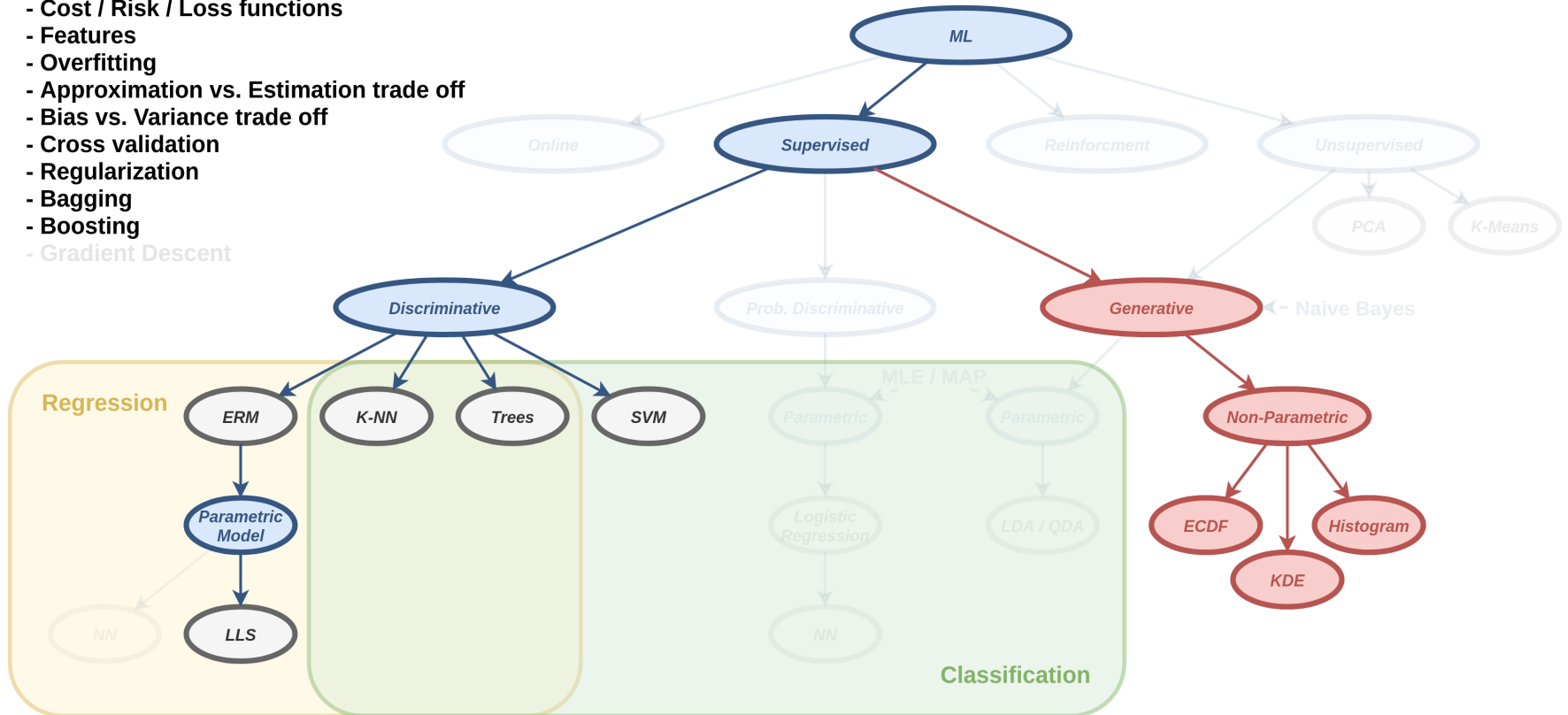


הרצאה 7 - שיערוך פילוג בשיטות לא פרמטריות

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



דיסקרימינטיבי vs. גנרטיבי

הגישה הדיסקרימינטיבית

מדגם



חזאי בעל ביצועים טובים על המדגם

הגישה הגנרטיבית

מדגם



פילוג על סמך המדגם



חזאי אופטימאלי בהינתן הפילוג

הקשר לבעיות **unsupervised learning**

- בקורס זה לא נעסוק כמעט בבעיות **unsupervised learning**.
- בבעיות **unsupervised learning** המדגם מכיל סוג אחד של משתנים x .
- ננסה ללמוד מהם התכונות שמאפיינות את הדגימות במדגם.
- אחת הדרכים הטובות ביותר לעשות זאת היא על ידי שיערוך הפילוג שלהם.

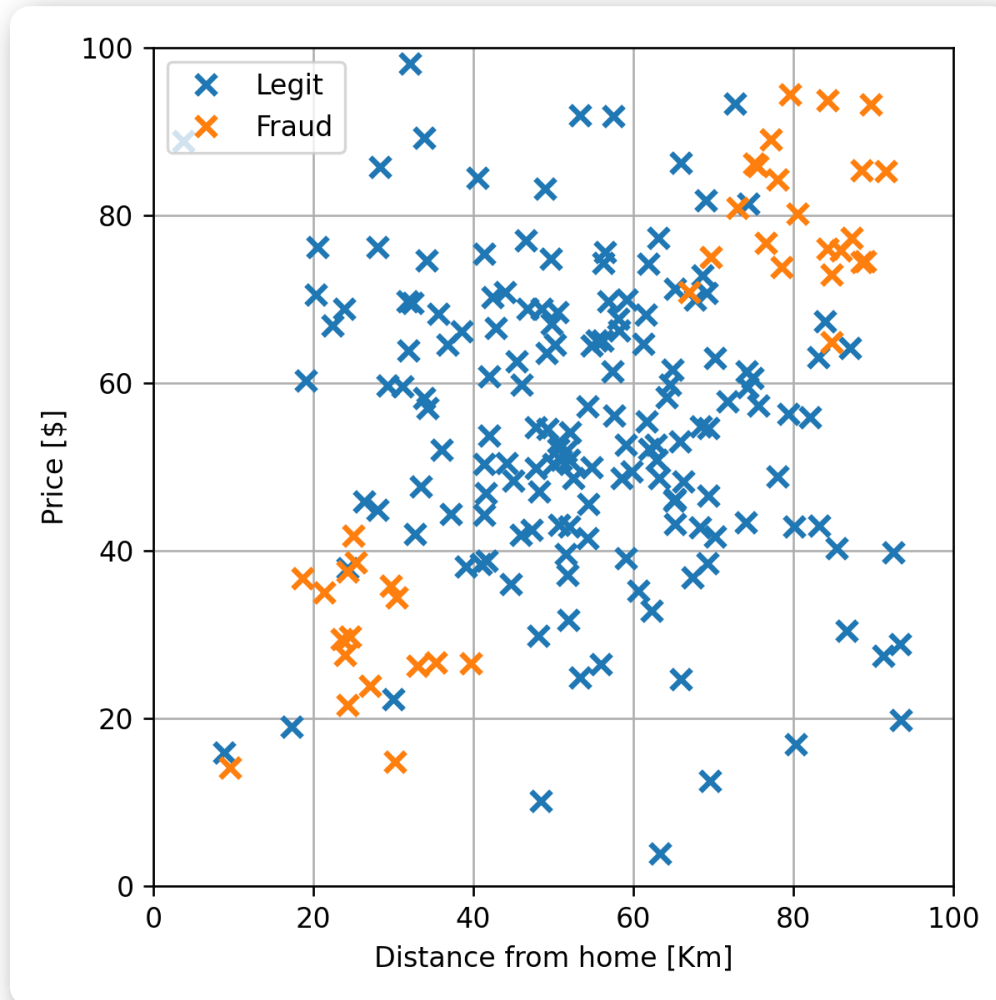
הבעיה של בניית מודל הסתברותי של משתנים אקראיים מתוך מדגם מכונה **בעיית שיערוך (estimation)**. את המודל ההסתברותי אנו נבטא בעזרת אחת מהפונקציות הבאות:

- פונקציית ההסתברות (probability mass function - PMF)
- פונקציית צפיפות ההסתברות (probability density function - PDF)
- פונקציית הפילוג המצרפית (cumulative distribution function CDF).

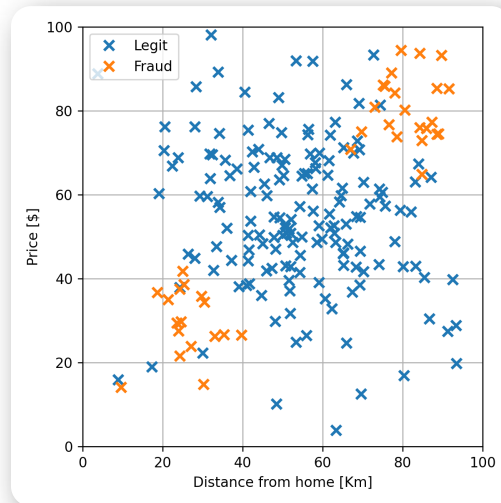
- **בבעיות חיזוי** אנו מועניינים לחזות את ערכו של **משתנה אקראי**, לרוב על סמך משתנה / וקטור אקראי בודד (**דגימה יחידה**).

- **בבעיות שיערוך** אנו מעוניינים לבנות **מודל הסתברותי** של משתנה / משתנים אקראיים לרוב על סמך **הרבה דגימות**.

