

---

## DSC 40A - Homework 7

Due: Friday, March 8 at 11:59pm

---

Write your solutions to the following problems by either typing them up or handwriting them on another piece of paper. Homeworks are due to Gradescope by 11:59pm on the due date. You can use a slip day to extend the deadline by 24 hours.


Homework will be evaluated not only on the correctness of your answers, but on your ability to present your ideas clearly and logically. You should **always explain and justify** your conclusions, using sound reasoning. Your goal should be to convince the reader of your assertions. If a question does not require explanation, it will be explicitly stated.

Homeworks should be written up and turned in by each student individually. You may talk to other students in the class about the problems and discuss solution strategies, but you should not share any written communication and you should not check answers with classmates. You can tell someone how to do a homework problem, but you cannot show them how to do it.

For each problem you submit, you should **cite your sources** by including a list of names of other students with whom you discussed the problem. Instructors do not need to be cited.

This homework will be graded out of 50 points. The point value of each problem or sub-problem is indicated by the number of avocados shown.

### Problem 1. Reflection and Feedback Form




 Make sure to fill out this [Reflection and Feedback Form, linked here](#) for two points on this homework! This form is primarily for your benefit; research shows that reflecting and summarizing knowledge helps you understand and remember it.

### Problem 2. Conditional Independence

Given the following three statements:

$$\begin{aligned}(1) \quad P(X, Y|E) &= P(X|E)P(Y|E) \\(2) \quad P(X|Y, E) &= P(X|E) \\(3) \quad P(Y|X, E) &= P(Y|E)\end{aligned}$$

In this problem, you will need to show the three of them are equivalent. You may assume none of the probabilities listed above can be zero.

- a)  Show that (1) and (2) are equivalent. That is,  $(1) \Rightarrow (2)$  and  $(2) \Rightarrow (1)$
- b)  Show that (2) and (3) are equivalent. That is,  $(2) \Rightarrow (3)$  and  $(3) \Rightarrow (2)$
- c)  Show that (1) and (3) are equivalent. That is,  $(1) \Rightarrow (3)$  and  $(3) \Rightarrow (1)$

### Problem 3. Independence, and Markov Property

A Markov Chain is a series of random variables that each one of them only depends on the preceding one. Don't get overwhelmed by its name, let's consider the following simple example:

$$X_1 \rightarrow X_2 \rightarrow X_3$$

which means  $X_2$  only depends on  $X_1$  and  $X_3$  only depends on  $X_2$ . That is, we have  $P(X_3|X_1, X_2) = P(X_3|X_2)$

- a) 🧐🧐🧐 Show that their joint distribution is

$$P(X_1, X_2, X_3) = P(X_1)P(X_2|X_1)P(X_3|X_2)$$

- b) 🧐🧐🧐🧐🧐 Show the Markov Property that given the present, the future and the past are independent in the Markov chain. That is,

$$P(X_1, X_3|X_2) = P(X_1|X_2)P(X_3|X_2)$$

#### Problem 4. Independence and Conditional Independence

Consider the sample space  $S = \{a, b, c, d, e, f, g\}$  with associated probabilities given in the table below.

outcome	$a$	$b$	$c$	$d$	$e$	$f$	$g$
probability	$\frac{5}{21}$	$\frac{2}{21}$	$\frac{1}{21}$	$\frac{4}{21}$	$\frac{2}{21}$	$\frac{4}{21}$	$\frac{3}{21}$

Let  $X = \{d, e\}$  and  $Y = \{e, f\}$ . Remember to show your work for all calculations.

- a) 🧐🧐🧐🧐 Are  $X$  and  $Y$  independent?
- b) 🧐🧐🧐🧐🧐 Determine if  $X$  and  $Y$  are conditionally independent given each of the following events  $Z$ .
1.  $Z = \{a, b, d, e, f, g\}$
  2.  $Z = \{a, d, e, f, g\}$
  3.  $Z = \{d, e, f, g\}$

#### Problem 5. Independence and Complements

Let  $E$  and  $F$  be two events in a sample space  $S$ , with  $0 < P(F) < 1$ .

