

# Mensch Maschine Interaktion

Übung 09  
Evaluation



Nicht vergessen:

Abgabe der ersten Online-  
Hausarbeit ist am 26.6.2020, 9:00

# Evaluation



# Klassifizierung von Evaluation

- Formativ vs. Summativ
- Qualitativ vs. Quantitativ
- Analytisch vs. Empirisch



# Evaluationsmethoden

	<b>Formativ vs. summativ</b>	<b>Qualitativ vs. quantitativ</b>	<b>Analytisch vs. empirisch</b>
Heuristische Evaluation			
GOMS KLM			
Beobachtungsstudie			
Kontrollierte Experimente			
Feldstudien			
Tagebuch-Studien			

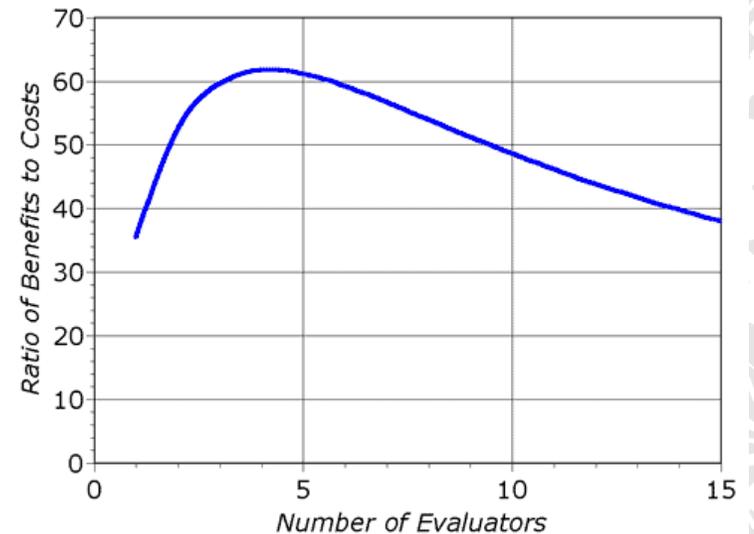
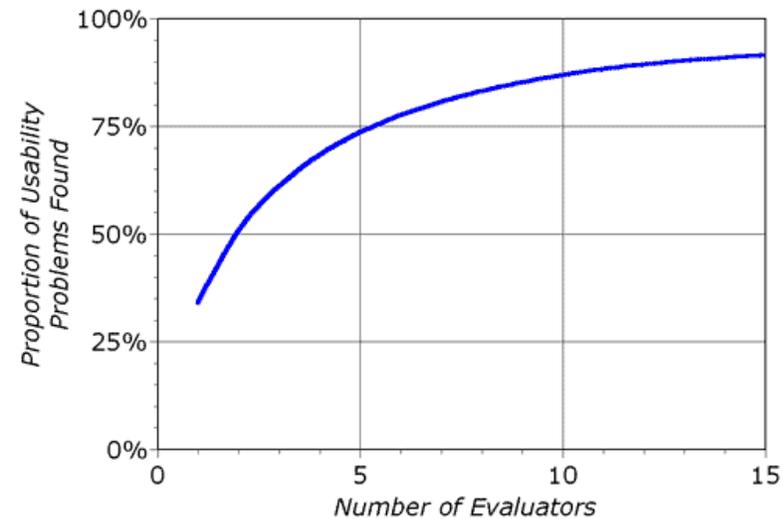
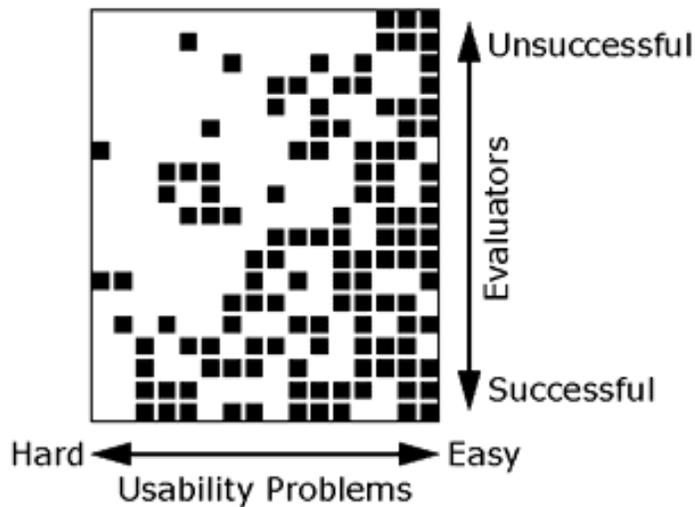
# Heuristische Evaluation

