

Lecture 24

Thursday, November 30, 2023 10:24 AM

Clustering evaluation

→ External measures

$$\mathcal{C} = \{C_1, C_2, \dots, C_r\}$$

↙ Clustering using algo A

$$|C_i| = n_i$$

$r = \#$ of clusters chosen by the user

$$\left\{ \begin{array}{l} T = \{T_1, T_2, \dots, T_k\} \\ \text{ground truth partitioning} \\ |T_i| = m_i \end{array} \right.$$

$k = \#$ of 'true' clusters

$$|C_i \cap T_j| = n_{ij}$$

γ

	T_1	T_2	\dots	T_j	T_k
C_1	n_{11}	n_{12}			
C_2					
\vdots					
C_i				n_{ij}	
\vdots					
C_r					

$$T = \left\{ \begin{array}{l} \text{Iris data} \\ T_1 = \text{setosa} \quad 50 \\ T_2 = \text{virginica} \quad 50 \\ T_3 = \text{versicolor} \quad 50 \end{array} \right\}$$

$$C = \left\{ \begin{array}{l} C_1 = 75 \\ C_2 = 30 \\ C_3 = 45 \end{array} \right\}$$

$$r = 3$$

Purity

Purity

for each cluster $C_i \rightarrow$ find the max over the true partitions

$$Purity_i = \frac{1}{n_i} \max_{j=1}^k \{ n_{ij} \}$$

Precision_i
(accuracy C_i)

$$Purity = \sum_{i=1}^r \left(\frac{n_i}{n} \right) Purity_i$$

$$Purity = \frac{1}{n} \sum_{i=1}^r \max_{j=1}^k \{ n_{ij} \}$$

C_1

C_2

C_3

	T_1	T_2	T_3	T_4	T_5
C_1	50	10	2	100	0
C_2				100	
C_3				50	

$n_1 = 162$

$n_3 = 1150$

$$Prec(C_1) = \frac{100}{162}$$

$$Recall(C_1) = \frac{100}{1150}$$

F-score: balance precision & recall

$$Prec_i = Purity_i = \frac{1}{n_i} \max_{j=1}^k \{ n_{ij} \}$$

$$Recall_i = \frac{1}{n_{j^*}} \max_{j=1}^k \{ n_{ij} \}$$

$$j^* = \arg \max_{j=1}^k \{ n_{ij} \}$$

$$F_i = \text{harmonic mean} (Prec_i, Recall_i)$$

$$= \frac{2 \times Prec_i \times Recall_i}{Prec_i + Recall_i}$$

prec_i + recall_i

$$F = \frac{1}{r} \sum_{i=1}^r F_i \quad (\text{macro})$$

Maximum matching

	T ₁	T ₂	T ₃	
C ₁	<u>50</u>	100	0	150 = n ₁
C ₂	10	<u>1000</u>	50	1060 = n ₂
C ₃	100	90	<u>150</u>	340 = n ₃
	160	1190	200	1550
	m ₁	m ₂	m ₃	n



Purity₁ = $\frac{100}{150}$

Purity₂ = $\frac{1000}{1060}$

Purity₃ = $\frac{150}{340}$

weighted avg = purity

one to many
or many to one

$\frac{50}{150}$

$\frac{1000}{1060}$

$\frac{150}{340}$

weighted avg =
maximum matching

one-to-one

Information Theoretic

Entropy

r

$$\frac{n_i}{n} = p_i$$

$$H(\mathcal{C}) = - \sum_{i=1} P(c_i) \log P(c_i)$$

$$H(\mathcal{T}) = - \sum_{j=1}^k P(T_j) \log P(T_j)$$

$$P_{T_j} = \frac{n_j}{n}$$

Mutual information

$$I(\mathcal{C}, \mathcal{T}) = \sum_{i=1}^r \sum_{j=1}^k \frac{p_{ij}}{P(c_i) \cdot P(T_j)} \log \left(\frac{p_{ij}}{P(c_i) \cdot P(T_j)} \right)$$

Observed
joint
prob

joint
prob if
 c_i is
independent
of T_j

$$I(\mathcal{C}, \mathcal{T}) \in [0, \infty]$$

$$p_{ij} = \frac{n_{ij}}{n} = \frac{|c_i \cap T_j|}{n}$$

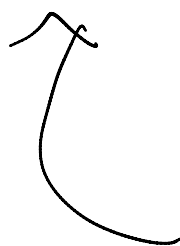
$$P(c_i) = \frac{n_i}{n} \quad P(T_j) = \frac{n_j}{n}$$

Normalized mutual information : $\in [0, 1]$

geometric mean of two ratios

$$\frac{I(\mathcal{C}, \mathcal{T})}{H(\mathcal{C})} \quad \text{and} \quad \frac{I(\mathcal{C}, \mathcal{T})}{H(\mathcal{T})}$$

$$NMI = \sqrt{\frac{I(\mathcal{C}, \mathcal{T})}{H(\mathcal{C})} \cdot \frac{I(\mathcal{C}, \mathcal{T})}{H(\mathcal{T})}} = \frac{I(\mathcal{C}, \mathcal{T})}{\sqrt{H(\mathcal{C}) \cdot H(\mathcal{T})}}$$



