

Seminar 2

Plotting many dimensions

- Tours
- Parallel coordinate plots
- Scatterplot matrices
- Multiple linked plots
- Using these together to explore data
- What we can learn about tennis!

Data Visualization

Discover, Explore and be Skeptical

Di Cook

Statistics, Iowa State University

soon to be Business Analytics, Monash University

LES DIABLERETS, FEB 1-4, 2015

2 -57

Notation

Data

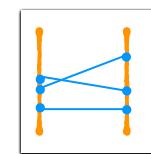
Crab	species	sex	frontal lobe	rear width	carapace length	carapace width	body depth
1	blue	male	8.1	6.7	16.1	19.0	7.0
2	blue	male	8.8	7.7	18.1	20.8	7.4
3	blue	male	9.2	7.8	19.0	22.4	7.7
4	blue	male	9.6	7.9	20.1	23.1	8.2
51	blue	female	7.2	6.5	14.7	17.1	6.1
52	blue	female	9.0	8.5	19.3	22.7	7.7
53	blue	female	9.1	8.1	18.5	21.6	7.7
101	orange	male	9.1	6.9	16.7	18.6	7.4
102	orange	male	10.2	8.2	20.2	22.2	9.0
151	orange	female	10.7	9.7	21.4	24.0	9.8
152	orange	female	11.4	9.2	21.7	24.1	9.7
153	orange	female	12.5	10.0	24.1	27.0	10.9

$$\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_p] = \begin{bmatrix} X_{11} & X_{12} & \dots & X_{1p} \\ X_{21} & X_{22} & \dots & X_{2p} \\ \vdots & \vdots & & \vdots \\ X_{n1} & X_{n2} & \dots & X_{np} \end{bmatrix}_{n \times p}$$

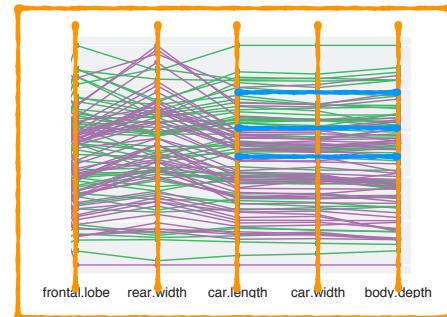
LES DIABLERETS, FEB 1-4, 2015

3 -57

Parallel coordinate plot



Cartesian to
parallel coords

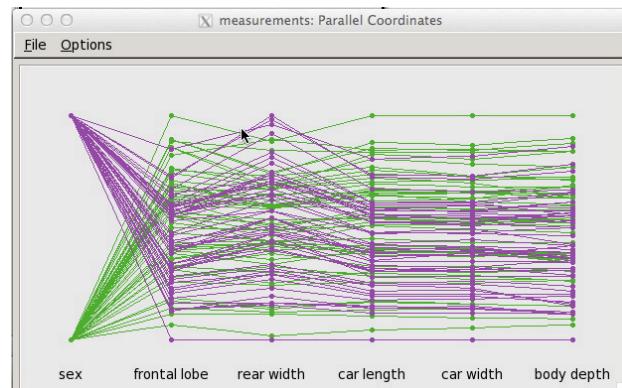
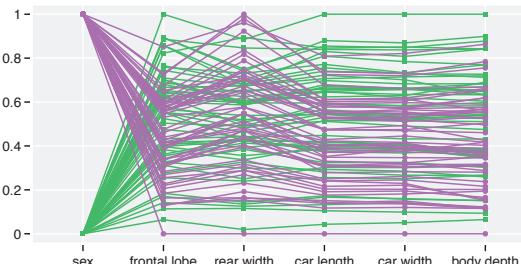


flat lines means positive association

Here, each variable is scaled individually by min/max

LES DIABLERETS, FEB 1-4, 2015

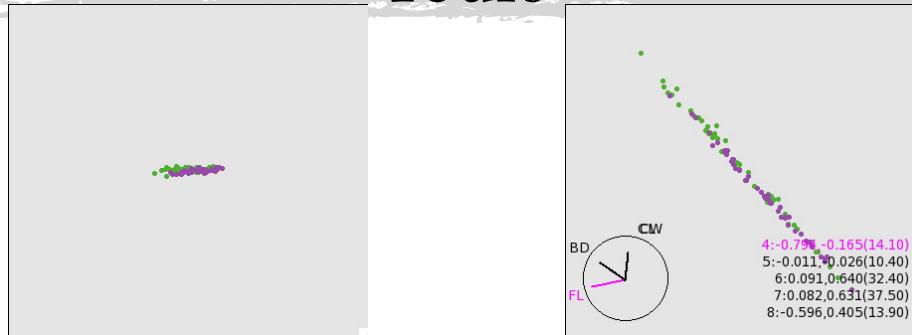
4 -57



LES DIABLERETS, FEB 1-4, 2015

5 -57

Tours



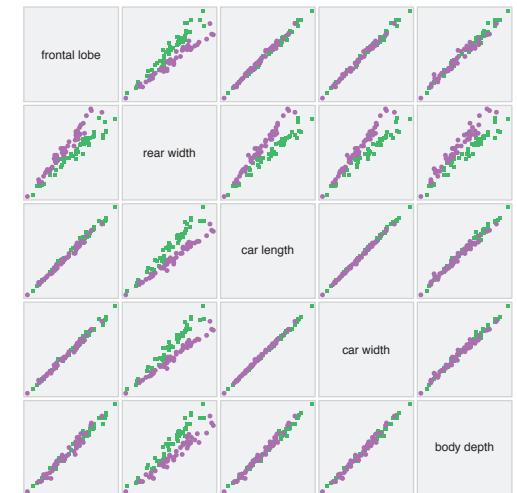
Motion graphic designed to study the joint distribution of multivariate data (Asimov 1985), in search of relationships that may involve several variables. It is created by generating a sequence of low-dimensional projections of high-dimensional data; these projections are typically 1D or 2D.