- **Experten** bewerten ein Interface mit einem Fragebogen entlang einer Aufgabe
- Fragebogen orientiert sich oft an den 10 Heuristiken von Nielsen

	Erfüllt	Nicht erfüllt	Kommentare
Das System hat eine Undo Funktion	<input type="radio"/>	<input type="radio"/>	

# Heuristische Evaluation

- Wie viele Bewertungen/Bewerter? Idealerweise mehrere!



# Heuristiken nach Jacob Nielsen

1. Visibility of System Status
2. Match between system and the real world
3. User control and freedom
4. Consistency and standards
5. Error prevention
6. Recognition rather than recall
7. Flexibility and efficiency of use
8. Aesthetic and minimalist design
9. Help users recognize, diagnose, and recover from errors
10. Help and documentation

<https://www.nngroup.com/articles/ten-usability-heuristics/>



Visibility of system status



Match between system and the  
real world



# User control and freedom



# Consistency and standards



# Error prevention



Recognition rather than recall



Flexibility and efficiency of use



# Aesthetic and minimalist design



Help users recognize, diagnose,  
and recover from errors



Help and documentation



# GOMS

Nach S. Card et al.

- **Goals:** Was will ich am Ende erreichen?
  - **Operators:** Welche Handlungen kann ich tätigen?
  - **Methods:** Welche Prozesse / Folgen von Operatoren führen zum Ziel?
  - **Selection Rules:** Welche Methode/n wähle ich zum Ziel?
- 
- "Top-down"-Methode

# GOMS-KLM

Nach S. Card et al.

- "Bottom-up"-Methode
- K – Keystroke  $\approx 0.28s$
- P – Pointing  $\approx 1.1s$
- H – Homing  $\approx 0.4s$
- M – Mental preparation  $\approx 1.35s$
- R – Response of the system



# GOMS KLM

- Berechnen Sie eine Abschätzung für folgendes Beispiel:

Sie haben aktuell Word geöffnet und ihre Übungsblatt-Abgabe fertig getippt. Nun wollen Sie diese als PDF Exportieren und auf Uni2Work hochladen. Sie können davon ausgehen, dass der Browser im Hintergrund geöffnet ist und die Uni2Work upload Seite bereits geöffnet ist.

# Beobachtungsstudie

- Quasi-Experiment
- **Kein Eingriff** in das Geschehen
- Keine Veränderung von Variablen
- Oft nachträgliches Einteilen der Teilnehmer in Gruppen
- Typische Auswertung: **Korrelationen**



# Kontrollierte Experimente

- Unter experimentellen, kontrollierten Bedingungen
- Variieren der **unabhängigen Variable**
- Messung der **abhängigen Variable**
  
- Jeder Teilnehmer durchläuft **jede Experimentalbedingung:**  
within-subject
- Jede Teilnehmer-Gruppe durchläuft **genau eine Bedingung:**  
between-groups

# Kontrollierte Experimente

- Experimente dienen dazu **Hypothesen** zu verifizieren
- Hypothese:  
*"Übungsteilnehmer erzielen bessere Noten in der Klausur"*
- Echt **verifizieren** ist schwer, **falsifizieren** ist oft leichter
- Ab wann sind die Noten denn "besser"?
- Deswegen **Null-Hypothese**:  
*"Übungsteilnehmer und Nichtteilnehmer erzielen im Mittel die gleichen Noten in der Klausur"*

# Studien-Design

- Teilnehmer in Studien sollten aus einer möglichst diversen Gruppe kommen
- Die Studie sollte nicht zu lang sein → Ermüdungseffekte
- Die Reihenfolge des Studienablaufs sollte die Ergebnisse nicht beeinflussen → Reihenfolgeeffekte, Lerneffekte
  - Randomisierung ist gut, aber keine Sicherheit
  - Deswegen am besten alle Permutationen testen
  - Oft nicht machbar da zu viele Möglichkeiten ( $n!$ )

