# Q1

## a)



| l2 distance mean | l2 distance SD |
| --- | --- |
| l1 distance mean | l1 distance SD |

## b)

$$E[R] = E[z_1 + \ldots + z_d]$$
$$= E[z_1] + \ldots + E[z_d]$$
$$= d \times \left(\frac{1}{6}\right)$$
$$= d/6$$

Since $X_i$, $Y_i$ are independently sampled for each $i$, $z_i$ and $z_j$ are independent for $i \neq j$

$$\Rightarrow Var[R] = Var[z_1 + \ldots + z_d]$$

$$= \text{Var}[z_1] + \ldots + \text{Var}[z_d]$$

$$= d \times \frac{7}{180}$$

$$= 7d/180$$

c) i) Let $R$ be the Euclidian distance

then $E$ is: $R - E[R] \leq d$

ii) $P(R - E[R] \leq d) = 1 - P(R - E[R] > d)$

$$= 1 - P(R - E[R] \geq d)$$

(since $E$ is a conti. random variable)

$$\geq 1 - \frac{\text{Var}[R]}{d^2}$$

iii) then $P(E) \geq 1 - \frac{7d/180}{d^2}$

$$= 1 - \frac{7}{180d}$$

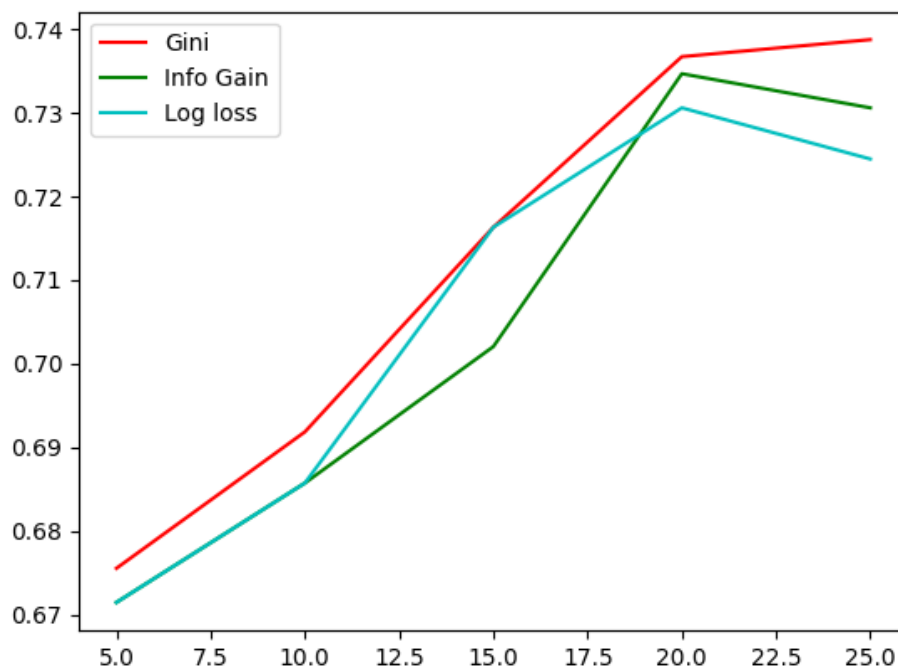$$\Rightarrow \lim_{d \to \infty} P(E) \geq \lim_{d \to \infty} 1 - \frac{7}{180d} = 1$$

