

# זיהוי תבניות

## זיהוי תבניות

לדוגמא: זיהוי פונמה; בהינתן אות דיבור, צ"ל איזו פונמה זו.

נפריד (ככל האפשר) את בעיית הסיווג לשתי בעיות נפרדות:

- חילוץ מאפיינים (feature extraction):

פעולה הממירה כל תצפית לווקטור במרחב המאפיינים.

- סיווג (classification):

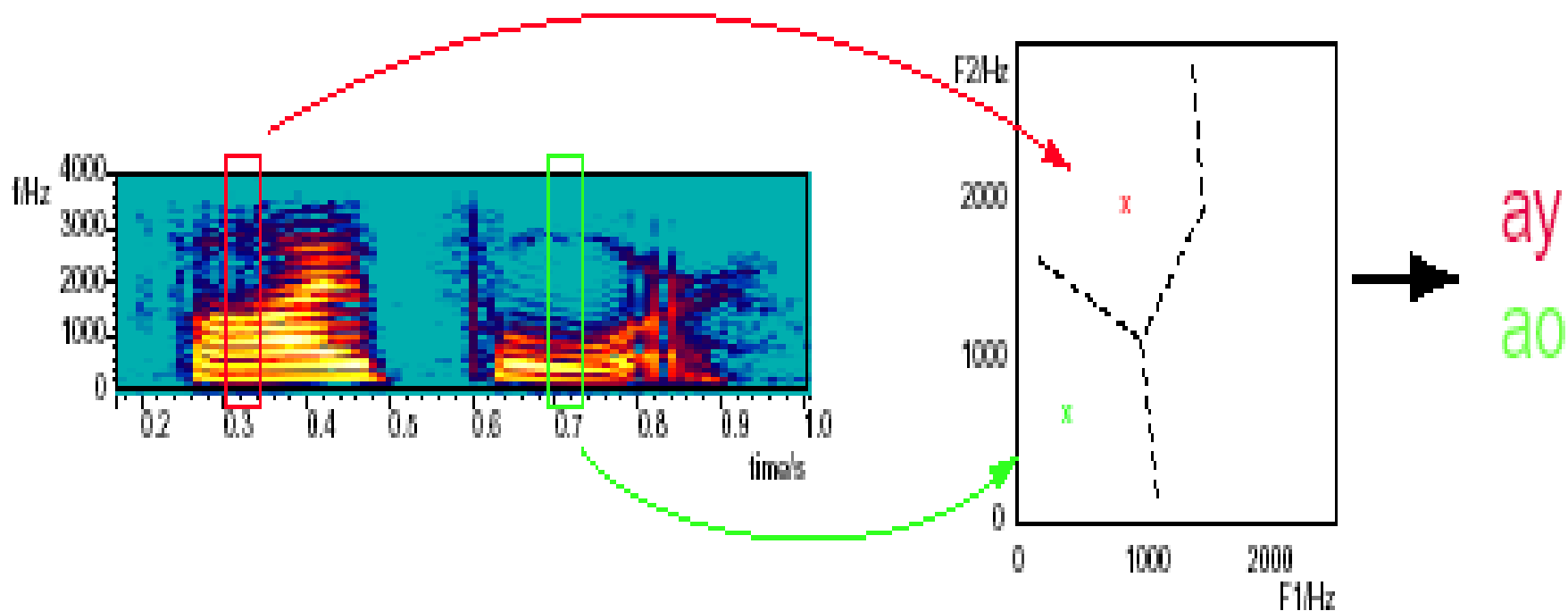
פעולה הממירה כל ווקטור במרחב המאפיינים למחלקה (

class

המתאימה ביותר.

# סיווג: מציאת תיוג דיסקרטי לתצפיות

## רציפות



## בניית מסווג

### הגדרת מחלקות (classes) – פלט רצוי

לדוגמא: הפונמות שונות.

### הגדרת מרחב המאפיינים (feature space) – קלט נתון

לדוגמא: התדרים של הפורמנט הראשון והשני.

### הגדרת אלגוריתם ההחלטה – הפונקציה המתאימה לכל פלט קלט

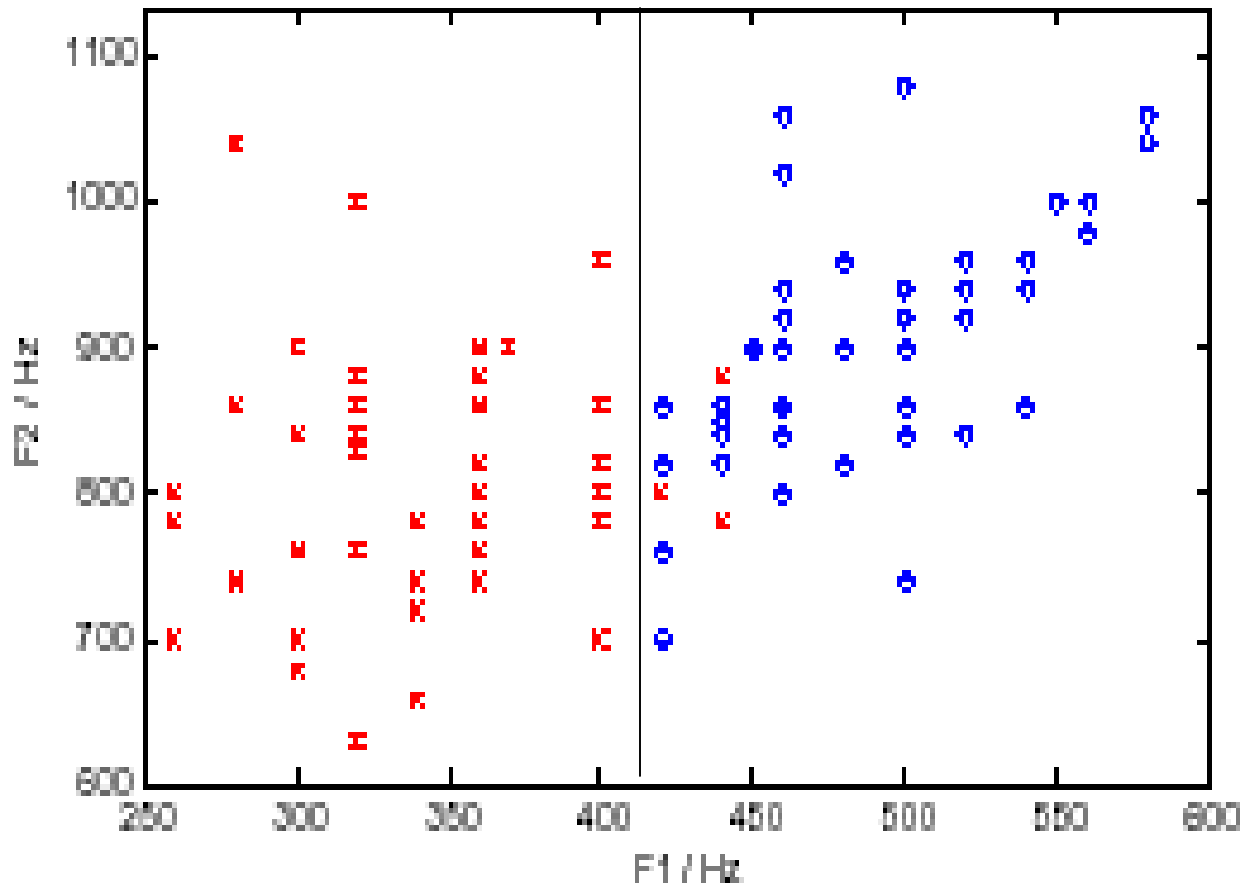
- אימון (עפ"י קבוצת אימון).

- סיווג.

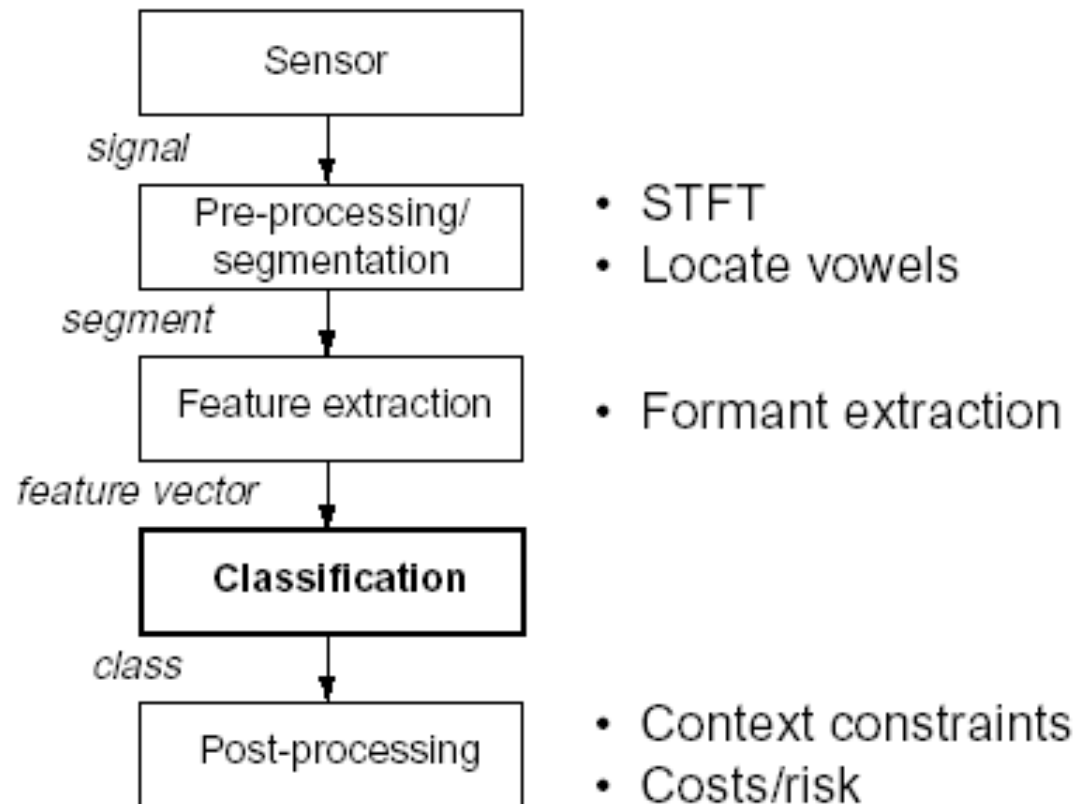
### מדידת ביצועים

4 מדידת אחוז שגיאה עבור קבוצת מבחן.