נסתכל לדוגמא על המדגם של הונאות אשראי מהרצאה הקודמת:

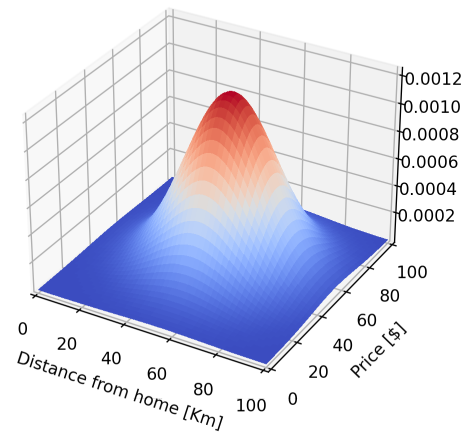


נרצה לשערך את הפילוג של המשתנים על פי מדגם זה

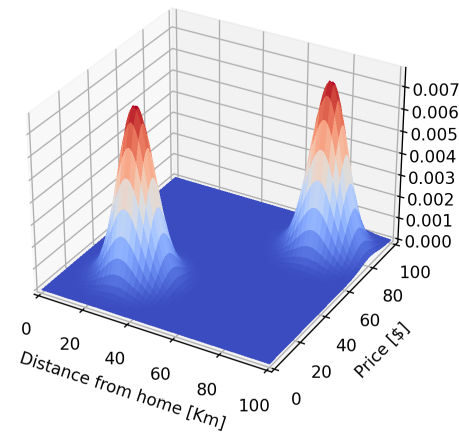


לדוגמא היינו רוצים למצוא פונקציות אשר יתארו את הפילוג של הדגימות החוקיות ושל ההונאות:

Legit PDF - $p_{x|y}(x|0)$



Fraud PDF - $p_{x|y}(x|1)$



שיערוך א-פרמטריות

בהרצאה הקרובה נעסוק בשיטות שיערוך אשר מכונות שיטות לא פרמטריות או א-פרמטריות, מהות השם תהיה ברורה יותר אחרי שנציג בהרצאה הבאה את הנושא של שיטות פרמטריות.

שיערוך ההסתברות של מאורע

דוגמא

נניח שיש בידינו את המדגם הבא של מדידות של זמני נסיעה (בדקות) מחיפה לתל אביב על כביש החוף:

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

ברצונינו לשערך את ההסתברות של המאורע שנסיעה מסויימת תיקח פחות משעה, $A = \{x < 60\}$.

שיערוך ההסתברות של מאורע

דוגמא

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

נשערך שהסתברות זו שווה למספר הפעמים היחסי שמאורע זה קרה במדגם הנתון:

$$\Pr(A) \approx \hat{p}_{A,\mathcal{D}} = 0.3$$

- **נשתמש בסימון "כובע" לציון גודל שאותו אנו חוזים / משערכים באופן אמפירי.**
- **נציין את העובדה שמשערך תלוי במדגם שבו השתמשנו על ידי הוספת \mathcal{D} מתחת למשערך.**

מדידה אמפירית (empirical measure) / משערך הצבה

בהינתן מדגם מסויים $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=0}^N$, המדידה האמפירית, $\hat{p}_{A,\mathcal{D}}$, הינה שיערוך של ההסתברות, $Pr(A)$, והיא מחושבת באופן הבא:

$$\hat{p}_{A,\mathcal{D}} = \frac{1}{N} \sum_{i=1}^N I\{\mathbf{x}^{(i)} \in A\}$$

נוכל כעת להשתמש בשיטה זו על מנת לנסות ולשערך את הפילוג של משתנים אקראיים.

משתנה אקראי דיסקרטי

דוגמא 1 - משתנה בינארי

- תוצאת הטלה של מטבע לא הוגן.
- הטלנו את המטבע 10 פעמים וקיבלנו:

$$\mathcal{D} = \{x^{(i)}\} = \{0, 0, 0, 0, 1, 0, 0, 1, 0, 0\}$$

מה ה PMF של x ?

משתנה אקראי דיסקרטי

דוגמא 1 - משתנה בינארי

$$\mathcal{D} = \{x^{(i)}\} = \{0, 0, 0, 0, 1, 0, 0, 1, 0, 0\}$$

גם כאן נשערך את ההסתברויות של הערכים ש x מקבל על פי השכיחות שלהם במדגם:

$$p_x(x) \approx \hat{p}_{x,\mathcal{D}}(x) = \begin{cases} 0.8 & 0 \\ 0.2 & 1 \end{cases}$$

• זו למעשה מדידה אמפירית של המאורע ש $\{x = x\}$.

משתנה אקראי דיסקרטי

דוגמא 2 - משתנה לא בינארי

- תוצאת הטלה של קוביה לא הוגנת.
- הטלנו את הקוביה 10 פעמים וקיבלנו:

$$\mathcal{D} = \{x^{(i)}\} = \{3, 2, 5, 1, 2, 6, 2, 5, 5, 3\}$$

מה ה PMF של x ?

משתנה אקראי דיסקרטי

דוגמא 2 - משתנה לא בינארי

$$\mathcal{D} = \{x^{(i)}\} = \{3, 2, 5, 1, 2, 6, 2, 5, 5, 3\}$$

בדיוק כמו קודם, נשערך את ההסתברות לקבל כל ערך לפי השכיחות שלו במדגם:

$$p_x(x) \approx \hat{p}_{x,\mathcal{D}}(x) = \begin{cases} 0.1 & 1 \\ 0.3 & 2 \\ 0.2 & 3 \\ 0 & 4 \\ 0.3 & 5 \\ 0.1 & 6 \end{cases}$$

בהינתן מדגם מסויים $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=0}^N$, נוכל לשערך את ה PMF של משתנה / וקטור אקראי דיסקרטי באופן הבא:

$$\hat{p}_{\mathbf{x},\mathcal{D}}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N I\{\mathbf{x}^{(i)} = \mathbf{x}\}$$

שימו לב שמובטח לנו שנקבל פונקציית הסתברות חוקית (חיובית שהסכום עליה שווה ל1).

שיערוך הפילוג המצרפי

נזכור כי פונקציית הפילוג המצרפי (ה CDF) מוגדרת באופן הבא:

$$F_{\mathbf{x}}(\mathbf{x}) = \Pr(\{x_j \leq x_j \forall j\})$$

נוכל אם כן לשערך גודל זה על ידי שימוש במדידה האמפירית בעבור המאורע של $A = \{x_j \leq x_j \forall j\}$ באופן הבא:

$$\hat{F}_{\mathbf{x},\mathcal{D}}(\mathbf{x}) = \hat{p}_{A,\mathcal{D}} = \frac{1}{N} \sum_{i=1}^N I\{x_j \leq x_j \forall j\}$$

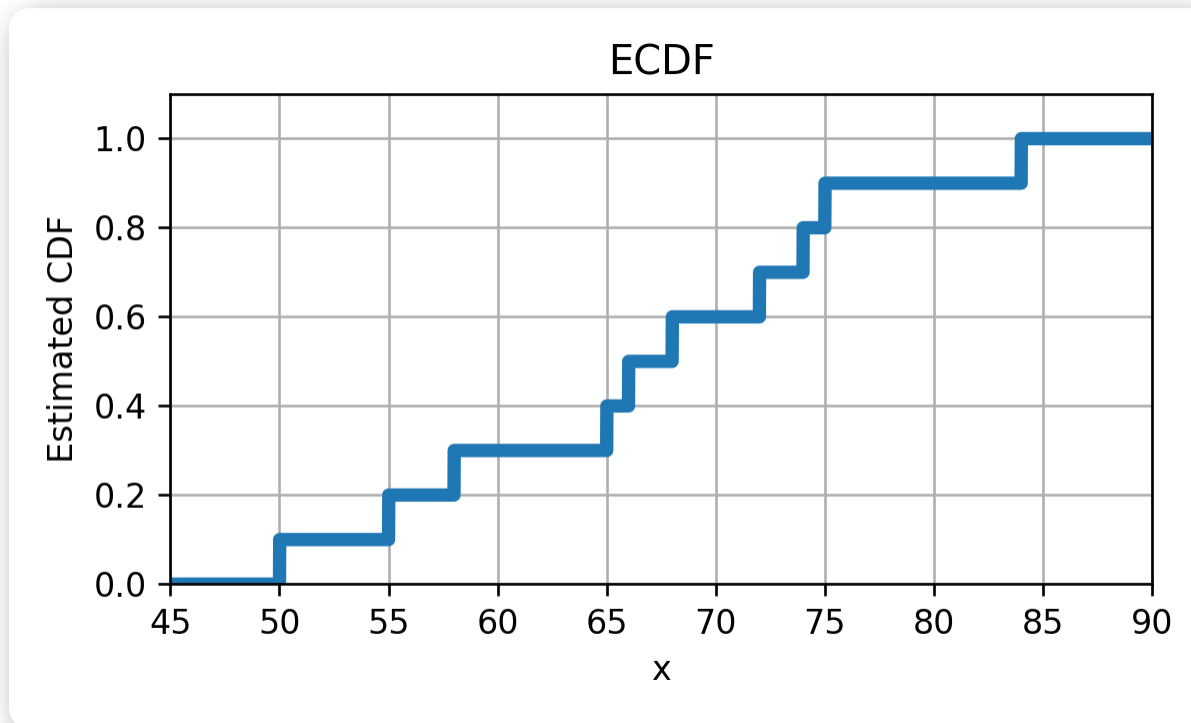
משערך זה נקרא empirical cumulative distribution function (ECDF).

