

Background Hypothesis-Error (HE) plots Reduced-rank displays Recent extensions Conclusions

Recent Advances in Visualizing Multivariate Linear Models

Michael Friendly Matthew Sigal

Statistical Society of Canada, May 26–29, 2013

Background Hypothesis-Error (HE) plots Reduced-rank displays Recent extensions Conclusions

Outline

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

Slides & R scripts: <http://datavis.ca/papers/ssc2013/>

Background Hypothesis-Error (HE) plots Reduced-rank displays Recent extensions Conclusions

Introduction: The LM family and friends

Models, graphical methods and opportunities

	Classical linear models	Generalized linear models
1	LM family: $E(y) = X\beta$, $V(y X) = \sigma^2 I$ 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, ...
2+	MLM: $E(Y) = X\beta$, $V(Y X) = \Theta\Sigma$ MANOVA, MMR, ... Graphical methods: ???	MGLM: ??? Graphical methods: ???

of response variables

Background Hypothesis-Error (HE) plots Reduced-rank displays Recent extensions Conclusions

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

What we know how to do well (almost)

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

Background Hypothesis-Error (HE) plots Reduced-rank displays Recent extensions Conclusions

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

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 \rightarrow Canonical HE plot

Background Hypothesis-Error (HE) plots Reduced-rank displays Recent extensions Conclusions

Visual overview: Recent extensions

Extending univariate methods to MLMs:

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

Background Hypothesis-Error (HE) plots Reduced-rank displays Recent extensions Conclusions

Data Ellipses: Galton's data

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

Background Hypothesis-Error (HE) plots Reduced-rank displays Recent extensions Conclusions

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 $\leq c^2$

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

S = sample variance-covariance matrix

- Radius: when y is \approx bivariate normal, $D^2(y)$ 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: $\bar{y} \pm 1s$
 - $c^2 = \chi^2_2(0.68) \approx 2.28$ 1 std. dev bivariate ellipse
 - $c^2 = \chi^2_2(0.95) \approx 6$ 95% data ellipse, 1D shadows: Scheffé intervals
- Construction: Transform the point circle, $U = (\sin \theta, \cos \theta)$,

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

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

- Robustify: Use robust estimate of S , e.g., MVE (minimum volume ellipsoid)
- p variables: Extends naturally to p-dimensional ellipsoids

Background Hypothesis-Error (HE) plots Reduced-rank displays Recent extensions Conclusions

The univariate linear model

- Model: $y_{n \times 1} = X_{n \times q} \beta_{q \times 1} + \epsilon_{n \times 1}$, with $\epsilon \sim N(0, \sigma^2 I_n)$
- LS estimates: $\hat{\beta} = (X^T X)^{-1} X^T y$
- General Linear Test: $H_0: C_{h \times q} \beta_{q \times 1} = 0$, 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 ?

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

Background Hypothesis-Error (HE) plots Reduced-rank displays Recent extensions Conclusions

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: $H_0: C_{h \times q} B_{q \times p} = 0_{h \times p}$
- Analog 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 $s = \min(p, df_e)$ orthogonal directions.
 - Test statistics (Wilks' Λ , Pillai trace criterion, Hotelling-Lawley trace criterion, Roy's maximum root) all combine info across these dimensions

Background Hypothesis-Error (HE) plots Reduced-rank displays Recent extensions Conclusions

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, Llanellyrn
- Variables: aluminum (Al), iron (Fe), magnesium (Mg), calcium (Ca) and sodium (Na)
- \rightarrow One-way MANOVA design, 4 groups, 5 responses

```
R> library(heplots)
R> Pottery
```

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

Background Hypothesis-Error (HE) plots Reduced-rank displays Recent extensions Conclusions

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> 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: ???

Background Hypothesis-Error (HE) plots Reduced-rank displays Recent extensions Conclusions

Motivating Example: Romano-British Pottery

Univariate plots are limited

- Do not show the relations of variables to each other

Background Hypothesis-Error (HE) plots Reduced-rank displays Recent extensions Conclusions

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.

