

text mining

表是一个矩阵 (matrix)

d_1						
d_2						

→ term vector

如果字很多, 維度會太大
可透過下二種方法來處理

Document vector = profile

主成分分析法 ⇒ 找特徵值
奇異值分解法

① words become features

資料 → 以什麼型式的文件讀入 → word segmentation →

去除停用字、保留詞幹 (stemming)

stop word (黃) walked, walking 雖形態不同 but same meaning

爬資料 ex: 八天