- a) 🧐🧐🧐🧐 If  $P(E|F) = P(E|\overline{F})$ , must it be true that  $E$  and  $F$  are independent? Provide a proof of independence, or give a counterexample by specifying a sample space  $S$  and two dependent events  $E$  and  $F$  that satisfy the given conditions.
- b) 🧐🧐🧐🧐 If  $P(E|F) = P(\overline{E}|F)$ , must it be true that  $E$  and  $F$  are independent? Provide a proof of independence, or give a counterexample by specifying a sample space  $S$  and two dependent events  $E$  and  $F$  that satisfy the given conditions.

#### Problem 6. Baby Avi

When Avi was a baby, he was a picky eater. When he was offered a food, he'd eat it only sometimes, depending on the type of food. Baby Avi ate

- bananas 95% of the time,
- crackers 60% of the time,
- meat 45% of the time, and
- zucchini 30% of the time.

- a) 🍌🍌🍌🍌🍌 Baby Avi's grandpa gave him one of the above four foods and Baby Avi ate it all up. You have no idea which of the four foods it was, so assume it was equally likely to be any of them.

Given that Baby Avi ate the food, what's the probability that it was a banana? A cracker? Meat? Zucchini? Show your work.

- b) 🍌🍌🍌🍌🍌 As before, Baby Avi's grandpa gave him one of the above four foods and Baby Avi ate it all up. This time, suppose you know that Baby Avi's grandpa offers

- bananas 15% of the time,
- crackers 40% of the time,
- meat 10% of the time, and
- zucchini 35% of the time.

Given that Baby Avi ate the food, what's the probability that it was a banana? A cracker? Meat? Zucchini? Show your work.

- c) 🍌🍌🍌🍌 Compare your answers to part (a) and part (b) above. Identify which of the four probabilities you computed increased and which decreased, and explain why this makes sense intuitively.

### Problem 7. Optional Extra Credit Challenge: Classifying Neutrino Signals in Majorana Demonstrator!

We are hosting a optional class-wide competition to see who can make the best predictions! Top 10 predictions can earn extra credit on Midterm 2! Please note that no point(🍌) is assigned to this problem.

High-Purity Germanium (HPGe) detector is one of the most sensitive detectors human beings have ever manufactured. It is sensitive in a sense that it measures the energy of elementary particles (electron, photon, etc) very accurately. Because of this, HPGe detectors have a wide range of applications, including the search for neutrinos and dark matters, medical imaging, as well as nuclear non-proliferation.

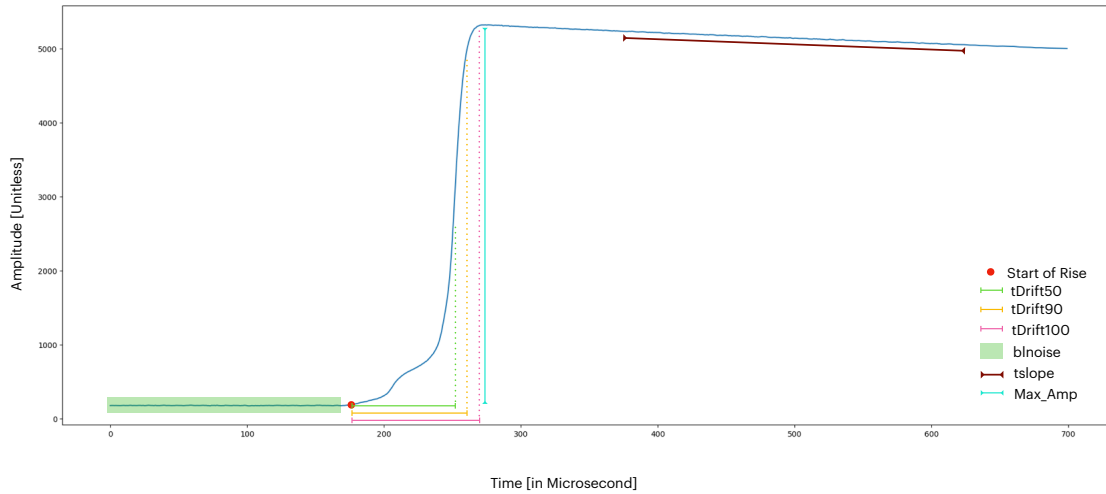
When a particle comes into HPGe detector, it produces a waveform, or time series data, as shown in the picture below. A time series is a sequence of data points that occur in successive order over some period of time. More formally, we can define time series this way: for each data point, a time series contains  $n$  pairs of  $t_i, a_i$  where  $t_i$  is the  $i^{th}$  time sample and  $a_i$  is the value at  $i^{th}$  time sample. To simplify this problem, we will not ask you to directly analyze the time series, but we extract certain features from the HPGe time series for you to build model.

