
מבוא למערכות לומדות

הוקלד בידי אייל צווכר

מבוא למערכות לומדות 67577 תשפ"ב סמסטר ב'

תוכן העניינים

2	1	כלים מתמטיים
2	1	הרצאה - תורת האמידה
13	2	תרגול - קצת לינארית

שבוע 1

כלים מתמטיים

1 הרצאה - תורת האמידה

מטרות הקורס

- היכרות עם בעיות למידה ואלגוריתמי למידה (חלק תיאורטי)
- הבנה של העקרונות מאחורי האלגוריתמים (חלק מתמטי)
- יכולת לממש אלגוריתמים של בעיות למידה (חלק תכנותי)

שאלה

נניח שאנחנו רוצים לדעת מה השעה, אבל אין לנו שום דרך לדעת. אז נשאל מישור מה השעה. נניח שהוא ענה 15 : 13 והשעה היא באמת 13 : 13.

- איך נדע מה השעה האמיתית?

- איך נדע כמה כל תשובה מדויקת?

ניתן לשאול הרבה אנשים ולחשב ממוצע של התשובות שלהם. התשובות שקיבלנו מהאנשים נקראות תצפיות (observation). נסמן ב- x_1, \dots, x_m הזמן המשוער שלנו הוא

$$\bar{x} = \frac{1}{m} \sum_i x_i$$

סט התצפיות שלנו נקרא המדגם (sample). נחשוב על x_1, \dots, x_n כערכים שהתקבלו (realized) במשתנים מקריים X_1, \dots, X_n .

הגדרה

נאמר כי x_1, \dots, x_m הם מתפלגים באופן זהה (identically distributed) אם כולם ערכים שנתקבלו מערכים משתנים המקריים המתפלגים עם אותה התפלגות, כלומר, $X_1, \dots, X_m \sim \mathcal{P}$, ו- $X_i = x_i$ לכל i .

הגדרה

נאמר שדגימות x_1, \dots, x_m הן בעלות אותה התפלגות ובלתי תלויים (independently identically distributed or i.i.d) אם הם מתפלגים באופן זהה ובלתי תלויים. במקרה זה נסמן

$$X_1, \dots, X_m \text{ i.i.d. } \mathcal{P}, \quad \forall i \ X_i = x_i$$

מהי \mathcal{P} ?

נניח כי \mathcal{P} היא התפלגות פרמטרית כלשהי שמוגדרת על ידי קבוצת פרמטרים $\theta \in \Theta$.

- $\vec{\theta}$ היא וקטור של פרמטרים.

- Θ היא קבוצת כל הערכים האפשריים לפרמטר.

- עבור $\mathcal{P} := \text{Pois}(\lambda)$ אזי $\vec{\theta} := \{\lambda\}$ ו- $\mathbb{R}_+ := \Theta$

- עבור $\mathcal{P} := N(\mu, \sigma^2)$ אזי $\vec{\theta} := \{\mu, \sigma^2\}$ ו- $\mathbb{R} \times \mathbb{R}_+ := \Theta$

בהינתן דגימות x_1, \dots, x_m שנדגמו i.i.d. לפי $\mathcal{P}(\vec{\theta})$ כאשר אנחנו לא יודעים את $\vec{\theta}$, נרצה לדעת מהו $\vec{\theta}^*$ שנותן את ההתאמה הכי טובה ל- $\vec{\theta}$ האמיתי.

(למשל, נניח ש- \mathcal{P} היא פואסונית ונחפש את הפרמטר λ הכי טוב).

פורמליזציה

• נגדיר פונקציית החלטה/כלל החלטה $\delta : \mathbb{R}^m \rightarrow \Theta$.

- עבור $\mathcal{P} := \text{Pois}(\lambda)$,

$$(x_1, \dots, x_m) \xrightarrow{\delta} \lambda$$

- עבור $\mathcal{P} := N(\mu, \sigma^2)$,

$$(x_1, \dots, x_m) \xrightarrow{\delta} (\mu, \sigma^2)$$

• את קבוצת כל פונקציות ההחלטה נסמן ב- $\Delta := \{\delta : \mathbb{R}^m \rightarrow \Theta\}$. היא מכונה מרחב ההשערות.

הגדרה

תהי $\delta \in \Delta$ פונקציית החלטה. אזי $\delta(X_1, \dots, X_m)$ נקראת אומד (estimator) של הפרמטר $\vec{\theta}$, או לפעמים רק אומד. $\vec{\theta}$ המתקבל נקרא אומדן.

אומדים של התפלגות גאוסיאנית

נניח שיש לנו דגימות $x_1, \dots, x_m \stackrel{\text{i.i.d.}}{\sim} \mathcal{P}$ עבור $\mathcal{P} := N(\mu, \sigma^2)$ ואנחנו רוצים להעריך את μ, σ^2 . איך נוכל לעשות את זה? איזה אומד נבחר עבור μ, σ^2 ?

• האומד ממוצע המדגם (sample mean) מוגדר על ידי

$$\hat{\mu}_X := \frac{1}{m} \sum x_i$$

הערה: הסימון $\hat{\mu}$ הוא האומד לפרמטר μ .

• האומד שונות המדגם (sample variance) מוגדר על ידי

$$\hat{\sigma}_X^2 := \frac{1}{m-1} \sum (x_i - \hat{\mu}_X)^2$$

תכונות של אומדים

נתבונן למשל בשונות המדגם. יכולנו לבחור בהרבה אפשרויות אחרות לפונקציות:

$$\hat{\sigma}_1 := \frac{1}{m-1} \sum (x_i - \hat{\mu})^2$$

$$\hat{\sigma}_2 := \frac{1}{m} \sum (x_i - \hat{\mu})^2$$

$$\hat{\sigma}_3 := \frac{1}{m} \sum |x_i - \hat{\mu}|$$

מה מהם "הכי טוב"? מה זה בכלל אומר "להיות הכי טוב"?

