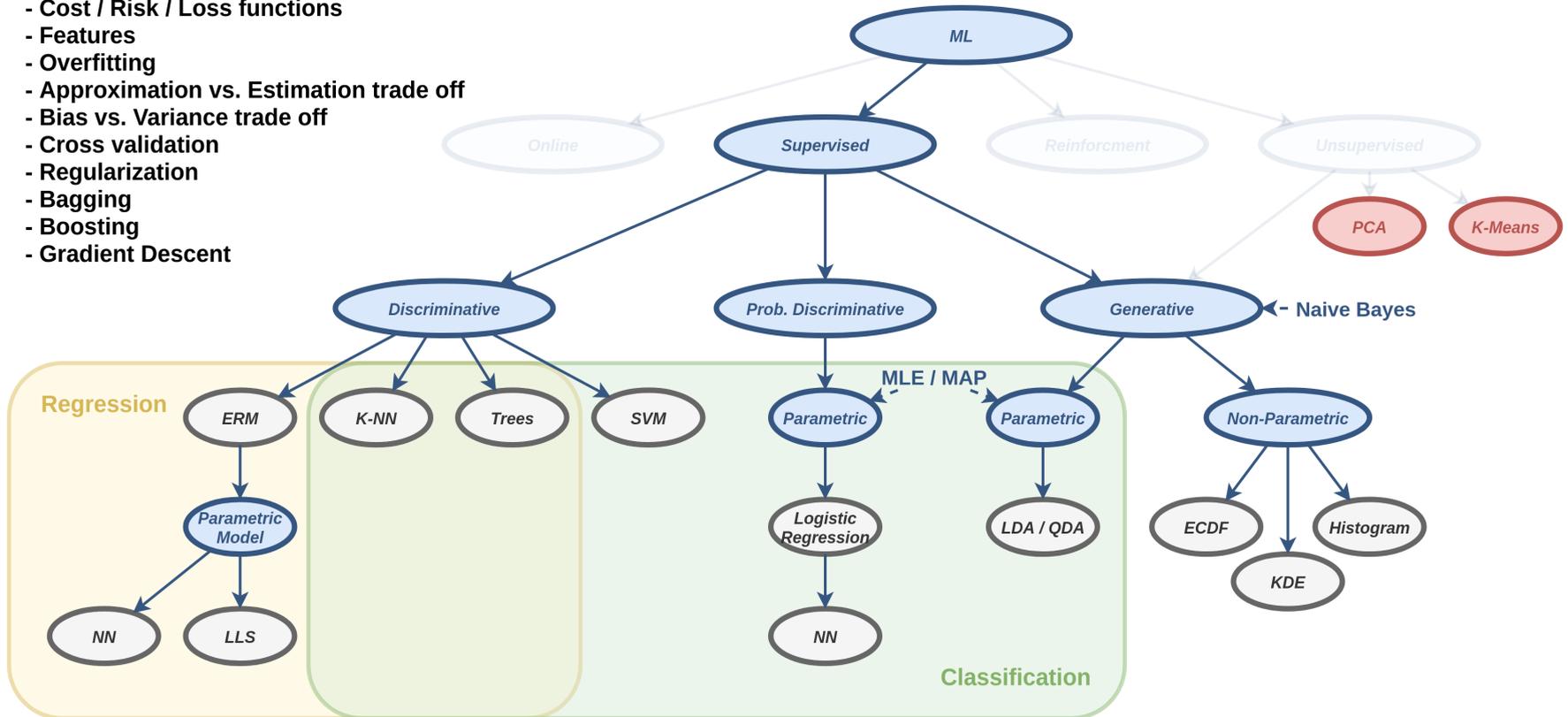


הרצאה 12 - PCA and K-Means

Subjects Covered in this Course

General concepts:

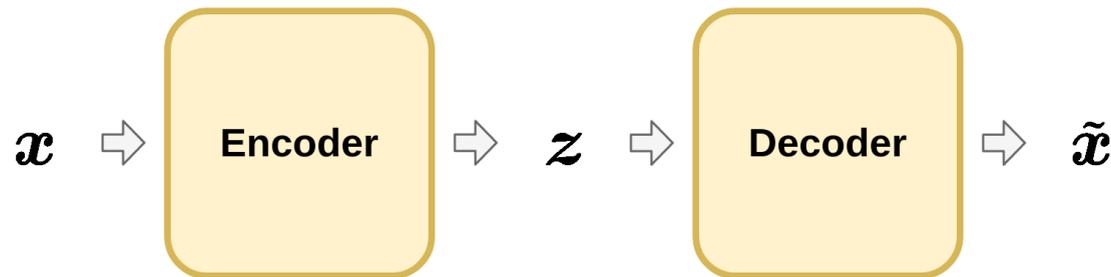
- Cost / Risk / Loss functions
- Features
- Overfitting
- Approximation vs. Estimation trade off
- Bias vs. Variance trade off
- Cross validation
- Regularization
- Bagging
- Boosting
- Gradient Descent



למידה לא מודרכת (Unsupervised Learning)

- שם כולל למגוון של בעיות בהם אנו מנסים בהינתן מדגם, ללמוד את התכונות של הדגימות או של המדגם כולו.
- המדגם יכול אוסף של דגימות (x) , ללא תווית (y) .
- דוגמאות:
 - אשכול (חלוקה לקבוצות).
 - מציאת ייצוג "נוח" יותר של הדגימות.
 - דחיסה.
 - זיהוי אנומליות.
 - למידת הפילוג של הדגימות.

מערכת Encoder-Decoder



דוגמאות לשימושים במערכת encoder-decoder הינם:

- דחיסה: נרצה ש z יהיה קטן ככל האפשר.
- תקשורת: נרצה ש z יהיה כמה שפחות רגיש לרעשים.
- הצפנה: נרצה שפעולת השחזור של x תהיה כמה שיותר קשה ללא ה decoder המתאים.
- \tilde{x} נקרא השחזור של x . בחלק מהמערכות ניתן להגיע לשיחזור מושלם, $\tilde{x} = x$, ובחלק מהמערכות לא.

ב PCA ננסה לבנות מערכת encoder-decoder שבה:

1. אנו מגבילים את האורך של הוקטור z .

2. אנו דורשים שה encoder וה decoder יהיו פונקציות אפיניות (affine = linear + offset).

3. התוחלת של שגיאת השחזור הריבועית $\mathbb{E} [\|\tilde{x} - x\|_2^2]$ היא מינימאלית.

נחליף את התוחלת בתוחלת אמפירית על מדגם.

D האורך של x ו K האורך של הוקטור z כאשר מתקיים כי $K \leq D$.

נרצה למצוא encoder:

$$z = T_1 x + b_1$$

ו decoder מהצורה של:

$$\tilde{x} = T_2 z + b_2$$

אשר ממצעים את התוחלת האמפירית של שגיאת השחזור הריבועית:

$$\arg \min_{T_1, T_2, b_1, b_2} \frac{1}{N} \sum_{i=1}^N \|\tilde{x}^{(i)} - x^{(i)}\|_2^2$$

דוגמאות למקרים שבהם נרצה לבצע הורדת מימד
:(dimensionality reduction)

1. בחירת מאפיינים לבעיות supervised learning
2. ויזואליזציה
3. דחיסה

הפתרון לבעיית האופטימיזציה

מסתבר שיש מספר רב של פתרונות. ניתן לבחור את הפרמטרים כך שיקיימו את האילוצים:

$$b_1 = -T_1 \mu$$

$$b_2 = \mu$$

$$T_1 = T_2^\top = T^\top$$

$$T^\top T = I$$

כאשר $\mu = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)}$.

הערה: שימו לב ש- $T \in \mathbb{R}^{D \times K}$ כך שמתקיים כי $T^\top T \in \mathbb{R}^{K \times K} =$

I_K כאשר I_K היא מטריצת היחידה. בנוסף, מתקיים $TT^\top \in \mathbb{R}^{D \times D}$ והיא לא שווה בהכרח ל- I_D .

