

1. Let's revisit the logic behind the voting method of classifier combination (used in Bagging, Random Forests, and Boosting to some extent). We are assuming that *the errors between the two classifiers are uncorrelated*

(a) First, let's assume our three independent classifiers both have an error rate of $e = 0.4$, calculated over 1000 instances with binary labels (500 A and 500 B).

(i) Build the confusion matrices for these classifiers, based on the assumptions above.

Actual class	#	P1	# for Sys1	P2	# for Sys 2	P3	# for Sys 3
A	500	A	(500 x 0.6=) 300	A	(300 x 0.6=) 180	A	(180*0.6=) 108
				B	(300 x 0.4=) 120	A	(120*0.6=) 72
		B	(500 x 0.4=) 200	A	(200 x 0.6=) 120	B	(120*0.4=) 48
		(200 x 0.4=) 80	A		(80*0.6=) 48		
		B	(500 x 0.4=) 200	B	(200 x 0.6=) 120	B	(80*0.4=) 32
					(200 x 0.4=) 80	A	(80*0.4=) 32
				B	(200 x 0.6=) 120	A	(120*0.4=) 48
B	500	A	(500 x 0.4=) 200	A	(200 x 0.4=) 80	B	(120*0.6=) 72
				B	(200 x 0.6=) 120	A	(120*0.4=) 48
		B	(500 x 0.6=) 300	A	(300 x 0.4=) 120	B	(120*0.6=) 72
				B	(300 x 0.6=) 180	A	(180*0.4=) 72
				B	(300 x 0.6=) 180	B	(180*0.6=) 108

A
A
A
B
A
B
B
A
A
A
B
A
B
B
B

(Same pred as above)

(ii) Using that the majority voting, what the expected error rate of the voting ensemble?

$$\text{For Actual A: } 48 + 48 + 48 + 32 = 176$$

$$\text{For Actual B: } 32 + 48 + 48 + 48 = 176$$

$$\text{Total error count} = 176 + 176 = 352$$

$$\text{Expected ER} = \frac{352}{1000} = 0.352 < 0.4 \text{ (individual ER)}$$

learners correct each others' mistake (assume errors uncorrelated)

Alternative :

$$\binom{3}{2} (0.4)^2 (0.6) + \binom{3}{3} (0.4)^3 = 0.352$$

two making mistakes all wrong

(b) Now consider three classifiers, first with $e_1 = 0.1$, the second and third with $e_2 = e_3 = 0.2$.

(i) Build the confusion matrices.

Actual		(e ₁ = 0.1)		(e ₂ = 0.2)		(e ₂ = 0.2)	
		Pred 1	Pred 2	Pred 1	Pred 2	Pred 3	
A	500	A	(500 x 0.9 =) 450	A	(450 x 0.8 =) 360	A	288
				B	(450 x 0.2 =) 90	B	72
		B	(500 x 0.1 =) 50	A	(50 x 0.8 =) 40	A	81
				B	(50 x 0.2 =) 10	B	18
	500	A	(500 x 0.1 =) 50	A	(50 x 0.2 =) 10	A	36
				B	(50 x 0.8 =) 40	B	8
		B	(500 x 0.9 =) 450	A	(450 x 0.2 =) 90	A	8
				B	(450 x 0.8 =) 360	B	2

(ii) Using the majority voting, what the expected error rate of the voting ensemble?

$$\text{Error count} = (18 + 8 + 8 + 2) + (2 + 8 + 8 + 18)$$

$$= 72$$

$$ER = \frac{72}{1000} = 0.072 < 0.1 \text{ (sys 1)}$$

Alternative :

$$ER = \underbrace{(0.1)(0.2)^2}_{\text{all wrong}} + \underbrace{(0.9)(0.2)^2}_{\text{S2 \& S3 wrong}} + \underbrace{2 \times (0.1)(0.2)(0.8)}_{\text{S1 \& S2/S3 wrong}}$$

$$= 0.072$$

- (iii) What if we relax our assumption of independent errors? In other words, what will happen if the errors between the systems were very highly correlated instead? (Systems make similar mistakes.)

