

מבוא ללמידה חישובית | סיכום הרצאה 13 (20942)

מנחה: ד"ר שי מימון
סמסטר: 2022'
נכתב על ידי: מתן כהן

שערוך - Estimation

1 מוטיבציה

במידה ויש לנו וקטור y של מדידות אשר קשורות לוקטור x כלשהו ונרצה להסיק מסקנות לגבי x דוגמה: נניח ש- x הוא וקטור של מיקומים ומהירויות של כלי טיס כלשהו ו- y הוא וקטור של מדידות של מכ"מים מכמה סנסורים.

היינו רוצים לאמוד את ערכי x על פי אותם נתונים בוקטור y .
ישנן 2 גישות להתמודדות עם שערוך:

- (1) הוקטור שמנסים לשערך הוא וקטור של פרמטרים דטרמיניסטיים אך הוא לא ידוע.
- (2) הוקטור שמנסים לשערך הוא וקטור אקראי.

2 Nonrandom Parameter Estimation

- בהרבה מצבים יותר נוח להניח שיש לנו פרמטרים דטרמיניסטיים אך לא ידועים מאשר פרמטרים אקראיים
 - לדוגמה - שערוך של תדר של גל סינוסי כלשהו שמעורבב ברעשים
- כאשר נרצה לשערך פרמטר או וקטור פרמטרים נסמן את השערוך ב- \hat{x} בתור פונקציה של y ונגדיר את השגיאה שלנו:

$$e(y) = \hat{x}(y) - x = \hat{x} - x$$

(הרבה פעמים נעלים את התלות ב- y מתוך נוחות)
יש לזכור נקודה קריטית מאוד:

אסור ליצור תלות בין \hat{x} ל- x !

כלומר - אסור שבפונקציית השערוך שלנו \hat{x} הפרמטרים של x יופיעו.

2.1 נרצה למדוד את טיב המשערוך שלנו בעזרת 2 מדדים

כיוון שידוע לנו ש x הוא לא אקראי אבל y כן (אלו המדידות שלנו) הרי שנוכל להשתמש במאפיינים סטטיסטיים כמו התוחלת של משתנה אקראי.

הטיה - bias

משערוך חסר הטיה מוגדרת בתור תוחלת השגיאה:

$$b_{\hat{x}}(x) = E[e(y)] = E[\hat{x}(y)] - x$$

משמע:

- אם תוחלת השגיאה היא אפס - נגיד שהמשערוך חסר הטיה
- אם המשערוך חסר הטיה אז התוחלת שלו שווה ערך לפרמטר אותו אנחנו מנסים לשערך!

Error-Covariance

נרצה, בנוסף להטיה נמוכה, שגם השונות המשותפת תהיה מאוד קטנה (מה שיצביע על שערך טוב). לשם כך נתבונן במטריצת השונות המשותפת שלנו על פני השגיאות (שונות בין כל 2 שגיאות) שמוגדרת:

$$\Lambda_e(x) = E \left[(e(y) - E(e(y))) \cdot (e(y) - E(e(y)))^T \right]$$

כמו כן, התוחלת של כפל השגיאות מוגדרת על ידי:

$$E(e(y) e^T(y)) = \Lambda_e(x) + b_{\hat{x}}(x) b_{\hat{x}}^T(x)$$

2.2 גישת הסבירות המרבית - Maximum Likelihood Estimation

- אחת מהגישות הרווחות ביותר לשערוך פרמטרים בהנחה שהפרמטרים דטרמיניסטיים לא ידועים.
- גישה זו ידוע בעיקר בשל התכונות האספימטוטיות שלה (כאשר מספר המדידות מאוד גדול) אשר משליכות על כך שככל שמספר הדגימות שלנו גדל המשערוך שלנו יהיה יותר קרוב לפרמטר - משמע ההטיה שלו תקטן ותשאף לאפס ביחד עם השונות.

- במילים אחרות ניתן להגיד שככל שישנן יותר מדידות המשערוך של הפרמטר ישאף לערך האמיתי של הפרמטר!
- במצב זה המשערוך ייקרא **עקבי**.

כיוון שאמרנו ש- y הוא וקטור של משתנים אקראיים אז יש לו פונקציית צפיפות פילוג אשר תלויה ב- x נגדיר את פונקציית ה-Likelihood:

$$L(x) = f_y(y; x)$$

הערה. כאשר פונקציה מוגדרת בעזרת $f_y(y; x)$ - ; $f_y(y; x)$ מבינים שמדובר בהנחה שהפרמטרים x הם דטרמיניסטיים לא ידועים. נזכור תמיד:

רוצים למצוא את סט הפרמטרים x אשר יסבירו את y בצורה הטובה ביותר

ולכן בפועל כאשר נקבל מצב בו ישנו y נתון ושני סטים של פרמטרים x_1 ו- x_2 המקיימים:

$$f_y(y; x_1) > f_y(y; x_2)$$

נרצה להשתמש ב- x_1 כיוון שהסבירות שלו הרבה יותר גבוהה!

