

אחזור מידע

תרגום חופשי של הספר Information Retrieval

של C. J. van RIJSBERGEN

פרק 6 - אחזור הסתברותי

תרגום: עירן גישס, נועם עמישר וזאב גטנר

עריכה: עפר דרורי

הקדמה

עד כה בספר זה עשינו שימוש מועט בתורת ההסתברות בבניית מודל של איזו שהיא תת מערכת ב-IR. הסיבה לכך היא פשוט העובדה שעיקר העבודה ב-IR היא בלתי הסתברותית ורק לאחרונה נעשתה פריצת דרך משמעותית כלשהי עם שיטות הסתברותיות. ההיסטוריה של השימוש בשיטות הסתברותיות מגיעה עד תחילת שנות השישים, אך מאיזו שהיא סיבה הרעיונות הראשוניים לא השתרשו. בפרק זה אתאר שעטות אחזור, כלומר חוקי חיפוש ועצירה המבוססים על שיקולים הסתברותיים. בפרק 2 עסקתי בקטלוג אוטומטי מבוסס על מודל הסתברותי של התפלגות אסימוני מלים במסמך (טקסט): כאן אעסוק בהתפלגות של מונחי קטלוג בקובץ מסמכים היוצרים אוסף של קובץ. אסתמך בעיקר על ההנחה המוכרת שהתפלגות של מונחי קטלוג באוסף או בתת קבוצה שלו יגלה לנו משהו אודות החשיבות הצפויה של כל מסמך נתון.

יתכן שכדאי להזהיר את הקורא שחלק מהחומר בפרק זה הוא מתמטי. אולם אני מאמין שהמסגרת של האחזור הנדון בפרק זה היא גם אלגנטית וגם בפוטנציאל חזקה ביותר. חרף העובדה שהעבודה על כך הנה חדשנית למדי ולפיכך אחדים עשויים לחשוב שעליה לעמוד במבחן הזמן, לדעתי היא מייצגת את פריצת הדרך החשובה ביותר ב-IR בשנים האחרונות. לפיכך, ללא בוושה, אעשה פרק זה תיאורטי היות והתיאוריה צריכה להיות מובנת היטב אם רוצים לעשות התקדמות כלשהי. ישנם מספר דרכים שונות להצגת התיאוריה הבסיסית. אני בחרתי להציג אותה בדרך כזו שהקשרים עם תחומים נוספים, כגון זיהוי תבנית, יעשו בקלות. יהיו לי דברים נוספים לומר אודות ניסוחים אחרים בהערות הביבליוגרפיות בסוף הפרק.

הכלי המתמטי הבסיסי בפרק זה הוא "משפט בייס". רוב המשוואות נובעות ישירות ממנו. למרות שהמתמטיקה הבסיסית עשויה בתחילה להראות מסובכת מעט התרגום הוא די פשוט. אז תנו לי לנסות ומייד לתת מעט תרגום למה שעתידי לבוא.

זכרו שהכלי הבסיסי שיש בידינו ע"מ להפריד את המסמכים הרלוונטיים מהבלתי רלוונטיים הוא פונקצית התאמה, הן אם אנו בסביבה מקובצת והן בסביבה בלתי מקובצת. הטעם בבחירת פונקצית התאמה מעולם לא הובהר, למעשה לרוב היא מבוססת על טענה אינטואיטיבית בהצלבה עם "OCCAM'S RAZOR". בפרק זה אנסה להשתמש בתיאוריה הסתברותית לשם קביעה כיצד פונקצית התאמה צריכה להראות וכיצד לעשות בה שימוש. הטענות הם בעיקרם תיאורטיים אך לדעתי די נחרצים. הספק היחיד שנוותר הוא על הקבילות של ההנחות שאנסה להציג בהסברי. המידע בו נעשה שימוש ע"מ לקבוע פונקצית התאמה נובע מהידע על התפלגות תנאי הקטלוג באוסף של מספר תת קבוצות שלה. אם היא מבוססת על מספר תת קבוצות של מסמכים אז תת קבוצה זו יכולה להיות מבוססת במספר דרכים: דגימה, קיבוץ או אחזור ניסיוני. המידע שנאסף משמש לקביעת ערכים של פרמטרים מסוימים הקשורים עם פונקצית ההתאמה. ברור שאם המידע מכיל מידע רלוונטי אז התהליך של הגדרת פונקצית ההתאמה יכול להיות מוזן ע"י מנגנון מישור דומה לזה הנובע מ-ROCCHIO שתואר בפרק הקודם. בדרך זו הפרמטרים של פונקצית ההתאמה יכולים "להילמד". אנו נתרכז בפונקצית התאמה הנובעת ממידע רלוונטי. יונח כי המסמכים מתוארים ע"י מאפייני מצב בינריים כלומר המצאות או היעדרות של תנאי קטלוג. זו אינה הגבלה על התיאוריה, בעקרון ההשלכה למאפיינים שרירותיים יכולה להימצא למרות שאין זה ברור שכדאי לעשות כן.

הערכה או חישוב של רלוונטיות

כשאנו מחפשים אוסף מסמכים או מנסים לאחזר מסמכים רלוונטיים מבלי לאחזר בלתי רלוונטיים. היות ואין עמנו נביא שיאמר לנו ללא טעות אילו מסמכים רלוונטיים ואילו בלתי רלוונטיים עלינו לעשות שימוש במידע שאינו מושלם ע"מ לנחש על כל מסמך שהוא אם הוא רלוונטי או בלתי רלוונטי. מבלי להיכנס לפרדוקס הפילוסופי הקשור לרלוונטיות אניח שאנו מסוגלים לנחש בלבד אודות רלוונטיות באמצעות מידע מתומצת על המסמך והיחסים שלו עם מסמכים אחרים. זו אינה הנחה מובנת במיוחד אם מאמינים שהדרך היחידה בה נקבעת רלוונטיות היא ע"י קריאת הטקסט כולו ע"י המשתמש. לכן דרך הגיונית לחישוב הניחוש שלנו היא לנסות ולהעריך לכל מסמך נתון את ההסתברות הרלוונטית שלו.

PQ(relevance/document)

הבה נניח (בעקבות רוברטסון [7]):

(1) הרלוונטיות של מסמך לבקשה היא בלתי תלויה במסמכים אחרים באוסף. עם הנחה זו אנו יכולים להציג עיקרון, במונחי הסתברות של רלוונטיות, שניתן להשתמש במידע הסתברותי באופן אופטימלי באחזור. רוברטסון מייחס עיקרון זה ל-ו.ס. קופר למרות ש-מרון כבר הדגיש את האופטימליות שלו ב-1964.

עקרון הדירוג ההסתברותי – אם תגובת מערכת שייכות אחזור לכל בקשה היא דרוג המסמכים באוסף בסדר יורד של רלוונטיות של הסתברות למשתמש שהגיש את הבקשה, בה הסתברות מוערכת בדיוק רב ככל הניתן על בסיס העובדה האם המידע היה זמין למערכת למטרותנו, היעילות הכוללת של המערכת למשתמש שלה תהיה הטובה ביותר הניתנת להשגה על בסיס נתונים אלו. כמוכן עקרון זה מעלה שאלות רבות באשר לתקפות ההנחות. לדוגמה "היפותזת הצבר" שמשמכים הקשורים בינם נוטים להיות רלוונטיים לאותן בקשות, במפורש מניחה את ההפך מהנחה (1). GOFFMAN, גם כן, בעבודתו טרח רבות ע"מ לבסס הנחת תלות מפורשת. אני מצטט "אם מסמך X הוערך כרלוונטי לשאלתה S הרלוונטיות של מסמכים אחרים בקובץ X עשויה להיות מושפעת היות והערך של המידע הנמצא במסמכים אלו עשוי לגדול או להצטמצם כתוצאה מהמידע הנמצא במסמך X". כמו כן קיימת השאלה על הדרך בה נמדדת היעילות הכוללת. רוברטסון, במאמרו מציג את עקרון הדירוג ההסתברותי למדידת יעילות במונחים של אחזור ונשורת. העיקרון גם נובע מהתיאוריה המתוארת בפרק זה. אך זהו אינו המקום להעלות את שאלות המחקר הללו, אולם, אני חושב כי זה הגיוני לאמץ את העיקרון כאחד שעליו ניתן לבנות מודל אחזור הסתברותי. מילת אזהרה, עקרון הדירוג ההסתברותי יראה נכון לשאלתה אחת בלבד. אין זה אומר שהביצועים על טווח של שאלות יהיה אופטימלי, לבסס תוצאה מסוג זה יש להיות ספציפי על הדרך למצע את הביצועים על פני שאלות.

עקרון הדירוג ההסתברותי מניח שביכולתנו לנחש F (רלוונטיות/מסמך), לא זו בלבד, הוא מניח שביכולתנו לעשות זאת בדיוק. זוהי הנחה בעייתית מאד והיא תעסיק אותנו בהמשך. הבעיה היא שאיננו יודעים אילו מסמכים רלוונטיים ואין אנו יודעים כמה ישנם כך שאין לנו דרך לחשב את F (רלוונטיות/מסמך). אך ביכולתנו באמצעות אחזור ניסיון, לנחש את F (רלוונטיות/מסמך) ובתקווה לשפר את הניחוש ע"י איטרציה. ע"מ לפשט את העניינים בדיון העוקב אניח שהסטטיסטיקות המתייחסות למסמכים רלוונטיים והבלתי רלוונטיים נגישות ואשתמש בהן ע"מ לבנות את המשוואות. אולם בכל זמן על הקורא להיות ער לעובדה שבכל מצב המידע היחסי צריך להיות מנוחש (או מוערך). נחזור לבעיה הנוכחית שהנה לחשב או להעריך את F (רלוונטיות/מסמך). לזאת נשתמש ב"משפט בייס" המייחס את ההסתברות הרלוונטית שאחרי ההסתברות הרלוונטית שלפני וסבירות הרלוונטיות לאחר בחינת מסמך. לפני שנקפוץ לביטוי פורמלי של הנושא עלי להציג מספר סמלים שיפשטו את הדברים בהמשך.