Highly correlated errors \Rightarrow voting unlikely to improve the results
 (making same mistakes)

E.g. S_1 & S_2 correlated, S_3 uncorrelated

$\Rightarrow S_1$ & S_2 "out-vote" S_3

If S_1 & S_2 wrong \Rightarrow ensemble wrong

\Rightarrow Choose uncorrelated classifiers for ensembling.

2. Consider the following dataset:

<i>id</i>	<i>apple</i>	<i>ibm</i>	<i>lemon</i>	<i>sun</i>	<i>label</i>
A	4	0	1	1	fruit
B	5	0	5	2	fruit
C	2	5	0	0	comp
D	1	2	1	7	comp
E	2	0	3	1	?
F	1	0	1	0	?

- (a) Treat the problem as an unsupervised machine learning problem (excluding the *id* and *label* attributes), and calculate the clusters according to **k-means** with $k = 2$, using the Manhattan distance:

- (i) Starting with seeds A and D.

(a) Seed A, D ($k=2$)

\Rightarrow two initial centroids:

$$C_1 : \langle 4, 0, 1, 1 \rangle$$

$$C_2 : \langle 1, 2, 1, 7 \rangle$$

- ① Calculate distances from each instance to the centroids (of clusters)

$$d_M(A, C_1) = 0$$

$$d_M(A, C_2) = |4-1| + |0-2| + |1-1| + |1-7| = 11$$

}

$$d_M(B, C_1) = |5-4| + |0-0| + |5-1| + |2-1| = \underline{6}$$

$$d_M(B, C_2) = \underline{15}$$

$$\begin{aligned} d_M(C, C_1) &= 9 \\ d_M(C, C_2) &= 12 \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\}$$

$$\begin{aligned} d_M(D, C_1) &= 11 \\ d_M(D, C_2) &= 0 \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\}$$

$$\begin{aligned} d_M(E, C_1) &= 4 \\ d_M(E, C_2) &= 11 \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\}$$

$$\begin{aligned} d_M(F, C_1) &= 4 \\ d_M(F, C_2) &= 9 \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\}$$

② Assign each instance to the nearest cluster

$$A \rightarrow C_1 \quad B \rightarrow C_1 \quad C \rightarrow C_1 \quad D \rightarrow C_2$$

$$E \rightarrow C_1 \quad F \rightarrow C_1$$

$$C_1 : A, B, C, E, F \quad C_2 : D$$

③ Update centroids (average the instances in that cluster)

$$C_1 = \left\langle \frac{4+5+2+2+1}{5}, \frac{5}{5}, \frac{1+5+3+1}{5} \right\rangle,$$

$$\frac{1+2+1}{5} \rangle$$

$$= \langle 2.8, 1, 2, 0.8 \rangle$$

id	apple	ibm	lemon	sun	label
A	4	0	1	1	fruit
B	5	0	5	2	fruit
C	2	5	0	0	comp
D	1	2	1	7	comp
E	2	0	3	1	?
F	1	0	1	0	?

$$C_2 = \langle 1, 2, 1, 7 \rangle \quad (\text{Just "D"})$$

⇒ repeat! (with new C_1 & C_2)

C_2 hasn't change : reuse dist from last iter.

$$d_M(A, C_2) = 11$$

$$d_M(B, C_2) = 15$$

$$d_M(C, C_2) = 12$$

$$d_M(D, C_2) = 0$$

$$d_M(E, C_2) = 11$$

$$d_M(F, C_2) = 9$$

$C_1: A, B, C, E, F$ } final assignment
 $C_2: D$

Same as last iter! \Rightarrow Converged! \Rightarrow Stop!

(ii) Starting with seeds A and F.

