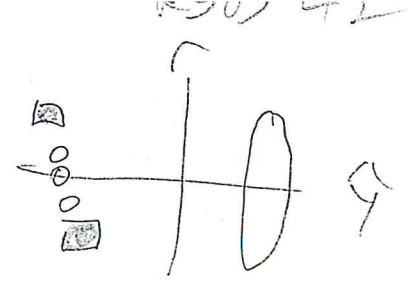
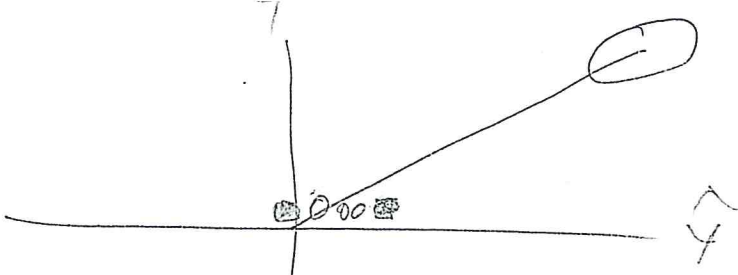
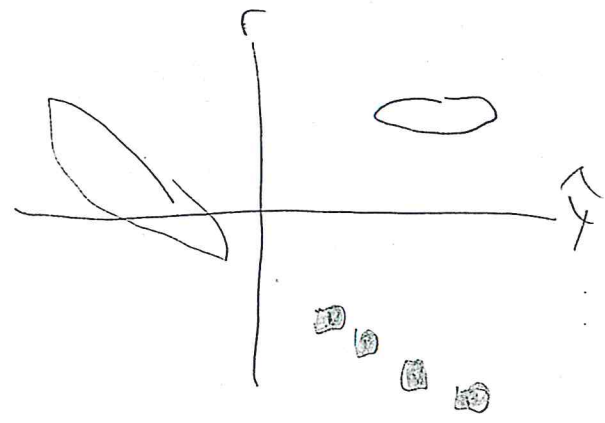
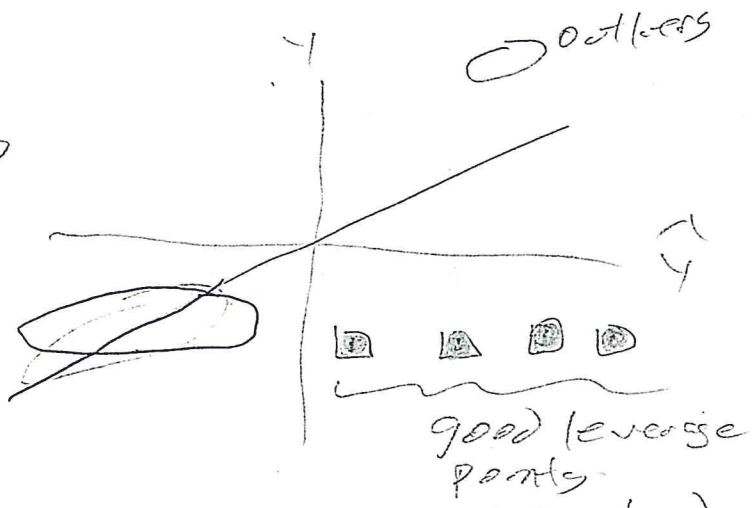


ex 3



\blacksquare = cook's distance is large so 3 outliers do not have large cook's distance. This is masking. wrt cook's distance
 None of the 3 outliers have large $|r_i|$; masking wrt residuals
 All of the outliers have (small) \hat{y}_i (the 3 outliers are correctly detected by \hat{y}).

ex 4



The good leverage points had large cook's distances but the outliers did not. masking for the outliers, swamping for the good leverage cases. wrt cook's distance.

see HW 7 5.10.

