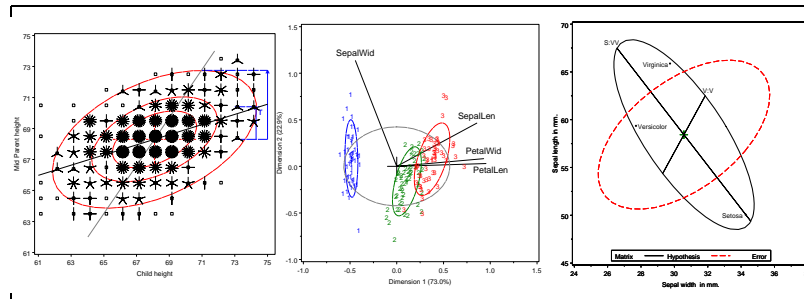


## Data ellipses, HE plots and Reduced-rank displays for MLMs: SAS software & examples



Michael Friendly

York University, <friendly@yorku.ca>

November, 2006

## Introduction: The LM family and friends

### ■ LM family

- Classical univariate models:  $y = X\beta + \epsilon$ , with  $\epsilon \sim \mathcal{N}(0, \sigma^2 I)$ ,
- $\rightarrow$  regression, ANOVA, ANCOVA, response surface models, ...
- Many graphical methods: effect plots, spread-level plots, influence plots, ...

### ■ GLM family

- Generalized univariate models:  $g(y) = X\beta + \epsilon$ , with general variance function
- $\rightarrow$  poisson regression, logistic regression, loglinear models
- Some graphical methods: effect plots, fourfold plots, mosaic plots, half-normal plots, ...

### ■ MLM family

- Classical multivariate models:  $Y = XB + E$ , with  $E \sim \mathcal{N}(0, I \otimes \Sigma)$ ,
- $\rightarrow$  MANOVA, MMRA, MANCOVA, ...
- Graphical methods: ???

### ■ MGLM family: ??

## Outline

### ■ Introduction

### ■ Bivariate plots for multivariate data

- Data ellipse
- Simpson's paradox: marginal vs. conditional relations
- Partial plots
- Robust data ellipse and outlier detection

### ■ Biplots: reduced-rank displays

### ■ HE plots

- HE plot matrices
- Showing contrasts
- HE plots for MMRA

### ■ Canonical discriminant plots: Reduced-rank HE plots

## SAS macros

Available at <http://www.math.yorku.ca/SCS/sasmac/>:

<code>biplot</code>	Generalized biplot display of variables and observations
<code>canplot</code>	Canonical discriminant structure plots
<code>ellipses</code>	Plot data ellipses (renamed from <code>contour</code> macro)
<code>heplot</code>	Plot H and E matrices for a bivariate MANOVA effect
<code>hemat</code>	HE plots for all pairs of response variables
<code>hemreg</code>	Extract H and E matrices for multivariate regression
<code>panels</code>	Display a set of plots in a rectangular layout
<code>outlier</code>	Robust multivariate outlier detection
<code>robcov</code>	Calculate robust covariance matrix via MCD or MVE
<code>scatmat</code>	Scatterplot matrices

Macros and named examples: <http://euclid.psych.yorku.ca/SCS/Papers/Private/hesoft-sas.zip>

### The Data Ellipse

#### Visual summary for bivariate marginal relations

- Shows: means, standard deviations, correlation, regression line(s)
- Defined: set of points whose squared Mahalanobis distance  $\leq c^2$ ,

$$D^2(\mathbf{y}) \equiv (\mathbf{y} - \bar{\mathbf{y}})^T \mathbf{S}^{-1} (\mathbf{y} - \bar{\mathbf{y}}) \leq c^2$$

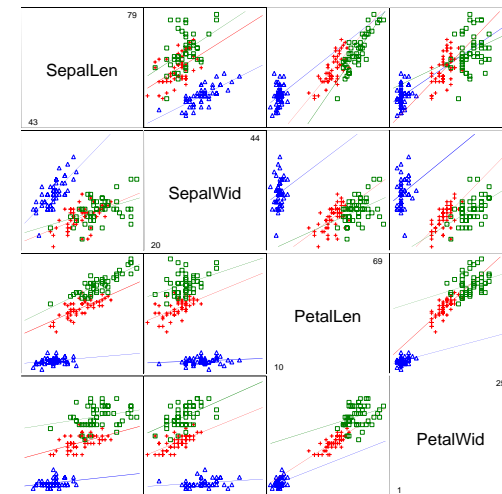
$\mathbf{S}$  = sample variance-covariance matrix