נשערך את הפילוג המצרפי של זמני הנסיעה בכביש החוף

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

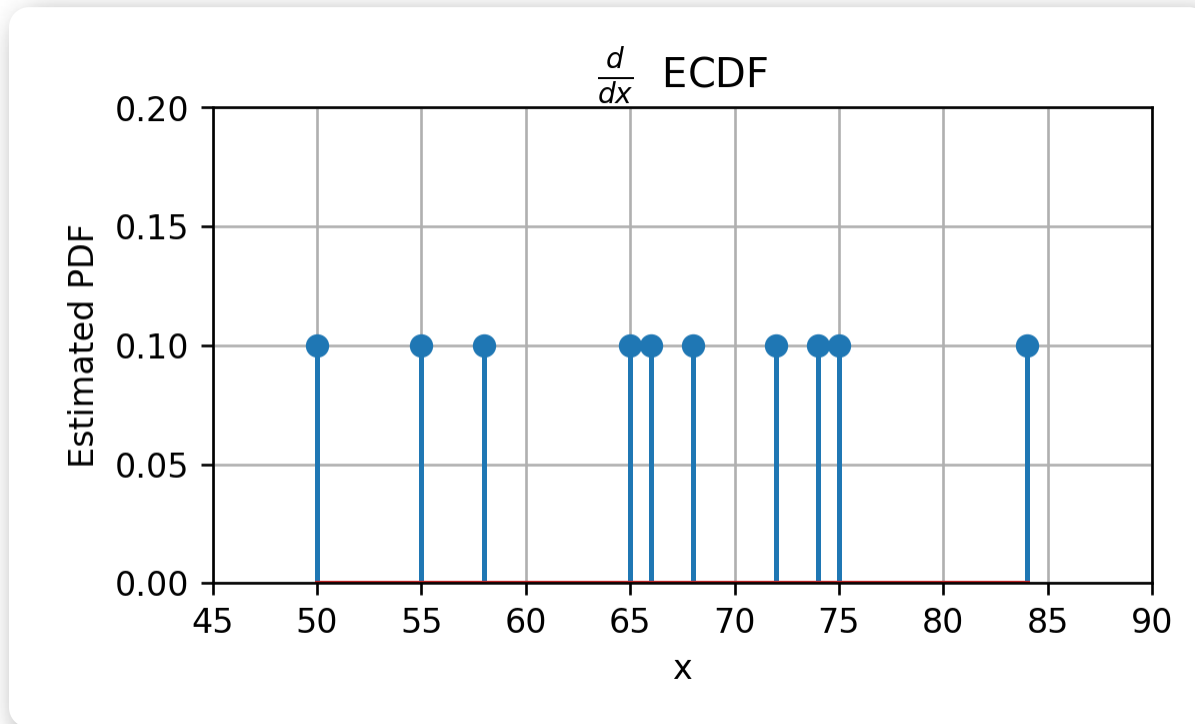
$$\hat{F}_{\mathbf{x}, \mathcal{D}}(\mathbf{x}) = \begin{cases} 0 & x < 50 \\ 0.1 & 50 \leq x < 55 \\ 0.2 & 55 \leq x < 58 \\ 0.3 & 58 \leq x < 65 \\ 0.4 & 65 \leq x < 66 \\ 0.5 & 66 \leq x < 68 \\ 0.6 & 68 \leq x < 72 \\ 0.7 & 72 \leq x < 74 \\ 0.8 & 74 \leq x < 75 \\ 0.9 & 75 \leq x < 84 \\ 1 & 84 \leq x \end{cases}$$

זוהי למעשה פונקציה קבועה למקוטעין אשר נראית כך:



בעיה: איך נראה ה PDF?

ככה:



פונקציה כזו היא לא מאד שימושית.

נסיון לשערך PDF על ידי קוונטיזציה של משתנה רציף.

- **נחלק את טווח הערכים למספר סופי של חלקים המכונים bins (תאים).**
- **נשתמש במדידה אמפירית על מנת לשערך את ההסתברות להימצא בכל תא.**

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

נחלק את התחום ל 5 קטעים:

$$[45, 54), [54, 63), [63, 72), [72, 81), [81, 90]$$

ההסתברות להיות בכל bin הינה:

$$\hat{p}_{\{45 \leq x < 54\}, \mathcal{D}} = 0.1$$

$$\hat{p}_{\{54 \leq x < 63\}, \mathcal{D}} = 0.2$$

$$\hat{p}_{\{63 \leq x < 72\}, \mathcal{D}} = 0.3$$

$$\hat{p}_{\{72 \leq x < 81\}, \mathcal{D}} = 0.3$$

$$\hat{p}_{\{81 \leq x \leq 90\}, \mathcal{D}} = 0.1$$

יש לבחור את ה bins כך שיכסו את התחום ולא יחפפו.

בכדי להפוך את ההסתברויות לצפיפות הסתברות נרצה "למרוח" את ההסתברות שקיבלנו באופן אחיד על פני ה bin.

$$\hat{p}_{x,\mathcal{D}}(x) = \begin{cases} \frac{1}{\text{size of bin } 1} \hat{p}_{\{x \text{ in bin } 1\},\mathcal{D}} & x \text{ in bin } 1 \\ \vdots \\ \frac{1}{\text{size of bin } B} \hat{p}_{\{x \text{ in bin } B\},\mathcal{D}} & x \text{ in bin } B \end{cases}$$

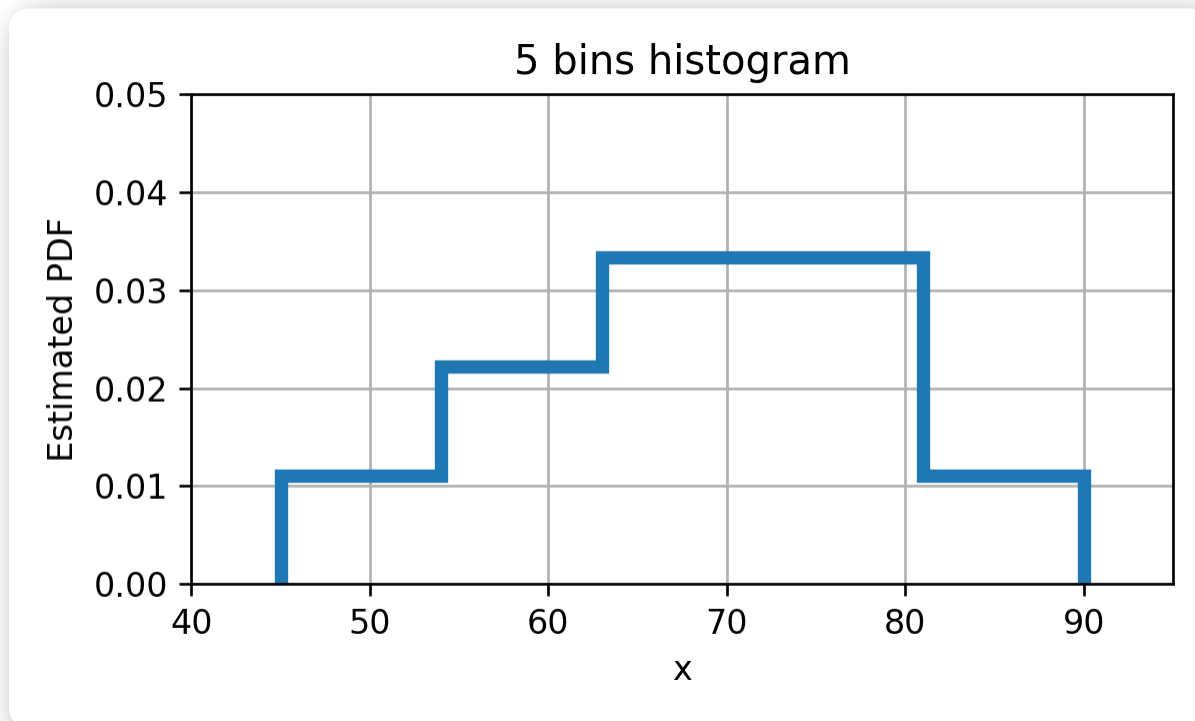
$$\hat{p}_{\{45 \leq x < 54\}, \mathcal{D}} = 0.1$$

$$\hat{p}_{\{54 \leq x < 63\}, \mathcal{D}} = 0.2$$

$$\hat{p}_{\{63 \leq x < 72\}, \mathcal{D}} = 0.3$$

$$\hat{p}_{\{72 \leq x < 81\}, \mathcal{D}} = 0.3$$

$$\hat{p}_{\{81 \leq x \leq 90\}, \mathcal{D}} = 0.1$$



היסטוגרמה - ניסוח פורמאלי

בהינתן מדגם מסויים $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=0}^N$, ההיסטוגרמה הינה שיערוך של ה PDF של משתנה / וקטור אקראי והיא מחושבת באופן הבא:

1. מחלקים את תחום הערכים ש \mathbf{x} יכול לקבל ל bins (תאים) לא חופפים אשר מכסים את כל התחום.
 2. לכל bin משערכים את ההסתברות של המאורע שבו \mathbf{x} יהיה בתוך התא.
 3. הערך של פונקציית הצפיפות בכל תא תהיה ההסתברות המשוערכת להיות בתא חלקי גודל התא.
- לבחירת ה bins יש השפעה גדולה על איכות השיערוך שנקבל. ננסה להבין את השיקולים בבחירת ה bins.