R> heplot3d(pottery.mod)

Background Hypothesis-Error (HE) plots Reduced-rank displays Recent extensions Conclusions

HE plots: Visualizing H and E (co) variation

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

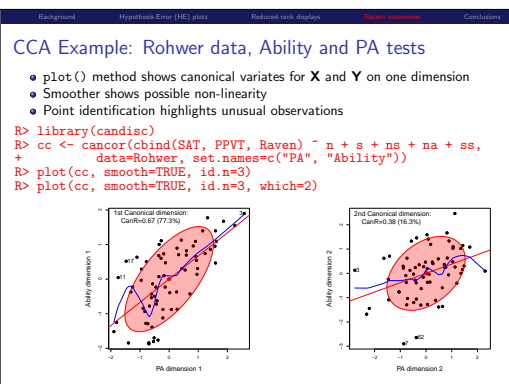
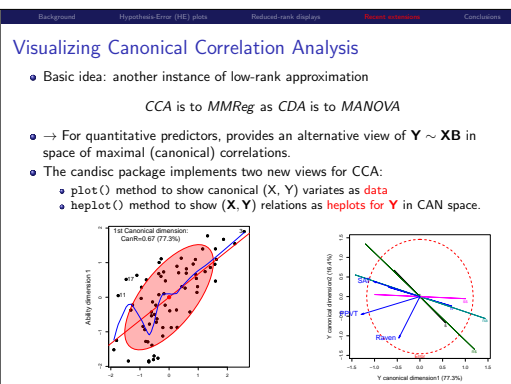
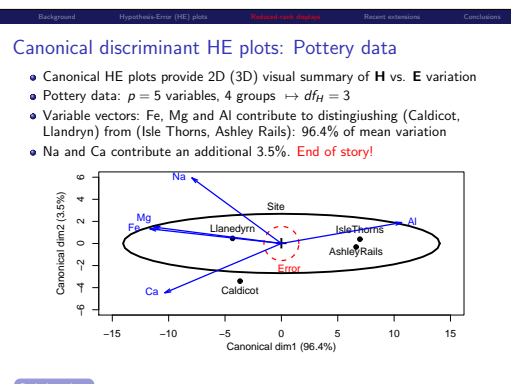
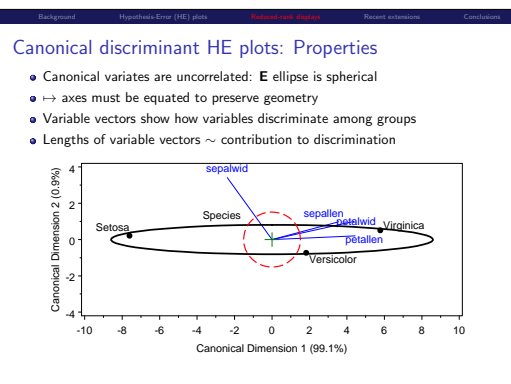
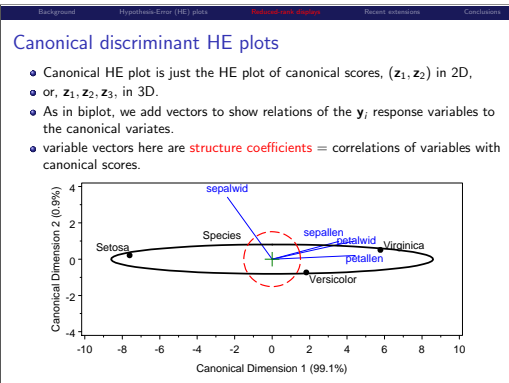
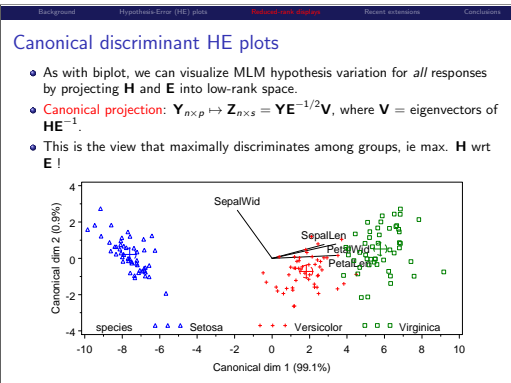
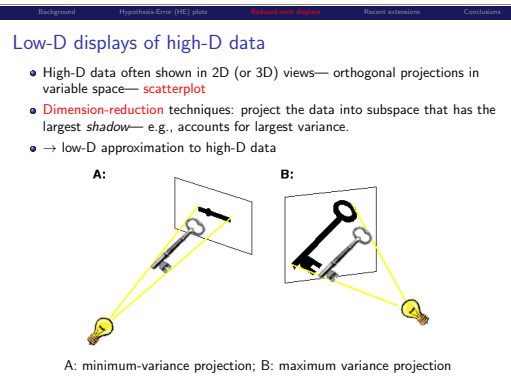