הפתרון לבעיית האופטימיזציה

לדוגמה, עבור המיפוי הבא

$$z = T_1 x + b_1 \quad \tilde{x} = T_2 z + b_2 \quad T_1 \in \mathbb{R}^{K \times D}, T_2 \in \mathbb{R}^{D \times K}$$

איברי ההטיה b_1 ו- b_2 יכולים להיקבע ע"י הדרישות

$$E[z] = 0 \quad \Rightarrow \quad b_1 = -T_1 \mu$$

-I

$$E[\tilde{x}] = E[x] \quad \Rightarrow \quad b_2 = E[x] = \mu$$

הפתרון לבעיית האופטימיזציה

הטרנספורמציות במקרה זה הופכות להיות:

$$\mathbf{z} = T^\top (\mathbf{x} - \boldsymbol{\mu})$$

$$\tilde{\mathbf{x}} = T\mathbf{z} + \boldsymbol{\mu}$$

ובעיית האופטימיזציה הינה:

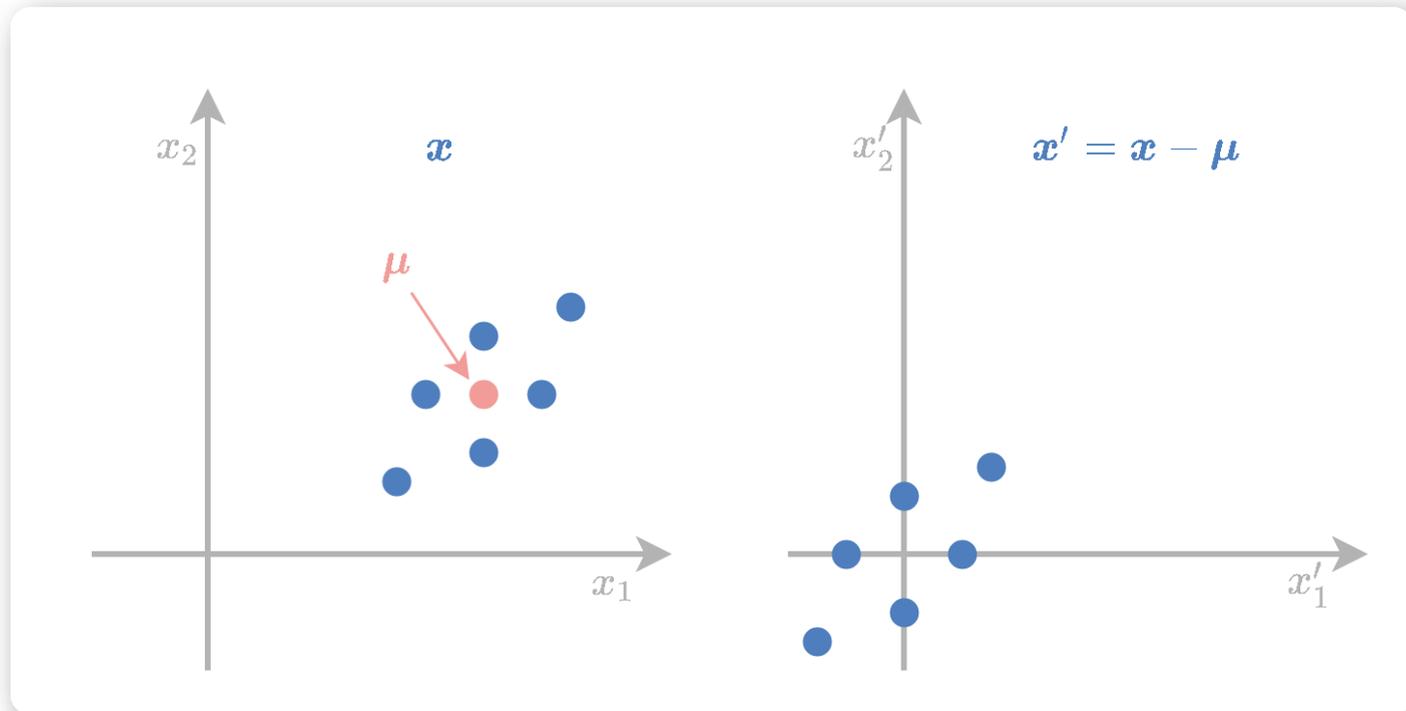
$$T^* = \arg \min_T \frac{1}{N} \sum_{i=1}^N \|\tilde{\mathbf{x}}^{(i)} - \mathbf{x}^{(i)}\|_2^2$$

s.t. $T^\top T = I$

$$T^* = \arg \min_T \frac{1}{N} \sum_{i=1}^N \|(TT^\top - I)(\mathbf{x}^{(i)} - \boldsymbol{\mu})\|_2^2$$

s.t. $T^\top T = I$

- ה encoder מחסר את הממוצע של x וה decoder מוסיף אותו בחזרה.
- נניח מעתה שהנתונים ממורכזים סביב האפס.



הפתרון לבעיית האופטימיזציה

הטרנספורמציות המתקבלות הינן:

$$\mathbf{z} = T^\top \mathbf{x}$$

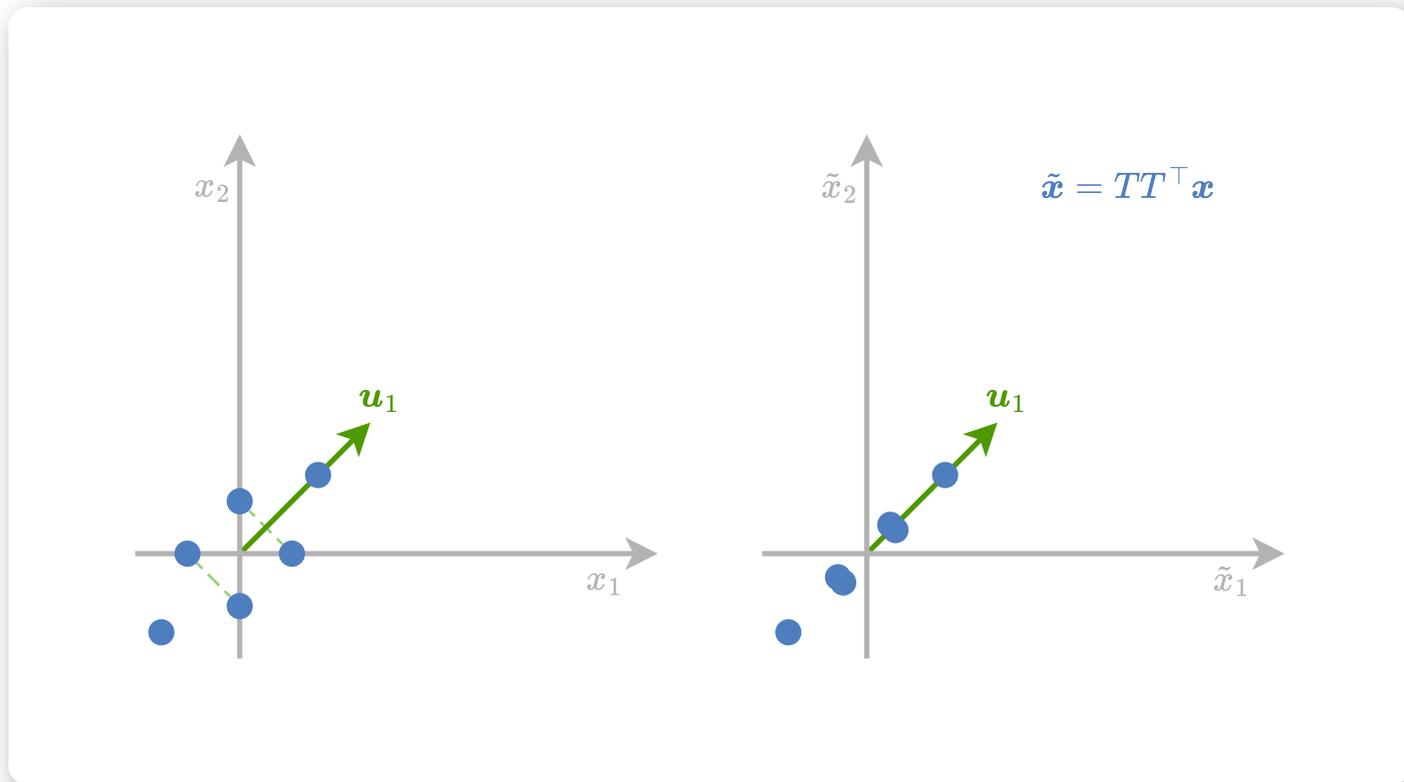
$$\tilde{\mathbf{x}} = T\mathbf{z} = TT^\top \mathbf{x}$$

