

כ' - 22/6/20

תרגול # 12

\* תרגול # 11 היה תרגול חציה לקראת המבחן.

דבר 4: תוכנת מני-3 - מוטיבציה

• PCA: 2 פרויקט, שינוי (נושקוים):

• ג- המחקר האפני הקורס קיור.

• ג המחקר שמחקרם את השוק.

מוטיבציה:

דאטה מנימז פיצורים גקור מליכ בפנינו 2 סוגים של אלגוריתם:

1. חישובי- כל שמד הפיצורים אלה, אולה סימבולי מחלקת הדיפוזציה,

ונצטרך יותר צדמור כז' לאמוז, ויקח לנו יותר זמן לעשות זאת

(כי דוד האלגוריתם שראנו הם פולינומליים קממן ו-  $m$ ).

2. data-exploration: נצטרך לחקור את הדאטה, אך זה קשה לחקור כמוה צדולה של פיצורים:

לחמין את ההתפלגות שלהם, מה היחסים ביניהם, איך הם משפיעים

4 המצאה... כחוס -

data-visualization: במימדים גקוריים לא ניתן לצייר את הדאטה.

$\Leftrightarrow$  נצטרך לעשות גרפיק תוכנת מנימז:  $\varphi: \mathbb{R}^d \rightarrow \mathbb{R}^k$  ( $k \leq d$ ),  $\{x_i\}_{i=1}^m \mapsto \{\varphi(x_i)\}_{i=1}^m$

נכצץ- טרנספורמציה  $\varphi$  הדאטה שלנו לגרפיק צדונו שלמים (אנאליים)

או לא אנאליים) של הפיצורים במרחב חדש, שבו נכצץ למיניה.

\* אלא אם הדאטה באמת יושב ב-  $\mathbb{R}^k$ , גרפיק תוכנת מנימז כדור

באפוז מניצץ. בתגלית זה נצטרך לשלוט בחלק שלמות אנתוני משמנים.

בוגלול:

• יש מעגל (טקסט) ב-  $\mathbb{R}^3$  שמקביל למישור  $xy$ . את ב"ו מתקוצות

היא נק' 4 המעגל + נחש שאוסט מליים אלס הכוונות.  $x \leftarrow y$

מכיוון שאנחנו יוצרים שהנק' יושקו 4 מעגל, וזה החלק ה"מעגל" בדאטה

(הצדקה), יכול להיות שנק' להפתיג מנימז למימז 2, וזה יקל עלינו את העבודה.

• נניח שאנו חוקרים מחלה מלימז, וחוקרים צדורה רכיב גנטי. רוצים למנוע

אנשים, שאוקים במחלה ובאוו שלא, ולומר האם יש אלמנט מוטציה גנטי

שצדומה למחלה.

אטם הוא כצ' אנון מל האלסקיה  $\Sigma = \{A, T, C, G\}$  באורך  $3.2 \times 10^9$ .

נשווה את הגנום של כמה אנשים, (מכאן את ההקבלים ונשאף את צדמנו האם

להקבלים האלו שמצונו יש קשר למחלה.



האם יש לנו כמות  $3.2 \times 10^9$  פיצויים?   
 נאמן 1: (הצמצם הפיזיקלי שלהם יש לנו הערכים בין אנליס שונים.

← הצבתו להצמצם ל-  $300 \times 10^6$  פיצויים.

נאמן 2: נצטרך שיהיה לנו אזור מסוים למה "reference", ולא כל יגד  
 השוני נשמה ל- "reference" הזה.

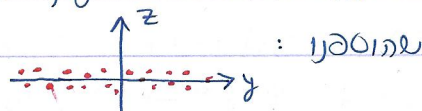
← הצבתו להצמצם ל-  $5 \times 10^6$  פיצויים.

★ הצגתו באותה אמת שיהיה לנצל יצחק קוצם על המצאה כדי להקטין את מיתות.  
 ★ זה לא feature selection. מוגנים על חלק מהמיתות.