The MAJORANA DEMONSTRATOR Experiment uses HPGe detector to search for neutrinoless double-beta decay. This kind of decay happens at least a trillion times slower than the age of our universe. Finding it will unravel the fundamental understanding of particle physics. For more information please read this [Press Release](#).

In the [supplementary Jupyter notebook \(linked\)](#), you are given access to a CSV containing training data with information about 3000 elementary particles which deposits their energy in a HPGe detector. Some of them are signal-like, i.e. they exhibit the same shape with neutrinoless double-beta decay, others are noise-like, i.e. they looks different from neutrinoless double-beta decay. This is a labeled dataset where signal-like data has a label of 1 and background-like data has a label of 0. We read this in as a DataFrame where the columns are different parameters. The columns are:

- **tDrift50:** Period from the Start of Rise ( $t_{SR}$ ) to when the waveform reaches 50% of Max\_Amp ( $t_{50}$ ), can also be written as  $t_{50} - t_{SR}$
- **tDrift90:** Period from the Start of Rise ( $t_{SR}$ ) to when the waveform reaches 90% of Max\_Amp ( $t_{90}$ ), can also be written as  $t_{90} - t_{SR}$

- **tDrift100:** Period from the Start of Rise ( $t_{SR}$ ) to when the waveform reaches Max\_Amp, can also be written as  $t_{Max\_Amp} - t_{SR}$
- **blnoise:** The standard deviation of amplitude values  $a_i$  in the green-colored region.
- **tslope:** The slope of the waveform tail.
- **Energy:** The energy of each waveform, i.e. the target of the previous challenge
- **Current\_Amplitude:** A new parameter extracted from the waveform, by taking a derivative of the waveform and read out the maximum of the derivative.



Your task for this problem is to produce a "classification score" for each HPGe detector waveforms. This score does not have to be between 0 and 1, but signal data point (data point with a label of 1) should have a higher score than noise data point (data point with a label of 0). This extra credit opportunity as minimal constraint: you are not allowed to import any additional python package, but you are allowed to use any model (including but not limited to the Naive Bayes classifier we will discuss in class) and any amount of parameters.

We've provided you with a function `calculate_AUC` to calculate the area under curve (AUC) of the Receiver Operating Characteristic (ROC) curve of your predictions on each waveform in the training data. For more information please read [this link](#). Your job is to fill in the body of the `predict` function. This function should take as input one row of the DataFrame (corresponding to one particular waveform) and return the predicted classification score corresponding to this waveform. How you make this prediction is up to you, subject to the rules above. Feel free to add more cells and functions, and to change the provided `predict` function, but do not change the provided `calculate_AUC` function.

The ten best prediction functions (as determined by the AUC on the hidden test data) will earn some extra credit on the upcoming midterm exam according to the following scheme: for  $n \leq 10$ , the  $n^{th}$  ranked prediction function in the class earns  $11 - n$  percentage points as extra credit on Midterm 1. But the total extra creit from both challenges (HW4 Q6 plus this one) is capped at 10%.