Carson McLean
CSC 2515
HW 1
2018-09-26

**①ⓐ**

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

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

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

$$= E(X^2) - (2)(E(X)E(Y)) + E Y^2 \qquad \text{2 is constant}$$
$$\text{X, Y independent}$$

$$= \int_0^1 x^2 \frac{1}{1-0} - \left(2 \cdot \int_0^1 x \frac{1}{1-0} \cdot \int_0^1 Y \frac{1}{1-0}\right) + \int_0^1 Y^2 \frac{1}{1-0}$$

$$= \frac{1}{3} - \left(2 \cdot \frac{1}{2} \cdot \frac{1}{2}\right) + \frac{1}{3} = \frac{2}{3} - \frac{1}{2} = \boxed{\frac{1}{6}}$$

$$E(X) = \int_a^b x \, f_x(x) \, dx \qquad \text{Continuous random variables}$$

$$E(g(X)) = \int_a^b g(x) f_x(x) \, dx$$

$$\text{Uniform PDF} = \frac{1}{b-a} = f(x)$$

$$Var(Z) = E(Z^2) - (E(Z))^2 = E(X^4 - 4X^3Y + 6X^2Y^2 - 4XY^3 + Y^4) - \left(\frac{1}{6}\right)^2$$

$$= \left[\int_0^1 x^4 \frac{1}{1-0} - \left(4 \cdot \int_0^1 x^3 \frac{1}{1-0} \cdot \int_0^1 Y \frac{1}{1-0}\right) + \left(6 \cdot \int_0^1 x^2 \frac{1}{1-0} \cdot \int_0^1 Y^2 \frac{1}{1-0}\right) - \left(4 \cdot \int_0^1 x \frac{1}{1-0} \int_0^1 Y^3 \frac{1}{1-0}\right) + \int_0^1 Y^4 \frac{1}{1-0}\right] - \frac{1}{36}$$

$$= \frac{1}{5} - \left(4 \cdot \frac{1}{4} \cdot \frac{1}{2}\right) + \left(6 \cdot \frac{1}{3} \cdot \frac{1}{3}\right) - \left(4 \cdot \frac{1}{2} \cdot \frac{1}{4}\right) + \frac{1}{5} - \frac{1}{36}$$

$$= \frac{1}{5} - \frac{1}{2} + \frac{2}{3} - \frac{1}{2} + \frac{1}{5} - \frac{1}{36} = \frac{1}{15} - \frac{1}{36} = \boxed{\frac{7}{180}}$$

(1)(6)

$$= \sum_{i=1}^{d} Z_i =$$

$$R = Z_1 + \cdots + Z_d = dZ \qquad \text{where } Z \text{ is a random variable as seen in } (1)(a) \; Z_i = (X_i - Y_i)^2$$

$$E(R) = E(dZ) = d\,E(Z) \quad \text{by removing the constant } d \qquad \text{From } (1)a, \; E(Z) = \tfrac{1}{6}$$

$$= \boxed{d/6}$$

$$Var(R) = Var(Z_1 + \cdots + Z_d) = Var\left(\sum_{i=1}^{d} Z_i\right)$$

By Bienaymé Formula

$$Var\left(\sum_{i=1}^{d} Z_i\right) = \sum_{i=1}^{d} Var(Z_i) = d\,Var(Z_i) \qquad \text{From } (1)(a), \; Var(Z) = 7/180$$

