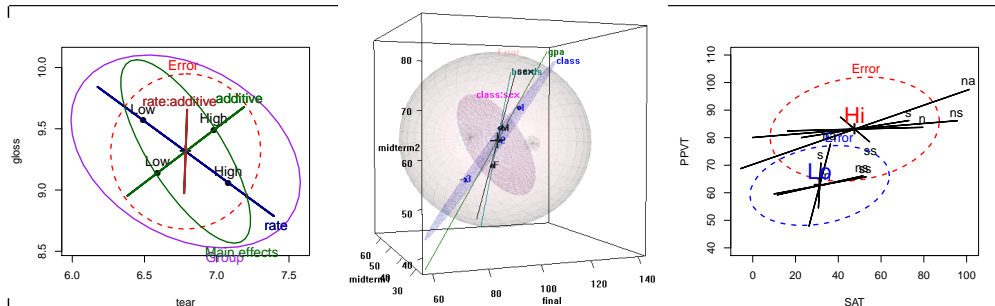


Visualizing multivariate linear models in R

Michael Friendly¹ Matthew Sigal¹

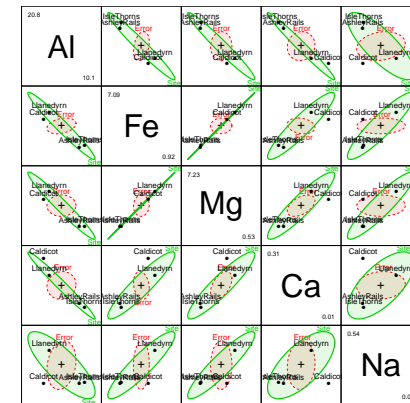
¹York University, Toronto

Cinquièmes Recontres R
Toulouse, June 22–24, 2016



Outline

- 1 Background
 - Overview
 - Visual overview
 - Data ellipses
 - The Multivariate Linear Model
- 2 Hypothesis Error (HE) plots
 - Motivating example
 - Visualizing H and E variation
 - MANOVA designs
 - MREG designs
- 3 Reduced-rank displays
 - Low-D displays of high-D data
 - Canonical discriminant HE plots
- 4 Recent extensions
 - Robust MLMs
 - Influence diagnostics for MLMs
- 5 Conclusions



Slides: <http://datavis.ca/papers/ToulouseR-2x2.pdf>

Overview: Research topics

- This talk outlines research on graphical methods for **multivariate** linear models (MLMs)— extending visualization for multiple regression, ANOVA, and ANCOVA designs to those with several response variables.
- The topics addressed include:
 - Visualizing multivariate tests with **Hypothesis–Error (HE) plots** in 2D and 3D
 - Low-D views: Generalized canonical discriminant analysis → canonical HE plots
 - Visualization methods for tests of equality of covariance matrices in MANOVA designs
 - Extending these methods to **robust** MLMs
 - Developing multivariate analogs of **influence measures** and **diagnostic plots** for MLMs.

Overview: R packages

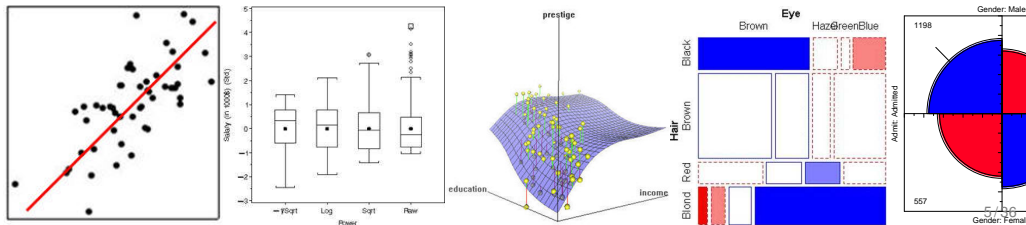
The following R packages implement these methods:

- **car** package: provides the infrastructure for hypothesis tests (**Anova()**) and tests of linear hypotheses (**linearHypothesis()**) in MLMs, including repeated measures designs.
- **heplots** package: implements the HE plot framework in 2D (**heplot()**), 3D (**heplot3d()**), and scatterplot matrix form (**pairs.mlm()**). Also provides:
 - **covEllipses()** for covariance ellipses, with optional robust estimation
 - **boxM()** and related methods for testing / visualizing equality of covariance matrices in MANOVA
 - Tutorial vignettes and many data set examples of use
- **candisc** package: generalized canonical discriminant analysis for an MLM, and associated plot methods.
- **mvinfluence** package: Multivariate extensions of leverage and influence (Cook's D) and **influencePlot.mlm()** in various forms.

