

Einführung in die statistische Auswertung von Datensätzen

Lehrgang in wissenschaftlicher Ornithologie
für Studierende und fortgeschrittene Amateure

Lorenz Gyax

Statistician & proximate Biologist

brunner und hess software ag, Zürich

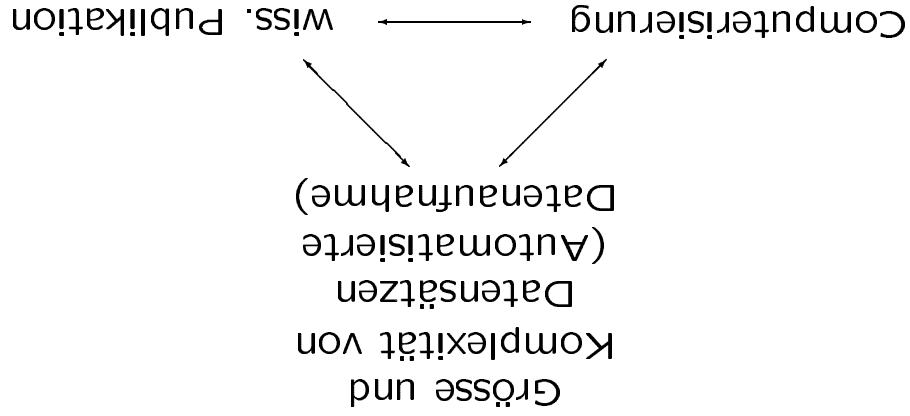
1. Juli 2000

1

Warum Statistik ?

- *saubere* Daten (Ausreisser finden)
- *Reduktion* von Daten
- eingängige und einfach zu interpretierende *graphische* Darstellungen
- quantitative *Beschreibung* von Daten
- quantitative Masse über *Unterschiede* und *Zusammenhänge*

2



3

Inhalt

Vielles nur angetönt, Statistik eigenes Studium

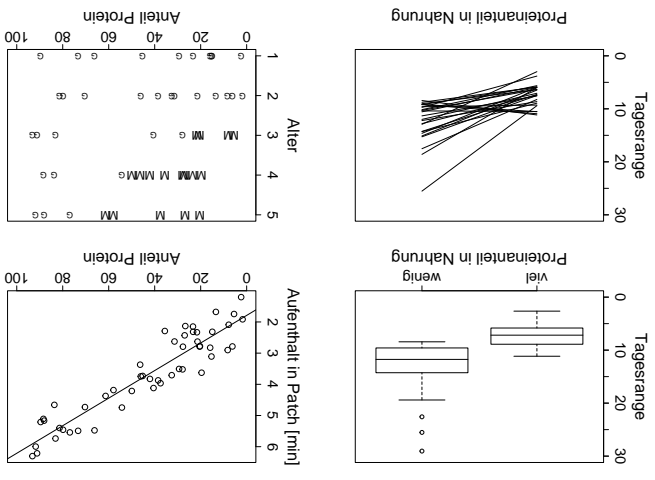
- Graphisches
- Wann? Überhaupt / in einem Projekt
- Wie? Grundsätze statistischer Auswertungen
- Vorwissen
- Wahl der Methode, Überblick

4

Graphisches

- Darstellen von (meist) hochdimensionalen Daten in 2D
- Info auf einen Blick erkennbar, Lesbarkeit
- ideal: Replikate noch immer einzeln einsehbar
- Variabilität und Verteilung sichtbar
- Fantasie für den Spezialfall (Zeit !)
- Statistik "überflüssig"
- Flexibilität des Programmes

5



6

Idee von Statistik

Wie wahrscheinlich sind die Daten unter Annahme des Zufalls (p-Wert von H_0).

Widerspruchstest: H_0 kann abgelehnt, aber weder H_0 noch H_1 bewiesen werden.