היסטוגרמה - המקרה הסקלרי

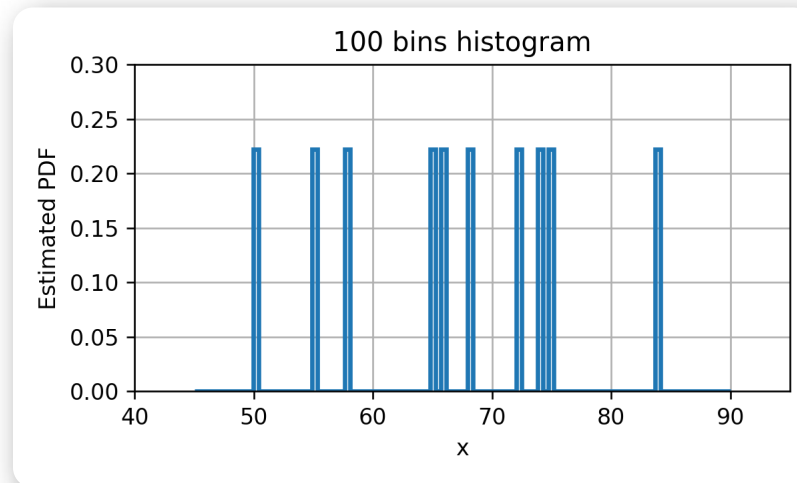
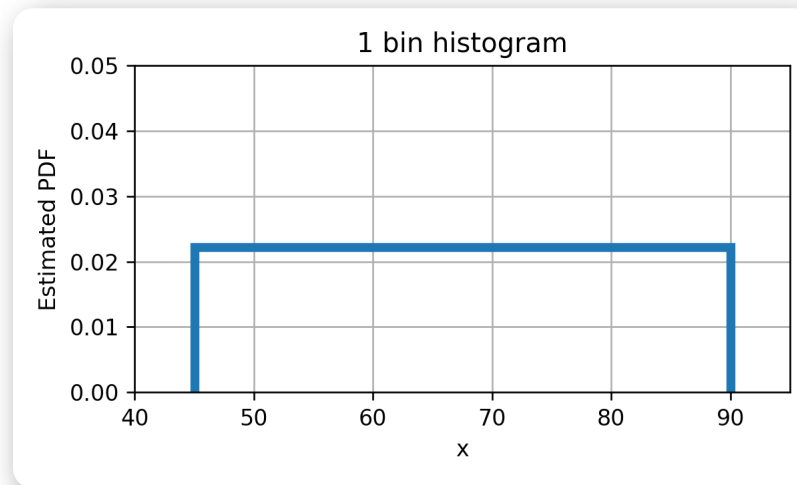
• מספר התאים B .

• l_b ו r_b את הגבול השמאלי והימני התא ה b .

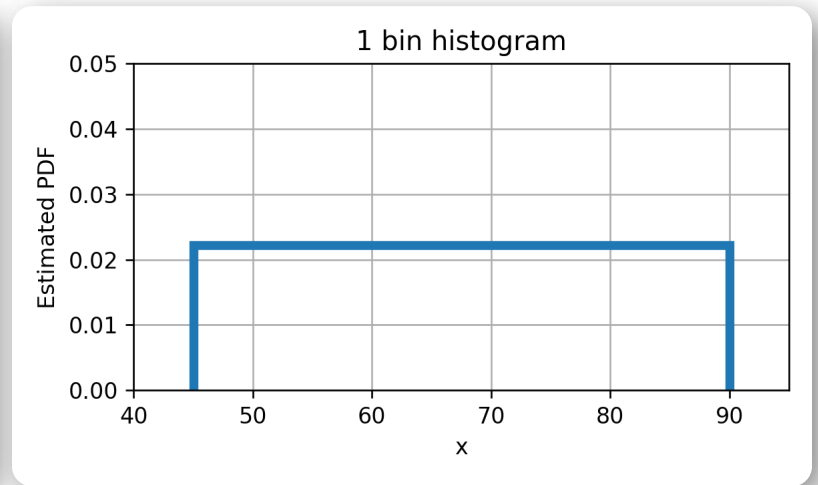
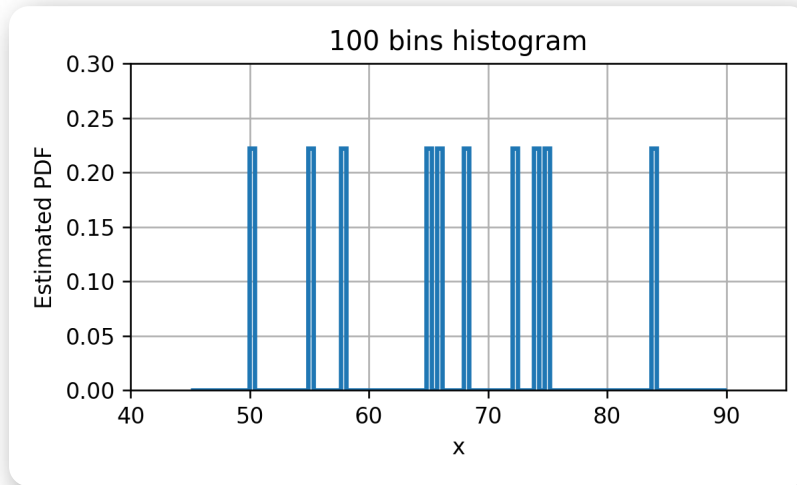
$$\hat{p}_{\mathbf{x}, \mathcal{D}}(\mathbf{x}) = \begin{cases} \frac{1}{N(r_1 - l_1)} \sum_{i=1}^N I\{l_1 \leq \mathbf{x}^{(i)} < r_1\} & l_1 \leq \mathbf{x} < r_1 \\ \vdots \\ \frac{1}{N(r_B - l_B)} \sum_{i=1}^N I\{l_B \leq \mathbf{x}^{(i)} < r_B\} & l_B \leq \mathbf{x} < r_B \end{cases}$$

של underfitting | Overfitting היסטוגרמה

דוגמא - שני מקרים קיצוניים



של underfitting | Overfitting היסטוגרמה



מספר תאים קטן

Underfitting: יכולת מוגבלת לקרב את ה PDF האמיתי.

מספר תאים גדול

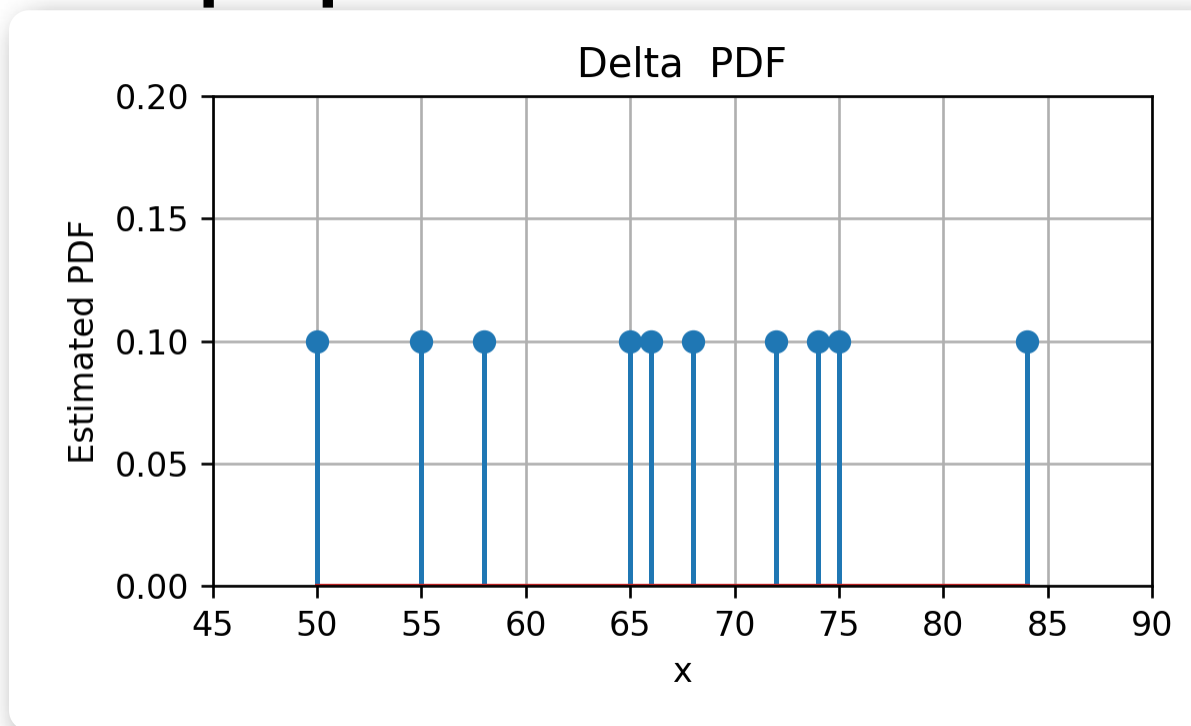
Overfitting: ההיסטוגרמה תתאר בצורה טובה את הדגימות אך לא את הפילוג האמיתי.

- מקובל לחלק ל k תאים אחידים בגודלם.
- מכיוון שה k האופטימאלי ישתנה מבעיה לבעיה, נאלץ לרוב לבחור אותו בעזרת ניסוי וטעיה.
- ישנם מספר כללי אצבע אשר במרבית המקרים יתנו תוצאה לא רעה.
- הכלל הנפוץ ביותר הינו לבחור את k להיות שורש מספר הדגימות במדגם (מעוגל כלפי מעלה):
$$k = \lceil \sqrt{N} \rceil$$

(Kernel Density Estimation (KDE

$\frac{1}{N}$ נתחיל מ PDF שבו אנו ממקמים פונקציית דלתא בגובה בכל נקודה אשר מופיעה במדגם.

