# Network Science

A.Y. 23/24

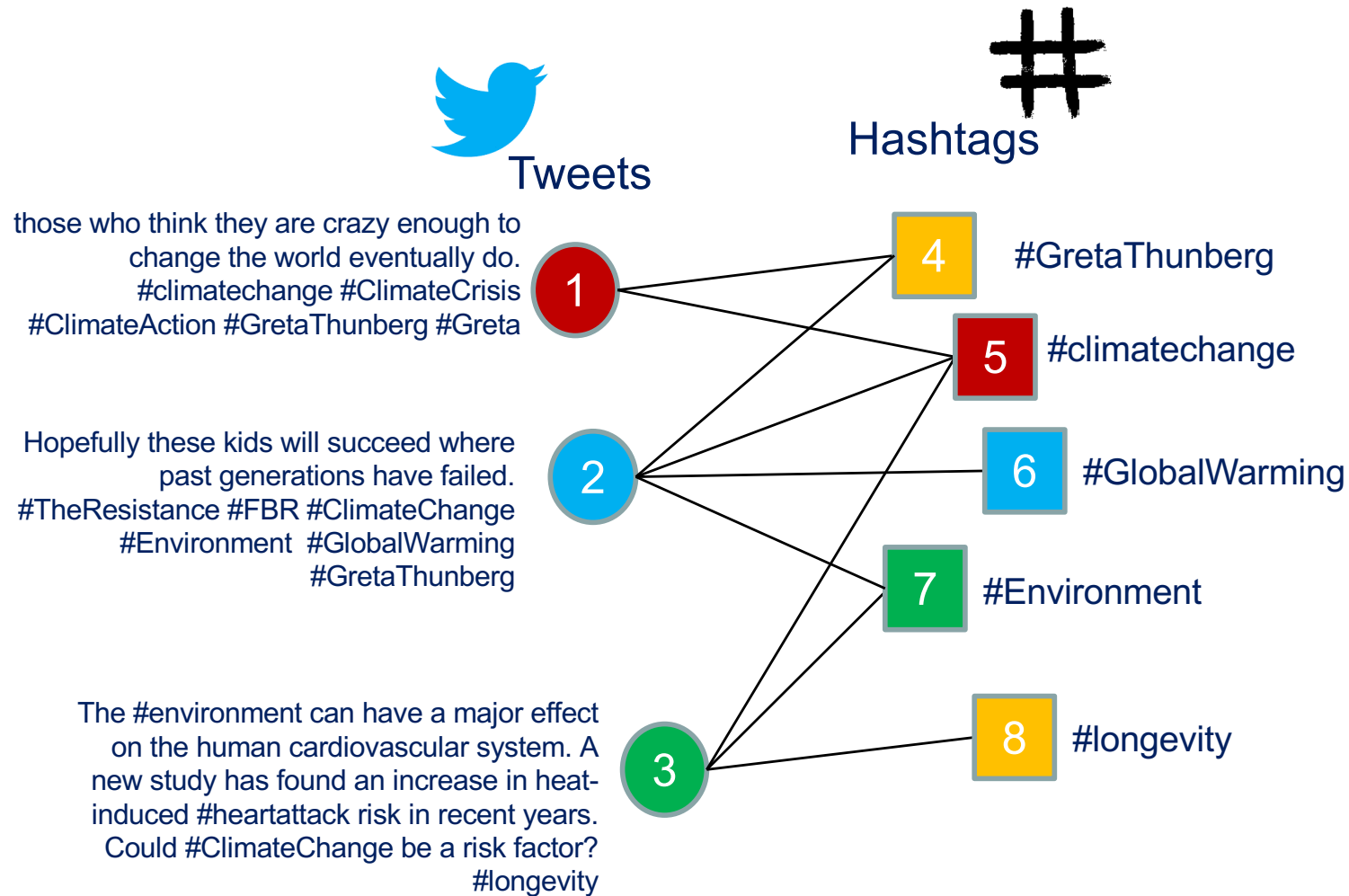ICT for Internet & multimedia, Data science, Physics of data

# Semantic networks

## network science tools for their study

UNIVERSITÀ
DEGLI STUDI
DI PADOVA

Tweets

Hashtags

those who think they are crazy enough to change the world eventually do. #climatechange #ClimateCrisis #ClimateAction #GretaThunberg #Greta

**1**

**4** #GretaThunberg

**5** #climatechange

Hopefully these kids will succeed where past generations have failed. #TheResistance #FBR #ClimateChange #Environment #GlobalWarming #GretaThunberg

**2**

**6** #GlobalWarming

**7** #Environment

The #environment can have a major effect on the human cardiovascular system. A new study has found an increase in heat-induced #heartattack risk in recent years. Could #ClimateChange be a risk factor? #longevity

**3**

**8** #longevity

❑ Data collection + polishing

❑ Building the semantic network (bipartite/projections)

❑ Topic (i.e., community) detection

   ❑ Modularity & InfoMap

   ❑ Non-negative matrix factorization (NMF)

   ❑ Latent Dirichlet allocation (LDA)

   ❑ Variational auto-encoders (VAE)

   ❑ Embeddings and BERTopic

# Data collection

how to get data from the Internet using APIs

**Twitter's plan to cut off free data access evokes 'fair amount of panic' among scientists**

Social media platform's intent to increase revenue could end or limit many research projects

8 FEB 2023 · 4:35 PM ET · BY KAI KUPFERSCHMIDT

**Twitter's plan to charge researchers for data access puts it in EU crosshairs**

Elon Musk's social media giant plans to charge academics to access its data – in potential violation of Europe's content rules

BY MARK SCOTT
MARCH 22, 2023

**Academic researchers blast Twitter's data paywall as 'outrageously expensive'**

By Brian Fung, CNN
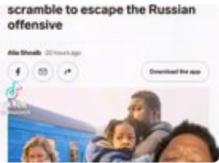Published 11:40 AM EDT, Wed April 5, 2023

7

# Reddit apps
https://www.reddit.com/prefs/apps

```
!pip install praw
```

```python
import pandas as pd
import praw
reddit = praw.Reddit(client_id='Qbdk-FkA9jSQB9T7drY8UQ',
                     client_secret="jtGPdqiaTj6hCWcvRPeS_nMNEnVkxw",
                     user_agent='reddit scraper 1.0 by u/Upbeat-Lychee-6630',
                     check_for_async=False)
print(reddit.read_only)
```

```
True
```

your own description of your app, including version and username

from Reddit apps

```python
df = pd.DataFrame([vars(post) for post in reddit.subreddit("all")
                   .search("#ukrainewar", sort='top', limit=10)])
```

"relevance", "hot", "top", "new", or "comments"

max 250 per call

can also add time_filter = "all", "day", "hour", "month", "week", or "year"

```python
df.to_excel('drive/MyDrive/Colab Notebooks/samples.xlsx',
            index=True)
```

10

there is a list of 116 entries per post,
on which you can choose!!!

from this you extract the date

| | title | created | score | upvote_ratio | ups | num_comments | selftext |
|---|---|---|---|---|---|---|---|
| 0 | Damn...we blinked and missed the T-34 stage of... | 1.666899e+09 | 10394 | 0.99 | 10394 | 738 | |
| 1 | Finnish🇫🇮 volunteer sends greetings home from ... | 1.680237e+09 | 2095 | 1.00 | 2095 | 57 | |
| 2 | Guess having 5 trucks fall into your office ca... | 1.663341e+09 | 1974 | 1.00 | 1974 | 88 | |
| 3 | [META] Important - Russia-Ukraine Crisis/War: ... | 1.645712e+09 | 1284 | 0.89 | 1284 | 1 | Hi, /u/Anonim97 here.\n\nWe - as a mods of 40k... |
| 4 | V*tniks coping hard Over the counter offensive... | 1.662956e+09 | 1081 | 1.00 | 1081 | 86 | |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 0 | Russia Ukraine War. | 1.695349e+09 | 1 | 1.00 | 1 | 0 | |

UNIVERSITÀ
DEGLI STUDI
DI PADOVA

♪ TikTok **for developers**

Products    Docs    Support    Blog    👤

Content Posting API ⌄

Display API ⌄

Research API ⌃

  About Research API

  Getting Started

  Frequently Asked Questions

  Codebook

  API Reference ⌃

    Query Videos

    Query User Info

    Query Video Comments

# Query Videos

## Request

| HTTP URL | https://open.tiktokapis.com/v2/research/video/query/ |
|---|---|
| HTTP Method | POST |
| Scopes | research.data.basic |

## Headers

| Key | Type | Description |
|---|---|---|

♪ TikTok **for developers**

### TikTok Research API application received

Thank you for applying to gain access to TikTok's research APIs. We have received your application and will proceed with the review process. You will be notified of the result via email within 3-4 weeks. If you have any questions, please contact our Support team.

12

UNIVERSITÀ
DEGLI STUDI
DI PADOVA

ENSEMBLEDATA   Home   Products ⌄   Documentation   Pricing

1 week trial period… register later on

# Social media scraping through simple APIs

Fetch data from TikTok, Instagram, YouTube through simple APIs.
Real-time, fast, reliable and easy to integrat

[>_]  Get started    [▤]  Documentatio

EnsembleData / **TikTokScraper**   Public

<> Code   ⊙ Issues   ⇄ Pull requests   ▷ Actions

⌥ main ⌄    ⑂ 1 branch    ⬠ 0 tags

fracogno Fix

📁 images                                      upload tt logo

📁 src                                         Fix

https://github.com/EnsembleData/TikTokScraper    Update README.md

13

# Data preprocessing

how to polish raw data from the Internet

1

## 1. Superficial cleaning

Removing website links
Removing accented characters
Removing text inside square brackets
Removing moderator messages
Removing double spaces
Removing non-text special words and characters
Removing extra-used new lines
Limiting all the repetitions to two characters and removing the extra characters
Removing punctuation except main sentence punctuation
Removing sentences that represent the rules of the community

Fixing **contractions**
Removing **emoji**
Removing **hashtags** and **mentions**
Removing **numbers**
**Lowercasing**
Correct **spellings**

the bare minimum to polish the text, useful as an input to sentiment analysis

## 2. Subsentence

Tokenise subsentences

useful for long text samples (e.g., Reddit)

## 3. Deep cleaning

Stop word removal
Word tokenization
POS tagging
Lemmatization

truly polished text, useful for building a semantic network

| PRP | VBZ | NNS | IN | DT | NN |
|---|---|---|---|---|---|
| She | sells | seashells | on | the | seashore |

