

למידה לא מפקחת - יורם לוזון, תשפ"ה

רגב יחזקאל אימרה

February 3, 2025

תוכן העניינים

3	I מנהלה.
3	II הקדמה.
5	III אלגוריתמי אשכול ב- \mathbb{R}^n .
5	1 אלגוריתם Expectation maximization.
8	2 אלגוריתם K-Means.
9	3 גרסה רכה של K-Means: אלגוריתם Fuzzy C Means.
9	IV מטריקות.
11	4 אלגוריתם אשכול ספקטרלי:
11	5 אלגוריתם Prim.
11	6 אשכול היררכי.
12	V אלגוריתמי Community detection.
12	7 אלגוריתם Newman Girvan.
12	8 אלגוריתם Louvain.
13	9 אלגוריתם Leiden.
13	10 הערכת צפיפות- היסטוגרמה, parzen windows ו-KDE.
14	11 אלגוריתם DBSCAN.
14	VI הורדת מימדים.
14	12 חזרה על תורת האינפורמציה.
14	12.1 מדד לאי וודאות.
15	13 בעיית PCA.
16	14 אלגוריתם ICA.
16	15 אלגוריתם C.M.D.S.

17	16 אלגוריתם ISOMAP.
17	17 אלגוריתם LLE.
17	18 אלגוריתם EIGENMAPS.
17	19 הולכים מקריים, צוואר בקבוק לאינפורמציה ואלגוריתם T-SNE.
18	20 אלגוריתם U-MAP.
18	21 מקודד אוטומטי.
18	VII למידה מפוקחת עצמאית.
19	VIII זיהוי חריגים.
20	IX נתונים דינאמיים.
20	22 מודלים מרקובים סמויים.
22	23 טיפול במשתנים קטגוריאליים- אלגוריתם MCA.
22	24 בונוס.

חלק I

מנהלה.

מרצה: יורם לוזון.

מייל מרצה: louzouy@math.biu.ac.il

אין מבחן בקורס, יש עבודה שמהותה היא שאי אפשר לייצר אותה דרך ChatGPT. העבודה בזוגות, תהיה כתובה ב-Latex ובאנגלית.

חלק II

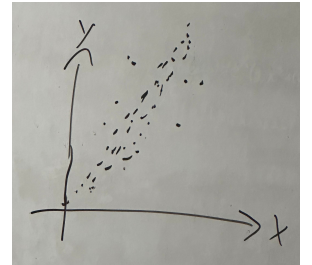
הקדמה.

בעיה: יש לנו אוסף תצפיות \bar{x}_i ולכל תצפית יש לנו תווית y_i וכל זוג כאלו מגיעים ממודל $\mathbb{P}(\{x_i, y_i\} | \theta)$ עם פרמטרים לא ידועים- θ .
בעיית סיווג לא ממוקדת זה בעיה שלא אכפת לי מ- y_i וכל מה שאני רוצה להעריך זה $\mathbb{P}(x_i | \hat{\theta})$.

בהסקה סטטיסטית המדד שלנו היה כמה $\hat{\theta}$ קרוב ל- θ .

כאן נניח שאני לא יודע את $\hat{\theta}$ אבל נניח שאני יודע להעריך את $\mathbb{P}(x_i | \hat{\theta})$. ככל שאני מעריך את $\mathbb{P}(x_i | \hat{\theta})$ יותר נכון ככה אני יותר שמח.

דוגמה 0.1. יש לי נקודות



וכעת יש לי x חדש ואני לא יודע מה הוא y , ואני רוצה לדעת מה ה- y שמתאים ל- x שלי.

אם אני רוצה למצוא קשר בין x ל- y בצורה של $y = f(x)$ זה יהיה בעייתי כי זה יהיה מאוד מסובך למצוא כזה. לכן אני אתקן ואומר כי $y = f(x) + \varepsilon$ כאשר ε הוא רעש שאני יכול להעריך את ההתפלגות שלו אבל לא אותו.

לכן נכתוב $y = f(x) + N(0, \sigma)$ כאשר אני לא יודע גם את σ . כעת, עוד בעיה היא שיש הרבה הרבה פונקציות בעולם ואני לא יכול לעבור על כולן ולראות מי הכי טובה לי, לכן אני רוצה לבחור את $f(x)$ להיות ממשפחת פונקציות מסויימת.

הנחה 1: נניח יש לי בסיס $\Phi_i(x)$ של פונקציות $i \in \{1, \dots, k\}$ וכן $f(x) = \sum_{i=1}^k w_i \Phi_i(x)$. איך אני יודע מהם ה- w_i שאני צריך?

הנחה 2: ערך ה- w שאני רוצה הוא ה- w שעבורו ההסתברות של הנתונים $\{(x_1, y_1), (x_2, y_2), \dots, (x_k, y_k)\}$ הינה מירבית, כלומר

$$L_{\{(x_i, y_i)\}}(w) = \mathbb{P}(\{(x_i, y_i)\} | w)$$

מקבל מקסימום.

הנחה 3: הדגימות שלי בלתי תלויות ומאותה התפלגות (i.i.d.).

לפיכך

$$L = \mathbb{P}(\{(x, y)\} | w) = \prod_{i=1}^k \mathbb{P}(\{(x_i, y_i)\} | w)$$

נפתור בעזרת נראות מירבית:

$$\begin{aligned} \log(L) &= \sum_{i=1}^k \log(\mathbb{P}(x_i, y_i | w)) = \sum_{i=1}^k \log(\mathbb{P}(y_i | x_i, w)) = \sum_{i=1}^k \log(\mathbb{P}(y_i | x_i, w) \cdot \mathbb{P}(x_i | w)) \\ &= \sum_{i=1}^k \log(\mathbb{P}(y_i | x_i, w)) + \underbrace{\log(\mathbb{P}(x_i))}_{\substack{\text{לא מעניין אותנו} \\ \text{כי גוזרים לפי } w}} = \sum_{i=1}^k \underbrace{\log\left(\frac{1}{\sqrt{2\pi}\sigma}\right)}_{\substack{\text{לא מעניין אותנו} \\ \text{כי גוזרים לפי } w}} - \frac{(y_i - f(x_i))^2}{2\sigma^2} \end{aligned}$$

לכן ללא קבועים סה"כ יש לי

$$\log(L) = - \sum_{i=1}^k (y_i - f(x_i))^2$$

אבל במקום למקסם את $\log(L)$ אני יכול למזער את $-\log(L)$ לכן נגדיר את ה-LOSS (הפסד) להיות $\text{Loss} = \sum_{i=1}^k (y_i - f(x_i))^2$ ואני רוצה למזער את ההפסד שלי (לגיטימי סה"כ).

הבעיה היא שלפעמים f היא מאוד קשה לחישוב או שאני לא יודע בכלל את w_i .

בעצם, אני לא רוצה למקסם את $\mathbb{P}(\{(x_i, y_i)\} | w)$, אלא את $\mathbb{P}(w | \{(x_i, y_i)\})$. למזלי הם קשורים אחד לשני לפי חוק בייס:

$$\mathbb{P}(w | \{(x_i, y_i)\}) = \frac{\mathbb{P}(\{(x_i, y_i)\} | w) \cdot \mathbb{P}_0(w)}{\mathbb{P}(\{(x_i, y_i)\})}$$

אבל אני רוצה למקסם את $\log(\mathbb{P}(w | \{(x_i, y_i)\}))$ ולכן אני ממקסם את

$$\log(\mathbb{P}(w | \{(x_i, y_i)\})) = \log(\mathbb{P}(\{(x_i, y_i)\} | w)) + \log(\mathbb{P}_0(w)) - \log(\mathbb{P}(\{(x_i, y_i)\}))$$

אבל נזכור שכשאני ממקסם את $\log(\mathbb{P}(w | \{(x_i, y_i)\}))$ אני גוזר לפי w ולכן אפשר להתעלם מ- $\log(\mathbb{P}(\{(x_i, y_i)\}))$ (כי זה קבוע ב- w ולכן כשאני אגזור את זה אני אקבל 0).

אפשר להגיד כי $\mathbb{P}_0(w) = \frac{1}{\sqrt{2\pi\beta}} e^{-\frac{\|w\|^2}{2\beta^2}}$ כי אני לא יודע את w , כי הכי נורמלי לנחש שהוא מתפלג נורמלי \odot .

ניקח לזה \log ואחרי ניקוי קבועים נקבל

$$\log(\mathbb{P}_0(w)) = -\frac{\|w\|^2}{2\beta^2}$$

ניזכר גם כי $\log(\mathbb{P}(\{(x_i, y_i)\} | w))$ הוא הנראות המירבית שלי ממקודם.

בסופו של יום אני רוצה לעשות

$$\max \left(-\frac{\|w\|^2}{2\beta^2} - \sum_{i=1}^k (y_i - f(x_i))^2 \right)$$

או באופן שקול

$$\min \left(\frac{\|w\|^2}{2\beta^2} + \sum_{i=1}^k (y_i - f(x_i))^2 \right)$$

ונקרא ל- $\frac{\|w\|^2}{2\beta^2}$ איבר רגולריזציה. הוא ביטוי מתמטי של ה- \mathbb{P}_0 .

סיכום מושגים:

מודל: $y = f(x) + N(0, \sigma)$

רעש: נורמלי.

נראות (LOSS): מסיקים מהמודל ומהרעש.

\mathbb{P}_0 : מסיקים מהנראות.

רגולריזציה: מסיקים מ- \mathbb{P}_0 .

דוגמה 0.2. אם אני בכיתה ומוודד את הגבהים של כל התלמידים ונכנס לי סטודנט חדש ואני רוצה לנחש את הגובה שלו, אני אנחש לפי הנראות המירבית, זה קלי קלות. אבל אם אני רוצה להפריד בין בנים ובנות הפעם אני יודע שההתפלגות שלי צריכה להיות מורכבת מ2 משתנים מקריים נורמלים שהיא $N(\mu_1, \sigma_1), N(\mu_2, \sigma_2)$. למודל כזה קוראים "מודל תערובת גאוסיאנים", או פשוט $G.M.M$. אם יש לי דגימה חדשה x אזי

$$\mathbb{P}(x) = \pi_1 N(x | \mu_1, \sigma_1) + \pi_2 N(x | \mu_2, \sigma_2)$$

כאשר π_i זה ההסתברות להיות באשכול ה- i .

$$\log(\mathbb{P}(x)) = \log \left(\pi_1 \cdot \frac{1}{\sqrt{2\pi\sigma_1}} e^{-\frac{(x-\mu_1)^2}{2\sigma_1^2}} + \pi_2 \cdot \frac{1}{\sqrt{2\pi\sigma_2}} e^{-\frac{(x-\mu_2)^2}{2\sigma_2^2}} \right)$$

איזה באסה, יש לי + בתוך ה- \log ובגזירה זה יוצא מזעזע. לגזור ולהשוות ל-0 זה מאוד לא פתיר, לכן חייבת להיות דרך טובה יותר לחשב את זה. הבעיה היא שהכנסנו משתנים סמויים וזה מה שהכניס לנו את הסימן + ל- \log .

$$\frac{\partial}{\partial \pi_1} \log(\mathbb{P}(x)) = \sum_{x \in \text{דגימות}} \frac{\frac{1}{\sqrt{2\pi\sigma_1}} e^{-\frac{(x-\mu_1)^2}{2\sigma_1^2}}}{\pi_1 \cdot \frac{1}{\sqrt{2\pi\sigma_1}} e^{-\frac{(x-\mu_1)^2}{2\sigma_1^2}} + \pi_2 \cdot \frac{1}{\sqrt{2\pi\sigma_2}} e^{-\frac{(x-\mu_2)^2}{2\sigma_2^2}}} = 0$$

בשביל זה הומצא אלגוריתם טוב שעוזר לנו מאוד.

