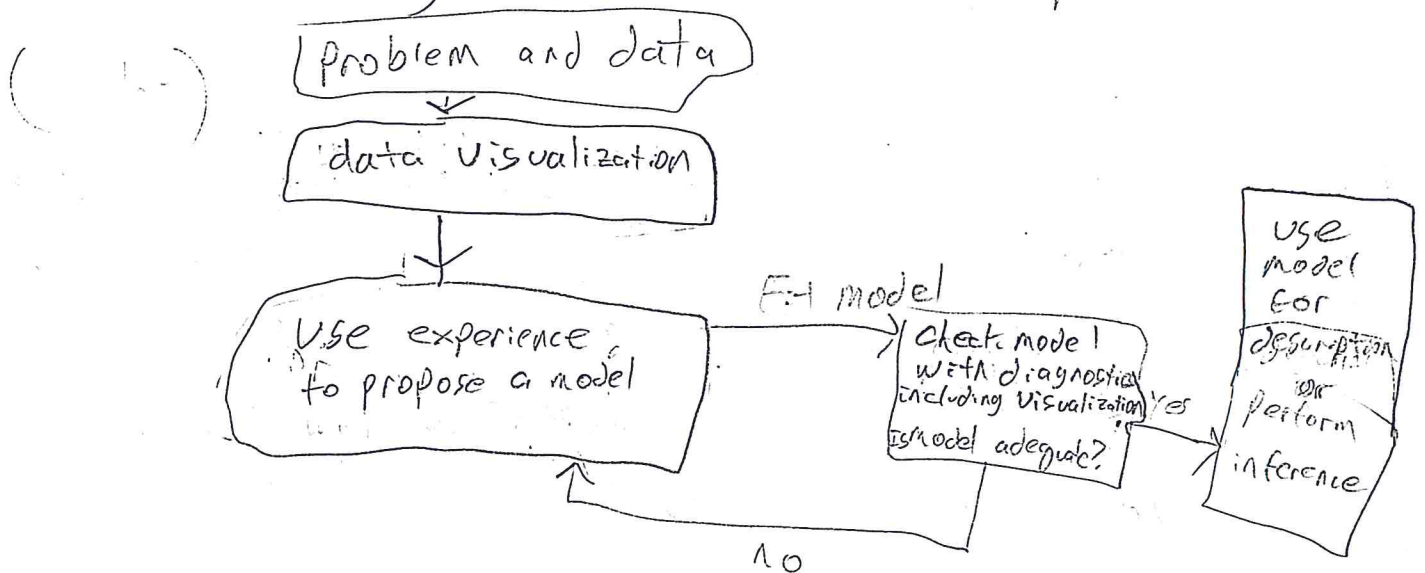1] <u>Statistics</u> is the science of obtaining useful information from data.

2] P.1 A <u>statistical model</u> is used to provide a useful approximation to the population that generated the data.

(A parametric location model)

[X] $Y_i = \mu + e_i$    $i = 1, ..., n$    where the $e_i$ are iid $N(0, \sigma^2)$. So the $Y_i$ are iid $N(\mu, \sigma^2)$. Confidence intervals for $\mu$ and hypothesis tests for $H_0: \mu = \mu_0$ vs $H_A: \mu \lessgtr \mu_0$ should be familiar from intro stat courses.

3] Model building is an iterative process.

( : )

Problem and data
↓
data visualization
↓
Use experience to propose a model → Fit model → Check model with diagnostics including visualization. Is model adequate? — yes → Use model for description or perform inference

no (loop back to propose a model)

