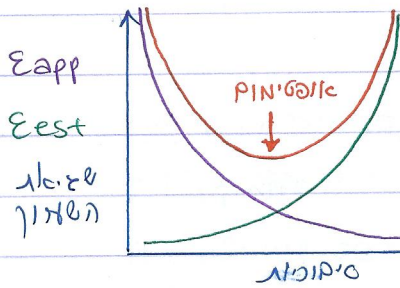


תאוריה:



מתייחסות: חזקים למצוא את הנק' האופטימלית

ה- bias-variance tradeoff

אם נבחר סימכות קטנה מדי

נעשה underfitting, ועצומה מדי

נעשה overfitting.

אפשר להתקדם לאופטימום בעזרת תאוריה.

אם כן במצוי אופטימיזציה ניסיון למצוא: $h^* = \argmin_{h \in H} F_S(h)$

(כאשר F_S היא פונק' החיזוי). החיזוי הוא שמצגשו נוסח למצוא:

אם $h^* = \argmin_{h \in H} (F_S(h) + \lambda R(h))$, כלומר נוסף איבר אלג' קבוצה

אלא רק בפרופורציה עצומה: בקטנה של שמצגיה כמה היא מסומכת.

איבר זה נקרא איבר הרגולרציה, ובפרקציה ליצור יציג כוונתן ששולטת

אם רמת הסימכות של ההיפוטזה שלנו מקבלים כן שניה קרובים לאופטימום.

(*) אם צרכים עצומים של λ : נכפול את $R(h)$ במקדם גדול, ולקטור

ה מסומכת נקבל גולגה גקורה לעצמך אומנו מלהיג ה- \argmin .

אם למצגה נעזף פונק' א סימכות נמוכה שק נשק אליה פחות.

אם נקבל שונת נמוכה ו- bias גקורה.

(*) אם צרכי ג קטנים: האיבר האומנו יהיה הצומיננטי (התנסו א

סימכות ה יהיה קטן משמעותי) ונעזף היפוטזה מסומכת גקורה.

אם נקבל שונת גקורה ו- bias נמוך.

דוגמא לאיבר הרגולרציה: נראה איך נימ' לתקום פונק' שונת:

ה- Mean Least Squares Setting יהיה הצרכי לנאליג: $F_S(h)$ יהיה

Ridge - ℓ_2 regularization: $\hat{w}_\lambda^{\text{ridge}} = \argmin_w RSS(\hat{y}; w) + \lambda \|w\|_2^2$ 1

* כל של- w הפגון יש נומה גקורה יגרי, כך ההיפוטזה מסומכת יגרי.

ההוכחה \Rightarrow למנה: גרא $u \in \mathbb{R}^{n \times n}$ מתייחסת אומנוג'ית, $D \in \mathbb{R}^{n \times n}$ מתייחסת

אלמנוג'ית. כל: $uDu^T + \lambda I = u(D + \lambda I)u^T$

הערה: גרא X design matrix ו- response vector y ק ש:

$u = X^T X$ פונק' ה- SVD של X , כל: ה- ridge estimator

הוא: $\hat{w}_\lambda^{\text{ridge}} = u \Sigma_\lambda^T v^T y$ כאשר $[\Sigma_\lambda^T]_{ii} = \frac{\sigma_i}{\sigma_i^2 + \lambda}$

הוכחה: $\lambda \geq 0$ פונקציה, $f_{\lambda_2}(w) = \text{RSS}(\beta; w) + \lambda \|w\|_2^2$: (10)

בשלב 1. נסתכל על w וננסה להפחית את λ עד שיהיה 0.

$$\frac{1}{2} \nabla_w f_{\lambda_2}(w) = \underbrace{XX^T w - Xy}_{\text{least squares}} + \underbrace{\lambda w}_{\text{regulation}}$$

אנחנו רוצים $XX^T w - Xy + \lambda w = 0 \Leftrightarrow (XX^T + \lambda I)w = Xy$: נשווה לאפס