Vorwissens → Problemwahrnehmung

„wir sehen, was wir glauben“

trotzdem: Wissen einsetzen → exakte Fragen

7

schliessende ↔ explorative Statistik

eingeschränkte ↔ offene Auswertung

Testen von Hypothesen ↔ „Data-mining“

gute Kontrolle ↔ viele Störeinflüsse

einfache Statistik ↔ multivariate Verfahren

Experimente ↔ Feldarbeit

Schätzungen

Voraussagen

8

Wann, wieviel

- Immer \leftrightarrow Studie entsprechend designen
- d. h. vor Studienbeginn Überlegungen machen
- sich beraten lassen (Kontrolle, Baseline)
- multivariates Auswerten
- Stärke des Einfluss wenn andere Variablen konstant gehalten
- kein Verdecken von Einflüssen
- multiples Testen (α)

6

Replikate

- zufällig ausgewählte Beobachtungseinheiten
- \rightarrow Stichprobengröße ($N = 5 \cdot \#V$?)
- course of dimensionality, Heterogenitäten
- statistische Unabhängigkeit (Individuen, Gebiete)
- 1 Datensatz pro Replik (1 Zeile in Datenfile)
- Investition/Kosten
- $\#$ Replikate \leftrightarrow Details für jedes Replikat

10

Skalen

- nominal/kategorial, binär
- geordnet (ordinal), binär
- Anzahlen: > 0 , Wurzeltransformation
- Beträge: > 0 , kontinuierlich, Log-transf.
- Intervall ($-\infty$ bis ∞)

12

Fehlerarten

- Fehler 1. Art (α , p-Wert):
P [H_0 fälschlich zu verwerfen]
- Multiples Testen
- Fehler 2. Art (β):
P [H_0 fälschlich beizubehalten]
- Macht ($1 - \beta$):
P [wahre Alternative anzunehmen]
- Power-Analysen

11

Vorgehen

- Daten- und Variablenreduktion
- Wahl der Auswertungsmethode, was passt zu Daten?
- Wahl des spezifischen Modelles
- Einfachheit versus Detailliertheit
- Statistisches Testen
- Prüfen der Voraussetzungen/Annahmen (meist: Verteilung der Residuen)

13

Daten- und Variablenreduktion

- oft riesige Datenmengen (autom. Verfahren) → *Zufallsauswahl* (Muster bereits sichtbar)
- Problem: *Kollinearität*
- Korrelierte Variablen (z. B. Wetter)
- nicht unabhängig variierend
- Lösung: z. B. *PCA*
- Interpretierbarkeit der Faktoren?
- Interpretation (Konstrukt), explorativ!

14

Wahl der Methoden

- Komplexität der Studie (KS)
- Skala der Zielvariable (ZV)
- Skala der erklärenden Variablen (EV)
- Verteilung der Fehler (VF)
- Schätzung der Fehler = Teilresultat
- Modelle rechnen ↔ Modellwahl

15

Rangtests

KS: einfach, ZV: ordinal/Intervall, EV: egal, VF: symmetrisch

RANGTESTS (nicht-parametrisch, resampling)

KS: (beliebig) komplex, ZV: > ordinal, EV: egal, VF: normalverteilt

Lineare Modelle

• ANOVA:

$$Y = \mu + E_{X_1} + E_{X_2} + \dots + \epsilon, \epsilon \sim N(0, \sigma^2)$$

• Regression:

$$Y = \alpha + \beta_1 X_1 + \beta_2 X_2 + \dots + \epsilon_i, \epsilon \sim N(0, \sigma^2)$$