Skip.

tie:

$d(E, C_1) = |2 - 4| + |0 - 0| + |3 - 1| + |1 - 1| = 4$
 $d(E, C_2) = |2 - 1| + |0 - 0| + |3 - 1| + |1 - 0| = 4$ } randomly pick one

(b) Perform agglomerative clustering of the above dataset (excluding the *id* and *label* attributes), using the Euclidean distance and calculating the group average as the cluster centroid.

In workshop slides.

3. Revise the concept of Unsupervised and Supervised evaluation for clustering evaluation

(a) Explain the two main concepts that we use to measure the how well do cluster labels match externally supplied class labels.

① Homogeneity

Measures if all the elements of each cluster have the same true label.

E.g. Entropy & Purity

If too many clusters with the same label \Rightarrow Homogeneity can't detect!

② Completeness

Measures if all the members of a class are assigned to the same cluster

(b) Explain the two main concepts that we use to measures the goodness of a clustering structure without respect to external information.

① Cohesiveness

Members of each cluster to be integrated and close to each other as possible.

e.g. W-SSE (within Cluster)

② Separability

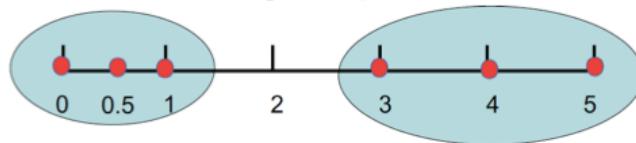
Clusters to be separate & independent as possible from the other clusters.

e.g. B-SSE (Between Cluster)

4. [OPTIONAL] Using the following dataset



- (a) If we develop 2 clusters: C1(A,B,C) and C2(D,E,F). Calculate the W-SSE and B-SSE for clusters, using Euclidean distance as the proximity function.



$$W_{SSE} = \sum_{i=1}^k \underbrace{\sum_{x \in C_i} (x - \mu_i)^2}_{\text{squared distance of the points to the centroid in } C_i}$$

$$B_{SSE} = \sum_{i=1}^k n_i (\bar{x} - \mu_i)^2$$

↑ overall mean ↑ mean of points in C_i

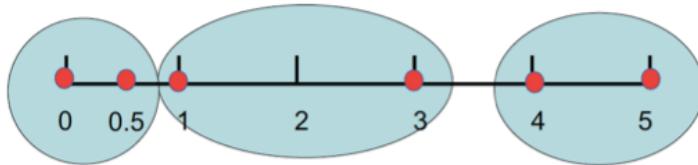
$$\Rightarrow W_{SSE} = [(0-0.5)^2 + (0.5-0.5)^2 + (1-0.5)^2] + \\ [(3-4)^2 + (4-4)^2 + (5-4)^2]$$

$$= 2.5$$

$$B_{SSE} = 3 \times (2.25 - 0.5)^2 + 3 \times (2.25 - 4)^2 \\ = 18.375$$

$$\frac{B_{SSE}}{W_{SSE}} = \frac{18.375}{2.5} = 7.35 \quad (\text{higher } \Rightarrow \text{better})$$

(b) Now consider developing 3 clusters as C1(A,B), C2(C,D) and C3(E,F). Calculate the W-SSE and B-SSE for these clusters.



$$\begin{aligned} W_{SSE} &= [(0 - 0.25)^2 + (0.5 - 0.25)^2] + [(1 - 2)^2 + (3 - 2)^2] \\ &\quad + [(4 - 4.5)^2 + (5 - 4.5)^2] \\ &= 2.625 \end{aligned}$$

$$\begin{aligned} B_{SSE} &= 2 \times (2.25 - 0.25)^2 + 2 \times (2.25 - 2)^2 + 2 \times (2.25 - 4.5)^2 \\ &= 18.25 \end{aligned}$$

$$\frac{B_{SSE}}{W_{SSE}} = \frac{18.25}{2.625} = 6.95$$

(c) Compare your results from parts (a) and (b) . Which one is a better clustering?

(a) better than (b) :

- ① Higher index
- ② Smaller W-SSE (cohesion)
- ③ Higher B-SSE (separability)