

confidence region is $\{0\}$, $n \geq 20p$,
 $B \geq 50p$ and the error dist is
 unimodal and not highly skewed.
 (This technique may be useful
 for data snooping.)

ex) $\sigma^2 V = \sigma^2 I$ $Y = \beta_1 + \beta_2 x_2 + \epsilon = 1 + x_2 + \epsilon$
 $= 1 + x_2 + 0x_3 + 0x_4 + \epsilon$ so $p=4$,
 $B_S = \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}$ $B_E = \begin{pmatrix} \beta_3 \\ \beta_4 \end{pmatrix}$.

Test $H_0: B_S = (1)$ with $p=1$ $B=1$ $n=6$

$H_0: B_E = (0)$ PRO BRO HYPD.

Nominal 95% 5000 runs. $B=1000$, $n=100=25p$
 Apply short CI for $\hat{\beta}_1, \dots, \hat{\beta}_B$ short undercoverage

	$\beta_1 = 1$	$\beta_2 = 1$	$\beta_3 = 0$	$\beta_4 = 0$	$H_0: B_E = 0$			$H_0: B_S = 1$		
					PRO	hypo	bro	p=1, n=6	br1	
reg	.946	.950	.947	.948	.940	.941	.941	.937	.936	.937
len	.396	.399	.394	.398	2.451	2.451	2.452	2.450	2.450	2.451
vs	.948	.950	.997	.996	.991	.974	.991	.938	.939	.940
len	.395	.398	.323	.323	2.699	2.699	3.002	2.450	2.450	2.457

vs and full
 similar

vs coverage higher
 volume likely lower

vs and full similar

Zero padding makes coverage higher

but volume (length CI) smaller

$D_{vs} \geq D_{ub}$

For β_1, β_2
 pop CI length = 0.392

For confidence region "length" is
 D_{vs}^2 or D_{ub}^2 the cutoff. $\sqrt{x^2_{2, .95} - 2.45}$
 can't compare vs vs reg since vol is not given.

43) Many regression models satisfy
 $Y | X \sim N(X^T \beta)$ with $\sqrt{n}(\hat{\beta} - \beta) \rightarrow N_p(0, W)$.

eg GLMs, survival regression, MLR.

For variable selection with AIC

Nearly the same argument as that

for MLR works with parametric
bootstrap.

34) if $Y | X \sim D(\beta^T X, \theta)$.

Poisson reg: $Y | X \sim \text{Poisson}(e^{\beta^T X})$

bin reg $Y | X_i \sim \text{bin}(m_i, \frac{e^{\beta^T X_i}}{1 + e^{\beta^T X_i}})$

weird regression for survival analysis
GLMs, etc.

44) Large sample PI for Y given X

Apply short large sample $100(1-\alpha)\%$ PI

to Y_1^*, \dots, Y_n^* iid $D(\beta^T X, \theta)$

obtained from the regression

close to 1.

Coverage can be higher than the nominal coverage for 2 reasons:

i) the bootstrap data cloud is more variable than the iid data cloud

ii) zero padding $\hat{\beta}_{-j, \min, 0}$.

45) If $\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{D} N_p(0, \sigma^2 V)$

where V is diagonal,

$$\sqrt{n}(\hat{\beta}_s - \beta_s) \xrightarrow{D} N_{g_s}(0, \sigma^2 V_s)$$

for both the full model $\hat{\beta}$ and $\hat{\beta}_{MFX}$ and maybe, for $\hat{\beta}_{VS}$.

If $\theta = A \beta_s$ expect similar coverage and volume for full model and variable selection confidence regions.

46) Let $\beta^T X = \beta_s^T X_s + \beta_E^T X_E = \beta_s^T X_s$.

Let $H_0: A \beta = \beta_0 = \underline{0}$ where

$\beta_0 = (\beta_{11}, \dots, \beta_{1g})^T$, $0 \in E$ so H_0 is true.

$\beta_E =$ other $p-g$ terms

a) Suppose a nominal 95% confidence region is used and $UB = 0.96B$.

