# Homework 1

Name: Shuyue Lai

NetID: shuyuel2

## Problem 1 Statistical Properties

**(a) Max and min:**

```python
# 1(a) Max and Min
# Max
def getMax(array):
    max = 0
    for item in array:
        if item > max:
            max = item
    # return np.max(array)
    return max

Midterm_Max = getMax(Midterm)
print("Midterm Max:", Midterm_Max)
Final_Max = getMax(Final)
print("Final Max:", Final_Max)

# Min
def getMin(array):
    min = 1000
    for item in array:
        if item < min:
            min = item
    # return np.min(array)
    return min

Midterm_Min = getMin(Midterm)
print("Midterm Min:", Midterm_Min)
Final_Min = getMin(Final)
print("Final Min:", Final_Min)
```

```
Midterm Max: 99
Final Max: 100
Midterm Min: 75
Final Min: 77
```

- Max:
    - Midterm = 99
    - Final = 100
- Minx:

- Midterm = 75
  - Final = 77

**(b) Mean, mode and median:**

```python
2  # Mean
3  def getMean(array):
4      sum = 0
5      for item in array:
6          sum = sum + item
7      amount = len(array)
8      mean = sum / amount
9      return mean
10
11 Midterm_Mean = getMean(Midterm)
12 print("Midterm Mean:", Midterm_Mean)
13 Final_Mean = getMean(Final)
14 print("Fianl Mean:", Final_Mean)
15
16 # Mode
17 def getMode(array):
18     counts = np.bincount(array)      # get frequncy
19     max_count = getMax(counts)       # get highest frequency
20
21     return np.where(counts==max_count)
22
23 Midterm_Mode = getMode(Midterm)
24 print("Midterm Mode:", Midterm_Mode)
25 Final_Mode = getMode(Final)
26 print("Final Mode:", Final_Mode)
27
28 # Median
29 def getMedian(array):
30     return np.median(array)
31
32 Midterm_Median = getMedian(Midterm)
33 print("Midterm Median:", Midterm_Median)
34 Final_Median = getMedian(Final)
35 print("Final Median:", Final_Median)
```

```
Midterm Mean: 88.7
Fianl Mean: 88.2
Midterm Mode: (array([96]),)
Final Mode: (array([80, 88]),)
Midterm Median: 89.0
Final Median: 88.0
```

- Mean:

  - Midterm: 88.7
  - Final: 88.2
- Mode:

  - Midterm: 96
  - Final: 80, 88
- Median:

  - Midterm: 89

- Final: 88

**(c) First quartile, third quartile and inter-quartile range:**

```
 1  # 1(c)First Quartile, Third Quartile and Inter-quartile Range
 2  # First Quartile
 3  def getFirstQuartile(array):
 4      array = np.sort(array)     # sort array
 5      length = len(array)
 6      mid = (int)(length / 2)
 7      first = (int)(mid / 2)
 8      third = mid + first
 9
10      if length >= 4:
11          if mid % 2== 0:
12              first_quartile = (array[first-1] + array[first]) / 2
13              third_quartile = (array[third-1] + array[third+1]) / 2
14          else:
15              first_quartile = array[first]
16              third_quartile = array[third]
17
18      inter_quartile = third_quartile - first_quartile
19
20      return array, first_quartile, third_quartile, inter_quartile
21
22  [Midterm_Sorted, Midterm_First_Quartile, Midterm_Third_Quartile, Mid
23  print("Midterm Sorted Array:", Midterm_Sorted)
24  print("Midterm First Quartile:", Midterm_First_Quartile)
25  print("Midterm Third Quartile:", Midterm_Third_Quartile)
26  print("Midterm Inter Quartile:", Midterm_Inter_Quartile)
27  [Final_Sorted, Final_First_Quartile, Final_Third_Quartile, Final_Int
28  print("Final Sorted Array:", Final_Sorted)
29  print("Final First Quartile:",Final_First_Quartile)
30  print("Final Third Quartile:", Final_Third_Quartile)
31  print("Final Inter Quartile:", Final_Inter_Quartile)
```

```
Midterm Sorted Array: [75 78 84 86 88 90 95 96 96 99]
Midterm First Quartile: 84
Midterm Third Quartile: 96
Midterm Inter Quartile: 12
Final Sorted Array: [ 77  80  80  85  88  88  90  95  99 100]
Final First Quartile: 80
Final Third Quartile: 95
Final Inter Quartile: 15
```