Context: The LM family and friends

Models, graphical methods and opportunities

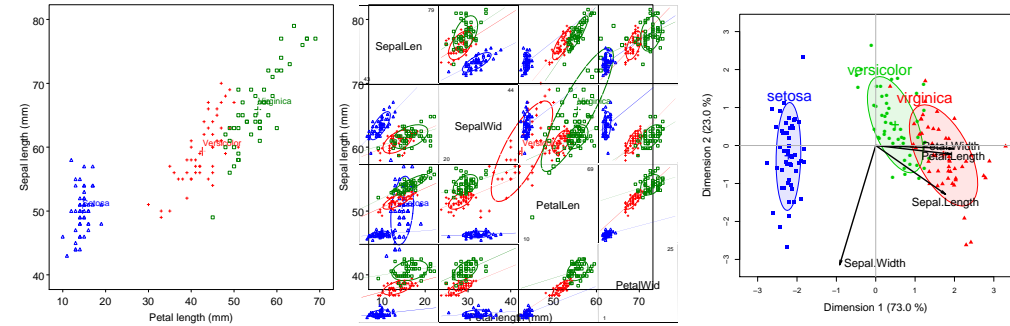
	Classical linear models	Generalized linear models	
# of response variables	1 LM family: $E(y)=X\beta$, $V(y X)=\sigma^2I$ ANOVA, regression, ... Many graphical methods: effect plots, spread-leverage, influence, ...	GLM: $E(y)=g^{-1}(X\beta)$, $V=V[g^{-1}(X\beta)]$ poisson, logistic, loglinear, ... Some graphical methods: mosaic plots, 4fold plots, diagnostic plots, ...	# of response variables
	2+ MLM: $E(Y)=X\beta$, $V(Y X)=I\otimes\Sigma$ MANOVA, MMRreg, ... Graphical methods: ???	MGLM: ??? Graphical methods: ???	



Visual overview: Multivariate data, $Y_{n \times p}$

What we know how to do well (almost)

- 2 vars: Scatterplot + annotations (data ellipses)
- p vars: Scatterplot matrix (all pairs)
- p vars: Reduced-rank display— show max. total variation \mapsto biplot

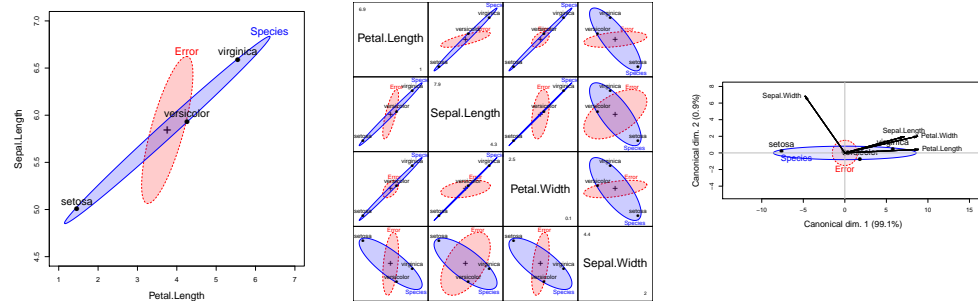


Visual overview: Multivariate linear model, $Y = XB + U$

$$Y = XB + U$$

What is new here?

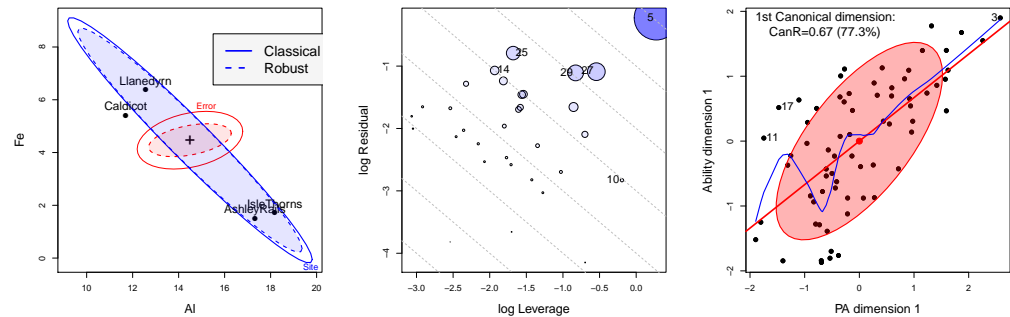
- 2 vars: HE plot— data ellipses of H (fitted) and E (residual) SSP matrices
- p vars: HE plot matrix (all pairs)
- p vars: Reduced-rank display— show max. H wrt. $E \mapsto$ Canonical HE plot



Visual overview: Recent extensions

Extending univariate methods to MLMs:

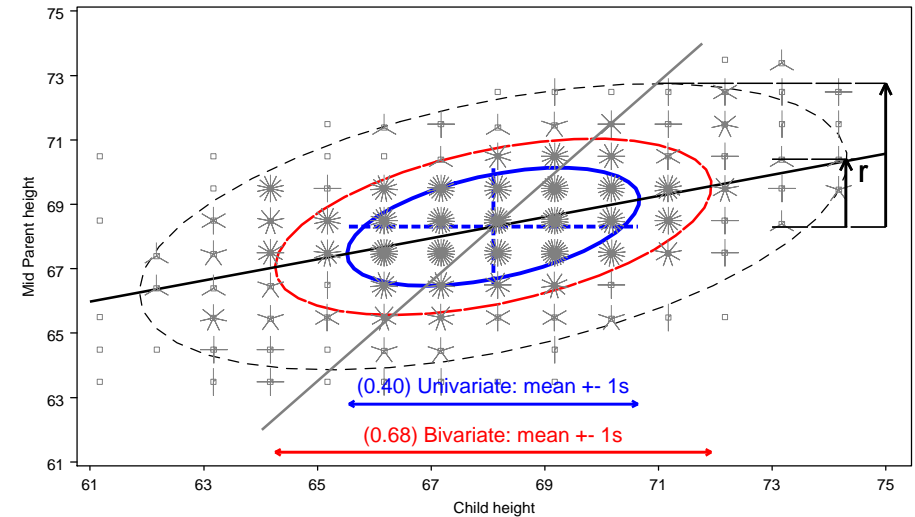
- Robust estimation for MLMs
- Influence measures and diagnostic plots for MLMs
- Visualizing canonical correlation analysis