מודל הסתברותי בסיסי

היות ואנו מניחים שכל מסמך מתואר ע"י תנאי קטלוג בהמצאות/היעדרות כל מסמך יכול להיות מיוצג ע"י וקטור בינארי

$$X = (x_1, x_2, \dots, x_n)$$

כש- $X_1=0$ או 1 מייצג העדות או המצאות של התנאי. אנו גם מניחים שישנם שני אירועים נבדלים

המסמך רלוונטי – W_1
 המסמך אינו רלוונטי – W_2

אז במונחים של סמלי אלו, מה שברצוננו לחשב לכל מסמך זה את $P(W_1/X)$ ואולי את $P(W_2/X)$ כדי שנוכל להחליט איזה רלוונטי ואיזה אינו רלוונטי. זהו שינוי קל במטרה, מפשוט לייצר דרוג ברצוננו שהתיאוריה תאמר לנו כיצד לצמצם את הדרוג. לכן אנו מנסחים את השאלה כשאלת בחירה. כמוכן שאין ביכולתנו להעריך את $P(W_1/X)$ ישירות לכן אנו חייבים למצוא דרך להעריך במושג של כמויות שאיננו יודעים עליהן. משפט בייס מלמד אותנו על התפלגויות בדידות

$$P(w_i/x) = \frac{P(x/w_i) P(w_i)}{P(x)} \quad i=1,2$$

כאן $P(w_i)$ הוא עיקר הסתברות לרלוונטיות $(i=1)$ או בלתי רלוונטיות, $P(x/w_i)$ פרופורציונלי למה שידוע כסבירות לרלוונטיות או בלתי רלוונטיות בהינתן X . במקרה הרציף זו תהיה פונקצית צפיפות ונרשום $p(x/w_i)$ לבסוף

$$P(x) = \sum_{w_i} P(x/w_i) P(w_i)$$

שזו ההסתברות לצפות ב- X על בסיס רנדומלי בהינתן שעשוי להיות רלוונטי או בלתי רלוונטי. שוב זה ירשם בפונקצית צפיפות $p(X)$ במקרה הרציף. למרות ש- $P(x)$ (או $p(x)$) בעיקר יופיע כגורם נרמול (כלומר הבטחה ש- $P(w_1/x) + P(w_2/x) = 1$) שזו הפונקציה עליה אנו יודעים את הדברים הרבים ביותר, היא אינה דורשת מידע של רלוונטיות ע"מ להיות מוגדרת. לפני שאדון כיצד מעריכים את האגף הימני של משפט בייס אראה כיצד ההחלטה בעד או נגד רלוונטיות נעשית.

כלל ההחלטה בו אנו משתמשים ידוע היטב כ"חוק ההחלטה של בייס". זהו
 $[P (w_1/x) > P(w_2/x) \rightarrow x \text{ is relevant, } x \text{ is non-relevant}] * D1$

הביטוי $D1$ הוא קיצור של: השווה את $P (w_1/x)$ עם $P (w_2/x)$. אם הראשון גדול מהשני אזי החלט X רלוונטי, אחרת החלט X בלתי רלוונטי. המקרה בו $P(w_1/x) = P(w_2/x)$ מטופל שרירותית ע"י קביעת בלתי רלוונטיות. הבסיס לחוק $D1$ הוא שהוא מצמצם את ההסתברות הממוצעת לטעות, הטעות שבקביעת מסמך רלוונטי כבלתי רלוונטי או להפך. ע"מ לראות זאת שים לב שלכל X ההסתברות לטעות היא

$$P_{\text{error}}(x) = \begin{cases} P(w_2/x) & \text{if we decide } w_1 \\ P(w_1/x) & \text{if we decide } w_2 \end{cases}$$

המשמעות של $[E \rightarrow p,q]$ היא שאם E נכון החלט P אחרת החלט Q .

במילים אחרות, מרגע שהחלטנו על דרך אחת (למשל רלוונטיות) אזי ההסתברות לטעות נתונה ע"י ההסתברות שהדרך ההפוכה היא הנכונה (כלומר בלתי רלוונטיות). ע"מ לעשות טעות זו קטנה ככל האפשר לכל X נתון עלינו לבחור ב- w_i לו $P (w_1/x)$ הוא הגדול ביותר ובהתאם עבורו ההסתברות לטעות היא הקטנה ביותר. ע"מ לצמצם את ההסתברות הממוצעת לטעות עלינו למזער

$$P(x) = \sum_i P(x|w_i) P(w_i)$$

סכום זה ימוזער ע"י מזעור $P(\text{error}/x)$ קטן ככל הניתן לכל X היות ו- $P(\text{error}/x)$ ו- $P(x)$ תמיד חיוביים. מטרה זו מושגת ע"י חוק ההחלטה $D1$ שכעת מקבל הצדקה.

כמוכן שטעות ממוצעת אינה הכמות היחידה שכדאי למזער. אם נקשר לכל סוג של טעות מחיר נוכל להסיר כלל החלטה שימזער ת הסיכון הכללי. הסיכון הכולל הוא ממוצע של תנאי הסיכון $R(w_i/x)$ שבפני עצמו מוגדר במושגים של פונקציה מחיר l_{ij} . כלומר l_{ij} הוא ההפסד הנגרם מההחלטה על w_i כש- w_j הוא הנכון. ההפסד הצפוי כשמחליטים w_i נקרא סיכון משתנה ומוגדר

$$R(w_i/x) = \sum_j l_{ij} P(w_j/x) \quad i = 1, 2$$

הסיכון הכולל הוא סכום באותה דרך שההסתברות הממוצעת לטעות היתה $R(w_i/x)$ כעת ממלא תפקיד של $P(w_i/x)$. הסיכון הכללי ממוזער ע"י

$$R(w_1/x) < R(w_2/x) \rightarrow x \text{ is relevant, } x \text{ is non-relevant] } D2]$$

$D1$ ו- $D2$ יכולים להראות חליפיים תחת תנאים מסוימים. ראשית נשכתב $D1$ ע"י שימוש במשפט בייס בדרך בה יעשה בו שימוש חליפי

$$[P(x/w_1) P(w_1) > P(x/w_2) P(w_2) \rightarrow x \text{ is relevant, } x \text{ is non-relevant] } D3$$

שים לב ש- $P(x)$ נעלם מהמשוואה היות שאין הוא משנה את התוצאה של ההחלטה. כעת ע"י שימוש בהגדרה $R(w_i/x)$ נקל להראות ש-

$$R(w_1/x) < R(w_2/x) \text{] [equivalence] [(121 - 111) P(x/w_1) P(w_1) > (112 - 122) P(x/w_2) P(w_2)]$$

כאשר פונקציה הפסד מיוחדת נבחרת

$$l_{ij} = \begin{cases} 0 & i=j \\ 1 & i \neq j \end{cases}$$

הרומז שלבחירה נכונה לא משתייך הפסד (הגיוני למדי) והפסד יחידה לטעות (פחות הגיוני) כאשר

$$[R(w_1/x) < R(w_2/x) \text{] [equivalence] P(x/w_1) P(w_1) > P(x/w_2) P(w_2)]$$

המציג את שווה הערך ל- $D2$ ו- $D3$ ולפיכך ל- $D1$ ו- $D2$ תחת פונקציה הפסד בינארית.

זה משלים את פיתוח חוק ההחלטה שישמש ע"מ להחליט על רלוונטיות או אי רלוונטיות, או להציג זאת באופן שונה, לשלוף או לא לשלוף. עד כה לא הושמו כל הגבלות על $P(x/w_1)$, לכן חוק ההחלטה הוא כללי למדי. הצגתי את הבעיה כבחירה בין שתי מחלקות ולפיכך התעלמתי מהבעיה של הדירוג. סיבה אחת לכך היא שכך הניתוח פשוט יותר והאחרת שברצוני שהניתוח יישאר למשתמש. בתוך המודל עד כה ניתן עדיין לדרג אך ערך החיתוך יהיה מפורש במונחים של הסתברויות קודמות ופונקציות מחיר. האופטימליות של עקרון הדרוג ההסתברותי נובע מהאופטימליות של חוק ההחלטה בכל נקודת חיתוך. עתה אבהיר את הצורה המדויקת של פונקציה ההסתברות בחוק ההחלטה.

פונקציית אחזור

החלק הקודם היה אבסטרקטי למדי והשאיר את הקשרים בין ההסתברויות השונות וה- IR פתוחים. למרות שזה הגיוני לרצות לחשב (רלוונטיות/מסמך) P אין זה כלל ברור כיצד לעשות זאת או האם ההמרה באמצעות משפט בייס היא הדרך הטובה ביותר. אף על פי כן נמשיך בהנחה ש- $P(x/w_i)$ היא הפונקציה הנכונה להנחה. פונקציה זו היא כמובן פונקציית הסתברות מצרפית והאינטראקציה בין רכיבי x עשויה להיות מורכבת באופן שרירותי. ע"מ להסיק חוק החלטה שניתן לעבוד עמן יש להניח הנחה מקלה על $P(x/w_i)$. הדרך המתמטית הקונבנציונלית להפשטת $P(x/w_i)$ היא להניח כי מרכיבי המשתנים x_i של X הם בלתי תלויים באופן סטוכסטי. באופן טכני זה מוביל להניח כי

$$P(x/w_i) = P(x_1/w_i) P(x_2/w_i) \dots P(x_n/w_i) \quad A1$$

מאוחר יותר אראה כיצד ניתן יהיה לפשט הנחה קשיחה זו. בשלב זה נתעלם מהעובדה שהנחת התניות בלתי תלויות ב- W_1 ו- W_2 באופן נפרד הינה בעלת משמעות על התניות התלות על W_1 ו- W_2 .

הבה עתה ניקח את הצורה המופשטת של $P(x/w_i)$ ונראה כיצד יראה חוק ההחלטה.

ראשית נגדיר מספר משתנים

$$p_i = \text{Prob}(x_i = 1/w_1)$$

$$q_i = \text{Prob}(x_i = 1/w_2)$$

במילים $p_i(q_i)$ היא ההסתברות שאם מסמך הוא רלוונטי (בלתי רלוונטי) אזי תנאי קטלוג ה- i יהי נוכח. ההסתברות המתאימות להעדרות מחושבת ע"י חיסור מ-1, כלומר $p_i = \text{Prob}(x_i = 1/w_1)$ ו- $q_i = \text{Prob}(x_i = 1/w_2)$. פונקציית הסבירות שנכנסת ל- D_3 תראה כעת

$$P(x/w) = \prod_{i=1}^n p_i^{x_i} (1-p_i)^{1-x_i}$$

$$P(x/w) = \prod_{i=1}^n q_i^{x_i} (1-q_i)^{1-x_i}$$

ע"מ להבין כיצד ביטויים אלו עובדים, על הקורא לבדוק ש-

$P((0,1,1,0,0,1)/w_1) = (1-p_1)p_2p_3(1-p_4)(1-p_5)p_6$. החלפת $P(x/w_i)$ ב- D_3 ורישום, חוק ההחלטה יומר לפונקציית הפרדה לינארית

$$c_i = \ln \frac{p_i(1-q_i)}{q_i(1-p_i)}$$

היכן שהקבועים a_i ו- b_i הם ברורים

$$\begin{aligned} \ln P(x/w) &= \sum_{i=1}^n (a_i x_i + b_i (1-x_i)) + c \\ &= \sum_{i=1}^n c_i x_i + c \end{aligned}$$

החשיבות שבכתיבה זו, מלבד פשטותה, היא שכל מסמך X ע"מ לחשב $g(x)$ עלינו פשוט לחבר את המקדמים c_i להם $x_i = 1$. C_i לרוב נקרא משקל. רוסרטסון וספרק גיונס מכנים את c_i כמשקל רלוונטיות וסלטון מכנה את $\exp(c_i)$ כקבוע או משל. מכאן השם פונקציית שקילה ל- $g(x)$.

