

Lasso i lineære modeller

Betragt n observationer $\{x_i, y_i\}_{i=1}^n$, hvor $x_i = (x_{i1}, \dots, x_{ip})$ er en p dimension vektor af fork-larende variable eller prediktorer og $y_i \in \mathbb{R}$ er den tilhørende respons variabel.

1.1 Mindste kvadraters metode

Den velkendte estimator for mindste kvadraters metode for (β_0, β) findes ud fra

$$\min_{\beta_0, \beta} \left\{ \frac{1}{2n} \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 \right\}$$

Løsningen hertil er givet ved

$$\hat{\beta}^{\text{OLS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$

Som bekendt er estimatoren unbiased, men ofte har den en høj varians. GRUNDE TIL AT PRØVE ANDET END OLS

1.2 Ridge regression

Ridge regression estimatoren findes ud fra

$$\min_{\beta_0, \beta} \left\{ \frac{1}{2n} \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 \right\}, \quad \text{underlagt at } \sum_{j=1}^p \beta_j^2 \leq t, \quad (1.1)$$

som kan omskrives til et Lagrange problem

$$\min_{\beta_0, \beta} \left\{ \frac{1}{2n} \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}, \quad (1.2)$$

hvor $\lambda \geq 0$ er en såkaldt strafparameter, som bestemmes separat. Der er en en-til-en korre-spondance mellem det betingede problem (??) og Lagrange problemet (??). Første led i (??) svarer til OLS, som finder de estimerede koefficienter ved at minimere SSR, mens sidste led mindsker de estimerede koefficienter. På matrix-vektor form er løsningen af ridge regression givet ved

$$\hat{\beta}^R = (\mathbf{X}^T \mathbf{X} + \lambda I_p)^{-1} \mathbf{X}^T \mathbf{y}$$

1.3 Lasso

Lasso finder løsningen til optimerings problemet

$$\min_{\beta_0, \beta} \left\{ \frac{1}{2n} \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 \right\}, \quad \text{ub} \sum_{j=1}^p |\beta_j| \leq t \quad (1.3)$$

Betingelsen $\sum_{j=1}^p |\beta_j| \leq t$ kan skrives mere kompakt ved $\|\beta\|_1 \leq t$. Dette kan udtrykkes på matrix-vektor notation. Lad $\mathbf{y} = (y_1, \dots, y_n)$ være en n dimensional vektor med responsvariable og \mathbf{X} være en $n \times p$ matrix med $x_i \in \mathbb{R}^p$ som den i 'te række, da kan (??) omskrives til

$$\min_{\beta_0, \beta} \left\{ \frac{1}{2n} \|\mathbf{y} - \beta_0 \mathbf{1} - \mathbf{X}\beta\|_2^2 \right\}, \quad \text{s.t. } \|\beta\|_1 \leq t,$$

hvor $\mathbf{1}$ er en n dimensionel vektor bestående af 1 og $\|\cdot\|_2$ betegner den Euklidiske norm af vektorer.

Grænsen t begrænser summen af de absolutte værdier af parameter estimererne. Denne skal specificeres ved en ekstern procedure kaldet *kryds validering*, som vil blive diskuteret i kap –.

Ofte standardiseres prediktorerne \mathbf{X} således at kolonnerne er centeret og har varians 1. Dvs $\frac{1}{n} \sum_{i=1}^n x_{ij} = 0$ og $\frac{1}{n} \sum_{i=1}^n x_{ij}^2 = 1$. Hvis ikke prediktorerne standardiseres da vil lasso estimererne afhænge af enhederne. Hvis prediktorerne er målt i samme enhed, da vil vi typisk ikke standardisere. For fuldstændigheden, antager vi også at responsvariablen y_i er centeret, dvs $\frac{1}{n} \sum_{i=1}^n y_i = 0$. Når data er centeret da kan vi se bort fra skæringen β_0 i lasso optimeringen. Given en optimal lasso løsning $\hat{\beta}$ på det centeret data, kan vi finde løsningen for det ikke-centeret data. Der gælder at

$$\begin{aligned} \hat{\beta}^{\text{ikke-centeret}} &= \hat{\beta}^{\text{centeret}} \\ \hat{\beta}_0^{\text{ikke-centeret}} &= \bar{y} - \sum_{j=1}^p \bar{x}_j \hat{\beta}_j \end{aligned}$$

Derfor ser vi bort fra skæringen resten af kapitlet.

