

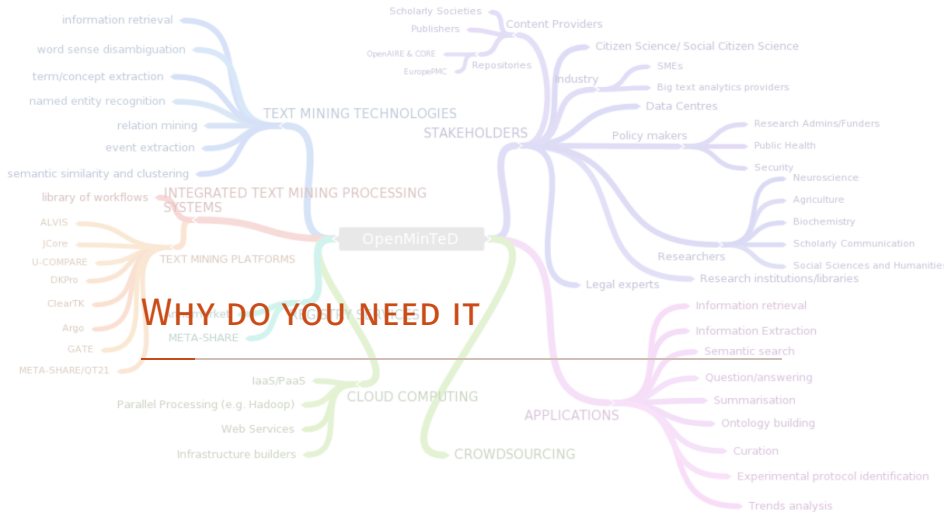
COMPUTATIONAL METHODS FOR TEXT ANALYSIS

BA PROGRAM “SOCIOLOGY AND SOCIAL INFORMATICS”

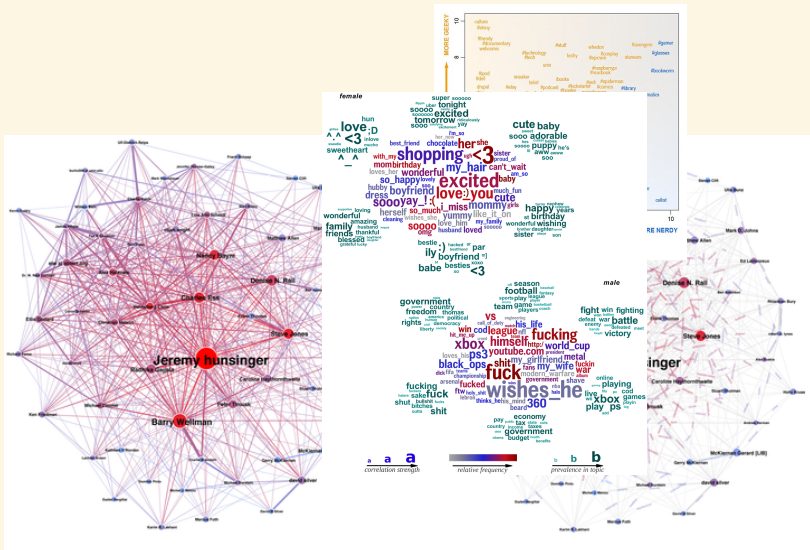
Kirill Maslinsky

2018

Higher School of Economics — Saint Petersburg



JUST TO LEARN TO MAKE THOSE PICTURES



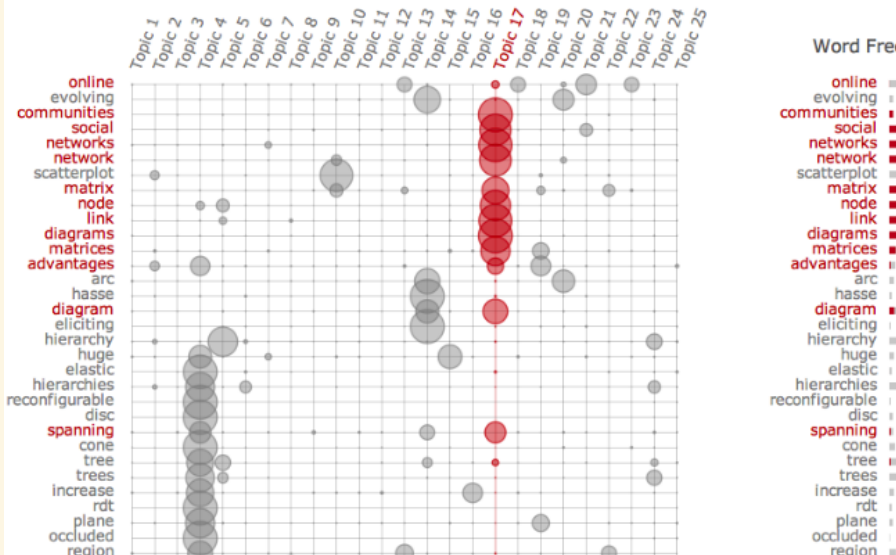
population studied	“all social media users of a town”
time spans	“all of the Post-Soviet history”
geographical scope	“all educational migration in Russia”