- Midterm:

| 75 | 78 | (84) | 86 | 88 | 90 | 95 | (96) | 96 | 99 |
|----|----|------|----|----|----|----|------|----|----|

  - First Quartile: 84
  - Third Quartile: 96
  - Inter-Quartile: 12
- Final:

| 77 | 80 | 80 | 85 | 88 | 88 | 90 | 95 | 99 | 100 |
|----|----|----|----|----|----|----|----|----|-----|

- 
  - First Quartile: 80
  - Third Quartile: 95
  - Inter-Quartile: 15

**(d) Variance(sample & population):**

```python
1   # 1(d)Variance (Sample & Population)
2   def getVariancePopulation(array):
3       length = len(array)
4       square_sum = 0
5       mean = getMean(array)
6
7       for item in array:
8           square_sum = square_sum + (item - mean) * (item - mean)
9
10      population_var = round(square_sum / length, 3)
11
12      # return np.var(array)
13      return population_var
14
15  def getVarianceSample(array):
16      length = len(array)
17      square_sum = 0
18      mean = getMean(array)
19
20      for item in array:
21          square_sum = square_sum + (item - mean) * (item - mean)
22
23      sample_var = round(square_sum / (length - 1) , 3)
24
25      # return np.var(array)
26      return sample_var
27
28  Midterm_Population_Variance = getVariancePopulation(Midterm)
29  print("Midterm Population Variance:", Midterm_Population_Variance)
30  Midterm_Sample_Variance = getVarianceSample(Midterm)
31  print("Midterm Sample Variance:", Midterm_Sample_Variance)
32  Final_Population_Variance = getVariancePopulation(Final)
33  Final_Sample_Variance = getVarianceSample(Final)
34  print("Final Population Variance:", Final_Population_Variance)
35  print("Final Sample Variance:", Final_Sample_Variance)
```

```
Midterm Population Variance: 58.61
Midterm Sample Variance: 65.122
Final Population Variance: 57.56
Final Sample Variance: 63.956
```

- Population:

$$Var = E((X - \bar{x})^2) = E(x^2) - [E(x)]^2$$

$$Var = \frac{1}{n} * \sum_{i=1}^{n}(x_i - \bar{x})^2$$

- Sample:

$$Var = \frac{1}{n-1} * \sum_{i=1}^{n}(x_i - \bar{x})^2$$

- Midterm:
  - Population Variance: 58.610
  - Sample Variance: 65.122
- Final:
  - Population Variance: 57.560
  - Sample Variance: 63.956

**(e) Standard Deviation (sample & population):**

```
1   # 1(e)Standard Deviation
2   def getStandardDeviationPopulation(array):
3       # return np.std(array)
4       return round(np.sqrt(getVariancePopulation(array)),3)
5
6   def getStandardDeviationSample(array):
7       # return np.std(array)
8       return round(np.sqrt(getVarianceSample(array)), 3)
9
10  Midterm_Standard_Population_Deviation = getStandardDeviationPopulati
11  print("Midterm Standard Population Deviation:", Midterm_Standard_Pop
12  Midterm_Standard_Sample_Deviation = getStandardDeviationSample(Midte
13  print("Midterm Standard Sample Deviation:", Midterm_Standard_Sample_
14  Final_Standard_Population_Deviation = getStandardDeviationPopulation
15  print("Final Standard Population Deviation:", Final_Standard_Populat
16  Final_Standard_Sample_Deviation = getStandardDeviationSample(Final)
17  print("Final Standard Sample Deviation:", Final_Standard_Sample_Devi
```

```
Midterm Standard Population Deviation: 7.656
Midterm Standard Sample Deviation: 8.07
Final Standard Population Deviation: 7.587
Final Standard Sample Deviation: 7.997
```