ch6 practical res. start and Robust MLR Estimators

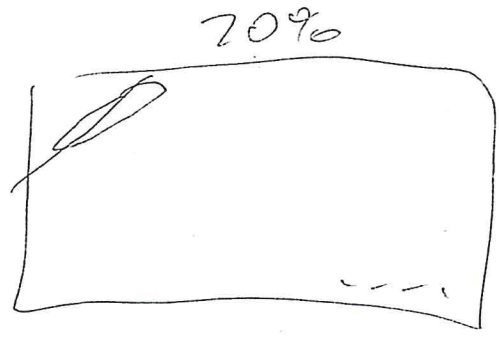
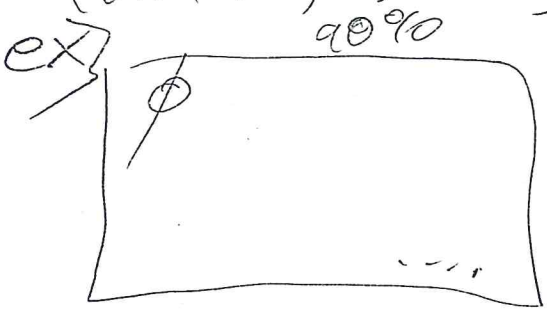
1) The MLD set MLR estimator $\hat{\beta}_D$ finds the convex set RMUN or RFCH set D applied to $U_i =$ continuous predictors. Then $\hat{\beta}_D$ is the OLS estimator MLR.

applied to the cases $(\underline{u}, \underline{u})$ corresponding to the $M \geq \frac{N}{2}$ cases in D . This estimator is \sqrt{n} consistent as long as the responses y_i were not used to select cases.

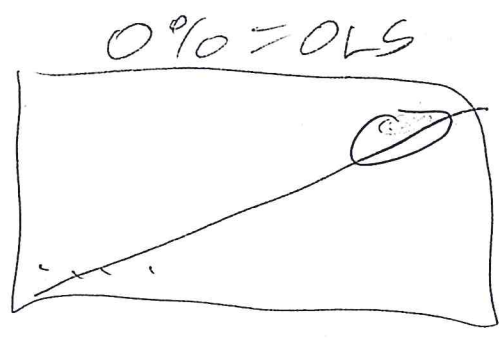
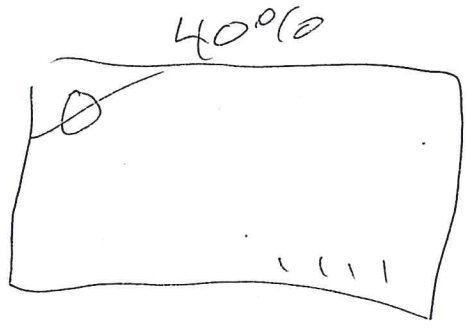
2) Compute (T, C) on the \underline{x}_i (or \underline{u}_i).

Trim the $M\%$ cases with the largest D_i^2 and compute $\hat{\beta}_M$ from the remaining cases. Use $M = 0, 10, \dots, 90$ to generate 10 response plots of $\hat{\beta}_M x_i$ vs y_i using all N cases.

The trimmed views (TV or TVreg) estimator $\hat{\beta}_{T, N}$ corresponds to the response plot (trimmed view) where the bulk of the plotted points follow the identity line with the smallest variance function, ignoring any outliers.



$\hat{\beta}_{T, N}$ uses 70% trimming.



identity line goes through the outliers bad

ex} data $(0,1), (1,2), (2,3), (3,4), (4,11)$ R58343
 lie on $Y = 1 + X$

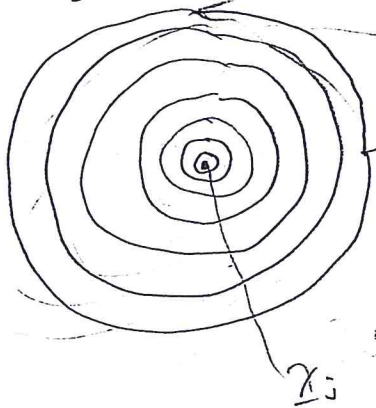
Let $J_1 = \{1,5\}$. The $p=2$ cases are $(0,1)$ and $(4,11)$.

