# PLOTTU – Identifying trending Topics

Based on Internet Search History

1

# Aim of the Project

- User narrow down answers to questions often by using internet research

- Knowing the knowlegde gap can help to take action, e.g., trainings or team building

- Interesting topics can be identifyied by
  - search engine key words
  - keywords of HTML pages visited
  - keywords and phrases extracted from webpages visited

2

## Proof of Concept Data Flow
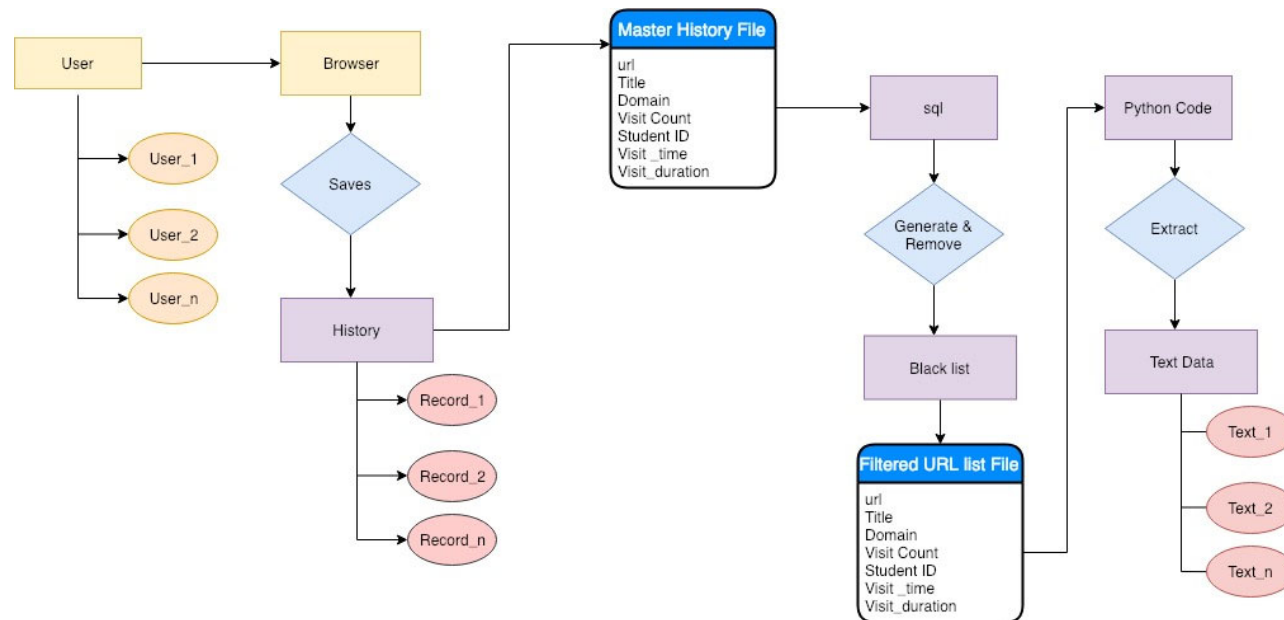
**Data Collection Scenario**
- no relevant downloads available
- students workgroup (n=34) with SPECIFIC data science question „How to assess a linear regression model?"
- Google search with Google Chorme browser
- duration approx. 15 minutes

