



**הטכניון**  
מכון טכנולוגי  
לישראל

**למידה עמוקה**

הכנת הנתונים למידול

# הכנת הנתונים למידול

- הנדסת מאפיינים - Feature engineering
- ניקוי והכנת הנתונים
- נרמול



**הטכניון**  
מכון טכנולוגי  
לישראל

**הנדסת מאפיינים**

# הנדסת מאפיינים

- הוספת והורדת משתנים
  - על בסיס הבנת משמעות הנתונים
  - על בסיס קורלציה
- הפיכת משתנים קטגוריאליים למספריים
  - עם סדר
  - ללא סדר
- הפיכת משתנים מספריים לקטגוריאליים
  - אינטראקציות ופולינומים

דוגמה	הסבר	מושג
אם $X_1$ הוא כמות שעות למידה ו- $X_2$ הוא רמת קושי של מבחן, ניתן ליצור משתנה חדש: $X_{\text{new}} = X_1 \times X_2$ .	יצירת תכונות חדשות על ידי הכפלה או שילוב של שני משתנים קיימים.	אינטראקציות בין משתנים
אם $X$ הוא גיל התלמיד, ניתן להוסיף משתנים כמו $X^2$ (גיל בריבוע) ו- $X^3$ (גיל בחזקה שלישית).	הוספת חזקות של משתנים קיימים כדי לאפשר למודל לזהות קשרים לא ליניאריים.	פולינומים (חזקות של משתנים)
ברגרסיה ליניארית, הוספת מאפיינים כאלה יכולה לשפר את דיוק התחזיות. במודלים כמו Random Forest או XGBoost זה פחות נחוץ.	כאשר יש קשרים לא ליניאריים או אינטראקציות משמעותיות בין משתנים.	מתי זה מועיל?

## מחברת ללימוד:

<https://colab.research.google.com/drive/1HQGJWUNgmb-McJaSuxIqH3pFK5KWq8yQ>

# תרגיל 1 – הנדסת מאפיינים

- טענו את בסיס נתונים titanic
- (הערה: ניתן גם בסיסי נתונים אחרים לבחירתכם/ן)
- המטרה: חיזוי השרדות
- בנו מחברת וכתבו קוד בהתאם להנחיות הבאות:
  - אילו מאפיינים תכלילו בסיווג ואילו לא? מדוע?
  - האם אפשר לחשוב על בניית מאפיינים נוספים?
  - אילו מאפיינים קטגוריאליים נדרש להמיר למספריים? כיצד לעשות זאת?



**הטכניון**  
מכון טכנולוגי  
לישראל

**ניקוי והכנת הנתונים**

# ניקוי והכנת הנתונים

- זיהוי ערכים חסרים
- זיהוי חריגים
- זיהוי כפילויות
- שיטות התמודדות

# תרגיל 2 – ניקוי והכנת הנתונים

- המשיכו את פיתוח המחברת בהתאם לשאלות הבאות:
  - האם קיימים ערכים חסרים? אם כן טפלו בהם
  - האם קיימים ערכים חריגים? אם כן טפלו בהם
  - האם קיימות כפילויות? אם כן – טפלו בהם



**הטכניון**  
מכון טכנולוגי  
לישראל

**נרמול**

נושא	הסבר
רקע	Scikit-learn היא ספרייה ללמידת מכונה ב-Python. אחת המשימות המרכזיות היא סטנדרטיזציה (התאמת הנתונים לממוצע 0 ושונות 1). בפעולה זו משתמשים ב- <code>fit_transform()</code> על נתוני האימון וב- <code>transform()</code> על נתוני הבדיקה.
מה עושה <code>fit_transform()</code> ?	מחושב הממוצע והשונות של כל תכונה, והנתונים משתנים בהתאם. נועד רק לקבוצת האימון.
מה עושה <code>transform()</code> ?	משתמש בממוצע והשונות שנלמדו מנתוני האימון וממיר בהתאם את נתוני הבדיקה.
למה זה חשוב?	אם נשתמש ב- <code>fit_transform()</code> גם על נתוני הבדיקה, המודל "יראה" אותם מראש ויהיה מוטת. מה שיפגע בהערכתו על נתונים חדשים.
דוגמה בקוד	שימוש ב- <code>StandardScaler</code> כדי לאמן מודל תוך שמירה על עקרון ההפרדה בין נתוני האימון לבדיקה.
סיכום השיטות	<code>fit_transform()</code> – מחושב ממוצע ושונות ומבצע שינוי (רק על נתוני האימון). <code>transform()</code> – משתמש בממוצע והשונות שנלמדו ומבצע שינוי (רק על נתוני הבדיקה).
מסקנה	השימוש הנכון מבטיח שהמודל ילמד רק מנתוני האימון ויבדוק את עצמו על נתוני הבדיקה כאילו הם חדשים לגמרי.

# תרגיל 3 – נרמול

- המשיכו את פיתוח המחברת בהתאם להנחיות הבאות:
  - בחנו את הסטטיסטיקה של הנתונים בעבודתכם.
  - נרמלו את הנתונים בשיטות שונות
  - איזו שיטה משיגה ביצועים טובים ביותר?
  - האם ניתן להסביר את שיפור הביצועים באמצעות שיטה זו?



**הטכניון**  
מכון טכנולוגי  
לישראל

**שאלות?**