

הסבר - Max Likelihood

תמונע קלינמן

כתבתי את הסיכום הזה לעצמי, ייתכן וכנראה שיש בו טעויות. מוזמנים לכתוב לי - timna.kleinman@mail.huji.ac.il

חלקים בכתום מסמנים קטע שהוא לא קריטי בשביל להבין, ולא קשור ישירות ל-Max Likelihood אבל לדעתי יעזרו לגבש תמונה כללית

1 טרמינולוגיה

נסמן -

• x_1, \dots, x_n הן הדגימות.

• θ היא המשתנה אותו אנו מנסים לאמוד.

נזכיר את חוק בייס -

$$\underbrace{P(\theta|x_1, \dots, x_n)}_{\text{posterior}} = \frac{\underbrace{P(x_1, \dots, x_n|\theta)}_{\text{likelihood}} \cdot \underbrace{P(\theta)}_{\text{prior}}}{\underbrace{P(x_1, \dots, x_n)}_{\text{(constant normalization)}}}$$

נסביר את ההיגיון בכל אחד מהשמות -

Prior - זה מידע מקדים שיש לנו על ההסתברות שהפרמטר הזה הוא האמיתי. לדוג' - יש לנו מטבע שקיבלנו מהסופר, ואנחנו מנסים לאמוד את ההסתברות שלו ליפול על עץ. כלומר $\hat{\theta}$ יהיה אומד להסתברות לקבל עץ. אזי, מאחר וזה סתם מטבע שקיבלנו מהסופר, ובלי קשר לדגימות שלנו (זריקות מטבע) בהסתברות גבוהה $\theta \simeq 0.5$, ובהסתברות ממש ממש נמוכה $\theta = 0$ ¹. במילים אחרות $P(\theta = 0.5)$ הוא גדול, ו- $P(\theta = 0)$ הוא קטן מאוד. זהו ה-prior.

Likelihood - עבור θ מסויימת מה הסיכוי (What is the likelihood) שהיינו מקבלים את הדגימות x_1, \dots, x_n ?

Posterior - אחרי (post) שראינו את הדגימות, מה הסיכוי שזו אכן ה- θ ?

¹כלומר לא נקבל עץ אף פעם

הערה 1.1. שימו לב להבדל בין השאלה של likelihood לבין posterior: likelihood מסתכל על θ ספציפית שכבר נבחרה, ואז שואל - האם בכלל סביר שהיינו מקבלים את הדגימות האלו? posterior אומר שנתונות לי כבר הדגימות, מה הסיכוי ש- θ היא הנכונה? מאחר ו- θ היא הפרמטר שמעניין אותנו, נוכל לכתוב

$$\underbrace{P(\theta|x_1, \dots, x_n)}_{\text{posterior}} \propto \underbrace{P(x_1, \dots, x_n|\theta)}_{\text{likelihood}} \cdot \underbrace{P(\theta)}_{\text{prior}}$$

כלומר, ה-posterior פרופורציונלי ל-prior, רק שמצטרפת אליו גם הסבירות מלכתחילה ש- θ היא הנכונה. נחזור לדוגמה שלנו עם זריקות המטבע - נניח וזרקנו מטבע 3 פעמים, ויצא לנו 3 פעמים פלי. אז, ה-likelihood יגיד לנו שהכי סביר שייצא לנו 3 פעמים פלי כאשר $\theta = 0$, כלומר כאשר הסיכוי לקבל עץ הוא 0. לכן, אם היינו בוחרים לפי ה-likelihood היינו בוחרים $\theta = 0$. ה-posterior לעומת זאת, אומנם כן "יבין" ש- $\theta = 0$ נשמע הכי סביר בהינתן הדגימות, אבל גם ישקלל לתוך זה את העובדה שממש לא סביר ש- $\theta = 0$ (ה-prior). לכן, יכול להיות שה-posterior ב- θ קטנה אומנם, אבל לא 0 כי זה ממש ממש לא סביר.

2 אמידה

מטרה: לאמוד את θ בהינתן x_1, \dots, x_n דגימות.

דרך א' (MAP): אינטואיטיבית, בהמשך לחלק של הטרמינולוגיה, מרגיש הכי הגיוני לבחור את $\hat{\theta}$ ע"י בחירת θ שהיא בעלת הסתברות posterior הכי גבוהה. כלומר

$$\begin{aligned} \hat{\theta}_{\text{MAP}} &= \arg \max_{\theta} (P(\theta|x_1, \dots, x_n)) \stackrel{\text{Bayes}}{=} \arg \max_{\theta} \left(\frac{P(x_1, \dots, x_n|\theta) \cdot P(\theta)}{P(x_1, \dots, x_n)} \right) \\ &\stackrel{\substack{P(x_1, \dots, x_n) \\ \theta \text{ on depend Doesn't}}}{=} \arg \max_{\theta} (P(x_1, \dots, x_n|\theta) \cdot P(\theta)) \stackrel{\text{בתי"ל}}{=} \arg \max_{\theta} \left(\prod_{i=1}^n P(x_i|\theta) \cdot P(\theta) \right) \\ &\stackrel{\text{מונוטוניית log}}{=} \arg \max_{\theta} \left(\log \left(\prod_{i=1}^n P(x_i|\theta) \cdot P(\theta) \right) \right) \stackrel{\text{תכונות log}}{=} \arg \max_{\theta} \left(\sum_{i=1}^n \log(P(x_i|\theta)) + \log(P(\theta)) \right) \end{aligned}$$

