \\ \implies \log \mathbb{P}[X|\boldsymbol{\mu}, W, \sigma] &= \sum_{i=1}^n -\frac{1}{2}(\mathbf{x}^i - \boldsymbol{\mu})C^{-1}(\mathbf{x}^i - \boldsymbol{\mu})^\top - \frac{1}{2} \log(2\pi \|C\|)\end{aligned}$$

2 Derivation for $\boldsymbol{\mu}^{\text{MLE}}$

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

Now we Differentiate the RHS term w.r.t. $\boldsymbol{\mu}$ so as to get the MLE estimate.

$$\begin{aligned}\frac{\partial \text{RHS}}{\partial \boldsymbol{\mu}} &= \frac{\partial}{\partial \boldsymbol{\mu}} \left(\sum_{i=1}^n -\frac{1}{2} (\mathbf{x}^i - \boldsymbol{\mu}) C^{-1} (\mathbf{x}^i - \boldsymbol{\mu})^\top \right) = 0 \\ &\quad \sum_{i=1}^n -C^{-1} (\mathbf{x}^i - \boldsymbol{\mu}) = 0 \\ &\quad \therefore \boldsymbol{\mu}^{\text{MLE}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^i\end{aligned}$$