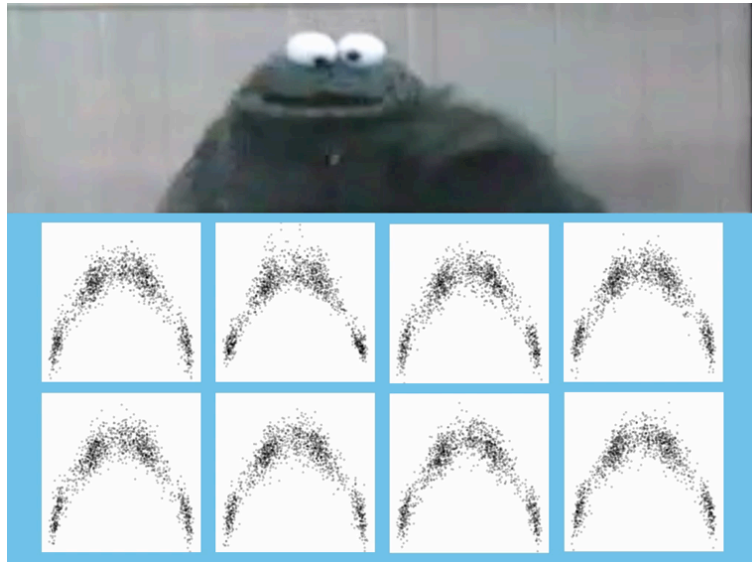


Data Visualization

Discover, Explore and be Skeptical

Di Cook
Statistics, Iowa State University
soon to be Business Analytics, Monash University

LES DIABLERETS, FEB 1-4, 2015



Video made by Hadley Wickham

LES DIABLERETS, FEB 1-4, 2015

3 -46

Seminar 3

Inference and Exploration

- Discoveries need to be calibrated by what might have been possible. Maintain a healthy skepticism.
- Underlying plots of data, are assumptions that implicitly specifying null hypotheses: what would you see if there really was nothing happening.
- Exploratory and inferential ARE NOT mutually exclusive.

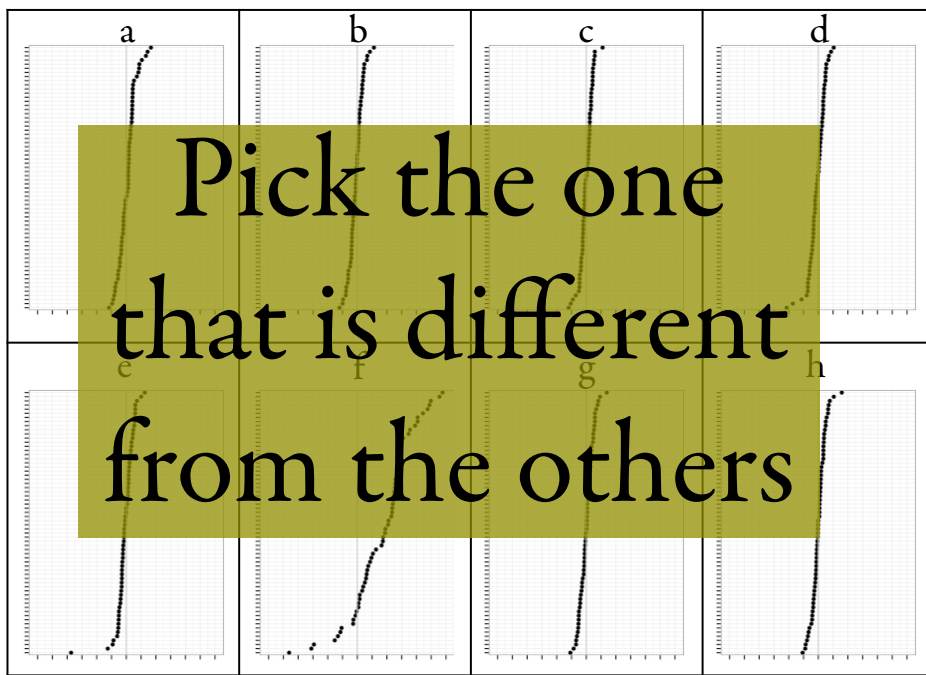
LES DIABLERETS, FEB 1-4, 2015

2 -46

Here is the math gap
exploration placed in the
CONTEXT of there being
NO MATH GAP ...

