



Hochschule Offenburg
offenburg.university

Natural Language Processing & Visual Analytics

Prof. Dr. Daniela Oelke



Elektrotechnik, Medizintechnik
und Informatik

Outline

- Text Classification
- Normalization of Natural Language Texts
- Feature Extraction from Natural Language Texts
- Information Extraction
- NLP with spaCy
- Visual Analytics

Data Science Job Postings

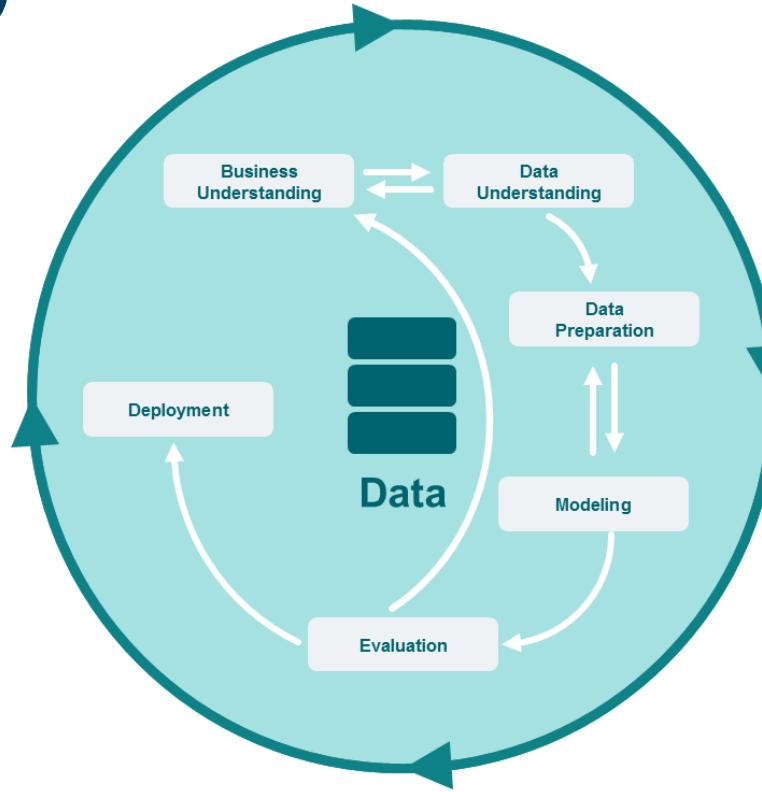
Data Science Job Postings

crawl_timestamp	job_title	category	company_name	city	state	country	inferred_city	post_date	job_description	job_type	job_board
06.02.2019 06:26	Enterprise Data Scientist I	Accounting/Finance	Farmers Insurance Group	Woodland Hills	CA	USA	Woodland hills	06.02.2019	Read what people are saying about working here. We are Farmers!Join a team of diverse pr...	undefined	indeed
06.02.2019 06:33	Data Scientist		Luxoft USA Inc	Middletown	NJ	USA	Middletown	05.02.2019	We have an immediate opening for a Sharp Data Scientist with a strong Mathematical/Statisti...	undefined	dice
06.02.2019 06:33	Data Scientist		Cincinnati Bell Technology Solutions	New York	NY	USA	New york	05.02.2019	Candidates should have the following background, skills and characteristics: Â· Experience d...	Full Time	dice
06.02.2019 06:33	Data Scientist, Aladdin Wealth Tech, Associate (M)	Accounting/Finance	BlackRock	New York	NY 10055 (Midt)	USA	New york	06.02.2019	Read what people are saying about working here. About BlackRockBlackRock helps investors...	undefined	indeed
06.02.2019 06:48	Senior Data Scientist	biotech	CyberCoders	Charlotte	NC	USA	Charlotte	05.02.2019	We are seeking an extraordinary Data Scientist in Charlotte to join our fast growing healthca...	Full Time	monster
06.02.2019 06:36	CIB â€“ Fixed Income Research â€“ Machine Learning	Accounting/Finance	JP Morgan Chase	New York	NY 10179 (Midt)	USA	New york	05.02.2019	Read what people are saying about working here. OpportunityThe opportunity is to join our...	undefined	indeed
06.02.2019 06:34	Data Scientist, Licensing Operations	Accounting/Finance	Spotify	New York	NY 10011 (Chels)	USA	New york	06.02.2019	Read what people are saying about working here. At Spotify our mission is to provide the w...	undefined	indeed
06.02.2019 06:52	Sr. Data Scientist (Can work on Xoriant W2)		Xoriant Corporation	Santa Clara	CA	USA	Santa clara	06.02.2019	Job Title:- Sr. Data Science Consultant Duration: 1Yrs (will get extend) Mode of interview:- Contract	Contract	dice
06.02.2019 06:34	Data Scientist, Aladdin Wealth Tech, Associate	Accounting/Finance	BlackRock	New York	NY 10055 (Midt)	USA	New york	06.02.2019	Read what people are saying about working here. About BlackRockBlackRock helps investors...	undefined	indeed
