

# Concordia University

Computer Science and Software Engineering Department

COMP 6321

Machine Learning

Assignment4

Student ID: 27162391

**Question: 1 (a)**

Answer: VC Dimension : 1

For one sample situation, there are two cases, both of which can be covered using one-sided interval  $[a, \infty)$ .

When there are two samples, the worst-case scenario happens when the minus sample is to the right of plus sample, which will not be covered by  $[a, \infty)$ . Thus its VC dimension is 1.

(b)

Answer: VC Dimension: 2

For one-sided intervals  $(-\infty, a]$  or  $[a, \infty)$ , we already know their VC dimension is at least 1 from (a). For two samples situation, there are two cases.

For three samples, the worst-case scenario happens when the distribution from left to right is plus, minus, plus; or minus, plus, minus.

(c)

Answer: VC Dimension: 2

From question (a) and (b) we can easily get the VC dimension for finite unions of one-sided intervals is at least 2. (All cases with 2 samples can be covered in (a) and (b), using only one one-sided interval)

For three samples, the worst-case scenario happens when the distribution from left to right is minus, plus, minus.

(d)

Answer: VC Dimension: 4

We should note that all the worst case scenario happens when the samples are distributed in a way that all neighbors of positive samples are negative and vice versa. Thus, we need only to consider this worst case of different number of samples.

For sample number equals to 4, below are the two worst cases, and both would be covered with two intervals:

For sample number equals to 4, worst-case happens when it's ordered in plus, minus, plus, minus, plus:

(e)

Answer: VC Dimension:  $2 \cdot k$

It would fail when sample numbers get to  $2 \cdot k + 1$ , where the samples are distributed with plus, minus, plus, ..., minus, plus. Just as the question 1(d).

**Question 2:**

Question 2 (a)

$\Rightarrow$  For function  $f(p_i) = \ln \frac{1}{p_i} = -\ln p_i$ ,  
its second derivative is :

$$\frac{\partial^2 f(p_i)}{\partial (p_i)^2} = \frac{1}{p_i^2} \geq 0,$$

Thus,  $f$  is a convex function and according to Jensen's inequality, we have the inequality function :

$$f\left(\sum_{i=1}^M \lambda_i x_i\right) \leq \sum_{i=1}^M \lambda_i f(x_i), \text{ where } \lambda_i \geq 0 \text{ and } \sum_{i=1}^M \lambda_i = 1$$

According to the question, we make  $\lambda_i = P(x_i)$ ,  $f(x_i) = -\log(x_i)$ , which is convex function, as proved above.

Thus, we can derive  $x_i = Q(x) / P(x)$  in the question. To adapt to the inequality function, we can adjust the KL-divergence formula to :

$$\begin{aligned} KL(P \parallel Q) &= \sum_x P(x) \log \frac{P(x)}{Q(x)} \\ &= \sum_x P(x) \log \frac{1}{\frac{Q(x)}{P(x)}} \end{aligned}$$

using Jensen's inequality as proved above, we can have :

$$\begin{aligned} \sum P(\alpha) \log \frac{1}{\frac{Q(\alpha)}{P(\alpha)}} &\geq \log \left( \sum_{\alpha} P(\alpha) \frac{Q(\alpha)}{P(\alpha)} \right) \\ &= \log \left( \sum_{\alpha} Q(\alpha) \right) \\ &= 0 \end{aligned}$$

combine the two formula together, we can get that :

$$KL(P||Q) \geq 0$$

Question 2 (b)

$\Rightarrow$  KL Divergence will be 0 when P is same distribution as Q, which is  $P(\alpha) = Q(\alpha)$ .

Thus,  $P(\alpha)/Q(\alpha) = 1$ .

Logarithm of it will be 0. Finally we would get a 0 for its summation.

Question 2 (c)

The maximum will be reached when any  $\alpha$  made  $Q(\alpha) = 0$ .

That is to say for some  $\alpha$ , Q is impossible event. In that case, we would get infinite as the maximum.



## Question 2 (d)

This is wrong. It can be proved with an example:

Below is an arbitrary case, where  $P$  and  $Q$  have 4 possible options and their probability of  $x$  equals to the 4 options are shown in the chart.

| $x$ | 0    | 1    | 2    | 3    |
|-----|------|------|------|------|
| $P$ | 0.25 | 0.25 | 0.25 | 0.25 |
| $Q$ | 0.1  | 0.4  | 0.1  | 0.4  |

Using the formula,

We can get:

$$KL(P||Q) = (0.5 \log 2.5 + 0.5 \log 0.625) // 0.10$$

$$KL(Q||P) = (0.2 \log 0.4 + 0.8 \log 1.6) // 0.08$$

The values of the two are not the same, which violates the equality.