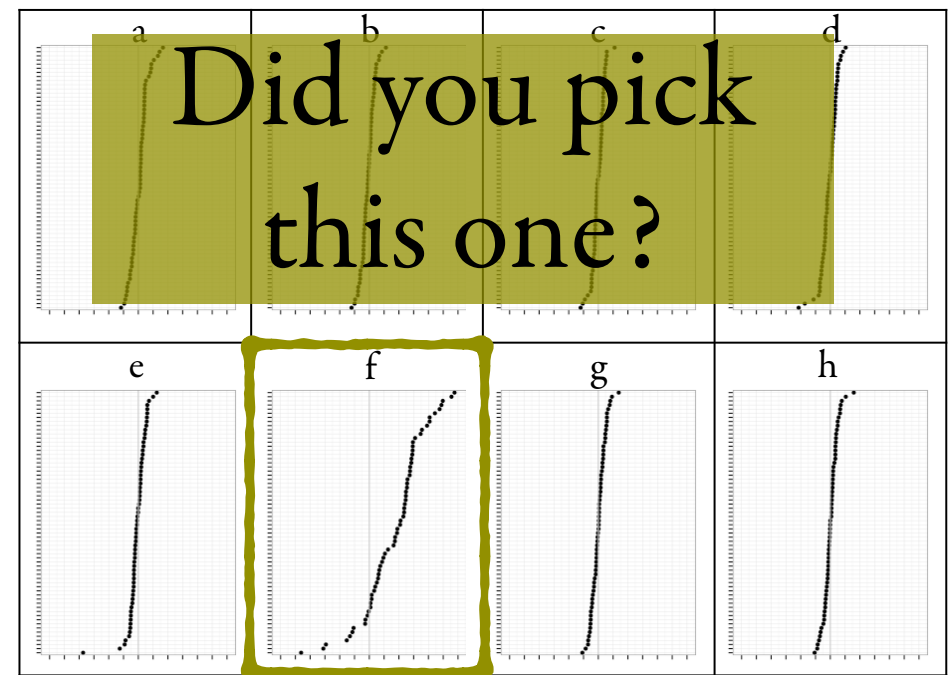
LES DIABLERETS, FEB 1-4, 2015

4 -46



LES DIABLERETS, FEB 1-4, 2015

5 -46



LES DIABLERETS, FEB 1-4, 2015

6 -46

Nulls by permutation

- Hold country fixed (subset by country)
- Permute the gender labels, so that the math scores are randomly assigned to a boy or girl
- Recalculate the difference between the means
- Plot the mean difference by country again

LES DIABLERETS, FEB 1-4, 2015

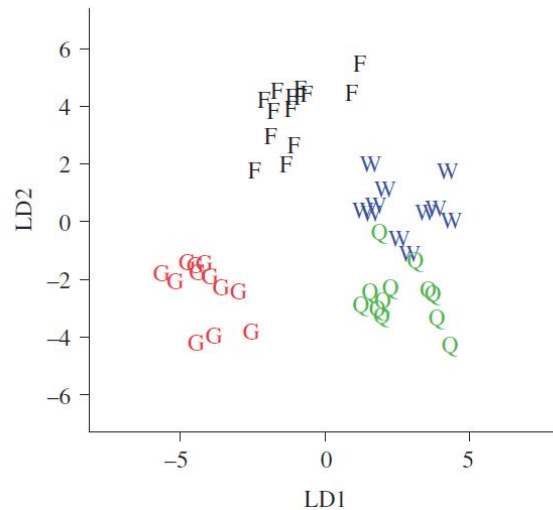
7 -46

LES DIABLERETS, FEB 1-4, 2015

8 -46

Let's do a real one

- 40 oligos (variables)
- 48 wasps (cases)
- 4 types of wasps
- Best LDA 2D separation of four groups
(Toth et al, 2010)



LES DIABLERETS, FEB 1-4, 2015

9 -46

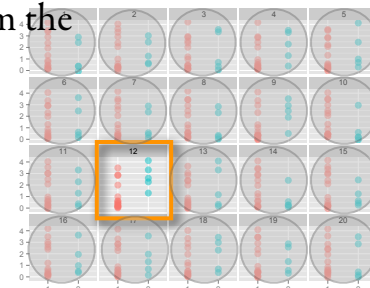
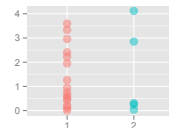
Really?

LES DIABLERETS, FEB 1-4, 2015

10 -46

Protocols

- Rorschach: Show many pictures of data with “nothing” happening, pictures from a null distribution
- Lineup: Embed the plot of the data among plots of data generated from the null distribution



Data plot
Null plots

Source: Buja et al (2009) RSPT(A)

LES DIABLERETS, FEB 1-4, 2015

11 -46

Hypothesis testing

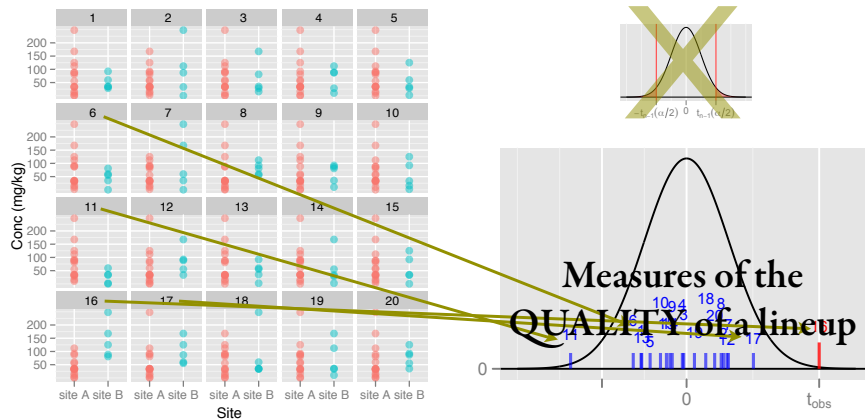
	Mathematical Inference	Visual Inference
Hypothesis	$H_0 : \mu_1 = \mu_2$ vs $H_a : \mu_1 \neq \mu_2$	$H_0 : \mu_1 = \mu_2$ vs $H_a : \mu_1 \neq \mu_2$
Test Statistic	$T(y) = \frac{\bar{y}_1 - \bar{y}_2}{s \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$	$T(y) =$
Sampling Distribution	$f_{T(y)}(t);$	$f_{T(y)}(t);$
Reject H_0 if	observed T is extreme	observed plot is identifiable