§} If $y_i = \underline{x}_i^T \underline{\beta} + e_i$, $\hat{\underline{\beta}}$ is a resistant estimator if $\hat{\underline{\beta}}$ is known to be useful for detecting certain types of outliers.

ex} $\text{lm}(\text{reg})$ and $\text{lm}(\text{sreg})$ in R are resistant estimators

¶} Let $D_i(\underline{x}_j) = \sqrt{(\underline{x}_i - \underline{x}_j)^T (\underline{x}_i - \underline{x}_j)}$ be the Euclidean distance of \underline{x}_i from \underline{x}_j for $i=1, \dots, n$. Let $D_{(1)}(\underline{x}_j) \leq \dots \leq D_{(n)}(\underline{x}_j)$ be the order statistics.

§} For a fixed \underline{x}_j consider the cases (y_i, \underline{x}_i) corresponding to the α percentage of \underline{x}_i closest to \underline{x}_j with $\alpha \in \{1, 2.5, 5, 10, 20, 33.3, \dots\}$

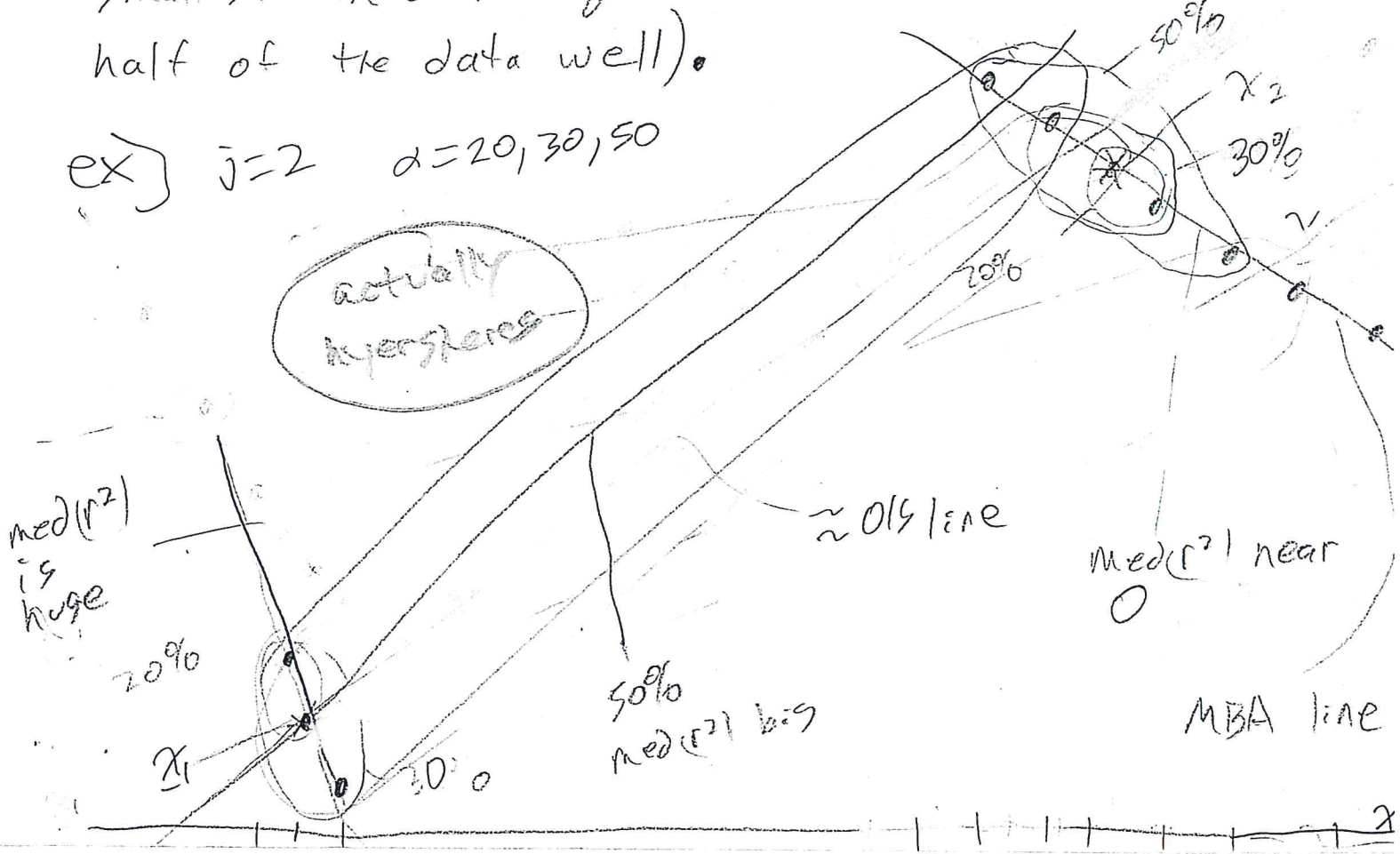


hypersphere $\left\{ \underline{z} \mid (\underline{z} - \underline{x}_j)^T (\underline{z} - \underline{x}_j) \leq D_{\left(\frac{\alpha n}{100}\right)}^2 \right\}$

if $\alpha = 50$ about 50% of the \underline{x}_i are in the sphere

$\hat{\beta}_{MBA}$ is a resistant estimator. Randomly select j \underline{x}_i from the n . Let $\hat{\beta}_{\alpha j}$ denote the OLS fit to the $\alpha\%$ cases with \underline{x}_i closest to \underline{x}_j for $\alpha \in \{1, 2.5, 5, 10, 20, 33, 50\} = A$. Also compute $\hat{\beta}$, the OLS fit to all of the data. This yields 50 OLS fits $\hat{\beta}_{\alpha j}$ for $\alpha \in A$ and $j=1, \dots, j$ and $\hat{\beta}$ ($= \hat{\beta}_{100, j} = \hat{\beta}_{100}$). Compute the median squared residual $\text{med} \left((y_i - \underline{x}_i^T \hat{\beta}_{\alpha j})^2 \right)$ for each of the 50 fits. $\hat{\beta}_{MBA}$ corresponds to the $\hat{\beta}_{\alpha j}$ with the smallest median squared residual (50 fits half of the data well).

