# text mining for poets
## DTL SU @ AU

kristoffer l nielbo
`kln@cas.au.dk`
`github.com/digitaltxtlab`

DAI|IMC|AARHUS UNIVERSITY

**text-mining for poets**

- text is available like never before
    - dictionaries, corpora, full-text DBs ...
    - information collection efforts [1]
    - `Information Super Highway Roadkill`, emails, blogs, websites ...
    - *billions & billions of words*

- what can we do with it all?

- it is better to do something simple, than nothing at all

- DIY is more satisfying than begging for "help" from a computer officer

**Pannang Curry 1-2-3!**

1. Open your right refrigerator door and remove ingredients from the following locations: Door shelf 2, Spot 1; Crisper drawer 1, Spot 3; Crisper drawer 1, Spot 5.

2. Open your third kitchen drawer and remove the utensils labeled "1", "3", "4", "9", and "12".

3. Use your arms to apply utensil 1 to ingredients 1-3. Place ingredients inside utensil 3.

*Note: This recipe uses ShaKL the Shared Kitchen Layout. To use ShaKL, you'll need to have installed ShaKL shelving, cabinetry, and utensils throughout your kitchen and pantry and have basic understanding of ShaKL managers. To learn more, read *Up and Running with ShaKL* (O'Billy Press, 2015). Want to improve ShaKL? Consider contributing to our team.

AARHUS UNIVERSITET

IMC
INTERACTING MINDS CENTRE

- a `shell` is a program whose primary purpose is to read commands and run other programs

- the `shell`'s main advantages are its high action-to-keystroke ratio, its support for automating repetitive tasks, and its capacity to access networked machines

- the `shell`'s main disadvantages are its primarily textual nature and how cryptic its commands and operation can be

`PS1='$ '` sets prompt string in console to

```
1  $
```

`prompt` indicates that the `shell` is waiting for input

```
1  $ whoami
2  kln
3  $
```

user ID or who the shell thinks you are

**whoami**

1. finds a program called `whoami`
2. runs that program
3. displays that program's output
4. displays a new prompt to tell us that it's ready for more commands

AARHUS UNIVERSITET

IMC
INTERACTING MINDS CENTRE

**unknown command**

```
1  $ somecommand
2  somecommand: command not found
3  $
```

- the **shell** runs other programs, so it does not work if the program does not exist
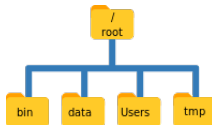
**print working directory** - current default directory

```
1  $ pwd
2  home/kln
3
4  $ a=$(pwd)
5  $ echo "current wd is: $a"
6  current wd is: /home/kln
```

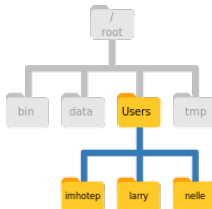the path to the home directory varies between operating systems:

- [linux] /home/yourname

- [mac] /Users/yourname

- [windows] C:\Users\yourname

organization of a filesystem (somewhat unix-centric view) #1



- at the top is the `root` directory that holds everything else (refer to by a / as the leading slash in /Users/kln

- the `root` directory contains multiple directories: `bin` for built-in programs, `data` for miscellaneous data files, `Users` for personal directories, `tmp` for temporary files &c

- so the current working directory /Users/kln is stored in `Users` which again is stored in `root` (because of /)

organization of a filesystem #2



- underneath `Users` each user with an account on the computer has a directory

`ls` - listing content of a directory

```
1  $ ls
2  Applications Documents    genesis       Movies        Pictures
3  Desktop      Downloads    Library       Music         Public
4
5  $ ls -F
6  Applications/ Documents/   genesis       Movies/       Pictures/
7  Desktop/      Downloads/   Library/      Music/        Public/
8
9  $ ls -F Desktop
10 creatures/         molecules/       notes.txt         solar.pdf
11 data/              north-pacific-gyre/ pizza.cfg       writing/
12
13 $ ls --help
```

## navigating a filesystem

### cd - change directory

```
1  $ cd Desktop
2  $ cd data
3  $ pwd
4  /Users/nelle/Desktop/data
5  $ cd ..
6  $ pwd
7  /Users/nelle/Desktop
```

### shortcuts

```
1   $ pwd
2   /Users/nelle/Desktop
3   $ cd
4   $ pwd
5   /Users/nelle
6   $ cd -
7   /Users/nelle
8   $ pwd
9   /Users/nelle
10  $ cd ~/Documents/books
11  $ pwd
12  /Users/nelle/Documents/books
```

AARHUS UNIVERSITET

IMC
INTERACTING MINDS CENTRE

but the `shell` can do more than navigation, it can manipulate (text) data →

# words|unstructured

**The First Book of Moses, called Genesis**

(1:1) In the beginning God created the heaven and the earth. (1:2) And the earth was without form, and void; and darkness [was] upon the face of the deep. And the Spirit of God moved upon the face of the waters.