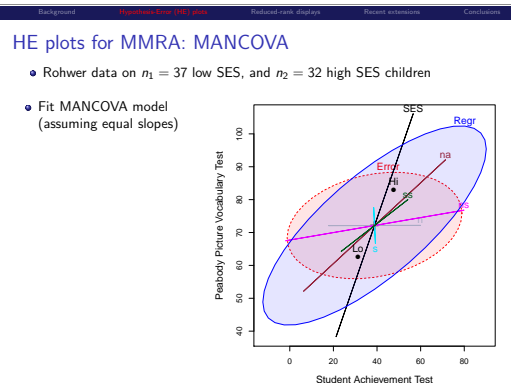
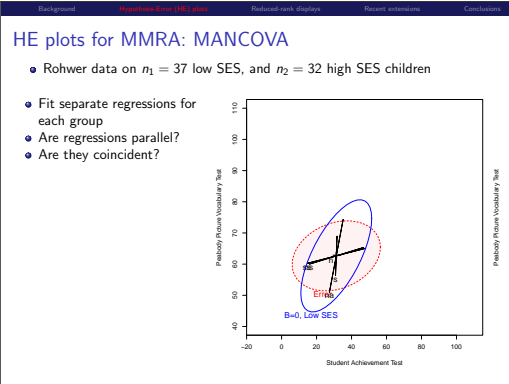
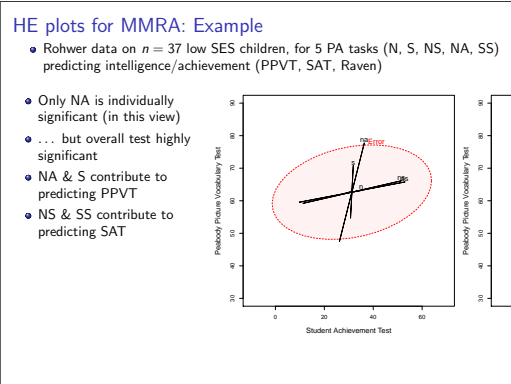
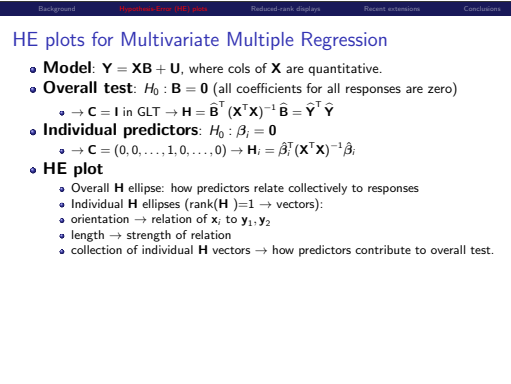
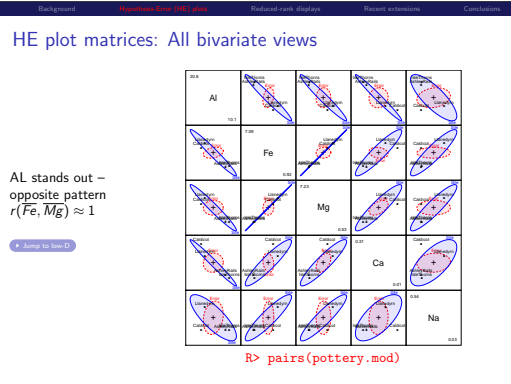
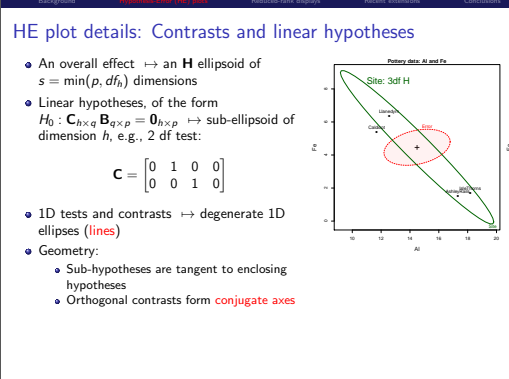
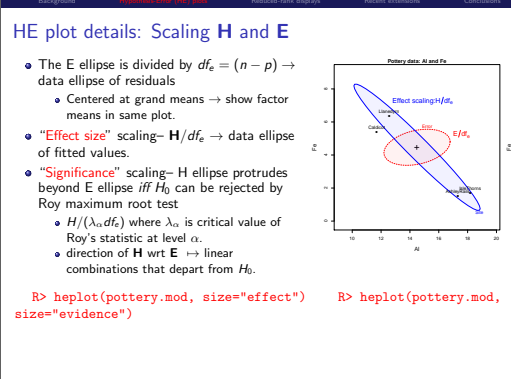
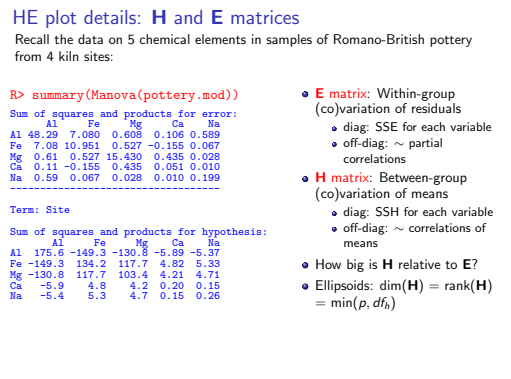
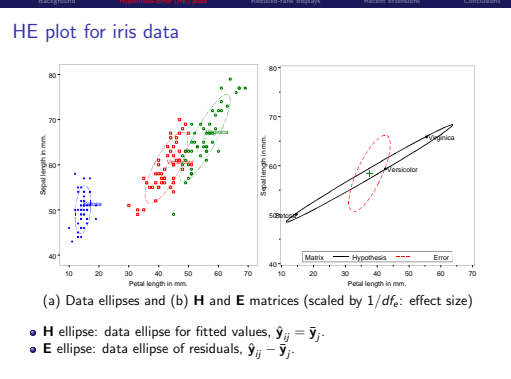
- H ellipse: data ellipse for fitted values, $\hat{y}_j = \bar{y}_j$
- E ellipse: data ellipse of residuals, $\hat{y}_j - \bar{y}_j$

