

# Comp 6321 - Machine Learning - Assignment 4

Federico O'Reilly Regueiro

December 4<sup>th</sup>, 2016

## Question 1: VC dimensions

### 1.a $[a, \infty)$

We can shatter a single point  $p_1, p_1 \in \mathbb{R}$ :

point	label	h
$p_1$	$\oplus$	$[a, \infty), a < p_1$
$p_1$	$\ominus$	$[a, \infty), a > p_1$

But if we have two points,  $p_1, p_2 \mid p_1 < p_2, p_1 \in \oplus, p_2 \in \ominus$ , then  $[a, \infty)$  cannot shatter them. Therefore, for this class of hypothesis:  $VC_{dim} = 1$

### 1.b $(-\infty, a]$ or $[a, \infty)$

Similarly to the previous question, we can shatter one point. Additionally, we can shatter two points,  $p_0, p_1 \mid p_0 < p_1, p_0$ :

point	label	h
$p_1$	$\ominus$	$(-\infty, a], a < p_1$
$p_2$	$\ominus$	
$p_1$	$\ominus$	$[a, \infty), p_1 < a < p_2$
$p_2$	$\oplus$	
$p_1$	$\oplus$	$(-\infty, a], p_1 < a < p_2$
$p_2$	$\ominus$	
$p_1$	$\oplus$	$[a, \infty), a < p_1$
$p_2$	$\oplus$	

However, three points  $p_1, p_2, p_3, \mid p_1 < p_2 < p_3, p_1 \in \ominus, p_2 \in \oplus, p_3 \in \ominus$  cannot be shattered. Therefore, for this class of hypothesis:  $VC_{dim} = 2$

### 1.c Finite unions of one-sided intervals

The union of more than one left-side interval  $(-\infty, a] \cup (-\infty, b] \dots \cup (-\infty, n]$  is equivalent to a single left-side interval  $(-\infty, \max(a, b, \dots n)]$ . The same applies for one or more right-side intervals being equivalent to  $[\min(a, b, \dots n), \infty)$ . Therefore, this hypothesis class is of the form  $(-\infty, a] \cup [b, \infty)$ .

Since  $\{(-\infty, a] \text{ or } [b, \infty)\} \subset \{(-\infty, a] \cup [b, \infty)\}$ , we know this class of hypothesis to be capable of shattering 2 points. But once again, three points  $p_1, p_2, p_3, \mid p_1 < p_2 < p_3, p_1 \in \ominus, p_2 \in \oplus, p_3 \in \ominus$  cannot be shattered with this class of hypothesis. Therefore, for this class:  $VC_{dim} = 2$

### 1.d $[a, b] \cup [c, d]$

This class of hypothesis can shatter four points due to the following:

- a Any four positives can be correctly classified by a single interval as can any labeling with a single positive.
- b Any two positives and two negatives can be classified with two intervals, given that a single interval is assigned to each positive.
- c Labeling three positives and one negative will always yield at most two groups of contiguous positive labels, each of which can be contained in one of the two intervals.

However, if we have five points  $p_1, p_2, p_3, p_4, p_5$ ,  $| p_1 < p_2 < p_3 < p_4 < p_5, p_1 \in \oplus, p_2 \in \ominus, p_3 \in \oplus, p_4 \in \ominus, p_5 \in \oplus$  cannot be shattered with this class of hypothesis. Therefore, for this class:  $VC_{dim} = 4$

### 1.e Unions of $k$ intervals

By induction:

Base step: One interval,  $k = 1, h = [a, b]$ , and two points,  $p_1, p_2 | p_1 < p_2, p_1$ :

point	label	h
$p_1$	$\ominus$	$[a, b], b < p_1$
$p_2$	$\ominus$	
$p_1$	$\ominus$	$[a, b], p_1 < a < p_2 < b$
$p_2$	$\oplus$	
$p_1$	$\oplus$	$[a, b], a < p_1 < b < p_2$
$p_2$	$\ominus$	
$p_1$	$\oplus$	$[a, b], a < p_1, p_2 < b$
$p_2$	$\oplus$	

We increase the set to three points with the following labels  $p_1, p_2, p_3, | p_1 < p_2 < p_3, p_1 \in \oplus, p_2 \in \ominus, p_3 \in \oplus$ , it cannot be shattered Therefore, for the base step  $VC_{dim} = 2 = 2k$ .

Now suppose that for the union of  $k$  intervals, VC dimension is  $2k^1$ , then we need to prove that with  $k + 1$  intervals we are able to shatter  $2(k + 1)$ .

Firstly we note that the most *difficult* configuration to classify would be an alternation of  $\oplus$  and  $\ominus$  points, since it would require using each one of the  $k$  intervals to classify a single point each; any other configuration would require less than  $k$  intervals and we would have some *leftover* intervals to be consumed in classifying newly inserted points.