Confusion matrix

Confusion matrix

Confusion matrix
Contingency table

Pairwise 'labels'

$$T = \left\{ \begin{matrix} T_1 & T_2 & \dots & T_k \\ y_1 & y_2 & & y_k \end{matrix} \right\} \quad \Bigg| \quad \mathcal{C} = \left\{ \begin{matrix} C_1 & C_2 & \dots & C_r \\ \hat{y}_1 & \hat{y}_2 & \dots & \hat{y}_r \end{matrix} \right\}$$

\vec{x}_i and \vec{x}_j True Predicted

1) True Positives : $y_i = y_j$ and $\hat{y}_i = \hat{y}_j$
TP

2) True negatives : $y_i \neq y_j$ and $\hat{y}_i \neq \hat{y}_j$

3) False Positive : $y_i \neq y_j$ and $\hat{y}_i = \hat{y}_j$

4) False negative : $y_i = y_j$ and $\hat{y}_i \neq \hat{y}_j$

		True	
		P	N
Predicted	P	TP	FP
	N	FN	TN

$\left(\begin{smallmatrix} n \\ 2 \end{smallmatrix} \right)$ pairs = $\frac{n(n-1)}{2}$

Jaccard : $\frac{TP}{TP+FN+FP}$

FM : Fowlkes - Mallows

geometric mean of prec & recall

$$\text{prec} = \frac{TP}{TP+FP}$$

$$\text{recall} = \frac{TP}{TP+FN}$$

$$\sqrt{\frac{TP}{(TP+FP)} \cdot \frac{TP}{(TP+FN)}}$$

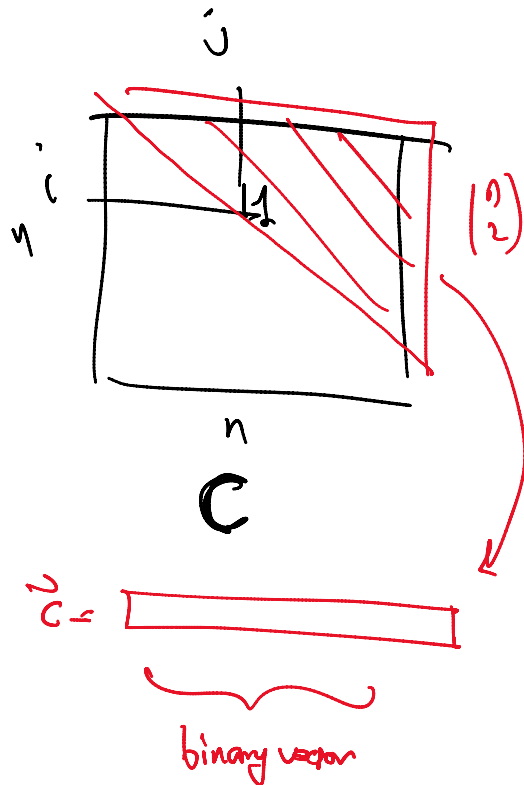
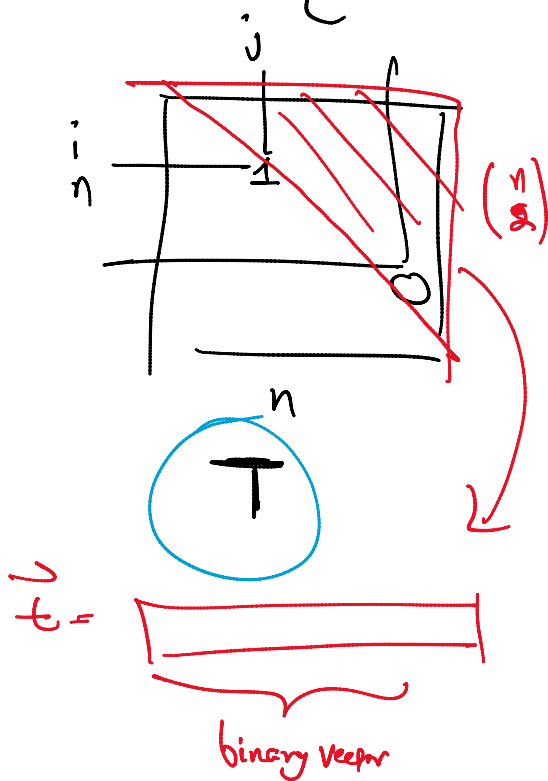
F-score: harmonic mean

