

כיצד נמדד - מידת אינפורמציה

Information gain = Info before - Info after

הפרש אינפורמציה

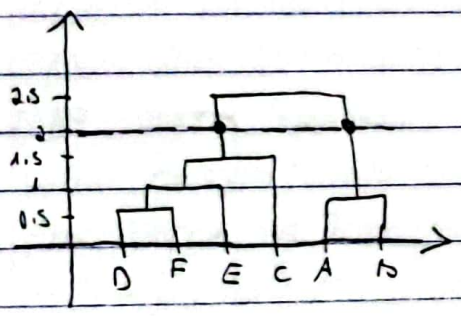
$$-\sum p \log p = \sum p \log \frac{1}{p}$$

אינפורמציה היא מידת אינפורמציה שיש לנו על המשתנה. ככל שיש לנו יותר אינפורמציה, כך יש לנו יותר אינפורמציה.

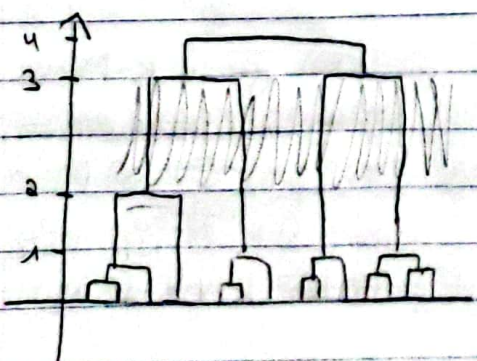
אינפורמציה היא מידת אינפורמציה שיש לנו על המשתנה.

אינפורמציה היא מידת אינפורמציה שיש לנו על המשתנה. אינפורמציה היא מידת אינפורמציה שיש לנו על המשתנה.

אינפורמציה היא מידת אינפורמציה שיש לנו על המשתנה. אינפורמציה היא מידת אינפורמציה שיש לנו על המשתנה.



אינפורמציה היא מידת אינפורמציה שיש לנו על המשתנה. אינפורמציה היא מידת אינפורמציה שיש לנו על המשתנה.



אינפורמציה היא מידת אינפורמציה שיש לנו על המשתנה. אינפורמציה היא מידת אינפורמציה שיש לנו על המשתנה.

אינפורמציה היא מידת אינפורמציה שיש לנו על המשתנה. אינפורמציה היא מידת אינפורמציה שיש לנו על המשתנה.

כ"ר - מ"ג - מ"ד - מ"ה

כ"ר - מ"ג - מ"ד - מ"ה
כ"ר - מ"ג - מ"ד - מ"ה

מ"ג - מ"ד - מ"ה
מ"ג - מ"ד - מ"ה

מ"ד - מ"ה
מ"ד - מ"ה

מ"ה
מ"ה

מ"ה
מ"ה

מ"ה
מ"ה

מ"ה
מ"ה

מ"ה
מ"ה

מ"ה
מ"ה

מ"ה
מ"ה

מ"ה
מ"ה

מ"ה
מ"ה

מ"ה
מ"ה

מבנה - פתרון

FT-Tree

(minsup=2) לנתון dataset הן נבנת FP-tree

{A,B} {A,C,D} {B,D} {A,B}

A: 3

B: 3

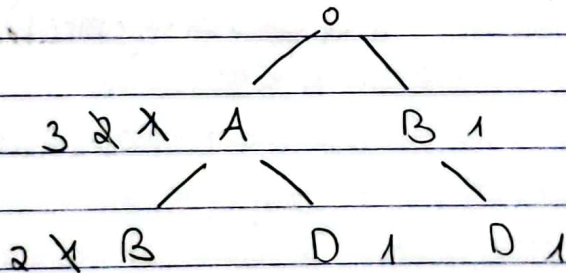
~~C: 1~~ minsup 2

D: 2

כל הנתונים שיש להם תמיכה לפחות 2 הם אלה שיש להם מניפוסט (c) (כל הנתונים) minsup 2

כל הנתונים שיש להם תמיכה לפחות 2 הם אלה שיש להם מניפוסט (c) (כל הנתונים) minsup 2

נתון dataset הן נבנת FP-tree. כל הנתונים שיש להם תמיכה לפחות 2 הם אלה שיש להם מניפוסט (c) (כל הנתונים) minsup 2. כל הנתונים שיש להם תמיכה לפחות 2 הם אלה שיש להם מניפוסט (c) (כל הנתונים) minsup 2.