# Latin-Square Permutation

Bed. 1	Bed. 2	Bed. 3	Bed. 4
Bed. 4	Bed. 1	Bed. 2	Bed. 3
Bed. 3	Bed. 4	Bed. 1	Bed. 2
Bed. 2	Bed. 3	Bed. 4	Bed. 1

# Latin-Square Permutation

Bed. 1	Bed. 3	Bed. 4	Bed. 2
Bed. 3	Bed. 1	Bed. 2	Bed. 4
Bed. 4	Bed. 2	Bed. 1	Bed. 3
Bed. 2	Bed. 4	Bed. 3	Bed. 1

# Feldstudien

- Laborstudien spiegeln oft nicht die wirklichen Situationen wieder
- Deswegen Feldstudien, beobachtung in echter Umgebung
- Geringere **interne Validität**
- Dafür größere **externe Validität**



# Tagebuch-Studien

- **Langzeit Studien:** Studien die Effekte über lange Zeiträumen untersuchen
- Dauerhafter Aufenthalt im Labor ist nicht machbar
- Die "Messung" kann dabei eigenverantwortlich als **Tagebuch** durch den Probanden erfasst werden
- Auch möglich über bspw. das Smartphone via **Experience Sampling** (s. auch das **PhoneStudy** Projekt der LMU)

# Evaluationsmethoden

	<b>Formativ vs. summativ</b>	<b>Qualitativ vs. quantitativ</b>	<b>Analytisch vs. empirisch</b>
Heuristische Evaluation			
GOMS KLM			
Beobachtungsstudie			
Kontrollierte Experimente			
Feldstudien			
Tagebuch-Studien			

# Auswertung von Studienergebnissen



# Daten

- Studien-Daten können in verschiedenen Formen kommen
- **Nominal:** bspw. Geschlecht, Herkunft, welches System genutzt wurde
- **Ordinal:** Reihenfolgen, ohne Aussage über Abstände, bspw. Likert-Skalen
- **Interval:** Reihenfolgen mit definierten Abständen, bspw. Temperatur in °C
- **Ratio:** Wie Interval aber mit absolutem Nullpunkt, bspw. Größe

# Auswertung

- Um belastbare Aussagen über Ihre Ergebnisse zu treffen, sollten Sie diese statistisch überprüfen
- Hierfür gibt es **statistische Hypothesentests**
- Diese geben u.a. eine Abschätzung, wie wahrscheinlich es ist, dass Ihre Daten, unter Annahme der Null-Hypothese, so auftreten würden, wie Sie sie beobachtet haben

$p = 0.05$  bedeutet in etwa *“Mit 5% Wahrscheinlichkeit würden wir die Daten so bekommen, wenn die Null-Hypothese gelten würde.”*

# Auswertung

- Bei ausreichend kleinem p-Wert können Sie die **Null-Hypothese ablehnen**. D.h. bei kleinen p ist ihre Hypothese wahrscheinlich richtig.
- Typische “klein-genug”-Werte:
  - $p < 0.05 \rightarrow$  signifikant \*
  - $p < 0.01 \rightarrow$  hoch signifikant \*\*
  - $p < 0.001 \rightarrow$  höchst signifikant \*\*\*