# SpaCy part-of-speech (POS) tags

| POS | description | example |
|---|---|---|
| **ADJ** | **adjective** | big, old, green, incomprehensible, first |
| ADP | adposition | in, to, during |
| **ADV** | **adverb** | very, tomorrow, down, where, there |
| AUX | auxiliary | is, has (done), will (do), should (do) |
| CONJ | conjunction | and, or, but |
| CCONJ | coordinating conjunction | and, or, but |
| DET | determiner | a, an, the |
| INTJ | interjection | psst, ouch, bravo, hello |
| **NOUN** | **noun** | girl, cat, tree, air, beauty |
| NUM | numeral | 1, 2017, one, seventy-seven, IV, MMXIV |

| POS | description | example |
|---|---|---|
| PART | particle | 's, not, |
| **PRON** | **pronoun** | I, you, he, she, myself, themselves, somebody |
| **PROPN** | **proper noun** | Mary, John, London, NATO, HBO |
| PUNCT | punctuation | ., (, ), ? |
| SCONJ | subordinating conjunction | if, while, that |
| SYM | symbol | $, %, §, ©, +, −, ×, ÷, =, :), 😝 |
| **VERB** | **verb** | run, runs, running, eat, ate, eating |
| X | other | sfpksdpsxmsa |
| SPACE | space |  |

spaCy

16

true date

superficial cleaning

deep cleaning

| title | created | score | upvote_ratio | selftext | title_sup_clean | title_deep_clean | title_deep_clean_pos |
|---|---|---|---|---|---|---|---|
| Damn...we blinked and missed the T-34 stage of... | 2022-10-27 | 10390 | 0.99 | NaN | damn we blinked and missed the t stage of the ... | damn blink miss t stage war | [damn ADV, blink VERB, miss VERB, t PROPN, sta... |
| Finnish 🇫🇮 volunteer sends greetings home from ... | 2023-03-31 | 2095 | 1.00 | NaN | finnish volunteer sends greetings home from so... | finnish volunteer send greeting home | [finnish ADJ, volunteer NOUN, send VERB, greet... |
| Guess having 5 trucks fall into your office ca... | 2022-09-16 | 1980 | 1.00 | NaN | guess having trucks fall into your office can ... | guess have truck fall office significant emoti... | [guess VERB, have VERB, truck NOUN, fall VERB,... |
| [META] Important - Russia-Ukraine Crisis/War: ... | 2022-02-24 | 1280 | 0.89 | Hi, /u/Anonim97 here.\n\nWe - as a mods of 40k... | important russia ukraine crisis war info and ... | important russia ukraine crisis war info way help | [important ADJ, russia PROPN, ukraine PROPN, c... |
| V*tniks coping hard Over the counter offensive... | 2022-09-12 | 1076 | 1.00 | NaN | v tniks coping hard over the counter offensive... | tnik cope hard counter offensive traitor pfp lmao | [tnik NOUN, cope VERB, hard ADJ, counter NOUN,... |
| ... | ... | ... | ... | ... | ... | ... | ... |

only ADJ, ADV, NOUN, PRON, PROPN, VERB kept

17

occurrence
(i.e., number of
times it appears in
the documents)

high-occurrence
words build the true
semantic network

might want to discard
words with very low
occurrence

word index

# Building the semantic network

bipartite and projected counterparts

number of occurrences
of words in documents

probability of words
given a documents

$$N_{wd} = \begin{array}{|cccc|}
\hline
0 & 1 & 1 & 1 \\
1 & 1 & 1 & 1 \\
1 & 0 & 1 & 0 \\
0 & 0 & 1 & 1 \\
1 & 0 & 1 & 1 \\
\hline
\end{array}$$

#globalwarming

#climatechange

#climateaction

#gretathunberg

#environment

$$P_{w|d} = \begin{array}{|cccc|}
\hline
0 & \frac{1}{2} & \frac{1}{5} & \frac{1}{4} \\
\frac{1}{3} & \frac{1}{2} & \frac{1}{5} & \frac{1}{4} \\
\frac{1}{3} & 0 & \frac{1}{5} & 0 \\
0 & 0 & \frac{1}{5} & \frac{1}{4} \\
\frac{1}{3} & 0 & \frac{1}{5} & \frac{1}{4} \\
\hline
\end{array}$$

we capture the statistical
properties by
normalizing by columns

we identify a
document
probability

$$p_d = \begin{cases} \dfrac{1}{D} & \text{equally likely} \\[2ex] \dfrac{n_d}{\sum_d n_d} & \text{custom} \end{cases}$$

**UNIVERSITÀ DEGLI STUDI DI PADOVA**

bipartite network

joint probability of words and documents

$$P_{wd} = P_{w|d} \, \text{diag}(p_d)$$

| | | | |
|---|---|---|---|
| 0 | $\frac{1}{8}$ | $\frac{1}{20}$ | $\frac{1}{16}$ |
| $\frac{1}{12}$ | $\frac{1}{8}$ | $\frac{1}{20}$ | $\frac{1}{16}$ |
| $\frac{1}{12}$ | 0 | $\frac{1}{20}$ | 0 |
| 0 | 0 | $\frac{1}{20}$ | $\frac{1}{16}$ |
| $\frac{1}{12}$ | 0 | $\frac{1}{20}$ | $\frac{1}{16}$ |

marginal probabilities

$$p_w = P_{wd} \, \mathbf{1} \qquad p_d = P_{wd}^{\mathsf{T}} \, \mathbf{1}$$

$$p_{w_1,w_2} = \sum_d p_{w_1|d, w_2} \, p_{w_2,d}$$

$$P_{ww} = P_{wd} \, \text{diag}(p_d)^{-1} \, P_{wd}^{\mathsf{T}}$$

$$p_w = P_{ww} \, \mathbf{1}$$

projection on words

projection on documents

$$p_{d_1,d_2} = \sum_w p_{d_1|w, d_2} \, p_{d_2,w}$$

$$P_{dd} = P_{wd}^{\mathsf{T}} \, \text{diag}(p_w)^{-1} \, P_{wd}$$

$$p_d = P_{dd} \, \mathbf{1}$$

21

bipartite network $P_{wd}$

Tweets

Hashtags

projection on hashtags $P_{ww}$

projection on tweets $P_{dd}$

term frequency =
frequency (probability) of
the word in the document

inverse document frequency =
(log) fraction of documents
that contain the word

$$\text{TF-IDF}_{w|d} = p_{w|d} \cdot -\log\left(\frac{\sum_d (nwd > 0)}{D}\right)$$

❑ An heuristic
❑ Punishes words that appear in many documents
❑ Enhances words that are document specific

$$T_{wd} = \text{TF-IDF}_{w|d} \, \text{diag}(\text{tf-idf}_d)^{-1} \cdot \text{diag}(p_d)$$

$T_{ww}$

$T_{dd}$

to be used in
place of $P_{wd}$

normalization vector $\text{tf-idf}_d = \text{TF-IDF}_{w|d}^T \, \mathbf{1}$
(to guarantee that columns sum up to 1)

23

# Topic detection

i.e., community detection in semantic networks

UNIVERSITÀ
DEGLI STUDI
DI PADOVA



bipartite network $P_{wd}$ or $T_{wd}$

a community identifies both
documents and words

a community
identifies words
need to use TopicRank to
identify documents

projection on words $P_{ww}$ or $T_{ww}$

projection on documents $P_{dd}$ or $T_{dd}$

a community
identifies documents
need to use
TopicRank to identify
words

25

UNIVERSITÀ
DEGLI STUDI
DI PADOVA



#metoo tweets

26

Projection on words
assigning documents to topics via TopicSpecific PageRank

Tweet 1 is assigned to Topic 1 !!!

Topic 1

those who think they are crazy enough to change the world eventually do. #climatechange #ClimateCrisis #ClimateAction #GretaThunberg #Greta

0.1234

4  #GretaThunberg

5  #climatechange

0.1221

Hopefully these kids will succeed where past generations have failed. #TheResistance #FBR #ClimateChange #Environment #GlobalWarming #GretaThunberg

6  #GlobalWarming

Topic 2

Tweets

7  #Environment

The #environment can have a major effect on the human cardiovascular system. A new study has found an increase in heat-induced #heartattack risk in recent years. Could #ClimateChange be a risk factor? #longevity

8  #longevity

Hashtags

27

statistical dependencies about words and topics

probability of a topic

$$\boldsymbol{P}_{\text{wt}} = \boldsymbol{P}_{wd}\,\boldsymbol{C}^{\mathsf{T}}$$

$$\boldsymbol{p}_{\text{t}} = \boldsymbol{P}_{w\text{t}}{}^{\mathsf{T}}\,\boldsymbol{1}$$

fraction of knowledge related to the topic that is explained by words (equal to 1 if topics use different words)

$$\text{NMI} = \frac{\text{I(W;T)}}{\text{H(T)}}$$

words
$W$

topics
$T$

**C** topic assignment to be assessed for quality

$$P_{tt} = C \, P_{dd} \, C^T$$

can be interpreted as a probability matrix linking topics, its entries are the sum of the links of **A** from topic i to topic j

| $P_{11}$ | $P_{12}$ | $P_{13}$ |
|---|---|---|
| $P_{21}$ | $P_{22}$ | $P_{23}$ |
| $P_{31}$ | $P_{32}$ | $P_{33}$ |

$$p_t = P_{tt} \, 1$$

can be interpreted as the probability vector of topics

modularity

$$Q = \sum_t (P_{tt} - p_t^2) < 1$$

to be maximized

normalized cut

normalized version

$$\text{Ncut} = 1 - \frac{\sum_t P_{tt}/pt}{\sum_t 1} > 0$$

to be minimized

PageRank vector (ranking of documents)

$$r = (1-c)\, P_{d|d}\, r + c\, 1/N$$

Here $c_i$ is the $i$th row of $C$

$$P_{d|d} = P_{dd}\, \text{diag}^{-1}(p_d)$$

$$q_i = \left(1 - (1-c)\frac{c_i 1}{N}\right) z_i 1 - c\, c_i P_{d|d} z_i^T$$

$$z_i = c_i\, \text{diag}(r)$$

$$\text{InfoMap} = f(q) + \sum_i f([q_i, z_i])$$

normalized version

$$\frac{\text{InfoMap}}{f(r)} - 1$$

to be minimized

entropy function

$$f(x) = -\sum_i x_i \log\left(\frac{x_i}{\sum_j x_j}\right)$$

31

InfoMap has either very big communities or not related to Louvain

❑ <u>Louvain Pdd</u> – provides the best results

   produces balanced clusters

❑ <u>Louvain soft</u> – slightly strengthens the result

❑ <u>Bipartite networks</u> – run much faster

   but performance deteriorates

❑ <u>InfoMap</u> – not robust ☹

   would be nice to see BigCLAM and SBMs
   … your task! ☺

# Non-negative Matrix Factorization

and its application to topic detection

$$P_{w|t}$$

$$P_{t|d} \simeq C$$

$$P_{w|d} \approx P_{w|t} \cdot P_{t|d}$$

the equivalence is
only approximate

underlying model

d → t → w

each document is
associated to a topic
distribution (one of the
colums of $C$)

each topic is
associated to a
word distribution
(one of the
columns of $P_{w|t}$)

$A = P_{w|d}$ is column stochastic

minimizing the Frobenius norm does not ensure a column stochastic product $W H$

$$\text{argmin}_{W \geq 0, H \geq 0} \sum_{ij} |A_{ij} - [WH]_{ij}|^2$$

$$\text{argmin}_{W \geq 0, H \geq 0} \sum_{ij} A_{ij} \log\left(\frac{A_{ij}}{[WH]_{ij}}\right) - A_{ij} + [WH]_{ij}$$

minimizing the generalized Kullback-Leibler divergence ensures a column stochastic product $W H$

$$f(y) = x \log\left(\frac{x}{y}\right) - x + y$$

$$f'(y) = -\frac{x}{y} + 1 = 0 \rightarrow y = x$$

Ho & Van Dooren. "Non-negative matrix factorization with fixed row and column sums." (2008)

39

```python
from sklearn.decomposition import NMF
Pwgd = Pwd/Pwd.sum(axis=0).flatten()
```

wisely initialize for best performance

run on different number of topics, then choose the best fit, e.g., according to modularity

choose generalized Kullback-Leibler divergence, and the related solver

```python
# fit nmf model X = W*H
model = NMF(n_components=i, init='nndsvd',
            solver='mu', beta_loss='kullback-leibler')
W = model.fit_transform(Pwgd)
H = sps.csr_matrix(model.components_)
# column normalized versions
H = sps.diags(W.sum(axis=0).flatten())*H # Ptgd
W = W/W.sum(axis=0).flatten() # Pwgt
# community assignment C
C = sps.csr_matrix(np.transpose(H/H.sum(axis=0).flatten()))
```

need to make W column stochastic, to have H column stochastic too

force column stochasticity in H (not needed though)

UNIVERSITÀ
DEGLI STUDI
DI PADOVA



42

❑ **Naturally provides a soft topic assignment**

❑ <u>NMF</u> – not strikingly good
   probably due to the fact that we want to express
   a sparse matrix through an eigenvector-like product
   with few eigenvectors (the fit is far from ideal)

❑  <u>Comparison</u> – with Louvain

   much weaker

❑ <u>Complexity</u> – generally slow
   need to test it for different numbers of topics ☹
   fast for fixed topic number

# Latent Dirichlet allocation

LDA = a stochastic model for topic detection

**UNIVERSITÀ DEGLI STUDI DI PADOVA**

topics vector
length $N_d$

words vector
length $N_d$

document

$d \longrightarrow t_d \longrightarrow w_d$

topics distribution
$p(t_{d,n} = j \mid \theta_d) = \theta_{d,j}$

words distribution
$p(w_{d,n} = i \mid \beta, t_{d,n} = j) = \beta_{i,j}$

topics probabilities
$\theta_d \sim$ Dirichlet($\alpha$)