וקיבלנו: $\hat{w}_{\lambda}^{\text{ridge}} = (XX^T + \lambda I)^{-1} Xy$

להלן הפירוק SVD:

$$\begin{aligned} \hat{w}_{\lambda}^{\text{ridge}} &= (XX^T + \lambda I)^{-1} Xy = (U \Sigma V^T V \Sigma^T U^T + \lambda I)^{-1} Xy = \\ &= (U \Sigma \Sigma^T U^T + \lambda I)^{-1} Xy = (U \Sigma^2 U^T + \lambda U U^T)^{-1} Xy = \\ &= U (\Sigma^2 + \lambda I)^{-1} U^T Xy = U (\Sigma^2 + \lambda I)^{-1} U^T U \Sigma V^T y = \\ &= U (\Sigma^2 + \lambda I)^{-1} \Sigma V^T y = U \Sigma_{\lambda} V^T y \end{aligned}$$

פונקציה: $[\Sigma_{\lambda}]_{ii} = \frac{\sigma_i^2}{\sigma_i^2 + \lambda}$, כאשר σ_i הם הערכים העצמיים של X .

המשוואה Ridge Regression היא משוואה

רגולרית (Ordinary Least Squares) OLS

עבור $X_t = [X | \sqrt{\lambda} I] \in \mathbb{R}^{d \times (m+d)}$, $y_{\lambda} = [y | 0] \in \mathbb{R}^{m \times d}$: שני

"הצבה"

"הצבה"

אם נסתכל על X_t, y_t של OLS פשוט (כפי שנקראו במקור) :

$$\hat{w}_{\lambda}^{\text{OLS}} = (X_{\lambda} X_{\lambda}^T)^{-1} X_{\lambda} y_{\lambda} = (X_{\lambda}^{\dagger})^T y_{\lambda}$$

כאשר X_{λ}^{\dagger} היא ה-pseudo-inverse של X_{λ} .

$$(XX^T)^{-1} X = (X^{\dagger})^T \text{ כאשר } p \text{ הוא מספר ה} \Sigma_{ii}^{\dagger} = \begin{cases} 1/\sigma_{ii} & : \Sigma_{ii} \neq 0 \\ 0 & : \Sigma_{ii} = 0 \end{cases}$$

אם X_{λ} היא ה-pseudo-inverse של X_{λ} .

הוא קרוי Ridge-Regression: X, y ו λ .

(*) גבולות חשיבות של ה-ridge estimator הוא שהוא מונע את התפרקות

של X (הוא מונע את (כאשר) λ קטן מדי) ופונקציה $(XX^T)^{-1}$ קטנה.

והוא קרוי רגולריות, כלומר שהוא מונע את התפרקות $(XX^T + \lambda I)^{-1}$.

והוא גם λ (ל- λ קטן) גורם את XX^T להיות מוגדרת, באופן

ואם לא נהיה מסוגלים להגדיר קטן מדי.

(*) גבולות חשיבות (נונים מוגדרים) : ה-estimator הוא biased

אך במקרה הזה זה גורם לשיפור הקצרים (אם כי זה גורם לשינוי λ).

2. Lasso - ℓ_1 Regularization

ממלאים למקרה שבו מעבר לסף העוצמה (למנוף מהפגות) של
 פהיו מסיקוסו קצונה מאוזן ולקצו (overfit) הוא מיצר פגונג
 sparse : הוא מאפס חלק מהפצרים.

* זה מועיל לנו, כי זה מקל את העלם האמיתי שבו ישנם פיצרים חלוקים
 פהיו מאחרים, אינם לא נכנה למנו. \Leftarrow דבר פגשו feature selection.

מה שהביא אמתנו ל- Lasso הוא מקרה ה- Best-subset Selection :
 במקרה זה איבר העוצמה הוא עם "נומג" אפס :

$$\textcircled{\star} \argmin_{\omega} f_{\ell_0}(\omega) = \underbrace{\argmin_{\omega} \|y - X^T \omega\|_2^2}_{\text{(Least Squares)}} + \lambda \|\omega\|_0$$

$\|\omega\|_0$: סופר את מס' הניכיוס שלקם אפס.