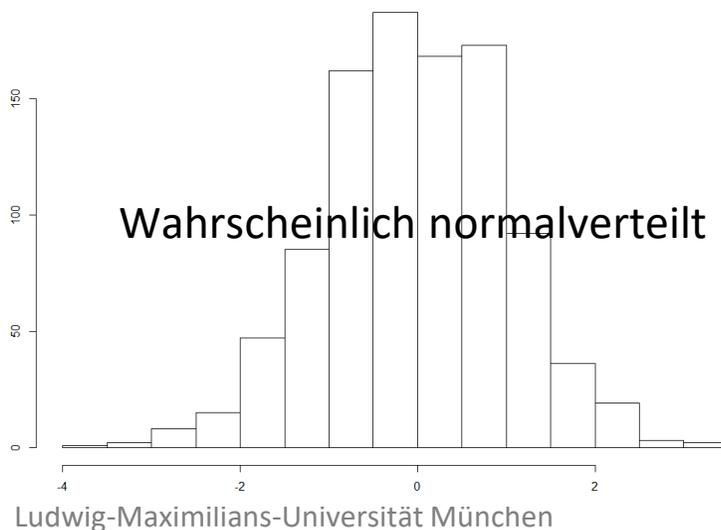
# Hypothesen-Tests

- Es gibt für verschiedene Situationen verschiedene Tests
- Nicht jeder Test ist für jede Situation geeignet oder hat Aussagekraft
- Wichtige Kriterien für die Test-Auswahl:
  - Form der Daten
  - Anzahl der Samples
  - Ziel der Analyse

<https://www.youtube.com/watch?v=rulIUAN0U3w>

# Normalverteilung

- Ein wichtiges Kriterium für die Test-Wahl: Sind die Daten normalverteilt?
- Auch hierfür gibt es statistische Tests, bspw. den Shapiro-Wilk-Test
- Hier ist die *Null-Hypothese*, dass die Daten normalverteilt sind



# t-Test als Beispiel für Hypothesentests

- Für normalverteilte Daten
- Mehrere Varianten:
  - Für 1 sample: Weicht der Mittelwert von einer Annahme ab
  - Für 2 samples unpaired/independent (bspw. zwei Gruppen, between groups): Sind die Mittelwerte signifikant unterschiedlich
  - Für 2 samples paired (bspw. Wenn Probanden zwei Systeme testen, within subject) : Sind die Mittelwerte signifikant unterschiedlich

Bei 2 samples sollten die Varianzen idealerweise ähnlich sein

<https://www.youtube.com/watch?v=rullUAN0U3w>

# Und wie macht man den Test jetzt?

- Hypothesentests sind in div. Software vor-implementiert, bspw.:
  - Excel (eingeschränkt, Microsoft, kostenpflichtig)
  - SPSS (IBM, kostenpflichtig)
  - R (open source, kostenlos)



# R-Studio

The screenshot displays the RStudio environment with the following components:

- Source Editor:** Contains R code for data manipulation and analysis. The code includes reading CSV files, combining them, filtering, and plotting. A boxplot is also shown.
- Environment Pane:** Lists objects in the workspace, such as 'data', 'agg', 'aggregated.by.group', etc., with their respective dimensions.
- Console:** Shows the R version (3.6.1) and the results of several commands, including a Shapiro-Wilk normality test for a random number generator.
- Plots Pane:** Displays a histogram titled "Histogram of runif(1000, min = 0, max = 100)". The x-axis is labeled "runif(1000, min = 0, max = 100)" and the y-axis is labeled "Frequency".

# R-Studio

**Skript-Editor um wiederholbare, automatische Auswertungen auszuführen**

```
1 data.group.1 <- read.csv(file.choose(), header = 1, sep = ";")
2 data.group.2 <- read.csv(file.choose(), header = 1, sep = ";")
3
4 data.group.1$collapse.first <- F
5 data.group.2$collapse.first <- F
6
7 # Combine the two data sets
8 data <- rbind(data.group.1, data.group.2)
9
10 # Filter those that completed the questionnaire
11 data <- data[data$completed == "1", ]
12
13 # This will attach the data to the environment
14 attach(data)
15
16 # Quick overview of the data
17 hist(Age.in.years)
18 plot(Gender)
19
20 # Sum
21
22 boxplot()
23 ZUI.SUS.I.think.that.I.would.like.to.use.this.system.frequently.Rating,
24 ZUI.SUS.I.found.the.system.unnecessarily.complex.Rating,
25 ZUI.SUS.I.thought.the.system.was.easy.to.use.Rating,
26 ZUI.SUS.I.think.that.I.would.need.the.support.of.a.technical.person.to.be.able.to.use.this.system.Rating,
27 ZUI.SUS.I.found.the.various.functions.in.this.system.difficult.to.use.Rating,
28 ZUI.SUS.I.thought.there.was.too.much.inconsistency.in.this.system.Rating,
29 ZUI.SUS.I.would.imagine.that.most.people.would.learn.to.use.this.system.very.quickly.Rating,
30 ZUI.SUS.I.found.the.system.very.difficult.to.use.Rating,
31 ZUI.SUS.I.felt.very.confident.using.the.system.Rating,
32 ZUI.SUS.I.needed.to.learn.a.lot.of.things.before.I.could.get.going.with.this.system.Rating
33
34 # calculate the aggregated sus score
35 ZUI.SUS <- (ZUI.SUS.I.think.that.I.would.like.to.use.this.system.frequently.Rating +
36 100 - ZUI.SUS.I.found.the.system.unnecessarily.complex.Rating +
37 ZUI.SUS.I.thought.the.system.was.easy.to.use.Rating +
38 100 - ZUI.SUS.I.think.that.I.would.need.the.support.of.a.technical.person.to.be.able.to.use.this.system.Rating +
39 ZUI.SUS.I.found.the.various.functions.in.this.system.difficult.to.use.Rating +
40 100 - ZUI.SUS.I.thought.there.was.too.much.inconsistency.in.this.system.Rating +
41 ZUI.SUS.I.would.imagine.that.most.people.would.learn.to.use.this.system.very.quickly.Rating +
42 ZUI.SUS.I.found.the.system.very.difficult.to.use.Rating +
43 ZUI.SUS.I.felt.very.confident.using.the.system.Rating +
44 ZUI.SUS.I.needed.to.learn.a.lot.of.things.before.I.could.get.going.with.this.system.Rating)
45
```