$\theta_d$ is the topic
assignment $C = P_{t|d}$

K topics

the number of
topics K is fixed

vector $\alpha$, to be
estimated, is the
topic probability
$p_t = \alpha/|\alpha|_1$

matrix $\beta$ , to be
estimated, is the
probability of
words given a
topic, $P_{w|t} = \beta$

45

topics assignment probability (Dirichlet)

$$p(\boldsymbol{\theta}_d | \boldsymbol{\alpha}) = \frac{\Gamma\left(\sum_{k=1}^{K} \alpha_k\right)}{\sum_{k=1}^{K} \Gamma(\alpha_k)} \prod_{k=1}^{K} [\theta_{d,k}]^{\alpha_k - 1}$$

words probability

$$p(\boldsymbol{w}_d | \boldsymbol{\beta}, \boldsymbol{\theta}_d) = \prod_{n=1}^{N_d} [\boldsymbol{\beta} \, \boldsymbol{\theta}_d]_{w_{d,n}}$$

this dependence between $\boldsymbol{\beta}$ and $\boldsymbol{\theta}$ is the trickiest part

overall probability

$$p(\text{corpus} | \boldsymbol{\alpha}, \boldsymbol{\beta}) = \prod_d \int p(\boldsymbol{w}_d | \boldsymbol{\beta}, \boldsymbol{\theta}_d) \, p(\boldsymbol{\theta}_d | \boldsymbol{\alpha}) d\boldsymbol{\theta}_d$$

target optimization

$$\operatorname{argmax}_{\boldsymbol{\alpha}, \boldsymbol{\beta}} p(\text{corpus} | \boldsymbol{\alpha}, \boldsymbol{\beta})$$

$$\boldsymbol{C} = \boldsymbol{P}_{t|d} = \boldsymbol{\theta}$$

$$\boldsymbol{P}_{wt} = \boldsymbol{\beta} \, \text{diag}(\boldsymbol{\alpha}/|\boldsymbol{\alpha}|_1)$$

this is what we get

46

```python
from sklearn.decomposition import LatentDirichletAllocation

# fit lda model
lda = LatentDirichletAllocation(n_components=i,
                                learning_method="batch")
lda.fit(Mwd.T)
# community assignment C = Ptgd'
C = sps.csr_matrix(lda.transform(Mwd.T))
```

initialise and fit model

LDA

extract topic assignment

NMF and LDA identify
few communities hardly
related to Louvain

❑ **Naturally provides a soft topic assignment**

❑ <u>**LDA**</u> **– not strikingly good**
  same eigenvector-like product as NMF
  worse than NMF … known issue ☹
  probably due to the Dirichlet assumption (questionable)
  and the variational inference (suboptimum approach)

❑ <u>**Comparison**</u> **– with Louvain**

  much weaker

❑ <u>**Complexity**</u> **– generally slow**
  need to test it for different numbers of topics ☹
  fast for fixed topic number

# Variational Auto Encoders

an application to topic analysis

the hidden prior in our case
is the topic (distribution)

the output in our case
is the document
(collection of words)

encoder

decoder

**d**

**z**

**z**

**d**

$p(\mathbf{z}|\mathbf{d})$

$p(\mathbf{d}|\mathbf{z})$

hidden prior
*embedding*

inference model
*posterior probabilities
to be estimated*

generative model
*given*

the (stochastic) model
explains how a
document is generated
from a topic (distribution)

but we are interested in the inverse link that, given
a document tells what topic it is associated with

$$p(\mathbf{z}|\mathbf{d}) = \frac{p(\mathbf{d}|\mathbf{z})\, p(\mathbf{z})}{p(\mathbf{d})} \cong q(\mathbf{z}|\mathbf{d})$$

impossible to know
in the closed form

needs an a-priori
model for the
embedding

is approximated
by a simple
alternative model

53

ELBO

$$\mathcal{L}_{\theta,\phi}(\boldsymbol{d}) \leq \log p_\theta(\boldsymbol{d})$$

encoder

$$\boldsymbol{d} \rightarrow q(\boldsymbol{z}|\boldsymbol{d}) \rightarrow \boldsymbol{z}$$

inference model
*estimate of p(z|d)*

hidden prior
*embedding p(z)*

decoder

$$\boldsymbol{z} \rightarrow p(\boldsymbol{d}|\boldsymbol{z}) \rightarrow \boldsymbol{d}$$

generative model
*given*

$$\mathcal{L}_{\theta,\phi}(\boldsymbol{d}) = \log p_\theta(\boldsymbol{d}) - D_{\mathrm{KL}}\left(q_\phi(\boldsymbol{z}|\boldsymbol{d})\middle\|p_\theta(\boldsymbol{z}|\boldsymbol{d})\right)$$

$$= \int dz\, q_\phi(\boldsymbol{z}|\boldsymbol{d}) \log\left(\frac{p_\theta(\boldsymbol{z},\boldsymbol{d})}{q_\phi(\boldsymbol{z}|\boldsymbol{d})}\right)$$

$$= \underbrace{\int dz\, q_\phi(\boldsymbol{z}|\boldsymbol{d}) \log\left(p_\theta(\boldsymbol{d}|\boldsymbol{z})\right)}_{\mathcal{L}_1} + \underbrace{\int dz\, q_\phi(\boldsymbol{z}|\boldsymbol{d}) \log\left(\frac{p_\theta(\boldsymbol{z})}{q_\phi(\boldsymbol{z}|\boldsymbol{d})}\right)}_{\mathcal{L}_2}$$

to be maximized wrt parameters $\theta$ and $\phi$ provides fitting on p(**z**), p(**d**|**z**), and q(**z**|**d**)

a-priori model (given)

inference model (approximate)

generative model (given)

54

$$\underbrace{\int dz\, q_\phi(z|d) \log\left(\frac{p_\theta(z)}{q_\phi(z|d)}\right)}_{\mathcal{L}_2}$$

both should have a simple parametrization on $\theta$ and $\phi$

e.g., the Gaussian case

$$p_\theta(z) = \frac{1}{\sqrt{\det\left(2\pi\, \mathrm{diag}(\boldsymbol{\sigma}_\theta^2)\right)}} \exp\left(-\tfrac{1}{2}(z - \boldsymbol{\mu}_\theta)^T \mathrm{diag}^{-1}(\boldsymbol{\sigma}_\theta^2)(z - \boldsymbol{\mu}_\theta)\right)$$

$$q_\phi(z|d) = \frac{1}{\sqrt{\det\left(2\pi\, \mathrm{diag}(\boldsymbol{\sigma}_\phi^2(d))\right)}} \exp\left(-\tfrac{1}{2}(z - \boldsymbol{\mu}_\phi(d))^T \mathrm{diag}^{-1}(\boldsymbol{\sigma}_\phi^2(d))(z - \boldsymbol{\mu}_\phi(d))\right)$$

$$\mathcal{L}_2(\theta, \phi) = \tfrac{1}{2} \sum_i 1 + \log\left(\frac{\sigma_{\phi,i}^2(d)}{\sigma_{\theta,i}^2}\right) - \frac{\sigma_{\phi,i}^2(d)}{\sigma_{\theta,i}^2} - \frac{\left(\mu_{\phi,i}(d) - \mu_{\theta,i}\right)^2}{\sigma_{\theta,i}^2}$$

55

$$\underbrace{\int dz \, q_\phi(z|d) \log\left(p_\theta(d|z)\right)}_{\mathcal{L}_1}$$

mostly too complex to be written in the closed form

solution: Monte Carlo approximation

$$\mathcal{L}_1(\theta, \phi) = \frac{1}{L}\sum_{\ell=1}^{L} \log\left(p_\theta(d|z_\ell)\right)$$

samples generated according to the correct distribution
$$z_\ell \sim q_\phi(z|d)$$

e.g., the Gaussian case
$$z_\ell = \mu_\phi(d) + \sigma_\phi(d)\, n_\ell$$

need to generate these once, then use them throughout the process

$$n_\ell \sim \mathcal{N}(0, I)$$

56

**UNIVERSITÀ DEGLI STUDI DI PADOVA**

here **W** builds a topic-to-word map that identifies the word distribution inside each document given the topic distribution

*word probabilities*

somehow related to a topic, but we do not know how

$$\text{softmax}(\mathbf{W}\mathbf{z}+\mathbf{b})$$

$\mathbf{z}$ → softmax(**Wz**+**b**) → $\mathbf{p}_w$ → $\mathbf{d}$

***hidden prior***
$p(\mathbf{z})$ **Gaussian** distributed
$\mathbf{z} = \mu_1 + \sigma_1\, \mathbf{n}$
with $\mathbf{n}$ normalized Gaussian

***document***
$p_\theta(\mathbf{d}|\mathbf{z})$ **multinomial**
distributed according to $\mathbf{p}_w(\mathbf{z})$

probability function for the multinomial
$\log(p_\theta(\mathbf{d}|\mathbf{z}))$
$= \mathbf{d}^T \log(\mathbf{p}_w(\mathbf{z})) + \text{const}$

*$\mathbf{d}$ one-hot-representation = number of occurrences of words in the document*

*different document regions are associated with different output values*

$\mathbf{d}$ → FFN → $\mu_0$ , $\log \sigma_0^2$ → $\mathbf{z}$

***document***
*one-hot-representation of appearing words*

*e.g., 2x 2-layer, ReLu activation, **linear** output*

***hidden prior***
$q_\phi(\mathbf{z}|\mathbf{d})$ **Gaussian** distributed
$\mathbf{z} = \mu_0 + \sigma_0\, \mathbf{n}$
with $\mathbf{n}$ normalized Gaussian

57

*one-hot-representation of a document = number of occurrences of words in the document*

normalized Gaussian samples

**decoder model**  **encoder model**

$$\mathcal{L}(\theta, \phi) = \frac{1}{L} \sum_{\ell=1}^{L} \sum_{m} \boldsymbol{d}_m^T \log \left( \underbrace{\text{softmax}\left( \boldsymbol{b} + \boldsymbol{W}}_{\text{decoder map}} (\boldsymbol{\mu}_0(\boldsymbol{d}_m) + \boldsymbol{\sigma}_0(\boldsymbol{d}_m)\, \boldsymbol{n}_{m,\ell})) \right)$$

$$+ \tfrac{1}{2} \sum_{m} \sum_{i} 1 + \log \left( \frac{\sigma_{0,i}^2(\boldsymbol{d}_m)}{\sigma_{1,i}^2} \right) - \frac{\sigma_{0,i}^2(\boldsymbol{d}_m)}{\sigma_{1,i}^2} - \frac{\left( \mu_{0,i}(\boldsymbol{d}_m) - \mu_{1,i} \right)^2}{\sigma_{1,i}^2}$$

**a-priori model**

Not very clear where the topic is, though!

59

UNIVERSITÀ
DEGLI STUDI
DI PADOVA



**a multigamma distribution is introduced**

$z$/sum($z$)

**true Dirichlet distribution (as in LDA)**

$z$ → softmax(log()) → $t$ → softmax($\mathbf{W}t+b$) → $\mathbf{p}_W$ → $d$

*hidden prior*
p($z$) *multigamma*
*distributed with parameters $\alpha_1$ and $\beta=1$*

*topic*
**Dirichlet**
*distributed*

*document*
$p_\theta(d|z)$ *multinomial*
*distributed*

$d$ → FFN → $\alpha_0>0$ → $z$

**shape vector**

*document*

*e.g., 2x 2-layer, ReLu activation,*
**softplus** $ln(1+e^x)$ *output*

*hidden prior*
$q_\phi(z|d)$ *multigamma*
*distributed with parameters $\alpha_0$ and $\beta=1$*

61

$$f(\boldsymbol{u}, \boldsymbol{\alpha}) = (\boldsymbol{u}\,\boldsymbol{\alpha}\,\Gamma(\boldsymbol{\alpha}_1))^{1/\alpha}$$

approx. uniform to multigamma map

normalized uniform samples

*one-hot-representation of a document = number of occurrences of words in the document*

**decoder model**

**encoder model**

$$\mathcal{L}(\theta, \phi) = \frac{1}{L} \sum_{\ell=1}^{L} \sum_{m} \boldsymbol{d}_m^T \log\left(\mathrm{softmax}\left(\boldsymbol{b} + \boldsymbol{W}\,\mathrm{softmaxlog}\left(\boldsymbol{f}(\boldsymbol{u}_{m,\ell}, \boldsymbol{\alpha}_0(\boldsymbol{d}_m))\right)\right)\right)$$

decoder map

$$+ \sum_{m} \sum_{i} \log\left(\frac{\Gamma(\alpha_{0,i}(\boldsymbol{d}_m))}{\Gamma(\alpha_{1,i})}\right) - \left(\alpha_{0,i}(\boldsymbol{d}_m) - \alpha_{1,i}\right)\psi(\alpha_{0,i}(\boldsymbol{d}_m))$$

**a-priori model**

digamma function

$$\psi(x) = \frac{\Gamma'(x)}{\Gamma(x)}$$

Now we know where the topic is!