Data ellipsoids: Visually sufficient summaries

- For any p -variable, multivariate normal $\mathbf{y} \sim \mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, the mean vector $\bar{\mathbf{y}}$ and sample covariance \mathbf{S} are **sufficient statistics**
- Geometrically, contours of constant density are **ellipsoids** centered at $\boldsymbol{\mu}$ with size and shape determined by $\boldsymbol{\Sigma}$
- \mapsto the **data** (concentration) ellipsoid, $\mathcal{E}(\bar{\mathbf{y}}, \mathbf{S})$ is a **sufficient visual summary**
- Easily robustified by using robust estimators of location and scatter

Data Ellipses: Galton's data



Galton's data on Parent & Child height: 40%, 68% and 95% data ellipses

9/36

10/36

The Data Ellipse: Details

• Visual summary for bivariate 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 covariance matrix

- **Radius:** when \mathbf{y} is \approx bivariate normal, $D^2(\mathbf{y})$ has a large-sample χ_2^2 distribution with 2 degrees of freedom.
 - $c^2 = \chi_2^2(0.40) \approx 1$: 1 std. dev univariate ellipse– 1D shadows: $\bar{y} \pm 1s$
 - $c^2 = \chi_2^2(0.68) = 2.28$: 1 std. dev bivariate ellipse
 - $c^2 = \chi_2^2(0.95) \approx 6$: 95% data ellipse, 1D shadows: Scheffé intervals
- **Construction:** Transform the unit circle, $\mathcal{U} = (\sin \theta, \cos \theta)$,

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

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

- **p variables:** Extends naturally to p -dimensional ellipsoids

The univariate linear model

- **Model:** $\mathbf{y}_{n \times 1} = \mathbf{X}_{n \times q} \boldsymbol{\beta}_{q \times 1} + \boldsymbol{\epsilon}_{n \times 1}$, with $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n)$
- **LS estimates:** $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$
- **General Linear Test:** $H_0 : \mathbf{C}_{h \times q} \boldsymbol{\beta}_{q \times 1} = \mathbf{0}$, where \mathbf{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 $y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \epsilon_i$

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

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

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

The multivariate linear model

- **Model:** $Y_{n \times p} = X_{n \times q} B_{q \times p} + U$, for p responses, $Y = (y_1, y_2, \dots, y_p)$
- **General Linear Test:** $H_0 : C_{h \times q} B_{q \times p} = 0_{h \times p}$
- Analogs of sums of squares, SS_H and SS_E are $(p \times p)$ matrices, H and E

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

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

- Analog of univariate F is

$$\det(H - \lambda E) = 0 ,$$

- 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 $s = \min(p, df_n)$ orthogonal directions.
 - Test statistics (Wilks' Λ , Pillai trace criterion, Hotelling-Lawley trace criterion, Roy's maximum root) all combine info across these dimensions

13/36

Motivating Example: Romano-British Pottery

Tubb, Parker & Nicholson analyzed the chemical composition of 26 samples of Romano-British pottery found at four kiln sites in Britain.

- **Sites:** Ashley Rails, Caldicot, Isle of Thorns, Llanedryn
- **Variables:** aluminum (Al), iron (Fe), magnesium (Mg), calcium (Ca) and sodium (Na)
- → One-way MANOVA design, 4 groups, 5 responses

```

1 R> library(heplots)
2 R> Pottery

1           Site  Al  Fe  Mg  Ca  Na
2 1  Llanedryn 14.4 7.00 4.30 0.15 0.51
3 2  Llanedryn 13.8 7.08 3.43 0.12 0.17
4 3  Llanedryn 14.6 7.09 3.88 0.13 0.20
5 . . .
6 25 AshleyRails 14.8 2.74 0.67 0.03 0.05
7 26 AshleyRails 19.1 1.64 0.60 0.10 0.03

```

14/36

HE plots Motivating example

HE plots Motivating example

Motivating Example: Romano-British Pottery

Questions:

- **Can** the content of Al, Fe, Mg, Ca and Na differentiate the sites?
- **How to understand** the contributions of chemical elements to discrimination?

Numerical answers:

```

R> pottery.mod <- lm(cbind(Al, Fe, Mg, Ca, Na) ~ Site)
R> car::Manova(pottery.mod)

