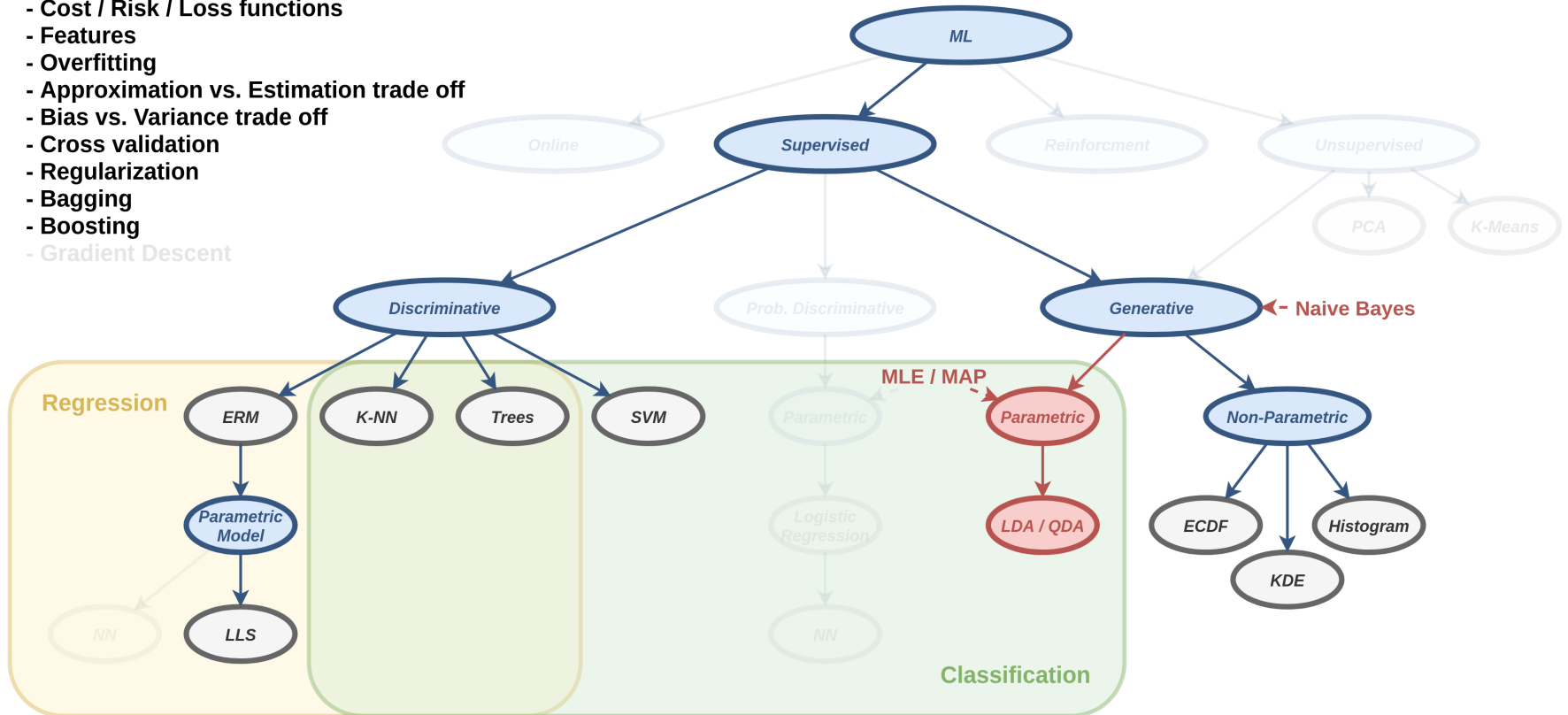


# הרצאה 8 - שיערוך פילוג בשיטות פרמטריות וסיווג גנרטיבי

## 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



**נניח שאנו רוצים לחזות האם אדם מסויים פרק את הכתף על פי הסימפטומים שלו. לשם כך נסתכל על המדגם הבא**

נזלת	נימול ביד	סימנים כחולים	נפיחות	כאב בכתף	פריקה	
-	+	+	+	+	+	1
-	-	+	+	+	+	2
-	+	-	+	+	+	3
+	-	+	+	+	+	4
-	-	-	-	-	-	5
+	-	-	-	+	-	6
-	-	+	-	-	-	7
-	-	-	-	+	-	8

**האם לאדם עם כל הסימפטומים יש פריקה של הכתף?**

# שיערוך נאיבי - חוסר תלות בין המשתנים

נניח שאנו רוצים לשערך צפיפות הסתברות של משתנה  $D$  ממדי. שיטה מאד נאיבית (לא מתוחכמת) לפתור את הבעיה היא להתעלם מהקשר בין המשתנים השונים ולהניח שהם בלתי תלויים סטטיסטית. זאת אומרת ש:

$$p_{\mathbf{x}}(\mathbf{x}) = p_{x_1}(x_1)p_{x_2}(x_2) \dots p_{x_D}(x_D) = \prod_{d=1}^D p_{x_d}(x_d)$$

• חיסרון: שהיא מגבילה מאד את הפילוגים שניתן ללמוד.

# מסוג ביס נאיבי - Naïve Bayes Classification

**תזכורת** במשימות סיווג קיים  $\mathbf{x} \in R^D$  ו- $y \in \{1, 2, \dots, C\}$ . בנוסף מוגדרת הסתברות משותפת  $P(\mathbf{x}, y)$ . נוכל להשתמש בשערוך הנאיבי לפתרון בעיות סיווג. נניח כי:

$$p_{\mathbf{x}|y}(\mathbf{x}|y) = p_{x_1|y}(x_1|y)p_{x_2|y}(x_2|y) \dots p_{x_D|y}(x_D|y) = \prod_{d=1}^D p_{x_d|y}(x_d|y)$$

- זו כמובן הנחה מאוד פשטנית שאינה מתקיימת במציאות.
- אנו **לא** נרצה להניח חוסר תלות בין  $\mathbf{x}$  ל  $y$ .

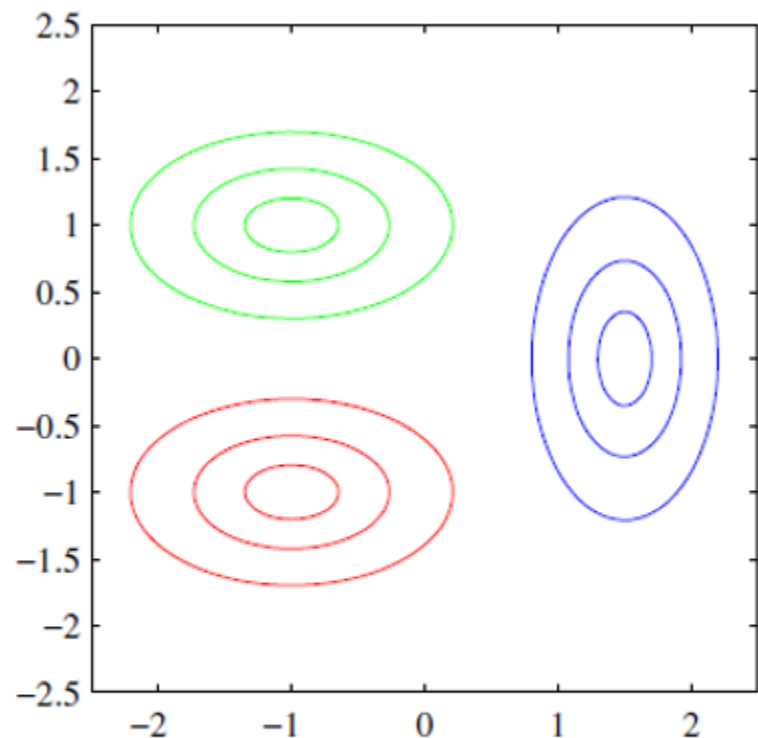
# מסוג ביים נאיבי - Naïve Bayes Classification

