# Visualizing Multivariate Uncertainty:

# Some Graphical Methods for Multivariable Spatial Data

## Michael Friendly

York University

`http://www.math.yorku.ca/SCS/friendly.html`

National Academy of Sciences

Washington, DC, Mar, 2005

**Outline**

- **Multivariate uncertainty and "moral statistics"**
  - A. M. Guerry's *Moral Statistics of France*
  - Guerry's data and analyses

- **Multivariate analyses**: Data-centric displays
  - Bivariate plots and data ellipses
  - Biplots
  - Canonical discriminant plots
  - HE plots for multivariate linear models

- **Multivariate mapping**: Map-centric displays
  - Star maps
  - Reduced-rank color maps

## Multivariate Uncertainty and "Moral Statistics" $\sim$ 1800

*It is a capital mistake to theorize before one has data.*     Sherlock Homes in *Scandal in Bohemia*

- **What to do about crime?**
  - Liberal view: increase education, literacy
  - Conservative view: build more prisons

- **What to do about poverty?**
  - Liberal view: increase social assistance
  - Conservative view: build more poor-houses

- **But**:
  - Little actual data – all armchair theorizing
  - No ways to understand or visualize *relationships* between variables
    - Statistical graphics just invented (Playfair)— line graph, bar chart, pie chart
    - All 1D or 1.5D (time series)

## The rise of "moral statistics" and modern social science

- **Political arithmetic**: William Petty (and others)
  - 1654— first attempt at scientific survey (on Irish estates)
  - 1687— idea that wealth and strength of a state depended on its subjects (number and characteristics)

- **Demography**: Johann Peter Süssmilch (1741)—
  - importance of measuring and analyzing population distributions
  - idea that ethical and state policies could encourage growth and wealth (increase birth rate, decrease death rate)
    - discourage alcohol, gambling, prostitution & priestly celibacy
    - encourage state support for medical care, distribution of land, lower taxes

- **Statistik**: Numbers of the state (1800–1820), Germany and France
  - collect data on imports, exports, transportation, ...

- **Guerry & Quetelet**
  - Quetelet: Concepts of "average man" and "social physics"
  - Guerry: First real social data analysis (Guerry, 1833)

## Guerry's data

- **Compte général** de l'administration de la justice criminelle en France
  - The first national compilation of official justice data (1825)
    - detailed data on all charges and disposition
    - collected quarterly in all 86 departments.
  - Other sources: Bureau de Longitudes (illegitimate births); Parent-Duchâtelet (prostitutes in Paris); Compte du ministere du guerre (military desertions); ...

- **Moral variables**: Scaled so 'more' is 'better'

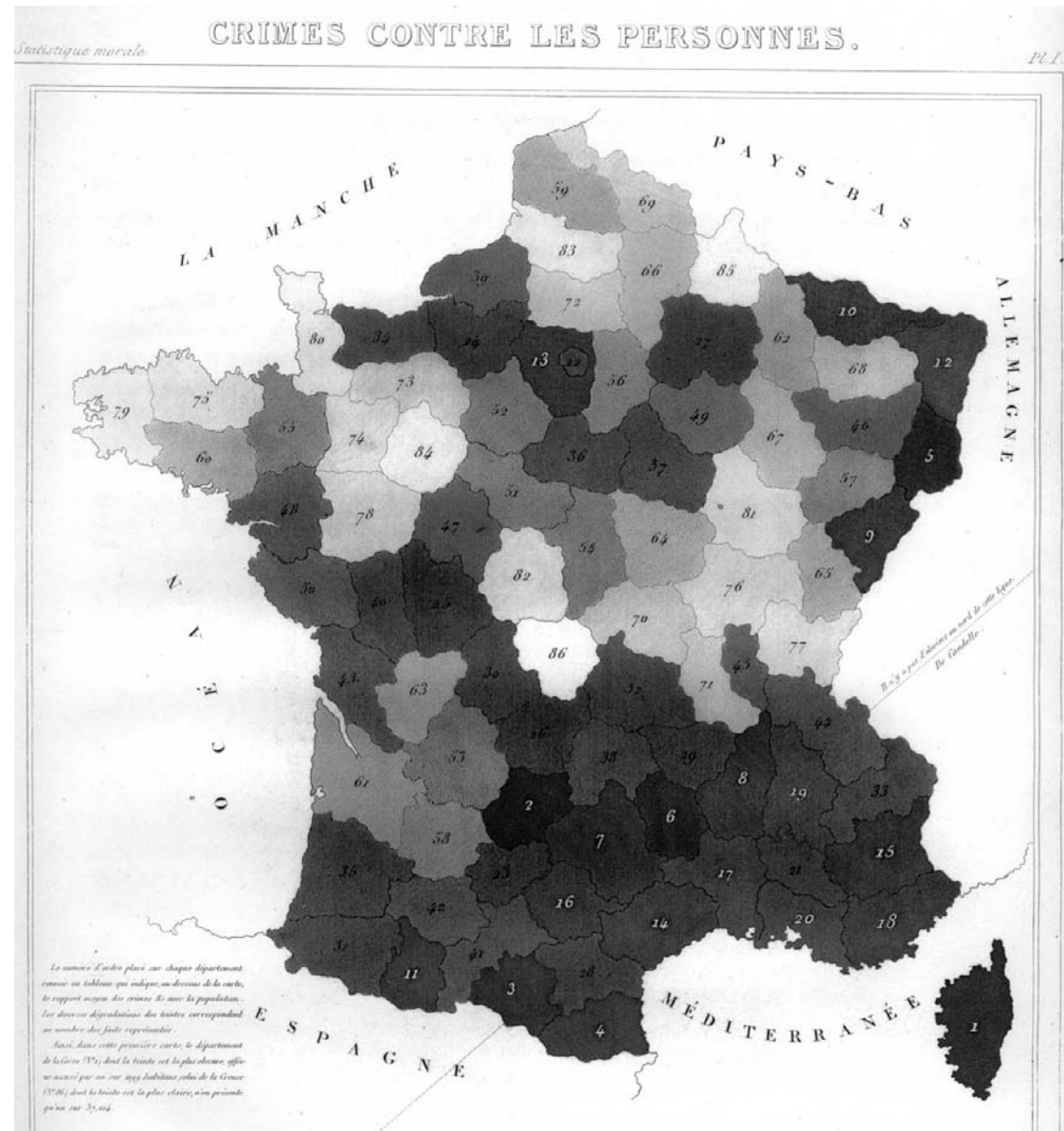  | | |
  |---|---|
  | `Crime_pers` | Population per Crime against persons |
  | `Crime_prop` | Population per Crime against property |
  | `Donations` | Donations to the poor |
  | `Infants` | Population per illegitimate birth |
  | `Literacy` | Percent who can read & write |
  | `Suicides` | Population per suicide |

  - Tried to define these to ensure comparability and representativeness
    - Crime: Use number of *accused* rather than *convicted*
    - Literacy: Reported levels of education unreliable; use data from military draft examinations (% of young men able to read and write)

- **Other variables**: Ranks by department: wealth, commerce, ...

## Guerry's Questions

■ Should crime and other moral variables be considered as structural, lawful characteristics of society, or simply as indicants of individual behavior?

- ■ Statistical regularity as the key to social science ("social physics") social equivalent of "law of large numbers")

- ■ Guerry showed that rates of crime had nearly invariant distributions over time (1825–1830) when classified by region, sex of accused, type of crime, etc. *"We would be forced to recognze that the facts of moral order, like those of physical order, obey invariant laws..."* (p.14)

■ Relations between crime and other moral variables

- ■ Do crimes against persons and crimes against property show the same or different trends?

- ■ How does crime relate to education and literacy?
  - ● Some "armchair" arguments had suggested increasing literacy to decrease crime: *"The definitive result shows that 67 out of 100 prisoners can neither read nor write. What stronger proof could there be that ignorance is the mother of all vices"* (A. Taillander, 1828)

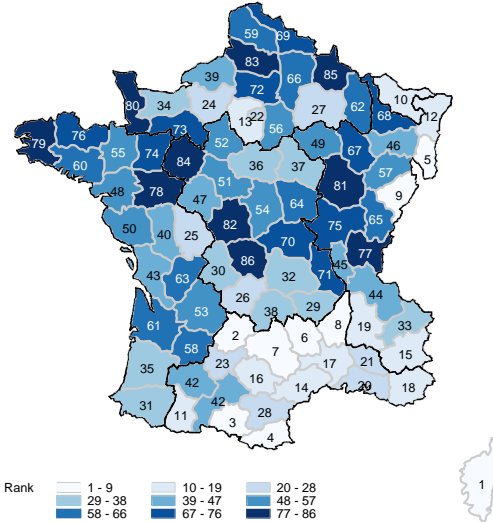- ■ Does crime vary coherently over regions of France (C, N, S, E, W)?
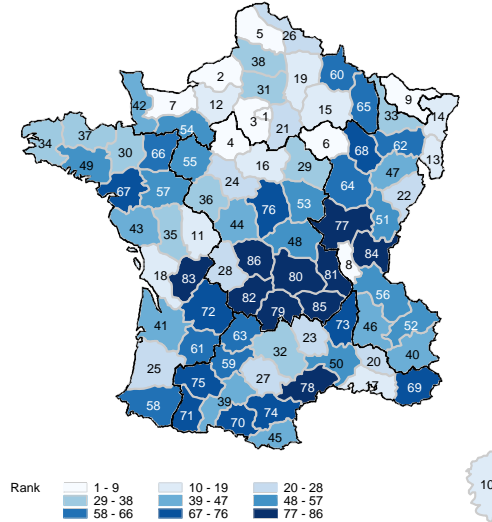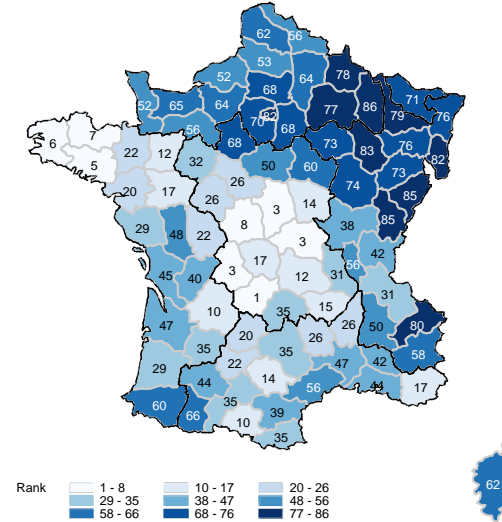
**Guerry's maps**

# Guerry's maps

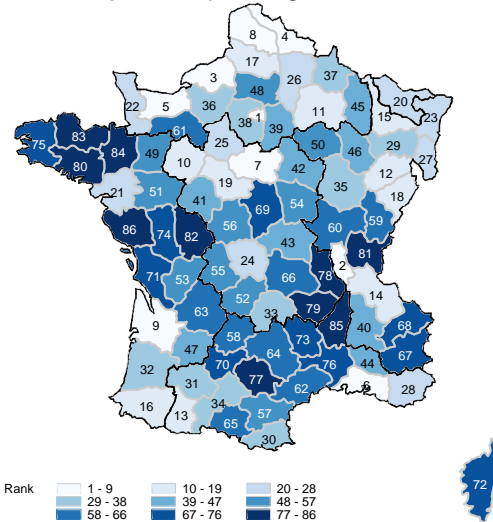### Population per Crime against persons



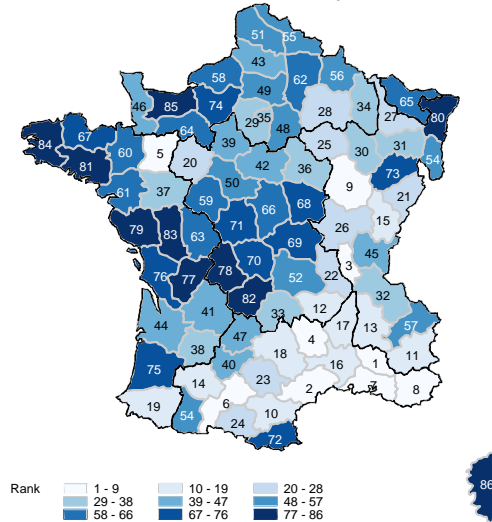### Population per Crime against property

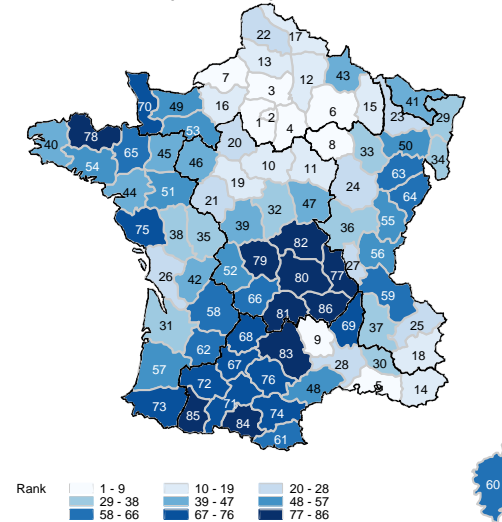

### Per cent who can Read and Write



### Population per Illegitimate birth



### Donations to the poor
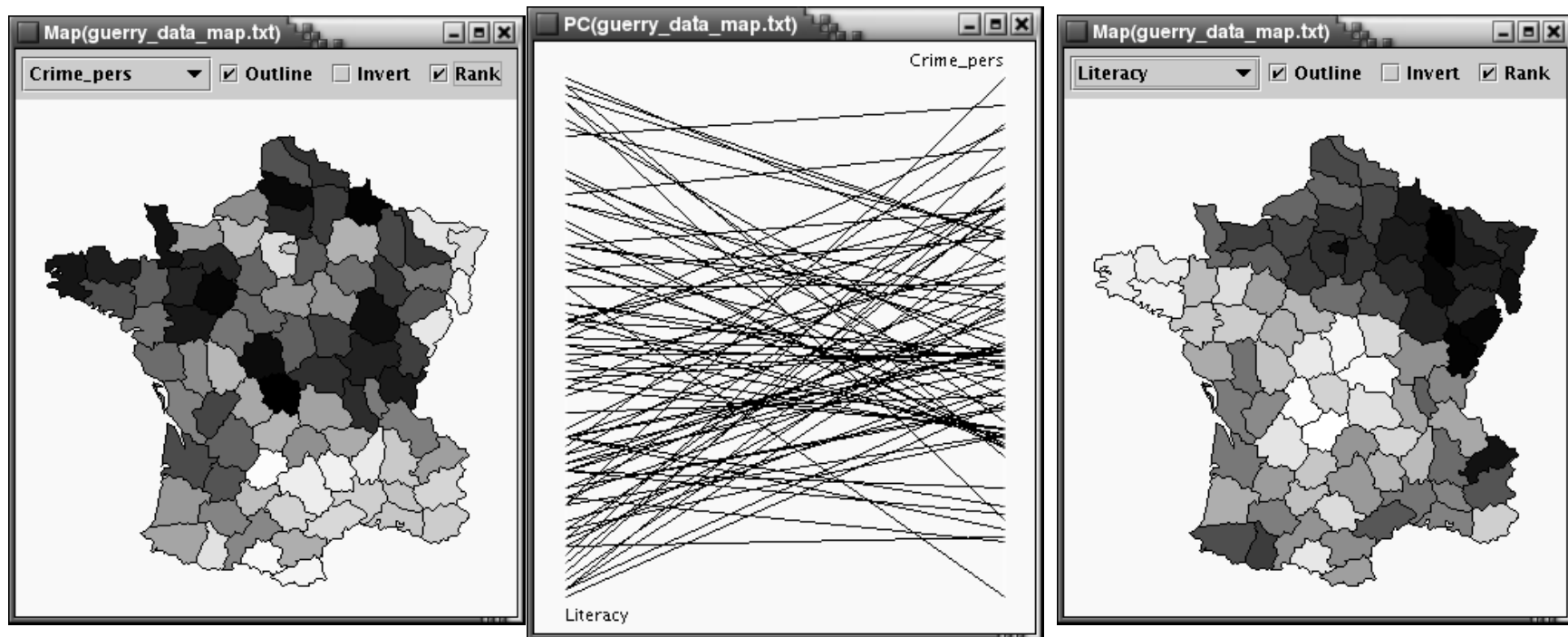


### Population per Suicide

# Guerry's analyses

Relate variables by comparing maps and ranked lists ($1^{st}$ || coordinate plot)