LES DIABLERETS, FEB 1-4, 2015

7 -57

Scatterplot matrix

- Same data: all pairs of the 5 variables displayed
- Strong association between all pairs.
- Difference between males and females on “rear.width”



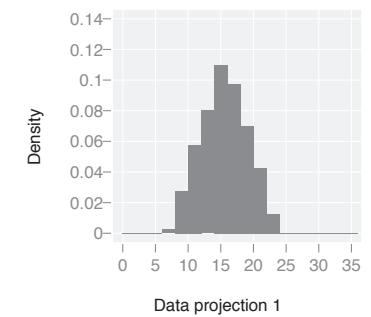
LES DIABLERETS, FEB 1-4, 2015

6 -57

Constructing a tour

	frontal lobe	rear width	carapace length	carapace width	body depth	
8.1	6.7	16.1	19.0	7.0		$\begin{bmatrix} 8.1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$
8.8	7.7	18.1	20.8	7.4		$\begin{bmatrix} 8.8 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$
9.2	7.8	19.0	22.4	7.7		$\begin{bmatrix} 9.2 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$
9.6	7.9	20.1	23.1	8.2		$\begin{bmatrix} 9.6 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$
7.2	6.5	14.7	17.1	6.1		$\begin{bmatrix} 7.2 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$
9.0	8.5	19.3	22.7	7.7		$\begin{bmatrix} 9.0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$
9.1	8.1	18.5	21.6	7.7		$\begin{bmatrix} 9.1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$
9.1	6.9	16.7	18.6	7.4		$\begin{bmatrix} 9.1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$
10.2	8.2	20.2	22.2	9.0		$\begin{bmatrix} 10.2 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$
10.7	9.7	21.4	24.0	9.8		$\begin{bmatrix} 10.7 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$
11.4	9.2	21.7	24.1	9.7		$\begin{bmatrix} 11.4 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$
12.5	10.0	24.1	27.0	10.9		$\begin{bmatrix} 12.5 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$

$\mathbf{A}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$, then the data projection is $\mathbf{X}\mathbf{A}_1 = \begin{bmatrix} 8.1 \\ 8.8 \\ 9.2 \\ 9.6 \\ 7.2 \\ \vdots \end{bmatrix}$



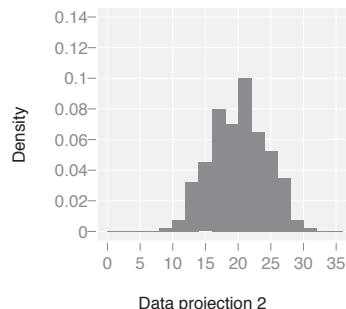
LES DIABLERETS, FEB 1-4, 2015

8 -57

Constructing a tour

frontal lobe	rear width	carapace length	carapace width	body depth
8.1	6.7	16.1	19.0	7.0
8.8	7.7	18.1	20.8	7.4
9.2	7.8	19.0	22.4	7.7
9.6	7.9	20.1	23.1	8.2
7.2	6.5	14.7	17.1	6.1
9.0	8.5	19.3	22.7	7.7
9.1	8.1	18.5	21.6	7.7
9.1	6.9	16.7	18.6	7.4
10.2	8.2	20.2	22.2	9.0
10.7	9.7	21.4	24.0	9.8
11.4	9.2	21.7	24.1	9.7
12.5	10.0	24.1	27.0	10.9

$$\mathbf{A}_2 = \begin{bmatrix} 0.707 \\ 0.707 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \text{ then } \mathbf{X}\mathbf{A}_2 = \begin{bmatrix} 0.707 \times 8.1 + 0.707 \times 6.7 = 10.5 \\ 0.707 \times 8.8 + 0.707 \times 7.7 = 11.7 \\ 0.707 \times 9.2 + 0.707 \times 7.8 = 12.0 \\ 0.707 \times 9.6 + 0.707 \times 7.9 = 12.4 \\ 0.707 \times 7.2 + 0.707 \times 6.5 = 9.7 \\ \vdots \end{bmatrix}$$



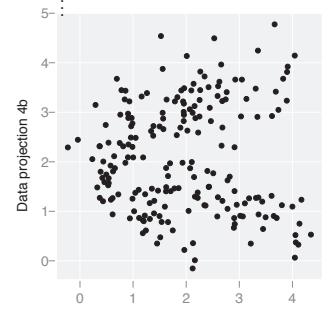
LES DIABLERETS, FEB 1-4, 2015

9 -57

Constructing a tour

frontal lobe	rear width	carapace length	carapace width	body depth
8.1	6.7	16.1	19.0	7.0
8.8	7.7	18.1	20.8	7.4
9.2	7.8	19.0	22.4	7.7
9.6	7.9	20.1	23.1	8.2
7.2	6.5	14.7	17.1	6.1
9.0	8.5	19.3	22.7	7.7
9.1	8.1	18.5	21.6	7.7
9.1	6.9	16.7	18.6	7.4
10.2	8.2	20.2	22.2	9.0
10.7	9.7	21.4	24.0	9.8
11.4	9.2	21.7	24.1	9.7
12.5	10.0	24.1	27.0	10.9

$$\mathbf{A}_4 = \begin{bmatrix} 0 & 0 \\ 0 & 0.950 \\ 0 & -0.312 \\ -0.312 & 0 \\ 0.950 & 0 \end{bmatrix} \text{ then } \mathbf{X}\mathbf{A}_4 = \begin{bmatrix} -0.312 \times 19.0 + 0.950 \times 7.0 = 0.72 & 0.950 \times 6.7 - 0.312 \times 16.1 = 1.34 \\ -0.312 \times 20.8 + 0.950 \times 7.4 = 0.54 & 0.950 \times 7.7 - 0.312 \times 18.1 = 1.67 \\ -0.312 \times 22.4 + 0.950 \times 7.7 = 0.33 & 0.950 \times 7.8 - 0.312 \times 19.0 = 1.48 \\ -0.312 \times 23.1 + 0.950 \times 8.2 = 0.58 & 0.950 \times 7.9 - 0.312 \times 20.1 = 1.23 \\ -0.312 \times 17.1 + 0.950 \times 6.1 = 0.46 & 0.950 \times 6.5 - 0.312 \times 14.7 = 1.59 \\ \vdots & \vdots \end{bmatrix}$$



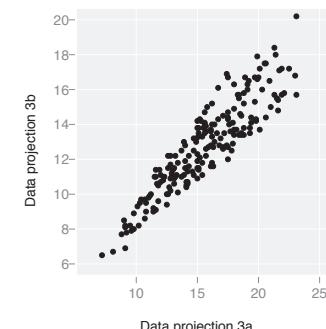
LES DIABLERETS, FEB 1-4, 2015

11 -57

Constructing a tour

frontal lobe	rear width	carapace length	carapace width	body depth
8.1	6.7	16.1	19.0	7.0
8.8	7.7	18.1	20.8	7.4
9.2	7.8	19.0	22.4	7.7
9.6	7.9	20.1	23.1	8.2
7.2	6.5	14.7	17.1	6.1
9.0	8.5	19.3	22.7	7.7
9.1	8.1	18.5	21.6	7.7
9.1	6.9	16.7	18.6	7.4
10.2	8.2	20.2	22.2	9.0
10.7	9.7	21.4	24.0	9.8
11.4	9.2	21.7	24.1	9.7
12.5	10.0	24.1	27.0	10.9

$$\mathbf{A}_3 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}, \text{ then the data projection is } \mathbf{X}\mathbf{A}_3 = \begin{bmatrix} 8.1 & 6.7 \\ 8.8 & 7.7 \\ 9.2 & 7.8 \\ 9.6 & 7.9 \\ 7.2 & 6.5 \\ \vdots & \vdots \end{bmatrix}$$