הערה

X_1, \dots, X_m הם משתנים מקריים, ולכן, כיוון ש- δ היא פונקציה של X_1, \dots, X_m אזי δ עצמו הוא משתנה מקרי. לכן, אנחנו יכולים לשאול מהי התוחלת של δ , מה השונות שלו, וכו'.

רעיון

כאשר נגדיר אומדן, תכונה שנרצה שתהיה לו היא שהוא יהיה בלתי מוטה (unbiased), כלומר, בממוצע הערך שהוא מחזיר שווה לערך האמיתי.

הגדרה

יהי δ אומדן עבור פרמטר $\vec{\theta}$. ההפרש $d := \delta(X_1, \dots, X_n) - \theta$ נקרא השגיאה של δ .

הגדרה

יהי δ אומדן עבור פרמטר $\vec{\theta}$. הגודל

$$\begin{aligned} \text{Bias}_{\vec{\theta}}[\delta(X_1, \dots, X_m)] &:= \mathbb{E}_{X_1, \dots, X_m | \vec{\theta}}[d] \\ &= \mathbb{E}_{X_1, \dots, X_m | \vec{\theta}}[\delta(X_1, \dots, X_m) - \vec{\theta}] \end{aligned}$$

נקראת ההטייה (bias או systemic error) של δ .
(התוחלת נלקחת על הרבה דגימות של X_1, \dots, X_m מההתפלגות עם פרמטר $\vec{\theta}$ מסויים)

הגדרה

יהי δ אומדן עבור פרמטר $\vec{\theta}$.
נאמר כי δ הוא בלתי מוטה אם לכל $\vec{\theta} \in \Theta$ מתקיים

$$\text{Bias}_{\vec{\theta}}[\delta(X_1, \dots, X_m)] = 0$$

טענה

האומדן ממוצע המדגם $\hat{\mu}$ הוא בלתי מוטה עבור ההתפלגות $\mathcal{P} := N(\mu, \sigma^2)$

הוכחה

לכן $\hat{\mu}$ הוא אומדן בלתי מוטה.

טענה

האומדן שונות המדגם $\hat{\sigma}^2$ הוא בלתי מוטה עבור ההתפלגות $\mathcal{P} := N(\mu, \sigma^2)$

$$\begin{aligned}
\mathbb{E}[\hat{\sigma}^2] &= \frac{1}{m-1} \sum_i \mathbb{E}[(X_i - \hat{\mu})^2] = \frac{1}{m-1} \mathbb{E} \left[\left(X_i - \frac{\sum_j X_j}{m} \right)^2 \right] \\
&= \frac{1}{m-1} \mathbb{E} \left[X_i^2 - 2X_i \cdot \underbrace{\left(\frac{1}{m} \sum_j X_j \right)}_{\hat{\mu}} + \underbrace{\frac{1}{m^2} \sum_{j,k} X_j X_k}_{\hat{\mu}^2} \right] = \\
&= \frac{1}{m-1} \mathbb{E} \left[\sum_i [X_i^2] - \frac{2}{m} \sum_{i,j} \mathbb{E}[X_i X_j] + \frac{1}{m} \sum_{j,k} \mathbb{E}[X_j X_k] \right] = \\
&= \frac{1}{m-1} \left(\left(1 - \frac{1}{m} \right) \sum_i \mathbb{E}[X_i^2] - \frac{1}{m} \sum_{i \neq j} \mathbb{E}[X_i X_j] \right) =
\end{aligned}$$

מכיוון שהמדידות הן i.i.d. אזי נגדיר $X \sim N(\mu, \sigma^2)$ ויתקיים:

- לכל i , $\mathbb{E}[X_i] = \mathbb{E}[X]$ ו- $\mathbb{E}[X_i^2] = \mathbb{E}[X^2]$
- לכל $i \neq j$, $\mathbb{E}[X_i X_j] = \mathbb{E}[X_i] \mathbb{E}[X_j] = (\mathbb{E}[X])^2$ (כי המ"מ ב"ת)

לכן

$$\begin{aligned}
&= \frac{1}{m-1} \left(\frac{m-1}{m} \cdot m \mathbb{E}[X^2] - \frac{m(m-1)}{m} (\mathbb{E}[X])^2 \right) = \\
&= \frac{1}{m-1} \left((m-1) \mathbb{E}[X^2] - (m-1) (\mathbb{E}[X])^2 \right) = \\
&= \mathbb{E}[X^2] - (\mathbb{E}[X])^2 = \text{Var}(X) = \sigma^2
\end{aligned}$$

כנדרש.

הגדרה

יהי δ אומד עבור פרמטר $\vec{\theta}$. השונות של δ היא

$$\text{Var}(\delta) := \mathbb{E}_{X_1, \dots, X_m | \vec{\theta}} \left[\left(\delta(X_1, \dots, X_m) - \mathbb{E}_{X_1, \dots, X_m | \vec{\theta}}[\delta(X_1, \dots, X_m)] \right)^2 \right]$$

דוגמא

תהי \mathcal{P} התפלגות עם שונות של σ^2 .
יהיו $\mathcal{P} \stackrel{\text{i.i.d.}}{\sim} X_1, \dots, X_m$. מהו $\text{Var}(\hat{\mu})$?

$$\begin{aligned}
\text{Var}(\hat{\mu}) &= \text{Var} \left(\frac{1}{m} \sum X_i \right) = \frac{1}{m^2} \text{Var} \left(\sum X_i \right) = \\
&= \frac{1}{m^2} \sum_{i=1}^m \text{Var}(X_i) \stackrel{\uparrow \text{i.i.d.}}{=} \frac{1}{m^2} \cdot m \cdot \sigma^2 = \frac{\sigma^2}{m}
\end{aligned}$$

משמעות

השונות של האומד נותן לנו מדד לכמה טובים הביצועים של האומד שלנו. אם השונות גדולה, האומד לא מאוד טוב.