כך ה- best-subset Selection בוחר את הפצרים החלוקים עם איפס
 הפצרים הנכונים. כך הקטנו את המימז d.

המקרה : מכלל Best subset of features (הא קצת אופטימיזציה NP-H
 ותא איעה קמורה. (כי $\|\omega\|_0$ אינה נורמה, ט נורמה אכן קמורה).

\Leftarrow זה מבקש אמתו ℓ_0 & לבנוח חסוקה בבחירה נורמה ℓ_1 :

1. Sparsity : פקור $q \leq 1$, אכן נקח כ- $\textcircled{\star}$ נורמה ℓ_1 , הפגונג
 (צולות) שנקב יהו צולוים כמו שהצרכנו מפי. ונקב
 active set של פיצרים ממימז קטן ℓ_1 .

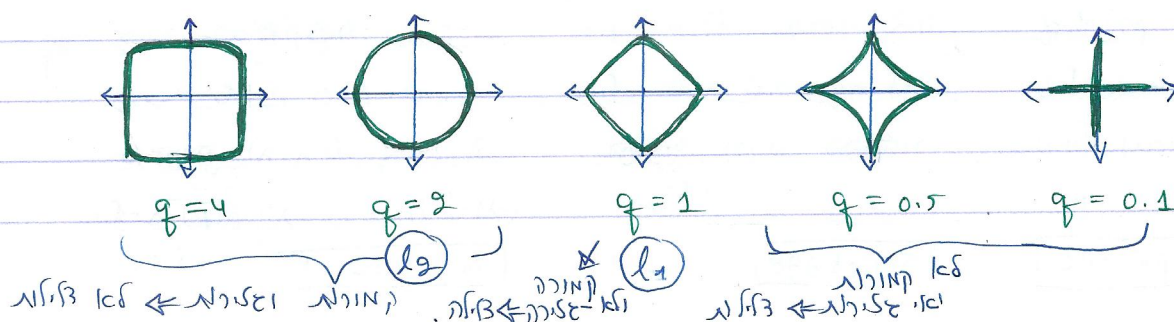
2. Convexity : פקור $q \geq 1$ קצת האופטימיזציה המועלנה עם החסר
 (קמיח) נורמה $\|\omega\|_q$ כ- ℓ_1 ויהיה קצת אופטימיזציה קמורה, $\textcircled{\star}$ (קמורה ואיבר
 הנאור, צציה זק)

פגמורה ראלו אלפי הפגורים אולג פיעלוג (קאמן פול).

קצת ה- Lasso Regression פוקת $q=1$ ומכוויחה אל & המנוח האול

$$\argmin_{\omega} f_{\ell_1}(\omega) = \argmin_{\omega} \|y - X^T \omega\|_2^2 + \lambda \|\omega\|_1$$

מה מקשים פגונג צולוים? ,במקן אל כצור פחציה פקור נורמה שונג :



אנטיאליזיה עבר

לשפרנו בן הסדרה לבנינו: $f_{\lambda_0}(\omega)$ היא קצת אופטימיזציה גור

אילונו. (כאילו כרזרסיה דוגמא עבר).

הפרמון האופטימלי של הקציה תוארה (עם איבר העצומיות)

מקרה קומון בין משטח האילוסטר (אבנו = תצורה) עם קו העקרה

(ה) - contour במרחב ω של הפוך $\|y - X^T \omega\|_2^2$ שיהא למטה

אינסוף. תצורה האליפטור נוקט מפרוק ה-SVD של X התיכני בילג.

תצורה מסר פניו, אין סיבה שהמחיר יגרוש על אחת הצירים, אך לעבור

על פניו על הצירים (שמוצגת על $q = 1 - \delta$ כפי שראו)

המחיר בן יקרים על אחת הצירים ולכן הפרמון יהיה על הצירים, כמות

חלק מהקובץ שלו יהיו 0.