4] P4 <u>Robust statistics</u> can give useful results when the model holds and when a certain specified model assumption is incorrect.
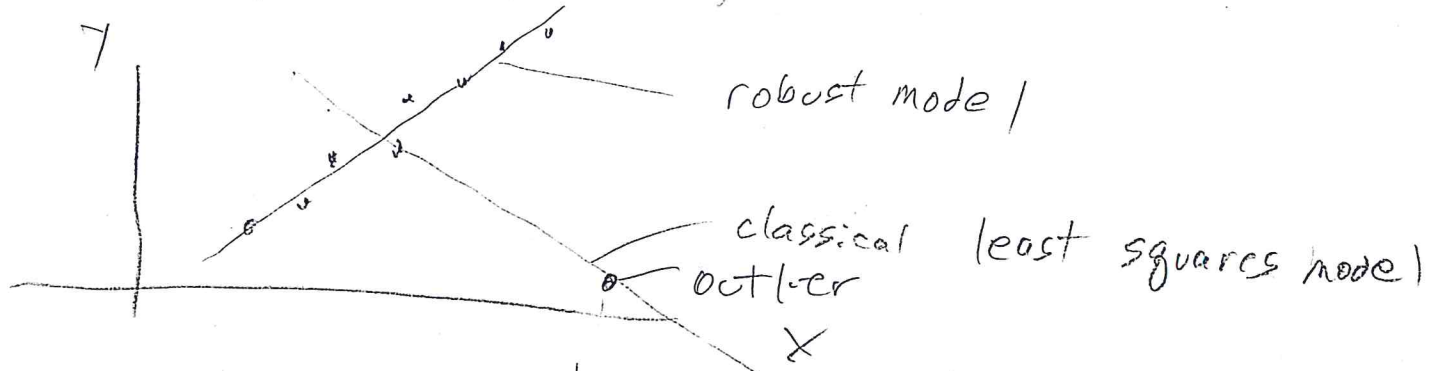
5) One assumption violation is the <u>(iv</u> presence of <u>outliers</u> : observations far from the bulk of the data. (often due to recording errors, not always bad eg spouse, good teachers doctors etc)

ex $Y = $ height $\qquad X = $ height at shoulder



we are also interested in methods that are <u>robust to the assumption of a parametric dist.</u> eg $\bar{Y} \sim N(\mu, \frac{\sigma^2}{n})$ if $Y_1, ..., Y_n$ are iid with $E(Y) = \mu$ $V(Y) = \sigma^2$.

6) PD In a <u>1D regression</u>, the response variable $Y$ that you want to predict with a vector $\underline{x} = (X_1, ..., X_p)^T$ of predictor variables is conditionally independent of $\underline{X}$ given $h(\underline{x})$ written $Y \perp\!\!\!\perp \underline{X} | h(\underline{x})$, where $h(\underline{x})$ is the sufficient predictor and $\hat{h}(\underline{x})$ is the estimated sufficient predictor. A <u>response plot</u> is a plot of ESP vs Y. $\underset{\text{ESP}}{\big\downarrow}$

7) The <u>single index model</u>

or $Y_i = m(\underline{x_i}^T \underline{\beta}) + e_i$ is a 1D model, where $e_i$ is an error, eg $e_1, ..., e_n$ are iid $N(0, \sigma^2)$.
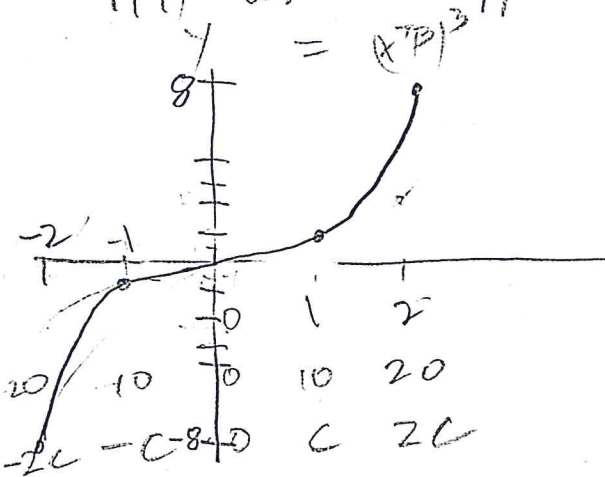
8) Another assumption violation is that $m$ is unknown or misspecified, see ch. 9

Suppose $Y = m(\underline{x}^T\underline{\beta})$

i) what happens if you plot $\underline{x}^T\underline{\beta}$ vs $y$?  (horiz axis) (vertical axis)

ii) what happens if you plot $c\underline{x}^T\underline{\beta}$ vs $y$ for $c>0$?

iii) what happens if you plot $c\underline{x}^T\underline{\beta}$ vs $y$ for $c<0$?



