# Question 1

a) $$E(X^2) = E(Y^2) = Var(X) + E(X)^2$$

$$= \frac{|b-a|^2}{12} + \frac{b+a}{2}$$

$$= \frac{1}{12} + \frac{1}{2} = \frac{13}{12}$$

$$E(Z) = E(X^2 - 2XY + Y^2)$$

$$= E(X^2) - 2E(X)E(Y) + E(Y)^2$$

$$= \frac{13}{2} - 2 + \frac{13}{2} = \frac{1}{6}$$

$$Var(X^2) = Var(Y^2)$$

$$= E(X^4) - E(X^2)^2$$

By fourth moment of mgf

$$E(X^4) = \frac{9}{5}$$

$$= \frac{9}{5} - \frac{13}{12}^2 = \frac{451}{720}$$

$$Var(X) = \frac{(b-a)^2}{12} = \frac{1}{12}$$

$$Var(Z) = Var\left((X-Y)^2\right)$$

$$= Var\left(X^2 - 2XY + Y^2\right)$$

$$= Var(X^2) + 4Var(XY) + Var(Y^2)$$

$\because$ X and Y are independent

$\therefore Var(XY) = Var(X)Var(Y)$

$$= Var(X^2) + 4Var(X)Var(Y) + Var(Y^2)$$

$$= \frac{451}{720} + 4 \cdot \frac{1}{12} \cdot \frac{1}{12} + \frac{451}{720}$$

$$= \frac{461}{360}$$

**b)**
$$E[R] = E[Z_1 + \ldots + Z_d]$$

$$\because \text{ since } X_1, Y_1 \ldots X_d, Y_d \sim \text{Uniform} (0, 1$$

$$\therefore E(|X_1 - Y_1|^2) = \ldots = E(|X_d - Y_d|^2)$$

$$\therefore E(Z_1) = \ldots = E(Z_d)$$

$$= d\, E(Z)$$

$$= d \cdot \frac{1}{6} = \frac{d}{6}$$

$$\text{Var}[R] = \text{Var}[Z_1 + \ldots + Z_d]$$

$$\because \text{Var}[Z_1] = \ldots = \text{Var}[Z_d]$$

$$\therefore \text{Var}[Z_1 + \ldots + Z_d] = \text{Var}(d Z_1)$$

$$= d^2 \text{Var}[Z_1]$$

)

$$= d^2 \frac{461}{720}$$

d)

$$\sigma = \sqrt{d^2 \frac{461}{720}} = d\sqrt{\frac{461}{720}}$$

$$\mu = d \cdot \frac{1}{6}$$

$$\because \frac{1}{6} < \sqrt{\frac{461}{720}} \approx 0.8$$

$$\therefore \mu < \sigma$$

As dimension $d$ increase, $\sigma$ increase faster than $\mu$.

maximum distance $= d \, E(z)$

Therefore as $d$ increase,

average distance distance increase,

so every points get further away.

But they get relative same distance,

which seems to be closer.

Question 2

a) $H(X) = \sum\limits_{x} P(x) \log_2\left(\frac{1}{P(x)}\right)$

$\because \ 0 \leq P(x) \leq 1 \ ,$

$\therefore \ \frac{1}{P(x)} \geq 1$

$\therefore \ \log_2\left(\frac{1}{P(x)}\right) \geq \log_2(1)$

$\Rightarrow \log_2\left(\frac{1}{P(x)}\right) \geq 0$

b) $\quad H(X,Y) = H(X) + H(Y|X)$

$\quad\quad$ by properties $\ H(Y|X) = H(Y)$

$\quad H(X,Y) = H(X) + H(Y)$

c) $\quad H(X,Y) = -\sum\limits_{x \in X} \sum\limits_{y \in Y} P(x,y) \log_2 P(x,y)$

$\quad\quad = -\sum\limits_{x \in X} \sum\limits_{y \in Y} P(x,y) \log_2 P(x) P(y|x)$

$\quad\quad = -\sum\limits_{x \in X} \sum\limits_{y \in Y} P(x,y) \log_2 P(x) - \sum\limits_{x \in X} \sum\limits_{y \in Y} P(x,y) \log_2 P(y|x)$

$\quad\quad = -\sum\limits_{x \in X} P(x) \log_2 P(x) - \sum\limits_{x \in X} \sum\limits_{y \in Y} P(x,y) \log_2 P(y|x)$