דוגמא

נוכל לומר שהאומד הכי טוב הוא $\delta \in \Delta$ כך ש- δ בלתי מוטה וגם $\text{Var}(\delta)$ מינימלי.

נראות מקסימלית

הגדרה

יהי $X \sim \mathcal{P}(\vec{\theta})$ ו- f פונקציית הצפיפות של \mathcal{P} . פונקציית הנראות היא

$$\mathcal{L}(\vec{\theta} | x) := f_{\vec{\theta}}(x)$$

כאשר x הוא realization של X .

דוגמא

נתבונן בהתפלגות גאוסיאנית עם $N(\mu, \sigma^2) \stackrel{\text{i.i.d.}}{\sim} x_1, \dots, x_m$ ונסמן $\vec{\theta} = \{\mu, \sigma^2\}$. פונקציית הנראות היא

$$\mathcal{L}(\vec{\theta} | x_i) = f_{\vec{\theta}}(x_i) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

לכל $i \in [m]$. מכיוון ש- x_1, \dots, x_m הן i.i.d. אזי

$$\begin{aligned} \mathcal{L}(\vec{\theta} | x_1, \dots, x_m) &= f_{\vec{\theta}}(x_1, \dots, x_m) = \\ &= \prod_{i=1}^m f_{\vec{\theta}}(x_i) = \\ &= \frac{1}{\sqrt{2\pi\sigma^2}^m} \exp\left(-\frac{1}{2\sigma^2} \sum (x_i - \mu)^2\right) \end{aligned}$$

הערה

נשתמש ב- \mathcal{L} על מנת להשוות בין דברים באופן יחסי! אם יהיו לנו הרבה דגימות, \mathcal{L} ישאף ל-0 בכל מקרה. אבל נוכל לשאול מי יותר סביר מבין $\vec{\theta}_1$ או $\vec{\theta}_2$ על ידי כך שנשווה בין $\mathcal{L}(\vec{\theta}_1 | x_1, \dots, x_m)$ ו- $\mathcal{L}(\vec{\theta}_2 | x_1, \dots, x_m)$.

הגדרה

תהי \mathcal{L} פונקציית הנראות עבור התפלגות \mathcal{P} כלשהי, המתואר על ידי $\vec{\theta} \in \Theta$. יהי $X \sim \mathcal{P}(\vec{\theta})$ משתנה מקרי ויהי x ערך שנדגם ממנו. אומד הנראות המרכזית (Maximum Likelihood Estimator or MLE) עבור $\vec{\theta}$ הוא

$$\hat{\theta}^{\text{MLE}} := \underset{\vec{\theta} \in \Theta}{\operatorname{argmax}} \mathcal{L}(\vec{\theta} | x)$$

דוגמא

נחשב את ה-MLE של הממוצע של ההתפלגות הגאוסיאנית $\hat{\mu}^{\text{MLE}}$, כאשר נתון לנו σ^2 .
יהיו $x_1, \dots, x_m \stackrel{\text{i.i.d.}}{\sim} N(\mu, \sigma^2)$ נרצה למצוא את

$$\hat{\mu}^{\text{MLE}} = \operatorname{argmax}_{\mu \in \mathbb{R}} \mathcal{L}(\mu | x_1, \dots, x_m, \sigma^2)$$

מתקיים

$$\begin{aligned} \hat{\mu}^{\text{MLE}} &= \operatorname{argmax}_{\mu \in \mathbb{R}} \mathcal{L}(\mu | x_1, \dots, x_m, \sigma^2) = \\ &\stackrel{\substack{\uparrow \\ \text{i.i.d.}}}{=} \operatorname{argmax}_{\mu \in \mathbb{R}} \prod_i \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right) = \\ &\stackrel{\substack{\uparrow \\ (*)}}{=} \operatorname{argmax}_{\mu \in \mathbb{R}} \prod_i \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right) = \\ &= \operatorname{argmax}_{\mu \in \mathbb{R}} \exp\left(-\frac{1}{2\sigma^2} \sum_i (x_i - \mu)^2\right) = \\ &\stackrel{\substack{\uparrow \\ \text{מונטוניית עולה}}}{=} \operatorname{argmax}_{\mu \in \mathbb{R}} \log\left(\exp\left(-\frac{1}{2\sigma^2} \sum_i (x_i - \mu)^2\right)\right) = \\ &= \operatorname{argmax}\left(-\frac{1}{2\sigma^2} \sum_i (x_i - \mu)^2\right) = \\ &= \operatorname{argmax}\left(-\sum_i (x_i - \mu)^2\right) \end{aligned}$$

(*) נבחין כי אנחנו מחפשים את מי שממקסם את הפונקציה ולא את הערך המקסימלי של הפונקציה! הקורס $\frac{1}{\sqrt{2\pi\sigma^2}}$ אולי משפיע על ערך הפונקציה אך הוא אינו משפיע על הערך שממקסם אותה

על מנת למקסם את הפונקציה, נגזור ונשווה ל-0:

$$\frac{\partial}{\partial \mu} \left(-\sum_i (x_i - \mu)^2 \right) = -\sum_{i=1}^m \frac{\partial (x_i - \mu)^2}{\partial \mu} = \sum_{i=1}^m 2(x_i - \mu) = 0$$

ולכן

$$\boxed{\hat{\mu}^{\text{MLE}} = \frac{1}{m} \sum_{i=1}^m x_i}$$

מסקנה

אומד ממוצע המדגם (sample mean), שראינו שהוא אומד בלתי מוטה, הוא גם האומד maximum likelihood!