- Conclusion: no clear relation between crime and literacy



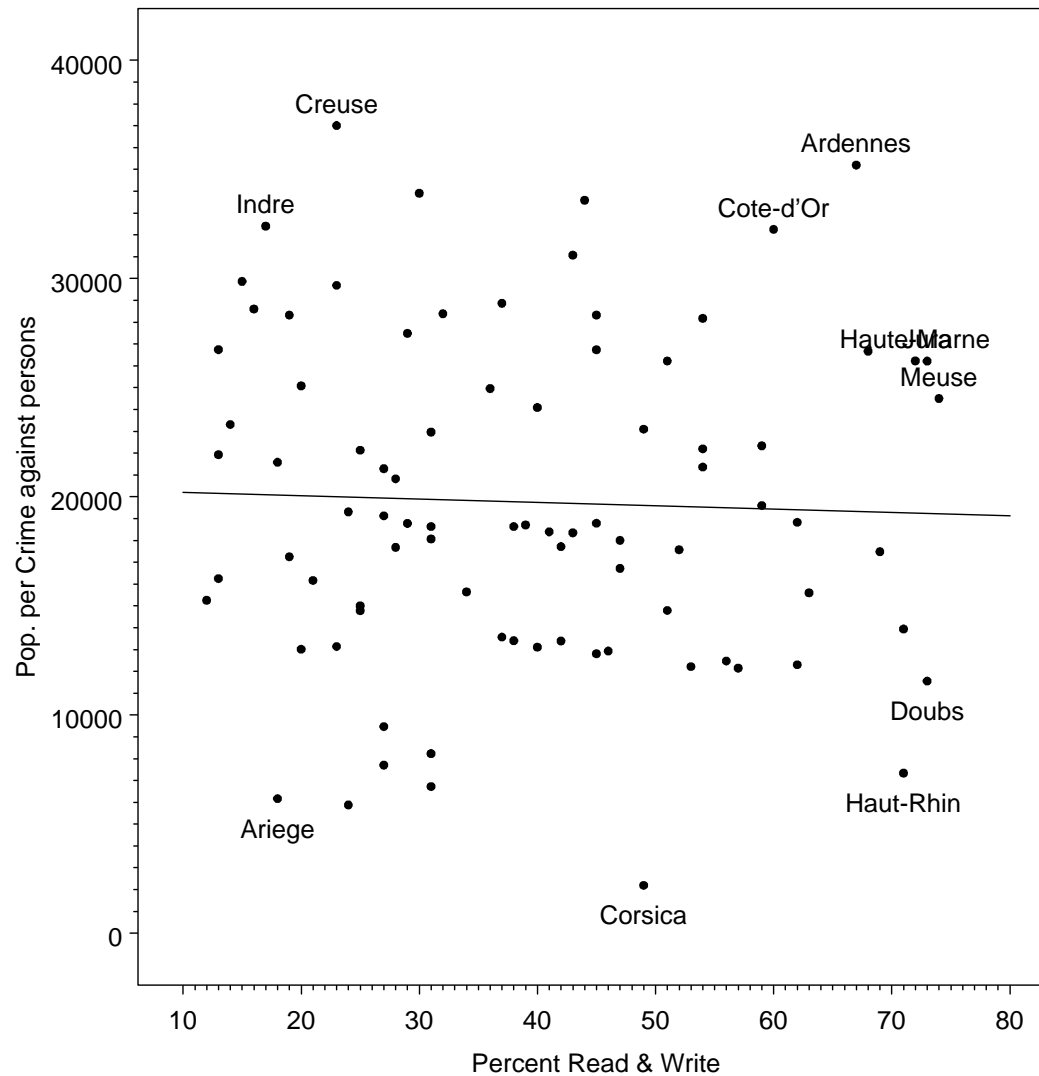Literacy                              Ranked lists                    Crimes against persons

- Similar analyses for other variables (suicide, illegitimate births, ...)
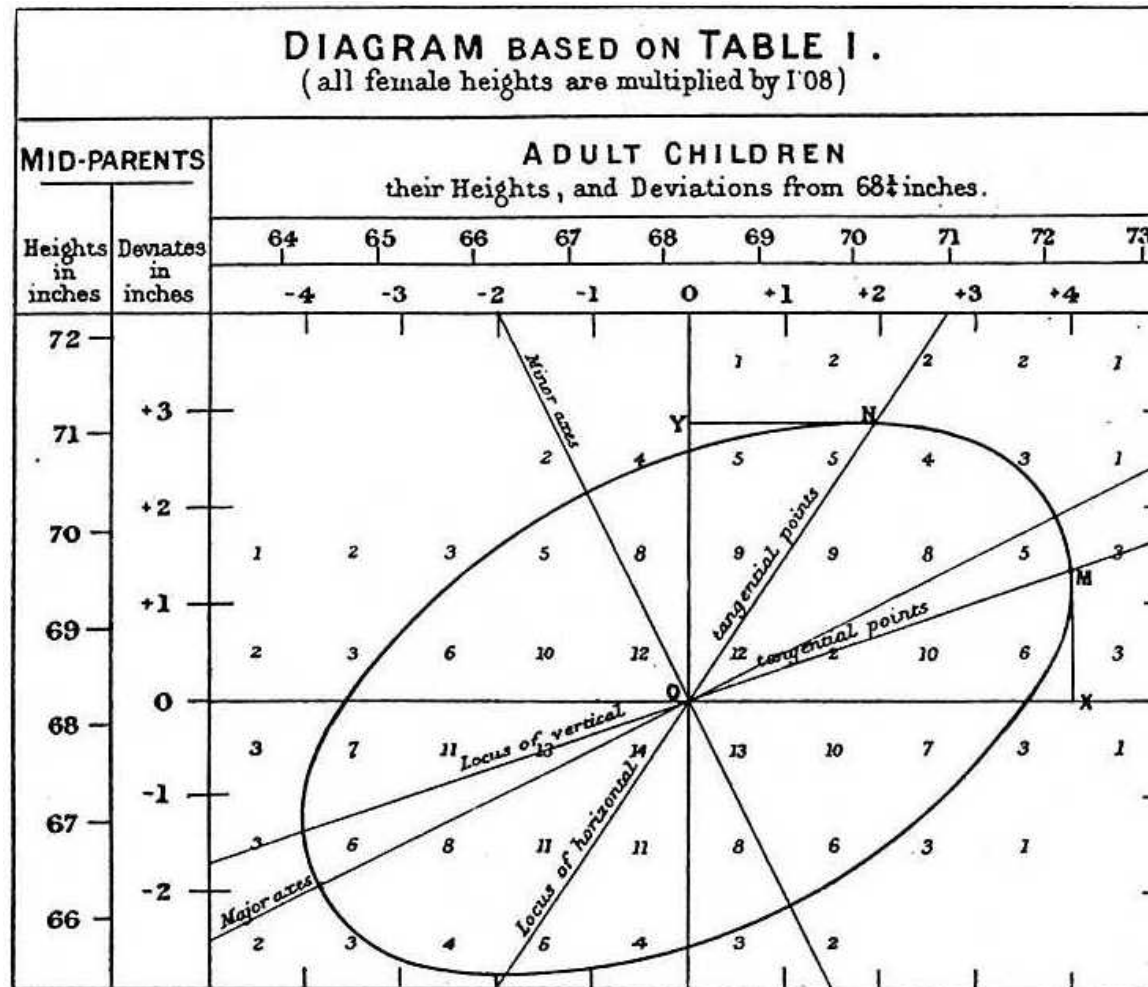
**Graphical methods for multivariate data**

■ **Bivariate displays**: Bivariate displays can be enhanced to show statistical relations more clearly and effectively

- ■ Scatterplots with data (concentration) ellipses and smoothed (loess) curves
- ■ Scatterplot matrices
- ■ Corrgrams and visual thinning

■ **Reduced-rank displays**: Multivariate visualization techniques can show the statistical data in simple ways, using dimension reduction techniques.

- ■ Biplots - show variables and observations in space accounting for greatest variance
- ■ Canonical discriminant plots - show variables and observations in space accounting for greatest between-group variation

■ **HE plots**: Visualization for Multivariate Linear Models

## Bivariate plots: Points and visual summaries



Scatterplot with linear regression line

## The Data Ellipse: Galton's Discovery



DIAGRAM BASED ON TABLE I.
(all female heights are multiplied by 1·08)

Pearson (1920): "... one of the most noteworthy scientific discoveries arising from pure analysis of observations."

11

© Michael Friendly

## The Data Ellipse: Details

■ **Visual summary for bivariate marginal relations**

■ **Shows**: means, standard deviations, correlation, regression line(s)

■ **Defined**: set of points whose squared Mahalanobis distance $\leq c^2$,

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

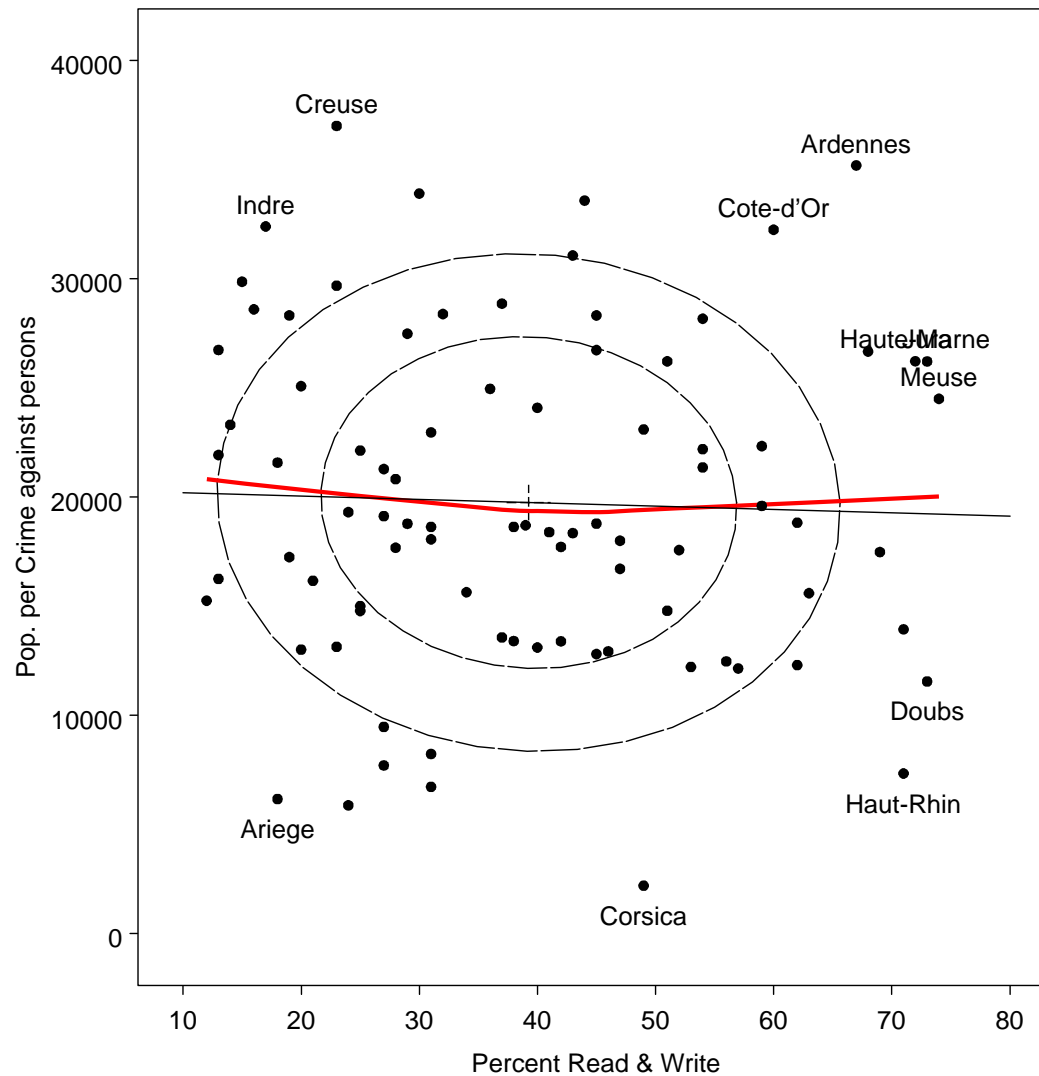$\boldsymbol{S}$ = sample variance-covariance matrix

■ **Radius**: when $\boldsymbol{y}$ is approx. bivariate normal, $D^2(\boldsymbol{y})$ 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: $\bar{y} \pm 1s$
● $c^2 = \chi^2_2(0.68) = 2.28$: 1 std. dev bivariate ellipse
● Small samples: $c^2 \approx 2F_{2,n-2}(1 - \alpha)$

■ **Construction**: Transform the unit circle, $\mathcal{U} = (\sin\boldsymbol{\theta}, \cos\boldsymbol{\theta})$,

