

FISH 604

*Module 3:*

Exploratory data analysis

---

Instructor: Franz Mueter

Lena Point, Rm 315

796-5448

[fmueter@alaska.edu](mailto:fmueter@alaska.edu)

# Objectives & Outcomes

MSL

FISH

604

## **You should appreciate...**

- ... the importance of visually analyzing your data
- .. the variety of methods available to displaying multi-dimensional data

## **You should know...**

- ...how to assess (approximate) normality
- ... how to detect outliers and what to do in the presence of outliers

## **You should be able to ...**

- ...quickly and efficiently explore the main features of simple and complex data sets
- ... identify and apply appropriate data transformations as needed



# Exploratory data analysis

---

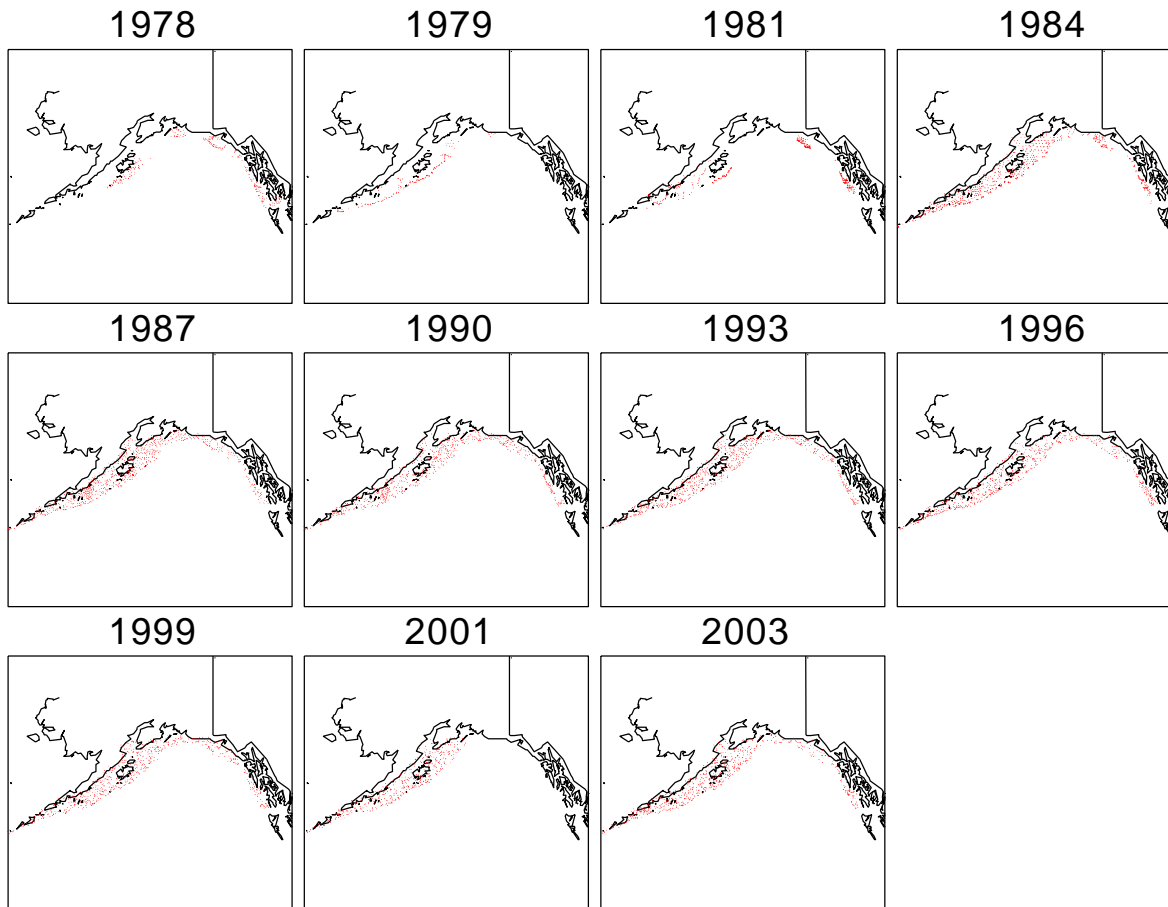
- Visualizing data
- Assessing distributions
- Outliers
- Standardization
- Transformations
- Correlations

# Visualizing data

- Explore spatial / temporal structure
- Detect relationships / correlations
  - Scatterplot
  - Scatterplot matrix
  - Parallel coordinates
- Explore grouped data (differences among groups)
  - Co-plots / Trellis graphics / ggplot 'aesthetics' & facets
- Assess distribution
  - Histograms, q-q plots, dotplots, boxplots
- Assess dependence
  - Serial correlation
  - Spatial correlation

# Spatial structure

Quick maps for exploring / mapping data locations & attributes



Spatial data are ubiquitous in ecology & environmental science  
→ GIS useful but not necessary

Figure 1 consists of three maps of the Bering Sea, labeled 2010, 2017, and 2018. Each map shows the distribution of juvenile Pacific halibut. The maps include density contours (black), 50% and 95% confidence intervals (red), and areas of high density (green and blue dots). The maps are plotted on a coordinate system with latitude (57 to 66) and longitude (-175 to -160). The word 'Alaska' is visible on the right side of each map.

- 6

# Spatial patterns

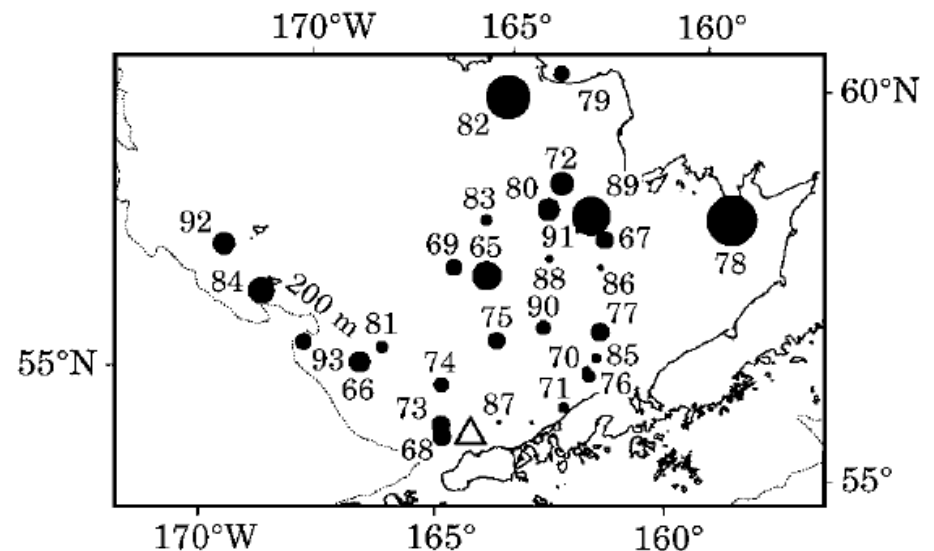
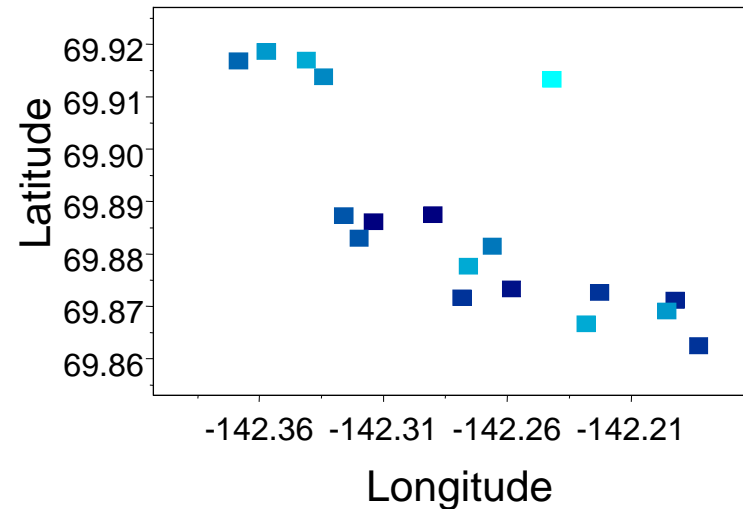
MSL

