

---

---

---

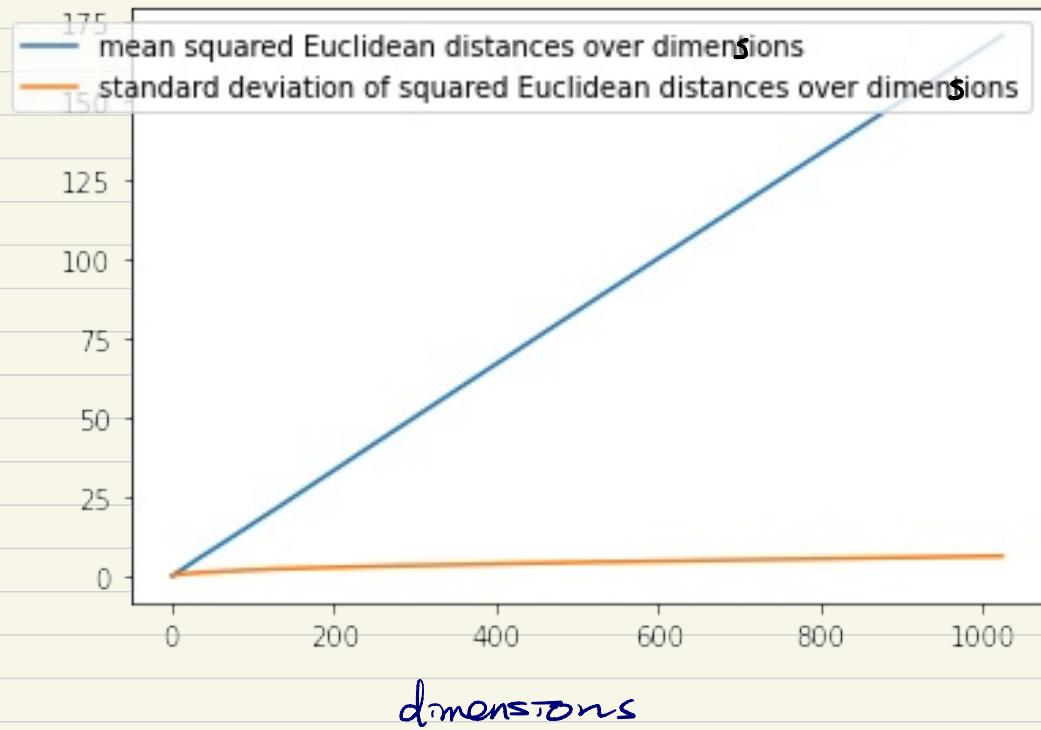
---

---



HW1

i) a)



b)

$$R = z_1 + \dots + z_d \quad ; \quad z_i = (x_i - y_i)^2$$

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$$

$$\mathbb{E}[R] = \mathbb{E}[z_1 + \dots + z_d]$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_d \end{bmatrix}$$

$$= \mathbb{E}[z_1] + \dots + \mathbb{E}[z_d]$$

$$= d \mathbb{E}[z_1] = \left(\frac{d}{6}\right)$$

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

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

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

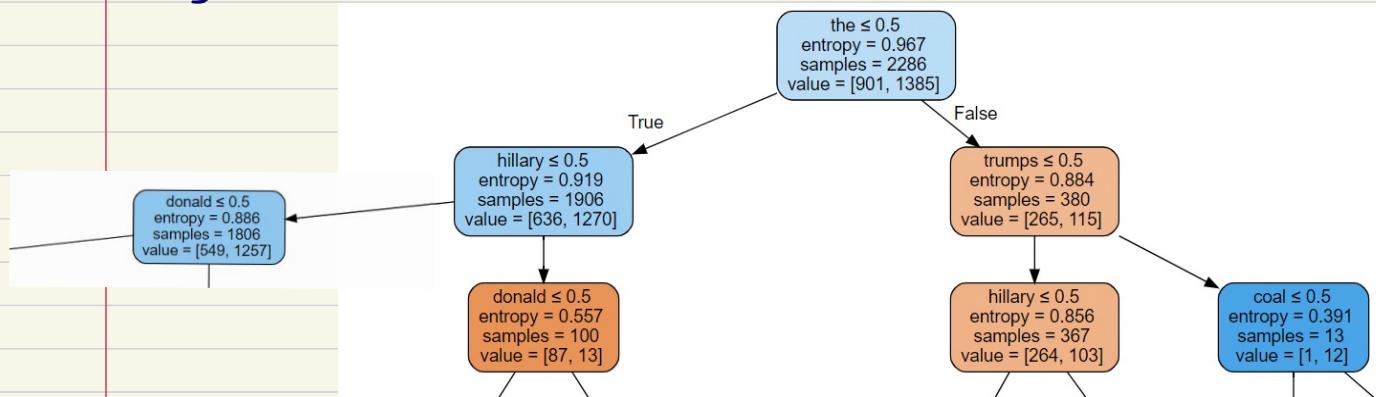
since  $z_i$ 's are independent

- 2) a) see submitted code .
- b) performance of mode with various  
max\_depth and split\_criteria .