62

UNIVERSITÀ
DEGLI STUDI
DI PADOVA

*each topic is a region in z … this also differentiates topic probabilities*

**z**

*hidden prior*
*p(**z**) normalized*
*Gaussian*
*distributed $N(\mathbf{0}_T, \mathbf{I}_T)$*

**FFN**

*e.g., 4-layer,*
*tanh activation,*
***softmax** output*

**t**

*topic*
*freely*
*distributed*

softmax(**Wt**+**b**)

**p**w

**d**

*document*
*$p_\theta(\mathbf{d}|\mathbf{z})$ multinomial*
*distributed*

**d**

*document*

**FFN**

*e.g., 2x 2-layer,*
*ReLu activation,*
***linear** output*

$\mu$

log $\boldsymbol{\sigma}^2$

**z**

*hidden prior*
*$q_\phi(\mathbf{z}|\mathbf{d})$ Gaussian*
*distributed*
*$\mathbf{z} = \mu + \boldsymbol{\sigma}\,\boldsymbol{n}$*

63

*one-hot-representation of a document = number of occurrences of words in the document*

normalized Gaussian samples

**decoder model**

**encoder model**

$$\mathcal{L}(\theta, \phi) = \frac{1}{L} \sum_{\ell=1}^{L} \sum_{m} \boldsymbol{d}_m^T \log\left(\underbrace{\mathrm{softmax}\left(\boldsymbol{b} + \boldsymbol{W}\,\mathrm{FFN}_1(\boldsymbol{\mu}_0(\boldsymbol{d}_m) + \boldsymbol{\sigma}_0(\boldsymbol{d}_m)\,\boldsymbol{n}_{m,\ell})\right)}_{\text{decoder map}}\right)$$

$$+ \tfrac{1}{2} \sum_{m} \sum_{i} 1 + \log\left(\sigma_{0,i}^2(\boldsymbol{d}_m)\right) - \sigma_{0,i}^2(\boldsymbol{d}_m) - \left(\mu_{0,i}(\boldsymbol{d}_m)\right)^2$$

Our estimate of the topic distribution for the *m*th document!

$$\boldsymbol{c}_m = \frac{1}{L} \sum_{\ell=1}^{L} \mathrm{FFN}_1(\boldsymbol{\mu}_0(\boldsymbol{d}_m) + \boldsymbol{\sigma}_0(\boldsymbol{d}_m)\,\boldsymbol{n}_{m,\ell})$$

64

each topic is a region in *z* … this also differentiates topic probabilities

**d**

*document*

**FFN**

*e.g., 4-layer, ReLu activation, **softmax** output*

**t**

*topic*
***freely***
*distributed*

softmax(**Wt+b**)

$p_w$

**d**

*document*

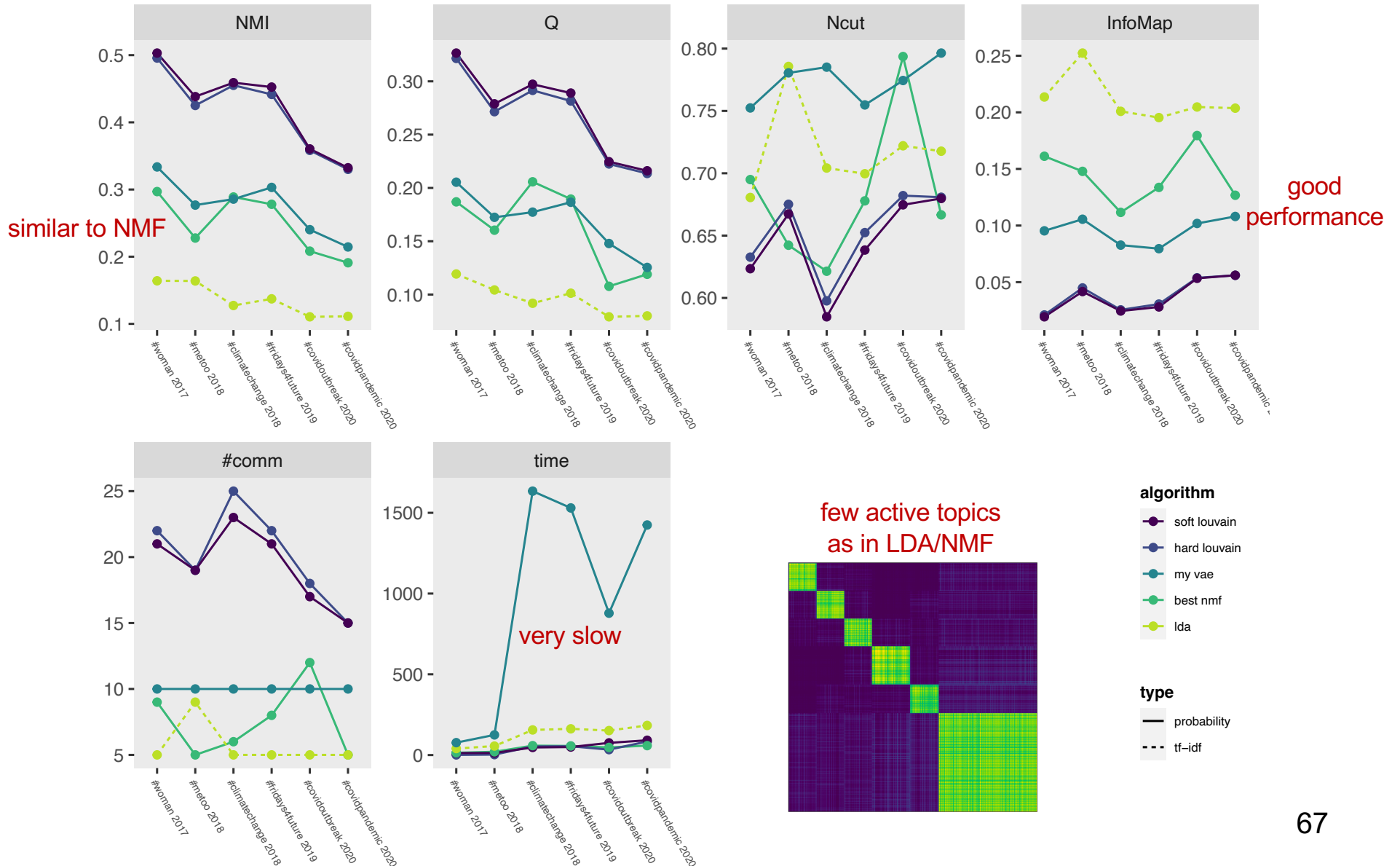$p_\theta$(**d|z**) ***multinomial*** *distributed*

*one-hot-representation of a document = number of occurrences of words in the document*

**decoder model**     **encoder model**

$$\mathcal{L}(\theta, \phi) = \sum_m \boldsymbol{d}_m^T \log\big(\mathrm{softmax}(\boldsymbol{b} + \boldsymbol{W}\,\mathrm{FFN}(\boldsymbol{d}_m))\big)$$

A comparison
with NMF, LDA, and Louvain

67

❑ Naturally provides a soft topic assignment

❑ <u>VAE</u> – interesting approach
   more flexible model than NMF or LDA
   gives improvements

❑ <u>Comparison</u> – with Louvain
   still far away
   would be nice to see other Deep Learning approaches
                                              … your task! ☺

# Transformer Architecture

with application to BERT, RoBERTa, OpenAI GPT

# ≔ Attention (machine learning)

Article    Talk

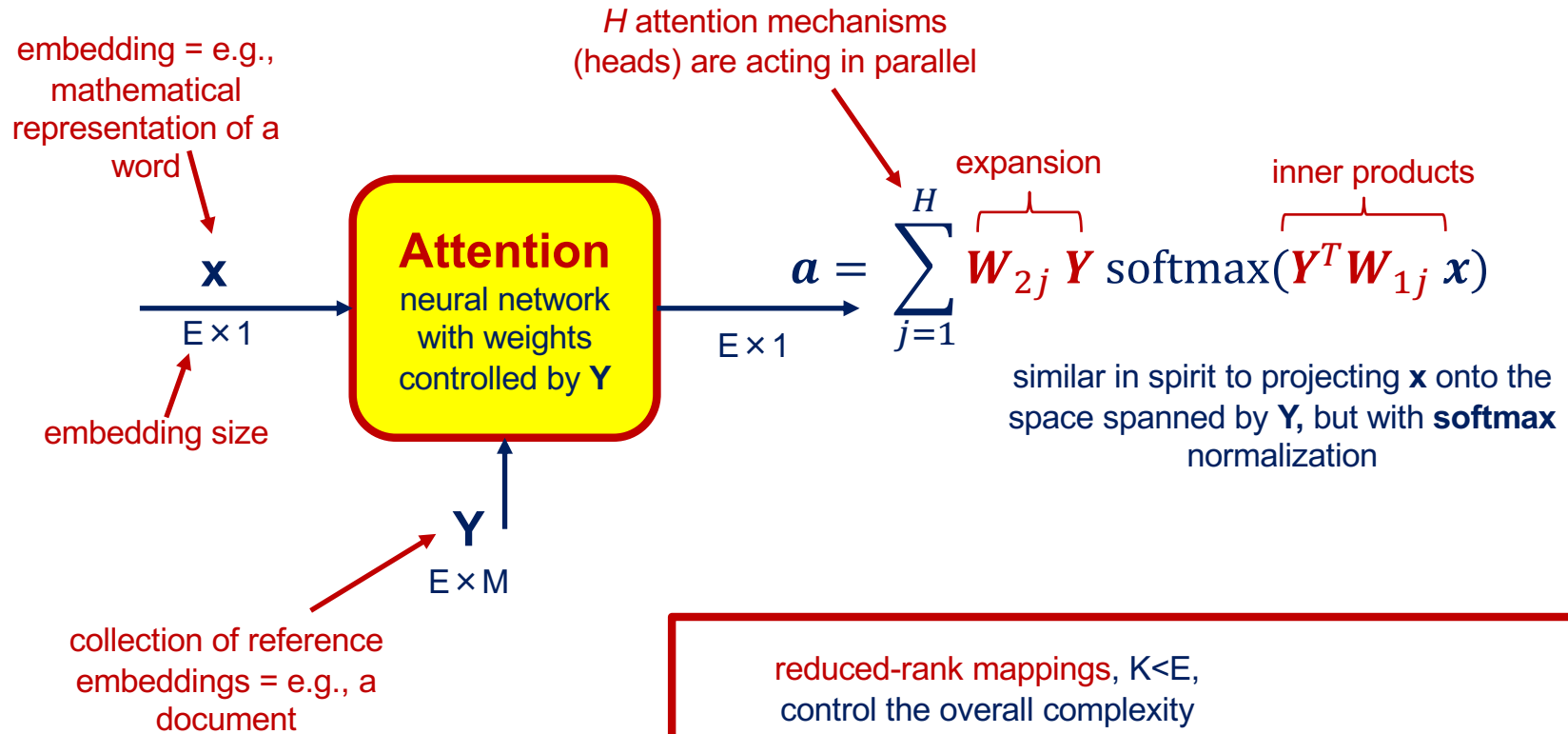From Wikipedia, the free encyclopedia

In artificial neural networks, **attention** is a technique that is meant to mimic cognitive attention. This effect enhances some parts of the input data while diminishing other parts — the motivation being that the network should devote more focus to the important parts of the data, even though they may be small portion of an image or sentence. Learning which part of the data is more important than another depends on the context, and this is trained by gradient descent.