Pols vowel formants: "u" (x) vs. "o" (o)



# Classification systems



## חילוץ מאפיינים – Feature Extraction

בחירת המאפיינים היא קריטית להצלחת הסיווג.

מה מגדיר מאפיין טוב ?

- מכיל את מירב האינפורמציה הרלוונטית.
  - מכיל את האינפורמציה הרלוונטית בצורה פשוטה (תחת מטריקת המרחב).
  - איזוריאנטי תחת טרנספורמציות לא רלוונטיות (תחת מטריקת המרחב).
  - לא רגיש לרעש (תחת מטריקת המרחב).
  - קל (או לפחות אפשרי) לחישוב.
- מאפיינים אפשריים: האות המקורי, הספקטרום, פורמנטים.

## חילוץ מאפיינים – קללת המימדיות

**Curse of dimensionality**: ככל שמימד מרחב המאפיינים גבוה יותר, כמות דוגמאות האימון הדרושה גדולה יותר (באופן אקספוננציאלי במימד).

מימד מרחב המאפיינים קטן מאפשר:

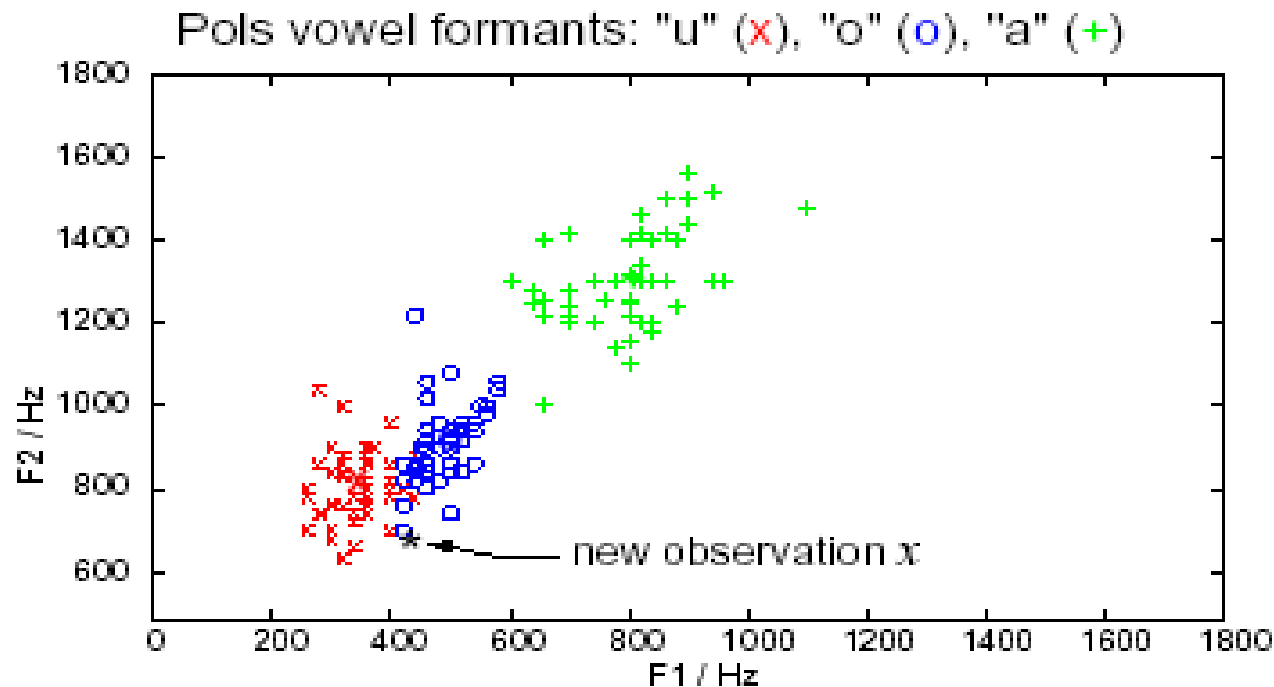
- מודלים פשוטים וקטנים יותר.
- פחות דוגמאות אימון דרושות.
- אימון מהיר יותר.

מסווג פשוט: Minimum distance classification.

האלגוריתם: מצא את דוגמת האימון הקרובה ביותר במרחב.

בחירת המטריקה חשובה ביותר.

הקטנת דרישות CPU וזכרון: השוואה רק עם נציגים מקבוצת האימון.

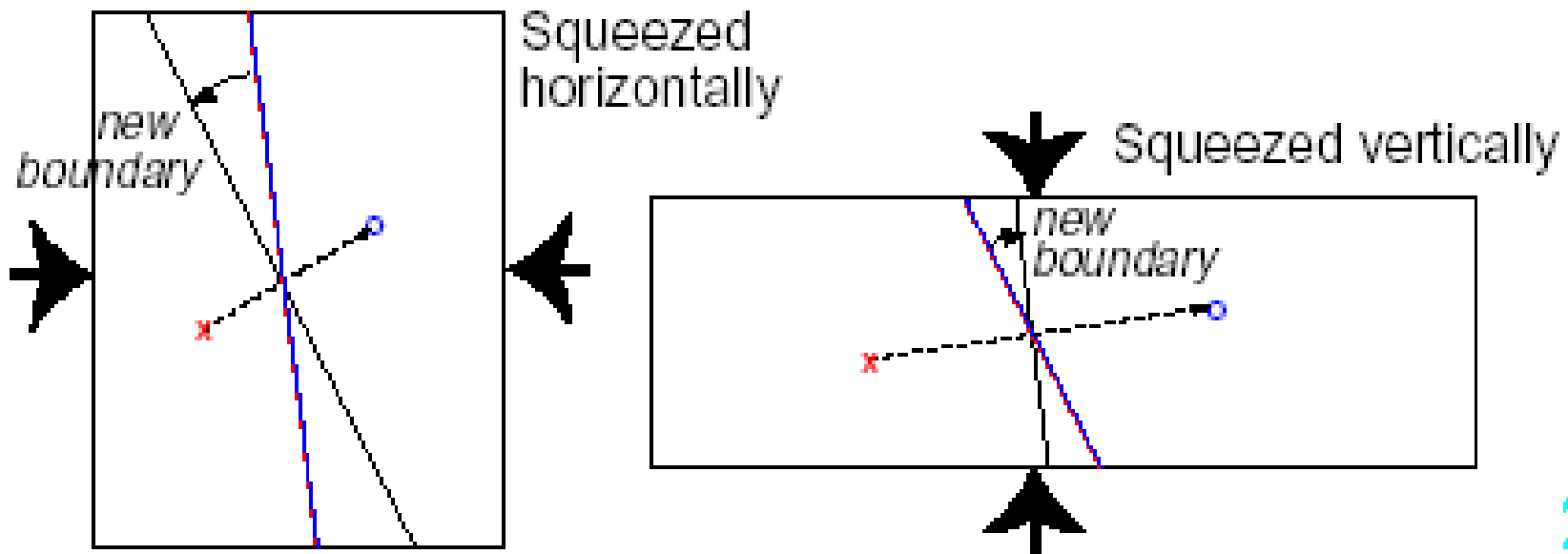


חסרונות: רגישות לרעש, אין הכללה. פתרון אפשרי: k-nearest neighbor

מטריקה אפשרית: מרחק אוקלידי.

חסרון מהותי: המרחק האוקלידי רגיש לטרנספורמציות מתיחה (scaling).

- **Scaling axes changes boundary:**



## קללת המימדיות: גישות לפתרון

הגישה הכללית: הורדת מימד תוך שמירה מירבית על האנפורציה הרלוונטית.

ההבדלים בין השיטות השונות:

- הגדרת המושג "אינפורמציה רלוונטית".
- השיטות להורדת המימד (בדר"כ טרנספורמציה ליניארית).