```

Type II MANOVA Tests: Pillai test statistic

	Df	test stat	approx F	num Df	den Df	Pr(>F)
Site	3	1.55	4.30	15	60	2.4e-05 ***

```

---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

What have we learned?

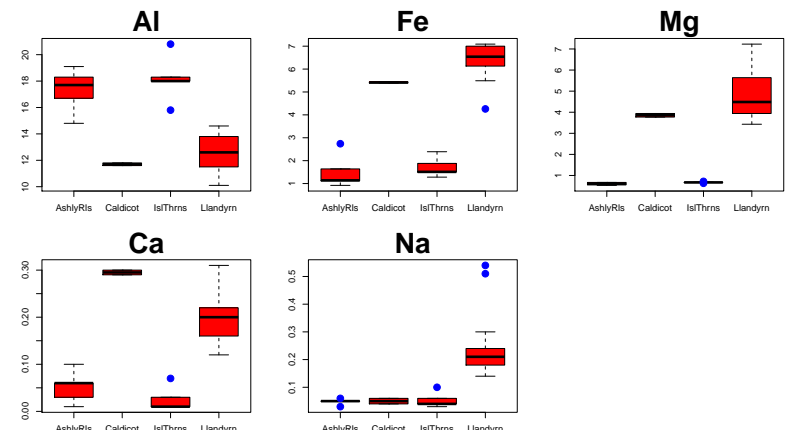
- **Can:** YES! We can discriminate sites.
- **But: How to understand** the pattern(s) of group differences: ???

15/36

Motivating Example: Romano-British Pottery

Univariate plots are limited

- Do not show the *relations* of response variables to each other
- Do not show *how* variables contribute to multivariate tests



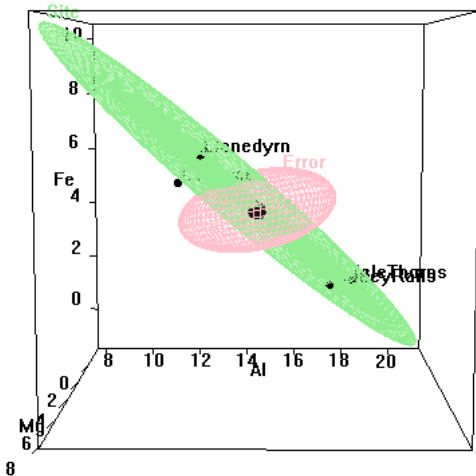
16/36

Motivating Example: Romano-British Pottery

Visual answer: HE plot

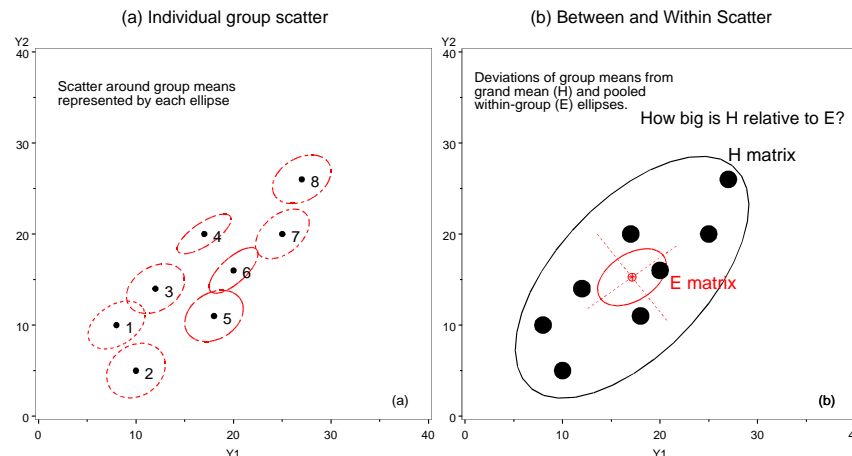
- Shows variation of means (H) relative to residual (E) variation
- Size and orientation of H wrt E : *how much* and *how* variables contribute to discrimination
- Evidence scaling: H is scaled so that it projects outside E iff null hypothesis is rejected.

Run heplot-movie.ppt



```
1 R> heplot3d(pottery.mod)
```

HE plots: Visualizing H and E variation



Ideas behind multivariate tests: (a) Data ellipses; (b) H and E matrices

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

HE plots: Visualizing multivariate hypothesis tests

HE plot details: H and E matrices

Recall the data on 5 chemical elements in samples of Romano-British pottery from 4 kiln sites:

```
1 R> summary(Manova(pottery.mod))
```

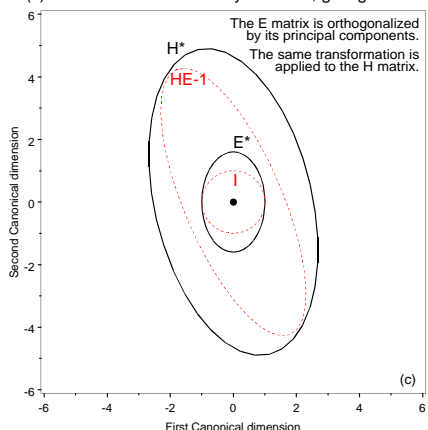
```
1 Sum of squares and products for error:
2   Al   Fe   Mg   Ca   Na
3 Al 48.29  7.080  0.608  0.106  0.589
4 Fe  7.08 10.951  0.527 -0.155  0.067
5 Mg  0.61  0.527 15.430  0.435  0.028
6 Ca  0.11 -0.155  0.435  0.051  0.010
7 Na  0.59  0.067  0.028  0.010  0.199
```

Term: Site

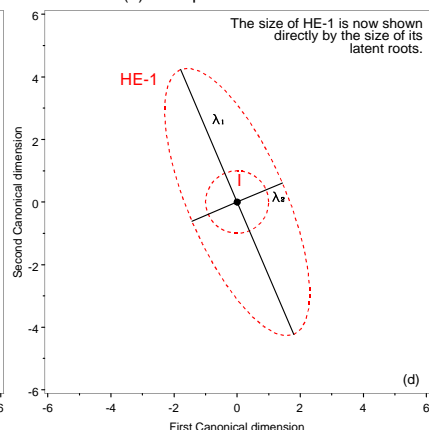
```
12 Sum of squares and products for hypothesis:
13   Al   Fe   Mg   Ca   Na
14 Al 175.6 -149.3 -130.8 -5.89 -5.37
15 Fe -149.3 134.2 117.7 4.82 5.33
16 Mg -130.8 117.7 103.4 4.21 4.71
17 Ca  -5.9   4.8   4.2  0.20 0.15
18 Na  -5.4   5.3   4.7  0.15 0.26
```

- E matrix: Within-group (co)variation of residuals
 - diag: SSE for each variable
 - off-diag: \sim partial correlations
- H matrix: Between-group (co)variation of means
 - diag: SSH for each variable
 - off-diag: \sim correlations of means
- How big is H relative to E ?
- Ellipsoids: $\dim(H) = \text{rank}(H) = \min(p, df_h)$

(c) H Matrix standardized by E matrix, giving HE^{-1}



(d) Principal axes of HE^{-1}

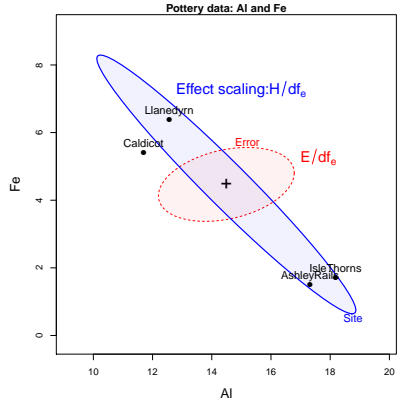


Ideas behind multivariate tests: latent roots & vectors of HE^{-1}

- $\lambda_i, i = 1, \dots, df_h$ show size(s) of H relative to E .
- latent vectors show canonical directions of maximal difference.

HE plot details: Scaling *H* and *E*

- The *E* ellipse is divided by $df_e = (n - p) \rightarrow$ data ellipse of residuals
 - Centered at grand means \rightarrow show factor means in same plot.
- “Effect size” scaling– $H/df_e \rightarrow$ data ellipse of fitted values.
 - Centered at grand means \rightarrow show factor means in same plot.
- “Significance” scaling– *H* ellipse protrudes beyond *E* ellipse *iff* H_0 can be rejected by Roy maximum root test
 - $H/(\lambda_\alpha df_e)$ where λ_α is critical value of Roy’s statistic at level α .
 - direction of *H* wrt *E* \rightarrow linear combinations that depart from H_0 .



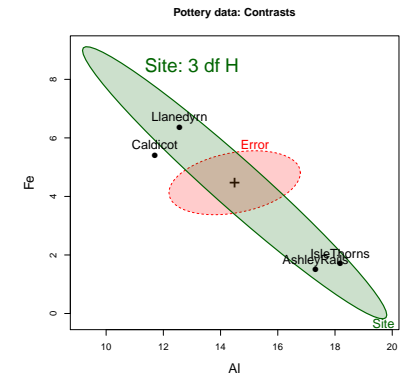
```
R> heplot(pottery.mod, size="effect")
heplot(pottery.mod, size="evidence")
```

R>

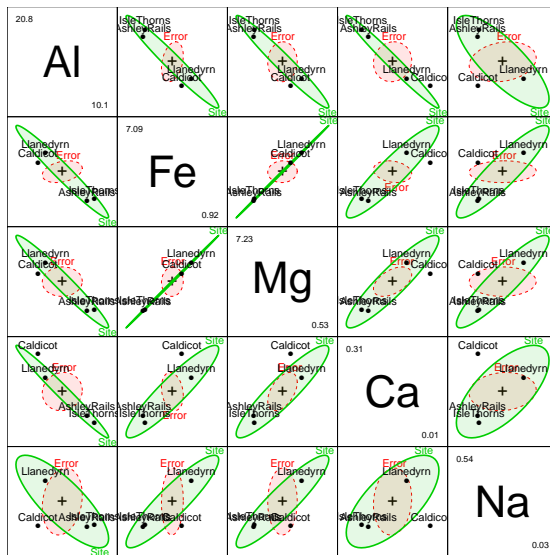
HE plot details: Contrasts and linear hypotheses

- An overall effect \rightarrow an *H* ellipsoid of $s = \min(p, df_h)$ dimensions
- Linear hypotheses, of rank *h*, $H_0 : C_{h \times q} B_{q \times p} = 0_{h \times p} \rightarrow$ sub-ellipsoid of dimension *h*

$$C = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$
- 1D tests and contrasts \rightarrow degenerate 1D ellipses (lines)
- Beautiful geometry:
 - Sub-hypotheses are **tangent** to enclosing hypotheses
 - Orthogonal contrasts form **conjugate axes**



HE plot matrices: All bivariate views



```
R> pairs(pottery.mod)
```

Two-way MANOVA– Plastic film data

- Data from an experiment to determine the optimal conditions for extruding plastic film.
 - Factors: **rate** of extrusion (low/high), amount of **additive** (low/high)
 - Responses: Tear resistance, film gloss, opacity
 - $\rightarrow 2 \times 2$ MANOVA design, 3 responses, $n = 5$ per cell.
- HE plots show main effects, interactions and linear hypotheses in relation to each other