אלגוריתמי אשכול ב- \mathbb{R}^n .

1 אלגוריתם Expectation maximization.

הרעיון: אם רק הייתי יודע את $\mu_1, \sigma_1, \pi_1, \mu_2, \sigma_2, \pi_2$, הייתי יכול בקלות להעריך את הסיכוי שכל דגימה שייכת לכל אשכול, ואז הייתי יכול לאמוד מחדש את $\mu_1, \sigma_1, \pi_1, \mu_2, \sigma_2, \pi_2$ וחוזר חלילה.

רקע: יש לי אוסף של j מטבעות, לכל מטבע יש סיכוי q_j להיות 1. בכל פעם אני לוקח מטבע מסויים ומטיל אותו n פעמים. לכן הסיכוי לתצפית k_i - מספר הפעמים שקיבלתי עץ בניסוי ה- i היא

$$\mathbb{P}(k_i|j) = \binom{n}{k_i} q_j^{k_i} (1 - q_j)^{n-k_i}$$

לכן

$$\mathbb{P}(k_i) = \sum_j \pi_j \mathbb{P}(k_i|j)$$

אני יודע כך הזמן למה k שווה, בלי לדעת מהו j . זה ניסוי אחד מתוך הרבה הרבה ניסויים שעשיתי. לכן הנראות שלי תהיה

$$L = \prod_i \left(\sum_j \pi_j \mathbb{P}(k_i|j) \right)$$

אני עכשיו עשיתי ניסוי ואני רוצה לדעת מהם π_j וכן q_j . נגדיר

$$\log(L) = \log \left(\prod_i \left(\sum_j \pi_j \mathbb{P}(k_i|j) \right) \right) = \sum_i \log \left(\sum_j \pi_j \mathbb{P}(k_i|j) \right)$$

וכמו בשיטת נראות מקסימלית אדרוש

$$\frac{\partial}{\partial q_i} \log(L), \frac{\partial}{\partial q_j} \log(L) = 0$$

הבעיה היא שהנגזרת יוצאת מאוד מסובכת

$$\frac{\partial}{\partial \pi_j} \log(L) = \sum_i \frac{1}{\sum_{j'} \pi_{j'} \mathbb{P}(k_i|j')} \cdot \pi_j \frac{\partial}{\partial q_j} \mathbb{P}(k_i|j) = 0$$

וכן

$$\begin{aligned} \frac{\partial}{\partial q_j} \mathbb{P}(k_i|j) &= \frac{k_i}{q_j} \mathbb{P}(k_i|j) - \frac{n - k_i}{(1 - q_j)} \mathbb{P}(k_i|j) \\ &= \underbrace{\left(\frac{k_i}{q_j} - \frac{n - k_i}{(1 - q_j)} \right)}_{:= f(k_i, q_j)} \mathbb{P}(k_i|j) \end{aligned}$$

לכן

$$\frac{\partial}{\partial q_i} \log(L) = \sum_i \frac{\pi_j \mathbb{P}(k_i|j) f(k_i, q_j)}{\sum_{j'} \pi_{j'} \mathbb{P}(k_i|j')} = 0$$

וזה רק המשוואה הראשונה.

מהמשוואה השנייה אני מקבל גם אילוץ: $\sum_j \pi_j = 1$ כלומר נפתור

$$\frac{\partial}{\partial \pi_j} \left(\sum_i \log \left(\sum_{j'} \pi_{j'} \mathbb{P}(k_i|j') \right) + \lambda \left(\sum_j \pi_j - 1 \right) \right) = 0$$

ונקבל

$$\sum_i \frac{\mathbb{P}(k_i|j)}{\sum_{j'} \pi_{j'} \mathbb{P}(k_i|j')} + \lambda = 0$$

ומכאן

$$(*) \sum_i \frac{\sum_j \pi_j \mathbb{P}(k_i|j)}{\sum_{j'} \pi_{j'} \mathbb{P}(k_i|j')} + \underbrace{\sum_j \pi_j \lambda}_{=1} = 0$$

כלומר

$$\lambda = -n$$

כעת מהמשוואה (*) נקבל

$$\sum_i \underbrace{\frac{\mathbb{P}(k_i|j)}{\sum_{j'} \pi_{j'} \mathbb{P}(k_i|j')}}_{:=\gamma_{ij}} = \pi_j n$$

לכן עבור π_j קיבלנו $\sum_i \gamma_{ij} = \pi_j n$ ועבור q_i נקבל $\sum_i \gamma_{ij} f(k_i, q_j) = 0$

מהו γ_{ij} ? זה המשקל היחסי של המטבע ה- j בניסוי i .

זה אומר, שאם הייתי מגדיר אינדקטור z_{ij} האירוע שבחרתי מטבע j בניסוי i אזי נקבל

$$\gamma_{ij} = \mathbb{P}(z_{ij}) = \mathbb{E}[z_{ij}]$$

אם רק הייתי יודע את π_j, q_j הייתי יודע גם את γ_{ij} לכל i ולכל j . אם גם הייתי יודע את γ_{ij} הייתי יודע את π_j, q_j .

לכן אני אנהש את π_j, q_j אחשב את γ_{ij} , מזה אני אחשב את π_j, q_j החדשים ומשם γ_{ij} וכו'.

בעיה כללית: בהינתן אוסף של משתנים סמויים z_{ij} כאשר ניסוי i מתבצע עם פרמטרים θ_j ויש לכל מצב j סיכוי π_j , ורוצים לאמוד אומדן נראות מירבית של θ בהינתן סדרת ניסויים x_1, \dots, x_n .

נכתוב:

$$\mathbb{P}(x|\theta) = \sum_z \mathbb{P}(x, z|\theta)$$

ולכן

$$\log(\mathbb{P}(x|\theta)) = \log\left(\sum_z \mathbb{P}(x, z|\theta)\right)$$

בעיה! יש לי + בתוך ה-log. אני אגדיר בשביל זה התפלגות חדשה $\mathbb{Q}(z)$ המקיימת $\mathbb{Q}(z) > 0$ וכן $\sum_z \mathbb{Q}(z) = 1$ ונקבל

$$\log(\mathbb{P}(x|\theta)) = \log\left(\sum_z \frac{\mathbb{P}(x, z|\theta)}{\mathbb{Q}(z)} q(z)\right) \geq \sum_z \mathbb{Q}(z) \log\left(\frac{\mathbb{P}(x, z|\theta)}{\mathbb{Q}(z)}\right)$$

* לפי אי שוויון ינסון.

נחשב את

$$\begin{aligned} \Delta &= \sum_z \mathbb{Q}(z) \log(\mathbb{P}(x|\theta)) - \sum_z \mathbb{Q}(z) \log\left(\frac{\mathbb{P}(x, z|\theta)}{\mathbb{Q}(z)}\right) \\ &= \sum_z \mathbb{Q}(z) \log\left(\frac{\mathbb{P}(x|\theta) q(z)}{\mathbb{P}(x, z|\theta)}\right) \\ &= \sum_z \mathbb{Q}(z) \log\left(\frac{\mathbb{Q}(z)}{\mathbb{P}(z|\theta, x)}\right) \end{aligned}$$

לכן בהינתן שתי התפלגויות \mathbb{Q}, \mathbb{P} נגדיר

$$D_{KL}(\mathbb{Q}, \mathbb{P}) = \sum_z \mathbb{Q}(z) \log\left(\frac{\mathbb{Q}(z)}{\mathbb{P}(z)}\right)$$

ונוכל לכתוב

$$\Delta = D_{KL}(\mathbb{Q}, \mathbb{P}(z|x, \theta))$$

ונקבל

$$\log(\mathbb{P}(x|\theta)) = \sum_z \mathbb{Q}(z) \log\left(\frac{\mathbb{P}(x, z|\theta)}{\mathbb{Q}(z)}\right) + D_{KL}(\mathbb{Q}, \mathbb{P}(z|x, \theta))$$

ואף נסמן $L := \sum_z \mathbb{Q}(z) \log \left(\frac{\mathbb{P}(x, z | \theta)}{\mathbb{Q}(z)} \right)$. כעת אני רוצה להעריך את $\log(\mathbb{P}(x | \theta))$ אבל זה מסובך לי. למזלי, להעריך את

$$\sum_z \mathbb{Q}(z) \log \left(\frac{\mathbb{P}(x, z | \theta)}{\mathbb{Q}(z)} \right) + D_{KL}(\mathbb{Q}, \mathbb{P}(z | x, \theta))$$

זה קל לי.

נאפס את $D_{KL}(\mathbb{Q}, \mathbb{P}(z | x, \theta))$, נרצה $\mathbb{Q}(z) = \mathbb{P}(z | x, \theta)$.

שלב א': בהינתן $\hat{\theta}_1$, נבחר $\mathbb{Q}(z) = \mathbb{P}(z | x, \hat{\theta}_1)$.

שלב ב': נכתוב $\log(\mathbb{P}(x | \theta)) = \sum_z \mathbb{P}(z | x, \hat{\theta}_1) \log \left(\frac{\mathbb{P}(x, z | \theta)}{\mathbb{P}(z | x, \hat{\theta}_1)} \right)$ ואחרי ניקוי קבועים וגזירה לפי θ נקבל אומד חדש ל- θ , נסמנו ב- $\hat{\theta}_2$.

שלב ג': נחזור על שלב א' עם $\hat{\theta}_2$.