נתון dataset הן נבנת FP-tree

* A *

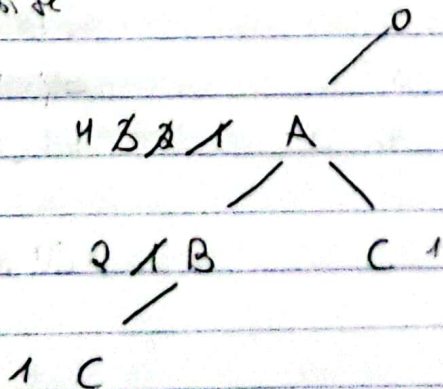
{A} {A,B} {A,C} {A,B,C}

minsup=2 נתון: FP-tree

A: 4

B: 2

C: 2



מיון - פתרון

יציבות בין הפתרון
: item sets ו הפתרון

$\{A\}$ $\{A,B\}$ $\{A,B,C\}$ $\{A,B,C,D\}$

? מהו confidence ו איך לרשום יציבות בין הפתרון

$$\text{conf} = \frac{1}{4}$$

$$\text{conf} = \frac{3}{4}$$

$$\text{conf} = \frac{1}{1}$$

$$\text{conf} = \frac{2}{3}$$

confidence = הסתברות הפתרון

הסתברות הפתרון

הסתברות הפתרון

$X \rightarrow Y$

$$\text{conf} = \frac{P(X \cup Y)}{P(X)}$$

$$\text{Confidence}(A \rightarrow B) = P(A|B) = \frac{P(A \cup B)}{P(A)}$$

(הסתברות הפתרון conf=1) . 3 הפתרון

כיצד ניתן לפרש את התוצאות

* את התוצאות של התהליך ניתן לפרש באמצעות שתי דרכים עיקריות:
1. שימוש בגרפים (למשל, גרף המרחק בין נקודות) כדי לראות כיצד המרחק בין נקודות משתנה ככל שהן מתקרבות.
2. שימוש בטבלה (למשל, טבלת המרחקים) כדי לראות את המרחקים בין נקודות באופן מסודר.

* נ"ל בציור הבא פירוט של

- 1. X-means (שיטה לזיהוי מספר קליסטרים)
- 2. RIPPER (שיטה לזיהוי קליסטרים)
- 3. EM (שיטה לזיהוי קליסטרים)
- 4. Cobweb (שיטה לזיהוי קליסטרים)

* שתי השיטות הבאות הן שיטות לזיהוי קליסטרים:

DBSCAN: שיטה לזיהוי קליסטרים המבוססת על מרחק בין נקודות. השיטה מחפשת נקודות "core" (נקודות שיש להן מספר מסוים של שכנים) ו"border" (נקודות שיש להן שכנים, אך אינן core).
 DBSCAN: שיטה לזיהוי קליסטרים המבוססת על מרחק בין נקודות. השיטה מחפשת נקודות "core" (נקודות שיש להן מספר מסוים של שכנים) ו"border" (נקודות שיש להן שכנים, אך אינן core).
 DBSCAN: שיטה לזיהוי קליסטרים המבוססת על מרחק בין נקודות. השיטה מחפשת נקודות "core" (נקודות שיש להן מספר מסוים של שכנים) ו"border" (נקודות שיש להן שכנים, אך אינן core).