הקבוע C שהנחנו כי הוא זהה לכל המסמכים X ישתנה משאילתא לשאילתא, אך ניתן לזהותו כחיתוך המיושם בפונקצית ההחזרה. החלק היחיד שניתן לשינוי בהתייחס לשאילתא נכונה היא פונקצית המחיר ושינוי זה הוא זה שיאפשר לנו לאחזר יותר או פחות מסמכים. ע"מ לראות זאת הבה נניח ש- $I_{11} = I_{22} = 0$ ושיש לנו בחירה בקביעת היחס I_{21}/I_{11} ע"י בחירת ערך לחשיבות היחסית שאנו מיחסים לפספוס מסמך רלוונטי בהשוואה לשליפת מסמל בלתי רלוונטי. בדרך זו ביכולתנו לייצר דרוג, כל דרגה מתייחסת ליחס שונה של I_{21}/I_{12} .

הבה נעבור עתה לחלק האחר של $g(x)$, כלומר c_i וננסה לפרט זאת במונחים של טבלת ה-'contingency' הקונבנציונלית.

	Relevant	Non-relevant	
$r_i=1$	r	$n-r$	n
$r_i=0$	$R-r$	$N-n-R+r$	$N-r$
	R	$N-n$	N

תהיה טבלה כזו לכל תנאי קטלוג. הראתי זאת לתנאי הקטלוג למרות שלא נעשה שימוש במקרא בתאים. אם יש לנו מידע מלא על המסמכים הרלוונטיים ובלתי רלוונטיים באוסף, אזי אנו נוכל להעריך את ע"י $p_i = r/R$ ואת q_i ע"י $(n-r)/(N-R)$. לכן $g(x)$ יכול להירשם

$$g(x) = \sum_{i=1}^n r_i \log \frac{r(R-r)}{(n-r)(N-n-R+r)} + C$$

זו למעשה פונקציה F_4 בה עשו שימוש רוברטסון וספרק ג'ונס בניסויי "המבט לאחור" שערכו. לנוחות עתידית הבה נקבע

$$K_i(N,n,R) = \log \frac{r(R-r)}{(n-r)(N-n-R+r)}$$

יש מספר דרכים להכיל את K_i . הפירוש המעניין ביותר ל- K_i הוא לומר שהוא מודד את המידה בה התנאי ה- i יכול להבדיל בין המסמכים הרלוונטיים והבלתי רלוונטיים.

באופן טיפוסי, המשקל, $K_i(N,r,n,R)$ מוערך מטבלת ה- contingency בה N אינו המספר הכולל של המסמכים במערכת אלא הוא תת מערכת כלשהי שנבחרה במיוחד ע"מ לאפשר להעריך את K_i . מאוחר יותר אעשה שימוש בפרוש דלעיל של K_i ע"מ להציג פונקציה נוספת הדומה ל- K_i בכדי למדוד את כוח ההפרדה של תנאי הקטלוג.

מונחי הקטלוג אינם בלתי תלויים

למרות שנוח מתמטית להניח שתנאי הקטלוג הם בלתי תלויים אין זה בהכרח ריאלי. ההתנגדות לחוסר תלות אינה חדשה. ב-1964 ה.ה. ויליאמס הביע זאת כך "ההנחה של חוסר תלות של מילים במסמך נעשית בד"כ כחלק מנוחות מתמטית. בלי ההנחה רבים מהיחסים המתמטיים העוקבים לא יוכלו להיות מבוטאים. עמה רבות מההנחות צריכות להתקבל הזהירות יתר. רק בגלל שהמתמטיקה נהית די בלתי ניתנת למעקב כאשר מניחים תלות אנשים נוטים להניח חוסר תלות. אום "תלות היא נורמה יותר מאשר ההפך". בציטוט מהתאורטיקן ההסתברות הידוע דה-פינטי. לכן התהליך הנכון הוא להניח תלות ולאפשר לניתוח לפשט לחוסר תלות רק אם אכן יהיה כך. כאשר מדברים על תלות אנו מתכוונים לתלות סטוכסטית. למידע IR תלות סטוכסטית פשוט נמדדת ע"י פונקצית קורלציה או ע"י דרך דומה אחרת.

הנחת התלות יכולה להיות הכרחית כאשר אנו מנסים להניח (רלוונטיות/מסמך) P במונחים של $P(x/w_i)$, היות שהדיוק בה מוערכת ההסתברות האחרונה תשפיע ללא ספק על ביצועי האחזור. לכן משימתנו המיידית היא לעשות שימוש בתלות (קורלציה) בין תנאי האינדקס בכדי לשפר את הערכתנו של $P(x/w_i)$ עליה נשען חוק ההחלטה שלנו.

באופן כללי התלות יכולה להיות מורכבת בצורה הבאה

$$P(x) = P(x_1)P(x_2/x_1)P(x_3/x_1, x_2) \dots P(x_n/x_1, x_2, \dots, x_{n-1})$$

לכן, בכדי להתייחס לכל המשתנים התלויים נצטרך לבטא כל משתנה באמצעות קבוצה הולכת וגדלה של המשתנים הנוספים. בעקרון זה אפשרי אך זה כנראה יהיה מאוד לא יעיל חישובית, ובמקרים מסויימים אף בלתי אפשרי כאשר אין מספיק מידע לחשב מקרים בסדר גבוה. במקום זאת נאמץ שיטת הערכה ל $P(x)$ שתיתן את מידע התלות המהותי. בצורה אינטואטיבית ניתן להסתכל בזאת כמישהו במסתכל על הפריסה הנ"ל ובוחר את המשתנה אם התלות הרבה ביותר. במילים אחרות אנו מחפשים תוצר מהצורה הבאה

$$P_j(x) = \prod_{i=1}^n P(x_i/x_{1..j}) \quad 0 \leq j \leq n \quad A_2$$

כאשר (m_1, m_2, \dots, m_n) מהווה פרמוטציה של השלמים 1 עד n ו $j(.)$ היא פונקציה הממפה את i לשלמים קטנים מ i ו $P(x_i/x_{m_0})$ שווה ל $P(x_i)$. להלן דוגמה לוקטור עם שישה מרכיבים $x = (x_1, \dots, x_6)$

$$P_t(x) = P(x_1)P(x_2/x_1)P(x_3/x_2)P(x_4/x_2)P(x_5/x_2)P(x_6/x_5)$$

שימו לב כמה דומה הנחה A_2 להנחה הבלתי תלויה A_1 , ההבדל היחיד הוא בכך שב A_2 לכל גורם קיים משתנה מתנה המזוהה עימו. בדוגמת הפרמוטציה (m_1, m_2, \dots, m_6) הם $(1, 2, \dots, 6)$ שזהו הסדר הטבעי, כמובן שהסיבה לכך שכתבנו את הפריסה כמו ב A_2 היא להראות שהפרמוטציה $(1, 2, \dots, 6)$ תיתן קירוב טוב. ברגע שניצור פרמוטציה זו ניתן לסווג את משתניה כבעלי סדר טבעי. הפרמוטציה ופונקציית $j(.)$ המגדירות ביחד עץ תלות ואת $P_t(x)$ המקביל נקראות התפלגות עץ התלות (מסדר ראשון). העץ המותאם לדוגמה עם ששת המשתנים מופיע באיור 6.1. העץ מראה משתנה המופיע בכל צד של פעימת ההתניה ב $P(./.)$. למרות שבחרנו לכתוב את הפונקציה $P_t(x)$ בצורה ש x_i הינו משתנה בלתי מותנה, ומכאן שורש העץ, למעשה כל עלה בעץ יכול לשמש כשורש העץ בתנאי שההתניה נעשית באופן רציף כלפי שורש העץ החדש.

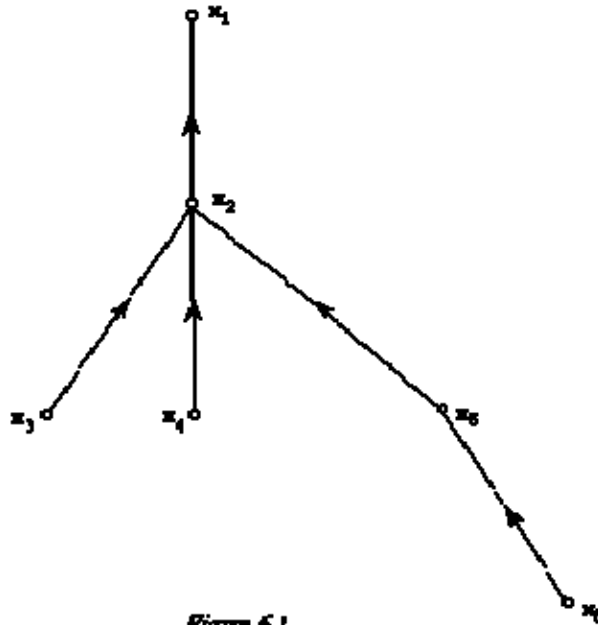


Figure 6.1.

$$P(x_1/x_2) = P(x_3/x_2)P(x_4/x_2)P(x_5/x_2)P(x_6/x_5)$$

אם נפעיל את המשוואה על הקשר בין השורש X_1 לבין הצאצא X_2 בדוגמה השורש יעבור ל X_2 והפריסה תשתנה להיות

$$P(x_1, x_2, \dots, x_6) = P(x_2)P(x_1/x_2)P(x_3/x_2)P(x_4/x_2)P(x_5/x_2)P(x_6/x_5)$$

כמובן, בשביל לעמוד בחוק של מתן תוויות מחדש נחליף בין השמות '1' ו '2'. כל הפריסות המוחלפות בצורה זו הינם קירוב טוב ל $P(x)$. ולכן העץ מייצג את קבוצת הפריסות המקבילות. ברור שקיימים מספר רב של עצי תלות, בעיית ההערכה שעלינו לפתור היא למצוא את העץ הטוב ביותר; כלומר למצוא את הפרמוטציה הטובה ביותר ואת המיפוי $j(\cdot)$ הטוב ביותר. החל מעכשיו נניח שנתינת התוויות נעשית בצורה $x_{mi} = x_i$.

בחירת עץ התלות הטוב ביותר

הבעיה העומדת בפנינו היא למצוא את פונקציית ההסתברות בצורת $P_t(x)$ על קבוצת מסמכים שהיא הקירוב הטוב ביותר לפונקציית ההסתברות המשותפת $P(x)$, וכמובן קירוב טוב יותר מזה הנובע מהנחה $A1^*$. הקבוצה עליה הקירוב מוגדר יכולה להיות שרירותית ומוגדרת כאוסף המסמכים הרלוונטיים או המסמכים הלא רלוונטיים. לעת עתה לא נגדיר את הקבוצה. כשנבנה חוק החלטה הדומה ל $D4$ נאלץ למצוא קירוב ל $P(x/w1)$ ו $P(x/w2)$.

