

FISH 604

Module 3:

Exploratory data analysis

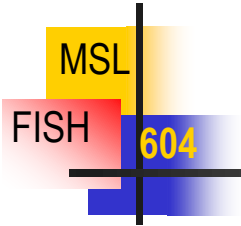
Instructor: Franz Mueter

Lena Point, Rm 315

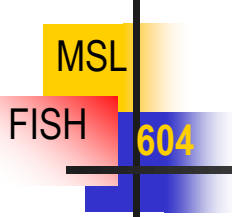
796-5448

fmueter@alaska.edu

Review/preview



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



Objectives & outcomes

Objectives

- Review standard methods for **Exploratory Data Analysis** to conduct prior to statistical modeling

Outcomes

- Know how to detect outliers and what to do in the presence of outliers
- Be able to identify when and how to apply appropriate data standardizations
- Be able to identify when and how to apply appropriate data transformations
- Be proficient in quickly exploring correlations among a set of variables (graphically & statistically)



Outlier detection

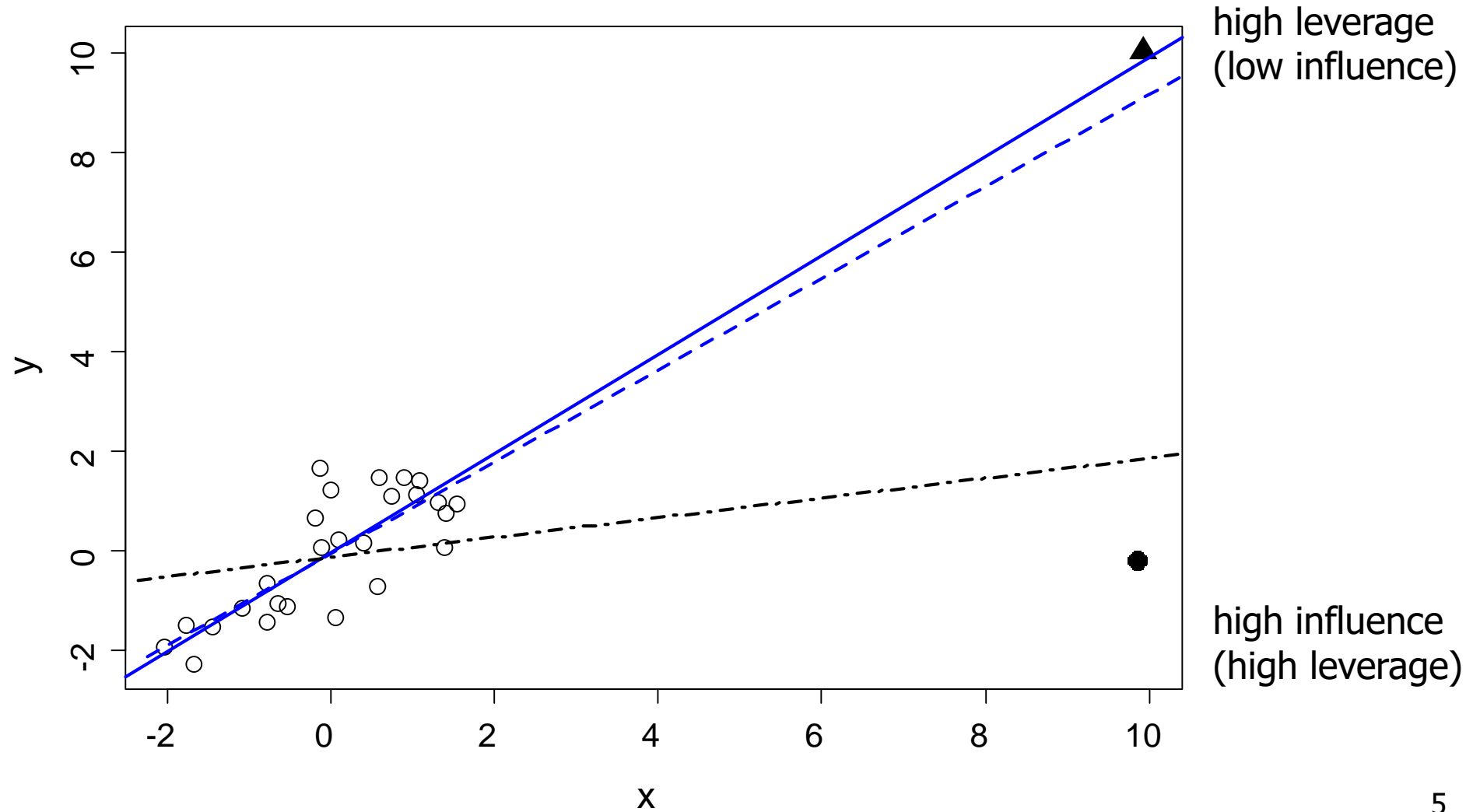
- Why?
 - Tests and model fitting algorithms (least-squares, likelihood) are often sensitive to outliers
- How?
 - Graphical, distance from mean
- What to do?
 - Ignore
 - Eliminate (Compare results)
 - Use robust measures / estimation