Vi kan omskrive lasso problemet til Lagrange form

$$\min_{\beta} \left\{ \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1 \right\} \quad (1.4)$$

for $\lambda \geq 0$. Af Lagrange dualiteten er der en bijektion mellem (??) og (??): for hver værdi af t hvor $\|\beta\|_1 \leq t$ er opfyldt, da findes en tilhørende værdi af λ som giver den samme løsning for (??). Mens løsningen $\hat{\beta}_\lambda$ til (??) løser grænse problemet med $t = \|\hat{\beta}_\lambda\|_1$

Variabel udvælgelsen for ridge regression og lasso illustreres på figur ??.

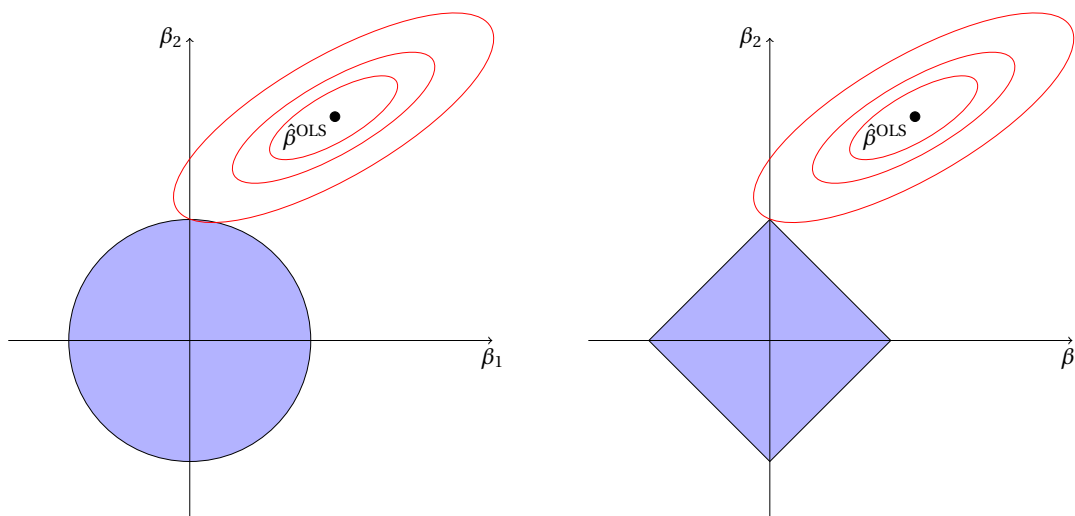


Figure 1.1: Konturer for SSR og betingelsesområderne for ridge regression (venstre) og lasso (højre). De blå arealer er betingelsesområderne $|\beta_1| + |\beta_2| \leq t$ og $\beta_1^2 + \beta_2^2 \leq t^2$, mens de røde ellipser er konturkurver for SSR. Konturkurverne har centrum i OLS estimatoren, $\hat{\beta}^{OLS}$.

For $p = 2$ underligges OLS betingelsen $\beta_1^2 + \beta_2^2 \leq t^2$ for ridge regression og betingelsen $|\beta_1| + |\beta_2| \leq t$ for lasso. Ellipserne omkring $\hat{\beta}^{OLS}$ er konturkurverne for SSR, dvs. SSR er konstant i en given ellipse. Værdien af SSR stiger, som ellipsen udvides fra $\hat{\beta}^{OLS}$. Ligningerne - og - indikerer at løsningen for ridge regression og lasso er givet ved det første punkt, hvor konturkurverne rammer betingelsesområdet. Siden ridge regression har et cirkulært betingelsesområde, vil skæringen med konturkurverne generelt ikke forekomme direkte på en akse. Modsat ridge regression har lassos betingelses område hjørner i hver akse, hvilket betyder, at hvis løsningen forekommer i et hjørne, da vil en af parametrene β_j være lig 0.