כשאנחנו מקבלים הורג מיתות אנחנו מאקרים נפח מהצמחה. חלם שמות  
 המיתות, נשאר את החלק שאנחנו רוצים. אלו גמול נוצר לשמור?

1. לשמור על השונות בצמחה, כדי שיהיה להחזיר את המופע אמת  
 אנו מציגים אלמנט.

צד: בצמחה המעלה, ארץ נשא את הצמחה על מישור  $xy$  נקרא נקודות  
 מוכנסות סביב ציר ה- $z$ , באינטרוול שקבע לפי המעלה המאוסני.



אם הצמחה לא מקבלת מישור, אלא נמצא באזור של  $\theta \in (0, 90^\circ)$

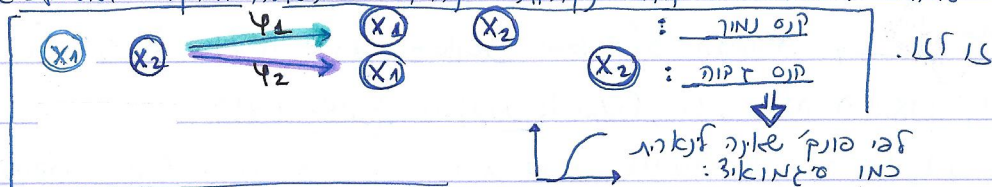
אז ההטחה יהיה מתבצעת על קואורדינטה שיהיה ל-3 הצירים, כך שיהיה  
 מ המרחק שנקרא סטור, יפיק לנו גמולה כמו ההטחה הקודמת על צד.

2. לשמור מרחקים קווקים. נצב את אלמנטים טיפניים שמבצעים הורג מיתות

לא לינארי: נאניש באופן שונה (ואולי לינארי) הטלה של נקודות

שמנואל קרום זו לא הצמחה המקורי אך מוקד זו מניח המיתות המופחת

לשמור הטלה שוקדמת נקודות שקדמות בצמחה המקורית ומטלה את קרום



3. לשמור צורה פנימית של הצמחה. נאציר:

המרחק האמפאלי - המרחק שליו הצמחה שלנו מיוצג בצורה המקורית.

למא קצמחה של המעלה המרחק האמפאלי הוא  $\mathbb{R}^3$ .

המרחק האנטרופי - מרחק המרחק עליו יושב הצמחה. בצמחה המעלה:  $\mathbb{R}^2$ .



## PCA

(1) PCA הוא פרוטקטור הנוצרת במינימום לא מרחק אפני בשימוש, הקונוס ביוגר  
 דנקונור הצאטה

$$w, u^* = \operatorname{argmin}_{w, u} \sum_{i=1}^m \|x_i - u w x_i\|_2^2$$

מחפשים מטריונר

(2) פרוטקטור של PCA הוא הנוצרת מינימום משמור שונור:  
 מחפשים:  $M^* = \operatorname{argmax}_M \operatorname{Var}(XM)$ , שומר הטלה למרחק מינימום  
 נמנך יותר שפטיאר  $M$  orthogonal הצאטה במ המרחק יהיה מקסימלי.

השונור שומר נרבה למרחק היא שונור ה- sample שמחושבר  $\bar{x}$ .

$$\text{Sample Variance} = \frac{1}{m} \sum_{i=1}^m \|x_i - \bar{x}\|_2^2$$

Sample Variance =  $\bar{x}$  -  $\bar{x}$  ממוצע  $\bar{x}$   $\{x_i\}_{i=1}^m$

★ יש אצב 2 פרוטקטור של PCA שומר לא נצב.

נצב כ"א מהפרוטקטור פלאו, ולקסול נראה שהם שקולים: (11)

### 1 ג-מרחק אפני:

כפי שהנאנו בתרבות, יש לנו פרוטקטור סאור לבציה זו.