LES DIABLERETS, FEB 1-4, 2015

12 -46

Consideration ONE

Sampling distribution comparison is against a finite



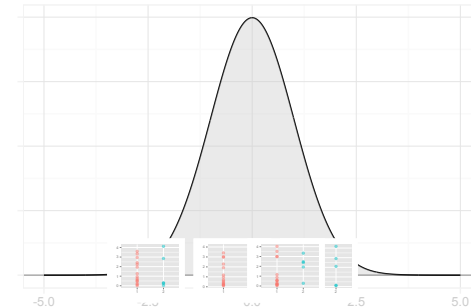
Source: Roy Chowdhury (2014)

LES DIABLERETS, FEB 1-4, 2015

13 -46

Consideration ONE

KEEP IN MIND: In practice, graphics is being used when there is no quantification of a sampling distribution. All we have is $(m-1)$ representatives from whatever



LES DIABLERETS, FEB 1-4, 2015

14 -46

Consideration TWO

- What is the p -value?
- For one observer, the probability of randomly selecting the data plot is $1/m$, where m is the number of plots in the lineup.
- With multiple observers, the p -value is estimated by

Number of independent observers

$$P(X \geq x) = 1 - \text{Binom}_{K, 1/m}(x) = \sum_{i=x}^K \binom{K}{i} \left(\frac{1}{m}\right)^i \left(\frac{m-1}{m}\right)^{K-i}$$

Number of observers choosing data plot

Source: Majumder et al (2013) To appear

LES DIABLERETS, FEB 1-4, 2015

15 -46

Consideration THREE

- What is the power of the test?
- There is a choice of type of plot to use. Some will be more optimal than others.
- Signal strength will be defined as “proportion of observers who identify the data plot”.
- Enables the comparison of different plot designs.
- Signal strength equals power, when only plot design changes.

Source: Majumder et al (2013) To appear

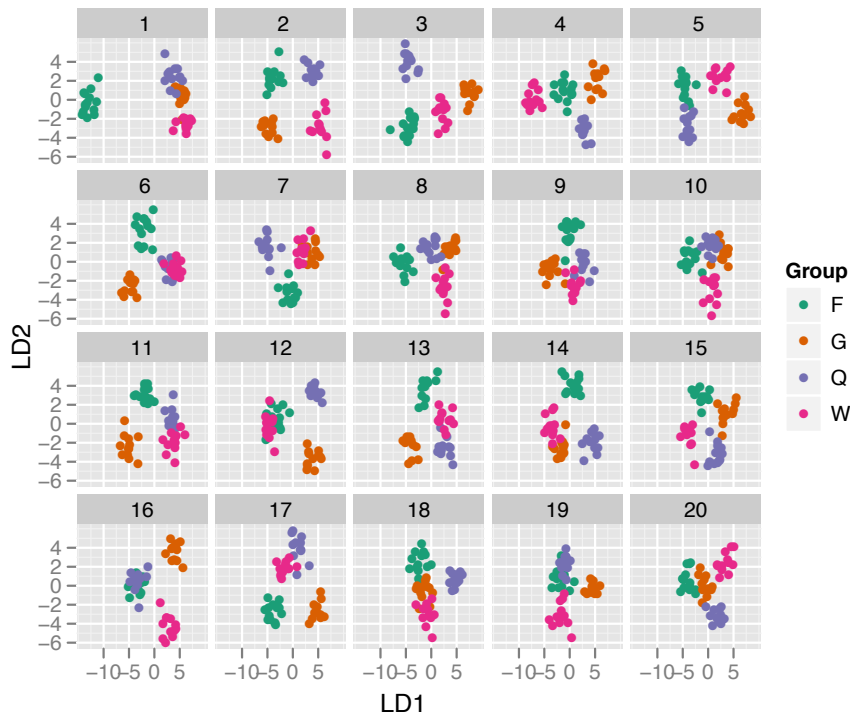
Hofmann et al (2012) InfoVis '12 Proc

LES DIABLERETS, FEB 1-4, 2015

16 -46

What we learn

- Wasps data is no different from random assignment of species label
- Difference between groups was due to sparseness of high-dimensional space



-46

Your Turn

For the following lineup of 20 plots, pick the ONE plot that is most different from the others

- Jot down the number of the plot

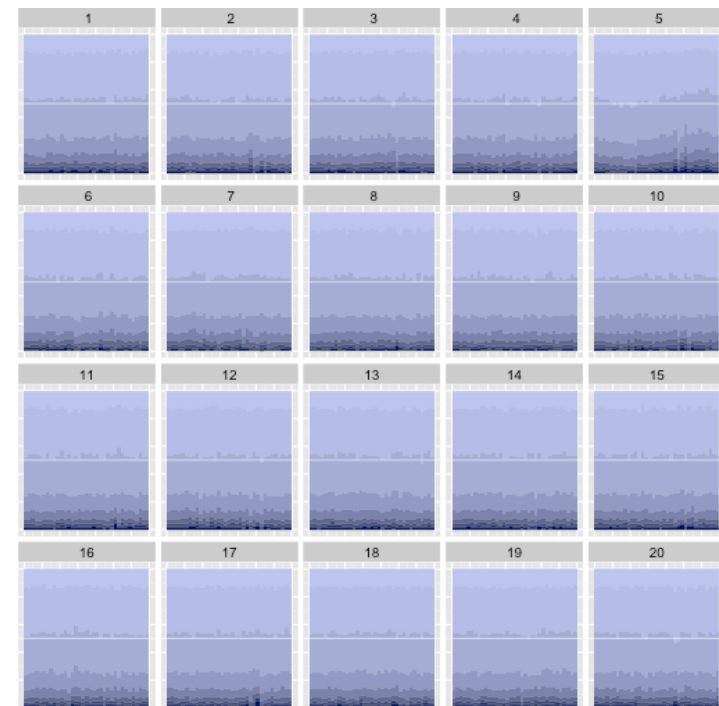
1	2	3	4	5
6	7	8	9	10
11	12	13	14	15
16	17	18	19	20

