# 

# Statistical Evaluation

Research Methods for Human-Centered Computing



Today's goal:

A (very) brief overview of statistical analysis

Outline:

- An overview of some stats
- An example



# **Statistics** A very brief overview



Covariance and correlation

Linear regression

t-tests, ANOVA, factorial ANOVA

Non-normal data

Multi-level data

Subjective data

Mediation analysis



A model is a way to explain or summarize the data The mean is a model

The quality of the model depends on how well it fits the data

We can measure the deviance between the model and the data





 $error_i = x_i - mean$ 

- $SS = \sum error_i^2$ SS = sum of squared errors
- s<sup>2</sup> = SS/(N-1)
  s<sup>2</sup> = variance
  s = standard deviation
  N-1 = degrees of freedom







Let's say you have 4 data points: 1, 3, 4, 8 Mean: 4

If you know the mean, how many data points are "free"?

Answer: Only three!

Once you know the first three, you will know the fourth one as well, because the mean needs to be 4!

 $(1+3+4+x)/4 = 4 \longrightarrow x$  has to be 8!



**Variance** is the variation of the data around a model (e.g. the mean)

$$s^2 = \sum (x_i - mean_x)^2 / (N-1)$$

It is the sum of the **error in x times the error in x**, divided by the degrees of freedom







**Covariance** measures the relationship between the variations of two variables, x and y

$$cov(x,y) = \sum(x_i - mean_x)(y_i - mean_y)/(N-1)$$

It is the sum of the **error in x times the error in y**, divided by the degrees of freedom



**Covariance** measures the relationship between the variations of two variables, x and y

 $cov(x,y) = \sum(x_i - mean_x)^*$  $(y_i - mean_y)/(N-1)$ 

It is the sum of the **error in x times the error in y**, divided by the degrees of freedom





# Standardization:

We can standardize any deviation by dividing it by the standard deviation of the measure ( $\sqrt{variance}$ )

If we want to standardize the covariance, we divide by **both** the standard deviation of x and the standard deviation of y.

The resulting metric is the **correlation coefficient**:

 $r = cov(x,y)/s_xs_y = \sum(x_i - mean_x)(y_i - mean_y)/(N-1)s_xs_y$ 



Which of these two graphs shows the strongest correlation?





# **Independent variables** (X): things that are manipulated (experiment) or innate (survey)

- Low vs. high diversity
- Number of search results
- Gender
- Age

They are outside the participants' control (in the experiment)



# **Dependent variables** (Y): things that are measured as an outcome of X

- Number of clicks
- Interaction time
- Facial expression
- Satisfaction\*



**Random variables** (also X): variables that are not of interest, but they may influence Y, so we measure them just in case.

**Control variables** (not X): variables that are not of interest, but they may influence Y, so we try to keep them stable



## More of X -> more of Y:

Does user satisfaction increase with the number of search results?

# More of X -> less of Y: Does Facebook usage satisfaction decrease with age?

#### User satisfaction 2 Ο -1 -2 5 15 25 0 10 20 30 Search results



Any type of model: outcome; = model + error;

Linear regression:

The model is a line with an intercept (a) and a slope (b)

$$Y_i$$
 = a + b $X_i$  +  $e_i$ 

#### **User satisfaction**





## How good is the **model**?

We can use deviation for this as well!

Compare against the deviation of the simplest model

In this case: the mean

#### **User satisfaction**





outcome<sub>i</sub> = model + error<sub>i</sub>

Multiple regression:

The model is a line with an intercept (a) and **several** slopes  $(b_1...b_n)$ 

$$Y_i = a + b_1 X_{1i} + b_2 X_{2i} + ... + b_n X_{ni} + e_i$$

This means you can predict satisfaction using usability **and** gender, in each case controlling for the other variable

Note: bs are partial correlations (not the same as r!)



E.g.: satisfaction<sub>i</sub> =  $1.00 + 2.00^*$  usability<sub>i</sub> +  $1.50^*$  gender<sub>i</sub> + e<sub>i</sub>

- For every 1 point increase in usability, satisfaction is expected to increase by 2 points, controlling for gender
- Controlling for usability, the satisfaction for males (1) is expected to be 1.5 points higher than for females (0)



Difference between two systems:

Do these two Uls (A and B) lead to a different level of usability?

Differences between two groups of people:

Do men (A) and women (B) perceive different levels of usability?



#### Usability



Differences between >2 systems / groups:

Are there differences in perceived system effectiveness between these 3 algorithms?

First do an omnibus test, then post-hoc tests or planned contrasts

Family-wise error!

#### Perceived system effectiveness





Two manipulations at the same time:

What is the combined effect of list diversity and list length on perceived recommendation quality?

Test for the interaction effect!

#### Perceived quality 0.6 low diversification 0.5 high diversification 0.4 0.3 0.2 O.1 10 items 5 items 20 items

Willemsen et al.: "Understanding the Role of Latent Feature Diversification on Choice Difficulty and Satisfaction", UMUAI



Standard tests assume that the dependent variable (Y) is an continuous, unbounded, normally distributed interval variable

- Continuous: variable can take on any value, e.g. 4.5 or 3.23 (not just whole numbers)
- Unbounded: range of values is unlimited (or at least does not stop abruptly)
- Interval: differences between values are comparable; is the difference between 1 and 2 the same as the difference between 3 and 4?