	max_depth	split_criteria	score
0	31	entropy	0.773469
1	41	gini	0.771429
2	46	gini	0.769388
3	16	gini	0.765306
4	36	gini	0.765306
5	46	entropy	0.763265
6	36	entropy	0.761224
7	41	entropy	0.759184
8	26	entropy	0.757143
9	31	gini	0.757143
10	26	gini	0.755102
11	21	entropy	0.753061
12	21	gini	0.742857
13	11	gini	0.742857
14	16	entropy	0.738776
15	6	gini	0.716327
16	11	entropy	0.714286
17	6	entropy	0.700000
18	1	entropy	0.661224
19	1	gini	0.661224

The best model with 31 max depth  
and split-criteria 'entropy' had a 77.3%  
accuracy tested on validation set .

c)



d)

The top-most split is 'the' .

Information Gain of 'the' : 0.05445

Information Gain of 'hillary' : 0.043199...

Information Gain of 'trumps' : 0.04057...

Information Gain of ' donald' : 0.04694 ...

3)

$$a) L = \frac{1}{2} (y - t)^2$$

$$\begin{aligned} J(\omega, b) &= \frac{1}{2N} \sum_{i=1}^N (y^{(i)} - t^{(i)})^2 + \frac{1}{2} \sum_{j=1}^D B_j - \omega_j^2 \\ &= \frac{1}{2N} \sum_{j=1}^D \left[ (w_j x_j^{(i)}) + b - t^{(i)} \right]^2 + \frac{1}{2} \sum_{j=1}^D B_j - \omega_j^2 \end{aligned}$$

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

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

$$w_j \leftarrow w_j - \alpha \frac{\partial J}{\partial w_j}$$

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

$$= w_j - \alpha \left[ \frac{1}{N} \sum_{i=1}^N (y^{(i)} - t^{(i)}) x_j^{(i)} + B_j w_j \right]$$

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

$$= (-\alpha B_j)(w_j) - \alpha \sum_{i=1}^N (y^{(i)} - t^{(i)})$$

so  $\boxed{w_j \leftarrow (-\alpha B_j)(w_j) - \alpha \sum_{i=1}^N (y^{(i)} - t^{(i)})}$

so at each step, we update  $w_j$  based on the amount by which the estimates

$y^{(i)}$  differ from the target. However, we also try to decrease  $w_j$  by  $\alpha B_j w_j$  each time towards 0, hence 'weight decay'.

$$\begin{aligned}
 \text{b) } \frac{\partial J^B}{\partial w_j} &= \frac{1}{N} \sum_{i=1}^N \left( \sum_{j'=1}^D (w_{j'} x_{j'}^{(i)}) - t^{(i)} \right) (x_{j'}^{(i)}) + B_j w_j \\
 &= \frac{1}{N} \sum_{i=1}^N \left[ (x_j^{(i)} \sum_{j'=1}^D (w_{j'} x_{j'}^{(i)})) - x_j^{(i)} t^{(i)} \right] + (B_j w_j) \\
 &= \frac{1}{N} \sum_{i=1}^N \left[ \sum_{j'=1}^D w_{j'} x_{j'}^{(i)} x_j^{(i)} - x_j^{(i)} t^{(i)} \right] + (B_j w_j) \\
 &= \sum_{j'=1}^D \left[ \underbrace{\frac{1}{N} \sum_{i=1}^N (x_j^{(i)} x_{j'}^{(i)})}_{A_{jj'}} + \underbrace{\mathbb{I}(j' = j) B_j}_{- \frac{1}{N} \sum_{i=1}^N x_j^{(i)} t^{(i)}} \right] w_{j'} + c_j
 \end{aligned}$$

$$\text{so } A_{jj'} = \frac{1}{N} \sum_{i=1}^N (x_j^{(i)} x_{j'}^{(i)}) + \mathbb{I}(j' = j) B_j$$

$$\text{and } \bar{c}_j = \frac{1}{N} \sum_{i=1}^N x_j^{(i)} t^{(i)}$$

where  $\mathbb{I}$  is the indicator function that returns 1 when  $j' = j$  and 0 otherwise.

$$c) \left( \underbrace{\frac{1}{N} \mathbf{x}^T \mathbf{x} + \text{Diag}(\mathbf{B})}_{\mathbf{A}} \right) \vec{\omega} - \underbrace{\frac{1}{N} \mathbf{x}^T \vec{\mathbf{t}}}_{\vec{\mathbf{c}}} \\ \text{set } = 0$$

assuming that  $\mathbf{A}$  is invertible,

$$\Rightarrow \mathbf{A}^{-1} \vec{\omega} = \vec{\mathbf{c}}$$

$$\Rightarrow \vec{\omega} = \mathbf{A}^{-1} \vec{\mathbf{c}}$$

$$\Rightarrow \boxed{\vec{\omega} = \left( \frac{1}{N} \mathbf{x}^T \mathbf{x} + \text{Diag}(\mathbf{B}) \right)^{-1} \frac{1}{N} \mathbf{x}^T \vec{\mathbf{t}}}$$