פרויקט סיום

מבוא למדע הנתונים

ד"ר אריה יעקבי

בפרויקט זה התלמידים יישמו באמצעות שפת התכנות פיתון את החומר שנלמד בקורס. הפרויקט יכלול קריאה של מקורות מידע מסוגים שונים (כגון CLSX ,CSV ועוד), ניקוי הנתונים והכנתם לקראת ניתוחים שונים, וביצוע ניתוחים בסיסיים על נתונים אלה.

הנחיות לפרויקט

1. בחירת מקורות המידע:

- 1.1. יש לבחור קובץ נתונים בכל תחום שנראה לכם. קראו לו input1_df. הקובץ יכול להיות בפורמט 2.1. עמודות נתונים. ניתן להוריד קבצים XML ,JSON ,XLSX בקובץ חייבים להיות לפחות 1000 רשומות ו- 5 עמודות נתונים. ניתן להוריד קבצים כאלה ברשת הווב. ישנן דוגמאות רבות של אתרים בהם ניתן להוריד קבצי נתונים באתר:
 https://www.kdnuggets.com/datasets/index.html
 - 1.2. מצא את אחוז הערכים החסרים. ניתן לבצע זאת על ידי מציאת מספר הערכים החסרים

df.isnull().sum()

- 1.3. אם אחוז הערכים החסרים קטן מ- 10% אז קובץ נתונים זה לא מתאים לפרויקט זה ותחפשו קובץ אחר.
 - 1.4. יש לשים לב שהקובץ אינו נקי מלכתחילה לדוגמה בדקו האם ישנם ערכים חסרים, ערכים שלא רשומים בפורמט נכון בעמודה (לדוגמה מספר במקום שם). אם הקובץ כבר נקי ומסודר פסלו אותו. לא אכפת לי שתיקחו קובץ שניקו אותו ופשוט תקלקלו אותו באופן משמעותי.
 - 1.5. לאחר קריאת הקובץ יש להדפיס את 5 השורות הראשונות, האחרונות ובאיזה שהוא מקום באמצע שתבחרו.
 - באמצעות הפונקציה data-frame באמצעות הפונקציה תיאורית של הנתונים ב- data-frame באמצעות הפונקציה .describe()
 - 1.7. יש לבחור קובץ נתונים <u>נוסף</u> בכל תחום שנראה לכם. קראו לו input2_df. הקובץ יכול להיות בפורמט. אוחור בכל תחום שנראה לכם. קראו להיות נתונים. וגם געודות נתונים. וגם בקובץ חייבים להיות לפחות 1000 רשומות ו- 5 עמודות נתונים. וגם עבורו בצעו את מה שכתוב בסעיפים 1.2 1.6
- 1.8. צרו data-frame ע"י פעולת Full Outer Join, קראו לו outer_join_df. אם אין להם מפתח משותף, תדאגו להוסיף עמודה מתאימה לאחד מה- data-frame כך שבשני הדטהפרמס תהיה עמודת מפתח עם ערכים דומים (כמו לדומה: מס מוצר, או מס זיהוי,...), כך שתוכלו להשתמש במתודה

df1.merge(df2, on='id', how='outer')

בדקו שוב האם ישנם ערכים חסרים בקובץ הממוזג.

2. צרו data-frame 2 חדשים ע"י בחירת עמודות מתוך ה- data-frame שנותר מפעולת ה- JOIN בסעיף data-frame 2 בסעיף. צרו data-frame 2 בחירת עמודות מתוך מלוץ, df2, df1 בהתאמה. אנחנו נעבוד על df2, df1 בסעיפים הבאים.

2. ניקוי וארגון קובץ הנתונים:

עבור ה - df2, df1 data-frames בצעו את הניקיון של הנתונים והפורמט שלהם לפי הדרישות הבאות:

- 2.1. אם בעמודה מספרית ישנם ערכים לא מתאימים כמו לדוגמה ערכים בולאניים, או סטרינגים, יש למצא אותם ולהוריד את השורות עם ערכים אלה (אם במקור df2, df1 frames אינם מכילים ערכים "לא מתאימים" תדאגו להכניס אותם). יש לוודא ששורות אלה נמחקו ולסדר את האינדקס של השורות באופן הנכון. כמו כן יש לדאוג שכל הערכים העמודה יהיו מסוג float64 ויש לוודא זאת.
- 2.2. אם בעמודה שיש בה סטרינגים (כמו לדוגמה עמודת שמות) ישנם ערכים לא מתאימים כמו לדוגמה ערכים לא מתאימים כמו לדוגמה ערכים בולאניים, או מספריים, יש למצא אותם ולהוריד את השורות עם ערכים אלה (אם במקור df2, df1 אינם מכילים ערכים "לא מתאימים" תדאגו להכניס אותם) . יש לוודא ששורות אלה נמחקו ולסדר את האינדקס של השורות באופן הנכון. כמו כן יש לדאוג שכל הערכים העמודה יהיו מסוג *string* ויש לוודא זאת.
 - 2.3. לאחר מכן ניתן לוודא שאם ישנם ערכים חסרים (תאים ריקים אן ערכי null) יש להחליף ערכים אלה באמצעות הממוצע, השכיח, או החציון של הערכים בעמודה. ניתן להשתמש במתודה: fillna . באמצעות הממוצע, השכיח, או החציון של הערכים בעמודה ניתן להשתמש במתודה: מי שיחליף ערכים אלה באופן שונה על יד לדוגמה מציאת הממוצע בין המספר הקודם למספר הבא. לדוגמה, אם באמצע העמודה יש את המספרים:

100

N\A

200

.150 = 2/(100+200) בממוצע של N\A- אז אפשר להחליף את ערך ה-N\A

שימו לב למה עושים במקרים חריגים כמו לדוגמה כיש רצף של ערכי null בעמודה אז אי אפשר לחשב את הממוצע של המספר הבא, או כשמדובר במספר האחרון בעמודה. במקרים אלה חשבו על פתרון יצירתי יותר.

- 2.4. $\frac{1}{1200}$ (df1 data-frames בה יש האחת מהעמודות באחד מה (df2, df1 data-frames). בה יש מספרים. עברו בלולאה על העמודה, החליפו כל מספר בעמודה במספר (בין $\frac{1}{1200}$) שהוא תוצאת חלוקת מספרים. עברו בלולאה על העמודת המספרים. לדוגמה, אם המספר הנכחי בעמודה הוא $\frac{750}{1200}$. אם המספר והמספר המקסימלי בעמודה הוא 1200, אז נחליף את 750 במספר $\frac{750}{1200}$. אם המספר הוא $\frac{750}{1200}$ (200) חלקי המספר המקסימלי בסדרה $\frac{1}{1200}$.
- 2.5. הדפיסו את כל השורות הכפולות (אם אין אז תכניסו שורות כאלה באמצעות קוד פייתון) והורידו אותן על ידי: (\drop duplicates.

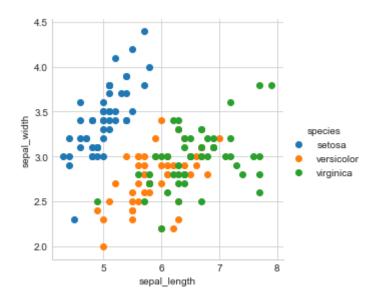
3. הצגה ויזואלית של הנתונים:

מספרים בפני עצמם אינם נוטים לעורר תגובה רגשית. אבל הדמיית נתונים יכולה לספר סיפור שנותן משמעות לנתונים. בעוד שהשוואת מספרים עשויה להרשים את הקוראים, חיזוק המספרים הללו באמצעות דיאגרמות עוזר להשפיע עוד יותר. קישור: איך לספר סיפור עם נתונים (מדריך למתחילים):

/https://venngage.com/blog/data-storytelling

עכשיו כשיש לנו df2, df1 data-frames שעברו ניקיון של הנתונים, בחרו אחד מהם או את שניהם כרצונכם, והציגו 10 דיאגרמות שונות של הנתונים עשו זאת כפי שעשינו באפליקציית הדוגמה של המניות בקורס (שיעור 11). ציירו כמה שיותר גרפים מכמה שיותר סוגים והסיקו מסקנות על הנתונים בהתאם ורשמו אותן. <u>גרפים מיוחדים ומסקנות יפות יזכו **בבונוס**.</u>

גרפים פשוטים כגון:



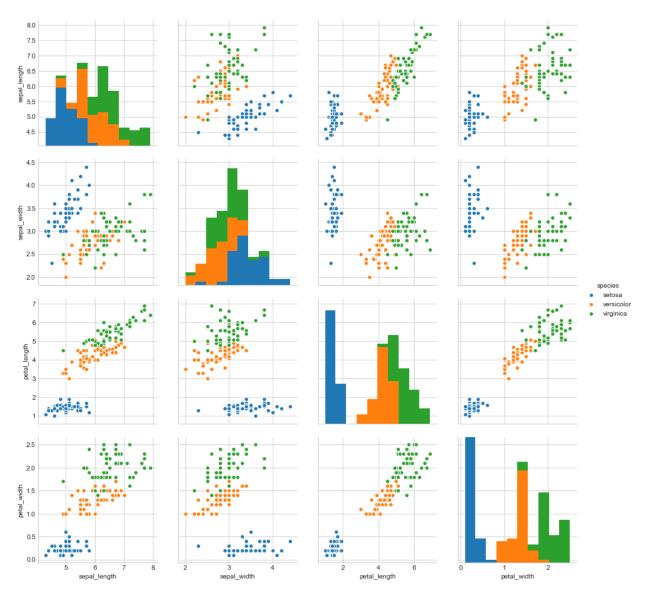
שהקוד הבא מציג:

```
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.datasets import load_iris

iris = load_iris()
print(iris.head())

plt.scatter(iris['sepal_length'],iris['sepal_width'],color=['r','b','g'])
plt.xlabel('Sepal length')
plt.ylabel('Sepal width')
plt.title('Scatter plot on Iris dataset')
```

ועד דוגמאות של גרפים מעניינים:



שמציגים 4 עמודות, אלה עם אלה. הקוד שמציג גרף זה:

```
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.datasets import load_iris

iris = load_iris()

sns.set_style("whitegrid");
sns.pairplot(iris, hue="species", size=3);
plt.show()
```

ממליץ מאוד על הספרי: /https://plotly.com/python/plotly-express שבה יש דוגמאות רבות וקוד מוכן.

https://towardsdatascience.com/data-visualization-for-machine-learning-and-data-:במסמך science-a45178970be7

יש דוגמאות שונות ויפות. דוגמאות מעניינות נוספות יש ב:

https://plotly.com/python/plotly-express/ https://matplotlib.org/stable/gallery/index

https://seaborn.pydata.org/generated/seaborn.jointplot.html

https://seaborn.pydata.org/

4. ניתוח הנתונים:

- 4.1. על אתם DF שבחרתם להשתמש בסעיף הקודם, בצעו חלוקה לאשכולות של הנתונים לפי חתכים שונים והציגו בגרף את התוצאות. כמובן שחשוב שתשתמשו במבחן המרפק (כולל גרף) על מנת לקבוע את מספר האשכולות הנכון. הסיקו מסקנות מתאימות (לפחות 5 מסקנות). בונוס יינתן ליותר מסקנות ולמסקנות מעניינות יותר.
 - 4.2. יש לבצע רגרסיה ליניארית (זה מופיע בשקפים) ולצייר את זה. הציגו והסבירו את מסקנות הרגרסיה.
 - 4.3. כל ניתוח נוסף מסוג אחר ומעניין יזכה **בבונוס**.

5. מודולים או אובייקטים שחובה להשתמש בהם בפרויקט הם:

- (רשות)Xmltodict .5.1
 - 5.2. האובייקט
 - Pandas .5.3
 - plotly.express .5.4
 - Matplotlib .5.5
 - Seaborn .5.6

6. מועד ואופן הגשה:

- 6.1. מועד הגשת הפרויקט יירשם באתר העבודות של הקורס.
- 6.2. כל דמיון חלקי או אחר בין התרגילים של תלמידים שונים יגרור פסילה של שתי העבודות ומכאן פסילת זכאות לקורס של צוותי התלמידים.
 - 6.3. הפרויקט יוגש באתר. בדקו במועד שהתרגיל הוגש כראוי. (בדקו שהפרויקט מופיע באתר לאחר העלאתו).
- 6.4. יש להגיש את הפרויקט במועדו . לא יתקבלו טענות מכל סוג כולל "חשבנו שהגשנו, טעינו בכתובת, יש לחץ או בחינות בסוף הסמסטר וכיו"ב". הערכו מראש ועמדו בקצב הזמנים המוצג.
 - 7. כיצד להגיש את התרגיל המסכם?
 - 7.1. יש לצרף לעבודה קובץ WORD שיכלול:
 - . כל הערה שברצונכם שאני אקרא טרם בדיקת התרגיל.

- שם הפרויקט, שמות המגישים ומספרי זיהוי, משפט המתאר את הפרויקט שלכם, הסברים וההערות לפי צורך.
 - . את שם קובץ הנתונים WORD ציינו ב
 - . כל גרף ומסקנה שבקשתי יכנסו בקובץ הוורד.
- 7.2. יש להכניס את קובץ הוורד כל קבצי הפיתון והנתונים לתיקיה ששמה מורכב משמות התלמידים ובתוכה יהיו כל הקבצים המוזכרים לעיל.
 - 7.3. יש לדחוס (לקווץ) את התיקייה על ידי תכנת WinRAR או ZIP. למי שאין שיתקין זה חינם.
 - 7.4. יש להעלות את התיקיה המכווצת לאתר כפתרון לעבודה.

8. כללי:

- 1.1. צוותי הפרויקט ימנו עד ארבעה סטודנטים.
 - 1.2. אי הגשת פרויקט תגרור כישלון בקורס.
- 1.3. התלמידים יצטרכו להגן בעל פה על הפרויקט. הגנה זו משמעותה הוא 40% מהציון הסופי של הפרויקט. בהגנה זו אנו יכולים לשאול כל שאלה הקשורה בפרויקט או בחומר (ההגנה היא עם חומר ומחשב פתוח וניתן להיעזר בהם ותתקיים בזום). הערות מצד סטודנטים כגון, "אני לא עשיתי את החלק הזה" יגררו הורדת נקודות. זמן ההגנה יכול לנוע מזמן קצר מאוד (דקות בודדות) עד הגנה ארוכה יותר (כרבע שעה).

<u>האחריות על למידת כל החומר והביצוע היא על כל הצוות.</u> לא ניתן לבצע משימה או ללמוד חומר מתקדם מבלי ללמוד היטב את החומר הקודם לכך.

הפרויקט יוגש באתר העבודות של הקורס. בדקו במועד שהפרויקט הוגש כראוי (בדקו שהתרגיל מופיע באתר לאחר העלאתו). לא נקבל טעויות בהגשה – בדקו בבקשה 5 פעמים לפני שאתם מגישים.

יש לשמור את העבודה אצלכם לגיבוי לפחות עד קבלת הציון.

<u>בהצלחה!</u>