LES DIABLERETS, FEB 1-4, 2015

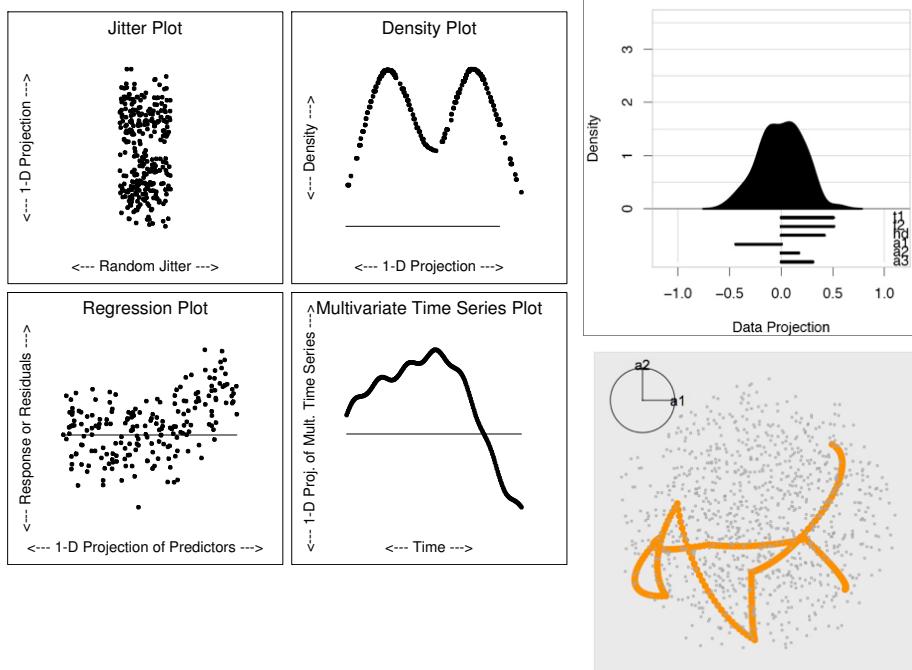
10 -57

Sequencing the projections

- Tours to display strictly real-valued multivariate data
- $\text{View}_i(t) = F(t)^T \mathbf{x}_i$, $F(t) = (\mathbf{f}_1(t), \dots, \mathbf{f}_d(t))$
- Render the view, and navigation info
- Method for choosing $F(t)$
- Interpolate between consecutive $F(t)$

LES DIABLERETS, FEB 1-4, 2015

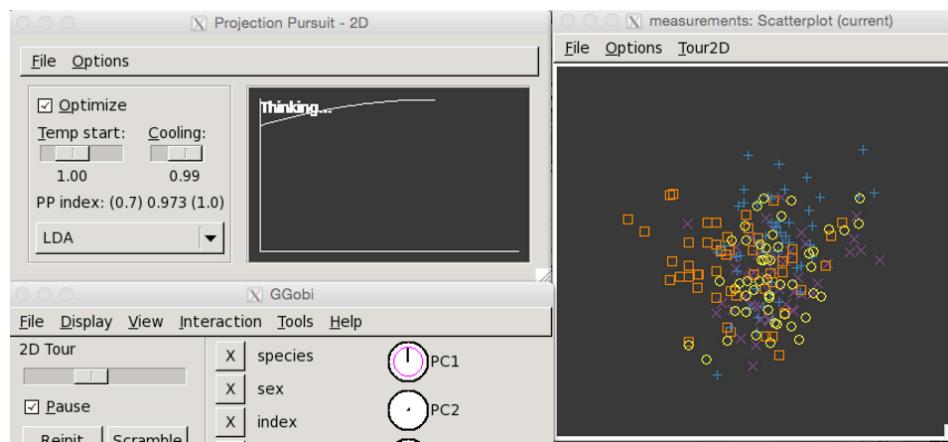
12 -57



LES DIABLERETS, FEB 1-4, 2015

13 -57

Guided tour: four groups of crabs, optimizing PP index



LES DIABLERETS, FEB 1-4, 2015

15 -57

Choosing F

- Grand tour: sample from a uniform on a sphere
- Guided: Optimize a projection pursuit index, eg

$$I_{LDA}(F) = 1 - \frac{|F^T W F|}{|F^T (W + B) F|}$$

$$y = F^T x$$

$$B = \sum_{i=1}^g n_i (\bar{y}_{i\cdot} - \bar{y}_{..})(\bar{y}_{i\cdot} - \bar{y}_{..})^T$$

$$W = \sum_{i=1}^g \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i\cdot})(y_{ij} - \bar{y}_{i\cdot})^T$$

- Manual: Choose a variable to control, allow user to interactively control coefficient, ranging between -1, 1, constrained on all other variables

LES DIABLERETS, FEB 1-4, 2015

14 -57

Why?

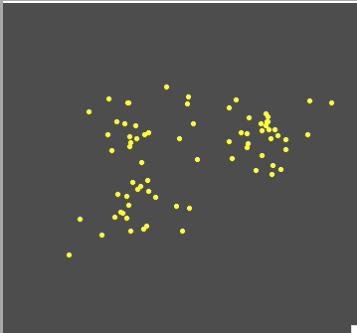
- Learn how several variables jointly vary
- Examine the multivariate distribution
- Check model fit
- Explore deviations from distribution: outliers, clusters, nonlinear relationships

LES DIABLERETS, FEB 1-4, 2015

16 -57

Your Turn

For each of the following videos answer these questions



How many clusters?

Anything else you see?

17

Your Turn

For each of the following videos answer these questions



Are the two groups different from each other?

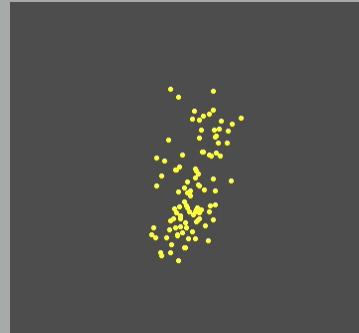
Do you see any outliers?

Any small clusters?

19

Your Turn

For each of the following videos answer these questions

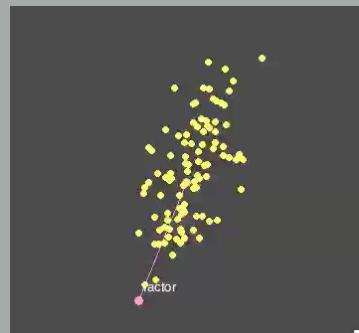


Linear dependence, or nonlinear dependence?