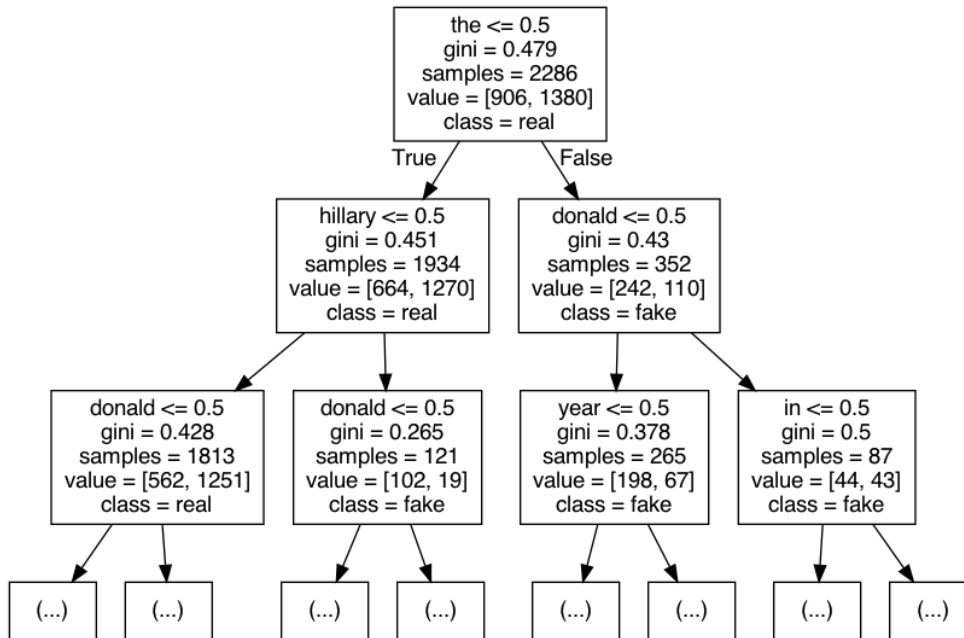so as $d \to \infty$, $P(E) = 1$, so any distance is $d$ away from its expectation

## Q2.  b) Function output:

```
Gini: score = 0.6755102040816326, depth = 5
Information gain: score = 0.6714285714285714, depth = 5
Log loss: score = 0.6714285714285714, depth = 5
Gini: score = 0.6918367346938775, depth = 10
Information gain: score = 0.6857142857142857, depth = 10
Log loss: score = 0.6857142857142857, depth = 10
Gini: score = 0.7163265306122449, depth = 15
Information gain: score = 0.7020408163265306, depth = 15
Log loss: score = 0.7163265306122449, depth = 15
Gini: score = 0.736734693877551, depth = 20
Information gain: score = 0.7346938775510204, depth = 20
Log loss: score = 0.7306122448979592, depth = 20
Gini: score = 0.7387755102040816, depth = 25
Information gain: score = 0.7306122448979592, depth = 25
Log loss: score = 0.7244897959183674, depth = 25
```

## Plot:

c) Gini w/ depth 25 achieved the highest accuracy

```
                        the <= 0.5
                        gini = 0.479
                        samples = 2286
                        value = [906, 1380]
                        class = real
                 True                    False
        hillary <= 0.5                        donald <= 0.5
        gini = 0.451                          gini = 0.43
        samples = 1934                        samples = 352
        value = [664, 1270]                   value = [242, 110]
        class = real                          class = fake

  donald <= 0.5      donald <= 0.5      year <= 0.5       in <= 0.5
  gini = 0.428       gini = 0.265       gini = 0.378      gini = 0.5
  samples = 1813     samples = 121      samples = 265     samples = 87
  value = [562,1251] value = [102, 19]  value = [198, 67] value = [44, 43]
  class = real       class = fake       class = fake      class = fake

  (...)    (...)     (...)    (...)     (...)    (...)     (...)    (...)
```

d) The keywords are selected from
    {"the", "hillary", "trumps", "donald"}

Their IG are as follows:

```
IG(Y|X) is 0.04570772617653496 for the keyword the
IG(Y|X) is 0.0426824963366705 for the keyword hillary
IG(Y|X) is 0.0371178553210577 for the keyword trumps
IG(Y|X) is 0.0419742322376332 for the keyword donald
```

## Q3.

**a)**

$$\frac{\partial J}{\partial w_{j'}} = \frac{1}{2N} \cdot \frac{\partial \left( \sum_{i=1}^{N} \left( y^{(i)} - t^{(i)} \right)^2 \right)}{\partial w_{j'}} \quad\longrightarrow\; = \sum_{j=1}^{D} w_j x_j^{(i)} + b$$

$$= \frac{1}{2N} \cdot \frac{\partial \left( \sum_{i=1}^{N} \left( \sum_{j=1}^{D} w_j x_j^{(i)} + b - t^{(i)} \right)^2 \right)}{\partial w_{j'}}$$

$$= \frac{1}{2N} \cdot 2 \sum_{i=1}^{N} \left( \sum_{j=1}^{D} w_j x_j^{(i)} + b - t^{(i)} \right) \left( x_{j'}^{(i)} \right)$$

$$= \frac{1}{N} \cdot \sum_{i=1}^{N} \underbrace{\left( \sum_{j=1}^{D} w_j x_j^{(i)} + b - t^{(i)} \right)}_{= \, y^{(i)}} \left( x_{j'}^{(i)} \right)$$

