



מדעי הנתונים ובינה עסקית, סמסטר ב' תש"פ

תרגיל בית 4 – שפת Python

קראו בעיון את כל ההוראות לפני ביצוע העבודה

הוראות כלליות:

- א. אי עמידה בכל אחת מההוראות יגרור הורדת ציון או פסילת העבודה.
- ב. הגשת העבודה בזוגות בלבד.
- ג. שפת תכנות – **Python 3.x**, סביבת פיתוח – **PyCharm** JetBrains גרסה 2019.x ומעלה. יש לוודא כי Anaconda ו-PyCharm מותקנים ע"פ קובץ ההתקנה של תוכנות מעבדה שנמצא ב-moodle.
- ד. יש להגיש את העבודה לתיקיית ההגשה הרלוונטית באתר הקורס (Moodle).
אחריותכם האישית לבדוק לפני הגשה כי כל הקבצים נפתחים כראוי.
- ה. יש להגיש קובץ zip - שם הקובץ יהיה מורכב משני מספרי תעודות הזהות של המגישים באופן הבא: ID1_ID2.zip
הקובץ יכיל את הקבצים הבאים:
 - הפרויקט המלא: קבצי קוד + GUI, חשוב : ללא קבצי הנתונים.
 - קובץ PDF (המכיל את שמות הסטודנטים ותעודות הזהות) של הדוח המתאר את מבנה הפרויקט שיצרתם ותפקיד של כל מחלקה ושיטה בפרויקט.
- ו. בנוסף, זוהי עבודה תכנותית ולפיכך יהיה משקל לכך בבדיקה. כלומר: יש לדאוג להערות בקוד, הסבר לפונקציות, חלוקה למחלקות, פונקציות קצרות וענייניות וכדומה.
- ז. שימו לב שתבצע בדיקה לאיתור עבודות מועתקות (גם אם חלקית). הקפידו לא לשתף קטעי קוד!
- ח. שאלות בנוגע לתרגיל יש לשאול אך ורק בפורום השאלות הרלוונטי המופיע ב-moodle (ולא במייל - שאלות במייל לא יענו).
- ט. מועד הגשת התרגיל – 27/6/2020 23:55



הוראות התרגיל:

בתרגיל זה עליכם להשתמש בספרייה Scikit-learn של Python על מנת לבצע clustering לקובץ נתונים. בנוסף, תטפלו ברשומות עם ערכים חסרים, כחלק מתהליך ניקוי הנתונים ותתרגלו עבודה עם מבני טבלה (Dataframe) שונים של Python. את פלט האלגוריתם תציגו באמצעות כלי הוויזואליזציה של הספרייה Plotly.

תיאור הקבצים שלרשותכם:

1. **Dataset info** – מידע כללי בנוגע לבסיס הנתונים ממנו לקוחים נתוני התרגיל. קובץ זה הינו לשימושכם בלבד ולא ישמש כנתון שעל תכניתכם לקרוא במהלך הריצה.
2. **Dataset** – קובץ האימון לאלגוריתם ה-clustering, בפורמט .xlsx.

תיאור המשימות שעליכם לממש במסגרת התרגיל:

1. ממשק משתמש פשוט שיוצג עם הרצת התכנית. הממשק יכיל:

- 1.1. הזנת ה-path לקובץ נתוני התרגיל (יש לממש אפשרות זו בעזרת browser).
- על הממשק להכיל תיבת טקסט אחת בלבד, אליה יוכנס הנתוב. במידת הצורך, יש לשמור פלטים בתיקיה זו. הטקסט אשר יופיע על הכפתור יהיה "Browse".
- 1.2. תיבת טקסט בה ניתן להזין את כמות ה-clusters אליהם יחולקו הנתונים. שם תיבת הטקסט יהיה "Number of clusters k".
- 1.3. תיבת טקסט בה ניתן להזין את כמות הריצות של האלגוריתם מ-seeds רנדומליים שונים. שם תיבת הטקסט יהיה "Number of runs".
- 1.4. לחצן לטעינת קובץ הנתונים, הכנתו וניקויו. הטקסט אשר יופיע על הכפתור יהיה "Pre-process".
- 1.5. לחצן לבניית מודל ה-Kmeans והצגת הוויזואליזציה של תוצאותיו. הטקסט אשר יופיע על הכפתור יהיה "Cluster".



2. הכנת הנתונים (תהליך 1.4):

2.1. עם לחיצה על הלחצן המתאים, התכנית תקרא את קובץ הנתונים הנתון (בפורמט xlsx) ותטען אותו למבנה נתונים מסוג Dataframe.