דוגמה 1.1. נבחר :

$$\mathbb{P}(x, z) = \left(\prod_i \prod_j \pi_j \binom{n}{k_i} q_j^{k_i} (1 - q_j)^{n - k_i} \right)^{z_{ij}}$$

נכניס משתנים סמויים: $z_{ij} = 1$ אם בניסוי ה- i בחרתי מטבע j ו- $z_{ij} = 0$ אם בניסוי ה- i לא בחרתי במטבע j .

$$\log(\mathbb{P}(x, z)) = \sum_i \sum_j z_{ij} \log \left(\pi_j q_j^{k_i} (1 - q_j)^{n - k_i} \right)$$

התעלמנו מ- $\binom{k_i}{n}$ כי הוא קבוע.

$$\mathbb{E}_z [\log(\mathbb{P}(x, z))] = \sum_i \sum_j \gamma_{ij} \log \left(\pi_j q_j^{k_i} (1 - q_j)^{n - k_i} \right)$$

נגזור לפי q_j :

$$\frac{\partial}{\partial q_j} \mathbb{E}_z [\log(\mathbb{P}(x, z))] = \sum_i \gamma_{ij} \left(\frac{k_i}{q_j} - \frac{n - k_i}{1 - q_j} \right) = 0$$

נקבל

$$\sum_i \gamma_{ij} \frac{k_i}{q_j} - \sum_i \gamma_{ij} \frac{n - k_i}{1 - q_j} = 0$$

ובסופו של יום נקבל

$$q_j = \frac{\sum_i \gamma_{ij} k_i}{\sum_i \gamma_{ij} n}$$

וכן

$$\pi_j = \sum_i \frac{\gamma_{ij}}{n}$$

כאשר

$$\gamma_{ij} = \frac{\prod_j \mathbb{P}(k_i | j)}{\sum_{j'} \pi_{j'} \mathbb{P}(k_i | j')}$$

דוגמה 1.2. עבור תערובת גאוסיאנים (Gaussian Mixture Model-G.M.M.) :

תזכורת: עבור חד מימד $\mathbb{P}(x | \theta) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$ וברב מימד $\mathbb{P}(\bar{x} | \theta) = \frac{1}{\sqrt{(2\pi)^D \det(\Sigma)}} e^{-\frac{(\bar{x}-\bar{\mu})^T \Sigma^{-1} (\bar{x}-\bar{\mu})}{2}}$ (זה בהנחה שיש לי משתנה גאוסיאני אחד).

נחשב \log לנראות :

$$\log(\mathbb{P}(\bar{x}_1, \dots, \bar{x}_n)) = \sum_{i=1}^n \frac{1}{\sqrt{(2\pi)^D \det(\Sigma)}} - \sum_{i=1}^n \frac{(\bar{x}_i - \bar{\mu})^T \Sigma^{-1} (\bar{x}_i - \bar{\mu})}{2}$$

נחשב את הגראדיאנט ונשווה ל-0 :

$$\sum_{i=1}^n \Sigma^{-1} (\bar{x}_i - \bar{\mu}) = 0$$

כלומר

$$\sum_{i=1}^n (\bar{x}_i - \bar{\mu}) = 0$$

ולכן

$$\hat{\mu} = \sum_{i=1}^n \frac{\bar{x}_i}{n}$$

ועם קצת אלגברה נגבל כי

$$\hat{\Sigma} = \sum_{i=1}^n \frac{(\bar{x}_i - \mu)(\bar{x}_i - \mu)^T}{n}$$

עכשיו החיים שלנו מסתבכים כי אנחנו רוצים לאמוד כמות מסויימת גדולה מ-1 של גאוסיאנים.

הבעיה:

$$\sum_{i=1}^n \sum_{j=1}^n \gamma_{ij} \log(\mathbb{P}(x_i | \theta_j)) = \sum_{i=1}^n \sum_{j=1}^n \gamma_{ij} \log \left(\pi_j c_j e^{-\frac{(x_i - \mu_j)^T \Sigma_j^{-1} (x_i - \mu_j)}{2}} \right)$$

כאשר

$$c_j = \frac{1}{\sqrt{(2\pi)^D \det(\Sigma_j)}}$$

לפי שלבי EM, נניח את $\pi_j, \bar{\mu}, \bar{\Sigma}_j$ וממנו נחשב את γ_{ij} ומוזה נשפר את $\pi_j, \bar{\mu}, \bar{\Sigma}_j$ ונשפר את γ_{ij}

$$\tilde{\gamma}_{ij} = \pi_j c_j e^{-\frac{(x_i - \bar{\mu}_j)^T \bar{\Sigma}_j^{-1} (x_i - \bar{\mu}_j)}{2}}$$

ננרמל את $\tilde{\gamma}_{ij}$ בשביל $\sum_j \tilde{\gamma}_{ij} = 1$ כלומר

$$\gamma_{ij} = \frac{\tilde{\gamma}_{ij}}{\sum_j \tilde{\gamma}_{ij}}$$

נגזור וכן הלאה ונקבל

$$\hat{\mu}_j = \sum_i \frac{\gamma_{ij} \bar{x}_i}{\sum_j \gamma_{ij}}$$

וגם

$$\hat{\Sigma}_j = \sum_{i=1}^n \frac{\gamma_{ij} (\bar{x}_i - \hat{\mu})(\bar{x}_i - \hat{\mu})^T}{\sum_i \gamma_{ij}}$$

וכן

$$\hat{\pi}_j = \sum_{i=1}^n \frac{\gamma_{ij}}{n}$$

2 אלגוריתם K-Means.

אשכול קשיח:

באשכול רך ראינו ב-G.M.M, רק נותנים סיכוי לכל אשכול.

באשכול קשיח- נחליף את השידוך להיות חד ערכי.

הנחות:

(א) נניח אשכול קשיח.

(ב) נניח שהשונויות בכל גאוסיאן שווה בכל מימד, ושווה בין הגאוסיאנים ואין שונות משותפת.

בהינתן ההנחות לעיל:

$$z_{ij} = \begin{cases} 0 \\ 1 \end{cases} \quad \text{א} \quad \gamma_{ij}$$

(ב) $z_{ij} = 1$ עבר המרכז הקרוב ביותר ל- x_i ו-0 אחרת.

לכן $\bar{\mu}$ הוא ממוצע של נקודות באשכול.

אלגוריתם 2.1: K-Means:

1. בחר k מרכזים באקראי.

2. שייך כל דגימה למרכז הקרוב אליה ביותר, וחשב מרכז מחדש.

3. אם חל שינוי, וכל האשכולות לא ריקים, חזור ל-2.

4. אם התרוקן אשכול, פצל אשכול ל-2 אשכולים באקראי, וחזור ל-2.

3 גרסה רכה של K-Means: אלגוריתם Fuzzy C Means.

אותן הנחות כמו K-Means אבל רך.

נגדיר את הטעות שלי להיות

$$LOSS = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \|\bar{x}_i - \bar{\mu}_j\|^2 \gamma_{ij}^m$$

כאשר

$$\sum_{j=1}^n \gamma_{ij} = 1, \gamma_{ij} \geq 0$$

אם $m \rightarrow 0$ אזי האשכול קשיח וכאשר $m \rightarrow \infty$ האשכול נהיה יותר ויותר קשיח.

נפתור באמצעות כופלי לגראנז' :

$$L = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \|\bar{x}_i - \bar{\mu}_j\|^2 \gamma_{ij}^m - \sum_{i=1}^n \lambda_i \left(\sum_{j=1}^n \gamma_{ij} - 1 \right)$$

נגזור ונקבל

$$\nabla \mu_j = \sum_{i=1}^n \|\bar{x}_i - \bar{\mu}_j\| \gamma_{ij}^m = 0 \Rightarrow \bar{\mu}_j = \frac{\sum_{i=1}^n \bar{x}_i \gamma_{ij}^m}{\sum_{i=1}^n \gamma_{ij}^m}$$

$$\frac{\partial L}{\partial \gamma_{ij}} = \frac{1}{2} \|\bar{x}_i - \bar{\mu}_j\|^2 m \gamma_{ij}^{m-1} - \lambda_i = 0 \Rightarrow \gamma_{ij} = \left(\frac{\lambda_i}{\frac{1}{2} \|\bar{x}_i - \bar{\mu}_j\|^2 m} \right)^{\frac{1}{m-1}}$$

$$\frac{\partial L}{\partial \lambda_i} = \sum_{j=1}^n \frac{\lambda_i^{\frac{1}{m-1}}}{\left(\frac{1}{2} \|\bar{x}_i - \bar{\mu}_j\|^2 m \right)^{\frac{1}{m-1}}} - 1 = 0 \Rightarrow \lambda_i = \left(\frac{1}{\sum_{j=1}^n \frac{1}{\left(\frac{1}{2} \|\bar{x}_i - \bar{\mu}_j\|^2 m \right)^{\frac{1}{m-1}}}} \right)^{m-1}$$

חלק IV

מטריקות.

כעת מה קורה אם יש לי רק מרחק בין עצמים d_{ij} ללא יכולת להגדיר הטלה ל- \mathbb{R}^n של העצמים? יש לי 2 פתרונות לבעיה הזאת:

1. למצוא לכל i ולעמוד \bar{x}_i ולהגדיר מטריקה $d(\bar{x}_i, \bar{x}_j) \approx d_{ij}$.

2. להגדיר גרף: התצפית ה- i תהפוך לקודקוד ה- i וכן $f(d_{ij})$ יהיה משקל הקשת דמיון ביניהם עבור f פונקציה מונוטונית יורדת. אם משתמשים בדמיון אז ניתן לזרוק קשתות עם ערך נמוך.

מטרה: אני רוצה לפרק את הגרף V ל-2 אשכולות (או יותר) A, B כך ש- $A \cup B = V$ וכן $A \cap B = \emptyset$.

האמונה שלי אומרת לי שאני רוצה $|A| \sim |B|$ וכן $\min \left(\sum_{\substack{ij \\ i \in A \\ j \in B}} w_{ij} \right)$ כלומר סכום משקלי הקשתות בין A ל- B יהיה מינימלי.

סימון: $z_i = \begin{cases} \frac{1}{2} & i \in A \\ -\frac{1}{2} & i \in B \end{cases}$ לכן לכל i, j באשכולות שונים $(z_i - z_j)^2 = 1$ וכן $(z_i - z_j)^2 = 0$ לכל i, j באותו אשכול. בנוסף לכל i מתקיים כי $z_i^2 = \frac{1}{4}$, לכן אני יכול עכשיו להגדיר פונקציית הפסד

$$LOSS = \sum_{i,j} w_{ij} (z_i - z_j)^2$$

טעות! אם אני אגדיר ככה את ההפסד אז אין בעיה לומר שפשוט כל הגרף שלי חלק מ- A וכן ש- B קבוצה ריקה וכביכול סיימתי אבל בעצם לא עשיתי פה כלום. לכן אני אגדיר את ההפסד שלי להיות

$$LOSS = \sum_{i,j} w_{ij} (z_i - z_j)^2 - \sum_i \underbrace{\left(\sum_j w_{ij} \right)}_{d_{ij}} (z_i - z_i)^2$$

כאשר יש לי אילוץ $\left| \sum_i z_i \right| < \varepsilon$.

