

פרק 13: רגולריזציה, פשטות וידע מוקדם

Ridge Regression ו- Lasso

מבוא

לפני שנדון ברגרסיה נתבונן בבעיית שערך ממוצע תחת אילוצים:

$$\text{Minimize } \sum_{i=1}^n (Y_i - \mu)^2 \text{ subject to } \mu^2 \leq C$$

בעזרת כופלי לגרנז' הבעיה שקולה ל:

$$\text{Minimize } \sum_{i=1}^n (Y_i - \mu)^2 + \lambda_c \mu^2$$

נגזור:

$$-2 \sum_{i=1}^n (Y_i - \hat{\mu}_c) + 2\lambda_c \hat{\mu}_c = 0$$

ונקבל:

$$\hat{\mu}_c = \frac{\sum_{i=1}^n Y_i}{n + \lambda_c} = K_c \bar{Y}, \quad K_c = \frac{n}{n + \lambda_c} \quad (*)$$

ורואים שהאפקט של C קריטי:

$$\begin{aligned} C \rightarrow 0, \quad \hat{\mu}_c &\rightarrow \bar{Y} \\ C \rightarrow \infty \quad \hat{\mu}_c &\rightarrow 0 \end{aligned}$$



מסקנה: האפקט של רגולריזציה הוא אפקט מקוץ (shrinking).

הערה: ניתן להוסיף באופן שקול מס' כלשהו n_c של נקודות מלאכותיות $Y_i = 0$ שהרי נקבל את הבעיה

$$\text{Minimize } \sum_{i=1}^n (Y_i - \mu)^2 + n_c (0 - \mu)^2$$

בתרגול נראה ש:

$$\text{Minimize } \sum_{i=1}^n (Y_i - \mu)^2 + \lambda_c \mu^2$$

שקול לשערך MAP (maximum a-posteriori) של μ עם פריור גאוס.

הסתכלות שקולה לפיכך היא להטות את הפתרון (bias) בכיוון רצוי. השאלה היא כיצד למצוא הטייה מועילה. האינטרפרטציות הן:

1. בייסיאנית: מציאת פריור מתאים
2. סיבוכיות: פתרון פשוט (נורמה קטנה) לעומת פתרון מסובך (נורמה גדולה)

Ridge Regression

נתבונן כעת ברגרסיה רב מימדית ונניח:

1. ל- X ממוצע 0, וקטור p מימדי.

2. ל- Y ממוצע 0

נחפש רגרסור מהצורה $Y \approx \beta^T X$. למודל כזה קוראים ה- "מודל הסטנדרטי".

נגדיר את הפונקציה הבאה:

$$SSE_\lambda(\beta) = \sum_{i=1}^n (Y_i - X_i^T \beta)^2 + \lambda \|\beta\|_2^2$$

ה- β שמביא למינימום את $SSE_\lambda(\beta)$ הוא הפתרון של בעיית ה- ridge regression:

$$\text{Minimize } \sum_{i=1}^n (Y_i - X_i^T \beta)^2 + \lambda \|\beta\|_2^2$$

קל לראות שהבעייה שקולה לבעיית אופטימיזציה עם אילוצים על הנורמה של הרגרסור β (מכופלי לגרנז') ולפיכך בעיית האופטימיזציה היא בעיה קמורה.

שאלה: מה אוסף הדוגמאות המלאכויות עבורן נקבל את הבעיה דלעיל ללא איבר הרגולריזציה?

נגזור לפי β ונקבל:

$$\frac{\partial}{\partial \beta(k)} SSE_\lambda(\beta) = 2 \sum_{i=1}^n (Y_i - X_i^T \beta) X_i^T(k) + 2\lambda \beta(k)$$

וברישום מטריצי לאחר השוואה ל-0 נקבל:

$$-Y^T X + \hat{\beta}_\lambda^T (X^T X + \lambda I) = 0$$

(נשים לב: β_λ היא וקטור p מימדי, Y וקטור n מימדי ו- X מטריצה $n \times p$).

ולכן:

$$\hat{\beta}_\lambda^T = Y^T X (X^T X + \lambda I)^{-1}$$

$$\hat{\beta}_\lambda = (X^T X + \lambda I)^{-1} X^T Y \quad \text{או}$$

(היפוך המטריצה מותר תמיד כי $\lambda > 0$).

זו הגרסה הוקטורית של משעריך הממוצע (*). נעיר שהאינטרפטציה הבייסיאנית עובדת גם כאן (ראה תרגול)

LASSO

LASSO היא בעיית האופטימיזציה בה נורמת 2 של איבר הרגולריזציה מוחלף ע"י נורמת 1, ז"א:

$$\text{Minimize} \quad \sum_{i=1}^n (Y_i - X_i^T \beta)^2 + \lambda \sum_{k=1}^p |\beta(k)|$$

קל לראות שזוהי בעיית קמורה עם אילוץ על נורמת 1 של β . למרבה הצער, הערך המוחלט מונע מאיתנו גזירה נוחה של פונקציית המטרה. למרות זאת ניתן לחשב את הפסאודו-גרדיאנט של פונקציית המטרה וקיימים אלגוריתמים יעילים לפיתרון בעיית האופטימיזציה.

אלגוריתם ה-LASSO עובד היטב אם קיימים פתרונות דלילים (ז"א רוב הקאורדינטות של β הם 0).

בחירת λ

אם היינו יודעים את ה-MSE של מסווג מסויים היינו בוחרים את המסווג עם ה-MSE המינימלי. ניתן לשערך את ה-MSE ע"י אימות צולב (Cross Validation).

שיטה אחרת היא לרשום חסם הכללה התלוי ב- λ , כפי שנעשה בהרצאה הדנה בתיאוריה.

נסיים בהבחנה מעניינת. ניתן להתבונן במטריצה הבאה:

$$S_\lambda = (X^T X + \lambda I)^{-1} X^T$$

המקיימת (עבור Ridge Regression): $\hat{\beta}_\lambda = S_\lambda Y$. ניתן לחשוב על S_λ כהטלה ממרחב n מימדי למרחב p מימדי. עבור $\lambda = 0$ מספר דרגות החופש הוא p , וככל שנגדיל את λ יקטן מספר דרגות החופש האפקטיביות שניתן להעריכו ע"י $n - TR(S_\lambda)$.