FISH

604

- Location
- Magnitude
  - symbol type
  - size
  - color

Iron (Beaufort Sea sediments)



# Displaying multiple attributes using 'aesthetics' (ggplot)

MSL

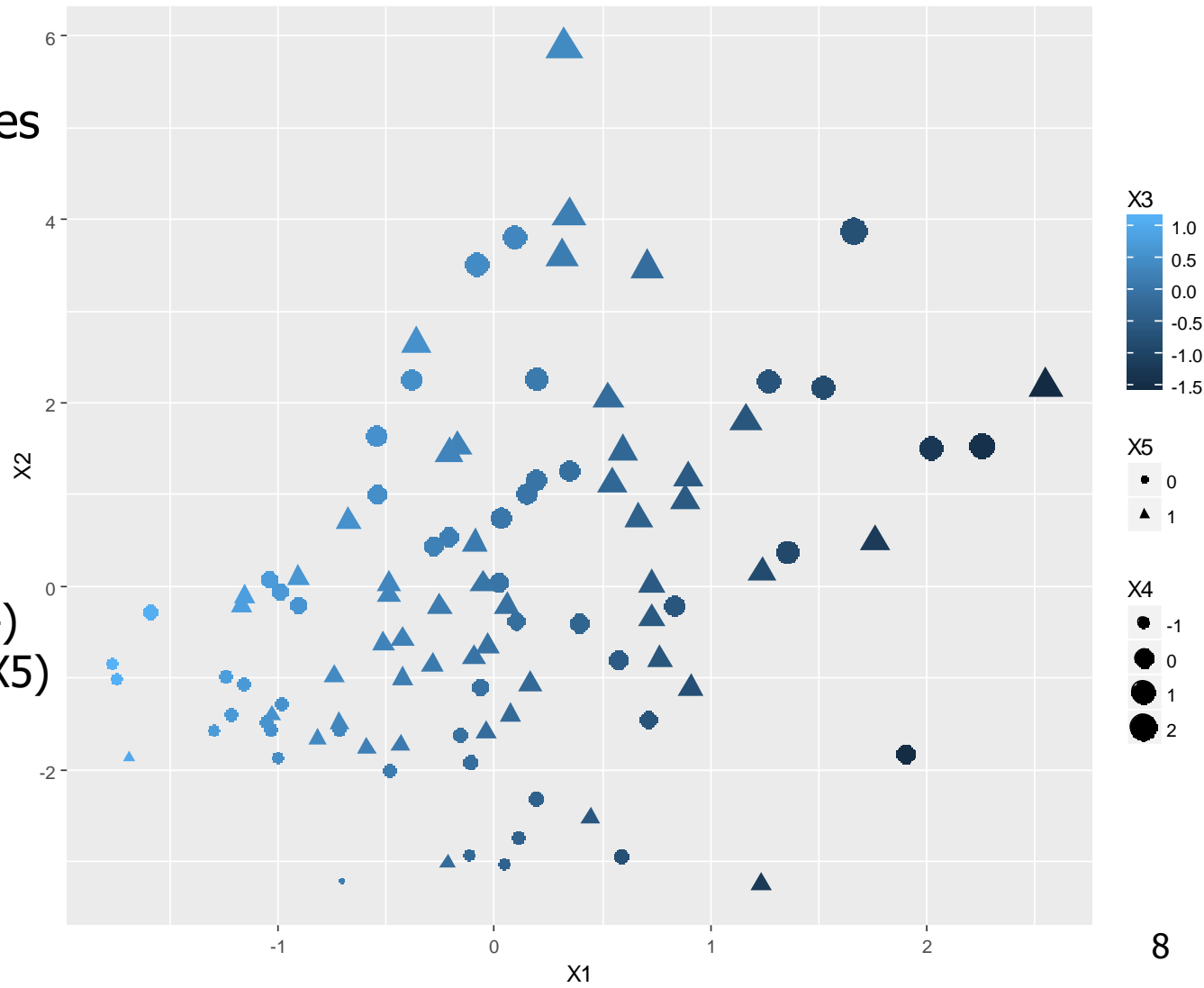
FISH

604

Displaying five variables  
in two dimensions:

- x-axis (X1)
- y-axis (X2)
- color (X3)
- size (X4)
- shape (X5)

Mix of continuous (X1-X4)  
& categorical variables (X5)