$\quad\quad = H(X) + H(Y|X)$

d)

$$KL(p||q) = \int p(x) \log_2\left(\frac{p(x)}{q(x)}\right) dx$$

$$= -\int p(x) \log_2\left(\frac{q(x)}{p(x)}\right) dx$$

$$= -E\left(\log_2 \frac{q}{p}\right)$$

by Jensen's inequality

$$E\left(\log_2 \frac{q}{p}\right) \geq \log_2\left(E\left(\frac{q}{p}\right)\right)$$

$$-E\left(\log_2 \frac{q}{p}\right) < -\log_2\left(E\left(\frac{q}{p}\right)\right)$$

∵ Since

∴ $-\log_2\left(\sum p(x) \frac{q(x)}{p(x)}\right) = 0$

∴ $-E\left(\log_2 \frac{q}{p}\right) < 0$

∴ $E\left(\log_2 \frac{q}{p}\right) \geq 0$

e)

$$KL(p(x,y) || p(x)p(y)) = \sum p(x,y) \log_2\left(\frac{p(x,y)}{p(x)p(y)}\right)$$

$$= \sum p(x,y) \log_2\left(\frac{p(y) p(x|y)}{p(x)p(y)}\right)$$

$$= \sum p(x,y) \log_2\left(\frac{p(x|y)}{p(x)}\right)$$

$$= \sum p(x,y) \log_2\left(\frac{1}{p(x)}\right) - \sum p(x,y) \log_2\left(\frac{1}{p(x|y)}\right)$$

$$= H(X) - H(X|Y)$$

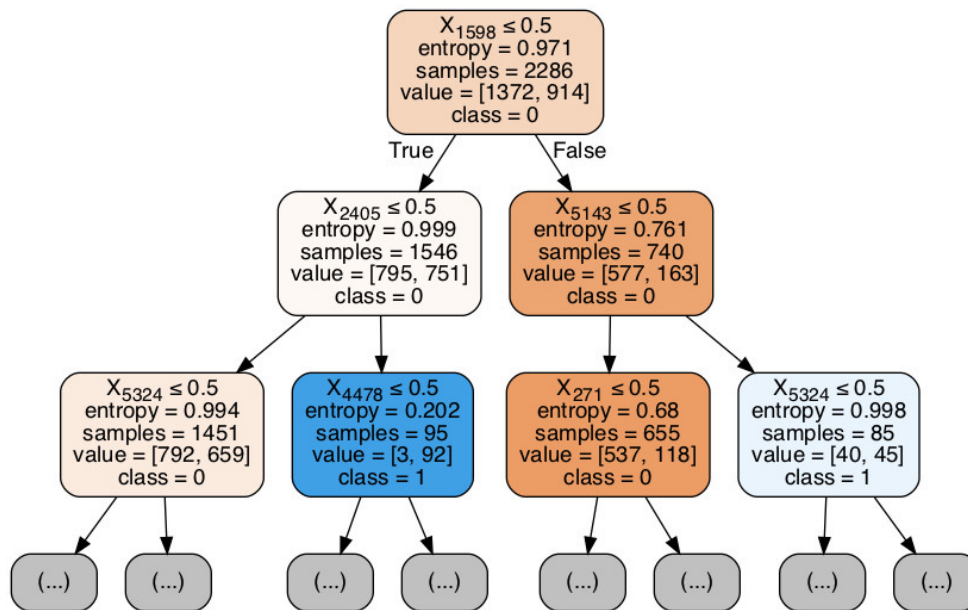$$= I(X;Y)$$

# Question 3

### b)

```
In [87]: models = [gini10, gini30, gini50, gini100, gini150, entro10, entro30, entro50, entro100, entro150]
         scores = []
         for model in models:
             score = model.score(x_validation_array, y_validation_array)
             scores.append(score)
             print(model.criterion, model.max_depth, 'Model Accuarcy:', score)
         best = max(scores)
         best_model = models[scores.index(best)]
         print('Best Model :', best_model)
         print("Best Model Accuracy ", test_accuracy(best_knn_model, x_test_array, y_test_array))
```

```
gini 10 Model Accuarcy: 0.7081632653061225
gini 30 Model Accuarcy: 0.7653061224489796
gini 50 Model Accuarcy: 0.7714285714285715
gini 100 Model Accuarcy: 0.753061224489796
gini 150 Model Accuarcy: 0.7714285714285715
entropy 10 Model Accuarcy: 0.7204081632653061
entropy 30 Model Accuarcy: 0.7612244897959184
entropy 50 Model Accuarcy: 0.7673469387755102
entropy 100 Model Accuarcy: 0.7775510204081633
entropy 150 Model Accuarcy: 0.7673469387755102
Best Model : DecisionTreeClassifier(class_weight=None, criterion='entropy', max_depth=100,
                       max_features=None, max_leaf_nodes=None,
                       min_impurity_decrease=0.0, min_impurity_split=None,
                       min_samples_leaf=1, min_samples_split=2,
                       min_weight_fraction_leaf=0.0, presort=False,
                       random_state=None, splitter='best')
Best Model Accuracy  0.5346938775510204
```