```
1 R> plastic.mod <- lm(cbind(tear, gloss, opacity) ~
2 rate*additive, data=Plastic)
3 R> Manova(plastic.mod, test.statistic="Roy")
```

1 Type II MANOVA Tests: Roy test statistic

	Df	test stat	approx F	num Df	den Df	Pr(>F)
2 rate	1	1.6188	7.5543	3	14	0.003034 **
3 additive	1	0.9119	4.2556	3	14	0.024745 *
4 rate:additive	1	0.2868	1.3385	3	14	0.301782
5 ---						
6 Signif. codes:	0	'***'	0.001	'**'	0.01	'*' 0.05 '.' 0.1 ' ' 1

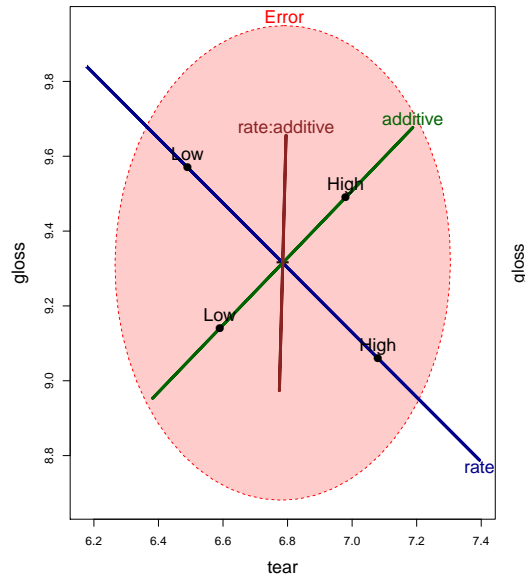
AL stands out – opposite pattern $r(Fe, Mg) \approx 1$

Jump to low-D

Two-way MANOVA– Plastic film data

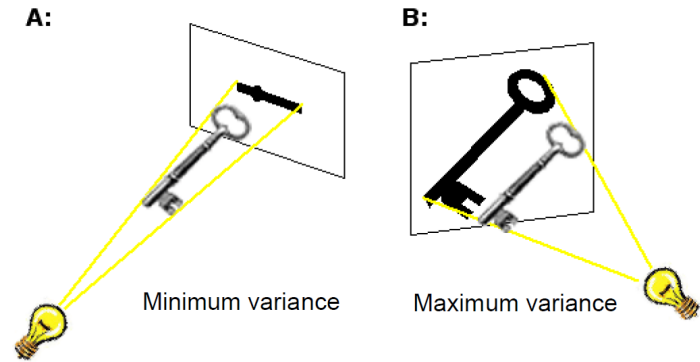
Visualizing tests of model terms and **composite** linear hypotheses

- Main effects & interaction term (1 df each)
- Balanced design: **Hs orthogonal** in 3D
- Add **H** ellipse for test of main effects, Main = rate + additive (2 df)
- Add **H** ellipse for test of group diffs, Group = rate * additive (3 df)



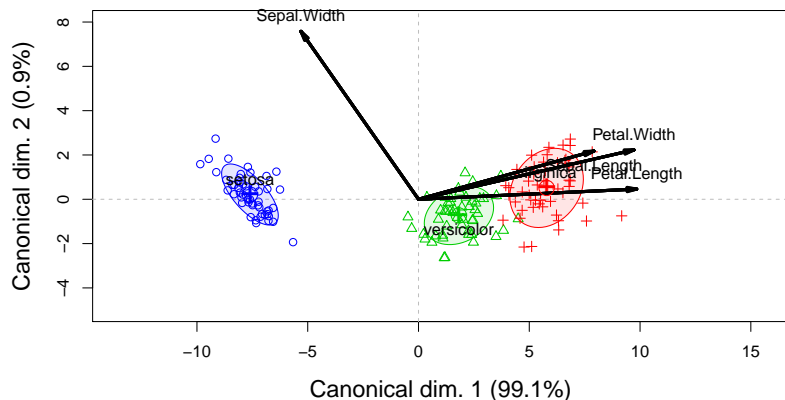
Low-D displays of high-D data

- High-D data often shown in 2D (or 3D) views— orthogonal projections in variable space— **scatterplot**
- **Dimension-reduction** techniques: project the data into subspace that has the largest **shadow**— e.g., accounts for largest variance.
- → low-D approximation to high-D data



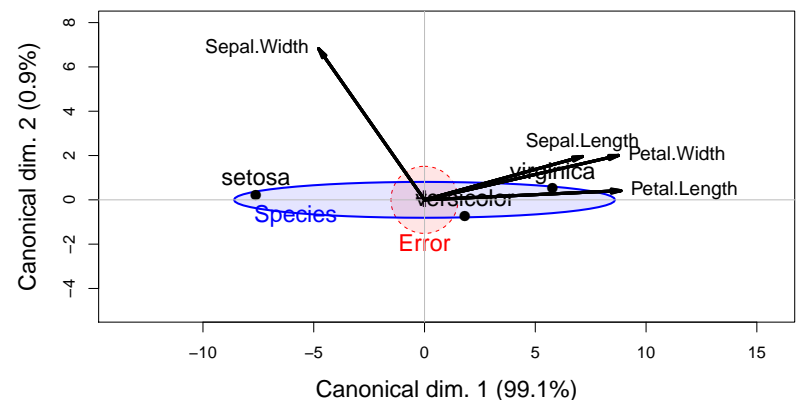
Canonical discriminant HE plots

- As with biplot, we can visualize MLM hypothesis variation for **all** responses by projecting **H** and **E** into low-rank space.
- **Canonical projection**: $Y_{n \times p} \mapsto Z_{n \times s} = YE^{-1/2}V$, where **V** = eigenvectors of HE^{-1} .
- This is the view that maximally discriminates among groups, ie max. **H** wrt **E** !



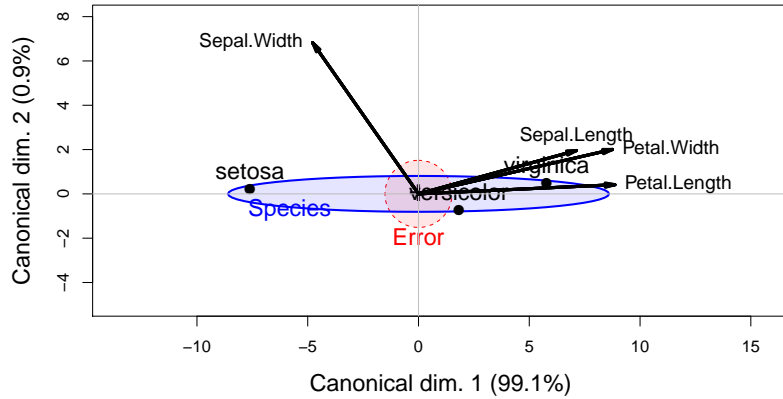
Canonical discriminant HE plots

- Canonical HE plot is just the HE plot of canonical scores, (z_1, z_2) in 2D, or z_1, z_2, z_3 , in 3D.
- As in biplot, we add vectors to show relations of the y_i response variables to the canonical variates.
- variable vectors here are **structure coefficients** = correlations of variables with canonical scores.



Canonical discriminant HE plots: Properties

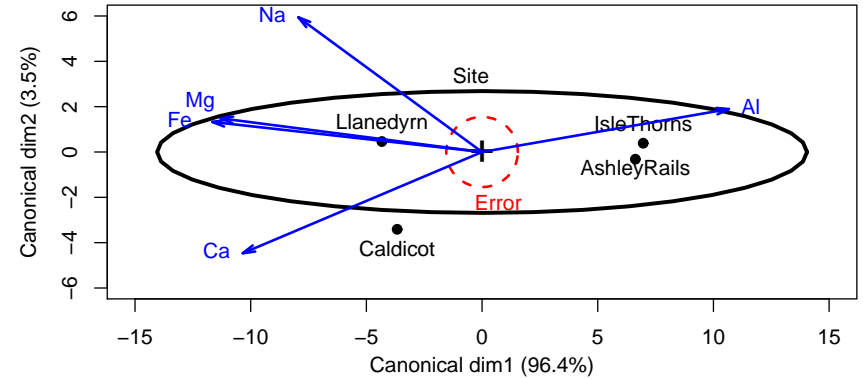
