

24) In high dimensions,
usually can't compute $\hat{\beta}^R$, β_{Full}
and S may not be a subset
of I_j for any $j = 1, \dots, J$.

25) Data splitting assumes
the data follows a regression model
and that the cases are ind.
Randomly divide the data into
two sets: the training set
or modeling set H has n_H cases
and the validation set V has
the remaining $n_V = n - n_H$ cases
 i_1, \dots, i_{n_V} . Use Set H to build a model I_H with
as predictors. For Set H , looking at the data
is allowed. variable selection can be

used. (On a full data set

35.5

on n cases, looking at the response invalidates hypothesis testing, and variable selection inference is difficult.)

Then fit model I_H with a predictors with the validation set. If

$n_v \geq J_a$ ($J \geq 5, 10$ etc), do

standard model checks and inference.

26] Drawbacks: i) n_v cases instead of n are used, so loss of efficiency compared to choosing model I_H without looking at the response.

ii) models that are much better than I_H may exist, especially in high dimensions

27] Data splitting is easy to do with high dimensional data, and iid cases are not needed.

28] * The R sample function (HD 36)

can be used to get a permutation

of $1, \dots, n$. Suppose a_1, \dots, a_n

represent the cases.

$n \leq 10$

sample(1:n, n) and sample(1:n) work.

see AWS

ex] $n=10$, a_1, \dots, a_{10} , $n_H=5$

($n_H \approx \frac{n}{2}$
is common)

sample(1:10)

4 1 9 6 10 | 3 5 2 8 7
└──────────┘ └──────────┘
H V

$a_4, a_1, a_9, a_6, a_{10}$ in H

could use sample(1:10, 5)

1 7 5 9 8

a_1, a_7, a_5, a_9, a_8 in H.

29] I like to use student names on
quizzes and exams.

30] §2.1 OLS MLR variable selection $n \gg p$ $\underline{y} = \underline{X}^T \underline{\beta} + \underline{\epsilon}$

Let $MSE(I) = \frac{SSE(I)}{n-k} = \frac{\sum (y_i - \hat{y}_i(I))^2}{n-k}$ where

$\underbrace{\text{contains constant}}_{\underline{\beta}_I \text{ is } k \times 1}$

Then $C_p(I) = \frac{SSE(I)}{MSE} + 2k - n$

$= (p-k)(F_I - 1) + k$ where F_I is the statistic for testing $H_0: \underline{\beta}_0 = \underline{0}$ out of the model

in the model $y = \underline{x}_I^T \underline{\beta}_I + \underline{x}_0^T \underline{\beta}_0 + e$.

Also $MSE = MSE(F) = \hat{\sigma}^2 = \text{full model MSE}$.

31] Let \underline{r} be the full model residuals

and \underline{r}_I be the submodel residuals:

$$\underline{y} = \underline{X}_I \underline{\beta}_I + \underline{e}, \quad r_{iI} = y_i - \hat{y}_{iI}, \quad i=1, \dots, n.$$

suppose every submodel I contains a constant and \underline{X} is a full rank matrix. For OLS, $=p$

$$\text{corr}(\underline{r}, \underline{r}_I) = \sqrt{\frac{SSE}{SSE(I)}} = \sqrt{\frac{n-p}{C_p(I) + n - 2k}} = \sqrt{\frac{n-p}{(p-k)F_I + n - p}}$$

$$F_I = \left[\frac{SSE(I) - SSE}{(n-k) - (n-p)} \right] / \frac{SSE}{n-p} = \frac{n-p}{p-k} \left[\frac{SSE(I)}{SSE} - 1 \right].$$

32] Consider all submodels I with $\underline{\beta}_I$ $k \times 1$, considered by the US method

The model that minimizes $C_p(I)$ maximizes $\text{corr}(\underline{r}, \underline{r}_I)$.

and $MSE(I)$

33} Of the J models considered, eg from forward selection, let I_{min} minimize $C_p(I)$.