- Radius: when  $\mathbf{y}$  is approx. bivariate normal,  $D^2(\mathbf{y})$  has a large-sample  $\chi_2^2$  distribution with 2 degrees of freedom.
  - $c^2 = \chi_2^2(0.40) \approx 1$ : 1 std. dev univariate ellipse– 1D shadows:  $\bar{y} \pm 1s$
  - $c^2 = \chi_2^2(0.68) \approx 2.28$ : 1 std. dev bivariate ellipse
  - $c^2 = \chi_2^2(0.95) \approx 6$ : 95% data ellipse, 1D shadows: Scheffé intervals
  - Small samples:  $c^2 \approx 2F_{2,n-2}(1 - \alpha)$
- Construction: Transform the unit circle,  $\mathbf{U} = (\sin \theta, \cos \theta)$ ,

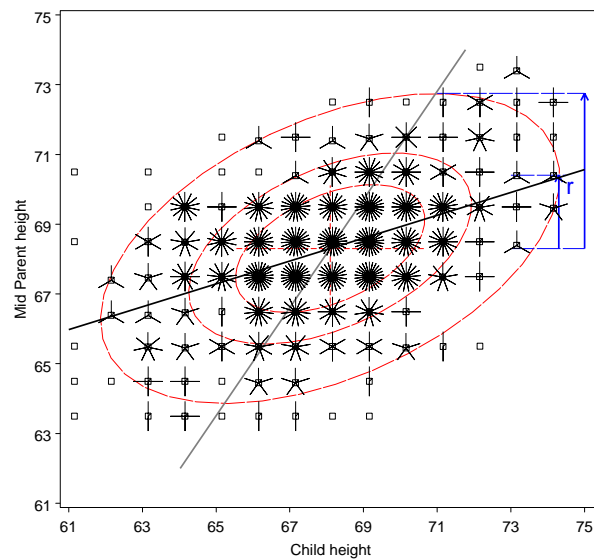
$$\mathcal{E}_c = \bar{\mathbf{y}} + c\mathbf{S}^{1/2}\mathbf{U}$$

$\mathbf{S}^{1/2}$  = any “square root” of  $\mathbf{S}$  (e.g., Cholesky)

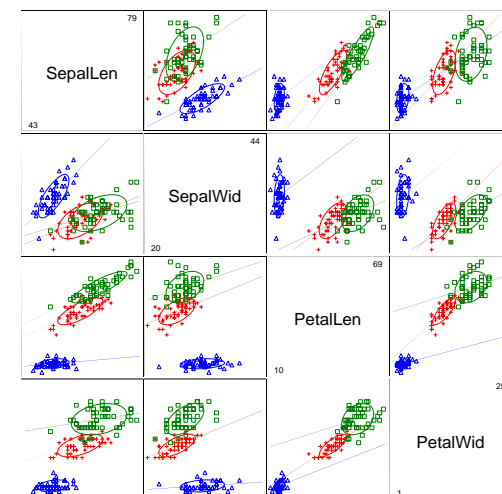
### Iris data: scatterplot matrix



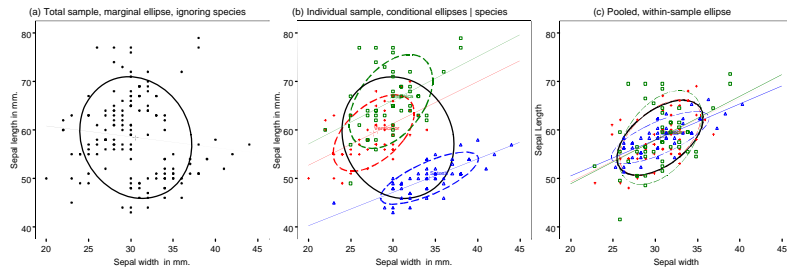
### Galton's data: Data Ellipses



### Iris data: scatterplot matrix + data ellipses

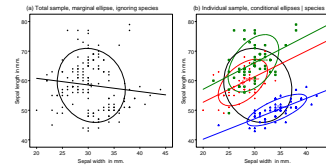


### Didactic displays: Simpson's paradox



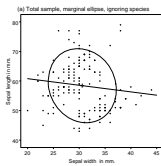
- Marginal relations: ignore other factors, covariates (species)
- Conditional relations: control or adjust for other factors, covariates
- Simpson's paradox: when these differ in direction.
- Pooled within-sample scatter:  $S_{\text{pooled}} = (N - g)^{-1} \sum_{i=1}^g (n_i - 1) S_i$
- Visual assessment of equal covariance matrices,  $\Sigma_1 = \Sigma_2 = \dots = \Sigma_g$