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

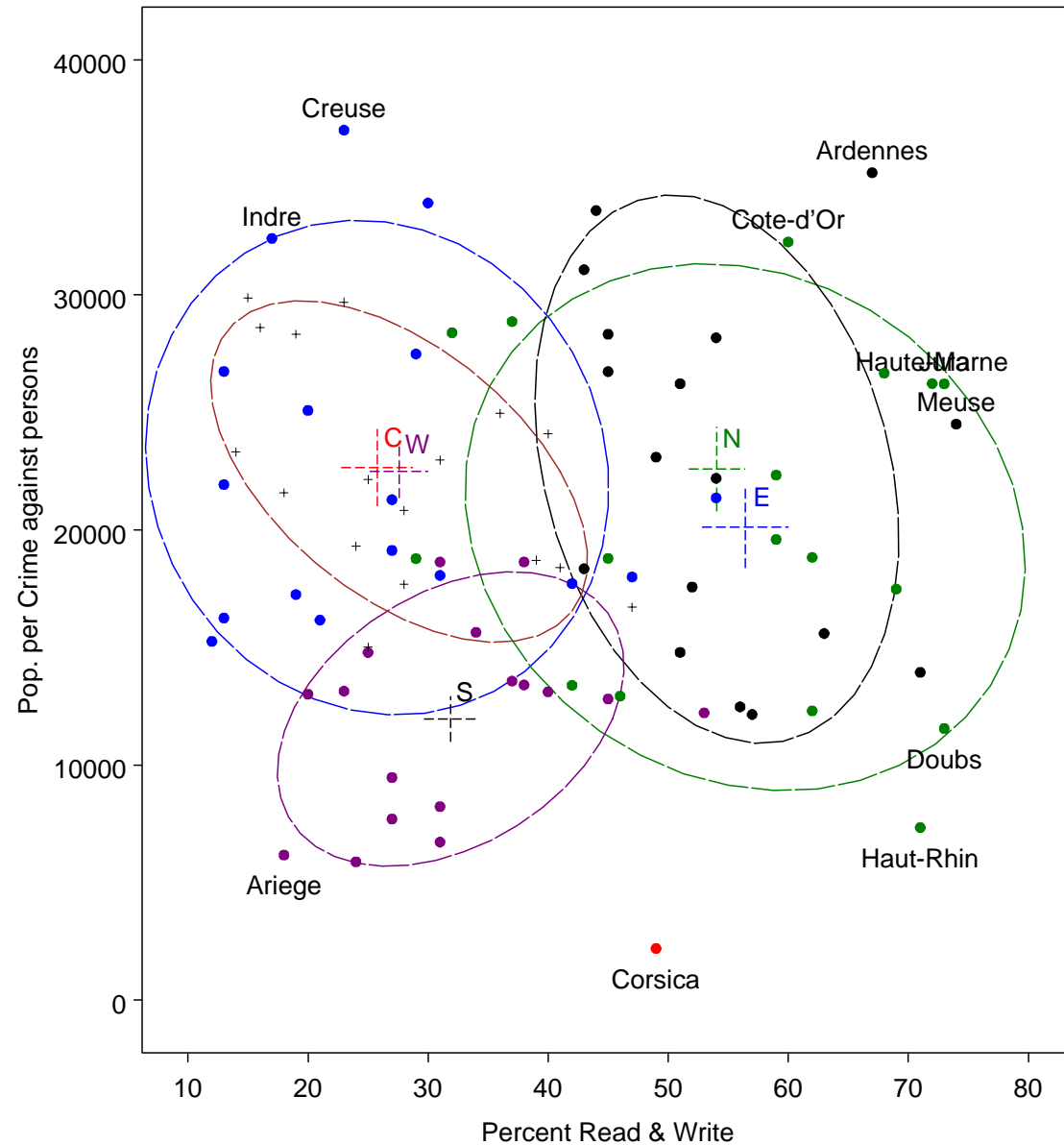$\boldsymbol{S}^{1/2}$ = any "square root" of $\boldsymbol{S}$ (e.g., Cholesky)

■ **Robust version**: Use robust covariance estimate (MCD, MVE)

■ **Nonparametric version**: Use kernel density estimation

© Michael Friendly

# Bivariate plots: Data ellipse and smoothing



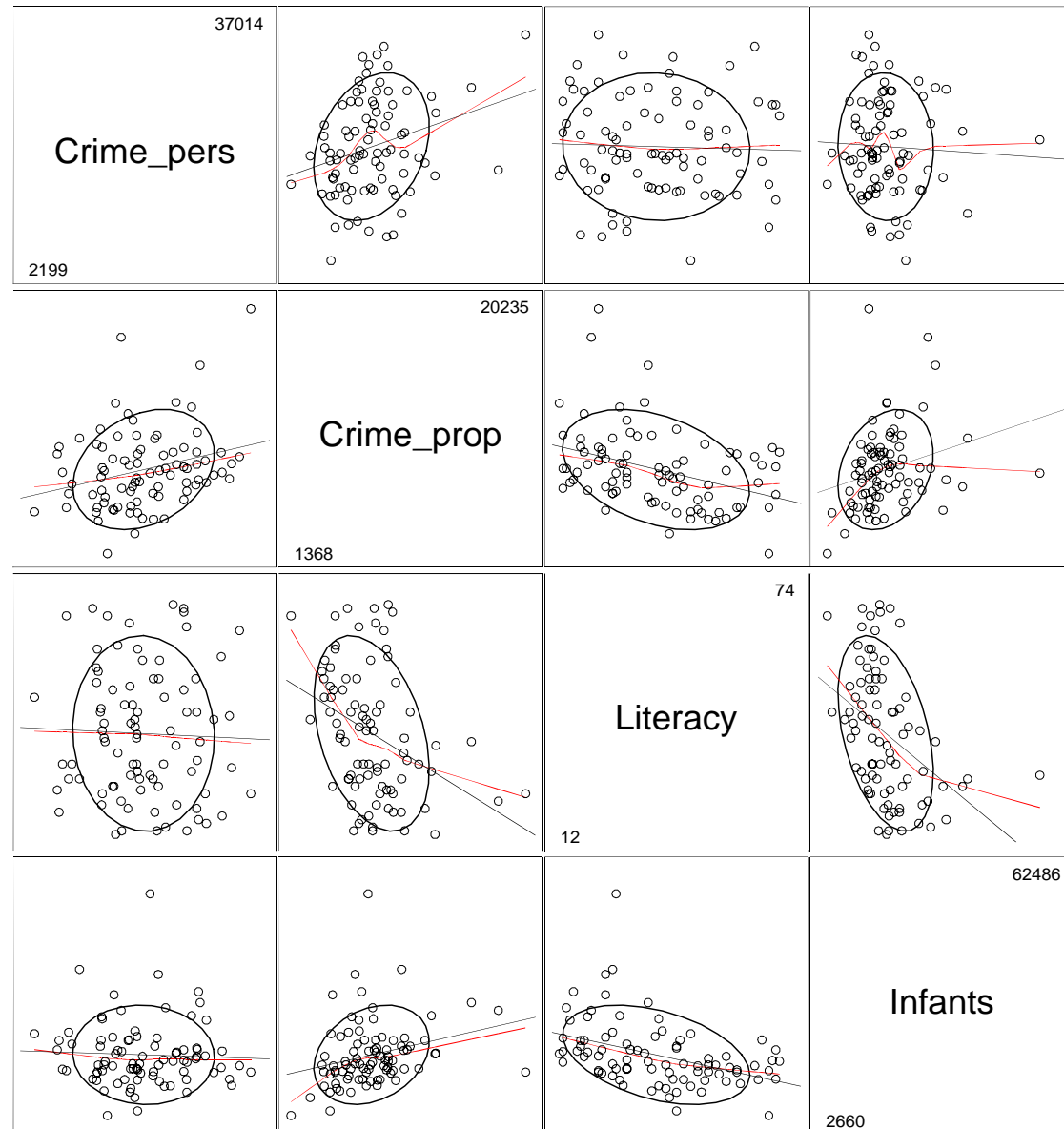Scatterplot with 68% data ellipse and smoothed (loess) curve
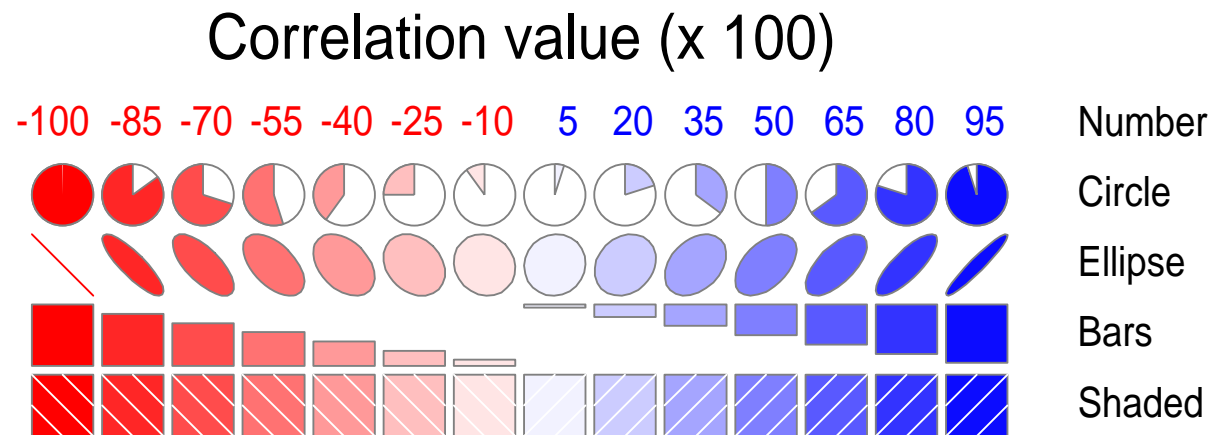
# Bivariate plots: Region differences

# Bivariate plots: Scatterplot matrices

## Corrgrams— Correlation matrix displays

- How to show a correlation matrix for different purposes? (Friendly, 2002)

- Render a correlation to depict sign and magnitude (tasks: lookup, comparison, detection)
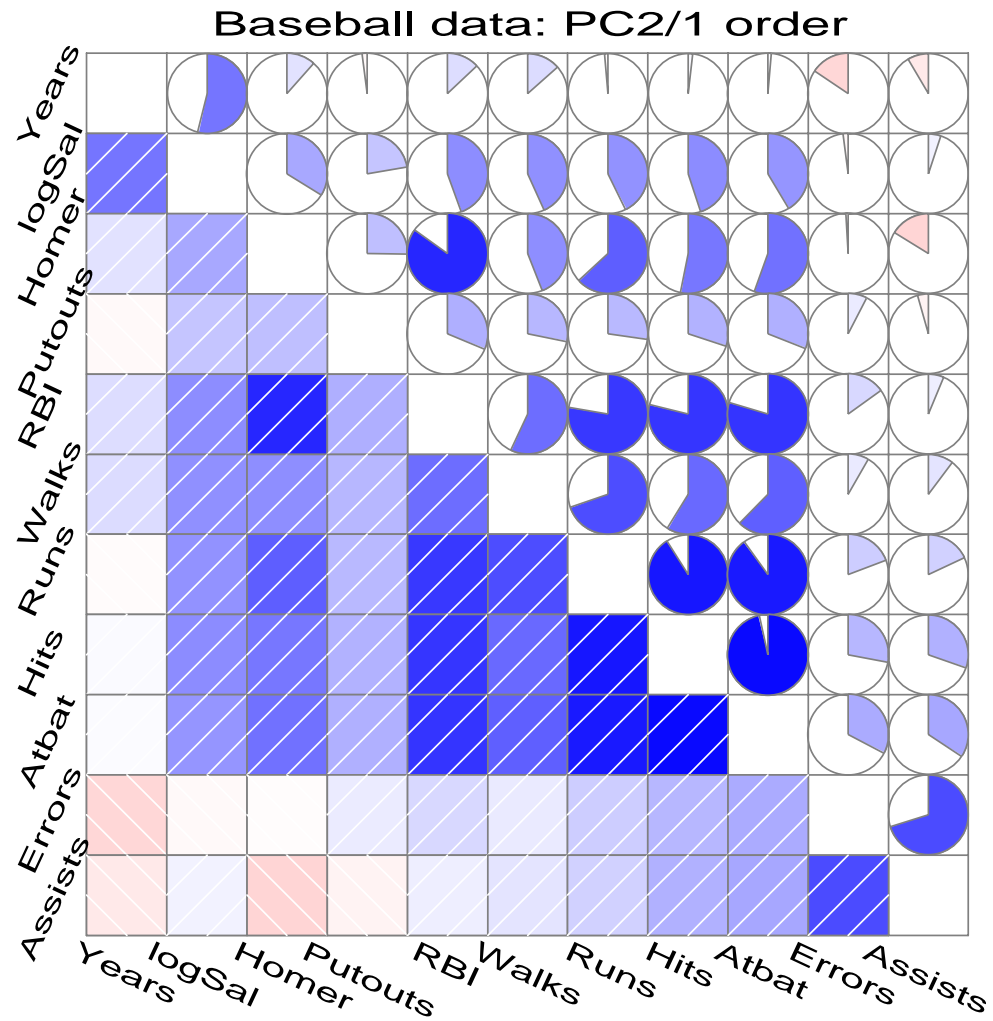
### Correlation value (x 100)

| -100 | -85 | -70 | -55 | -40 | -25 | -10 | 5 | 20 | 35 | 50 | 65 | 80 | 95 | |
|------|-----|-----|-----|-----|-----|-----|---|----|----|----|----|----|----| |



Number
Circle
Ellipse
Bars
Shaded

Task-specific renderings:

| **Task** | Lookup | Comparison | Detection |
|----------|--------|------------|-----------|
| **Rendering** | Number | Circle | Shading |

# Corrgrams— Rendering

Baseball data: (lower) Patterns vs. (upper) comparison



Baseball data: PC2/1 order

# Corrgrams— Variable ordering

Baseball data: (a) alpha vs. (b) correlation ordering (Friendly and Kwan, 2003)