Inductive step: We add points  $p_{2k+1}, p_{2k+2}$ , with no inequality constraints, to the  $2k$  points shattered with  $k$  intervals. Without loss of generality, we suppose the previous points to be in an alternating configuration of labels as we mentioned above. We can contemplate three possible scenarios for the added points:

i  $p_{2k+1}, p_{2k+2} \in \ominus$

ii  $p_{2k+1} \in \oplus, p_{2k+2} \in \ominus$ , note<sup>2</sup>

iii  $p_{2k+1}, p_{2k+2} \in \oplus$

#### case i

Since the previous  $2k$  points could be shattered and there are no two contiguous  $\oplus$  labels in the previous set of  $2k$  points, introducing two  $\ominus$  labels anywhere will not disrupt prior labeling if the intervals capturing the adjacent  $\oplus$  points are adjusted accordingly.

<sup>1</sup>ie, we can shatter  $2k$  points but not  $(2k) + 1$  points.

<sup>2</sup>Equivalent to  $p_{2k+2} \in \oplus, p_{2k+1} \in \ominus$

**case ii**

As above, the  $\ominus$  point will not disrupt prior labeling. The  $\oplus$  point will either fall beside another  $\oplus$  point where it can be included in the interval<sup>3</sup> capturing the adjacent  $\oplus$ , or at either end of the set, besides an  $\ominus$  point, in which case the  $k + 1^{th}$  interval will correctly classify it.

**case iii**

If the previous  $2k$  points are labeled with alternating  $\ominus$  and  $\oplus$ , then one end of the set will have  $\ominus$  and the other  $\oplus$ . Thus on inserting points  $p_{2k+1}$  and  $p_{2k+2}$  one of them will necessarily fall beside another  $\oplus$  and, in the worst case, the other point could be placed at the end of the interval on the end with the  $\ominus$ , in which case the  $k + 1^{th}$  interval would correctly classify it.

Thus  $k + 1$  intervals shatter  $2(k + 1)$  points. Conversely, with  $2(k + 1) + 1$  points and the following configuration  $\oplus, \ominus, \dots, \oplus$  we would not be able to shatter the set of points with  $k + 1$  intervals.

Thus the inductive step holds.

Then, for this class with  $k$  intervals,  $VC_{dim} = 2k$ .

## Question 2: KL Divergence

### 2.a $KL(P||Q) \geq 0, \forall P, Q$

Since  $\log(x)$  is a concave function, in order to use Jensen's inequality as stated for convex functions, we make the expression convex by proving  $-KL(P||Q) \leq 0, \forall P, Q$ . We use  $P(x_i) > 0, \forall P(x_i)$  for convenience.

$$\begin{aligned} -KL(P||Q) &= -\sum_{i=1}^m P(x_i) \log \left( \frac{P(x_i)}{Q(x_i)} \right) \\ &= \sum_{i=1}^m P(x_i) \log \left( \frac{Q(x_i)}{P(x_i)} \right) \end{aligned}$$

And by Jensen's inequality, taking  $\frac{P(x_i)}{Q(x_i)}$  to be a random variable uniformly distributed over  $i$ , we can then write:

$$\begin{aligned} -KL(P||Q) &\leq \log \left( \sum_{i=1}^m P(x_i) \frac{Q(x_i)}{P(x_i)} \right) \\ &\leq \log \left( \sum_{i=1}^m Q(x_i) \right) \\ &\leq \log(1) \\ &\leq 0 \end{aligned}$$

### 2.b $KL(P||Q) = 0?$

When both distributions are equal, i.e.  $P(x_i) = Q(x_i), \forall i$ ,  $KL(P||Q)$  becomes:

$$\begin{aligned} KL(P||Q) &= -\sum_i P(x_i) \log \left( \frac{Q(x_i)}{P(x_i)} \right) \\ &= -\sum_i P(x_i) \log(1) \\ &= 0 \end{aligned}$$

Which makes sense since the divergence of two equal distributions should be zero.

---

<sup>3</sup>Once the bounds of said interval have been adjusted

## 2.c Max $KL(P||Q)$ ?

$$KL(P||Q) = \sum_i P(x_i) \log \left( \frac{P(x_i)}{Q(x_i)} \right)$$

$$\lim_{Q(x_i) \rightarrow 0} (KL(P||Q)) = \infty, \text{ for some } i, | P(x_i) \not\rightarrow 0$$

Which can be interpreted as the divergence between the true distribution  $P(x)$  and the modeling distribution  $Q(x)$  approaching infinite if  $Q(x)$  cannot represent an event  $x_i$  with a non-zero probability in  $P(x)$ .