embedding = e.g., mathematical representation of a word

embedding size

**Attention**
neural network with weights controlled by **Y**

$\mathbf{x}$
$E \times 1$

$E \times 1$

**Y**
$E \times M$

collection of reference embeddings = e.g., a document

*H* attention mechanisms (heads) are acting in parallel

expansion

inner products

$$a = \sum_{j=1}^{H} W_{2j}\, Y\, \mathrm{softmax}(Y^T W_{1j}\, x)$$

similar in spirit to projecting **x** onto the space spanned by **Y**, but with **softmax** normalization

reduced-rank mappings, K<E, control the overall complexity

$W_{ij}$ = $V_{ij}$ $U_{ij}$

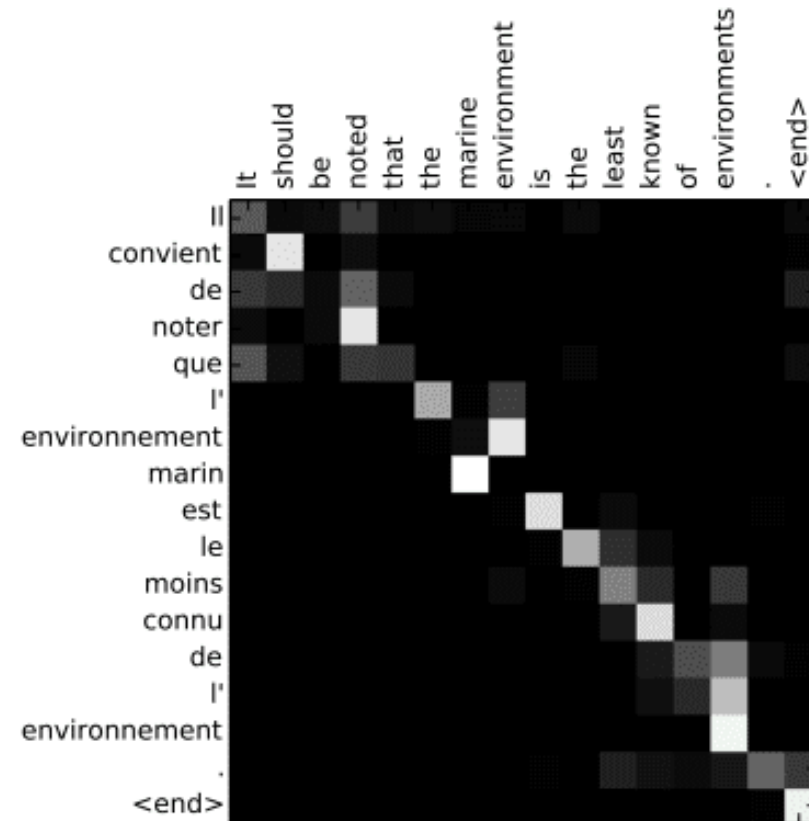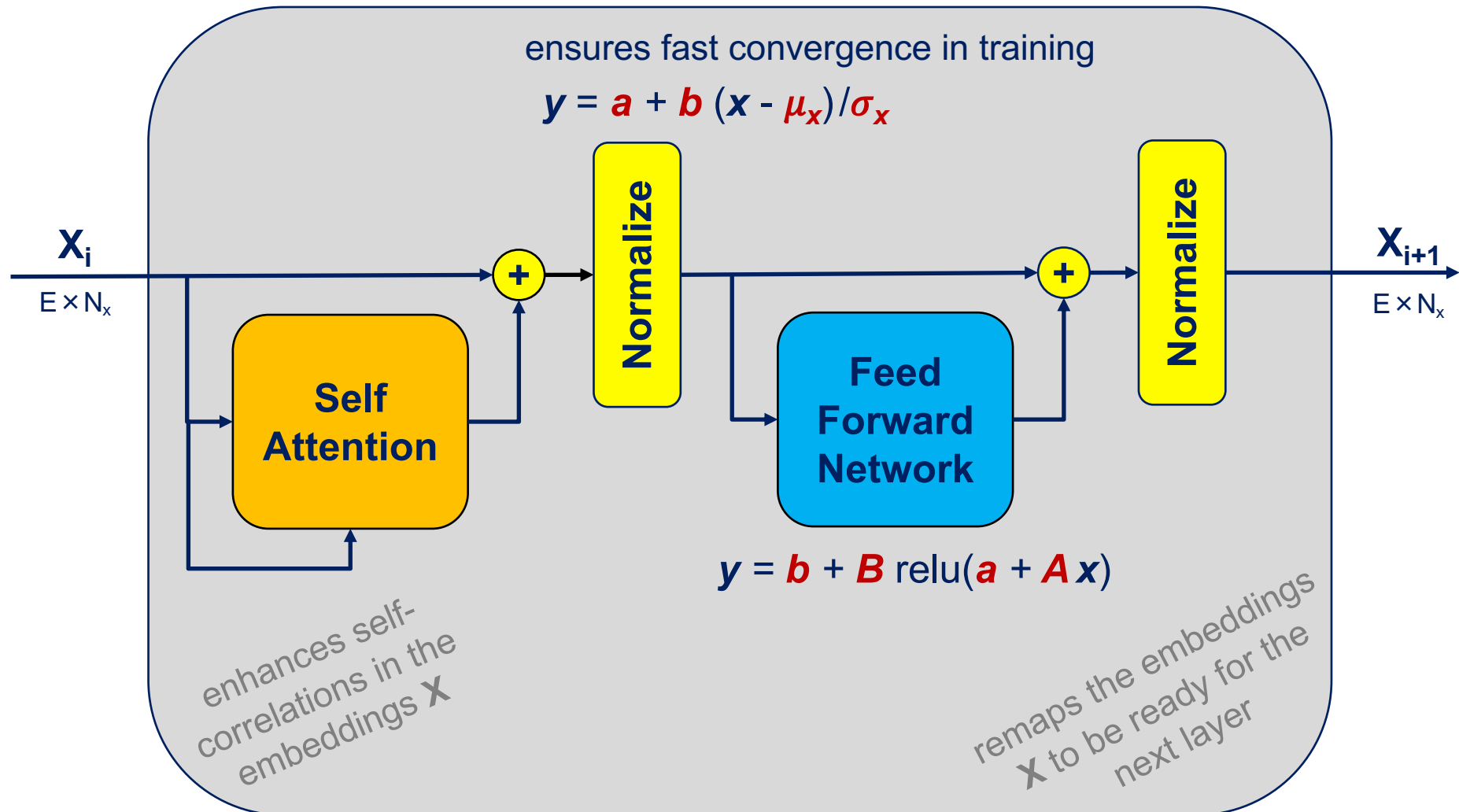$E \times E$        $E \times K$     $K \times E$

overall complexity 4EKH

71

$$\text{softmax}(Y^T W_{1j} X)$$

ensures fast convergence in training

$$y = a + b\,(x - \mu_x)/\sigma_x$$

$X_i$

$E \times N_x$

**Self Attention**

**Normalize**

**Feed Forward Network**

**Normalize**

$X_{i+1}$

$E \times N_x$

$$y = b + B\,\text{relu}(a + A\,x)$$

enhances self-correlations in the embeddings $X$

remaps the embeddings $X$ to be ready for the next layer

73

# Transformer Architecture

Vaswani, Ashish, et al. "Attention is all you need" (2017)
Google's patent https://patents.google.com/patent/US10452978B2/en

estimate $\hat{Y}$,
given the alternative representation $X$,
and the past information $Y_s$
(i.e., $Y$ shifted right by one position)

word-to-embedding map

Decoder

output embeddings $E \times N_y$

$V$
$E \times D$

$V^T$

$Y_s$
$D \times N_y$

Linear

$Y_0$

$P_0$

$Y_1$

$Y_{L-1}$

$Y_L$

Linear

softmax

$\hat{Y}$
$D \times N_y$

embedding-to-word
(probability) map

$X$
$D \times N_x$

Linear

$X_0$

$P_0$

positional
encoding

$X_1$

$X_{L-1}$

$X_L$

input embeddings $E \times N_x$

Encoder

75

UNIVERSITÀ
DEGLI STUDI
DI PADOVA

harvardnlp

Members    PI    Code    Publications

# The Annotated Transformer
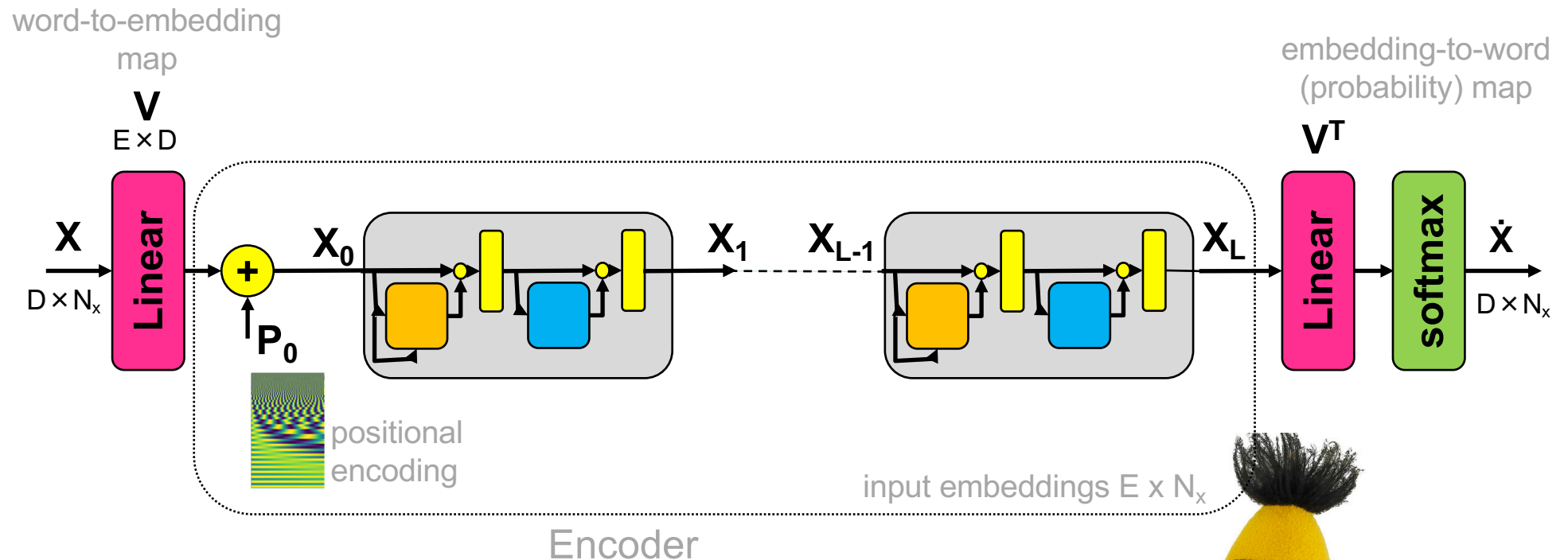
Apr 3, 2018

— — — — — — — — -

**There is now a new version of this blog post updated for modern PyTorch.**

— — — — — — — — -

```
from IPython.display import Image
Image(filename='images/aiayn.png')
```

## Attention Is All You Need

# BERT

Devlin, Jacob, et al. "Bert: Pre-training of deep bidirectional transformers for language understanding" (2018)