## Question 2 (c)

According to the definition, we can derive the left part as:

$$KL(p(x,y) \parallel Q(x,y)) = \sum_x \sum_y p(x,y) \log \frac{p(x,y)}{Q(x,y)}$$

.... (1)

Using the product rule of joint probability, we can have:

$$p(x,y) = p(y|x) p(x)$$

$$Q(x,y) = Q(y|x) Q(x)$$

.... (2)

Using (2) into (1), the formula can be derived into:

$$\begin{aligned} & \sum_x \sum_y p(x,y) \log \frac{p(x,y)}{Q(x,y)} \\ &= \sum_x \sum_y p(y|x) p(x) \log \frac{p(y|x) p(x)}{Q(y|x) Q(x)} \\ &= \sum_x \sum_y p(y|x) p(x) \log \frac{p(y|x)}{Q(y|x)} + \\ & \quad \sum_x \sum_y p(y|x) p(x) \log \frac{p(x)}{Q(x)} \\ & \quad \dots (3) \end{aligned}$$



$$= \sum_{\alpha} P(\alpha) \left( \sum_y P(y|\alpha) \log \frac{P(y|\alpha)}{Q(y|\alpha)} \right) + \sum_{\alpha} P(\alpha) \log \frac{P(\alpha)}{Q(\alpha)} \left( \sum_y P(y|\alpha) \right)$$

For the last sum, the first part is KL-divergence between two conditional probability distribution,  $P(y|\alpha)$  and  $Q(y|\alpha)$ , which is:

$$\sum_{\alpha} P(\alpha) \left( \sum_y P(y|\alpha) \log \frac{P(y|\alpha)}{Q(y|\alpha)} \right) = KL(P(y|\alpha) || Q(y|\alpha))$$

For the right half part, since

$$\sum_y P(y|\alpha) = 1$$

Thus we can further get:

$$\begin{aligned} & \sum_{\alpha} P(\alpha) \log \frac{P(\alpha)}{Q(\alpha)} \left( \sum_y P(y|\alpha) \right) \\ &= \sum_{\alpha} P(\alpha) \log \frac{P(\alpha)}{Q(\alpha)} \\ &= KL(P(\alpha) || Q(\alpha)) \end{aligned}$$



So far, we've proved that the left part is equal to the summation of right part.

It's proved correct.

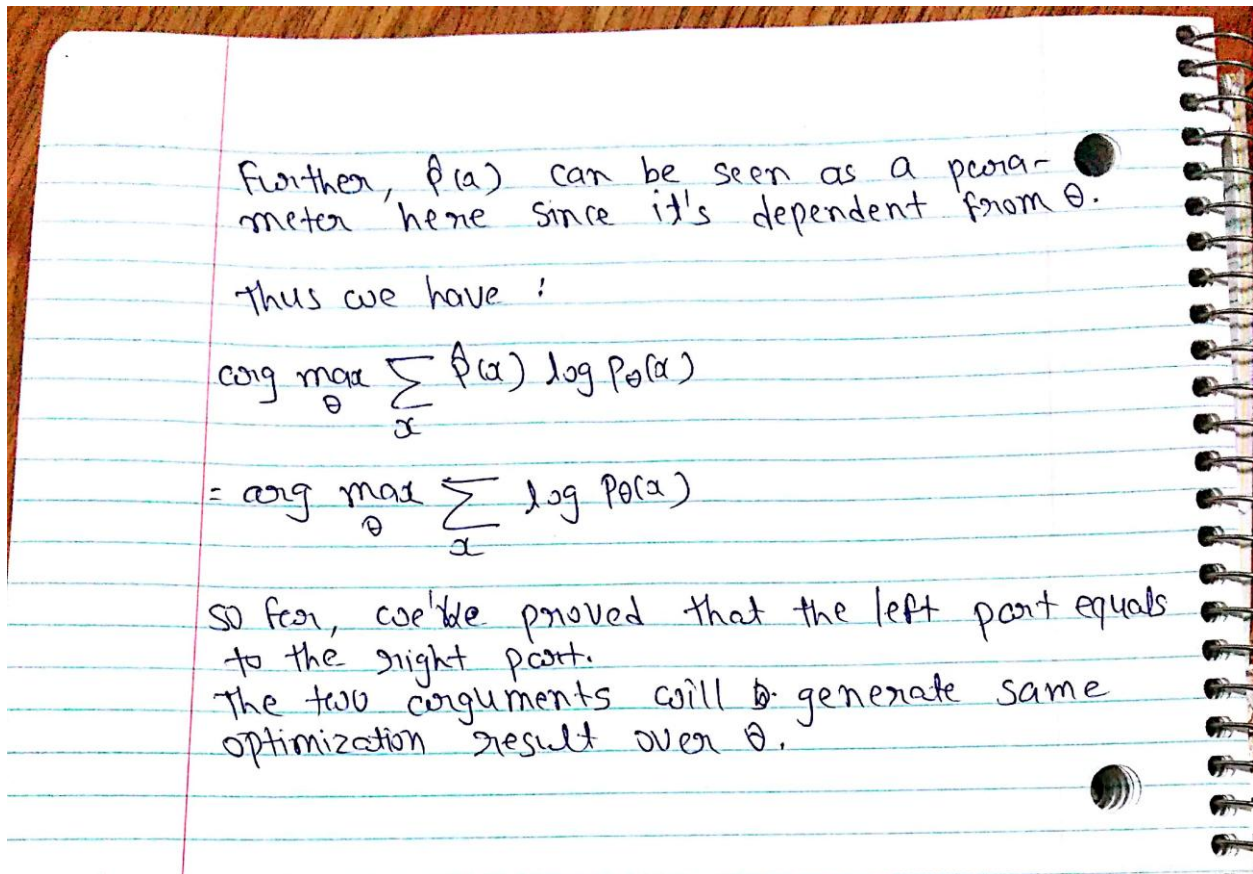
Question 2 (F)