הערה ממוצע המדגם ממזער את ריבועי המרחקים ממנו.

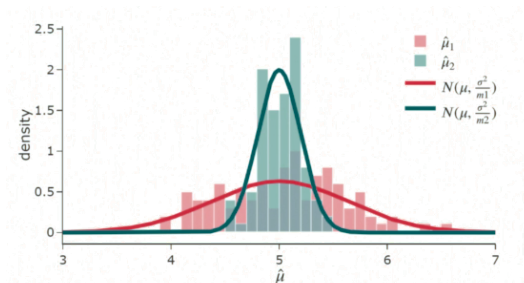
כמה טוב השיערוך?

מכיוון שהאומד הוא משתנה מקרי, יש לו התפלגות. למשל, במקרה של ה-sample mean, אנחנו יודעים ש- $x_1, \dots, x_m \stackrel{\text{i.i.d.}}{\sim} N(\mu, \sigma^2)$ לכן,

$$\mathbb{E}[\hat{\mu}] = \mu, \quad \text{Var}(\hat{\mu}) = \frac{\sigma^2}{m}$$

ניתן גם להראות (אבל לא נראה) שמתקיים $\hat{\mu} \sim N\left(\mu, \frac{\sigma^2}{m}\right)$ בפרט:

- האומד נותן התפלגות גאוסיאנית סביב הערך האמיתי μ (זה הגיוני כי הוא unbiased)
- השונות פרופורציונלית לשונות של המידע המקורי, והיא דועכת באופן לינארי ביחס למספר הדגימות.



הגרף מתאר מצב בו הגרלנו הרבה פעמים m_1 ו- m_2 ערכים מהתפלגות גאוסיאנית עם μ, σ^2 , וראינו אילו ערכי $\hat{\mu}_1, \hat{\mu}_2$ התקבלו לנו לפי ממוצע המדגם. לכן נוכל לשאול שאלות כגון:

- מה הסיכוי שנשערך ערך מסוים?
- כמה אנחנו בטוחים בתוצאה שהתקבלה לנו?
- מה הסיכוי לסטייה מהערך האמיתי, וכיצד זה תלוי במספר הדגימות?

דוגמא

נניח שאנחנו רוצים להעריך את ההטייה של מטבע לא הוגן כלשהו. כלומר, נתון מטבע שיש לו סיכוי של p לתת Heads ואנחנו רוצים לדעת מהו p . נחשוב על הטלת מטבע בתור משתנה מקרי ברנולי

$$\mathcal{D}_p(X) := \begin{cases} p & X = 1 \\ 1-p & X = 0 \end{cases}, \quad p \in [0, 1]$$

בהינתן דגימה של m הטלות $S := \{x_1, \dots, x_m\}$, נסמן את התפלגות של m הדגימות על ידי \mathcal{D}_p^m ונסמן את ההסתברות לקבל את הדגימות S ב- $\mathcal{D}_p^m(S)$.

בהינתן מדגם כזה, נרצה להפיק אלגוריתם למידה \mathcal{A} שיעריך את p ונשאל כמה האלגוריתם שלנו מדויק. האלגוריתם מקבל את S כקלט, שנבחר i.i.d. לפי \mathcal{D} ומוציא כפלט את ההערכה ל- p . נסמן אותה ב- $\hat{p}(S)$ או \hat{p} או בקיצור נמרץ \hat{p} .

אלגוריתם הלמידה בו נשתמש יהיה פשוט השכיחות האמפירית (כלומר, כמה פעמים מופיע Heads במדגם), שזה בדיוק ממוצע המדגם.

$$\hat{p}(S) = \frac{1}{m} \sum_i x_i$$

אנחנו יודעים שהאומד הזה הוא בלתי מוטה. לכן, בתוחלת נקבל $\mathbb{E}_S(\hat{p}) = p$ בעיות:

- המדגם שלנו סופי, ולכן \hat{p} שלנו לא יכול מדויק. לכן נסתפק ב- \mathcal{A} שיהיה מדויק מספיק עבורנו, כלומר, $|\hat{p} - p| \leq \varepsilon$ עבור $\varepsilon \in (0, 1)$.
- אלא אם כן $p \in \{0, 1\}$, תמיד יש סיכוי כלשהו לקבל S שלא מייצג את \mathcal{D} . לכן, יהיה עלינו לשאול מה הסיכוי ש- \mathcal{A} יחזיר תוצאה טובה מספיק, כלומר, נחשב $\Pr(|\hat{p} - p| \leq \varepsilon)$.