---

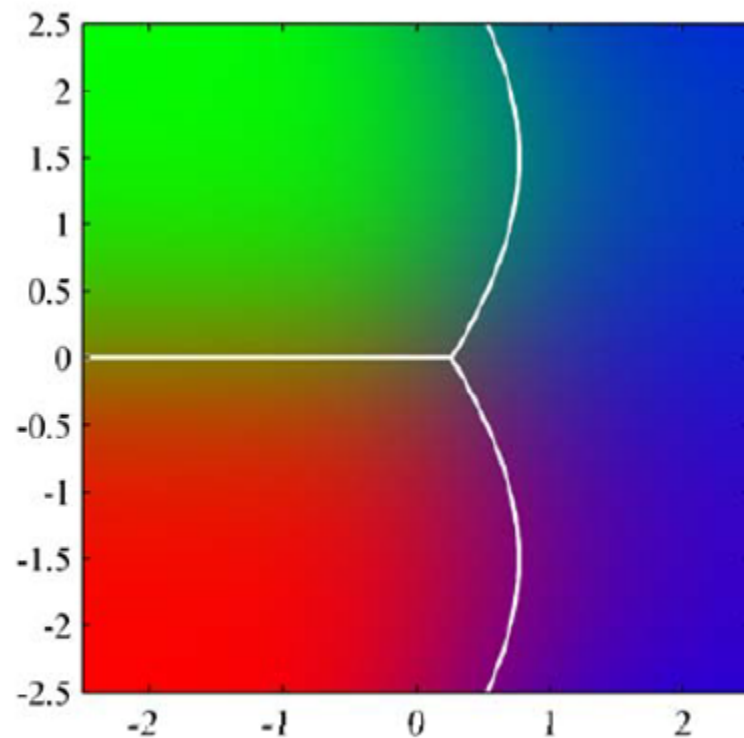
החזאי אשר ממזער את הסיכוי לטעות יהיה:

$$\begin{aligned}\hat{y} = h(\mathbf{x}) &= \arg \max_y p_{y|\mathbf{x}}(y|\mathbf{x}) \\ &= \arg \max_y p_{\mathbf{x}|y}(\mathbf{x}|y)p_y(y) \\ &= \arg \max_y p_y(y) \prod_{d=1}^D p_{x_d|y}(x_d|y)\end{aligned}$$

# מסוג ביים נאיבי - Naïve Bayes Classification



$$p(x|y), y = 1, 2, 3$$



$$p(y|x), y = 1, 2, 3$$

האיור מתוך, C.M. Bishop, Pattern Recognition and Machine Learning

# דוגמא 1 - זיהוי פריקה של הכתף

תחת הנחת חוסר התלות נשערך בנפרד את ההסתברות המותנית של כל אחד מהרכיבים  $p_{x_d|y}(x_d|y)$ ,  $x_d \in \{0, 1\}$ . לדוגמא:

$$p_{x_{\text{pain}}|y}(x_{\text{pain}}|1) = \begin{cases} \frac{4}{4} = 1 & 1 \\ \frac{0}{4} = 0 & 0 \end{cases}$$

$$p_{x_{\text{pain}}|y}(x_{\text{pain}}|0) = \begin{cases} \frac{2}{4} = 0.5 & 1 \\ \frac{2}{4} = 0.5 & 0 \end{cases}$$

$$p_{x_{\text{swelling}}|y}(x_{\text{swelling}}|1) = \begin{cases} \frac{4}{4} = 1 & 1 \\ \frac{0}{4} = 0 & 0 \end{cases}$$

$$p_{x_{\text{swelling}}|y}(x_{\text{swelling}}|0) = \begin{cases} \frac{0}{4} = 0 & 1 \\ \frac{1}{4} = 1 & 0 \end{cases}$$

ובאופן דומה לשאר הרכיבים.



# דוגמא 1 - זיהוי פריקה של הכתף

החיזוי עבור המקרה בו מופיעים כל הסימפטומים יהיה

$$\hat{y} = \arg \max_y p_y(y) \prod_{d=1}^D p_{x_d|y}(1|y)$$

זאת אומרת שעלינו לבדוק האם:

$$p_y(1) \prod_{d=1}^D p_{x_d|y}(1|1) \stackrel{?}{>} p_y(0) \prod_{d=1}^D p_{x_d|y}(1|0)$$

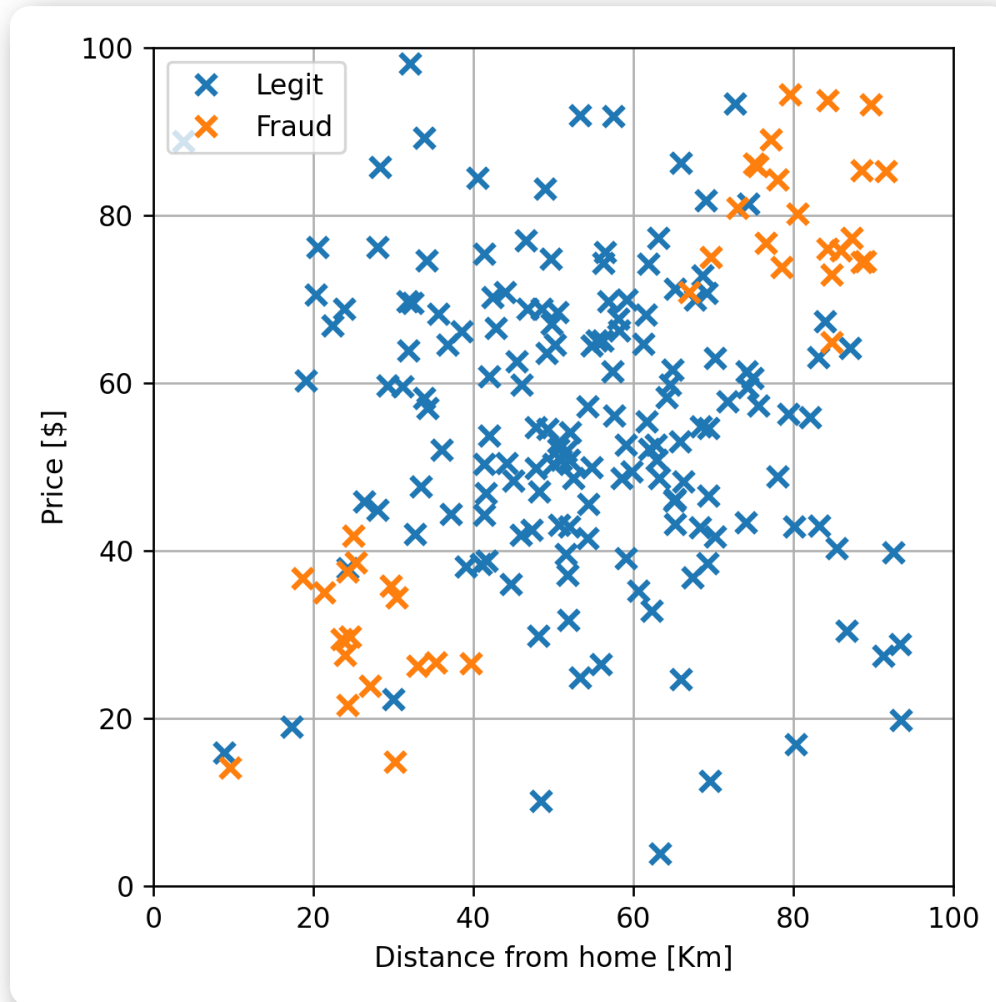
$$0.5 \times 1 \times 1 \times 0.75 \times 0.5 \times 0.25 \stackrel{?}{>} 0.5 \times 0.5 \times 0 \times 0.25 \times 0 \times 0.25$$

$$0.09375 \stackrel{?}{>} 0$$

מכיוון שתנאי זה מתקיים, החיזוי יהיה שישנה פריקה של כתף.

