

## מדעי הנתונים ובינה עסקית

## תרגיל מס' 1

הוראות כלליות:

- הגשת העבודה היא בזוגות.
- יש להגיש את העבודה דרך המודל (שותף אחד מכל זוג צריך להגיש).
- יש להגיש את העבודה עד לתאריך 9.4.2020.
- יש לכתוב תשובות סופיות במקום המיועד ובסוף המסמך לצרף את כל החישובים המתאימים.
- יש להקפיד על 3 ספרות אחרי הנקודה העשרונית, **אי הקפדה תגרור הורדת נקודות.**
- תשובה סופית, גם אם היא נכונה, ללא הצגת דרך החישוב המלאה בנספחים, **לא תזכה בניקוד.**

## שאלה 1 (30 נק')

כאשר: 1931-1968 הנתונים הבאים מתארים את שיעורי הפשיעה בקנדה בשנים

Tfr	Total fertility rate per 1000 women.
Partic	Women's labor-force participation rate per 1000.
Degrees	Women's post-secondary degree rate per 10,000.
Fconvict	Female indictable-offense conviction rate per 100,000.
Ftheft	Female theft conviction rate per 100,000.
Mconvict	Male indictable-offense conviction rate per 100,000.
Mtheft	Male theft conviction rate per 100,000.

year	tfr	partic	degrees	fconvict	ftheft	mconvict	mtheft
1935	2755	238	13.2	79.4	20.4	765.7	247.1
1936	2696	240	13.2	91	22.1	816.5	254.9
1937	2646	241	12.2	100.4	22.4	821.8	272.4
1938	2701	242	12.6	108.9	21.8	956.8	285.8
1939	2654	244	12.3	123.6	21.1	1035.7	292.2
1940	2766	245	12	157.3	21.4	951.6	256
1941	2832	246	11.7	154.3	25.3	850.9	205.8
1942	2964	268	11.2	143.9	27.1	769.7	188
1943	3041	333	11.5	147.5	29	811.2	205.8
1944	3010	335	11.1	97.3	24.2	865.4	207.9

1945	3018	331	12.5	76.6	24.7	866.8	197.8
1946	3374	253	15	72.8	20.5	968	195.9
1947	3545	243	17.6	68.9	20.7	894	198.9
1948	3441	241	21.2	66.7	22.8	830.8	198.6
1949	3456	242	22.7	55	18.5	811.2	216.1
1950	3455	237	21.4	55	19.3	775.4	212
1951	3503	242	20.7	55	20	739.8	208
1952	3641	240	19.5	54.7	19.4	740.3	199.3
1953	3721	233	36.6	47.9	16.1	725.1	176.4

א. יש לנרמל את המשתנה **partic** באמצעות Min-Max normalization כאשר  $New\ min = 1, New\ max = 10$ .

1. מהו המינימום המקורי של המשתנה? 233
2. מהו המקסימום המקורי של המשתנה? 335
3. מהו הערך המנורמל של המשתנה בשנת 1952? 1.618
4. יש לצרף טבלה המכילה את המשתנה המנורמל בנספחים.

ב. יש לבצע החלקה למשתנה **mconvict** באמצעות Weighted moving average כאשר:

t	Weight
-3	180.
-2	220.
-1	60.

1. מהו הערך של המשתנה **mconvict** המוחלק בשנת 1952? 760.484
2. מהו הערך של המשתנה **mconvict** המוחלק בשנת 1946? 856.484
3. נא לצרף את המשתנה המוחלק בנספחים.
4. נא לצרף גרף עם המשתנה המקורי והמשתנה לאחר ההחלקה בנספחים.