Then the prediction region and
Bickel and Ren confidence regions
contain at least 96% of the
bootstrap sample. If

$\beta_{0j}^* = \underline{0}$ for more than 4% of
the $(\hat{\beta}_{01}^*, \dots, \hat{\beta}_{0p}^*)$ then $\underline{0}$ is

in the confidence region and the
bootstrap test fails to reject H_0 .
If this occurs for each run of
the simulation (5000 runs), then
the observed coverage is 100% > 95%.

b) Now suppose $\hat{\beta}_{0j}^* = \underline{0}$ for

$j = 1, \dots, p$. Then S_T^* is singular

but $\{\underline{0}\}$ is the large sample

100(1- δ)% confidence region

for each of the 3 bootstrap methods:
prediction region, Bickel and Ren, hybrid.

Then $p_{\text{val}} = 1$ estimates the pop p-value.

For the Twin Model for forward selection,
there may be strong evidence
that X_0 is not needed in the model
given X_T is in the model (if the "100%")

assume

$\beta_0 = \tau = 0$

for
forward
selection

12/5/1)

St.P to M reg Lasso, Relaxed lasso, vs LM 82
ridge regression & elastic net

are other ways to fit $Y = X\beta + \epsilon$

They are especially useful if $X'X$ is singular, or if $(X'X)^{-1}$ is ill conditioned, eg if $p \gg n$.

2) Lasso, Relaxed lasso vs and elastic net do variable selection, some $\hat{\beta}_i = 0$ are possible.

3) Let the nontrivial predictors $U_j = (x_{1j}, \dots, x_{nj})^T$. Let the $n \times (p-1)$ matrix of standardized nontrivial predictors be $W = (w_{ij})$ where

$$\sum_{i=1}^n w_{ij} = 0 \quad \text{and} \quad \sum_{i=1}^n w_{ij}^2 = n$$

$\underbrace{\sum_{i=1}^n w_{ij} = 0}_{j\text{th standardized predictor has}}$ $\underbrace{\sum_{i=1}^n w_{ij}^2 = n}_{\hat{\sigma}_j^2 = \frac{1}{n} \sum_{i=1}^n w_{ij}^2 = 1}$

$\bar{w}_j = 0$

Then the sample correlation matrix of the nontrivial predictors is

$R_U = \frac{W^T W}{n}$. Let $\underline{z} = \underline{y} - \bar{y}$

where $\bar{y} = \bar{Y} \mathbf{1}$. The MLR methods in i) fit $\underline{z} = W\alpha + \epsilon$ and then find $\hat{\beta}$ from $\hat{\underline{\mu}} = \alpha \mathbf{1}$ and $\hat{Y} = \hat{\underline{\mu}} + \bar{Y}$.

$$4) \text{ Let } Q(\underline{\beta}) = \frac{1}{a} (\underline{z} - \underline{w}\underline{\beta})^T (\underline{z} - \underline{w}\underline{\beta}) + \frac{\lambda_{in}}{a} \sum_{j=1}^{p-1} |\beta_j|$$

where $\lambda_{in} \geq 0$, $a > 0$ and $j > 0$ are known.

Then $j=2$ gives ridge regression, $j=1$ gives lasso and $a=1, 2, \dots$ and $2n$ are common.

The residual sum of squares

$$RSS(\underline{\beta}) = (\underline{z} - \underline{w}\underline{\beta})^T (\underline{z} - \underline{w}\underline{\beta}), \text{ and}$$

$\lambda_{in} = 0$ corresponds to the OLS estimator $\hat{\underline{\beta}}_{OLS} = (\underline{w}^T \underline{w})^{-1} \underline{w}^T \underline{z}$

5) Lasso and ridge regression use a grid of m λ values $0 \leq \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_{m-1} \leq \lambda_m$.

For lasso, λ_m is the smallest value of λ such that $\hat{\underline{\beta}}_{\lambda_m} = \underline{0}$.

Hence $\hat{\beta}_i \neq 0$ for $i \leq m$.

6) The elastic net estimator minimizes the criterion

$$Q(\underline{\beta}) = RSS(\underline{\beta}) + \lambda_1 \|\underline{\beta}\|_2^2 + \lambda_2 \|\underline{\beta}\|_1$$

where $\lambda_1 = (1-\alpha)\lambda_{in}$ and $\lambda_2 = 2\alpha\lambda_{in}$

with $0 \leq \alpha \leq 1$.