באופן כללי, נרצה למקסם את L :

Maximum Likelihood Estimator

$$\hat{x}_{ML} = \underset{x}{\operatorname{argmax}} L(x) = \underset{x}{\operatorname{argmax}} f_y(y; x)$$

בהרבה מאוד מקרים נראה שיהיה לנו הרבה יותר קל למקסם פונקציה מונוטונית של הסבירות מאשר את פונקציית הסבירות עצמה, לדוגמה:

$$\log L(x) = \log f_y(y; x)$$

ולמקסם:

$$\hat{x}_{ML} = \underset{x}{\operatorname{argmax}} \log L(x) = \underset{x}{\operatorname{argmax}} \log f_y(y; x)$$

2.2.1 מציאת המקסימום

- במקרים בהם פונקציית הסבירות תהיה גזירה נוכל להשתמש בנגזרת על מנת למצוא את המקסימום בו אנו חושקים:

$$\frac{\partial f_{\mathbf{y}}(\mathbf{y}; \mathbf{x})}{\partial \mathbf{x}} = 0$$

$$\frac{\partial \log f_{\mathbf{y}}(\mathbf{y}; \mathbf{x})}{\partial \mathbf{x}} = 0$$

- לא תמיד יהיו לנו פתרונות סגורים או גלובליים ונוכל להשתמש באלגוריתמי גרדיאנט על מנת לחפש נקודות מקסימום
- גישה מוכרת שנתפרה עבור בעיית ה Maximum Likelihood היא גישת ה EM - Estimate Maximize

דוגמה. התפלגות ברנולי - הטלת מטבע

נתון סט של N הטלות מטבע בת"ס (iid) עם ההסתברות θ עבור עץ (1) ו- $1 - \theta$ עבור פאלי (0)

$$y_n = \begin{cases} 1 & \theta \\ 0 & 1 - \theta \end{cases}, \quad n = 1, 2, \dots, N, \quad iid$$

נשתמש בגישת הסבירות המירבית על מנת לשערך את פרמטר θ - ההסתברות להופעת "1".
בשלב הראשון נרשם את פונקציית הצפיפות בעזרת θ , תחילה עבור n יחיד כלשהו:

$$f_{y_n}(y_n; \theta) = \begin{cases} \theta & y_n = 1 \\ 1 - \theta & y_n = 0 \end{cases} = \theta^{y_n} \cdot (1 - \theta)^{1 - y_n}, \quad y_n \in \{0, 1\}$$

באופן כללי, כיוון שכלל המשתנים בת"ס נוכל להשתמש במכפלה ולקבל:

$$L(\theta) = f_{\mathbf{y}}(\mathbf{y}; \theta) = \prod_{n=1}^N f_{y_n}(y_n; \theta) = \prod_{n=1}^N \theta^{y_n} \cdot (1 - \theta)^{1 - y_n} = \theta^{\sum_{n=1}^N y_n} \cdot (1 - \theta)^{N - \sum_{n=1}^N y_n}$$

קיבלנו פונקציה די "מפחידה" - ברור שבגלל שמדובר בהמון מעריכים נוכל להשתמש בפונקציית הלוגריתם ולפשט את הדברים:

$$\log f_{\mathbf{y}}(\mathbf{y}; \theta) = \left(\sum_{n=1}^N y_n \right) \cdot \log \theta + \left(N - \sum_{n=1}^N y_n \right) \cdot \log(1 - \theta)$$

נעת נוכל למצוא מקסימום על ידי גזירה והשוואה לאפס ונקבל שמדובר בממוצע (חייב לגזור שוב כדי לוודא שזו נק' מקס):

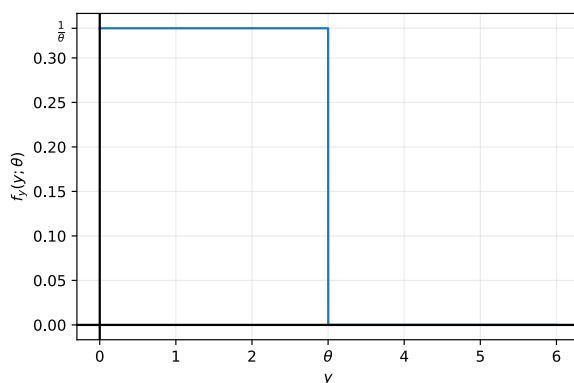
$$\begin{aligned} \frac{\partial \log f_{\mathbf{y}}(\mathbf{y}; \theta)}{\partial \theta} &= \left(\sum_{n=1}^N y_n \right) \cdot \frac{1}{\theta} + \left(N - \sum_{n=1}^N y_n \right) \cdot \frac{-1}{1 - \theta} = 0 \\ \Rightarrow \frac{1}{\theta(1 - \theta)} &\cdot \left[\left(\sum_{n=1}^N y_n \right) \cdot (1 - \theta) - \left(N - \sum_{n=1}^N y_n \right) \cdot \theta \right] = 0 \\ \Rightarrow \frac{1}{\theta(1 - \theta)} &\cdot \left[\left(\sum_{n=1}^N y_n \right) - \left(\theta \sum_{n=1}^N y_n \right) - N\theta + \left(\theta \sum_{n=1}^N y_n \right) \right] = 0 \\ \Rightarrow \frac{1}{\theta(1 - \theta)} &\cdot \left[\left(\sum_{n=1}^N y_n \right) - N\theta \right] = \frac{N}{\theta(1 - \theta)} \cdot \left[\frac{1}{N} \sum_{n=1}^N y_n - \theta \right] = 0 \\ \Rightarrow \hat{\theta}_{ML} &= \bar{y} = \frac{1}{N} \sum_{n=1}^N y_n \end{aligned}$$