- provide basic understanding of how to properly use collections of texts as **quantitative evidence**,
- and to make this knowledge **practical**

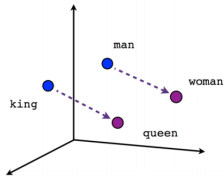


COURSE CONTENT

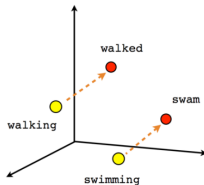
BREAD AND BUTTER: TOPIC MODELING



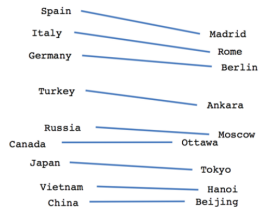
KILLER FEATURE: WORD EMBEDDINGS



Male-Female

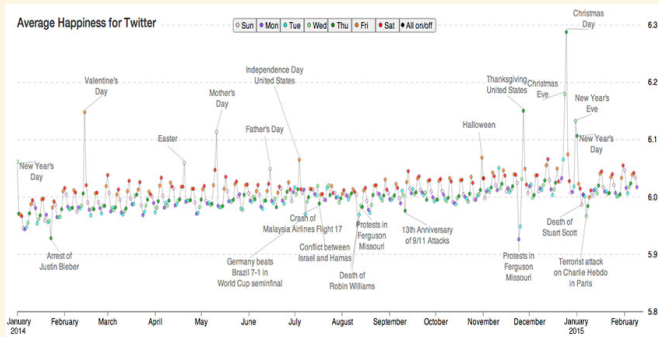


Verb tense

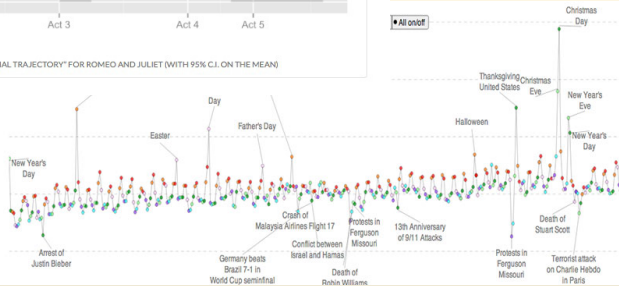
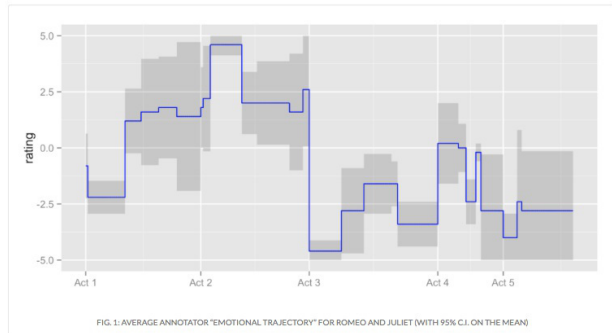


Country-Capital

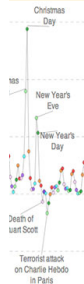
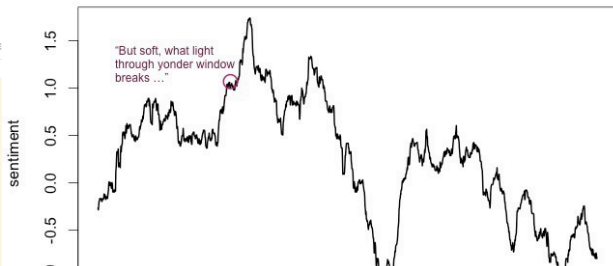
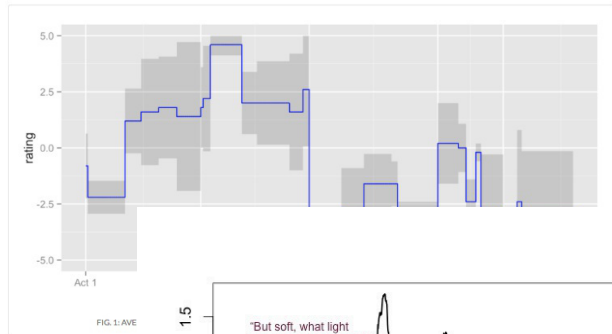
THE ICING ON THE CAKE: SENTIMENT ANALYSIS



THE ICING ON THE CAKE: SENTIMENT ANALYSIS



THE ICING ON THE CAKE: SENTIMENT ANALYSIS



COURSE TOPICS

- Basic word statistics:
 - lexical statistics (word frequency distributions),
 - distributive semantics (word co-occurrence patterns),
 - vector representation of text.
- Methods for supervised and unsupervised modeling:
 - dictionary methods,
 - document classification and clusterization,
 - topic modeling,
 - word embeddings,
 - sequence modeling.
- Applied tasks:
 - automating content analysis (extracting theme and topic),
 - sentiment analysis,
 - information extraction from unstructured text.