נתייחס כעת לאילוץ של $T^\top T = I$. אילוץ זה אומר שהעמודות של T צריכות להיות אורתו-נורמאליות.

נסמן את העמודות של T ב \mathbf{u}_j :

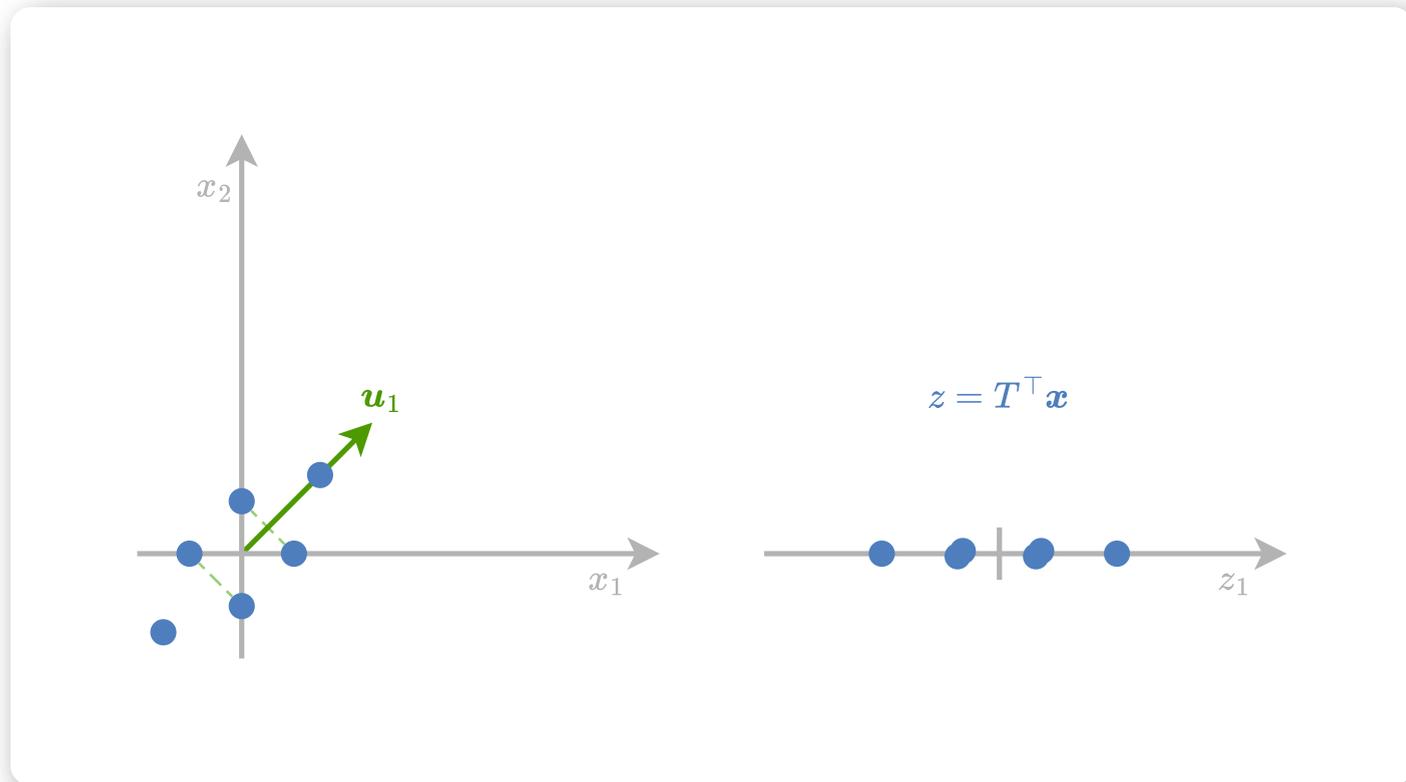
$$T = \begin{pmatrix} | & | & \dots & | \\ \mathbf{u}_1 & \mathbf{u}_2 & \dots & \mathbf{u}_K \\ | & | & & | \end{pmatrix}$$

הפעולה של $\tilde{x} = TT^T x$ מטילה את הוקטור x על תת-המרחב הליניארי הנפרס על ידי הוקטורים u_j .



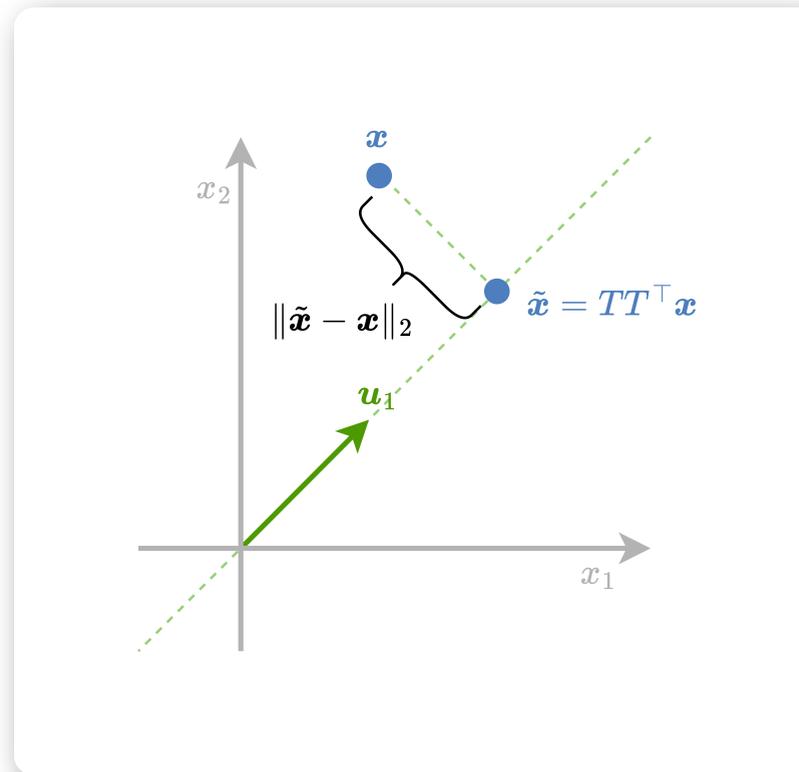
פרשנות גיאומטרית

הפעולה של $z = T^T x$ גם מטילה את x על אותו תת-מרחב, היא רק משאירה אותו במערכת הצירים של u_j :



נסתכל כעת על המשמעות הגיאומטרית של שגיאת השחזור

$$\|\tilde{x} - x\|_2^2$$



בעיית האופטימיזציה היא הבעיה של מציאת תת-המרחב ממימד K אשר ההטלה של נקודות המדגם עליו הם הקרובות ביותר לנקודות המקוריות.