משמע $\hat{\theta} \rightarrow \theta$, $N \rightarrow \infty$, בנוסף נראה שההטיה אפסית - $E(\hat{\theta}_{ML}) = \frac{1}{N} \sum_{n=1}^N E(y_n) = \theta$ והשונות קטנה ככל ש- N גדל.

דוגמה. התפלגות אחידה

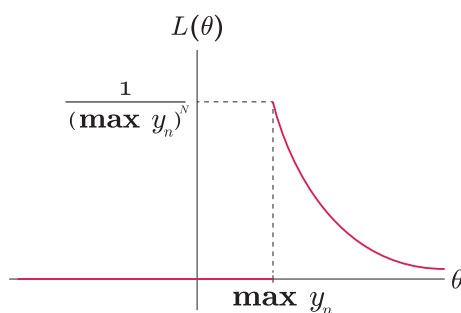
ניזכר בפונקציית צפיפות הפילוג של התפלגות האחידה עם פרמטר θ של מדידות בת"ס:

$$f_{y_n}(y_n; \theta) = \begin{cases} \frac{1}{\theta} & 0 \leq y_n \leq \theta \\ 0 & \text{else} \end{cases} \quad iid.$$



נחשב את פונקציית צפיפות הפילוג המשותפת כמו בדוגמה הקודמת:

$$L(\theta) = f_{\mathbf{y}}(\mathbf{y}; \theta) = \prod_{n=1}^N f_{y_n}(y_n; \theta) = \begin{cases} \frac{1}{\theta^N} & \max(y_n) \leq \theta \\ 0 & \text{else} \end{cases}$$



וכעת, קל להבחין על פי הגדרת הפונקציה ש:

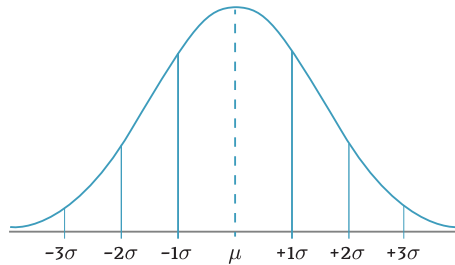
$$\hat{\theta}_{ML} = \max \{y_1, y_2, \dots, y_N\}$$

ועל פי הגדרה זו ומכך שראינו שהפונקציה מגדירה $\max(y_n) \leq \theta$ הרי שבהכרח אנחנו נתקרב אל θ **מלמטה** - משמע אנחנו תמיד נהיה עם הטיה מסוימת.

הנקודה חשובה כאן היא להבחין שככל ש- N יגדל ככה המקסימום שנקבל עבור $\hat{\theta}_{ML}$ יהיה קרוב יותר ל- θ עצמו!

דוגמה. התפלגות נורמלית (גאוסית)

$$f_{y_n}(y_n; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp \left\{ -\frac{1}{2\sigma^2} (y_n - \mu)^2 \right\} \quad iid$$



פונקצית צפיפות הפילוג המשותפת:

$$\begin{aligned} L(\theta) &= f_{\mathbf{y}}(y_n; \mu, \sigma^2) = \prod_{n=1}^N f_{y_n}(y_n; \mu, \sigma^2) \\ &= (2\pi\sigma^2)^{-N/2} \cdot \exp \left\{ -\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \mu)^2 \right\} \end{aligned}$$

מתבקש להשתמש ב \log :

$$\log f_{\mathbf{y}}(\mathbf{y}; \mu, \sigma^2) = \frac{N}{2} \cdot \log(2\pi) - \frac{N}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \mu)^2$$

נגזור לפי μ ונשווה לאפס (פונקציה ריבועית אז נקבל מקסימום גלובלי):

$$\begin{aligned} \frac{\partial \log f_{\mathbf{y}}(y_n; \mu, \sigma^2)}{\partial \mu} &= \frac{1}{\sigma^2} \sum_{n=1}^N (y_n - \mu) = 0 \\ \Rightarrow \hat{\mu}_{ML} &= \frac{1}{N} \sum_{n=1}^N y_n = \bar{y} \end{aligned}$$

באותה צורה נוכל **לגזור** על פי σ^2 כדי לקבל את התוצאה הרצויה:

$$\begin{aligned} \frac{\partial \log f_{\mathbf{y}}(y_n; \mu, \sigma^2)}{\partial \sigma^2} &= -\frac{N}{2} \cdot \frac{1}{\sigma^2} + \frac{1}{2(\sigma^2)^2} \cdot \sum_{n=1}^N (y_n - \mu)^2 = 0 \\ \Rightarrow \hat{\sigma}_{ML}^2 &= \frac{1}{N} \sum_{n=1}^N \left(y_n - \underbrace{\bar{y}}_{\hat{\mu}_{ML}} \right)^2 \end{aligned}$$

Maximum Likelihood and Least-Squares 2.2.2

נרצה לקשור בין הגישה לבין שיטת הריבועים הפחותים.