According to KL-divergence definition, the left part can be re-written as:

$$\begin{aligned} & \arg \min_{\theta} KL(\hat{P} \parallel P_{\theta}) \\ &= \arg \min_{\theta} \sum_x \hat{P}(x) \log \frac{\hat{P}(x)}{P_{\theta}(x)} \\ &= \arg \min_{\theta} \left[ \sum_x \hat{P}(x) \log \hat{P}(x) - \sum_x \hat{P}(x) \cdot \log P_{\theta}(x) \right] \end{aligned}$$

Since  $\hat{P}$  stands for the empirical distribution, which is based on counts for each value of  $x$  in this data, it has fixed value / distribution over  $x$ , and is dependent from  $\theta$ , thus we have,

$$\begin{aligned} & \arg \min_{\theta} \left[ \sum_x \hat{P}(x) \log \hat{P}(x) - \sum_x \hat{P}(x) \log P_{\theta}(x) \right] \\ &= \arg \max_{\theta} \sum_x \hat{P}(x) \cdot \log P_{\theta}(x) \end{aligned}$$



### Question 3:

Please check the file **Kmeans.m**. The three input arguments are: **testData**- data loaded from **hw4.dat**; **K** the initial number of clusters, and is 8 in this question; **initialCenter**- the initial centroids given in the question.

The three output: **newCenter**- the final centers of different clusters; **nearestVec**- this is a vector of size 210012\*1, each of it stores the nearest center for the corresponding pixel; **distancelter**- a vector that stores the sum of squared distance from each pixel to the nearest centroid of each iteration.

MATLAB CODE:

```
Data=load('hw4-image.txt');
```

```
K=8;
```

```
initialCenter = [255,255,255;255,0,0;128,0,0;0,255,0;0,128,0;0,0,255;0,0,128;0,0,0];
```

```
[newCenter, nearestVec] = Kmeans(testData, K , initialCenter);
```

- **How many clusters there are in the end?**

There are 6 clusters in the end, which can be checked in the chart below.

- **The final centroids of each cluster:**

The chart records all the centroids information.

| R        | G        | B        |
|----------|----------|----------|
| 241.2296 | 238.6252 | 233.8629 |
| 194.4116 | 136.3331 | 90.94365 |
| 136.2656 | 61.08973 | 10.10385 |
| NaN      | NaN      | NaN      |
| 157.2917 | 97.59398 | 51.4333  |
| NaN      | NaN      | NaN      |
| 78.92744 | 37.10829 | 13.0707  |
| 25.978   | 23.23575 | 23.60599 |

- **The number of pixels associated to each cluster:**

**MATLAB CODE:**

```
noVec=zeros(8,1);
```

```
for i=1:K
```

```
    noVec(i)=sum(nearestVec==i);
```

```
end
```

The result is shown in the chart below:

| index | # of pixels |
|-------|-------------|
| 1     | 4930        |
| 2     | 15190       |
| 3     | 52535       |
| 4     | 0           |



