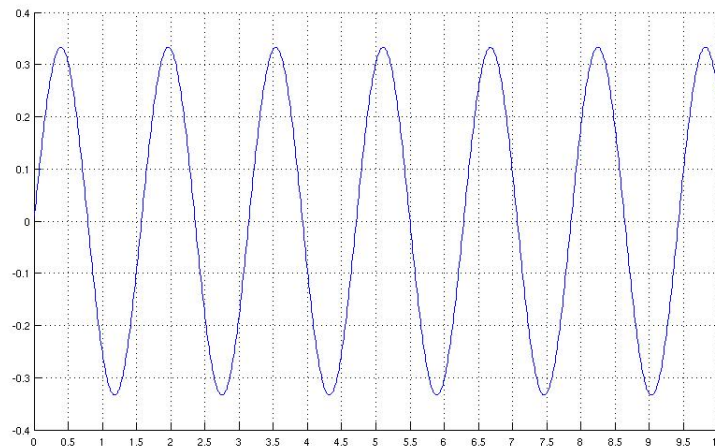


COMS20011 - Data-Driven Computer Science

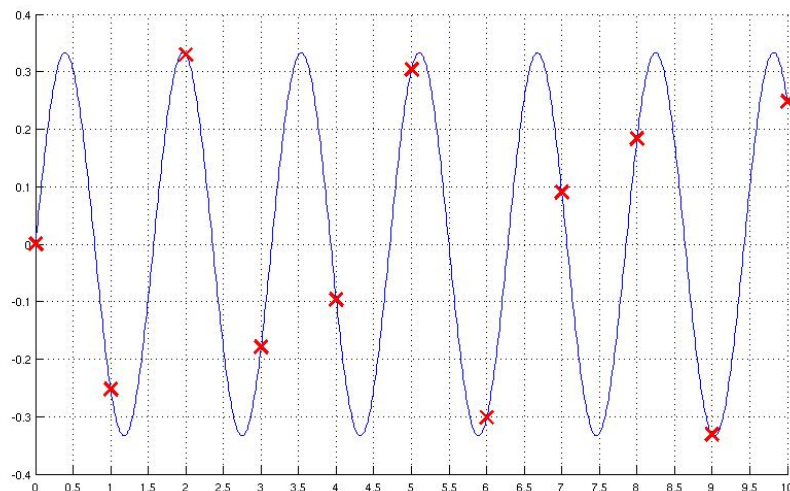
Problem Sheet 1 - Data Acquisition and Distances

January 2023

1. On the $\sin(x)$ signal below, label the following terms and approximate their values: period, frequency and amplitude

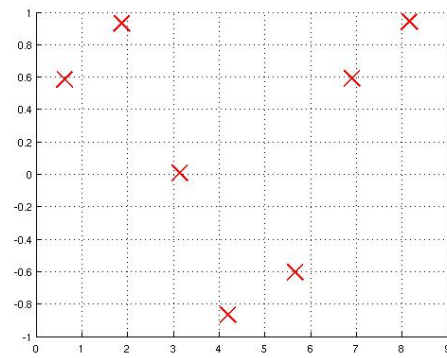
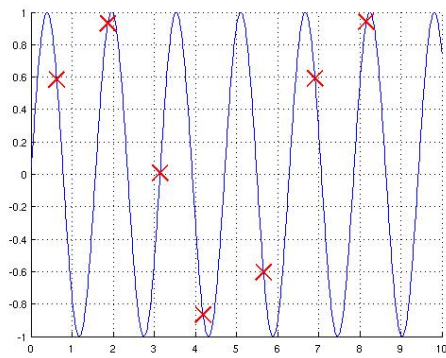


2. For the signal above, convert it into its digital representation using the sampled points. You need to think about the number of bits you would represent each sample as. This is referred to as **Quantisation**. Example, if you need 8 different levels of sound, then 3 bits are sufficient ($2^3 = 8$).