34} When $H_0: \beta_0 = 0$ is true,

$$(p-k) F_I \xrightarrow{D} \chi^2_{p-k} \quad \text{and}$$

$$(p-k)(F_I - 1) + k = C_p(I) \xrightarrow{D} \chi^2_{p-k} + 2k - p$$

if the errors are iid with constant variance σ^2 .

35} Suppose the full model F is one of the models considered so $S \subseteq F$. Let $\underline{y} = \underline{X} \underline{\beta} + \underline{e}$,

Forward sel and backward elim

$$C_p(I_{min}) \leq C_p(F) = p \quad \text{and} \quad \text{corr}(r, r_{I_{min}}) \rightarrow 1 \quad \text{as } n \rightarrow \infty$$

If $\underline{x}^T \underline{\beta} = \underline{x}_S^T \underline{\beta}_S$ and $P(S \subseteq I_{min})$ does not go to one, then $\text{corr}(r, r_{I_{min}})$ would not go to one. Hence $P(S \subseteq I_{min}) \rightarrow 1$ as $n \rightarrow \infty$

(assuming S is unique and has minimal dimension).

ex} Assume $S = I_i$ is the model that deletes predictor x_i . Then $\underline{\beta}_{I_i}$ is $(p-1) \times 1$. Consider all subsets selection. Then $S \subseteq I_i$ and $S \subseteq F$ so only π_S and π_F are positive since $P(S \subseteq I_{min}) \rightarrow 1$ as $n \rightarrow \infty$. Since $C_p(F) = p$, $S = I_i$ is selected if $C_p(S) < p$.

By 34) $C_p(s) \xrightarrow{D} \chi_{p-k}^2 + 2k - p \underset{k=p-1}{=} \chi_1^2 + p - 2.$ (37.5)

Hence $\pi_g = P(\chi_1^2 + p - 2 < p) = P(\chi_1^2 < 2)$

and $\pi_F = 1 - \pi_g$. This result will also hold for backward elimination. For forward selection, expect $\pi_g < P(\chi_1^2 < 2)$ and $\pi_F > 1 - P(\chi_1^2 < 2)$.

(S will often not be 1 of the 3 models for forward sel.)

37.4 A Data Splitting Prediction Region

36] Suppose there is training data $\underline{x}_1, \dots, \underline{x}_n$ and test data \underline{x}_f (a future value).

A large sample $100(1-\delta)\%$ prediction region is a set A_n such that $P(\underline{x}_f \in A_n)$ is eventually bounded below by $1-\delta$ as $n \rightarrow \infty$.

α is often used

37] If \underline{x}_f has a pdf, we often want $P(\underline{x}_f \in A_n) \rightarrow 1-\delta$ as $n \rightarrow \infty$. A prediction interval $[L_n, U_n] = A_n$ is a prediction region where $p=1$ ($\underline{x}_f = X_f$).

38] If $\underline{x}_1, \dots, \underline{x}_n, \underline{x}_f$ are iid $N_p(\underline{\mu}, \Sigma)$, Σ^{-1} exists, then the large sample $100(1-\delta)\%$ classical prediction region is $\{ \underline{z} \geq D_{\frac{\delta}{2}}^2(\underline{\bar{x}}, \hat{\Sigma}) \leq \chi_{p, 1-\delta}^2 \}$.

$\hat{\Sigma} = \hat{\Sigma}_n$.

39] Let $g_n = \min(1-\delta+0.05, 1-\delta+\frac{p}{n})$, $\delta > 0.1$ (HD 38)

$$g_n = \min\left(1-\frac{\delta}{2}, 1-\delta+10\frac{\delta p}{n}\right), \delta \leq 0.1.$$

If $1-\delta < 0.999$ and $g_n \leq 1-\delta+0.001$, set $g_n = 100(1-\delta)$. Let $D_{(n)}$ be the

$100g_n$ th sample quantile of the $D_i = D_{x_i}$

where $D_{x_i}^2 = (x_i - \bar{x})^T S^{-1} (x_i - \bar{x})$, $i=1, \dots, n$.