7) Assume $(W^T W)^{-1}$ exists,

These estimators minimize convex Q ,
Karush-Kuhn-Tucker (KKT)
conditions for convex optimality
exist and are like normal
equations.

$$\text{KKT for lasso! } \frac{1}{n} W^T (z - W \hat{\beta}_L) + \frac{\lambda_1 n}{2n} \underline{s}_n = \underline{0}$$

$$\text{or } -W^T (z - W \hat{\beta}_L) + \frac{\lambda_1 n}{2} \underline{s}_n = \underline{0}$$

where $s_{in} \in \{-1, 1\}$ and $s_{in} = \text{sign}(\hat{\beta}_{iL})$
if $\hat{\beta}_{iL} \neq 0$, $\text{sign}(\hat{\beta}_{iL}) = \begin{cases} 1 & \beta_i > 0 \\ -1 & \beta_i < 0 \end{cases}$

So $\underline{s}_n = \underline{s}_n(\hat{\beta}_L)$ depends on $\hat{\beta}_L$. Thus

$$\hat{\beta}_L = (W^T W)^{-1} W^T z - \frac{\lambda_1 n}{2n} (W^T W)^{-1} \underline{s}_n =$$

$$\hat{\beta}_{OLS} - \frac{\lambda_1 n}{2n} (W^T W)^{-1} \underline{s}_n$$

KKT for EN: $2W^T W \hat{\beta}_{EN} - 2W^T z + 2\lambda_1 \hat{\beta}_{EN} + \lambda_2 \underline{s}_n = \underline{0}$
After algebra,

$$\hat{\beta}_{EN} = \hat{\beta}_{OLS} - n (W^T W + \lambda_1 I_{p-1})^{-1} \left[\frac{\lambda_1}{n} \hat{\beta}_{OLS} + \frac{\lambda_2}{2n} \underline{s}_n \right]$$

a) If $\frac{\lambda_1 n}{\sqrt{n}} \rightarrow 0$, then $\sqrt{n} (\hat{\beta}_A - \beta) \rightarrow N_p(0, \sigma^2 V)$

(probably
not kkt)

$$\text{identity EOP RR: } \hat{\underline{\mu}}_R = (\underline{W}^T \underline{W} + \lambda_n \underline{I}_{p-1})^{-1} \underline{W}^T \underline{z}$$

$$= (\underline{W}^T \underline{W} + \lambda_n \underline{I}_{p-1})^{-1} \underline{W}^T \underline{W} (\underline{W}^T \underline{W})^{-1} \underline{W}^T \underline{z}$$

$$= (\underline{W}^T \underline{W} + \lambda_n \underline{I}_{p-1})^{-1} \underline{W}^T \underline{W} \hat{\underline{\mu}}_{OLS} = \underline{A}_n \hat{\underline{\mu}}_{OLS} =$$

$$\left[\underline{I}_{p-1} - \lambda_n (\underline{W}^T \underline{W} + \lambda_n \underline{I}_{p-1})^{-1} \right] \hat{\underline{\mu}}_{OLS} =$$

$$\underline{B}_n \hat{\underline{\mu}}_{OLS} = \hat{\underline{\mu}}_{OLS} - \frac{\lambda_n}{n} (\underline{W}^T \underline{W} + \lambda_n \underline{I}_{p-1})^{-1} \hat{\underline{\mu}}_{OLS}$$

Since $\underline{A}_n - \underline{B}_n = \underline{O}$.

8) Assume the sample correlation matrix
(*) $\underline{R}_U = \frac{\underline{W}^T \underline{W}}{n} \xrightarrow{P} \underline{V}^{-1}$ where

$\underline{V}^{-1} = \underline{\rho}_U$ = the pop correlation matrix

of the nontrivial predictors U_i are iid from a pop.