### Using the ellipses macro



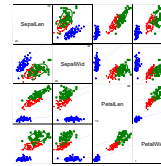
```
[... contiris3.sas ...]
1 title '(b) Individual sample, conditional ellipses | species';
2 %ellipses(data=iris,
3 x=SepalWid, y=SepalLen,
4 group=species, /* grouping variable */
5 colors=blue red green black,
6 interp=rl,
7 all=YES, /* include total sample ellipse */
8 symbols=triangle plus square,
9 line = 5 5 5 1, /* line styles */
10 width= 3 3 3 1, /* line widths */
11 haxis=axis2,
12 pvalue=.68);
```

### Using the ellipses macro



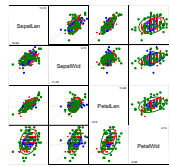
```
[contiris3.sas ...]
1 title '(a) Total sample, marginal ellipse, ignoring species';
2 axis2 order=(20 to 45 by 5) offset=(3);
3 %ellipses(data=iris,
4 x=SepalWid, y=SepalLen,
5 colors=black,
6 interp=rl, /* add regression line */
7 width=3, /* line width */
8 symbols=dot, /* plotting symbol */
9 haxis=axis2, /* horizontal axis */
10 pvalue=.68); /* ellipse size(s) */
```

### Using the scatmat macro



```
[scatirisd.sas ...]
1 %include data(iris);
2 %include macros(scatmat); /*-- or place scatmat.sas in SASAUTOS;
3 %scatmat(data=iris,
4 var=SepalLen SepalWid PetalLen PetalWid,
5 group=Species,
6 symbols=triangle plus square,
7 colors= blue red green,
8 hsym=4, htitle=9,
9 interp=rl, /* draw regression lines */
10 anno=ellipse); /* and add data ellipses */
```

### Partial plots



- View covariation after some effects “adjusted for,”  $\text{VAR}(\mathbf{Y} | \mathbf{X}) = \text{VAR}(\mathbf{U})$ .
- model SepalLen--PetalWid = Species → within-cell error ( $\mathbf{E}$  matrix)

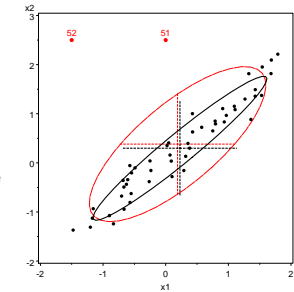
[... scatirisd.sas]

```
*-- remove group means to view within-cell relations;
proc glm data=iris noprint;
  class species;
  model SepalLen SepalWid PetalLen PetalWid = Species /nouni;
  output out=resids
         r=seplen sepwid petlen petwid;
%scatmat(data=resids,
  var=SepLen SepWid PetLen PetWid, group=Species,
  names=SepalLen SepalWid PetalLen PetalWid,
  symbols=triangle plus square,
  colors= blue red green,
  hsym=4, httitle=9,
  interp=rl,
  anno=ellipse);
```

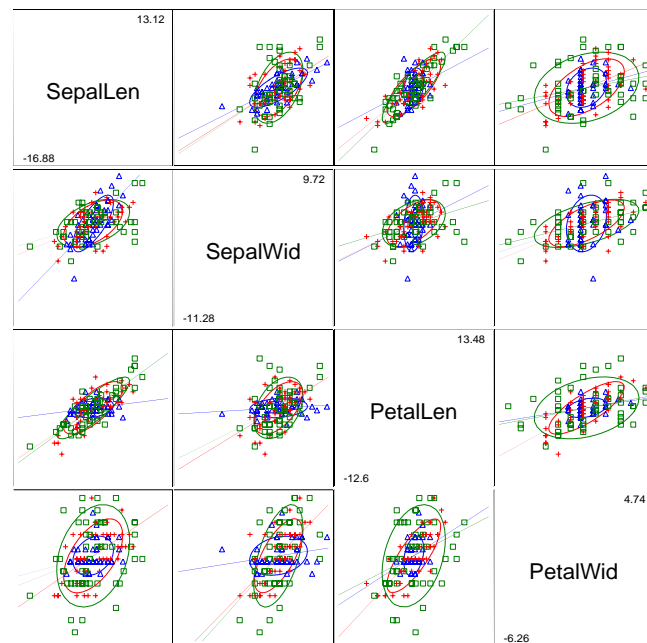
### Robust data ellipse and outlier detection

- Multivariate outliers:

- Shift the mean vector
- Inflate  $\mathbf{S}$  → makes them look less extreme

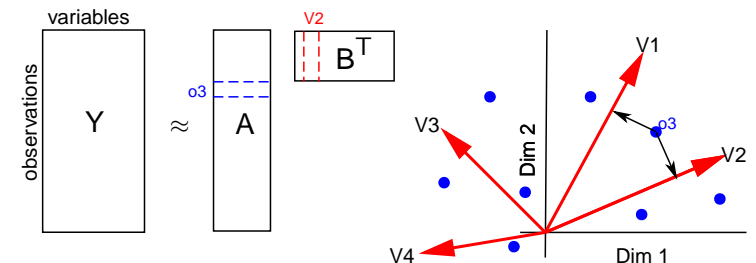


- Robust covariance estimation
  - Multivariate trimming : ignore observations where  $\text{Pr}(D^2(\mathbf{y}_i) > \chi_p^2) < \alpha$ . [implemented in `outlier` macro.]
  - MVE and MCD methods minimize volume or determinant of  $\mathbf{S}$  containing a fraction  $\gamma$  of points. [implemented in `robcov` macro.]
  - use data ellipse based on robust  $\mathbf{S}$  [implemented via `WEIGHT=` in `ellipses` macro.].



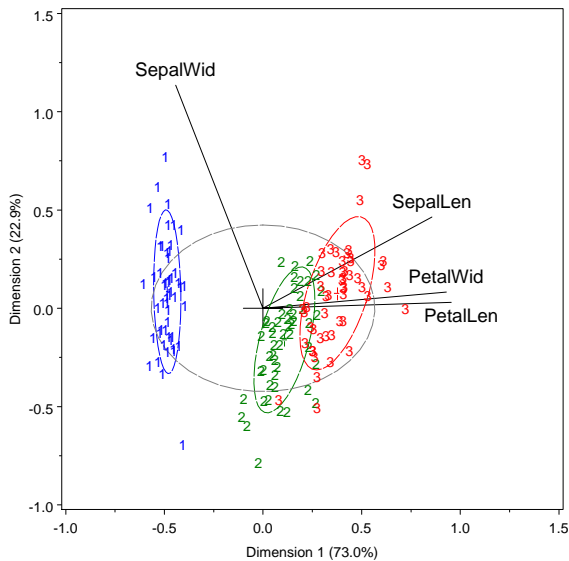
### Biplots: Low-D views of multivariate data

- Display variables *and* observations in a reduced-rank space of  $d$  ( $=2$  or  $3$ ) dimensions,



- Biplot properties:

- Plot observations as points, variables as vectors from origin (mean)
- Angles between vectors show correlations ( $r \approx \cos(\theta)$ )
- $y_{ij} \approx \mathbf{a}_i^T \mathbf{b}_j$ : projection of observation on variable vector
- Observations are uncorrelated overall (but not necessarily within group)

Biplot of Iris data: *setosa*: blue (1); *versicolor*: green (2); *virginica*: red (3)

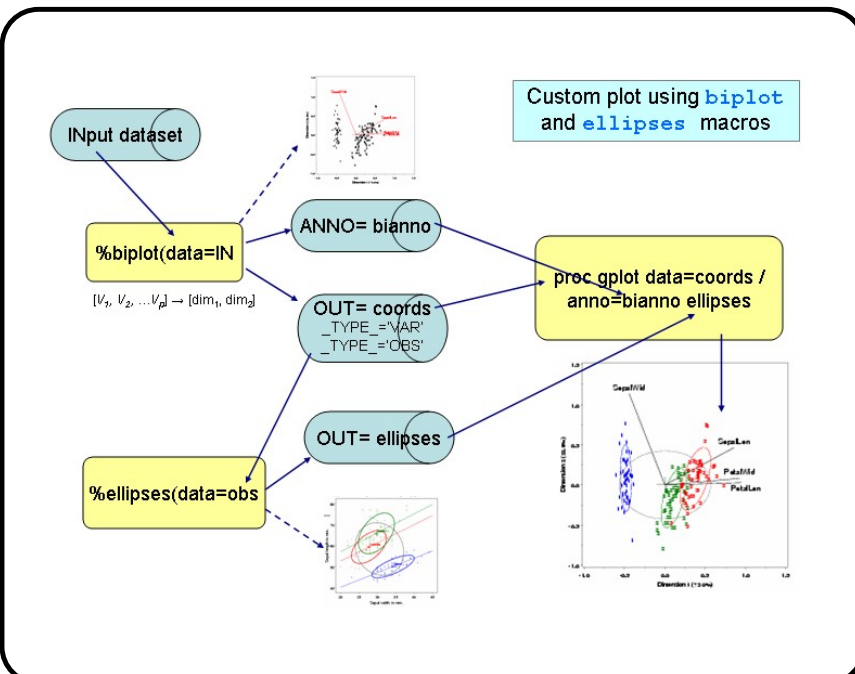
### Custom plot using `biplot` macro and `ellipses` macro