06.02.2019 07:03	Data Scientist		Adroit Resources	San Francisco	CA	USA	San francisco	05.02.2019	Â¢C 3+ years related a professional experienceÂ¢C Proven achievements resulting from pr...	Contract	dice
06.02.2019 07:19	Data Scientist	Computer/Internet	Northrop Grumman	Monterey	CA 93940	USA	Monterey	06.02.2019	Read what people are saying about working here. At Northrop Grumman, innovation isn't ju...	undefined	indeed
06.02.2019 07:24	Data Scientist		Envoy Consulting Group, Inc.	Reston	VA	USA	Reston	05.02.2019	Data Scientist EnvoyIT is looking for a Data Scientist for a Full-time position with one of our...	Full Time	dice
06.02.2019 07:21	ETL Developer / Data Scientist	Computer/Internet	Nobilis	Reston	VA 20191	USA	Reston	06.02.2019	Read what people are saying about working here. RESPONSIBILITIESThe Citizen Services Mis...	undefined	indeed
06.02.2019 07:22	Research Data Scientist	Computer/Internet	ARUP Laboratories	Salt Lake City	UT	USA	Salt lake city	06.02.2019	Read what people are saying about working here. RESPONSIBILITIES	undefined	indeed
06.02.2019 07:33	Data Scientist		Perspecta	Washington	DC	USA	Washington	05.02.2019	Every day at Perspecta, we en...	RESPONSIBILITIES	
06.02.2019 07:38	Senior Data Scientist :: VK (9)		Pyramid Consulting, Inc.	McLean	VA	USA	McLean	05.02.2019	Immediate need for Senior Data...	The Citizen Service Mission Area (CS) provides guidance and support to its clien...	
06.02.2019 07:39	Senior Data Scientist - W2		Orpine.com	Boston	MA	USA	Boston	05.02.2019	Work Authorization: Only US C...	CS is building a behavioral analytics team to support a wide variety of client ne...	
06.02.2019 07:39	Data Scientist		Apeiro Technologies	McLean	VA	USA	McLean	05.02.2019	Hi, Hope you are doing great. V...	This position is for a Software Developer with exposure to primary to SAS, Extract...	
06.02.2019 07:27	Data Scientist	Computer/Internet	Bank of America	Seattle	WA 98104 (First)	USA	Seattle	06.02.2019	Read what people are saying about working here. This position is for a Software Develop...	RESPONSIBILITIES	
06.02.2019 07:27	Data Scientist	Computer/Internet	Sallie Mae	Newark	DE	USA	Newark	06.02.2019	Read what people are saying about working here. This position is for a Software Develop...	QUALIFICATIONS	
06.02.2019 07:31	Sr. Data Scientist	Arts/Entertainment/Publishing	Taboola	Los Angeles	CA	USA	Los angeles	06.02.2019	Read what people are saying about working here. This position is for a Software Develop...	RESPONSIBILITIES	
06.02.2019 07:27	Data Scientist	Computer/Internet	VideoAmp	Santa Monica	CA	USA	Santa monica	06.02.2019	Read what people are saying about working here. This position is for a Software Develop...	QUALIFICATIONS	