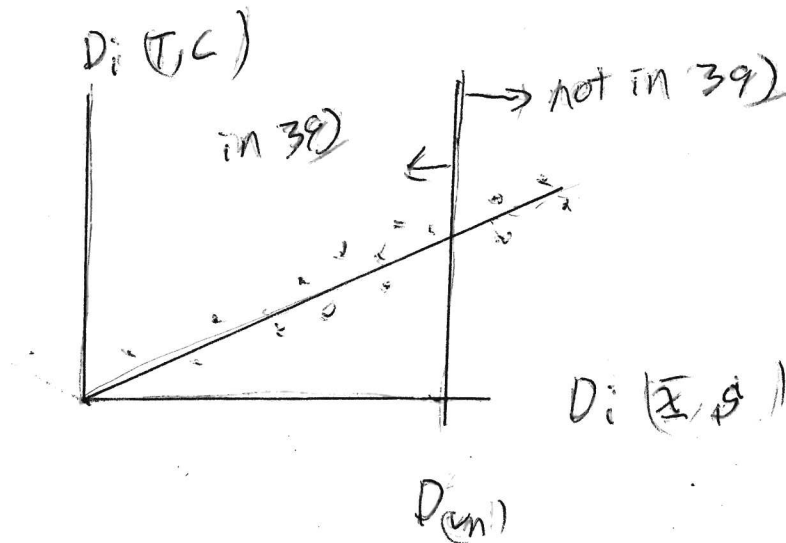
If x_1, \dots, x_n, x_f are iid and Σ_x^{-1} exists,

then the large sample $100(1-\delta)\%$ nonparametric prediction region for x_f is

$$\left\{ \underline{z} : D_{\underline{z}}^2(\bar{x}, S) \leq D_{(n)}^2 \right\}.$$

40] 39) uses $D_{(n)}^2$ while 38) uses $\chi_{p, 1-\delta}^2$.

39) contains $\approx 100g_n\%$ of the training data.



If $\delta = 0.95$, $n \leq 20p$, then $g_n = 0.975$.

Methods usually fit training data better than test data.

41) If x_1 and $x_2 = x_f$ are iid random variables, 38.5

then $P(x_1 \leq x_2) = P(x_2 \leq x_1) \geq 0.5$.

So $[x_1, \infty)$ and $(-\infty, x_1]$ are 50% prediction intervals for $x_2 = x_f$.

42) Use data splitting to divide the iid training data $\underline{x}_1, \dots, \underline{x}_n$ into 2 sets

H and V . Compute (T_H, G_H) from the n_H cases in H where G_H^{-1} exists eg

$G_H = I_p$. Then compute the squared validation

$$\text{distances } D_j^2 = D_{\underline{x}_{ij}}^2 = (\underline{x}_{ij} - T_H)^T G_H^{-1} (\underline{x}_{ij} - T_H)$$

for the cases $\underline{x}_{i_1}, \dots, \underline{x}_{i_{n_V}}$ in V . Let

$(D_{(U_V)}^2)$ be the U_V th order statistic of the D_j^2 where $U_V = \min(n_V, \lceil (n_V+1)(1-\delta) \rceil)$

where the ceiling function $\lceil x \rceil$ rounds x up to the nearest integer: $\lceil 7.7 \rceil = \lceil 8 \rceil = 8$.

Then the large sample 100(1- δ)%
data splitting prediction region for

HD 39

$$\underline{x} \text{ is } \left\{ \underline{z}: D_{\underline{z}}^2(T_H, G_H) \leq D_{(UV)}^2 \right\}$$