- Use `biplot` macro to obtain biplots scores and variable vectors
- Use `ellipses` macro to obtain 68% data ellipses
- SAS "Annotate" programming to produce custom display

```

----- [bipliris.sas ...] -----
1  *-- Get biplot scores (_type_='OBS') and var vectors (_type_='VAR');
2  %biplot(data=iris,
3     var=SepalLen SepalWid PetalLen PetalWid,
4     id=id,           /* observation ID, here, species number */
5     std=std,        /* standardize to mean=0, var=1 */
6     scale=0.36,    /* Scale factor for variable vectors */
7     htext=1.5 2,   /* text heights for obs. and var labels */
8     xextra=0 1,    /* extra tick mark for labels */
9     gplot=no,      /* suppress the plot */
10    colors=black,   /* we change these later */
11    out=biplot,     /* output coordinates data set */
12    anno=bianno);  /* output Annotate data set

```



### Custom plot using `biplot` macro and `ellipses` macro

```

----- [... bipliris.sas ...] -----
1  *-- Customize the Annotate data set;
2  data bianno;
3  set bianno;
4  if _type_='OBS' then do;
5     *-- change colors;
6     select(_name_);
7         when ('3') color='RED';
8         when ('2') color='GREEN';
9         when ('1') color='BLUE';
10        otherwise;
11    end;
12    end;
13    *-- adjust label position to avoid overplotting;
14    else do; /* _type_='VAR' */
15        if _name_='PetalLen' and function='LABEL' then position='E';
16    end;

```

```

1 _____ [... bipliris.sas ...] _____
2 *-- Select just the observation scores;
3 data biplobs;
4   set biplot;
5   where (_type_='OBS');
6
7 *-- Obtain data ellipses for each group and total sample;
8 %ellipses(data=biplobs,
9   x=dim1, y=dim2, /* data ellipses for biplot dimension */
10  group=_name_,
11  all=yes, /* include total sample ellipse */
12  colors=blue green red gray,
13  plot=no, /* suppress the plot */
14  pvalue=0.68,
15  vaxis=axis98, /* use AXIS statements generated by %biplot */
16  haxis=axis99,
17  out=ellipses /* output Annotate data set */
18 );

```

### HE plots: Visualization for the MLM

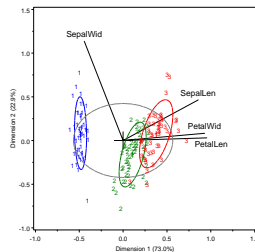
#### ■ LM: $y = X\beta + \epsilon$

- Test *any* hypothesis by *General Linear Test*:  $H_0 : C\beta = 0$ , where  $C$  = matrix of constants; rows specify  $h$  linear combinations or contrasts of parameters.
- e.g., Test of  $H_0 : \beta_1 = \beta_2 = 0$  in model  $y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \epsilon_i$

$$C\beta = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

- All  $\rightarrow$  F-test: How big is  $SS_H$  relative to  $SS_E$ ?

$$F = \frac{SS_H/df_h}{SS_E/df_e} = \frac{MS_H}{MS_E} \rightarrow (MS_H - F MS_E) = 0$$



```

1 _____ [... bipliris.sas] _____
2 *-- Join the biplot annotations and the data ellipses;
3 data bianno;
4   set bianno ellipses;
5
6 *-- Plot the observations along with annotations;
7 symbol1 v=none;
8 proc gplot data=biplot;
9   plot dim2 * dim1 /
10     anno=bianno frame
11     vaxis=axis98 haxis=axis99
12     vminor=1 hminor=1;

```

### HE plots: Visualization for the MLM

#### ■ MLM: $Y = XB + E$ , for $p$ responses, $Y = (y_1, y_2, \dots, y_p)$

- *General Linear Test*:  $H_0 : CB = 0$
- Analogs of sums of squares,  $SS_H$  and  $SS_E$  are  $(p \times p)$  matrices,  $H$  and  $E$ ,