Hvis t er tilstrækkelig stor, da vil betingelsesområderne indeholde $\hat{\beta}^{OLS}$ og derfor vil ridge regression og lasso estimatorene være lig OLS estimatoren.

På figur ?? har vi blot betragtet det simple tilfælde hvor $p = 2$. Når $p = 3$ vil betingelsesområdet for ridge regression være en kugle, mens betingelsesområdet for lasso vil være en polydron.

Da lasso penalty ikke er differentialbel, findes der ikke en eksplicit løsning til lasso problemet. Vi antager at responsvariablerne y_i og prediktorene x_{ij} er standardiseret således at $\frac{1}{n} \sum_{i=1}^n y_i = 0$, $\frac{1}{n} \sum_{i=1}^n x_{ij} = 0$ og $\frac{1}{n} \sum_{i=1}^n x_{ij}^2 = 1$. Da kan vi se bort fra skæringen β_0 . Lagrange formen er nyttig for numerisk udregning af løsningen som findes vha en simpel procedure kaldet *coordinate descent*.

Vi kan opskrive objektfunktionen i (??) som

$$\frac{1}{2n} \sum_{i=1}^n \left(y_i - \sum_{k \neq j} x_{ik} \beta_k - x_{ij} \beta_j \right)^2 + \lambda \sum_{k \neq j} |\beta_k| + \lambda |\beta_j|$$

Vi kan se at løsningen for hver β_j kan udtrykkes ved den partial residual $r_i^{(j)} = y_i - \sum_{k \neq j} x_{ik} \hat{\beta}_k$,

som fjerner ... Da er den j 'te koefficient opdateret ved

$$\hat{\beta}_j = S_\lambda \left(\frac{1}{n} \langle \mathbf{x}_j, \mathbf{r}^{(j)} \rangle \right), \quad (1.5)$$

hvor $r_i = y_i - \sum_{j=1}^p x_{ij} \hat{\beta}_j$ er de fulde residualer. Den beskrevne algoritme svarer til metoden *cyclical coordinate descent*, som minimerer en konveks objektfunktion langs hver koordinat af gangen. Under milde regularitets betingelser, konvergerer løsningen til et global optimum. Fra opdateringen (??) ser vi at algoritmen foretager en univariat regression af den partial residual på hver prediktor, cycling gennem prediktorerne indtil konvergens. *pathwise coordinate descent*

Coordinate descent er særlig hurtig til at løse lasso problemet da ...

Homotopy metoder er en alternativ teknisk til at løse lasso problemet. Disse producerer en helt sti af løsninger i en frekventiel sekvens, ved at starte med nul. Denne sti er faktisk piecewise lineær. Algoritmen kaldet *least angle regression* (LARS) er en homotopy metode som effektivt konstruerer piecewise lineære stier. En mere teoretisk gennemgang af coordinate descent og LARS algoritmen er givet i kapitel –.

1.4 Generaliseringer af lasso

I denne sektion introduceres generaliseringen af lasso. Alle har de to essentielle egenskaber af standard lasso, nemlig shrinkage og udvælgelse af variable.

Empirisk studier viser at lasso ikke er godt til højt korreleret variable

1.4.1 Elastic net

Som nævnt er lasso ikke god til at håndtere højt korreleret variable. Dette ses ved at koefficienter stierne er uregelmæssige.

Med ved at kombinere et kvadreret ℓ_2 strafled med ℓ_1 strafled fås en metode kaldet elastiske net, som er bedre til korreleret grupper og vælger de korreleret variater (eller ikke) sammen.

Det elastiske net løser det konvekse problem

$$\min_{\beta_0, \beta} \left\{ \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - x_i^T \beta)^2 + \lambda \left[\frac{1}{2} (1 - \alpha) \|\beta\|_2^2 + \alpha \|\beta\|_1 \right] \right\}, \quad (1.6)$$

hvor $\alpha \in [0, 1]$ er en parameter som kan varieres.

Hvis $\alpha = 1$, da reduceres strafledet til ℓ_1 -normen eller strafledet for lasso og hvis $\alpha = 0$ reduceres det til den kvadrerede ℓ_2 -norm, svarende til strafledet for ridge regression.