מקרה פרטי: Orthogonal design: $XX^T = I_d$, X - design matrix

במקרה זה לפרמון ש צורה סגורה (כמו שהיתה ב-Ridge), כולל אין.

נראה את 3 הצורות שלפניו עלינו:

$$(i) \hat{\omega}_{\lambda}^{ridge} = \frac{\hat{\omega}^{OLS}}{1+\lambda}$$

\Leftrightarrow

במקרה הזה, קוא כל ילד פשוט מהבאלי

$$(ii) \hat{\omega}_{\lambda}^{lasso} = \eta_{\lambda}^{soft}(\hat{\omega}^{OLS})$$

where

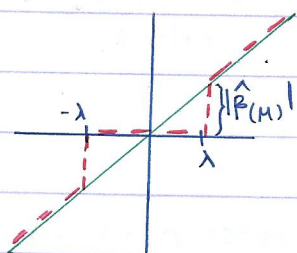
$$\eta_{\lambda}^{soft}(x) = \text{sign}(x) [|x| - \lambda]_+$$

$$= \begin{cases} x - \lambda & : x > \lambda \\ 0 & : |x| < \lambda \\ x + \lambda & : x \leq -\lambda \end{cases}$$

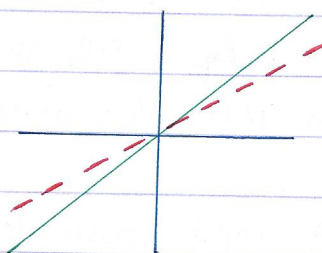
$$(iii) \hat{\omega}_{\lambda}^{subset} = \eta_{\lambda}^{hard}(\hat{\omega}^{OLS})$$

where

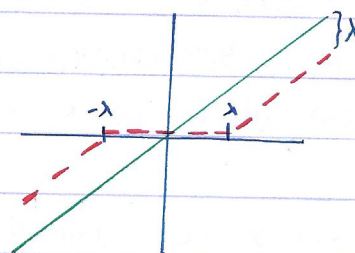
$$\eta_{\lambda}^{hard}(x) = \mathbb{1}_{[|x| \geq \lambda]} \cdot x$$



Best Subset



Ridge



Lasso



הפרמון
העצומי
הפסגון
על פני
OLS

עשה רק סבסטיקציה:

לא עשה סבסטיקציה,

עשה סבסטיקציה: שוח

שוח עניס לאנס 0-5.

לא שוח לים עבר 0-5.

עניס לאנס 0-5.

מחול למדוס לא מלטה

עצבים קטנים בקטור

לא תלבים, מוקטנים ק-ג.

פוס

$$\frac{1}{1+\lambda}$$

כומר מכולל לא, הקצורה.

סבסטיקציה

Shrinkage

Shrinkage + סבסטיקציה

response vector y - X design matrix

OLS - \hat{w} - $\eta_{\lambda}^{\text{soft}}(x)$: $\hat{w}^{\text{lasso}}(\lambda) = \eta_{\lambda}^{\text{soft}}(x)$

$\hat{w} = (X^T X)^{-1} X^T y$: \hat{w} - $\eta_{\lambda}^{\text{soft}}(x)$

(\hat{w} - $\eta_{\lambda}^{\text{soft}}(x)$)

: P - P λ

$$\begin{aligned} f_{\lambda}(w) &= \frac{1}{2} \|y - X^T w\|_2^2 + \lambda \|w\|_1 = \frac{1}{2} (\|y\|_2^2 - 2y^T X^T w + w^T X X^T w) + \lambda \|w\|_1 \\ &= \frac{1}{2} \|y\|_2^2 + \sum_{j=1}^d (\frac{1}{2} w_j^2 - \hat{w}_j w_j + \lambda |w_j|) \\ &= \frac{1}{2} \|y\|_2^2 + \sum_{j=1}^d (\frac{1}{2} w_j^2 - \hat{w}_j w_j + \lambda \text{sign}(w_j) w_j) \end{aligned}$$