Outliers: Influence & leverage

MSL

FISH

604



Outliers: Leverage & Influence

MSL

FISH

604

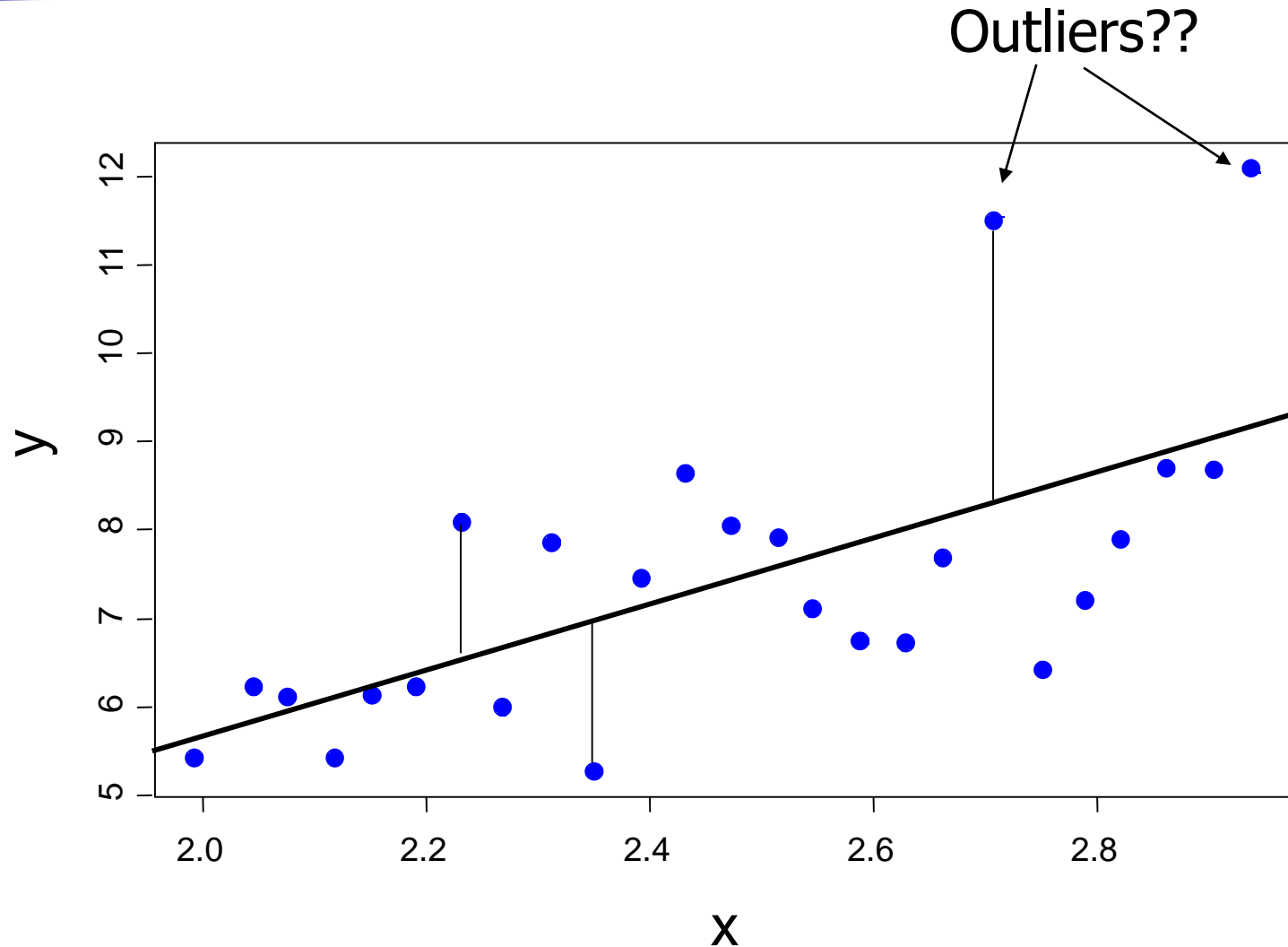
Leverage

- Defined as diagonals of "hat matrix": H_{ii}
- Large values often due to extreme values in X
- Rule of thumb: If leverage $> 2p/n$, look at the data point more closely
- Point with high leverage may or may not also be influential!

Influence

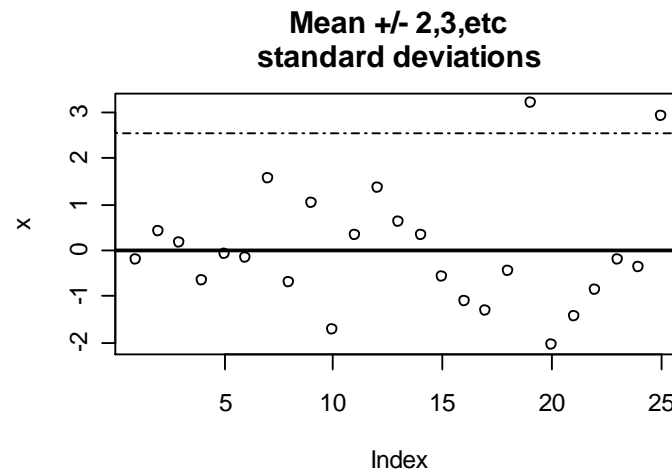
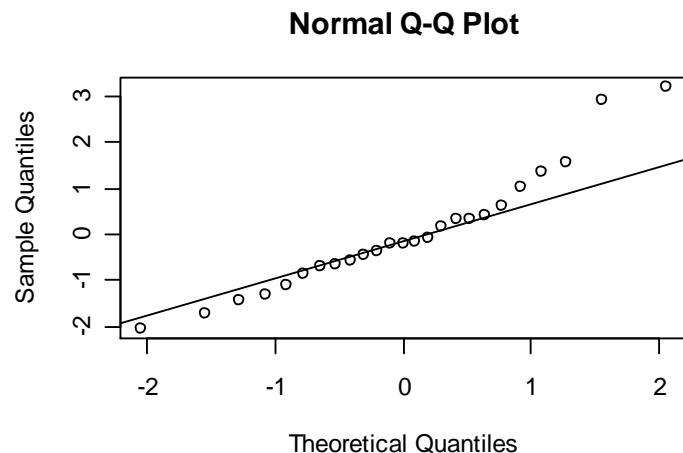
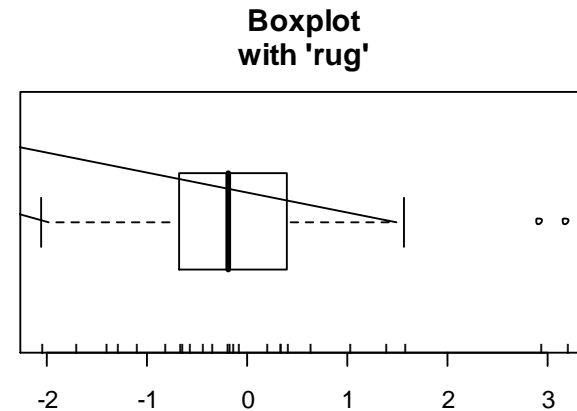
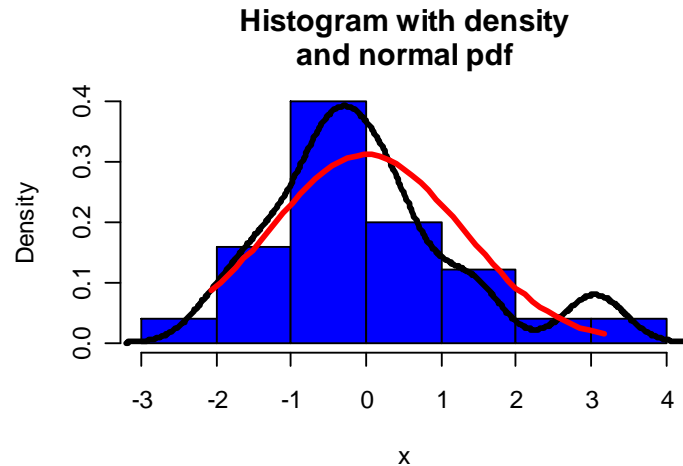
- Removal of influential point results in large change in fit!
- Usually measured by change in predicted values or regression coefficient(s) when i^{th} observation is removed,
- Cook's distance is a common measure of influence