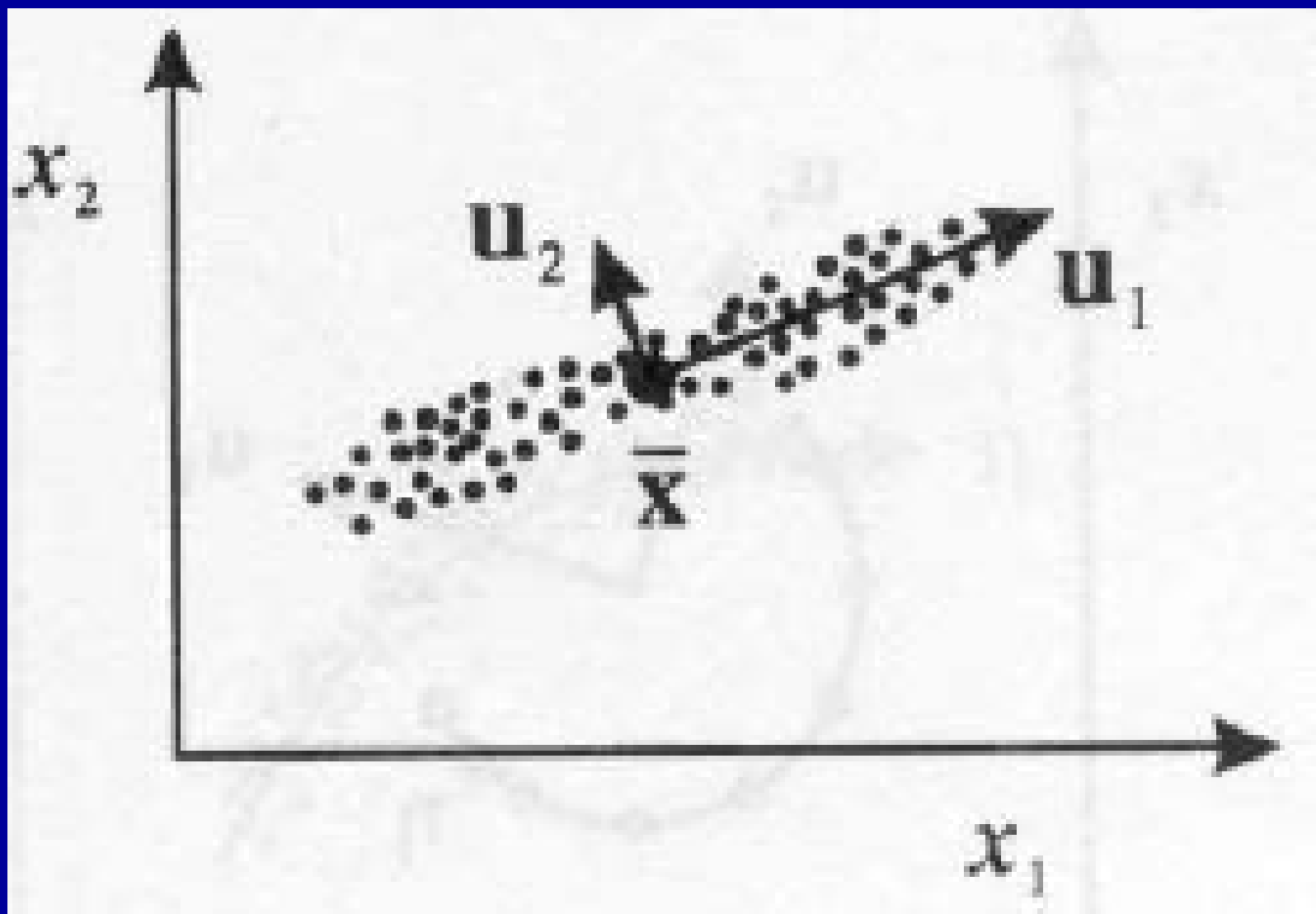
## Principal Component Analysis – PCA

### מוטיבציה:

בהינתן מרחב  $m$ -מימדי, נרצה למיין את המימדים השונים לפי סדר "החשיבות", ונוותר על  $m$  המימדים ה-"פחות חשובים".

כדי שהתהליך יהיה יעיל, נרצה תחילה להחליף את הבסיס של המרחב לבסיס אחר כך ש- $m$  המימדים שנוותר עליהם יהיו "לא חשובים" ככל האפשר.

התהליך של החלפת הבסיס מתבצע ע"י מציאת ווקטורים עצמיים למטריצת הקובריאנס של ווקטורי האימון, וקביעת הבסיס כאוסף הווקטורים העצמיים.



נתונים הווקטורים  $x_1, \dots, x_n$  במרחב ווקטורי  $d$ -מימדי.

נרצה למצוא טרנספורמציה אורתונורמלית  $W$  המעבירה ווקטור  $x_i$  ל-  $y_i = Wx_i$ , המקיימת תנאי אופטימליות עבור הקריטריון הבא:

לכל ווקטור  $y_i$  נמחק את המימדים  $m+1, \dots, d$  ונציב במקומם קבועים  $b_{m+1}, \dots, b_d$ , כלומר:

$$x_i \rightarrow y_i = \sum_{j=1}^d z_{i,j} u_j, \quad y'_i = \sum_{j=1}^m z_{i,j} u_j + \sum_{j=m+1}^d b_j u_j$$

כאשר  $u_j = We_j$ ,  $\{e_j\}$  הבסיס סטנדרטי.

נחשב את שגיאת הקירוב עבור  $x_i$ :

$$\text{err}_i = x_i - y'_i = \sum_{j=m+1}^d (z_{i,j} - b_j) u_j$$

נחשב את ממוצע ריבוע השגיאה:

$$E = \frac{1}{n} \sum_{i=1}^n \|x_i - y'_i\|^2 = \frac{1}{n} \sum_{i=1}^n \left\| \sum_{j=m+1}^d (z_{i,j} - b_j) u_j \right\|^2 = \frac{1}{n} \sum_{i=1}^n \sum_{j=m+1}^d (z_{i,j} - b_j)^2$$

נמצא את ערכי  $b_j$  המביאים את  $E$  למינימום:

$$\frac{\partial E}{\partial b_j} = 0$$

$$\frac{\partial E}{\partial b_j} = -2 \sum_{i=1}^n (z_{i,j} - b_j)$$

$$b_j = \frac{1}{n} \sum_{i=1}^n z_{i,j} = \frac{1}{n} \sum_{i=1}^n u_j^t x_i = u_j^t \bar{x}$$

כלומר  $b_j$  הנו הערך הממוצע עבור קואורדינטה  $j$  במרחב החדש.

כעת נמצא את  $W$ :

$$E = \frac{1}{n} \sum_{i=1}^n \|x_i - y'_i\|^2 = \frac{1}{n} \sum_{i=1}^n \sum_{j=m+1}^d (z_{i,j} - b_j)^2 = \frac{1}{n} \sum_{i=1}^n \sum_{j=m+1}^d (u_j^t x_i - u_j^t \bar{x})^2 =$$
$$= \sum_{j=m+1}^d u_j^t \Sigma u_j$$

$$\Sigma = \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^t$$

$W$  הרצוי יביא למינימום את  $E$ .

פתרון (ללא הוכחה):  $\{u_j\}$  הם הווקטורים העצמיים של  $\Sigma$ .

לפיכך, השגיאה תהיה:

$$E = \sum_{j=m+1}^d u_j^t \Sigma u_j = \sum_{j=m+1}^d u_j^t \lambda_j u_j = \sum_{j=m+1}^d \lambda_j u_j^t u_j = \sum_{j=m+1}^d \lambda_j$$

לפיכך נוותר על  $d-m$  הווקטורים העצמיים בעלי הערכים-עצמיים הנמוכים ביותר.

### לסיכום, האלגוריתם הוא כדלהלן:

קלט: סידרת ווקטורים  $x_i$ .

פלט: סידרת ווקטורים  $z_i$  שהיא היטל של הקלט על מרחב במימד נמוך יותר.

1. נחשב את ווקטור ממוצע הקלט  $\leftarrow Ex$ .

2. נחשב את מטריצת הקובריאנס  $\Sigma$ .

3. נמצא את הווקטורים העצמיים של  $\Sigma$ , ונמיינם עפ"י ערכיהם העצמיים.

4. נבחר את  $m$  ( $m < d$ ) הווקטורים העצמיים בעלי ע"ע מקסימליים –  $\{u_j\}$ .

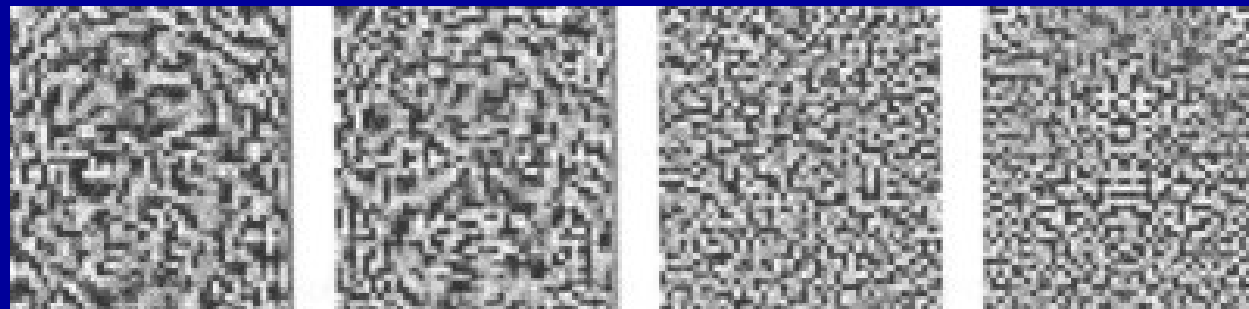
5. לכל ווקטור  $x_i$  נחשב את ההיטל שלו על המרחב הנפרט  $z_i = [(x_i^t - Ex)u_1 \dots (x_i^t - Ex)u_m]^t$ .



הפנים הממוצעות (בצד שמאל) וארבעת ה-Eigenfaces הראשונים



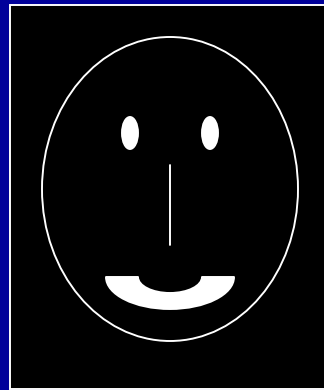
Eigenfaces 15,100,200,250,300 (מצד שמאל)



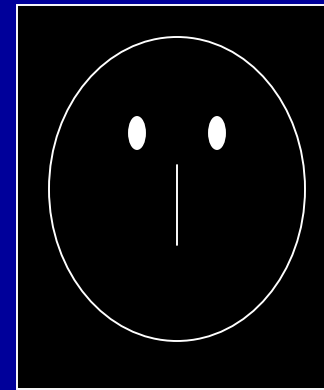
אינטרפטציה לשימוש ב-PCA (eigenfaces) עבור זיהוי פנים:

שיטת הרכבת קלסטרון ע"י המשטרה.