See: http://www.math.yorku.ca/SCS/sasmac/corrgram.html
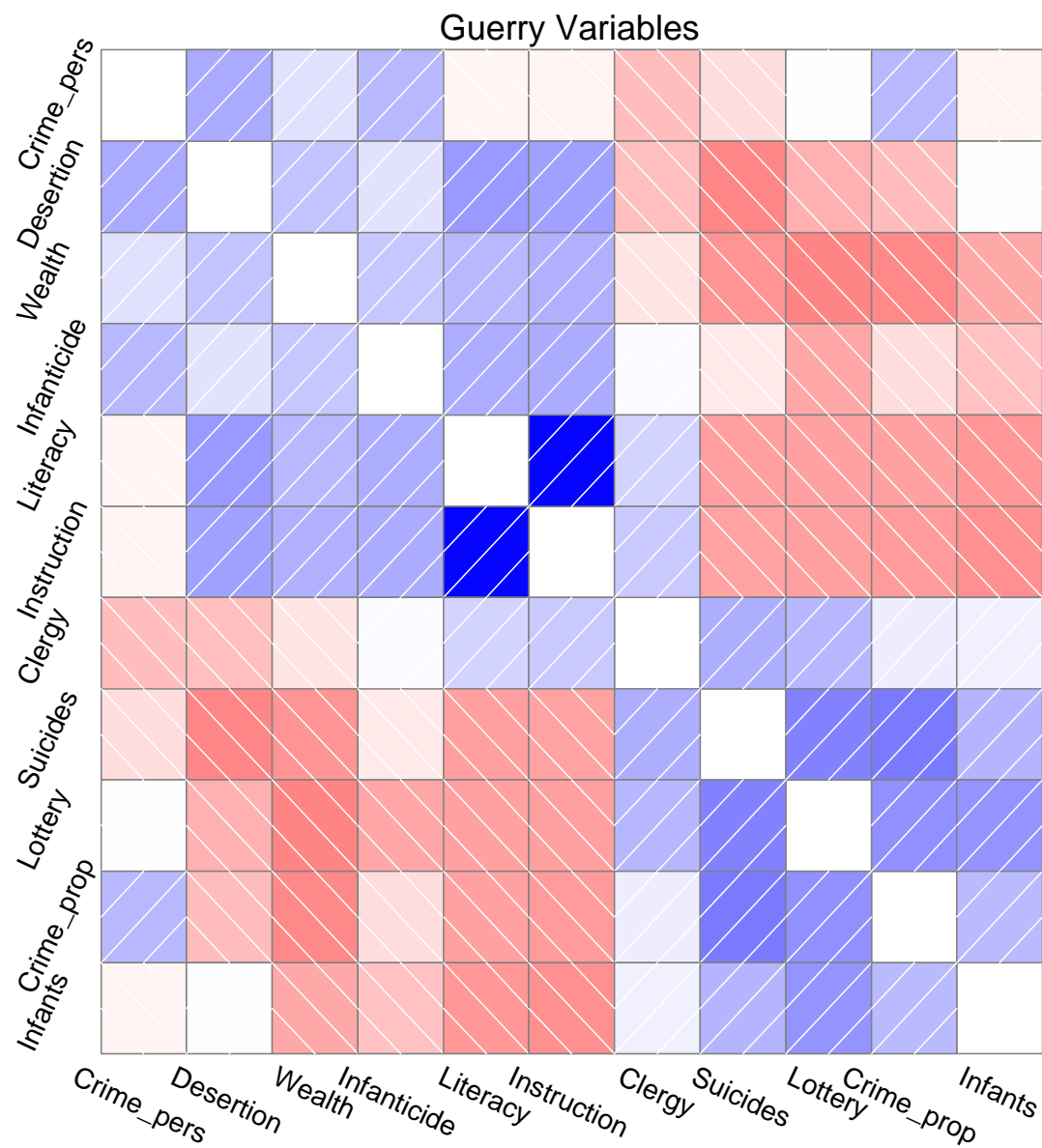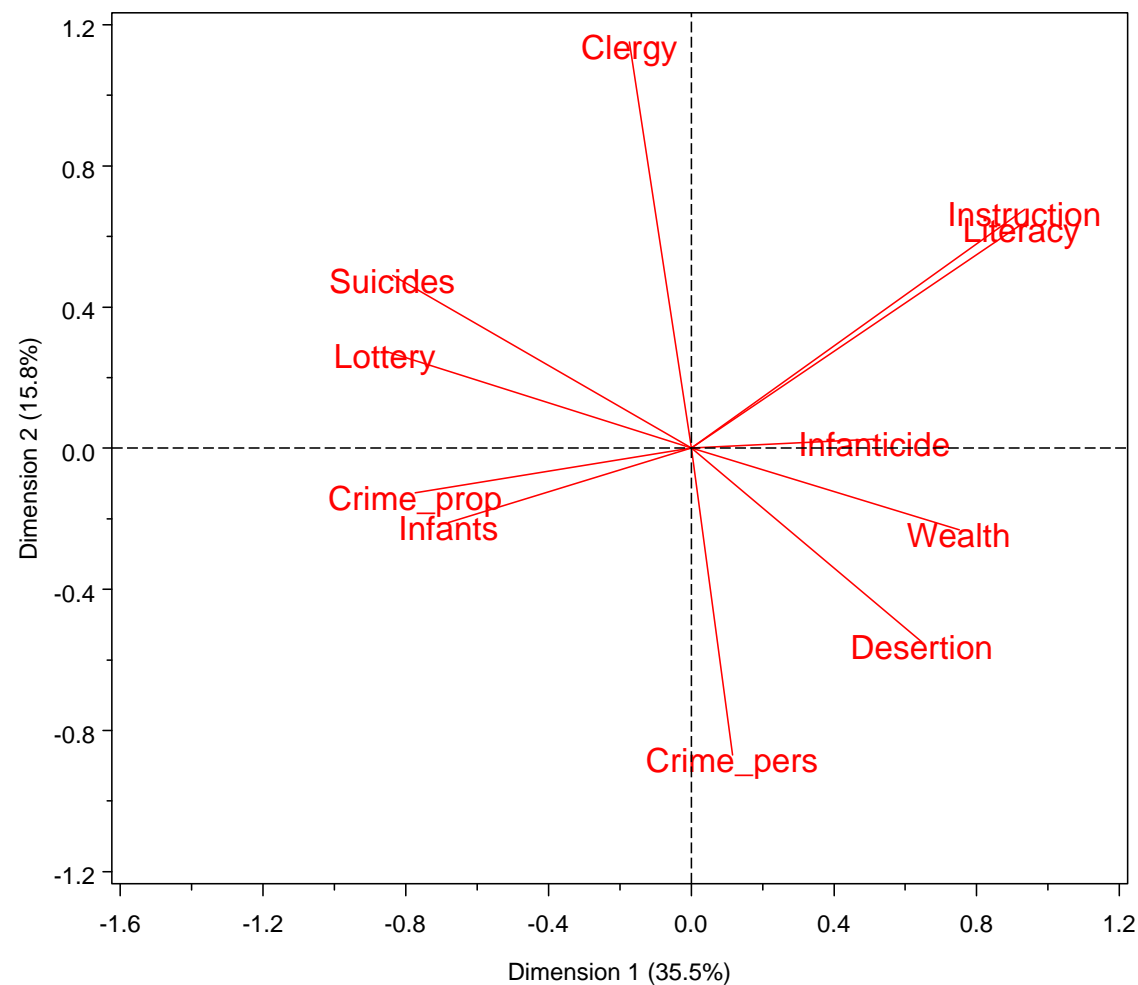
## Corrgrams— Variable ordering

■ Reorder variables to show similarities: PC1 or angles (PC2/PC1)
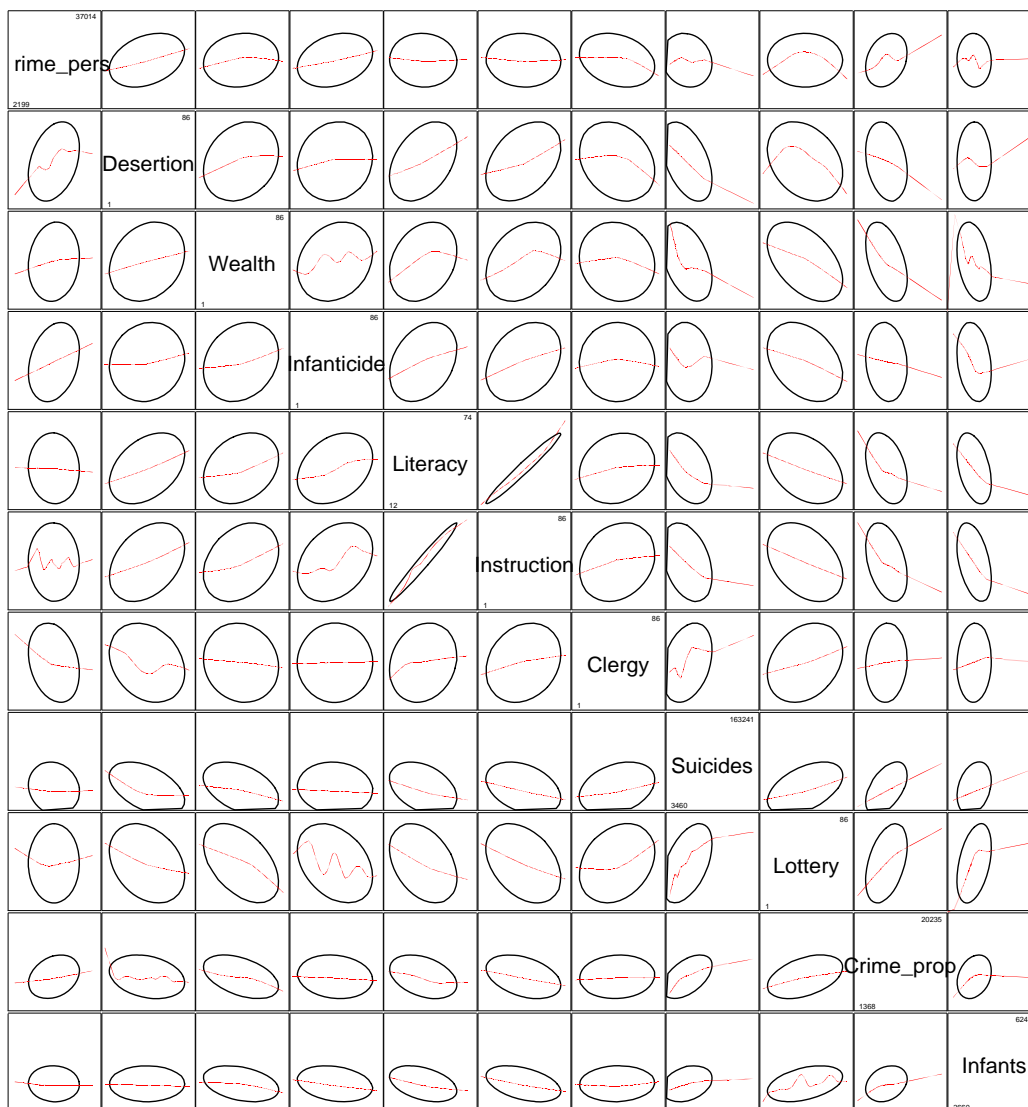
# Corrgrams— Guerry data



Guerry Variables

**Guerry data— Variable ordering**

# Visual thinning: Minimal summaries for large data sets

Guerry data: schematic scatterplot matrix: 68% data ellipse $+$ loess smooth

## Multivariate analyses: Reduced rank displays

- Multivariate visualization techniques can show the statistical data in simple ways, using dimension reduction techniques.

  - Biplots - show variables and departments in space accounting for greatest variance
  - Canonical discriminant plots - show variables and departments in space accounting for greatest between-region variation

- Can try to show geographic location by color coding or other visual attributes.

  - Color code by region
  - Show data ellipse to summarize regions

- $\rightarrow$ **Data-centric displays**: The multivariate data is shown directly; geographic relations indirectly
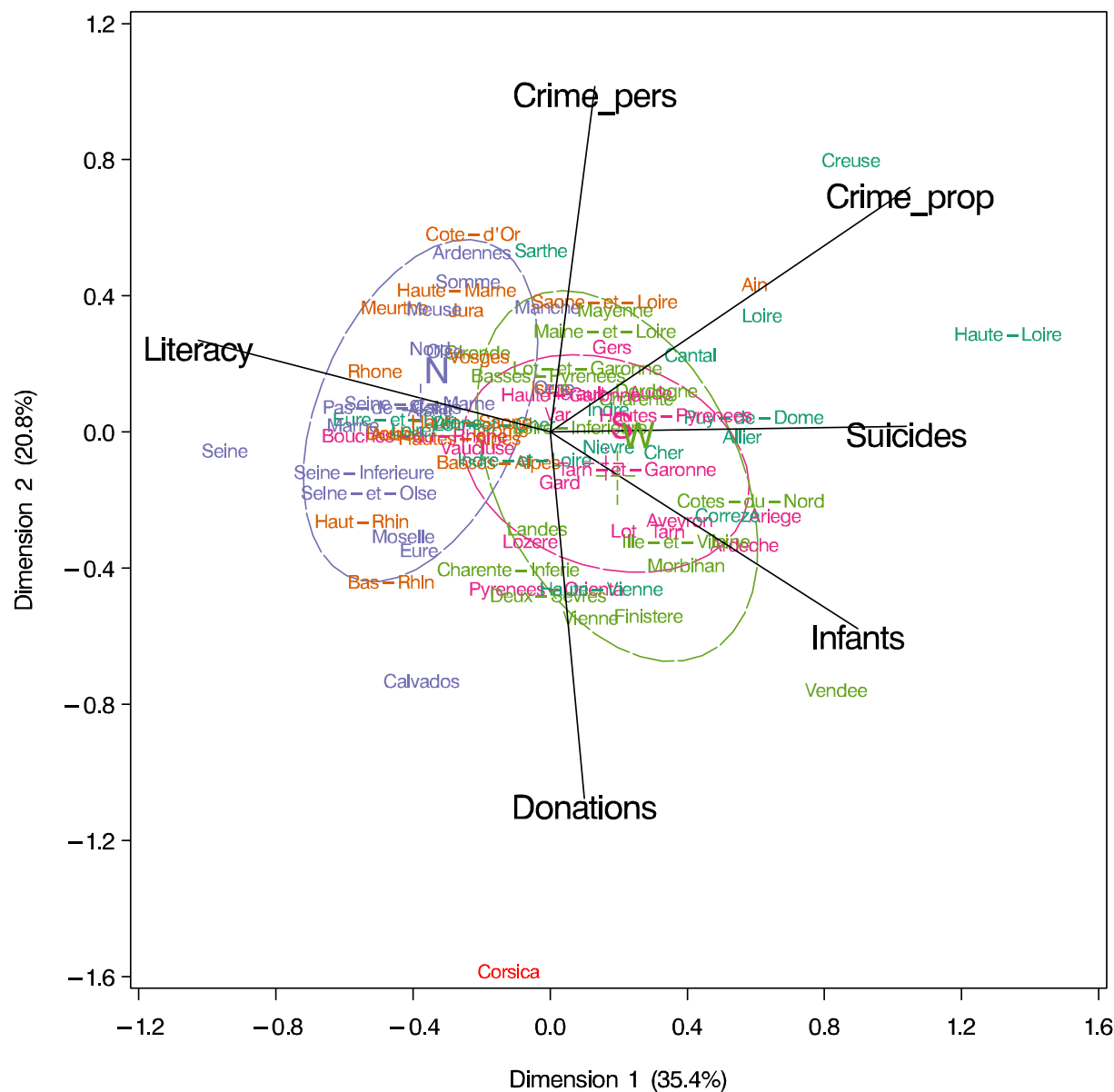
## Biplots