**Dataset reduction**
- 104 queries extracted
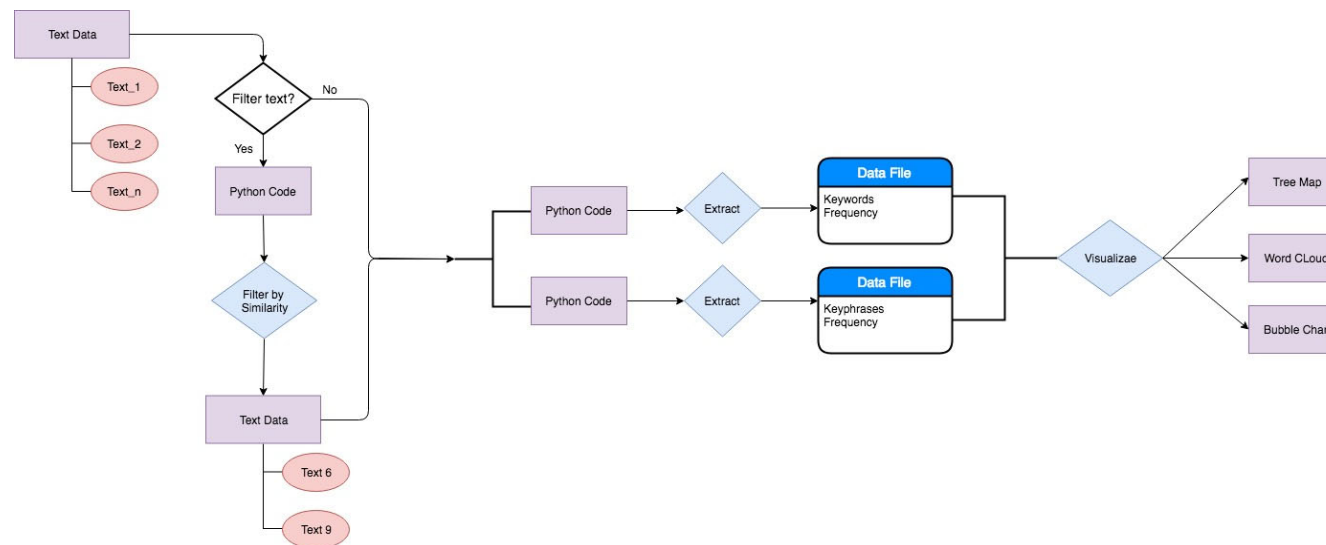- preprocessing and keyword analysis using Python

3

HTW Berlin | Amuj Dixit | Aelee Im | Ting Shen | Cornelia Niklas | Prof. Dr. Tilo Wendler

3

# Proof of Concept Data Flow

- Data Flow Part 1

HTW Berlin | Anuj Dixit | Aelee Im | Ting Shen | Cornelia Niklas | Prof. Dr. Tilo Wendler

# Proof of Concept Data Flow

- Data Flow Part 2

# Steps to Analyse Data

Extracting Google Chrome Data

HTW Berlin | Anuj Dixit | Aelee Im | Ting Shen | Cornelia Niklas | Prof. Dr. Thilo Wendler

# Steps to Analyse Data

Extracting Google Chrome Data

HTW Berlin | Anuj Dixit | Aelee Im | Ting Shen | Cornelia Niklas | Prof. Dr. Tilo Wendler