- Standard Deviation:

$$Std = \sqrt{Var}$$

- Midterm:
  - Population Standard Deviation: 7.656
  - Sample Standard Deviation: 8.070
- Final:
  - Population Standard Deviation: 7.587
  - Sample Standard Deviation: 7.997

# Problem 2 Normalize Data

**(a) Min-Max Normalization:**

```python
# 2(a) Min-Max Normalization
def getMinMaxNormalization(array):
    old_min = getMin(array)
    old_max = getMax(array)
    new_min = 0
    new_max = 1

    new_list = []
    for item in array:
        new_list.append(round((item-old_min)/(old_max-old_min)*(new_
    
    new_array = np.array(new_list)
    return new_array

Midterm_MinMax_Normalization = getMinMaxNormalization(Midterm)

print("Midterm MinMax Normalization Student 1:", Midterm_MinMax_Norm
print("Midterm MinMax Normalization Student 2:", Midterm_MinMax_Norm
print("Midterm MinMax Normalization Student 3:", Midterm_MinMax_Norm
```

```
Midterm MinMax Normalization Student 1: 0.833
Midterm MinMax Normalization Student 2: 0.458
Midterm MinMax Normalization Student 3: 0.125
```

-

$$v' = \frac{v - min}{max - min} * (newMax - newMin) + newMin$$

- Student 1: 0.833

- Student 2: 0.458

- Student 3: 0.125

**(b) Variance(population) of min-max normalized midterm scored:**

```python
# 2(b) Variance(Population) of Min-Max Normalization
def getVarianceMinMaxNormalization(array):
    new_array = getMinMaxNormalization(array)
    return getVariancePopulation(new_array)

Midterm_Variance_MinMax_Normalization = getVarianceMinMaxNormalizati
print("Midterm Variance MinMax Normalization:", Midterm_Variance_Min
```

```
Midterm Variance MinMax Normalization: 0.102
```

- Variance (Population):

$$Var = \frac{1}{n} * \sum_{i=1}^{n}(x_i - \bar{x})^2$$

- Min-Max Normalization:

$$v' = \frac{v - min}{max - min} * (newMax - newMin) + newMin$$

- Midterm Variance Min-Max Normalization:0.102

**(c) Z-Score Normalization of final scores:**

```python
# 2(c) Z-Score Normalizaion
def getZScoreNormalization(array):
    mean = getMean(array)
    std = getStandardDeviationPopulation(array)

    new_list = []
    for item in array:
        new_list.append(round((item-mean)/std,3))

    new_array = np.array(new_list)
    return new_array

Final_ZScore_Normalization = getZScoreNormalization(Final)
print("Final ZScore Normalization Student 4:", Final_ZScore_Normaliz
print("Final ZScore Normalization Student 5:", Final_ZScore_Normaliz
print("Final ZScore Normalization Student 6:", Final_ZScore_Normaliz
```

```
Final ZScore Normalization Student 4: 0.896
Final ZScore Normalization Student 5: -0.422
Final ZScore Normalization Student 6: -1.476
```

- 

$$v' = \frac{v - \mu}{\sigma}$$

- Student 4: 0.896

- Student 5: -0.422

- Student 6: -1.476

**(d) Variance(population) of z-score normalized final score:**

```python
# 2(d) Variance(Population) of Z-Score Normalization
def getVarianceZScoreNormalization(array):
    new_array = getZScoreNormalization(array)
    return getVariancePopulation(new_array)

Final_Variance_ZScore_Normalization = getVarianceZScoreNormalization
print("Final Variance ZScore Normalization:", Final_Variance_ZScore_
```

```
Final Variance ZScore Normalization: 1.0
```