מתוך העובדה ש $T^\top T = I$ ניתן להראות ש:

$$\|\tilde{\mathbf{x}} - \mathbf{x}\|_2^2 = \|\mathbf{x}\|_2^2 - \|\tilde{\mathbf{x}}\|_2^2 = \|\mathbf{x}\|_2^2 - \|\mathbf{z}\|_2^2$$

שכן, עבור $T^\top T = I$ מתקיים כי $(I - TT^\top)^2 = (I - TT^\top)$ לכן,

$$\begin{aligned}\|\mathbf{x} - \tilde{\mathbf{x}}\|_2^2 &= \|\mathbf{x} - TT^\top \mathbf{x}\|_2^2 \\ &= \|(I - TT^\top) \mathbf{x}\|_2^2 \\ &= \mathbf{x}^\top (I - TT^\top) \mathbf{x} \\ &= \|\mathbf{x}\|_2^2 - \|\mathbf{z}\|_2^2\end{aligned}$$

ובנוסף

$$\|\tilde{\mathbf{x}}\|_2^2 = \|T\mathbf{z}\|_2^2 = \mathbf{z}^\top T^\top T \mathbf{z} = \|\mathbf{z}\|_2^2$$

מכאן שנוכל לרשום את בעיית האופטימיזציה באופן הבא:

$$T^* = \arg \min_T \frac{1}{N} \sum_{i=1}^N \left(\|\mathbf{x}^{(i)}\|_2^2 - \|\mathbf{z}^{(i)}\|_2^2 \right)$$

s.t. $T^\top T = I$

נזכור ש $\|\mathbf{x}\|_2^2$ הוא תכונה של הוקטורים במדגם והם אינם תלויים ב T ולכן:

$$T^* = \arg \min_T - \frac{1}{N} \sum_{i=1}^N \|\mathbf{z}^{(i)}\|_2^2$$

s.t. $T^\top T = I$

$$T^* = \arg \min_T - \frac{1}{N} \sum_{i=1}^N \|z^{(i)}\|_2^2$$

s.t. $T^\top T = I$

הבעיה של מזעור שגיאת השחזור הריבועית שקולה לבעיה של מקסום הגודל $\sum_{i=1}^N \|z^{(i)}\|_2^2$.

גדול זה מכונה ה variance של אוסף הוקטורים $\{z^{(i)}\}_{i=1}^N$.

נגדיר:

$$X = \begin{pmatrix} - & \mathbf{x}'^{(1)} & - \\ - & \mathbf{x}'^{(2)} & - \\ & \vdots & \\ - & \mathbf{x}'^{(N)} & - \end{pmatrix} = \begin{pmatrix} - & (\mathbf{x}^{(1)} - \boldsymbol{\mu})^\top & - \\ - & (\mathbf{x}^{(2)} - \boldsymbol{\mu})^\top & - \\ & \vdots & \\ - & (\mathbf{x}^{(N)} - \boldsymbol{\mu})^\top & - \end{pmatrix}$$

ומטריצת ה covariance האמפירית של \mathbf{x} תהיה:

$$P = X^\top X$$

$$P = X^\top X$$

P ממשית וסימטרית ולכן מובטח כי ניתן לפרק אותה באופן הבא:

$$P = U \Lambda U^\top$$

כאשר U היא מטריצה הוקטורים עצמיים:

$$U = \begin{pmatrix} | & | & \dots & | \\ u_1 & u_2 & & u_D \\ | & | & & | \end{pmatrix}$$

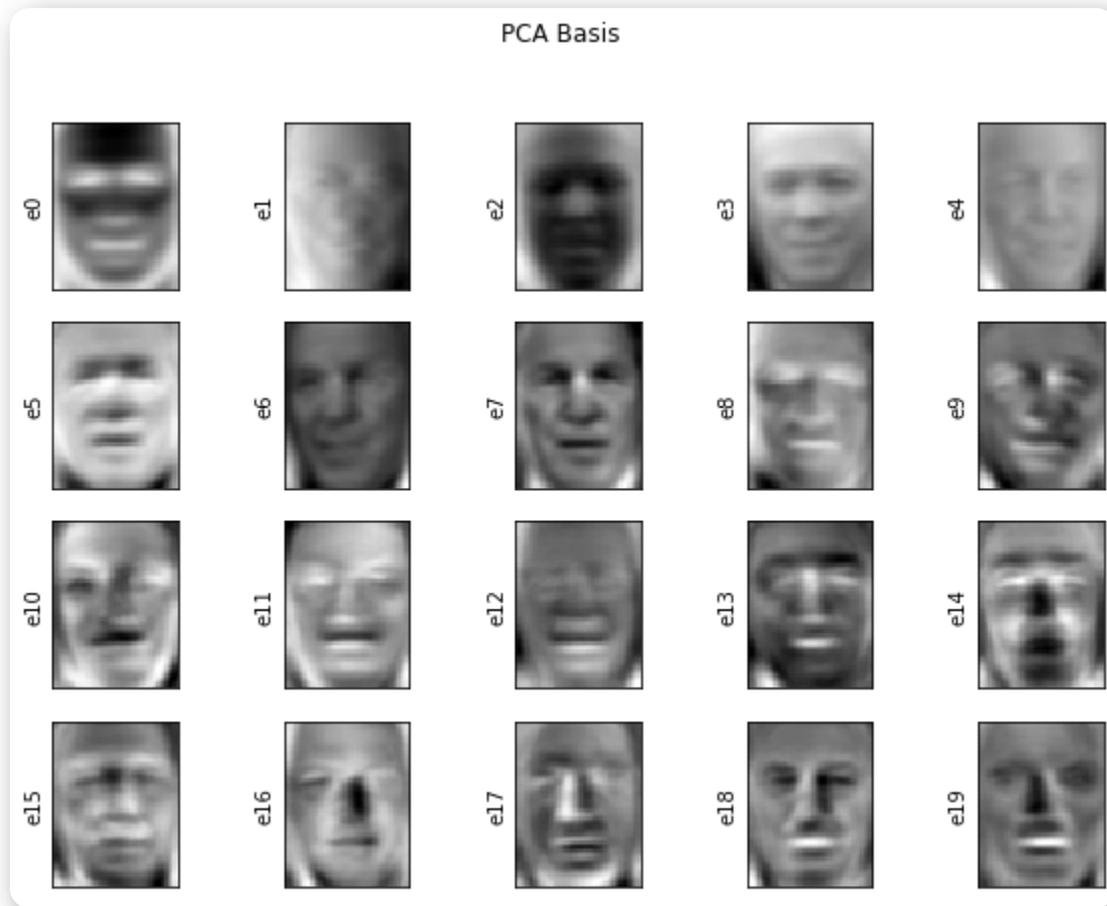
Λ היא מטריצה הערכים העצמיים:

T תהיה מטריצה אשר העמודות שלה הם K העמודות הראשונות במטריצה U :