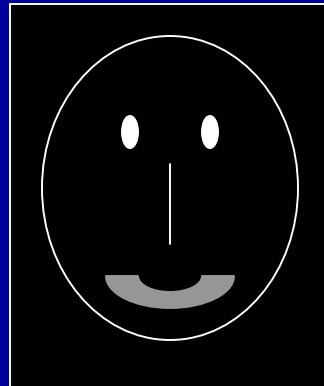
$$x_1 =$$



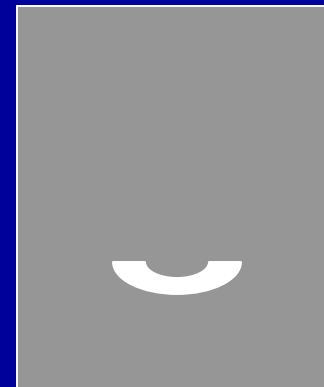
$$x_2 =$$



$$Ex = (x_1 + x_2) / 2 =$$



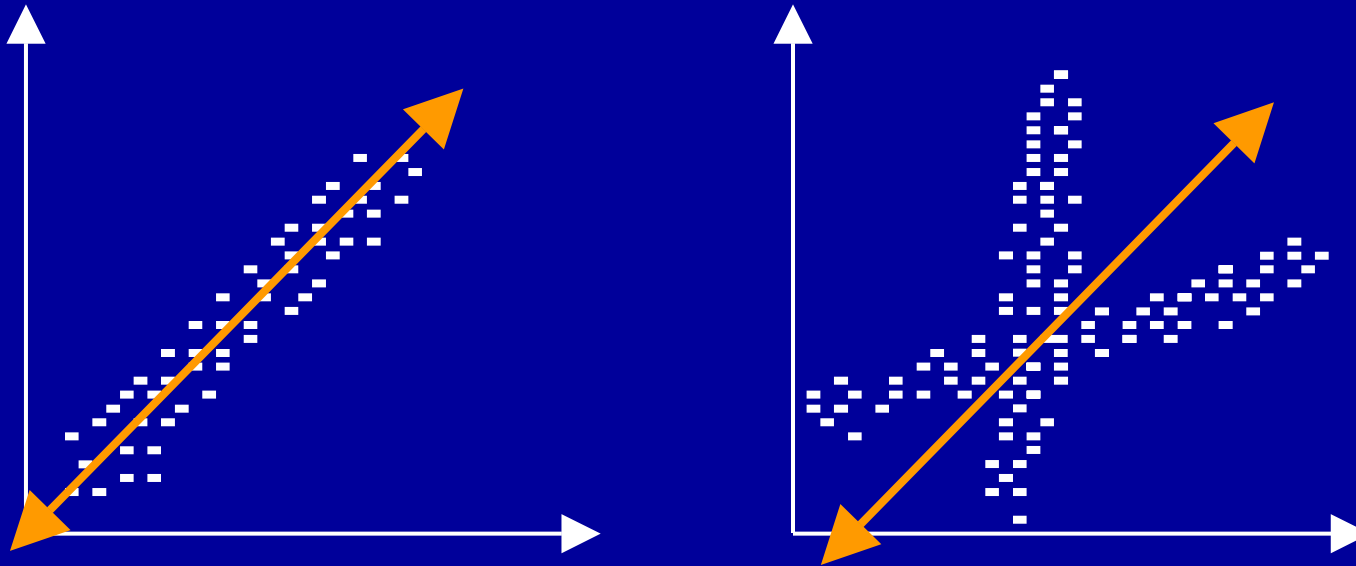
$$u_1 =$$



## מגבלות שיטת PCA:

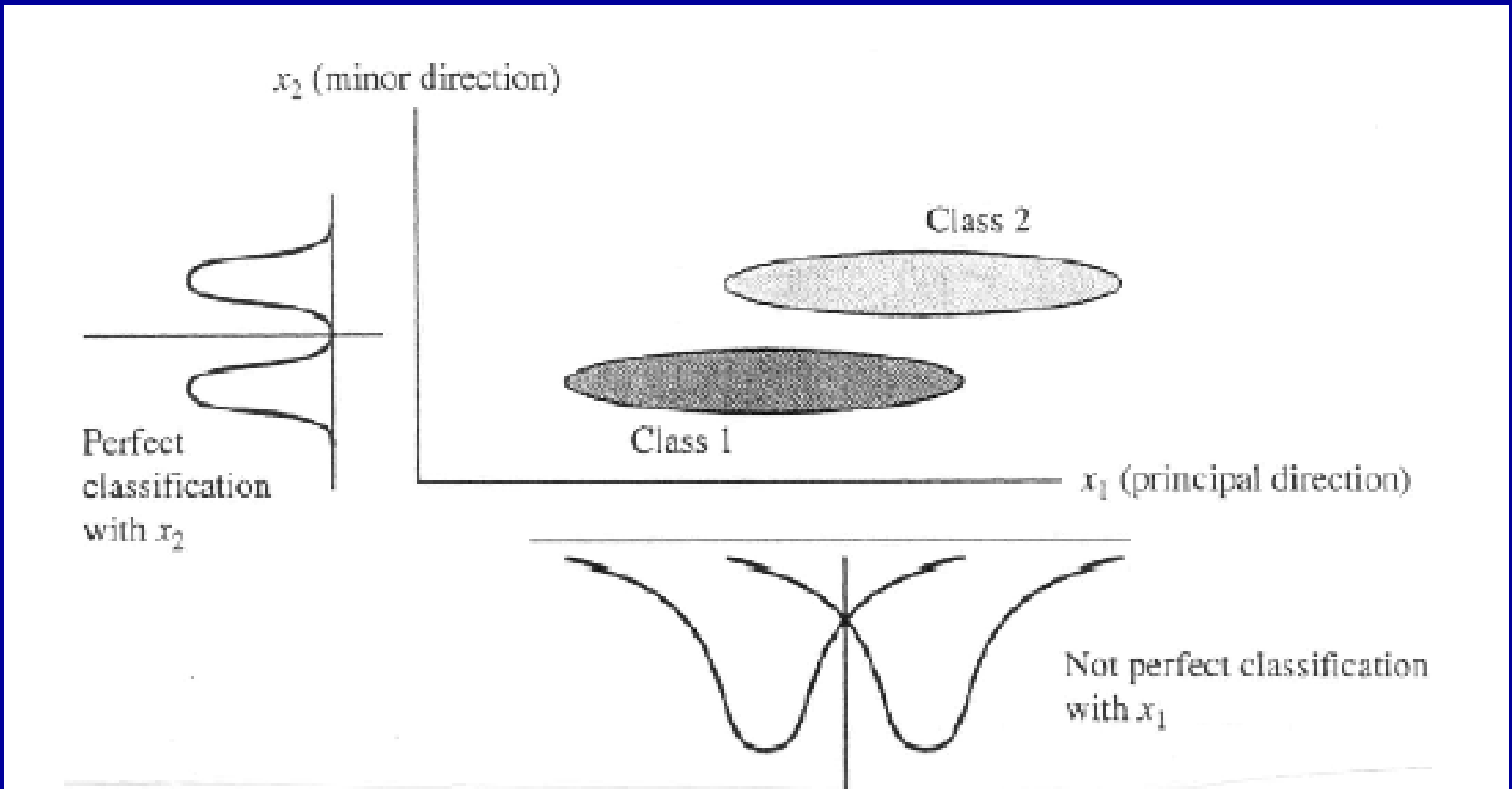
1. PCA הנה שיטה ליניארית, ואינה יכולה לאפיין תת-מרחבים לא ליניאריים. (לדוגמא קואורדינטות פולריות).

2. PCA אינה מתאימה לתאר צבירים.



3. PCA רגישה לנרמול הצירים השונים.

4. PCA אינה מנסה להפריד בין מחלקות שונות.



# זיהוי תבניות: מסווגים בייסיאנים

# סיווג בייסיאני

- חוק Bayes
- בחירה בין השערות שונות – ML, MAP
- מסווג בייסיאני נאיבי
- למידה פרמטרית / א-פרמטרית
- למידת MAP, ML

# חוק בייס (Bayes)

$$\Pr(w|X) = \frac{\Pr(X|w) \Pr(w)}{\Pr(X)}$$

$\Pr(w)$ : הסתברות אפריורית של השערה  $w$ .

$\Pr(X)$ : הסתברות אפריורית של התצפיות  $X$ .

$\Pr(w|X)$ : הסתברות מותנית של השערה  $w$  בהינתן תצפיות  $X$ .

$\Pr(X|w)$ : הסתברות מותנית של תצפיות  $X$  בהינתן ההשערה  $w$ .

# בחירת השערות

$$\Pr(w|X) = \frac{\Pr(X|w) \Pr(w)}{\Pr(X)}$$

נרצה לבחור את ההשערה המסתברת ביותר בהינתן התצפיות:

.Maximum a posteriori (MAP)

$$w_{\text{MAP}} = \operatorname{argmax}_{w \in \{W\}} \Pr(w|X) = \operatorname{argmax}_{w \in \{W\}} \frac{\Pr(X|w) \Pr(w)}{\Pr(X)} = \operatorname{argmax}_{w \in \{W\}} \Pr(X|w) \Pr(w)$$

הגדרה: נראות (likelihood) היא ההסתברות של התצפיות בהינתן המודל  $w$ .

לעיתים נניח כי  $\Pr(w)$  מתפלג אחיד, ואז נבחר את ההשערה בעלת

$$w_{\text{ML}} = \operatorname{argmax}_{w \in \{W\}} \Pr(X|w) \quad \text{: (Maximum Likelihood)}$$

## דוגמא לבחירת השערות

האם לנבדק יש סרטן?

נבדק מבצע בדיקת מעבדה, והבדיקה יוצאת חיובית.

סטטיסטיקה ידועה על הבדיקה:

• אם יש סרטן, 98% שהבדיקה תהיה חיובית.

• אם אין סרטן, 97% שהבדיקה תהיה שלילית.

בנוסף: ל- 0.8% מן האוכלוסייה יש בסרטן.

### נגדיר:

$w$  – יש סרטן / אין סרטן.