- Canonical variates are uncorrelated: E ellipse is spherical
- \mapsto axes must be equated to preserve geometry
- Variable vectors show how variables discriminate among groups
- Lengths of variable vectors \sim contribution to discrimination



29/36

Canonical discriminant HE plots: Pottery data

- Canonical HE plots provide 2D (3D) visual summary of H vs. E variation
- Pottery data: $p = 5$ variables, 4 groups $\mapsto df_H = 3$
- Variable vectors: Fe, Mg and Al contribute to distinguishing (Caldicot, Llanedryn) from (Isle Thorns, Ashley Rails): 96.4% of mean variation
- Na and Ca contribute an additional 3.5%. **End of story!**



Run heplot-movie.ppt

30/36

Robust MLMs

- R has a large collection of packages dealing with robust estimation:
 - `robust::lmrob()`, `MASS::rlm()`, for *univariate* LMs
 - `robust::glmrob()` for univariate *generalized* LMs
 - **High breakdown-bound** methods for robust *PCA* and robust covariance estimation
 - However, none of these handle the **fully general MLM**
- `heplots` now provides `robmlm()` for robust MLMs:
 - Uses a simple M-estimator via iteratively re-weighted LS.
 - Weights: calculated from Mahalanobis squared distances, using a simple robust covariance estimator, `MASS::cov.trob()` and a weight function, $\psi(D^2)$.

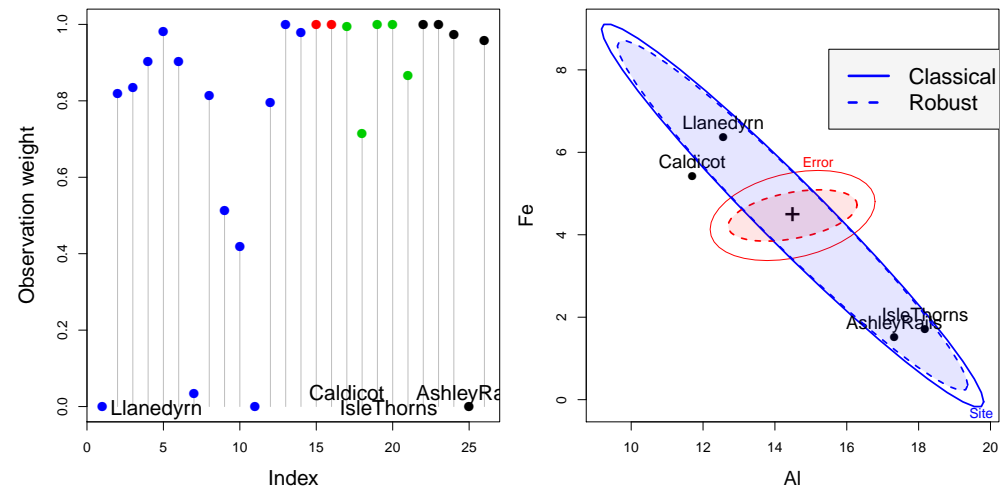
$$D^2 = (\mathbf{Y} - \hat{\mathbf{Y}})^T \mathbf{S}_{\text{trob}}^{-1} (\mathbf{Y} - \hat{\mathbf{Y}}) \sim \chi_p^2 \quad (1)$$

- This fully extends the "mlm" class
- Compatible with other `mlm` extensions: `car::Anova()` and `heplot()`.

31/36

Robust MLMs: Example

For the Pottery data:



- Some observations are given weights ~ 0
- The E ellipse is considerably reduced, enhancing apparent significance

32/36

Influence diagnostics for MLMs

- Influence measures & diagnostic plots well-developed for *univariate* LMs
 - Influence measures: Cook's D, DFFITS, dfbetas, etc.
 - Diagnostic plots: Index plots, `car::influencePlot()` for LMs
 - However, these have been unavailable for MLMs
- The `mvinfluence` package now provides:
 - Calculation for multivariate analogs of univariate influence measures (following Barrett & Ling, 1992), e.g., Hat values & Cook's D :