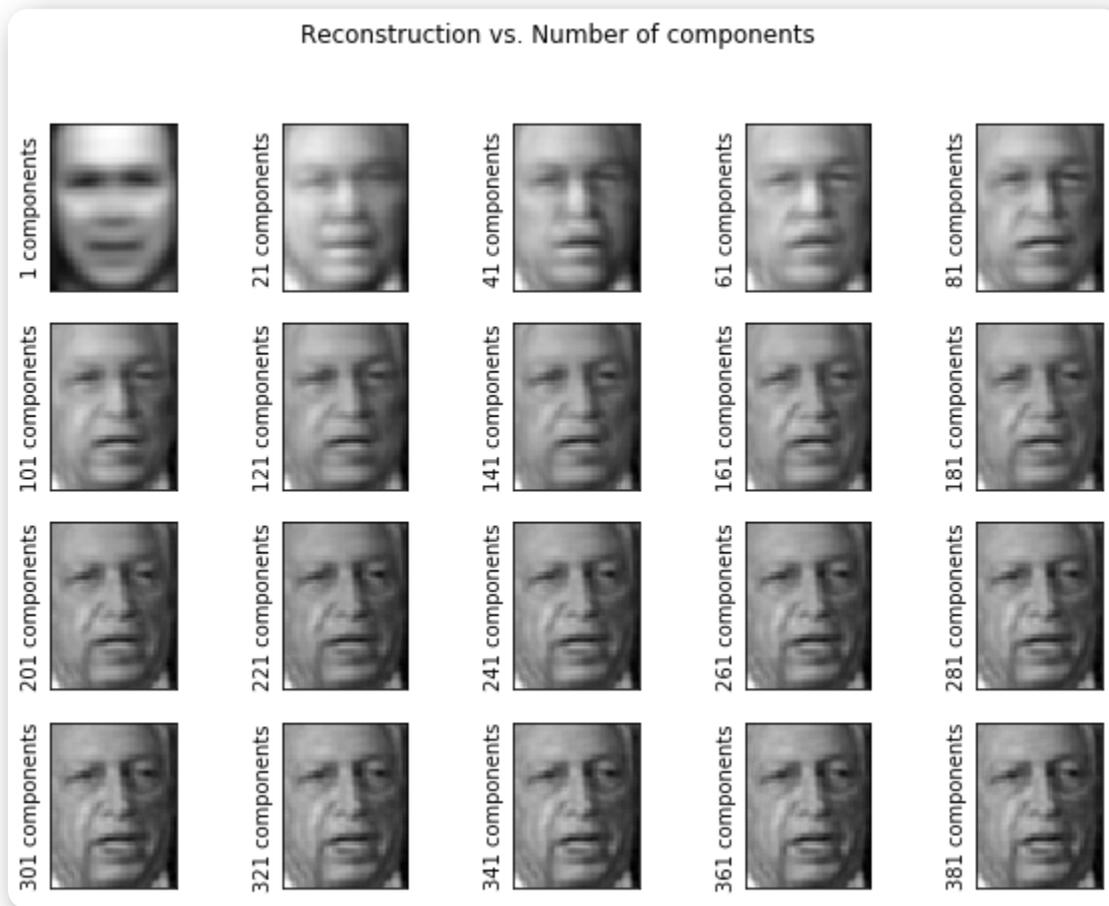
$$T = \begin{pmatrix} | & | & \dots & | \\ u_1 & u_2 & \dots & u_K \\ | & | & \dots & | \end{pmatrix}$$

- הכיוונים $u^{(j)}$ מכונים **כיוונים העיקריים**
- הרכיבים של הוקטור z מכונים **הרכיבים העיקריים (principal components)**.

פירוק תמונות של פנים לכיוונים העיקריים:



תמונה המשוחזרת עבור ערכים שונים של K :

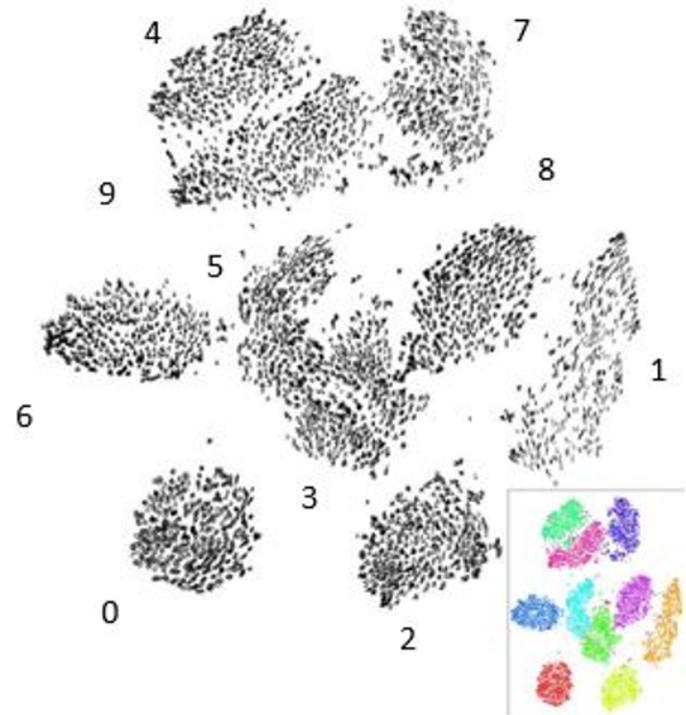


הרחבות לא לינאריות

קיימות הרחבות לא לינאריות רבות ל-PCA. נידונות בקורס עיבוד וניתוח מידע (ענ"ם).



0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9 9 9 9 9 9

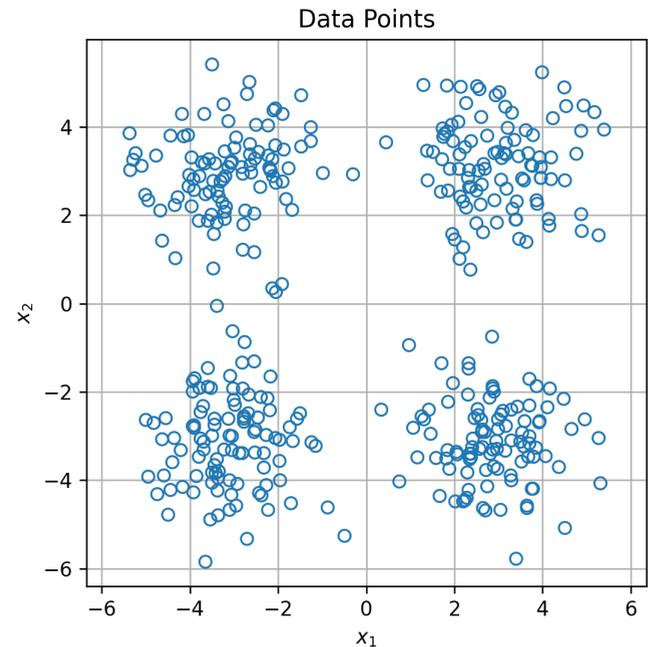
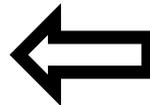
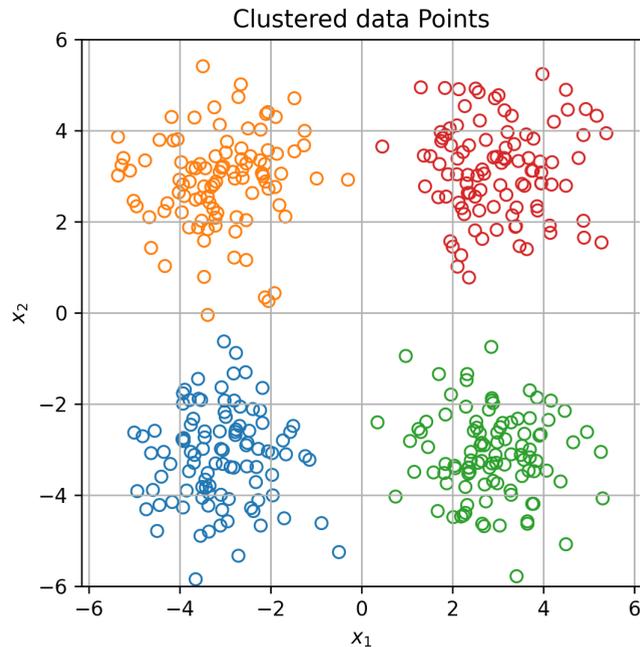


$$D = 28 \times 28 = 784$$

$$K = 2$$

הפעלת אלגוריתם tSNE על MNIST.

באלגוריתמי אשכול ננסה לחלק אוסף של פרטים לקבוצות המכונים אשכולות (clusters), כאשר לכל קבוצה איזשהן תכונות דומות. כמובן, בממדים גבוהים לא רואים זאת בעין.



