

Clustering – K-means

Clustering – K-means

אתם עובדים כחברי צוות אנליסטים בחברת – "CardWise" חברת כרטיסי אשראי מובילה בישראל. בשנים האחרונות חל גידול משמעותי בבסיס הלקוחות של החברה, אך יחד עם ההתרחבות ניצבה החברה בפני אתגר: להבין טוב יותר את התנהגות הלקוחות השונים כדי לייעל שיווק, למזער סיכונים ולאפשר הצעות ממוקדות.

בהתבסס על מערך הנתונים הפנימי של החברה (קובץ, Customer Data.csv) תנסו לזהות קבוצות לקוחות דומות זו לזו (segments) באמצעות שיטות clustering, בעיקר KMeans. כך תוכלו בתום התרגיל להציע לקהלי יעד מותאמים: למשל, מבצעים ללקוחות שמתמשים בתשלומים חד-פעמיים, או הצעות אשראי ללקוחות בעלי מסגרות גבוהות אך משתמשים בתדירות נמוכה.

לאחר שתבצעו ניתוח חקר נתונים, קדם-עיבוד ובחירת מספר הקלאסטרים, המטרה היא:

1. לספק פרופיל תמציתי לכל קבוצה

2. להמליץ לצוות השיווק על מסרים ומבצעים

להלן פירוט ה-features:

משתנה	תיאור
CUST_ID	מזהה ייחודי של הלקוח.
BALANCE	סך היתרה (במטבע) בחשבון הכרטיס, ממוצע חודשי לאורך תקופת ה-TENURE.
BALANCE_FREQUENCY	תדירות העדכון של היתרה (0-1): ערך קרוב ל-1 אומר שהיתרה מתעדכנת כמעט בכל חודש; ערך קרוב ל-0 – לעתים רחוקות.
PURCHASES	סכום כל הרכישות שבוצעו בכרטיס (מבטא את היקף ההוצאות החודשיות הכולל).
ONEOFF_PURCHASES	סכום הרכישות שבוצעו כעסקאות חד-פעמיות (ללא תשלום בתשלומים).
INSTALLMENTS_PURCHASES	סכום הרכישות שבוצעו בתשלומי פיצול (installments).
CASH_ADVANCE	סכום המשיכות במזומן (Cash Advance) שבוצעו בכרטיס.
CASH_ADVANCE_TRX	מספר המשיכות במזומן שבוצעו.

משתנה	תיאור
CASH_ADVANCE_FREQUENCY	תדירות המשיכות במזומן (0-1): ערך גבוה מצביע על קיום משיכות קרובות זו לזו.
PURCHASES_FREQUENCY	תדירות כללית של רכישות (0-1): כמה קרובות זו לזו העסקאות ביחס לאורך ה-TENURE.
ONEOFF_PURCHASES_FREQUENCY	תדירות עסקאות חד-פעמיות. (0-1)
PURCHASES_INSTALLMENTS_FREQUENCY	תדירות עסקאות בתשלומים. (0-1)
PURCHASES_TRX	מספר העסקאות הכולל שבוצעו ברכישות (כולל חד-פעמיות ובתשלומים).
CREDIT_LIMIT	מסגרת האשראי המוקצית ללקוח (Credit Card Limit).
PAYMENTS	סך התשלומים שבוצעו על ידי הלקוח (Pay-off Amount), לאורך תקופת ה-TENURE.
MINIMUM_PAYMENTS	סכום התשלום המינימלי שדרש הבנק בכל חודש, ממוצע לאורך כל ה-TENURE.
PRC_FULL_PAYMENT	אחוז החודשים שבהם הלקוח שילם את היתרה במלואה. (0-1)
TENURE	משך ההחזקה של הכרטיס (במספר חודשים), לאורך נמדדו שאר המשתנים.

חלק א' – טעינת הנתונים ובדיקה ראשונית

1. טען את הקובץ ל- DataFrame
2. הצג head(), info() ו- describe() לקבלת תמונה על סוגי העמודות, ערכים חסרים וסטטיסטיקות כלליות.
3. בדוק האם קיימים ערכי חסר בעמודות, וכיצד תטפל בהם (הסבר שיטת טיפול – מחיקה או השמה בערך ממוצע/מדיאני).

חלק ב' – ניתוח חקר נתונים (EDA)

1. חשב ופרש את מטריצת הקורלציה (correlation) בין התכונות העיקריות (למשל: BALANCE, PURCHASES, CASH_ADVANCE, CREDIT_LIMIT, PAYMENTS).
2. ויזואליזציה:
 - היסטוגרמות והתפלגות של BALANCE ו- PURCHASES כדי לזהות א-סימטריות או אוטליירים.
 - Scatter plot בין PURCHASES לבין CASH_ADVANCE

3. זהה נקודות קיצון (Outliers) בעזרת תיבות whisker (boxplot) עבור עמודת CREDIT_LIMIT ועמודת CASH_ADVANCE.

חלק ג' – קדם-עיבוד והתאמה למודל

1. בחר את קבוצת העמודות הרלבנטיות להצבה במודל (לדוגמה, תכונות כספיות ותדירות רכישות). בצעו מחקר קטן באמצעות deep research לראות מה יכול להיות רלוונטי ומה פחות.
2. בצע Scale/Standardization (StandardScaler) כדי לאזן את סקאלת התכונות (במידת הצורך).

חלק ד' – מציאת מספר הקלאסטרים Elbow

1. בחר את הערך האופטימלי של k והסבר את הבחירה.

חלק ה' – אימון מודל KMeans ופרופילינג של קלאסטרים

1. אמן את KMeans עם k הנבחר ו-`random_state=42`.
2. הוסף ל-dataframe עמודת Cluster עם התוויות שהוקצו לכל לקוח.
3. חשב את המאפיינים הממוצעים (centroids) של כל קלאסטר והצג בטבלה (אפשר להשתמש ב-groupby):
 - גודל קלאסטר (מספר לקוחות)
 - ממוצע, CREDIT_LIMIT, CASH_ADVANCE, PURCHASES, BALANCE, PAYMENTS
4. פרש כל קלאסטר:
 - למשל: איזה סגנון הוצאה ותדירות רכישות מאפיין אותו? האם זה High-Value “Cash-Advance Seekers”, “Frequent Spenders”, “Customers וכד'?

חלק ו' – סיכום ומסקנות עסקיות

1. אילו תובנות עסקיות ניתן להסיק מהתפלגות הקלאסטרים?
2. כיצד ניתן ליישם את המודל לשיפור שיווק ממוקד, ניהול סיכונים אשראי או מבצעים מותאמים?

בהצלחה!