## שאלה 2 (35 נק')

מעוניינים לחזות נוכחות של מחלת לב בחולים. לשם כך, Cleveland Clinic Foundation רופאים מ-  
משתנים מועמדים ומשתנה (3) מכיל **hw1\_dataset** המצורף (CSV) נאסף מידע על 100 חולים. קובץ ה-  
prediction. אחד ( )

הגדרות המשתנים:

Attribute	Domain	Type
gender	1=male, 0=female	Nominal
fbs	0=false, 1=true	Nominal
slope	1=upsloping, 2=flat, 3=downsloping	Nominal
Prediction	0,1,2,3,4	Nominal

- א. מהי האנטרופיה של המשתנה gender? 0.869
- ב. מהי האנטרופיה של המשתנה fbs? (10 נק') 0.577
- ג. מהי האנטרופיה המותנית של משתנה המטרה בהינתן המשתנה slope? 1.577
- ד. מהי האנטרופיה המקסימלית של המשתנה slope? 1.585
- ה. יש לחשב את ה-mutual information של משתנה המטרה ו-gender בשתי דרכים:
- לפי ההפרש בין האנטרופיות הרלוונטיות
  - לפי הנוסחה של mutual information  $I(X;Y)$

**שאלה 3 (35 נק')**

יש להשתמש באותו קובץ נתונים משאלה 2.

- א. יש לבנות מודל ID3 (עד רמה 2, קודקוד השורש הוא רמה 0). יש להראות חישוב ידני של לפחות משתנה אחד בכל קדקוד, לשאר המשתנים ניתן לכתוב רק את התוצאה הסופית. יש להציג את העץ שהתקבל.
- ב. מהו ה-majority rule accuracy?
- ג. מהו ה-training accuracy של המודל?
- ד. 5 חולים חדשים הגיעו לבית החולים. ע"פ המודל שבניתם, סווגו את הרשומות הבאות:

Gender	fbs	Slope	Prediction
0	1	1	<b>0</b>
1	0	2	<b>1</b>
1	1	3	<b>0</b>
0	0	3	<b>3 (random choice between 3 and 0)</b>
0	0	2	<b>0</b>

**בהצלחה!****נספחים****א1.**

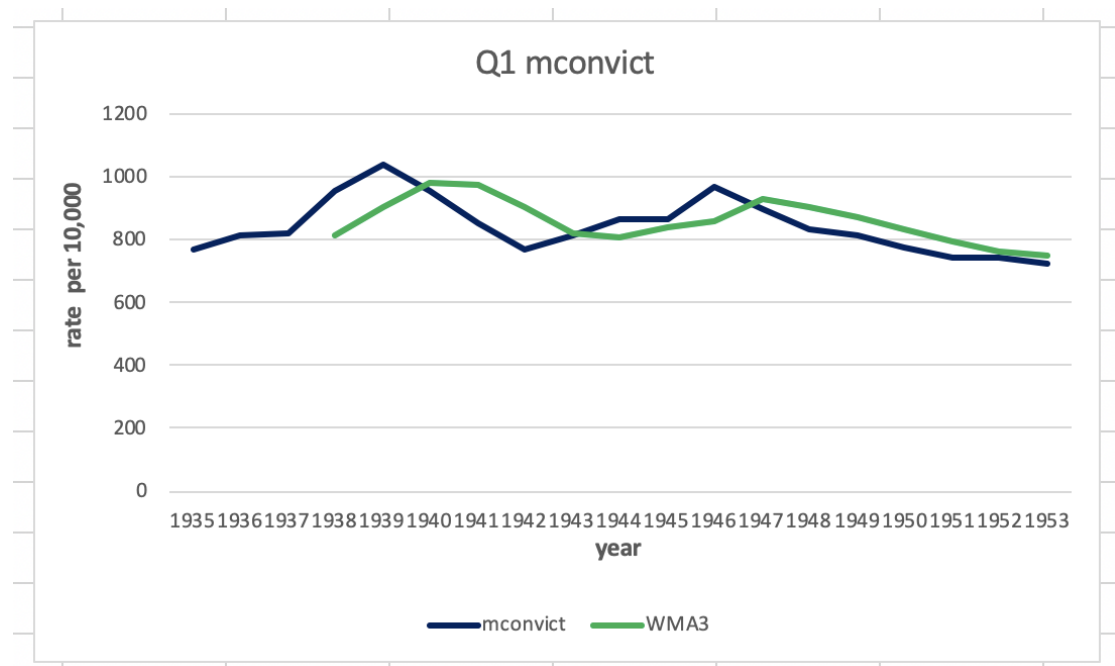
<b>year</b>	<b>normalized partic</b>
1935	1.441
1936	1.618
1937	1.706
1938	1.794
1939	1.971
1940	2.059
1941	2.147
1942	4.088
1943	9.824
1944	10.000
1945	9.647
1946	2.765
1947	1.882
1948	1.706
1949	1.794
1950	1.353
1951	1.794
1952	1.618
1953	1.000