דוגמה 3.1. יש לי מטריצת דרגות

$$W = \begin{pmatrix} 0 & 2 & 0 & 5 & 0 & 3 \\ 2 & 0 & 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 2 & 1 & 0 \\ 5 & 1 & 2 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 & 4 \\ 3 & 0 & 0 & 1 & 4 & 0 \end{pmatrix}$$

כלומר מחוץ לאלכסון $d_{ij} = 0$ וכן $d_{ii} = \sum_j w_{ij}$.

מטריצת הדרגות שלי תהיה

$$D = \begin{pmatrix} 10 & & & & & \\ & 4 & & & & \\ & & 4 & & & \\ & & & 9 & & \\ & & & & 5 & \\ & & & & & 8 \end{pmatrix}$$

כעת,

$$LOSS = \sum_{i,j} w_{ij} (z_i - z_j)^2 - \sum_{i,j} d_{ij} (z_i - z_j)^2$$

לכן אני מקבל השראה להגדיר את מטריצת הלפסיאן L_a להיות $L_a = W - D$.

במקרה שלנו נקבל

$$L_a = \begin{pmatrix} -10 & 2 & 0 & 5 & 0 & 3 \\ 2 & -4 & 1 & 1 & 0 & 0 \\ 0 & 1 & -4 & 2 & 1 & 0 \\ 5 & 1 & 2 & -9 & 0 & 1 \\ 0 & 0 & 1 & 0 & -5 & 4 \\ 3 & 0 & 0 & 1 & 4 & -8 \end{pmatrix}$$

וכן

$$LOSS = \sum_{i,j} L_{a_{ij}} (z_i - z_j)^2 \\ = -2\bar{z}^T L_a \bar{z}$$

כלומר אני רוצה לפתור

$$\begin{aligned} \max z^T L_a z \\ \text{s.t. } z_i \in \{-0.5, 0.5\} \\ \left| \sum_i z_i \right| < \varepsilon \end{aligned}$$

כיוון ש- z_i בדיד זה ממש קשה לפתור, אז נרשה ל- z_i להיות רציף אבל עדיין נבקש $\sum_i z_i^2 = \frac{n}{4}$ כאשר כמובן n זה מספר הדגימות שלי.

$$\text{נגדיר } \bar{v} = z \cdot \frac{2}{\sqrt{n}} \text{ ולכן } \|v\|^2 = 1 \text{ כלומר } v \text{ וקטור יחידה ולכן } \underbrace{\sum_i v_i^2}_{:=\beta} \leq \frac{2}{\sqrt{n}} \varepsilon$$

תכונות של מטריצת הלפסיאן:

1. ע"ע ממשיים- כי היא סימטרית.
2. ו"ע $(1, \dots, 1)^T$ כי סכום כל שורה הוא 0.

נפרק את v לפירוק לפי בסיס של וקטורים עצמיים של הלפסיאן:

$$v = \sum_i \alpha_i \bar{u}_i$$

כאשר \bar{u}_i ו"ע של L_a . לכן

$$L_a \cdot v = \sum_i \alpha_i \lambda_i \bar{u}_i$$

$$v^T L_a v = \sum_i \alpha_i^2 \lambda_i$$

כעת, העי"ע של L_a הם כולם אי חיוביים, כלומר אני ממקסם את זה

$$\begin{aligned} \max \sum_i \alpha_i^2 \lambda_i \\ \text{s.t. } \sum_i \alpha_i^2 &= 1 \\ \lambda_j &\leq \lambda_i = 0 \end{aligned}$$

לכן אני ארצה $\alpha_1 = 1, \alpha_{i \neq 1} = 0$ אבל להזכירנו יש לי אילוץ $\sum_i v_i^2 \leq \beta < 1$, לכן הפתרון שלי יהיה $\alpha_1 = \beta, \alpha_2 = \sqrt{1 - \beta^2}$ ואם חוזרים כל הדרך אחורה נקבל

$$z = \frac{\sqrt{n}}{2} \left(\beta \bar{u}_1 + \sqrt{1 - \beta^2} \bar{u}_2 \right)$$

לכן מהדוגמה הזאת בא לי הרעיון לאלגוריתם אשכול ספקטרלי:

4 אלגוריתם אשכול ספקטרלי:

אלגוריתם 4.1. צעדים לאשכול ספקטרלי:

$$1. \text{ חשב דרגה } d_{ij} = \sum_j w_{ij}$$

$$2. \text{ חשב } L_a = W - D \text{ כאשר } D = \begin{cases} d_{ij} & i = j \\ 0 & i \neq j \end{cases}$$

$$3. \text{ חשב וי"ע שני של } L_a : \bar{u}_2$$

$$\bar{z} = \frac{\sqrt{n}}{2} \left(\beta \bar{u}_1 + \sqrt{1 - \beta^2} \bar{u}_2 \right)$$

$$4. \text{ אם } z_i > 0 \text{ שייך ל-} A \text{ אחרת שייך ל-} B.$$

5 אלגוריתם Prim.

נניח יש לי המון המון נקודות וגרף ביניהם שבו קשת מסמנת קרבה בין 2 קודקודים.

האלגוריתם מוצא לי עץ"מ.

אלגוריתם 5.1. אלגוריתם Prim.

1. בכל שלב חבר לעץ הקיים את הקשת הכי זולה מחוץ לעץ.

2. חזור ל-1.

התוצאה של אלגוריתם Prim הוא עץ, לכן אם אני אסנן את ה- k קשתות הכי יקרות שלי אני אקבל k אשכולות כמו שרציתי באשכול ספקטרלי.

6 אשכול היררכי.

אלגוריתם 6.1. אלגוריתם לאשכול היררכי:

1. שים כל דגימה באשכול משלה.

2. חבר שני אשכולות קרובים ביותר לאשכול חדש וחשב את המרחק שלו מאשר האשכולות.

3. כל עוד לא כולם מחוברים חזור ל-2.

שאלה מעולה: איך אני מגדיר מרחק בין אשכולות?

תשובה מעולה: זה תלוי.

כל אחד יכול להגדיר בעצמו ולראות מה הכי טוב לו.

דוגמה 6.2. דוגמאות למרחקים בין אשכולות: $d(A, (B, C))$

1. מרחק single link:

$$\min(d(A, B), d(A, C))$$

2. מרחק full link:

$$\max(d(A, B), d(A, C))$$

3. מרחק average link :

$$\frac{d(A, B) + d(A, C)}{2}$$

4. מרחק UPGMA :

$$\frac{|B|}{|B| + |C|} d(A, B) + \frac{|C|}{|B| + |C|} d(A, C)$$

5. מרחק energy distance :

$$d(\{x_i\}, \{y_j\}) = 2 \sum_{i,j} d(x_i, y_j) - \sum_i d(x_i, x'_i) - \sum_j d(y_j, y'_j)$$

חלק V

אלגוריתמי Community detection.

השאלה שלנו : בהינתן רשת (גרף), איך מחלקים אותה לקהילות (אשכולות)?
 תזכורת: יש לנו A_{ij} -דמיון בין i ל- j . אם אין קשת בין i ל- j , בנוסף, $d_i = \sum_j A_{ij}$ הוא דרגת קודקוד.
 אם הגרף כיווני דרגת כניסה היא $\sum_j A_{ji}$ ודרגת יציאה היא $\sum_j A_{ij}$.

7 אלגוריתם Newman Girvan.

7.1. הגדרה. המרכזיות של קשת / קודקוד x הוא מספר המסלולים המינימליים שעוברים דרך הקשת / קודקוד x .
 נסמן על ידי $\text{Centrality}(x)$.

7.2. אלגוריתם Newman Girvan :

1. חשב מרכזיות לכל קשת.
2. הוצא קשתות עם המרכזיות הכי גבוהה לפי סדר המרכזיות.
3. חזור על 2 עד שהגעת למספר הקשתות הרצוי.

8 אלגוריתם Louvain.

נגדיר את הרווח להיות

$$Q = \frac{1}{2m} \sum_{i,j} \left(A_{ij} - \frac{k_i k_j}{2m} \right) C(i, j)$$

כאשר :

$$k_i \text{ דרגה של קודקוד } i. \quad \text{מספר הקשתות} = \frac{1}{2} \sum_{i,j} A_{ij} = m = \frac{1}{2} \sum_i k_i$$

$$C(i, j) = \begin{cases} 1 & i, j \text{ באותו אשכול} \\ 0 & \text{אחרת} \end{cases}$$

שאלה חשובה היא כמה קשתות יש בקבוצה מול כמה ציפית.

אם $A_{ij} - \frac{k_i k_j}{2m}$ גבוה, רוצים $C(i, j) = 1$, אחרת רוצים $C(i, j) = 0$. זו בעיה כי לכל קודקוד אני צריך לקבוע באיזה אשכול הוא ויש מלא אשכולים וזה מאוד יקר. לכן הפתרון הוא להפוך את $C(i, j)$ לרצף תוך כדי שנייצר אילוצים ונוכל לגזור אותו.

8.1. אלגוריתם Louvain.

1. נסדר קודקודים בסדר אקראי.
2. נחבר כל קודקוד לאשכול הקרוב אליו, אם החיבור מעלה את Q (פורמלית, בהינתן 2 קודקודים שכנים i, j חשב את

$$\Delta(i, j) = Q(j \text{ עובר לאשכול } i) - Q(i, j \text{ הנוכחי})$$

ואם $\Delta(i, j) > 0$ תעביר את i לאשכול j .

3. אם אין איך להתקדם, חבר כל אשכול לקודקוד אחד וצור קשת פנימית עם פעמיים סכום הקשתות בכל אשכול $w_{k,k} = \sum_{\substack{i \in C_k \\ j \in C_k}} A_{ij}$ וקשת

$$w_{k,l} = \sum_{\substack{i \in C_k \\ j \in C_l}} A_{ij}$$

4. חזור עד שאין שינוי.

9 אלגוריתם Leiden

$$Q = \frac{1}{2m} \sum_{i,j} \left(A_{ij} - \frac{k_i k_j}{2m} \right) C(i, j)$$

הבדלים בין לוביין לליידן:

1. מאפשר לפצל אשכולות.