# Dotplots: Show each datapoint (whenever possible)!

MSL

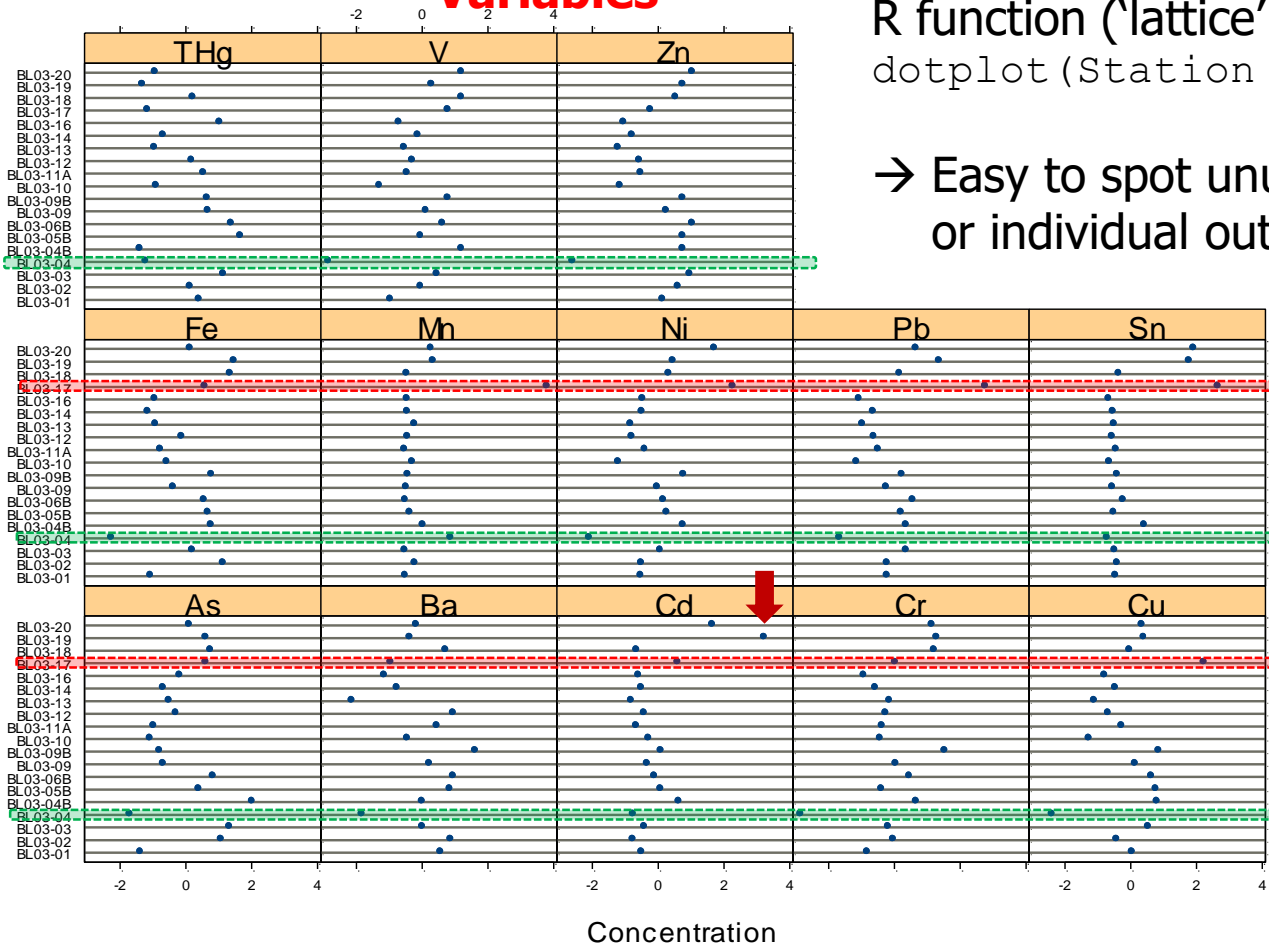
FISH

604

## Metal concentrations, Beaufort Lagoon

**Variables**

**Samples**

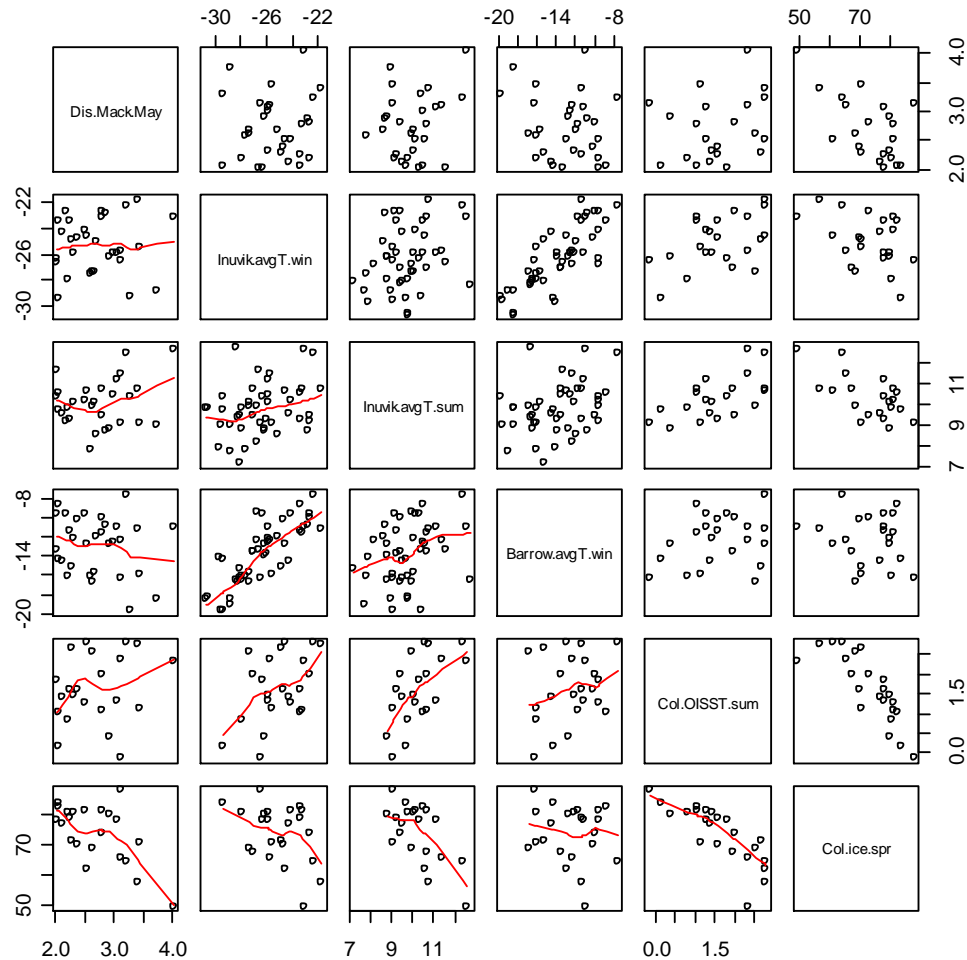


R function ('lattice' package):

```
dotplot(Station ~ Conc | Metal)
```

→ Easy to spot unusual samples or individual outliers

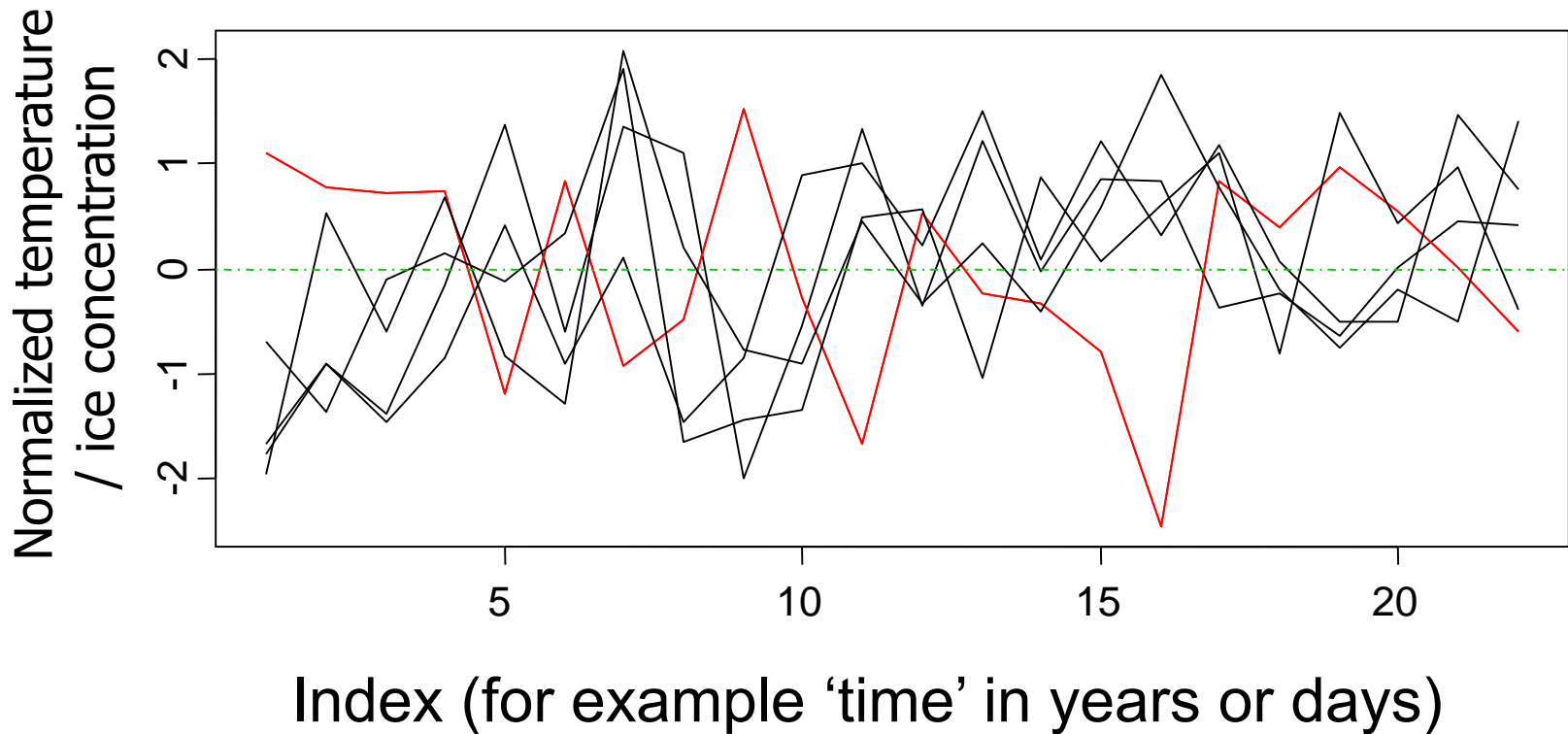
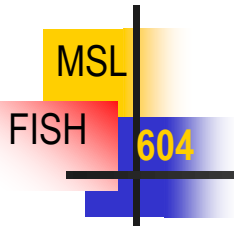
# Relationships among variables: Scatterplot matrix