# דוגמא 2 - זיהוי הונאות אשראי

ננסה להשתמש בשיטה זו לבעיית חיזוי הונאות האשראי



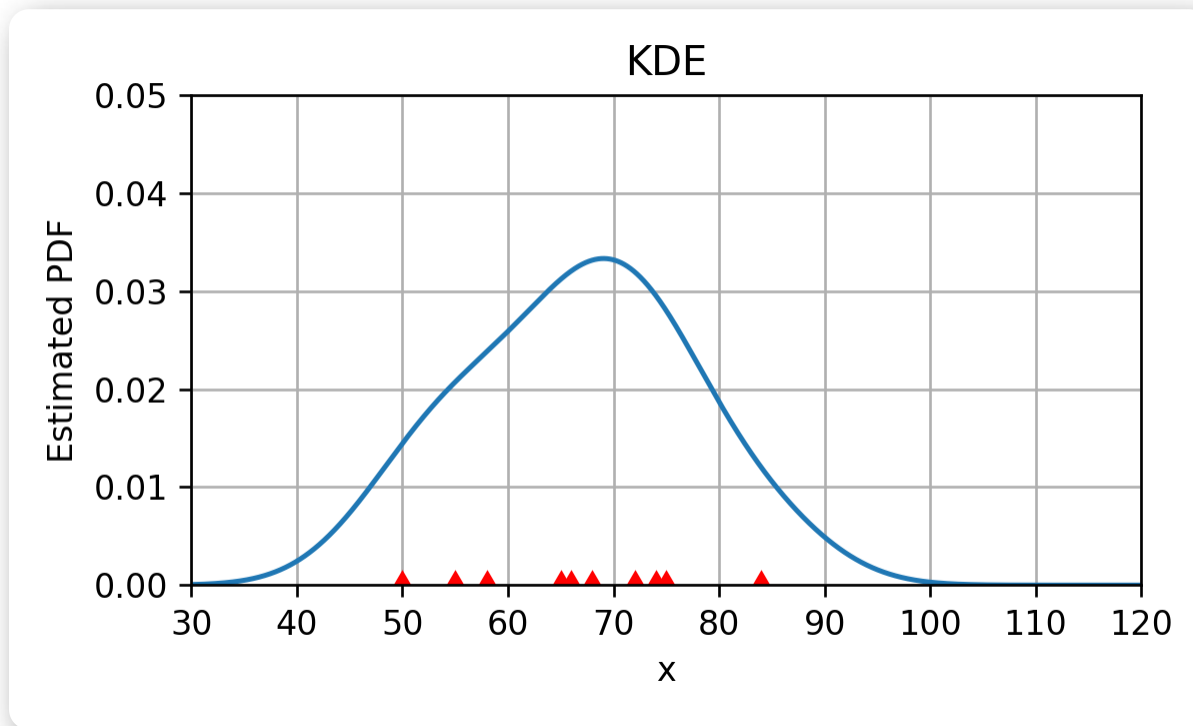
- דומה לשימוש שעשינו במודלים פרמטריים בגישה הדטרמיניסטית.
- נגביל את הצורה של הפונקציה שאותה אנו רוצים לשערך למודל פרמטרי.
- נסמן את וקטור הפרמטרים ב  $\theta$ .
- כאן המודל חייב לייצר פילוג חוקי עבור כל בחירה של פרמטרים.
- מגבלה קשה אשר מצמצמת מאד את המודלים הפרמטריים שאיתם ניתן לעבוד.

- נרצה למצוא דרך לתת "ציון" לכל בחירה של פרמטרים ולחפש את הפרמטרים אשר מניבים את הציון הטוב ביותר.
- נציג שתי גישות שונות להתייחס לפרמטרים של המודל.
- כל גישה מובילה לדרך מעט שונה לבחירה של הפרמטרים.

# דוגמא: שיערוך הפילוג של זמן הנסיעה

$$\mathcal{D} = \{x^{(i)}\} = \{55, 68, 75, 50, 72, 84, 65, 58, 74, 66\}$$

משערך ה KDE של הפילוג (לא דנו הסמסטר) הינו:



נרצה לשערך פרמטרים של פילוג נורמלי שיתאר בצורה טובה את הדגימות במדגם.

# הגישה הלא-בייסיאנית (קלאסית או תדירותית) (**Frequentist**)

---

$$p_{\mathbf{x}}(\mathbf{x}; \theta)$$

- נתייחס לפרמטרים כאל מספרים שאותם יש לקבוע על מנת שהמודל יתאר בצורה טובה את המדגם.
- ההנחה היא כי יש ערך לא ידוע של הפרמטר שהוא ה"טוב" ביותר.

# Maximum Likelihood Estimator (MLE)

---

- נסמן ב  $p_{\mathcal{D}}(\mathcal{D}; \theta)$  את ההסתברות לקבלת המדגם  $\mathcal{D} = \{x^{(i)}\}$ .
- גודל זה מכונה **הסתברות (likelihood)** של המדגם.
- אנו מעוניינים למצוא את הפרמטרים  $\theta$  אשר מניבים את הסבירות הכי גבוהה.
- מקובל לסמן את פונקציית ה **likelihood** באופן הבא:

$$\mathcal{L}(\theta; \mathcal{D}) \triangleq p_{\mathcal{D}}(\mathcal{D}; \theta)$$

משערך ה MLE של  $\theta$  הוא הערך אשר ממקסם את פונקציית ה **likelihood**:

$$\hat{\theta}_{\text{MLE}} = \arg \max_{\theta} \mathcal{L}(\theta; \mathcal{D})$$

# Maximum Likelihood Estimator (MLE)

---

$$\hat{\theta}_{\text{MLE}} = \arg \max_{\theta} \mathcal{L}(\theta; \mathcal{D})$$

- נרשום את בעיית האופטימיזציה כבעיית מינימיזציה:

$$\hat{\theta}_{\text{MLE}} = \arg \min_{\theta} -\mathcal{L}(\theta; \mathcal{D})$$

- כאשר הדגימות במדגם הם i.i.d:

$$p_{\mathcal{D}}(\mathcal{D}; \theta) = \prod_i p_{\mathbf{x}}(\mathbf{x}^{(i)}; \theta)$$

ולכן:



# Maximum Likelihood Estimator (MLE)

$$\hat{\theta}_{\text{MLE}} = \arg \min_{\theta} -\mathcal{L}(\theta; \mathcal{D}) = \arg \min_{\theta} - \prod_i p_{\mathbf{x}}(\mathbf{x}^{(i)}; \theta)$$

במקרים רבים נוכל להחליף את המכפלה על כל הדגימות בסכום, על ידי החלפת פונקציית ה-likelihood ב- $\log$ -likelihood:

$$\hat{\theta}_{\text{MLE}} = \arg \min_{\theta} -\log \mathcal{L}(\theta; \mathcal{D}) = \arg \min_{\theta} - \sum_i \log \left( p_{\mathbf{x}}(\mathbf{x}^{(i)}; \theta) \right)$$

**הערה:** בקורסים "עיבוד אותות אקראיים" ו-"הסקה סטטיסטית" מרחיבים הרבה בנושא תכונות משערך זה ואחרים.

**ננסה להתאים למדגם מודל של פילוג נורמלי:**

$$p_x(x; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

**וקטור הפרמטרים הינו  $\theta = [\mu, \sigma]^T$ .**

$$p_x(x; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

**נרשום את בעיית האופטימיזציה של מציאת משעריך ה MLE:**

$$\begin{aligned}\hat{\theta}_{\text{MLE}} &= \arg \min_{\theta} - \sum_i \log\left(p_x(x^{(i)}; \theta)\right) \\ &= \arg \min_{\theta} - \sum_i \log\left(\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x^{(i)} - \mu)^2}{2\sigma^2}\right)\right) \\ &= \arg \min_{\theta} \sum_i \log(\sigma) + \frac{1}{2} \log(2\pi) + \frac{(x^{(i)} - \mu)^2}{2\sigma^2} \\ &= \arg \min_{\theta} N \log(\sigma) + \frac{1}{2\sigma^2} \sum_i (x^{(i)} - \mu)^2\end{aligned}$$