# Steps to Analyse Data

Extracting Google Chrome Data

| | A | B | C | D | E | F | G |
|---|---|---|---|---|---|---|---|
| | id | url | title | visit_count | StudentID | visit_time | visit_duration |
| | 1 | http://www.f3.htw-berlin.de/ | Fachbereich 3 | 1 | s01 | 45:06.4 | 0 |
| | 2 | https://www.f3.htw-berlin.de/ | Fachbereich 3 | 1 | s01 | 45:06.4 | 0 |
| | 3 | https://www.google.de/search?q=dropbox+login&oq=drop&aqs=chrome.2.69i5 | dropbox login - Google-Suche | 1 | s01 | 46:24.5 | 0 |
| | 4 | https://www.dropbox.com/en_GB/login | Login - Dropbox | 1 | s01 | 46:30.4 | 0 |
| | 5 | https://www.dropbox.com/profile_services/redirect_to_identity_provider?actic | Anmelden ?€? Google Konten | 1 | s01 | 46:33.9 | 0 |
| | 6 | https://accounts.google.com/o/oauth2/auth?access_type=offline&client_id=801 | Anmelden ?€? Google Konten | 1 | s01 | 46:33.9 | 0. |
| | 7 | https://accounts.google.com/signin/oauth?client_id=801668726815.apps.google | Anmelden ?€? Google Konten | 6 | s01 | 46:34.1 | 0. |
| | 7 | https://accounts.google.com/signin/oauth?client_id=801668726815.apps.google | Anmelden ?€? Google Konten | 6 | s01 | 46:34.4 | 0. |
| | 7 | https://accounts.google.com/signin/oauth?client_id=801668726815.apps.google | Anmelden ?€? Google Konten | 6 | s01 | 46:34.4 | 0. |
| | 7 | https://accounts.google.com/signin/oauth?client_id=801668726815.apps.google | Anmelden ?€? Google Konten | 6 | s01 | 46:34.4 | 0. |
| | 7 | https://accounts.google.com/signin/oauth?client_id=801668726815.apps.google | Anmelden ?€? Google Konten | 6 | s01 | 46:34.4 | 0. |
| | 7 | https://accounts.google.com/signin/oauth?client_id=801668726815.apps.google | Anmelden ?€? Google Konten | 6 | s01 | 46:34.6 | 0. |
| | 8 | https://accounts.google.com/signin/oauth/identifier?client_id=801668726815.a | Anmelden ?€? Google Konten | 2 | s01 | 46:34.6 | 0. |
| | 8 | https://accounts.google.com/signin/oauth/identifier?client_id=801668726815.a | Anmelden ?€? Google Konten | 2 | s01 | 47:28.4 | 0. |
| | 9 | https://accounts.google.com/signin/v2/challenge/pwd?client_id=801668726815 | Anmelden ?€? Google Konten | 1 | s01 | 47:28.4 | 0. |
| | 10 | https://accounts.google.com/CheckCookie?hl=de&checkedDomains=youtube&c | Weiterleitung | 1 | s01 | 48:08.3 | 0. |
| | 11 | https://accounts.youtube.com/accounts/SetSID?ssdc=1&sidt=ALWU2csfkg0Lpih | Weiterleitung | 1 | s01 | 48:08.3 | 0. |
| | 12 | https://accounts.google.de/accounts/SetSID?ssdc=1&sidt=ALWU2cs6YpBTfFvRyl | Weiterleitung | 2 | s01 | 48:08.3 | 0. |
| | 12 | https://accounts.google.de/accounts/SetSID?ssdc=1&sidt=ALWU2cs6YpBTfFvRyl | Weiterleitung | 2 | s01 | 48:08.3 | 0. |
| | 13 | https://accounts.google.de/accounts/SetSID | Weiterleitung | 1 | s01 | 48:08.3 | 0. |
| | 14 | https://accounts.google.com/signin/oauth/consent?authuser=0&part=AJi8hAM3WwrP8cpvAs6X8nXsDejSciGYwiNO | | 1 | s01 | 48:10.2 | 0. |
| | 15 | https://accounts.google.com/signin/oauth/consent?authuser=0&part=AJi8hAM3WwrP8cpvAs6X8nXsDejSciGYwiNO | | 1 | s01 | 48:10.2 | 0 |
| | 16 | https://www.dropbox.com/google/authcallback?state=ABzxVboKkpADlbt_i3lIU6SR5wh_tiDwq95u41WY6RP_1CKF-C | | 1 | s01 | 48:10.2 | 0.029705 |
| | 17 | https://www.dropbox.com/ | ?????? - Dropbox | 2 | s01 | 48:10.6 | 0 |
| | 17 | https://www.dropbox.com/ | ?????? - Dropbox | 2 | s01 | 48:10.9 | 0 |

8

HTW Berlin | Anuj Dixit | Aelee Im | Ting Shen | Cornelia Niklas | Prof. Dr. Tilo Wendler

8

# Steps to Analyse Data

Extracting Google Chrome Data

| domain | path | params | query | fragment | scheme |
|---|---|---|---|---|---|
| www.f3.htw-berlin.de | | NaN | NaN | NaN | http |
| www.f3.htw-berlin.de | | NaN | NaN | NaN | https |
| www.google.de | search | NaN | q=dropbox+login&oq=drop&aqs=chrome.2.69i57j0l5... | NaN | https |
| www.dropbox.com | en_GB login | NaN | | NaN | https |
| www.dropbox.com | profile_services redirect_to_identity_provider | NaN | action=login_user&... | | |
| accounts.google.com | o oauth2 auth | NaN | access_type=offline... | | |
| accounts.google.com | signin oauth | NaN | client_id=8016687... | | |
| accounts.google.com | signin oauth | NaN | client_id=8016687... | | |

| Attribute | Index | Value | Value if not present |
|---|---|---|---|
| scheme | 0 | URL scheme specifier | *scheme* parameter |
| netloc | 1 | Network location part | empty string |
| path | 2 | Hierarchical path | empty string |
| params | 3 | Parameters for last path element | empty string |
| query | 4 | Query component | empty string |
| fragment | 5 | Fragment identifier | empty string |
| username | | User name | None |
| password | | Password | None |
| hostname | | Host name (lower case) | None |
| port | | Port number as integer, if present | None |

