

# M5-L1-P1

October 3, 2023

## 1 Problem 1 (6 points)

In this problem, you will implement a function to calculate gini impurity on an arbitrary input vector.

For reference, the formula for Gini impurity is:

$$\text{Gini}(D) = 1 - \sum_{i=1}^k p_i^2$$

where  $D$  is the dataset containing samples from  $k$  classes and  $p_i$  is the probability of a data point belonging to class  $i$ .

### 1.1 Gini Impurity Function

Complete the function `gini(D)` below. It should take as input a 1-D array, where is the number of samples corresponding to each output class.

For example, consider the input array `D = np.array([4, 9, 7, 0, 3])` In this example, there are 5 input classes and 23 total samples. For this input, your function should return 0.707.

Your function should work regardless of the length of the input vector.

```
[ ]: import numpy as np

def gini(D):
    D_adjusted = (D/np.sum(D))**2
    return (1 - np.sum(D_adjusted))

D = np.array([4, 9, 7, 0, 3])
g = gini(D)
print(f"gini([4,9,7,0,3]) = {g:.3f} (should be about {0.707})")
```

`gini([4,9,7,0,3]) = 0.707 (should be about 0.707)`

### 1.2 More test cases

Compute and print the gini impurity for D1, D2, D3, and D4, defined below:

```
[ ]: D1 = np.array([1,0,0])
      D2 = np.array([0,0,4])
      D3 = np.array([0, 20, 0, 0, 0, 3])
      D4 = np.array([6, 6, 6, 6])

      i = 1
      for D in [D1, D2, D3, D4]:
          print("Gini Impurity for D%i: " % i + "%f" % gini(D))
          i += 1
```

```
Gini Impurity for D1: 0.000000
Gini Impurity for D2: 0.000000
Gini Impurity for D3: 0.226843
Gini Impurity for D4: 0.750000
```