השיטה הבאה היא שיטה לזיהוי קליסטרים המבוססת על מרחק בין נקודות. השיטה מחפשת נקודות "core" (נקודות שיש להן מספר מסוים של שכנים) ו"border" (נקודות שיש להן שכנים, אך אינן core).
 השיטה הבאה היא שיטה לזיהוי קליסטרים המבוססת על מרחק בין נקודות. השיטה מחפשת נקודות "core" (נקודות שיש להן מספר מסוים של שכנים) ו"border" (נקודות שיש להן שכנים, אך אינן core).
 השיטה הבאה היא שיטה לזיהוי קליסטרים המבוססת על מרחק בין נקודות. השיטה מחפשת נקודות "core" (נקודות שיש להן מספר מסוים של שכנים) ו"border" (נקודות שיש להן שכנים, אך אינן core).

margin - SVM - linear

hard margin - SVM - linear - margin - SVM - linear - margin - SVM - linear

soft margin - SVM - linear - margin - SVM - linear - margin - SVM - linear

polynomial kernel - SVM - linear - margin - SVM - linear - margin - SVM - linear

kernel - SVM - linear - margin - SVM - linear - margin - SVM - linear

support - SVM - linear - margin - SVM - linear - margin - SVM - linear

support - SVM - linear

$support(A) \leq support(B)$

dataset, items in A & B

1. $support(A) > support(B)$

2. $support(A) = support(B)$

3. $support(A) < support(B)$

4. $support(A) > support(B)$

$support(A) > support(B) \iff support(A) > support(B)$

margin - SVM - linear - margin - SVM - linear - margin - SVM - linear

margin - SVM - linear - margin - SVM - linear - margin - SVM - linear

SVM - linear

כיתה ראשונה - מנתח

* ZeroR מנתח class אחת בלבד וכל הנתונים שייכים אליה
 - Overfitting - מנתח את הנתונים ומחזיר את accuracy של 100%
 * ZeroR מנתח dataset של WEKA, accuracy של 100%
 * מנתח dataset של WEKA ומחזיר accuracy של 100%
 * $\frac{100}{25} > 5$ (כלומר 4% מהנתונים הם שייכים לאותה class)

* DBSCAN - מנתח את הנתונים ומחזיר core points ו-borders
 - Core points - מנתח את הנתונים ומחזיר core points
 - border - מנתח את הנתונים ומחזיר border points
 - noise - מנתח את הנתונים ומחזיר noise points
 - מנתח את הנתונים ומחזיר core points ו-borders

* מנתח את הנתונים ומחזיר PCA
 - PCA מנתח את הנתונים ומחזיר PCA
 - PCA מנתח את הנתונים ומחזיר PCA
 - PCA מנתח את הנתונים ומחזיר PCA

* PCA מנתח את הנתונים ומחזיר PCA
 - PCA מנתח את הנתונים ומחזיר PCA
 - PCA מנתח את הנתונים ומחזיר PCA
 - PCA מנתח את הנתונים ומחזיר PCA

* מנתח את הנתונים ומחזיר PCA
 - PCA מנתח את הנתונים ומחזיר PCA
 - PCA מנתח את הנתונים ומחזיר PCA

DBSCAN - DBSCAN

DBSCAN

(0,0), (2,1), (2,3), (4,2)

DBSCAN

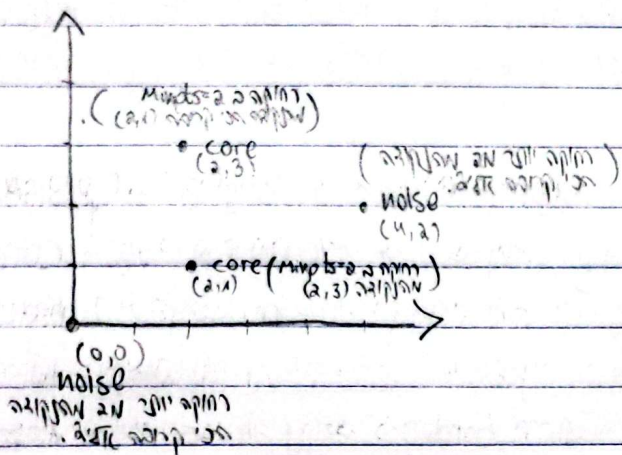
Minpts=2, E=2

noise

border

core

Minpts



DBSCAN

noise

border

core

Minpts

DBSCAN

* * *

(0,0), (2,1), (2,3), (4,2)