$X$  - תוצאת הבדיקה.

$$\Pr(w = \text{cancer}) = 0.008,$$

$$\Pr(w = \text{!cancer}) = 0.992$$

$$\Pr(X=\text{positive} \mid w=\text{cancer}) = 0.98, \quad \Pr(X=\text{negative} \mid w = \text{cancer}) = 0.02$$

$$\Pr(X=\text{positive} \mid w=\text{!cancer})= 0.03, \quad \Pr(X=\text{negative} \mid w = \text{!cancer})= 0.97$$

## דוגמא לבחירת השערות

$$W_{\text{MAP}} = \operatorname{argmax}_{w \in \{W\}} \Pr(X|w)\Pr(w)$$

$$\Pr(X=\text{positive} \mid w=\text{cancer}) * \Pr(w=\text{cancer}) = 0.98 * 0.008 = 0.00784$$

$$\Pr(X=\text{positive} \mid w=\text{!cancer}) * \Pr(w=\text{!cancer}) = 0.03 * 0.992 = 0.02976$$

$$W_{\text{MAP}} = \text{!cancer}$$

$$\Pr(w=\text{cancer} \mid X=\text{positive}) = 0.00784 / (0.00784 + 0.02976) = 20.9\%$$

$$\Pr(w=\text{!cancer} \mid X=\text{positive}) = 0.02976 / (0.00784 + 0.02976) = 79.1\%$$

... אבל 20.9% היא הסתברות מספיק גבוהה בשביל לנקוט בפעולות נוספות.

← כדאי לשקלל את ה"עלות" של קבלת כל השערה.

$$W_{\text{MAP-cost}} = \operatorname{argmax}_{w \in \{W\}} \Pr(w|X) * \text{cost}(\text{decision} = !w \mid w)$$

נניח שהעלות של בחירה "cancer" בטעות היא 1,  
והעלות של בחירה "!cancer" בטעות היא 10.

$$\begin{aligned} & \Pr(X=\text{positive} \mid w=\text{cancer}) * \Pr(w=\text{cancer}) * \text{cost}(\text{decision} = !\text{cancer} \mid \text{cancer}) \\ & = 0.98 * 0.008 * 10 = 0.0784 \end{aligned}$$

$$\begin{aligned} & \Pr(X=\text{positive} \mid w=! \text{cancer}) * \Pr(w=! \text{cancer}) * \text{cost}(\text{decision} = \text{cancer} \mid ! \text{cancer}) \\ & = 0.03 * 0.992 * 1 = 0.02976 \end{aligned}$$

$$W_{\text{MAP-cost}} = \text{cancer}$$

## נוסחאות רלוונטיות בהסתברות

$$1. \Pr(A, B) = \Pr(A|B) \Pr(B) = \Pr(B|A) \Pr(A)$$

$$2. \Pr(A \vee B) = \Pr(A) + \Pr(B) - \Pr(A, B)$$

$$3. \Pr(B) = \sum_{i=1}^n \Pr(B|A_i) \Pr(A_i) \quad \text{if } \Pr(A_i \cap A_j) = 0 \quad \forall i \neq j,$$
$$\Pr(\cup A_i) = 1$$

# Naïve Bayes מודל

בהינתן ווקטור תצפיות  $X=x_1, x_2, \dots, x_n$

הנחת "naïve bayes" - קיימת אי תלות בין התצפיות השונות.

$$\Pr(x_1, \dots, x_n | w) \cong \prod_{i=1}^n \Pr(x_i | w)$$

כך שמקבלים:

$$W_{NB} = \operatorname{argmax}_{w \in \{W\}} \Pr(w) \prod_{i=1}^n \Pr(x_i | w)$$

## נחזור לדוגמא עם בדיקת הסרטן:

נניח שישנה בדיקה נוספת ( $Y$ ) בלתי-תלויה שניתן לבצע.

הנתונים לגבי הבדיקה:

$$\Pr(Y=\text{positive} \mid w=\text{cancer}) = 0.5, \quad \Pr(Y=\text{negative} \mid w = \text{cancer}) = 0.5$$

$$\Pr(Y=\text{positive} \mid w=\!\text{cancer}) = 0.1, \quad \Pr(Y=\text{negative} \mid w = \!\text{cancer}) = 0.9$$

אפשר לראות שהבדיקה  $Y$  פחות מובהקת מן הבדיקה  $X$ , אך אי-תלותן משפרת את יכולת החיזוי.

נניח שהבדיקה  $Y$  אף היא חיובית עבור הנבדק שלנו.

$$\Pr(X, Y=\text{positive} \mid w=\text{cancer}) * \Pr(w=\text{cancer}) = 0.98 * 0.5 * 0.008 = 0.00392$$

$$\Pr(X, Y=\text{positive} \mid w=\!\text{cancer}) * \Pr(w=\!\text{cancer}) = 0.03 * 0.1 * 0.992 = 0.002976$$

$$W_{\text{MAP}} = \text{cancer}$$

$$\Pr(w=\text{cancer} \mid X, Y=\text{positive}) = 0.00392 / (0.00392 + 0.00298) = 56.8\%$$

$$\Pr(w=\!\text{cancer} \mid X, Y=\text{positive}) = 0.00298 / (0.00392 + 0.00298) = 43.2\%$$

חגי אהרונוביץ

## מסקנה:

אם נבצע הרבה בדיקות בלתי תלויות, נוכל להגיע ליכולת חיזוי גבוהה יותר.

## ביקורת:

1. ההסתברות של סרטן באוכלוסייה  $\Pr(\text{cancer})=0.008$  פחות רלוונטית מההסתברות של סרטן בקרב אוכלוסיית הנבדקים (שהיא כנראה גבוהה יותר).
2. האם אכן הבדיקות בלתי תלויות?

# שערוך מודל

- במציאות, ההתפלגויות אינן ידועות מראש.
- יש צורך ללמוד אותן.
- כאשר המשתנים הם בדידים (כמו בדוגמת ה-cancer), יש צורך לחשב אוסף סופי של הסתברויות.
- לדוגמא, נרצה לדעת מהו  $P(\text{cancer})$ .  
בהינתן מדגם (קבוצת אימון) של  $n$  אנשים, ידוע על  $m$  מתוכם החולים בסרטן.  
נניח התפלגות בינומית.

# שערוך מודל

השיטה הבייסיאנית נותנת מענה אופטימלי אך דורשת מידע מוקדם על כל ההתפלגויות הרלוונטיות.

במציאות ההתפלגויות אינן ידועות מראש – צריך לשערך (estimate) אותן.

לפעולה זו של שערוך ההסתברויות נקרא גם שערוך מודל, או "אימון מודל" (model training, model estimation).

את האימון נבצע באמצעות "דוגמאות אימון".

בקורס נדבר בעיקר על למידה מונחית: כאשר נתון הסיווג של דוגמאות האימון.

## גישות לשערוך פרמטרים

### הגישה הפרמטרית:

1. נניח שפונקצית ההתפלגות היא בעלת צורת התפלגות מוכרת (בינומית, נורמלית, אקספוננציאלית).
2. ננסה "לנחש" את הפרמטרים של ההתפלגות.

### הגישה הא-פרמטרית:

1. לא נניח (כמעט) דבר בנוגע לצורת פונקצית ההפלגות.
2. נלמד את פונקצית ההתפלגות ישירות מדוגמאות האימון.

## הגישה הפרמטרית

נניח שרוצים לשערך את ההתפלגות  $P(x|w)$ .

לדוגמא  $w = \{\text{male, female}\}$  ו- $x = \text{height}$ .

נציג את  $P(x|w)$  כך:

$$P(x|w_i) = f(x, \Theta_i)$$

לדוגמא:  $\Theta = \{\mu, \sigma^2\}$

$$P(x|w_i) = \frac{1}{\sqrt{2\pi\sigma_i}} e^{-\frac{(x-\mu_i)^2}{2\sigma_i^2}}$$

כל מה שנותר הוא לשערך את ערך הפרמטרים  $\Theta$ .

## הגישה הא-פרמטרית

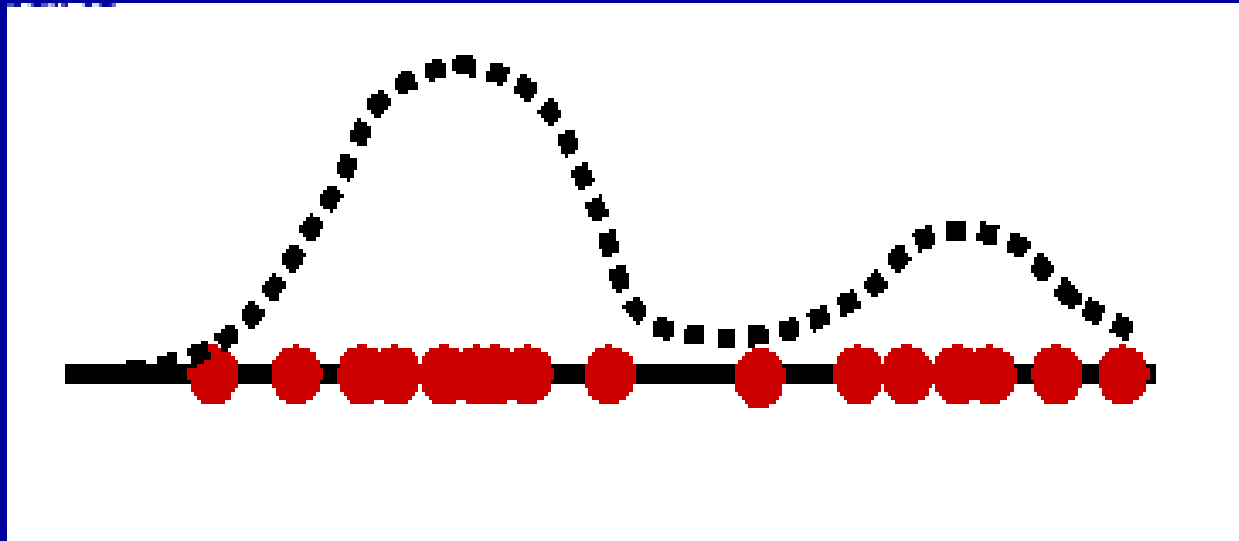
נניח שלא ידועה לנו אופי פונקצית ההתפלגות  $P(x|w)$ .