2. מוסיף רעש.

סיכום:

אם יש לי גרף סביר שאני רוצה לחלק ל2 קבוצות- אשכול ספקטרלי.

גרף ענק- Leiden\ Louvain.

אם חשוב לי היררכיה של אשכולות- אשכול היררכי.

10 הערכת צפיפות- היסטוגרמה, parzen windows ו-KDE.

נניח יש לי סדרה של מספרים 1, 2, 3, 1, 4, 5, 1, 8, 2, 7, 1, 4, 3, 2, 7, 5, 8, 10, 9, 7, 4, 3, 8, 10

דרך אחת להעריך כמה מופעים יש לי מכל מספר היא באמצעות ספירה של ממש (ונירמול):

מספר	1	2	3	4	5	6	7	8	9	10
צפיפות	$\frac{4}{24}$	$\frac{3}{24}$	$\frac{3}{24}$	$\frac{3}{24}$	$\frac{2}{24}$	$\frac{0}{24}$	$\frac{4}{24}$	$\frac{3}{24}$	$\frac{1}{24}$	$\frac{2}{24}$

לפעולה זו קוראים היסטוגרמה.

אני אקח אינטרוולים מרווחים שווה בשווה (כמה שבא לי) בין המקסימום ועוד קצת לבין המינימום פחות קצת ובין כל מרווח (bin) אני אשים את כל המספרים שמתאימים לשם. כלומר במקרה הזה $\rho_i = \frac{n_i}{\Delta x}$. הבעיה היא שאם יש לי עוד מספר 545435453 אז יהיה לי היסטוגרמה דפוקה. לכן נגדיר $\rho = \frac{n_0}{\Delta x_i}$ כלומר כל מרווח (bin) להכיל n_0 ערכים. באופן טיפוסי $n_0 = \sqrt{n}$.

כעת אם אני בשני מימדים אני יכול לעשות מקביל לזה ויצא לי $\rho_i = \frac{n_i}{\Delta x \Delta y}$.

אם יש לי k בינים בכל מימד ב- l מימדים נקבל שיש לי k^l בינים וזה מאוד מאוד יקר. נניח שיש לי ∞ זיכרון אז ברוב התאים יהיו לי 0 איברים ובחלק מהתאים יהיה לי קצת מאוד מאוד איברים, כלומר ההיסטוגרמה שלי תהיה מאוד שטוחה.

פתרון אלטרנטיבי: במקום לעשות קוביות נעשה כדור $B(x, \epsilon) = \{z : \|x - z\| \leq \epsilon\}$. אני אסמן $V =$ נפח כדור ב- l מימדים עם רדיוס ϵ ו- $n(x) =$ כמות הנקודות ברדיוס סביב x . לכן נגדיר $\rho(x) = \frac{n(x)}{V}$. הבעיה היא שיכול להיות שאני אדגום כדור ריק ואפספס דגימות. אני רוצה איכשהו להציל את התכונה של צפיפות בלי להיות תלוי במזל שלי באיפה אני דוגם.

פתרון: אם יש לי ϵ רדיוס בתוכו יש k נקודות סביב x אזי $\rho(x) = \frac{k}{V_\epsilon}$ כאשר V_ϵ נפח כדור יחידה עם רדיוס ϵ . באופן טיפוסי נבחר $k = \sqrt{n}$.

נשאלת השאלה איך אפשר להפוך את זה לרציף? נסמן $\rho(\bar{x}) = \sum_i \frac{\mathbb{I}(\|\bar{y}_i - \bar{x}\| \leq \epsilon)}{V_\epsilon}$. במקום אינדיקטור ניקח פונקציה יותר יפה שגם יהיה לי קל לגזור f ונקבל

$$\rho(\bar{x}) = \sum_i \frac{f\left(\frac{\|\bar{y}_i - \bar{x}\|}{\epsilon}\right)}{V_\epsilon}$$

נבחר את $f(x) = \frac{1}{(\sqrt{2\pi\epsilon})^l} e^{-\frac{\|\bar{y}_i - \bar{x}\|^2}{2\epsilon^2}}$. זה יוצא נורא יפה, לכן נגדיר $KDE(x) = \sum_i \frac{\mathbb{P}(y_i|x)}{V}$ סיכוי ל- y_i תחת הנחת להיות צפיפות נורמלית סביב x ברדיוס ϵ . לכאורה העלות של זה היא $O(|V|)$ אבל אפשרי גם ב- $O(\log |V|)$.

איך מזה אני בונה מזה אשכולות?

(א) נקודות בצפיפות גבוהה שייכות לאשכול.

(ב) נקודות בצפיפות נמוכה הן רעש

מבחינתנו אשכול חייב להיות מספיק צפוף בצורה רציפה.

11 אלגוריתם DBSCAN.

11.1 אלגוריתם DBSCAN.

- (א) בהינתן ε נגדיר צפיפות של נקודה כתור מסתר הנקודות האחרות שהו ברדיוס ε ממנה.
 (ב) נגדיר $\rho(x) = \frac{k}{V_\varepsilon}$ כאשר k הוא מספר הנקודות במרחק ε מ- x וכן V_ε נפח כדור l מימדי ברדיוס ε .
 (ג) נקודת גרעין היא נקודה עם צפיפות $\rho(x) \geq \rho_0$ וכן $k(x) \geq k_0$. באופן טיפוסי נבחר $k_0 = 5, \varepsilon = 0.1$.
 (ד) נגדיר אשכול כתור אוסף של נקודות גרעין שמחוברות אחת לשניה ברדיוס ε וכל הנקודות שסביבן ברדיוס ε .
 (ה) כל הנקודות שלא חלק מהאשכול הן רעש.

חלק VI

הורדת מימדים.

נניח אני מסתכל על רשימת ציונים \bar{x} מ-48 קורסים. אין 48 תכונות שקובעות את הציון, יש בערך 5 תכונות- \bar{y} . נאמר $\bar{x} \sim f(\bar{y})$ אבל

$$y \in \mathbb{R}^L, x \in \mathbb{R}^K, K \gg L$$

אני לא מכיר את f ולא את y , אבל אני רוצה מודל שבהינתן \bar{x} יביא לי \bar{y} . מה יכול לאפיין לי את y ?
 יש לנו מלא נתונים x_1, \dots, x_n שמגיעים האיזושהו $x_i = f(y_i)$. המטרה שלנו:

(א) מהו f ?

(ב) מהו y_i ?

אנחנו רוצים להוריד את המימד של הנתונים שלנו תוך כדי שאנחנו שומרים על כמה שיותר אינפורמציה. איך נעשה את זה?

(א) רוצים למזער טעות: $\tilde{x}_i = f(y_i)$ שיביא לנו $\|\tilde{x}_i - x_i\|^2$ מינימלית PCA.

(ב) רוצים לשמר אינפורמציה ICA.

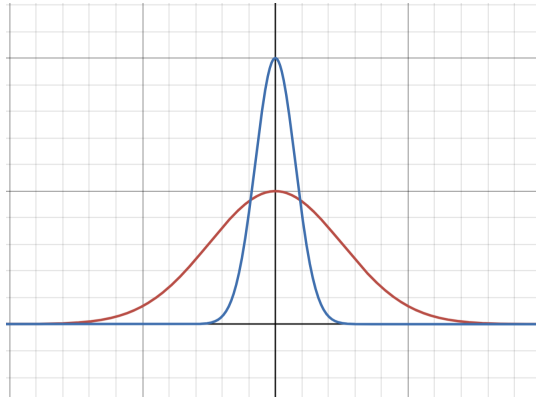
(ג) מרחק M.D.SISOMAP $d(y_i, y_j) \sim d(x_i, x_j)$.

(ד) צפיפות: מסלול של הולך אקראי TSNE\UMAP.

12 חזרה על תורת האינפורמציה.

12.1 מדד לאי וודאות.

אם יש לי התפלגות ρ אזי ככל שהסטיית תקן שלי יותר גבוהה ככה יש לי אי וודאות יותר גבוהה.



לדוגמה, יש פה 2 התפלגויות, האדומה עם סטיית תקן יותר גבוהה ולכן יותר מפוזרת.

נפרמל את הרעיון הזה:

הגדרה 12.1. בהינתן התפלגות ρ נגדיר את האנטרופיה של ρ להיות

$$\mathbb{H}(\rho) = -\sum_i \mathbb{P}_i \log(\mathbb{P}_i) = \int_{\mathbb{R}} \mathbb{P}(x) \log(\mathbb{P}(x)) dx$$

הגדרה 12.2. בהינתן שני משתנים מקריים X, Y נגדיר את האינפורמציה המשותפת שלהם להיות

$$\mathbb{H}(X, Y) = -\sum_{x,y} \mathbb{P}(x, y) \log(\mathbb{P}(x, y))$$

הגדרה 12.3. בהינתן משתנה מקרי X ומשתנה מקרה Y אני רוצה לדעת כמה הידיעה של Y חידש לי על הידיעה על X . לכן נגדיר את **האינפורמציה המותנית להיות**

$$\begin{aligned} MI(X, Y) &= \mathbb{H}(X) - \mathbb{H}(X|Y) \\ &= \mathbb{H}(Y) + \mathbb{H}(X) - \mathbb{H}(X, Y) \\ &= MI(Y, X) \end{aligned}$$

13 בעיית PCA

נניח שיש לי וקטור $\bar{x}_i = \sum_{j=1}^m \alpha_{i,j} \bar{U}_j$ כאשר \bar{U}_i בסיס אורתונורמלי. אזי $\alpha_{ij} = \langle \bar{x}_i, \bar{U}_j \rangle$ ואני רוצה להחליט איזה $\alpha_{i,j}$ אני שומר ואיזה אני זורק. יהיה לי $\tilde{x}_i = \sum_{j=1}^k \alpha_{i,j} \bar{U}_j$ נגדיר

$$\begin{aligned} L &= \sum_i \|x_i - \tilde{x}_i\|^2 \\ &= \sum_i \left\| \sum_{j=k+1}^m \alpha_{i,j} \bar{U}_j \right\|^2 \\ &\vdots \\ &= \sum_i \sum_{j=k+1}^m \alpha_{i,j}^2 \\ &= \sum_i \sum_{j=k+1}^m (\bar{x}_i^T \bar{U}_j)^T (\bar{x}_i^T \bar{U}_j) \\ &= \sum_{j=k+1}^m \bar{U}_j^T \left(\underbrace{\sum_i \bar{x}_i \bar{x}_i^T}_{:=C} \right) \bar{U}_j \\ &= \sum_{j=k+1}^m \bar{U}_j^T C \bar{U}_j \end{aligned}$$

נזכור כי וקטורי יחידה $U_j^T U_j = 1$ ולכן נגדיר את הלגרנז'יאן להיות

$$La = \sum_{j=k+1}^m \bar{U}_j^T C \bar{U}_j - \sum_j (U_j^T U_j - 1) \lambda_j$$

כאשר λ_j מסמנים לי פה כופלי לגראנז'. נגזור:

$$\nabla_{U_i} La = 0 \Rightarrow 2\bar{C}\bar{U}_j - 2\lambda_j \bar{U}_j = 0$$

נקבל

$$\bar{C}\bar{U}_j = \lambda_j \bar{U}_j$$

כלומר \bar{U}_j ו"ע של \bar{C} . לכן $\lambda_j = \sum_{j=k+1}^m \bar{U}_j^T C \bar{U}_j$ כאשר כאן λ_j מסמן לי ע"ע מתאים ל- \bar{U}_j . לכן נבחר λ_j להיות הקטנים ביותר.