For ethvert $\alpha < 1$ og $\lambda > 0$ da er det elastiske net problem (??) streng konveks, dvs der eksisterer en entydig løsning uanset korrelations strukturen af X_j .

Figur – viser betingelsesområdet for henholdsvis det elastiske net og standard lasso for tre variable. Heraf ses at det elastiske net deler egenskaber af ℓ_1 kuglen og ℓ_2 kuglen: de skarpe hjørner og kanter opfordre til variable udvælgelse og de kurvede konturer opfordre til deling af koefficienterne.

Det elastiske net har en ekstra tuning parameter α som skal bestemmes. I praksis kan den ses som en højere-level parameter, og kan sættes på subjektiv grunder. Alternativt, kan man inkludere en sekvens af værdier for α vha krydsvalidering.

Det elastiske net problem (??) er konveks for $(\beta_0, \beta) \in \mathbb{R} \times \mathbb{R}^p$ og vi kan anvende en række algoritmer til at løse det. Coordinate descent er særlig effektiv, og opdateringer er blot en simpel udvidelse af dem for standard lasso i -. Igen centreres kovariaterne, således at skæringen findes til sidst.

Coordinate descent opdateringen for j 'te koefficient er givet ved

$$\hat{\beta}_j = \frac{S_{\lambda\alpha} \left(\sum_{i=1}^n r_{ij} x_{ij} \right)}{\sum_{i=1}^n x_{ij}^2 + \lambda(1-\alpha)},$$

hvor $S_\mu(z) = \text{sign}(z)(z - \mu)_+$ er soft-thresholding operatoren og $r_{ij} = y_i - \hat{\beta}_0 - \sum_{k \neq j} x_{ik} \hat{\beta}_k$ er den partial residual.

1.4.2 Grouped lasso

For mange regressions problemer har kovariaterne en naturlig grupperet struktur, og da foretrækkes det at alle koefficienter indenfor en gruppe er ikke-nul (eller nul) samtidig. Betragt en lineær regressions model som har J grupper af kovariater, hvor vektoren $Z_j \in \mathbb{R}^{p_j}$ for $j = 1, \dots, J$ repræsenterer kovariaterne i gruppe j . Formålet er da at prædiktere responsvariablen $Y \in \mathbb{R}$ baseret på en samling af kovariater (Z_1, \dots, Z_J) . En lineær model for regressions funktionen $\mathbb{E}[Y|Z]$ er givet ved $\theta_0 + \sum_{j=1}^J Z_j^T \theta_j$, hvor $\theta_j \in \mathbb{R}^{p_j}$ repræsenterer en gruppe af p_j regressions koefficienter.

Given en samling af n samples $\{(y_i, z_{i,1}, z_{i,2}, \dots, z_{i,J})\}_{i=1}^n$ løser group lasso følgende konveks problem

$$\min_{\theta_0 \in \mathbb{R}, \theta_j \in \mathbb{R}^{p_j}} \left\{ \frac{1}{2} \sum_{i=1}^n \left(y_i - \theta_0 - \sum_{j=1}^J z_{ij}^T \theta_j \right)^2 + \lambda \sum_{j=1}^J \|\theta_j\|_2 \right\}, \quad (1.7)$$

hvor $\|\theta_j\|_2$ er den euklidiske norm af vektoren θ_j . Dette er en grupperet generalisering af lasso, som har følgende egenskaber:

- Afhængig af λ , vil enten alle indgange i vektoren $\hat{\theta}_j$ være nul eller ikke-nul
- Når $p_j = 1$, da har vi at $\|\theta_j\|_2 = |\theta_j|$, således at alle grupper er singletons, dermed reduceres optimerings problemet (??) til standard lasso problemet.

På figur – sammenlignes betingelsesområdet for den grupperet lasso med lasso for tre variable. Vi ser at den grupperet lasso deler egenskaber med både ℓ_1 og ℓ_2 kuglen.

I (??), straffes alle grupper ligeligt, hvilket betyder at større grupper vil have en tendens til at blive valgt.

