

The ABC of Computational Text Analysis

#6 LEARNING REGULAR EXPRESSIONS

Alex Flückiger

Faculty of Humanities and Social Sciences
University of Lucerne

April 7, 2023

Recap last Lecture

- well-solved assignment #1 
example solution
- counting words 
particular words or entire vocabulary
- preprocessing and cleaning 

Outline

- introducing regular expression 
- practicing the writing of patterns 

Text as Pattern

Formal Search Patterns

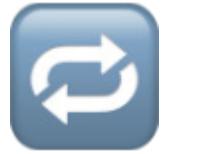
How to extract all email addresses in a text collection?

```
1 Please contact us via info@organization.org.  
2 ---  
3 For specific questions ask Mrs. Green \(.a.green@mail.com\).  
4 ---  
5 Reach out to support@me.ch
```

👉 **Solution:** Write a single pattern to match any valid email address

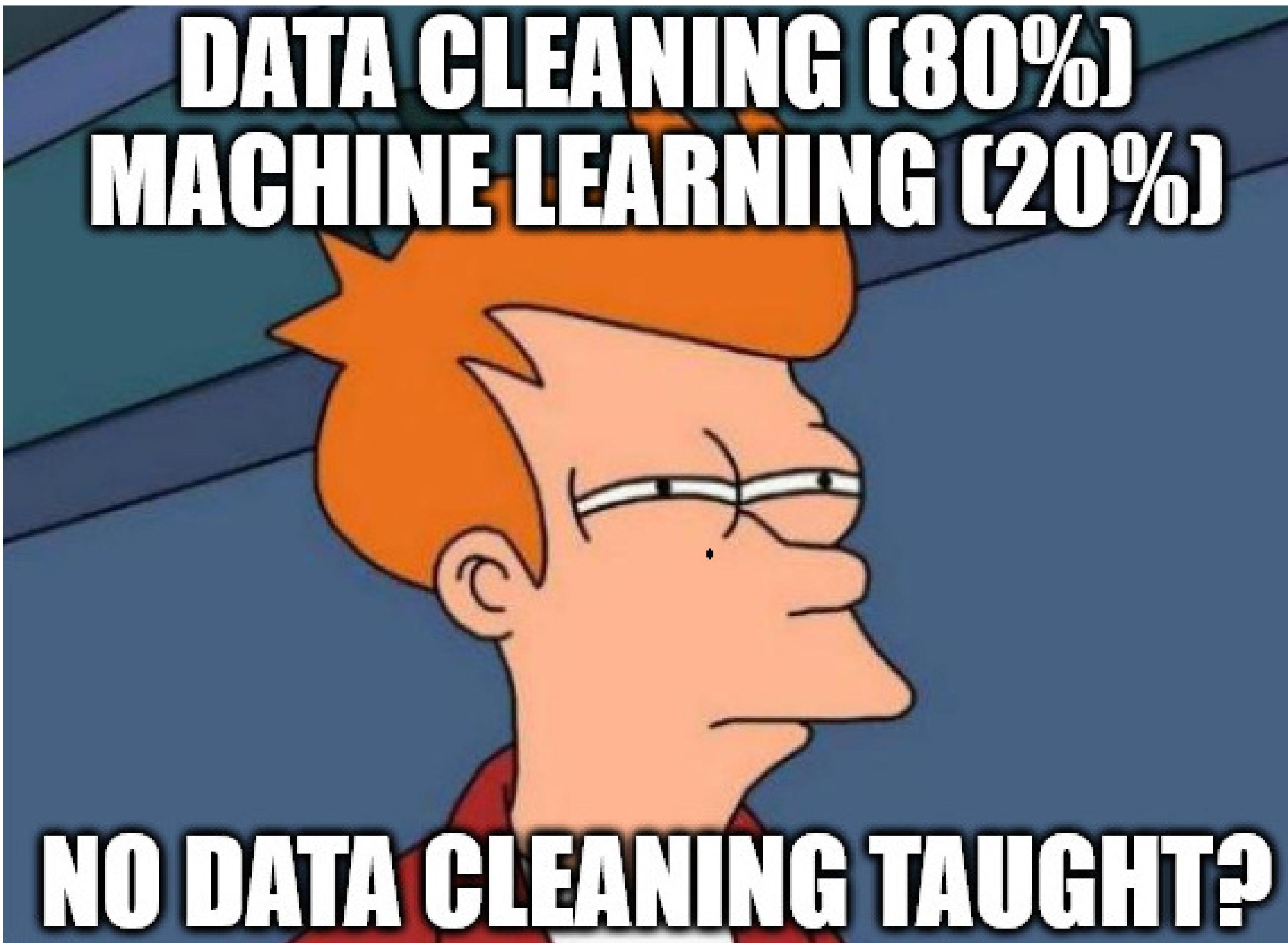
```
1 [A-Z0-9._%+-]+@[A-Z0-9.-]+\.[A-Z]{2,} # match any email address
```

What are patterns for?