2.2. לאחר קריאת הקובץ, יתבצע תהליך ניקוי הנתונים:

- א. יש להשלים ערכים נומריים חסרים בערך הממוצע של כל ערכי התכונה.
 - ב. יש לנרמל את כל ערכי קובץ הנתונים לערך הסטנדרטי שלהם (חיסור הממוצע וחלוקה בסטיית התקן). פעולה זו מכונה Standardization.
 - ג. קיבוצ הנתונים לפי התכונה "country". כחלק מתהליך זה, יש ליצור רשומה אחת עבור כל מדינה, כך שתמצע את ערכי התכונות על פני השנים (תכונה year). יש ליצור dataframe חדש שבו כל רשומה מסמלת מדינה וערכי הפיצ'רים יהיו הממוצע על פני כל הרשומות השייכות לאותה מדינה (לכל מדינה יש X רשומות שמתארות את המדינה על פני השנים). אחרי הקיבוצ אפשר להתעלם מהפיצ'ר year.
- 2.3. בשלב זה יופיע dialog אשר יכיל את הודעה "Preprocessing completed successfully!" המודיעה על סיום הכנת הנתונים ואפשר למשתמש ללחוץ על "OK" להמשך.

✓ ניתן להשתמש בספרייה הייעודית של Scikit-learn עבור ניקוי הנתונים.

3. חלוקת הנתונים לאשכולות (תהליך 1.5):

- 3.1. קובץ הנתונים ישמש לבניית מודל k-means באמצעות הספרייה הייעודית של Scikit-learn. קראו את התיעוד של המחלקה `sklearn.cluster.KMeans` באתר הרשמי של הספרייה: <http://scikit-learn.org/stable/modules/generated/sklearn.cluster.KMeans.html>
- 3.2. יש להפעיל את האלגוריתם Kmeans על סט האימון עם הפרמטרים הבאים:
 - א. מספר האשכולות (`n_clusters`) שהוכנס ע"י המשתמש בתיבת הטקסט המתאימה (תהליך 1.2 בממשק המשתמש).
 - ב. מספר הריצות (`n_init`) שהוכנס ע"י המשתמש בתיבת הטקסט המתאימה (תהליך 1.3 בממשק המשתמש). כל ריצה מתחילה כברירת מחדל עם ערכי centroid רנדומליים.
- 3.3. יש להצמיד את פלט האלגוריתם לכל רשומה (מדינה) ב-dataframe שיצרתם לאחר שלב ניקוי הנתונים.
- 3.4. יש ליצור שני פלטים המסכמים את תוצאות הריצה של האלגוריתם:
 - א. יש ליצור תרשים פיזור (scatter) של ערכי התכונה Generosity כתלות בערכי התכונה social_support. יש לצבוע את הנקודות על פי ערכי הפלט המתאים מהאלגוריתם.



- ❖ יש להקפיד להוסיף כותרות מתאימות לצירים ולתרשים כולו.
- ❖ יש להשתמש בפונקציה `scatter` של הספרייה `matplotlib`.
- ב. יש ליצור `horopleth map` (מפת מדינות) המדגימה את פלט האלגוריתם (חלוקה של המדינות לאשכולות) עבור המדינות בקובץ הנתונים.
- ❖ יש להשתמש בספרייה `Plotly` ע"י התקנתה (בהנחה שפלטפורמת ה- `Anaconda` מותקנת על גבי `python`):
 - יש להתקין את הספרייה באמצעות הפקודה `pip install plotly` בחלון ה-`Terminal`.
 - יש להירשם (חד פעמי עם אי-מייל כלשהו) באתר <https://plot.ly/accounts/login/?action=login>
 - יש להיכנס למסך ה-`Settings` אחרי `login` לאתר `Plotly` וללחוץ על תת-תפריט ה-`API Keys` ואז על `Regenerate Key` כדי לקבל את ה-`APIKEY`.
 - שם המשתמש בתוספת המפתח יאפשר לכם לשמור את פלט מפת המדינות בחינם (עד 25 תרשימים) כקובץ על המחשב.
- ❖ יש להיעזר בדוגמא הנתונה בקישור הבא: <https://plot.ly/python/choropleth-maps>
- ❖ יש לשמור את תרשימים המדינות כתמונה סטטית תוך שימוש בפונקציה הבאה:


```
import plotly.plotly as py
py.sign_in(username, API key)
py.image.save_as(choromap, filename='name.png')
```
- 3.5. יש להציג את שני הפלטים זה לצד זה במסך ה-`GUI`.
- 3.6. יש להציג `dialog` נוסף שיעדכן שתהליך ה-`clustering` הסתיים. לחיצה על "OK" תסיים את ריצת התכנית.