9

# Extracted Relevant Content

Result of Python scripts:

| | A | B | C | D |
|---|---|---|---|---|
| 1 | urls | key_word | key_phrase | |
| 2 | blog.minitab.com/blog/adventures-in-statistics-2/how-high-should-r- | squared,regression,question,prediction,minitab,model,statistic,analy: | - t;wrong question;continuous improvement;data analysis;main goal | |
| 3 | blog.minitab.com/blog/adventures-in-statistics/how-high-should-r-sq | squared,regression,question,prediction,minitab,model,statistic,analy: | - t;wrong question;continuous improvement;data analysis;main goal | |
| 4 | courses.lumenlearning.com/introstats1/chapter/testing-the-significar | value,line,correlation,coefficient,significant,population,linear,sample, | y values;standard deviation;best-fit line;correlation coefficient;critical value | |
| 5 | de.wikipedia.org/wiki/Anzahl_der_Freiheitsgrade_(Statistik) | displaystyle,mathbf,freiheitsgrade,top,anzahl,boldsymbol,varepsilon,: | anzahl der freiheitsgrade;sum _;sim chi ^;beta _;displaystyle beta _ | |
| 6 | de.wikipedia.org/wiki/Chi-Quadrat-Test | displaystyle,chi,bearbeiten,cdot,test,quadrat,quelltext,frac,verteilung | taschenbuch der statistik;sum _;der ablehnungsbereich f??r;die pr??fgr????e;displaystyle alpha | |
| 7 | de.wikipedia.org/wiki/F-Test | displaystyle,test,sigma,mathrm,bearbeiten,wert,stichprobe,frac,geq,\ | sigma _;worden w??re;displaystyle =p;displaystyle alpha _;displaystyle f_ | |

**Extracted Data:**

- frequency of keywords -> word cloud
- key word -> related key phrase

| | A | B |
|---|---|---|
| | word_name | fre_num |
| 2 | regression | 88 |
| 3 | test | 79 |
| 4 | displaystyle | 79 |
| 5 | statistic | 76 |
| 6 | data | 68 |

```python
import nltk
from nltk.corpus import stopwords
a = df['titel'].str.cat(sep=' ')
words = nltk.tokenize.word_tokenize(a)
delete = [',', '.', ':', ';', '?', '(', ')', '[', ']', '&', '!', '*', '@', '#', '$', '%','-','+','=','1m1','1m5','2m2',
words = [word for word in words if word not in delete]
stops = set(stopwords.words("english"))
words = [word for word in words if word not in stops]
word_dist = nltk.FreqDist(words)
rslt = pd.DataFrame(word_dist.most_common(12),
                    columns=['Word', 'Frequency'])
print(rslt)
```

```
      Word  Frequency
0     spss        149
1     maps        149
2     mail        120
3     wiki        103
4      chi         95
5   square         89
6  regression      85
7  calendar        84
8     test         80
9   Berlin         74
10 altenberg       73
```

10

HTW Berlin | Anuj Dixit | Aelee Im | Ting Shen | Cornelia Niklas | Prof. Dr. Thomas Wendler

10

# Extracted Relevant Content

Result of Python scripts:

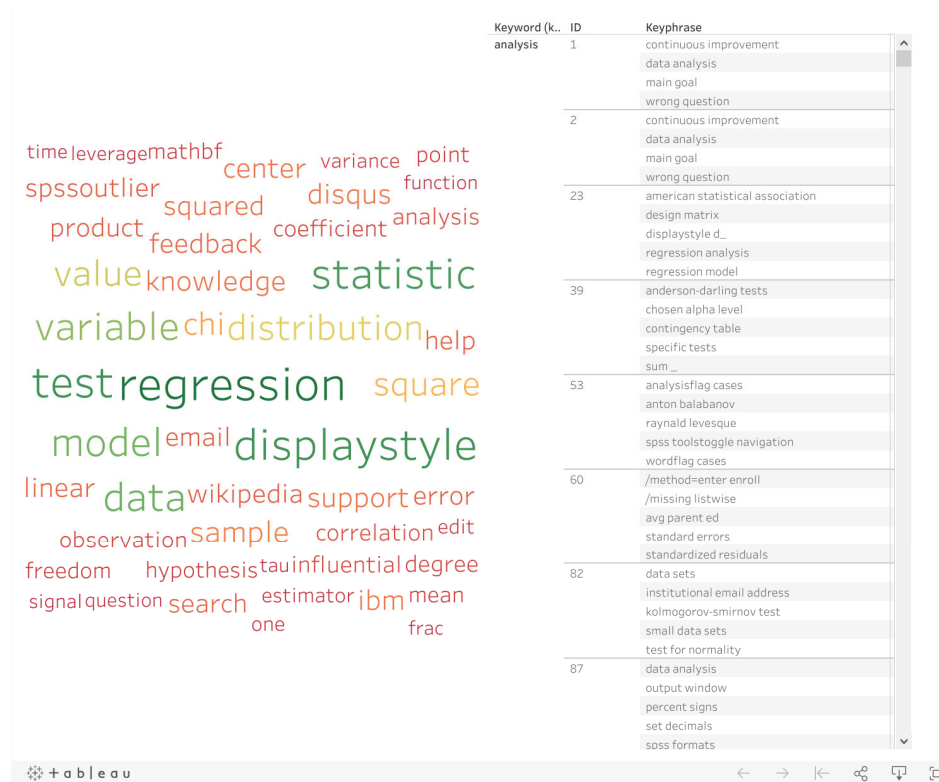| | A | B | C | D |
|---|---|---|---|---|
| 1 | urls | key_word | key_phrase | |
| 2 | blog.minitab.com/blog/adventures-in-statistics-2/how-high-should-r- | squared,regression,question,prediction,minitab,model,statistic,analy: | - t;wrong question;continuous improvement;data analysis;main goal | |
| 3 | blog.minitab.com/blog/adventures-in-statistics/how-high-should-r-sq | squared,regression,question,prediction,minitab,model,statistic,analy: | - t;wrong question;continuous improvement;data analysis;main goal | |
| 4 | courses.lumenlearning.com/introstats1/chapter/testing-the-significar | value,line,correlation,coefficient,significant,population,linear,sample, | y values;standard deviation;best-fit line;correlation coefficient;critical value | |
| 5 | de.wikipedia.org/wiki/Anzahl_der_Freiheitsgrade_(Statistik) | displaystyle,mathbf,freiheitsgrade,top,anzahl,boldsymbol,varepsilon,; | anzahl der freiheitsgrade;sum _;sim chi ^;beta _;displaystyle beta _ | |
| 6 | de.wikipedia.org/wiki/Chi-Quadrat-Test | displaystyle,chi,bearbeiten,cdot,test,quadrat,quelltext,frac,verteilung | taschenbuch der statistik;sum _;der ablehnungsbereich f??r;die pr??fgr????e;displaystyle alpha | |
| 7 | de.wikipedia.org/wiki/F-Test | displaystyle,test,sigma,mathrm,bearbeiten,wert,stichprobe,frac,geq,\| | sigma _;worden w??re;displaystyle =p;displaystyle alpha _;displaystyle f_ | |

**Extracted Data:**

- frequency of keywords -> word cloud
- key word -> related key phrase

| | A | B |
|---|---|---|
| 1 | word_name | fre_num |
| 2 | regression | 88 |
| 3 | test | 79 |
| 4 | displaystyle | 79 |
| 5 | statistic | 76 |
| 6 | data | 68 |

| | A | B |
|---|---|---|
| 1 | ID | key_word |
| 2 | 1 | squared |
| 3 | 1 | regression |
| 4 | 1 | question |
| 5 | 1 | prediction |
| 6 | 1 | minitab |

| | A | B |
|---|---|---|
| 1 | key_phrase | ID |
| 2 | wrong question | 1 |
| 3 | continuous improvement | 1 |
| 4 | data analysis | 1 |
| 5 | main goal | 1 |
| 6 | wrong question | 2 |
| 7 | continuous improvement | 2 |

11

# Identifying Trending Topics - VISUALIZATION



https://public.tableau.com/profile/anuj.dixit#!
/vizhome/shared/J7PW25C8G

# Further Improvements

Possible improvements

- implementing a black list of pages visited helps to block not relevant content, e.g., weather, news pages, ...

- company structure can provide additional context:
  searches based on department structure or job descriptions should be similar

- After tracking down certain content, similarity between pages can help to provide relevant similar pages to look at

- Twitter introduced also an algorithm to identify trending topics.

  - Get trends near a location

  - Get locations with trending topics

  - Basically the algorithms works as follows:

    - extract keywords

    - determine the number of occurrences in tweeds

    - number of occurrences -> time series

    - if slope in time series is large -> arising / trending topic

13

HTW Berlin | Anuj Dixit | Aelee Im | Ting Shen | Cornelia Niklas | Prof. Dr. Tilo Wendler          13