Background Hypothesis-Error (HE) plots Reduced-rank displays Recent extensions Conclusions

HE plots: Visualizing multivariate hypothesis tests

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.



Background

Hypothesis-Error (HE) plots

Reduced-rank displays

Recent extensions

Conclusions

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**
- The `heplots` package 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.rob()` and a weight function, $\psi(D^2)$.

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

- This fully extends the "mlm" class
- Compatible with other mlm extensions: `car::Anova` and `heplots::heplot`.
- Downside: Does not incorporate modern consistency factors or other niceties.

Background

Hypothesis-Error (HE) plots

Reduced-rank displays

Recent extensions

Conclusions

Robust MLMs: Example

For the Pottery data:

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

Background

Hypothesis-Error (HE) plots

Reduced-rank displays

Recent extensions

Conclusions

Influence diagnostics for MLMs

- Influence measures and diagnostic plots are 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:

$$H_I = \mathbf{X}_I (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}_I^T \tag{2}$$
$$D_I = [\text{vec}(\mathbf{B} - \mathbf{B}_{(I)})]^T [\mathbf{S}^{-1} \otimes (\mathbf{X}^T \mathbf{X})] [\text{vec}(\mathbf{B} - \mathbf{B}_{(I)})] \tag{3}$$

- 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)

Background

Hypothesis-Error (HE) plots

Reduced-rank displays

Recent extensions

Conclusions

Influence diagnostics for MLMs: Example

For the Rohrer data:

Cook's D vs. generalized Hat value

Leverage - Residual (LR) plot

Background

Hypothesis-Error (HE) plots

Reduced-rank displays

Recent extensions

Conclusions

Influence diagnostics for MLMs: LR plots

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

Background

Hypothesis-Error (HE) plots

Reduced-rank displays

Recent extensions

Conclusions

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 —