R function:

`pairs(data, lower.panel=panel.smooth)`

# Relationships among variables: Parallel coordinates



# Multivariate exploration (e.g. "brushing") (e.g. using GGobi)

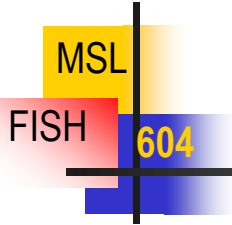
MSL

FISH

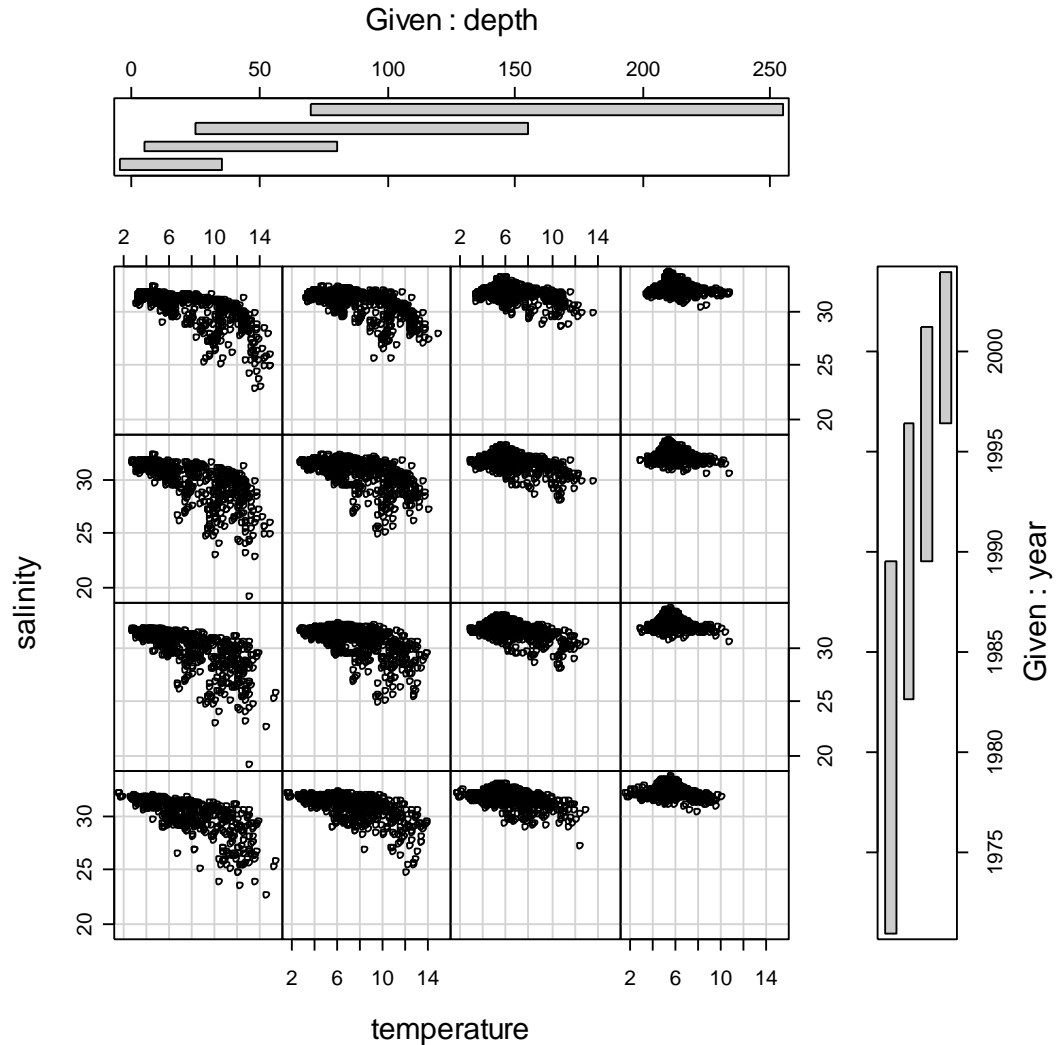
604

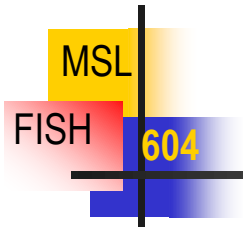