$y = (x^T\beta)^3$

| $\underline{x}^T\underline{\beta}$ | $y$ | $10\underline{x}^T\underline{\beta}$ |
|---|---|---|
| -2 | -8 | -20 |
| -1 | -1 | -10 |
| 0 | 0 | 0 |
| 1 | 1 | 10 |
| 2 | 8 | 20 |

$\underline{x}^T\underline{\beta} = w$

If $Y = m(\underline{x}^T\underline{\beta}) + e^i$ and $\underline{\beta}$ was known a plot of $c\underline{x}^T\underline{\beta}$ vs $y$ lets you "see $m$" up to error.

EX) $Y = \underline{x}^T\underline{\beta} + e$     let $w = \underline{x}^T\underline{\beta}$

$Y = \underbrace{0 + 1w}_{} + e$

$m$ is the line through the origin with unit slope



$Y_i - \underline{x}^T\underline{\beta} = e_i$ is the vertical deviation $\underline{x}^T\underline{\beta} =$

Note that $m(\underline{x}^T\underline{\beta}) = m\left(\dfrac{a + c\underline{x}^T\underline{\beta} - a}{c}\right) = m_{a,c}(a + c\underline{x}^T\underline{\beta})$

where $m_{a,c}(v) = m\left(\dfrac{v-a}{c}\right)$,

So a plot of $a + c\underline{x}^T\underline{\beta}$ vs $y$ "shows $m$."

9] Idea: if $\hat{\beta}$ is a good estimator of $c\beta$ for some $c \neq 0$, then a plot of $x^T\hat{\beta}$ vs $y$ is almost a plot of $cx^T\beta$ vs $y$.

10] often the least squares estimator

$$\hat{\beta} = c\beta + \underline{u}$$ where the bias vector

$\underline{u} = 0$ or is small, often $\underline{u}$ can be made small by computing least squares on a subset of the data.

11] There are "useful models" but no "true model."

With a single predictor, the model
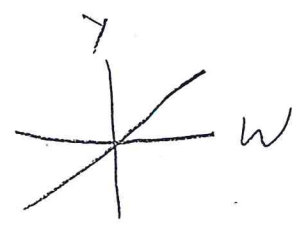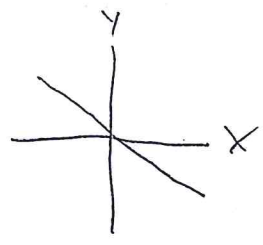$$Y = m(x) + e$$ can be visualized with a scatterplot of $x$ vs $y$. For a 1D model
$$Y = m(x^T\beta) + e,$$ $m$ can be visualized with a scatterplot of $x^T\hat{\beta}$ vs $y$ if $\hat{\beta}$ is a good estimator of $c\beta$ for $c \neq 0$.

ex) $y = -x$     so $m(x) = -x$

| | | | | |
|---|---|---|---|---|
| y | 1 | 0 | -1 | 2 |
| w=-x | 1 | 0 | -1 | 2 |
| x | -1 | 0 | 1 | 2 |

for predicting $Y$, $w$ is about as good as $x$ although the "true $m$" is flipped about the $y$ axis.

1) p19* The location model is $Y_i = \mu + e_i$, $i = 1, ..., n$.

2) p19 $\underline{know}$ $\quad$ An important robust technique for the location model is to make a plot of the data.