Correlation Measure

define two matrices

$$T(i,j) = \begin{cases} 1 & \text{if } y_i = y_j \\ 0 & \text{otherwise} \end{cases}$$

$$C(i,j) = \begin{cases} 1 & \text{if } \hat{y}_i = \hat{y}_j \\ 0 & \text{otherwise} \end{cases}$$



Correlation?

$$\cos \theta = \frac{\vec{t}}{\|\vec{t}\|} \cdot \frac{\vec{c}}{\|\vec{c}\|}$$

Hubert statistic

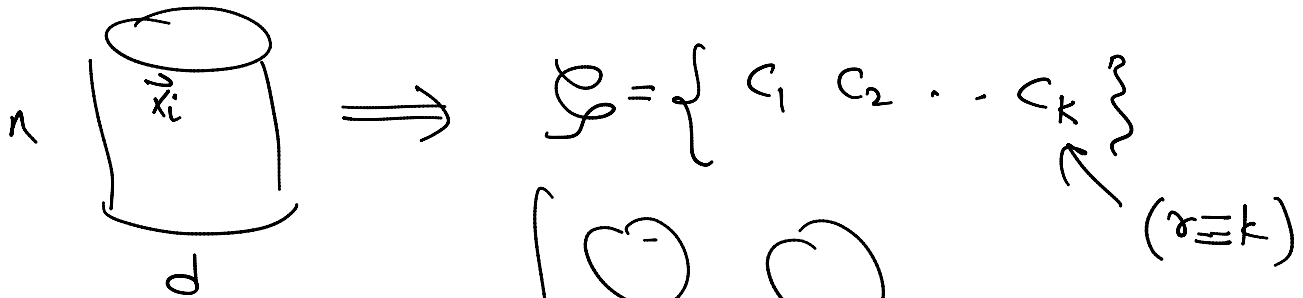
→ Compute using 2×2 Confusion matrix
(TP FP TN FN)

- Compute using 2×2 Confusion matrix

$$(T_P, F_P, T_N, F_N)$$

Internal/Intrinsic

No T: no ground truth labels/partitions



Intuition

- points that are
- 1) similar should be in the same cluster
 - 2) diff clusters should be far apart
- intra-cluster distances should be small
" " " " " " high
- inter-cluster distance should be high / large

$W \equiv$ proximity matrix \equiv distance matrix

$n \times n$ of distances

$w_{ij} = \|\vec{x}_i - \vec{x}_j\|$ \leftarrow L_2 norm for points $\left(\begin{array}{l} \text{shortest path length} \\ \text{in a graph} \end{array} \right)$

$$W = \{w_{ij}\}_{i,j=1 \dots n}$$

$$W = \{ w_{ij} \}_{i,j=1 \dots n}$$

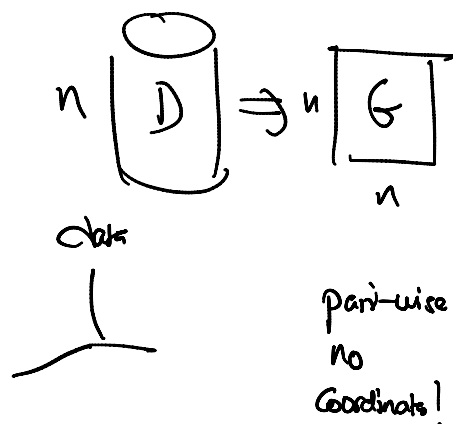
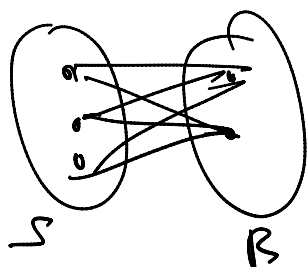
symmetric matrix

$$w_{ij} \geq 0$$

clustering \equiv partition \equiv k-way graph cut

$$W(S, R) = \sum_{i \in S} \sum_{j \in R} w_{ij}$$

$S, R \subseteq D$



$$D \Rightarrow \mathcal{C} = \{C_1, C_2, \dots, C_k\}$$

$$W(C_i, C_i) = \sum_{j \in C_i} \sum_{j \in C_i} w_{ij} \equiv \text{intra-cluster distance for } C_i$$

$$W_{in} = \frac{1}{2} \sum_{i=1}^k W(C_i, C_i)$$

over all clusters / intra-cluster

$$W_{out} = \frac{1}{2} \sum_{i=1}^k W(C_i, \bar{C}_i) \quad \text{inter-cluster}$$

Normalized cut

$$N_c = \sum_{i=1}^k \frac{W(C_i, \bar{C}_i)}{n} = \sum_{i=1}^k \frac{W(C_i, \bar{C}_i)}{n}$$

$$\max_{\{c_1, \dots, c_k\}} N_c = \sum_{i=1}^n \frac{W(c_i, \bar{c}_i)}{w(c_i)} = \sum_{i=1}^n \frac{W(c_i, \bar{c}_i)}{w(c_i, v)}$$

\downarrow
 W is a distance matrix

min if W were a kernel matrix