- Write down a reason for your choice
- And rate how confident that you picked the data on a scale of 1 (very sure) to 5 (don't have a clue)

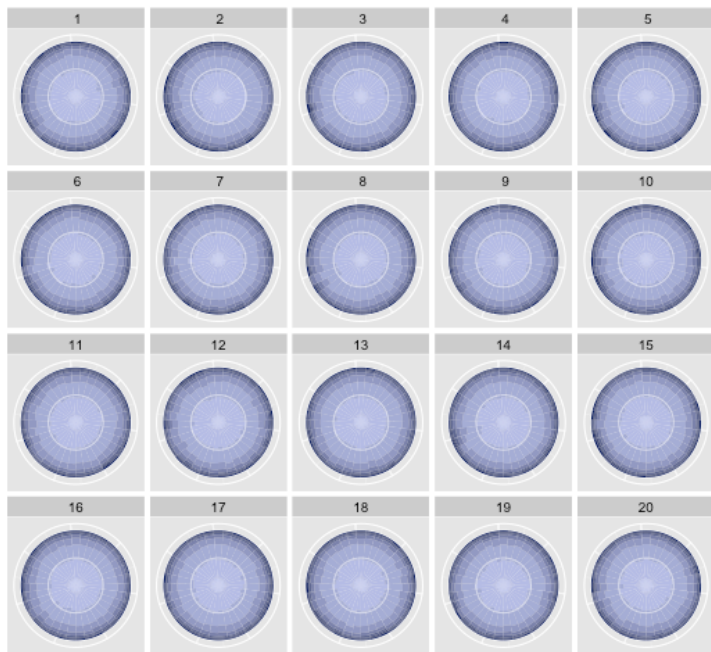
19

LES DIABLERETS, FEB 1-4, 2015

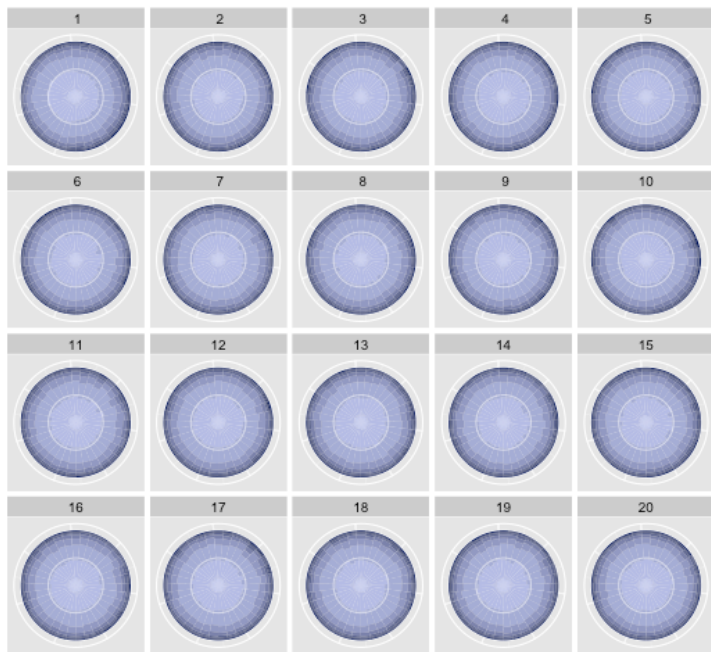
18 -46



20 -46



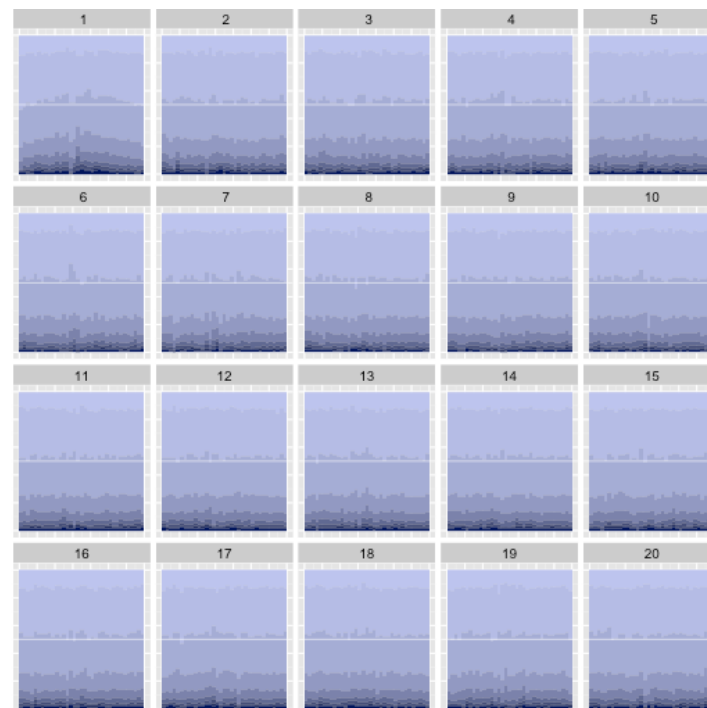
21 -46



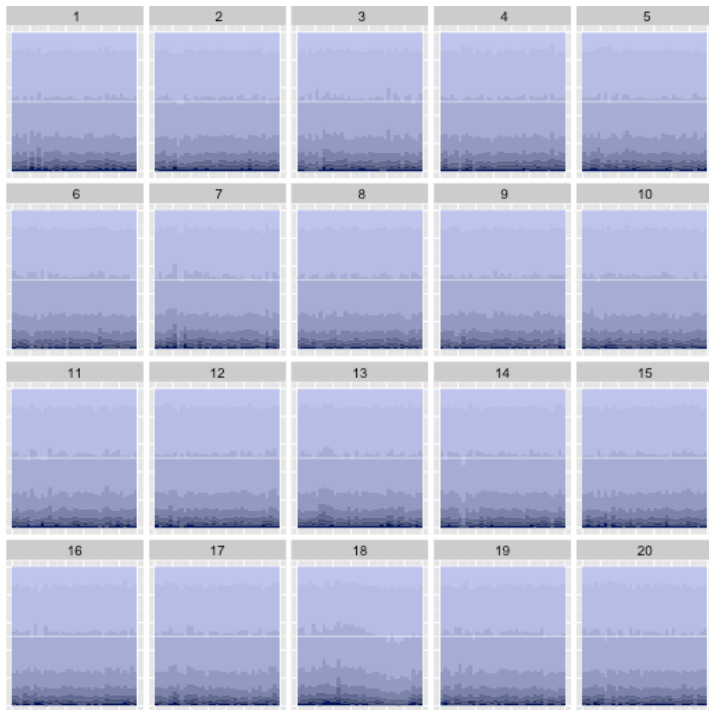