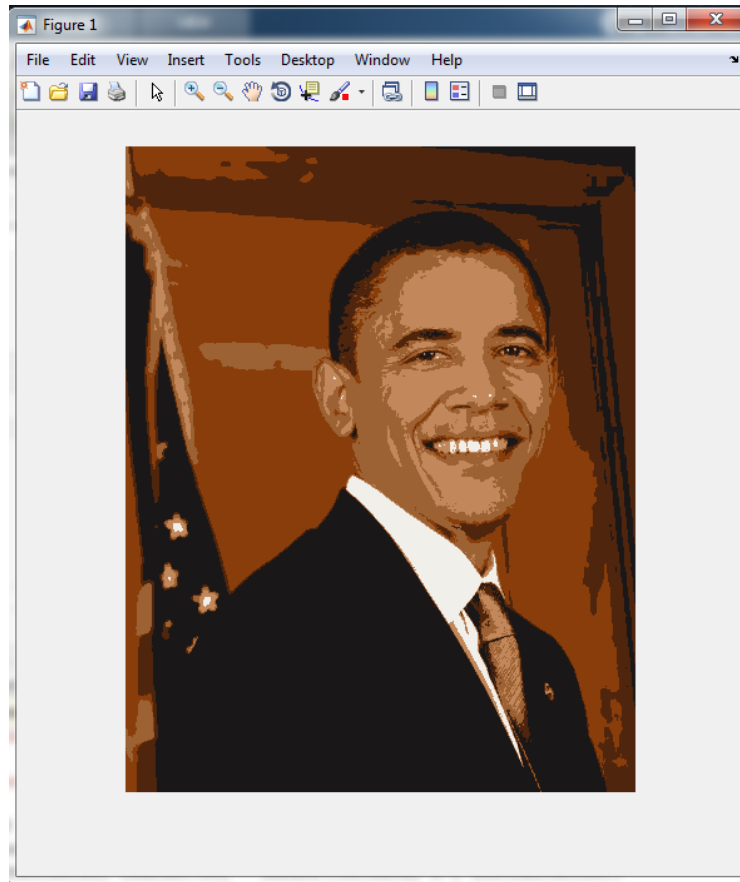
|   |       |
|---|-------|
| 5 | 22075 |
| 6 | 0     |
| 7 | 40365 |
| 8 | 74917 |

- The sum of squared distances from each pixel to the nearest centroid after every iteration of the algorithm:

Below is a chart showing the sum of squared distances for all the pixels to their nearest centroid in total 50 iterations.

| index | value    | index | value   | index | value   | index | value   | index | value   |
|-------|----------|-------|---------|-------|---------|-------|---------|-------|---------|
| 1     | 14442779 | 11    | 4393300 | 21    | 4390808 | 31    | 4135942 | 41    | 4135029 |
| 2     | 5799305  | 12    | 4394124 | 22    | 4387059 | 32    | 4135806 | 42    | 4135105 |
| 3     | 5347163  | 13    | 4395501 | 23    | 4378179 | 33    | 4135453 | 43    | 4135148 |
| 4     | 5187489  | 14    | 4395484 | 24    | 4356133 | 34    | 4135970 | 44    | 4135168 |
| 5     | 4765128  | 15    | 4395365 | 25    | 4292430 | 35    | 4135470 | 45    | 4135189 |
| 6     | 4423282  | 16    | 4395379 | 26    | 4222993 | 36    | 4135199 | 46    | 4135182 |
| 7     | 4385986  | 17    | 4395338 | 27    | 4184325 | 37    | 4134699 | 47    | 4135174 |
| 8     | 4386233  | 18    | 4394996 | 28    | 4159897 | 38    | 4134667 | 48    | 4135182 |
| 9     | 4388704  | 19    | 4394155 | 29    | 4144723 | 39    | 4134824 | 49    | 4135182 |
| 10    | 4391991  | 20    | 4393061 | 30    | 4137972 | 40    | 4134930 | 50    | 4135182 |

The visualized image is shown below:

**Question 4:**

|            |                | K-medoids  | K-means  |
|------------|----------------|--|--|
| Similarity |                | The K-medoids algorithm shares the properties of K-means that each iteration decreases the error, will always converges, different initial center will give different local optimization, instead of global minimum. |  |
| Difference | Speed          | Converge slow  | Converge fast, by comparison                   |
|            | Noise          | Good to lower variance   | Bad for outliers                               |
|            | Cluster number | Number of clusters will not change after calling this algorithm  | Some cluster might disappear during clustering |

The algorithm for K-medoids can be checked in file ***kMedoids.m***. Note that all the arguments mean the same as in method Kmeans.

References:

- Christopher M. Bishop, *Pattern Recognition and Machine Learning*, Springer, 2006
- Kevin P. Murphy, *Machine Learning. A Probabilistic Perspective*, MIT Press, 2012.
- [https://en.wikipedia.org/wiki/Kullback%E2%80%93Leibler\\_divergence](https://en.wikipedia.org/wiki/Kullback%E2%80%93Leibler_divergence)
- <http://www.csse.monash.edu.au/~lloyd/tildeMML/KL/>
- Pattern Recognition and Machine Learning, Bishop, Page 56
- A note from friend for the equations.