$$H = (C\hat{B})^T [C(X^T X)^{-1} C^T]^{-1} (C\hat{B}),$$

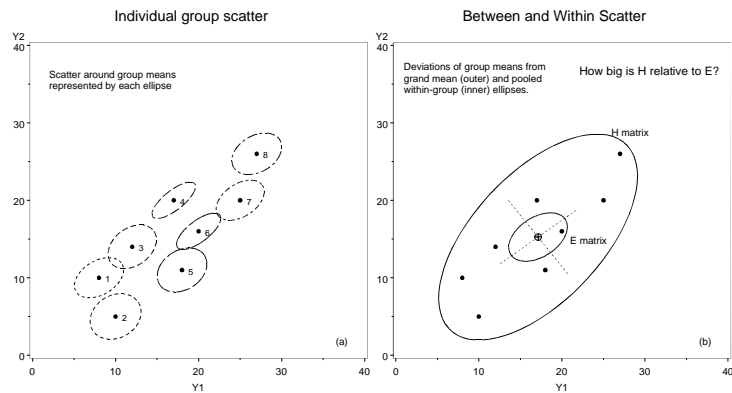
$$E = Y^T [I - H] Y.$$

- Analog of univariate  $F$  is

$$|H - \lambda E| = 0,$$

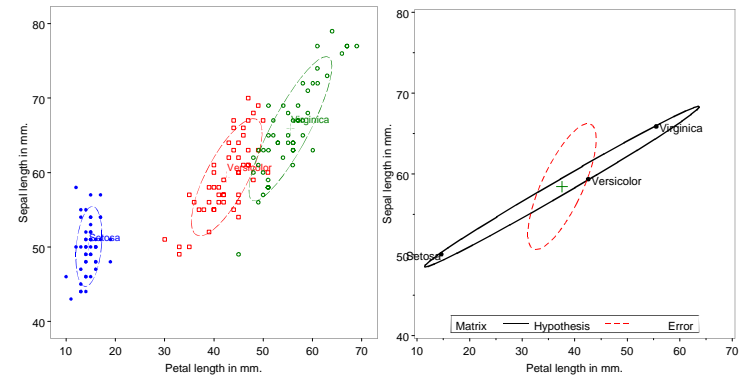
- How big is  $H$  relative to  $E$ ?
  - Latent roots  $\lambda_1, \lambda_2, \dots, \lambda_s$  measure the "size" of  $H$  relative to  $E$  in  $s = \min(p, df_h)$  orthogonal directions.
  - Test statistics (Wilks'  $\Lambda$ , Pillai trace criterion, Hotelling-Lawley trace criterion, Roy's maximum root) all combine info across these dimensions
- HE plot: Shows size, dimensionality, and effect-correlation of  $H$  relative to  $E$ .

### Didactic displays for MANOVA tests



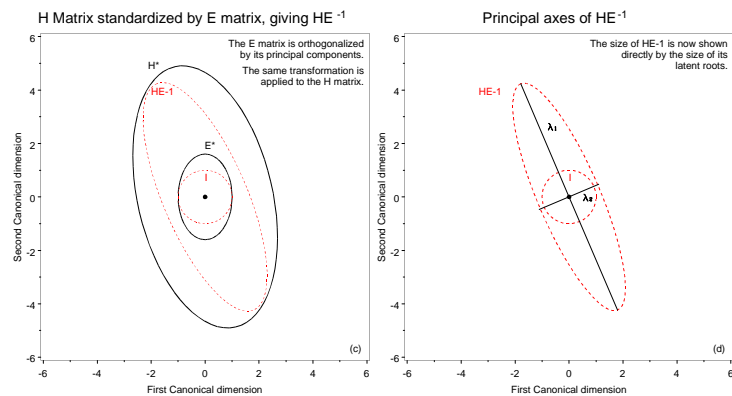
Essential ideas behind multivariate tests: (a) Data ellipses; (b)  $H$  and  $E$  matrices

### HE plot for iris data



(a) Data ellipses and (b)  $H$  and  $E$  matrices (scaled by  $1/df_e$ : effect size)

- $H$  ellipse: data ellipse for fitted values,  $\hat{y}_{ij} = \bar{y}_j$ .
- $E$  ellipse: data ellipse of residuals,  $\hat{y}_{ij} - \bar{y}_j$ .



Essential ideas behind multivariate tests: latent roots of  $HE^{-1}$

### HE plots: heplot macro

- Use PROC GLM to obtain outstat= data set w/  $H$  and  $E$
- (One  $H$  matrix for each effect in model)

```

_____ [heplot3a.sas ...] _____
1 proc glm data=iris outstat=stats noprint;
2   class species;
3   model SepalLen SepalWid PetalLen PetalWid = species / nouni ss3;
_____

■ heplot macro draws  $H$  and  $E$  for specified effect

_____ [... heplot3a.sas ...] _____
1 axis1 label=(a=90) order=(40 to 80 by 10);
2 %heplot(data=iris,
3   stat=stats,           /* Data set containing H & E matrices */
4   var=PetalLen SepalLen, /* Variables to plot */
5   effect=species,       /* Effect to plot */
6   vaxis=axis1);
_____

```

## PROC GLM outstat= dataset

```

1 proc glm data=iris outstat=stats noprint;
2   class species;
3   model SepalLen SepalWid PetalLen PetalWid = species / nouri ss3;
4   * contrast 'S:VV' species -2 1 1;
5   * contrast 'V:V' species 0 -1 1;

```

The outstat= dataset contains:

- $E$  matrix for full model
- One  $H$  matrix for each effect in model
- One  $H$  matrix for each contrast statement

