

מבוא לבינה מלאכותית: מבחן בית - הפורמט הקצר

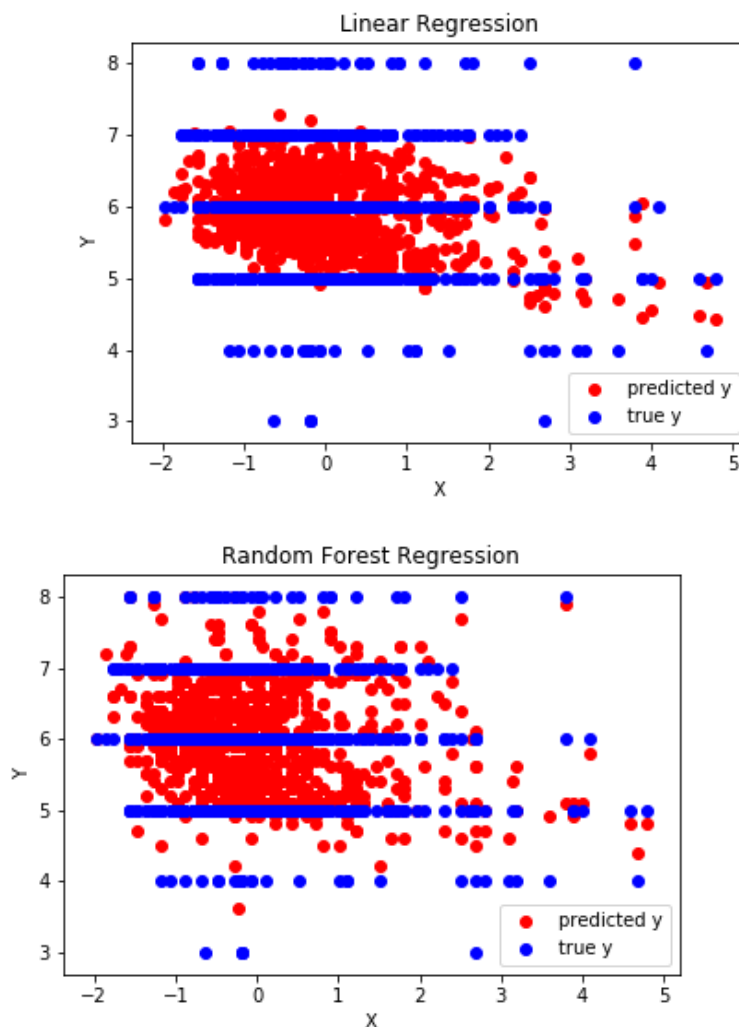
יוסי לוי ת.ז. 037597028

13.08.2020

1

כל ערכי הפיצ'רים הם ערכים מנורמלים. בכל חלקי השאלה

1. בחלק זה נתבקשנו להריץ את הנתונים על שני מודלים ולחזות את איכות היין. המודל הראשון שנבחר הוא רגרסיה ליניארית והשני הוא Random Forest.

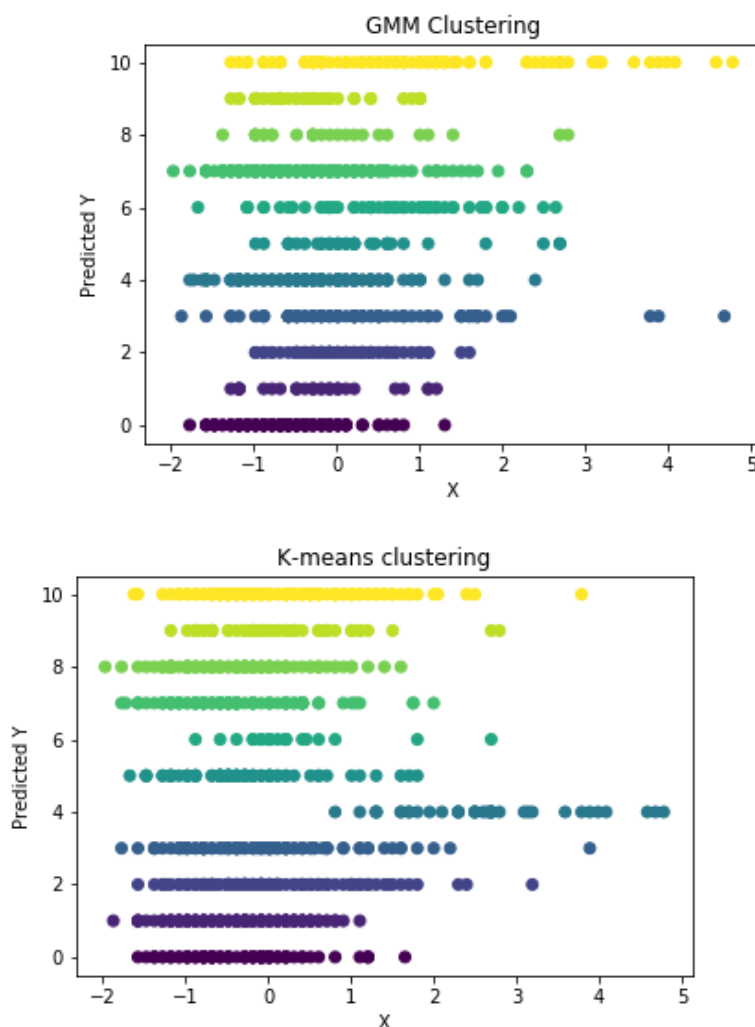


הגרפים מציגים את היכולת לחזות את איכות היין (PREDICTED Y) אל מול הערך האמיתי של איכות היין (TRUE Y).