233				1	1
335				2	
10		new max		3	
1		new min			



## 1.ב.

year	mconvict	WMA3
1935	765.7	
1936	816.5	
1937	821.8	
1938	956.8	810.536
1939	1035.7	901.846
1940	951.6	979.84
1941	850.9	971.038
1942	769.7	906.318
1943	811.2	820.306
1944	865.4	809.216
1945	866.8	836.25
1946	968	856.484
1947	894	927.268
1948	830.8	905.384
1949	811.2	869.4
1950	775.4	830.416
1951	739.8	793.248
1952	740.3	760.484
1953	725.1	746.508
Weight	t	
0.18	-3	
0.22	-2	
0.6	-1	



## 2.א.

א.		total	probability				
	male	71	0.710				
	female	29	0.290				
	total	100	1				
	H(X)	0.869					
ב.		total	probability				
	TRUE	13	0.130				
	FALSE	87	0.870				
	total	100	1				
	H(X)	0.557					
	predict\slope	1	2	3	total		
	0	39	12	6	57		
	1	7	10	2	19		
	2	2	7	1	10		
	3	2	6	2	10		
	4	0	4	0	4		
	total	50	39	11	100		
	P(predict, slope)	1	2	3	total		
	0	0.390	0.120	0.060	0.570		
	1	0.070	0.100	0.020	0.190		
	2	0.020	0.070	0.010	0.100		
	3	0.020	0.060	0.020	0.100		
	4	0.000	0.040	0.000	0.040		
	total	0.500	0.390	0.110	1.000		
	P(predict slope)	1	2	3			
	0	0.780	0.308	0.545			
	1	0.140	0.256	0.182			
	2	0.040	0.179	0.091			
	3	0.040	0.154	0.182			
	4	0.000	0.103	0.000			
	total	1.000	1.000	1.000			
	p(slope, prediction) * log(p(prediction slope))	1	2	3	total		ג.
	0	0.140	0.204	0.052	0.396		
	1	0.199	0.196	0.049	0.444		
	2	0.093	0.173	0.035	0.301		
	3	0.093	0.162	0.049	0.304		
	4	0.000	0.131	0.000	0.131		
	total	0.524	0.867	0.185	1.577		
	$H(\text{prediction}   \text{slope}) = - \sum p(\text{slope}, \text{prediction}) * \log(p(\text{prediction}   \text{slope}))$						
			max H(slope)	1.585			ד.
	הערה: לפי המקסימום היא:						
	$\log(\text{number of total possible values})$						



[illegible]

## 3.8

root	0	1	2	3	4	sum
	57.000	19.000	10.000	10.000	4.000	100
Pi	0.570	0.190	0.100	0.100	0.040	
log(Pi)	-0.811	-2.396	-3.322	-3.322	-4.644	
Pi*log(Pi)	-0.462	-0.455	-0.332	-0.332	-0.186	
$Info(D) = -\sum_{i=1}^m p_i \log_2(p_i)$						
Info(root)	1.768					