Regression with outliers?



Outlier detection (linear models or 'normalized' residuals)

Examine distribution of residuals from a model fit

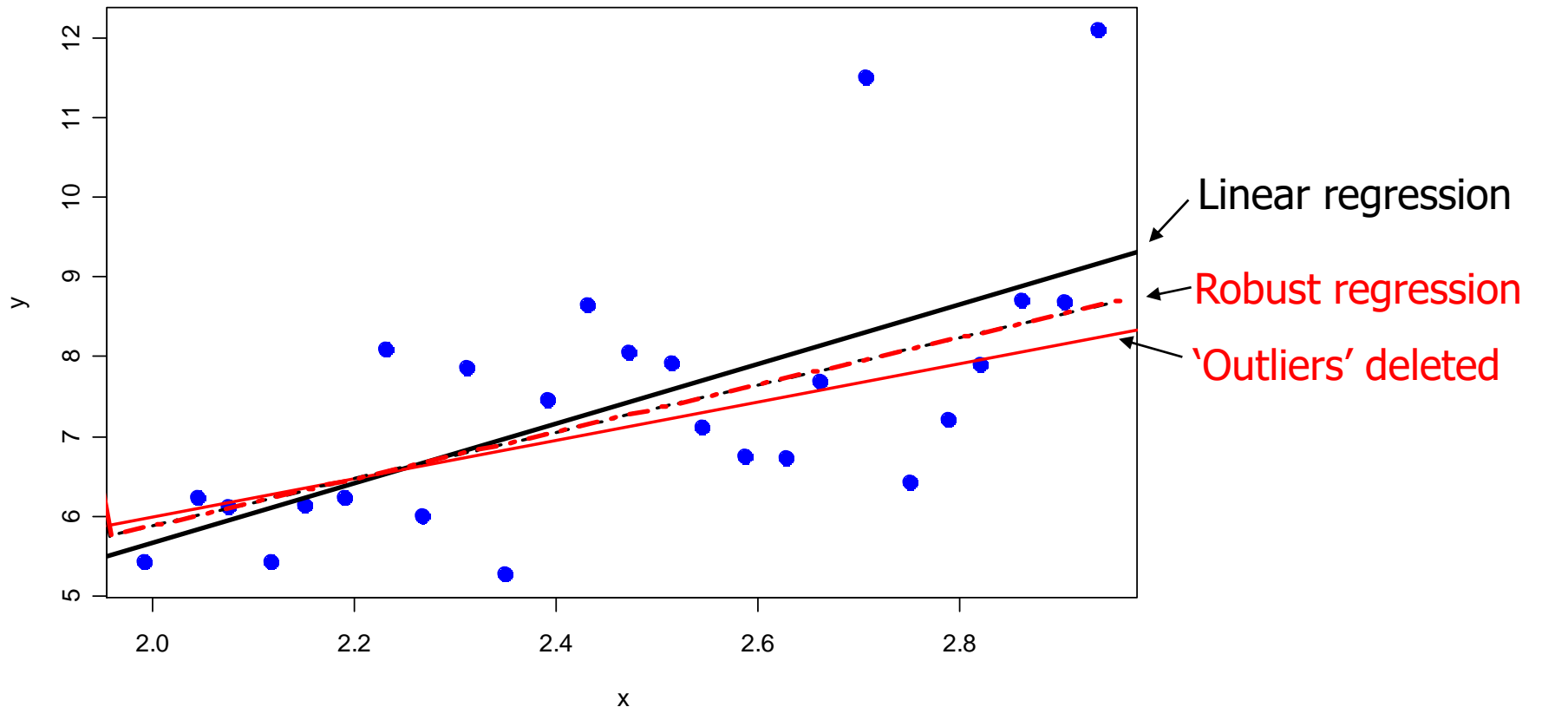


Robust regression

MSL

FISH

604



R code: `plot(x, y)`
`fit <- lm(y~x); abline(fit)`
`fit.rob <- rlm(y~x); abline(fit.rob)`

(`rlm` in MASS package)

Dealing with outliers:

Re-run analysis with outliers removed

Do the conclusions change when the case is deleted?