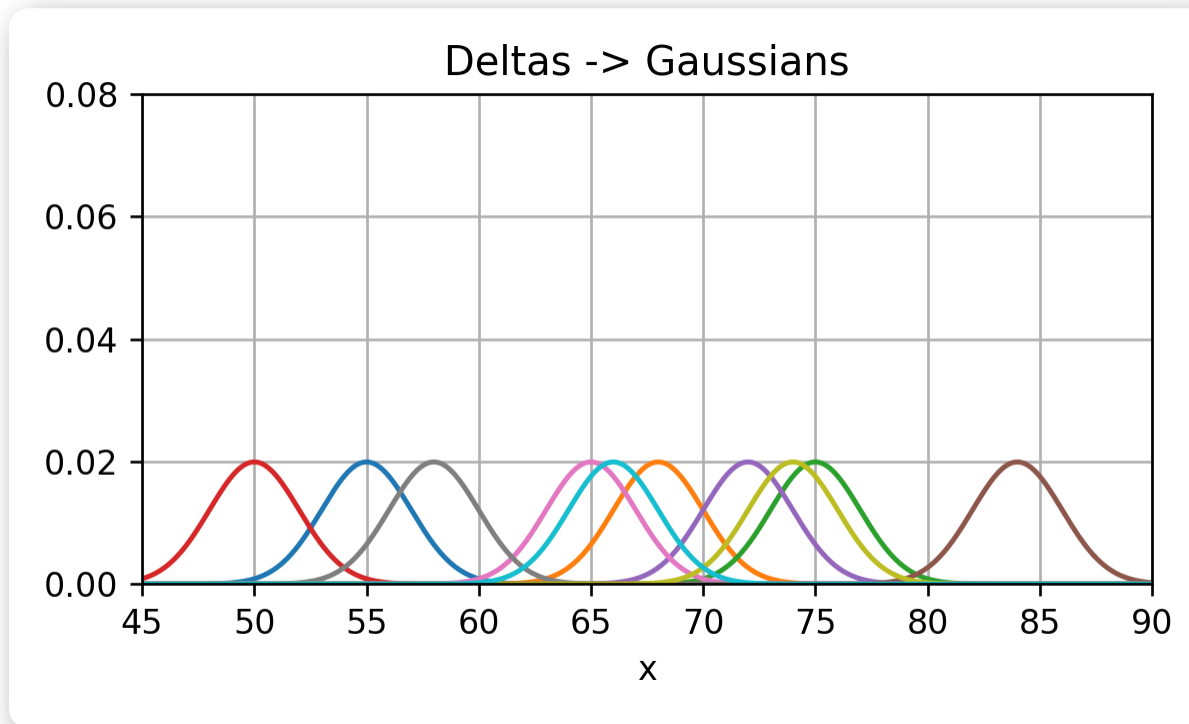
לדוגמא, בעבור זמני הנסיעה בכביש החוף נקבל:



(Kernel Density Estimation (KDE

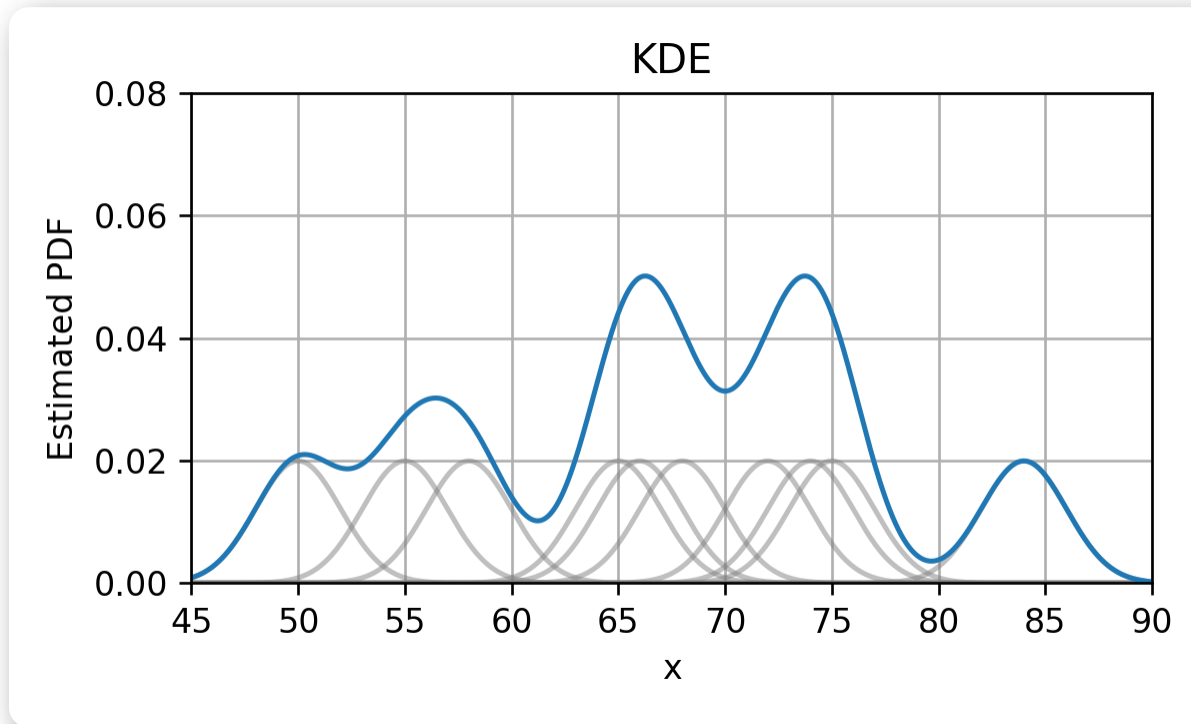
נחליף כל דלתא בפונקציית גרעין בעלת רוחב גדול מ-0.

לדוגמא גאוסיאנים:



(Kernel Density Estimation (KDE

נסכום את כל פונקציות הגרעין לקבלת ה PDF המשוער:



(Kernel Density Estimation (KDE

- פונקציות הגרעין (kernel) מכונות גם **Parzen window**.
- ומקובל לסמנם ב $\phi(\mathbf{x})$.

אם כן, משערך ה KDE נתון על ידי:

$$\hat{p}_{\mathbf{x},\phi,D}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \phi(\mathbf{x} - \mathbf{x}^{(i)})$$

הערה: תנאי מספיק והכרחי בכדי שנקבל PDF חוקי, הינו שפונקציית הגרעין תהיה בעצמה PDF חוקי.

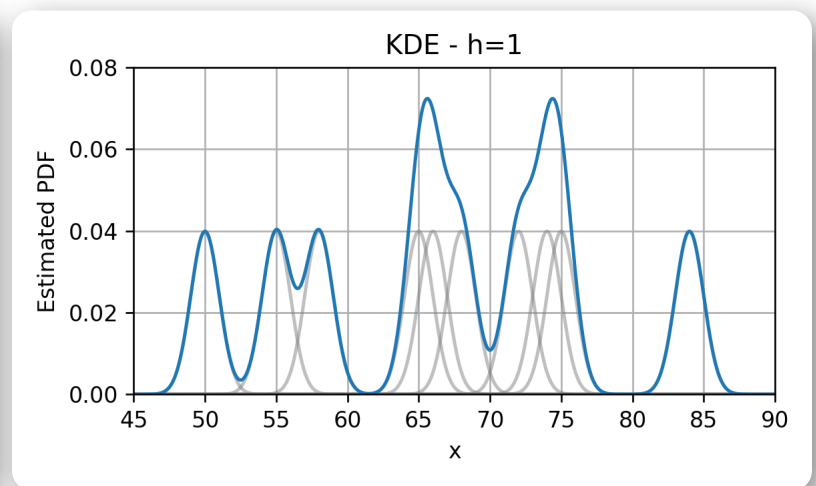
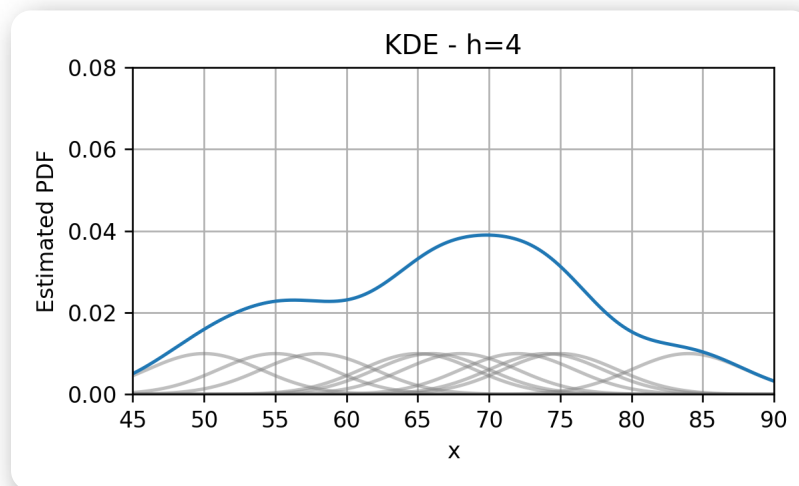
בהקשר של עיבוד אותות: למעשה אנו מבצעים קונבולוציה בין פונקציית הדלתאות לבין פונקציית הגרעין. נרצה שהגרעין ישמש כמעין low pass filter.