$$H_I = \mathbf{X}_I(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_I^T \quad (2)$$

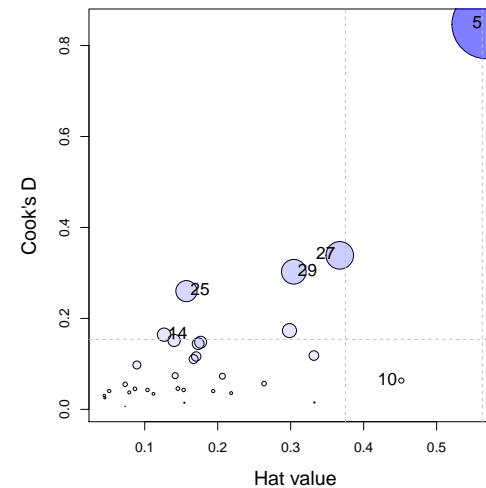
$$D_I = [\text{vec}(\mathbf{B} - \mathbf{B}_{(I)})]^T [\mathbf{S}^{-1} \otimes (\mathbf{X}^T \mathbf{X})] [\text{vec}(\mathbf{B} - \mathbf{B}_{(I)})] \quad (3)$$

- Provides deletion diagnostics for *subsets* (I) of size $m \geq 1$.
- e.g., $m = 2$ can reveal cases of **masking** or **joint influence**.
- Extension of `influencePlot()` to the multivariate case.
- A new plot format: **leverage-residual (LR) plots** (McCulloch & Meeter, 1983)

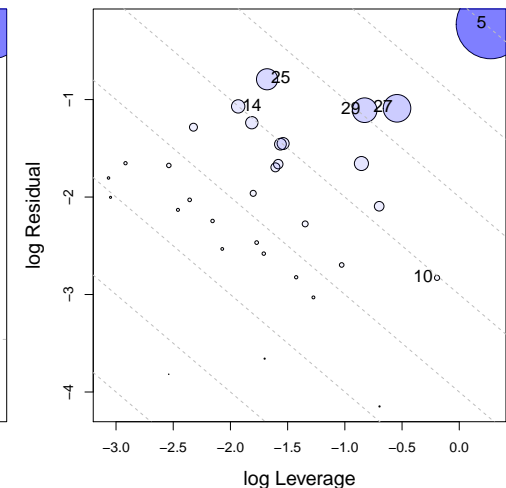
33/36

Influence diagnostics for MLMs: Example

For the Rohwer data:



Cook's D vs. generalized Hat value

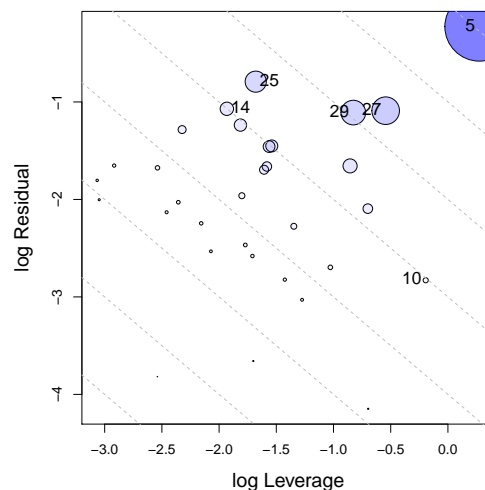


Leverage - Residual (LR) plot

34/36

Influence diagnostics for MLMs: LR plots

- Main idea: Influence \sim Leverage (L) \times Residual (R)
- $\mapsto \log(\text{Infl}) = \log(L) + \log(R)$
- \mapsto contours of constant influence lie on lines with slope = -1.
- Bubble size \sim influence (Cook's D)
- This simplifies interpretation of influence measures



35/36

Conclusions: Graphical methods for MLMs

Summary & Opportunities

- **Data ellipse**: visual summary of bivariate relations
 - Useful for multiple-group, MANOVA data
 - Embed in scatterplot matrix: pairwise, bivariate relations
 - Easily extend to show partial relations, robust estimators, etc.
- **HE plots**: visual summary of multivariate tests for MANOVA and MMRA
 - Group means (MANOVA) or 1-df H vectors (MMRA) aid interpretation
 - Embed in HE plot matrix: all pairwise, bivariate relations
 - Extend to show partial relations: HE plot of "adjusted responses"
- **Dimension-reduction techniques**: low-rank (2D) visual summaries
 - Biplot: Observations, group means, biplot data ellipses, variable vectors
 - Canonical HE plots: Similar, but for dimensions of maximal discrimination
- **Beautiful and useful geometries**:
 - Ellipses everywhere; eigenvector-ellipse geometries!
 - Visual representation of significance in MLM
 - Opportunities for other extensions

— FIN et Merci —

36/36