(1:3) And God said, Let there be light: and there was light. (1:4) And God saw the light, that [it was] good: and God divided the light from the darkness. (1:5) And God called the light Day, and the darkness he called Night. And the evening and the morning were the first day.

(1:6) And God said, Let there be a firmament in the midst of the waters, and let it divide the waters from the waters. (1:7) And God made the firmament, and divided the waters which [were] under the firmament from the waters which [were] above the firmament: and it was so. (1:8) And God called the firmament Heaven. And the evening and the morning were the second day.

(1:9) And God said, Let the waters under the heaven be gathered together unto one place, and let the dry [land] appear: and it was so. (1:10) And God called the dry [land] Earth; and the gathering together of the waters called he Seas: and God saw that [it was] good. (1:11) And God said, Let the earth bring forth grass, the herb yielding seed, [and] the fruit tree yielding fruit after his kind, whose seed [is] in itself, upon the earth: and it was so. (1:12) And the earth brought forth grass, [and] herb yielding seed after his kind, and the tree yielding fruit, whose seed [was] in itself, after his kind: and God saw that [it was] good. (1:13) And the evening and the morning were the third day.

(1:14) And God said, Let there be lights in the firmament of the heaven to divide the day from the night; and let them be for signs, and for seasons, and for days, and years: (1:15) And let them be for lights in the firmament of the heaven to give light upon the earth: and it was so. (1:16) And God made two great lights; the greater light to rule the day, and the lesser light to rule the night: [he made] the stars also. (1:17) And God set them in the firmament of the heaven to give light upon the earth, (1:18) And to rule over the day and over the night, and to divide the light from the darkness: and God saw that [it was] good. (1:19) And the evening and the morning were the fourth day. (1:20) And God said, Let the waters bring forth abundantly the moving creature that hath life, and fowl [that] may fly above the earth in the open firmament of heaven. (1:21) And God created great whales, and every living creature that moveth, which the waters brought forth abundantly, after their kind, and every winged fowl after his kind: and God saw that [it was] good. (1:22) And God blessed them, saying, Be fruitful, and multiply, and fill the waters in the seas, and let fowl multiply in the earth. (1:23) And the evening and the morning were the fifth day.

(1:24) And God said, Let the earth bring forth the living creature after his kind, cattle, and creeping thing, and beast of the earth after his kind: and it was so. (1:25) And God made the beast of the earth after his kind, and cattle after their kind, and every thing that creepeth upon the earth after his kind: and God saw that [it was] good.

(1:26) And God said, Let us make man in our image, after our likeness: and let them have dominion over the fish of the sea, and over the fowl of the air, and over the cattle, and over all the earth, and over every creeping thing that creepeth upon the earth. (1:27) So God created man in his [own] image, in the image of God created he him; male and female created he them. (1:28) And God blessed ...

AARHUS UNIVERSITET

IMC
INTERACTING MINDS CENTRE

### case-folding

```
1  $ tr A-Z a-z
2  In the beginning God created the heaven and the earth
3  in the beginning god created the heaven and the earth
```

### translate from and to file

```
1  $ tr "God" "Joe" < genesis > jenesis
2  $ head genesis | tr "God" "Joe"
3
4  $ tr "{}" " " < genesis > genesis_nobrace
```

### remove (squeeze) repetitions

```
1  echo  "  In the beginning God created the heaven and the earth" | tr -s [:space:] " "
2  : In the beginning God created the heaven and the earth
```

### remove characters

```
1  $ echo "{1:1} In the beginning God created the heaven and the earth" | tr -d "0-9"
2  {:} In the beginning God created the heaven and the earth
```

### complement

```
1  $ echo "{1:1} In the beginning God created the heaven and the earth" | tr -cd "0-9:"
2  1:1
```

### collapse lines

```
1  $ tr -s "\r\n" " " < genesis > genesis_line
2  $ tr -s "\n" " " < filename > new_filename
```

### man(ual) page for more information

```
1  $ man tr
```

#### `tokenization` - unigrams

```
1  $ tr -sc "A-Za-z" "\n" < genesis
```

#### `sort` in alphabetic order

```
1  $ tr -sc "A-Za-z" "\n" < genesis | sort
```

#### `uniq` - lexicon of document

```
1  $ tr -sc "A-Za-z" "\n" < genesis | sort | uniq
2  $ tr -sc "A-Za-z" "\n" < genesis | sort | uniq -c
3  $ tr -sc "A-Za-z" "\n" < genesis | sort | uniq -c > lexicon.txt
```

### file conversion

```
1  man pdftotext
2  man catdoc
```

### script for char extraction from pdf files in directory

```
1  #!/bin/bash
2
3  for f in *.pdf
4  do
5    echo "converting: - $f"
6    pdftotext $f $f.txt
7  done
```

tired of cryptic commands and operations from the command line?

luckily we have:

```
1  if questions:
2      try:
3          answer()
4      except RunTimeError:
5          pass
6      else:
7          print "thank you"
```