._SOURCE_.	._TYPE_.	._NAME_.	DF	sepalen	sepalwid	petallen	petalwid
ERROR	ERROR	sepalen	147	3896	1363	2462	564
		sepalwid	147	1363	1696	812	481
		petallen	147	2462	812	2722	627
		petalwid	147	564	481	627	616
species	SS3	sepalen	2	6321	-1995	16525	7128
		sepalwid	2	-1995	1134	-5724	-2293
		petallen	2	16525	-5724	43710	18677
		petalwid	2	7128	-2293	18677	8041

## MANOVA: Contrasts

- As in ANOVA, significant effects in MANOVA may be interpreted by using contrasts
- In the GLT, any contrast(s) can be expressed as a  $(h_i \times q)$   $C$  matrix, whose rows sum to zero, e.g.,

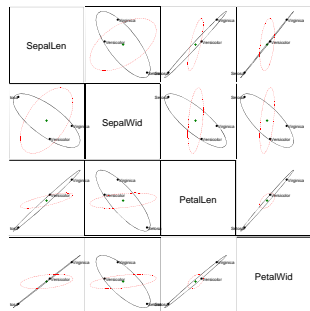
$$C = \begin{bmatrix} -2 & 1 & 1 \\ 0 & 1 & -1 \end{bmatrix}$$

- MANOVA: Pairwise orthogonal 1 df contrasts ( $C_i^T C_j = 0$ )  $\rightarrow$  rank 1  $H_i$  matrices that decompose the overall hypothesis SSCP matrix,

$$H = H_1 + H_2 + \dots + H_{df_h} .$$

- Each  $H_i$  plots as a vector
  - collection  $\rightarrow$  visual summary of the overall test
  - orthogonal in the  $HE^{-1}$  vs.  $I$  plot

## HE plot matrices: hemat macro



[... hematiris.sas]

```

1 %hemat(data=iris,
2   stat=stats,
3   var=SepalLen SepalWid PetalLen PetalWid,
4   effect=species);

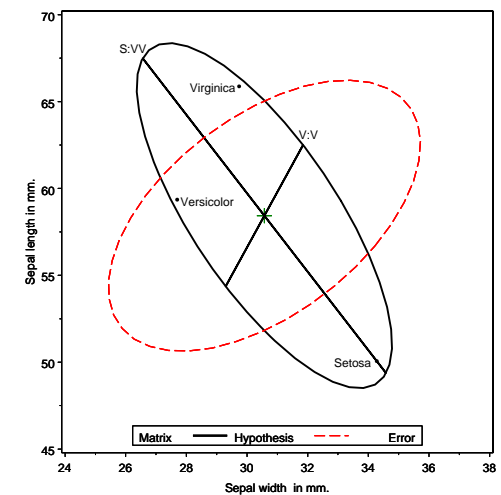
```

## ■ Example: Iris data

- $C_1$ : setosa vs. average of versicolor and virginica ("S:VV")
- $C_2$ : versicolor vs. virginica ("V:V")

$$C_1 = \begin{pmatrix} -2 & 1 & 1 \\ 0 & 1 & -1 \end{pmatrix}$$

$$C_2 = \begin{pmatrix} -2 & 1 & 1 \\ 0 & 1 & -1 \end{pmatrix}$$





### Overlaying multiple HE plots with the `panels` macro

```

1 _____ [heplot4.sas ...] _____
2 proc glm data=iris outstat=stats;
3   class species;
4   model SepalLen sepalwid PetalLen petalwid = species / nouni ss3;
5   contrast 'S:VV' species -2 1 1;
6   contrast 'V:V' species 0 -1 1;
7   manova H=species /short summary;
8
9 _____ [heplot4.sas] _____
10 goptions nodisplay;
11 %heplot(data=iris,stat=stats, var=SepalWid SepalLen, effect=species,
12   vaxis=axis1, haxis=axis2);
13 *-- Contrasts;
14 %heplot(data=iris,stat=stats, var= SepalWid SepalLen,
15   effect=S:VV, ss=contrast, class=, efflab=S:VV,
16   vaxis=axis1, haxis=axis2);
17 %heplot(data=iris,stat=stats, var= SepalWid SepalLen,
18   effect=V:V, ss=contrast, class=, efflab=V:V,
19   vaxis=axis1, haxis=axis2);
20 goptions display;
21 *-- Overlay panels;
22 %panels(rows=1, cols=1, replay=1:1 1:2 1:3);