- Biplots represent both variables (attributes) and observations (departments) in the same plot— a low-rank (2D) approximation to a data matrix (Gabriel, 1971)

$$\boldsymbol{Y}^{\star} = \boldsymbol{Y} - Y_{..} \approx \boldsymbol{A}\boldsymbol{B}^{\mathsf{T}} = \sum_{k=1}^{d} \boldsymbol{a}_k \boldsymbol{b}_k^{\mathsf{T}}$$
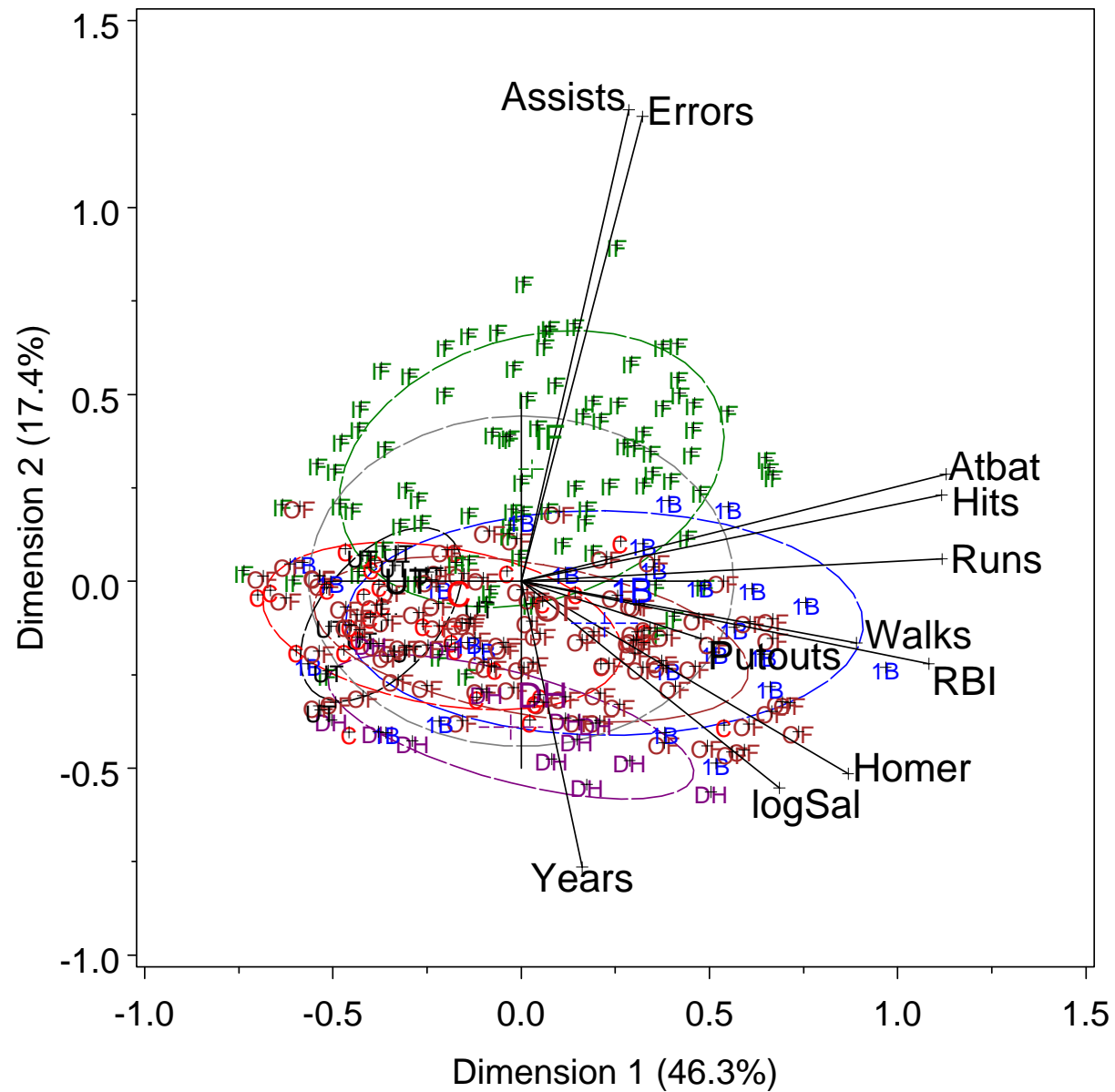
  - Variables are usually represented by vectors from origin (mean)
  - Observations are usually represented by points
  - Can show clusters of observations by data ellipses

- Properties:

  - Angles between vectors show correlations ($r \approx \cos(\theta)$)
  - Length of variable vectors $\sim$ % variance accounted for
  - $y_{ij} \approx \boldsymbol{a}_i^{\mathsf{T}} \boldsymbol{b}_j$: projection of observation on variable vector
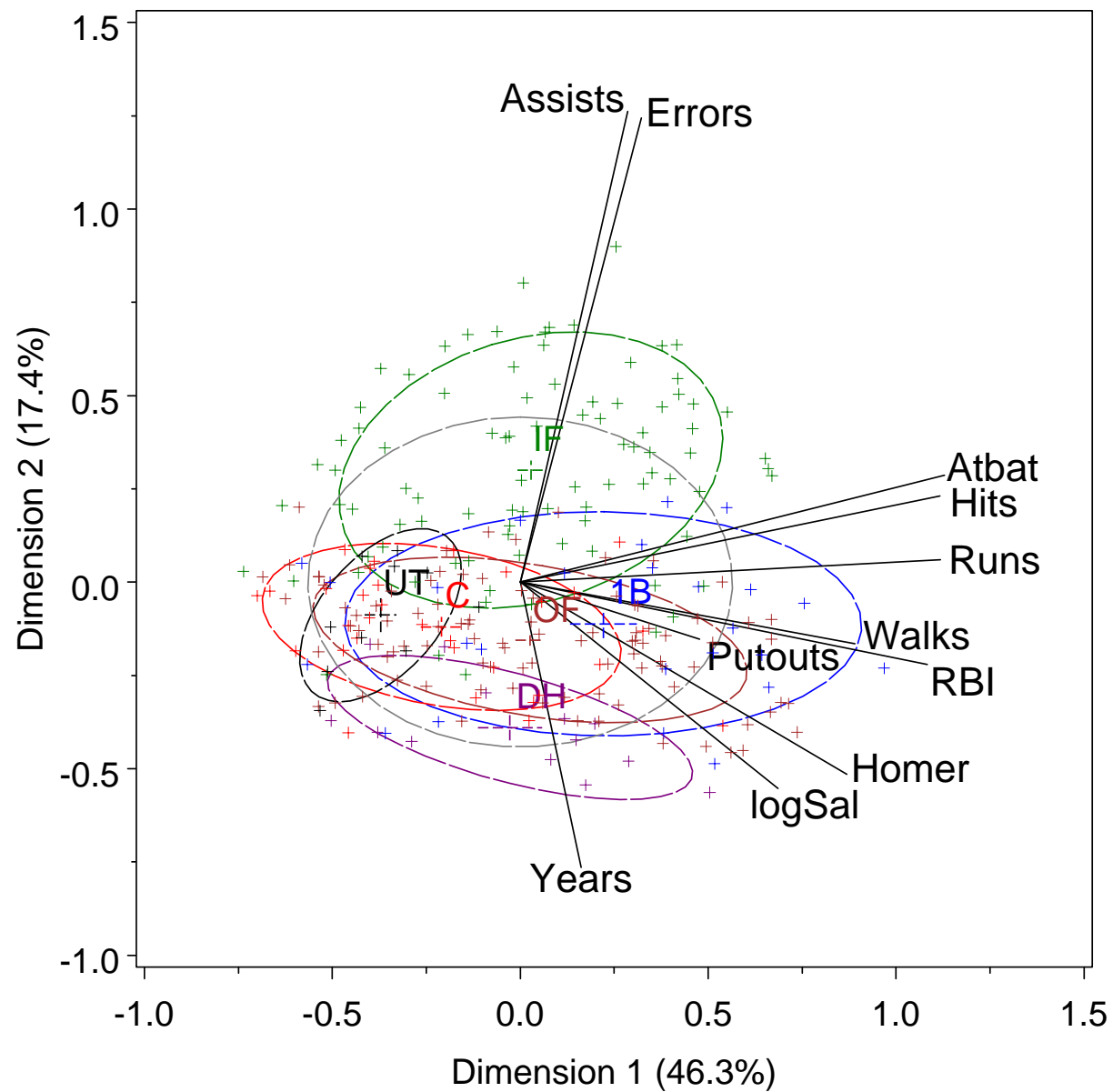  - Dimensions are uncorrelated overall (but not necessarily within group)

ⓒ Michael Friendly

# Biplots: Guerry data

# Biplots: Baseball data



Biplot of baseball data, players labeled by position

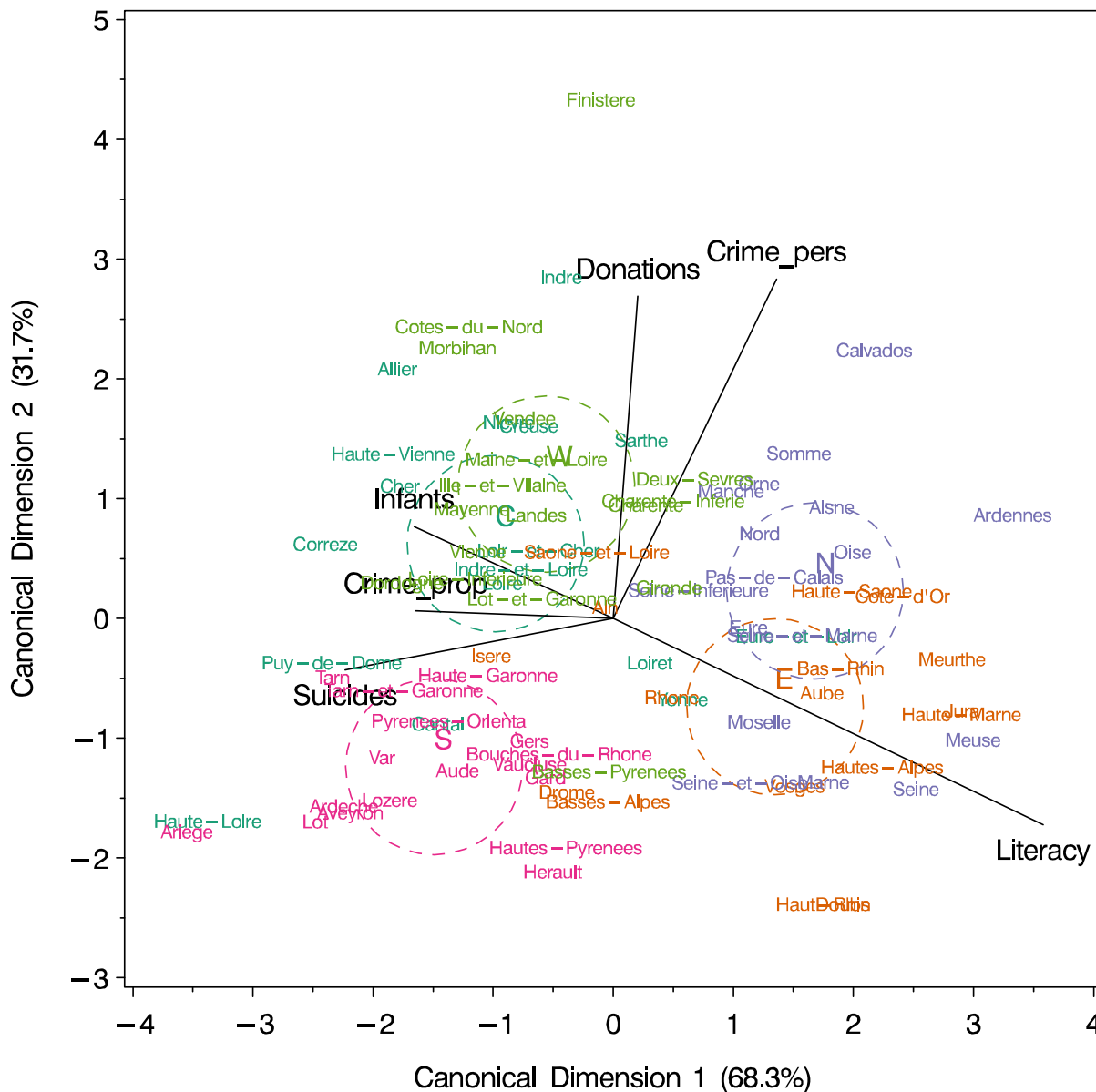Biplot of baseball data, with data ellipses by position

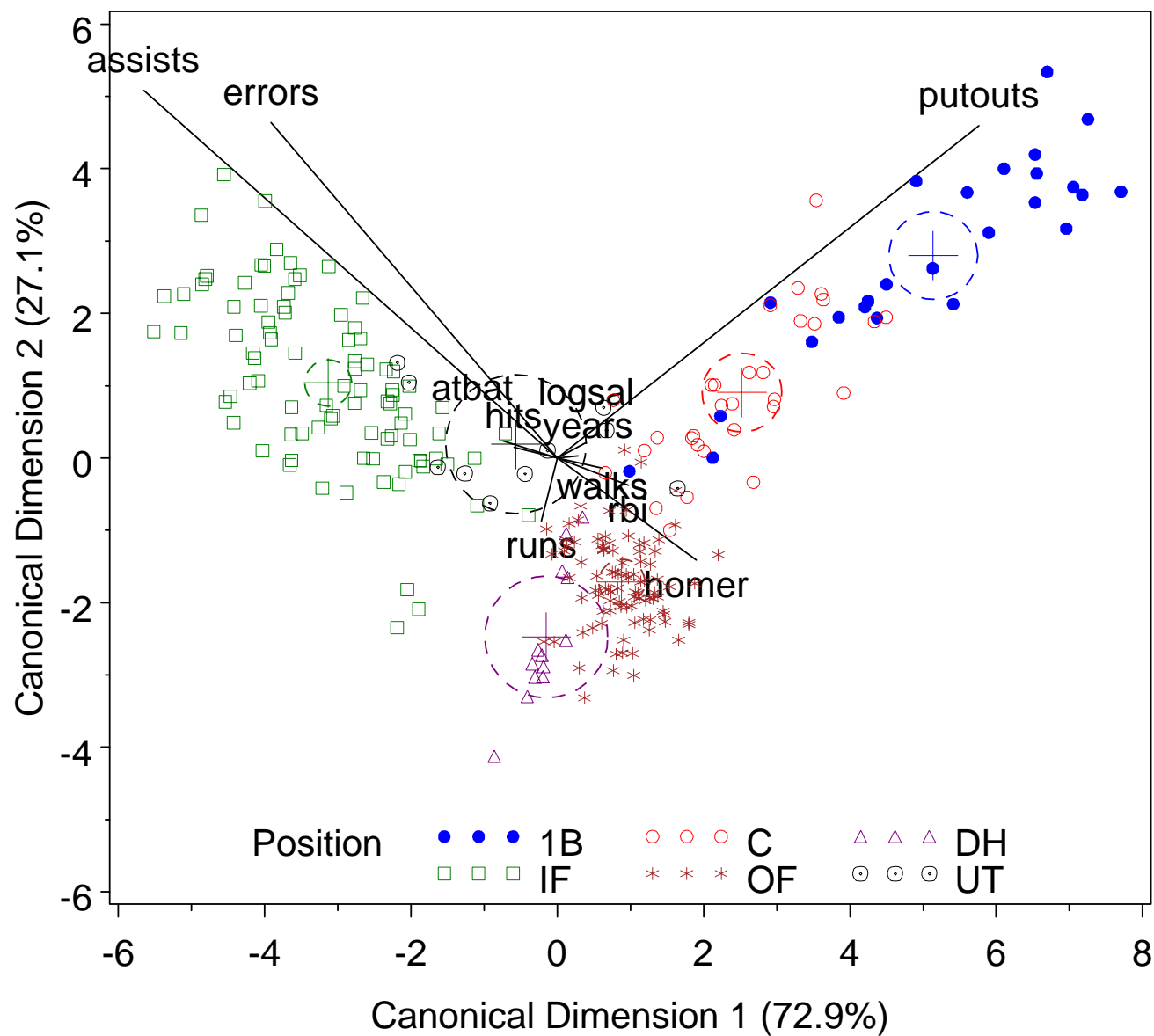Biplot of baseball data, with data ellipses by position (thinned)