NO

YES

Proceed with case included.
Examine case to see if anything
can be learned from it.

(acknowledge outlier)

*Is there a reason to believe that
outlier comes from a different
population (or is an error)?*

YES

NO

Omit case and proceed!

*Does the case have unusually
"distant" explanatory variables?*

YES

NO

Omit the case and proceed!
Report conclusions for the reduced
range of explanatory variables

Need more data or info to resolve!
Report conclusions for analysis w/
and w/out the influential case!



Standardizations

- Why?
 - Plot / compare variables on common scale
 - Variables in multivariate analyses
 - To make sure each variable has the same influence
 - Multiple regression
 - Stabilize numerical estimation methods
 - **Make coefficients comparable**
- How?
 - **Standardize to $N(0,1)$**
 - Standardize to common maximum
 - Standardize relative to mean

Standardization / scaling



- Normalize to $N(0,1)$: (any continuous variable)

$$x' = (x_i - \bar{x}) / sd(x) \quad \rightarrow \text{R function } \text{scale}()$$

- Relative to mean ($x > 0$):

$$x' = (x_i - \bar{x}) / \bar{x} \quad \Rightarrow \quad \bar{x}' = 0$$

$$x' = x_i / \bar{x} \quad \Rightarrow \quad \bar{x}' = 1$$

- Same maximum ($x > 0$):

$$x' = x_i / \max(x) \quad \Rightarrow \quad 0 < x' < 1$$

Transformations

MSL

FISH

604

- Why?
 - To meet regression assumptions:
 - Normalize residuals
 - Eliminate / reduce heteroscedasticity
 - Ensure additivity of effects
 - Reduce leverage of extreme values
 - Linearize relationships (e.g. log-log plots)
- How?
 - Logarithmic (natural log) transformation
 - Square-root transformation
 - Arc-sine transformation
 - Power transformations
 - Box-Cox transformations

Log-transformation

MSL

FISH

604

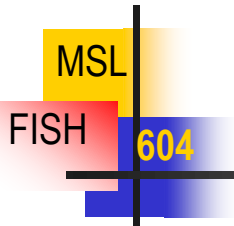
→ Add small constant if x may include 0

- $x' = \log(x + 1)$ OR $x' = \log(x)$
(typically base 10 or **natural logarithm**)
- Log-transform dependent variable in ANOVA when
 - Multiplicative effects are present (Non-additive)
- Log-transform dependent variable in regression when:
 - standard deviations are proportional to mean, i.e. CV is constant (heteroscedasticity)
- Try log-transforming both dependent (y) and independent (x) variable when
 - Errors are multiplicative: $y = f(x) * e^{\varepsilon}$

$$\rightarrow \log(y) = \log\{f(x)\} + \varepsilon$$

where $\varepsilon \sim N(0, \sigma^2)$

Log-transform for additivity



Example: two-way ANOVA

1. Additive effects

		Factor A			
		Level 1	Level 2	Level 3	
Factor B	Level 1	10	→ +10 → 20	→ +5 → 25	↓ +10
	Level 2	20	→ 30	→ 35	

2. Multiplicative effects

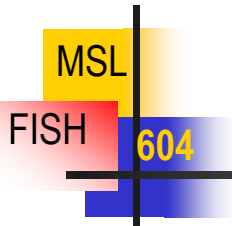
		Factor A			
		Level 1	Level 2	Level 3	
Factor B	Level 1	10	→ *3 → 30	→ *2 → 60	violates ANOVA assumption
	Level 2	20	→ 60	→ 120	

2. Log-transformed (log10)

-> additivity restored!

		Factor A			
		Level 1	Level 2	Level 3	
Factor B	Level 1	1.00	→ +.48 → 1.48	→ +.3 → 1.78	meets ANOVA assumption
	Level 2	1.30	→ 1.78	→ 2.08	

Log-transform for homoscedasticity

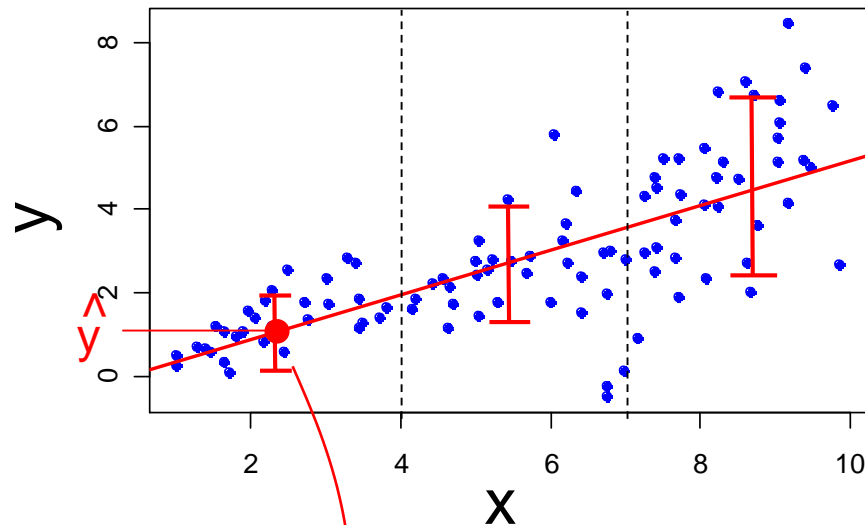


Example: simple linear regression

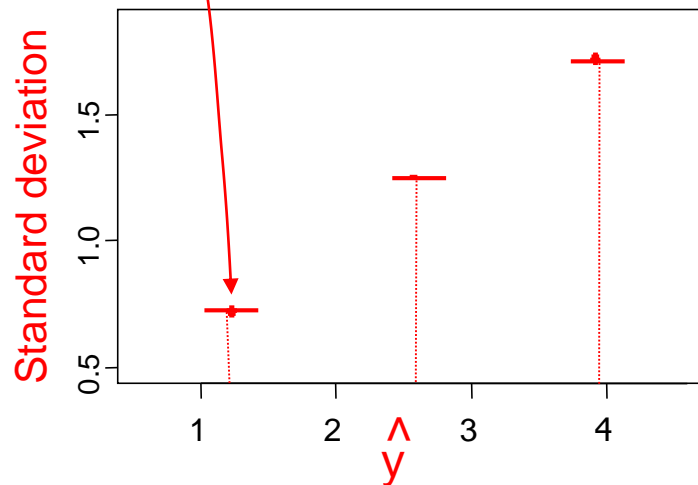
Standard deviation of y increases in proportion to mean of y