כל שפטיאר הנוצרת מהשימור ועמקצ בחלקים מסוימים קרה:

פרוטקטור ה- PCA: שפטיאר הצאטה ונמנכט אלו (פחוג אלו  $\bar{x}$ ),

$$A = \sum_{i=1}^m x_i x_i^T$$

וניצור אלו המטריונר: נמנכט אלו הטלה ותוץ

ש  $A$  ונוכים מהם אלו  $w, u$ : צמיונר  $u$  יפיו א הוקטורים העיקריים

(עם צוב)  $A$ ,  $u = u^T$ .

למה צב ככה? כעניציה להקציר אלו שומר נחפש מטריונר  $u$  שפטיאר א שפטיאר

אלו  $x$  למו המרחק  $\|x - u w x\|_2$ , שומר מחפשים העקרה  $\varphi: X \rightarrow \mathbb{R}$ ,

אך נרבה שפטיאר זו גמאצ אלו סכום המרחקים, שומר אלו  $\|x - u w x\|_2^2$  (\*).

מסבר שפטיאר המטריונר המקימט שלא קן המטריונר שפטיאר אלו הטלה

האורמאטולג  $u$  מ המרחק.

$$\left[ \begin{array}{l} \text{הטלה אורמאטולג - יפיו } R^d \subset \mathbb{R}^d \text{ מינימום א, ויפיו קסים אורמאטולג} \\ \text{ש-} V: \{v_1, \dots, v_k\} \text{ אלו: } P = \sum_{i=1}^k v_i v_i^T \text{ מטריונר הטלה אורמאטולג} \\ \text{מטריונר מטריונר הטלה:} \\ 1. P \text{ סימטרי} \\ 2. \text{ש } x \in \mathbb{R}^d, u \in V: \|x - P x\|_2 \geq \|x - u\|_2 \text{ אלו} \\ P x \in \mathbb{R}^d \text{ היא הטלה של } x \text{ למו המרחק.} \end{array} \right]$$

הנק הקרובה ביותר  $x \in \mathbb{R}^d$  למו המרחק הוא הטלה שו  $u$

המרחק:  $P_x$ , ולק כפי למאצ אלו (\*) נצטרך למאצ מטריונר הטלה אורמאטולג.

כלומר בעל האופטימיזציה היא מנסה  
 (ע"י  $w = u^T$  טאנחן צייכט מטרית הטלה אורטוגונלית)  $(\sum_{i=1}^k u_i u_i^T)$   
 מבוזר, לאחר מזה, זה שקול למקסימום:  $\|x\|^2 - \text{trace}[u^T (\sum_{i=1}^m x_i x_i^T) u]$   
 נגזר:  $:= A$

(בחין ש- $A$  סימטרית ו- $\text{positive semidefinite (PSD)}$  :  $G_N$

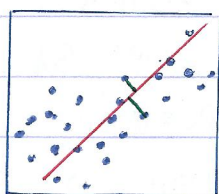
$\mathbb{R}^{k \times k}$  סימטרית שנקראת  $x^T A x \geq 0 : x \in \mathbb{R}^k$   
 נניח ש- $A$  היא  $\text{PSD}$ , ו- $v \in V$  :  $v^T A v = v^T (\sum_{i=1}^k x_i x_i^T) v = \sum_{i=1}^k v^T x_i x_i^T v = \sum_{i=1}^k \langle v, x_i \rangle \langle x_i, v \rangle = \sum_{i=1}^k \langle v, x_i \rangle^2 \geq 0$   
 סימטרית מ"מ

מטון ש- $A$  סימטרית ו- $\text{PSD}$  קיימת פירוק  $A = V D V^T$  :  $V$  אורטוגונלית,  $D$  אלכסונית של אלכסון של מופע ה"ץ של  $A$  (שהם חיוביים) (הערכים העיגוליים של  $A$ ), ו- $V$  מופע ה"ץ המאליני.  
 ש"ץ מופע  $D$

