

HW2 SEUNG EUN LEE

1. (a)

K	Accuracy
1	0.751847
5	0.754889
11	0.764885
21	0.746632
41	0.752282
61	0.737505
81	0.726641
101	0.728814
201	0.731421
401	0.719687

1.(b)

K	Accuracy
1	0.856150
5	0.870056
11	0.878748
21	0.884398
41	0.885267
61	0.882660
81	0.877445
101	0.875272
201	0.860061
401	0.839635

1.(c)

[illegible]

[illegible]

31	t32	spam	spam	spam	spam	spam	spam	spam	spam	no	no
32	t33	spam	spam	spam	spam	spam	no	no	no	no	no
33	t34	spam	spam	spam	spam	spam	no	no	no	no	no
34	t35	spam	spam	spam	spam	spam	spam	spam	spam	spam	spam
35	t36	spam	spam	spam	spam	spam	spam	spam	spam	spam	spam
36	t37	spam	spam	spam	spam	spam	spam	spam	spam	spam	spam
37	t38	spam	spam	spam	spam	spam	spam	spam	spam	spam	spam
38	t39	spam	spam	spam	spam	spam	spam	spam	spam	spam	spam
39	t40	no	no	no	no	no	no	no	no	no	no
40	t41	no	no	no	no	no	no	no	no	no	no
41	t42	spam	spam	spam	spam	spam	spam	spam	spam	no	no
42	t43	no	no	no	no	no	no	no	no	no	no
43	t44	no	no	no	no	no	no	no	no	no	no
44	t45	spam	spam	spam	spam	spam	spam	spam	spam	spam	spam
45	t46	spam	spam	spam	spam	spam	spam	spam	spam	spam	spam
46	t47	spam	spam	spam	spam	spam	spam	spam	spam	spam	spam
47	t48	spam	spam	spam	spam	spam	spam	spam	spam	spam	spam
48	t49	spam	spam	spam	spam	spam	spam	spam	spam	spam	spam
49	t50	spam	spam	spam	spam	spam	spam	spam	spam	spam	spam

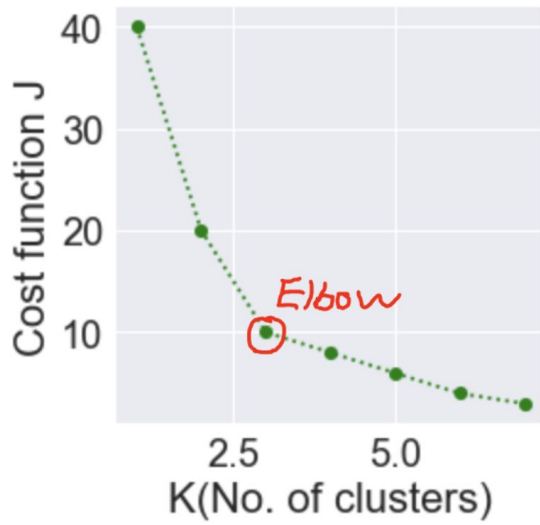
1.(d)

- Z-score normalization increased the accuracy than the not-normalized one.
- For (a) K=11, for (b) K=41 have the best accuracies

1.(e)