תזכורת (א"ש מרקוב)

עבור מ"מ X אי שלילי עם תוחלת סופית מתקיים:

$$\Pr(X \geq a) \leq \frac{\mathbb{E}(X)}{a}$$

במקרה שלנו, $|\hat{p} - p|$ הוא מ"מ אי שלילי. לכל רמת דיוק $\varepsilon \in (0, 1)$ נוכל לחסום את הסיכוי לסטות מ- p ביותר מ- ε על ידי

$$\Pr[|\hat{p} - p| \geq \varepsilon] \leq \frac{\mathbb{E}(|\hat{p} - p|)}{\varepsilon}$$

נותר לנו לחסום מלמעלה את $\mathbb{E}(|\hat{p} - p|)$. ניזכר כי $\text{Var}(A) = \mathbb{E}(A) - (\mathbb{E}(A))^2$ ולכן

$$\text{Var}(|\hat{p} - p|) = \mathbb{E}(|\hat{p} - p|) - \mathbb{E}^2(|\hat{p} - p|)$$

ולכן, מהיות $\text{Var}(|\hat{p} - p|) \geq 0$,

$$\mathbb{E}^2(|\hat{p} - p|) \leq \mathbb{E}(|\hat{p} - p|^2)$$

לכן

$$\begin{aligned} \mathbb{E}^2(|\hat{p} - p|) &\leq \mathbb{E}(|\hat{p} - p|^2) = \mathbb{E}((\hat{p} - p)^2) = \\ &= \mathbb{E}((\hat{p} - \mathbb{E}(\hat{p}))^2) = \text{Var}(\hat{p}) = \end{aligned}$$

(*)
(*)

(*) כיוון ש- \hat{p} הוא אומד בלתי מוטה, $\mathbb{E}(\hat{p}) = p$. נזכור כי \hat{p} הוא ממוצע של m הטלות מטבע, ולכן

$$= \text{Var}\left(\frac{1}{m} \sum X_i\right) = \frac{1}{m^2} \text{Var}\left(\sum_{\text{i.i.d.}} x_i\right) = \frac{p(1-p)}{m} \stackrel{(*)}{\leq} \frac{1}{4m}$$

(*) שכן הפונקציה $f(p) := p(1-p)$ מקבלת את הערך המקסימלי שלה בקטע $[0, 1]$ עבור $p = \frac{1}{2}$ ו- $f(p) = \frac{1}{4}$.

לכן

$$\mathbb{E}(|\hat{p} - p|) \leq \frac{1}{\sqrt{4m}}$$

אם כן

$$\Pr[|\hat{p} - p| \geq \varepsilon] \leq \frac{\mathbb{E}(|\hat{p} - p|)}{\varepsilon} \leq \frac{1}{\sqrt{4m\varepsilon^2}}$$

כלומר

$$\Pr[|\hat{p} - p| \leq \varepsilon] \geq 1 - \frac{1}{\sqrt{4m\varepsilon^2}}$$

בתור שלב אחרות, נסמן $\delta = \frac{1}{\sqrt{4m\varepsilon^2}}$. נוכל לבדוד את m ונקבל $m = \frac{1}{4\varepsilon^2\delta^2}$. לכן, לכל $\varepsilon, \delta \in (0, 1)$, אם נקבל $m \geq \frac{1}{4\varepsilon^2\delta^2}$ דגימות, אזי

$$\Pr(|\hat{p} - p| \leq \varepsilon) \geq 1 - \delta$$

אפשר לשפר את התוצאה הזו! למשל, אם נשתמש באי-שוויון צ'בישב.

תזכורת (א"ש צ'בישב)

יהי X מ"מ בעל תוחלת ושונות. אזי לכל $\varepsilon > 0$ מתקיים

$$\Pr[|X - \mathbb{E}[X]| \geq \varepsilon] \leq \frac{\text{Var}(X)}{\varepsilon^2}$$

לכן, בהינתן $x_1, \dots, x_m \stackrel{\text{i.i.d.}}{\sim} \text{Ber}(p)$ כמו שהגדרנו נקבל

$$\begin{aligned} \Pr[|\hat{p} - p| \geq \varepsilon] &= \Pr[|\hat{p} - \mathbb{E}(\hat{p})| \geq \varepsilon] \leq \frac{\text{Var}(\hat{p})}{\varepsilon^2} = \\ &= \frac{1}{\varepsilon^2} \text{Var}\left(\frac{1}{m} \sum_{i=1}^m x_i\right) \stackrel{\text{i.i.d.}}{=} \frac{1}{m^2 \varepsilon^2} \sum \text{Var}(x_i) = \\ &= \frac{p(1-p)}{m\varepsilon^2} \leq \frac{1}{4m\varepsilon^2} \end{aligned}$$

החסם יתנהג כמו $\frac{1}{\sqrt{m}}$, וזה יותר טוב מהחסם של מרקוב, שנותן $\frac{1}{\sqrt{m}}$.

לכן, נגדיר $\delta = \frac{1}{4m\varepsilon^2}$ ונקבל ש- $m = \frac{1}{4\delta\varepsilon^2}$. אם כן, הסקנו כי לכל $\varepsilon, \delta \in (0, 1)$, אם ניתן $m > \frac{1}{4\delta\varepsilon^2}$ דגימות אזי

$$\Pr(|\hat{p} - p| \leq \varepsilon) \geq 1 - \delta$$

התפלגויות רבות משתנים

עד עכשיו התמודדנו עם התפלגויות של משתנה אחד.

אבל מה אם יש כמה דגימות של משתנים שתלויים זה בזה?

למשל, כיצד נמדל גובה ומשקל של אנשים?

אפשר להחליט למשל ש- $w_1, \dots, w_m \stackrel{\text{i.i.d.}}{\sim} N(75, 3)$ ו- $h_1, \dots, h_m \stackrel{\text{i.i.d.}}{\sim} N(170, 5)$.

אבל הגובה והמשקל תלויים זה בזה!

לכן נשתמש בוקטורים מקריים.

