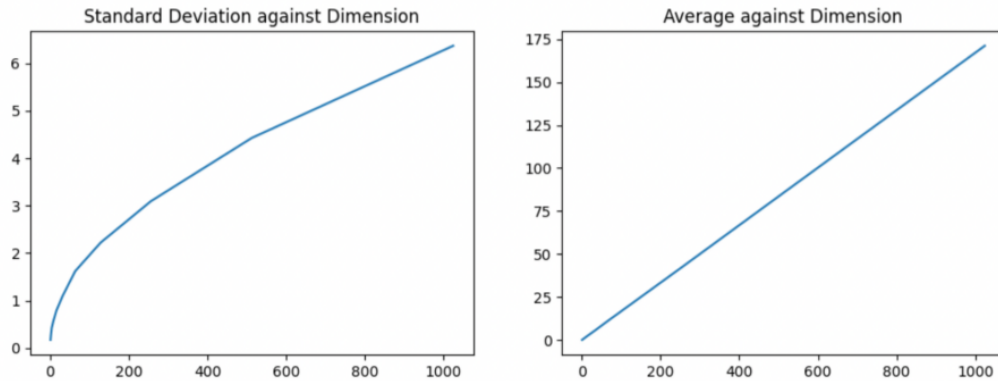


Question 1 a)



Question 1 b)

let $R = Z_1 + Z_2 + \dots + Z_n$ where $Z_i = (X_i - Y_i)^2$

notice that for Z_i, Z_j for $1 \leq i, j \leq n$
 Z_i and Z_j are independent, because

the difference of distance on one dimension
is independent on another dimension

we know $E[Z_i] = \frac{1}{6}$, $\text{Var}[Z_i] = \frac{7}{180}$

by linearity: $E[R] = E[Z_1 + \dots + Z_d] = E[Z_1] + \dots + E[Z_d]$

$$= 6 E[Z_i] = \frac{d}{6}$$

by independence: $\text{Var}[R] = \text{Var}[Z_1 + \dots + Z_d]$

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

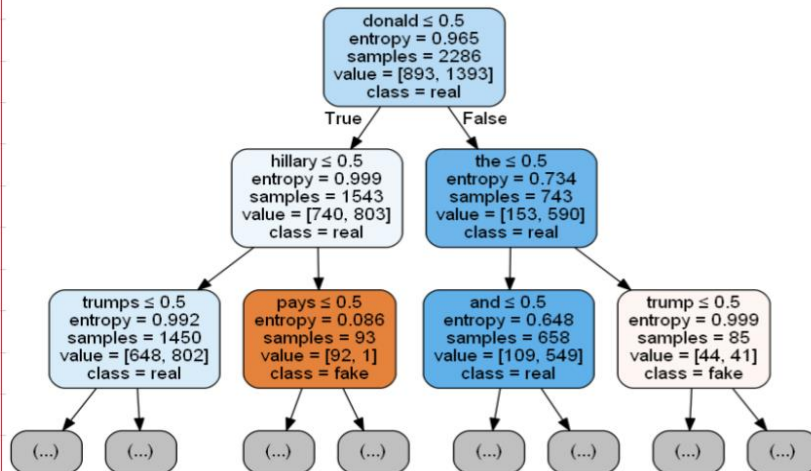
$$= d \text{Var}[Z_i]$$

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

Question 2b)

decision tree classifier using max_depth 5 and split criteria information gain has accuracy 0.7
decision tree classifier using max_depth 5 and split criteria Gini coefficient has accuracy 0.7061224489795919
decision tree classifier using max_depth 10 and split criteria information gain has accuracy 0.7081632653061225
decision tree classifier using max_depth 10 and split criteria Gini coefficient has accuracy 0.7183673469387755
decision tree classifier using max_depth 50 and split criteria information gain has accuracy 0.7591836734693878
decision tree classifier using max_depth 50 and split criteria Gini coefficient has accuracy 0.7489795918367347
decision tree classifier using max_depth 100 and split criteria information gain has accuracy 0.7612244897959184
decision tree classifier using max_depth 100 and split criteria Gini coefficient has accuracy 0.7428571428571429
decision tree classifier using max_depth 150 and split criteria information gain has accuracy 0.7755102040816326
decision tree classifier using max_depth 150 and split criteria Gini coefficient has accuracy 0.7673469387755102

2c)