If $\frac{\lambda_n}{n} \rightarrow 0$, then $\frac{\underline{W}^T \underline{W} + \lambda_n \underline{I}_{p-1}}{n} \xrightarrow{P} \underline{V}^{-1}$

and $n (\underline{W}^T \underline{W} + \lambda_n \underline{I}_{p-1})^{-1} \xrightarrow{P} \underline{V}$.

9) Let $\hat{\underline{\mu}}_A$ be $\hat{\underline{\mu}}_{EU}$, $\hat{\underline{\mu}}_L$ or $\hat{\underline{\mu}}_R$.

Assume $\sqrt{n} (\hat{\underline{\mu}}_{OLS} - \underline{\mu}) \xrightarrow{D} N_{p-1}(\underline{0}, \sigma^2 \underline{V})$.

a) If $\frac{\lambda_{1n}}{\sqrt{n}} \xrightarrow{P} 0$, $\sqrt{n}(\hat{\beta}_1 - \beta_1) \xrightarrow{D} N_{p_1}(0, \sigma^2 V)$ LM 84

b) If $\frac{\lambda_{1n}}{\sqrt{n}} \xrightarrow{P} \tau \geq 0$, $\lambda_{1n} \xrightarrow{P} \psi \in [0, 1]$ and $\underline{s}_n \xrightarrow{P} \underline{s} = \underline{s}_m$

then $\sqrt{n}(\hat{\beta}_1 - \beta_1) \xrightarrow{D} N_{p_1}(-\sigma^2[(1-\psi)/\tau \underline{m} + \psi \underline{s}], \sigma^2 V)$.

c) If $\frac{\lambda_{1n}}{\sqrt{n}} \xrightarrow{P} \tau \geq 0$, then

$\sqrt{n}(\hat{\beta}_R - \beta) \xrightarrow{D} N_{p_1}(-\tau \underline{V} \underline{m}, \sigma^2 V)$.

d) If $\frac{\lambda_{1n}}{\sqrt{n}} \xrightarrow{P} \tau \geq 0$ and $\underline{s}_n \xrightarrow{P} \underline{s} > \underline{s}_m$,

then $\sqrt{n}(\hat{\beta}_L - \beta) \xrightarrow{D} N_{p_1}(-\frac{\tau}{2} \underline{V} \underline{s}, \sigma^2 V)$.

10) Lasso and ridge regression are consistent estimators of $\underline{\beta}$

if $\frac{\lambda_{1n}}{n} \rightarrow 0$ as $n \rightarrow \infty$,

\sqrt{n} consistent if $\frac{\lambda_{1n}}{\sqrt{n}} \xrightarrow{P} \tau \geq 0$

and asymptotically equivalent to OLS if $\frac{\lambda_{1n}}{\sqrt{n}} \xrightarrow{P} 0$.

11) Roughly, lasso sets about one $\hat{\beta}_i = 0$ if $\frac{\lambda_{1n}}{\sqrt{n}} \xrightarrow{P} \tau > 0$.

12) Relaxed lasso is OLS applied to the x_i that have lasso $\hat{\beta}_i \neq 0$.

Let T_{min} be this set. If $P(S \subseteq T_{min}) \rightarrow 1$,
 then relaxed lasso theory is
 like that of forward selection
 except the π_j differ.

Relaxed elastic net is OLS
 applied to the X_j that have
 $EN \hat{\beta}_j \neq 0$. Expect relaxed versions
 perform better unless $(X_{T_{min}}' X_{T_{min}})^{-1}$ is ill conditioned.

13} Proof Sketch for 9}

If $\frac{\hat{\lambda}_{1n}}{\sqrt{n}} \xrightarrow{P} \gamma$ and $\frac{\hat{\lambda}_{2n}}{\sqrt{n}} \xrightarrow{P} \psi$, then