אלגוריתם 13.1. אלגוריתם PCA.

1. העבר את ממוצע הנקודות ל-0.

2. חשב את מטריצת ה- $(\bar{x}_i - \bar{\mu}_i)(\bar{x}_i - \bar{\mu}_i)^T$. $C = \sum_j$

3. חשב ו"ע של C .

4. חשב הטלה של k ו"ע עם ע"ע גבוהים ביותר $\alpha_{k_1}, \dots, \alpha_{k_n}$.

5. היצוג של \bar{x}_i הוא $\begin{pmatrix} \alpha_{i1} \\ \vdots \\ \alpha_{ik} \end{pmatrix}$.

זה הייצוג החדש.

14 אלגוריתם ICA

נניח אני עכשיו בים ושומע מלא מלא גלים. מאחורי אני שומע אנשים שמדברים עלי ואני מעוניין בהינתן הגלי קול שלהם להבין מה הם אמרו עלי. הבעיה היא שאין לי רק את הגלי קול שלהם אלא זה מתערבב עם הגלי קול של הים. איך אני יכול להפריד בין הגלי קול של אלו שהיו מאחורי ובין הגלי קול של הים? יש לי $\bar{x} = A\bar{y}$ כאשר \bar{y} זה מה שאני מחפש. מה הבעיה? כיוון ש- \bar{x} שלי יהיה משתנה מקרי לא נורמלי, כשאני מכפיל את \bar{y} (וקטור משתנים מקריים) במטריצה כלשהי B אני מקבל $B\bar{y}$ וקטור נורמלי. אבל אז איך זה יתכן ש- \bar{x} וקטור נורמלי? הרי אני יודע שהוא קיים אז איך אני מגיע אליו? הרעיון שלי הוא שבהינתן \bar{x}_i נורמלים נחפש A כך ש- $y_i = A^{-1}\bar{x}_i$ כמה שפחות נורמלי.

מה מאפיין התפלגות נורמלית משאר ההתפלגויות? להתפלגות נורמלית יש רק מומנטים 1 ו-2 שונים מ-0, כל שאר המומנטים שווים ל-0.

אלגוריתם 14.1. ICA

נגדיר $K = \mathbb{E} \left[\left(\frac{x-\mu}{\sigma} \right)^4 \right]$. בהתפלגות נורמלית אנחנו רוצים שהוא יהיה כמה שיותר מינימלי.

בהינתן x_i ממימד n בחר מטריצה $C_{k \times n}$. חשב $y_i = Cx_i$ ובצע צעד בכיוון שממקסם

$$\sum_{\substack{\text{מימדים} \\ y}}^k K \left(\begin{array}{c} \text{כל מימד} \\ y \end{array} \right)$$

עוד דרך:

אלגוריתם 14.2. ICA בשיטה אחרת:

בהינתן \bar{x}_i נרצה להטיל אותו ל- y_i במימד 1 $\bar{U}^T \bar{x}_i$ כך ש- $MI(y_i, \bar{x}_i)$. פתרון:

(א) בחר \bar{U} אקראי.

(ב) חשב $MI(y_i, \bar{x}_i)$.

(ג) בצע צעד $\Delta U = \nabla_U MI(y_i, \bar{x}_i)$.

15 אלגוריתם C.M.D.S

בהינתן d_{ij} בין כל הזוגות המטרה היא למצוא הטלה \bar{x}_i כך ש-

$$\|x_i - x_j\| = d_{ij}$$

כלומר

$$(x_i - x_j)(x_i - x_j)^T = d_{ij}^2$$

כלומר

$$x_i^T x_i - 2x_i^T x_j + x_j^T x_j = d_{ij}^2$$

לכן אם נסמן $b_{ij} = x_i^T x_j$ נקבל

$$d_{ij}^2 = b_{ii} - 2b_{ij} + b_{jj}$$

כעת בה"כ נניח $\sum_i \bar{x}_i = 0$. d_{ij}^2 ידוע אבל b_{ij} לא ידוע כי x_i לא ידועים. אבל:

$$b_{ij} = \left(\sum_i x_i^T \right) x_j = 0$$

לכן

$$\sum_i (b_{ii} - 2b_{ij} + b_{jj}) = \sum_i d_{ij}^2$$

נסמן $\bar{x} \bar{x}^T = B$ ונסמן $T = \text{trace}(B)$. לכן

$$T = \sum_j \sum_i \frac{d_{ij}^2}{2n}$$

ולכן

$$\sum_i T + nb_{jj} = \sum_j \sum_i d_{ij}^2$$

כלומר

$$b_{jj} = \frac{\sum_i d_{ij}^2 - T}{n}$$

ואז נקבל כבר

$$b_{ij} = \frac{b_{ii} + b_{jj} - d_{ij}^2}{2}$$

אלגוריתם 15.1. C.M.D.S. אלגוריתם

1. חשב $T = \sum_j \sum_i \frac{d_{ij}^2}{2n}$
2. חשב $b_{jj} = \frac{\sum_i d_{ij}^2 - T}{n}$
3. חשב $b_{ij} = \frac{b_{ii} + b_{jj} - d_{ij}^2}{2}$
4. פרק $B = U^T \Sigma U$
5. $X = \Sigma^{\frac{1}{2}} U$

16 אלגוריתם ISOMAP

טוב לדאטה על מניפה.

בהינתן גרף לא כיווני של מרחקים w_{ij} המרחק בין i, j נגדיר $K.N.N$ בתור גרף כיווני שכל קודקוד מחובר ל- k השכנים הקרובים אליו.

אלגוריתם 16.1 ISOMAP

1. חשב גרף $K.N.N$.
2. חשב מרחקים בגרף $K.N.N$ (דייקסטרה).
3. בצע $M.D.S$.

17 אלגוריתם LLE

אני רוצה לשמר את המבנה הלינארי הלוקאלי של כל הנקודות.

אלגוריתם 17.1 LLE

1. בנה גרף $K.N.N$.
2. קרב כל נקודה להיות צירוף לינארי של השכנים: $x_i \sim \sum_j w_{ij} x_j$ על ידי מציאת w_{ij} עבורם $\|x_i - \sum_j w_{ij} x_j\|^2$ מינימלי. שימו לב שזה נפתר בנפרד. עבור \tilde{x}_i פותרים $\frac{1}{2} \sum_i \|\tilde{x}_i - \sum_j w_{ij} \tilde{x}_j\|^2$ מינימלי- פתרון גלובאלי. w_{ij} נתון- מחפשים את \tilde{x}_i כדי לא להגיע לפתרון ה-0 דורשים על \tilde{x}_i מעגל היחידה.

18 אלגוריתם EIGENMAPS

שוב עובדים על גרף $K.N.N$. היינו רוצים ש- d_{ij}^2 יהיה נמוך $\Leftrightarrow \|\tilde{x}_i - \tilde{x}_j\|$ יהיה נמוך. נגדיר דימיון $r_{ij} = f(d_{ij})$ פונקציה יורדת מ- ∞ ל-0. ההפסד יהיה

$$L = \sum_{i,j} r_{ij} \|\tilde{x}_i - \tilde{x}_j\|$$

ונחפש $\min(L)$. צריך אילוץ על \tilde{x}_i נגדיר

$$D_I = \begin{pmatrix} \sum_j w_{1j} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sum_j w_{nj} \end{pmatrix}$$

ולכן האילוץ יהיה $\tilde{x}^T D_I \tilde{x} = I$ כאשר \tilde{x}_i צמודים ל- D_I . לכן נקבל את הבעיה

$$\begin{aligned} \min \sum_{i,j} r_{ij} \|\tilde{x}_i - \tilde{x}_j\| \\ \text{s.t. } \tilde{x}^T D_I \tilde{x} = I \end{aligned}$$

נקבל שהפתרון הוא $L = W - D_I$.

19 הולכים מקריים, צוואר בקבוק לאינפורמציה ואלגוריתם T-SNE

צוואר בקבוק: נעשה את הניסוי המשעשע הבא:

אליס ובוב שניהם רשומים לקורס כלשהו באוניברסיטה. רצה הגורל ויום אחד בוב לא יכול להגיע להרצאה ואליס ישבה וסבלה שעה וחצי לבדה בהרצאה. בסוף השיעור הגיע בוב אחרי שרץ בטורבו 5 בלוקים וביקש מהמרצה שיסכם לו את מה שקרה בשיעור, ובאורך פלא 3 משפטים המרצה סיכם לבוב את מה שהלך בשיעור. נרגש ונפעם בוב הלך לאליס וסיפר לה את ה-3 משפטים שהמרצה אמר לו והמומה אליס אמרה לו שזה סיכום מאוד מדויק של כל מה שהלך בשיעור. מה המסקנה? היה אפשר לדחוס את השעה וחצי של סבל שחווה אליס לכדי 3 משפטים וסתם היה בזבז של זמן.

הילוך מקרי: נניח יש לי עכשיו גרף של נקודות במרחב כלשהו. נסמן $\mathbb{P}(j|i) = \mathbb{P}(\text{הסיכוי שאקפוץ מ-} i \text{ ל-} j) = \frac{f(\|x_i - x_j\|)}{\sum_{i \neq j} f(\|x_i - x_j\|)}$ כאשר $\sum_j \mathbb{P}(j|i) = 1$ וכן ניקח את $f \sim \mathcal{N}(0, \sigma)$ לכן

$$\mathbb{P}_{j|i} = \frac{e^{-\frac{\|x_i - x_j\|^2}{2\sigma_i^2}}}{\sum_{i \neq j} e^{-\frac{\|x_i - x_j\|^2}{2\sigma_i^2}}}$$

מתקיים

$$\mathbb{P}_{j,i} = \frac{\mathbb{P}_{j|i} + \mathbb{P}_{i|j}}{2n}$$

ואכן מתקיים $\sum_i \sum_j \mathbb{P}_{i,j} = 1$ אני רוצה \tilde{x}_i בו הסיכוי לקפוץ מ- \tilde{x}_j דומה לסיכוי לקפוץ מ- i ל- j .

