

(3) למה נשתמש ב-feature selection?
כי לפעמים כמות המופעים שלנו יכולה להיות
לחיצה למיליון וחמושים מסווגים לא יכולים
להתמודד עם כמות כל כך גדולה של features.
אנו צריכים בדרך כלל להוריד כמות איתן
של המוצר, כדי להוריד כמות ריבוי, כדי
להפחית מידע בלתי רלוונטי. ערכים ממוצעים לא קשורים
ויש להם חשיבות מועטה מוצר ממוצע שלא יהיה
מקצועי לא נכונה.

למצוא בהרצאה שלוש סוגים של feature selection:
Mutual Information, Chi-square, Frequency.

בשיטה Chi-square היא שיטה המבוססת
על נדונים סטטיסטיים, ויכולה להבחין
מונחים בנות שיתופיים לפעולה סיווג.
מחשבים על ידי הסתברות של הופעת מילה
מסוימת ביחד עם class מסווג.

בשיטה Mutual Information היא שיטה המבוססת
על מידע יותר מנצח ופרסומי מנצח. אבל
יכולה להשתמש במאפיינים לא שיתופיים נצחיים.
מקבלים מוצר רב ובנות אם מונח נמצא במסמך
ומסתמך שייך ל-class מסווג.
ושיטה ה' פשוטה היא Frequency. מתחשבים
בה במאפיינים ה' נכונים. אלא שה' מסווגים אבל
זה עובד כי עוקבים עם נדונים ה' שיתופיים ובלתי
המקרים זו שיטה ה' יעילה.

אנו נשתמש ב-feature selection ב-MB מודל ווריאנטי כי
אחרת נסבוק מידע בלתי רלוונטי והרבה סבירות.