**נפתור על ידי גזירה והשוואה ל-0.**

נסמן את ה objective ב  $f$ :

$$f(\boldsymbol{\theta}; \mathbf{x}) = N \log(\sigma) + \frac{1}{2\sigma^2} \sum_i (x^{(i)} - \mu)^2$$

$$\begin{cases} \frac{\partial f(\boldsymbol{\theta}; \mathbf{x})}{\partial \mu} = 0 \\ \frac{\partial f(\boldsymbol{\theta}; \mathbf{x})}{\partial \sigma} = 0 \end{cases}$$

$$\Leftrightarrow \begin{cases} -\frac{1}{\sigma^2} \sum_i (x^{(i)} - \mu) = 0 \\ \frac{N}{\sigma} - \frac{1}{\sigma^3} \sum_i (x^{(i)} - \mu)^2 = 0 \end{cases}$$

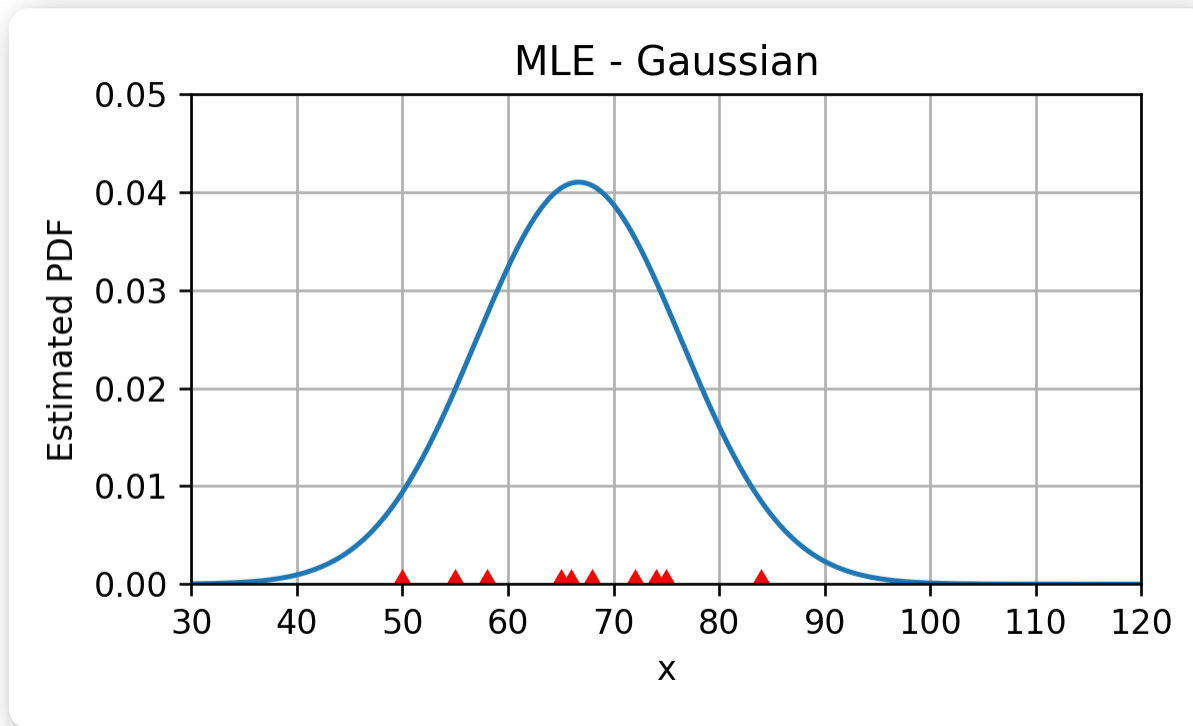
$$\Leftrightarrow \begin{cases} N\mu - \sum_i x^{(i)} = 0 \\ N\sigma^2 - \sum_i (x^{(i)} - \mu)^2 = 0 \end{cases}$$

$$\Leftrightarrow \begin{cases} \mu = \frac{1}{N} \sum_i x^{(i)} \\ \sigma^2 = \frac{1}{N} \sum_i (x^{(i)} - \mu)^2 \end{cases}$$

במקרה של הנסיעות בכביש החוף נקבל:

$$\mu = 66.7 \text{ [min]}$$

$$\sigma = 9.7 \text{ [min]}$$



- הפרמטרים של המודל הם ריאליזציות (הגרלות) של משתנה אקראי.
- גישה זו מניחה שיש בידינו מודל לפילוג המשותף של הפרמטרים והמדגם.

$$p_{\mathcal{D},\theta}(\mathcal{D}, \theta) = p_{\mathcal{D}|\theta}(\mathcal{D}|\theta)p_{\theta}(\theta)$$

**תחת ההנחה שבהינתן הפרמטרים הדגימות במדגם הם i.i.d:**

$$p_{\mathcal{D},\theta}(\mathcal{D}, \theta) = p_{\theta}(\theta) \prod_i p_{\mathbf{x}|\theta}(\mathbf{x}^{(i)}|\theta)$$

- עלינו לקבוע את  $p_{\mathbf{x}|\theta}(\mathbf{x}|\theta)$  ואת  $p_{\theta}(\theta)$ .
- זכרו בשבגישה הקודמת הנחנו כי  $\theta$  קבוע אבל לא ידוע.

## A Prioiri Distribution

הפילוג השולי של הפרמטרים  $p_{\theta}(\theta)$ , מכונה לרוב הפילוג הפריורי (prior distribution) או הא-פריורי (a priori distribution), זאת אומרת הפילוג של  $\theta$  לפני שראינו את המדגם.

## A Posteriori Distribution

הפילוג של הפרמטרים בהינתן המדגם  $p_{\theta|\mathcal{D}}(\theta|\mathcal{D})$  מכונה הפילוג הפוסטריורי (posterior distribution) או א-פוסטריורי (a posteriori distribution) (או הפילוג בדיעבד), זאת אומרת, הפילוג אחרי שראינו את המדגם.

# Maximum A-posteriori Probability (MAP)

MAP משערך את הערך אשר ממקסם את הפילוג הא-פוסטריורי (הערך הכי סביר של  $\theta$  בהינתן המדגם  $p_{\theta|\mathcal{D}}(\theta|\mathcal{D})$ ):

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} p_{\theta|\mathcal{D}}(\theta|\mathcal{D}) = \arg \min_{\theta} - \log p_{\theta|\mathcal{D}}(\theta|\mathcal{D})$$

על פי חוק בייס, נוכל לכתוב זאת כ:

$$\hat{\theta}_{\text{MAP}} = \arg \min_{\theta} - \log \frac{p_{\mathcal{D}|\theta}(\mathcal{D}|\theta)p_{\theta}(\theta)}{p_{\mathcal{D}}(\mathcal{D})} = \arg \min_{\theta} - \log p_{\mathcal{D}|\theta}(\mathcal{D}|\theta) - \log p_{\theta}(\theta)$$

כאשר הדגימות במדגם **בהינתן**  $\theta$  הן **i.i.d**, מתקיים:

$$\hat{\theta}_{\text{MAP}} = \arg \min_{\theta} - \log p_{\theta}(\theta) - \sum_i \log p_{\mathbf{x}|\theta}(\mathbf{x}^{(i)}|\theta)$$



# Maximum A-posteriori Probability (MAP)

---

$$\hat{\theta}_{\text{MAP}} = \arg \min_{\theta} - p_{\theta}(\theta) \prod_i p_{\mathbf{x}|\theta}(\mathbf{x}^{(i)}|\theta)$$

גם כאן נוכל להפוך את המכפלה לסכום על ידי מזעור מינוס הלוג של הפונקציה:

$$\hat{\theta}_{\text{MAP}} = \arg \min_{\theta} - \log(p_{\theta}(\theta)) - \sum_i \log(p_{\mathbf{x}|\theta}(\mathbf{x}^{(i)}|\theta))$$

# ההבדל בין MLE ל MAP

$$\hat{\theta}_{\text{MLE}} = \arg \min_{\theta} - \sum_i \log \left( p_{\mathbf{x}}(\mathbf{x}^{(i)}; \theta) \right)$$

$$\hat{\theta}_{\text{MAP}} = \arg \min_{\theta} - \log (p_{\theta}(\theta)) - \sum_i \log \left( p_{\mathbf{x}|\theta}(\mathbf{x}^{(i)} | \theta) \right)$$

- האיבר  $-\log(p_{\theta}(\theta))$  מוסיף את הידע שיש לנו לגבי איזה ערכים של  $\theta$  יותר סבירים.
- ראינו תוספת שכזו כאשר דיברנו על רגולריזציה.
- ניתן לחשוב על בעיית ה MAP כעל בעיית MLE עם רגולריזציה.
- בתרגיל הבית אתם תראו את השקילות שבין בעיית MAP לבין לבעיית MLE עם רגולריזציה.