ex) $j=2$ $\alpha=20, 30, 50$



RSB3 44

7) The Robust estimator $\hat{\beta}_R$ is the OLS

estimator applied to the RMVN set
(computed) from $\underline{u}_i = (x_{i1}, \dots, x_{ip}, y_i)^T$

where $x_{i1} \equiv 1$. $\hat{\beta}_R$ is \sqrt{n} consistent if
the \underline{u}_i come from a large class of EC dist's
which is a much stronger assumption than
iid errors e_i .

§6.2

8) An estimator $\hat{\beta}$ is high breakdown (HB) if the
median absolute or squared residual stays
bounded even if nearly half (but $< \frac{n}{2}$) of the
cases are contaminated

9) If $\hat{\beta}_0$ is high breakdown, so is the
attractor after k FLTS concentration steps.

10) If a constant is in the model
 $\hat{\beta}_M = (\overbrace{\text{MED}(n)}^{\text{of the } y_i}, 0, \dots, 0)^T$ has median absolute
residual = $\text{MAD}(n)$. Hence $\hat{\beta}_M$ is HB.

Also $\|\hat{\beta}_M\| = |\text{MAD}(n)|$ and $\text{MAD}(n)$ is HB.

11) Let $\underline{\beta}_M \geq \underline{\beta} - 1$. Let the OLS list $\underline{\beta}_O$ to the $\underline{\beta}_M$ whose γ values are closest to MEDM), let $\hat{\underline{\beta}}_B$ (LMS) be the attractor after k concentration steps. Let $\underline{\beta}_{HB} = 0.9999 \hat{\underline{\beta}}_B$. Then $\underline{\beta}_{HB}$ and $\hat{\underline{\beta}}_B$ are HB MLR estimators.

12) The hbreg estimator $\hat{\underline{\beta}}_H$ uses 3 attractors.

- i) $\hat{\underline{\beta}}_C = \hat{\underline{\beta}}_{OLS}$ a consistent estimator
- ii) $\hat{\underline{\beta}}_A = \hat{\underline{\beta}}_{MBA}$ or $\hat{\underline{\beta}}_{rmreg}$ a practical outlier resistant estimator.
- iii) $\hat{\underline{\beta}}_B$ a practical HB attractor.