נניח שוקטור המדידות שלנו $\mathbf{y} \in \mathbb{R}^N$ מוגדר על ידי המדידות X , וקטור \mathbf{w} ווקטור אקראי $\mathbf{v} \in \mathbb{R}^N$ כך ש $v_n \sim N(0, \sigma^2)$ בת"ס:

$$\mathbf{y} = X\mathbf{w} + \mathbf{v}, \quad \mathbf{w} \in \mathbb{R}^{d+1}, \quad \mathbf{v} = \begin{bmatrix} v_1 \\ \vdots \\ v_N \end{bmatrix} \sim N(\mathbf{0}, \sigma^2 \cdot I)$$

נוכל לרשום את הפילוג של \mathbf{y} (כיוון שהזזה של משתנה אקראי פשוט מסיטה את הממוצע, נבחין כי השונות לא תשתנה אך הממוצע יזוז על פי אותו היסט - במקרה שלנו $X\mathbf{w}$):

$$(2.1) \quad \mathbf{y} \sim N(X \cdot \mathbf{w}, \sigma^2 \cdot I)$$

נכליל את המקרה הסקלרי של ההתפלגות הנורמלית: עבור משתנה אקראי $\mathbf{y} \sim N(\boldsymbol{\mu}, C)$ פונקציית צפיפות הפילוג שלו היא:

$$f_{\mathbf{y}}(\mathbf{y}; \boldsymbol{\mu}, C) = \frac{1}{\sqrt{\det(2\pi C)}} \cdot \exp \left\{ -\frac{1}{2} \cdot (\mathbf{y} - \boldsymbol{\mu})^T \cdot C^{-1} \cdot (\mathbf{y} - \boldsymbol{\mu}) \right\}$$

נעת נציב את 2.1:

$$\begin{aligned} f_{\mathbf{y}}(\mathbf{y}; X \cdot \mathbf{w}, \sigma^2 \cdot I) &= \frac{1}{\sqrt{\det(2\pi \cdot \sigma^2 I)}} \cdot \exp \left\{ -\frac{1}{2} \cdot (\mathbf{y} - X \cdot \mathbf{w})^T \cdot \underbrace{C^{-1}}_{\frac{1}{\sigma^2} I} \cdot (\mathbf{y} - X \cdot \mathbf{w}) \right\} \\ &= \frac{1}{\sqrt{(2\pi \cdot \sigma^2)^N \cdot \underbrace{\det(I)}_1}} \cdot \exp \left\{ -\frac{1}{2\sigma^2} \cdot \|\mathbf{y} - X \cdot \mathbf{w}\|^2 \right\} \end{aligned}$$

נשתמש ב- \log על מנת לפשט את החישובים שלנו:

$$\log f_{\mathbf{y}}(\mathbf{y}; X \cdot \mathbf{w}, \sigma^2 \cdot I) = -\frac{N}{2} \cdot \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|\mathbf{y} - X \cdot \mathbf{w}\|^2$$

ונקבל:

$$\hat{\mathbf{w}}_{ML} = \underset{\mathbf{w}}{\operatorname{argmax}} \log L(\mathbf{w})$$

כיוון שהחלק הראשון של הביטוי לא תלוי כלל ב \mathbf{w} ונוכל להתעלם מקבועים - נוכל פשוט למצוא ערך \mathbf{w} אשר יביא למינימום (החלפת הסימן) את הביטוי ונקבל בדיוק את גישת ה-Least Squares:

$$\hat{\mathbf{w}}_{ML} = \underset{\mathbf{w}}{\operatorname{argmin}} \|\mathbf{y} - X \cdot \mathbf{w}\|^2 = \hat{\mathbf{w}}_{LS}$$

Maximum Likelihood and Logistic Regression 2.2.3

נזכר שרגרסיה לוגיסטית קיבלנו פונקציה הממודלת על ידי $\theta(s) = \frac{1}{1 + e^{-s}}$

$$P(y_n = 1 | \mathbf{x}_n) = \theta(\mathbf{w}^T \mathbf{x}_n)$$

$$P(y_n = -1 | \mathbf{x}_n) = 1 - \theta(\mathbf{w}^T \mathbf{x}_n) = \theta(-\mathbf{w}^T \mathbf{x}_n)$$

נעת נשתמש ב-Maximum Likelihood על מנת למדל את אותה בעיה.
תחילה ננסה לרשום את הפונקציה בצורה פשוטה:

$$P(y_n | \mathbf{x}_n) = \theta(y_n \cdot \mathbf{w}^T \mathbf{x}_n), \quad y_n \in \{-1, 1\}$$