23 -46



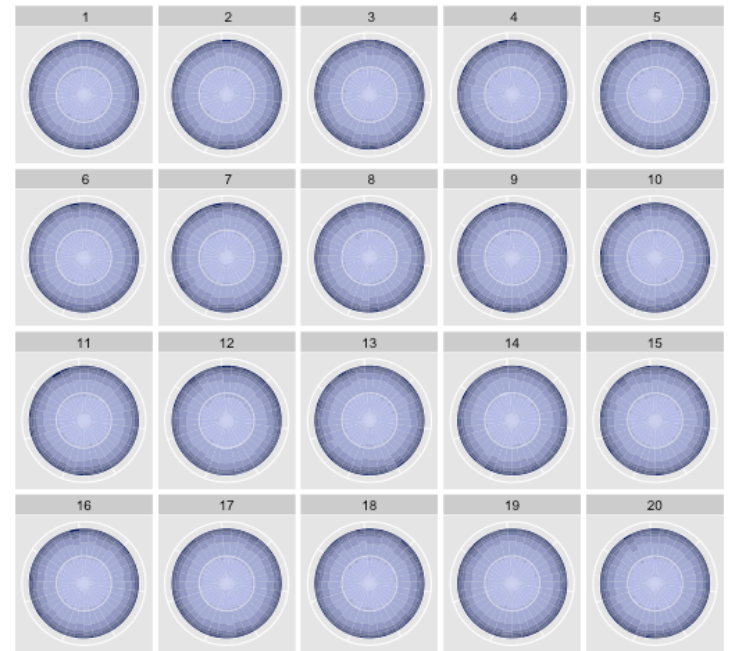
22 -46



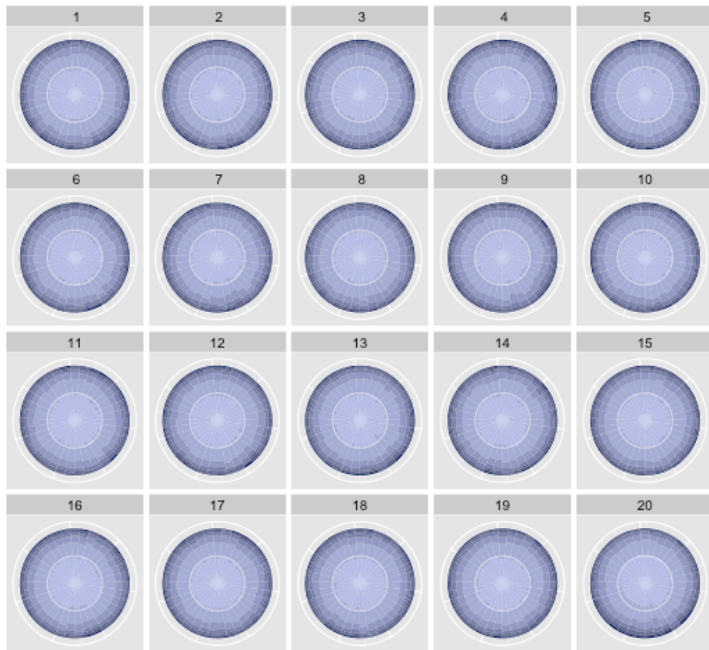
24 -46



25 -46



26 -46



27 -46

Lineup	# Correct	Reason	Confidence
1			
2			
3			
4			
5			
6			
7			
8			

LES DIABLERETS, FEB 1-4, 2015

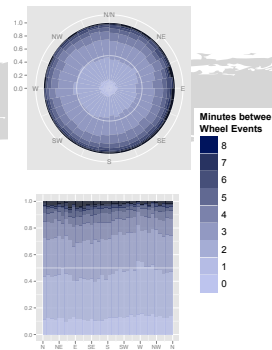
28 -46

Study

- Examine wind direction and airport efficiency.