הגדרה

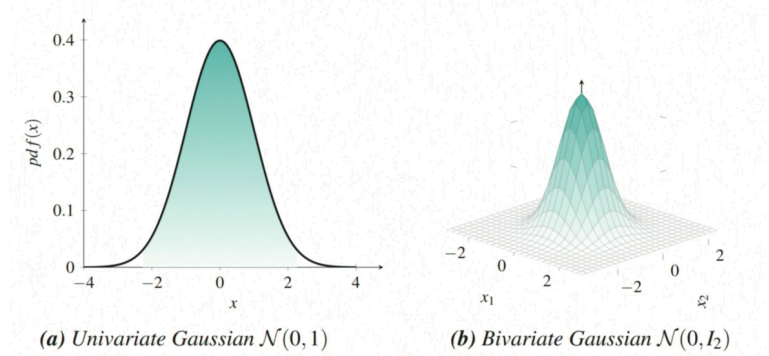
יהיו X_1, \dots, X_d קבוצה של m מוגדרים מעל אותו מרחב הסתברות.

אזי $X := \begin{pmatrix} X_1 \\ \vdots \\ X_d \end{pmatrix}$ נקרא וקטור מקרי.

וקטור מקרי הוא פונקציה ממרחב המדגם ל- \mathbb{R}^d .

הגדרה

יהי $X := \begin{pmatrix} X_1 \\ \vdots \\ X_d \end{pmatrix}$ וקטור מקרי. ההתפלגות המשותפת של X_1, \dots, X_d הוא פונקציה ההתפלגות על כל הערכים האפשריים של X_1, \dots, X_d .



נשים לב כי קווי הגובה הם אליפסות יפות.

הגדרה

יהי $X := \begin{pmatrix} X_1 \\ \vdots \\ X_d \end{pmatrix}$ וקטור מקרי. מטריצת השונות המשותפת (Σ covariance matrix) היא מטריצה $d \times d$ שהערך ה- (i, j) שלה הוא השונות המשותפת:

$$\Sigma_{ij} := \text{Cov}(X_i, X_j) = \mathbb{E}[(X_i - \mathbb{E}(X_i))(X_j - \mathbb{E}(X_j))]$$

הערה המטריצה היא ריבועית וסימטרית, והיא גם מטריצת psd.

הגדרה

יהי $X := \begin{pmatrix} X_1 \\ \vdots \\ X_d \end{pmatrix}$ וקטור מקרי עם התפלגות נורמלית של כמה משתנים עם תוחלת $\mu \in \mathbb{R}^d$ ומטריצת שונות משותפת $\Sigma \in \mathbb{R}^{d \times d}$. אזי פונקציית ההתפלגות המשותפת היא

$$f(X) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2}(X - \mu)^T \Sigma^{-1}(X - \mu)\right)$$

(כאשר $|\Sigma| = \det(\Sigma)$)
במקרה זה נסמן $X \sim \mathcal{N}(\mu, \Sigma)$.

הערה אם $d = 1$ נקבל את ההתפלגות הגאוסית.

יהי $X \sim N(\mu, \Sigma)$ וקטור מקרי גאוסיאני דו מימדי, עם מטריצת שונותיות משותפות אלכסונית, כלומר,

$$X \sim N\left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{pmatrix}\right)$$

אזי ההתפלגויות השוליות הן

$$f(X_i) = \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{1}{2} \left(\frac{X_i - \mu_i}{\sigma_i}\right)^2\right)$$

בחזרה לדוגמא של גובה ומשקל.

נוכל למדל כל דגימה כוקטור $\vec{x} \in \mathbb{R}^d$, שהוא realization של הוקטור $X := (X_{weight}, X_{height})^T$. לפיכך, המדגם שלנו הוא $S := \{\vec{x}_1, \dots, \vec{x}_n\}$, כאשר כל \vec{x}_i הוא realization של X . הקואורדינאטה ה- j של הדגימה ה- i היא $\vec{x}_i(j)$ או x_{ij} . הקואורדינאטה ה- j היא realization של הקואורדינאטה ה- j של X_i , שהיא $X_i^{(j)}$.

נרצה לשערך את ההתפלגות של המשתנים

$$\vec{x}_1, \dots, \vec{x}_m \stackrel{\text{i.i.d.}}{\sim} N\left(\begin{pmatrix} 75 \\ 170 \end{pmatrix}, \begin{pmatrix} 3 & 1 \\ 1 & 5 \end{pmatrix}\right)$$

האומד ממוצע המדגם רב המשתנים שלנו הוא פשוט האומד ממוצע המדגם לכל משתנה בנפרד:

$$\hat{\mu} = \begin{pmatrix} \vdots \\ \hat{\mu}_j \\ \vdots \end{pmatrix} = \begin{pmatrix} \vdots \\ \frac{1}{m} \sum_i x_{i,j} \\ \vdots \end{pmatrix}$$

האומדן של ה-Cov מעט יותר מורכב, שכן יש תלויות בין המשתנים. האומד שונות המדגם המשותפת הבלתי מוטה בין X_i ל- X_j נתון על ידי

$$\hat{\sigma}(X_i, X_j) = \frac{1}{m-1} \sum_k (x_{ki} - \hat{\mu}_i)(x_{kj} - \hat{\mu}_j)$$

לבסוף, האומד מטריצת השונות המשותפת של המדגם היא מטריצה $\hat{\Sigma}$ מסדר $d \times d$ כך ש-

$$\hat{\Sigma}_{ij} := \hat{\sigma}(X_i, X_j)$$

בכתיב וקטורי, עבור $X \in \mathbb{R}^{m \times d}$ מטריצה ששורותיה הן הדגימות $\vec{x}_1, \dots, \vec{x}_m$