ונוכל לרשום את המכפלה עבור N המדידות בת"ס:

$$\prod_{n=1}^N P(y_n | \mathbf{x}_n) = \prod_{n=1}^N \theta(y_n \cdot \mathbf{w}^T \mathbf{x}_n)$$

נעבור שוב ל-log ונשתמש בהגדרת θ :

$$\begin{aligned} \log L &= \prod_{n=1}^N P(y_n | \mathbf{x}_n) = \sum_{n=1}^N \log \theta(y_n \cdot \mathbf{w}^T \mathbf{x}_n) \\ &= - \sum_{n=1}^N \log(1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n}) \end{aligned}$$

ועל פי הצבה פשוטה נבחין כי זו פונקצית המחיר שפיתחנו עבור cross entropy כאשר דיברנו על Logistic Regression:

$$\hat{\mathbf{w}}_{ML} = \underset{\mathbf{w}}{\operatorname{argmax}} \log L(\mathbf{w}) = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{n=1}^N \log(1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n})$$

זוהי פונקציה קונבקסית של \mathbf{w} ולכן ההתכנסות שלנו עם אלגוריתם איטרטיבי תביא אותנו על מינימום גלובלי.

3 שערורך פרמטר אקראי

בגישה זו אנחנו מתייחסים ל \mathbf{x} (הפרמטר שנרצה לשערך) בתור משתנה אקראי כלשהו שיש לנו מידע מקדים עליו (נניח אי שלילי) - משמע ישנו פילוג $p_{\mathbf{x}}(\mathbf{x})$ אשר מתאר את הפילוג הא-פריורי (A-priori) של \mathbf{x} עוד לפני שצפינו בתצפיות! בצורה זו, במידה ובידינו הפילוג $p_{\mathbf{x}}(\mathbf{x})$ והפילוג $p_{\mathbf{y}|\mathbf{x}}(\mathbf{y} | \mathbf{x})$ נוכל לקבל את הפילוג המשותף (הסטטיסטיקה המלאה):

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

וכמובן שנוכל להשתמש בנוסחאת בייס ולקבל את הפילוג הא-פוסטריורי (A-posteriori) אשר משמעו "איך הפילוג של \mathbf{x} מושפע לאחר שצפינו במדידות \mathbf{y} :"

Bayes' theorem

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

לגישה זו קוראים Bayesian Estimation

3.1 שערורך בייסיאני - Bayesian Estimation

בשיטה זו נרצה להגדיר פונקצית מחיר בין פרמטר \mathbf{x} לבין המשערך שלנו שתלוי ב- \mathbf{y} : $\mathbf{f}(\mathbf{y})$ ועל מנת למצוא את המשערך הטוב ביותר נרצה לפתור בעיה אופטימיזציה מהצורה הבאה:

$$\hat{\mathbf{x}}(\cdot) = \underset{\mathbf{f}(\cdot)}{\operatorname{argmin}} E(C(\mathbf{x}, \mathbf{f}(\mathbf{y})))$$

נתבונן בתוחלת שבבעיה האופטימיזציה:

$$\underbrace{\iint C(\mathbf{x}, \mathbf{f}(\mathbf{y})) \cdot p_{\mathbf{x}\mathbf{y}}(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y}}_{E(C(\mathbf{x}, \mathbf{f}(\mathbf{y})))} \stackrel{\text{Bayes}}{=} \int \left[\underbrace{\int C(\mathbf{x}, \mathbf{f}(\mathbf{y})) p_{\mathbf{x}|\mathbf{y}}(\mathbf{x} | \mathbf{y}) d\mathbf{x}}_{(1)} \right] p_{\mathbf{y}}(\mathbf{y}) d\mathbf{y}$$

ונוכל להתבונן ב- (1) ולרשום את $\mathbf{f}(\mathbf{y})$ בתור \mathbf{a} ולקבל:

$$(3.1) \quad \hat{\mathbf{x}}(\mathbf{y}) = \underset{\mathbf{a}}{\operatorname{argmin}} \int C(\mathbf{x}, \mathbf{a}) p_{\mathbf{x}|\mathbf{y}}(\mathbf{x} | \mathbf{y}) d\mathbf{x}$$

נשתמש שוב בנוסחאת בייס על הביטוי: $p_{\mathbf{x}|\mathbf{y}}(\mathbf{x} | \mathbf{y})$, נשמיט את המכנה $p_{\mathbf{y}}$ (כיוון שלא קשור לבעיה) ונקבל את בעיה האופטימיזציה:

$$\hat{\mathbf{x}}(\mathbf{y}) = \underset{\mathbf{a}}{\operatorname{argmin}} \int C(\mathbf{x}, \mathbf{a}) p_{\mathbf{y}|\mathbf{x}}(\mathbf{y} | \mathbf{x}) p_{\mathbf{x}}(\mathbf{x}) d\mathbf{x}$$

בדוגמה הבאה נדבר על מקרה בו המשתנה האקראי הוא סקלרי, קל לעבור למצב וקטורי בצורה של component wise.