מקובל להוסיף פרמטר h אשר שולט ברוחב של הגרעין:

$$\phi_h(\mathbf{x}) = \frac{1}{h^D} \phi\left(\frac{\mathbf{x}}{h}\right)$$

בתוספת פרמטר זה המשערך יהיה:

$$\hat{p}_{\mathbf{x},\phi,h,\mathcal{D}}(\mathbf{x}) = \frac{1}{Nh^D} \sum_{i=1}^N \phi\left(\frac{\mathbf{x} - \mathbf{x}^{(i)}}{h}\right)$$



פונקציות גרעין נפוצות

שתי הבחירות הנפוצות ביותר לפונקציית הגרעין הינן:

1. חלון מרובע:

$$\phi_h(\mathbf{x}) = \frac{1}{h^D} I\{|x_j| \leq \frac{h}{2} \quad \forall j\}$$

כלל אצבע עבור חלון ריבועי: נבחר את גודל החלון אדפטיבית כך שיכלול \sqrt{N} דגימות מסביב לנקודה הנחקרת.

2. גאוסיאן:

$$\phi_\sigma(x) = \frac{1}{\sqrt{2\pi}\sigma^D} \exp\left(-\frac{\|x\|_2^2}{2\sigma^2}\right)$$

כלל אצבע לבחירת רוחב הגרעין במקרה הגאואסי הסקלרי:

$$\sigma = \left(\frac{4 \cdot \text{std}(x)^5}{3N}\right)^{\frac{1}{5}} \approx 1.06 \text{std}(x) N^{-\frac{1}{5}}$$

שיערוך של פילוגים מעורבים

- נניח שאנו רוצים לשערך את הפילוג המשותף של x ו y כאשר x הוא משתנה רציף ו y הוא משתנה בדיד.
- במקרים כאלה נוח לפרק את פונקציית הפילוג המשותף באופן הבא:

$$p_{x,y}(x, y) = p_{x|y}(x|y)p_y(y)$$

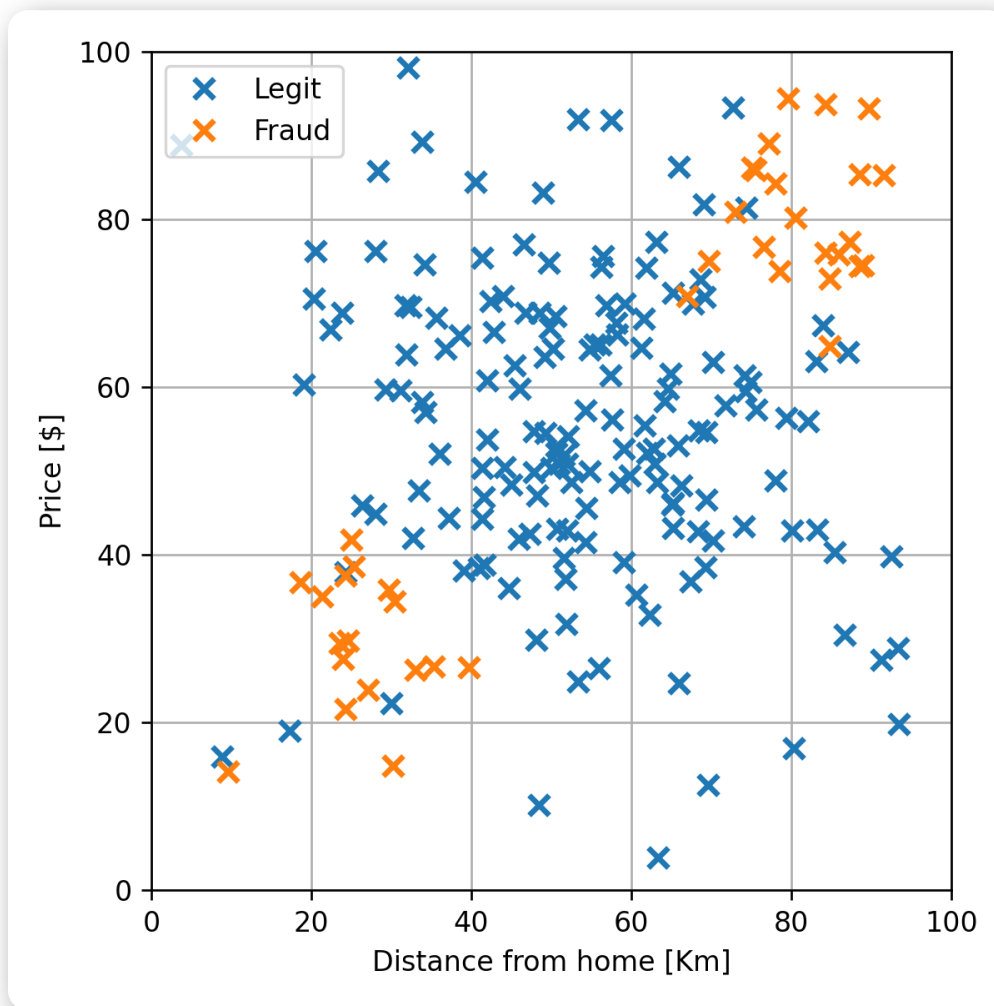
ולהפריד את בעיית השיערוך לשני חלקים:

1. השיערוך של $p_y(y)$

2. השיערוך של $p_{x|y}(x|y)$ - כאן נשערך את הפילוג בנפרד לכל ערך של y .

שיערוך של פילוגים מעורבים - דוגמא

נחזור לדוגמא של הונאות האשראי:



שיערוך של פילוגים מעורבים - דוגמא

נתחיל בשיערוך של y .

- בדיד ולכן נוכל לשערך את ה PMF שלו על פי השכיחות של הערכים במדגם.
- מתוך ה 200 עסקאות ישנם 160 עסקאות חוקיות ו 40 עסקאות שחשודות כהונאה. לכן:

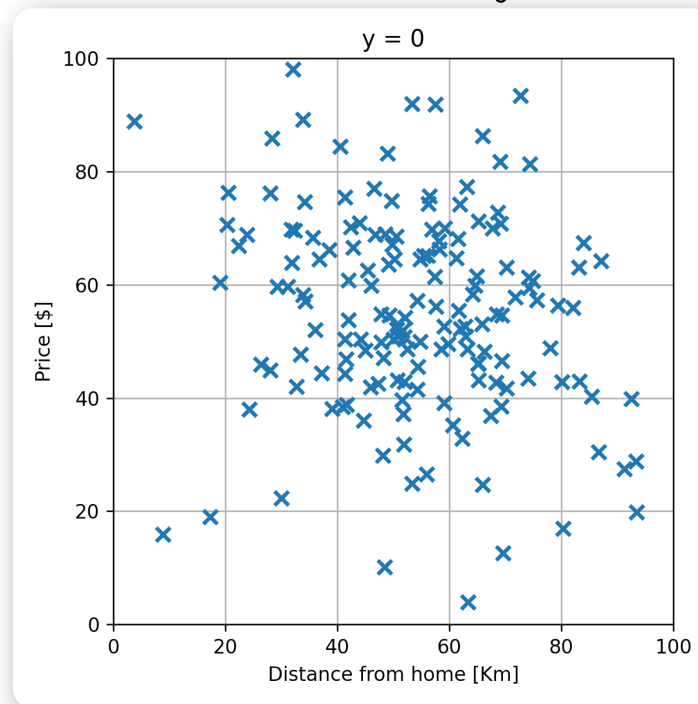
$$\hat{p}_{y,D}(y) = \begin{cases} \frac{160}{200} & 0 \\ \frac{40}{200} & 1 \end{cases} = \begin{cases} 0.8 & 0 \\ 0.2 & 1 \end{cases}$$

שיערוך של פילוגים מעורבים - דוגמא

נמשיך לשיערוך של $p_{x|y}(x|y)$.

