

3.

7. The table below provides a training data set containing six observations, three predictors, and one qualitative response variable.

Obs.	X_1	X_2	X_3	Y
1	0	3	0	Red
2	2	0	0	Red
3	0	1	3	Red
4	0	1	2	Green
5	-1	0	1	Green
6	1	1	1	Red

Suppose we wish to use this data set to make a prediction for Y when $X_1 = X_2 = X_3 = 0$ using K -nearest neighbors.

- Compute the Euclidean distance between each observation and the test point, $X_1 = X_2 = X_3 = 0$.
- What is our prediction with $K = 1$? Why?
- What is our prediction with $K = 3$? Why?
- If the Bayes decision boundary in this problem is highly non-linear, then would we expect the *best* value for K to be large or small? Why?

.....

a) Euclidean Distances

1- 0, 3, 0

$$d = \sqrt{(0-0)^2 + (3-0)^2 + (0-0)^2}$$

$$= \sqrt{3^2} = \underline{\underline{3}}$$

2. $2, 0, 0$

$$d = \sqrt{2^2 + 0^2 + 0^2} = \sqrt{2^2} = \underline{\underline{2}}$$

3. $0, 1, 3$

$$d = \sqrt{0^2 + 1^2 + 3^2} = \sqrt{1+9} = \sqrt{10}$$

$$\therefore d = \underline{\underline{3.16}}$$

4. $0, 1, 2$

$$d = \sqrt{0^2 + 1^2 + 2^2} = \sqrt{1+4} = \sqrt{5}$$

$$\therefore d = \underline{\underline{2.24}}$$

5. $-1, 0, 1$

$$d = \sqrt{(-1)^2 + 0^2 + 1^2} = \sqrt{1+1} = \sqrt{2}$$

$$\therefore d = \underline{\underline{1.414}}$$

6. $1, 1, 1$

$$d = \sqrt{1^2 + 1^2 + 1^2} = \sqrt{3}$$

$$\therefore d = \underline{\underline{1.732}}$$

b) When $K=1$, our prediction will be

Green since the point $(-1, 0, 1)$ is the

closest to the test point $(0, 0, 0)$; and

the point $(-1, 0, 1)$ belongs to class 'Green'.

c) When $K=3$, the points that we'll take into consideration will be:

$(-1, 0, 1)$ Green

$(1, 1, 1)$ Red

$(2, 0, 0)$ Red

Since $2/3$ of the closest points belong to class 'Red', our prediction for $(0,0,0)$ with $k=3$ will be 'Red'.

d) If the Bayes decision boundary is highly non-linear, we can expect the best k value to be small as a model with a smaller k value would only look at few of the nearest neighbours of a test point for prediction, thus allowing the model to estimate a highly non-linear Bayes decision boundary.