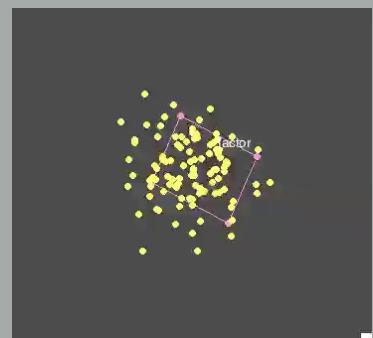
18

Your Turn

For each of the following videos answer these questions

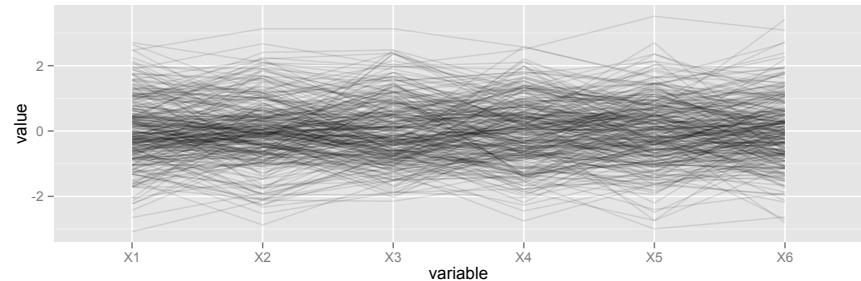


One- or two-dimensional?



20

Normal data

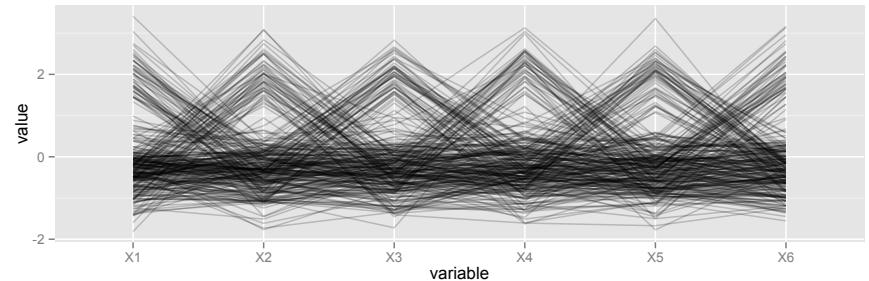


Nothing interesting! All a little moderate correlation.
Modelling is going to be easy!

LES DIABLERETS, FEB 1-4, 2015

21 -57

Clustered data

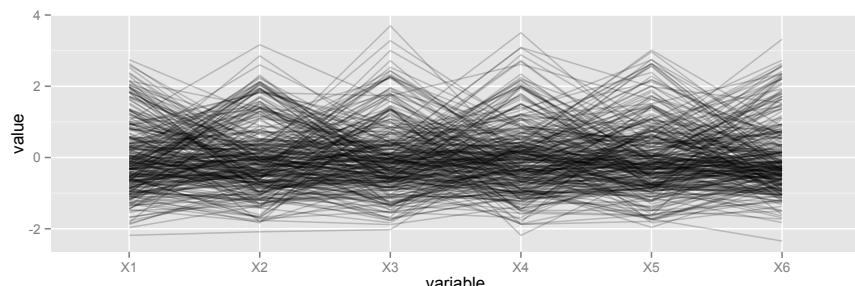


See the criss-crossing, gaps between lines.
Will need to extract the clusters before doing any other modeling, otherwise pretty regular data

LES DIABLERETS, FEB 1-4, 2015

22 -57

(Less) Clustered data

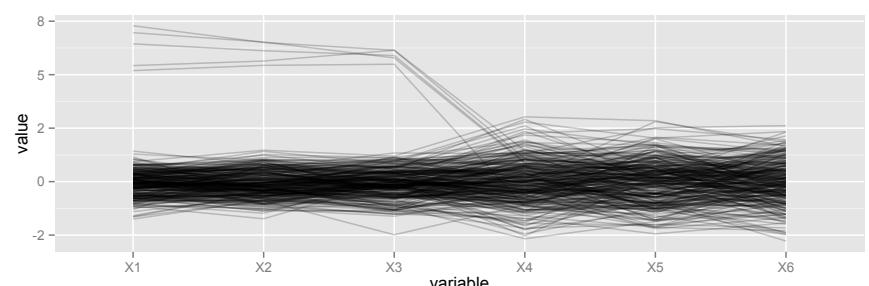


Still see the criss-crossing, gaps between lines, but less prominent.
Will need to deal with the multi-modality

LES DIABLERETS, FEB 1-4, 2015

23 -57

Outliers in the data

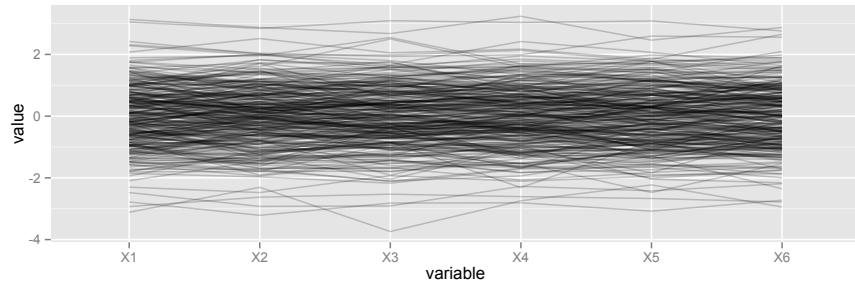


Small group of observations that are outliers on X1-X3.
Need to do something with these cases, remove with justification, or fix

LES DIABLERETS, FEB 1-4, 2015

24 -57

Strong association

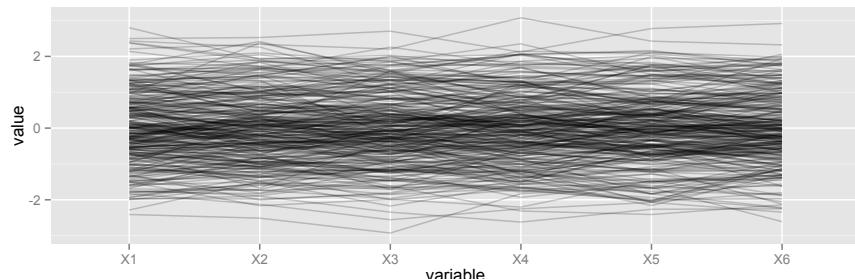


Very flat lines indicate strong positive association between all variables.

LES DIABLERETS, FEB 1-4, 2015

25 ·57

Strong negative association - fixed

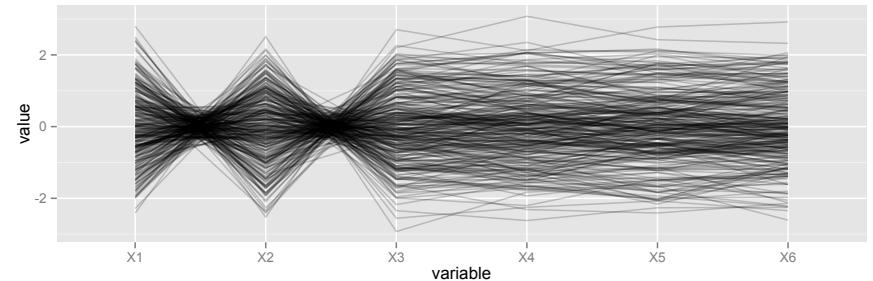


X2 is multiplied by -1, then it is positively associated with other variables.

