

## מבוא ללמידת מכונה

### תרגיל 3

#### תאריך הגשה: באתר הקורס

##### שאלה 1

העזרו במצגת שהועברה בקורס וממשו מחלקה הנקראת LogisticRegression שתכיל בנוסף לשאר גם את המתודות `fit`, `predict`, `score`:  
`fit(self, X, y)` – מקבלת את  $X$  מטריצה המכילה את ה `features`,  $y$  מכילה את ה `labels`, ומאמנת מודל בעזרת אלגוריתם רגרסיה לוגיסטית עם גראדינט דיסנט. את וקטור המשקולות  $w$  שמרו בשדה `weights_` של המחלקה כך שתהיה אפשרות לגשת אליו.

`predict(self, X)` - מחשבת את הפלט של המודל על הדוגמאות ב- $X$ .

`predict_proba(self, X)` - מחשבת את ההסתברויות שהמודל מתאים לדוגמאות ב- $X$ .

`score(self, X, y)` - מחשבת את אחוז הדוגמאות ב- $X$  שהרשת המאומנת מסווגת בצורה נכונה, מחזיר תוצאה בין 0 ל-1.

##### שאלה 2

ממשו פונקציית `main` שמשתמשת במחלקה LogisticRegression שכתבתם על ה- `ham/spam dataset` שמופיע ב-`mama2`. המטרה היא להעריך את ההסתברות שהודעת מייל הינה `spam`.

הדפיסו את התוצאות ואת וקטור המשקולות בסוף התהליך.

שימו לב: יש לעשות עיבוד מקדים על ה-`dataset` כך שתוכלו להשתמש בו.

##### שאלה 3

כתבו פונקציה לשרטוט ה-`ROC curve` עבור המחלקה LogisticRegression על ה-`ham/spam dataset`.  
בשימוש במחלקה שלכם לצורך סיווג דואר זבל, איזה `threshold` הייתם בוחרים לצורך תרגום ההסתברויות לסיווג?

##### שאלה 4

- הוסיפו למחלקה LogisticRegression שלכם אפשרות ל-`multiclass classification`. השתמשו ב-`one vs rest` לצורך כך.
- ממשו פונקציית `main` שמשתמשת במחלקה LogisticRegression שכתבתם על ה- `iris dataset`. שימו לב שבדאטהסט הזה יש יותר משתי מחלקות.

אפשר להוריד את הדאטהסט מ-`sklearn`, בעזרת הקוד הבא:

```
from sklearn import datasets
```

```
iris = datasets.load_iris()
```

תזכרו: בעבודה על הדאטהסטים השונים (ham/spam, iris) בחירת הייצוג תמיד בשליטתכם. בין הכלים שניתן להשתמש לצורך הייצוג: כלים שאת מכירים (PCA, scaling, polynomial features, ועוד), וכל רעיון יצירתי אחר שיש לכם.