06.02.2019 07:28	Data Scientist	Computer/Internet	Apple	Cambridge	MA	USA	Cambridge	06.02.2019	Read what people are saying about working here. This position is for a Software Develop...	RESPONSIBILITIES	
06.02.2019 07:50	Data Scientist		Amtex Enterprises	Plano	TX	USA	Plano	06.02.2019	Data Scientist 2 years Qualificationsdatabase development and application integrati...	RESPONSIBILITIES	
06.02.2019 07:53	Senior Data Scientist		Sysmind, LLC	Dallas	TX	USA	Dallas	05.02.2019	Position : Senior Data Scientist, advanced knowledge of SAS, C #, Java, Perl, Python, Scala and r programming	RESPONSIBILITIES	
06.02.2019 07:51	Data Scientist with Sagemaker Experience NO OF		For All	Orlando	FL	USA	Orlando	05.02.2019	Position : Senior Data Scientist, advanced knowledge of SAS, C #, Java, Perl, Python, Scala and r programming	RESPONSIBILITIES	
06.02.2019 08:06	Data Scientist / Data Engineer									SAS or equivalent training in SAS or other similar tool	
06.02.2019 08:03	Data Scientist with Exp in Data Modeling									Hadoop specifically Spark and Spark Streaming	
06.02.2019 08:06	Data Scientist									HDFS to Hadoop File System integration	
06.02.2019 08:04	Data Scientist needed with solid R and Python skill										
06.02.2019 08:08	Level1 Solutions Engineer - Data Scientist										
06.02.2019 08:40	Data Scientist	military									
06.02.2019 08:30	Data Scientist	business and financial operation									
06.02.2019 08:30	Senior Data Scientist - Tallahassee, FL - \$150K-\$170K	business and financial operation									
06.02.2019 08:28	Data Scientist	business and financial operation									
06.02.2019 08:27	Data Scientist	business and financial operation									
06.02.2019 08:27	Data Scientist	business and financial operation									
06.02.2019 08:28	Data Scientist	business and financial operation									
06.02.2019 08:31	Data Scientist Tampa, FL \$110-130K	business and financial operation									
06.02.2019 07:48	Data Scientist										
06.02.2019 08:29	Clinical Data Scientist	business and financial operation									
06.02.2019 09:37	Sr. Mgr, Data Scientist	Engineering/Architecture									
06.02.2019 09:37	Data Scientist Intern, Engineering - Software Dev	Engineering/Architecture									
06.02.2019 09:49	Data Scientist	Manufacturing/Mechanical	Comcast	Philadelphia	PA 19103	USA	Philadelphia	06.02.2019	Read what people are saying about working here. Comcast brings together the best in medi...	undefined	indeed
06.02.2019 09:49	Principal Data Scientist (REMOTE)	Engineering/Architecture	Fiserv, Inc.	Boston	VA 20100	USA	Boston	05.02.2019	Read what people are saying about working here. Job DescriptionWe dream of a job unde...	undefined	indeed