חישוב אנטרופיה של השורש:

$$Info(prediction) = -0.57 * \log_2(0.57) - 0.19 * \log_2(0.19) - 0.1 * \log_2(0.1) - 0.1 * \log_2(0.1) - 0.04 * \log_2(0.04) = \mathbf{1.768}$$

חישוב information gain של פיצול לפי gender:

$$IG(prediction, Gender) = Info(prediction) - Info_{gender}(prediction)$$

נקבע וקטור לערכי הרשומות במשתנה המטרה:

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

Gender/prediction	0	1	2	3	4	sum
0	23.000	1.000	1.000	3.000	1.000	29.000
1	34.000	18.000	9.000	7.000	3.000	71.000
Pi	0.793	0.034	0.034	0.103	0.034	
	0.479	0.254	0.127	0.099	0.042	

$$\begin{aligned} Info(Prediction_0) &= -P_{0|0} * \log_2(P_{0|0}) - P_{1|0} * \log_2(P_{1|0}) - P_{2|0} * \log_2(P_{2|0}) - \\ &P_{3|0} * \log_2(P_{3|0}) - P_{4|0} * \log_2(P_{4|0}) = \\ &-0.793 * \log_2(0.793) - 0.034 * \log_2(0.034) - 0.034 * \log_2(0.034) - 0.103 \\ &* \log_2(0.103) - 0.034 * \log_2(0.034) = \mathbf{1.106} \end{aligned}$$

$$\begin{aligned} Info(Prediction_1) &= -P_{0|1} * \log_2(P_{0|1}) - P_{1|1} * \log_2(P_{1|1}) - P_{2|1} * \log_2(P_{2|1}) - \\ &P_{3|1} * \log_2(P_{3|1}) - P_{4|1} * \log_2(P_{4|1}) = \\ &-0.479 * \log_2(0.479) - 0.254 * \log_2(0.254) - 0.127 * \log_2(0.127) - 0.099 \\ &* \log_2(0.099) - 0.042 * \log_2(0.042) = \mathbf{1.911} \end{aligned}$$

$$\begin{aligned} \text{Info}_{\text{gender}}(\text{prediction}) &= \sum \left( \frac{|prediction_i|}{|prediction|} \cdot \text{Info}(prediction_i) \right) \\ &= 0.29 * 1.106 + 0.71 * 1.911 = 0.321 + 1.357 = \mathbf{1.678} \end{aligned}$$

$$\begin{aligned} \text{IG}(\text{prediction}, \text{Gender}) &= \text{Info}(\text{prediction}) - \text{Info}_{\text{gender}}(\text{prediction}) = \\ &1.768 - 1.678 = \mathbf{0.09} \end{aligned}$$

חישוב information gain של פיצול לפי fbs:

$$\begin{aligned} \text{IG}(\text{prediction}, \text{fbs}) &= \text{Info}(\text{prediction}) - \text{Info}_{\text{fbs}}(\text{prediction}) = \\ &1.768 - 1.745 = \mathbf{0.023} \end{aligned}$$

חישוב information gain של פיצול לפי slope:

$$\begin{aligned} \text{IG}(\text{prediction}, \text{fbs}) &= \text{Info}(\text{prediction}) - \text{Info}_{\text{slope}}(\text{prediction}) = \\ &1.768 - 1.577 = \mathbf{0.191} \end{aligned}$$

נבחר לפצל לפי המשתנה slope מכיוון שהוא מספק את רווח הידע הגדול ביותר (הקטנת אי הוודאות הגדולה ביותר).