דוגמה. Minimum Uniform Cost (MUC)

נניח שפונקציית המחיר שלנו היא:

$$C(a, \hat{a}) = \begin{cases} 1 & |a - \hat{a}| > \varepsilon \\ 0 & \text{otherwise} \end{cases}$$

נציב את C בבעיית האופטימיזציה שהגדרנו ב-3.1 ונשתמש בהסתברות המשלימה:

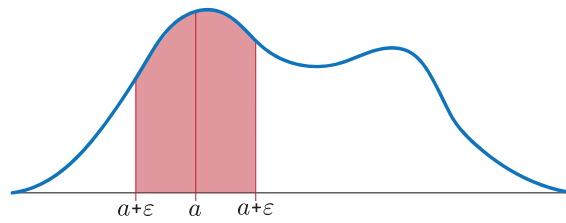
$$\hat{x}_{MUC}(\mathbf{y}) = \underset{a}{\operatorname{argmin}} \left[\int_{\{\mathbf{x}|C=1\}=|\mathbf{x}-a|>\varepsilon} p_{\mathbf{x}|\mathbf{y}}(\mathbf{x} | \mathbf{y}) d\mathbf{x} \right] = \underset{a}{\operatorname{argmin}} \left[1 - \int_{a-\varepsilon}^{a+\varepsilon} p_{\mathbf{x}|\mathbf{y}}(\mathbf{x} | \mathbf{y}) d\mathbf{x} \right]$$

כיוון שאנחנו רוצים להביא למינימום את הביטוי, אנחנו בעצם רוצים להביא למקסימום את הביטוי:

$$\underset{a}{\operatorname{argmax}} \int_{a-\varepsilon}^{a+\varepsilon} p_{\mathbf{x}|\mathbf{y}}(\mathbf{x} | \mathbf{y}) d\mathbf{x}$$

בעצם מה שעשינו כאן זה למצוא נקודת מקסימום בתחום מסוים של פונקציית ההתפלגות שלנו על מנת להחזיר לנו את

ההסתברות הגבוהה ביותר:



אם נשאיף את ε לאפס נמצא את הנקודה המקסימלית של ה a-posteriori distribution ונקבל המשעור הא-פוסטריורי

המקסימלי:

Maximum A posteriori Estimation (MAP)

$$\hat{x}_{MAP}(\mathbf{y}) = \underset{a}{\operatorname{argmax}} p_{\mathbf{x}|\mathbf{y}}(a | \mathbf{y}) = \lim_{\varepsilon \rightarrow 0} \hat{x}_{MUC}(\mathbf{y})$$

ואם נתבונן בהבדלים בין Maximum Likelihood לבין Maximum A-posteriori Estimation נראה שהעיקר הוא המידע

המקדים על x :

$$\hat{x}_{ML}(y) = \underset{a}{\operatorname{argmax}} p(y; a) = \underset{a}{\operatorname{argmax}} \log p(y; a)$$

$$\hat{x}_{MAP}(y) \stackrel{\text{Bayes}}{=} \underset{a}{\operatorname{argmax}} \frac{p_{y|x}(y | a) \cdot p_x(a)}{p_y(y)} \stackrel{\log \text{ is monotonic}}{=} \underset{a}{\operatorname{argmax}} \log p_{y|x}(y | a) + \log p_x(a)$$

דוגמה. Minimum Absolute Error (MAE)

ניתן להשתמש בפונקציית המחיר של הערך המוחלט:

$$C(a, \hat{a}) = |a - \hat{a}|$$

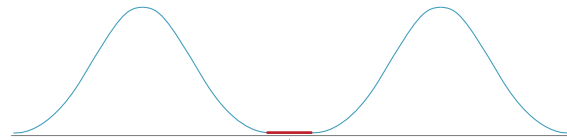
נקבל בעזרת הצבה ב- 3.1 את:

$$\hat{\mathbf{x}}(\mathbf{y}) = \operatorname{argmin}_a \int_{-\infty}^{\infty} |x - a| p_{\mathbf{x}|\mathbf{y}}(\mathbf{x} | \mathbf{y}) d\mathbf{x}$$

ובעצם את החציון של פונקציית ההתפלגות - הערך שבו ההסתברות של ההתפלגות האפוסטריורית היא בדיוק חצי:

$$\int_{-\infty}^{\hat{x}_{MAE}(\mathbf{y})} p_{\mathbf{x}|\mathbf{y}}(\mathbf{x} | \mathbf{y}) d\mathbf{x} = \int_{\hat{x}_{MAE}(\mathbf{y})}^{\infty} p_{\mathbf{x}|\mathbf{y}}(\mathbf{x} | \mathbf{y}) d\mathbf{x} = \frac{1}{2}$$

הערה. יכולים להיות כמה משערכים עם אותו ערך (החציון הוא לא בהכרח משערך יחיד) - חשוב לזכור!