# בגישה בייסיאנית השערוך הוא בעיית חיזוי

---

- **אנו מתייחסים גם למדגם וגם לפרמטרים בריאליזציות של משתנים אקראיים.**
- **אנו מניחים שאנו יודעים את הפילוג המשותף שלהם.**
- **ואנו מנסים למצוא את הערך של הפרמטרים בהינתן המדגם.**
- **זוהי בדיוק בעיית חיזוי קלאסית של משתנה אקראי אחד בהינתן משתנה אקראי אחר על סמך הפילוג המשותף.**

# דוגמא - הוספת prior

- נחזור לדוגמא של התאמת מודל של פילוג נורמלי לפילוג של זמן הנסיעה בכביש החוף.
- לשם הפשטות נקבע את סטיית התקן של המודל ל  $\sigma = 10$ .
- הפרמטר היחיד של המודל יהיה  $\mu$ :

$$p_{x|\mu}(x|\mu) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

# דוגמא - הוספת prior

$$p_{x|\mu}(x|\mu) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

• נניח שיש לנו ידע קודם על פילוג הצפוי של  $\mu$ .

• נניח שהפילוג הא-פריורי של  $\mu$  הוא גם פילוג נורמלי עם תוחלת  $\mu_\mu = 60$  וסטיית תקן של  $\sigma_\mu = 5$ :

$$p_\mu(\mu) = \frac{1}{\sqrt{2\pi}\sigma_\mu} \exp\left(-\frac{(\mu-\mu_\mu)^2}{2\cdot\sigma_\mu^2}\right)$$

נרשום את משעריך ה-MAP של  $\mu$ :

$$\hat{\mu}_{\text{MAP}} = \arg \min_{\mu} -\log(p_\mu(\mu)) - \sum_i \log(p_{x|\mu}(x^{(i)}|\mu))$$

# דוגמא - הוספת prior

$$\hat{\mu}_{\text{MAP}} = \arg \min_{\mu} -\log(p_{\mu}(\mu)) - \sum_i \log(p_{\mathbf{x}|\mu}(\mathbf{x}^{(i)}|\mu))$$

גזירה והשוואה ל-0 נותנת את התוצאה הבאה:

$$\begin{aligned} \frac{\partial f(\theta; x)}{\partial \mu} &= 0 \\ \Leftrightarrow \frac{1}{\sigma_{\mu}^2}(\mu - \mu_{\mu}) - \frac{1}{\sigma^2} \sum_i (x^{(i)} - \mu) &= 0 \\ \Leftrightarrow \left( \frac{1}{\sigma_{\mu}^2} + \frac{N}{\sigma^2} \right) \mu &= \frac{\mu_{\mu}}{\sigma_{\mu}^2} + \frac{1}{\sigma^2} \sum_i x^{(i)} \\ \Leftrightarrow \mu &= \frac{\frac{\sigma^2}{N\sigma_{\mu}^2} \mu_{\mu} + \frac{1}{N} \sum_i x^{(i)}}{\frac{\sigma^2}{N\sigma_{\mu}^2} + 1} \end{aligned}$$

זו למעשה ממוצע ממושקל בין הממוצע של  $x$  במדגם לבין  $\mu_{\mu}$ .

# דוגמא - הוספת prior

---

עבור הדוגמא שלנו נקבל:

$$\mu = 64.8 \text{ [min]}$$

ערך זה מעט יותר קרוב ל 60 מאשר התוצאה שקיבלנו בשיערוך ה MLE. זאת משום ה prior ש"מושך" את הפרמטרים לאיזורים הסבירים יותר ולכן הוא מקרב אותו ל  $\mu = 60$ .

# שימוש בשיערוך פרמטרי לפתרון בעיות supervised learning

---

נראה עכשיו איך להשתמש בשיערוך הצפיפות שתארנו צורך פתרון בעיות סיווג בלמידה מפוקחת. נציג שיטה אשר משתמשת במודל של פילוג נורמלי וב MLE לפתרון בעיות סיווג.



# Quadratic Discriminant Analysis ((QDA

- נשתמש בפילוג נורמלי וב MLE על מנת לשערך את  $p_{\mathbf{x}|y}(\mathbf{x}|y)$

- אנו צריכים לשערך מודל עבור כל אחת מ  $C$  המחלקות של  $y$ :

- וקטור תוחלת  $\mu_c$  עבור כל אחד מהמחלקות.

- מטריצת קווריאנס  $\Sigma_c$  עבור כל אחד מהמחלקות.

$$p_{\mathbf{x}|y}(\mathbf{x}|c; \mu_c, \Sigma_c) = \frac{1}{\sqrt{(2\pi)^D |\Sigma_c|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu_c)^\top \Sigma_c^{-1}(\mathbf{x} - \mu_c)\right)$$

הפילוג המשותף של  $\mathbf{x}$  ו  $y$  יהיה:

# Quadratic Discriminant Analysis ((QDA

---

בעיית האופטימיזציה של MLE תהיה:

$$\begin{aligned}\hat{\theta}_{\text{MLE}} &= \arg \min_{\theta} -\log \mathcal{L}(\theta; \mathcal{D}) \\ &= \arg \min_{\theta} -\sum_i \log \left( p_{\mathbf{x}|y}(\mathbf{x}^{(i)} | y^{(i)}; \boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y) p_y(y^{(i)}) \right) \\ &= \arg \min_{\theta} -\sum_i \log \left( p_{\mathbf{x}|y}(\mathbf{x}^{(i)} | y^{(i)}; \boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y) \right) + \log \left( p_y(y^{(i)}) \right) \\ &= \arg \min_{\theta} -\sum_i \log \left( p_{\mathbf{x}|y}(\mathbf{x}^{(i)} | y^{(i)}; \boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y) \right)\end{aligned}$$

בהינתן ש- $p_y$  ידוע.

# Quadratic Discriminant Analysis (QDA)

---

$$\hat{\theta}_{\text{MLE}} = \arg \min_{\theta} - \sum_i \log \left( p_{\mathbf{x}|y}(\mathbf{x}^{(i)} | y^{(i)}; \mu_y, \Sigma_y) \right)$$

• נחלק את הסכימה לסכימות נפרדות על כל אחת מהמחלקות.

# Quadratic Discriminant Analysis ((QDA

---

• נגדיר לשם כך את הסימונים הבאים:

◦  $\mathcal{I}_c = \{i : y^{(i)} = c\}$  - זאת אומרת, אוסף האינדקסים של הדגמים במדגם שמקיימים  $y^{(i)} = c$ .

◦  $|\mathcal{I}_c|$  - מספר האינדקסים ב  $\mathcal{I}_c$

$$\hat{\theta}_{\text{MLE}} = \arg \min_{\theta} - \sum_{i \in \mathcal{I}_1} \log \left( p_{\mathbf{x}|y}(\mathbf{x}^{(i)} | 1; \mu_1, \Sigma_1) \right) - \sum_{i \in \mathcal{I}_2} \log \left( p_{\mathbf{x}|y}(\mathbf{x}^{(i)} | 2; \mu_2, \Sigma_2) \right) - \dots$$

# Quadratic Discriminant Analysis ((QDA

עבור המחלקה  $c$  נקבל את בעיית האופטימיזציה הבאה:

$$\begin{aligned}\hat{\boldsymbol{\mu}}_{c,\text{MLE}}, \hat{\boldsymbol{\Sigma}}_{c,\text{MLE}} &= \arg \min_{\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c} - \sum_{i \in \mathcal{I}_c} \log \left( p_{\mathbf{x}|y}(\mathbf{x}^{(i)} | c; \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) \right) \\ &= \arg \min_{\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c} \sum_{i \in \mathcal{I}_c} \log \left( \sqrt{|\boldsymbol{\Sigma}_c|} \right) + \frac{1}{2} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_c)^\top \boldsymbol{\Sigma}_c^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_c)\end{aligned}$$

• ניתן לפתור את הבעיה הזו על ידי גזירה והשוואה ל-0.