# Exploring grouped data: Co-plots

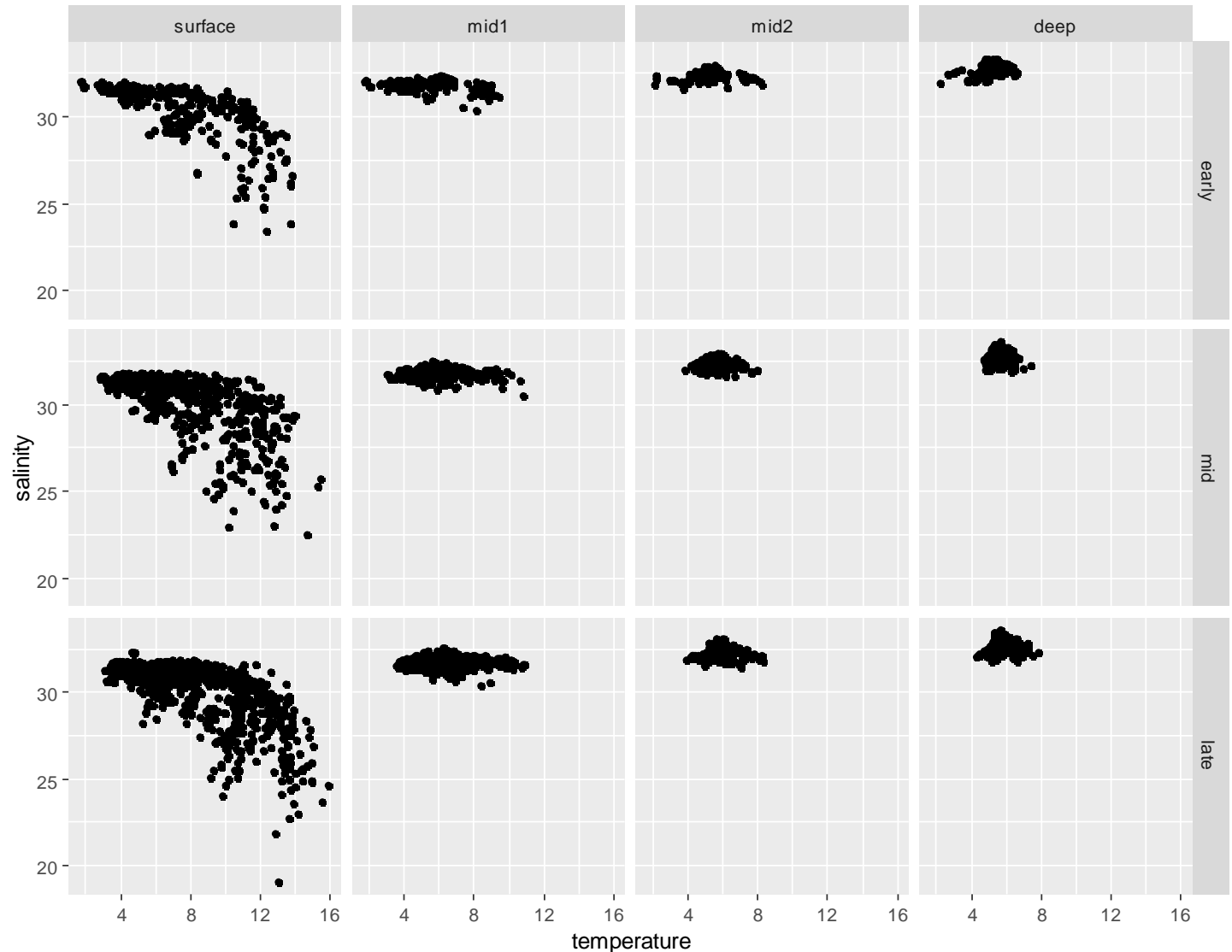


- GAK 1 data
- T-S plots by
  - depth and
  - time period





# Exploring grouped data: ggplot (non-overlapping groups)



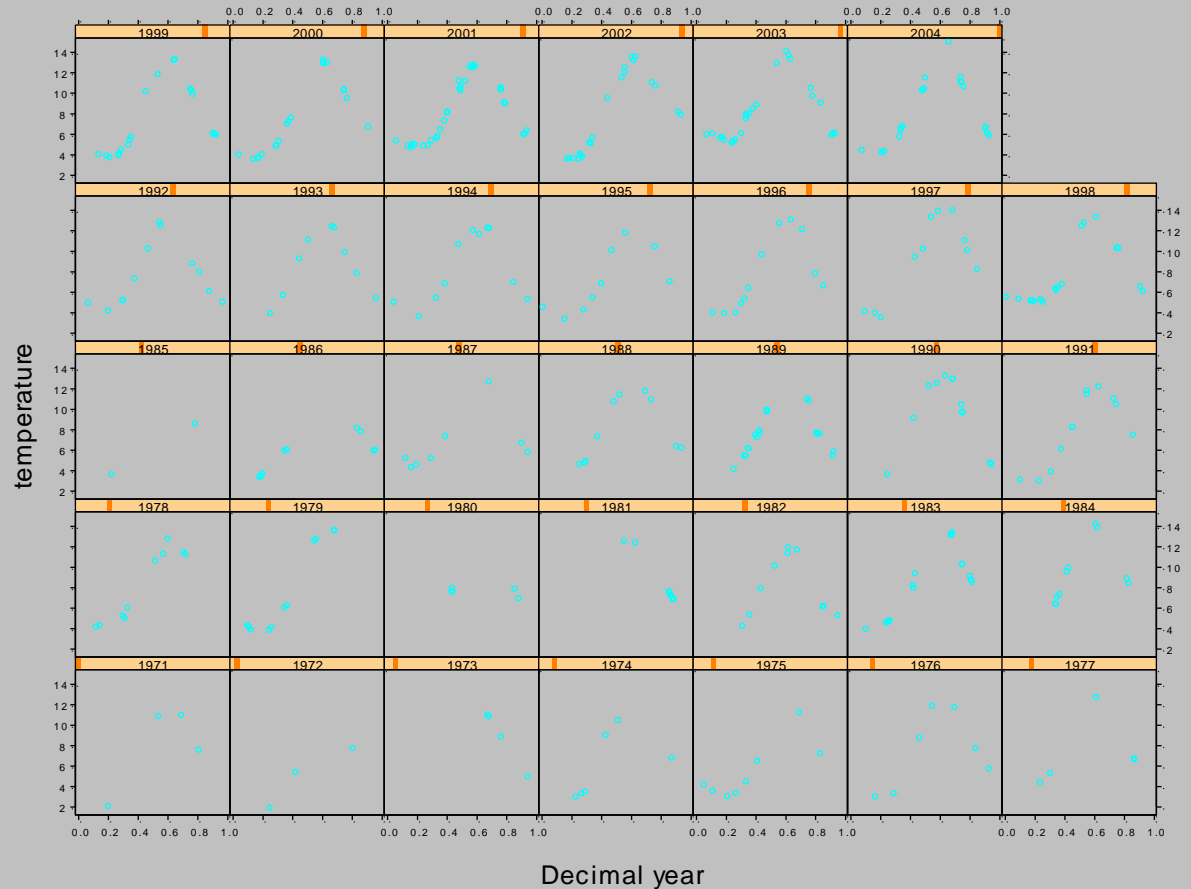
# Exploring grouped data: Trellis graphics