טיב הקירוב נמדד ע"י פונקציה המוגדרת היטב. אם $P(x)$ ו $Pa(x)$ הם שני הסתברויות דיסקרטיות אזי

$$I(P, Pa) = \sum_i P_i \log \frac{P_i}{P_{i,j}}$$

הינו מדד לטיב הקירוב $Pa(x)$. במונחים של פונקציה זו נרצה למצוא התפלגות של עצים תלויים $P_t(x)$ כך ש $I(P, P_t)$ מהווה מינימום, במילים אחרות למצוא עץ תלות מכלל העצים התלויים שיגרום ל $I(P, P_t)$ להיות קטן ככל האפשר. במקרה וסטיית האינדקסים i ו j מאי תלות נמדדת באמצעות (EMIM) expected mutual information measure (ראה פרק 3)

$$I(x_i, x_j) = \sum_{i,j} P_{i,j} \log \frac{P_{i,j}}{P_i P_j}$$

אזי הקירוב הכי טוב $P_t(x)$, במונח שצמצום $I(P, P_t)$, ניתן באמצעות עץ הפריסה המקסימלי (MST) (ראה פרק 3) על המשתנים x_1 עד x_n . עץ הפריסה נגזר מהגרף שצמתיו הם אינדקס 1 עד n וקצותיו משוקללים באמצעות $I(x_i, x_j)$. ה MST הוא למעשה העץ ששקלול כל צמתיו

$$\sum_{i,j} I(x_i, x_j)$$

מהווה מקסימום. זוהי הגדרה מאוד כללית לעץ, אך הגדרה טובה יותר תהיה טכנית. הוכחה מדויקת של האופטימיזציה ניתן למצוא ב[13] Chow and Liu. אנו מעוניינים ביישום של מבנה העץ.

דרך אחת להסתכל על ה MST היא בכך שהוא כולל בתוכו את מרבית התלויות המשמעותיות בין המשתנים ומותנה בכך שסכום המשתנים צריך להיות המקסימום. לדוגמה באיור 6.1 הקשר בין המשתנים (צמתים x_1 עד x_6) הוכנס רק בגלל שהסכום

$$I(x_1, x_2) + I(x_2, x_3) + I(x_2, x_4) + I(x_2, x_5) + I(x_5/x_6)$$

הוא מקסימום. כל סכום יהיה קטן או שווה לסכום זה. שימו לב שאין זה אומר שכל משקל בפני עצמו המשוך לקצה יהיה גדול מקצה שאינו בעץ, למרות שכך זה ברוב במקרים. ברגע שמצאנו את עץ התלויות קירוב ההתפלגות יכול להירשם בפורמט $A2$. מכאן ניתן לגזור פונקציה מאבחנת כמו שעשינו במקרה הבלתי תלוי

$$t_i = \text{Prob}(x_i = 1/x_j(i) = 1)$$

$$r_i = \text{Prob}(x_i = 1/x_j(i) = 0) \text{ ו } r_1 = \text{Prob}(x_1 = 1)$$

$$P(x_i/x_j(i)) = [t_i^{x_i} (1-t_i)^{1-x_i}]^{x_j(i)} [r_i^{x_i} (1-r_i)^{1-x_i}]^{1-x_j(i)}$$

אז

$$\begin{aligned} \log P(x) &= \sum_{i=1}^n [r_i \log r_i + (1-r_i) \log (1-r_i)] + \\ &+ \sum_{i=1}^n \left[r_{i,j} \log \frac{1-r_i}{1-r_{i,j}} + r_{i,j} \log \frac{r_i(1-r_i)}{(1-r_{i,j})} \right] + \text{constant} \end{aligned}$$

זוהי פונקצית שקלול לא ליניארית הפשוטה יותר מזו הנגזרת מ $A1$ כאשר מניחים שהמשתנים בלתי תלויים, כאשר $t_i = r_i$. הקבוע הוא בעל אותה פרשנות במונחים של הסתברות אפריורית ופונקצית הפסד. פונקצית ההחלטה השלמה היא

$$g(x) = \log P(x/w1) - \log P(x/w2)$$

שעכשיו מערבת חישובים בסדר גודל כפול מבמקרה הליניארי. רק הסכומים הכוללים $x_j(i)$ גורמים לפונקצית השקלול להיות שונה מבמקרה הליניארי, וזהו החלק המאפשר לקחת בחשבון ש x_i תלוי ב x_j . כשמשמשים בפונקצית השקלול במסמך המכיל את $x_j(i)$ או את $x(j)i$ וגם את x_i יקבל את התרומה מחלק זה של פונקצית השקלול.

קל יותר לראות איך $g(x)$ משלב משקלות שונות לתנאים שונים, אם מסתכלים במשקל שנתרם על ידי $g(x)$ למסמך נתון x בשביל מצבים שונים של הזוג $x_i x_j(i)$. כאשר $x_i = 1$ וגם $x_j(i) = 0$ המשקל הנתרם הוא

$$\log \frac{\text{Prob}(x_i = 1 | x_j = 0) \cdot A_{ij}}{\text{Prob}(x_i = 1 | x_j = 0) \cdot A_{ij}}$$

ובאופן דומה בשביל שלושת המצבים x_i וגם $x_j(i)$

כעת ראינו עד כמה פשוטה פונקצית השקלול הלא ליניארית. לדוגמה בהינתן מסמך ש i מופיע אבל $j(i)$ לא מופיע, אז המשקל הנתרם באמצעות $g(x)$ מבוסס על היחס בין שתי הסתברויות. הראשונה היא המופע של i בקבוצת מסמכים רלוונטים כשנתון ש $j(i)$ לא מופיע, השניה היא האנלוגיה על קבוצת מסמכים לא רלוונטית. על בסיס יחס זה אנו מחליטים אם יש מספיק נתונים בשביל להציב x עבור מסמכים רלוונטים או לא רלוונטים. חשוב לזכור בשלב זה שהנתונים בשביל לבצע את ההצבה בדרך כלל מתבססת על ההערכה של זוג הסתברויות.

הערכת משתנים

השימוש בפונקציה שקלול מהסוג הנ"ל באחזור דורש הערכה של פרמטרים מקושרים. כעת נבצע הערכה של t_i ו r_i במקרה הלא ליניארי, המקרה הליניארי יוצג לאחר מכן. ניתן כעת דוגמה של תהליך ההערכה המשתמש במקסימום פשטני באמצעות דמיון. הבסיס להערכות הוא הטבלה הבאה

	$r_1 = 1$	$r_1 = 0$	
$t_1 = 1$	11	01	01
$t_1 = 0$	01	11	01
	01	11	01

כאן אימצנו את הסימון לתאים כך ש $[x]$ משמעות מספר המופעים בתא x . אם נתעלם מכך שהקבוצה מבוססת טבלה; הערכותינו יהי כדלקמן:

$$f(t_1 = 1, r_1 = 1) = r_1 = \frac{01}{01}$$

$$f(t_1 = 1, r_1 = 0) = r_1 = \frac{01}{01}$$

במקרה הכללי יהיו בידינו שתי טבלאות מסוג זה היוצרות את הפונקציה $g(x)$, האחת בשביל להעריך פרמטרים המקושרים עם $P(x/w1)$ והשניה בעבור פרמטרים המקושרים עם $P(x/w2)$. בגבול יהיה לנו מידע מלא איזה מסמכים שנאספו רלוונטים ואיזה לא. אם היינו מחשבים את ההערכות למקרה הגבולי היינו מקבלים את הגבול העליון למקרה האחזור במקרה זה. באופן מציאותי יותר יהיה בידינו מדגם של מסמכים, סביר להניח קטן (לא בהכרח מקרי), כך שבעבור כל מסמך מידת הרלוונטיות הייתה ידועה. המקבץ הקטן יהווה מקור נתונים לכל טבלת 2 על 2 שנרצה לבנות. ההערכות לכן יהיו משוחדות בצורה בלתי נמנעת. הערכות הנ"ל הינן דוגמאות להערכות בנקודה. קיימות מספר דרכים בשביל לקבל חוקיות המתאימה להערכה בנקודה. לצערנו חוק ההערכה המוצלח ביותר הינו עדיין בעיה פתוחה [14]. למעשה סטטיסטיקנים מאמינים שאין לבצע כלל הערכה בנקודה [15]. למרות זאת באחזור מידע קשה לראות כיצד נמנעים מהערכה בנקודה. בכל חוק הערכה בנקודה נעשות מספר הנחות שרירותיות. למזלנו באחזור מידע קיים סיכוי לתקן הנחות אלו באמצעות מידע הנאסף ממערכות אחזור, ובכך מסירים חלק מהשרירותיות.

שתי הנחות מרכזיות נעשות כאשר יוצרים חוקי ההערכה על פי תיאוריית בייס:

(1) הצורה של ההתפלגות הקודמת במרחב המשתנים וגם

(2) הצורה של פונקציה ההפסד המשמשת למדידת הטעות הנעשית בעת הערכת המשתנה. ברגע ששתי הנחות אלו קיימות ומוגדרות היטב על ידי הגדרתם כהתפלגות ופונקציה הפסד, אז, ביחד עם עיקרון בייס המחפש לצמצם את הטעות בהינתן התצפית, ניתן להסיק מספר חוקי ההערכה. הספרות הסטטיסטית לא נותנת העדפות לחוקים מסוימים. (לפרטים ניתן להסתכל ב van Rijsbergen[2]). החוקים החשובים של הערכת a ביחס ל p כולם באים בצורה

$$\hat{p} = \frac{x + a}{n + a + 1}$$

כאשר x הוא מספר ההצלחות מתוך n ניסיונות, a ו b הם משתנים המוכתבים על ידי השילוב בין פונקציה הקדימות וההפסד. כך יש בידינו קבוצה שלמה של חוקי ההערכה. לדוגמה כאשר a=b=0 יש לנו את ההערכה x/n, וכאשר a=b=0.5 יש לנו את החוק של Good [16]. החוק האחרון הוא למעשה החוק של Robertson and Sparck Jones [1] השתמשו בו בהערכותיהם. כל מצב של a ו b ניתן להצדיק בתנאים של סבירות תוצאות התפלגות אפריורית.

מאחר ומה שסביר בעיני אדם אחד אינו סביר בעיני אדם שני, ההחלטה הסופית צריכה להתבסס על הביצועים בניסויים. למזלנו באחזור מידע אנו במצב ייחודי לבצע ניסויים אלו.

