

1. כי לcold, אני يولיה והפרויקט שלי עוסק בחיזוי רמת הפופולריות של שירים בספוטיפי.
2. שאלת המחקר שלי הייתה: מה גורם לשיר להיות פופולרי - הסאונד או החשיפה שלו? המטרה הייתה לחזות את ממד הפופולריות שהוא משתנה רציף ונמדד בסקללה בין 0-100 **באטען מאפייני שמע ומידע על פלילייטים, אלבומים ואמנים**. עבדתי על DATA של כ-33 אלף שירים, מתוך קגל, כולל 23 פיצרים ראשוניים.
3. תהליך העבודה כלל הכתנת הנתונים לfile Flat כאשר הורדתי כל מיניAIMAGES כמו להבה, ניתוחEDA, ניקוי חרגים והשלמת MISNIG, יצירת פיצרים חדשים, בחירת פיצרים חשובים ולבסוף בניית מודלים וfine tuning. כל השלבים נעשו בPIPELINE דרך גוגול קולאבר.
4. בשלב-הEDA ראייתי שרוב השירים יוצאו ביום שישי ובחודש ינואר. רוב השירים נמצאים בטוווח ציוני פופולריות של 60-20. הפיקים ב-2008 ו-2014ocab מצביעים על % השירים עם פופולריות 0 כאשר שנת 2008 היא שנת השקתה הספוטיפי. חלק מהמשתנים כמו ספיינס, אקוסטיות, אינסטרומנטליות הרואו זנב ימני, יכולمر רוב השירים לא אקוסטיים או דיבוריים, אבל יש קבוצה קטנה שכן.
5. כאן בדקתי את הקשרים בין כל הפיצרים המספריים. רואים למשל קשר חזק בין אנרגיה ורעש - לאודנס. לעומת זאת קשר שלילי בין אקוסטיות לאנרגיה- ככל מר שירים יותר אקוסטיים הם פחות אנרגטיים. קשר חזק מתון בין עד כמה השיר נשמע שמח או אופטימי לבין עד כמה השיר מתאים לירוק. הקשרים של המשתנים עם הפופולריות עצמה חלשים. קורלציה 1 זה משך השיר – שזה אותו פיצ'ר ובסוף השארתי 1.
6. ניקוי החרגים במשתנים רציפים נעשה בשיטה של QI כולל קורלציה והתפלגות. מה שינוי התפלגות אף לא קורלציה הוסר והושלם באמצעות מודל MICE.
7. כאן יצרתי משתנים חדשים. התווסףו עוד מהשלב הקודם המשתנים טמפרליים כמו שנה, חודש ויום בשבוע. אחר כך משתני חשיפה- כמה שירים יש לאמן, כמה באלבום, וכמה פלילייטים כוללים את השיר. בנוסף יצרתי משתנים מבוססי טקסט מהמלילים הנפוצות בשם השיר. **למשל המילה "original"** היא דוגא קשורה לשירים פחות פופולריים מהמומצע. בסוף שלב Feature engineering, לאחר סינון משתנים חופפים, נותרו 57 פיצרים.
8. ניסיתי שלוש גישות לבחירת פיצרים:
 1. – **One-Hot Encoding** – אחרי שזרקתי את עמודות שמות האמן והפלילייט.
 2. – **Target Encoding** – חד עם הפיצרים החזקים מהגישה הראשונה.
 3. **גישה משולבת Target Encoding** – חד עם הפיצרים נבחרו לפי הסכמה של לפחות שני מודלים בין : Lasso, Ridge, Gradient Boosting, Random Forest.
9. בסוף נבחרה הגישה של **Target Encoding** כי היא גם הסבירה הכי הרבה מהשונות בפופולריות — בערך 55%. זה היה אחרי Fine-Tuning XGBoost ל-216 **שילובים ב-GridSearchCV**.
10. אפשר לראות בתרשימים שהגורמים המשפיעים ביותר הם החשיפה בפלילייטים וביצוע האמן, הרבה יותר מאשר מאפייני הסאונד.
11. ראייתי שהפקטור הכי משמעותי לפופולריות של שיר הוא החשיפה שלו — דרך האמן או הפלילייטים שהוא מופיע בהם.
12. **המאפיינים של הסאונד עצמו**, כמו קצב או אנרגיה, כן משפיעים, אבל הרבה פחות. המודל מצליח להסביר בערך 55% מהשונות בפופולריות, מה שאומר שיש עוד גורמים שלא נכללו – כמו שיווק, פרסום או טרנדים חברתיים.
13. יש גם מגבלות שצורך לזכור:
 1. ממד הפופולריות עצמו נקבע על ידי ספוטיפי והוא **לא ממד אובייקטיבי**.
 2. שני שירים עם אותה כמות השמעות יכולים לקבל ציוני פופולריות שונים, לפי פרמטרים כמו **מתי האזינו לאחרונה, כמה זמן שמעו, וכמה שיתפו**.
 3. בנוסף, פופולריות היא **משהו דינמי** — מה שפופולרי היום לא בהכרח יהיה פופולרימחר.
 4. לסיכון אפשר לראות שפופולריות היא **לא רק עניין של סאונד, אלא הרבה יותר עניין של חשיפה**.