$$= \frac{1}{N} \cdot \sum_{i=1}^{N} \left( y^{(i)} - t^{(i)} \right) \left( x_{j'}^{(i)} \right) \qquad (1)$$

$$\frac{\partial R}{\partial w_{j'}} = \frac{\partial \left( \sum_{j=1}^{D} \alpha_j |w_j| + \frac{1}{2} \sum_{j=1}^{D} \beta_j \cdot w_j^2 \right)}{\partial w_{j'}} \quad\longrightarrow\; := f(w_j) = \begin{cases} w_j, & w_j > 0 \\ 0, & w_j = 0 \\ -w_j, & w_j < 0 \end{cases}$$

$$= \begin{cases} \text{if } w_{j'} > 0: & \dfrac{\partial \left( \sum_{j=1}^{D} \alpha_j w_j + \frac{1}{2} \sum_{j=1}^{D} \beta_j w_j^2 \right)}{\partial w_{j'}} \\[4mm] & \quad = \alpha_{j'} + \beta_{j'} w_{j'} \qquad (2) \\[4mm] \text{if } w_{j'} = 0: & \quad = 0 \qquad (3) \\[4mm] \text{if } w_{j} < 0: & \partial \left( -\sum_{j=1}^{D} \alpha_j w_j + \frac{1}{2} \sum_{j=1}^{D} \beta_j w_j^2 \right) \Big/ \partial w_{j'} \end{cases}$$

$$= -d_{j'} + \beta_{j'} w_{j'} \qquad (4)$$

So $\dfrac{\partial J_{reg}^{\alpha/\beta}(w)}{\partial w_{j'}}$ is:

if $w_{j'} = 0$ : $\dfrac{1}{N} \cdot \sum\limits_{i=1}^{N} (y^{(i)} - t^{(i)})(x_{j'}^{(i)})$

if $w_{j'} > 0$ : $\dfrac{1}{N} \cdot \sum\limits_{i=1}^{N} (y^{(i)} - t^{(i)})(x_{j'}^{(i)}) + d_{j'} + \beta_{j'} w_{j'}$

if $w_{j'} < 0$ : $\dfrac{1}{N} \cdot \sum\limits_{i=1}^{N} (y^{(i)} - t^{(i)})(x_{j'}^{(i)}) - d_{j'} + \beta_{j'} w_{j'}$

$$\frac{\partial J}{\partial b} = \frac{1}{2N} \cdot \partial \left[ \sum_{i=1}^{N} \left( y^{(i)} - t^{(i)} \right)^2 \right] / \partial b$$

$$\searrow \sum_{j=1}^{D} w_j x_j + b$$

$$= \frac{1}{2N} \cdot \sum_{i=1}^{N} 2 \left( \sum_{j=1}^{D} w_j x_j + b - t^{(i)} \right)$$

$$= \frac{1}{N} \cdot \sum_{i=1}^{N} (y^{(i)} - t^{(i)})$$

let $a > 0$ be the learning rate, so overall:

if $w_{j'} > 0$ : $w_{j'} \leftarrow w_{j'} - \alpha \left( \dfrac{1}{N} \cdot \sum\limits_{i=1}^{N} (y^{(i)} - t^{(i)})(x_{j'}^{(i)}) + d_{j'} + \beta_{j'} w_{j'} \right)$

$\Leftrightarrow w_{j'} \leftarrow w_{j'} - \dfrac{\alpha}{N} \sum\limits_{i=1}^{N} (y^{(i)} - t^{(i)})(x_{j'}^{(i)}) - d_{j'} \alpha - \alpha \beta_{j'} w_{j'}$

$$\Leftrightarrow \quad w_j{}' \leftarrow w_j{}'(1-\alpha\beta_j{}') - \frac{\alpha}{N}\sum_{j=1}^{N}(y^{(i)}-t^{(i)})(x_j{}'^{(i)}) - \alpha_j{}'\alpha$$

$$b \leftarrow b - \frac{\alpha}{N}\sum_{j=1}^{N}(y^{(i)}-t^{(i)})$$

if $w_j{}'=0$: $\quad w_j{}' \leftarrow w_j{}' - \alpha\left(\frac{1}{N}\cdot\sum_{j=1}^{N}(y^{(i)}-t^{(i)})(x_j{}'^{(i)})\right)$