LES DIABLERETS, FEB 1-4, 2015

27 ·57

Strong negative association

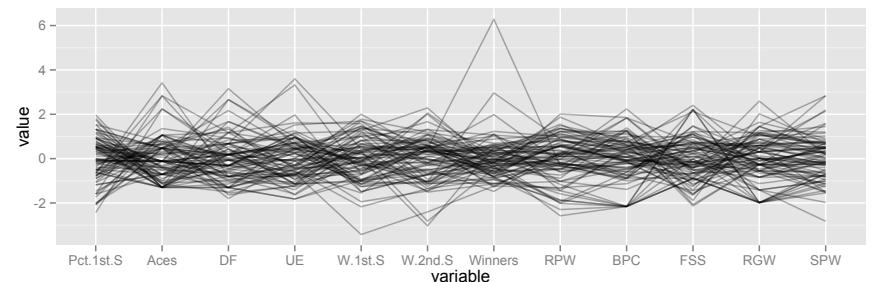


Crossed lines in first three variables indicate X2 is strongly negatively correlated with other vars.

LES DIABLERETS, FEB 1-4, 2015

26 ·57

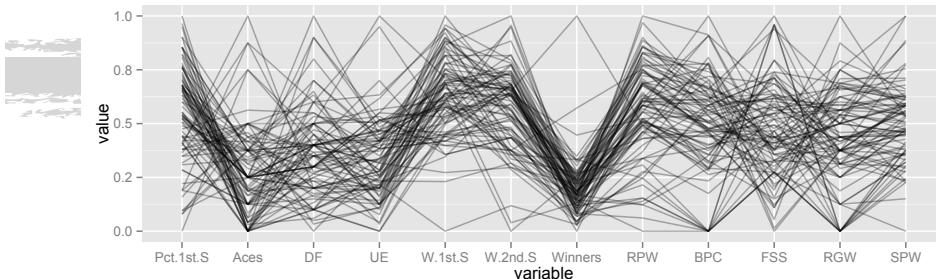
Tennis statistics



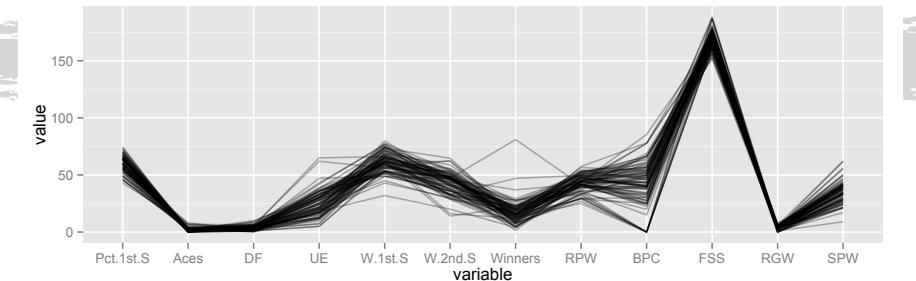
Scale matters: mean/sd, 0/1, individual/global
Enables correlation to be seen better, and outliers

LES DIABLERETS, FEB 1-4, 2015

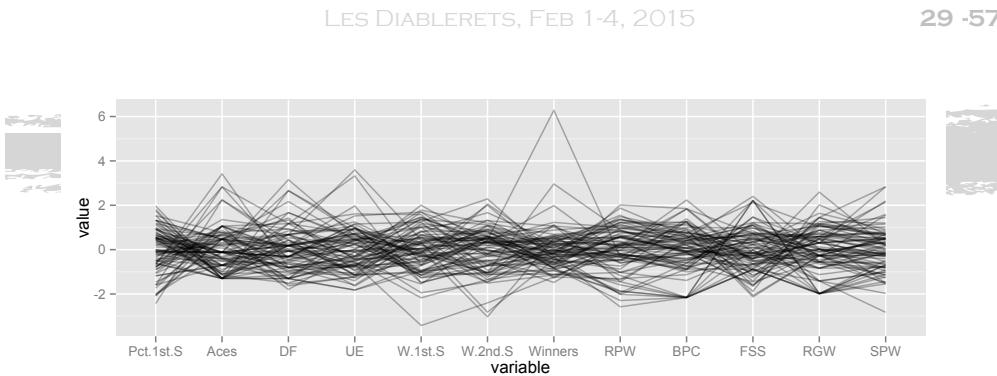
28 ·57



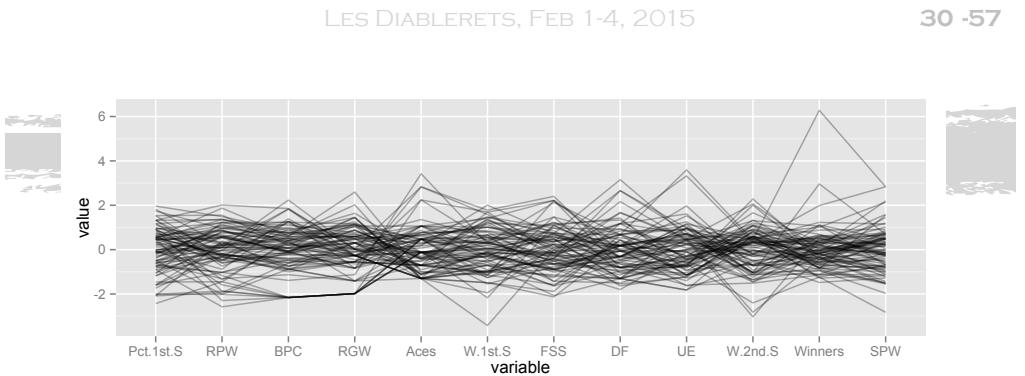
Scale matters: mean/sd, 0/1, individual/global
Emphasizes the univariate distributions



Scale matters: mean/sd, 0/1, individual/global



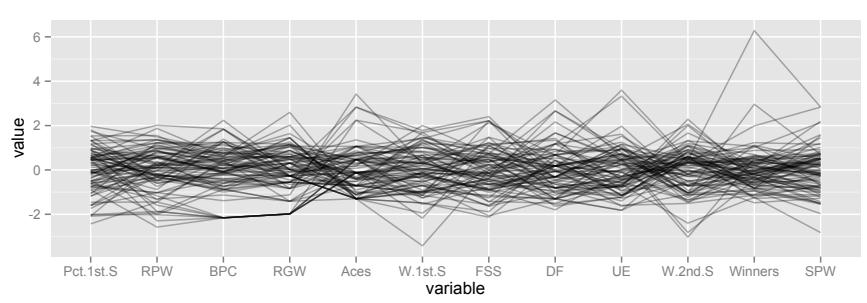
Order matters: place variables that are highly correlated close to each other