$\frac{\hat{\lambda}_1}{\sqrt{n}} \xrightarrow{P} (1-\psi)/\tau$ and $\frac{\hat{\lambda}_2}{\sqrt{n}} \xrightarrow{P} 2\psi\tau$. Then

$$\sqrt{n} \begin{pmatrix} \hat{\beta} \\ EN \hat{\beta} - \eta \end{pmatrix} = \underbrace{\sqrt{n} (\hat{\beta}_{OLS} - \eta)}_{\xrightarrow{P} N_{p-1}(0, \sigma^2 V)} - \underbrace{n (W'W + \hat{\lambda}_1 I_p)}_{\xrightarrow{P} -V} \underbrace{\begin{pmatrix} \frac{\hat{\lambda}_1}{\sqrt{n}} \hat{\beta}_{OLS} \\ \frac{\hat{\lambda}_2}{2\sqrt{n}} \underline{s} \end{pmatrix}}_{\xrightarrow{P} (1-\psi)\tau\eta + 4\tau\underline{s}}$$

$$\text{so } \sqrt{n} \begin{pmatrix} \hat{\beta}_{EN} \\ EN \hat{\beta}_{EN} - \eta \end{pmatrix} \xrightarrow{D} N_{p-1}(\underline{0}, \sigma^2 V) = \underline{V} \underline{[(1-\psi)\tau\eta + 4\tau\underline{s}]}$$

$$\sim N_{p-1} \left(\underline{V} \underline{[(1-\psi)\tau\eta + 4\tau\underline{s}]}, \sigma^2 V \right)$$

$\underline{0}$ if $\tau=0$

$$(14) \hat{m}_R = (W^T W + \lambda_{in} I_{p-1})^{-1} W^T \underline{z}$$

$W^T W$ is symmetric and square

so $W^T W \geq 0$ with eigenvalues

$$\psi_1 \geq \dots \geq \psi_p \geq 0, \text{ if } \psi_p = 0$$

then $(W^T W)^{-1}$ does not exist.

Let (ψ, \underline{g}) be an eigenvector
eigenvector pair of $W^T W = n R_0$.

$$\text{Then } [W^T W + \lambda_{in} I_{p-1}] \underline{g} =$$

$$W^T W \underline{g} + \lambda_{in} \underline{g} = \psi \underline{g} + \lambda_{in} \underline{g} =$$

$$(\psi + \lambda_{in}) \underline{g}. \text{ So } (\underbrace{\psi + \lambda_{in}}_{> 0 \text{ if } \psi_{in} > 0}, \underline{g})$$

> 0 if $\psi_{in} > 0$

is an eigenvalue eigenvector pair
of $\underbrace{W^T W + \lambda_{in} I_{p-1}}_{\text{positive definite}}$ > 0 if $\lambda_{in} > 0$

Hence $(W^T W + \lambda_{in} I_{p-1})^{-1}$ exists $\forall \lambda_{in} > 0$,
even if $W^T W$ is singular or ill conditioned.

12) Let $A = [\underline{a}_1 \dots \underline{a}_k]$. Then the vec operator stacks the columns of A so

$$\text{vec}(A) = \begin{pmatrix} \underline{a}_1 \\ \underline{a}_2 \\ \vdots \\ \underline{a}_k \end{pmatrix}.$$

13) Let $A = (a_{ij})$ and B be $p \times q$. The

$m \times n$

Kronecker product of A and B is

$$\underbrace{A \otimes B}_{m \times n \times p \times q} = \begin{pmatrix} a_{11}B & a_{12}B & \dots & a_{1n}B \\ a_{21}B & a_{22}B & \dots & a_{2n}B \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1}B & a_{m2}B & \dots & a_{mn}B \end{pmatrix}.$$

If A and B are non singular, then

$$[A \otimes B]^{-1} = A^{-1} \otimes B^{-1}.$$

$A \otimes B$ is a compact method for writing down certain block matrices.

14) The multivariate least squares central

limit theorem (MLS CLT). Suppose

$$\max_{1 \leq i \leq n} (x_{i1}, \dots, x_{in}) \rightarrow 0, \quad \frac{1}{n} X^T X \rightarrow W^{-1} \text{ and}$$

the zero mean iid ξ_i have finite 4th moments,

Then $\sqrt{n} \text{vec}(\hat{B} - B) \xrightarrow{D} N_{pm} \left(0, \underset{\substack{\uparrow \\ m \times m}}{\Sigma_\epsilon} \otimes W \right)$.