Pick $a = 1.4 > 1$ and set $\hat{\underline{\beta}}_H = \hat{\underline{\beta}}_C$. If

$a Q_L(\hat{\underline{\beta}}_A) < Q_L(\hat{\underline{\beta}}_C)$, set $\hat{\underline{\beta}}_H = \hat{\underline{\beta}}_A$. If

$a Q_L(\hat{\underline{\beta}}_B) < \min\{Q_L(\hat{\underline{\beta}}_C), a Q_L(\hat{\underline{\beta}}_A)\}$, set

$\hat{\underline{\beta}}_H = \hat{\underline{\beta}}_B$. Q_L is the LMS, \underline{LTA} or LTS criterion
↑
default.

13) Let $\hat{\underline{\beta}}_L$ be the LMS, LTS or LTA estimator.

If $\hat{\underline{\beta}}_C$ and $\hat{\underline{\beta}}_L$ are both consistent estimators of $\underline{\beta}$, then $\hat{\underline{\beta}}_H$ is a HB MLR estimator asymptotically equivalent to $\hat{\underline{\beta}}_C$.

14) Typically $\hat{\underline{\beta}}_C$ and $\hat{\underline{\beta}}_L$ are both consistent estimators (for $\underline{\beta}$ if the iid ϵ_i are from a large class

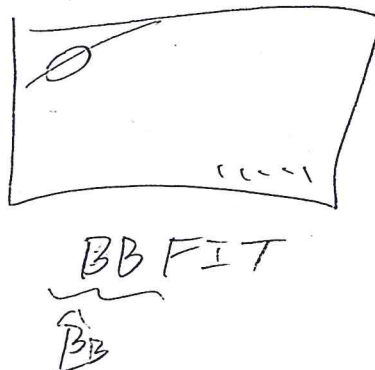
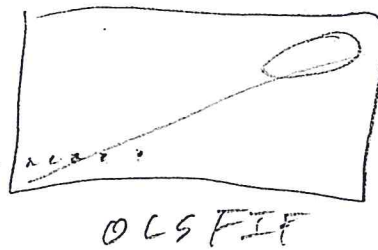
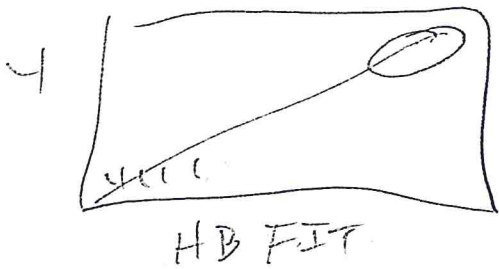
of zero mean finite variance symmetric distributions. K909 4)

If the distribution is skewed the constant β_1 estimated by $\hat{\beta}_L$ and $\hat{\beta}_C$ are different. Then $\hat{\beta}_H$ tends to select $\hat{\beta}_B$ too often.

15) $\hat{\beta}_B$ tends to be poor for e_i symmetric.

If the e_i dist is skewed, $\hat{\beta}_B$ seems to estimate the slopes ($\beta_i; i > 1$) fairly well.

16) Just as the MCD criterion sometimes selected (T_{DBH}, C_{DBH}) instead of (T_{MB}, C_{MB}) when outliers make (T_{DBH}, C_{DBH}) bad, the LTR criterion sometimes select $\hat{\beta}_C$ when the response plots look as follows:



end exam 2 material

ch 7 Variable Selection and Lasso (450)

1) Variable selection is a search for a subset of predictors that can be deleted without important loss of information if $n \geq 10p$ or so that model I is good for prediction if $n < 5p$.
 κ remaining predictors in the model

2) Let I_{min} be the model selected by the variable selection method. For forward selection and backward elimination, use the model that minimizes $C_p(I)$ if $n \geq 10p$ EBIC(I) if $n < 10p$.