ברמה הרעיונית איך אלגוריתם $T - SNE$ יפעל:

1. חשב סיכויי קפיצה ל- k שכנים.
2. הפוך לסימטרי.
3. מצא \tilde{x}_i עבורם סיכויי הקפיצה דומים לאלו של x_i אבל החלף את ההתפלגות הנורמלית בהתפלגות t :

$$g(\|\tilde{x}_i - \tilde{x}_j\|) = \frac{1}{1 + \|\tilde{x}_i - \tilde{x}_j\|}$$

ללא שום פרמטר.

4. מצא \tilde{x}_i שממזערים את $D_{KL}(\mathbb{P}, \tilde{\mathbb{P}})$.

כאשר בחרתי את σ_i בחרתי אותו לקבל אנטרופיה של $\mathbb{P}_{j|i}$ כרצוני $-\sum_j \mathbb{P}_{j|i} \log(\mathbb{P}_{j|i})$.

כעת יש לי $\mathbb{P}_{j,i} = \frac{\mathbb{P}_{j|i} + \mathbb{P}_{i|j}}{2n}$. אני מחפש \bar{y}_i עבורם הסיכוי לקפוץ מ- i ל- j הוא \mathbb{Q}_{ij} עם התפלגות (ב). נגדיר $\mathbb{Q}_{ij} = \frac{\frac{1}{1 + \|\bar{y}_i - \bar{y}_j\|}}{\sum_{i,j} \frac{1}{1 + \|\bar{y}_i - \bar{y}_j\|}}$ ואני רוצה להעריך

את הדמיון בין ההתפלגות הנצפית שלי להתפלגות האמיתית, לכן נמדוד את הדמיון לפי $D_{KL}(\mathbb{P}_{ij}, \mathbb{Q}_{ij})$. איך מזה נמצא את \bar{y} ? נעשה ירידה במורד הגראדיאנט.

20 אלגוריתם U-MAP

זה $sne - t$ לקטני אמונה. נגדיר $\mathbb{P}_{ij} = e^{-\frac{\|x_i - x_j\|^2}{2\sigma_i^2}}$ וכן $\mathbb{Q}_{ij} = \frac{1}{1 + e^{-\frac{\|y_i - y_j\|^2 + 2\rho}{\sigma_y}}}$ המרחק בין \mathbb{P} ל- \mathbb{Q} יהיה

$$D(\mathbb{P}, \mathbb{Q}) = \sum_{ij} \mathbb{P}_{ij} \ln \mathbb{Q}_{ij} + \sum_{ij} (1 - \mathbb{P}_{ij}) \ln(1 - \mathbb{Q}_{ij})$$

21 מקודד אוטומטי.

יש לי תצפיות \bar{x}_i . יש לי מקודד $y_i = f(\bar{x}_i)$ וכן מפענח $\tilde{x}_i = g(y_i)$. אני רוצה $\tilde{x}_i \sim \bar{x}_i$. נגדיר את ההפסד להיות

$$L = \frac{1}{2} \sum_i D(g(f(\bar{x}_i)), \bar{x}_i)^2$$

ואני רוצה למזער את L . בדרך נבחר $D(x, y) = \|x - y\|^2$.

אני רוצה שהמקודד יהיה רציף מקומית, כלומר אם יש לי $x_i \rightarrow y_i \rightarrow \tilde{x}_i$ וכן $\tilde{x}_i \sim \bar{x}_i$ אזי $\tilde{x}_i + \bar{\epsilon}_1 \sim \tilde{x}_i + \bar{\epsilon}_2$.

כעת נניח אני רוצה ללמוד שני ערכים שונים לחלוטין מ- $f(x) : \mu_{\bar{x}}$ להיות ממוצע של x וכן $\sigma_{\bar{x}}$ להיות השונות של x . נגדיר $y(x) = \mu_x + \sigma_x \mathcal{N}(0, 1)$ ואני רוצה למצוא $\hat{x} = g(y(\hat{x}))$ והפעם אני עובר מסביבה של \bar{x} ל- y ל- \tilde{x} .

נגיד \bar{x} במימד 1000, μ_x, σ_x, \bar{y} במימד 15. הבעיה היא שהמודל יקטין לי את σ עד לאיזה רמת דיוק שאני רוצה כדי לדייק לכן אני צריך להוסיף אילוץ $\bar{y}_i \sim \mathcal{N}(0, 1)$ בכל מימד. לכן נעדכן את L להיות $L = \frac{1}{2} \sum_i d(\bar{x}_i, \tilde{x}_i)^2 + D_{KL}(\bar{y}_i, \mathcal{N}(0, 1))$. אם אני לא רוצה לחשוש מזה שחלק מהנתונים שלי

נעלמים בזמן שאני מוצא את f ו- g אני יכול למצוא אותן כשכל פעם אני משתמש גם בחלק מהמידע שלי.

חלק VII

למידה מפוקחת עצמאית.

נניח יש לי ספר, ואת אותו ספר מתורגם לרוסית. איך אני יכול להבין שמדובר על אותו ספר?

נניח הספר בעברית הוא x_i והספר ברוסית הוא x_j . אני אייצר מלא מלא הפרעות על x_i ועם x_j :

$$\begin{aligned} x_i &\rightarrow \{\tilde{x}_i\} \\ x_j &\rightarrow \{\tilde{x}_j\} \end{aligned}$$

כעת אני לוקח זוג $[x_1, x_2]$. ניצור הטלה שלהם $z_1, z_2 \in \mathbb{R}^n$ הטלות $x_1, x_2 \rightarrow$. אני רוצה להגדיר הפסד שיביא לי ערך גבוה (מרחק נמוך) אם x_1, x_2 מאותה תמונה וערך נמוך (מרחק גבוה) אם x_1, x_2 לא מאותה תמונה. נגדיר SELF SUPERVISED LOSS להיות

$$L = \frac{\sum_{i,j} e^{-d(z_i, z_j)}}{\sum_{i,j} e^{-d(z_i, z_j)}}$$

VIII חלק

זיהוי חריגים.

כיום יש 3 סוגים של חריגים:

1. Point Anomaly

2. Contextual Anomaly

3. Collective Anomaly

יש כמה שיטות לזיהוי אנומליות:

- אשכול.
- צפיפות.
- מרחק.
- מודל סטטיסטי.
- Embedding.

מודל סטטיסטי לזיהוי אנומליות:

אנומליה = אירוע עם סיכוי נמוך.

בהינתן אוסף נקודות \bar{x}_i ומודל שתלוי בפרמטרים סמויים $\bar{\theta}$, נחשב באמצעות EM את $\bar{\theta}$, אבל לא נכניס לחישוב ערכים עם סיכוי מתחת ל- ε (לבחירת).

1. ניקח את הנקודות, נבנה מודל סטטיסטי.

2. נחשב סיכוי של נק' ונזהה אנומליות.

3. נוציא אנומליות מהנתונים.

4. נחזור ל-1).

אשכול לזיהוי אנומליות:

נוכל לעשות DBSCAN ואם יש לנו נקודה רחוקה מנקודות גרעין אזי היא חריגה.

נוכל גם לעשות FUZZY - C - MEANS ונקודה ש"שייכת לכמה אשכולות" תהיה חריגה.

צפיפות לזיהוי אנומליות:

אפשר להפעיל KDE על הנתונים שלי ולהתבונן בערכי ה-KDE של הנקודות.

מרחק לזיהוי אנומליות:

נחשב מרחקים של כל הנקודות ונשאל מה המרחק של השכן ה-k.

בנוסף: SVM לזיהוי אנומליות:

בהינתן מרכז $\bar{\mu}$ ורדיוס R נגדיר הפסד

$$L(x_i) = \begin{cases} 0 & \|x_i - \bar{\mu}\| < R \\ \|x_i - \bar{\mu}\| - R & \text{else} \end{cases}$$

לכן אני רוצה לפתור

$$\min_{\mu, R} \sum_i L(x_i) + f(R)$$

ואז אם $\|x_i - \bar{\mu}\| > R$ אזי x_i חריג. בדרך כלל: $f(R) = \alpha R^2$.

בחיים האמיתיים זה לא עובד.

נטיל למימד יותר גבוה: $x \rightarrow \phi(x_i)$ ונעשה

$$L(x_i) = \begin{cases} 0 & \|\phi(x_i) - \bar{\mu}\| < R \\ \|\phi(x_i) - \bar{\mu}\| - R & \text{else} \end{cases}$$

$$\min_{\mu, R} \sum_i L(x_i) + f(R)$$

חלק IX

נתונים דינאמיים.

יש 2 סוגים של נתונים דינאמיים:

הראשון הוא בדידים בזמן: O_1, O_2, O_3, \dots כאשר O_i הוא תצפית מזמן i .

הנחות שניח לאורך הדרך:

(א) מרקוביות: העתיד תלוי בהווה בלבד ולא בעבר:

$$\mathbb{P}(x_{t+1}|x_t, \dots, x_1) = \mathbb{P}(x_{t+1}|x_t)$$

(ב) המציאות סמויה.

מקרה הכי פשוט:

(א) הזמן בדיד.

(ב) המצבים הסמויים סופיים ובדידים.

$$a_{ij} = \mathbb{P}(j|i)$$

$$\mathbb{P}_i = \mathbb{P}(i)$$

$$b_{ik} = \mathbb{P}(i \text{ מצב } k \text{ תצפיות})$$

שאלות:

(א) בהינתן מודל (\mathbb{P}, A, B) וסדרת תצפיות o_1, \dots, o_T , מה ההסתברות לסדרת התצפיות?

(ב) בהינתן מודל ותצפיות כנ"ל, מה רצף האירועים הכי סביר (רצף המצבים הסמויים)?

(ג) בהינתן אוסף של תצפיות- מה המודל הכי סביר?

22 מודלים מרקוביים סמויים.

נגדיר:

סיכוי התחלתי לכל מצב: P .

סיכוי לעבור ממצב i למצב j : A .