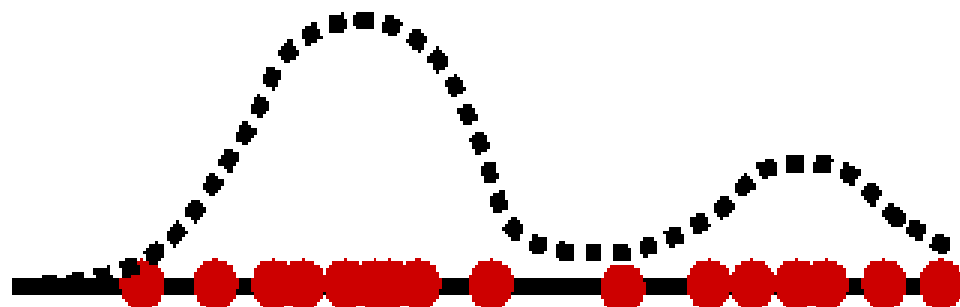
### שיטה מס' 1: שיטת ההיסטוגרמה

נניח ש- $X$  הוא חד מימדי.

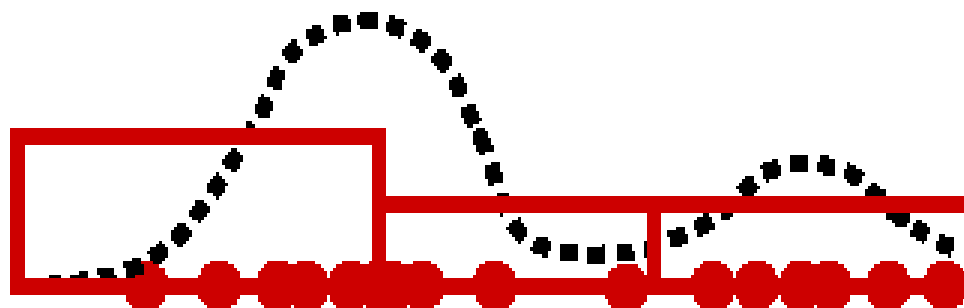
1. נחלק את הטווח של ערכי  $X$  ל- $M$  תחומים שווים.

2. נתייחס להתפלגות של  $X$  כהטלת מטבע עם  $M$  ערכים אפשריים.

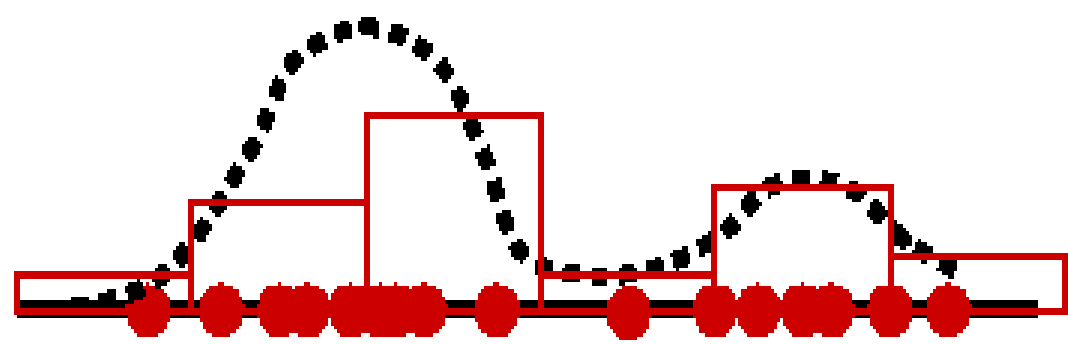




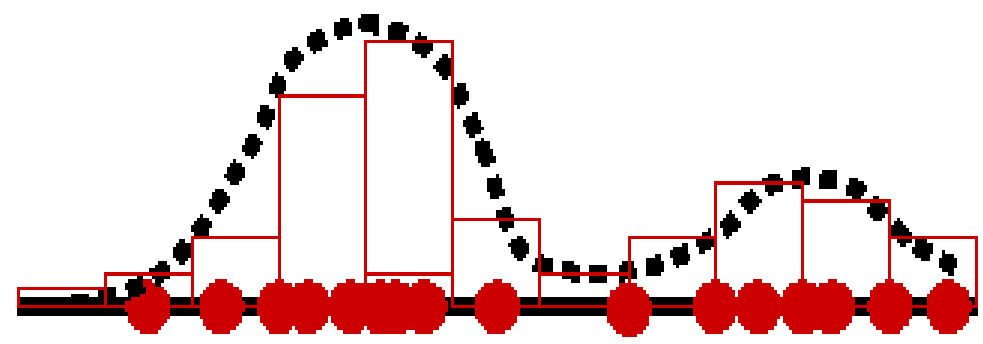
3 bins



7 bins



11 bins



## שיטת ההיסטוגרמה

מספר התחומים  $M$  משפיע מאוד על איכות השיטה:

$M$  גדול: ההיסטוגרמה "רועשת".

$M$  קטן: מאבדים הרבה אינפורמציה.

גודל ה- $M$  האופטימלי תלוי ב:

- גודל קבוצת האימון.
- מבנה פונקצית ההתפלגות.
- רגישות הבעיה לחוסר דיוק / רעש.

חסרונות עיקריים: במימדים גבוהים יותר נצטרך קבוצת אימון בגודל  $M^{\dim}$ .

פונקציות ההתפלגות שמקבלים איננה רציפה.

דוגמא:

$w = \{w_1 = \text{male}, w_2 = \text{female}\}$  ו- $x$  הינו הגובה.

נחלק את, תחום  $x$  לשני חלקים:  $[x > 1.80]$ ,  $[x \leq 1.80]$

$\Theta = \{p_i\}$

$$P(x | w_i) = \begin{cases} p_i & \text{iff } x > 1.80 \\ 1-p_i & \text{else} \end{cases}$$

במקרה זה, כל מה שנותר לשערך הוא שתי הסתברויות  $(p_1, p_2)$ .

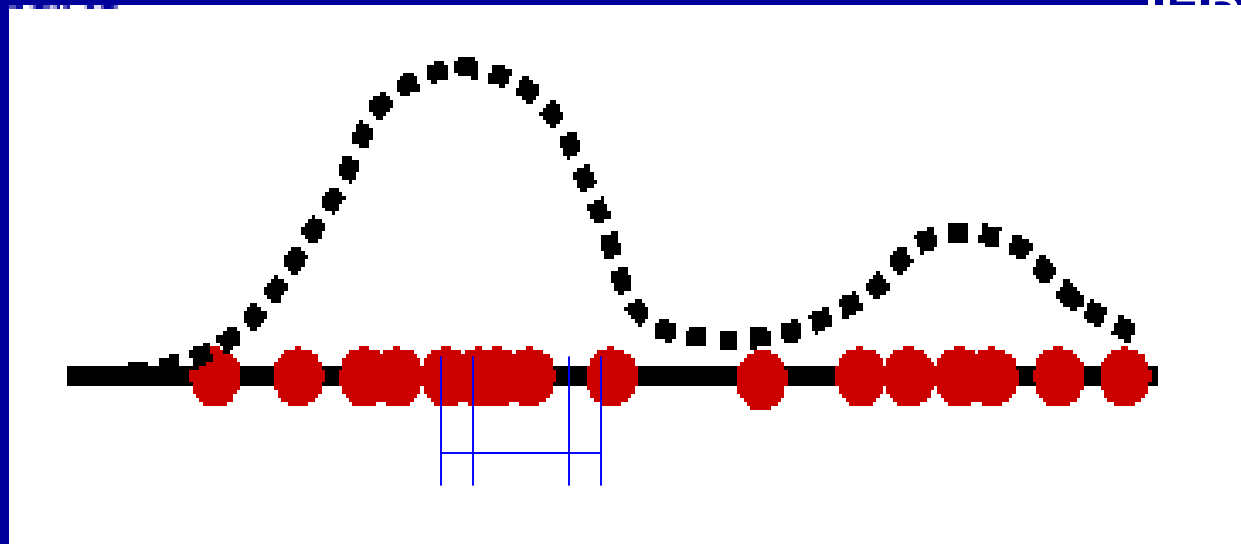
## שיפור החסרונות של גישת ההסטוגרמה

### גישת החלון

כדי לחשב את צפיפות ההסתברות  $P(X|w)$  בנקודה  $x_0$ , נחשב את מספר

המקרים מקבוצת האימון שבהם ערך  $x$  הוא בחלון כלשהו סביב  $x_0$ , ונחלק

בנפח החלון



דוגמא:

$w = \{w_1 = \text{male}, w_2 = \text{female}\}$  ו- $x$  הינו הגובה.

נחלק את, תחום  $x$  לטווחים החופפים הבאים:

לכל  $x_0$ , נתאים טווח  $[x_0 - d, x_0 + d]$ .

$$\Pr(x = x_0 | w_i) = \Pr(x \in [x_0 - d, x_0 + d] | w_i) / 2d$$

איך נבחר את גודל החלון  $d$ ?

- שרירותית

- נבחר גודל משתנה (כפונקציה של מספר הנקודות סביב  $x_0$ ).

## שיטה מס' 2: Vector Quantization

נרצה לבחור קבוצה של  $k$  "נציגים" במרחב המאפיינים,

כך שכל ווקטור במרחב יוחלף ע"י אחד הנציגים.

היתרון: עברנו מבעיית סיווג רציפה במימד גבוה לבעיית סיווג בדידה.

בעיית השערוך של ההתפלגות  $\Pr(X|w)$  הופכת לבעיית שערוך של

ההתפלגות  $\Pr(i|w)$ , כש- $i$  הוא מספר טבעי  $1 \leq i \leq k$ .

החסרון: איבוד רזולוציה (דיוק).