2d)

```
feature_names, data_train, label_train, data_valid, label_valid, data_test, label_test = load_data()
select_model(feature_names, data_train, label_train, data_valid, label_valid, data_test, label_test)
print(
    f'information gain at root split donald is {compute_information_gain(data_train, label_train, "donald", feature_names)}')
print(
    f'information gain at other keyword the is {compute_information_gain(data_train, label_train, "the", feature_names)}')
print(
    f'information gain at other keyword hillary is {compute_information_gain(data_train, label_train, "hillary", feature_names)}')
print(
    f'information gain at other keyword trumps is {compute_information_gain(data_train, label_train, "trumps", feature_names)}')
```

```
C:\Users\wangy\anaconda3\python.exe "C:/uoft-academic/2021 fall/CSC311/hw1/a1.py"
```

```
information gain at root split donald is 0.052603317747226375
```

```
information gain at other keyword the is 0.047175200175596954
```

```
information gain at other keyword hillary is 0.04127665563376459
```

```
information gain at other keyword trumps is 0.04340669920550855
```

Question 3

a) know that

$$J_{\text{reg}}^B(w) = \frac{1}{2N} \sum_{i=1}^N (y^{(i)} - t^{(i)})^2 + \frac{1}{2} \sum_{j=1}^D \beta_j w_j^2$$

and $y = \sum_{j=1}^D w_j x_j^{(i)} + b$ (2)

substitute $y^{(i)}$ in $J_{\text{reg}}^B(w)$ with (2)

$$J_{\text{reg}}^B(w) = \frac{1}{2N} \sum_{i=1}^N \left(\left(\sum_{j=1}^D w_j x_j^{(i)} + b \right) - t^{(i)} \right)^2 + \frac{1}{2} \sum_{j=1}^D \beta_j w_j^2$$

break $J_{\text{reg}}^B(w) = J(w) + R(w)$,

then,

$$\frac{d}{dw_j} J_{\text{reg}}^B(w) = \frac{d}{dw_j} J(w) + \frac{d}{dw_j} R(w)$$

calculate the partial derivative w.r.p w_j for $1 \leq j \leq D$ for $J(w)$

$$\begin{aligned} \frac{d}{dw_j} J &= \frac{d}{dw_j} \left(\frac{1}{2N} \sum_{i=1}^N \left(\left(\sum_{j=1}^D w_j x_j^{(i)} + b \right) - t^{(i)} \right)^2 \right) \\ &= \frac{1}{N} \sum_{i=1}^N \left[\left(\sum_{j=1}^D w_j x_j^{(i)} + b - t^{(i)} \right) \left(\frac{d}{dw_j} \sum_{j=1}^D w_j x_j^{(i)} + \frac{d}{dw_j} b - \frac{d}{dw_j} t^{(i)} \right) \right] \\ &= \frac{1}{N} \sum_{i=1}^N \left[\left(\sum_{j=1}^D w_j x_j^{(i)} + b - t^{(i)} \right) \cdot \frac{d}{dw_j} \sum_{j=1}^D w_j x_j^{(i)} \right] \\ &= \frac{1}{N} \sum_{i=1}^N \left[\left(\sum_{j=1}^D w_j x_j^{(i)} + b - t^{(i)} \right) \cdot x_j^{(i)} \right] \end{aligned}$$

substitute $y^{(i)}$ back:

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

calculate the partial derivative w.r.p w_j for $1 \leq j \leq D$ for $R(w)$

$$\frac{d}{dw_j} R(w) = \frac{d}{dw_j} \left(\frac{1}{2} \sum_{j=1}^D \beta_j w_j^2 \right) = \beta_j w_j$$

then,

$$\frac{dJ_{\text{reg}}}{dw_j} = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - t^{(i)}) x_j^{(i)} + \beta_j w_j$$

According to the gradient descent

$$w_j \leftarrow w_j - \alpha \frac{dJ_{\text{reg}}}{dw_j}$$

$$w_j \leftarrow w_j - \alpha \frac{1}{N} \sum_{i=1}^N (y^{(i)} - t^{(i)}) x_j^{(i)} - \alpha \beta_j w_j$$

$$w_j \leftarrow (1 - \alpha \beta_j) w_j - \frac{\alpha}{N} \sum_{i=1}^N (y^{(i)} - t^{(i)}) x_j^{(i)}$$

calculate $\frac{dJ_{reg}}{db}$

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

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

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

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

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

According to the gradient descent

$$b \leftarrow b - \alpha \frac{dJ_{reg}}{db}$$

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

In conclusion: $w_j \leftarrow (1 - \alpha \beta_j) w_j - \frac{\alpha}{N} \sum_{i=1}^N (y^{(i)} - t^{(i)}) x_j^{(i)}$

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

Weight decay: According to weight update form, we rescale

the weights w_j by $(1 - \alpha \beta_j) w_j$ on each gradient descent.