$$\Leftrightarrow \quad w_j{}' \leftarrow w_j{}' - \frac{\alpha}{N}\sum_{j=1}^{N}(y^{(i)}-t^{(i)})(x_j{}'^{(i)})$$

$$b \leftarrow b - \frac{\alpha}{N}\sum_{j=1}^{N}(y^{(i)}-t^{(i)})$$

if $w_j{}'<0$: $\quad w_j{}' \leftarrow w_j{}' - \alpha\left(\frac{1}{N}\cdot\sum_{j=1}^{N}(y^{(i)}-t^{(i)})(x_j{}'^{(i)}) - \alpha_j{}' + \beta_j{}'w_j{}'\right)$

$$\Leftrightarrow \quad w_j{}' \leftarrow (1-\alpha\beta_j{}')w_j{}' - \frac{\alpha}{N}\cdot\sum_{j=1}^{N}(y^{(i)}-t^{(i)})(x_j{}'^{(i)}) + \alpha\alpha_j{}'$$

$$b \leftarrow b - \frac{\alpha}{N}\sum_{j=1}^{N}(y^{(i)}-t^{(i)})$$

This is called weight decay possibly because, for cases $w_j{}'<0$ and $w_j{}'>0$, the update rule for $w_j{}'$ contains the term $(1-\alpha\beta_j{}')\,w_j{}'$.

$\alpha>0$ and $\beta_j{}' \geq 0$, so $(1-\alpha\beta_j{}')\,w_j{}' \leq w_j{}'$

$\Rightarrow$ within the update rule, the weight $w_j{}'$ decays to a

lesser term.

b) $\lambda_1 = 0$

$$\Rightarrow J_{reg}^{\beta}(\omega) = \frac{1}{2N} \cdot \sum_{i=1}^{N} \left( y^{(i)} - t^{(i)} \right)^2 + \frac{1}{2} \sum_{j=1}^{D} \beta_j \omega_j^2$$

Note: for the sake of consistency $\overset{\text{to } 3\alpha)}{\smile}$, in my notations I
swapped $j$ and $j'$ defined in the question

$$\Rightarrow \frac{\partial J_{reg}^{\beta}(\omega)}{\partial \omega_{j'}} = \frac{1}{2N} \cdot \sum_{i=1}^{N} 2 \left( y^{(i)} - t^{(i)} \right) \left( x_{j'}^{(i)} \right) + \beta_{j'} \omega_{j'}$$

$$= \frac{1}{N} \cdot \sum_{i=1}^{N} \underbrace{\left( \sum_{j=1}^{D} \left( \omega_j x_j^{(i)} \right) - t^{(i)} \right) \left( x_{j'}^{(i)} \right)}_{} + \beta_{j'} \omega_{j'}$$

$$= \sum_{i=1}^{N} \sum_{j=1}^{D} \left( \omega_j x_j^{(i)} \right) x_{j'}^{(i)} - \sum_{i=1}^{N} t^{(i)} x_{j'}^{(i)}$$

$$= \sum_{j=1}^{D} \sum_{i=1}^{N} \frac{1}{N} \omega_j x_j^{(i)} x_{j'}^{(i)} - \frac{1}{N} \sum_{i=1}^{N} t^{(i)} x_{j'}^{(i)} + \beta_{j'} \omega_{j'}$$

Define the indicator function $I(j) = \begin{cases} 0 & \text{if } j \neq j' \\ 1 & \text{if } j = j' \end{cases}$

$$= \sum_{j=1}^{D} \sum_{i=1}^{N} \frac{1}{N} \omega_j x_j^{(i)} x_{j'}^{(i)} + I(j) \beta_j \omega_j - \frac{1}{N} \sum_{j=1}^{N} t^{(i)} x_{j'}^{(i)}$$

$$= \sum_{j=1}^{D} \left( \frac{1}{N} \sum_{i=1}^{N} x_j^{(i)} x_{j'}^{(i)} + I(j) \beta_j \right) \omega_j - \frac{1}{N} \sum_{j=1}^{N} t^{(i)} x_{j'}^{(i)}$$

So $A_{jj'} = \frac{1}{N} \sum\limits_{i=1}^{N} x_j^{(i)} x_{j'}^{(i)} + I(j) B_j$