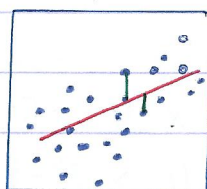
נסור להסוי, אלא אנתן מסיק למקסימום :  $\|x\|^2 - \text{trace}[u^T A u] \stackrel{\text{SVD}}{=} \text{trace}[u^T V D V^T u] \stackrel{B := u^T V}{=} \text{trace}[B^T D B] = \sum_{i=1}^d [B^T D B]_{ii} = \sum_{i=1}^d D_{ii} \langle B_i, B_i \rangle = \sum_{i=1}^d D_{ii} \sum_{j=1}^m B_{ij}^2 \leq \sum_{i=1}^d D_{ii}$   
 קסר מלי מלי  $B_i$  :  $\sum_{j=1}^m B_{ij}^2 \leq 1$

כלומר קיבלנו ש- $\text{trace}[u^T A u] \leq \sum_{i=1}^d D_{ii}$  :  
 וזה הערך המקסימלי שאפשר לקבל. מהו ה- $\text{argmax}$ ?

$u$  שמוביזה הם א ה"ץ העזרים ביותר ב- $A$ , מקרה מרקי:  
 $\text{trace}[u^T A u] = \sum_{i=1}^k D_{ii}$



PCA



OLS

אנטלציה זאמלונג:  
 יש לך זאמל ו- $z$   
 למקור המוצגים ע"י:



אך המונח, באיזה שנים? ההסבר נעץ בתחילת ש"ס מהמזלזל -

• סל: התחיל שמש"ל צ"ע. וספר בקירוב מ מרחק א צ"ע ואז

רוב האנשים למדו בקיבוץ א ה-י. שומר:

$y = \beta_0 + x^T \beta_1 + \varepsilon$ ,  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$

$$\varepsilon \sim \mathcal{N}(0, \sigma^2) \quad , \quad y = \beta_0 + x^T \beta_1 + \varepsilon$$

• PCA: (למיזוג אנרגיה מפורקת,  $x, y$ , נובים אל המטלה האוטומטית  $U$  ו- $V$  המרחב. מפני שלצדית  $\mu$  נוסף הצדית שווה  $\sigma^2 I$ ,  $\sigma^2$  נק' יטלה  $\sigma^2$  הכיוונים באופן שווה ולכן נוצר למשל אוטומטית. הנבטה  $\mu$ ,  $\sigma^2 I$ :  $(x, y) \sim \mathcal{N}(\mu, \sigma^2 I)$ .

$$(x, y) \sim \mathcal{N}(\mu, \sigma^2 I_d) : \mu, \sigma \in \mathbb{R}$$

כל נחשב את  $M$  המינימלית הקרובה ביותר ל- $X$  מתוך  $M$  אורתוגונלית,  $M^* = \underset{M \text{ orthogonal}}{\operatorname{argmax}} \operatorname{Var}(XM)$  : את פיצול הקוליות  $f$  מ- $M$  המינימלית.

$M^* = \underset{M \text{ orthogonal}}{\operatorname{argmax}} \operatorname{Var}(XM)$  : מצא את המטריצה  $M$  ה-orthogonal המקסימלית את התנודות של  $XM$

גנרליזר ה-PCA: באשיר נמיכס אל תצטח (נחסר  $\bar{x}$ )

שקלא והכרעא :

מיוחק (תק' מג-מחב) (שאג סיבו) ננסה לעדיר, ואלא מיוחק ההטלה האורגניזם

$$a^2 + b^2 = c^2 \quad : \text{Pythagoras Theorem}$$
$$b = \sum \frac{\|x_i - \mu_W x_i\|_2^2}{2} \quad \text{w/c}$$

פונק' המטרה כפרוש ראשוני,  $C = \text{Var}(XM)$  : פונ' המטרה כפרוש ראשוני.  
שנויים למעלה → ולפי זה נבדוק שיהיו. (ii) ← שנויים למטה