Using this dataset:
 Which analysis goals could a data science student have?
 Which analysis goals could the HR department of a company have?

crawl_timestamp	job_title	category	company_name	city	state	country	inferred_city	post_date	job_description	job_type	job_board
06.02.2019 06:26	Enterprise Data Scientist I	Accounting/Finance	Farmers Insurance Group	Woodland Hills	CA	USA	Woodland hills	06.02.2019	Read what people are saying about working here. We are Farmers!Join a team of diverse pr	undefined	indeed
06.02.2019 06:33	Data Scientist	luxoft USA Inc	Middletown	NJ	USA	USA	Middletown	05.02.2019	We have an immediate opening for a Sharp Data Scientist with a strong Mathematical/Statisti	undefined	dice
06.02.2019 06:33	Data Scientist	Cincinnati Bell Technology Solutions	New York	NY	USA	USA	New York	05.02.2019	Candidates should have the following background, skills and characteristics: Â· Experience in Full Time	Full Time	dice
06.02.2019 06:33	Data Scientist, Aladdin Wealth Tech, Associate (M)	Accounting/Finance	BlackRock	New York	NY 10055 (Midtown)	USA	New York	06.02.2019	Read what people are saying about working here. About BlackRockBlackRock helps investors	undefined	indeed
06.02.2019 06:48	Senior Data Scientist	biotech	CyberCoders	Charlotte	NC	USA	Charlotte	05.02.2019	We are seeking an extraordinary Data Scientist in Charlotte to join our fast growing healthca	Full Time	monster
06.02.2019 06:56	CIB â€“ Fixed Income Research â€“ Machine Learn	Accounting/Finance	JP Morgan Chase	New York	NY 10179 (Midtown)	USA	New York	05.02.2019	Read what people are saying about working here. OpportunityThe opportunity is to join our	undefined	indeed
06.02.2019 06:34	Data Scientist, Licensing Operations	Accounting/Finance	Spotify	New York	NY 10011 (Chelsea)	USA	New York	06.02.2019	Read what people are saying about working here. At Spotify our mission is to provide the wo	undefined	indeed
06.02.2019 06:52	Sr. Data Scientist (Can work on Xoriant W2)	Xoriant Corporation	Santa Clara	CA	USA	USA	Santa clara	06.02.2019	Job Title: - Sr. Data Science Consultant Duration: 1+ Yrs (will get extend) Mode of Interview - Contract	Contract	dice
06.02.2019 06:34	Data Scientist, Aladdin Wealth Tech, Associate	Accounting/Finance	BlackRock	New York	NY 10055 (Midtown)	USA	New York	06.02.2019	Read what people are saying about working here. About BlackRockBlackRock helps investors	undefined	indeed
06.02.2019 07:03	Data Scientist	Computer/Internet	Adroit Resources	San Francisco	CA	USA	San francisco	05.02.2019	â€¢ 3+ years related a professional experienceÂ· Proven achievements resulting from pr	Contract	dice
06.02.2019 07:19	Data Scientist	Computer/Internet	Northrop Grumman	Monterey	CA 93940	USA	Monterey	06.02.2019	Read what people are saying about working here. At Northrop Grumman, innovation isn't ju	undefined	indeed
06.02.2019 07:24	Data Scientist	Computer/Internet	Envoy Consulting Group, Inc.	Reston	VA	USA	Reston	05.02.2019	Data Scientist EnvoyIT is looking for a Data Scientist for a Full-time position with one of our	Full Time	dice
06.02.2019 07:21	ETL Developer / Data Scientist	Computer/Internet	Noblis	Reston	VA 20191	USA	Reston	06.02.2019	Read what people are saying about working here. RESPONSIBILITIESThe Citizen Services Mis	undefined	indeed
06.02.2019 07:22	Research Data Scientist	Computer/Internet	ARUP Laboratories	Salt Lake City	UT	USA	Salt lake city	06.02.2019	Read what people are saying about working here. Job DetailsDescription:Schedule:Monday -	Monday	indeed
06.02.2019 07:33	Data Scientist	Computer/Internet	Perspecta	Washington	DC	USA	Washington	05.02.2019	tion s Full Time	Full Time	dice
			Pyramid Consulting, Inc.	McLean	VA	USA	McLean	05.02.2019	Mont undefined	Mont	dice
			Orpine.com	Boston	MA	USA	Boston	05.02.2019	Business In Contract	Contract	dice
				WA 98104 (First)	WA	USA	Seattle	06.02.2019	of our Contract	Contract	dice
				Inc	WA 98104 (First)	USA	Newark	06.02.2019	Data Sci Undefined	indeed	dice
									and we will	indeed	dice

Cross Industry Standard Process for Data Mining (CRISP-DM)