• linear? X_2 , $\log(X)$, \sqrt{X} , $\sin(X)$ möglich

• Residuenanalyse; Tests: Levene,

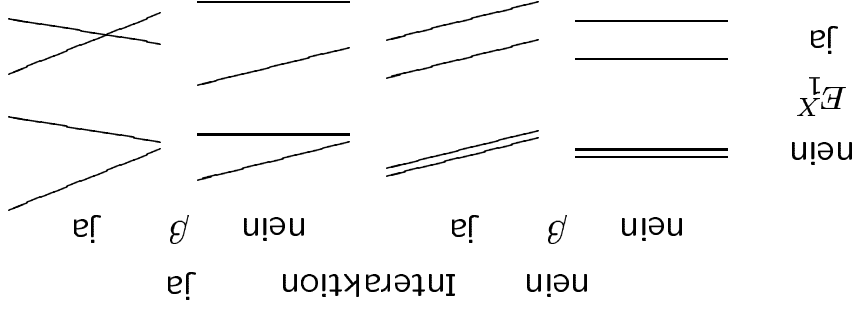
Graphiken: qq, TA, LL, RVE, RC, RI plots

17

Interaktionen

• multiplikative Effekte (alle Skalen!):

$$Y = \alpha + E_{X_1} + \beta X_2 + \beta^I E_{X_1} X_2$$

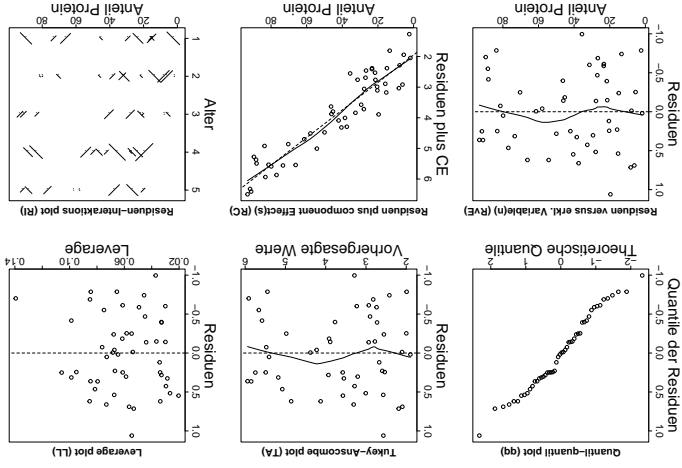


18

Wiederholte Messung ("geparte Daten")

- mehrere Messungen am selben Replikat
- Verlaufskurven
- Repeated Measure; Annahmen
- einige Kennzahlen schätzen pro Replikat
- gute Charakteristika des Verlaufes
- Maximum, Zeit bis . . . , Krümmung
- Auswertung pro Kennzahl

19



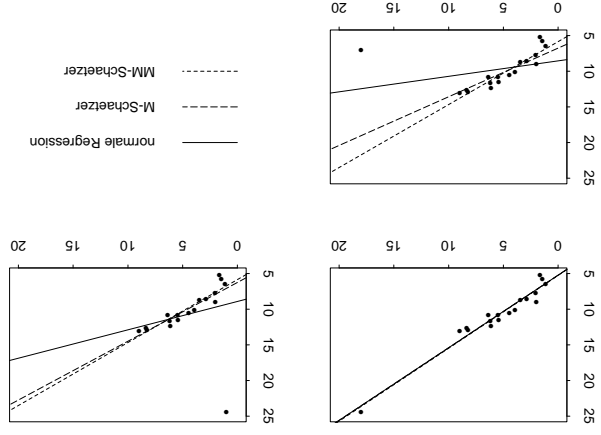
20

KS: komplex, ZV: > ordinal, EV: egal, VF: fast normalverteilt

Robuste Methoden

- Normalverteilung mit langen Schwänzen
- . . . oder kleiner Anzahl Ausreisser
- hat man eigentlich *immer*
- spezielle Methoden (eingesch. Implement.)
- z. B. Mittelwert versus Median
- z. B. Standardabweichung versus MAD
- $\text{mad} = C \cdot \text{med}(X_i - \text{med}(X))$