UNIVERSITÀ
DEGLI STUDI
DI PADOVA

https://github.com/google-research/bert

word-to-embedding map

$V$

$E \times D$

embedding-to-word (probability) map

$V^T$

$X$

$D \times N_x$

Linear

$+$

$P_0$

$X_0$

$X_1$

$X_{L-1}$

$X_L$

Linear

softmax

$\dot{X}$

$D \times N_x$

positional encoding

input embeddings $E \times N_x$

Encoder

| | Embeddings size E | Self-attention heads H | Head dimension K = E/H | FFN inner size I = 4E | Parameters per layer $12E^2+9E$ | Layers L | Dictionary size D | Total parameters |
|---|---|---|---|---|---|---|---|---|
| BERT base | 768 | 12 | 64 | 3072 | 7.1M | 12 | 30.5K | 110M |
| BERT large | 1024 | 16 | 64 | 4096 | 12.6M | 24 | 30.5K | 340M |

max tokens $N_x = 512$

Created by researchers at Google AI Language

78

**Masked Language Model**

15% masked tokens replaced with:
- [MASK] token (80% of the times)
- Original token (10%)
- Random token (10%)

**Next Sequence Prediction**
- Next sequence (50% of the times)
- Random sequence (50%)



Output [CLS] fed into an additional output layer for softmax classification (of correct/wrong next sequence)

Output masked tokens fed into the output layer $V^T$ and evaluated for probability of correct estimate

79

# RoBERTa

Liu, Yinhan, et al. "Roberta: A robustly optimized BERT pretraining approach" (2019)

### Larger training corpora (10x larger)
*training on BookCorpus + Wikipedia and also CC-News, OpenWebText, Stories*

### Dynamic masking
*training data was duplicated 10 times so that each sequence is masked in 10 different ways over the 40 epochs of training*

### Full-sentences without NSP loss
*full sentences sampled contiguously from one or more documents, such that the total length is at most 512 tokens*

### Large mini-batches

### A larger byte-level BPE (byte pair encoding) of 50K subword units
*a hybrid between character- and word-level representations that allows handling the large vocabularies common in natural language corpora*

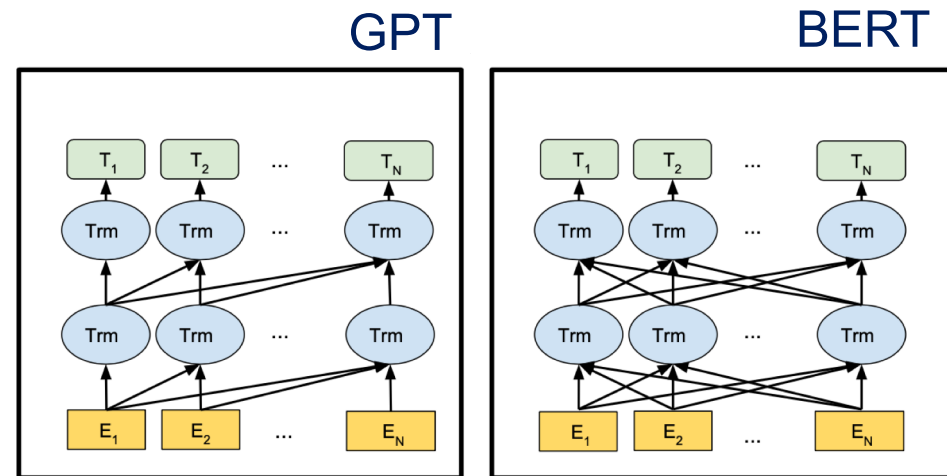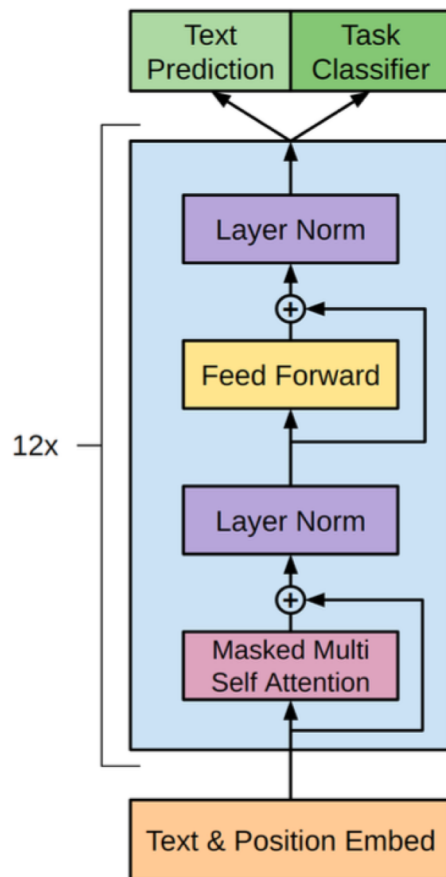Optimized by University of Washington & Facebook AI    80

# Generative Pre-Training (GPT)

Radford, Alec, et al. "Improving language understanding by generative pre-training." (2018)

(unsupervised) pre-training on Language Modelling (no mask)

$$L_1(\mathcal{U}) = \sum_i \log P(u_i | u_{i-k}, \ldots, u_{i-1}; \Theta)$$



same parameters of BERT-base, but with Masked Attention
trained on BookCorpus only

# GPT-2

Radford, Alec, et al. "Language models are unsupervised multitask learners" (2019)

## McCann et al. (2018)

language provides a flexible way to specify tasks, inputs, and outputs all as a sequence of symbols… it is therfore possible to **train a single model** with **sufficient capacity** to infer and perform many **different tasks**

model gets complex!

data gets larger!

## WebText

| Parameters | Layers | $d_{model}$ | |
|---|---|---|---|
| 117M | 12 | 768 | GPT, BERT-base |
| 345M | 24 | 1024 | BERT-large |
| 762M | 36 | 1280 | |
| 1542M | 48 | 1600 | GPT-2 |

scraping all outbound links (45M links) from Reddit, a social media platform, which received at least 3 karma – exclude WikiPedia

# GPT-3

Brown, Tom, et al. "Language models are few-shot learners" (2020)

increasingly larger data and model!

| Model Name | $n_{\text{params}}$ | $n_{\text{layers}}$ | $d_{\text{model}}$ | $n_{\text{heads}}$ | $d_{\text{head}}$ | Batch Size | Learning Rate |
|---|---|---|---|---|---|---|---|
| GPT-3 Small | 125M | 12 | 768 | 12 | 64 | 0.5M | $6.0 \times 10^{-4}$ |
| GPT-3 Medium | 350M | 24 | 1024 | 16 | 64 | 0.5M | $3.0 \times 10^{-4}$ |
| GPT-3 Large | 760M | 24 | 1536 | 16 | 96 | 0.5M | $2.5 \times 10^{-4}$ |
| GPT-3 XL | 1.3B | 24 | 2048 | 24 | 128 | 1M | $2.0 \times 10^{-4}$ |
| GPT-3 2.7B | 2.7B | 32 | 2560 | 32 | 80 | 1M | $1.6 \times 10^{-4}$ |
| GPT-3 6.7B | 6.7B | 32 | 4096 | 32 | 128 | 2M | $1.2 \times 10^{-4}$ |
| GPT-3 13B | 13.0B | 40 | 5140 | 40 | 128 | 2M | $1.0 \times 10^{-4}$ |
| GPT-3 175B or "GPT-3" | 175.0B | 96 | 12288 | 96 | 128 | 3.2M | $0.6 \times 10^{-4}$ |

**Layer normalization at the input** (plus one at the output)

**Sparse attention patterns**
*alternating dense and locally banded sparse attention patterns in the layers*
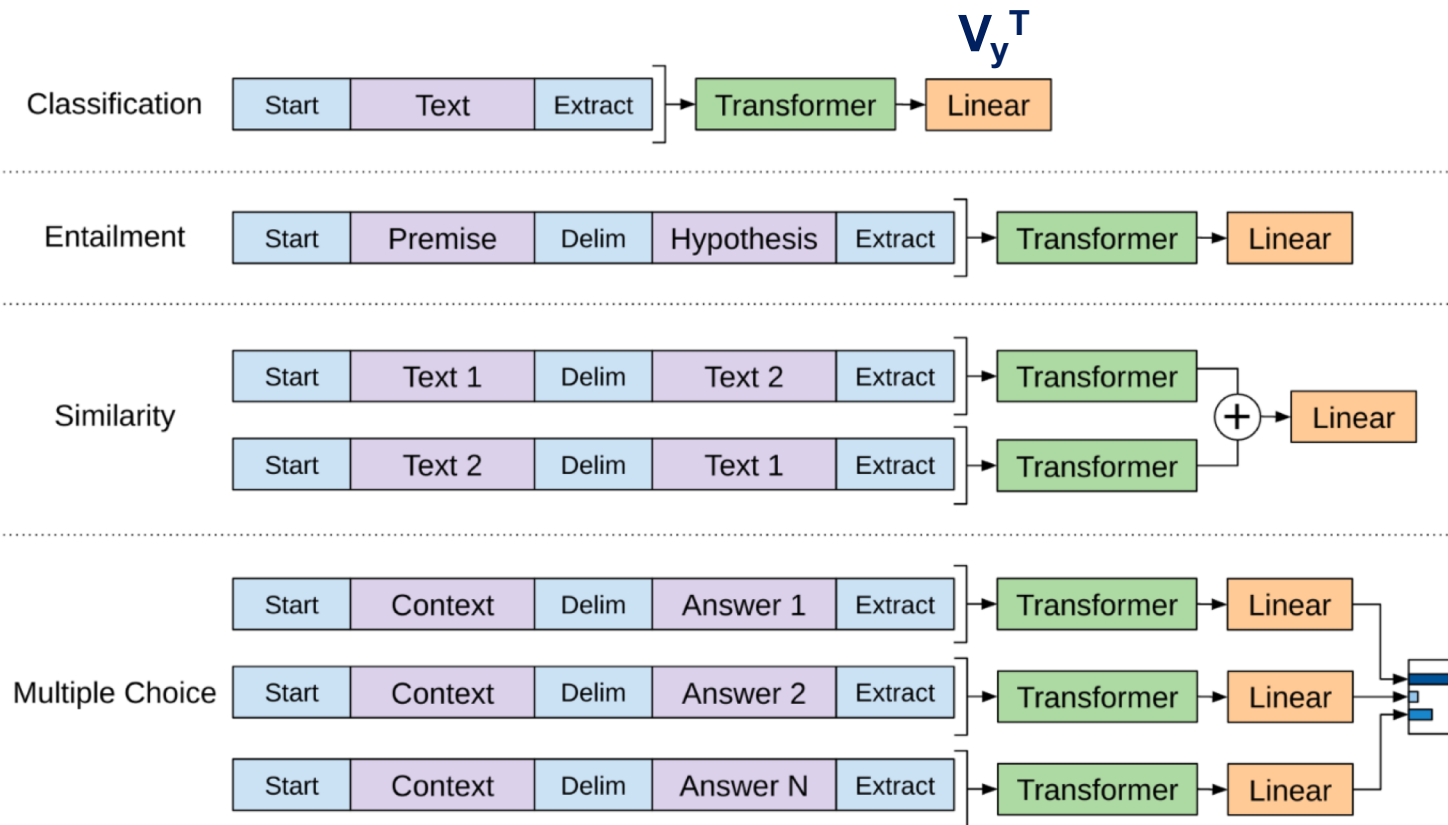
**Byte-level BPE (byte pair encoding)** of 50K subword units
*also prevent BPE from merging across character categories (to avoid dog, dog!, dog?)*

**Modified initialization**

Created by researchers at OpenAI   83

$$L_2(\mathcal{C}) + \lambda * L_1(\mathcal{C})$$

log softmax($\mathbf{V_y^T X_L}$)

Language Modelling loss

$\mathbf{V_y^T}$



84

| Task | Description | Possible approach |
|---|---|---|
| Masked language prediction | predict masked words in a text | This is what BERT model is pre-trained for |
| Text classification or Sentiment analysis | assign a label to a given sequence of text | Apply linear transform+softmax on K classes, and train the model for the specific classification task |
| Text translation | translate a text | Need to pre-train a full Transfomer Architecture for this task |
| Summarization | generate a summary of a document | GPT example: context given by a document; then generate 100 tokens by top-2 random sampling (Fan et al., 2018), i.e., take at each step the most likely next word at random among the top-2 candidates; finally select first 3 sentences as abstract |
| Question answering | answer a question | GPT example: the context of the language model is seeded with example question answer pairs which helps the model infer the short answer style of the dataset |
| Document question answering | answer a question on a given text | GPT example: context seeded by a text; then as for question answering |
| Conversational | ChatBot | InstructGPT/ChatGPT: Fine-tuned models using reinforcement learning from human feedback |