Not true for most behaviors!

- Number of clicks
- Time, money
- 1-5 ratings
- Decisions





Linear regression:

 $Y_i = a + b_1 X_{1i} + b_2 X_{2i} + ... + b_k X_{ki} + e_i$ 

What if Y is binary (0 or 1)?

We can try to predict the **probability** of Y=1 - P(Y)

However, this probability is a number between 0 and 1

For linear regression, we want an unbounded linear Y!

Can we find some transformation that allows us to do this? Yes:  $P(Y) = 1/(1+e^{-U})$ 



$$P(Y) = 1 / (1 + e^{-U})$$

Conversely: U = ln(P(Y)/(1-P(Y)))

Interpretation:

P(Y)/(1-P(Y)) is the **odds** of Y

Therefore, U is the log odds, or **logit** of Y





Standard regression requires **uncorrelated errors** 

This is not the case when...

...you have repeated measurements of the same participant (e.g. you measured 5 task performance times per participant, for 60 participants)

...participants are somehow related (e.g. you measured the performance of 5 group members, for 60 groups)



Count variables often look like this

> Examples: # of purchases, # of clicks, time\*, price\*

Not normal, heteroscedastic!

Can we find some transformation that makes this work?

Yes: Y = 
$$e^{\bigcup}$$





How to interpret the b coefficients?

- b is the increase in U for each increase of  $\boldsymbol{X}$
- b is the increase in the **log rate** of Y for each increase in X
- e<sup>b</sup> is the ratio of rate Y for each increase in X
- e<sup>b</sup> is the **rate ratio**

Why the ratio?

 $b = \log(rate_{x+1}) - \log(rate_x) = \log(rate_{x+1} / rate_x)$ therefore,  $e^b = rate_{x+1} / rate_x$ 



Question: "I only act to satisfy immediate concerns, figuring the future will take care of itself."

Answer categories:

- 1=extremely uncharacteristic
- 2=somewhat uncharacteristic
- 3=uncertain
- 4=somewhat characteristic
- 5=extremely characteristic



This is ordinal, not interval!

Is the difference between "extremely uncharacteristic" and "somewhat uncharacteristic" the same as the difference between "uncertain" and "somewhat characteristic"?

Also, likely not very normally distributed!

How can we solve these problems?











The model estimates intercepts for each threshold 1|2, 2|3, 3|4, 4|5

These thresholds are the **log odds** of any person having **at least** this value

How to interpret the b coefficients?

 $e^b$  is the **odds ratio** for a 1pt increase in X

e.g. if the odds ratio is 1.40, then the odds of a higher value increase by 40% if X is 1 higher



Repeated measurements

- e.g. participants make 30 decisions
- (Partially) within-subjects design
  - e.g. participants are randomly assigned to 1 of 3 games, and tested once with sound on and once with sound off

Grouped data

e.g. participants perform tasks in groups of 5

A combination of the above



Consequence: errors are correlated

There will be a user-bias (and maybe an task-bias)

Golden rule: data-points should be **independent** 





Take the average of the repeated measurements

- Reduces the number of observations
- It becomes impossible to make inferences about individual tasks/users/etc.





## In regression:

- define a random intercept for each user
- impose an error
   covariance structure





Behavior is an "observed" variable

- Relatively easy to quantify
- E.g. time, money spent, click count, yes/no decision

Perceptions, attitudes, and intentions (subjective valuations) are "unobserved" variables

They happen in the user's mind

How can we quantify them?



Psychometrics:

Ask multiple questions on a 5- or 7-point scale

# E.g. perceived system effectiveness:

- "Using the system is annoying"
- "The system is useful"
- "Using the system makes me happy"
- "Overall, I am satisfied with the system"
- "I would recommend the system to others"
- "I would quickly abandon using this system"

# Use **factor analysis** to validate the scales





Why would the new system (X) have a higher usability (Y)?



To learn something from a study, we need a **theory** behind the effect

- This makes the work generalizable
- This may suggest future work
- Measure **mediating** variables
  - Measure subjective system aspects
  - Find out how they mediate the effect on user experience

Statistical method: **structural equation modeling** (SEM)



> Does the system influence usability via understandability?





> Does the system influence usability via understandability?





> Does the system influence usability via understandability?





> Does the system influence usability via understandability?





> Does the system influence usability via understandability?

Types of mediation Partial mediation **Full mediation** 

Negative mediation





> Does the system influence usability via understandability?









Less attractive 30% sales Higher choice satisfaction

More attractive 3% sales Lower choice satisfaction



# Satisfaction = benefit - cost

- Benefit of more options: easier to find the right option
- Cost of more options: more comparisons, higher potential regret

Is this also true for **recommendations**?





# **Example** from Bollen et al.: "Choice Overload"

What is the effect of the number of recommendations? What about the composition of the recommendation list?

# Tested with **3 conditions**:

- **–** Top 5:
  - recs: 1 2 3 4 5
- **–** Top 20:
  - recs: 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20
- Lin 20:

- recs: 1 2 3 4 5 99 199 299 399 499 599 699 799 899 999 1099 1199 1299 1399 1499











Bollen et al.: "Understanding Choice Overload in Recommender Systems", RecSys 2010





Bollen et al.: "Understanding Choice Overload in Recommender Systems", RecSys 2010