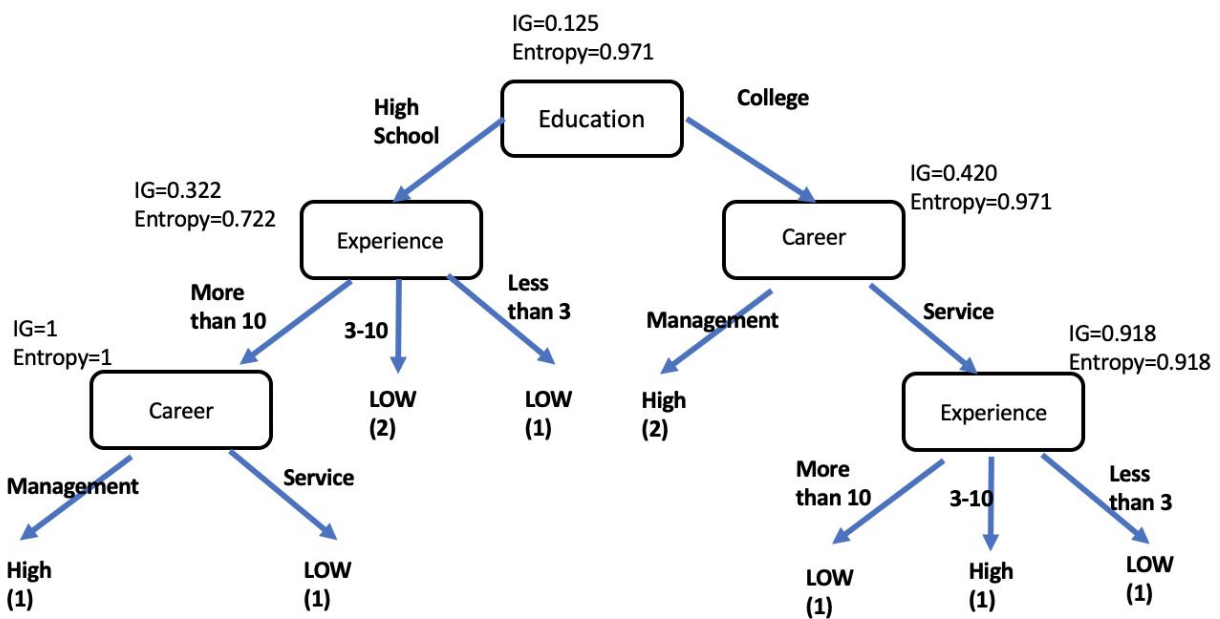
<Elbow method>

: Try various values of K, and if there is a point where the value changes significantly, where a point bends like an elbow bend. This point is judged as the optimal K value.



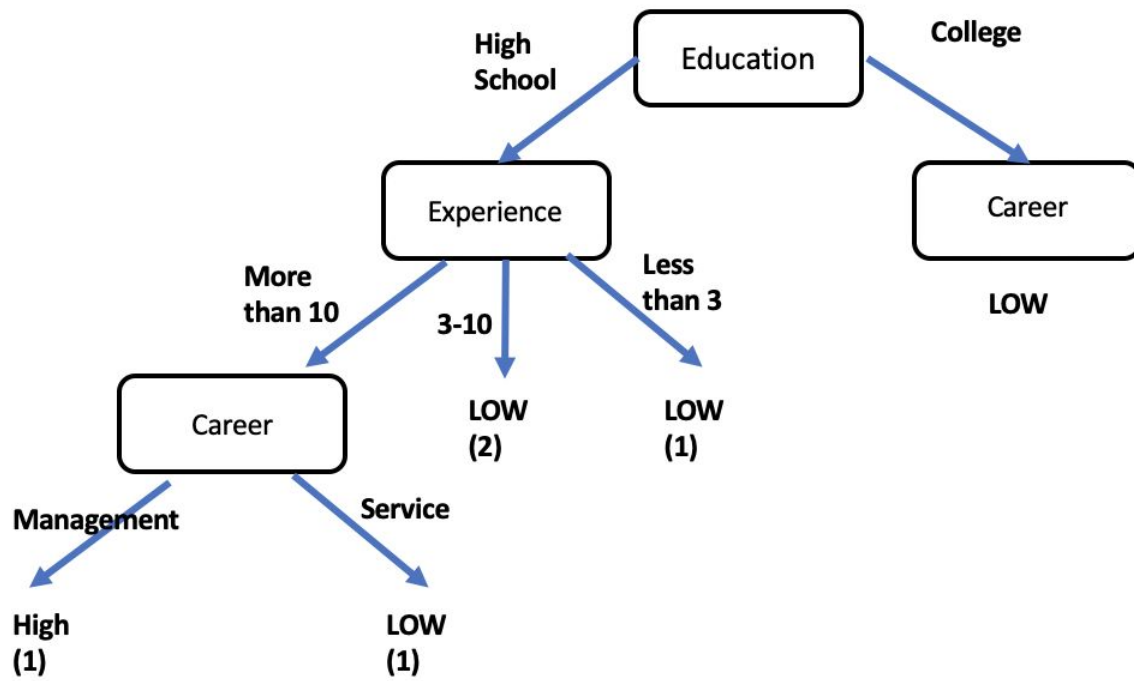
2.

<pre-prune>



<post-prune>

College >> low



3.

3.

i)

$$P(y = \text{Low} \mid x = \text{HighSchool, Service, } < 3)$$

$$= P(x = \text{H.S.} \mid y = \text{Low}) * P(x = \text{Service} \mid y = \text{Low}) * P(x = < 3 \mid y = \text{Low}) * P(y = \text{Low})$$

$$= \frac{4}{6} * \frac{4}{6} * \frac{2}{6} * \frac{6}{10}$$

< Laplace smoothing > = add - 1 smoothing

$$\frac{4+1}{6+2} * \frac{4+1}{6+2} * \frac{2+1}{6+3} * \frac{6}{10} = \underline{0.78} //$$

$$P(y = \text{High} \mid x = \text{H.S., Service, } < 3) = P(x = \text{H.S.} \mid y = \text{High}) * P(x = \text{Service} \mid y = \text{High}) * P(x = < 3 \mid y = \text{High}) * P(y = \text{High})$$

$$= \frac{1}{4} * \frac{1}{4} * \frac{1}{4} * \frac{4}{10}$$

< Laplace smoothing >

$$\Rightarrow \frac{1+1}{4+2} * \frac{1+1}{4+2} * \frac{1+1}{4+3} * \frac{4}{10} = \underline{0.13} //$$