ככל שמספר הנציגים יהיה גדול יותר, איבוד הדיוק יהיה קטן יותר, אבל

## אלגוריתם למציאת k הנציגים: k-means

נרצה למצוא k נציגים אשר עבורם איבוד הדיוק הממוצע עבור קבוצת

האימון הוא מינימלי.

נגדיר את העיוות של ווקטור  $x$  כריבוע המרחק האוקלידי בינו לבין הנציג

הקרוב ביותר אליו.

הערה: מציאת k נציגים כאלו היא NP-קשה.

אלגוריתם k-means:

1. אתחל את קבוצת k הנציגים (למשל אקראית).

בצע עד להתכנסות:

2. סווג כל ווקטור בקבוצת האימון לנציג הקרוב ביותר.

חגי אהרונוביץ

3. חשב מחדש את הנציגים עם "י" תוספת הווקטורים המשוייכים

## ניתוח אלגוריתם k-means

- אלגוריתם k-means מתכנס תמיד.
- האלגוריתם מתכנס למינימום מקומי.
- אתחול אקראי מספק בדרך-כלל.
- בחירת  $k$  משמעותית – קטן מדי: אובדן דיוק. גדול מדי: over-fitting.
- בדר"כ נבחר את  $k$  אמפירית.

שימוש נוסף לאלגוריתם: דחיסה (לא משמרת).

## הגישה הפרמטרית: שיטות לשערוך פרמטרים

כדי לשערך את הפרמטרים של מודל כלשהו, יש צורך בשני מרכיבים:

1. קריטריון טיב (= פונקצית ציון) למודל.
2. אלגוריתם חיפוש במרחב הפרמטרים.

קריטריון הטיב למודל מאפשר להשוות בין שני מודלים ולענות על השאלה:

"איזה מודל טוב יותר?"

אלגוריתם החיפוש משתמש בקריטריון הטיב כדי למצוא מודל אופטימלי

## קריטריוני טיב

הסתברות אפוסטריורית מירבית (MAP - Maximum a-posteriori)

(

נרצה לבחור את המודל המסתבר ביותר בהינתן תצפיות האימון.  
$$w_{MAP} = \operatorname{argmax}_{w \in \{W\}} \Pr(w|X) = \operatorname{argmax}_{w \in \{W\}} \frac{\Pr(X|w) \Pr(w)}{\Pr(X)} = \operatorname{argmax}_{w \in \{W\}} \Pr(X|w) \Pr(w)$$

כזכור,  $P(x | w) = f(x, \theta)$ , ולכן:

$$\Theta_{MAP} = \operatorname{argmax}_{\Theta \in \{\Theta\}} \Pr(\Theta|X) = \operatorname{argmax}_{\Theta \in \{\Theta\}} \frac{f(X, \Theta) \Pr(\Theta)}{\Pr(X)} = \operatorname{argmax}_{\Theta \in \{\Theta\}} f(X, \Theta) \Pr(\Theta)$$

## קריטריוני טיב (2)

נראות מקסימלית (ML - Maximum Likelihood)

נניח שההתפלגות  $\Pr(\Theta)$  היא אחידה, כלומר כל המודלים האפשריים הם

סבירים אפריורית באותה מידה.

$$\Theta_{\text{MAP}} = \operatorname{argmax}_{\Theta \in \{\Theta\}} \Pr(\Theta|X) = \operatorname{argmax}_{\Theta \in \{\Theta\}} \frac{f(X, \Theta)\Pr(\Theta)}{\Pr(X)} = \operatorname{argmax}_{\Theta \in \{\Theta\}} f(X, \Theta)\Pr(\Theta)$$

$$\Theta_{\text{ML}} = \operatorname{argmax}_{\Theta \in \{\Theta\}} f(X, \Theta)\Pr(\Theta) = \operatorname{argmax}_{\Theta \in \{\Theta\}} f(X, \Theta)$$

## נראות מקסימלית (ML - Maximum Likelihood)

דוגמא 1: התפלגות בינומית.

דוגמא 2: התפלגות נורמלית.

## הסתברות אפוסטריורית מירבית (MAP - Maximum a-posteriori)

דוגמא 1: התפלגות בינומית עם מידע אפריורי.

## שיערוך התפלגות בינומית: בעיית שכיחות אפס

נחזור לדוגמת מידול טקסט.

נניח הנחה מרקובית מסדר 2:

$$\Pr(x_i | x_1, \dots, x_{i-1}, w) = \Pr(x_i | x_{i-2}, x_{i-1}, w)$$

נצטרך לשערך מתוך קבוצת האימון עבור כל שלשה אפשרית של אותיות

$$\Pr(x_i = a | x_{i-2} = b, x_{i-1} = c, w)$$

עפ"י שיטת שיערוך Maximum Likelihood, נספור לכל שלשה את

מספר המופעים  $k$  של השלשה בקבוצת האימון, ונחלק בגודל קבוצת האימון:

$$\Pr(x_i = a | x_{i-2} = b, x_{i-1} = c, w) = k(a,b,c) / n$$

הבעיה: סביר להניח שתהיינה שלשות שעבורן נקבל:  $k(a,b,c) = 0$  חגי אפרוביץ

## שיערוך התפלגות בינומית: בעיית שכיחות אפס (המשך)

מה יקרה עם שלשה כזו תופיע בטקסט המבחן?

... נקבל:

$$\Pr(X|w) = 0$$

פתרון אפשרי:

נעדכן את פונקצית שיערוך הפרמטרים לפונקציה הבאה:

$$\Pr(x_i = a | x_{i-2} = b, x_{i-1} = c, w) = \frac{k(a,b,c) + \varepsilon}{n + \varepsilon * m^3}$$

m – גודל האלף-בית

n – גודל טקסט האימון

$\varepsilon$  – פרמטר שרירותי (בסדר גודל 0.1-1).

## שיערוך פרמטרים: הגישה הבייסיאנית

במקום לחפש השמה לפרמטרים (מודל) הממקסמת קריטריון כלשהו,

נחשב את ההסתברות של כל השמה (מודל).

כשנרצה להשתמש במודל לצורך חישוב כלשהו, נשתמש בכל המודלים,

ונמשקלם בהסתברותם.

כלומר נחשב את  $\Pr(\theta|X)$  לכל  $\theta$  ונחשב את  $\Pr(Y|X)$  כך:

$$\Pr(Y|X) = \int_{\Theta} \Pr(Y|\theta) \Pr(\theta|X) d\theta$$

## הגישה הפרמטרית: תערובת גאוסיינים (GMM)

ראינו שקיים פתרון סגור לשערוך פרמטרים של התפלגות נורמלית.

הבעיה: במציאות, רוב ההתפלגויות אינם נורמליות.

הפתרון: נקרב כל התפלגות ע"י תערובת של התפלגויות נורמליות

תערובת גאוסיינים (Gaussian Mixture Model).

$$\Pr(x) = \sum_{k=1} w_k \Pr(x | N(\mu_k, \sigma_k^2))$$

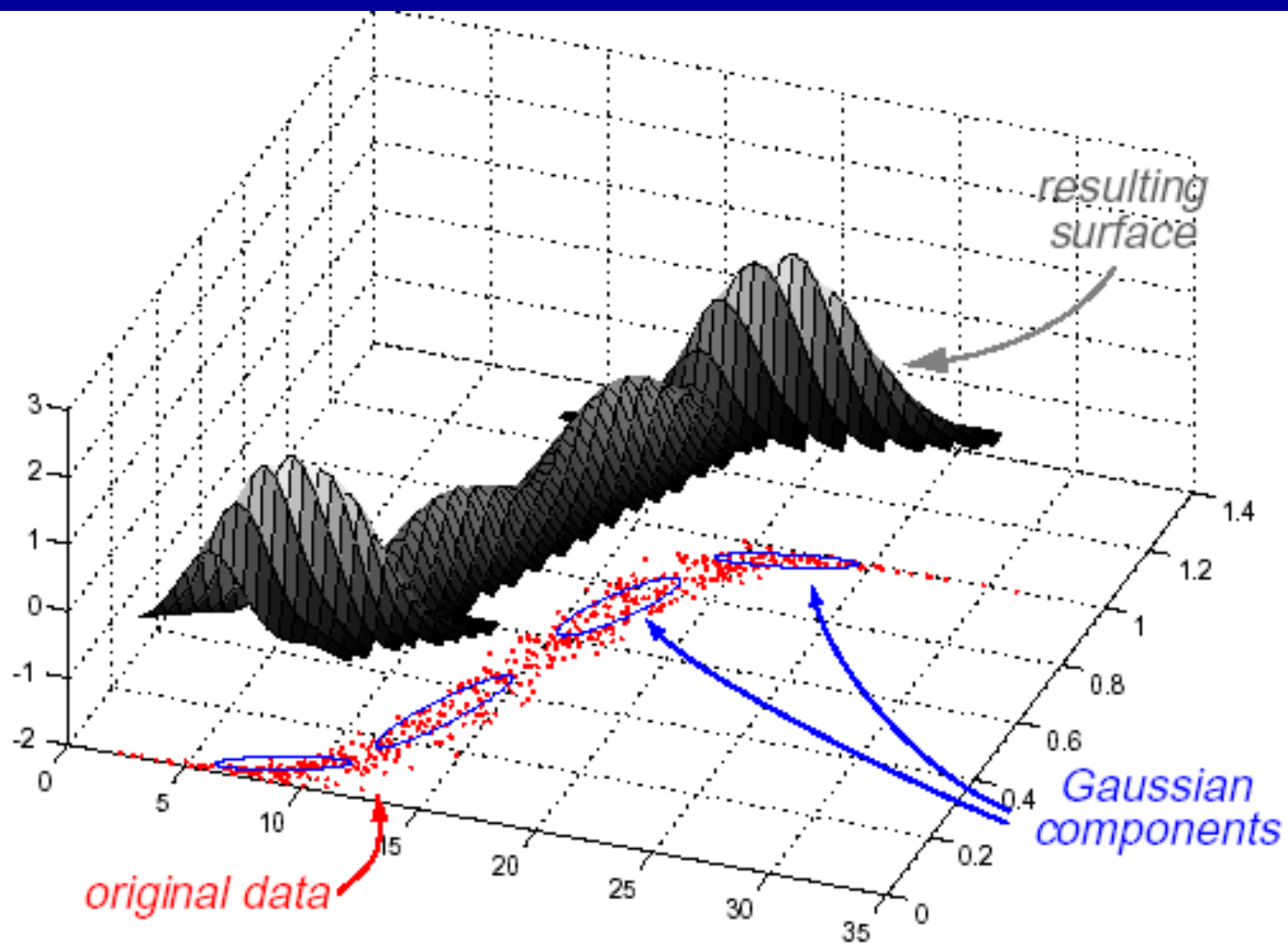
$$\sum_{k=1}^n w_k = 1$$

טענה: אפשר להראות של כל התפלגות (סבירה) אפשר לקרב ע"י

תערובת של מספר קטן יחסית התפלגויות נורמליות.