Udregning af group lasso

Lad os omskrive optimerings problemet (??) på matrix-vektor form

$$\min_{\theta_1, \dots, \theta_J} \left\{ \frac{1}{2} \left\| \mathbf{y} - \sum_{j=1}^J \mathbf{Z}_j \theta_j \right\|_2^2 + \lambda \sum_{j=1}^J \|\theta_j\|_2 \right\}.$$

Vi ignorerer skæringen θ_0 , da vi centrerer variableerne og responsvariablen. For dette problem er nul subgradient ligningerne givet ved

$$-\mathbf{Z}_j^T \left(\mathbf{y} - \sum_{\ell=1}^J \mathbf{Z}_\ell \hat{\theta}_\ell \right) + \lambda \hat{s}_j = 0, \quad j = 1, \dots, J,$$

hvor $\hat{s}_j \in \mathbb{R}^{p_j}$ er et element af subdifferential af normen $\|\cdot\|_2$ evalueret i $\hat{\theta}_j$. Når $\hat{\theta}_j \neq 0$ da har vi, at $\hat{s}_j = \frac{\hat{\theta}_j}{\|\hat{\theta}_j\|_2}$, og når $\hat{\theta}_j = 0$ da har vi, at \hat{s}_j er enhver vektor hvor $\|\hat{s}_j\|_2 \leq 1$. En metode at løse nul subgradient ligningerne er ved at fastholde alle block vektorer $\{\hat{\theta}_k, k \neq j\}$, og da løse for $\hat{\theta}_j$. Hermed udføres block coordinate descent på objektfunktionen af group lasso. Da problemet er konveks, og straffledet kan separeres efter block, er det garanteret at konvergere til en optimal løsning. Med $\{\hat{\theta}_k, k \neq j\}$ fastholdt, kan vi skrive

$$-\mathbf{Z}_j^T (\mathbf{r}_j - \mathbf{Z}_j \hat{\theta}_j) + \lambda \hat{s}_j = 0,$$

hvor $\mathbf{r}_j = \mathbf{y} - \sum_{k \neq j} \mathbf{Z}_k \hat{\theta}_k$ er den j 'te partial residual. Fra betingelserne opfyldt af subgradienten \hat{s}_j , må vi have at $\hat{\theta}_j = 0$ hvis $\|\mathbf{Z}_j^T \mathbf{r}_j\|_2 < \lambda$, og ellers må $\hat{\theta}_j$ opfylde

$$\hat{\theta}_j = \left(\mathbf{Z}_j^T \mathbf{Z}_j + \frac{\lambda}{\|\hat{\theta}_j\|_2} \mathbf{I} \right)^{-1} \mathbf{Z}_j^T \mathbf{r}_j. \quad (1.8)$$

Denne opdatering er ens med løsningen af ridge regression, bortset fra at den underliggende straf parameter afhænger af $\|\hat{\theta}_j\|_2$. Desværre har ligning (??) ikke en lukket løsning for $\hat{\theta}_j$ medmindre at \mathbf{Z}_j er ortonormal. I dette special tilfælde har vi, at

$$\hat{\theta}_j = \left(1 - \frac{\lambda}{\|\mathbf{Z}_j^T \mathbf{r}_j\|_2} \right)_+ \mathbf{Z}_j^T \mathbf{r}_j.$$

1.4.3 Ikke-konvekse straffled

tilpasse modeller som er mere sparse end lasso Den vægtede lasso løser

$$\min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \sum_{j=1}^p w_j |\beta_j| \right\},$$

hvor $w_j = \frac{1}{|\beta_j|^v}$. Straffledet for den vægtede lasso kan ses som en approksimation til ℓ_q straffledene med $q = 1 - v$.

1.5 Generaliseret lineære modeller

Generaliseret lineære modeller er en generalisering af den sædvanlige *lineære model* som tillader at responsvariablen kan have andre fordelingen end normalfordelingen.

Der er tre hovedkomponenter i GLM:

Stokastisk komponent

Systematisk komponent Den specificerer den lineære kombination af de forklarende variable

$$\eta_i = x_i^T \beta$$

Link funktion Denne forbinder den stokastiske og systematiske komponent. Den viser hvordan den forventede værdi af responsvariablen er forbundet med den lineære prediktor af forklarende variable

$$\eta_i = g(\mu_i),$$

hvor μ er den forventede værdi for Y og g er en streng monoton link funktion.

1.5.1 Logistisk regressions model