## Canonical discriminant plots

- Project the variables into a low-rank (2D) space that maximally discrimates among regions (Friendly, 1991)

  - Visual summary of a MANOVA
  - Canonical dimensions are linear combinations of the variables with maximum univariate $F$-statistics.
  - Vectors from the origin (grand mean) for the observed variables show the correlations with the canonical dimensions

- Properties:

  - Canonical variates are uncorrelated
  - Circles of radius $\sqrt{\chi_2^2(1-\alpha)/n_i}$ give confidence regions for group means.
  - Variable vectors show how variables discriminate among groups
  - Lengths of variable vectors $\sim$ contribution to discrimination

# Canonical discriminant plots: Guerry data, by Region
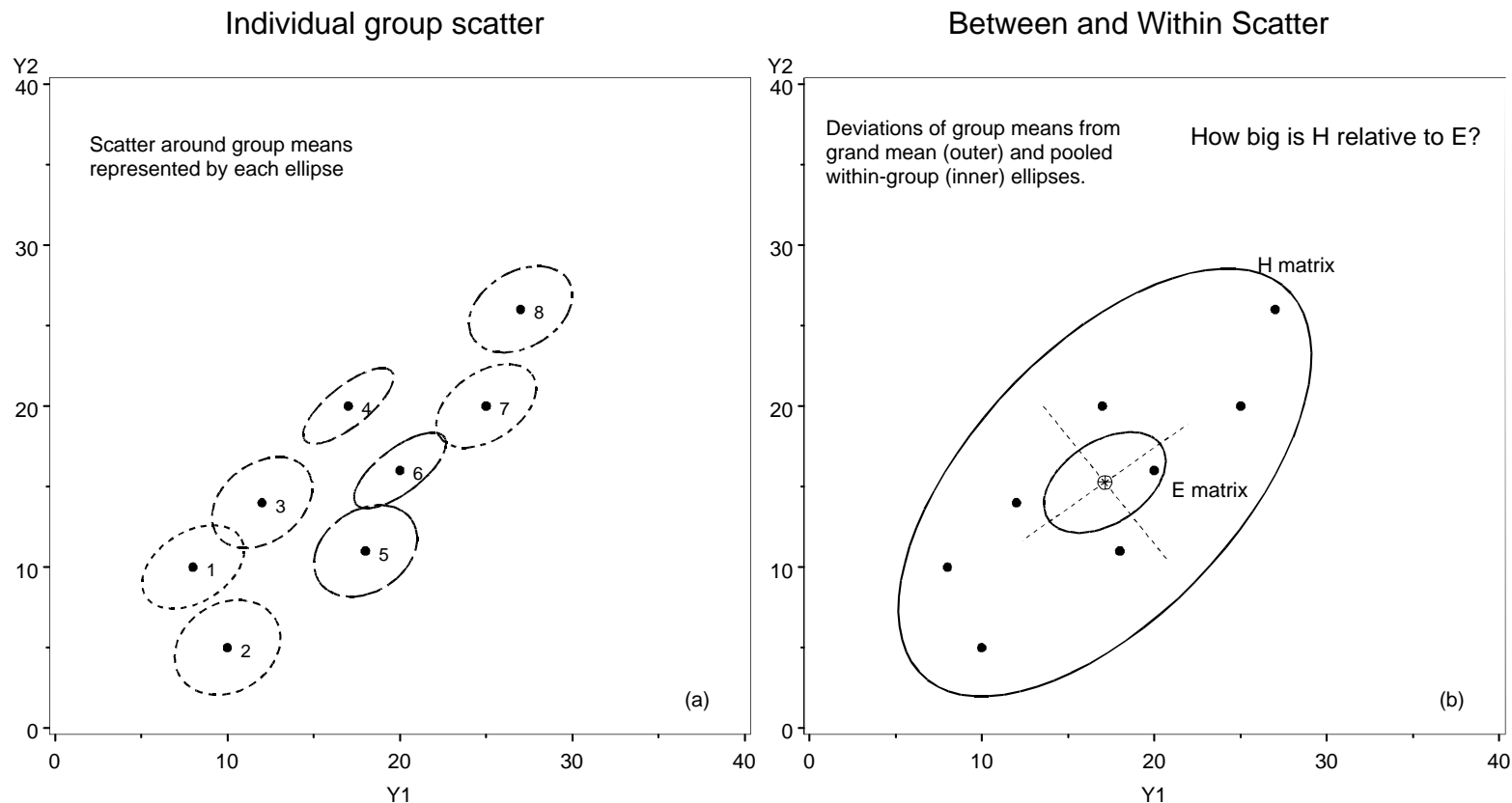
# CDA plots: Baseball data, by player position

## HE plots: Visualization for Multivariate Linear Models

■ How are $p$ responses, $\boldsymbol{Y} = (\boldsymbol{y}_1, \boldsymbol{y}_2, \ldots, \boldsymbol{y}_p)$ related to $q$ predictors, $\boldsymbol{X} = (\boldsymbol{x}_1, \boldsymbol{x}_2, \ldots, \boldsymbol{x}_q)$? (Friendly, 2004a)

■ MANOVA: $\boldsymbol{X} \sim$ discrete factors
■ MMRA: $\boldsymbol{X} \sim$ quantitative predictors
■ MANCOVA, response surface models,

All the same MLM:

$$\underset{(n \times p)}{\boldsymbol{Y}} = \underset{(n \times q)}{\boldsymbol{X}} \underset{(q \times p)}{\boldsymbol{B}} + \underset{(n \times p)}{\boldsymbol{E}}$$

■ Analogs of univariate tests:

■ Explained variation: $MS_H \longmapsto (p \times p)$ covariance matrix, $\boldsymbol{H}$
■ Residual variation: $MS_E \longmapsto (p \times p)$ covariance matrix, $\boldsymbol{E}$
■ Test statistics: $F \longmapsto |\boldsymbol{H} - \lambda \boldsymbol{E}| = 0 \mapsto \lambda_1, \lambda_2, \ldots \lambda_s$

■ How big is $\boldsymbol{H}$ relative to $\boldsymbol{E}$ ?

■ Latent roots $\lambda_1, \lambda_2, \ldots \lambda_s$ measure the "size" of $\boldsymbol{H}$ relative to $\boldsymbol{E}$ in $s = \min(p, df_h)$ orthogonal directions.
■ Test statistics: Wilks' $\Lambda$, Pillai trace, Hotelling-Lawley trace, Roy's maximum root combine these into a single number
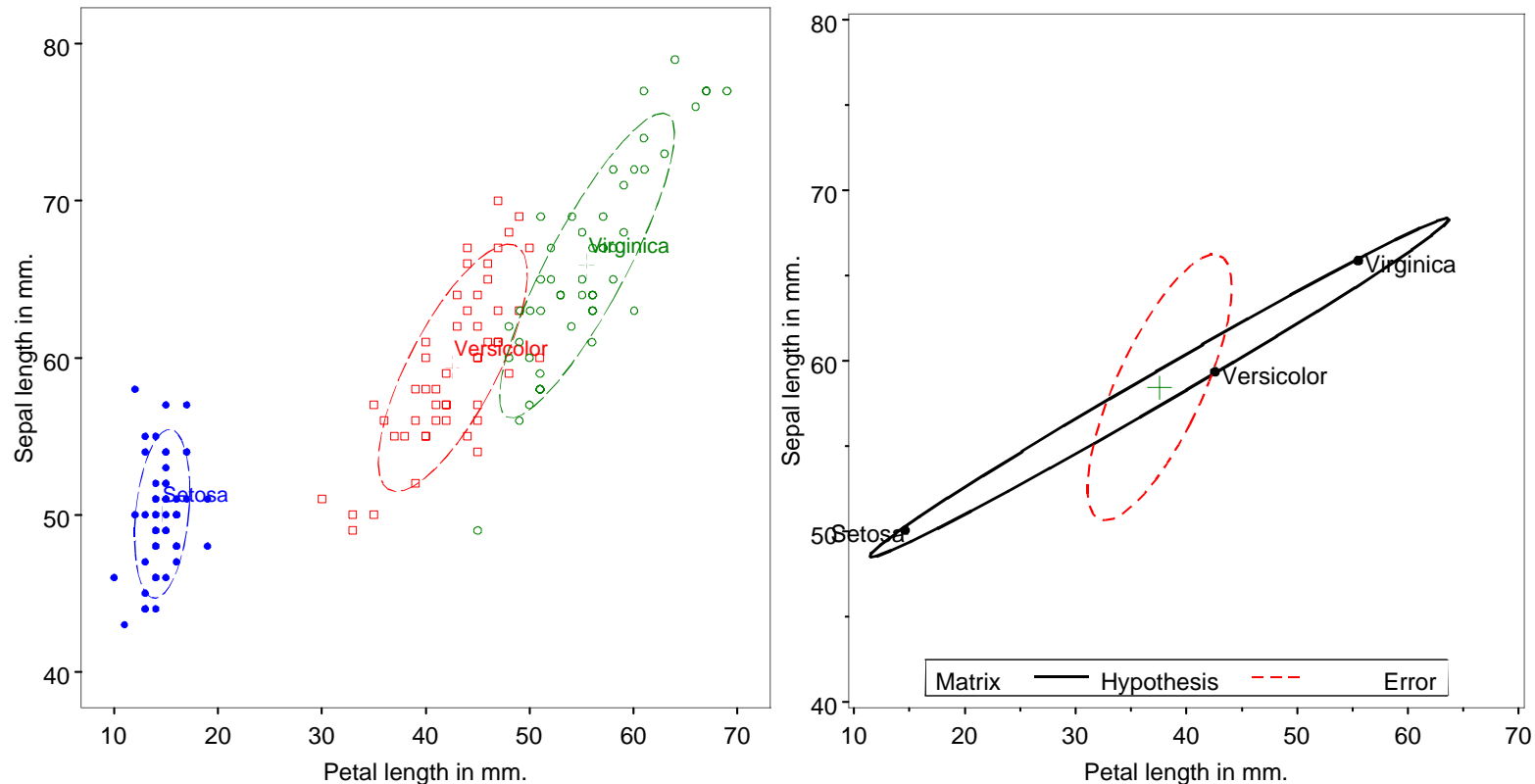
## HE plots: Visualization for Multivariate Linear Models

- **HE plot**: for two response variables, $(y_1, y_2)$, plot a $H$ ellipse and $E$ ellipse

- **HE plot matrices**: For all $p$ responses, plot an HE scatterplot matrix

- $\rightarrow$ **Shows**: size, dimensionality, and effect-correlation of $H$ relative to $E$ .

Individual group scatter

Between and Within Scatter

Scatter around group means represented by each ellipse

Deviations of group means from grand mean (outer) and pooled within-group (inner) ellipses.

How big is H relative to E?

H matrix

E matrix

(a)

(b)

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

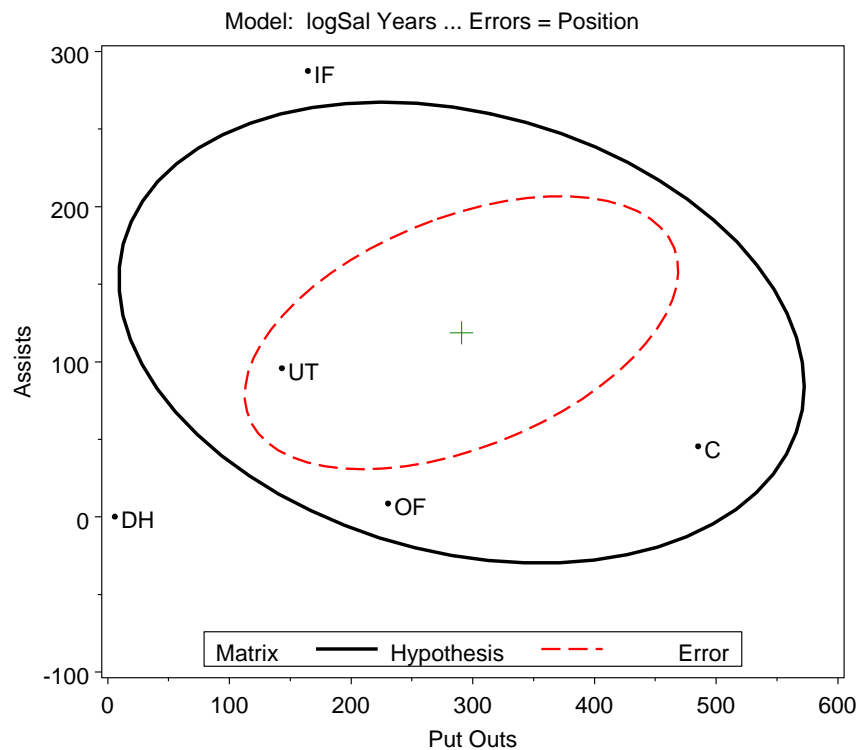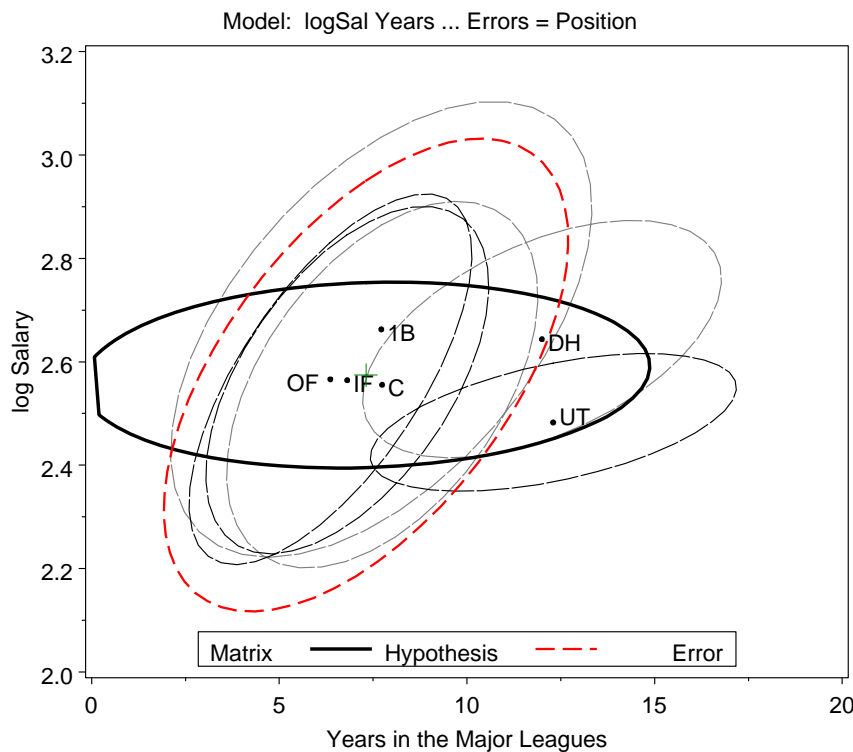## Simple example: Iris data
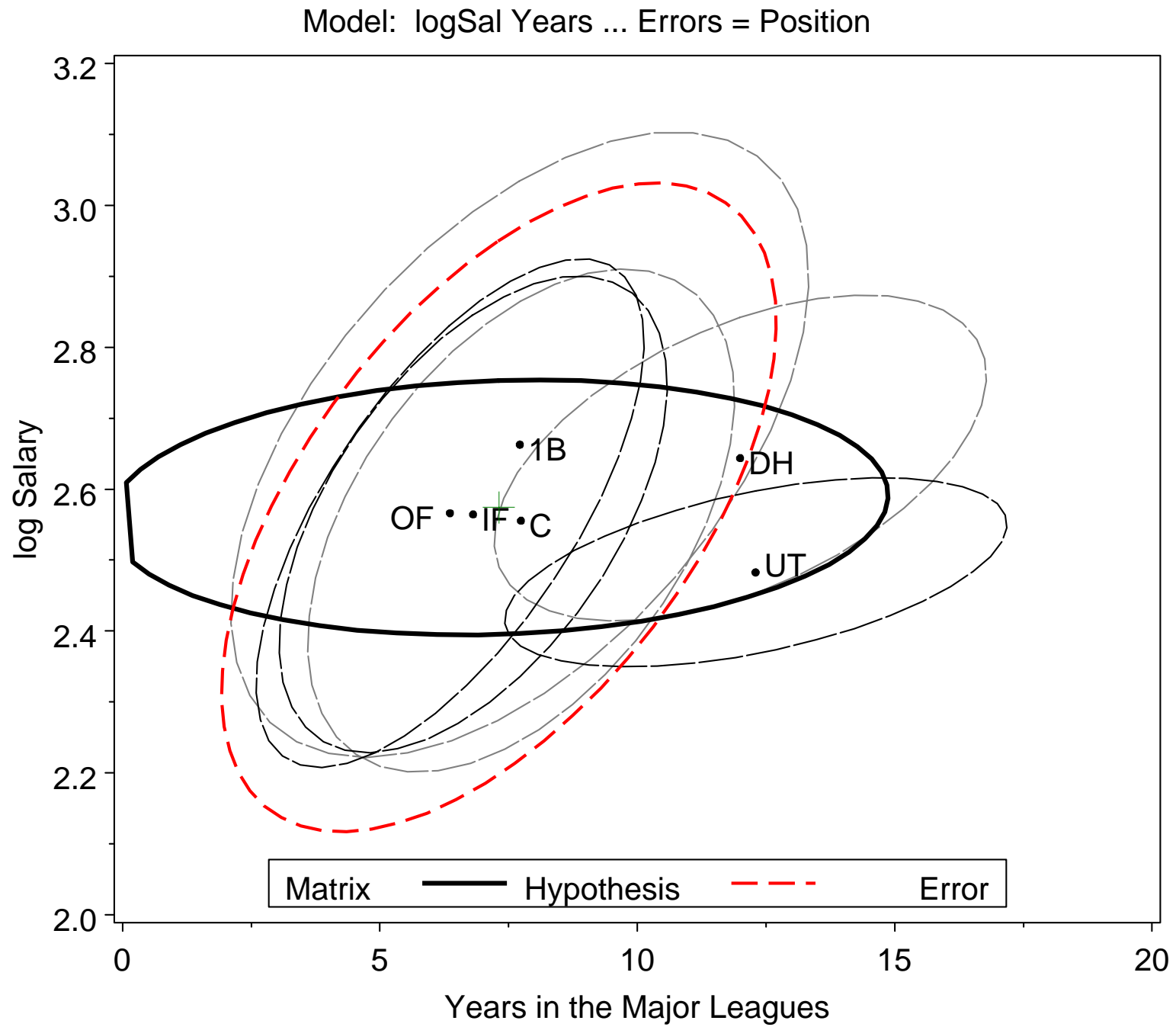