סיבה חשובה אחת כדי שחוקי הערכה יהיו שונים מהצורה הפשטנית x/n, היא בכך שזה אינו סביר בעבור מדגמים קטנים. לדוגמה מדגם עם דגימה אחת (n=1) ותוצאה (x=0) או (x=1) שיגרום לתוצאה בעבור p (p=0) או (p=1). ברור שזהו מגווח, מאחר וברוב המקרים אנחנו יודעים בהסתברות גבוהה ש $0 < p < 1$.

בשביל להתגבר על קושי זה אנו יכולים לשלב מידע מוקדם בהתפלגות על הערכים האפשריים אותם אנו מנסים להעריך. ברגע שנקבל את הפיסביליות של זה ונגדיר דרך למדוד את טעות ההערכה, עיקרון בייס (או עקרונות אחרים) לרוב יובילו לתוצאה שונה מ x/n . לא נוסיף יותר בנושא זה, קיים פירוט רב יותר בספרות סטטיסטית.

סיכום

בשלב זה נסכם בשפה פשוטה. הנקודה הראשונה היא שעד כה ניסינו להעריך את $P(\text{relevance}/\text{document})$ כלומר ההסתברות שמסמך יהיה רלוונטי. למרות שקל לכתוב משפט זה לא ברור שיש לו משמעות. רלוונטיות של מסמך זהו מונח שאינו מוגדר היטב, יש שיטענו שזה יכול להיות 1 או 0. אחרים יטענו שניתן להסיק זאת על ידי ההסתברות $P(x/\text{relevance})$ לפי נוסחת בייס. לא נטיל ספק בשימוש בנוסחת בייס השאלה אם ל $P(\text{relevance}/x)$ יש משמעות בהקשר של אחזור מידע. נניח שכן כדי שנוכל לקשר בין $P(\text{relevance}/x)$ לבין התפלגות המידע בעבור אינדקסים.

גישה זו הייתה חסרת משמעות אלמלא היינו מניחים שההתפלגות להרבה אינדקסים למסמכים רלוונטים שונה מזו למסמכים לא רלוונטים. אם היינו מניחים את ההפך אזי $P(\text{relevance}/\text{document})$ יתפלג כמו $P(\text{relevance})$, קבוע לכל המסמכים ומכאן לא מסוגל להחליט מי מהמסמכים הוא חסר ערך לאחזור. מכאן אנו מניחים שמידע עקיף נובע מההתפלגות המשותפת של אינדקסים מעל שתי קבוצות שיעזור לנו להבחין בינם. ברגע שנקבל השקפה זו אנו מחויבים לניסוח הנגזר מן הנאמר. המחויבות היא שאנו חייבים לנחש את $P(\text{relevance}/\text{document})$ בדיוק מרבי, או לחילופין לנחש את $P(\text{document}/\text{relevance})$ וגם את $P(\text{relevance})$ באמצעות התפלגות המידע הנובעת מהמסמך.

ההסבר במונחי דירוג במקום אבחנה פשוטה הוא טריוויאלי: ההפחתה הנובעת מהקבוע $g(x)$ בהדרגתיות קטנה ובאמצעות זאת מגדילה את מספר המסמכים המאוחרים. תוצאות הדירוג הם אופטימאליות עקב העובדה שעבור כל הפחתה אנו ממזערים את הסיכון הכללי. צריך להתייחס לכך בזהירות כי זה מניח שיש בידינו את $P(x/w_i)$ בצורה נכונה והערכות שלנו הם הטובות ביותר שאפשר, לא סביר שיתממש בניסוי.

אם נסכים שהמשתמש יקבע את גודל ההפחתה לאחר האחזור אזי הצורך בתיאוריה בנושא ההפחתה בטל. ההשלכה הנובעת מכך שבמקום לעבוד אם היחס

$$\frac{R_i(\text{relevance})}{R_i(\text{non-relevance})}$$

נעבוד אם היחס

$$\frac{R_i(\text{relevance})}{R_i}$$

במקרה האחרון איננו מסתכלים על בעיית האחזור כאבחנה בין מסמכים רלוונטיים ללא רלוונטיים, במקום אנו מנסים לחשב את ההסתברות $P(\text{relevance}/x)$ לכל מסמך x ולהציג גודל זה בסדר יורד למשתמש. בכל צורה שנסתכל בזה עדיין נדרש הערכה של שתי פונקציות הסתברות משותפת.

חוקי ההחלטה הנובעים מהני"ל מבוטאים בצורה של $P(x/w_i)$. מכאן אדם יכול להסיק שהערכה של הסתברויות אלו הכרחית לביצועי האחזור, וכמובן העובדה שרק ניתן להעריך הסתברויות אלו היא הסבר לחלק מחוסר האופטימאליות של הביצועים. בכדי להעריך צריך להניח מספר דברים לגבי הצורה של $P(x/w_i)$. הנחה ברורה היא להניח חוסר תלות סטוכסטי לגבי מרכיבי x . באופן כללי זה לא ריאלי מאחר וטבע פונקצית האחזור שאינדקסים יהיו תלויים אחד בשני. נצטט מאמר מוקדם של "מרון" בנקודה זו: "בשביל לעשות זאת אדם יצטרך ליצוא תוכנת מחשב לבצע ניתוח אנליטי של אינדקסים כך שהמחשב "יידע" מי מהתנאים קרוב ביותר אל תנאי אחר ויכול להצביע על הכיוון הוודאי ביותר שבו ניתן להגדיל בקשה". לכן גישה ריאלית יותר היא להניח

תלות כלשהי בין התנאים כאשר מעריכים את $P(x/w_1)$ ואת $P(x/w_2)$ או $P(x)$.

כעת נמשיך ונדון בדרכים להשתמש במודל זה של אחזור ובוזמנית נדון בבעיות הצעות. בהתחלה לא נשנה כמעט את המודל אבל לאחר מכן הדון בדרכים מעשיות למימוש שלא בהכרח ייצמדו להנחות על פיהם בנינו מודל זה. באופן טבעי ההצדקות לכך טמונות בניסויים שרבים עדיין נותר לבצע [17]. נתחיל בהסבר לשינוי קטן הנובע בצורך לצמצם את ממדי הבעיה.

מגרעת הממדים

בהסקת חוקי ההחלטה הנחנו שמסמך מוצג באמצעות ווקטור בעל n ממדים כאשר n הוא גודל האינדקסים. בדרך כלל n יהיה גדול למדי, ולכן ממדי וקטורי המסמך יהיה לרוב גדול ממספר הדגימות המשמש להעריך את הפרמטרים בפונקציית ההחלטה. דבר זה יוביל לבעיות שצוינו באופן מתמיד בספרות בנושא זיהוי דפוסים. למרות שניתוח הבעיה בזיהוי דפוסים רלוונטי לאחזור מידע גם כן, הפתרון אינו ישים מיידית. בזיהוי דפוסים הבעיה היא: בהינתן מספר הדגימות ששימשו "לאמן" את פונקציית ההחלטה (פונקציית השקלול שלנו), האם קיים מספר אופטימאלי של מדידות שניתן ליצור מדפוס בלתי ידוע כך שממוצע ההסתברויות של השמה נכונה ניתן למקסמו? במקרה שלנו בכמה אינדקסים ניתן להשתמש כדי לקבוע רלוונטיות. היוז[18] מראה שבעבור מבנה הסתברותי כללי מספר המדידות הוא קטן באופן מפתיע למרות שדגימות בגודל סביר משמשות בכדי "לאמן" את פונקציית ההחלטה.

באופן אידיאלי נרצה לבחור תת קבוצה (קטנה) של אינדקסים שבעבורם פונקציית השקלול $g()$ תוגבל ובכך תמקסם את ממוצע ההסתברויות של השמה נכונה. בזיהוי דפוסים יש טכניקות מסובכות לעשות זאת לבעיות מקבילות. באחזור מידע התמזל מזלנו וקיימת דרך טבעית שבאמצעותה ניתן לצמצם את ממדי הבעיה. אנו מקבלים שתנאי השאלה הם מדריך טוב לתכונות המוצלחות ביותר ביישום של פונקציית $g()$ בשביל להחליט בין רלוונטיות ואי רלוונטיות של מסמכים. לכן במקום לחשב את פונקציית השקלול לכל התנאים אנו מגבילים את $g()$ לתנאים המצוינים בשאלה ואפשרי שגם לתנאים דומים. המשמעות של הנ"ל היא שבזמן תהליך האחזור כל המסמכים מוקרנים ממרחב רב ממדי לתת מרחב הנפרש באמצעות מספר תנאים קטן יותר.

פרטים חישוביים

נעבור כעת לחלקים המעשיים יותר בחישוב $g(x)$ לכל x כאשר מניחים שהמשתנים x_i תלויים סטוכסטית. המטרה המרכזית של חלק זה היא להוכיח שהחישובים הנם פיזיביליים. הדרך הבהירה ביותר לעשות זאת היא לחשב עבור כל אובייקט את ה EMIM את ה MST ואת ה $g(\cdot)$ בנפרד ובסדר זה.

1. חישוב ה EMIM

ניתן לפשט את חישוב ה EMIM וכן ניתן להעריך אותו וכך להפחית את זמן החישוב. נבחן תחילה את הדרך הראשונה:

כשמחשבים את $I(x_i, x_j)$ לצורך בניית ה MST אנו צריכים לדעת רק את סדר הדירוג עבור כל ה הערכים עצמם אינם משנים. לכן ניתן להשתמש במידע על ההסתברויות בטבלה הבאה:

	$x_i = 1$	$x_i = 0$	
$x_j = 1$	[4]	[2]	[7]
$x_j = 0$	[5]	[6]	[8]
	[5]	[6]	[8]

במקרה זה תהיה הפונקציה $I(x_i, x_j)$ מונוטונית לגמרי עם הפונקציה:

$$[1] \log \frac{[1]}{[5][7]} + [2] \log \frac{[2]}{[6][7]} + [3] \log \frac{[3]}{[5][8]} + [4] \log \frac{[4]}{[6][8]}$$

וזהו אכן דרך פשוטה לחישוב ה EMIM. נחשוב כעת על מקרה שבו אנו רוצים לחשב את $p(x)$ במקרה זה, ה MST מבוסס על מידע המופיע מספר פעמים והנגזר מכל האוסף. ברגע שיש לנו את זה ([1]) ואנו יודעים את מספר המסמכים בקובץ ([9]) אז כל קובץ הפוך יכול את יתר המידע על התדירות שאנו צריכים כדי למלא את יתר התאים בטבלה. כלומר, בעזרת ([5]) ו ([7]) הנתונים על ידי הקובץ ההפוך אנו יכולים להסיק את תוצאתם של [2] [3] [4] [6] ו [8].

ישנה בעיה של טיפול בערכי 0 בתאים 0-4 והיא נפתרת אם קובעים שהתוצאה של $0 \log 0$ שווה ל אפס. שאר התאים אינם יכולים לקבל ערכי אפס.