c)

```
X_1598 ≤ 0.5
entropy = 0.971
samples = 2286
value = [1372, 914]
class = 0
```

True                    False

```
X_2405 ≤ 0.5
entropy = 0.999
samples = 1546
value = [795, 751]
class = 0
```

```
X_5143 ≤ 0.5
entropy = 0.761
samples = 740
value = [577, 163]
class = 0
```

```
X_5324 ≤ 0.5
entropy = 0.994
samples = 1451
value = [792, 659]
class = 0
```

```
X_4478 ≤ 0.5
entropy = 0.202
samples = 95
value = [3, 92]
class = 1
```

```
X_271 ≤ 0.5
entropy = 0.68
samples = 655
value = [537, 118]
class = 0
```

```
X_5324 ≤ 0.5
entropy = 0.998
samples = 85
value = [40, 45]
class = 1
```

(...)  (...)    (...)  (...)    (...)  (...)    (...)  (...)

```
In [83]: models = [gini10, gini30, gini50, gini100, gini150, entro10, entro30, entro50, entro100, entro150]
         scores = []
         for model in models:
             score = model.score(x_validation_array, y_validation_array)
             scores.append(score)
         best = max(scores)
         best_model = models[scores.index(best)]
         print(best_model)
         print("Best Model Accuracy ", test_accuracy(best_model, x_test_array, y_test_array))

DecisionTreeClassifier(class_weight=None, criterion='entropy', max_depth=100,
                       max_features=None, max_leaf_nodes=None,
                       min_impurity_decrease=0.0, min_impurity_split=None,
                       min_samples_leaf=1, min_samples_split=2,
                       min_weight_fraction_leaf=0.0, presort=False,
                       random_state=None, splitter='best')
Best Model Accuracy  0.5102040816326531
```
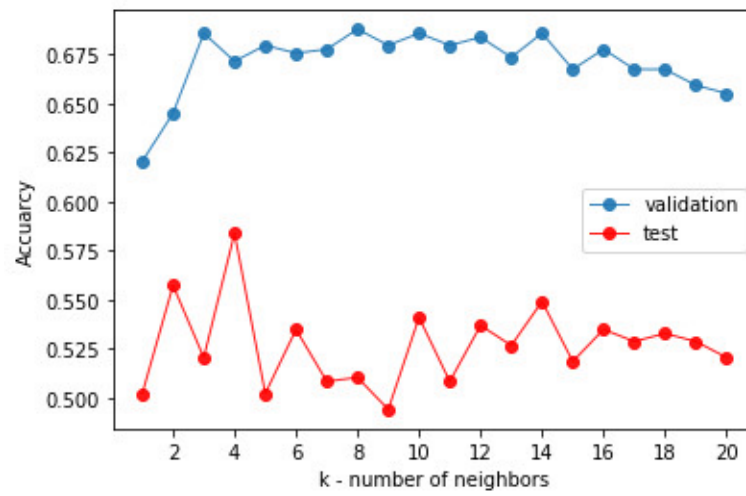
d)

```
In [213]: print(compute_info_gain(vectorizer.inverse_transform(x_train),y_train_array, 'the')) #Topmost
          print(compute_info_gain(vectorizer.inverse_transform(x_train),y_train_array, 'hillary'))
          print(compute_info_gain(vectorizer.inverse_transform(x_train),y_train_array, 'donald'))

          0.05330262393898977
          0.040765700073505995
          0.05010238863021219
```

In [ ]:

e)



```
In [74]: print(best_knn_model)
         print("Best Model Accuracy ", test_accuracy(best_knn_model, x_test_array, y_test_array))

         KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
                              metric_params=None, n_jobs=None, n_neighbors=17, p=2,
                              weights='uniform')
         Best Model Accuracy  0.5285714285714286
```