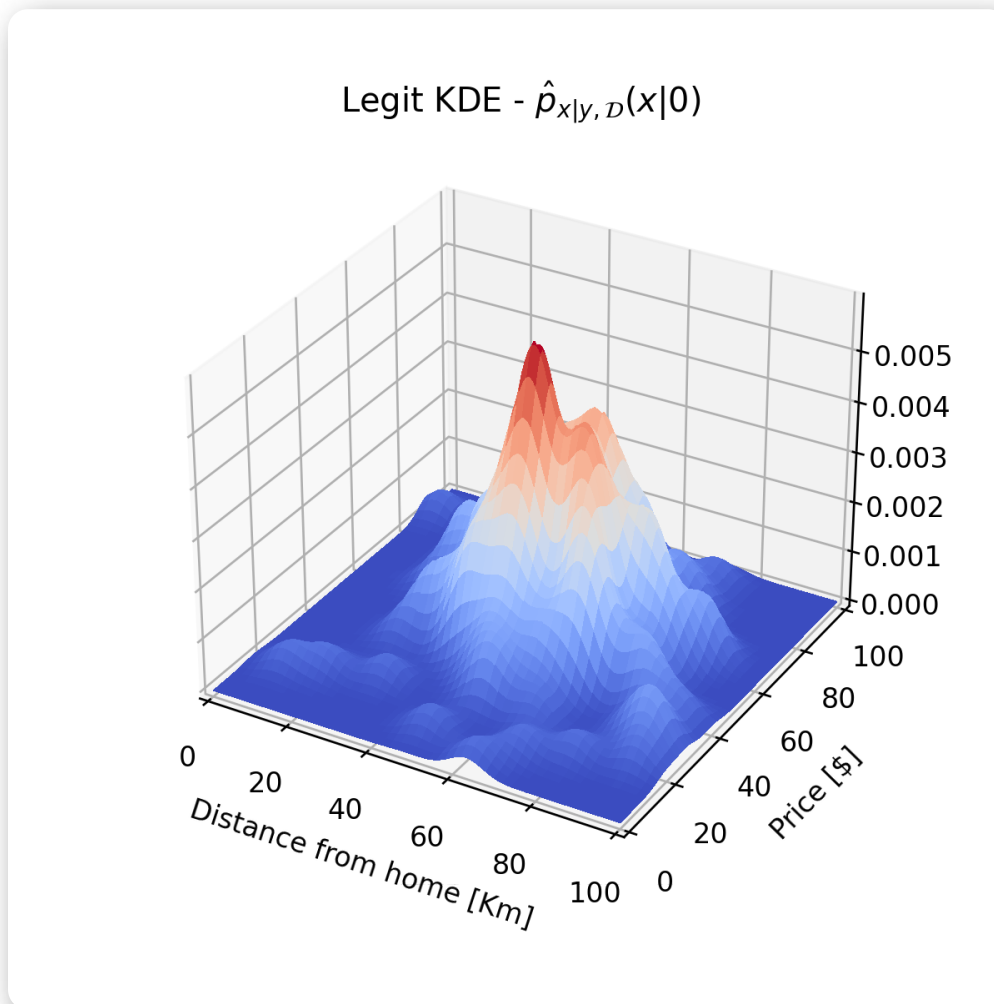
• נשערך בנפרד את $p_{x|y}(x|0)$ ואת $p_{x|y}(x|1)$.

נתחיל מ $p_{x|y}(x|0)$. בשביל לשערך פילוג זה נסתכל רק על הדגימות השייכות של $y = 0$:



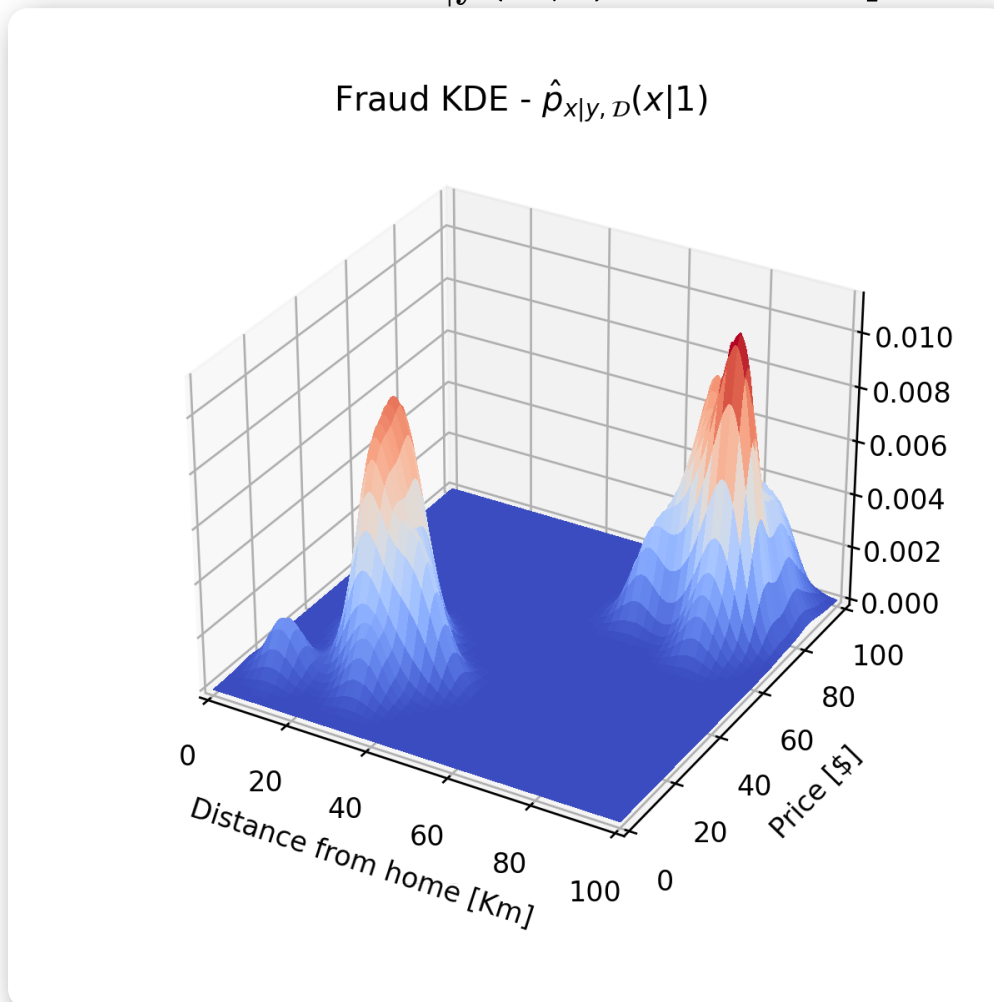
שיערוך של פילוגים מעורבים - דוגמא

נשתמש ב KDE על מנת לשערך את $p_{x|y}(x|0)$:



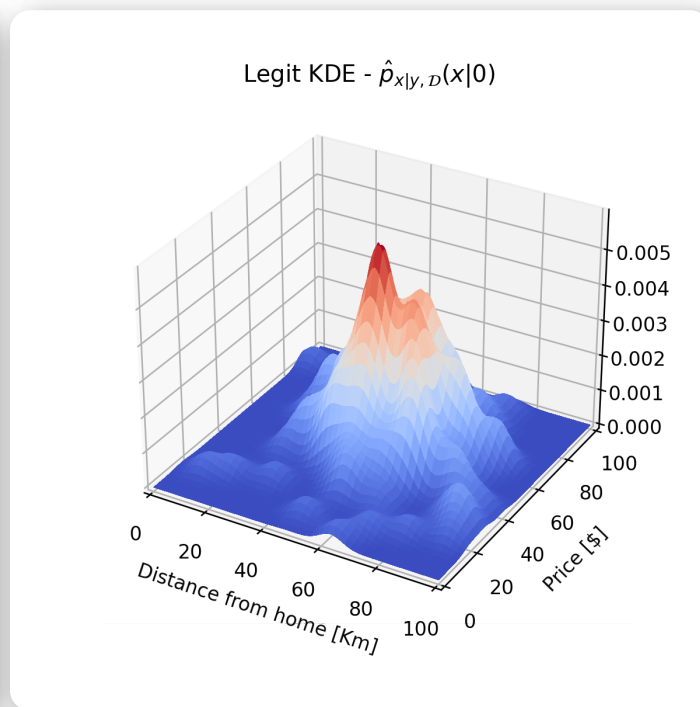
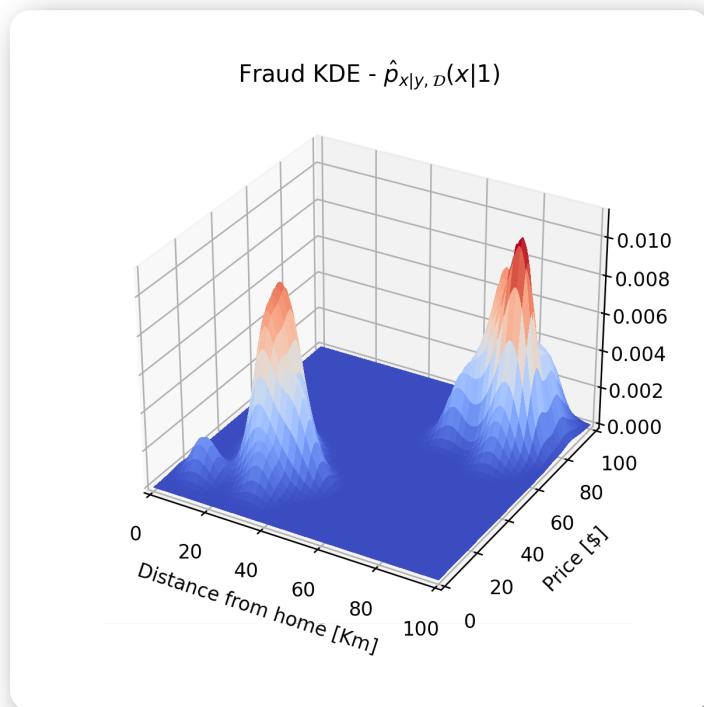
שיערוך של פילוגים מעורבים - דוגמא

באופן דומה נשערך גם את $p_{x|y}(x|1)$:



שיערוך של פילוגים מעורבים - דוגמא

$$\hat{p}_{y,D}(y) = \begin{cases} 0.8 & 0 \\ 0.2 & 1 \end{cases}$$



שלושת הפילוגים ששיערכנו מרכיבים את הפילוג המשותף על פי:

$$p_{\mathbf{x},y}(\mathbf{x}, y) = p_{\mathbf{x}|y}(\mathbf{x}|y)p_y(y)$$

שימוש בפילוג המשוערך לפתרון בעיות supervised learning

הגישה הגנרטיבית

מדגם



פילוג על סמך המדגם



חזאי אופטימאלי בהינתן הפילוג

עשינו את השלב הראשון, נעשה כעת את השלב השני.