🤗 Hugging Face

https://huggingface.co/docs/transformers/v4.29.1/en/index

State-of-the-art Machine Learning
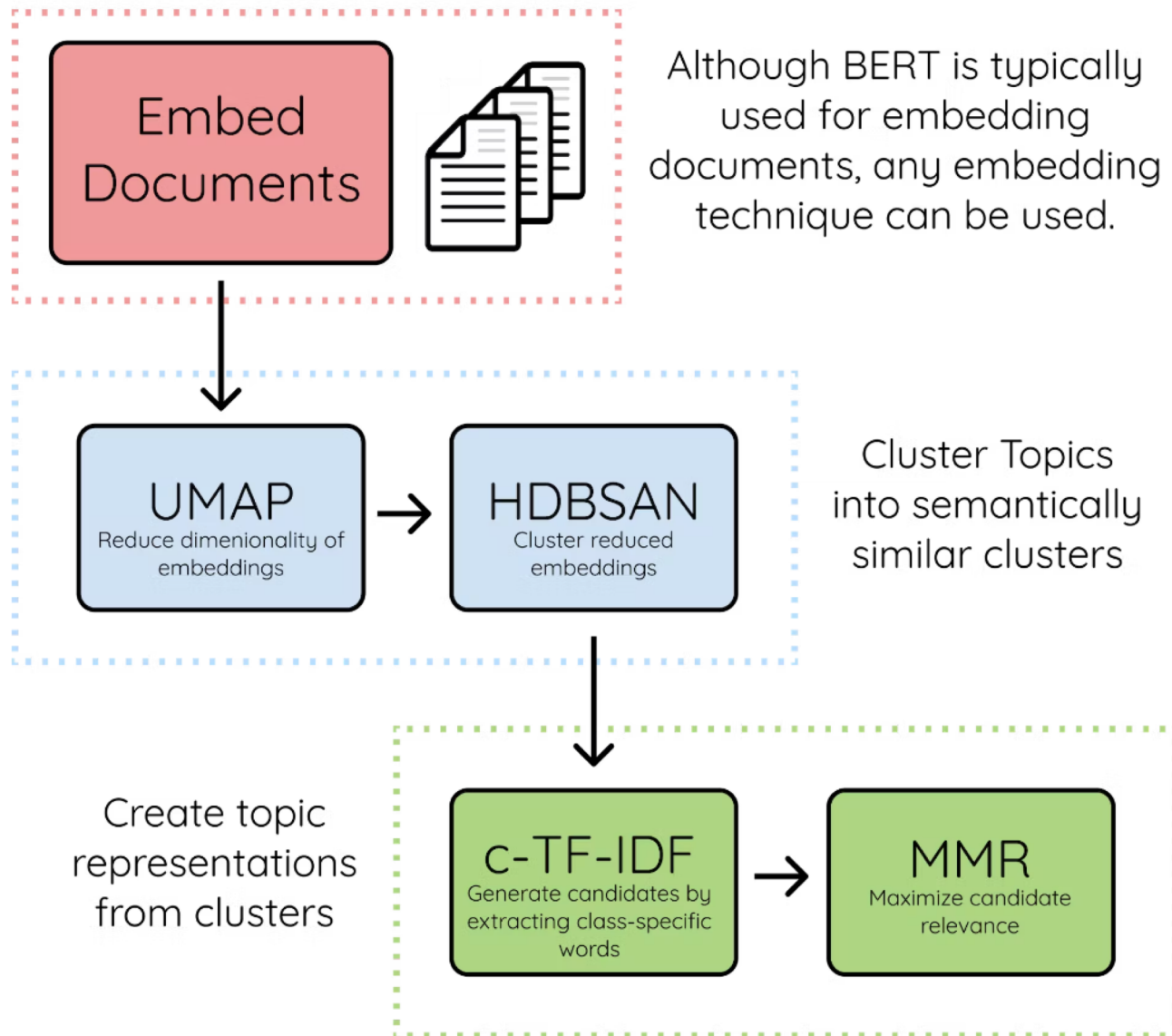for PyTorch, TensorFlow, and JAX

○ PyTorch      ⬆ TensorFlow

ALBERT, BART, **BERT**, BigBird, BigBird-Pegasus, BioGpt, BLOOM, CamemBERT, CANINE,
ConvBERT, CTRL, Data2VecText, DeBERTa, DeBERTa-v2, DistilBERT, ELECTRA, ERNIE,
ErnieM, ESM, FlauBERT, FNet, Funnel Transformer, GPT-Sw3, **OpenAI GPT-2**, GPTBigCode,
GPT Neo, GPT NeoX, GPT-J, I-BERT, LayoutLM, LayoutLMv2, LayoutLMv3, LED, LiLT, LLaMA,
Longformer, LUKE, MarkupLM, mBART, MEGA, Megatron-BERT, MobileBERT, MPNet, MVP,
Nezha, Nyströmformer, OpenLlama, **OpenAI GPT**, OPT, Perceiver, PLBart, QDQBert, Reformer,
RemBERT, **RoBERTa**, RoBERTa-PreLayerNorm, RoCBert, RoFormer, SqueezeBERT, TAPAS,
Transformer-XL, XLM, XLM-RoBERTa, XLM-RoBERTa-XL, XLNet, X-MOD, YOSO

# BERT Topic

exploiting embeddings for topic detection

87

**Embed Documents**

Although BERT is typically used for embedding documents, any embedding technique can be used.

**UMAP**
Reduce dimenionality of embeddings

**HDBSAN**
Cluster reduced embeddings

Cluster Topics into semantically similar clusters

Create topic representations from clusters

**c-TF-IDF**
Generate candidates by extracting class-specific words

**MMR**
Maximize candidate relevance

88

UNIVERSITÀ
DEGLI STUDI
DI PADOVA

```python
!pip install bertopic
from bertopic import BERTopic
from sentence_transformers import SentenceTransformer
```

initialise model

```python
sentence_model = SentenceTransformer("all-MiniLM-L6-v2")
bert_model = BERTopic(embedding_model=sentence_model,
                      min_topic_size=20,nr_topics='auto')
```

```python
docs = list(df2["text_sup_clean"])
topics, probabilities = bert_model.fit_transform(docs)
```

fit model

```python
topics = bert_model.reduce_outliers(docs, topics)
```

reduce outliers

```python
# extract community assignments
C = sps.csr_matrix((len(topics),max(topics)+2))
for i in range(C.shape[1]):
  C[np.array(topics)==(i-1),i] = 1

# remove zero assignments
C = C[:,np.unique(scipy.sparse.find(C)[1])]
```
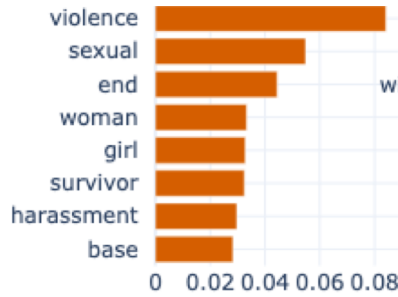
extract C from topic assignment

89

UNIVERSITÀ
DEGLI STUDI
DI PADOVA



similar to Louvain

poor performance

poor performance

poorer performance

**algorithm**
- bertopic reduce
- bertopic
- soft louvain
- hard louvain
- best nmf

**type**
- — probability
- --- embeddings

90

BERTopic and Louvain mostly identify different communities, more balanced in Louvain

91

**sexual violence**

| term | value |
|---|---|
| violence | |
| sexual | |
| end | |
| woman | |
| girl | |
| survivor | |
| harassment | |
| base | |

0  0.02 0.04 0.06 0.08

**refugees**

| term | value |
|---|---|
| refugee | |
| syria | |
| withrefugees | |
| rohingya | |
| woman | |
| syrian | |
| story | |
| flee | |

0  0.05  0.1

**peace**

| term | value |
|---|---|
| peace | |
| peacebuilding | |
| pk70 | |
| peacekeeper | |
| woman | |
| role | |
| peace72 | |
| sustainingpeace | |

0  0.05  0.1

**girlsinict**

| term | value |
|---|---|
| girl | |
| digital | |
| girlsinict | |
| education | |
| close | |
| life | |
| woman | |
| world | |

0  0.02 0.04 0.06 0.08

**executive director**

| term | value |
|---|---|
| ap | |
| rt | |
| director | |
| sweden | |
| executive | |
| congratulation | |
| thank | |
| general | |

0  0.02 0.04 0.06 0.08

**mothersday**

| term | value |
|---|---|
| fathersday | |
| child | |
| mothersday | |
| dad | |
| family | |
| mom | |
| parentsday | |
| daughter | |

0  0.02 0.04 0.06

**sustainability**

| term | value |
|---|---|
| woman | |
| sustainable | |
| climate | |
| ocean | |
| change | |
| india | |
| community | |
| affect | |

0  0.02  0.04

**equal pay**

| term | value |
|---|---|
| pay | |
| work | |
| man | |
| woman | |
| more | |
| do | |
| less | |
| equalpayday | |

0  0.02 0.04 0.06

**discrimination**

| term | value |
|---|---|
| human | |
| right | |
| love | |
| day | |
| discrimination | |
| today | |
| standup4humanrights | |
| humanright | |

0  0.02 0.04 0.06

**gender equality**

| term | value |
|---|---|
| equality | |
| ac | |
| gender | |
| right | |
| speak | |
| woman | |
| feminist | |
| equal | |

0  0.02  0.04

**whatwomenwant**

| term | value |
|---|---|
| health | |
| care | |
| quality | |
| whatwomenwant | |
| access | |
| reproductive | |
| maternal | |
| healthcare | |

0  0.05  0.1

**politics**

| term | value |
|---|---|
| gender | |
| council | |
| equality | |
| advisory | |
| g7 | |
| un | |
| ed | |
| system | |

0  0.02 0.04 0.06

93

Same barchart with Louvain
#metoo2018

```python
def bertopic_overwrite(bert_model_in,docs,C):
  bert_model = copy.deepcopy(bert_model_in)

  # build the documents dataframe: 'Document' + "Topic"
  documents = pd.DataFrame(docs,columns=['Document'])
  tmp = np.array([C[i].argmax() for i in range(C.shape[0])])
  documents["Topic"] = tmp

  # update topic assignment
  bert_model.topics_ = tmp.tolist()

  # build cf-idf values
  documents_per_topic = documents.groupby(['Topic'],
                  as_index=False).agg({'Document': ' '.join})
  c_tf_idf_, words = bert_model._c_tf_idf(documents_per_topic)
  bert_model.c_tf_idf_ = c_tf_idf_

  # extract words representations
  topic_representations_ = bert_model._extract_words_per_topic(words, documents)
  bert_model.topic_representations_ = topic_representations_
  bert_model.topic_labels_ = {key: f"{key}_" + "_".join([word[0] for word in values[:4]])
                              for key, values in
                              bert_model.topic_representations_.items()}

  # exit
  return bert_model
```

❑ Naturally provides a hard topic assignment

❑ Useful tool

❑ More readable output with deep cleaned text
    but same performance

❑ Comparison – with Louvain
    weaker in general, especially in modularity
    equivalent NMI = relevant topics
    lower modularity = the documents that identify the
                                    topics are less distinguishable
    higher complexity involved
    less balanced topics, but generally meaningful
    topics correlated with Louvain

96

# Sentiment analysis

adding useful insights to your data

❑ <u>Sentiment</u> – e.g., positive, negative, neutral

   enduring cognitive content that defines the affective state

❑ <u>Emotion</u> – e.g., anger, disgust, fear, joy, sadness

   intense affective state of short duration with a precise cause

❑ <u>Ingroup bias</u> – e.g., use of pronouns I, we, us

   tendency to favor one's own group over other groups

❑ <u>Outgroup bias</u> – e.g., use of pronoun they

   tendency to dislike members of groups we don't identify with

❑ <u>Agency</u> – e.g., use of action verbs do, take, make

   perception that an individual is able to contribute to/a group
   can collectively reach a social change