2 דוגמאות למקרים שבהם נרצה לאשכול:

1. על מנת לבצע הנחות על אחד מהפרטים באשכול על סמך פרטים אחרים באשכול.
לדוגמא: להציע ללקוח מסויים בחנות אינטרנט מוצרים על סמך מוצרים שקנו לקוחות אחרים באשכול שלו.
2. לתת טיפול שונה לכל אשכול.
לדוגמא: משרד ממשלתי שרוצה להפנות קבוצות שונות באוכלוסיה לערוצי מתן שירות שונים: אפליקציה, אתר אינטרנט, נציג טלפוני או הפניה פיסית למוקד שירות.

K-Means הוא אלגוריתם אשכול אשר מנסה לחלק את הדגימות במדגם ל K קבוצות על סמך המרחק בין הדגימות.

סימונים

- K - מספר האשכולות (גודל אשר נקבע מראש).
- \mathcal{I}_k - אוסף האינדקסים של האשכול ה- k .
לדוגמא: $\mathcal{I}_5 = \{3, 6, 9, 13\}$
- $|\mathcal{I}_k|$ - גודל האשכול ה- k (מספר הפרטים בקבוצה)
- $\{\mathcal{I}_k\}_{k=1}^K$ - חלוקה מסוימת לאשכולות

K-Means מנסה למצוא את החלוקה לאשכולות אשר תמצער את המרחק הריבועי הממוצע בין כל דגימה לכל שאר הדגימות שאיתו באותו האשכול:

$$\arg \min_{\{\mathcal{I}_j\}_{k=1}^K} \frac{1}{N} \sum_{k=1}^K \frac{1}{2|\mathcal{I}_k|} \sum_{i,j \in \mathcal{I}_k} \|\mathbf{x}^{(j)} - \mathbf{x}^{(i)}\|_2^2$$

שאלה: האם פונקציית מרחק ריבועית תמיד מתאימה?

הבעיה השקולה

נגדיר את מרכז המסה:

$$\mu_k = \frac{1}{|\mathcal{I}_k|} \sum_{i \in \mathcal{I}_k} \mathbf{x}^{(i)}$$

ניתן להראות כי בעיית האופטימיזציה המקורית, שקולה לבעיה של מיזעור המרחק הממוצע של הדגימות ממרכז המסה של האשכול:

$$\arg \min_{\{\mathcal{I}_j\}_{k=1}^K} \frac{1}{N} \sum_{k=1}^K \sum_{i \in \mathcal{I}_k} \|\mathbf{x}^{(i)} - \mu_k\|_2^2$$

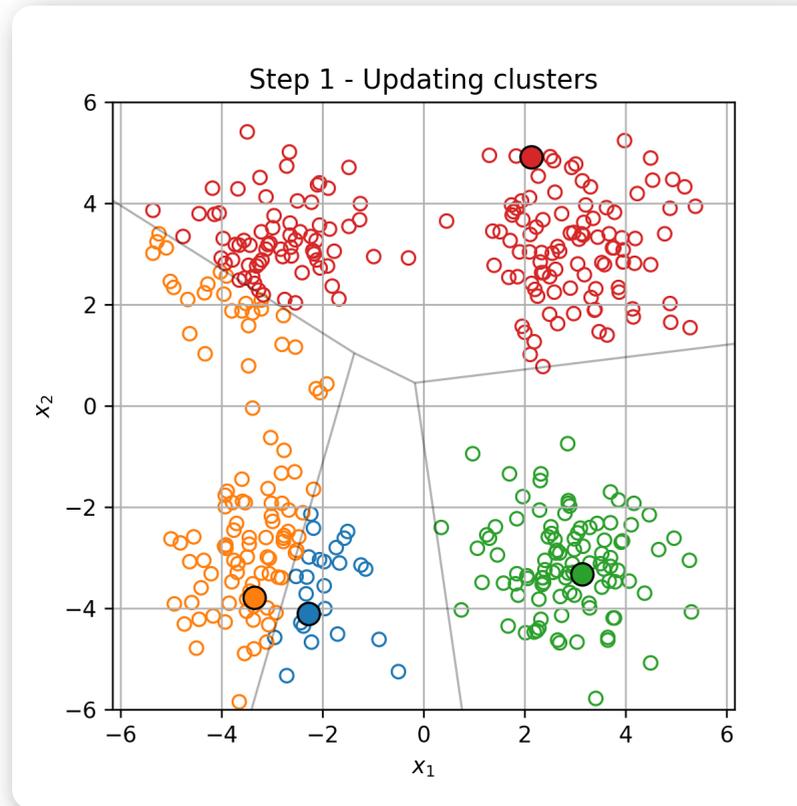
$$\begin{aligned} \sum_{i,j \in \mathcal{I}_k}^K \left\| \mathbf{x}^{(i)} - \mathbf{x}^{(j)} \right\|_2^2 &= \sum_{i,j \in \mathcal{I}_k} \left\| \mathbf{x}^{(i)} - \boldsymbol{\mu}_k + \boldsymbol{\mu}_k - \mathbf{x}^{(j)} \right\|_2^2 \\ &= \sum_{i,j \in \mathcal{I}_k} \left\| \mathbf{x}^{(i)} - \boldsymbol{\mu}_k \right\|_2^2 + \sum_{i,j \in \mathcal{I}_k} \left\| \mathbf{x}^{(j)} - \boldsymbol{\mu}_k \right\|_2^2 - 2 \sum_{i,j \in \mathcal{I}_k} \left(\mathbf{x}^{(i)} - \boldsymbol{\mu}_k \right)^\top \left(\mathbf{x}^{(j)} - \boldsymbol{\mu}_k \right) \\ &= 2 |\mathcal{I}_k| \sum_{i \in \mathcal{I}_k} \left\| \mathbf{x}^{(i)} - \boldsymbol{\mu}_k \right\|_2^2 - 2 \sum_{i \in \mathcal{I}_k} \left(\mathbf{x}^{(i)} - \boldsymbol{\mu}_k \right)^\top \sum_{j \in \mathcal{I}_k} \left(\mathbf{x}^{(j)} - \boldsymbol{\mu}_k \right) \\ &= 2 |\mathcal{I}_k| \sum_{i \in \mathcal{I}_k} \left\| \mathbf{x}^{(i)} - \boldsymbol{\mu}_k \right\|_2^2 \end{aligned}$$

שכן:

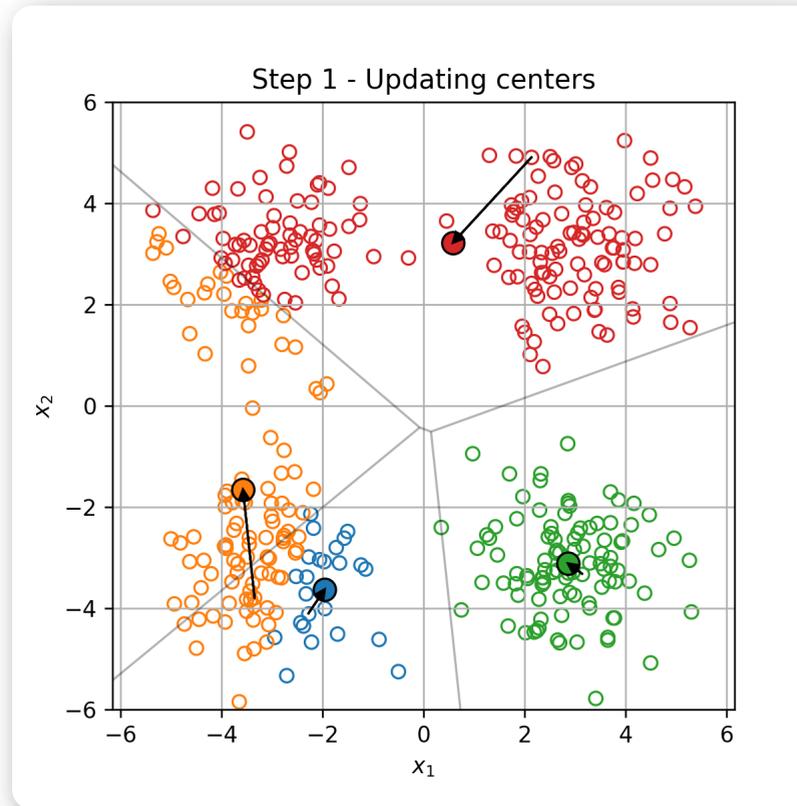
$$\sum_{i \in \mathcal{I}_k} \left(\mathbf{x}^{(i)} - \boldsymbol{\mu}_k \right) = |\mathcal{I}_k| \cdot \frac{1}{|\mathcal{I}_k|} \sum_{i \in \mathcal{I}_k} \mathbf{x}^{(i)} - |\mathcal{I}_k| \boldsymbol{\mu}_k = 0$$