**Übersicht über aktuelle Daten**

Object	Class	Attributes
data	data.frame	39 obs. of 3 variables
agg	data.frame	3 obs. of 2 variables
agg.by.group	data.frame	36 obs. of 4 variables
agg.by.user	data.frame	30 obs. of 4 variables
agg.by.user.not.stem	data.frame	30 obs. of 4 variables
agg.by.user.stem	data.frame	36 obs. of 4 variables
data	data.frame	9 obs. of 133 variables
data.group.1	data.frame	6 obs. of 133 variables
data.group.2	data.frame	4 obs. of 133 variables
daten.teil.1	data.frame	6 obs. of 132 variables
daten.teil.2	data.frame	4 obs. of 132 variables
demographics	data.frame	39 obs. of 11 variables
F.it	data.frame	39 obs. of 11 variables
fitted.norm	data.frame	39 obs. of 11 variables
meta.data	data.frame	28 obs. of 2 variables
meta.data.1.2	data.frame	20 obs. of 17 variables
meta.data.1.3	data.frame	18 obs. of 17 variables
meta.data.2.3	data.frame	18 obs. of 17 variables
study.data	data.frame	1092 obs. of 7 variables
study.data.filtered	data.frame	908 obs. of 8 variables
sus.test.data	data.frame	18 obs. of 2 variables
test.data	data.frame	216 obs. of 2 variables
tlx.data	data.frame	39 obs. of 16 variables
tlx.not.stem	data.frame	30 obs. of 16 variables
tlx.stem	data.frame	36 obs. of 16 variables
tmp	data.frame	36 obs. of 16 variables
tmp.ex.1	data.frame	36 obs. of 18 variables
tmp.ex.2	data.frame	36 obs. of 18 variables

**Terminal um Einzel-Befehle einzugeben**

```
R version 3.6.1 (2019-07-05) -- "Action of the Toes"
Copyright (C) 2019 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

> workspace loaded from ~/RStudio/
> plot(rnorm(1000))
> hist(rnorm(1000))
> hist(runif(1000, min=0, max=100))
> shapiro.test(rnorm(1000))

      Shapiro-Wilk normality test

data:  rnorm(1000)
W = 0.9988, p-value = 0.7321

> shapiro.test(runif(1000, min=0, max=100))
Error in shapiro.test(runif(1000, min = 0, max = 100)) :
could not find function "shapiro.test"
> shapiro.test(runif(1000, min=0, max=100))

      Shapiro-Wilk normality test

data:  runif(1000, min = 0, max = 100)
W = 0.95169, p-value = 2.2e-16
```

**Ausgabe von bspw. Grafiken**

# R-Studio

- Kommende Woche wollen wir mit R-Studio arbeiten
- Deswegen bitte bis dahin R-Studio installieren



# Nächste Übung: Evaluation II & IxD

22.06.20 - 26.06.20