H_0 : wind direction has no effect on efficiency
against the alternative H_a : wind direction does have an effect.

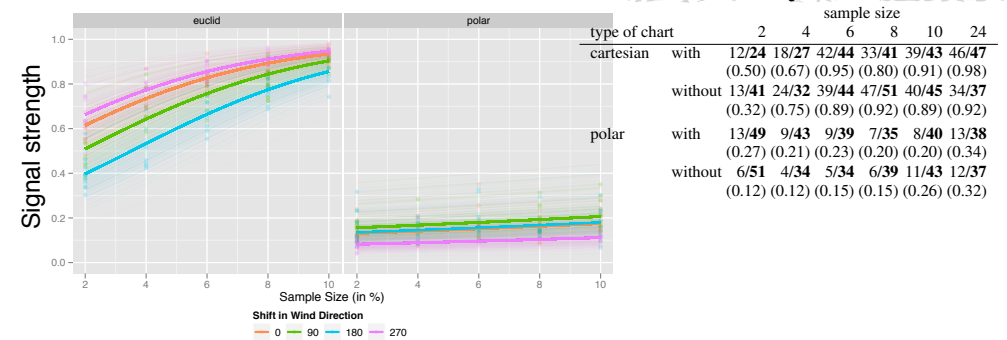
- Decide on best display: (conditional) wind rose charts or bar charts, where each of 36 wind directions the percentage of flights falling into one-minute intervals between successive flights, from zero minutes to eight minutes or more is shown in color scale.



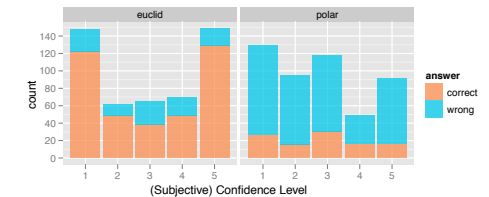
LES DIABLERETS, FEB 1-4, 2015

29 -46

Results of full study



Reported confidence in their choice on plot



LES DIABLERETS, FEB 1-4, 2015

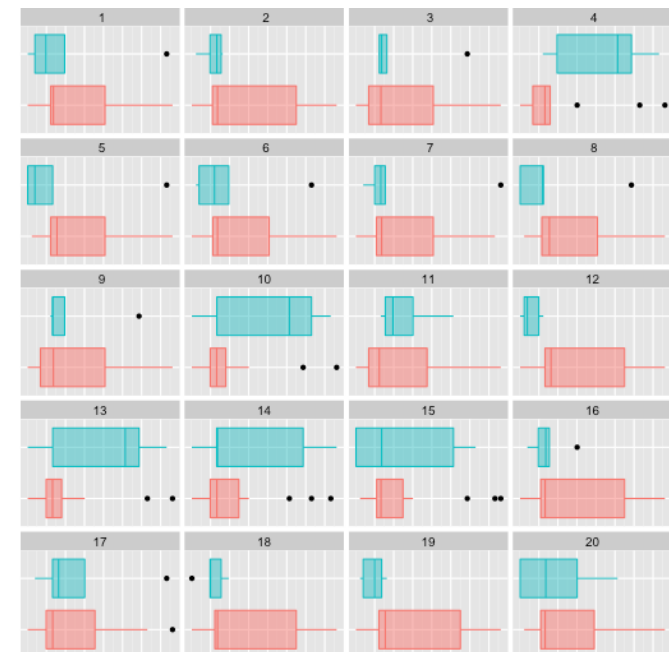
30 -46

A similar study

MOTIVATION: A very small data set of chemical concentrations taken from a superfund clean up site (5 values), compared with samples taken from a normal site (15 values).

Can we see a difference between the two groups, using a side-by-side dotplot? Are side-by-side dotplots better for comparing groups, or side-by-side boxplots, or stacked histograms or density plots?

In which group is the blue group further to the right?