3) * A model for variable selection is

$$\underline{x}^T \underline{\beta} = \underline{x}_S^T \underline{\beta}_S + \underline{x}_E^T \underline{\beta}_E = \underline{x}_S^T \underline{\beta}_S$$

where $\underline{x} = (\underline{x}_S^T \underline{x}_E^T)^T$, \underline{x}_S is $a_S \times 1$, and \underline{x}_E is $(p - a_S) \times 1$.

Given that \underline{x}_S is in the model, $\underline{\beta}_E = \underline{0}$ and E denotes the subset of terms that can be eliminated given S is in the model. (S is the model you would like to select.)

4) Let \underline{x}_I denote an $a \times 1$ vector of terms in the candidate model. If $S \subseteq I$ and \exists

holds, then $\underline{x}^T \underline{\beta} = \underline{x}_I^T \underline{\beta}_I = \underbrace{\underline{x}_I^T \underline{\beta}_I}_{\text{in model}} + \underbrace{\underline{x}_0^T \underline{\beta}_0}_{\text{out of model}} = \underline{x}_I^T \underline{\beta}_I$ RS83 46

ex) $p=4$ a constant β_1 is always in the model

$\underline{\beta} = (\beta_1, \beta_2, 0, 0)^T$. There are $J = 2^p - 1 = 8$ possible subsets of $\{1, \dots, p\}$ that always contain 1 are $I_1 = \{1\}$, $I_2 = \{1, 2\}$,

$I_3 = \{1, 3\}$, $I_4 = \{1, 4\}$, $I_5 = \{1, 2, 3\}$, $I_6 = \{1, 2, 4\}$,

$I_7 = \{1, 3, 4\}$ and $I_8 = \{1, 2, 3, 4\}$. There are

$2^{p-1} = 4$ subsets I_2, I_5, I_6 and I_8

such that $S \subseteq I_j$. Let $\hat{\underline{\beta}}_{I_j} = (\hat{\beta}_1, \hat{\beta}_3, \hat{\beta}_4)^T$ and

$$\underline{x}_{I_j} = (x_{1j}, x_{3j}, x_{4j})^T$$

*) If $\hat{\underline{\beta}}_I$ is $a \times 1$, use zero padding

to form a $p \times 1$ vector $\hat{\underline{\beta}}_{I,0}$.

ex) $p=4$ $\hat{\underline{\beta}}_{I_{j=1}} = (\hat{\beta}_1, \hat{\beta}_3)^T \Rightarrow$

$$\hat{\underline{\beta}}_{I_{j=1},0} = (\hat{\beta}_1, 0, \hat{\beta}_3, 0)^T$$

I_j	Model	x_2	x_3	x_4	x_5	$\beta_{I_j,0}$	(460)
I_2	1		*			$(\hat{\beta}_1, 0, \hat{\beta}_3, 0, 0)^T$	
I_3	2		*	*		$(\hat{\beta}_1, 0, \hat{\beta}_3, \hat{\beta}_4, 0)^T$	
I_4	3	*	*	*		$(\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3, \hat{\beta}_4, 0)^T$	
I_5	4	*	*	*	*	$(\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3, \hat{\beta}_4, \hat{\beta}_5)^T$	Full model

Assume $x_1 \equiv 1$ is always in the model and $p \geq 5$.

Given $C_p(I_j)$ if $I_{\min} = I_3$

$$\text{Then } \hat{\beta}_{\text{us}} = \hat{\beta}_{I_{\min},0} = (\hat{\beta}_1, 0, \hat{\beta}_3, \hat{\beta}_4, 0)^T$$

$$6) \quad Y = X^T \beta + e = \text{full model}$$

$$Y = X_I^T \beta_I + e = \text{submodel.}$$

The full model is a submodel.

7) Underfitting occurs if $S \not\subseteq I$ so

X_I is missing important predictors. Underfitting occurs if X_I is $a \times 1$ with $a < a_s$ where β_s is $a_s \times 1$. Overfitting occurs if

$n < sa$ or if $S \subseteq I$ but $S \neq I$.

not enough data to estimate a parameters well.

Overfitting is serious if $n < sa$, but

"not much of a problem" if $n > 10p$ or $n > 20p$.