DBSCAN

noise

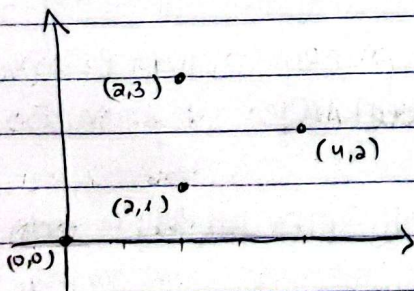
border

core

Minpts=2, E=1.2

Minpts=3, E=1.3

Minpts=2, E=1.5, 4

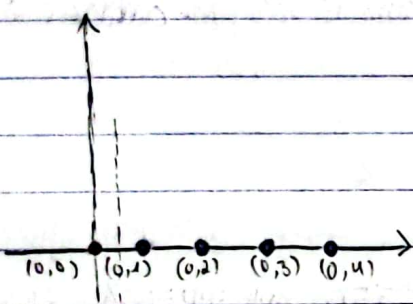


מיון נתונים - SVM

SVM עמית

$x_1 = (0, 0)$	$y_1 = -1$	} α_i
$x_2 = (0, 1)$	$y_2 = 1$	
$x_3 = (0, 2)$	$y_3 = 1$	
$x_4 = (0, 3)$	$y_4 = 1$	
$x_5 = (0, 4)$	$y_5 = 1$	

המטרה היא למצוא את ה-SVM המקסימלי



- הצורה הכללית של ה-SVM היא:
- 1. $\alpha_1 = 2, \alpha_2 = -2, \alpha_3 = 0, \alpha_4 = 0, \alpha_5 = 0$
 - 2. $\alpha_1 = 1, \alpha_2 = 1, \alpha_3 = 0, \alpha_4 = 0, \alpha_5 = 0$
 - 3. $\alpha_1 = 2, \alpha_2 = 2, \alpha_3 = 0, \alpha_4 = 0, \alpha_5 = 0$
 - 4. $\alpha_1 = 1, \alpha_2 = 2, \alpha_3 = 3, \alpha_4 = 4, \alpha_5 = 5$

הצורה הכללית של ה-SVM היא $w \cdot x + b = 0$.
 המטרה היא למצוא את ה-SVM המקסימלי.
 ה-SVM המקסימלי הוא זה שיש לו המרווח המרבי.
 ה-SVM המקסימלי הוא זה שיש לו המרווח המרבי.
 ה-SVM המקסימלי הוא זה שיש לו המרווח המרבי.

ה-Lagrangian function היא:

$$L = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \cdot \alpha_j \cdot y_i \cdot y_j \cdot (x_i \cdot x_j)$$

למשל:

$$L = 2\alpha - \frac{1}{2} \cdot \alpha^2 \cdot 1 \cdot 1 \cdot (0,1)(0,1) = 2\alpha - \frac{1}{2} \alpha^2$$

למשל:

$$L = 2\alpha - \frac{1}{2} \alpha^2$$

$$L' = 2 - \alpha = 0$$

$$\alpha = 2$$

ה-SVM המקסימלי הוא זה שיש לו המרווח המרבי

כ"מ - Naive Bayes

Naive Bayes מניח שכל הfeatures הם independent זהו ה"naive" שבו נקראת המודל.
 Naive Bayes מניח שכל הfeatures הם independent זהו ה"naive" שבו נקראת המודל.
 Naive Bayes מניח שכל הfeatures הם independent זהו ה"naive" שבו נקראת המודל.

Naive Bayes מניח שכל הfeatures הם independent זהו ה"naive" שבו נקראת המודל.
 Naive Bayes מניח שכל הfeatures הם independent זהו ה"naive" שבו נקראת המודל.
 Naive Bayes מניח שכל הfeatures הם independent זהו ה"naive" שבו נקראת המודל.