הרבה נקודות שיקיימו את המשוואה

דוגמה. Minimum Mean Square Error (MMSE)

המשערך הפופולרי ביותר, קובע את פונקציית המחיר הריבועית המוכרת (במקרה הווקטורי):

$$C(\mathbf{a}, \hat{\mathbf{a}}) = \|\mathbf{a} - \hat{\mathbf{a}}\|^2 = (\mathbf{a} - \hat{\mathbf{a}})^T (\mathbf{a} - \hat{\mathbf{a}}) = \sum_{i=1}^N (a_i - \hat{a}_i)^2$$

נציב ב- 3.1 ונקבל:

$$\hat{\mathbf{x}}_{MMSE}(\mathbf{y}) = \operatorname{argmin}_a \int_{-\infty}^{\infty} (\mathbf{a} - \mathbf{x})^T (\mathbf{a} - \mathbf{x}) p_{\mathbf{x}|\mathbf{y}}(\mathbf{x} | \mathbf{y}) d\mathbf{x}$$

במקרה הסקלרי:

$$\hat{x}_{MMSE}(\mathbf{y}) = \operatorname{argmin}_a \int_{-\infty}^{\infty} (a - x)^2 p_{x|\mathbf{y}}(x | \mathbf{y}) dx$$

נמצא נקודת מינימום על ידי גזירה על פי a ונקבל את משערך התוחלת המותנית:

$$\hat{x}_{MMSE}(\mathbf{y}) = \int_{-\infty}^{\infty} x p(x | \mathbf{y}) dx = E(x | \mathbf{y})$$

משערך התוחלת המותנית

Minimum Mean Square Error (MMSE)

$$\hat{\mathbf{x}}_{MMSE}(\mathbf{y}) = E(\mathbf{x} | \mathbf{y})$$

תכונות ה-MMSE

הטיה

המשערור חסר הטיה:

$$b_{MMSE} = E(\mathbf{e}(\mathbf{x}, \mathbf{y})) = E(\hat{\mathbf{x}}_{MMSE}(\mathbf{y}) - \mathbf{x}) = E_{\mathbf{y}}(E_{\mathbf{x}|\mathbf{y}}(\mathbf{x} | \mathbf{y})) - E(\mathbf{x}) = 0$$

covariance - שונות משותפת

$$\begin{aligned}\Lambda_{MMSE} &= E(\mathbf{e}\mathbf{e}^T) = E[(\mathbf{x} - E(\mathbf{x} | \mathbf{y}))(\mathbf{x} - E(\mathbf{x} | \mathbf{y}))^T] \\ &= E_{\mathbf{y}}\{E_{\mathbf{x}|\mathbf{y}}[(\mathbf{x} - E(\mathbf{x} | \mathbf{y}))(\mathbf{x} - E(\mathbf{x} | \mathbf{y}))^T | \mathbf{y}]\} \\ &= E(\Lambda_{\mathbf{x}|\mathbf{y}}(\mathbf{y}))\end{aligned}$$

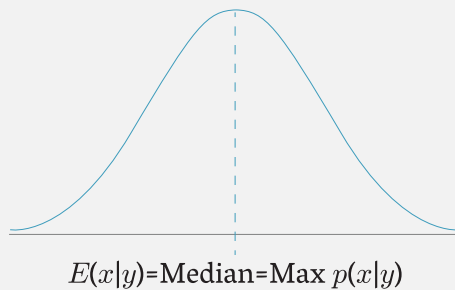
שגיאת השערור אורתוגונלית לכל פונקציה (לינארית או לא לינארית) של המדידות שלנו

$$E[(\hat{\mathbf{x}}(\mathbf{y}) - \mathbf{x})\mathbf{g}^T(\mathbf{y})] = 0$$

סיכום הדוגמאות:

- MAP מחפש את נקודת המקסימום של הפילוג
- MAE מחפש את נקודת (נקודות) החציון של הפילוג
- MMSE מחפש את נקודת התוחלת של הפילוג

במידה ו- \mathbf{x} ו- \mathbf{y} מפולגים גאוסית במשותף - שלושת המשערים יתכנסו לאותה נקודה



4 דוגמה לשערך

נניח שיש בידנו סט מדידות y_n המוגדר על ידי המודל:

$$y_n = w_0 + \mathbf{x}_n^T \cdot \mathbf{w} + V_n, \quad n = 1, 2, \dots, N$$

כאשר:

- w_0 - פרמטר
- \mathbf{x}_n - מדידות ידועות
- \mathbf{w} - סט משתנים אקראיים בת"ס אחד בשני וב- V_n המתפלגים: $w_j \sim N(0, \tau^2)$, $j = 1, 2, \dots, d$
- V_n - רעשים המתפלגים: $V_n \sim N(0, \sigma^2)$, iid

MAP

נכתוב את פונקציית צפיפות הפילוג של \mathbf{w} , תוך שימוש בנתון שהמשתנים הם בת"ס:

$$f_{\mathbf{w}}(\mathbf{w}) = \prod_{j=1}^d \frac{1}{\sqrt{2\pi\tau^2}} \cdot \exp \left\{ -\frac{1}{2\tau^2} \cdot (w_j)^2 \right\}$$