כעת נדון בדרך השנייה שבה מעריכים את ה EMIM מרון וקונס השתמשו בעבודתם המוקדמת בנוסחה:

$$d(x_i, x_j) = P(x_i = 1, x_j = 1) - P(x_i = 1) P(x_j = 1) (*)$$

כדי למדוד את הסטייה מעצמאות עבור שני אברי אינדקס i, j . לבד מה \log זוהי הדרך הראשונה של הרחבת ה EMIM. אם מחשבים את ה MST (עץ תלות) לפי הנוסחה הקודמת, תוצאות החישוב של $p(x|w_i)$ לא יחיו מושלמות אבל ההערכה תהיה קרובה מספיק מחישוב ותחסוך זמן רב לעומת חישוב מלא של ה EMIM. באופן דומה, לוי השתמש בנוסחה:

$$\log \frac{P(x_i = 1, x_j = 1)}{P(x_i = 1) P(x_j = 1)}$$

להערכת ה EMIM. ללא ספק, ישנן דרכים נוספות להערכת ה EMIM אך על השאלה אם ניתן להשתמש בהן למציאת עץ תלות שבעזרתו ניתן להגיע להערכה טובה של פונקציה ההסתברות המשותפת ניתן לענות לאחר ניסויים אמפיריים.

2. חישוב ה-MST :

ישנם מספר אלגוריתמים לחישוב ה-MST. כנראה שהטוב שבהם הוא האלגוריתם החדש של וויטני שהסיבוכיות שלו היא $O(n^2)$ כאשר n הוא מספר איברי האינדקס בעץ התלות. זהו אינו מכשול לשימוש באלגוריתם עבור קבוצות גדולות היות וקל לחלק את המידע באמצעות טכניקות שנידונו קודם ואחר כך, ניתן לבנות את העץ על ידי הפעלת האלגוריתם על כל אשכול של איברי אינדקס. דבר זה מוריד את הסיבוכיות ל $O(k^2)$ כאשר k קטן בצורה משמעותית מ n .

לפי עקרונות אלו, הוכיחו בנטלי ופרידמן ששימוש במרחב הנדסי שבו איברי האינדקס הם נקודות במרחב, הסיבוכיות של האלגוריתם לחישוב ה-MST יורדת ל $O(n \log n)$ בקרוב ויותר מכך ניתן לבנות עץ מפתח שהוא כמעט MST שגם הסיבוכיות שלו היא $O(n \log n)$.

ישנו חוסר יעלות גדול, בכך שכל $n(n-1)/2$ האסוציאציות מחושבות ורק בחלק קטן מהן יש משמעות בהקשר שהן אינן ערך אפס וניתן להשתמש בהן בעץ המפתח. רוב האסוציאציות מקבלות ערכי אפס וניתן להתעלם מהן, אך לרוע המזל ניתן לדעת איזו אסוציאציה תקבל ערך אפס רק לאחר חישוב. קרופט, גילה לאחרונה להתעלם מאסוציאציות מבלי לחשבן תחילה. דרך זאת מניחה שהקובץ והצורה ההפוכה שלו קיימים ולכן יתכן שבאלגוריתם של קרופט נדרש זמן להפוך את הקובץ.

3. חישוב $g(x)$

יש להדגיש שבמקרים לא ליניאריים, הערכת הפרמטרים של $g(x)$ תגרום באופן אידיאלי ל MST שונה עבור כל $p(x|w1)$ ו $p(x|w2)$. כמוכן ישנו מידע מלא על הופעת איברי האינדקס בקבוצות ה שייכות/אחוסר שייכות במצב הניסויי. חישוב $g(x)$ באמצעות מידע מלא עשוי להיות נחוץ כשגוזרים את הגבולות העליונים של יעילות האחזור לפי המודל, כפי שנעשה במקרה העצמאי אצל רוברטסון, ספארק וג'ונס. במקרה שבו אין מסמכים רלוונטיים שידועים מראש, יהיה צורך להשתמש בטכניקה של משוב רלוונטיות כדי להעריך את הפרמטרים ואז הביצועים ישאפו לגבול העליון. עובדה זאת, מובטחת לפי תיאוריית השאיפה שנדונה בעמוד 106 בפרק 5.

ישנם מספר דרכים לפרש את המודל לגבי $g(x)$ והדבר תלוי בשאלה האם רוצים לקבוע את ההסתברות אפריורית או אפוסטריורית. במקרה הראשון יש לקבוע את ה MST הן עבור איברי האינדקס במסמכים הרלוונטיים והן במסמכים הלא רלוונטיים. היות וניתן לעשות זאת רק עם מידע מדגמי ולא מלא, עץ התלות יהיה רחוק מן האופטימאלי. דרך היוריסטית להתגבר על הבעיה היא לבנות עץ תלות עבור כל האוסף. מבנה העץ צפוי להיות במקרה זה זהה למבנה של שני עצי התלות המבוססים על המסמכים הרלוונטיים והלא רלוונטיים. $p(x|w1)$ ו $p(x|w2)$ מחושבים על ידי חישוב ההסתברות המותנית עבור הצמתים בעץ התלות. עד כמה חישוב זה מדויק, ניתן לקבוע באמצעות ניסויים בלבד.

אם מניחים שההסתברויות נקבעות אפוסטריורית, הרי ניתן לדרג את המסמכים לפי $p(w1|x)$ ולהשאיר למשתמש את ההחלטה מתי הוא ראה מספיק. כלומר, אנו משתמשים בנוסחה:

$$P(x|w1) = \frac{P(x|w1) P(w1)}{P(x)}$$

להערכת ההסתברות של רלוונטיות עבור כל מסמך X . כאן אנו זקוקים להערכה בלבד של $p(x|w1)$ ולאחר חישוב $p(x)$ אנו משתמשים בעץ מפתח עבור כל האוסף מבלי להתייחס למידע שייכות. לגישה זאת מספר יתרונות. יתרון אחד הוא, שאם יש הנחת עצמאות עבור כל האוסף הרי יש קונסיסטנטיות עם הנחת העצמאות עבור המסמכים הרלוונטיים. הנחות אלו יגבירו את החישובים בצורה משמעותית. למרות שהנחת העצמאות עבור $W1$ איננה סבירה, היא עשויה להיות כורך המציאות, היות וקשה למדוד את התלות המדויקת בניסויים או במדגמים.

דרך אלטרנטיבית לשימוש בעץ תלות (היפוטיזת האסוציאציות)

לפי הטיעונים בפסקה הקודמת ניתן להסיק שעץ התלות היחיד האפשרי הוא עץ תלות עבור כל האוסף. נמשיך הלאה לפי קו זה. בניית עץ תלות עבור איברי האינדקס מבלי להשתמש במידע על תלות בניהם דומה לבניית טבלת סיווג לאיברי אינדקס. בפרק 3 ציינתי את הקשר בין ה-MST והקשר היחיד והוכחתי שהאחד אינו כה שונה מהשני. עובדה זאת מובילה לרעיון ששימוש בעץ תלות דומה לשימוש בקיבוץ איברים.

הרעיון של קיבוץ איברים הוסבר בפרק 2. ניתן לסכם זאת באמירה שאפשר ליישם מספר אסטרטגיות של הוספת ומחיקת איברים לפי השיטה של קיבוץ איברים. אם נתעלם לרגע מהמחיקה, ברור כיצד ניתן להשתמש בעץ התלות כדי להוסיף איברים או להרחבת השאילתה. הסיבה לכך ניתנה ב-1964 ע"י מרון ששאל: "כיצד ניתן להגדיל את ההסתברות לאחזור מספר מסמכים המכילים מידע ייחודי". דרך ברורה מאלה היא להרחיב את הבקשה המקורית באיברי אינדקס בעלי משמעות קרובה לאינדקסים בבקשה המקורית. ההנחה כאן היא שמשמעות קשורה יכולה להתגלות ע"י אסוציאציה סטטיסטית. לכן אני מציע שבהינתן שאילתה שבה הקריטריונים אינם ברורים, אנו משתמשים בעץ התלות כדי לגלות אילו איברים נוספים הקשורים באיברים בשאילתה עשויים לעזור לנו במציאת מסמכים רלוונטיים. כלומר אני טוען שאיברי אינדקס הקשורים ישירות לאיבר מתוך שאילתה בעץ תלות, עשויים לעזור באחזור. למעשה ניסחתי מחדש את ההיפוטיזה עליה מבוססת תיאורית קיבוץ האיברים. ברצוני כעת לנסח זאת פורמאלית ולקרוא לכך היפוטיזת האסוציאציות: אם איבר אינדקס כלשהו יכול להועיל לסינון מסמכים רלוונטיים ומסמכים לא רלוונטיים, אז איברי אינדקס בעלי אסוציאציות לאיבר זה יועילו ככל הנראה גם הם.

אנו מפרשים היפוטיזה זאת בכך שאם איבר בשאילתה נבחר ע"י המשתמש, אזי הוא מסנן טוב בין מסמכים רלוונטיים ולא רלוונטיים ולכן אנו מעוניינים באסוציאציות הקרובות שלו. ההיפוטיזה אינה דנה בדרכים שבהם ניתן למדוד את האסוציאציה בין איברי אינדקס אך בפרק זה דנתי בשימוש ב-EMIM. היא גם לא דנה בכימות הסינון, בכך אדון בפסקה הבאה. במידה מסוימת בת זוג של תיאוריית קיבוץ האיברים וניתן לבחון אותן באותה דרך.

כוח הסינון של מונח אינדקס:

בעמוד 120 הגדרתי:

$$K_i(N, n, P) = \log \frac{P^i(P - 1)}{(n - 1)!(N - n - P + 1)}$$

ולמעשה הערתי שנוסחה זאת מודדת את הכוח של איבר i לסנן בין מסמכים רלוונטיים ולא רלוונטיים. המשקולות בפונקציה נגזרים מהנחת העצמאות. אם נתעלם לרגע מכך שמשקלים אלו הם תוצאה של מודל מסוים, ולחלופין נתייחס לכוח הסינון של איבר אינדקס כדבר העומד בפני עצמו. זהו בבדאי אינו דבר קל לעשותו. סלטון בעבודתו חיפש דרכים אפקטיביות למדוד את כוח הסינון של איבר אינדקס. זה נראה הגיוני להצמיד לכל איבר אינדקס הנכנס לתהליך אחזור, משקל המתייחס לכוח הסינון שלו. שימוש ב K_i כדרך למדוד כוח זה בעייתי במקצת היות והוא לא מוגדר כאשר ארגומנטים של פונקציה הלוגריתם הופכים לאפס. לכן אנו מחפשים דרך חלקה יותר למדידת כוח הסינון. הפונקציה עליה אני עומד להמליץ היא אכן חלקה יותר ומאפשרת לי להסיק מסקנה כללית בעלת חשיבות רבה. יש להדגיש עם זאת שהפונקציה המוצעת היא רק אפשרות אחת, אך כל פונקציה מתאימה תהייה דומה לזאת המוצעת כאן.