$C_{j'} = \frac{1}{N} \sum\limits_{i=1}^{N} t^{(i)} x_{j'}^{(i)}$

c) Note that $X = \begin{pmatrix} x_1^{(1)} & x_2^{(1)} & \cdots & x_D^{(1)} \\ \vdots & \vdots & & \vdots \\ x_1^{(N)} & x_2^{(N)} & \cdots & x_D^{(N)} \end{pmatrix}$

Then $A = \begin{pmatrix} \frac{1}{N}\sum\limits_{i=1}^{N} x_1^{(i)} x_1^{(i)} & \cdots & \frac{1}{N}\sum\limits_{i=1}^{N} x_1^{(i)} x_D^{(i)} \\ \vdots & \ddots & \\ \frac{1}{N}\sum\limits_{i=1}^{N} x_D^{(i)} x_1^{(i)} & & \frac{1}{N}\sum\limits_{i=1}^{N} x_D^{(i)} x_D^{(i)} \end{pmatrix} + \begin{pmatrix} B_1 & & O \\ & \ddots & \\ O & & B_D \end{pmatrix}$

$= \frac{1}{N} \begin{pmatrix} x_1^{(1)} x_1^{(1)} + \cdots + x_1^{(N)} x_1^{(N)} & \cdots & x_1^{(1)} x_D^{(1)} + \cdots + x_1^{(N)} x_D^{(N)} \\ \vdots & \ddots & \\ x_D^{(1)} x_1^{(1)} + \cdots + x_D^{(N)} x_1^{(N)} & \cdots & x_D^{(1)} x_D^{(1)} + \cdots + x_D^{(N)} x_D^{(N)} \end{pmatrix}$

Let $\vec{x_i} = \begin{pmatrix} x_i^{(1)} \\ \vdots \\ x_i^{(N)} \end{pmatrix}$:

$= \frac{1}{N} \begin{pmatrix} \vec{x_1} \cdot \vec{x_1} & \cdots & \vec{x_1} \cdot \vec{x_D} \\ \vdots & \ddots & \\ \vec{x_D} \cdot \vec{x_1} & & \vec{x_D} \cdot \vec{x_D} \end{pmatrix}$

$$= \frac{1}{N} \begin{pmatrix} -\vec{x}_1- \\ \vdots \\ -\vec{x}_D- \end{pmatrix} \begin{pmatrix} | & & | \\ \vec{x}_1 & \cdots & \vec{x}_D \\ | & & | \end{pmatrix}$$

$$\underbrace{\phantom{\begin{pmatrix} -\vec{x}_1- \\ \vdots \\ -\vec{x}_D- \end{pmatrix}}}_{X^T} \underbrace{\phantom{\begin{pmatrix} | & & | \\ \vec{x}_1 & \cdots & \vec{x}_D \\ | & & | \end{pmatrix}}}_{X}$$

$$= \frac{1}{N} X^T X + \begin{pmatrix} \beta_1 & & 0 \\ & \ddots & \\ 0 & & \beta_D \end{pmatrix}$$

$$C = \frac{1}{N} \begin{pmatrix} x_1^{(1)} & \cdots & x_1^{(N)} \\ \vdots & \ddots & \vdots \\ x_D^{(1)} & \cdots & x_D^{(N)} \end{pmatrix} \begin{pmatrix} t^{(1)} \\ \vdots \\ t^{(N)} \end{pmatrix}$$

$$= \frac{1}{N} X^T \vec{t} \quad \text{for target vector } \vec{t} = \begin{pmatrix} t^{(1)} \\ \vdots \\ t^{(N)} \end{pmatrix}$$

then $A\vec{w} - C = 0$

$$\Rightarrow \left[ \frac{1}{N} X^T X + \begin{pmatrix} \beta_1 & & 0 \\ & \ddots & \\ 0 & & \beta_D \end{pmatrix} \right] \vec{w} - \frac{1}{N} X^T \vec{t} = 0$$

$$\Rightarrow \left[ X^T X + N \begin{pmatrix} \beta_1 & & 0 \\ & \ddots & \\ 0 & & \beta_D \end{pmatrix} \right] \vec{w} = X^T \vec{t}$$

$$\Rightarrow \vec{w} = \left[ X^T X + N \begin{pmatrix} \beta_1 & & 0 \\ & \ddots & \\ 0 & & \beta_D \end{pmatrix} \right]^{-1} X^T \vec{t}$$

assuming that $X^T X + N \begin{pmatrix} \beta_1 & & 0 \\ & \ddots & \\ 0 & & \beta_D \end{pmatrix}$ is invertible