∴ Low probability is high

∴ Predicted Class = Low //

$$3C(ii) \quad P(y = \text{low} | X = \text{College, Retail}, <3) = P(x = \text{College} | y = \text{low}) * \\ P(x = \text{Retail} | y = \text{low}) * P(x < 3 | y = \text{low}) * P(y = \text{low})$$

$$\Rightarrow \frac{2}{6} * \frac{0}{6} * \frac{2}{6} * \frac{6}{10}$$

<Laplace smoothing>

$$\frac{2+1}{6+2} * \frac{0+1}{6+3} * \frac{2+1}{6+3} * \frac{6}{10} = \underline{0.0083}$$

$$P(y = \text{high} | x = \text{College, Retail}, <3) = P(x = \text{College} | y = \text{high}) * P(x = \text{Retail} | y = \text{high}) \\ * P(x < 3 | y = \text{high}) * P(y = \text{high})$$

$$\Rightarrow \frac{3}{4} * \frac{0}{4} * \frac{1}{4} * \frac{4}{10}$$

<Laplace smoothing>

$$\frac{3+1}{4+2} * \frac{0+1}{4+3} * \frac{1+1}{4+3} * \frac{4}{10} = \underline{0.011}$$

\therefore High //

$$3(iii) P(y=low | x=graduate, service, 3-10) = P(x=graduate | y=low) * P(x=service | y=low) * P(x=3-10 | y=low) * P(y=low)$$

$$= \frac{0}{6} * \frac{4}{6} * \frac{2}{6} * \frac{6}{10}$$

<Laplace smoothing>

$$\frac{0+1}{6+3} * \frac{4+1}{6+2} * \frac{2+1}{6+3} * \frac{6}{10} = \underline{0.014},,$$

$$P(y=high | x=graduate, service, 3-10)$$

$$= P(x=graduate | y=high) * P(x=service | y=high) * P(x=3-10 | y=high) * P(y=high)$$

$$= \frac{0}{4} * \frac{1}{4} * \frac{1}{4} * \frac{4}{10}$$

<Laplace smoothing>

$$\frac{0+1}{4+3} * \frac{1+1}{4+2} * \frac{1+1}{4+3} * \frac{4}{10} = \underline{0.0054},,$$

∴ low //

4.

4. 25 students \rightarrow 60% accuracy

(a) take 3 models, M.V.C C3
C3's accuracy

error $\Rightarrow 0.4 \rightarrow p$

$$\sum_{k \geq n/2}^n \binom{n}{k} p^k (1-p)^{n-k}$$

$$\sum_{k \geq 1.5}^3 \binom{3}{k} (0.4)^k (1-0.4)^{3-k}$$

$$\Rightarrow \sum_{k=2}^3 \binom{3}{k} (0.4)^k (0.6)^{3-k}$$

$$\Rightarrow \binom{3}{2} 0.4^2 0.6^1 + \binom{3}{3} (0.4)^3 (0.6)^0$$

$$\Rightarrow 3(2 \cdot 0.4^2 \cdot 0.6^1 + 3 \cdot 0.4^3 \cdot 0.6^0)$$

$$\Rightarrow 3 \times (0.4)^2 (0.6) + (0.4)^3 (0.6)^0$$

$$= 0.352$$

$$\therefore \text{accuracy} = 0.648$$

$$(b) \sum_{k \geq 2.5}^5 \binom{5}{k} (0.4)^k (0.6)^{5-k}$$

$$\Rightarrow \sum_{k=3}^5 \binom{5}{k} (0.4)^k (0.6)^{5-k}$$

$$= 5(3(0.4)^3 (0.6)^2 + 5(4(0.4)^4 (0.6)^1 + 5(5(0.4)^5 (0.6)^0)$$

$$= 0.317$$

$$\therefore \text{accuracy} = 0.683$$

$$(c) \sum_{k \geq 25/2}^{25} \binom{25}{k} (0.4)^k (0.6)^{25-k} = 0.154$$

$$\therefore \text{accuracy} = 0.846$$

(d)

Assumption of every classifier has the same probability will cause difference from the reality.
Another assumption is that each classifier must be independent of each other.

(e)

Handwritten calculation on lined paper:

(c) accuracy = 45% \rightarrow $e = 0.55$

$$\sum_{k=12}^{25} \binom{25}{k} (0.55)^k (0.45)^{25-k}$$

$\Rightarrow 0.694$

accuracy = $\frac{0.306}{11}$

(simple program to calculate the error probability)

```
from scipy.special import comb
import math

def ensemble_error(n_classifier, error):
    k_start = int(math.ceil(n_classifier / 2.))
    probs = [comb(n_classifier, k) * error**k * (1-error)**(n_classifier - k)
              for k in range(k_start, n_classifier + 1)]
    return sum(probs)

x=ensemble_error(n_classifier=25, error=0.55)
print(x)
```