סיכוי לראות תצפית k בהינתן מצב i : B .

$$B = \begin{pmatrix} 0.2 & 0.1 & 0.1 & 0.1 & 0.5 \\ 0.7 & 0 & 0 & 0.3 & 0 \\ 0 & 0.5 & 0.5 & 0 & 0 \end{pmatrix} \begin{matrix} 1 \\ 2 \\ 3 \end{matrix}, A = \begin{pmatrix} 0.8 & 0 & 0.2 \\ 0 & 0.4 & 0.6 \\ 0.1 & 0 & 0.9 \end{pmatrix}, P = \begin{pmatrix} 0.5 \\ 0.3 \\ 0.2 \end{pmatrix}$$

ה ט ג ב א

אני רוצה לשאול מה הסיכוי של רצף תצפיות (O_1, \dots, O_T) . לפי משפט ההתפלגות השלמה

$$\mathbb{P}(O_1, \dots, O_T) = \sum_i \mathbb{P}(O_1, \dots, O_T, T \text{ בזמן } i \text{ מצב}) = \sum_i \alpha_i(T)$$

כאשר

$$\alpha_i(1) = \mathbb{P}_i b_{i, O_1}$$

ואז

$$\alpha_i(t) = \sum_j \mathbb{P}(O_1, \dots, O_{t-1}, t-1 \text{ בזמן } j \text{ מצב}) a_{ji} b_{i, O(t)}$$

בעת, בהינתן אוסף תצפיות אני רוצה לשאול מה הרצף הכי סביר. ביותר פורמלי: מה הרצף S_1, \dots, S_t הכי סביר לתצפיות O_1, \dots, O_t . אני רוצה

$$\delta_i(t) = \mathbb{P}(S_1, \dots, S_t, O_1, \dots, O_t)$$

כתור התחלה :

$$\delta_i(1) = \mathbb{P}_i b_{i,O_1}$$

וכן

$$\delta_i(t+1) = \mathbb{P}(S_1, \dots, S_t, O_1, \dots, O_t) = \max(\mathbb{P}_i(S_1, \dots, S_{t-1}, O_1, \dots, O_{t-1}) a_{ji} b_{i,o_t})$$

בהינתן אוסף רצפי זמן של תצפיות, מה המודל הכי סביר? נעשה EM :

(א) נניח מודל (P_0, A_0, B_0) .

(ב) נחשב רצף מצבים הכי סביר לכל רצף תצפיות.

(ג) נחשב מודל חדש לפי רצף מצבים.

(ד) כל עוד לא התכנסנו נחזור ל-(ב').

יש לי כמה שאלות כעת :

(א) מה הסיכוי לרצף תצפיות? אם יש לנו

$$\alpha_i(T) = \mathbb{P}(O_1, \dots, O_T, i \text{ להיות במצב } i)$$

$$\alpha_i(1) = \mathbb{P}_i b_{i,O_1}$$

$$\alpha_i(t+1) = \sum_j \alpha_j(t) a_{ji} b_{i,O_{t+1}}$$

אז

$$\mathbb{P}(O_1, \dots, O_T) = \sum_j \alpha_j(t)$$

(ב) מה רצף המצבים הכי סביר בהינתן רצף תצפיות? זה מה שעשינו :

$$\delta_i(t) = \mathbb{P}(O_1, \dots, O_t, i \text{ במסלול הכי סביר})$$

וכן

$$\delta_i(1) = \mathbb{P}_i b_{i,O_1}$$

אז

$$\delta_i(t+1) = \max(\delta_i(t) a_{ji} b_{i,o_t})$$

נמצא את מקסימום $\delta_i(T)$ וחוזרים אחורה במסלול הכי סביר- כמו בתכנון דינמי.

(ג) בהינתן הרבה רצפי תצפיות, מה הוא המודל הכי סביר (לוקאלי)? נגדיר

$$\gamma_i(t) = \mathbb{P}(O_1, \dots, O_T, t \text{ בזמן } i \text{ להיות במצב } i)$$

$$\beta_i(t) = \mathbb{P}(O_1, \dots, O_T | t \text{ בזמן } i \text{ להיות במצב } i)$$

ואז

$$\beta_i(T) = 1$$

$$\beta_j(t+1) = \sum_i a_{ji} \beta_i(t) b_{i,O_t}$$

וכן גם

$$\alpha_i(t) = \sum_j a_{ji} \alpha_j(t-1) b_{i,O_t}$$

אז נקבל

$$\gamma_i(t) = \mathbb{P}(O_1, \dots, O_t | t \text{ בזמן } i \text{ להיות במצב } i) \mathbb{P}(O_{t+1}, \dots, O_T | t \text{ בזמן } i \text{ הייתי במצב } i) = \alpha_i(t) \beta_i(t)$$

עכשיו אנחנו רוצים לדעת את הסיכוי של i בזמן t :

$$\mathbb{P}(i, t \text{ בזמן } i) = \frac{\mathbb{P}(i, t \text{ בזמן } i, O_1, \dots, O_T)}{\mathbb{P}(O_1, \dots, O_T)} = \frac{\alpha_i(t) \beta_i(t)}{\sum_{i'} \alpha_{i'}(t) \beta_{i'}(t)} \Rightarrow \sum_i \mathbb{P}(i, t \text{ בזמן } i) = 1$$

מכאן שסיכוי ההתחלה החדש

$$new \mathbb{P}_i(1) = \frac{\alpha_i(1) \beta_i(1)}{\sum_{i'} \alpha_{i'}(1) \beta_{i'}(1)}$$

$$new a_{ij} = \frac{\sum_{\text{זמנים}} \text{מעברים מ-} i \text{ ל-} j}{\sum_{\text{זמנים}} i \text{ סיכוי להיות במצב } i} = \frac{\sum_t \alpha_i(t) a_{ij} b_{i, O_{t+1}} \beta_j(t+1)}{\sum_t \gamma_i(t)}$$

וכן

$$new b_{ik} = \frac{\sum_t i \text{ ראייתי תצפית } k \cdot \text{הייתי במצב } i}{\sum_t i \text{ הייתי במצב } i} = \frac{\sum_t \gamma_{it} \mathbb{I}_{O_k, t}}{\sum_t \gamma_{it}}$$

וקיבלתי מודל חדש, כלומר זה שלב ה-EM. כמו בכל EM את המודל הזה אני שוב פעם אפשר עד שאני מגיע לקיצון מקומי.

23 טיפול במשתנים קטגוריאליים - אלגוריתם MCA.

יהי $x_i \in \{z_1, \dots, z_k\}$ זה משתנה בדיד. איך אני עובר ממנו למשנה רציף?

אם זה בינארי עוברים לוקטור 0, 1:

כלומר אם לא $x_i = 0$, אם כן $x_i = 1$ ואם לא ידוע עד שמים חציין/ממוצע. אם יש יותר מ-2 אפשרויות נעשה one-hot: אם יש לי k קטגוריות אני אצור k עמודות וכל דגימה תקבל 1 בקטגוריה שלה.

לדוגמה: אם יש לי רשימה {אמא, אח, אמא, אבא, אבא} אז נקבל

אח	אבא	אמא
0	1	0
0	1	0
0	0	1
1	0	0
0	0	1

אם אני רוצה לעשות PCA עם one hot לא נקבל תוצאות טובות- לשם כך יש אלגוריתם MCA:

נמרכז את x : נוריד מכל עמודה את הממוצע שלה. אחרי זה נחשב $B = x^T x$, מטריצה בגודל מספר העמודות \times מספר העמודות. נטיל על הו"ע של B כמו ב-PCA ואם יש לי שורה של לא ידוע נחליף בשורה אפשים.

אם יש לי דאטא לא קטגוריאלי משולב עם כן קטגוריאלי, נפצל אותם, נעשה לכל אחד את ההורדת מימדים שלו ואז נאחד אותם ועל התואצה נעשה עוד הורדת מימדים.

אחרי MCA: קיבלתי וקטור רציף, אפשר לשלב עם משתנים רציפים אחרים אם רוצים להמשיך להוריד מימד.

24 בנוס.

כעת יש לי שאלה אחרת: בהינתן נדב ורשימת הסרטים שהוא ראה, אני רוצה לחזות איזה סרטים הוא רוצה לראות. כלומר אם יש לי מטריצה של ראה סרט לא ראה סרט אני רוצה

מסקולות		סרטים																					
אנשים	<table><tr><td></td><td></td><td></td><td></td></tr><tr><td></td><td></td><td></td><td></td></tr><tr><td></td><td></td><td></td><td></td></tr><tr><td></td><td></td><td></td><td></td></tr><tr><td></td><td></td><td></td><td></td></tr></table>																					=	
	V																						

מסקולות		סרטים											
אנשים	<table><tr><td>□□</td><td></td></tr><tr><td></td><td></td></tr><tr><td></td><td></td></tr><tr><td></td><td></td></tr><tr><td></td><td></td></tr></table>	□□											
□□													
	W												

מסקולות		סרטים																					
אנשים	<table><tr><td></td><td></td><td></td><td></td></tr><tr><td></td><td></td><td></td><td></td></tr><tr><td></td><td></td><td></td><td></td></tr><tr><td></td><td></td><td></td><td></td></tr><tr><td></td><td></td><td></td><td></td></tr></table>																						
	K																						

מסקולות		סרטים																					
אנשים	<table><tr><td></td><td></td><td></td><td></td></tr><tr><td></td><td></td><td></td><td></td></tr><tr><td></td><td></td><td></td><td></td></tr><tr><td></td><td></td><td></td><td></td></tr><tr><td></td><td></td><td></td><td></td></tr></table>																						
	K																						

אני רוצה $\tilde{V} = V \sim WH$. אני רוצה \tilde{V} עם ערכים שהם לא שליליים. איך אני עושה את זה? הכי פשוט: נכריח את H, W להיות רק עם ערכים חיוביים. אני לא יכול לעשות $W^{t+1} = W^t + ?$ כי W^{t+1} יכול להרוס לי את הכל ולהכניס מינוסים. לכן אני אעשה (משתנה חיובי) $W_{ij}^{t+1} = W_{ij}^t$. איך אני עושה את זה? בעזרת מורד הגרדיאנט עם טוויסט:

בהינתן $F(h)$ אני אגדיר $G(h, h')$ כך ש- $G(h, h') \geq F(h)$ וכן $G(h, h) = F(h)$. כדי למזער כל פונקציה $F(h)$:

(א) נגדיר $G(h, h')$

(ב) נגדיר $h^0 = h$ עבור h כרצוני וכן $h^{t+1} = \arg \min G(h, h^t)$. כעת, $F(h^{t+1}) = G(h^{t+1}, h^{t+1}) \leq G(h^{t+1}, h^t) \leq G(h^t, h^t) = F(h^t)$. נדלג על המון המון שלבים ואז נקבל את העדכון

$$new H_{kj} = \frac{H_{kj} [W^T V]_{kj}}{[W^T W H]_{kj}}$$

$$new W_{ik} = \frac{H_{ik} [V H^T]_{ik}}{[W H H^T]_{ik}}$$

וסיימנו את הקורס.