$$= \boxed{\dfrac{d7}{180}}$$

# Q2b

```
########################################
Max Depth: 1 // Split Criteria: entropy
0.6959183673469388
Max Depth: 1 // Split Criteria: gini
0.6959183673469388
Max Depth: 4 // Split Criteria: entropy
0.7877551020408163
Max Depth: 4 // Split Criteria: gini
0.7816326530612245
Max Depth: 8 // Split Criteria: entropy
0.8244897959183674
Max Depth: 8 // Split Criteria: gini
0.8163265306122449
Max Depth: 12 // Split Criteria: entropy
0.8183673469387756
Max Depth: 12 // Split Criteria: gini
0.8040816326530612
Max Depth: 16 // Split Criteria: entropy
0.810204081632653
Max Depth: 16 // Split Criteria: gini
0.789795918367347


########################################
Best Results:
Max Depth: 8 // Split Criteria: entropy
0.8244897959183674
```
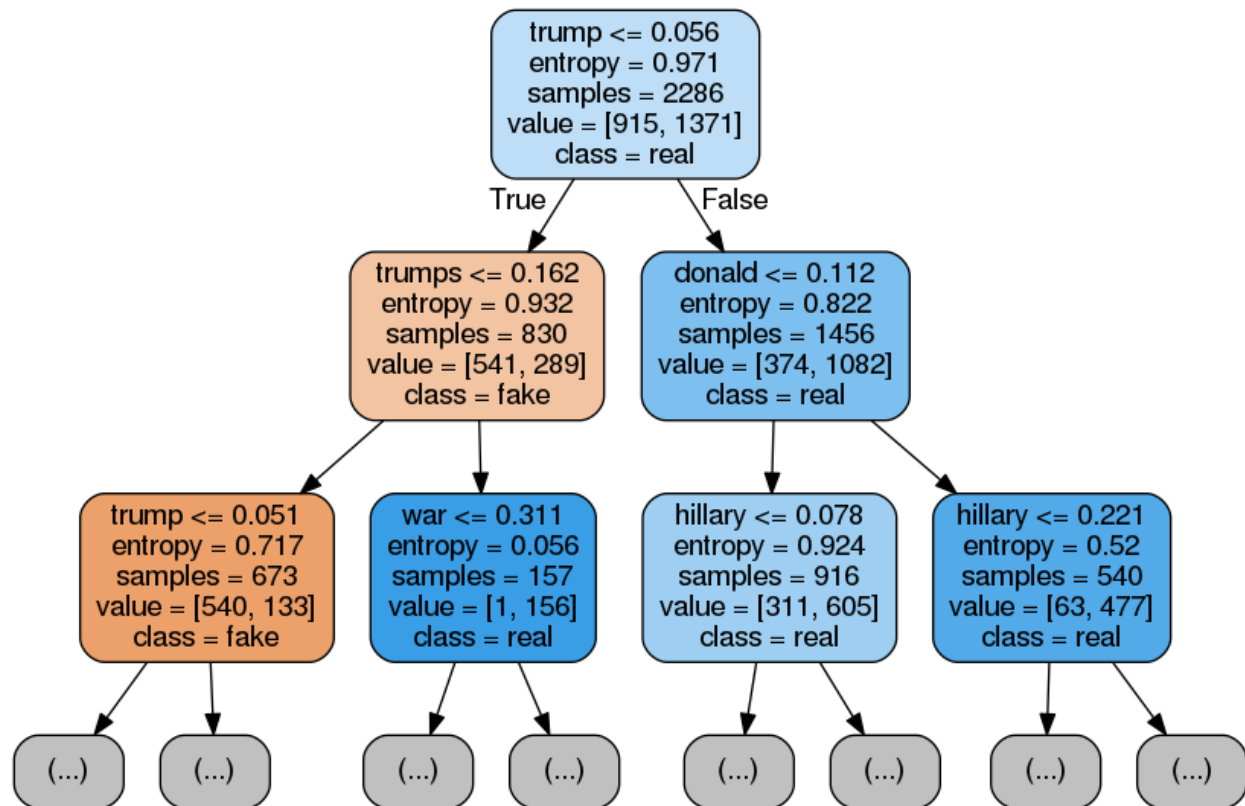
# Q2c