Order matters: place variables that are highly correlated close to each other
Less line crossing, easier to digest positive correlation, and then negative correlation

Your Turn

Take two minutes and discuss with your neighbor what you see

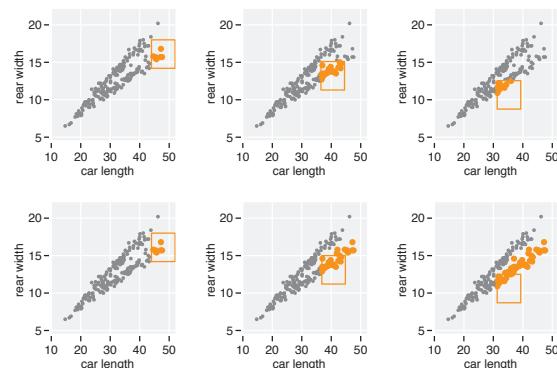


Clusters? Outliers? (In one variable or multiple variables?) Association?

33

Brushing

Persistent painting vs transient brushing



LES DIABLERETS, FEB 1-4, 2015

35 -57

Interactivity

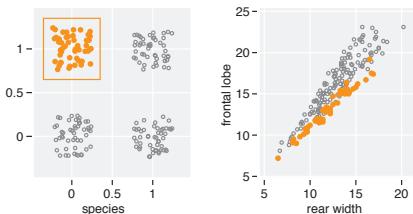
- Brushing and linking between multiple plots
- Identifying
- Scaling

LES DIABLERETS, FEB 1-4, 2015

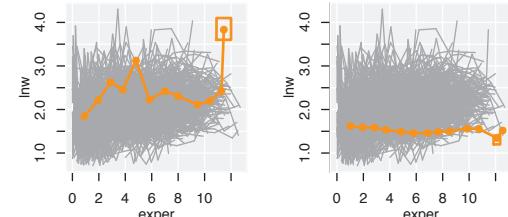
34 -57

Linking

- One-to-one



- Using a categorical variable

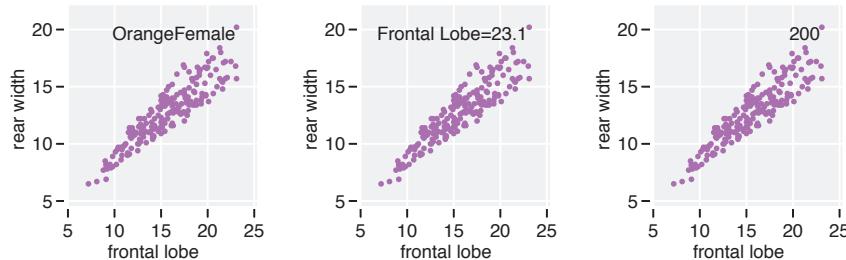


LES DIABLERETS, FEB 1-4, 2015

36 -57

Identification

Lookup label information on data element



LES DIABLERETS, FEB 1-4, 2015

37 -57

Exploring tennis statistics

- 2014 was a great year for Swiss tennis
- Stan Wawrinka surprised everyone and defeated Nadal in the final, after defeating Berdych and Djokovic to get there
- Switzerland won its first ever Davis Cup

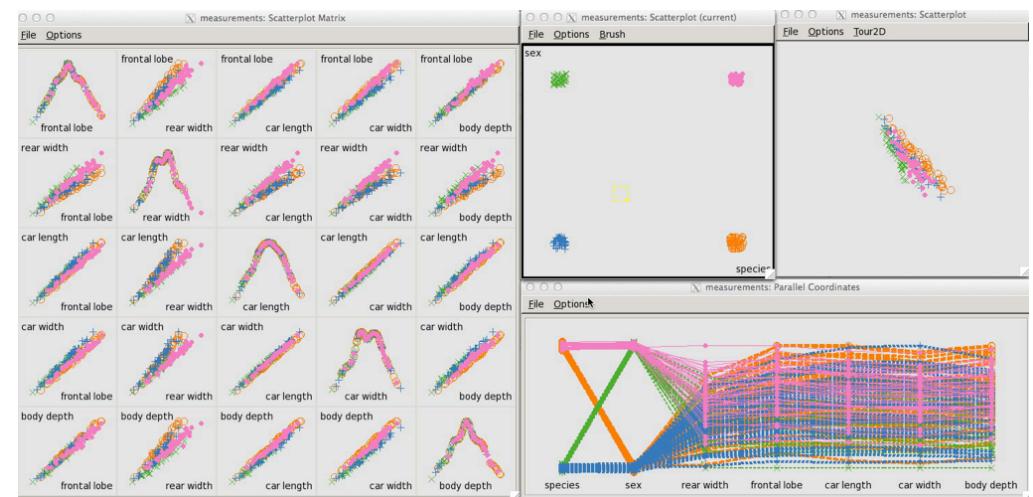


Statistics for 2014 pulled from http://www.ausopen.com/en_AU/players/overview/atpw367.html

LES DIABLERETS, FEB 1-4, 2015

39 -57

Putting it together



LES DIABLERETS, FEB 1-4, 2015

38 -57

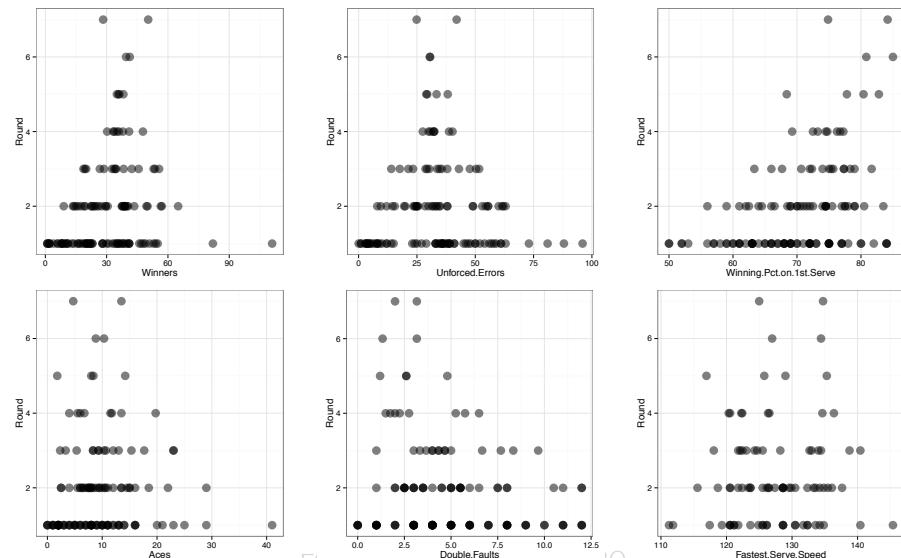
Questions

- What performance statistics suggest advancing in the tournament?
- How did Stan Wawrinka win?