Plot standard deviation of residuals against predicted y

→ Linear increase



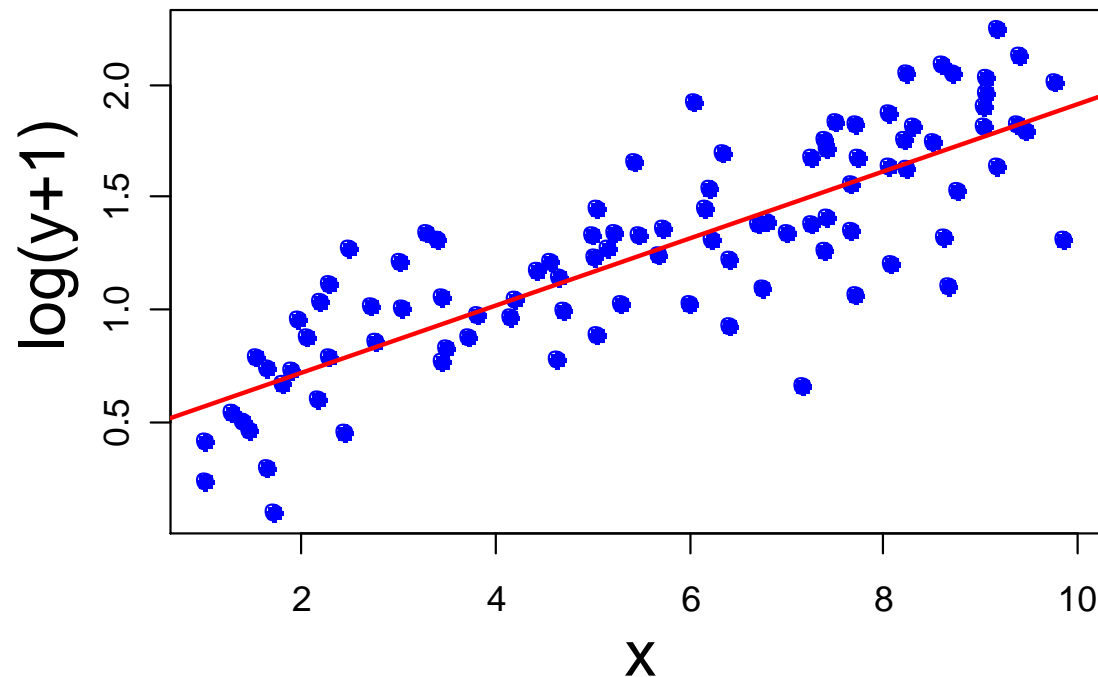
Heteroscedasticity violates regression assumption! (affects CIs & p-values)



Visualize distribution of responses in vertical 'slices' to determine relationship between SD of response and predicted values