```

### Reduced-rank HE plots: `canplot` macro and `heplot` macro

- `canplot` macro: get scores `can1` and `can2` on 1st two canonical dimensions

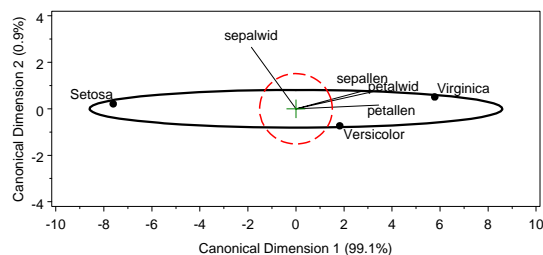
```

1 _____ [hecaniris.sas ...] _____
2 %canplot(
3   data=iris,
4   class=species,
5   var=SepalLen SepalWid PetalLen PetalWid,
6   plot=N0,
7   scale=3.5, /* scale factor for variable vectors */
8   out=canscores, /* output data set containing discrim scores */
9   anno=cananno); /* output data set containing annotations */

```

### Reduced-rank HE plots: `canplot` macro and `heplot` macro

- HE plot in 2D space that maximally discriminates among groups
- → HE plot of canonical scores: weights given by the latent vectors of  $HE^{-1}$



Canonical HE plot for the iris data.

### Reduced-rank HE plots: `canplot` macro and `heplot` macro

- HE plot of canonical scores

```

1 _____ [... hecaniris.sas] _____
2 *-- Get H and E matrices for canonical scores;
3 proc glm data=canscores outstat=stats;
4   class species;
5   model can1 can2 = species / nouni ss3;
6   manova h=species;
7   run;
8
9 *-- Axis statements to equate axis units;
10 axis1 length=2.6 IN order=(-4 to 4 by 2) label=(a=90);
11 axis2 length=6.5 IN order=(-10 to 10 by 2);
12 %heplot(data=canscores, stat=stats,
13   x=Can1, y=Can2,
14   effect=species,
15   haxis=axis2, vaxis=axis1,
16   legend=none, hsym=1.6,
17   anno=cananno);

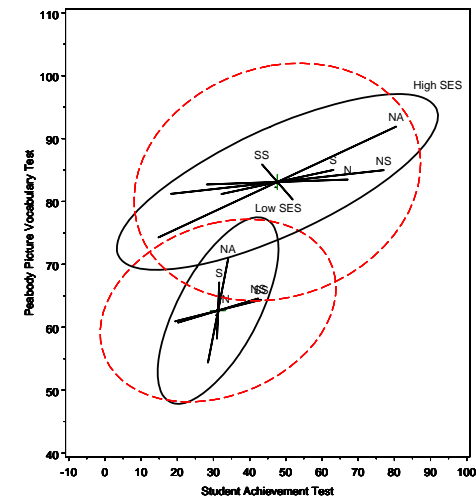
```

### HE plots for MMRA

- **Overall test:**  $H_0 : B = 0$  (all coefficients for all responses are zero)
  - $\rightarrow C = I$  in GLT  $\rightarrow H = \hat{B}^T (X^T X)^{-1} \hat{B}$
- **Individual predictors:**  $H_0 : \beta_i = 0$ 
  - $\rightarrow C = (0, 0, \dots, 1, 0, \dots, 0) \rightarrow H_i = \hat{\beta}_i^T (X^T X)^{-1} \hat{\beta}_i$
- **HE plot**
  - Overall  $H$  ellipse: how predictors relate collectively to responses
  - Individual  $H$  ellipses (vectors):
    - orientation  $\rightarrow$  relation of  $x_i$  to  $y_1, y_2$
    - length  $\rightarrow$  strength of relation

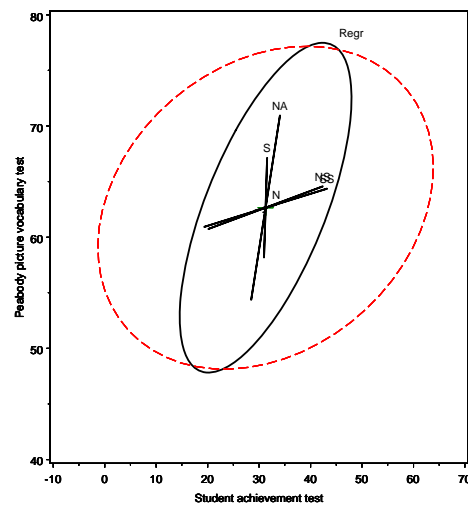
### HE plots for MMRA: ANCOVA

- Rohwer data on  $n_1 = 37$  low SES children, and  $n_2 = 32$  high SES children
  - Fit separate regressions for each group



### HE plots for MMRA: Example

- Rohwer data on  $n = 37$  low SES children, for 5 PA tasks (N, S, NS, NA, SS) predicting intelligence/achievement (PPVT, SAT, Raven): Type III SSCPs



### HE plots for MMRA: ANCOVA

- Rohwer data on  $n_1 = 37$  low SES children, and  $n_2 = 32$  high SES children
  - Fit ANCOVA model (assuming equal slopes)