• מטריצת שונות המדגם המוטה נתונה על ידי

$$\hat{\Sigma} := \frac{1}{m} \sum_{i=1}^m (\vec{x}_i - \hat{\mu})(\vec{x}_i - \hat{\mu})^T = \frac{1}{m} \tilde{X}^T \tilde{X}$$

כאשר \tilde{X} היא המטריצה הממורכזת: $\tilde{X}_{:,i} = X_{:,i} - \hat{\mu}$

2 תרגול - קצת לינארית

הגדרה

פונקציה $d : X \times X \rightarrow \mathbb{R}$ (עבור קבוצה X כלשהי) תקרא מטריקה אם

$$1. \quad d(x, y) = 0 \text{ אם } x = y$$

$$2. \quad d(x, y) = d(y, x)$$

$$3. \quad d(x, y) \leq d(x, z) + d(z, y)$$

דוגמא

$$d(u, v) = \sum_i |u_i - v_i| \text{ היא מטריקה.}$$

הוכחה

$$1. \quad \sum_i |u_i - v_i| = 0 \text{ אם כל אחד מהאיברים בסכום הוא } 0 \text{ אם } u_i = v_i \text{ לכל } i \text{ אם } u = v$$

$$2. \quad d(u, v) = \sum_i |u_i - v_i| = \sum_i |v_i - u_i| = d(v, u)$$

3.

$$\begin{aligned} d(u, v) &= \sum_i |u_i - v_i| = \sum_i |(u_i - w_i) + (w_i - v_i)| \leq \\ &\leq \sum_i |u_i - w_i| + \sum_i |w_i - v_i| = d(u, w) + d(w, v) \end{aligned}$$

הגדרה

פונקציה $\|\cdot\| : \mathbb{R}^d \rightarrow \mathbb{R}$ תיקרא נורמה אם היא מקיימת את התכונות הבאות:

$$1. \quad \|v\| \geq 0 \text{ לכל } v, \text{ ו-} \|v\| = 0 \text{ אם } v = 0$$

$$2. \quad \|\alpha v\| = |\alpha| \|v\|$$

$$3. \quad \|v + u\| \leq \|v\| + \|u\|$$

דוגמאות

$$1. \quad \text{נורמת } \ell_1 \text{ המוגדרת על ידי } \|v\| = \sum_i |v_i|$$

$$2. \quad \text{נורמת } \ell_2 \text{ המוגדרת על ידי } \|v\|_2 = \left(\sum_i |v_i|^2 \right)^{\frac{1}{2}}$$

$$3. \quad \text{נורמת } \ell_\infty \text{ המוגדרת על ידי } \|v\|_\infty = \max_i |v_i|$$

הכללה

הנורמת הללו הן חלק ממשפחת נורמות ℓ_p , המוגדרת על ידי

$$\|v\|_p = \left(\sum_i |v_i|^p \right)^{\frac{1}{p}}$$

הגדרה

כדור היחידה לפי הנורמה $\|\cdot\|$ מוגדר על ידי

$$B_{\|\cdot\|} = \{x \in V \mid \|x\| \leq 1\}$$

הגדרה

פונקציה $\langle \cdot, \cdot \rangle : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ נקראת מכפלה פנימית אם היא מקיימת את התכונות הבאות:

1. סימטריות: $\langle u, v \rangle = \langle v, u \rangle$
2. לינאריות: $\langle \alpha v, u + w \rangle = \alpha \langle v, w \rangle + \alpha \langle v, u \rangle$
3. אי-ניוון: $\langle u, u \rangle \geq 0$ ו- $\langle u, u \rangle = 0$ אם ורק אם $u = 0$.

הגדרה

נגדיר נורמה על ידי $\|v\| = \langle v, v \rangle^{\frac{1}{2}}$

הגדרה

$$\cos \theta = \frac{\langle v, u \rangle}{\|v\| \|u\|}$$

הגדרה

נאמר כי v מאונך ל- u ונסמן $v \perp u$ אם $\langle u, v \rangle = 0$.

הגדרה

הטלה אורתוגונית של u על v הוא

$$\vec{P} = p\hat{u}$$

כאשר

$$p = \langle v, \hat{u} \rangle$$

$$\hat{u} = \frac{u}{\|u\|}$$

כלומר,

$$\vec{P} = \langle v, \hat{u} \rangle \hat{u}$$

הגדרה

אם $v \in \mathbb{R}^n$ ו- $u \in \mathbb{R}^m$, אזי המכפלה החיצונית של v, u היא מטריצה מסדר $n \times n$ כך ש-

$$(v \cdot u)_{ij} = v_i u_j$$

הדרגה של מטריצה זו היא 1.

הגדרה

יהי $V \subseteq \mathbb{R}^d$, $\dim(V) = k$, v_1, \dots, v_k בסיס אורתונורמלי ל- V . אזי מטריצת ההטלה האורתוגונלית על המרחב V היא

$$P = \sum_i v_i \otimes v_i = \sum_i v_i v_i^t$$

(P היא מטריצה מגודל $d \times d$ מדרגה k)

מסקנה

אם $x \in \mathbb{R}^d$, אזי יתקיים

$$Px \in V$$

למשל, אם $k = 1$ נקבל

$$Px = (v_1 v_1^t) x = v_1 v_1^t x = v_1 \langle v_1, x \rangle$$

שזו בדיוק ההטלה של x על $\text{span}\{v_1\}$ (כי $\|v_1\| = 1$).