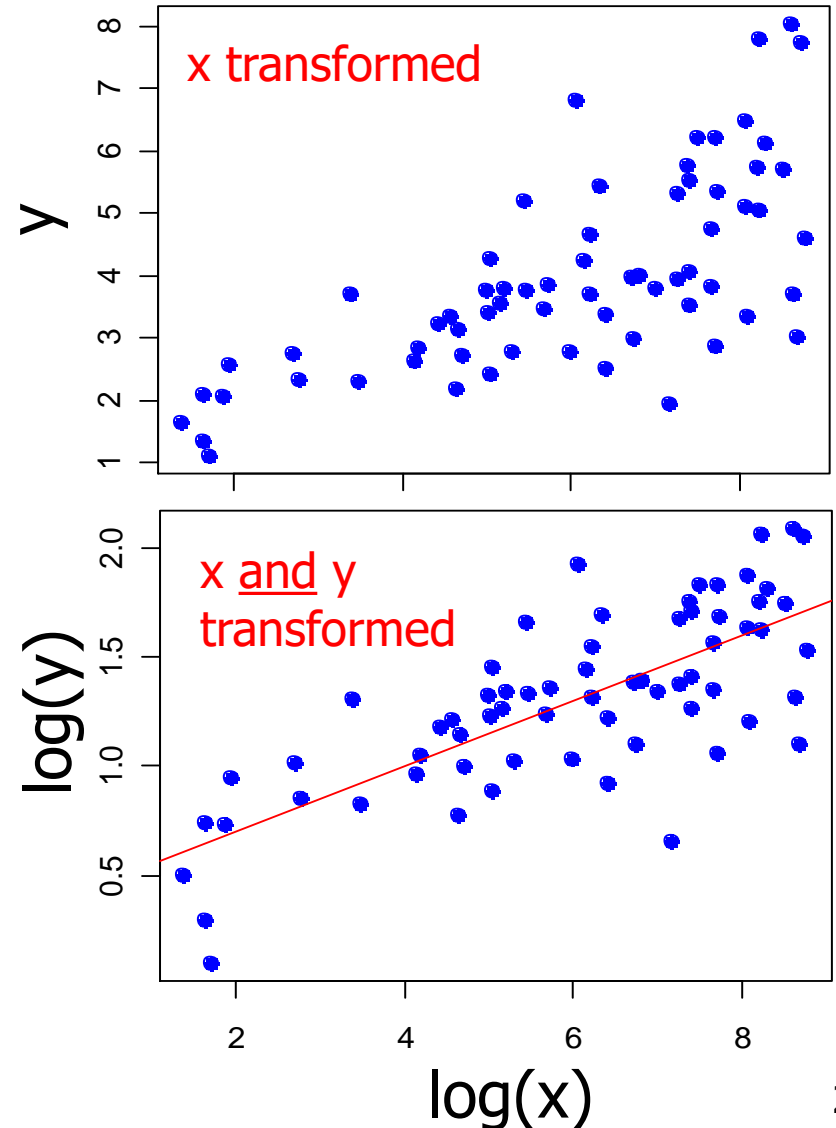
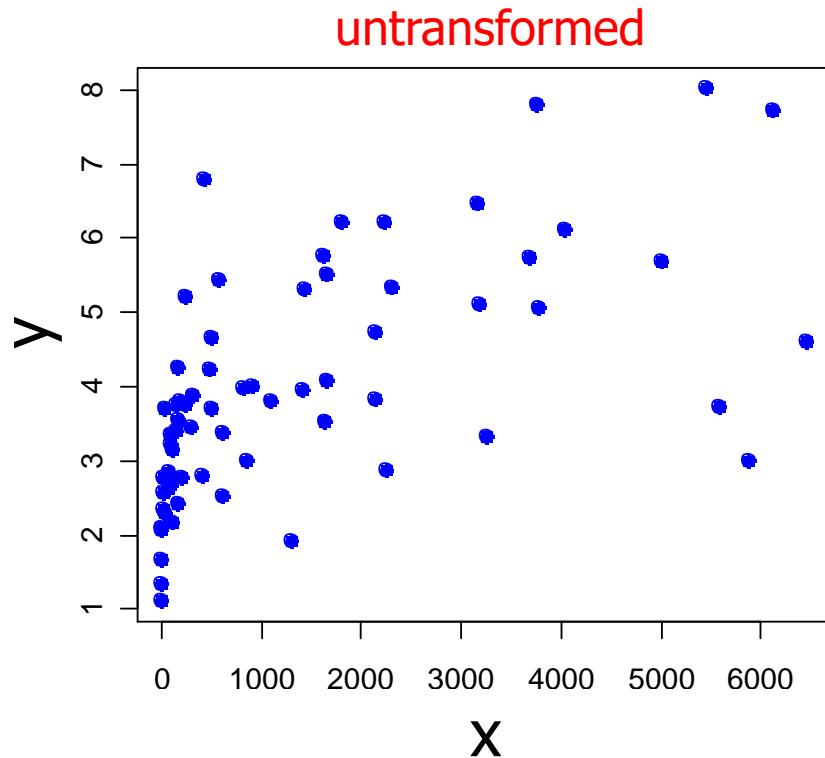
Log-transform for homoscedasticity

Example: simple linear regression with log-transformed y-variable



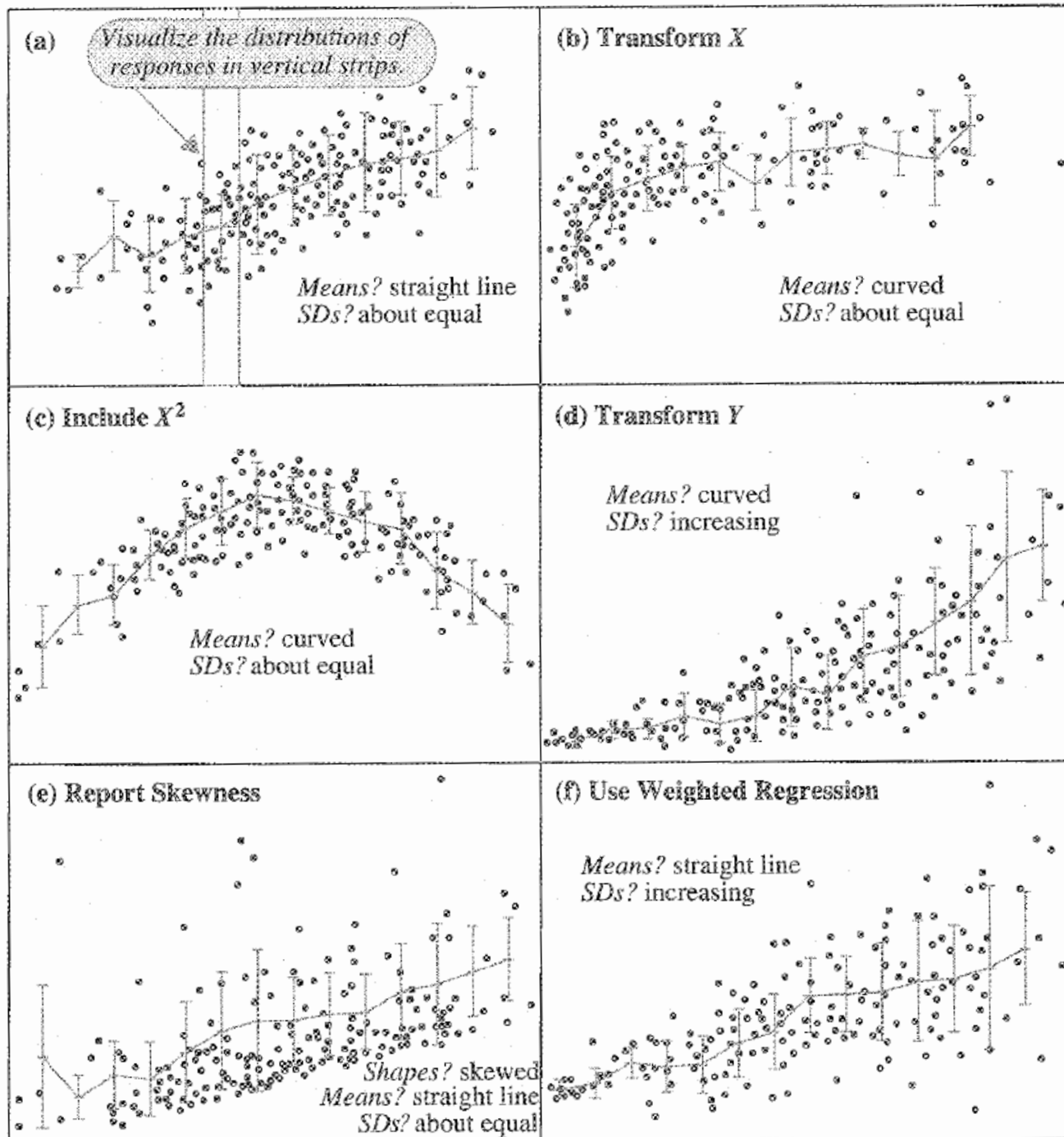
Transformed variable has constant standard deviation (Homoscedasticity)