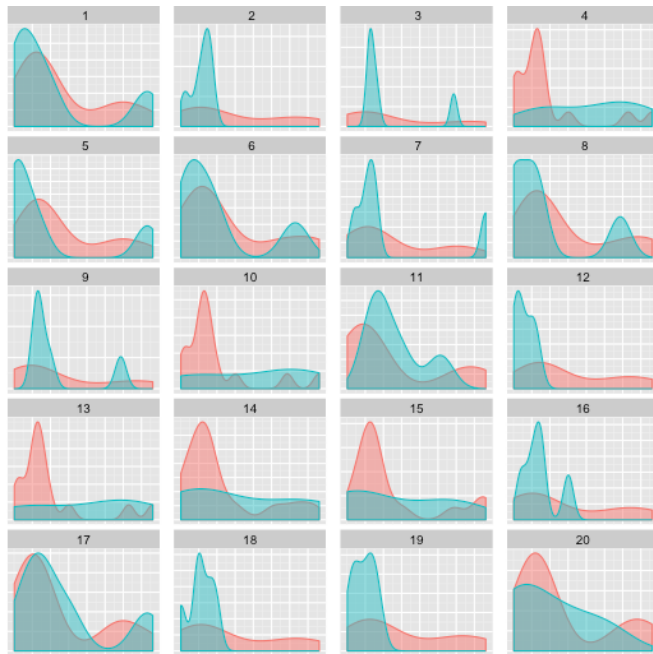


LES DIABLERETS, FEB 1-4, 2015

31 -46

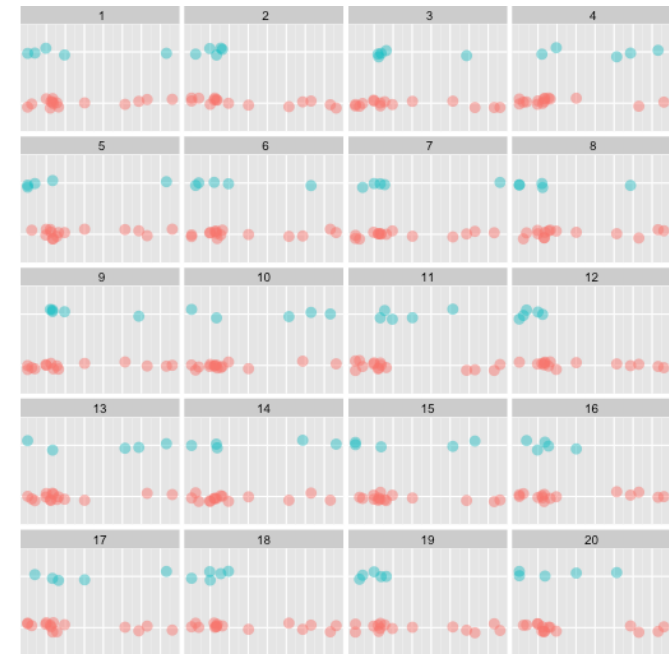
LES DIABLERETS, FEB 1-4, 2015

32 -46



LES DIABLERETS, FEB 1-4, 2015

33 -46



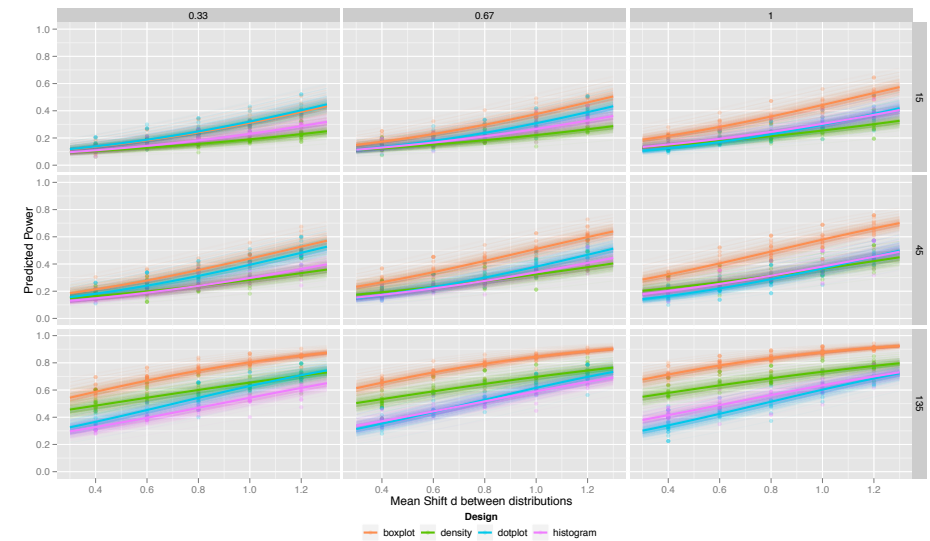
LES DIABLERETS, FEB 1-4, 2015

34 -46



LES DIABLERETS, FEB 1-4, 2015

35 -46



Boxplots beats all other plot designs,
except for really small data sets.

LES DIABLERETS, FEB 1-4, 2015

36 -46

Process

1. Decide on appropriate plot of the data, using good graphical principles.
2. Make the lineup before you have seen the actual data plot.
3. Pick the plot that is different from the rest.
4. If you have already seen the plot of the data, you can show the lineup to someone who hasn't, and use their results.
5. Services like Amazon's Mechanical Turk allow employing independent observers, from a broad cross-section of society.
6. (We are not doing invalid post-hoc testing.)

LES DIABLERETS, FEB 1-4, 2015

37 -46

LES DIABLERETS, FEB 1-4, 2015

38 -46

nullabor package

- Builds on the ggplot2 package for making data plots.
- Generate the lineups automatically, so that you see this before you see the plot of the data.
- Encrypts the location of the actual plot, for you to decrypt when you're ready.

