

מבחן במסדי נתונים

מר נתנאל שקרון

7029210-2,6

סמסטר קיץ מועד א' ח תשרי התשפ"ב, 14.09.2021

הנחיות כלליות:

- משך הבחינה: 150 דקות.
- יש לענות בגוף השאלון! המחברת תשמש כטיטא בלבד, מענה במחברת עלול לגרור ציון 0.
- אין להכניס שום חומר עזר.
- השימוש במחשבון אסור.
- בשאלות האמריקאיות רק תשובה אחת נכונה.
- בסיום הבחינה - נא למסור את השאלון ואת המחברת.

	1	2	3	4	5	6	Total
Max points	28	19	18	15	10	10	100
Grade							

ב ה צ ל ח ה !

ב. כתבו שאילתה המחזירה את רשימת המחלקות בהן מספר העובדים גדול מ-3 ומספר העובדים בכל אחת מהן (10 נק)

[illegible]

ג. הסבירו איך ניתן להבטיח שלא יימחקו מחלקות שיש בהן עובדים (4 נק)

שאלה 2 – Normalization (17 נק')

נתונה הטבלה הבאה :

מסרד	מנחה	התמחות	ת"ז סטודנט
11.1.3	ד"ר אורי כהנה	Physics	111
9.4.7	ד"ר ישי קפלן	Music	111
9.3.1	ד"ר נוימן - חזון	Math	320
11.2.6	ד"ר ישראל בוחבוס	Physics	671
10.2.2	פרופסור עומר אלפרוביץ	Physics	803

תחת ההנחות :

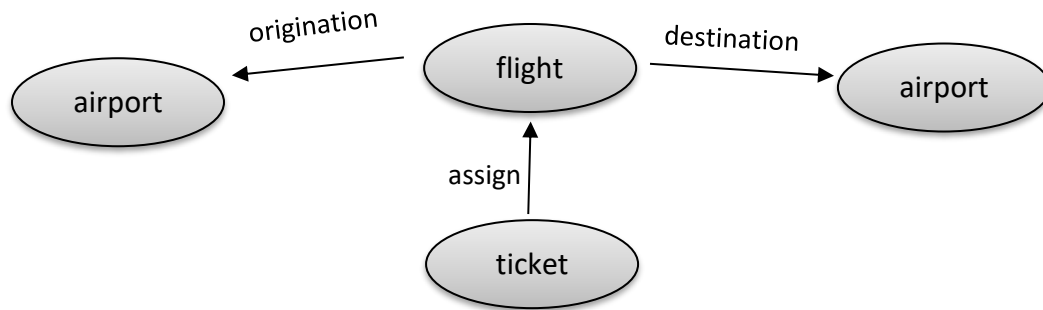
- כל סטודנט יכול להתמחות ביותר מהתמחות אחת
 - בכל התמחות שנבחרה לסטודנט יכול להיות רק מנחה אחד
 - לכל התמחות יש מספר מנחים
 - כל מנחה מנחה תחום התמחות אחת בלבד
 - כל מנחה עשוי להנחות יותר מסטודנט אחד
- א. מה/ם המפתחות האפשריים של הטבלה ? (5 נק)

ב. באיזו רמת נרמול נמצאת הרלציה הנ"ל ? (5 נק)

ג. במידה והרלציה דורשת נרמול - תאר את שלבי הנרמול הנדרשים בפירוט והצג את מבנה הטבלאות המתקבל בכל שלב . (7 נק)

שאלה 3 – Neo4J | NoSQL (18 נק')

א. נתון בסיס נתונים גרפי של neo4j, המתאר רשת של מידע על טיסות, שדות תעופה וכרטיסי טיסה.



סוגי הצמתים והמידע על כל הצומת:

צומת **Flight** - מכילה (בין השאר) את התכונות – תאריך ושעת טיסה, קוד חברת תעופה, משך הטיסה, מרחק הטיסה.

צומת **Airport** - מכילה (בין השאר) את התכונות- קוד שדה התעופה, מיקום שדה התעופה.

צומת **Ticket** - מכילה (בין השאר) את התכונות- מס' מושב בטיסה, שם המחלקה בטיסה (first/ business/ economy), מחיר הכרטיס

קשרים בין הצמתים:

בין צומת Ticket לצומת Flight יכול להיות מוגדר הקשר- **Assign** – בשביל לתאר שיוך בין כרטיס לטיסה.

בין צומת Flight לצומת Airport יכול להיות מוגדרים הקשרים – 1. **Origination** – בשביל לתאר את הקשר בין שדה תעופה ממנו יוצאת הטיסה לטיסה עצמה. 2. **Destination** – בשביל לתאר את הקשר בין שדה התעופה אליו מגיע הטיסה לטיסה עצמה.

עליכם לכתוב שאילתה ב- cypher המחזירה את המחירים של טיסות **ישירות** שיוצאות משדה התעופה באורלנדו (**ORD**), עבור כרטיסים **במחלקת עסקים**.

עבור כל תוצאה יש להציג את יעד הטיסה, מחיר הכרטיס וקוד חברת התעופה. (12 נק')
לדוגמא:

destiantion	price	airlineCode
TLV	1000	LY
CDG	850	AF
JFK	300	AA

[illegible]

שאלה 4 – spark (15 נק')

נתון מסמך במבנה json המתאר פרטי טיסות ברחבי העולם.

להלן תיאור טיסה אחת מתוך מסמך הג'ייסון:

```
{
  "id": 123,
  "airlineCode": "LY",
  "origin": "TLV",
  "destination": "ORD",
  "duration": 5.5,
  "departureTime": "01/01/2021, 22:00",
  "arrivalTime": "02/01/2021, 04:00"
}
```

כתבו קטע קוד ב-spark (pySpark) המקבל את המסמך ומחזיר את מס' הטיסות המגיעות לכל יעד , עבור טיסות בהן משך הטיסה קטן מ-7 שעות.

(ניתן להניח שכבר קיים rdd שנוצר ממסמך ה-json)

This image shows a single sheet of white paper with horizontal blue or grey ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

שאלה 6 - Linear regression (13 נק')

נתון הדאטה הבא בעל פיצ'ר אחד

$$x = [10, 17, 23, 40]$$

$$y = [15, 30, 44, 66]$$

סטודנט א' הציע את פונקציית הניבוי הבאה: $h(x) = x + 6$ סטודנט ב' הציע את פונקציית הניבוי הבאה: $h(x) = 2x + 1$

א. חשבו את ערך ה Loss בכל אחד מהאפשרויות

איזו פונקציית ניבוי עשויה לנבא ערך טוב יותר ל- y עבור דאטה חדש? נמק (5 נק)

ב. שפר את פונקציית הניבוי של סטודנט א' באמצעות הפעלת איטרציה אחת של שיטת מורד הגרדיאנט, חשב עבורה את ה- loss החדש (5 נק)

ג. האם נכון להגיד כי ככל שמוספים יותר פיצ'רים ה- test error או train error קטנים באותה מידה, נמקו (3 נק)