Text Classification

Finding the right category

crawl_timestamp	job_title	category	company_name	city	state	country	inferred_city	post_date	job_description	job_type	job_board
06.02.2019 06:26	Enterprise Data Scientist I	Accounting/Finance	Farmers Insurance Group	Woodland Hills	CA	Usa	Woodland hills	06.02.2019	We have an immediate opening for a Sharp Data Scientist with a strong Mathematical/Statistical background. We are looking for someone who can work well in a team and independently. We are looking for someone who can work well in a team and independently.	Full Time	indeed
06.02.2019 06:33	Data Scientist	Luxoft USA Inc	Middletown	NJ	Usa	Middletown	New York	05.02.2019	Candidates should have the following background, skills and characteristics: Experience in Big Data, Machine Learning, Python, Java, C/C++, R, SQL, Hadoop, NoSQL, and distributed systems. Strong communication skills, ability to work in a team and independently.	Full Time	dice
06.02.2019 06:33	Data Scientist	Cincinnati Bell Technology Solutions	New York	NY	Usa	New York	New York	06.02.2019	We are seeking an experienced Data Scientist with a strong background in machine learning, data mining, and statistical analysis. The ideal candidate will have experience working with large datasets and developing predictive models. A solid understanding of Python, R, and SQL is required.	Full Time	indeed
06.02.2019 06:33	Data Scientist, Aladdin Wealth Tech, Associate (M)	BlackRock	JP Morgan Chase	New York	NY	Usa	New York	05.02.2019	Read what people are saying about working here. About BlackRockBlackRock helps investors and companies manage trillions of dollars in assets and liabilities across nearly every asset class. Our clients include more than 400 of the world's leading financial institutions, pension funds, sovereign wealth funds, governments and corporations in more than 35 countries.	Full Time	indeed
06.02.2019 06:48	Senior Data Scientist	bioTech	CyberCoders	Charlotte	NC	Usa	Charlotte	05.02.2019	We are seeking an extraordinary Data Scientist in Charlotte to join our fast growing health care company. The ideal candidate will have a strong background in data science, machine learning, and statistical analysis. A solid understanding of Python, R, and SQL is required.	Full Time	monster
06.02.2019 06:36	CIB â€“ Fixed Income Research â€“ Machine Learning	Accounting/Finance	Spotify	New York	NY	Usa	New York	06.02.2019	Read what people are saying about working here. OpportunityThe opportunity is to join our fast growing music company. The ideal candidate will have a strong background in data science, machine learning, and statistical analysis. A solid understanding of Python, R, and SQL is required.	Full Time	indeed
06.02.2019 06:34	Data Scientist, Licensing Operations	Accounting/Finance	Xoriant Corporation	Santa Clara	CA	Usa	Santa clara	06.02.2019	Read what people are saying about working here. At Spotify our mission is to provide the world's best music. The ideal candidate will have a strong background in data science, machine learning, and statistical analysis. A solid understanding of Python, R, and SQL is required.	Full Time	indeed
06.02.2019 06:52	Sr. Data Scientist (Can work on Xoriant W2)	BlackRock	New York	NY	10055 (Midt)	Usa	New York	06.02.2019	Job Title: Sr. Data Science Consultant Duration: 1-Yrs (will get extend) Mode of interview: Contract	Contract	dice
06.02.2019 06:34	Data Scientist, Aladdin Wealth Tech, Associate	Accounting/Finance	Adroit Resources	San Francisco	CA	Usa	San francisco	05.02.2019	Read what people are saying about working here. About BlackRockBlackRock helps investors and companies manage trillions of dollars in assets and liabilities across nearly every asset class. Our clients include more than 400 of the world's leading financial institutions, pension funds, sovereign wealth funds, governments and corporations in more than 35 countries.	Contract	indeed
06.02.2019 07:03	Data Scientist	Northrop Grumman	Monterey	CA	93940	Usa	Monterey	06.02.2019	Read what people are saying about working here. At Northrop Grumman, innovation isn't just a buzzword - it's how we do business. We're looking for individuals who are excited to join us in our mission to defend, protect, and serve our nation and our allies.	Full Time	indeed
06.02.2019 07:19	Data Scientist	Envoy Consulting Group, Inc	Reston	VA	Usa	Reston	Reston	05.02.2019	Data Scientist EnvoyIT is looking for a Data Scientist for a Full-time position with one of our clients. The ideal candidate will have a strong background in data science, machine learning, and statistical analysis. A solid understanding of Python, R, and SQL is required.	Full Time	dice
06.02.2019 07:24	Data Scientist	Nobilis	Reston	VA	20191	Usa	Reston	06.02.2019	Read what people are saying about working here. RESPONSIBILITIESThe Citizen Services Mission	Full Time	indeed
06.02.2019 07:21	ETL Developer / Data Scientist	ARUP Laboratories	Salt Lake City	UT	Usa	Salt lake city	Salt Lake City	06.02.2019	Read what people are saying about working here. Job DetailsDescriptionSchedule:Monday - Friday	Full Time	indeed