- Variance (Population):

$$Var = \frac{1}{n} * \sum_{i=1}^{n}(x_i - \overline{x})^2$$

- Z-Score Normalization:

$$v' = \frac{v - \mu}{\sigma}$$

- Variance(population) of z-score normalized final score: 1.000

# Problem 3 Relationship between Midterm and Final

**(a) Covariance(Population):**

```
1   # 3(a) Covariance(population)
2   def getCovariancePopulation(array1, array2):
3       mean1 = getMean(array1)
4       mean2 = getMean(array2)
5       length = len(array1)
6
7       covariance = sum(array1 * array2) / length - mean1 * mean2
8       return round(covariance, 3)
9
10  Covariance = getCovariancePopulation(Midterm, Final)
11  print("Covariance:", Covariance)
```

Covariance: 18.16

$$Cov = E[(x_1 - \mu_1) * (x_2 = \mu_2)] = E(x_1 * x_2) - \mu_1 * \mu_2$$

- Covariance: 18.160

**(b)Pearson's Correlation Coefficient(population standard deviation and covariance):**

```
1   # 3(b) Pearson's Correlation Coefficient
2   def getPearsonCorrelation(array1, array2):
3       covariance = getCovariancePopulation(array1, array2)
4       std1 = getStandardDeviationPopulation(array1)
5       std2 = getStandardDeviationPopulation(array2)
6
7       pearson = covariance / (std1 * std2)
8       return round(pearson, 3)
9
10  PearsonCorrelation = getPearsonCorrelation(Midterm, Final)
11  print("Peasrson Correlation Coefficient:", PearsonCorrelation)
```

Peasrson Correlation Coefficient: 0.313

$$Pearson = \frac{cov(x_1, x_2)}{\sigma_{x_1} * \sigma_{x_2}}$$

- Pearson Correlation Coefficient: 0.313

**(c) Independent?**

- Since the Pearson Correlation Coefficient is 0.313, which means M which denotes midterm scores and F denotes final scores are not independent and are positive related.

**(d) Distance**

```python
# 3(d) Distance
# Manhattan Distance
def getManhattan(array1, array2):
    length = len(array1)

    sum = 0
    for iter in range(length):
        sum = sum + abs(array1[iter] - array2[iter])

    return sum

Manhattan = getManhattan(Midterm, Final)
print("Manhattan Distance:", Manhattan)

# Euclidean Distance
def getEuclidean(array1, array2):
    length = len(array1)

    sum = 0
    for iter in range(length):
        sum = sum + pow((array1[iter] - array2[iter]),2)

    return round(pow(sum, 0.5), 3)

Euclidean = getEuclidean(Midterm, Final)
print("Euclidean Distance:", Euclidean)
```

```
28  # Supremum Distance
29  def getSupremum(array1, array2):
30      length = len(array1)
31
32      max = 0
33      for iter in range(length):
34          temp = abs(array1[iter] - array2[iter])
35          if max < temp:
36              max = temp
37
38      return round(temp, 3)
39
40  Supremum = getSupremum(Midterm, Final)
41  print("Supremum Distance:", Supremum)
42
43  # Cosine Similarity
44  def getCosineSimilarity(array1, array2):
45      cosine = np.dot(array1, array2) / (pow(sum(pow(array1, 2)), 0.5)
46
47      #return np.dot(array1, array2) / (norm(array1)*norm(array2))
48      return round(cosine, 3)
49
50  Cosine = getCosineSimilarity(Midterm, Final)
51  print("Cosine Similarity:", Cosine)
```

```
Manhattan Distance: 75
Euclidean Distance: 28.302
Supremum Distance: 16
Cosine Similarity: 0.995
```

- Minkowski Distance:

$$d_{ij} = \sqrt[p]{|x_{i1} - x_{j1}|^p + |x_{i2} - x_{j2}|^p + L + |x_{il} - x_{jl}|^p}$$

- Manhattan Distance:
  -
$$d_{ij} = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + L + |x_{il} - x_{jl}|$$

  - Midterm & Final: 75
- Euclidean Distance:
  -
$$d_{ij} = \sqrt[2]{|x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + L + |x_{il} - x_{jl}|^2}$$

  - Midterm & Final: 28.302
- Supremum Distance:
  -
$$d_{ij} = lim_{p \to \infty} \sqrt[p]{|x_{i1} - x_{j1}|^p + |x_{i2} - x_{j2}|^p + L + |x_{il} - x_{jl}|^p} = max_{f=1}^{l}|x_{if} - x_{jf}|$$

- - Midterm & Final: 16
  - Cosine Similarity of m and f:
    -

$$CosineSimilarity = \frac{A \cdot B}{||A|| * ||B||} = \frac{\sum_{i=1}^{n} A_i * B_i}{\sqrt{\sum_{i=1}^{n}(A_i)^2} * \sqrt{\sum_{i=1}^{n}(B_i)^2}}$$

  - Midterm & Final: 0.995

**(e)**

Manhattan Distance: Means the difference of each student between midterm and final. The distance is always positive and the larger value means the difference is larger. In the scenario, the larger value means students perform different in midterm and final.

**(f) Kullback-Leibler Divergence / Jaccard Coefficient**

I don't think using Kullback-Leibler Divergence is a good choice since the model is discrete and some of the p(x) will be zero which will lead mistake when taking log. Besides, I also don't think using Kullback-Leibler Divergence is a good choice since the attribute is binary while it's not appropriate to separate Midterm and Final in this way.

- KL:

$$D_{kL}(p(x)||q(x)) = \sum q(x) * ln\frac{p(x)}{q(x)}$$

- Jaccard:

$$Jaccard = \frac{|a \cap b|}{|A| + |B| - |A \cap B|)}$$

# Problem 4

**(a) X^2 correlation value:**

$$X^2 = \sum_{i=1}^{n} \frac{(O_i - E_i)^2}{E_i}$$

```
 1  # 4(a) x^2 correlation
 2  def getX2Correlation(array):
 3      length = array.shape[0]
 4
 5      sum_row_0 = array[0][0] + array[0][1]
 6      sum_row_1 = array[1][0] + array[1][1]
 7      sum_col_0 = array[0][0] + array[1][0]
 8      sum_col_1 = array[0][1] + array[1][1]
 9      sum_total = sum_row_0 + sum_row_1
10
11      new_array = np.zeros([2,2])
12      new_array[0][0] = sum_row_0 * sum_col_0 / sum_total
13      new_array[0][1] = sum_row_0 * sum_col_1 / sum_total
14      new_array[1][0] = sum_row_1 * sum_col_0 / sum_total
15      new_array[1][1] = sum_row_1 * sum_col_1 / sum_total
16
17      sum = 0
18      for i in range(length):
19          for j in range(length):
20              sum = sum + pow((array[i][j] - new_array[i][j]), 2) / ne
21      return round(sum, 3)
22
23  X2Correlation = getX2Correlation(Purchase)
24  print("X2 Correlation:", X2Correlation)
```

X2 Correlation: 2062.333

**(b)**

We can reject the null hypothesis of independence at a confidence level of 0.001.
Therefore, "purchasing beer" and "purchasing diaper" are correlated.

**(c)**

p = [200 : 80+20 : 3000] = [2 : 1 : 30] = [0.061 : 0.030 : 0.909]

**(d) Kullback-Leibler Divergence**

```
1  # 4(d) Kullback-Leibler Divergence
2  def getKL(p, q):
3      length = len(p)
4      sum = 0
5      for iter in range(length):
6          sum = sum + q[iter] *np.log(q[iter]/p[iter])
7
8      return round(sum, 3)
9
10 p = np.array([0.061, 0.030, 0.909])
11 q = np.array([0.5, 0.3, 0.2])
12 KL = getKL(p, q)
13 print("Kullback Leibler Divergence:", KL)
```

Kullback Leibler Divergence: 1.44

$$D_{kL}(p(x)\|q(x)) = \sum q(x) * ln\frac{p(x)}{q(x)}$$

p = [0.061 0.030 0.909]

```
q = [0.5 0.3 0.2 ]

KL = 1.440
```