433 42) can be used for high dimensional data, if \underline{x}^{-1} does not exist, and for estimators with little theory such as

$$(T_H, G_H) = (T_{\text{covmb2}}, G_{\text{covmb2}}), \text{ The}$$

$$\text{"actual coverage"} \approx \frac{UV}{n_U+1} \approx \frac{\sqrt{(n_U+1)(1-\delta)}}{n_U+1}$$

n_U	δ	$UV = \min(n_U, \sqrt{(n_U+1)(1-\delta)})$	$\frac{UV}{n_U+1}$
1	1/2	$1 = \sqrt{2 \cdot \frac{1}{2}}$	$\frac{1}{2} = 0.5$
2	1/3	$2 = 3 \cdot \frac{2}{3}$	$\frac{2}{3}$
3	1/4	$3 = 4 \cdot \frac{3}{4}$	$\frac{3}{4}$
\vdots			
9	1/10	$9 = 10 \cdot \frac{9}{10}$	$\frac{9}{10}$
10	1/10	$10 = \min(10, \sqrt{11 \cdot \frac{9}{10}})$	$\frac{10}{11} \approx 0.9091$
19	1/20	$19 = 20 \cdot \frac{19}{20}$	$\frac{19}{20}$
20	1/20	$20 = \min(20, \sqrt{21 \cdot \frac{19}{20}})$	$\frac{20}{21} \approx 0.9524$
99	1/20	$95 = 100 \cdot \frac{95}{100}$	$\frac{95}{100}$
100	1/20	$96 = \sqrt{101 \cdot \frac{95}{20}}$	$\frac{96}{101} \approx 0.9505$
999	1/20	$950 = 1000 \cdot \frac{950}{1000}$	$\frac{950}{1000}$
1000	1/20	$951 = \sqrt{1001 \cdot \frac{950}{20}}$	$\frac{951}{1000} = 0.951$

N_V	$\frac{1}{20}$	UV	cov		39.5
600		571	571/600	.9501	

If $N_V \geq 19$, can get "Coverage" close to 0.95.

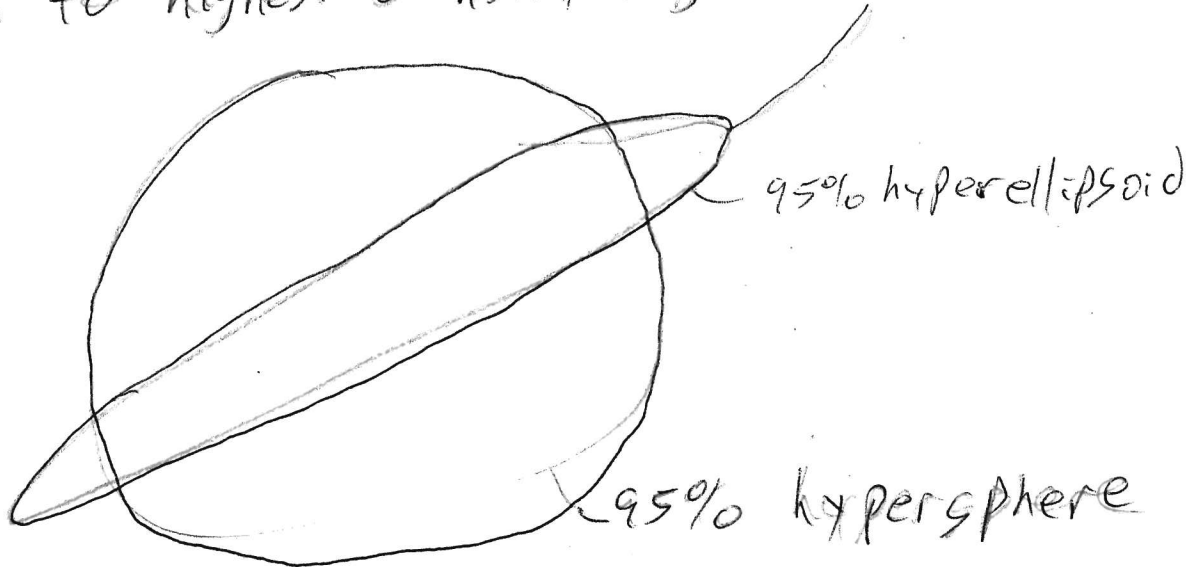
ex) $x_1 = \dots = x_n = \underline{0} = \underline{x}_f = T$, Δ nonsingular
 Coverage = 100% closed region
 $D_i^2 = 0$ for $i=1, \dots, n$ $D_{(uv)}^2 = 0$

pred region = $\{z : z = \underline{0}\}$.