דוגמא לאלגוריתם למידה

נניח כי יש d משוואות בת"ל ב- d נעלמים מהצורה $y_i = x_i^t w = \sum_j x_{ij} w_j$ לכל $i \in [d]$.

נתונים x, y . נרצה למצוא את w .
כלומר, נתונים $y = \begin{pmatrix} y_1 \\ \vdots \\ y_d \end{pmatrix}$ ו- $x = \begin{pmatrix} \cdots & x_1 & \cdots \\ & \vdots & \\ \cdots & x_n & \cdots \end{pmatrix}$ ומטרנו היא למצוא w כך ש- $y = xw$.
נרצה לבנות את הפונקציה הבאה:

```
def lig_reg(x, y):
    ...
    return w
```

המשוואות בת"ל, ולכן דרגתה של x מלאה ולכן קיימת לה מטריצה הופכית ולכן נוכל להשתמש באלגוריתם המחזיר את

$$w = x^{-1} x w = x^{-1} y$$

הגדרה

בהינתן מטריצה A , ללכסן את A משמעותו למצוא מטריצה אלכסונית D ומטריצה הפיכה P כך ש-

$$D = PAP^{-1}$$

$$P^{-1}DP = A$$

פירוק לערכים עצמיים (Eigen Value Decomposition)

תהי $A \in \mathbb{R}^{d \times d}$ סימטרית. אזי קיימות $U, D \in \mathbb{R}^{d \times d}$ כך ש- $A = UDU^t$ המקיימות

1. U אורתוגונלית שעמודותיה הן ו"ע של A

2. D אלכסונית שערכי האלכסון שלהם הם הע"ע של A

דוגמא

נתבונן במטריצה $A = \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}$. נבצע לה פירוק EVD. נמצא ע"ע על ידי מציאת שורשים לפולינום:

$$\det(A - \lambda I) = 0$$

השורשים הם $\lambda_1 = 3, \lambda_2 = 1$, ולכן $D = \begin{pmatrix} 3 & 0 \\ 0 & 1 \end{pmatrix}$. כעת נותר למצוא את הוקטורים העצמיים, ונקבל

$$U = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{pmatrix}$$

אם יש לנו וקטור $x \in \mathbb{R}^d$ המקיים $\|x\| = 1$, אזי

$$Ax = (UDU^t)x = UD \begin{pmatrix} x^t u_1 \\ \vdots \\ x^t u_d \end{pmatrix} = U \begin{pmatrix} \lambda_1 x^t u_1 \\ \vdots \\ \lambda_d x^t u_d \end{pmatrix} = \sum \lambda_i \langle x, u_i \rangle u_i$$

פירוק לערכים סינגולריים (Singular Value Decomposition)

הגדרה

u, v הם וקטורים סינגולריים של A המתאימים לערך סינגולרי σ אם $Av = \sigma u$. במקרה זה v יקרא וקטור סינגולרי ימני ו- u יקרא וקטור סינגולרי שמאלי.

משפט ה-SVD

תהי $A \in \mathbb{R}^{m \times d}$. אזי $A = U\Sigma V^T$ כאשר

- $U \in \mathbb{R}^{m \times m}$ אורתוגונלית שעמודותיה u_1, \dots, u_m ו"ס שמאליים של A .
- $V \in \mathbb{R}^{d \times d}$ אורתוגונלית שעמודותיה v_1, \dots, v_d ו"ס ימניים של A .
- $\Sigma \in \mathbb{R}^{m \times d}$ "אלכסונית" שערכי ע"ס של A .

הערה

המשמעות של מטריצה אלכסונית היא ש- Σ היא מהצורה $\begin{pmatrix} D & 0 \end{pmatrix}$ כאשר D אלכסונית.

משמעות

יהי $x \in \mathbb{R}^d$. נזכור כי $x = \sum_i \langle x, v_i \rangle v_i$.

$$Ax = A \sum_i \langle x, v_i \rangle v_i = \sum_i \langle x, v_i \rangle U \Sigma V^T v_i =$$

נבחין כי

$$[V^T v_i]_j = \delta_{ij}$$

ולכן

$$V^T v_i = e_i$$

(כאשר e_i הוקטור ה- i בבסיס הסטנדרטי)
ולכן

$$= \sum_i \langle x, v_i \rangle U \Sigma e_i = \sum_i \sigma_i \langle x, v_i \rangle U e_i = \sum_i \sigma_i \langle x, v_i \rangle u_i$$

שאלה

כיצד נמצא את V, U ?

תשובה

נניח שבחרנו בסיס אורתונורמלי כלשהו v_1, \dots, v_d . נבחין כי אם נפעיל את A על הבסיס נקבל

$$Av_1, \dots, Av_d$$

אבל אין שום סיבה ש- Av_1, \dots, Av_d יהיו אורתונורמליים!
נבחין כי $A^T A$ היא מטריצה סימטרית, ולכן קיים לה פירוק EVD.
נניח שיש לנו פירוק SVD.

$$\begin{aligned} A^T A &= (U \Sigma V^T)^T U \Sigma V^T = V \Sigma^T \underbrace{U^T U}_I \Sigma V^T = \\ &= V \underbrace{\Sigma^T \Sigma}_D V^T \end{aligned}$$

לכן, בהינתן A נוכל לבצע EVD על $A^T A$, ונקבל v_1, \dots, v_d בסיס אורתונורמלי ו- $\lambda_1, \dots, \lambda_d$.
נוכל לבצע תהליך דומה על AA^T ונקבל u_1, \dots, u_m בסיס אורתונורמלי ו- $\lambda_1, \dots, \lambda_n$.
ויתקיים

$$Av_i = \sqrt{\lambda_i} u_i$$