21



22

Generalisierte lineare Modelle

- kumulative Logits
KS: komplex, ZV: ordinal, EV: egal, VF: multinominal
- Diskriminanzanalyse
KS: komplex, ZV: nominal/binär, EV: egal, VF: multinominal
- Logistische Regression (siehe später)
KS: komplex, ZV: Andere, EV: egal, VF: Andere
- log-lineare Modelle
- Poissonregression

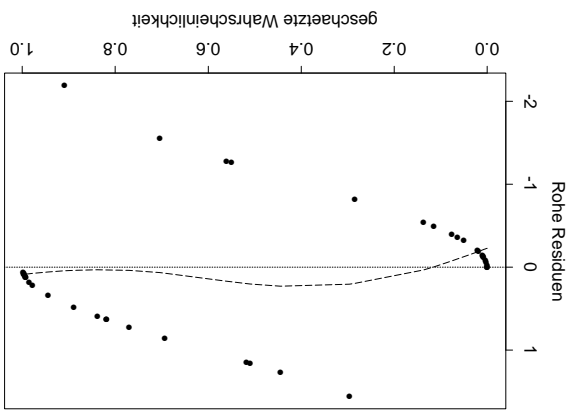
23

Logistische Regression

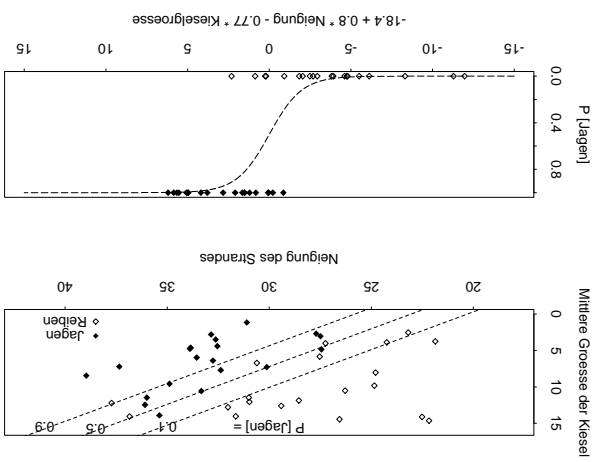
- binäre Zielvariable
- Modell, das passt
- Interaktionen möglich
- Schätzwerte $\hat{=}$ Wahrscheinlichkeiten

KS: komplex, ZV: binär, EV: egal, VF: binomial

24



26



25

Arbeiten mit %en

- Prozente nahe an 0, 1 nicht normalverteilt
- ev. $\arcsin \sqrt{Y}$
- "Sum 1 constraint" ↔ unabhängige Analyse
- ev. eine oder mehrere Anteile weglassen
- nur alle Ergebnisse zusammen interpretieren
- Compositional Analysis

27

Spezielles / Modernes

- nichtlineare Regression
- nichtparametrische Regression (z. B. lowess)
- resampling Verfahren: Jackknife / Bootstrap
- Fehlerverteilung aus Daten geschätzt
- so gut wie ursprüngliche Daten
- Zeitreihen (wenig Replikate, lange Reihen)
- neuronale Netze

28

Statistische Software

- an Anforderungen angepasst (Forscher und Problem)

- was ist vorhanden

- SPSS (www.spss.com)

- Glim, Genstat, SAS

- S-plus (www.mathsoft.com)

- R (www.ci.tuwien.ac.at/R/)

31

Es gibt immer verschiedene Wege ein statistisches Problem anzugehen, verschiedene Methoden sollten bei starker Struktur zur gleichen Interpretation führen!

Statistik ist ein Hilfsmittel zur Entdeckung und quantitativen Beschreibung von Mustern und Zusammenhängen in hochdimensionalen Daten.

32

“Lehre jemanden den t-Test und er wird für
einen Tag glücklich sein; lehre jemanden die
Regression und er wird für eine Woche lang
glücklich sein; lehre jemanden Statistik und
er wird sein ganzes Leben lang Probleme
haben.”
(unbekannter Statistiker)