44) may need $N_H \geq 20P$ to get a good estimate of Δ_H (if $\Delta_H \neq I_P$).

may need $N_V \geq 50$ for $D_{(uv)}$ to be a good estimator of the "percentile" $D_{(uv)} \approx D_{(n(1-p))}$.

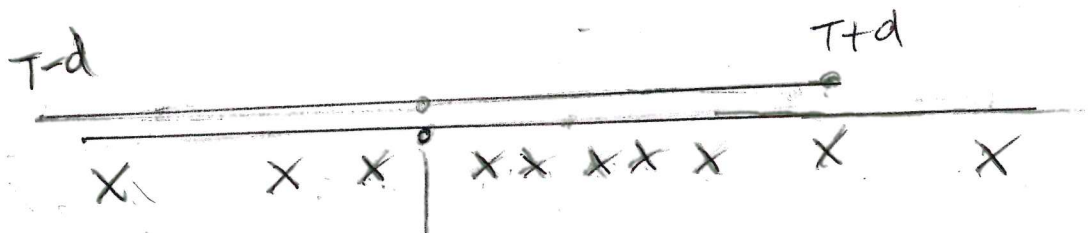
Drawback: volume of 42) may be huge compared to "highest density region volume."



45) If $P=1$, the 42) prediction interval

is the closed interval $[T-d, T+d]$ (HD 40)

where $d \geq 0$ is the smallest number such that the PI contains $\geq U_V$ cases.



$$T = \text{MED}(N_H)$$

$$N_V = 10, U_V = 9$$

(cases shown from V , not H)

46] Proof that 42] is a prediction region.

Assume $\underline{x}_1, \dots, \underline{x}_n, \underline{x}_t$ are iid where

$\underline{x}_t = \underline{x}_{n+1}$. Compute (T, d_H) from the

cases in H . Consider the squared

validation distances D_j^2 for $j=1, \dots, N_V$

and D_{n+1}^2 for case \underline{x}_t (unobserved).

Since the N_V+1 cases are iid, the

prob that D_t^2 has rank K for $K=1, \dots, N_V+1$

is $\frac{1}{N_V+1}$ for each t (i.e. the ranks

follow a discrete uniform dist).

$$t=1, \dots, N_V+1$$

Let $D_{(j)}^2$ be the order statistics (40.5)
 without using the unknown squared
 "validation distance" $D_{n_{v+1}}^2 = D_{x_f}^2$. Then

$D_{(i)}^2$ has rank i if $D_{(i)}^2 < D_{n_{v+1}}^2$

but rank $i+1$ if $D_{(i)}^2 > D_{n_{v+1}}^2$.

$$D_{n_{v+1}}^2 = D_{x_f}^2$$

take $x = n_{v+1}$

↓

	$D_{(1)}^2$	$D_{(2)}^2$...	$D_{(j)}^2$	$D_{(j+1)}^2$...	$D_{(n_v)}^2$
rank	1	2	...	j	$j+2$...	n_{v+1}

Thus $D_{(uv)}^2$ has rank $uv+1$ if $D_{x_f}^2 < D_{(uv)}^2$

$$\text{and } P\left(x_f \in \left\{z: D_z^2(T_H, Q_H) \leq D_{(uv)}^2\right\}\right)$$

$$= P\left(D_{x_f}^2 \leq D_{(uv)}^2\right) \geq \frac{uv}{n_{v+1}} \rightarrow 1-\delta \text{ as } n_v \rightarrow \infty$$

If there are no tied ranks,

$$P\left(D_{x_f}^2 \leq D_{(uv)}^2\right) = P\left(D_{x_f}^2 < D_{(uv)}^2\right) =$$