נשתמש בלוגריתם על מנת לפשט את הדברים:

$$\begin{aligned} \log f_{\mathbf{w}}(\mathbf{w}) &= -\frac{d}{2} \log(2\pi\tau^2) - \frac{1}{2\tau^2} \sum_{j=1}^d w_j^2 \\ &= -\frac{d}{2} \log(2\pi\tau^2) - \frac{1}{2\tau^2} \|\mathbf{w}\|^2 \end{aligned}$$

כעת נבחין כי הביטוי $w_0 + \mathbf{x}_n^T \cdot \mathbf{w}$ משמש לנו מעין קבוע, ולכן נוכל להבחין כי ההתפלגות של y_n היא נורמלית כהתפלגות V_n אשר מוסטת בקבוע $w_0 + \mathbf{x}_n^T \cdot \mathbf{w}$:

$$y_n | \mathbf{w} \sim N(w_0 + \mathbf{x}_n^T \cdot \mathbf{w}, \sigma^2)$$

ונוכל למצוא את פונקציית ההתפלגות:

$$f_{\mathbf{y}|\mathbf{w}} = \prod_{n=1}^N f_{y_n|\mathbf{w}}(y_n | \mathbf{w})$$

כל אחד מהביטויים $f_{y_n|\mathbf{w}}(y_n | \mathbf{w})$ מוגדר:

$$f_{y_n|\mathbf{w}}(y_n | \mathbf{w}) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp \left\{ -\frac{1}{2\sigma^2} \cdot (y_n - w_0 - \mathbf{x}_n^T \cdot \mathbf{w})^2 \right\}$$

נשתמש שוב ב- \log :

$$\log f_{y_n|\mathbf{w}}(y_n | \mathbf{w}) = -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \cdot (y_n - w_0 - \mathbf{x}_n^T \cdot \mathbf{w})^2$$

ולכן:

$$\log f_{\mathbf{y}|\mathbf{w}} = -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - w_0 - \mathbf{x}_n^T \cdot \mathbf{w})^2$$

מצאנו את שני המרכיבים של MAP וכעת צריך למקסם:

$$\hat{\mathbf{w}}_{MAP} = \underset{\mathbf{w}}{\operatorname{argmax}} \log f_{\mathbf{y}|\mathbf{w}}(\mathbf{y} | \mathbf{w}) + \log f_{\mathbf{w}}(\mathbf{w})$$

קל לראות ש:

• ב- $\log f_{\mathbf{y}|\mathbf{w}}(\mathbf{y} | \mathbf{w})$ הביטוי התלוי ב- \mathbf{w} הוא:

$$-\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - w_0 - \mathbf{x}_n^T \cdot \mathbf{w})^2$$

○ ביטוי זה מוכפל ב- $-\frac{1}{2}$, לכן נוכל להשמיט אותו

• ב- $\log f_{\mathbf{w}}(\mathbf{w})$ הביטוי התלוי ב- \mathbf{w} הוא:

$$-\frac{1}{2\tau^2} \|\mathbf{w}\|^2$$

○ ביטוי זה גם מוכפל ב- $-\frac{1}{2}$ ונוכל להשמיט את המכפלה גם כאן

נקבל את בעיית האופטימיזציה הפשוטה יותר (לאחר מכפלה ב- σ^2 על מנת לקבל ביטוי נוח יותר):

$$\hat{\mathbf{w}}_{MAP} = \underset{\mathbf{w}}{\operatorname{argmin}} \left\{ \sum_{n=1}^N (y_n - w_0 - \mathbf{x}_n^T \cdot \mathbf{w})^2 + \frac{\sigma^2}{\tau^2} \cdot \|\mathbf{w}\|^2 \right\}$$

ונבחין כי זה ביטוי זהה לביטוי שקיבלנו ב-Ridge Regression (אם נסמן $\lambda = \frac{\sigma^2}{\tau^2}$), רק שכאן λ נקבע על פי שונות הרעש V_n

והשונות של סט המשתנים האקראיים \mathbf{w} .

- השונות σ^2 ו- τ^2 יהוו מדד לעד כמה אנו מאמינים במידע המקדים (\mathbf{w}) לעומת המדידות (y_n)
- אם τ^2 מאוד קטן - נסיק שיש לנו וודאות מאוד גדולה שהערכים מתפלגים סביב 0 (פעמון מאוד צר)
- משמע נקבל פרמטר רגולריזציה λ מאוד גדול - מה שיגרור אילוף של פרמטרים מאוד קטנים
- אם σ^2 מאוד קטנה (ביחס ל τ^2 כמובן) - מדובר במדידות שהן כמעט ונטולות ברעש
- משמע נקבל פרמטר רגולריזציה מאוד קטן ונסתמך על ערכי \mathbf{w} בשביל לאפסם את הבעיה שלנו

הערה. במידה ונשנה את הפלגות w_j כך ש:

$$w_j \sim \text{Laplace}(0, \tau^2), \quad j = 1, 2, \dots, d$$

עם אותן הנחות של חוסר תלות - נקבל את Lasso Regression.