- מובטח כי פונקציית המטרה תקטן בכל צעד.
- מובטח כי האלגוריתם יעצר לאחר מספר סופי של צעדים.
- **לא** מובטח כי האלגוריתם יתכנס לפתרון האופטימאלי. בפועל במרבית מתכנס לפתרון קרוב מאד לאופטימאלי.
- אתחולים שונים יכולים להוביל לתוצאות שונות.

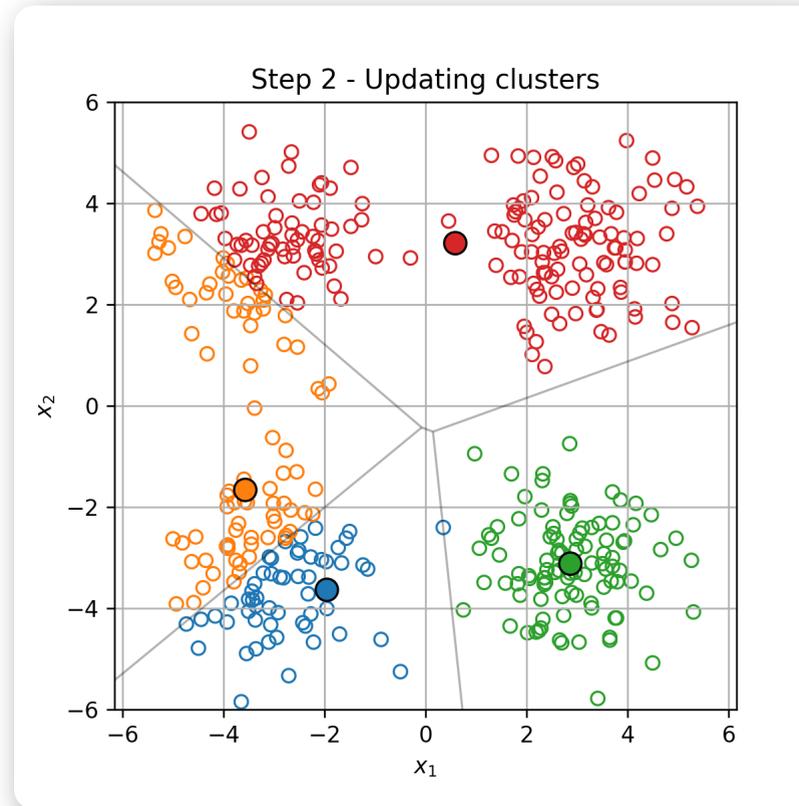
אתחול (וחלוקה ראשונית לאשכולות):



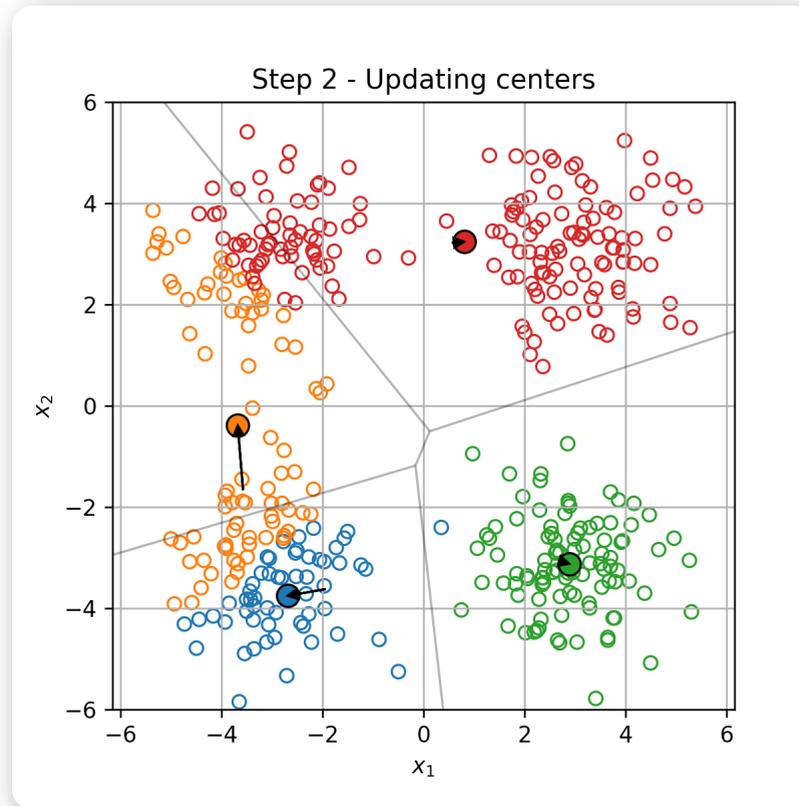
עדכון המרכזים:



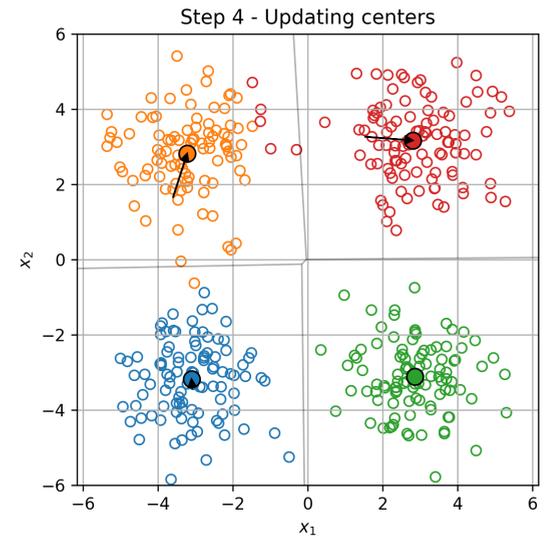
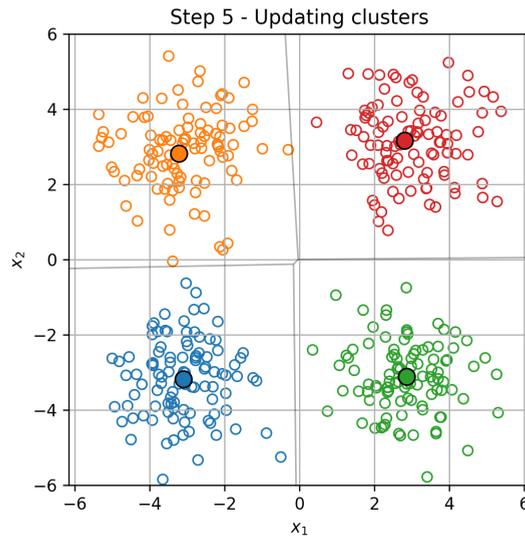
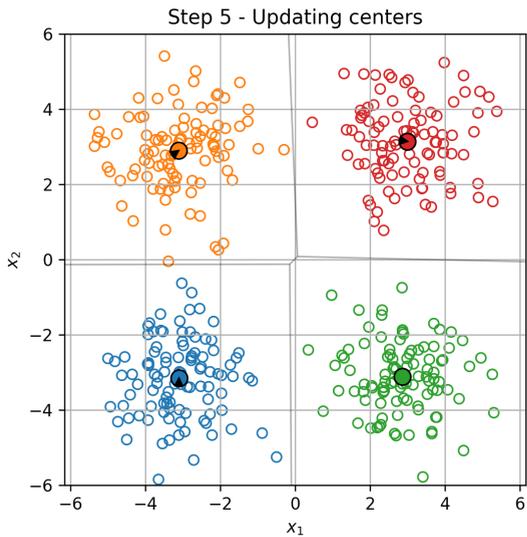
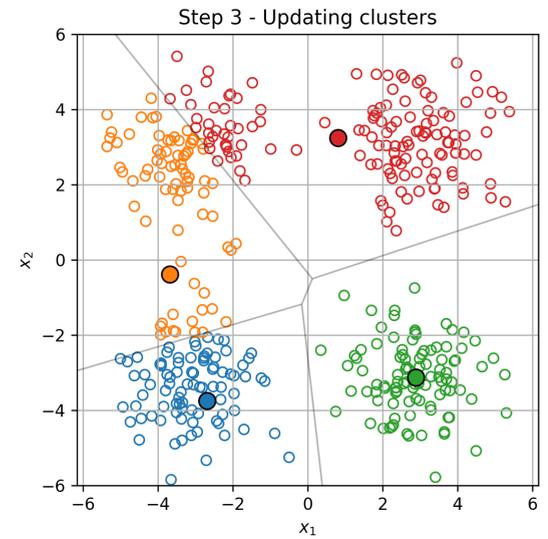
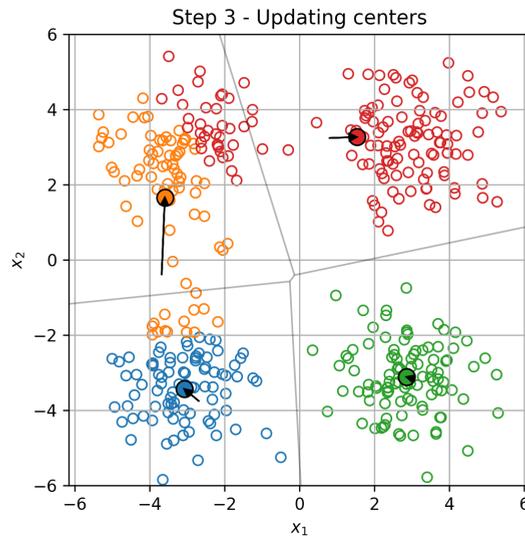
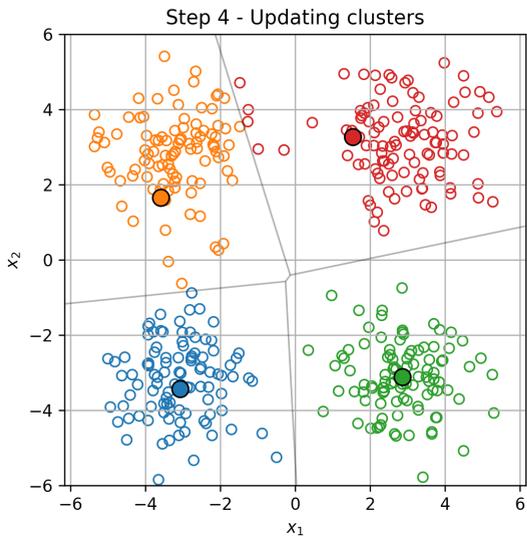
עדכון האשכולות:



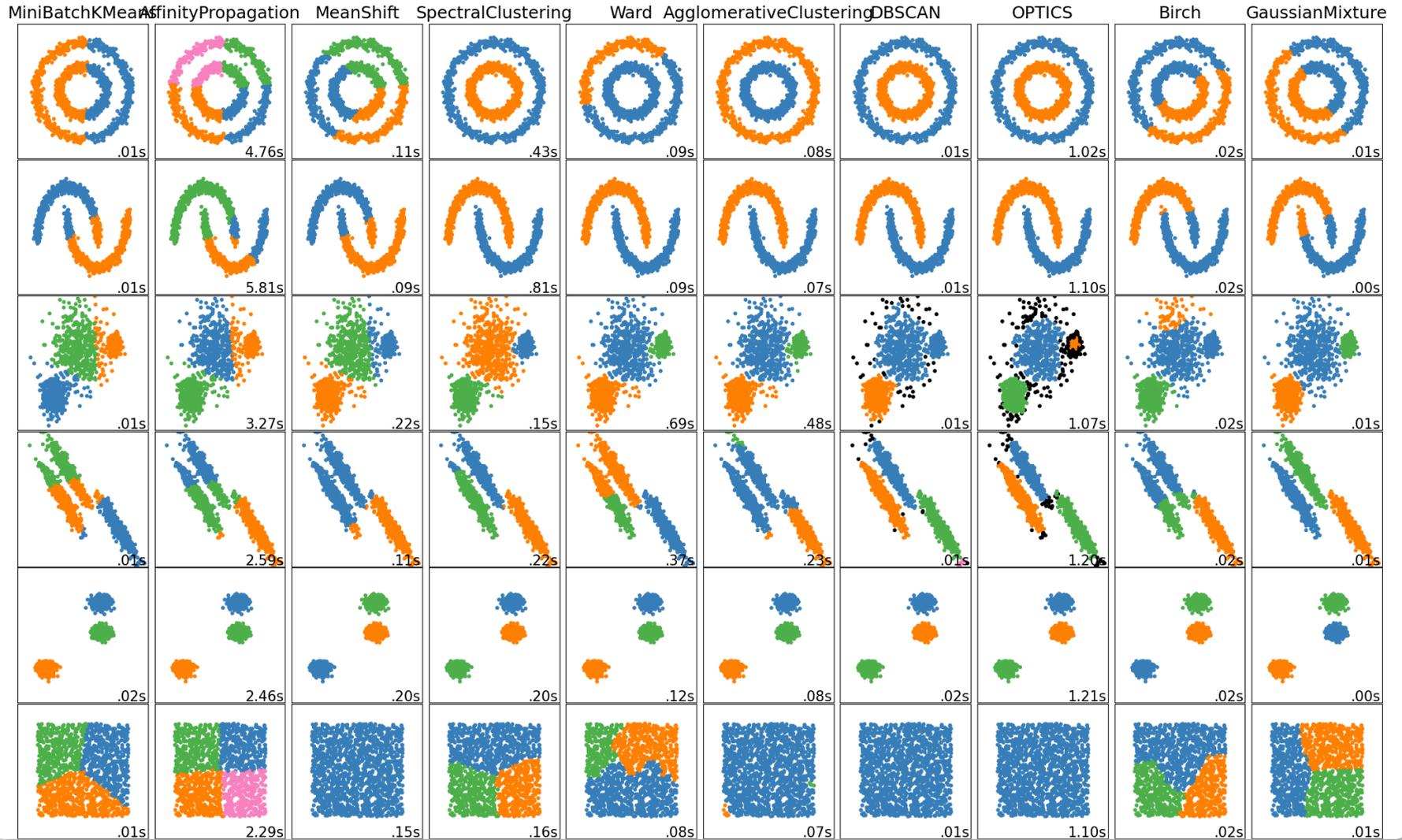
עדכון המרכזים:



וּחֹזֵר חֲלִילָה (הַסֵּדֵר הוּא מִיָּמִין לַשְּׂמָאל):



אלגוריתמי אשכול שונים



לרוב לא נוכל לצייר את האשכולות בשני ממדים.