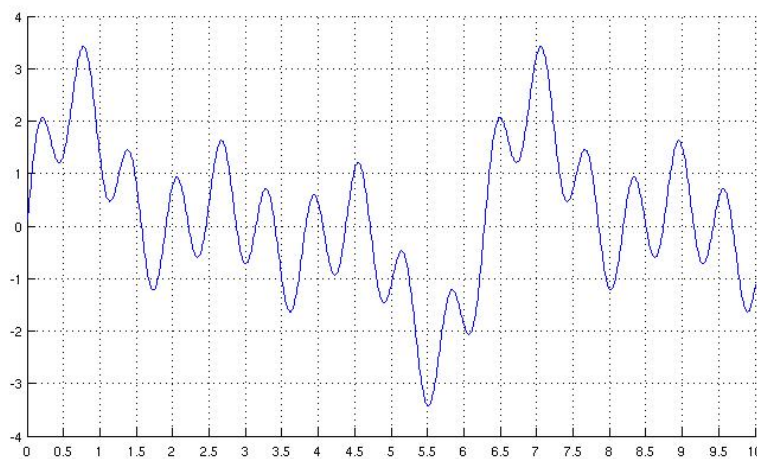
What is the sampling rate in this case??

3. Repeat the digitization and reconstruction step for this data below, can you notice any difference?



4. Based on your understanding of the **Nyquist Sampling Rate** theorem, what is a sufficient sampling rate for the signal below?

Note: You might want to look at Fourier Analysis (ahead of our deeper look later in the course) to understand how this sinusoidal wave was constructed.



5. **Refreshing your memory:**

For the set of measurements:

-3, 2, 4, 6, -2, 0, 5

calculate:

mean

median

variance

standard deviation

6. **Distance measures:** Calculate the following distance measures for the data provided:

- $A = (4, 5, 6)$, $B = (2, -1, 3)$ - Distance Measure Manhattan Distance L_1
- $P = (4, 5, 6)$, $Q = (2, -1, 3)$ - Distance Measure 3-norm L_3
- $E = (4, 5, 6)$, $F = (2, -1, 3)$ - Distance Measure Chebyshev Distance L_∞
- $A1 = \text{'Shot'}$, $A2 = \text{'Chop'}$ - Distance Measure Hamming Distance
- $A1 = \text{'weather'}$, $A2 = \text{'further'}$ - Distance Measure Hamming Distance
- $A1 = \text{'Tank'}$, $A2 = \text{'Thanks'}$ - Distance Measure Edit Distance
- $A1 = \text{'water'}$, $A2 = \text{'further'}$ - Distance Measure Edit Distance
- $A1 = \text{'plankton'}$, $A2 = \text{'plants'}$ - Distance Measure Edit Distance
- *** OPTIONAL *** Order, ascendingly, the following words { 'tap', 'river', 'liquid', 'ice' } based on their WUP relatedness to: 'water'. Use 1-WUP as the distance measure and the online <http://ws4jdemo.appspot.com>

7. **Distance measures:** Assume you were given a set of whatsapp messages, each with a timestamp (yy-mm-dd hh:mm) and text content (word, word, ...). Propose a distance measure for:

- calculating whether one message is an exact copy of the other message
- calculating whether one message was sent before the other message
- calculating whether one message contains the same set of words as the other message
- calculating whether one message contains the other message (with potential extras at the start and the end)
- calculating whether both messages discuss the same topic

Check your distance measures satisfy: non-negativity, reflexive, symmetric and triangle inequality.

8. You collected a four dimensional dataset of values $\mathbf{x} = (x_1, x_2, x_3, x_4)$ and calculated the mean to be $(3, 2.6, -0.4, 2.6)$, and the covariance matrix to be

$$\begin{bmatrix} 4 & 0.1 & -4 & -0.1 \\ 0.1 & 0.01 & -0.1 & 0 \\ -4 & -0.1 & 4 & 0.1 \\ -0.1 & 0 & 0.1 & 9 \end{bmatrix}$$

- You are asked to only select two variables, x_1 and another variable, to take forward for a machine learning algorithm that predicts future values of the variable \mathbf{x} . Which other variable would you pick: x_2 , x_3 or x_4 and why?
- Calculate the eigenvalues and eigenvectors for your chosen covariance matrix
- Using the probability density function of the normal distribution in two dimensions, calculate the probability that the following new data $(3, 2.61, 0, 3)$ belongs to the dataset \mathbf{x} [Note: only use the two variables you picked in (a)]