MSL

FISH

604

Ocean  
temperature  
data collected  
irregularly over  
multiple years

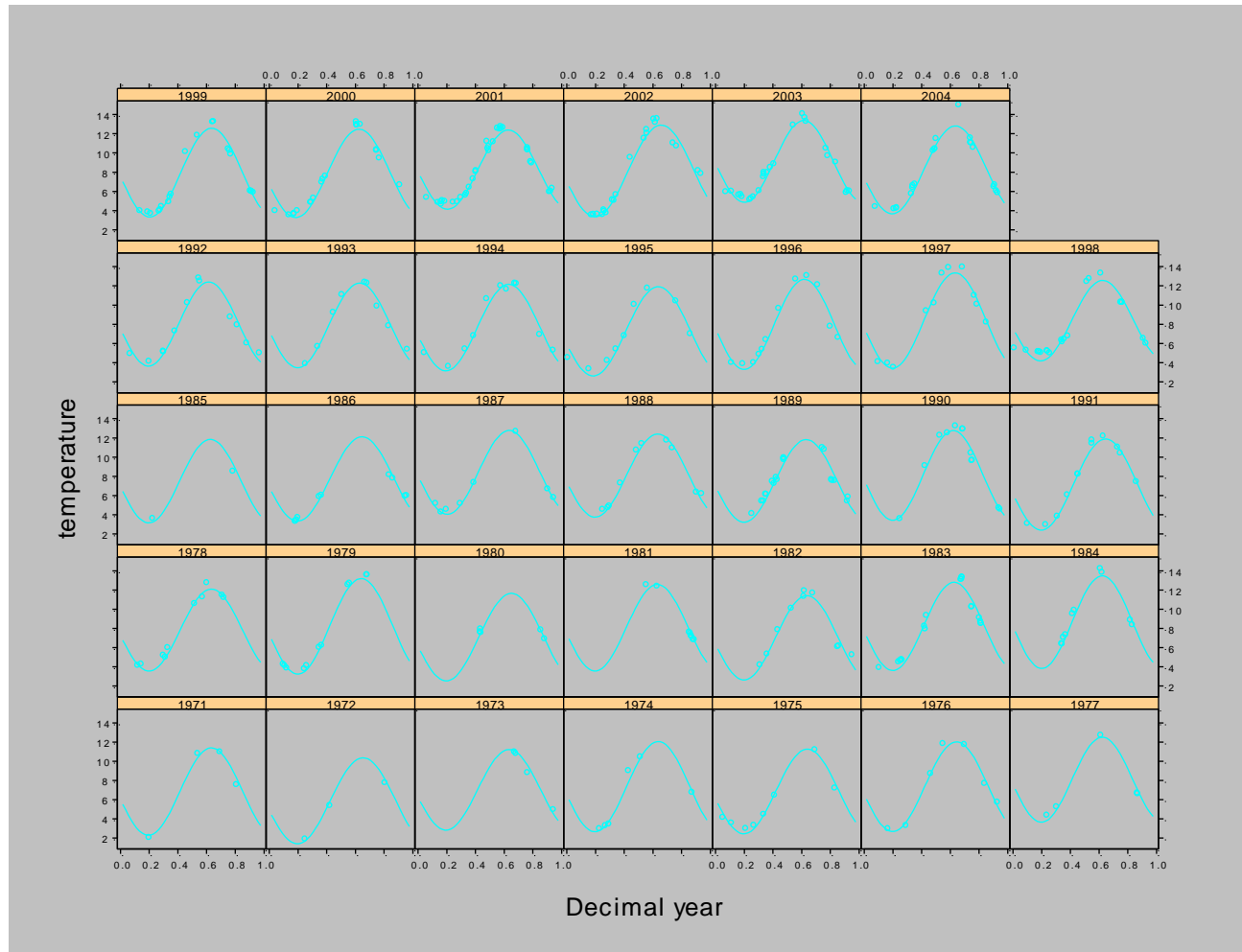


R code:  
`xyplot(temperature ~ dec.year | Year, data = GAK1, subset = depth==0)` 15

# Visualizing model fits

Fitted non-linear  
mixed-effects  
model

→ Learn how to  
visualize  
simple and  
complex  
model fits



R code (requires library nlme):  
`plot(augPred(fitted.object))`



# Visualizing model fits

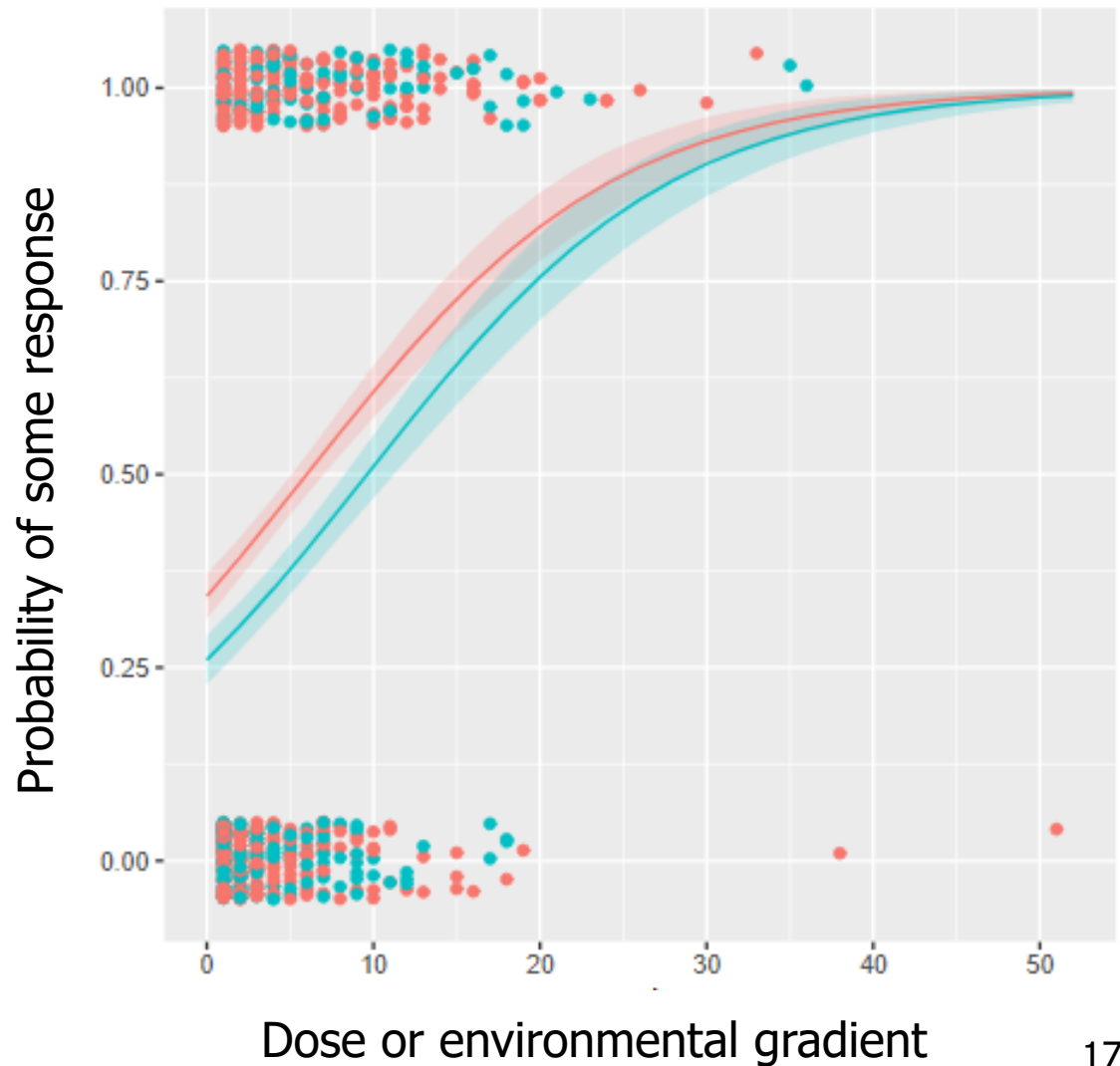
MSL

FISH

604

`ggplot()` is great  
for visualizing model  
fits for many  
“standard” models,  
but be aware of  
what you are doing!