- finding 
- extracting 
- removing/cleaning 
- replacing 

... specific parts in texts

Data Cleaning is paramount!



What are Regular Expressions (RegEx)?

RegEx builds on two classes of symbols

- **literal** characters and strings
 - letters, digits, words, phrases, dates etc.
- **meta** expressions with special meaning
 - e.g., `\w` represents alphanumeric characters
 - `[Cc]o+l` → Col, col, Cool, coool ...
- akin to regular languages

Finding + Extracting

extended globally search for regular expression and print (egrep)

- tool to filter/keep matching lines only

```
1 # check a regular expression quickly
2 echo "check this pattern" | egrep "pattern"
3
4 egrep "yes" file.txt          # search in a specific file
5 egrep -r "yes" folder        # search recursively within folder
6
7 egrep "yes" *.txt            # keep lines containing pattern (yes)
8 egrep -i "yes" *.txt        # ditto, ignore casing (Yes, yes, YES . .
9 egrep -v "noisy" *.txt      # do NOT keep lines containing noisy
10
11 # extract raw match only to allow for subsequent counting
12 egrep -o "only" *.txt       # print match only instead of entire line
13 egrep -h "only" *.txt       # suppress file name
```

Quantifiers

repeat preceding character X times

- $?$ zero or one
- $+$ one or more
- $*$ zero or any number
- $\{n\}$, $\{m, n\}$ a specified number of times

```
1 egrep -r "Bundesrath?es"          # match old and new spelling
2 egrep -r "a+"                      # match one or more "a"
3 egrep -r "e{2}"                    # match sequence of two "e"
```



Do not confuse regex with Bash wildcards!

Character Sets

- [. . .] any of the characters between brackets
 - any vowel: [auoei]
 - any digit: [0-9]
 - any letter: [A-Za-z]
- [^ . . .] any character but none of these (negation)
 - anything but the vowels: [^auoei]

```
1 # match the capitalized and non-capitalized form
2 egrep -r "[Gg]rüne"
3
4 # match sequences of 3 vowels
5 egrep -r [aeiou]{3}
6
7 # extract all bigrams (sequence of two words)
8 egrep -rohi "[a-zA-Z]+ [a-zA-Z]+"
```

Special Symbols

- `.` matches any character (excl. newline)
- `\.` escapes to match literal
 - `\.` means the literal `.` instead of “any symbol”
- `\w` matches any alpha-numeric character
 - same as `[A-Za-z0-9_]`
- `\s` matches any whitespace (space, newline, tab)
 - same as `[\t\n]`

```
1 # match anything between brackets
2 egrep -r "\(.*\)"
```

The power of . * ...

matches *any character any times*

More Complex Examples

```
1 # extract basename of URLs
2 egrep -ro "www\.\w+\.[a-z]{2,}"
3
4 # extract valid email addresses (case-insensitive)
5 egrep -iro "[A-Z0-9._%+-]+@[A-Z0-9.-]+\.[A-Z]{2,}" **/*.txt
```

Combining RegEx with Frequency Analysis

something actually useful

```
1 # count political areas by looking up words ending with "politik"
2 egrep -rioh "\w*politik" */*.txt | sort | uniq -c | sort -h
3
4 # count ideologies/concepts by looking up words ending with "ismus"
5 egrep -rioh "\w*ismus" */*.txt | sort | uniq -c | sort -h
```



Start simple,
add complexity
subsequently.

In-class: Exercise

1. Use the command line to navigate to the local copy of the Github repository KED2023 and make sure it is up-to-date with git pull. Change in to the directory materials/data/swiss_party_programmes/.txt.
2. Use egrep to extract all uppercased words like UNO, OECD, SP and count their frequency.
3. Use egrep to extract all plural nouns with female endings e.g. Schweizerinnen (starting with an uppercase letter, ending with innen, and any letter in between). Do the same for the male forms. Is there a qualitative or a quantitative difference between the gendered forms?

```
1 # Some not so random hints
2 piping with |
3 sort
4 uniq -c
5 egrep -roh **/* .txt
```

In-class: Solution

1. Use egrep to extract all uppercased words like UNO, OECD, SP and count their frequency.

```
egrep -roh "[A-Z]{2,}" */*.txt | sort | uniq -c | sort -h
```

2. Use egrep to extract all plural nouns with female endings

e.g. Schweizerinnen (starting with an uppercase letter, ending with innen, and any letter in between). Do the same for the male forms. Is there a qualitative or a quantitative difference between the gendered forms?

```
egrep -roh "[A-Z][a-z]+innen\b" */*.txt | sort | uniq -c | sort -h
```

```
egrep -roh "[A-Z][a-z]+er\b" */*.txt | sort | uniq -c | sort -h
```

(there is no way with regular expression to extract only nouns of the male form but not Wasser and the like. For this, you have to use some kind of machine learning.)