אינטרפרטציה: כל תצפית  $x$  נוצרת ע"י אחד מן הגאוסיינים הנבחר עפ"י חגי אהרונוביץ

## דוגמא להתפלגות המתוארת ע"י GMM



אם היינו יודעים עבור כל תצפית לאיזה גאוסין היא "שייכת",  
יכולנו לשערך את הפרמטרים של ה-GMM בצורה דומה לשערוך  
פרמטרי  
גאוסין בודד.

הבעיה: אנו לא יודעים עבור כל תצפית לאיזה גאוסין היא "שייכת".  
מעבר לכך, לא ידוע על שיטה כלשהי המאפשרת למצוא את הערכים  
ה-ML  
של תערובת גאוסיינים.

הפתרון (החלקי): נפעיל שיטה אשר מאפשרת למצוא מקסימום  
מקומי

לפונקצית ה-likelihood (במקום המקסימום גלובאלי הדרוש).

לשיטה זו קוראים EM.

## אלגוריתם EM

אלגוריתם EM הוא אלגוריתם כללי (לאו דווקא ל-GMM) המטפל בבעיה הבאה:

- רוצים לשערך את הפרמטרים  $\theta$  של מודל פרמטרי.
- קיימים משתנים חבויים  $Y$  אשר באמצעותם קל לשערך את המודל הפרמטרי מתוך התצפיות  $X$ .

במקרה שלנו,  $Y$  מציין עבור כל תצפית מאיזה גאוסין היא "נוצרה".

אלגוריתם ה-EM הנו האלגוריתם האיטרטיבי הבא:

1. אתחל את המודל  $\theta$  בערכים כלשהם.

חזור עד להתכנסות:

2. חשב את ערכי  $Y$  עפ"י  $\theta$  ו- $X$ .

3. שערך את  $\theta$  עפ"י  $Y$  ו- $X$ .

ניתן להוכיח עבור אלגוריתם EM:

1. בכל איטרציה ה-likelihood משתפר.

2. האלגוריתם מתכנס.

אבל... לאו דווקא מגיעים למקסימום גלובאלי.

## אלגוריתם EM עבור שיערוך GMM

1. אתחל בצורה כלשהי את הפרמטרים של המודל, למשל ע"י:  
בחירה אקראית של ווקטורים כפרמטרי התוחלת  $\mu_k$ ,  
שימוש בקובריאנס הכללי לאיתחול פרמטרי הקובריאנס  $\sigma_k$ ,  
ערך אחיד לפרמטרי המשקל  $w_k$ .

2. בצע מס' איטרציות (עד להתכנסות):

E: חשב את ההסתברות של כל אחד מן הגאוסיינים עפ"י כל

ווקטור  $x$ :

$$(I) \quad \Pr(x|k) = \Pr(x | N(\mu_k, \sigma_k^2))$$

$$(II) \quad \Pr(k|x) = \frac{\Pr(x|k)\Pr(k)}{\Pr(x)} = \frac{\Pr(x|k)w_k}{\Pr(x)}$$

$$(III) \quad \sum_j \Pr(j|x) = 1$$

$$\Rightarrow \Pr(k|x) = \frac{\Pr(x|k)w_k}{\sum_j \Pr(x|j)w_k}$$

M: שערך את הפרמטרים של ה-GMM ע"י המשוואות הבאות:

$$\bar{\mu}_k = \frac{\sum_{x \in X} \Pr(k|x) \cdot \bar{x}}{\sum_{x \in X} \Pr(k|x)}$$

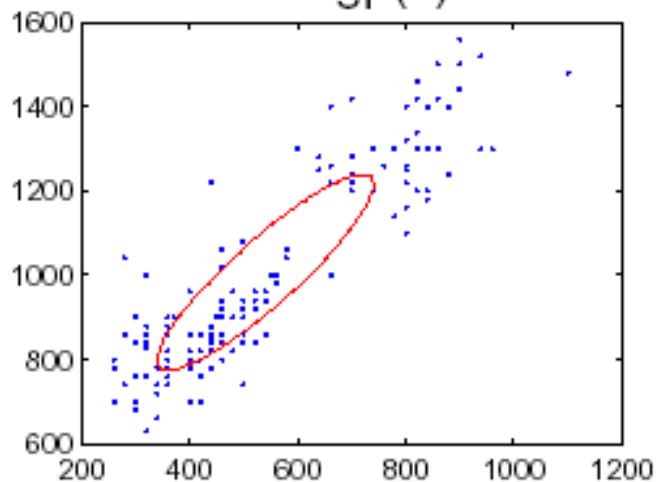
$$\sigma_{k_i,j} = \frac{\sum_{x \in X} \Pr(k|x) \cdot (x_i - \mu_{k_i}) \cdot (x_j - \mu_{k_j})}{\sum_{x \in X} \Pr(k|x)}$$

$$w_k = \frac{1}{|X|} \sum_{x \in X} \Pr(k|x)$$

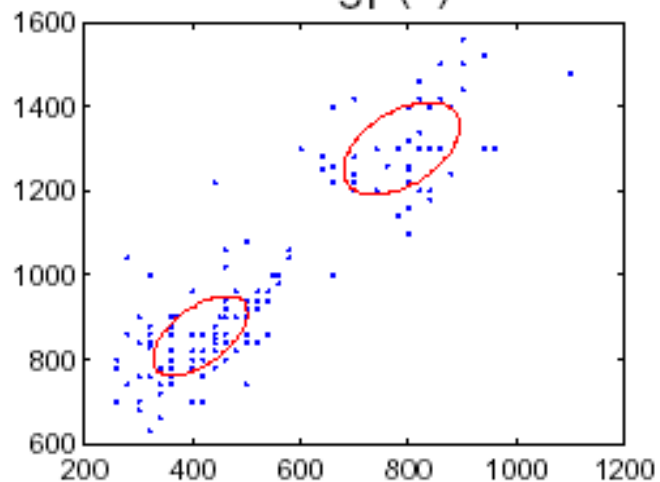
# GMM examples

- **Vowel data fit with different mixture counts:**

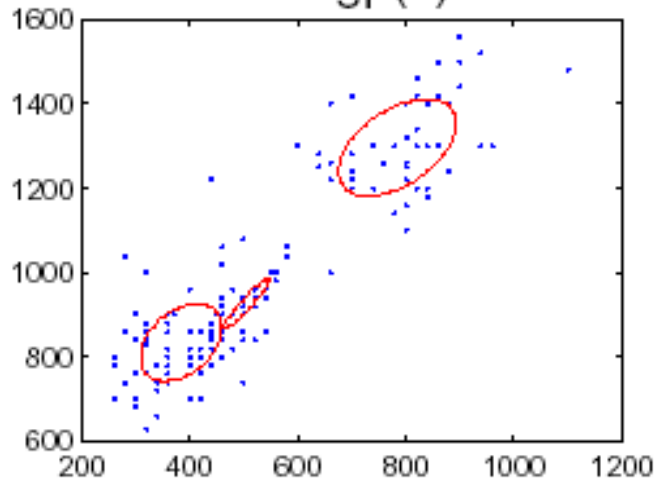
1 Gauss  $\log p(x) = -1911$



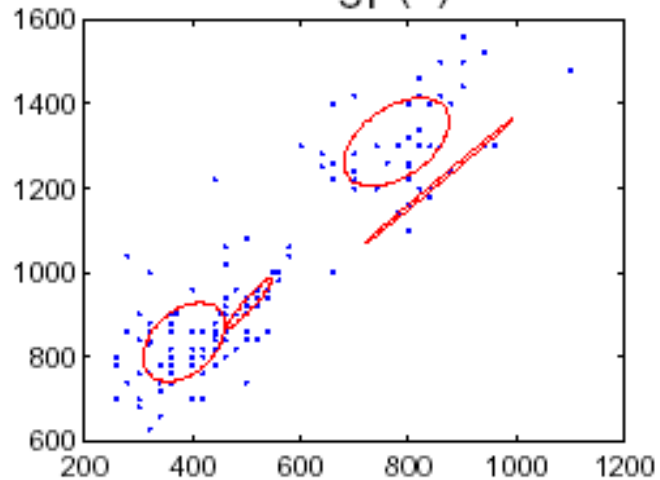
2 Gauss  $\log p(x) = -1864$



3 Gauss  $\log p(x) = -1849$



4 Gauss  $\log p(x) = -1840$



## הערות:

1. מודל עשיר מדי (יותר מדי גאוסיינים) ילמד טוב "מדי" את דוגמאות האימון, ולא יכליל.
2. את מס' הגאוסיינים ניתן לקבוע באמצעות קבוצת בקרה (validation set).
3. את מס' הגאוסיינים ניתן לקבוע עפ"י גודל קבוצת האימון: 100-20 דוגמאות/גאוסייין.
4. לעיתים נניח כי מטריצת הקובריאנס אלכסונית.
5. לעיתים נקשור (tie) פרמטרים שונים יחדיו – בעיקר קובריאנס.