3) Common assumption: $e_1, ..., e_n$ are iid from a distribution with 0 mean and variance $\sigma^2$ and $n$ is large enough so that the sample mean $\bar{y} = \frac{1}{n} \sum_{i=1}^{n} Y_i \approx N(\mu, \frac{\sigma^2}{n})$, ie the central limit theorem CLT holds.
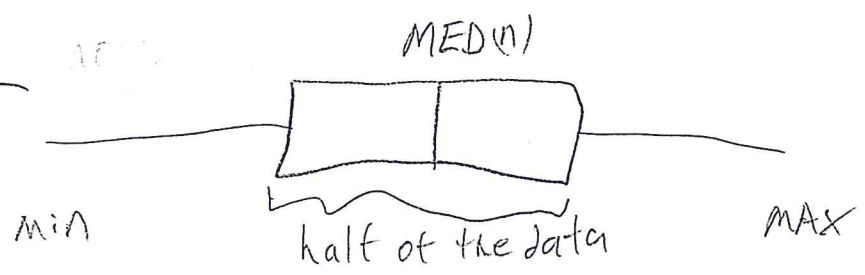
4) A $\underline{dot\ plot}$ is a plot of $i$ vs $Y_i$

A $\underline{histogram}$ tries to approximate the probability density function (pdf) $f(y)$ of a continuous random variable (RV) $Y$ and to approx the probability mass function (pmf) $P(Y=y)$ of a discrete RV $Y$.

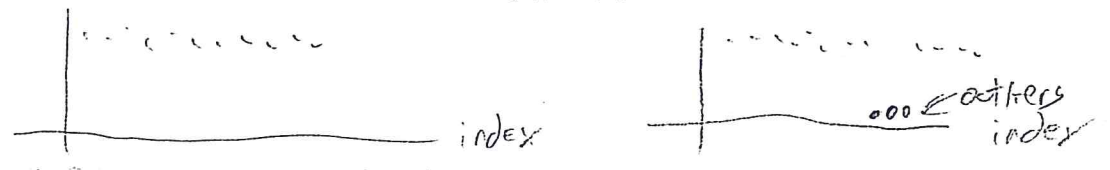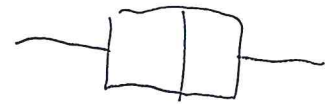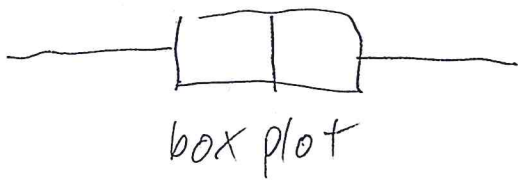A $\underline{density\ estimate}$ is a smoother approx for $f(y)$.

A typical $\underline{boxplot}$ $\quad$ summarized the dot plot.

MED(n)

MIN $\qquad$ half of the data $\qquad$ MAX

ex) Y = height
dot plot

location ....... index

.... outliers
index

histograms normal data

+ outliers

↑
outliers

density estimate (normal)

outliers

box plot

0.0
outliers

**5)** Know P20-27 the <u>sample mean</u> $\overline{Y} = \frac{1}{n}\sum_{i=1}^{n} Y_i$
If their data is $Y_{(1)}, ..., Y_n$, then the
<u>order statistics</u> are $Y_{(1)}, ..., Y_{(n)}$ where
$Y_{(1)} \leq Y_{(2)} \leq \cdots \leq Y_{(n)}$ are the $Y_i$'s written in
min ↗ ascending order. $\overset{\uparrow}{max}$

The <u>sample variance</u> $S^2 = Var(n) = \frac{\sum_{i=1}^{n}(Y_i - \overline{Y})^2}{n-1} = \frac{\sum_{i=1}^{n} Y_i^2 - n(\overline{Y})^2}{n-1}$
The <u>sample standard deviation</u> $S = \sqrt{Var(n)}$
The <u>sample median</u> $MED(n) = \begin{cases} Y_{(\frac{n+1}{2})}, & n \text{ odd} \\ \frac{Y_{(\frac{n}{2})} + Y_{(\frac{n}{2}+1)}}{2}, & n \text{ even} \end{cases}$

The <u>sample median absolute deviation</u> $MAD(n) = MED|Y_i - MED(n)|$
That is, let $D_i = |Y_i - MED(n)|$
Then $MAD(n)$ is the median of $D_1, ..., D_n$.

Review for Final.  given data compute

$\bar{Y}, S, MED(n)$ and $MAD(n)$

ex) Also see HW1 problem 2.10

Consider the data set 66, 3, 8, 5, 2.