בוצע מבחן TSET-T לבלתי תלויים על סמך ריבוע השגיאות של ערך החיזוי של כל דגימה אל מול הערך האמיתי של איכות היין (PREDICTED Y - TRUE Y). ערך ה-P VALUE

של מבחן זה יצא $1.3067221214172147 \times 10^{-5}$ כשמוצג ריבועי השגיאות של הרגרסיה הליניארית יצא 0.5690247717229263 ושל המודל RANDOM FOREST יצא 0.382377551020408 ולכן מבין שתי השיטות הללו מודל הרגרסיה הטוב יותר הוא מודל ה-RANDOM FOREST. בו סך ממוצעי ריבועי השגיאות קטן באופן מובהק מזה של מודל הרגרסיה הליניארית.

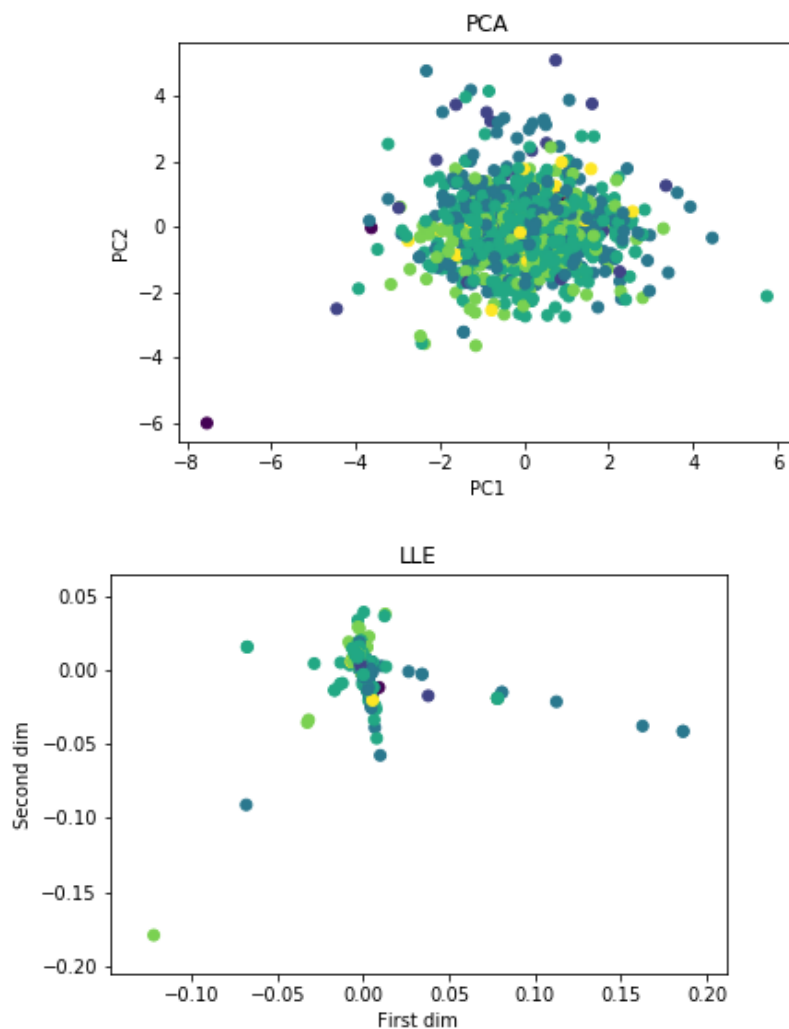
2. בחלק זה הדרישה היתה לבצע אשכול של הנתונים על סמך שני מודלים שונים. המודל הראשון שנבחר הוא K-MEANS והמודל השני הוא GMM (GAUSSIAN MIXTURE MODEL). בחלק זה טיב האשכול הוערך על ידי SCORE SILHOUETTE. ההרצה של כל אחת משיטות ה-CLUSTERING הנ"ל בוצעה 30 פעמים ולכל אחת חולץ מדד זה. הממוצע K-MEAN היה 0.11178440152783391 וממוצע GMM היה 0.0016525120223483622 . כבר מהשוואה זו ניתן לראות כי ששיטת K-MEAN עדיפה לחלוקה לאשכולות מאחר והיא בעלת הערך הגבוה יותר והקרובה ל-1 (הערך הגבוה ביותר למדד). מעבר לכך בוצע מבחן T-TEST גם כאן וערך ה-P-VALUE יצא $1.260634256993649 \times 10^{-56}$, גם כאן ההבדלים יצאו מובהקים.



הגרף העליון מראה את שיטת החלוקה לאשכולות על סמך GMM והגרף התחתון לפי שיטת K-MEANS. הגרפים מציגים את החלוקה לאשכולות על סמך הפיצ'ר הראשון כאשר ציר ה-Y מסמן את האשכול של כל דגימה ונע בטווח 10^{-0} .

3. בחלק זה הדרישה היתה ליישם שתי שיטות להורדת מימדים. השיטה הראשונה שנבחרה היא שיטת LLE (LOCALLY LINEAR EMBEDDING), והשיטה השניה היא PCA (PRINCIPAL COMPONENT ANALYSIS). המדד להשוואה בין השיטות הוא RECONSTRUCTION MINIMUM ERROR. מדד זה ב-LLE הוא $9.209806608190625 \times 10^{-19}$ וב-PCA הוא 0.46563534972005655 ולכן מבין שניהם ההטלה של LLE היא קורלטיבית יותר לאיכות היי. המשמעות היא שהורדת המימדים שמרה על עיקר המידע ונוכל להסתמך על

המודל בכדי לחזות בהצלחה גבוהה יותר על סמך LLE.



שני הגרפים הנ"ל מראים את האשכולות לפי שיטת PCA ושיטת LLE.