it refers to weight decay as it makes weights smaller

b) from part a), we know $\frac{dJ_{reg}^B}{dw_j} = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - t^{(i)}) x_j^{(i)} + \beta_j w_j$

substitute $y^{(i)} = \sum_{j'=1}^D w_{j'} x_{j'}^{(i)}$

$$\frac{dJ_{reg}^B}{dw_j} = \frac{1}{N} \sum_{i=1}^N \left(\left(\sum_{j'=1}^D w_{j'} x_{j'}^{(i)} \right) - t^{(i)} \right) x_j^{(i)} + \beta_j w_j$$

multiply both side by N :

$$\sum_{i=1}^N \left[\left(\sum_{j'=1}^D w_{j'} x_{j'}^{(i)} \right) - t^{(i)} \right] x_j^{(i)} + \beta_j w_j = N \frac{dJ_{reg}^B}{dw_j}$$

$$\sum_{i=1}^N \left[\sum_{j'=1}^D w_{j'} x_{j'}^{(i)} x_j^{(i)} + \beta_j w_j \right] - \sum_{i=1}^N t^{(i)} x_j^{(i)} = N \frac{dJ_{reg}^B}{dw_j}$$

Expressing $\beta_j w_j$ as $\sum_{j'=1}^D w_{j'} \beta_j I_{j'=j}$: Here $I_{j'=j}$ is the indication function

$$\sum_{i=1}^N \left[\sum_{j'=1}^D \left(x_{j'}^{(i)} x_j^{(i)} + \beta_j I_{j'=j} \right) w_{j'} \right] - \sum_{i=1}^N t^{(i)} x_j^{(i)} = N \frac{dJ_{reg}^B}{dw_j}$$

exchange order of $\sum_{i=1}^N$ and $\sum_{j'=1}^D$:

$$\sum_{j'=1}^D \left[\sum_{i=1}^N \left(x_{j'}^{(i)} x_j^{(i)} + \beta_j I_{j'=j} \right) w_{j'} \right] - \sum_{i=1}^N t^{(i)} x_j^{(i)} = N \frac{dJ_{reg}^B}{dw_j}$$

$$\text{According to the form } \sum_{j'=1}^D A_{jj'} w_{j'} - c_j = \frac{dJ_{reg}^B}{dw_j}$$

We can see that

$$A_{jj'} = \frac{1}{N} \sum_{i=1}^N \left(x_{j'}^{(i)} x_j^{(i)} + \beta_j I_{j'=j} \right)$$

$$c_j = \frac{1}{N} \sum_{i=1}^N t^{(i)} x_j^{(i)}$$

c) since $C = \begin{bmatrix} \frac{1}{N} \sum_{i=1}^N t^{(1)} x_i^{(1)} \\ \vdots \\ \frac{1}{N} \sum_{i=1}^N t^{(i)} x_j^{(i)} \\ \vdots \\ \frac{1}{N} \sum_{i=1}^N t^{(i)} x_p^{(i)} \end{bmatrix}$ $X = \begin{bmatrix} x_1^{(1)} & \dots & x_j^{(1)} & \dots & x_p^{(1)} \\ \vdots & & \vdots & & \vdots \\ x_1^{(i)} & & x_j^{(i)} & & x_p^{(i)} \\ \vdots & & \vdots & & \vdots \\ x_1^{(N)} & & x_j^{(N)} & & x_p^{(N)} \end{bmatrix}$ $t = \begin{bmatrix} t^{(1)} \\ \vdots \\ t^{(i)} \\ \vdots \\ t^{(N)} \end{bmatrix}$

$$C = \frac{1}{N} X^T t \quad (\text{verified by matrix product})$$

similarly, first break $A_{jj'} = \frac{1}{N} \sum_{i=1}^N (x_j^{(i)} x_{j'}^{(i)})$
 $+ \frac{1}{N} \sum_{i=1}^N \beta_{j'} I_{j'=j}$

First notice that $\frac{1}{N} \sum_{i=1}^N (x_j^{(i)} x_{j'}^{(i)})$ is the dot product of j^{th} and j'^{th} col by aggregate all the i and j' which equal to $X^T X$

for the second part, $\frac{1}{N} \sum_{i=1}^N \beta_{j'} I_{j'=j} = \beta_{j'} I_{j'=j}$ if $j'=j$ by aggregate all i and j' which equal to $\text{diag}(\beta)$ i.e. a diagonal matrix on the j^{th} position in the diagonal is β_j .

add the two parts.

$$A = \frac{1}{N} X^T X + \text{diag}(\beta) \quad \text{and} \quad C = \frac{1}{N} X^T t$$

According to the form $\sum_{j=1}^n A_{jj'} w_{j'} - c_j = \frac{dJ_{reg}}{dw_j} = 0$

we know $Aw = C$

$$w = A^{-1}C$$

$$w = (X^T X + \lambda I \text{diag}(\beta))^{-1} X^T \cdot t$$