(a) Data ellipses and (b) $H$ and $E$ ellipses

- $H$ **ellipse**: Shows 2D covariation of *predicted* values (means)
- $E$ **ellipse**: Shows 2D covariation of *residuals*
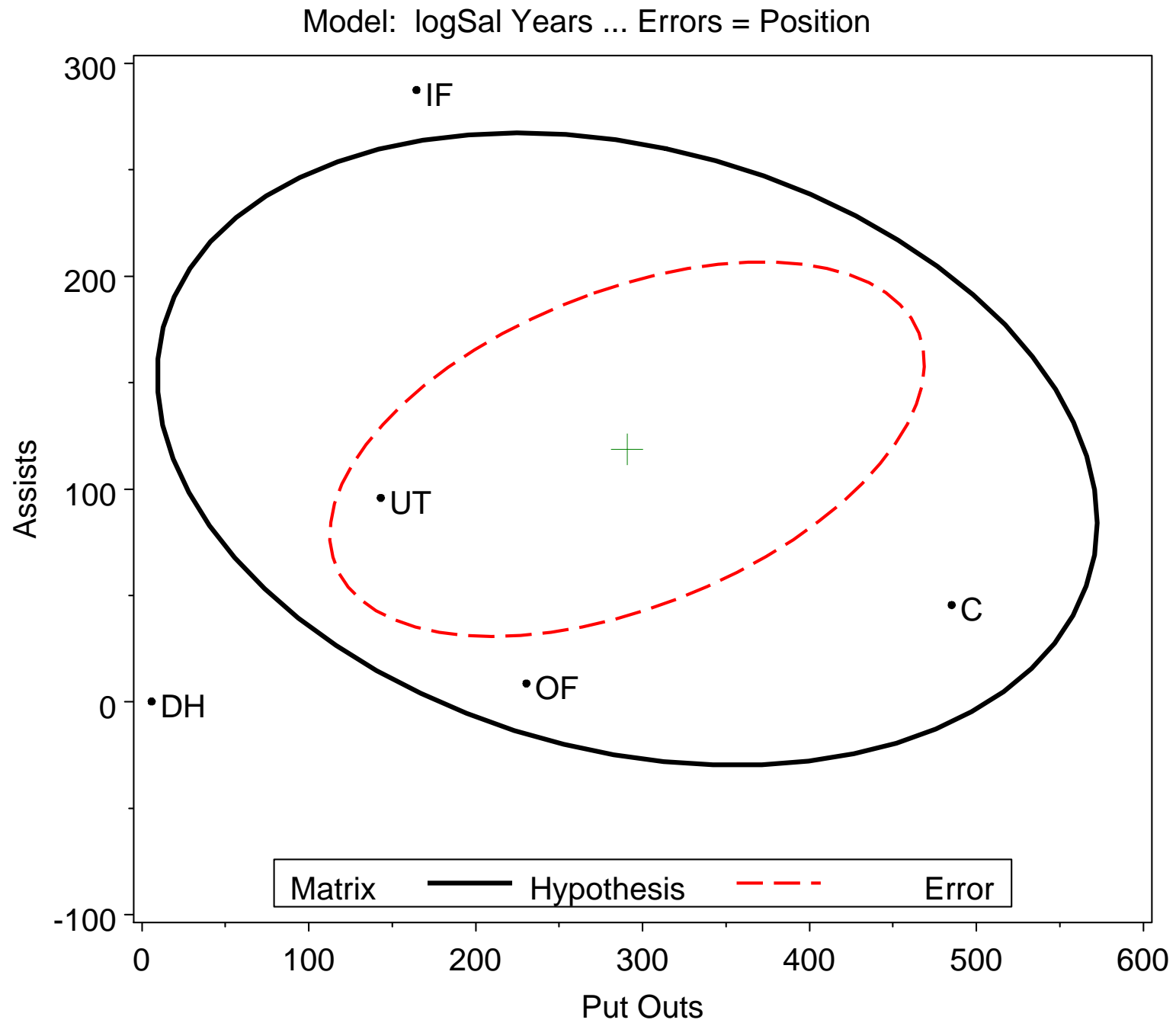- **points**: show group means on both variables

## Baseball data: Variation by position

■ How do relations among variables vary with player's position?

   ■ Fit MANOVA model,

     `(logSal Years Homer Runs Hits RBI Atbat Walks Putouts`
     `Assists Errors) = Position`

   ■ HE plots for selected pairs: (Years, logSal), (Putouts, Assists)

Model: logSal Years ... Errors = Position

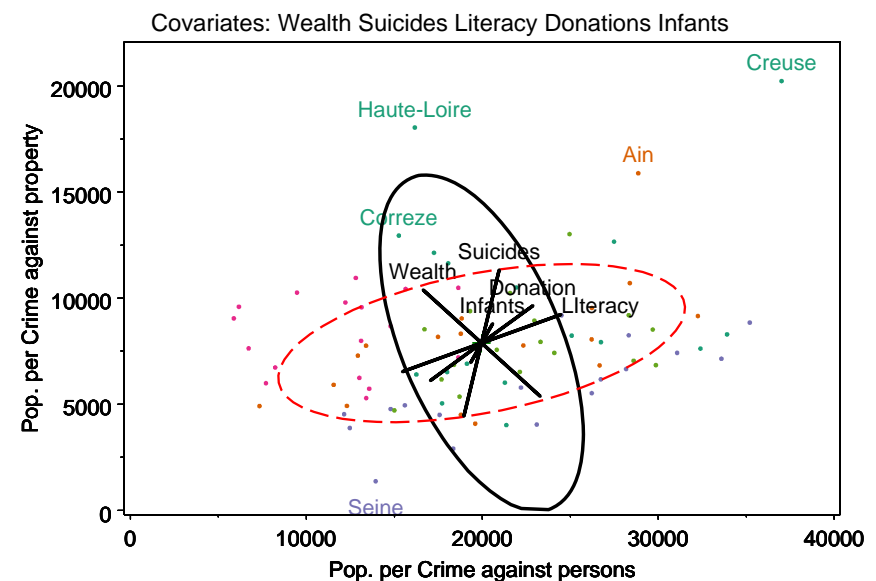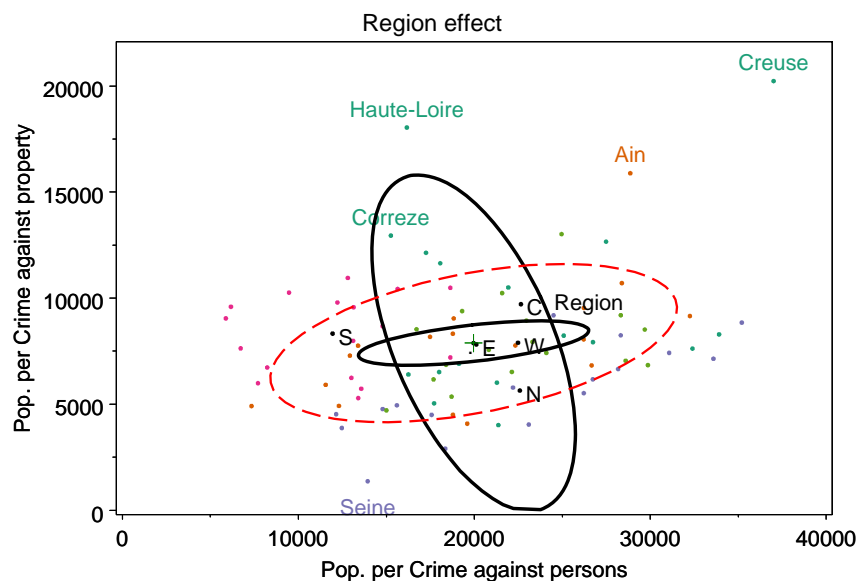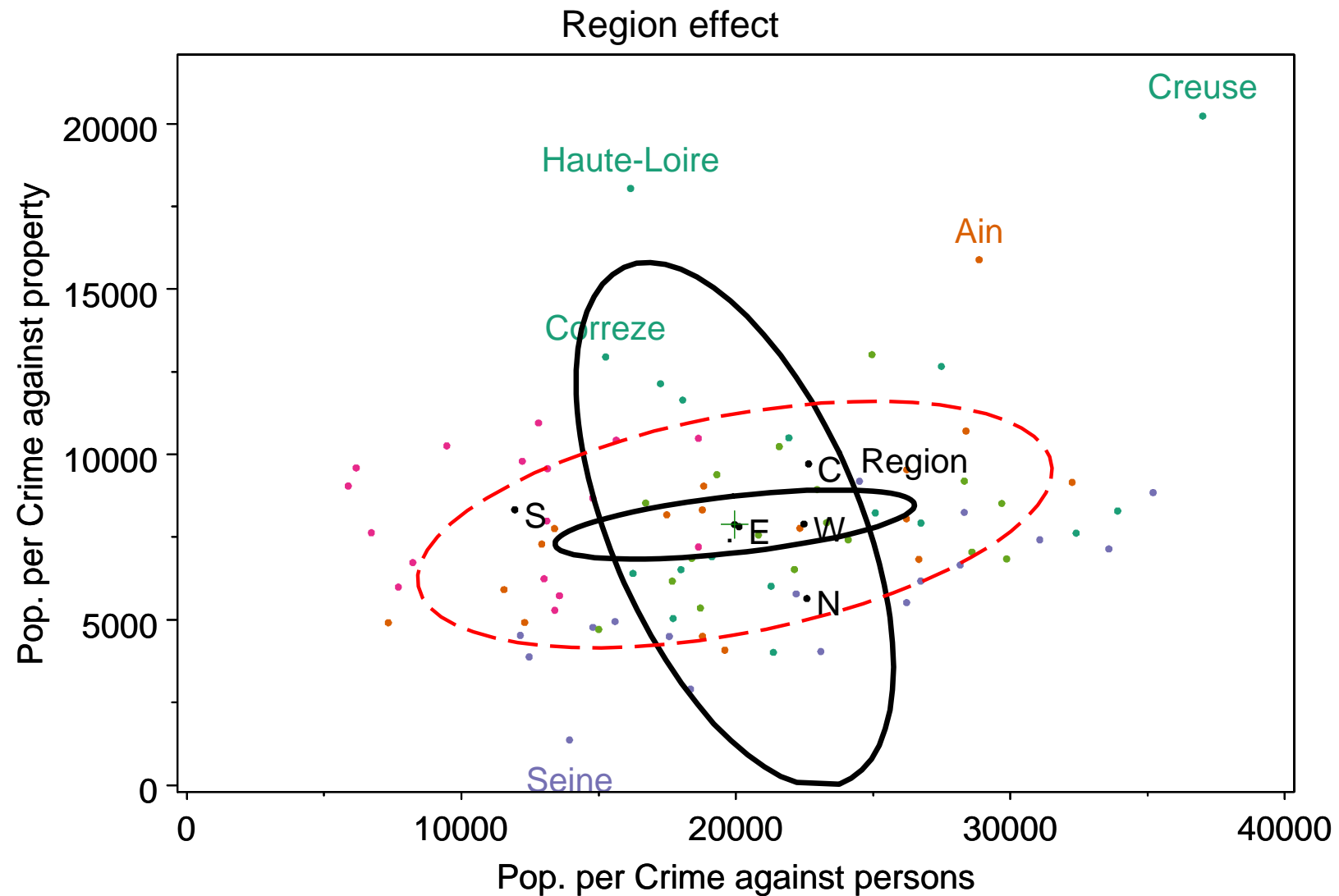Model: logSal Years ... Errors = Position

## Guerry data: Predicting crime

- How do rates of crime vary with other variables?