# LIWC linguistic inquiry and word count

Tausczik, Pennebaker. "The psychological meaning of words:
LIWC and computerized text analysis methods."  (2010)

UNIVERSITÀ
DEGLI STUDI
DI PADOVA

https://citeseerx.ist.psu.edu/document?repid=rep1&type=pdf&doi=79d2494cc10a9633c42115df84bb74ed447080f6

**LIWC**  HOME  TRY IT NOW

## INTRODUCING LIWC-22

### A NEW SET OF TEXT ANALYSIS TOOLS AT YOUR FINGERTIPS

https://www.liwc.app/

❑ word count (or dictionary) methodology

❑ over 60 dictionaries coded and validated for their accuracy in reflecting psychological content

❑ simplicity of implementation and usage

❑ state-of-the-art in psychology

❑ one licence available in the instructor's PC ☺

99

UNIVERSITÀ
DEGLI STUDI
DI PADOVA

| Category | Examples | Words in Category | Psychological Correlates |
|---|---|---|---|
| *Linguistic processes* | | | |
| Word count | | | Talkativeness, verbal fluency |
| Words/sentence | | | Verbal fluency, cognitive complexity |
| Dictionary words | (Percentage of all words captured by the program) | | Informal, nontechnical language |
| Words >6 letters | (Percentage of all words longer than 6 letters) | | Education, social class |
| Total function words | | 464 | |
| Total pronouns | I, them, itself | 116 | Informal, personal |
| Personal pronouns | I, them, her | 70 | Personal, social |
| First-person singular | I, me, mine | 12 | Honest, depressed, low status, personal, emotional, informal |
| First-person plural | We, us, our | 12 | Detached, high status, socially connected to group (sometimes) |
| Second person | You, your, thou | 20 | Social, elevated status |
| Third-person singular | She, her, him | 17 | Social interests, social support |
| Third-person plural | They, their, they'd | 10 | Social interests, out-group awareness (sometimes) |

ingroup

outgroup

100

| Category | Examples | Words in Category | Psychological Correlates |
|---|---|---|---|
| Indefinite pronouns | It, it's, those | 46 | |
| Articles | A, an, the | 3 | Use of concrete nouns, interest in objects and things |
| Common verbs | Walk, went, see | 383 | |
| Auxiliary verbs | Am, will, have | 144 | Informal, passive voice |
| Past tense | Went, ran, had | 145 | Focus on the past |
| Present tense | Is, does, hear | 169 | Living in the here and now |
| Future tense | Will, gonna | 48 | Future and goal oriented |
| Adverbs | Very, really, quickly | 69 | |
| Prepositions | To, with, above | 60 | Education, concern with precision |
| Conjunctions | And, but, whereas | 28 | |
| Negations | No, not, never | 57 | Inhibition |
| Quantifiers | Few, many, much | 89 | |
| Numbers | Second, thousand | 34 | |
| Swear words | Damn, piss, fuck | 53 | Informal, aggression, |
| Psychological processes | | | |
| Social processes | Mate, talk, they, child | 455 | Social concerns, social support |
| Family | Daughter, husband | 64 | |
| Friends | Buddy, friend, neighbor | 37 | |
| Humans | Adult, baby, boy | 61 | |
| Affective processes | Happy, cried, abandon | 915 | Emotionality |

focus on past, present or future

101

| WC | Analytic | Clout | Authentic | Tone | WPS | Sixltr | Dic | function | pronoun |
|---|---|---|---|---|---|---|---|---|---|
| ppron | **i** | **we** | **you** | **shehe** | **they** | ipron | article | prep | auxverb |
| adverb | conj | negate | verb | adj | compare | interrog | number | quant | affect |
| **posemo** | **negemo** | **anx** | **anger** | **sad** | social | family | friend | female | male |
| insight | cause | discrep | tentat | certain | differ | percept | see | hear | feel |
| bio | **body** | **health** | **sexual** | ingest | drives | **affiliation** | **achieve** | **power** | **reward** |
| risk | **focus past** | **focus present** | **focus future** | relativ | motion | space | time | work | leisure |
| home | money | **relig** | death | **informal** | **swear** | netspeak | assent | nonflu | filler |
| AllPunc | Period | Comma | Colon | SemiC | QMark | Exclam | Dash | Quote | Apostro |
| Parenth | cogproc | | | | | | | | |

## Choose the ones of interest to your project!

**Stage 1** agency seed dictionary (ASD) assembly

**Stage 2** agency lexicographic dataset (ALD) extraction

**Stage 3** agency evaluation of ALD by human coders

**Stage 4** model fine-tuning

**Stage 5** gold standard dictionary (GSD) assembly

**Stage 6** agency evaluation of GSD by human coders

**Stage 7** model selection

BERTAgent
https://pypi.org/project/bertagent/

SWPS University

- ❑ agency in text
- ❑ uses deep learning
- ❑ based on BERTA
- ❑ a validated tool
- ❑ available on Python

103

| Variable | M | SD | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|---|---|
| 1. HumEval | 0.12 | 1.54 | | | | | | | | |
| 2. PietA | 0.05 | 0.05 | .17** [.06, .28] | | -1.25 | 0.28 | 0.05 | 5.35** | -1.78 | -10.95** |
| 3. PietB | 0.02 | 0.03 | .25** [.14, .35] | .40** [.30, .49] | | 1.27 | 1.16 | 6.58** | -0.70 | -10.00** |
| 4. PietC | 0.05 | 0.05 | .17** [.06, .28] | .99** [.99, 1.00] | .40** [.30, .49] | | 0.03 | 5.34** | -1.80 | -10.93** |
| 5. NicoPos | 0.03 | 0.04 | .17** [.05, .27] | .18** [.07, .29] | .23** [.12, .34] | .17** [.06, .28] | | 5.49** | -3.81** | -11.08** |
| 6. NicoNeg | 0.01 | 0.03 | -.28** [-.38, -.17] | -.10 [-.21, .01] | -.01 [-.12, .11] | -.10 [-.21, .02] | -.03 [-.14, .09] | | -5.73** | -13.40** |
| 7. NicoCom | 0.02 | 0.05 | .30** [.19, .40] | .20** [.09, .31] | .19** [.08, .30] | .19** [.08, .30] | .82** [.78, .85] | -.60** [-.67, -.52] | | -10.38** |
| 8. BATot | 0.09 | 0.35 | .78** [.73, .82] | .21** [.10, .31] | .24** [.13, .34] | .20** [.09, .31] | .22** [.11, .33] | -.42** [-.51, -.33] | .42** [.33, .51] | |

Human evaluation

BERTAgent

best correlation with Human evaluation

Z-statistics: correlation is statistically more relevant that DWC

104

# Agency in US elections
Twitter, 2020-2021
by Jan Nikadon @ swps

UNIVERSITÀ
DEGLI STUDI
DI PADOVA



Election day

Capitol Hill

Women's day

SWPS
University

Death of George Floyd 2020/05/25

Independence Day 2020/07/04

Death of Jacob Blake 2020/08/23

Election Day 2020/11/03

Thanksgiving 2020/11/26

Christmas 2020/12/25

New Year 2021/01/01

Storm on the Capitol 2021/01/06

Inauguration 2021/01/20

International Women's Day 2021/03/08

Easter 2021/04/04

Derek Chauvin found guilty 2021/04/20

agency raises before
elections than drops

twitter

105

# Agency in postpartum depression
## Reddit Posts 2021
### by Selen Arslan @ unipd/swps

UNIVERSITÀ DEGLI STUDI DI PADOVA

SWPS University



composite emotion score
(LIWC, EmoPos - EmoNeg)

semantic agency
(BERTAgent)

106

**UNIVERSITÀ DEGLI STUDI DI PADOVA**

```
[1]  !pip install bertagent
```

install, import,
set instance

```
from bertagent import BERTAgent
ba0 = BERTAgent()
```

input sentences must
be superficially cleaned

**SWPS University**

```
# provide example sentences
sents = ["hardly wo
         "hard work
         "striving
         "struggling
         "struggling
         "unable to
         "this car
         "this car
         "this poli
         ]
```

```
'hardly working individual' : -0.57
'hard working individual' : 0.44
'striving to achieve my goals' : 0.73
'struggling to achieve my goals' : -0.67
'struggling to survive' : -0.52
'unable to survive' : -0.57
'this car runs on gasoline with lead' : -0.03
'this car runs on gasoline and it will lead us' : 0.09
'this politician runs for office and he will lead us' : 0.58
```

BERTAgent output

```
# assign agency
vals = ba0.predict(sents)
# print results
for item in zip(sents, vals):
    print(f"  {item[0]!r} : {item[1]:.2f}")
```

run BERTAgent

107

# Using sentiment analysis

an overview on how it can be useful in your projects

Socio-psychological linguistic markers
a view on the entire tweets corpus

# Student's *t*-test

Article   Talk

From Wikipedia, the free encyclopedia

A **_t_-test** is a type of statistical analysis used to compare the averages of two groups and determine whether the differences between them are more likely to arise from random chance. It is any statistical hypothesis test in which the test statistic follows a Student's *t*-distribution under the null hypothesis. It is most commonly applied when the test statistic would follow a normal distribution if the value of a scaling term in the test statistic were known (typically, the scaling term is unknown and is therefore a nuisance parameter). When the scaling term is estimated based on the data, the test statistic—under certain conditions—follows a Student's *t* distribution. The *t*-test's most common application is to test whether the means of two populations are different.

111

## Assumption:

the average values of the two populations being compared should follow a normal distribution (e.g., with many samples)

## Hypothesis (to be tested):

H0 - the two sets have equal average value, $m_x = m_y$

H1 - the two sets have different average value, $m_x \neq m_y$

sample averages

Student's t distribution
with $\nu$ degrees of freedom
(an approximation under H0)

Test statistic: $t = \dfrac{\bar{x} - \bar{y}}{\sqrt{\dfrac{s_x^2}{N_x} + \dfrac{s_y^2}{N_y}}} \sim t_\nu$

$\hat{\nu} = \dfrac{\left(\dfrac{s_x^2}{N_x} + \dfrac{s_y^2}{N_y}\right)^2}{\dfrac{\left(s_x^2/N_x\right)^2}{N_x - 1} + \dfrac{\left(s_y^2/N_y\right)^2}{N_y - 1}}$

unbiased estimator
of variance

samples
sizes

estimate of $\nu$
(degrees of freedom)
used for calculations

112

we declare different
average values H1

| p-value | evidence |
|---------|----------|
| < .01 | very strong evidence against $H_0$ |
| .01 − .05 | strong evidence against $H_0$ |
| .05 − .10 | weak evidence against $H_0$ |
| > .1 | little or no evidence against $H_0$ |

Student-t PDF with $\nu$
degrees of freedom
under H0 (same averages)

½ p

½ p

-|t|

|t|

the *p* value
informs on
whether a
difference exists
(statistically)

it is a false
positive rate

probability p that, by
confirming H1, we are
making errors on H0

(conservative) region choice to
declare H1 from the t value

113

UNIVERSITÀ
DEGLI STUDI
DI PADOVA

$$\text{Cohen's } d = \frac{\bar{x} - \bar{y}}{\sqrt{a \, s_x^2 + (1-a) \, s_y^2}} \qquad a = \frac{N_x - 1}{N_x + N_y - 2}$$

the *d* value informs on the size of the effect

| Relative size | Effect size |
| --- | --- |
| | 0.0 |
| Small | 0.2 |
| Medium | 0.5 |
| Large | 0.8 |
| | 1.4 |

we confirm H1

114

115

projecting the adjacency matrix on topics

(b) We

(d) Future focus

relevant statistically changes of **we-future** only in the climate action community

# Wrap-up
on topic detection

❑ **What available tools should be used**

      Louvain & BERTopic

      compare their performance through NMI, modularity, etc.

      LIWC & BERTAgent

      to enrich your analysis under a socio-psycological lens

❑ **What available tools should NOT be used**

      InfoMap, NMF & LDA

      they show poor performance

❑ **What would be nice to see implemented**

      soft Louvain made fast

      performance of BigCLAM and SMBs

      NFTM VAE and its performance