Find a) $\bar{Y}$     b) $S$     c) $MED(n)$   d) $MAD(n)$

Soln) a) $\bar{Y} = \dfrac{\sum\limits_{i=1}^{n} Y_i}{n} = \dfrac{84}{5} = \boxed{16.8} = \dfrac{66+3+8+5+2}{5}$

b) $S^2 = \dfrac{\sum Y_i^2 - n(\bar{Y})^2}{n-1} = \dfrac{4458 - 5(16.8)^2}{4} = \dfrac{3046.8}{4}$

$$S^2 = 761.7$$

$$S = \sqrt{S^2} = \sqrt{761.7} = \boxed{27.5989 = S}$$

(Don't forget to square $\bar{Y} = 16.8$.)

c)    Sort data    2, 3, 5, 8, 66

$\boxed{MED(n) = 5}$

d)    $Y_i - MED(n)$ :   $-3, -2, 0, 3, 61$

Sort $|Y_i - MED(n)|$    0, 2, 3, 3, 61

$\boxed{MAD(n) = 3}$

So $\overline{Y}$ is the value such that the sum of the distances of the $Y_i$'s $< \overline{Y}$ $\underset{>}{=}$ "

$(Y_i$'s $= \overline{Y}$ contributed 0 to the sum)

So outliers affect $\overline{Y}$

ex   $\underbrace{0,0,0,\ 0,0,0,\ 0,0,0,}_{9\ 0\text{'s}}$   1000

$\overset{9(100)}{\downarrow}$   $\overline{\phantom{1000-\overline{Y}=900}}$ $1000-\overline{Y}=900$

|————————————————|
0   100                              1000

$\overline{Y}=1000$ is not a typical value

$\underbrace{1000-\overline{Y}}_{\text{deviations from }Y_i > \overline{Y}} = \underbrace{9\,(\overline{Y}-0)}_{Y_i\text{'s} < \overline{Y}} = 900$

9) MED($n$) is such that

at least half of the $Y_i$'s $\leq$ MED($n$)   and

at least half of the $Y_i$'s $\geq$ MED($n$)

ex   In the last ex, MED($n$)$=0$, a typical value

1, 2, 3                          1, 2, 3, 4
    $\uparrow$
MED($n$)$=2$                    MED($n$)$=2.5$

1, 2, 3, 4, 5

can replace these by any numbers greater than 3, and MED($n$) does not change

replace this by any number $x$

|  | MED($n$) |
| --- | --- |
| $-\infty < x \leq 3$ | 3 |
| $3 < x \leq 4$ | $x$ |
| $4 < x < \infty$ | 4 |

MED($n$) is the "most" outlier resistant estimator of location

(MCD max)

9) $S^2 = \frac{1}{n-1} \sum (Y_i - \bar{Y})^2 \approx$ sample mean of $(Y_i - \bar{Y})^2$

not robust! a single outlier greatly changes $S^2$
and $S$

b) $MAD(n) =$ median $|Y_i - MED(n)|$ is the "most" outlier resistant measure of spread.

10) pop quantities   $E(Y)$   $Var(Y)$   $MED(Y)$   $MAD(Y)$
$\phantom{pop quantities}$  $\phantom{E(Y)}$ || $\phantom{Var(Y)}$ ||
$\phantom{pop quantities aaaaa}$  $\mu$ $\phantom{aaaaaa}$ $\sigma^2$

sample analog   $\bar{Y}$   $S^2$   $MED(n)$   $MAD(n)$

11) *p.23 The population median $MED(Y)$

is any value such that $P(Y \leq MED(Y)) \geq 0.5$
and $P(Y \geq MED(Y)) \geq 0.5$

12) * $\overset{pop}{MAD(Y)} = MED\left(|Y - MED(Y)|\right)$

13) p 23 2   Let $f_Y(y)$ be the pdf of Y.

a) The family of pdf's $f_W(w) = f_Y(w - \mu)$

is the <u>location family</u> for $W = \mu + Y$, $\mu \in \mathbb{R}$.

b) The family of pdfs $f_W(w) = \frac{1}{\sigma} f_Y\left(\frac{w}{\sigma}\right)$

is the <u>scale family</u> for $W = \sigma Y$, $\sigma > 0$.

c) The family of pdfs $f_W(w) = \frac{1}{\sigma} f_Y\left(\frac{w-\mu}{\sigma}\right)$

is the location scale family for $W = \mu + \sigma Y$ where $\mu \in \mathbb{R}$ and $\sigma > 0$.

$$\text{note: } F_W(w) = P(W \le w) = P(\mu + \sigma Y \le w)$$

$$= P\left(Y \le \frac{w-\mu}{\sigma}\right) = F_Y\left(\frac{w-\mu}{\sigma}\right) = \tau$$

So $f_W(w) = \frac{d}{dw} F_Y\left(\frac{w-\mu}{\sigma}\right) = \frac{1}{\sigma} f_Y\left(\frac{w-\mu}{\sigma}\right).$

14] p33 know for E1    The cdf $F_Y(y) = P(Y \le y).$

a) Let $M = MED(Y)$. To find $M$, solve $F_Y(M) = 0.5.$

b) Let $D = MAD(Y)$. After finding $M$, find $D$

by solving $F_Y(M+D) - F_Y(M-D) = 0.5,$
often numerically.

c) If $W = \mu + \sigma Y$, then $MED(W) = \mu + \sigma M$

and $MAD(W) = \sigma D.$

d) If $Y$ has a pdf that is symmetric

about $\mu$, then $MED(Y) = \mu$ and

$MAD(Y) = y_{0.75} - MED(Y)$

where $P(Y \le y_\alpha) = \alpha$, ie $y_{0.75}$ is the

75th percentile of $Y$.



50%

M-D  M  M+D

50%    50%

M-D  M  M+D
    50%

Suppose $Y$ is a RV with a symmetric

pdf $f_y$ and cdf $F_y(y) = \begin{cases} \dfrac{y - \theta_1}{\theta_2 - \theta_1} & \begin{array}{l} y \leq \theta_1 \\ \text{for } \theta_1 \leq y \leq \theta_2 \end{array} \\[4pt] 1 & y \geq \theta_2 \end{cases}$

Find a) MED(Y)    b) MAD(Y)

Soln   a) $= F_y(M) = \dfrac{M - \theta_1}{\theta_2 - \theta_1} \overset{\text{set}}{=} 0.5$

or  $M = \dfrac{\theta_2 - \theta_1}{2} + \theta_1 = \dfrac{\theta_2 - \theta_1 + 2\theta_1}{2} = \dfrac{\theta_1 + \theta_2}{2}$

b) Let $U = y_{.75}$

So $F_y(U) = \dfrac{U - \theta_1}{\theta_2 - \theta_1} = 0.75$

or   $U = (\theta_2 - \theta_1)\dfrac{3}{4} + \theta_1$

and $MAD(Y) = U - M = \dfrac{3}{4}(\theta_2 - \theta_1) + \dfrac{4\theta_1}{4} + \dfrac{-2\theta_1 - 2\theta_2}{4}$

$= \dfrac{3\theta_2 - 3\theta_1 + 4\theta_1 - 2\theta_1 - 2\theta_2}{4} = \dfrac{\theta_2 - \theta_1}{4}$

15) If $Y$ is from a 2 parameter family [0⁰]

$$\mu = c_1 E(Y) \quad \text{and} \quad \gamma = c_2 VAR(Y),$$

then the method of moments estimator

is $\left(\hat{\mu} = c_1 \bar{Y}, \quad \hat{\gamma} = c_2 \frac{n-1}{n} S^2\right)$.

16) know for E1! **The MAD Method!**
P27

if $Y$ is a 2 parameter family

with $\theta = g_1(MED(Y), MAD(Y))$ and

$\quad \lambda = g_2(MED(Y), MAD(Y))$ then

$\hat{\theta} = g_1(MED(n), MAD(n))$

$\hat{\lambda} = g_2(MED(n), MAD(n))$.

ex) $Y \sim N(\mu, \sigma^2)$

$\quad MED(Y) = \mu \qquad\qquad MAD(Y) \approx 0.6745\sigma$

$\quad$ So $\hat{\mu} = MED(n)$ and $\hat{\sigma} \approx \dfrac{MAD(n)}{0.6745} \approx 1.483 MAD(n)$.

ex) $Y \sim C(\mu, \sigma)$ $\qquad$ Cauchy $E(Y)$ and $Var(Y)$ do
$\quad$ not exist, but $MED(Y) = \mu$ and $MAD(Y) = \sigma$
$\qquad$ So $\hat{\mu} = MED(n)$ and $\hat{\sigma} = MAD(n)$.

17) 'Consider intervals' that contain $c$ cases $[Y_{(1)}, Y_{(c)}]$, $[Y_{(2)}, Y_{(c+1)}]$, ..., $[Y_{(n-c+1)}, Y_{(n)}]$. Compute $Y_{(c)} - Y_{(1)}$, $Y_{(c+1)} - Y_{(2)}$, ..., $Y_{(n)} - Y_{(n-c+1)}$. Then $shorth(c) = [Y_{(s)}, Y_{(s+c-1)}]$ is the closed interval with the shortest length.

ex} know for E1

Let $c = 4$.  Data below has $n = 7$.

$$0, 1, \quad 3, \qquad 6, \qquad 9, 10, 11$$

Intervals containing $c=4$ cases {

$6 = 6 - 0$

$8 = 9 - 1$

$7 = 10 - 3$

$5 = 11 - 6$

↑

5 is shortest length

$\boxed{[6, 11] = shorth(4)}$

18) The highest $100(1-\delta)\%$ density region of a pdf is found by moving a horizontal line down from the top of the pdf so that the line intersects the pdf at one or more intervals and the sum of the areas under the pdf corresponding to the intervals $= (1-\delta)$.  The pdf can't have a positive flat interval eg $u(a,b)$.



unimodal

2.5%   95%   2.5%

highest 95% region

bimodal

2.5%  40%  15%  40%  2.5%

$I_1$   $I_2$   $I_1 \cup I_2 =$ highest 80% region

estimates the highest density $100(1-\delta)\%$ region if that region is an interval. Then the shorth(c) estimator can be used as a $100(1-\delta)\%$ large sample prediction interval (PI).

If $Y_1, \ldots, Y_n, Y_f$ are iid where $Y_f$ is a future observation, $[L_n, U_n]$ is a $100(1-\delta)\%$ large sample PI if $P[Y_f \in [L_n, U_n]]$ is eventually bounded below by $1-\delta$ as $n \to \infty$. Often want $P(Y_f \in [L_n, U_n]) \to 1-\delta$ as $n \to \infty$. The most used statistical PI's assume $Y_i \sim N(\mu, \sigma^2)$ and have coverage much smaller than $100(1-\delta)\%$ eg 95% of 50%. Let $k_1 = \lceil \frac{n\delta}{2} \rceil$, $k_2 = \lceil n(1-\frac{\delta}{2}) \rceil$. $[Y_{(k_1)}, Y_{(k_2)}]$ is the nonparametric PI.

**24] Concentration:** Start with estimator $T_0$. (ch9, ch10 b4?)

move Find the "half set" of cases closest to $T_0$ $Y_i$ such that $|Y_i - T_0| \leq MED|Y_i - T_0|$. Let $(T_1, S_1^2)$ the sample mean and variance of these cases. $(T_1, \cdot)$.

Iterate to obtain $(T_1, S_1^2), \ldots, (T_k, S_k^2)$, could iterate until convergence. If $k$ is fixed, eg $k=10$, don't need to compute $S_i^2$.

Fact: $\leq S_i^2 \leq S_{i-1}^2$ and convergence occurs when $S_i^2 = S_{i-1}^2$.

The MB estimator uses $T_0 = MED(n)$.

The DGK estimator uses $T_0 = \bar{Y}$

When iterated to convergence, the halfset seems to be estimating the same thing as the shorth ($c \approx \frac{n}{2}$) estimator.