## 2.d $KL(P||Q) = KL(Q||P)$ ? Justify

No, suppose  $x_i = \{0, 1\}$  with  $P(0) = \frac{1}{3}, P(1) = \frac{2}{3}$  modelled by  $Q(0) = \frac{1}{2}, Q(1) = \frac{1}{2}$ . Then:

$$\begin{aligned} KL(P||Q) &= \frac{1}{3} \log \left( \frac{\frac{1}{3}}{\frac{1}{2}} \right) + \frac{2}{3} \log \left( \frac{\frac{2}{3}}{\frac{1}{2}} \right) \\ &= \frac{1}{3} \log \left( \frac{2}{3} \right) + \frac{2}{3} \log \left( \frac{4}{3} \right) \\ &= 0.056633 \\ KL(Q||P) &= \frac{1}{2} \log \left( \frac{\frac{1}{2}}{\frac{1}{3}} \right) + \frac{1}{2} \log \left( \frac{\frac{1}{2}}{\frac{2}{3}} \right) \\ &= \frac{1}{2} \log \left( \frac{3}{2} \right) + \frac{1}{2} \log \left( \frac{3}{4} \right) \\ &= 0.058891 \\ KL(P||Q) &\neq KL(Q||P) \end{aligned}$$

**2.e Prove**  $KL(P(Y, X)||Q(Y, X)) = KL(P(X)||Q(X)) + KL(P(Y|X)||Q(Y|X))$

$$\begin{aligned}
KL(P(Y, X)||Q(Y, X)) &= \sum_x \sum_y P(x, y) \log \left( \frac{P(x, y)}{Q(x, y)} \right) \\
&= \sum_x \sum_y P(x, y) \log \left( \frac{P(y|x)P(x)}{Q(y|x)Q(x)} \right) \\
&= \sum_x \sum_y P(x, y) \left( \log \left( \frac{P(x)}{Q(x)} \right) + \log \left( \frac{P(y|x)}{Q(y|x)} \right) \right) \\
&= \sum_x \sum_y P(x, y) \log \left( \frac{P(x)}{Q(x)} \right) + \sum_x \sum_y P(x, y) \log \left( \frac{P(y|x)}{Q(y|x)} \right) \\
&= \sum_x \sum_y P(x|y)P(y) \log \left( \frac{P(x)}{Q(x)} \right) + \sum_y \sum_x P(y|x)P(x) \log \left( \frac{P(y|x)}{Q(y|x)} \right) \\
&= \sum_x \sum_y \left[ P(x|y)P(y) \log \left( \frac{P(x)}{Q(x)} \right) \right] + \sum_y \left[ P(y|x)P(x) \log \left( \frac{P(y|x)}{Q(y|x)} \right) \right] \\
&= \sum_x \log \left( \frac{P(x)}{Q(x)} \right) \sum_y [P(x|y)P(y)] + \sum_y \left[ P(y|x) \log \left( \frac{P(y|x)}{Q(y|x)} \right) \right] \\
&= \sum_x P(x) \log \left( \frac{P(x)}{Q(x)} \right) + \sum_y P(y|x) \log \left( \frac{P(y|x)}{Q(y|x)} \right) \\
&= KL(P(X)||Q(X)) + KL(P(Y|X)||Q(Y|X))
\end{aligned}$$

**2.f Prove**  $\arg \min_{\hat{P}} KL(\hat{P}||P) = \arg \max_{\theta} \sum_{i=1}^m \log P_{\theta}(x_i)$

First, we develop the *LHS* and we note that:

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

And since  $\sum_x \hat{P}(x) \log(\hat{P}(x))$  depends solely on the observations and is fixed w.r.t.  $\theta$ , we can then equivalently write:

$$\begin{aligned}
\arg \min_{\hat{P}} KL(\hat{P}||P_{\theta}) &= \arg \min_{\theta} - \sum_x \hat{P}(x) \log(P_{\theta}(x)) \\
&= \arg \max_{\theta} \sum_x \hat{P}(x) \log(P_{\theta}(x))
\end{aligned} \tag{1}$$

Then, for a given set of  $m$  observations  $x_i, x \in X, i \in \{1, 2, \dots, m\}$  we can have  $n \leq m$  unique values, which we index as  $x_j, x \in X, j \in \{1, 2, \dots, n\}$  to avoid confusion with observation indexing  $i \in \{1, 2, \dots, m\}$ . Thus we can substitute  $\hat{P}(x_j) = \frac{|x_j|}{m}$  into equation 1:

$$\arg \min_{\theta} KL(\hat{P}||P_{\theta}) = \arg \max_{\theta} \sum_{j=1}^n \frac{|x_j|}{m} \log(P_{\theta}(x_j)) \tag{2}$$

Now we develop the *RHS*:

$$\begin{aligned} M.L.E(x_i, \theta) &= \arg \max_{\theta} \sum_{i=1}^m \log(P_{\theta}(x_i)) \\ &= \arg \max_{\theta} \sum_{i=1}^m \frac{1}{m} \log(P_{\theta}(x_i)) \\ &= \arg \max_{\theta} \sum_{j=1}^n \frac{|x_j|}{n} \log(P_{\theta}(x_j)) \end{aligned}$$

Which is precisely equal to equation 2.

### **Question 3: Implementation: K-means**

### **Question 4: K-medoids - advantages and disadvantages vs K-means**

K-medoids has two main advantages: it can use any measure of similarity between points, which means it is more flexible than K-means which uses euclidean distance between vectors. The use of medoids as opposed to means also makes the partitioning more robust towards outliers.

Conversely, a disadvantage of K-medoids is that its runtime cost is  $O(n^2)$  whereas K-means' is  $O(n)$ .