LES DIABLERETS, FEB 1-4, 2015

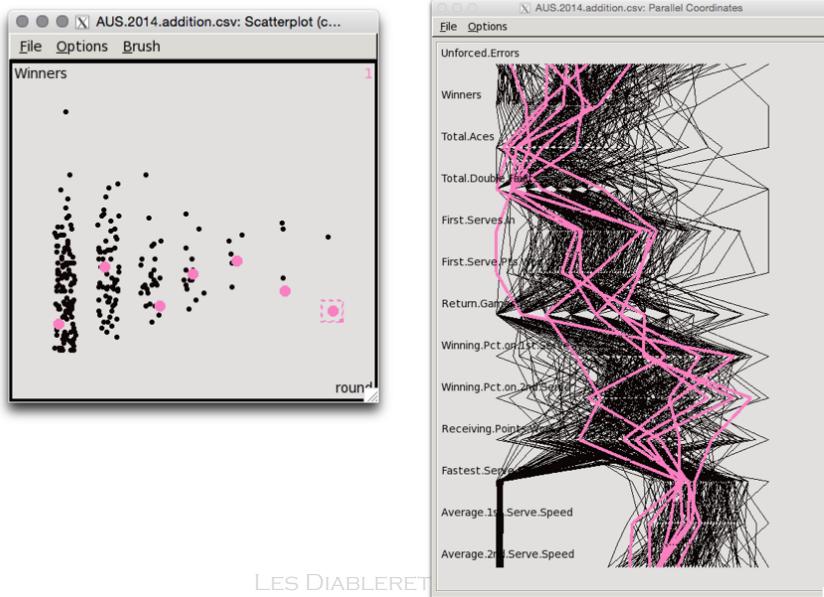
40 -57

Performance



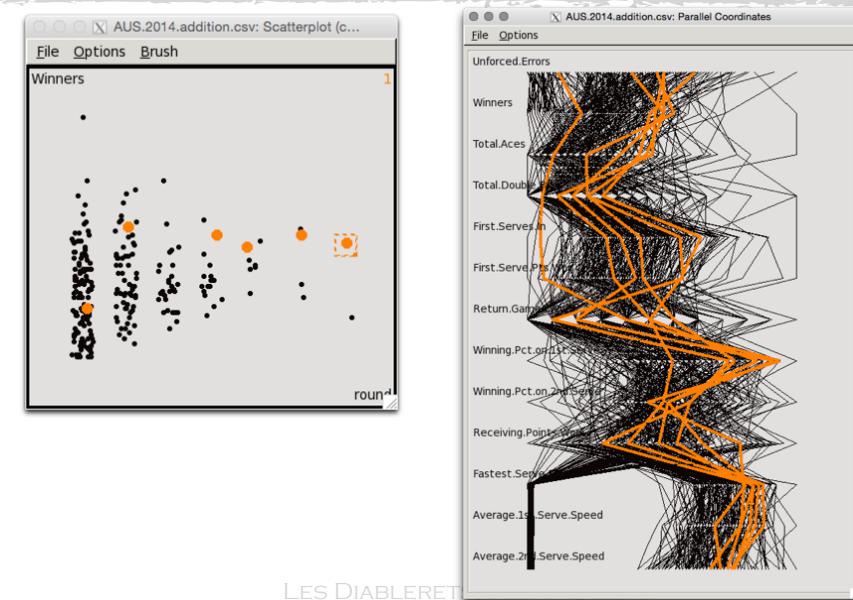
-57

Wawrinka vs Nadal



43 -57

Wawrinka vs Nadal



LES DIABLERETS

42 -57

What we learn

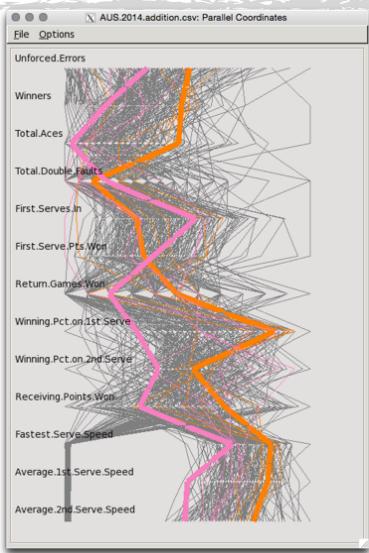
- ➊ It is important to keep control in the game, not too many winners, not too few, and errors same.
- ➋ Serve speed is not that important.
- ➌ Winning your serve is important.
- ➍ Stan beat Rafa on winners in final!
- ➎ Nadal serves slowly, and not so much slower in final.

LES DIABLERETS

LES DIABLERETS, FEB 1-4, 2015

44 -57

Stan vs Rafa



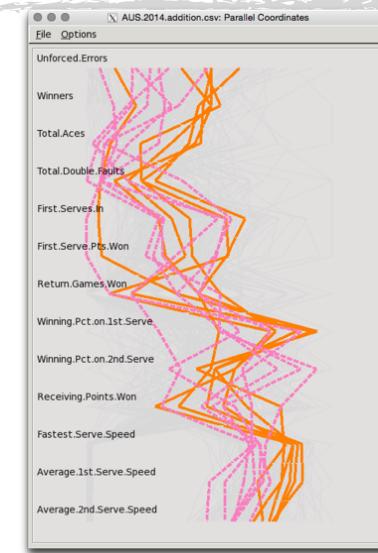
Stan won the game on winners, unforced errors, aces, serving % and speed.

It looks like he played aggressively, he was in the zone, and the gamble paid off.

LES DIABLERETS, FEB 1-4, 2015

45 -57

Stan vs Rafa



Generally, throughout tournament, Stan had more winners, errors, first serve %, serve speed.

He had a better tournament performance.

LES DIABLERETS, FEB 1-4, 2015

46 -57

Your Turn

Take two minutes to come up with some more questions

- ⌚ What other things would you like to investigate in the game of tennis?
- ⌚ What calculations, tables, plots would you make to tackle these questions?