Note: LS CLT had the same assumptions except the zero mean iid ϵ_i had 2nd moments, 4th moments are likely used because a bigger theorem also showed $\hat{\Sigma}_\epsilon$ is \sqrt{n} consistent and asymptotically normal. If $m=1$,

$$\hat{\Sigma}_\epsilon = \sigma^2 \quad \text{and} \quad \hat{\Sigma}_\epsilon \otimes W = \sigma^2 W.$$

15) To test $H_0: LB=0$ vs $H_1: LB \neq 0$, let

$$W_e = \hat{E}^T \hat{E} = (n-p) \hat{\Sigma}_\epsilon, \quad \text{let}$$

$$H = \hat{B}^T L^T \left[L (X^T X)^{-1} L^T \right]^{-1} L \hat{B}, \quad \text{and let}$$

$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$ be the eigenvalues of

$$W_e^{-1} H, \quad \text{Here } H \text{ is not the hat matrix,}$$

$$\text{Wilks' } \Lambda(L) = \prod_{i=1}^m \frac{\lambda_i}{1 + \lambda_i}$$

$$\text{Pillai's } V(L) = \sum_{i=1}^m \frac{\lambda_i}{1 + \lambda_i} \quad \text{and}$$

$$\text{Hotelling Lawley } U(L) = \sum_{i=1}^m \lambda_i =$$

$$\frac{1}{n-p} [\text{vec}(L\hat{\beta})]^T \left[\hat{\Sigma}^{-1} \otimes (L(X^T X)^{-1} L^T)^{-1} \right] \text{vec}(L\hat{\beta}).$$

16) Under regularity conditions, if H_0 is true

$$-\left[\frac{n-p+1}{2} - 0.5(m-r+3) \right] \log [V(L)] \xrightarrow{D} \chi_{rm}^2$$

$$(n-p) V(L) \xrightarrow{D} \chi_{rm}^2 \quad \text{and}$$

$$(n-p) U(L) \xrightarrow{D} \chi_{rm}^2.$$

17) For $m=1$, the Seber test statistic for $H_0: L\beta = 0$ was

$$F_R = \frac{[L\hat{\beta}]^T [L(X^T X)^{-1} L^T]^{-1} [L\hat{\beta}]}{r \hat{\sigma}^2}$$

$$= \frac{(n-p) U(L)}{r} \quad \text{since } \hat{\Sigma}^{-1} = \frac{1}{\hat{\sigma}^2}.$$

So the Hotelling Lawley test statistic is the MLR partial F test statistic extended to $m > 1$ response variables.

18) Use $F_{rm, n-rm}$ approx instead of χ_{rm}^2 !

$$p_{val} = P \left[F_{r, m, n-rm} > \frac{-(n-p+1 - 0.5(m-r+3)) \log(RL)}{rm} \right]$$

$$p_{val} = P \left[F_{r, m, n-rm} > \frac{(n-p) V(L)}{rm} \right]$$

$$p_{val} = P \left[F_{r, m, n-rm} > \frac{(n-p) U(L)}{rm} \right].$$

These are large sample tests.

19) The Mahalanobis distance of a $p \times 1$ vector \underline{w}_i with respect to a multivariate location and dispersion estimator $(\underline{T}, \underline{C})$ is

$$D_i(\underline{T}, \underline{C}) = D_{\underline{w}_i}(\underline{T}, \underline{C}) = \sqrt{(\underline{w}_i - \underline{T})^T \underline{C}^{-1} (\underline{w}_i - \underline{T})}$$

The classical Mahalanobis distance uses

$(\underline{T}, \underline{C}) = (\underline{\bar{w}}, \underline{S}_{\underline{w}})$ where $\underline{\bar{w}} = \frac{1}{n} \sum_{i=1}^n \underline{w}_i$ and

the sample covariance matrix estimator

$$\underline{S}_{\underline{w}} = \frac{1}{n-1} \sum_{i=1}^n (\underline{w}_i - \underline{\bar{w}}) (\underline{w}_i - \underline{\bar{w}})^T.$$

20) Let $MD_i = D_i(\underline{\bar{w}}, \underline{S}_{\underline{w}})$ and let $RD_i = D_i(\underline{T}, \underline{C})$