$[X, P] = i\hbar$ $\langle E_1 | \psi \rangle = \sqrt{2} S_1 \dots$ $\frac{d}{dt} \langle \psi | \psi \rangle = 0$ $\frac{d}{dt} \langle \psi | \hat{H} | \psi \rangle = \langle \psi | [\hat{H}, \hat{H}] | \psi \rangle = 0$
 $H = \frac{p^2}{2m} + V(x)$ $\langle E_1 | \psi \rangle = \sqrt{2} S_1 \dots$ $\frac{d}{dt} \langle \psi | \hat{H} | \psi \rangle = \langle \psi | \frac{d\hat{H}}{dt} | \psi \rangle$
 $H|\psi\rangle = E|\psi\rangle$ $\langle \psi | \hat{H} | \psi \rangle = \int \psi^* (\hat{H} \psi) dx = E \int \psi^* \psi dx = E$ $\frac{d}{dt} \langle \psi | \hat{H} | \psi \rangle = \langle \psi | \frac{d\hat{H}}{dt} | \psi \rangle$

$\hat{X} = \frac{m\omega}{\hbar} X$ $\hat{P} = \frac{\hbar}{i} \frac{d}{dx}$
 $\langle \hat{X}^2 \rangle = \frac{\hbar^2}{m^2 \omega^2} \langle \psi | \frac{d^2}{dx^2} | \psi \rangle$ $\langle \hat{P}^2 \rangle = \hbar^2 \langle \psi | \frac{d^2}{dx^2} | \psi \rangle$
 $\langle \hat{X} \hat{P} \rangle = \frac{\hbar}{i} \langle \psi | X \frac{d}{dx} | \psi \rangle$ $\langle \hat{P} \hat{X} \rangle = \frac{\hbar}{i} \langle \psi | \frac{d}{dx} X | \psi \rangle$

$a = \frac{1}{\sqrt{2}} (\hat{X} + i\hat{P})$ $a^\dagger = \frac{1}{\sqrt{2}} (\hat{X} - i\hat{P})$
 $a^\dagger a = \hat{N}$ $a a^\dagger = \hat{N} + 1$
 $[a, a^\dagger] = 1$

WHAT TO EXPECT

$\hat{H} = a^\dagger a + \frac{1}{2} E = \hbar \omega \left(\hat{N} + \frac{1}{2} \right)$ $E_n = \hbar \omega \left(n + \frac{1}{2} \right)$
 $\langle \hat{X} \rangle = 0$ $\langle \hat{P} \rangle = 0$ $\langle \hat{X}^2 \rangle = \frac{\hbar}{m\omega} \left(n + \frac{1}{2} \right)$ $\langle \hat{P}^2 \rangle = \hbar m \omega \left(n + \frac{1}{2} \right)$

$\langle \hat{X}^4 \rangle = \frac{3\hbar^2}{m^2 \omega^2} \left(n + \frac{1}{2} \right) \left(n + \frac{3}{2} \right)$ $\langle \hat{P}^4 \rangle = 3\hbar^2 m^2 \omega^2 \left(n + \frac{1}{2} \right) \left(n + \frac{3}{2} \right)$
 $\langle \hat{X}^2 \hat{P}^2 \rangle = \frac{5\hbar^2}{m^2 \omega^2} \left(n + \frac{1}{2} \right) \left(n + \frac{3}{2} \right) \left(n + \frac{1}{2} \right) \left(n + \frac{3}{2} \right)$

$\langle \hat{X}^2 \rangle = \frac{\hbar}{m\omega} \left(n + \frac{1}{2} \right)$ $\langle \hat{P}^2 \rangle = \hbar m \omega \left(n + \frac{1}{2} \right)$
 $\langle \hat{X}^4 \rangle = \frac{3\hbar^2}{m^2 \omega^2} \left(n + \frac{1}{2} \right) \left(n + \frac{3}{2} \right)$ $\langle \hat{P}^4 \rangle = 3\hbar^2 m^2 \omega^2 \left(n + \frac{1}{2} \right) \left(n + \frac{3}{2} \right)$

$\langle \hat{X}^2 \hat{P}^2 \rangle = \frac{5\hbar^2}{m^2 \omega^2} \left(n + \frac{1}{2} \right) \left(n + \frac{3}{2} \right) \left(n + \frac{1}{2} \right) \left(n + \frac{3}{2} \right)$



HOW COURSEWORK WILL BE ORGANIZED

- An interesting recent article
- with an explanation of the necessary concepts and methods during lecture
- followed by detailed analysis of the method in class
- concluded by the task to reproduce the method with your own data

Practical work with real texts in class and at home.

- command line
- mining your own text collection
- R scripts
- bugs in scripts, googling, bugs in scripts again
- seeking and getting help from your peers and course instructor
- **happy end**

WORK IN GROUPS



WHAT YOU CAN LEARN

- State-of-the-art of natural language processing:
 - solved problems
 - topical issues and unsolved problems
- Terms:
 - a minimal vocabulary of necessary linguistic terms (**with meanings! :)**)
 - appropriate keywords to search for current research and tools
- Tools:
 - Where to apply methods for computational text analysis and how to interpret their results
 - Existing software for text analysis (for Russian and English)
 - Existing linguistic resources — dictionaries, corpora, pre-trained models (for Russian and English)