Log-transform both y and x



→ In some cases, this can result in both linearity & homoscedasticity (and simplify analyses)

If possible,
aim for equal
variances
and a 'linear'
relationship



Other transformations

- Square-root transformation

- When variances are proportional to the mean (e.g. in biology: samples from a Poisson distribution)

$$y' = \sqrt{y + 3/8}$$

- Arc-sine

- To normalize proportions, as long as they are not too close to zero (~ 0.2 - 0.8)

$$p' = \arcsin \sqrt{p}$$

Other transformations

MSL

FISH

604

- Power transformation
 - When standard deviations decrease with the mean and/or if distribution is left-skewed:

$$y' = y^2$$

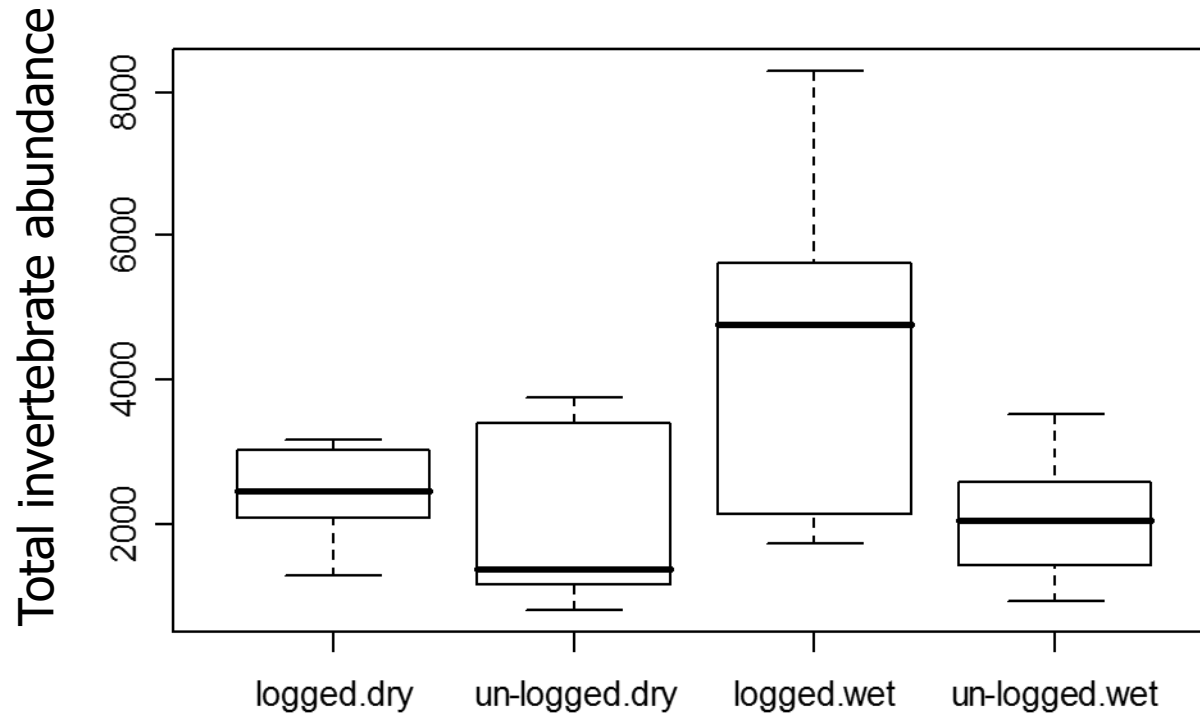
(rare!)

- Box-cox transformation
 - General approach to determine optimal transformation for normalizing data

$$y' = \begin{cases} y^\lambda - 1/\lambda & \lambda \neq 0 \\ \log(y) & \lambda = 0 \end{cases}$$

Box-Cox transformation

- Which transformation, if any?



R-code: `boxplot(totalN ~ treatment * ecoregion, data=logging)`

Box-Cox transformation

MSL

FISH

604

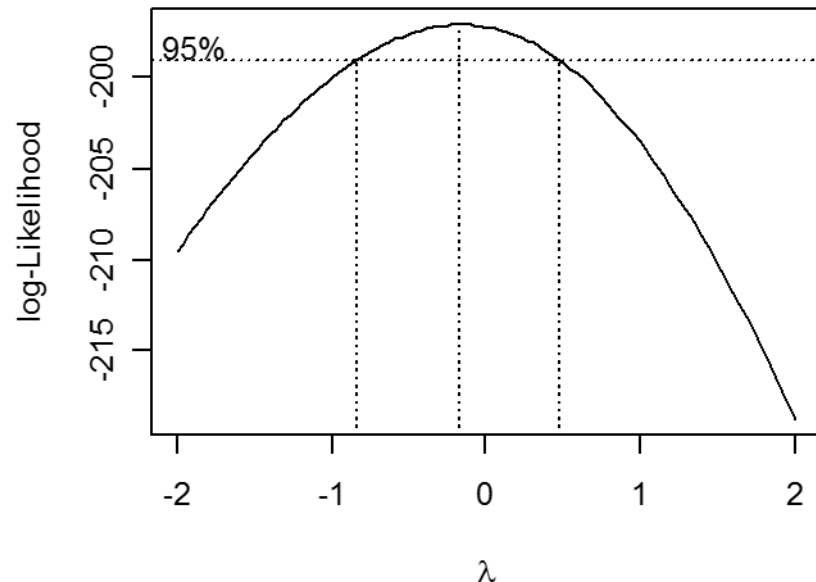
- Fit desired model:

```
fit <- aov(totalN ~ treatment * ecoregion,  
          data=logging)
```