במקום להשתמש ב K_i אני מציע להשתמש ברדיוס המידע כפי שמוגדר בפרק 3 בעמוד 42, כדרך למדידת כוח הסינון של איבר אינדקס. זהו קרוב משפחה של "אמצעי המידע המשותף הצפוי" שנשתמש בו בהמשך. נשתמש ב U וב V בתור משקולות חיוביות כאשר $1=V+U$ ובסימונים הרגילים של פונקציות הסתברות. אנו יכולים לנסח את רדיוס המידע כך:

$$\begin{aligned} & u^{P(x_i=1|\mathcal{X}_1)} \log \frac{P(x_i=1|\mathcal{X}_1)}{u^{P(x_i=1|\mathcal{X}_1)} + v^{P(x_i=1|\mathcal{X}_2)}} + \\ & + v^{P(x_i=1|\mathcal{X}_2)} \log \frac{P(x_i=1|\mathcal{X}_2)}{u^{P(x_i=1|\mathcal{X}_1)} + v^{P(x_i=1|\mathcal{X}_2)}} + \\ & + u^{P(x_i=0|\mathcal{X}_1)} \log \frac{P(x_i=0|\mathcal{X}_1)}{u^{P(x_i=0|\mathcal{X}_1)} + v^{P(x_i=0|\mathcal{X}_2)}} + \\ & + v^{P(x_i=0|\mathcal{X}_2)} \log \frac{P(x_i=0|\mathcal{X}_2)}{u^{P(x_i=0|\mathcal{X}_1)} + v^{P(x_i=0|\mathcal{X}_2)}} \end{aligned}$$

הפירוש המעניין של רדיוס המידע מוצג בצורה בהירה יותר במונחים של פונקציות הסתברות רציפות. במקום להשתמש בצפיפויות של $p(,|w2$ ו $p(,|w1$ אשתמש בהסתברויות התואמות של $U1$ ו $U2$. נגדיר תחילה את הממוצע של שתי הסתעפויות ישירות:

$$R(u1, u2/v) = uI(u1/v) + vI(u2/v)$$

$I(ui|v)$ מודד את כמות המידע שמרוויחים מהידיעה שיש לדחות את V לטובת U_i . כעת רדיוס

$$\frac{u1v}{v} \mathcal{R}(\mu_1, \mu_2 | \mathcal{P})$$

המידע הוא המינימום של:

וכך מורידים את ה V השרירותי. למעשה, מסתבר שהמינימום מושג כאשר:

$$v = u u1 + v u2$$

זהו למעשה ממוצע של שני הערכים שיש לסנן. אם אנו מקבלים את V ו U כהסתברויות פריוריות אז V מובע בצפיפות:

$$p(x) = p(x/w1) P(w1) + p(x/w2) P(w2)$$

המוגדרות עבור כל האוסף ללא התייחסות לרלוונטיות. בערך זה אנו בטוחים אך את הערכים $U1$ ו $U2$ יש לנחש. לכן זה הגיוני שבעת מדידת ההפרש בין $U1$ ו $U2$ ימזג את כל המידע הקיים. רדיוס המידע בדיוק עושה זאת.

ישנה בעיה טכנית אחת בשימוש ברדיוס המידע או בכל דרך מדידת סינון אחרת המבוססת על ארבעת התאים בטבלת השכיחות. הבעיה היא ששיטות אלו אינן מבחינות בין התרומות השונות של התאים השונים בטבלה כך שלמשל, איבר אינדקס יכול להיות מסנן טוב היות והוא מופיע באופן שוטף במסמכים לא רלוונטיים וכמעט ואינו מופיע במסמכים רלוונטיים. ולכן שקלול איבר אינדקס באופן פרופורציונאלי לכוח הסינון שלו בכל עת שהוא מופיע במסמך הוא הדבר השגוי לעשותו במקרה זה. כלומר יש להשתמש במידע שבטבלת השכיחות כשיוצרים תוכנית שקלול.

היפותזת תוספת הסינון

בפסקה הקודמת, הנחתי הנחה על עצמאות או תלות בצורה פשוטה. או שהנחתי עצמאות של $W1$ או $W2$ או תלות של שניהם. אך כפי שנאמר קודם, זאת אינה הדרך היחידה של הנחת ההנחות הללו. רוברטסון וספארק גיונס טוענים שהנחת תלות או עצמאות במסמכים כלשהם יכולה להוביל למסקנה בנוגע לאוסף המסמכים כולו. כדי לראות זאת, ניקח שני אינדקסים i, j כאשר :

$$P(x_i, x_j) = P(x_i, x_j / w_1)P(w_1) + P(x_i, x_j / w_2) P(w_2)$$

$$P(x_i) P(x_j) = [P(x_i / w_1)P(w_1) + P(x_i, w_2) P(w_2)] [P(x_j / w_1) P(w_1) + P(x_j, w_2) P(w_2)]$$

לעצמאות בלתי מותנית, צריך שיתקיים :

$$P(x_i, x_j) = P(x_i, /w_1) P(x_j, w_1) P(w_1) + P(x_i /w_2) P(x_j/ w_2) P(w_2)$$

וזה קורה רק כאשר $p(w_1)=0$ או כאשר $p(x_i|w_1)=p(x_i|w_2)$ או כאשר $p(x_j|w_1)=p(x_j|w_2)$. ואם ננסח זאת במילים, כאשר לפחות איבר אינדקס אחד הנו חסר תועלת בסינון בין מסמכים רלוונטיים ולא רלוונטיים. ובאופן כללי עצמאות מותנית תוביל לתלות בלתי מותנית. כעת, נניח שאיברי האינדקס הם אכן עצמאים בצורה מותנית, נקבל את התוצאות יוצאות הדופן הבאות :

קנדל וסטוארט מגדירים מקדם מתאם חלקי בין שני ערכים כך :

$$\rho(X, Y|W) = \frac{\rho(X, Y) - \rho(X, W)\rho(Y, W)}{(1 - \rho(X, W)^2)^{1/2} (1 - \rho(Y, W)^2)^{1/2}}$$

כאשר $[\rho](..|w)$ ו $[\rho](..)$ הנם בהתאמה מקדמי המתאם. כעת אם X ו Y הנם עצמאים בצורה מותנית, אז :

$$[\rho](X, Y|W) = 0$$

שמוכיח, כאשר משתמשים בביטוי על המתאם החלקי, ש :

$$[\rho](X, Y) = [\rho](X, W) [\rho](Y, W)$$

היות ו :

$$|[\rho](X, Y)| \leq 1, |[\rho](X, W)| \leq 1, |[\rho](Y, W)| \leq 1$$

זה כשלעצמו, מוכיח שכשמשמשים בהיפותזת העצמאות המותנית :

$$|[\rho](X, Y)| < |[\rho](X, W)| \text{ or } |[\rho](Y, W)| (**)$$

ולכן, אם W הנו משתנה מקרי המייצג את הרלוונטיות, אזי המתאם בינו לבין כל אחד מאיברי האינדקסים גדול מהמתאם בין איברי האינדקסים.

באופן איכותי, אנסה להכליל זאת לפונקציות נוספות בנוסף למדדי מתאם. לינפוט, מגדיר סוג של אמצעי של מתאם התורם מידע :

$$rij = (1 - \exp(-2I(x_i, x_j)))^{1/2} \leq rij \leq 1$$

או :

$$I(x_i, x_j) = -\log \frac{(1 - r_{ij}^2)}{2}$$

כאשר $I(x_i, x_j)$ הנו האמצעי המוכר למדידת המידע ההדדי. אבל מוריד את מקדם המתאם הסטנדרטי $[\rho](..)$ אם (x_i, x_j) מתפלג נורמאלית. לכן זה הגיוני להניח שעבור התפלגות שאינה נורמאלית rij יתנהג בצורה דומה ל $[\rho](..)$ ויהיה גם הוא מספק. אבל rij מונוטונית עבור $I(x_i, x_j)$ ולכן גם היא מספקת. לכן אנו יכולים לומר כעת שתחת עצמאות מותנית, המידע

הגלום באיבר אינדקס אחד על איבר אינדקס אחר הוא פחות מהמידע שיש בכל אחד משני האיברים לגבי המשתנה המותנה W . או בניסוח מתמטי:

$$I(x_i, x_j) < I(x_i, W) \text{ or } I(x_j, W),$$

כאשר $i(\cdot, w)$ הוא רדיוס המידע, עם המשקלים שלו המבוטאים כהסתברויות פריאוריות. ונזכור כי $I(\cdot, w)$ הוצע כמדד לכוח הסינון. אני סבור שיש לבטא מסקנה זאת כהיפוטזה, כאשר W מבוטה כתלות.

היפותזת תוספת הסינון: תחת ההיפוטזה של עצמאות מותנית, המידע הסטטיסטי הגלום באיבר אינדקס אחד על איבר אינדקס אחר פחות מהמידע המגולם בכל אחד מהם על הרלוונטיות.

עלי להדגיש שהטענה שהובילה להיפוטזה אינה מוכחת. הטענה היא טענה איכותית בלבד. מלבד זאת, הטענה מספקת, בתוספת להיפוטזה, צידוק מסוים ובסיס תיאורטי לשימוש ב MST המבוסס על $I(x_i, x_j)$ לשיפור אחזור המידע. היפוטזת הסינון היא דרך לחזק את היפוטזת האסוציאציה תחת עצמאות מותנית.

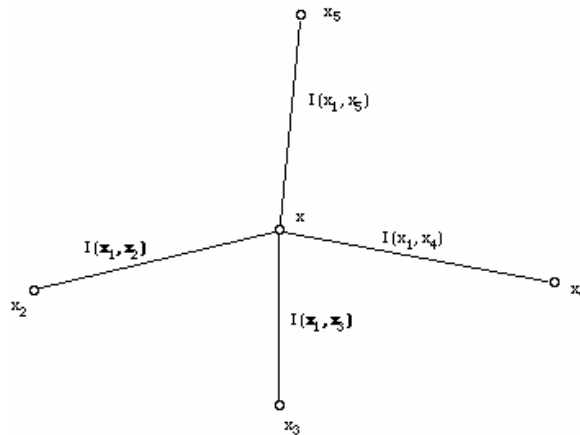


Figure 6.2

תוצאה אחת של היפותזת הסינון, היא שהיא מספקת בסיס הגיוני לדירוג איברי אינדקס הקשורים לשאילתה, בעץ תלות בסדר של I (איבר אינדקס, איבר שאילתה) ערכים שישקפו את סדר ערכי כוח הסינון. הבסיס של זה הוא, שכל שאיבר אינדקס קשור לאיבר שאילתה (שנמדד ע"י EMIM) כך יש סיכוי שכוח הסינון שלו יהיה גדול יותר. כדי להבהיר את הנקודה, נשתמש בדוגמה שנידונה בתרשים 6.2 שלמעלה. נניח ש X_1 מייצג איבר שאילתה וש:

$$I(x_1, x_2) < I(x_1, x_3) < I(x_1, x_4) < I(x_1, x_5)$$

אז, אומרת ההיפוטזה שלנו, שמבלי לדעת מראש מה כוח הסינון של כל אחד מאיברי האינדקס 2,3,4,5 זה הגיוני להניח ש:

$$I(x_2, W) < I(x_3, W) < I(x_4, W) < I(x_5, W).$$

ברור שאיננו יכולים להבטיח סדר זה, אך זהו הניחוש הטוב ביותר בהתחשב בבורות שלנו.