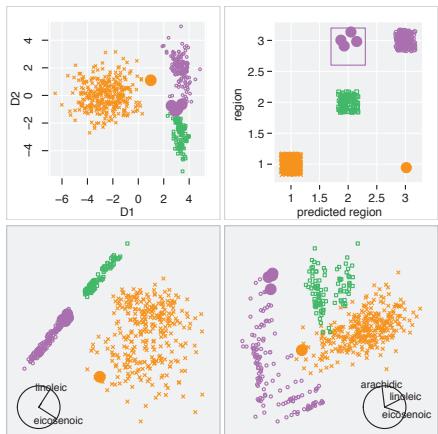
Exploring model fits

- Classification, examining boundary rules and misfits
- Clustering, exploring self-organizing maps

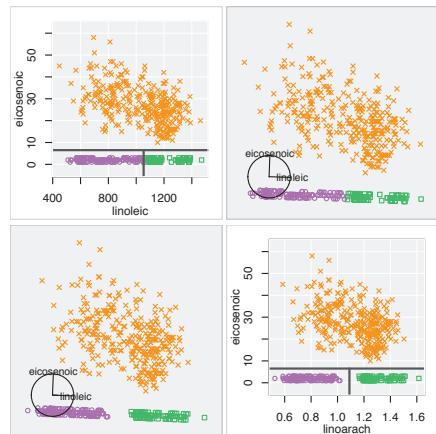
LES DIABLERETS, FEB 1-4, 2015

49 -57

LDA



Trees



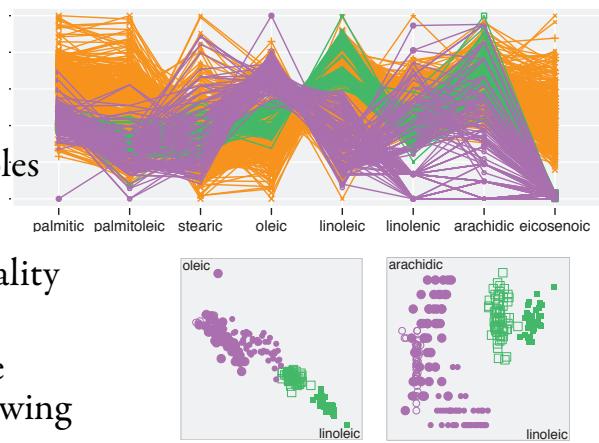
LES DIABLERETS, FEB 1-4, 2015

51 -57

Italian Olive oils



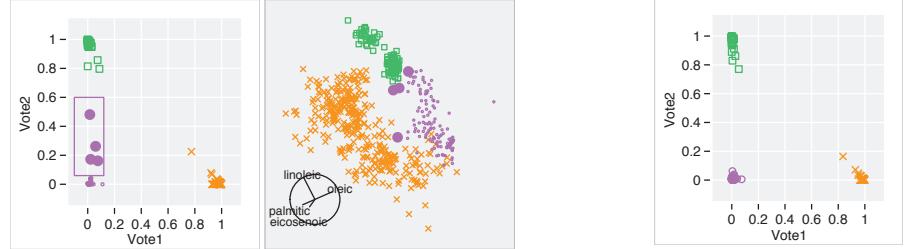
- Ancient statistics textbook data
- 572 samples
- 8 fatty acid composition variables
- 9 classes
- Related to food quality and pricing
- Fatty acid signature associated with growing region



LES DIABLERETS, FEB 1-4, 2015

50 -57

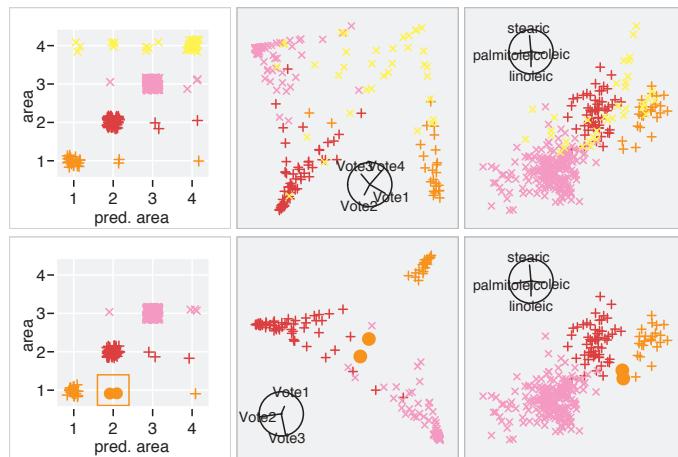
Random forests



LES DIABLERETS, FEB 1-4, 2015

52 -57

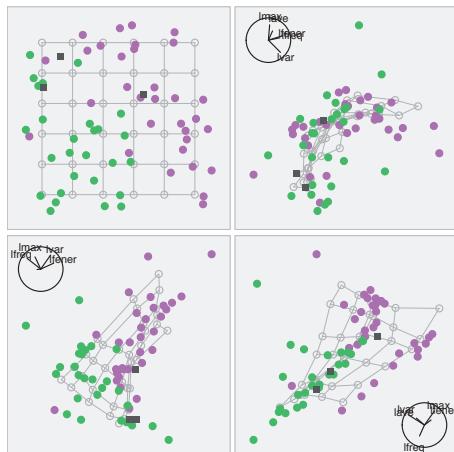
Random forests



LES DIABLERETS, FEB 1-4, 2015

53 -57

Self-organizing maps

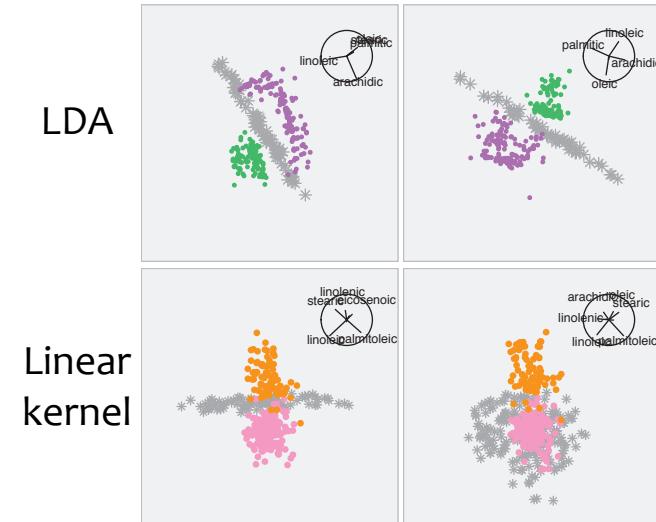


- Model is fit by warping a sheet through high-d
- Lay out sheet to see proximities
- Instead, use the tour to examine how sheet fits the data

LES DIABLERETS, FEB 1-4, 2015

55 -57

Support vector machines



LES DIABLERETS, FEB 1-4, 2015

54 -57

Summary

- We can see beyond 3D, with a combination of dynamic graphics and linking between multiple plots.
- Statistical graphics explores abstract relationships between variables, and enables building a conceptual map of structure in data
- It is important to examine the model fit IN THE DATA SPACE.

LES DIABLERETS, FEB 1-4, 2015

56 -57

Acknowledgements & Resources

- Used R packages ggplot2, cranvas (<http://cranvas.org>), and also ggobi (<http://www.ggobi.org>), knitr/rmarkdown
- New tools: ggviz, shiny, animint