חישוב אנטרופיה של הקדקוד בו slope=1:

$$\text{info}(prediction_{\text{slope}=1}) = 1.048 \text{ האנטרופיה של קדקוד זה חושבה בשלב הקודם:}$$

חישוב information gain של פיצול לפי gender:

$$\text{IG}(prediction_{\text{slope}=1}, \text{Gender}) = \text{info}(prediction_{\text{slope}=1}) - \text{Info}_{\text{gender}}(prediction_{\text{slope}=1})$$

$$\begin{aligned} \text{Info}(\text{Prediction}_{\text{slope}=1, \text{gender}=0}) &= -P_{0|0} * \log_2(P_{0|0}) - P_{1|0} * \log_2(P_{1|0}) - P_{2|0} * \\ &\log_2(P_{2|0}) - P_{3|0} * \log_2(P_{3|0}) - P_{4|0} * \log_2(P_{4|0}) = \\ &-0.944 * \log_2(0.944) - 0.056 * \log_2(0.056) - 0 - 0 - 0 = \mathbf{0.31} \end{aligned}$$

$$\begin{aligned} \text{Info}(\text{Prediction}_{\text{slope}=1, \text{gender}=1}) &= -P_{0|1} * \log_2(P_{0|1}) - P_{1|1} * \log_2(P_{1|1}) - P_{2|1} * \\ &\log_2(P_{2|1}) - P_{3|1} * \log_2(P_{3|1}) - P_{4|1} * \log_2(P_{4|1}) = \\ &-0.688 * \log_2(0.688) - 0.188 * \log_2(0.188) - 0.063 * \log_2(0.063) - 0.063 \\ &* \log_2(0.063) - 0 = \mathbf{1.324} \end{aligned}$$

$$\begin{aligned} \text{Info}_{\text{gender}}(prediction_{\text{slope}=1}) &= \sum \left( \frac{|prediction_{\text{slope}=1,i}|}{|prediction_{\text{slope}=1}|} \cdot \text{Info}(prediction_{\text{slope}=1,i}) \right) \\ &= \frac{18}{50} * 0.31 + \frac{32}{50} * 1.324 = \mathbf{0.959} \end{aligned}$$

$$IG(prediction_{slope=1}, Gender) = Info(prediction_{slope=1}) - Info_{gender}(prediction_{slope=1}) = 1.048 - 0.959 = \mathbf{0.089}$$

חישוב information gain של פיצול לפי fbs:

$$IG(prediction_{slope=1}, fbs) = Info(prediction_{slope=1}) - Info_{fbs}(prediction_{slope=1}) = 1.048 - 0.99 = \mathbf{0.058}$$

נבחר לפצל לפי המשתנה fbs מכיוון שהוא מספק את רווח הידע הגדול ביותר (הקטנת אי הוודאות הגדולה ביותר).

חישוב אנטרופיה של הקדקוד בו slope=2:

$$info(prediction_{slope=2}) = 2.224$$

חישוב information gain של פיצול לפי gender:

$$IG(prediction_{slope=2}, Gender) = info(prediction_{slope=2}) - Info_{gender}(prediction_{slope=2})$$

$$\begin{aligned} Info(Prediction_{slope=2, gender=0}) &= -P_{0|0} * \log_2(P_{0|0}) - P_{1|0} * \log_2(P_{1|0}) - P_{2|0} * \log_2(P_{2|0}) - P_{3|0} * \log_2(P_{3|0}) - P_{4|0} * \log_2(P_{4|0}) = \\ &= -0.571 * \log_2(0.571) - 0 - 0.143 * \log_2(0.143) - 0.143 * \log_2(0.143) - 0.143 * \log_2(0.143) = \mathbf{1.664} \end{aligned}$$

$$\begin{aligned} Info(Prediction_{slope=2, gender=1}) &= -P_{0|1} * \log_2(P_{0|1}) - P_{1|1} * \log_2(P_{1|1}) - P_{2|1} * \log_2(P_{2|1}) - P_{3|1} * \log_2(P_{3|1}) - P_{4|1} * \log_2(P_{4|1}) = \\ &= -0.25 * \log_2(0.25) - 0.313 * \log_2(0.313) - 0.188 * \log_2(0.188) - 0.156 * \log_2(0.156) - 0.094 * \log_2(0.094) = \mathbf{2.216} \end{aligned}$$

$$\begin{aligned} Info_{gender}(prediction_{slope=2}) &= \sum \left( \frac{|prediction_{slope=2,i}|}{|prediction_{slope=2}|} * Info(prediction_{slope=2,i}) \right) \\ &= \frac{7}{39} * 1.664 + \frac{32}{39} * 2.216 = \mathbf{2.117} \end{aligned}$$

$$IG(prediction_{slope=1}, Gender) = Info(prediction_{slope=1}) - Info_{gender}(prediction_{slope=1}) = 2.224 - 2.117 = \mathbf{0.107}$$