עם זאת, ישנן שתי סיבות עיקריות (שאני מכירה) לכך שנעדיף להשתמש ב-likelihood -

- ה-prior הרבה פעמים לא ידוע לנו (וכך גם ה-posterior), ולכן לא נוכל כלל לחשב את הביטוי הרצוי.
- אם נסתכל על החלק האחרון בפיתוח, נשים לב כי אנו סוכמים על הדגימות. לכן, אם ניקח "ממש הרבה" דגימות, ה-prior יהפוך לזניח ונוכל להשתמש ב-likelihood בלבד.

הערה 2.1. כל הנושא הזה של כן או לא prior מתקשר לגישות שונות - bayesian vs. frequentist שאליה לא ניכנס כאן.

דרך ב' (Max Likelihood): כאמור למעלה, ה-likelihood הוא $P(x_1, \dots, x_n|\theta)$. **חשוב לשים לב!** כי ה-likelihood הוא פונקציה של θ - הדגימות נתונות לנו, ואנחנו שואלים "עבור θ ספציפית, כמה סביר שהיינו מקבלים את הדגימות הנתונות?". לכן נסמן

$$\mathcal{L}(\theta|x_1, \dots, x_n) = P(x_1, \dots, x_n|\theta) \left(\underbrace{= f_{\theta}(x_1, \dots, x_n) = f(x_1, \dots, x_n|\theta)}_{\text{במקרה הרצוי}} \right)$$

כעת, נרצה לבחור את θ בעלת ה-likelihood הגבוה ביותר -

$$\begin{aligned}\hat{\theta}_{ML} &= \arg \max_{\theta} (P(x_1, \dots, x_n | \theta)) \stackrel{\text{i.i.d}}{=} \arg \max_{\theta} \left(\prod_{i=1}^n P(x_i | \theta) \right) \stackrel{\text{מונטונייט}}{=} \arg \max_{\theta} \left(\log \left(\prod_{i=1}^n P(x_i | \theta) \right) \right) \\ &\stackrel{\text{מונטונייט}}{=} \arg \max_{\theta} \left(\sum_{i=1}^n \log P(x_i | \theta) \right)\end{aligned}$$

הערה 2.2. לא היינו חייבים להוסיף את ה- \log , אך פעמים רבות הוא עוזר ומקל מאוד על החישובים.

לבסוף, נזכיר כי על מנת למצוא את המקסימום של הביטוי נרצה לגזור אותו ולהשוות ל-0 (ואם אנחנו פדנטיים אז גם לוודא שהוא מקסימום).

דוגמה 2.3. (מהשעור) נניח שמשפחת ההתפלגויות ידועה לנו והיא ההתפלגות הנורמלית. אנחנו מחפשים את μ . כלומר $x_1, \dots, x_n \stackrel{\text{i.i.d}}{\sim} \mathcal{N}(\mu, \sigma^2)$ עבור איזושהי σ^2 . נזכיר כי ההתפלגות הנורמלית היא מהצורה $-\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right)$. עתה, נשתמש ב- $\hat{\theta}_{ML}$ שחישבנו קודם (עבור $\theta := \mu$) -

$$\begin{aligned}\hat{\mu}_{ML} &= \arg \max_{\theta} \left(\sum_{i=1}^n \log f_{\theta}(x_i) \right) = \arg \max_{\theta} \left(\sum_{i=1}^n \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{(x_i - \mu)^2}{2\sigma^2} \right) \right) \right) \\ &\stackrel{\frac{1}{\sqrt{2\pi\sigma^2}} = \text{constant}}{=} \arg \max_{\theta} \left(\sum_{i=1}^n \log \exp \left(-\frac{(x_i - \mu)^2}{2\sigma^2} \right) \right) = \arg \max_{\theta} \left(-\sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2} \right) \\ &\stackrel{\frac{1}{2\sigma^2} = \text{constant}}{=} \arg \max_{\theta} \left(-\sum_{i=1}^n (x_i - \mu)^2 \right) = \arg \min_{\theta} \left(\sum_{i=1}^n (x_i - \mu)^2 \right)\end{aligned}$$

עכשיו, על מנת למצוא את המינימום של הפונקציה ביחס ל- μ נוכל לגזור אותה (ביחס ל- μ) ונקבל

$$\begin{aligned}\frac{\partial \left(\sum_{i=1}^n (x_i - \mu)^2 \right)}{\partial \mu} &= \sum_{i=1}^n \frac{\partial (x_i - \mu)^2}{\partial \mu} = \sum_{i=1}^n 2(x_i - \mu) \stackrel{!}{=} 0 \\ \sum_{i=1}^n (x_i - \mu) &= \sum_{i=1}^n x_i - \sum_{i=1}^n \mu = \sum_{i=1}^n x_i - n\mu \stackrel{!}{=} 0 \\ \hat{\mu}_{ML} &= \frac{1}{n} \sum_{i=1}^n x_i\end{aligned}$$