• הפיתוח עבור  $\boldsymbol{\Sigma}_c$  הוא מעט מורכב ואנו לא נראה אותו בקורס זה ונקפוץ ישר לפתרון. הפיתוח מודגם בקורס "עיבוד וניתוח מידע".

# Quadratic Discriminant Analysis ((QDA

---

החישוב של  $\mu_c$

$$f(\boldsymbol{\theta}; \mathbf{x}) = \sum_{i \in \mathcal{I}_c} \log \left( \sqrt{|\Sigma_c|} \right) + \frac{1}{2} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_c)^\top \Sigma_c^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_c)$$

$$\frac{\partial f}{\partial \boldsymbol{\mu}_c} = 0$$

$$\Leftrightarrow - \sum_{i \in \mathcal{I}_c} \Sigma_c^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_c) = 0$$

$$\Leftrightarrow |\mathcal{I}_c| \boldsymbol{\mu}_c - \sum_{i \in \mathcal{I}_c} \mathbf{x}^{(i)} = 0$$

$$\Leftrightarrow \boldsymbol{\mu}_c = \frac{1}{|\mathcal{I}_c|} \sum_{i \in \mathcal{I}_c} \mathbf{x}^{(i)}$$

# Quadratic Discriminant Analysis ((QDA

---

הפרמטרים של המודל יהיו:

$$p_y(c) = \frac{|\mathcal{I}_c|}{N}$$

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

$$\Sigma_c = \frac{1}{|\mathcal{I}_c|} \sum_{i \in \mathcal{I}_c} \left( \mathbf{x}^{(i)} - \mu_c \right) \left( \mathbf{x}^{(i)} - \mu_c \right)^T$$

הצגנו כאן משערך אינטואיטבי להסתברות השיוך לכל מחלקה. אפשר להראות שזו התוצאה המתקבלת גם מהנחת מודל פילוג מולטינומיאלי למשתנה הקטגורי של המחלקה.

# Quadratic Discriminant Analysis ((QDA

---

עבור פונקציית מחיר של סיכוי הטעות, החזאי האופטימלי יהיה:

$$\begin{aligned}\hat{y} = h(\mathbf{x}) &= \arg \max_y p_{\mathbf{x}|y}(\mathbf{x}|y; \boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y) p_y(y) \\ &= \arg \max_y -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_y)^\top \boldsymbol{\Sigma}_y^{-1}(\mathbf{x} - \boldsymbol{\mu}_y) + \log \left( \frac{p_y(y)}{\sqrt{|\boldsymbol{\Sigma}_y|}} \right)\end{aligned}$$



# המקרה הבינארי - משטח הפרדה ריבועי

עבור המקרה של סיווג בינארי (סיווג לשתי מחלקות) מתקבל החזאי הבא:

$$h(x) = \begin{cases} 1 & \mathbf{x}^T C \mathbf{x} + \mathbf{a}^T \mathbf{x} + b > 0 \\ 0 & \text{otherwise} \end{cases}$$

כאשר:

$$C = \frac{1}{2}(\Sigma_0^{-1} - \Sigma_1^{-1})$$

$$\mathbf{a} = \Sigma_1^{-1} \boldsymbol{\mu}_1 - \Sigma_0^{-1} \boldsymbol{\mu}_0$$

$$b = \frac{1}{2} (\boldsymbol{\mu}_0^T \Sigma_0^{-1} \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1^T \Sigma_1^{-1} \boldsymbol{\mu}_1) + \log \left( \frac{\sqrt{|\Sigma_0|} p_y(1)}{\sqrt{|\Sigma_1|} p_y(0)} \right)$$

התנאי שקיבלנו  $\mathbf{x}^T C \mathbf{x} + \mathbf{a}^T \mathbf{x} + b > 0$  הוא ריבועי ב  $x$  ומכאן מקבל האלגוריתם את שמו.

# (Linear Discriminant Analysis (LDA

- מניח שלפונקציות הפילוג של המחלקות השונות יש את אותה מטריצת הקווריאנס.
- שהפרמטרים של המודל יהיו כעת:
  - וקטור תוחלת  $\mu_c$  עבור כל אחד מהמחלקות.
  - מטריצת covariance אחת  $\Sigma$  משותפת לכל המחלקות.

$$p_{\mathbf{x}|y}(\mathbf{x}|c; \mu_c, \Sigma) = \frac{1}{\sqrt{(2\pi)^D |\Sigma|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu_c)^\top \Sigma^{-1}(\mathbf{x} - \mu_c)\right)$$

# (Linear Discriminant Analysis (LDA

---

פתרון בעיית ה MLE:

$$p_y(c) = \frac{|\mathcal{I}_c|}{N}$$

$$\boldsymbol{\mu}_c = \frac{1}{|\mathcal{I}_c|} \sum_{i \in \mathcal{I}_c} \mathbf{x}^{(i)}$$

$$\boldsymbol{\Sigma} = \frac{1}{N} \sum_i \left( \mathbf{x}^{(i)} - \boldsymbol{\mu}_{y^{(i)}} \right) \left( \mathbf{x}^{(i)} - \boldsymbol{\mu}_{y^{(i)}} \right)^T$$

# (Linear Discriminant Analysis (LDA

---

עבור פונקציית מחיר של misclassification rate, החזאי האופטימאלי המתקבל ממודל זה הינו:

$$\begin{aligned}\hat{y} = h(\mathbf{x}) &= \arg \max_y p_{\mathbf{x}|y}(\mathbf{x}|y; \boldsymbol{\mu}_c, \Sigma) p_y(y) \\ &= \arg \max_y -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_y)^\top \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_y) + \log(p_y(y)) \\ &= \arg \min_y \mathbf{x}^\top \Sigma^{-1} \boldsymbol{\mu}_y - \frac{1}{2} \boldsymbol{\mu}_y^\top \Sigma^{-1} \boldsymbol{\mu}_y - \log(p_y(y))\end{aligned}$$

עבור המקרה של סיווג בינארי (סיווג לשני מחלקות) מתקבל החזאי הבא:

$$h(x) = \begin{cases} 1 & \mathbf{a}^T \mathbf{x} + b > 0 \\ 0 & \text{otherwise} \end{cases}$$

כאשר:

$$\mathbf{a} = \Sigma^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)$$

$$b = \frac{1}{2} (\boldsymbol{\mu}_0^T \Sigma^{-1} \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1) + \log \left( \frac{p_y(1)}{p_y(0)} \right)$$