חישוב information gain של פיצול לפי fbs:

$$IG(prediction_{slope=2}, fbs) = Info(prediction_{slope=2}) - Info_{fbs}(prediction_{slope=2}) = 2.224 - 2.122 = \mathbf{0.102}$$

נבחר לפצל לפי המשתנה gender מכיוון שהוא מספק את רווח הידע הגדול ביותר (הקטנת אי הוודאות הגדולה ביותר).

חישוב אנטרופיה של הקדקוד בו slope=3:

$$info(prediction_{slope=3}) = 1.686$$

חישוב information gain של פיצול לפי gender:

$$IG(prediction_{slope=3}, Gender) = info(prediction_{slope=3}) - Info_{gender}(prediction_{slope=3})$$

$$\begin{aligned} \text{Info}(\text{Prediction}_{\text{slope}=3, \text{gender}=0}) &= -P_{0|0} * \log_2(P_{0|0}) - P_{1|0} * \log_2(P_{1|0}) - P_{2|0} * \\ &\log_2(P_{2|0}) - P_{3|0} * \log_2(P_{3|0}) - P_{4|0} * \log_2(P_{4|0}) = \\ &-0.5 * \log_2(0.5) - 0 - 0 - 0.5 * \log_2(0.5) - 0 = 1 \end{aligned}$$

$$\begin{aligned} \text{Info}(\text{Prediction}_{\text{slope}=3, \text{gender}=1}) &= -P_{0|1} * \log_2(P_{0|1}) - P_{1|1} * \log_2(P_{1|1}) - P_{2|1} * \\ &\log_2(P_{2|1}) - P_{3|1} * \log_2(P_{3|1}) - P_{4|1} * \log_2(P_{4|1}) = \\ &-0.571 * \log_2(0.571) - 0.286 * \log_2(0.286) - 0.143 * \log_2(0.143) - 0 - 0 \\ &= 1.379 \end{aligned}$$

$$\begin{aligned} \text{Info}_{\text{gender}}(\text{prediction}_{\text{slope}=3}) &= \sum \left( \frac{|\text{prediction}_{\text{slope}=3, i}|}{|\text{prediction}_{\text{slope}=3}|} \cdot \text{Info}(\text{prediction}_{\text{slope}=3, i}) \right) \\ &= \frac{4}{11} * 1 + \frac{7}{11} * 1.379 = 1.241 \end{aligned}$$

$$\begin{aligned} \text{IG}(\text{prediction}_{\text{slope}=1}, \text{Gender}) &= \text{Info}(\text{prediction}_{\text{slope}=1}) - \text{Info}_{\text{gender}}(\text{prediction}_{\text{slope}=1}) = \\ &1.686 - 1.241 = 0.445 \end{aligned}$$

**חישוב information gain של פיצול לפי fbs:**

$$\begin{aligned} \text{IG}(\text{prediction}_{\text{slope}=2}, \text{fbs}) &= \text{Info}(\text{prediction}_{\text{slope}=2}) - \text{Info}_{\text{fbs}}(\text{prediction}_{\text{slope}=2}) = \\ &1.686 - 1.523 = 0.163 \end{aligned}$$

**נבחר לפצל לפי המשתנה gender מכיוון שהוא מספק את רווח הידע הגדול ביותר (הקטנת אי הוודאות הגדולה ביותר).**