

## HW4

$$1. \text{ dist} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2 + (z_2 - z_1)^2}$$

$$\text{Obs}_1 = \sqrt{(0-0)^2 + (3-0)^2 + (0-0)^2} = \sqrt{9} = 3$$

$$\text{Obs}_2 = \sqrt{(2-0)^2 + (0-0)^2 + (0-0)^2} = \sqrt{4} = 2$$

$$\text{Obs}_3 = \sqrt{(0-0)^2 + (1-0)^2 + (3-0)^2} = \sqrt{1+9} = \sqrt{10}$$

$$\text{Obs}_4 = \sqrt{(0-0)^2 + (-1-0)^2 + (2-0)^2} = \sqrt{1+4} = \sqrt{5}$$

$$\text{Obs}_5 = \sqrt{(-1-0)^2 + (0-0)^2 + (1-0)^2} = \sqrt{1+1} = \sqrt{2}$$

$$\text{Obs}_6 = \sqrt{(1-0)^2 + (1-0)^2 + (1-0)^2} = \sqrt{1+1+1} = \sqrt{3}$$

2. Since  $k=1$ , select the nearest neighbor to  $(0,0,0)$  which is  $\text{Obs}_5$ .  $\therefore$  the prediction is green.

3. Since  $k=3$ , select the three closest neighbors to  $(0,0,0)$ , which are  $\text{Obs}_5$ ,  $\text{Obs}_6$ , and  $\text{Obs}_2$ . That is green, red, and red respectively.  $\therefore$  the prediction is red.

Each neighbor contributes equally to the prediction when using uniform weights.

$$\text{When } k=3 \quad P(\text{Red}) = \frac{2}{3} \text{ and } P(\text{Green}) = \frac{1}{3}$$

$$\text{All Obs} \quad P(\text{Red}) = \frac{4}{6} = \frac{2}{3} \text{ and } P(\text{Green}) = \frac{2}{6} = \frac{1}{3}$$

Distance weight:

$$w_i = \frac{\exp(-d_i^2)}{\sum_j \exp(-d_j^2)} \quad \text{for } k=3$$

$$\text{Obs}_5 \quad w_5 = \frac{\exp(-1.41^2)}{\exp(-1.41^2) + \exp(-1.73^2) + \exp(-2.05^2)} \approx 0.413$$

$$\text{Obs}_6 \quad w_6 = \frac{\exp(-1.73^2)}{\exp(-1.41^2) + \exp(-1.73^2) + \exp(-2.05^2)} \approx 0.288$$

$$P(\text{Green}) \approx 0.413$$

$$P(\text{Red}) \approx 0.288 + 0.299 = 0.587$$

$$\text{Obs}_2 \quad w_2 = \frac{\exp(-2.00^2)}{\exp(-1.41^2) + \exp(-1.73^2) + \exp(-2.05^2)} \approx 0.299$$

# HWA Continued

Steps of k-means algorithm:

1. Assign each data point to the nearest cluster centroid
2. Update the cluster centroid to be the mean of the data points assigned to it.

Proof:

Assign: Since the data points are being assigned to the nearest cluster centroid, the sum of squared distances between the data points and their cluster centroids,  $\therefore$  the objective function will stay the same or decrease.

Update: Since the cluster centroid is being updated to be the mean of its data points, the average distance between the points will decrease or stay the same.

Since both steps decrease the objective function or keep it the same, the objective function is non-increasing.

Steps of medians of data points w/ Manhattan distance

1. Initialize  $k$  centroids randomly
2. Assign data points to the nearest centroid based on Manhattan distance
3. Update centroids to be the medians of the data
4. Repeat steps 2 and 3 until convergence

Proof:

Assign: Assigning the data points to the nearest centroid based on the Manhattan distance minimizes the sum of absolute distances from each data point to its centroid.  $\therefore$  the next iteration ( $J_{t+1}$ ) will always be less than or equal to the current iteration ( $J_t$ )

Update: Updating the positions of the centroids to the medians of their assigned points ensures the centroids are positioned to minimize the sum of absolute distances.  $\therefore$  the next iteration ( $J_{t+1}$ ) will always be less than or equal to the current iteration ( $J_t$ )

Since both steps decrease the objective function or keep it the same, the objective function is non-increasing.