חזאים אופטימאליים של פונקציות מחיר מוכרות - תזכורת

• **MSE**: התוחלת המותנית:

$$h^*(\mathbf{x}) = \mathbb{E}[y|\mathbf{x}]$$

• **MAE**: החציון של הפילוג המותנה:

$$h^*(\mathbf{x}) = y_{\text{median}} \quad \text{s.t.} \quad F_{y|\mathbf{x}}(y_{\text{median}}|\mathbf{x}) = 0.5$$

(כאשר $F_{y|\mathbf{x}}$ היא פונקציית הפילוג המצרפי של y בהינתן \mathbf{x}).

• **Misclassification rate**: הערך הכי סביר (ה mode):

$$h^*(\mathbf{x}) = \arg \max_y p_{y|\mathbf{x}}(y|\mathbf{x})$$

בעבור הפילוג שמצאנו נחפש את החזאי אשר ממזער את ה **misclassification rate**.

$$h(\mathbf{x}) = \arg \max_y p_{y|\mathbf{x}}(y|\mathbf{x})$$

במקרה הבנארי חזאי זה שווה ל:

$$h(\mathbf{x}) = \begin{cases} 1 & p_{y|\mathbf{x}}(1|\mathbf{x}) > p_{y|\mathbf{x}}(0|\mathbf{x}) \\ 0 & \text{else} \end{cases}$$

את $p_{y|\mathbf{x}}(y|\mathbf{x})$ נוכל לחשב מתוך הפילוג המשותף באופן הבא:

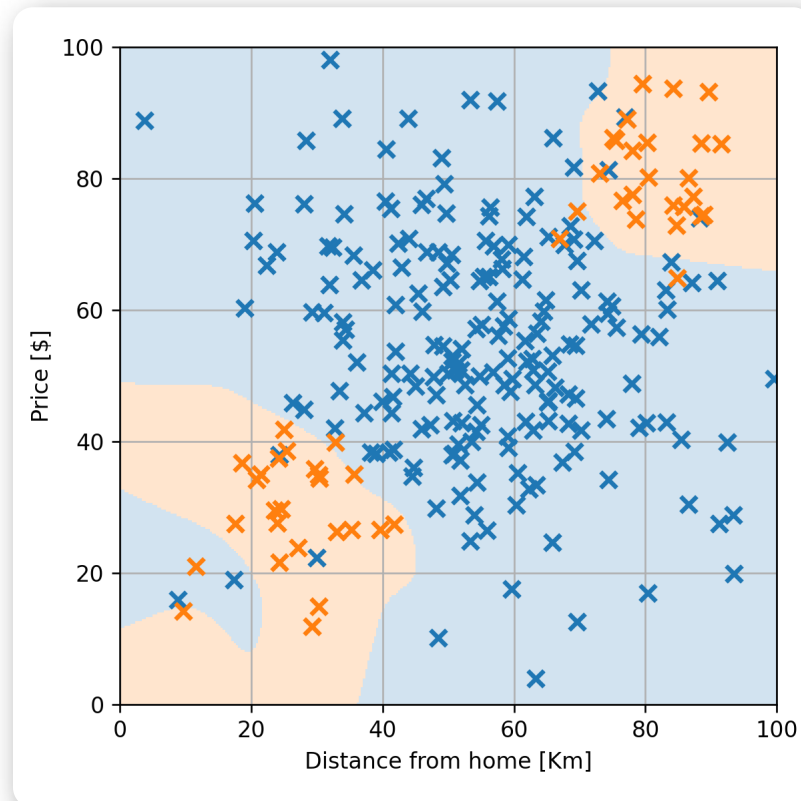
$$p_{y|\mathbf{x}}(y|\mathbf{x}) = \frac{p_{\mathbf{x},y}(\mathbf{x}, y)}{p_{\mathbf{x}}(\mathbf{x})} = \frac{p_{\mathbf{x}|y}(\mathbf{x}|y)p_y(y)}{p_{\mathbf{x}}(\mathbf{x})}$$

אם כן, בכדי לבדוק האם עסקה מסויימת הינה הונאה או לא, עלינו לבדוק האם:

$$\begin{aligned}
 & p_{y|x}(1|x) > p_{y|x}(0|x) \\
 \Leftrightarrow & \frac{p_{x|y}(x|1)p_y(1)}{p_x(x)} > \frac{p_{x|y}(x|0)p_y(0)}{p_x(x)} \\
 \Leftrightarrow & p_{x|y}(x|1)p_y(1) > p_{x|y}(x|0)p_y(0)
 \end{aligned}$$

$$p_{\mathbf{x}|y}(\mathbf{x}|1)p_y(1) > p_{\mathbf{x}|y}(\mathbf{x}|0)p_y(0)$$

נציב את פונקציות הפילוג ששיערכנו קודם לכן ונקבל את החזאי הבא:



ה misclassification rate של חזאי זה על ה test set הינו .0.12

ה bias וה variance של משערך

- המשערכים תלויים בצורה חזקה במדגם שאיתו אנו עובדים.
- נסתכל על האקראיות של השיערוך הנובעת מהאקראיות של המדגם.
- נשתמש בסימון $\mathbb{E}_{\mathcal{D}}$ בכדי לסמן תוחלת על פני הפילוג של המדגם.
- נגדיר bias ו variance של משערך

ה bias וה variance של משערך

Bias

בעבור שיערוך של גודל כל שהוא z בעזרת משערך \hat{z}_D , ה bias (היסט) של השיערוך מוגדר כ:

$$\text{Bias}(\hat{z}) = \mathbb{E}_D[\hat{z}_D] - z$$

כאשר ההטיה שווה ל-0, אנו אומרים שהמשערך אינו מוטה (Unbiased).

Variance

ה variance (שונות) של המשערך יהיה:

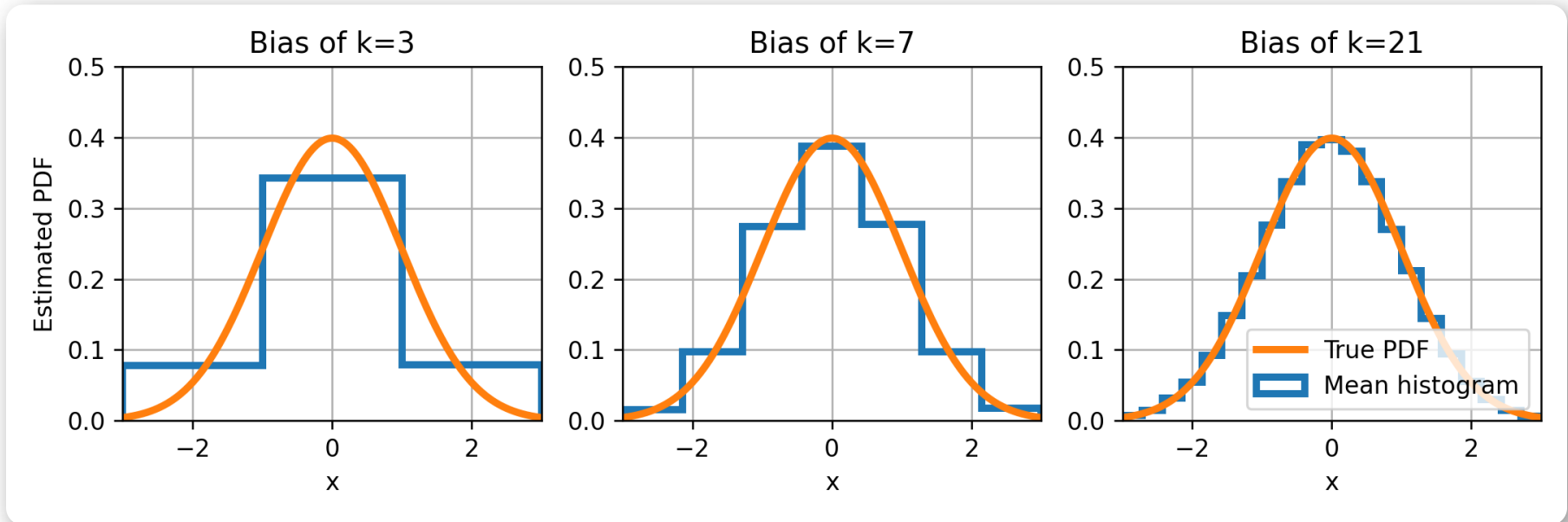
$$\text{Var}(\hat{z}) = \mathbb{E}_D \left[(\hat{z}_D - \mathbb{E}_D[\hat{z}_D])^2 \right] = \mathbb{E}_D[\hat{z}_D^2] - \mathbb{E}_D[\hat{z}_D]^2$$

מספר ה bins במונחים של bias ו variance

- ננסה לשערך את ה PDF של משתנה אקראי נורמאלי בעזרת היסטוגרמות בעלות 3, 7 ו 21 bins.

ה bias

נשרטט את ההיסטוגרמה הממוצעת לצד ה PDF האמיתי.



מספר ה bins במונחים של bias ו variance

מספר ה bins במונחים של bias ו variance

ה variance

- בכל שורה בגרף הקודם מגרילים שלושה מדגמים ומחשבים להם את ההיסטוגרמה.
 - אנו מצפים שבעבור מקרים שבהם ה variance נמוך השינויים יהיו קטנים ובעבור variance גבוה השינויים יהיו גדולים.
 - ה variance גדל ככל שאנו מגדילים את כמות ה bins.
- בדומה לחזאים בגישה הדיסקרימינטיבית, גם בהיסטוגרמה ישנו **bias-variance tradeoff**.