# Q2d

**Top Most Split**

xi = trump
Information Gain: 0.032391364179196636

**Other Keywords**

xi = trumps
Information Gain: 0.04736671065028819

xi = donald
Information Gain: 0.05227450569883951

xi = energy
Information Gain: 0.0006377983877460247

xi = hillary
Information Gain: 0.04058134341841224

xi = canada
Information Gain: 0.00027362801976738016

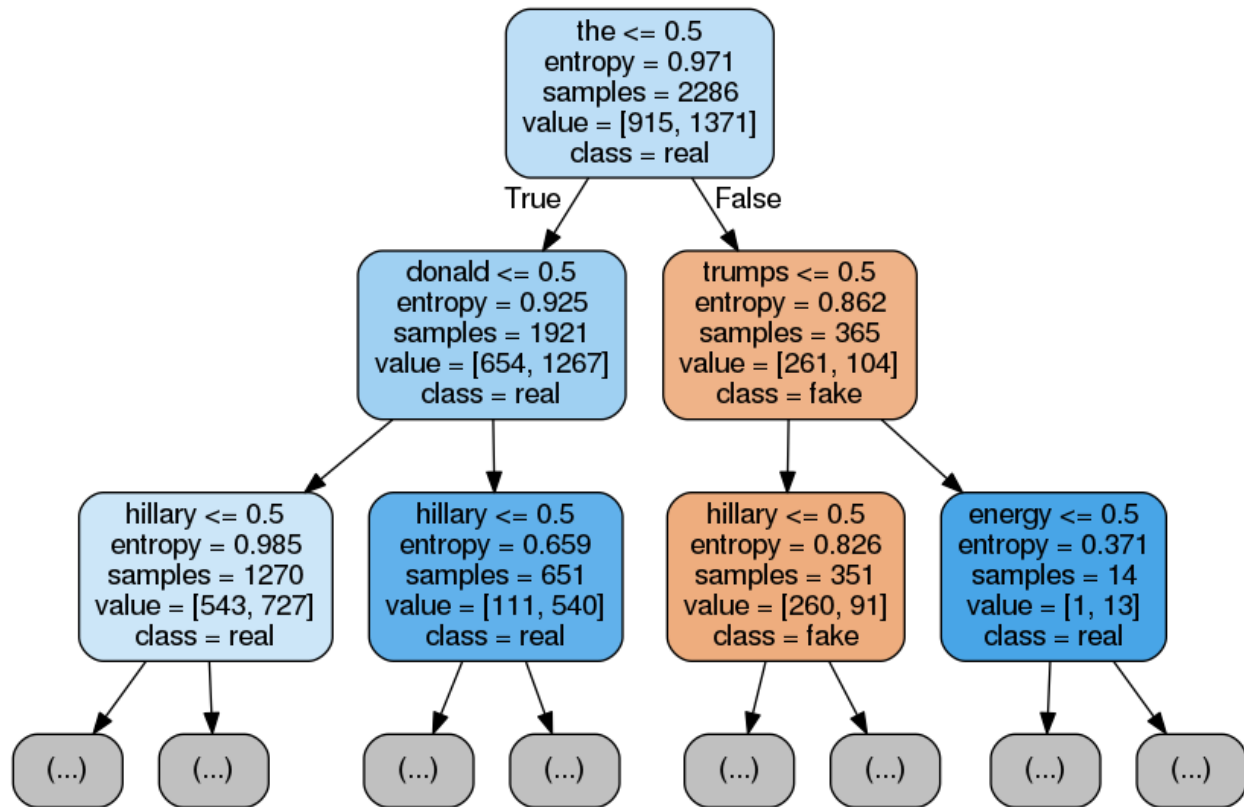xi = somemadeupword
Information Gain: 0.0

xi = the
Information Gain: 0.055942512911258624

xi = cat
Information Gain: 0.0

**Explanation**

In my code, I use a TfidfVectorizer() for the Classifier which in turn is used to create the Tree diagram. Yet, IG is calculated in Q2d using whether or not the $x_i$ keyword is present or not. Therefore, there is a divergence in how the graph diagram gets generated, and which node appears at the top, versus how IG is calculated. This results in the root node word for Tfidf (ie, 'trump') to appear to not have the highest IG score. In a Decision Tree, we would expect that the root node have the highest IG so that we maximize the split and better fit the data.

So, in my code as provided, simply switching to a CountVectorizer() [line 32-33], which is more closely in line to the IG keyword present vs absent style calculation, we see the following new Tree graph:

**Top Most Split**

xi = the
Information Gain: 0.055942512911258624

**Other Keywords**

xi = trump
Information Gain: 0.032391364179196636

xi = trumps
Information Gain: 0.04736671065028819

xi = donald
Information Gain: 0.05227450569883951

xi = energy
Information Gain: 0.0006377983877460247

xi = hillary

Information Gain: 0.04058134341841224

xi = canada
Information Gain: 0.00027362801976738016

xi = somemadeupword
Information Gain: 0.0

xi = cat
Information Gain: 0.0

So, with a CountVectorizer(), the IG calculation is in line with what is seen in the decision tree graph. However, I have decided to leave Tfidf as the default. There was no requirement to use one or the other, and Tfidf achieves ~80% accuracy, while Count is 75%. Further, 'the' is a stop word that perhaps should have been removed.