מ"ו - אלגוריתמים לבידול

1. OneR - קטגוריאל בלבד
2. SVM - SVM
3. XMeans - clustering
4. Random Forest - SVM

max/min - מרחק בין קטגוריה אחת לשנייה

- Average distance - מרחק בין כל הקטגוריות, כלומר מרחק בין כל הקטגוריות
- Manhattan distance - מרחק בין כל הקטגוריות, כלומר מרחק בין כל הקטגוריות
- Euclidean distance - מרחק בין כל הקטגוריות, כלומר מרחק בין כל הקטגוריות

DBSCAN - זיהוי קבוצות

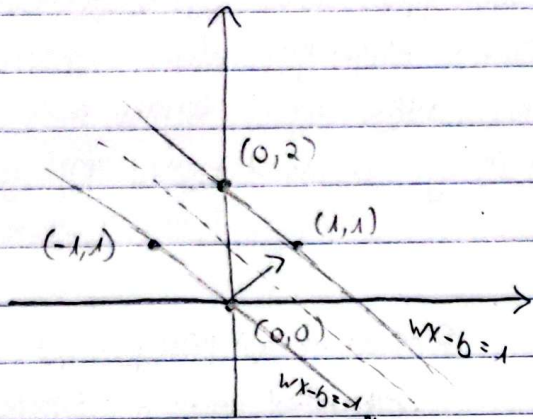
- EM - זיהוי קבוצות, כלומר זיהוי קבוצות
- K-means - זיהוי קבוצות, כלומר זיהוי קבוצות

Find the optimal solution

!!!!!!!

$w = (1, 1) \in \mathbb{R}^2$, $x_1 = (0, 0)$ $x_2 = (1, 1)$ $x_3 = (0, 2)$ $x_4 = (-1, 1)$

? b for each row



0	1
1	2
-1	3
2	4

$w x - b = 1 \Rightarrow (1, 1)(1, 1) - b = 1 \Rightarrow 2 - b = 1 \Rightarrow b = 1$

$w x - b = -1 \Rightarrow (1, 1)(0, 0) - b = -1 \Rightarrow b = 1$

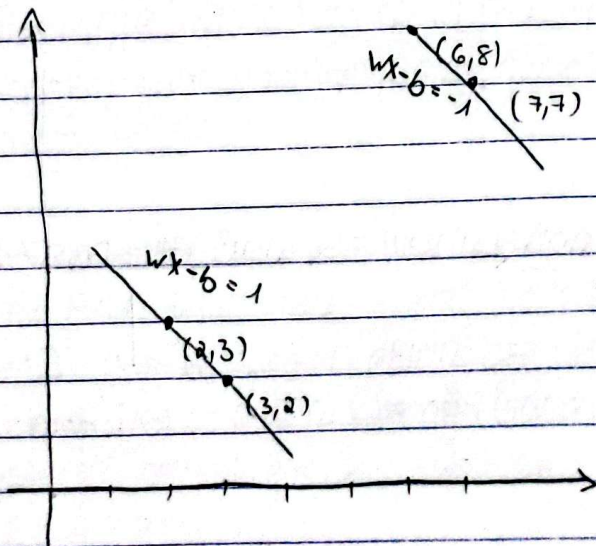
$(0, 2)(1, 1) - b = 1 \Rightarrow 2 - b = 1 \Rightarrow b = 1$

$(-1, 1)(1, 1) - b = -1 \Rightarrow -1 + 1 - b = -1 \Rightarrow b = 1$

and then we will find

* * *

!!!!!!



$x_1 = (2, 3)$ $y = 1$

$x_2 = (3, 2)$ $y = 1$

$x_3 = (6, 8)$ $y = -1$

$x_4 = (7, 7)$ $y = -1$

!!!!!!

? b 1 w can

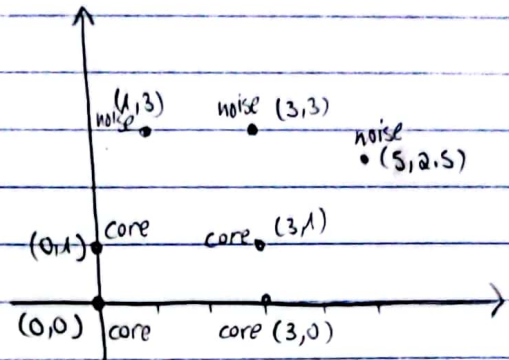
$w x - b = 1$

$w x - b = -1$

DBSCAN algorithm - DBSCAN algorithm

x	y
0	1
1	3
0	0
3	1
3	3
3	0
5	2.5

minpts = 2
E = 1.5



DBSCAN algorithm
border = ...

minpts < E - core
minpts < E - border
minpts < E - noise