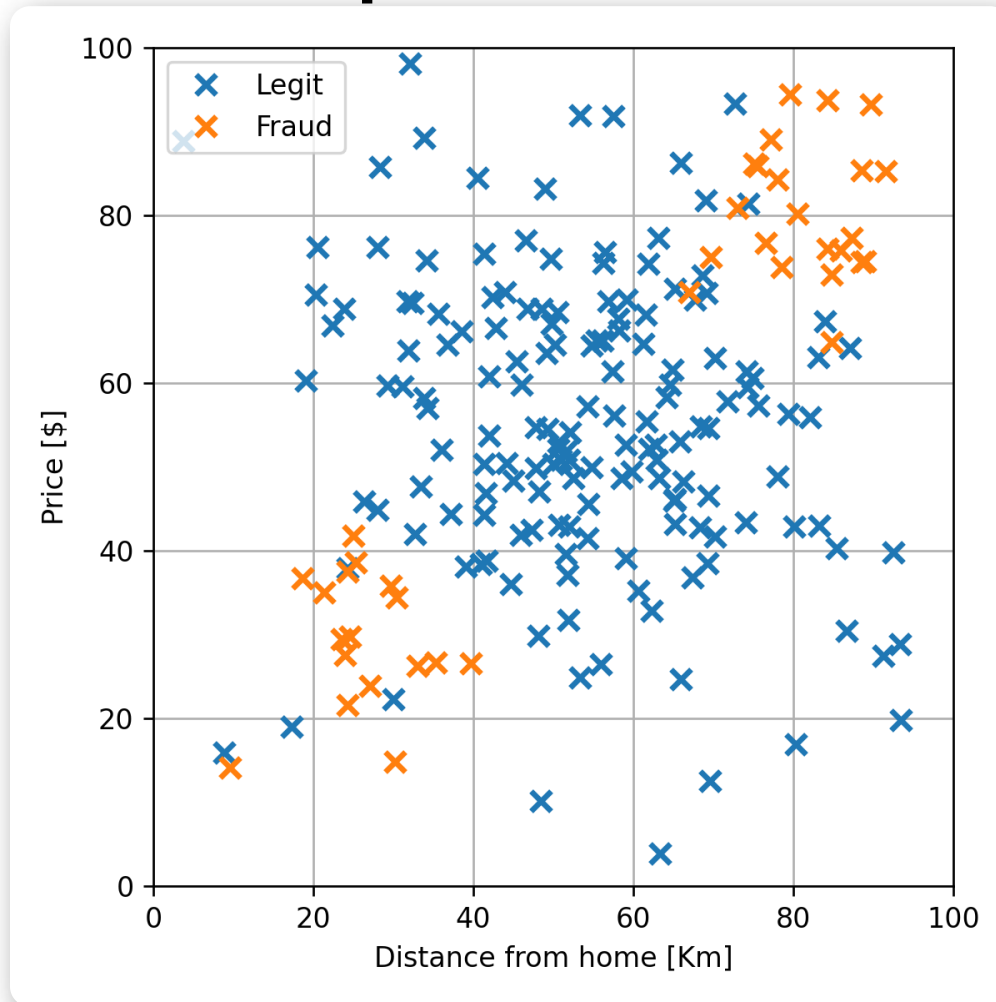
התנאי שקיבלנו  $\mathbf{a}^T \mathbf{x} + b > 0$  הוא לינארי ב  $\mathbf{x}$  ומכאן מקבל האלגוריתם את שמו.

# המקרה הכללי (לא בינארי)

---

במקרה הכללי המרחב יהיה מחולק ל  $C$  איזורים שהשפות שלהם יהיו מורכבות מהמישורים המתקבלים מהשפות שבין כל זוג מחלקות. דוגמא למקרה עם 3 מחלקות תופיע בתרגול.

נסתכל שוב על הבעיה של חיזוי עסקאות שחשודות כהונאות:



# התאמה של מודל QDA

$$p_y(0) = \frac{|\mathcal{I}_0|}{N} = 0.81$$

$$p_y(1) = \frac{|\mathcal{I}_1|}{N} = 0.19$$

$$\boldsymbol{\mu}_0 = \frac{1}{|\mathcal{I}_0|} \sum_{i \in \mathcal{I}_0} \mathbf{x}^{(i)} = [55.1, 54.6]^\top$$

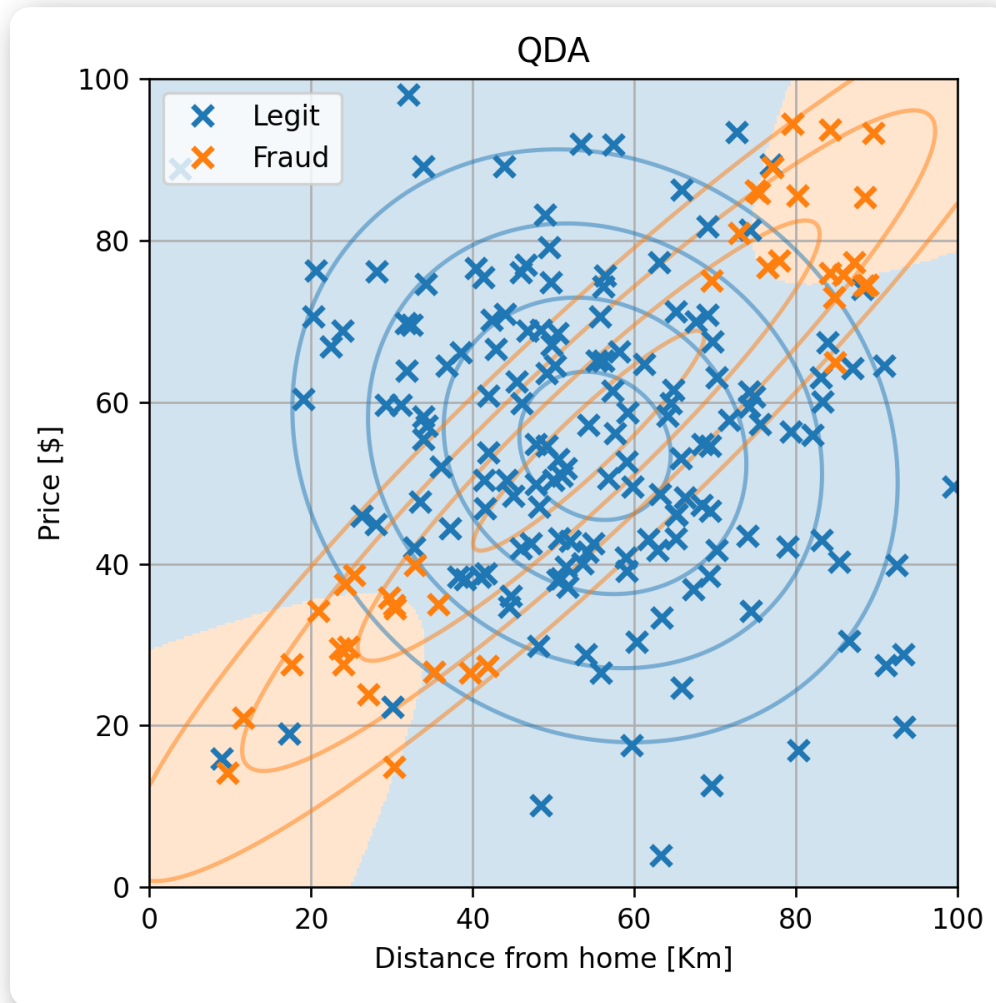
$$\boldsymbol{\mu}_1 = \frac{1}{|\mathcal{I}_1|} \sum_{i \in \mathcal{I}_1} \mathbf{x}^{(i)} = [54.4, 55.2]^\top$$

$$\boldsymbol{\Sigma}_0 = \frac{1}{|\mathcal{I}_0|} \sum_i \left( \mathbf{x}^{(i)} - \boldsymbol{\mu}_{y^{(i)}} \right) \left( \mathbf{x}^{(i)} - \boldsymbol{\mu}_{y^{(i)}} \right)^\top = \begin{bmatrix} 350.9 & -42.9 \\ -42.9 & 336 \end{bmatrix}$$

$$\boldsymbol{\Sigma}_1 = \frac{1}{|\mathcal{I}_1|} \sum_i \left( \mathbf{x}^{(i)} - \boldsymbol{\mu}_{y^{(i)}} \right) \left( \mathbf{x}^{(i)} - \boldsymbol{\mu}_{y^{(i)}} \right)^\top = \begin{bmatrix} 817.9 & 730.5 \\ 730.5 & 741.7 \end{bmatrix}$$

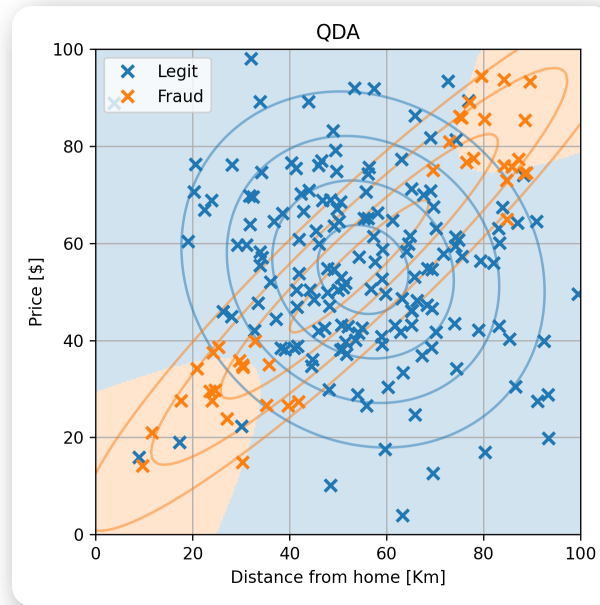


# התאמה של מודל QDA



שגיאת החיזוי (miscalssification rate) על ה test set הינה 0.08.