Example: Logistic  
regression fit with  
two groups



# Assessing distributions

## Visually

- histograms / density plots / boxplots
- quantile-quantile plots (q-q plots)

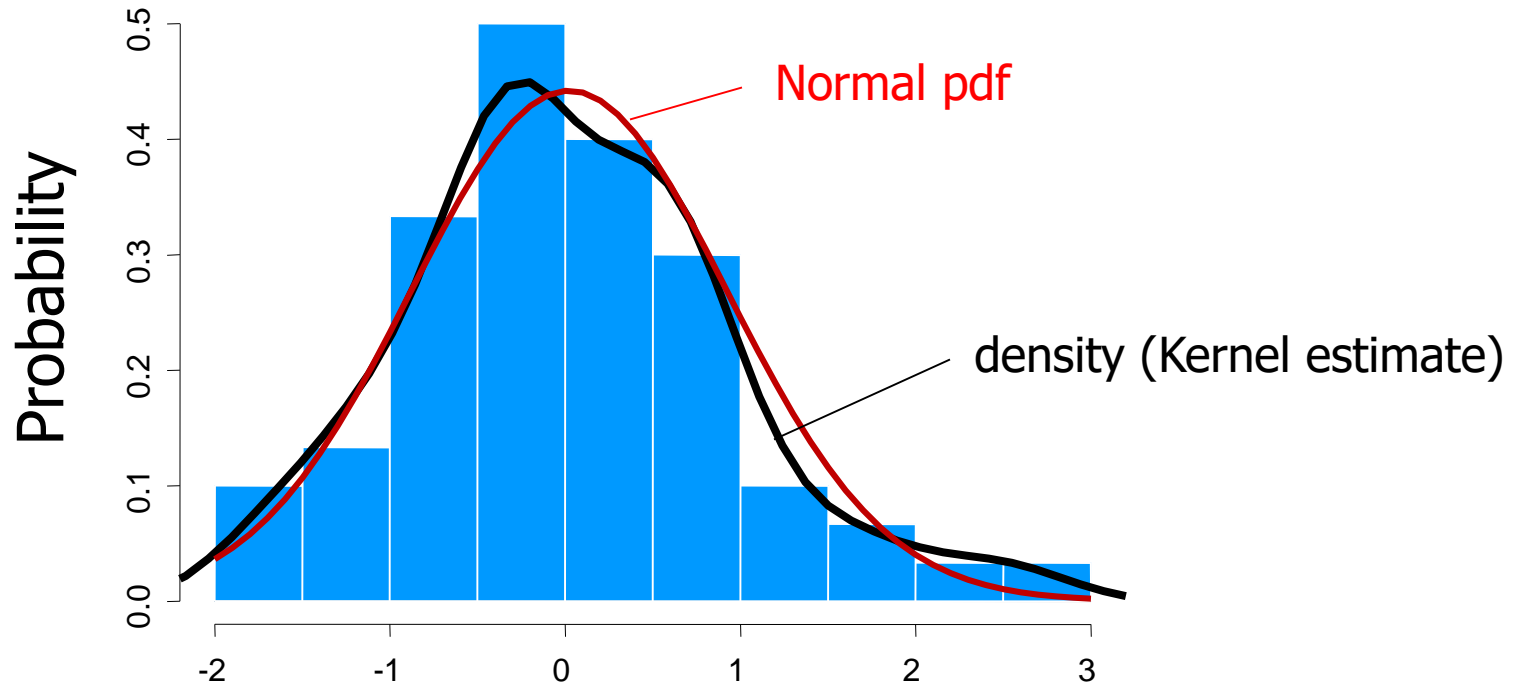
## Statistical

- tests for normality
- testing against any distribution
  - Kolmogorov-Smirnov test

## R functions

```
hist(); histogram(); density(); boxplot()  
qqplot(); qqnorm(); qqline()  
shapiro.test(); ks.test()
```

# Histogram with density



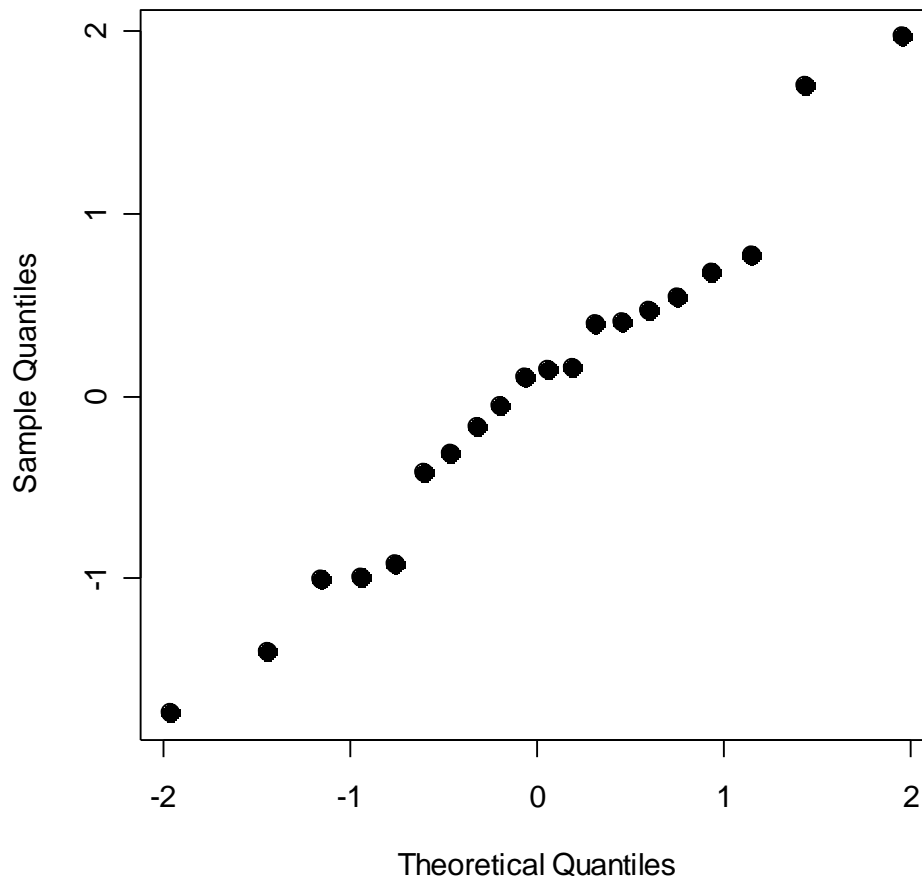
Algorithms to determine bin size:

Scott (1979): "optimal" bin width  $= 3.49 \cdot \sigma \cdot N^{-1/3}$

Freedman and Diaconis (1981): "robust" bin width  $= 2 \cdot IQR \cdot N^{-1/3}$

# quantile-quantile plots

Normal q-q plot



- Sample quantiles
  - Sort data:  $x_{(i)}$
  - $f_i = (i-0.5)/n$
  - $x_{(i)}$  is the  $f_i$  quantile of the data
- Theoretical quantiles
  - $q_{\mu,\sigma}(f_i)$  is normal  $f_i$  quantile
  - $q_{\mu,\sigma}(f) = \mu + \sigma q_{0,1}(f)$
- Normal q-q plot
  - $x_{(i)}$  vs.  $q_{0,1}(f_i)$