LES DIABLERETS, FEB 1-4, 2015

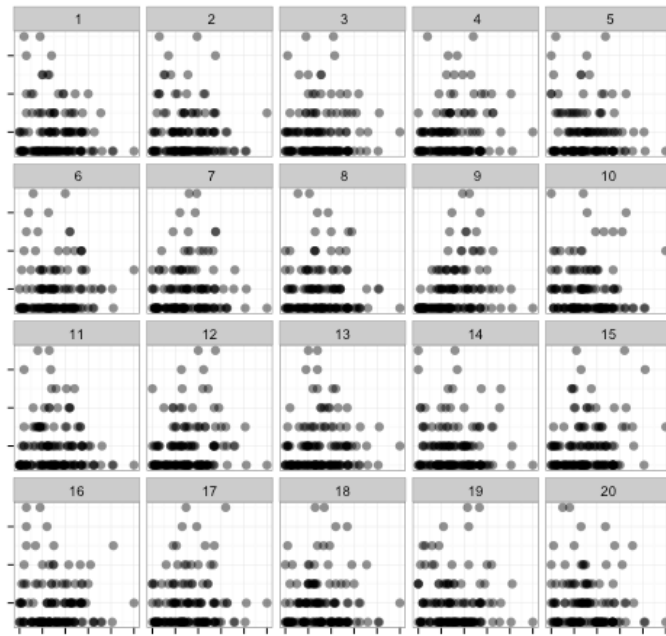
39 -46

Tennis statistics

- The relationships between round and performance statistic were not regular.
 - Simple linear model may not pick up if there is a relationship between the variables.
 - Lineups can be used to determine if there is a **real** relationship.
 - Permutation of "round" label is used to generate nulls
- Ready?

LES DIABLERETS, FEB 1-4, 2015

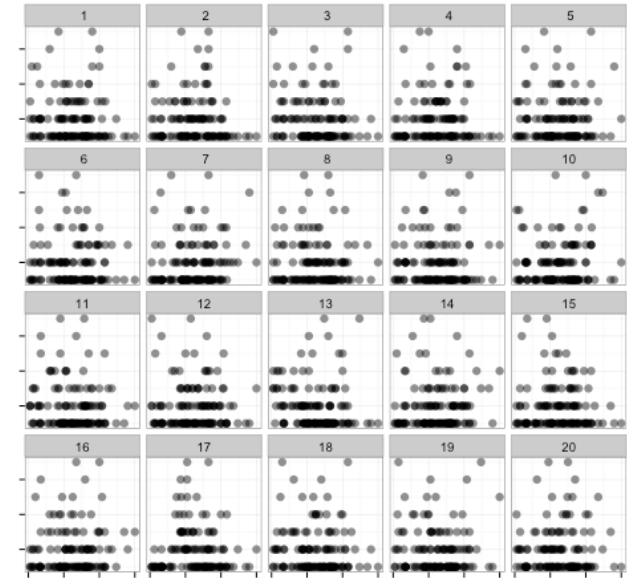
40 -46



decrypt("Y25b yGKG Uu I1OUKU1u Xj")

LES DIABLERETS, FEB 1-4, 2015

41 -46



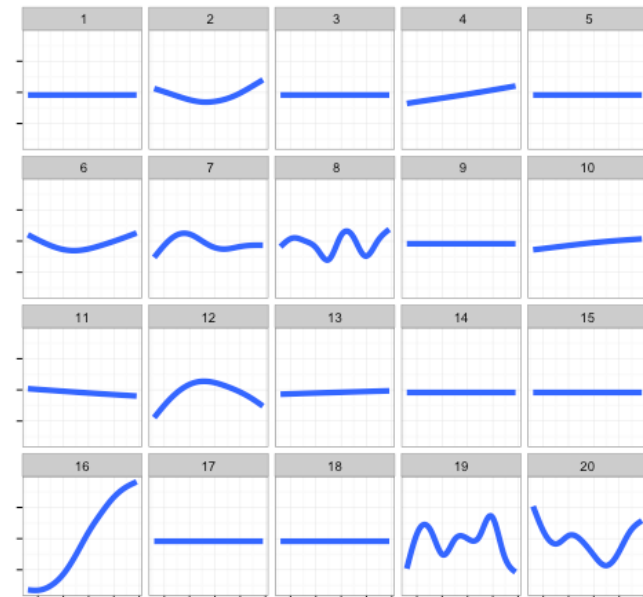
decrypt("Y25b yGKG Uu I1OUKU1u Xj")

LES DIABLERETS, FEB 1-4, 2015

42 -46

Climate change

- Look at neighboring stations, to that examined in first day's slides.
- Produce lineups of temperature measurements.
- Is there a temperature trend?



LES DIABLERETS, FEB 1-4, 2015

43 -46

decrypt("Y25b yGKG Uu I1OUKU1u Xj")

LES DIABLERETS, FEB 1-4, 2015

44 -46

Summary

- ☞ This is an exciting time to be a statistician, if you are interested in data analysis.
- ☞ Data is so widely available and accessible, that it is easy to extract, analyze and learn about our world.
- ☞ It is possible to both explore and yet maintain a healthy skepticism.

EDA & Inference

If the plot that is picked is the plot of the real data, this is statistical significance, and a p -value can be placed on the discovery.

- ☞ Buja et al (2009) RSPT A (econ eg)
- ☞ Wickham et al (2010) InfoVis/TVCG
- ☞ Hofmann et al (2012) InfoVis/TVCG
- ☞ Majumder et al (2013) JASA
- ☞ Roy Chowdhury et al (2014) Comput. Stat.
- ☞ Zhao et al (2014) IJITAS