הערות ביבליוגרפיות

הרקע התיאורטי של פרק זה נמצא רק בכמה מאמרים. האחד עוסק בגישה של נתינת משקל הסתברותית המבוססת על מידע על רלוונטיות שנגזר מעבודתו של יו ושותפיו. האחר הוא מאמר המפורסם של רוברטסון וספארק ג'ונס. לרוע המזל שני המאמרים הללו מבוססים על ההנחה של עצמאות סטוכסטית. המאמר שלי ומאמרם של בוקשטיין וקראפט הם שני המאמרים היחידים שידועים לי שמנסים לבנות מודל בלי הנחה זאת. יתכן שיש לציין גם את מאמרו המוקדם של נגוטיה שמנסה לנקוט בגישה לא ליניארית. סיכומו של רוברטסון נותן מידע על כמה התפתחויות חדשות בנושא.

לפי דויל, היו מרון וקונס הראשונים שתארו שימוש באסוציאציות של איברי האינדקס לצורך ייעול החיפוש. עם זאת, דויל עצמו עבד בשנות החמישים על רעיונות דומים, וכתב כמה מאמרים בנושא. בשנת 1961 סיפק סטילס, שכבר היה מודע לעבודתם של מרון וקונס, פרוצדורה של שימוש בחזרה כפולה בחיפוש, בצורה דומה לשיטה המבוססת על עץ התלות. הוא גם השתמש ב[[chi]] כדי למדוד אסוציאציה בין איברי אינדקס, שמבחינה מתמטית דומה לשימוש במידת חזרת המידע המשותף הנצפית. אם זאת, יש להעדיף את השיטה האחרונה, כשמודדים תלות (ראה גודמן וקרוסקל). סטילס עצמו היה מאוד נחרץ לגבי התועלת שבשימוש באסוציאציות בין איברי אינדקס. הוא ראה שבאמצעותם, ניתן לאתר מסמך הרלוונטי לבקשה גם אם למסמך עצמו לא היה אינדקס זהה לזה שבבקשה.

המודל בפרק זה קשור לעוד שני רעיונות שנידונו במחקרים מוקדמים. האחד הוא משקל הצפיפות של מסמכים הפוכים שכבר נידון בפרק 2 והשני הוא קיבוץ איברים. אם נתייחס תחילה לרעיון הראשון, שמגיע עד למחקרם המוקדם של אדמונסון וויליס, ניתן לכתוב:

$$P(\text{relevance/document}) \propto \frac{1}{P(\text{document})}$$

במילים: עבור כל מסמך, ההסתברות לרלוונטיות קשורה ביחס הפוך להסתברות של חזרה אקראית. אם מניחים ש $p(\text{document})$ הוא המכפלה של איברי אינדקס בודדים שקיימים או חסרים במסמך, הרי שלאחר חישובים מתאימים מקבלים את כלל המשקל של צפיפות המסמך. הוא מניח את הסבירות הגבוהה ש $p(\text{document relevance})$ קבוע עבור כל המסמכים. השאלה מדוע עקרון זה עובד בצורה כה טובה אינה ברורה אך יתכן שיש לה תשובה בעבודה האחרונה של יו וסלטון.

הקשר עם קיבוץ האיברים נעשה עוד בתחילת הפרק. ניתן להתייחס לעץ המפתח בתור סיווג לאיברי אינדקס. אחת ההשלכות החשובות של המודל המתואר בפרק זה היא שיש פירות ברור כיצד להשתמש בעץ באחזור מידע. עבודות קודמות בתחום התייחסו למקרים פרטיים ולא הובילו למסקנות ממצות.

זה כבר זריך להיות ברור, שהמודל הכמותי מכיל בתוכו מספר נושאים כגון קיבוץ איברים, ניתוח אסוציאציות, משקל הצפיפויות במסמך ומשקל רלוונטיות.

References

1. ROBERTSON, S.E. and SPARCK JONES, K., 'Relevance weighting of search terms', *Journal of the American Society for Information Science*, **27**, 129-146 (1976)
2. van RIJSBERGEN, C.J., 'A theoretical basis for the use of co-occurrence data in information retrieval', *Journal of Documentation*, **33**, 106-119 (1977).
3. BOOKSTEIN, A. and KRAFT, D., 'Operations research applied to document indexing and retrieval decisions', *Journal of the ACM*, **24**, 410-427 (1977).

4. MARON, M.E., 'Mechanized documentation: The logic behind a probabilistic interpretation', In: *Statistical Association Methods for Mechanized Documentation* (Edited by Stevens et al.) National Bureau of Standards, Washington, 9-13 (1965).
5. OSBORNE, M.L., 'A Modification of Veto Logic for a Committee of Threshold Logic Units and the Use of 2-class Classifiers for Function Estimation', Ph.D. Thesis, Oregon State University (1975).
6. GOOD, I.J., *Probability and the Weighting of Evidence*, Charles Griffin and Co.Ltd., London (1950).
7. ROBERTSON, S.E., 'The probability ranking principle in IR', *Journal of Documentation*, **33**, 294-304 (1977).
8. GOFFMAN, W., 'A searching procedure for information retrieval', *Information Storage and Retrieval*, **2**, 294-304 (1977).
9. WILLIAMS, J.H., 'Results of classifying documents with multiple discriminant functions', In : *Statistical Association Methods for Mechanized Documentation* (Edited by Stevens et al.) National Bureau of Standards, Washington, 217-224 (1965).
10. DE FINETTI, B., *Theory of Probability*, **Vol. 1**, 146-161, Wiley, London (1974).
11. KU, H.H. and KULLBACK, S., 'Approximating discrete probability distributions', *IEEE Transactions on Information Theory*, **IT-15**, 444-447 (1969).
12. KULLBACK, S., *Information Theory and Statistics*, Dover, New York (1968).
13. CHOW, C.K. and LIU, C.N., 'Approximating discrete probability distributions with dependence trees', *IEEE Transactions on Information Theory*, **IT-14**, 462-467 (1968).
14. COX, D.R., 'The analysis of multivariate binary data', *Applied Statistics*, **21**, 113-120 (1972).
15. BOX, G.E.P. and TIAO, G.C., *Bayesian Inference in Statistical Analysis*, 304-315, Addison-Wesley, Reading, Mass. (1973).
16. GOOD, I.J., *The Estimation of Probabilities: An Essay on Modern Bayesian Methods*, The M.I.T. Press, Cambridge, Mass. (1965).
17. HARPER, D. and van RIJSBERGEN, C.J., 'An evaluation of feedback in document retrieval using co-occurrence data', *Journal of Documentation*, **34**, 189-216 (1978).
18. HUGHES, G.F., 'On the mean accuracy of statistical pattern recognizers', *IEEE Transactions on Information Theory*, **IT-14**, 55-63 (1968).
19. MARON, M.E. and KUHN, J.L., 'On relevance, probabilistic indexing and information retrieval', *Journal of the ACM*, **7**, 216-244 (1960).
20. IVIE, E.L., 'Search Procedures Based on Measures of Relatedness Between Documents', Ph.D. Thesis, M.I.T., Report MAC-TR-29 (1966).
21. WHITNEY, V.K.M., 'Minimal spanning tree, Algorithm 422', *Communications of the ACM*, **15**, 273-274 (1972).

22. BENTLEY, J.L. and FRIEDMAN, J.H., *Fast Algorithm for Constructing Minimal Spanning Trees in Coordinate Spaces*, Stanford Report, STAN-CS-75-529 (1975).
23. CROFT, W.B., 'Clustering large files of documents using single link', *Journal of the American Society for Information Science*, **28**, 341-344 (1977).
24. SALTON, G., *Dynamic Information and Library Processing*, Prentice-Hall, Englewoods Cliffs, NJ., 441-445 (1975).
25. JARDINE, N. and SIBSON, R., *Mathematical Taxonomy*, pp. 12-15, Wiley, London and New York (1971).
26. KENDALL, M.G. and STUART, A., *Advanced Theory of Statistics*, **Vol. 2**, 2nd ed., Griffin, London (1967).
27. LINFOOT, E.H., 'An informational measure of correlation', *Information and Control*, **1**, 85-89 (1957).
28. YU, C.T. and SALTON, G., "Precision Weighting - An effective automatic indexing method", *Journal of the ACM*, **23**, 76-85 (1976).
29. YU, C.T., LUK, W.S. and CHEUNG, T.Y., 'A statistical model for relevance feedback in information retrieval', *Journal of the ACM*, **23**, 273-286 (1976).
30. NEGOITA, C.V., 'On the decision process in information retrieval', *Studii si cercetari de documentare*, **15**, 269-281 (1973).
31. ROBERTSON, S.E., 'Theories and models in information retrieval', *Journal of Documentation*, **33**, 126-148 (1977).
32. DOYLE, L.B., *Information Retrieval and Processing*, Melville Publishing Co., Los Angeles, California (1975).
33. DOYLE, L.B., 'Programmed interpretation of text as a basis for information retrieval systems', In: *Proceedings of the Western Joint Computer Conference*, San Francisco, 60-63 (1959).
34. DOYLE, L.B., 'Semantic road maps for literature searchers', *Journal of the ACM*, **8**, 553-578 (1961).
35. DOYLE, L.B., 'Some compromises between word grouping and document grouping', In: *Statistical Association Methods for Mechanized Documentation* (Edited by Stevens *et al.*) National Bureau of Standards, Washington, 15-24 (1965).
36. STILES, H.F., 'The association factor in information retrieval', *Journal of the ACM*, **8**, 271-279 (1961).
37. GOODMAN, L. and KRUSKAL, W., 'Measures of association for cross-classifications', *Journal of the American Statistical Association*, **49**, 732-764 (1954).
38. EDMUNDSON, H.P. and WYLLYS, R.E., 'Automatic abstracting and indexing - Survey and recommendations', *Communications of the ACM*, **4**, 226-234 (1961).
39. YU, C.T. and SALTON, G., 'Effective information retrieval using term accuracy', *Communications of the ACM*, **20**, 135-142 (1977).

40. SPARCK JONES, K., *Automatic Keyword Classification for Information Retrieval*, Butterworths, London (1971).