: P - P λ \hat{w} - $\eta_{\lambda}^{\text{soft}}(x)$

$$0 = \frac{\partial}{\partial w_j} \left(\frac{1}{2} \|y - X^T w\|_2^2 + \lambda \|w\|_1 \right) = \frac{\partial}{\partial w_j} \left(\frac{1}{2} \|y\|_2^2 + \sum_{j=1}^d (\frac{1}{2} w_j^2 - \hat{w}_j w_j + \lambda \text{sign}(w_j) w_j) \right) =$$

$$= \frac{\partial}{\partial w_j} (\frac{1}{2} w_j^2 - \hat{w}_j w_j + \lambda \text{sign}(w_j) w_j) = \frac{\partial}{\partial w_j} (\frac{1}{2} w_j^2 - \hat{w}_j w_j + \lambda \text{sign}(w_j) w_j) =$$

$$= w_j - \hat{w}_j + \lambda \text{sign}(w_j) \iff w_j = \hat{w}_j - \lambda \text{sign}(w_j)$$

: P - P λ \hat{w} - $\eta_{\lambda}^{\text{soft}}(x)$

: λ $\hat{w}_j < \lambda$ P λ $\hat{w}_j > \lambda$

($\text{sign}(w_j) = -1$)

- $\hat{w}_j + \lambda > \lambda - \lambda = 0$: $w_j < 0$: $\hat{w}_j - \lambda \text{sign}(w_j) < 0$: $w_j < 0$: $\hat{w}_j - \lambda \text{sign}(w_j) < 0$

- $\hat{w}_j - \lambda > \lambda - \lambda = 0$: $w_j > 0$: $\hat{w}_j - \lambda \text{sign}(w_j) > 0$: $w_j > 0$: $\hat{w}_j - \lambda \text{sign}(w_j) > 0$

- $\hat{w}_j - \lambda < \lambda - \lambda = 0$: $w_j = 0$: $\hat{w}_j - \lambda \text{sign}(w_j) = 0$: $w_j = 0$: $\hat{w}_j - \lambda \text{sign}(w_j) = 0$

: λ $\hat{w}_j = 0$ P λ $\hat{w}_j = 0$

$w_j \geq 0$: $\hat{w}_j \geq \lambda$: $w_j = \hat{w}_j - \lambda \text{sign}(w_j)$: $w_j = \hat{w}_j - \lambda$: $\hat{w}_j \geq \lambda$: $w_j = \hat{w}_j - \lambda$

$w_j \leq 0$: $\hat{w}_j \leq -\lambda$: $w_j = \hat{w}_j - \lambda \text{sign}(w_j)$: $w_j = \hat{w}_j + \lambda$: $\hat{w}_j \leq -\lambda$: $w_j = \hat{w}_j + \lambda$

$w_j \geq 0$: $\hat{w}_j \geq \lambda$: $w_j = \hat{w}_j - \lambda \text{sign}(w_j)$: $w_j = \hat{w}_j - \lambda$: $\hat{w}_j \geq \lambda$: $w_j = \hat{w}_j - \lambda$

$w_j \leq 0$: $\hat{w}_j \leq -\lambda$: $w_j = \hat{w}_j - \lambda \text{sign}(w_j)$: $w_j = \hat{w}_j + \lambda$: $\hat{w}_j \leq -\lambda$: $w_j = \hat{w}_j + \lambda$

$w_j \geq 0$: $\hat{w}_j \geq \lambda$: $w_j = \hat{w}_j - \lambda \text{sign}(w_j)$: $w_j = \hat{w}_j - \lambda$: $\hat{w}_j \geq \lambda$: $w_j = \hat{w}_j - \lambda$

$w_j \leq 0$: $\hat{w}_j \leq -\lambda$: $w_j = \hat{w}_j - \lambda \text{sign}(w_j)$: $w_j = \hat{w}_j + \lambda$: $\hat{w}_j \leq -\lambda$: $w_j = \hat{w}_j + \lambda$

$\hat{w}^{\text{lasso}}(x) = \eta_{\lambda}^{\text{soft}}(x)$

: λ \hat{w} - $\eta_{\lambda}^{\text{soft}}(x)$