Replacing + Removing

stream editor (sed)

- advanced find + replace using regex
`sed "s/WHAT/WITH/g" file.txt`
- `sed` replaces any sequence, `tr` only single symbols

```
1 echo "hello" | sed "s/llo/y/g"          # replace "llo" with a "y"
2
3 # by setting the g flag in "s/llo/y/g",
4 # sed replaces all occurrences, not only the first one
```

Contextual Replacing

reuse match with grouping

- define a group with parentheses (`group_pattern`)
- `\1` equals the expression inside first pair of parentheses
- `\2` expression of second pair
- ...

```
1 # swap order of name (last first -> first last)
2 echo "Lastname Firstname" | sed -E "s/(.+)\ (.+)/\2 \1/"
3
4 # matching also supports grouping
5 # match any pair of two identical digits
6 egrep -r "([0-9])\1"
```

More Meta-Symbols

- `\b` matches word boundary
`word\b` does not match `words`
- `^` matches begin of line and `$` end of line
`^A` matches only `A` at line start
- `|` is a disjunction (OR)
`(Mr|Mrs|Mr\.|Mrs\.) Green` matches alternatives

Greediness Trap

- greedy ~ match the longest string possible
- quantifiers `*` or `+` are greedy
- non-greedy by excluding some symbols
[`^EXCLUDE_SYMBOLS`] instead of `.*`

```
1 # greedy: an apple, other apple
2 echo "an apple, other apple" | egrep "a.*apple"
3
4 # non-greedy: an apple
5 echo "an apple, other apple" | egrep "a[^,]*apple"
```

Assignment #2



- get/submit via OLAT
 - starting tomorrow
 - deadline 15 April 2023, 23:59
- use forum on OLAT
 - subscribe to get notifications
- ask friends for support, not solutions

In-class: Exercises I

1. Use egrep to extract capitalized words and count them. What are the most frequent nouns?
2. Use egrep to extract words following any of these strings: der die das.
Hint: Use a disjunction.
3. Do the self-check on the next slide.
4. Use sed -E to remove the table of content, the footer and the page number in the programme of the Green Party. Check the corresponding PDF to get a visual impression and test your regular expression with egrep first to see if you match the correct parts in the document.

In-class: Solution I

1. Use egrep to extract capitalized words and count them. What are the most frequent nouns?

```
egrep -roh "[A-Z][a-z]+" **/*.txt | sort | uniq -c | sort -h
```

2. Use egrep to extract words following any of these strings: der die das.
Hint: Use a disjunction.

```
egrep -roh "(der|die|das) \w+" **/*.txt
```

3. Use sed -E to remove the table of content, the footer and the page number in the programme of the Green Party. Check the corresponding PDF to get a visual impression and test your regular expression with egrep first to see if you match the correct parts in the document.

```
cat gruene_programme_2019.txt | sed "1,192d" | sed -E  
"s/^Wahlplattform.*2023$/g" | sed -E "s/^[\t\n]+$/g"
```

In-class: Self-Check

equivalent patterns

```
1 a+ == aa*                      # "a" once or more than once
2 a? == _(a|_)                   # "a" once or nothing
3 a{3} == aaa                     # three "a"
4 a{2,3} == _(aa|aaa)            # two or three "a"
5 [ab] == _(a|b)                  # "a" or "b"
6 [0-9] == _(0|1|2|3|4|5|6|7|8|9) #any digit
```

In-class: Exercise II

1. Count all the bigrams (sequence of two words) using character sets and quantifiers. What about trigrams (three words)?
2. Extract the words following numbers (also consider numbers like: 1 '000, 1,000 or 5%). Then, count all the words while excluding the numbers themselves. Hint: Pipe another grep to remove the digits.
3. You are ready to come up with your own patterns...

In-class: Solution II

1. Count all the bigrams (sequence of two words) using character sets and quantifiers. What about trigrams (three words)?

```
egrep -hoir "\b[a-z]+ [a-z]+\b" | sort | uniq -c | sort -h
```

```
egrep -hoir "\b[a-z]+ [a-z]+\b" | sort | uniq -c | sort -h
```

2. Extract the words following numbers (also consider numbers like: 1'000, 1,000 or 5%). Then, count all the words while excluding the numbers themselves. Hint: Pipe another grep to remove the digits.

```
egrep -rhoi "[0-9][0-9, '%]+ [a-z]+" | egrep -io "[a-z]+" | sort | uniq -c | sort -h
```

Alternative: `egrep -rhoi "[0-9][0-9, '%]+ [a-z]+" | sed -E "s/[0-9][0-9, '%]+//g" | sort | uniq -c | sort -h`

In-class: Exercise III

1. Since you know about RegEx, we can use a more sophisticated tokenizer to split a text into words. What is the difference between the old and new approach? Test it and check the helper page with man.

```
1 # new, improved approach
2 cat text.txt | tr -sc "[a-zäöüA-ZÄÖÜĞ-9-]" "\n"
3
4 # old approach
5 cat text.txt | tr " " "\n"
```

More Resources

required

- Ben Schmidt. 2019. Regular Expressions.
- Cheatsheet of this course

highly recommended

- Nikolaj Lindberg. egrep for Linguists.

online regular expression editor



Questions?