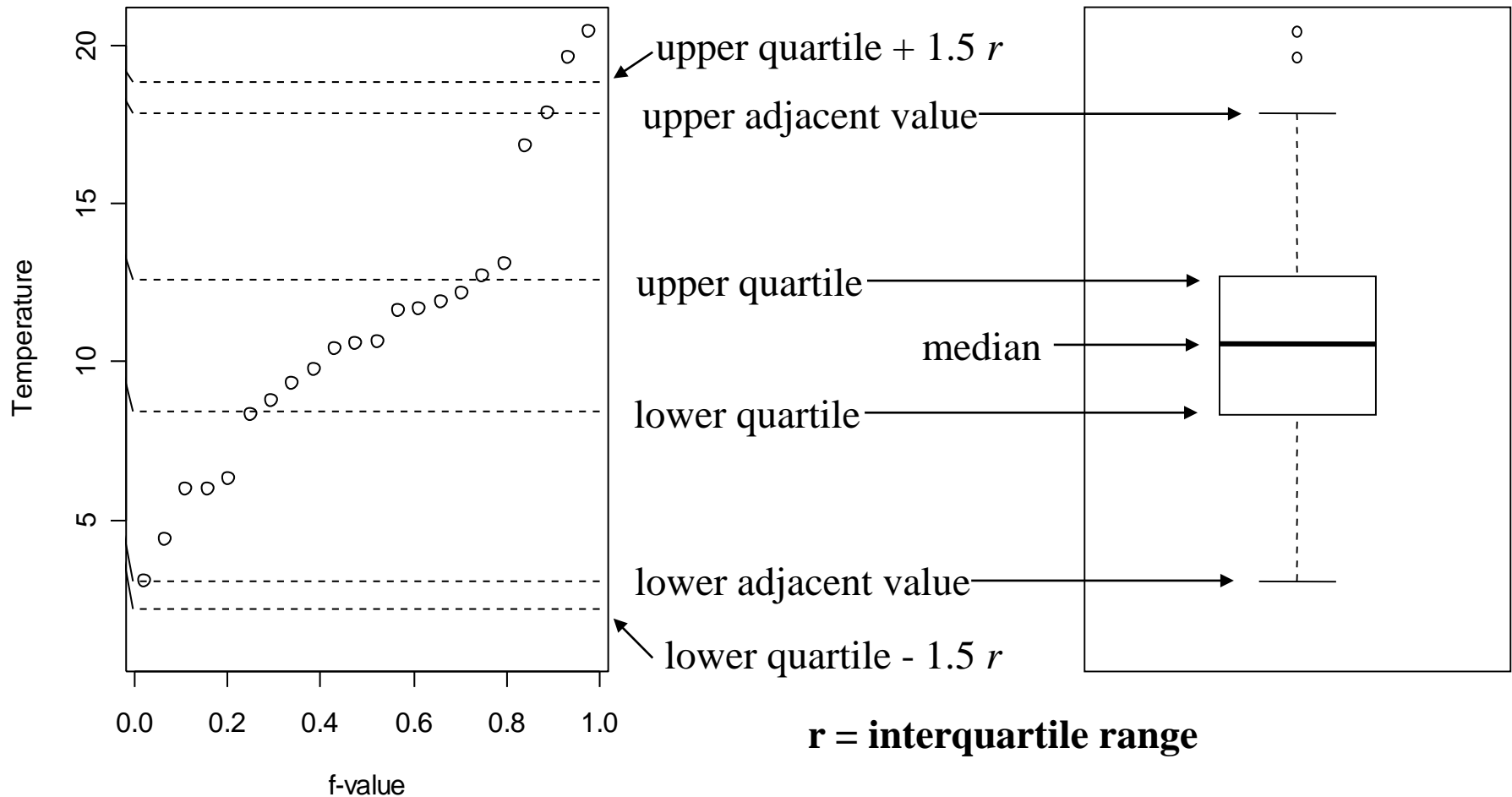
→ usually sufficient as a diagnostic check of normality assumption!

# Boxplot

MSL

FISH

604



- Extremely effective for visualizing / comparing many distributions
- Be aware of what boxes & whiskers show and explain it in figure captions



# Tests for normality

---

- Shapiro-Wilk test
- Kolmogorov-Smirnov test
- Many others!!!

## R functions

```
shapiro.test(); ks.test()
```

- Remember that we only care about normality of residuals from a linear model, NOT normality of the response variable (or the independent variables)
- Residuals from other models (e.g. GLM) often transformed to 'approximate' normality for visual assessment only

# Shapiro-Wilk test

MSL

FISH

604

- Tests the null hypothesis that a sample  $x_1, \dots, x_n$  came from a normally distributed population

- Test statistic: 
$$W = \frac{\left( \sum_{i=1}^n a_i x_{(i)} \right)^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

where  $x_{(i)}$  is  $i^{\text{th}}$  order statistic

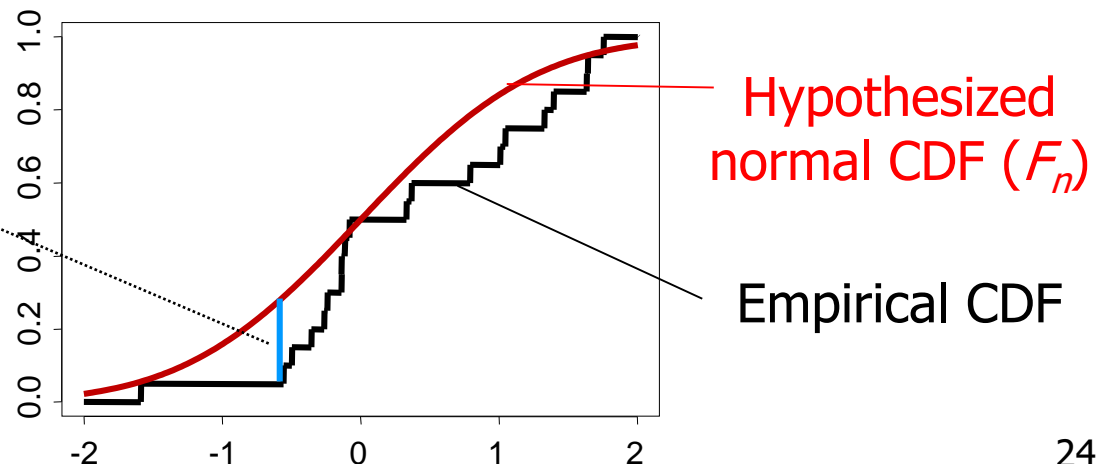
$a_i$  are constants that depend on expected values and covariance of order statistics

- Null hypothesis is rejected if  $W$  is too small

# Kolmogorov-Smirnov test

- Null hypotheses:
  - Two samples have the same underlying probability distributions
  - A single sample comes from a hypothesized distribution (e.g. normal)
- Test is based on maximum distance between two cumulative distribution functions:

$$D = \max \left( F_n - \frac{i-1}{N}, \frac{i}{N} - F_n \right)$$





The logo consists of a yellow square with 'MSL' in black, a red square with 'FISH' in black, and a blue square with '604' in yellow. A black crosshair is overlaid on the squares.

# Reading assignment

---

- Zuur et al. (2007). **Analyzing Ecological Data**. Springer.  
*Chapter 4: Exploration*  
(see pdf posted on Canvas)



# Further reading

---

## ***Graphical analysis***

- Cleveland, W.S., 1993. Visualizing Data. AT&T Bell Laboratories, Murray Hill, NJ.
- Wilkinson, L. 1999. The Grammar of Graphics. Springer, New York
- Yau, Nathan 2013. Data Points: Visualization that means something. Wiley.

## ***Multivariate data***

- Cook, D., Swayne, D.F., and Buja, A. 2007. Interactive and Dynamic Graphics for Data Analysis: With R and GGobi. Springer, New York.

## ***General exploratory analyses***

- Zar, J.H., 1984. Biostatistical Analysis. Prentice-Hall, Englewood Cliffs, NJ.
- Barnett V, 2004. Environmental Statistics: methods and applications, John Wiley & Sons, Chichester, England.