# הבעיה של הגישה הגנרטיבית הפרמטרית

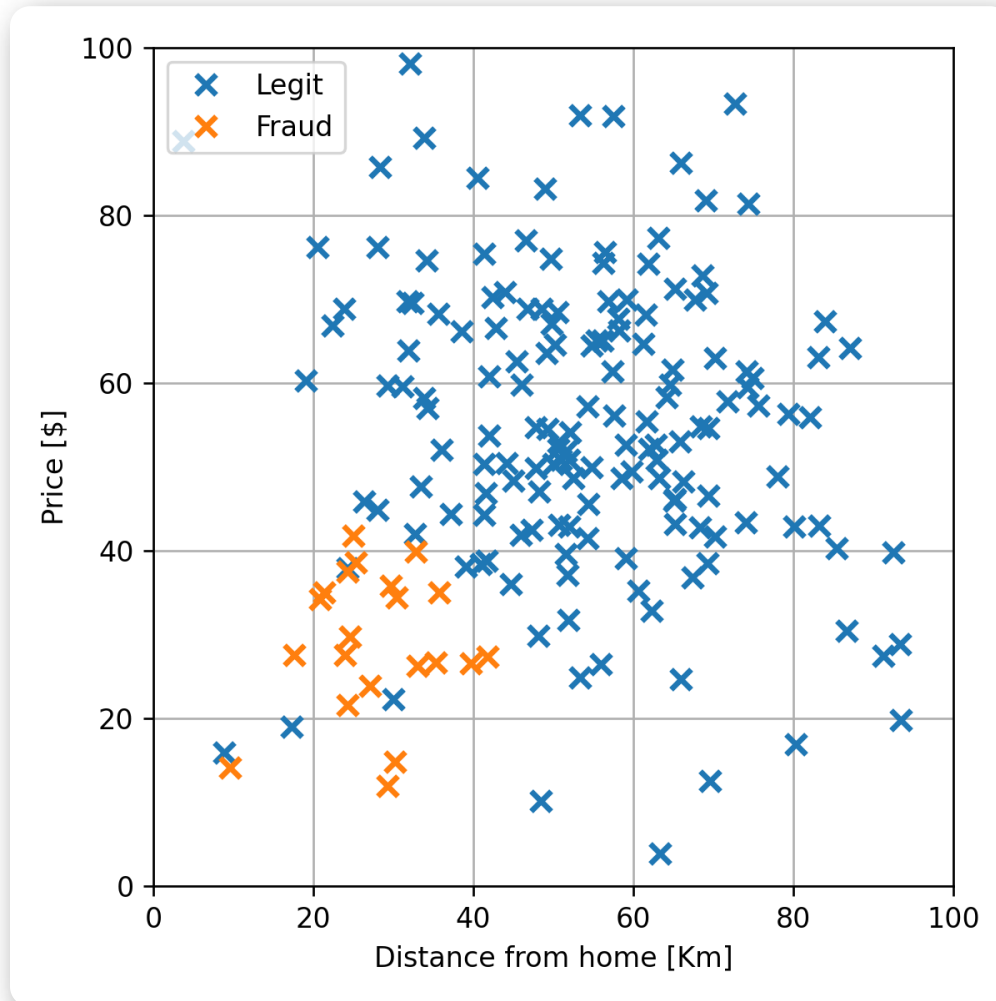


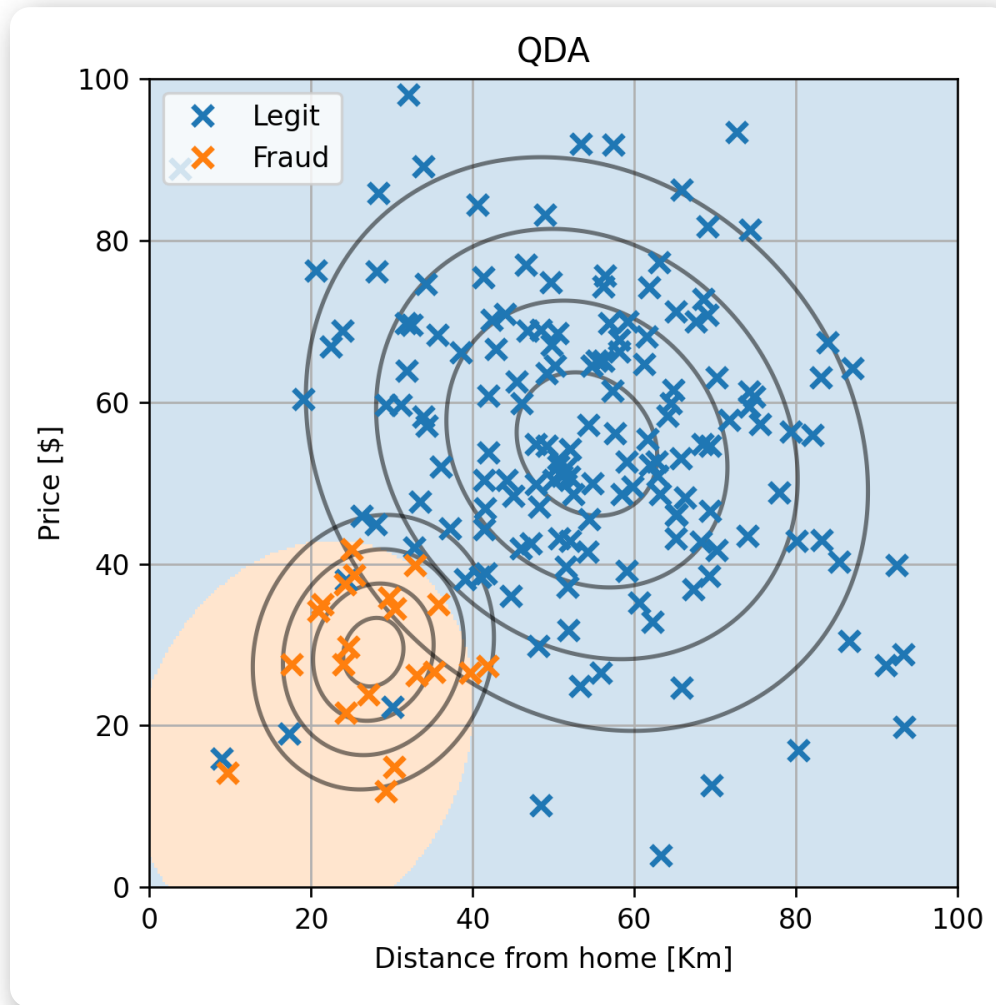
- הינו רוצים מודל אשר יכול לייצג בנפרד את שני האיזורים.
- לצערינו המבחר של המודלים בהם אנו יכולים לא גדול.
- המגבלה הזו נובעת מהצורך שהמודל ייצג פילוג חוקיים.

**הערה:** קיימות הרבה הרחבות שנידונות במקורסים מתקדמים יותר, למשל, תערובת של גאוסיאניים.

# דוגמא למדגם שמתאים למודל של QDA

לצורך הדגמה נסתכל על גירסא של המדגם שבה יש רק איזור אחד של ההונאות:





## נציג גם את התוצאה המתקבלת ממודל ה-LDA:

