25)

take the sample mean of each variable and collect in a vector.

know for E2

| | length $X_1$ | Nasal ht $X_2$ | bizonal $X_3$ | |
|---|---|---|---|---|
| $\Sigma X$ | 168 | 51 | 104 | 1st person |
| | 176 | 55 | 106 | |
| | 179 | 55 | 104 | $n = 5$ |
| | 176 | 46 | 114 | |
| | 177 | 48 | 101 | 5th person |

see HW3

$\Sigma X_{1i} = 876 \quad \Sigma X_{2i} = 255 \quad \Sigma X_{3i} = 529$

$$\bar{X} = \frac{1}{5}\begin{pmatrix} 876 \\ 255 \\ 529 \end{pmatrix} = \begin{pmatrix} 175.2 \\ 51.0 \\ 105.8 \end{pmatrix} \leftarrow \text{sample mean of length}$$

know for E2

ex) Let T be the coordinate wise median, MED($W$).

$$T = \begin{pmatrix} 176 \\ 51 \\ 104 \end{pmatrix}.$$

168 176 176 177 179

46  48  51  55  55

101  104  104  106  114

26) The R function cov.mcd gives estimator $(\hat{\mu}, \hat{\Sigma})$ that is often useful when outliers are present but RFCH and RMVN are faster and backed by theory,

§13.6

27) PIO7 The jth start $(T_{-1,j}, C_{-1,j})$ is an initial MLD

estimator. Then $(T_{0,j}, C_{0,j}) = (\bar{X}_{0,j}, S_{0,j})$ is the

classical estimator computed from the $C_n \simeq \frac{n}{2}$ cases

with the smallest $D_i (T_{-1,j}, C_{-1,j})$. Repeat the iteration

for $t_i$ steps resulting in the sequence of estimators

$(T_{-1,j}, C_{-1,j}) (T_{0,j}, C_{0,j}), \ldots, (T_{t_i,j}, C_{t_i,j}) = (T_{k,j}, C_{k,j}) = (\bar{X}_{k,j}, S_{k,j})$

the jth attractor. The __concentration estimator__

$(T_A, C_A)$ is the attractor used to obtain the

final estimators, $j = 1, \ldots, E$.

28) P107 $\det(C_{t+1,j}) \leq \det(C_{t,j})$ with equality

iff $(T_{t,j}, C_{t,j}) = (T_{t+1,j}, C_{t+1,j})$ for $t \geq 0$. So

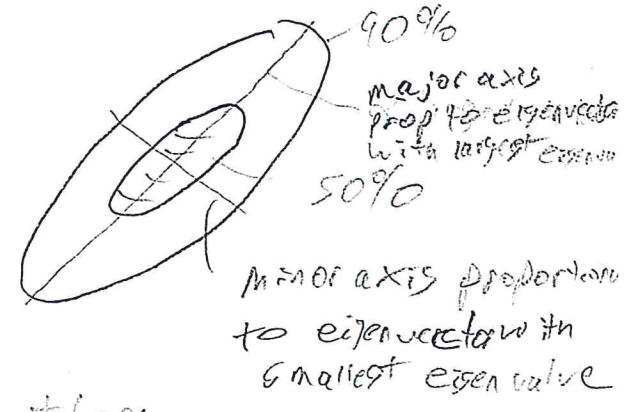the determinant decreases until convergence.

29) P100 The volume of the hyperellipsoid

$$\left\{ z : (z - \bar{x}_{k,j})^T S_{k,j}^{-1} (z - \bar{x}_{k,j}) \leq h^2 \right\} =$$

$$\underbrace{\frac{2\pi^{p/2}}{p \, \Gamma(\frac{p}{2})}}_{\text{constant}} h^p \sqrt{\det(S_{k,j})}.$$ So small volume goes with

small determinant.

30) For EC distributions, the regions of highest density

are $\left\{ z : (z - \mu)^T \Sigma^{-1} (z - \mu) \leq D^2_{1-\alpha}(\mu, \Sigma) \right\}$ are hyperellipsoid

of $\not\equiv$. If $(T, C)$ is a consistent estimator of $(\mu, s\not\equiv)$ for some $s > 0$, then

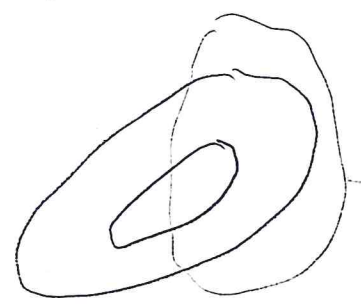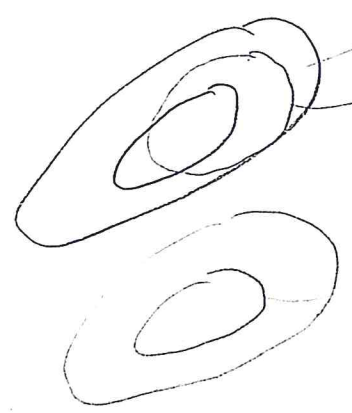$$\{i : D_i^2 (T, C) \leq MED(D_i^2 (T, C))\}$$ estimates the highest 50% density region.

31] Using the concentrated attractor iterated to convergence seems to estimate a highest density region for EC data + outliers if the attractor region does not contain any outliers. (If there are 40% outliers, then the highest density region containing $\frac{5}{6}$ of the clean data is estimated.)

$$\frac{5}{6} 60\% = 50\%$$



90%

major axis prop to eigenvector with largest eigenvalue

50%

minor axis proportional to eigenvector with smallest eigenvalue

Step 0 Scale ellipsoid) to contain half the data take classical estimator

cover half the data

Step 1 new region is a lot better

at convergence, region for attractor seems pretty good)

32] $\not\equiv$ determines the shape of the highest density region which is a hyper ellipsoid for EC distributions. (with $g\downarrow$),

w) An elemental start for MLD uses $p+1$ randomly selected cases $J_j$. Then

$$(T_{oj}, C_{oj}) = (\bar{X}_{J_j}, S_{J_j}) = \text{classical estimator}$$

applied to the $p+1$ cases. COV.MCD = FMCD uses 500 randomly chosen elemental starts. For each of the $K \div 500$ attractors, find $\det(C_{1,K}), \ldots, \det(C_{500,K})$. Suppose $J = M$ corresponds to the attractor with the minimum determinant. Then $(T_{FMCD}, C_{FMCD})$ uses the attractor $(T_{MK}, C_{MK})$.

B4) The DGK estimator is very simple. Let $(T_{01}, C_{01}) = (\bar{X}, S) = $ classical estimator

be the start. Then $(T_{DGK}, C_{DGK})$ is the attractor. Using $K = 10$ concentration steps works well. The DGK estimator has much more outlier resistance than the classical estimator.
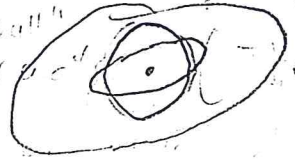
B5) PIID The MB (median ball) estimator

$$(T_0, C_0) = (MED(W), I)$$ as the start where MED(W)

is the coordinatewise median. Hence $(T_0, C_0)$ is the classical estimator applied to the $C_n \approx \frac{n}{2}$ cases closest to MED(W) in Euclidean distance, use $K = 5$ concentration steps. This estimator is a highly outlier resistant. can also be old

**36)** the FCH estimator uses the MB attractor if $T_{DGK}$ has a greater Euclidean distance from $MED(W)$ than half the data (so $T_{DGK}$ is outside of the median ball that contains half of the data). Otherwise FCH MB uses the DGK or MB estimator that has the smallest determinant $(T_A, C_A)$. Then $T_{FCH} = T_A$ and

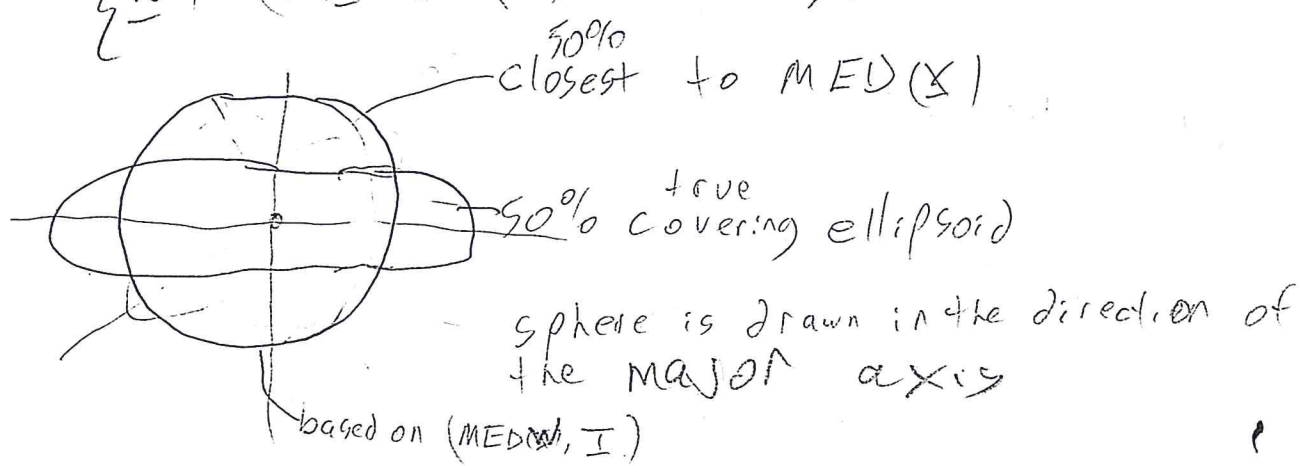$$C_{FCH} = \frac{MED(D_i^2(T_A, C_A))}{\chi^2_{p,.5}} C_A.$$

[DGK if det DGK = determi ... otherwise MB]

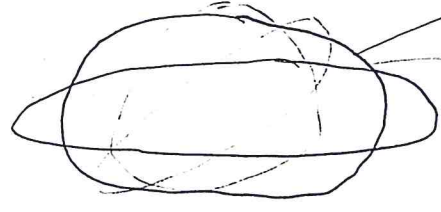**37)** On a large class of EC distributions, the prob that $(T_A, C_A) = (T_{DGK}, C_{DGK})$ goes to one as $n \to \infty$, and $(T_{FCH}, C_{FCH})$ is a $\sqrt{n}$ consistent estimator of $(\mu, c\Sigma)$ where $c > 0$ and $c = 1$ for MVN data.

**38)** The median ball attractor results in ellipsoids $\{\underline{x} | (\underline{x} - \underline{T}_{MB})^T C^{-1}_{MB} (\underline{x} - \underline{T}_{MB}) \le d^2\}$ that are "too short" in the major axis and "too fat" in the minor axis compared to $\{\underline{x} | (x - \underline{\mu}) \Sigma^{-1} (x - \mu) \le d^2\}$.

50%
closest to $MED(\underline{X})$

true
50% covering ellipsoid

sphere is drawn in the direction of the major axis

based on $(MED_W, I)$

$(\overline{x}_{OLD}, C_{OLD})$

too fat on the minor axis
too short on the major axis

based on

$$(\overline{x}_{FCH}, C_{FCH}) \overset{?}{=} (\overline{x}_{MB}, C_{MB})$$

better but still biased, if
(k concentration steps are used, ...)

(It is not known whether MB is biased or not) if concentration
2) is iterated to convergence.)

... the ... estimator is ... ,

By 19... , high ... is ... ,

39) When outliers are present that the DGK estimator
can't detect, usually $\det(C_{MB}) < \det(C_{DGK})$ so
the FCH estimator is highly outlier resistant,

40) Even if the 50% of cases with the smallest distances
based on the start contains outliers, the attractor
may use a half set without outliers,

44) Idea

$\overline{x}_{J0}$

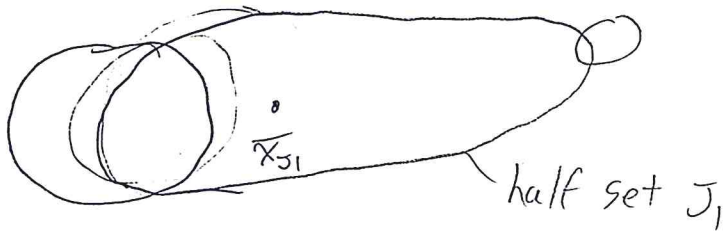Outliers

bulk of data
= clean data

half set $J_0$ based on start

If the number of outliers in the
half set $J < 25\% n$, the $\overline{x}_J$ is closer to the bulk

of the data than to the outliers. After a concentration step, fewer outliers and more clean cases will be used ( if the outliers are far enough away). After K concentration steps, the half set may be clean.



$\bar{x}_{J_1}$

half set $J_1$



$\bar{x}_{J_k}$

half set $J_k$ contains no outliers.

clean data " won the tug of war! HW 7C illustrates this phenomenon.

42) The RFCH estimator can give good results even when nearly 50% of the cases are outliers.

43) A scatterplot matrix of the distances from cov.mcd, FCH, DGk and the median ball estimator can be useful. HW 7B makes DD plots based on these 4 estimators.

44) P1V6 RFCH and RMVN are reweighted versions of FCH. They are $\sqrt{n}$ consistent estimators of $(\mu, c_i \Sigma)$ where $c_i = 1$ for gMVN data, $c_i = c_{RFCH}$, $c_2 = c_{RMVN}$.