הטקסט על הכפתורים לא ניתן לשינוי, וחשוב שיהיה זהה למוגדר לעיל.
ניתן להניח שהפעולות יבוצעו בסדר הנכון – הכנת הנתונים ואז חלוקה לאשכולות.

כותרת כל החלונות (כולל הדיאלוגים שפורטו לעיל) צריכה להיות
"K Means Clustering"

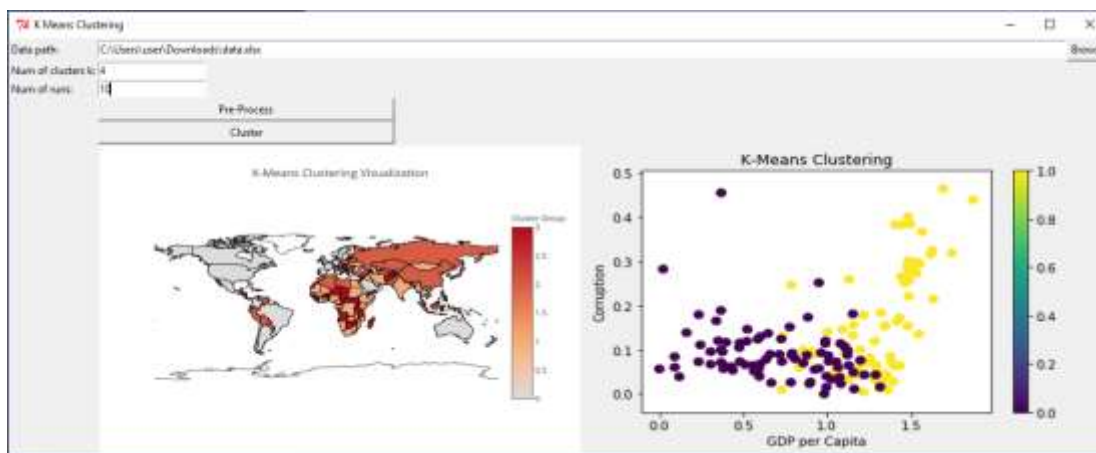


4. דוח:

יש לצרף דוח המתאר את המחלקות והשיטות שכתבתם בקוד.

הערות חשובות נוספות:

- מצורף תצלום חלון ה-GUI הנדרש לתרגיל זה. מומלץ להשתמש בממשק פשוט ביותר של Tkinter (from Tkinter import *). התמונה להמחשה בלבד. ממשק המשתמש לא חייב להראות כבתמונה, אך צריך להיות אינטואיטיבי, נוח ופשוט.



- את חלון ה-GUI אין צורך לבנות מאפס, אפשר להשתמש בקובץ calculator.py שהוצג במעבדה 6 (introduction python) ולשנות את ה class הנתון (להוסיף/ לשנות שדות וכפתורים על פי הנדרש).
- ההנחה היחידה בעבודה היא כי סדר הפעולות הנדרשות יתבצע בסדר הנכון, כלומר קודם הכנת הנתונים ("Pre-process") ורק אז החלוקה לאשכולות ("Cluster").
- על התכנית לדעת להתמודד עם שגיאות כמו למשל קובץ נתונים ריק, מספר לא תקין בתיבות הטקסט השונות. במקרה של נתון לא תקין, יש להציג הודעת שגיאה מתאימה (המעידה על סוג השגיאה) ולא לאפשר לחיצה על כפתור ה-"Cluster".
- יש לבצע בדיקות קלט מלאות לכל השדות בממשק המשתמש ולכל הפרמטרים של האלגוריתם. היעזרו בתיעוד כדי לדעת מה טווח הערכים שמקבל כל פרמטר.
- הקוד ייבדק על קובץ נתונים דומה במבנה (אותם עמודות) אך עם ערכים שונים מהקובץ לדוגמא שניתן לכם.
- אינכם נדרשים להתקין חבילות תוכנה נוספות חוץ מ-plotly (בהנחה ש-Anaconda מותקנת). אין להשתמש בחבילות שאינן קיימות בפלטפורמת ה-Anaconda וב-Plotly (כי הן מחייבות התקנה נפרדת).