  - Fit MANCOVA model,
    ```
    (Crime_pers Crime_prop) = Region + Wealth + Suicides
                            + Literacy + Donations + Infants
    ```

  - HE plots: Overall, plus for Region and covariate effects

Region effect

- Overall: Predicted crimes against persons and property are negatively correlated
- Larger variation in crimes against property
- Region variation greater in crimes against persons

Covariates: Wealth Suicides Literacy Donations Infants
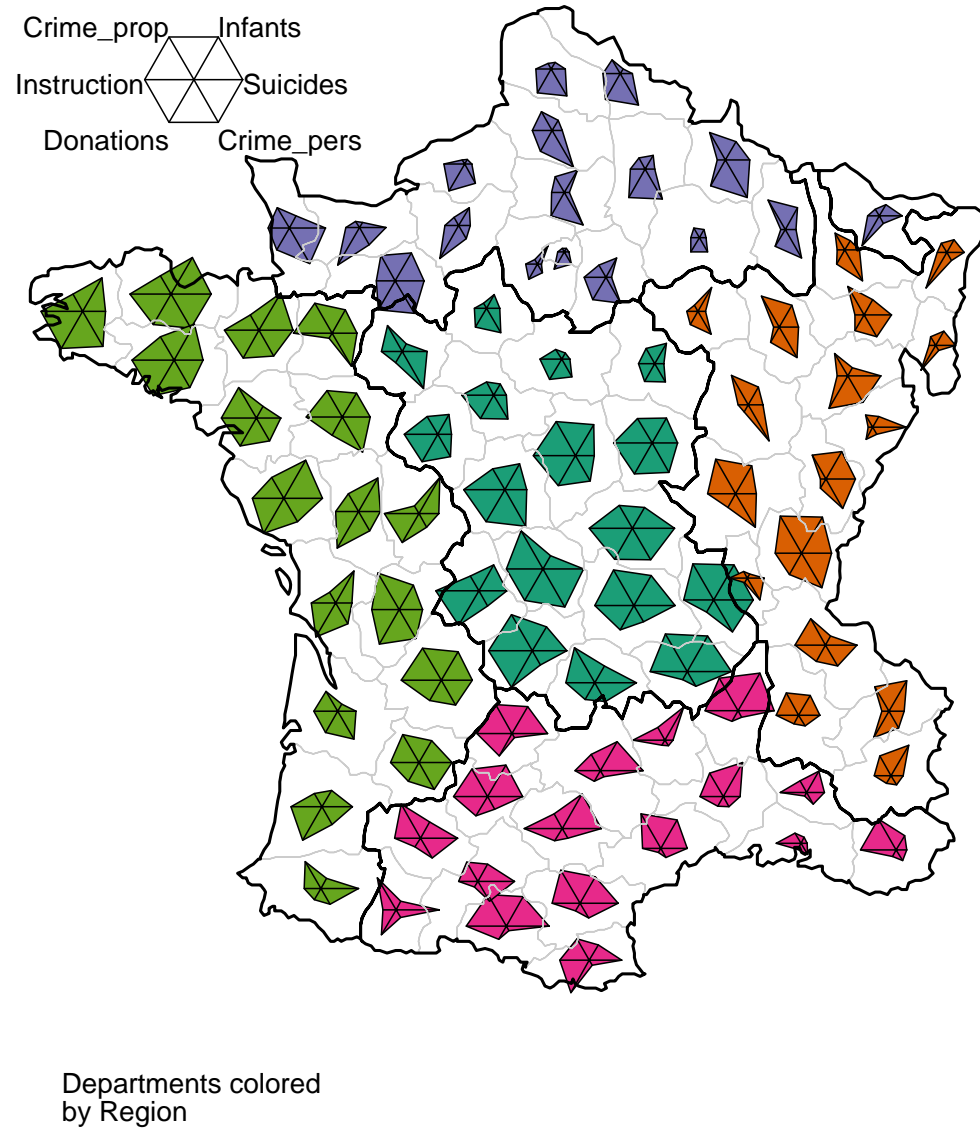
- Each quantitative variable (covariate) plots as a 1D ellipse (vector)
  - **Orientation**: relation of $x_i$ to $y_1, y_2$
  - **Length**: strength of relation

**Multivariate mapping: Map-centric displays**

- How to generalize choropleth maps to many variables?

- **Star maps**: Show multivariate data on the map using star icons, variable $\sim$ length of ray

- **Reduced-rank RGB displays**: Factor analysis $\rightarrow$ (F1, F2, F3) factor scores $\mapsto$ (R, G, B) shading

- **PREFMAP (x, y) maps**: Fit data variables to (Long, Lat) map coordinates. Display variables as vectors in map coordinates.

# Star maps

Star map of Guerry Variables (Ranks)

Crime_prop ——— Infants

Instruction ⬡ Suicides

Donations        Crime_pers

Departments colored
by Region

# Star maps: Medians by region

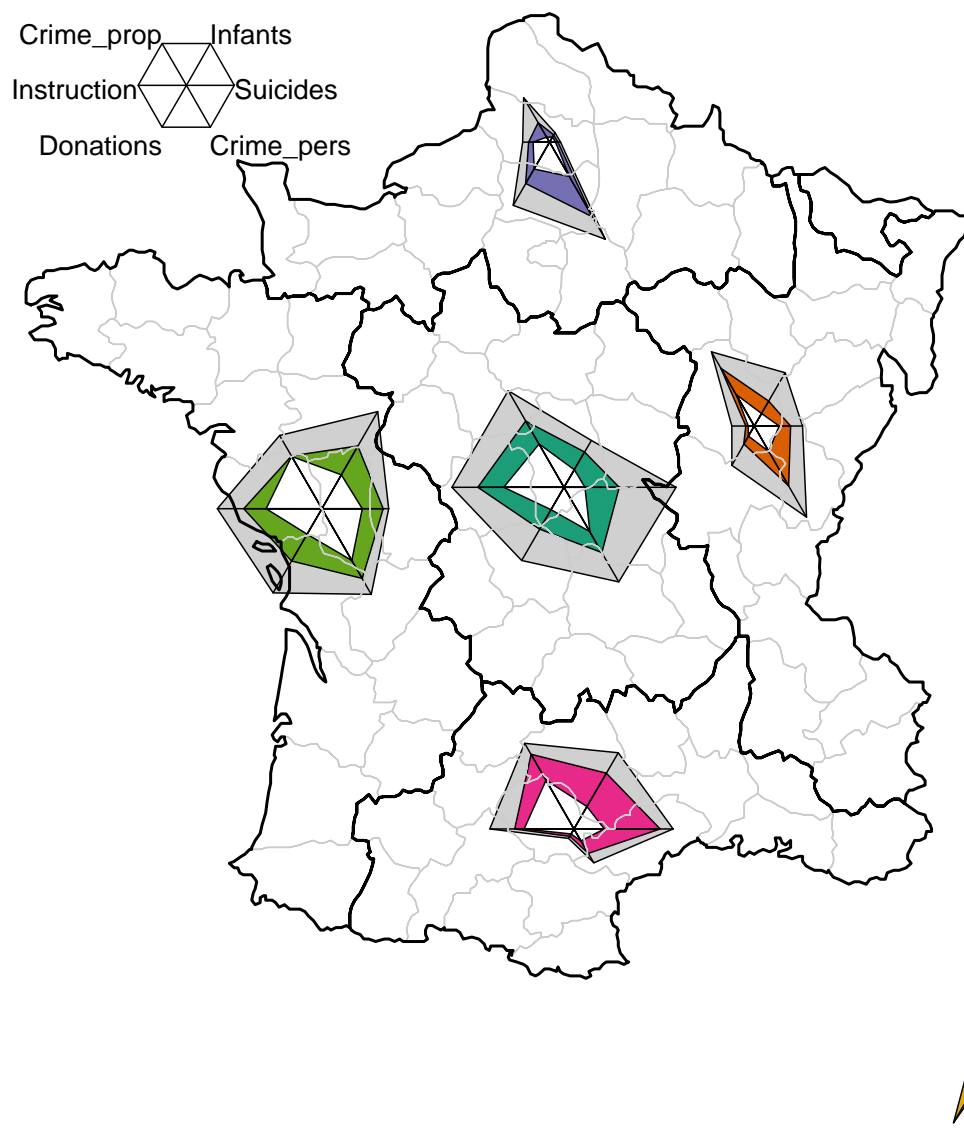Median Ranks by Region

# Star maps: Multivariate boxplots by region



- stars for Q1, Median, Q3
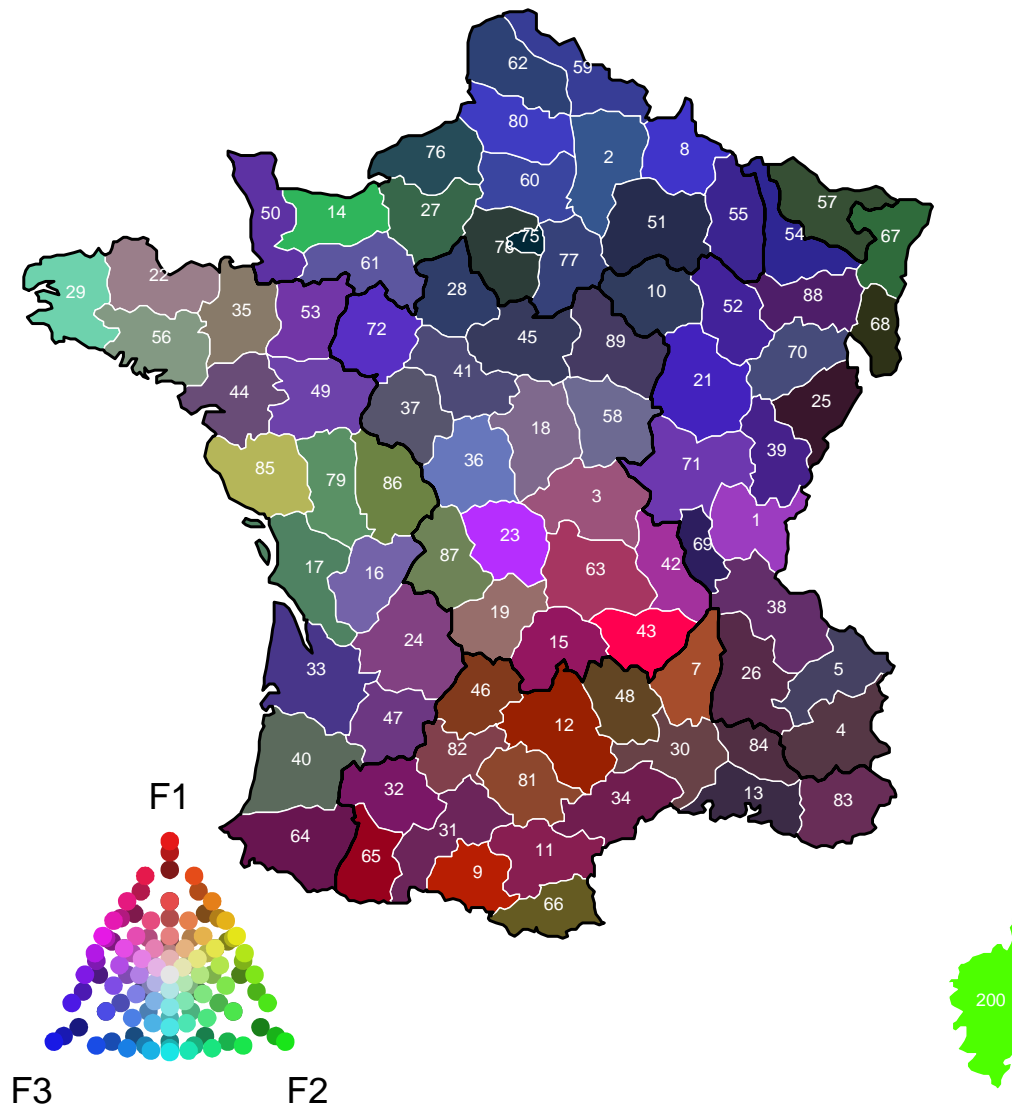- How to show unusual depts?

## Reduced-rank color-coded displays

■ Use dimension-reduction technique (PCA, Factor analysis, ...) to produce scores for observations (departments) on 3 dimensions ($F_1, F_2, F_3$)

| Variable | Factor1 Civil society | Factor2 Moral values | Factor3 Crime |
|---|---|---|---|
| Pop per Crime against persons | | | 0 97 |
| Pop per Crime against property | 0 75 | | 0 39 |
| Percent Read & Write | -0 72 | | |
| Pop per illegitimate birth | 0 62 | 0 42 | |
| Donations to the poor | | 0 89 | |
| Pop per suicide | 0 80 | | |

■ Scale $(F_1, F_2, F_3) \rightarrow [0,1]$

■ Color mapping function, e.g., $\mathcal{C}(F_1, F_2, F_3) \mapsto \mathrm{rgb}(F_i, F_j, F_k)$

# Reduced-rank color-coded displays

RGB 3-factor map: R=f1, G=f2, B=f3
Variables: Crime_pers Crime_prop Literacy Infants Donations Suicides

## Summary and future directions

- **Guerry's challenge**: Understanding uncertainty in multivariate, spatial data
  - How visualize and understand relations among many variables?
  - How to relate these to geographic information?

- **Understanding multivariate variation**
  - Visual summaries (data ellipses, smoothings) can show statistical relations more clearly and effectively
  - Reduced-rank visualization methods can show simpler, approximate views, based on several criteria.
  - Multivariate statistical models need their own visualization methods, just beginning – HE plots as an example.

- **Understanding multivariate, spatial variation**
  - Multivariate visualizations applied to spatial data can be revealing, but still need work
  - Statistical methods for spatial data need to be extended to the multivariate setting.

© Michael Friendly

# References

Friendly, M. (1991). *SAS System for Statistical Graphics*. Cary, NC: SAS Institute, 1st edn.

Friendly, M. (2002). Corrgrams: Exploratory displays for correlation matrices. *The American Statistician*, 56(4), 316–324.

Friendly, M. (2004a). Graphical methods for the multivariate general linear model. *Journal of Computational and Graphical Statistics*. (submitted 2/17/04; resubmitted: 10/02/04).

Friendly, M. (2004b). Milestones in the history of data visualization: A case study in statistical historiography. In W. Gaul and C. Weihs, eds., *Studies in Classification, Data Analysis, and Knowledge Organization*. New York: Springer. (In press).

Friendly, M. and Kwan, E. (2003). Effect ordering for data displays. *Computational Statistics and Data Analysis*, 43(4), 509–539.

Gabriel, K. R. (1971). The biplot graphic display of matrices with application to principal components analysis. *Biometrics*, 58(3), 453–467.

Guerry, A.-M. (1833). Essai sur la statistique morale de la France. Paris. English translation: Hugh P. Whitt and Victor W. Reinking, Lewiston, N.Y. : Edwin Mellen Press, c2002.

Pearson, K. (1920). Notes on the history of correlation. *Biometrika*, 13(1), 25–45.