- Determine "best" transformation:

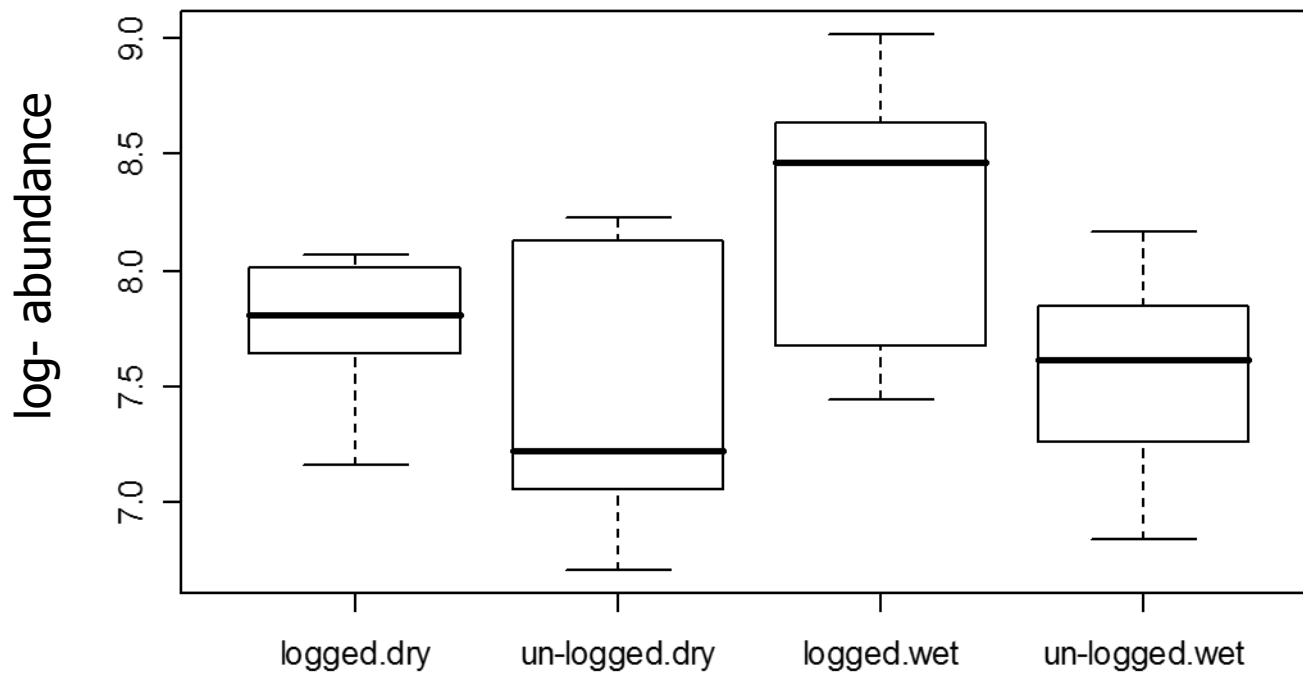
```
boxcox(fit)
```

- log-transform!?



Box-Cox transformation

- Approximately normal distribution and equal variances:



R-code: `boxplot(log(totalN) ~ treatment * ecoregion, data=logging)`



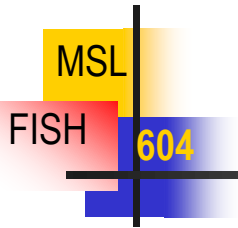
Correlations

- Correlation matrix
- Pearson's product moment correlation
- Robust correlations (rank based)
 - Spearman's rho
 - Kendall's tau
- Testing significance of correlations

Some useful R functions

```
cor(); cor.test(x,y,method="pearson")  
cor.test(x,y,method="spearman")  
rcorr() # library(Hmisc)]  
corrplot() # library(corrplot)  
ggcorrplot() # library(ggcorrplot)
```

Correlation matrix



p < 0.05
p < 0.01

	airT.Kodiak	airT.Sitka	GAK1.SST	GAK1.SSS	GAK1.BS	FW.spring	FW.annual	DW.win	DW.sum
airT.Kodiak	1	0.68	0.69	-0.58	0.08	0.50	0.27	0.04	-0.17
airT.Sitka	0.68	1	0.65	-0.83	0.19	0.69	0.60	-0.30	-0.10
GAK1.SST	0.69	0.65	1	-0.55	-0.20	0.70	0.36	0.01	-0.20
GAK1.SSS	-0.58	-0.83	-0.55	1	-0.35	-0.68	-0.38	0.44	0.16
GAK1.BS	0.08	0.19	-0.20	-0.35	1	-0.08	0.01	-0.39	0.31
FW.spring	0.50	0.69	0.70	-0.68	-0.08	1	0.40	-0.04	-0.13
FW.annual	0.27	0.60	0.36	-0.38	0.01	0.40	1	-0.26	-0.30
DW.win	0.04	-0.30	0.01	0.44	-0.39	-0.04	-0.26	1	0.04
DW.sum	-0.17	-0.10	-0.20	0.16	0.31	-0.13	-0.30	0.04	1

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 vertical line is to the right of the yellow square, and a black horizontal line is below the red square.

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.