06.02.2019 07:22	Research Data Scientist	Perspecta	Washington	DC	Usa	Washington	Washington	05.02.2019	Every day at Perspecta, we enable hundreds of thousands of people to take on our nation's most complex challenges. We're looking for individuals who are excited to join us in our mission to defend, protect, and serve our nation and our allies.	Full Time	dice
06.02.2019 07:33	Data Scientist	Pyramid Consulting, Inc.	McLean	VA	Usa	McLean	McLean	05.02.2019	Immediate need for Senior Data Scientist with experience in the IT Industry. This is 12 Month contract position.	Full Time	indeed
06.02.2019 07:38	Senior Data Scientist :: VK (9)	Orpine.com	Boston	MA	Usa	Boston	Boston	05.02.2019	Work Authorization: Only US Citizen's & Green Card. A. Job description: A. â€¢ Business intelligence	Contract	dice
06.02.2019 07:39	Senior Data Scientist - W2	Apeiro Technologies	McLean	VA	Usa	McLean	McLean	05.02.2019	H, I, Hope you are doing great. We are having requirement for Data Scientist with one of our clients. The ideal candidate will have a strong background in data science, machine learning, and statistical analysis. A solid understanding of Python, R, and SQL is required.	Contract	dice
06.02.2019 07:39	Data Scientist	Bank of America	Seattle	WA	98104 (First)	Usa	Seattle	06.02.2019	Read what people are saying about working here. Job Description:Position SummaryData Scientist	Full Time	indeed
06.02.2019 07:27	Data Scientist	Taboola	Newark	DE	Usa	Newark	Newark	06.02.2019	Read what people are saying about working here. At Operations strategy and Analytics, we are looking for a Data Scientist to join our team.	Full Time	indeed
06.02.2019 07:31	Sr. Data Scientist	VideoAmp	Los Angeles	CA	Usa	Los angeles	Los angeles	06.02.2019	Read what people are saying about working here. Read something interesting online today?	Full Time	indeed
06.02.2019 07:27	Data Scientist	Apple	Seattle	WA	20191	Usa	Seattle	06.02.2019	Read what people are saying about working here. About UsVideoAmpâ€™s mission is to bring video to life.	Full Time	indeed
06.02.2019 07:28	Data Scientist	Amtex Enterprises	Cambridge	MA	Usa	Cambridge	Cambridge	06.02.2019	Read what people are saying about working here. SummaryPosted: Feb 5, 2019Weekly Hour:	Full Time	indeed
06.02.2019 07:50	Data Scientist	Sysmind, LLC	Plano	TX	Usa	Plano	Plano	06.02.2019	Data Scientist 2 years Qualifications and Skills: A. â€¢ A minimum of 3+ years work experience in a business environment.	Contract	dice
06.02.2019 07:53	Senior Data Scientist	FocuzMindz	Dallas	TX	Usa	Dallas	Dallas	05.02.2019	Position : Senior Data Scientist Location : Dallas, TX Job Description: 8-10 Years ExperienceA. â€¢	Full Time	indeed
06.02.2019 07:51	Data Scientist with Sagemaker Experience NO OF	Digital Intelligence Systems, LLC	Dallas	TX	Usa	Dallas	Dallas	05.02.2019	Job Title: Data Scientist with Sagemaker Experience NO OPTâ€™s please Work Location & F	Contract	dice
06.02.2019 08:06	Data Scientist / Data Engineer	Advent Global Solutions, Inc.	Bothell	WA	Usa	Bothell	Bothell	05.02.2019	Incorporated in 1994, DISYS is one of the largest IT staffing firms in the US. DISYS has nearly 200 offices worldwide.	Contract	dice
06.02.2019 08:03	Data Scientist with Exp in Data Modelling	Resource Informatics Group	Houston	TX	Usa	Houston	Houston	05.02.2019	Urgent need for a Data Scientist. Essentially, the guidance was, send resumes and weâ€™ll call you.	Contract	dice
06.02.2019 08:06	Data Scientist	Viri Technology	Seattle	WA	Usa	Seattle	Seattle	06.02.2019	We have below urgent positon in TX. Below are the more details on it. Let me know your interest.	Full Time	dice
06.02.2019 08:04	Data Scientist needed with solid R and Python skill	Kforce Technology Staffing	Irvine	CA	Usa	Irvine	Irvine	07.02.2019	Viri Technology is seeking an up and coming Data Scientist for an immediate FTE role with a client in CA.	Full Time	dice
06.02.2019 08:08	Level 1 Solutions Engineer - Data Scientist	L3 Technologies	Chantilly	VA	Usa	Chantilly	Chantilly	03.02.2019	RESPONSIBILITIES: Kforce has a client seeking a Level 1 Solutions Engineer - Data Scientist in VA.	Full Time	indeed
06.02.2019 08:40	Data Scientist	military	West Chester Towns	OH	Usa	West chester	West chester	05.02.2019	Description Building large systematic reports and one-off small pieces of exploratory analysis.	Full Time	monster
06.02.2019 08:30	Data Scientist	Partners, Inc.	Tallahassee	FL	Usa	Tallahassee	Tallahassee	05.02.2019	Has operations in multiple states across the country. National motor vehicle industry leader.	Full Time	careerbuilder
06.02.2019 08:28	Data Scientist	Dallas	Dallas	TX	Usa	Dallas	Dallas	05.02.2019	an Data Science Companyâ€, focused on the development of advanced digital technology and analytical solutions.	Full Time	careerbuilder
06.02.2019 08:27	Data Scientist	Plymouth Meeting	PA	Usa	Plymouth Meeting	PA	Plymouth Meeting	05.02.2019	Advanced Digital Technology and Analytics.	Full Time	careerbuilder
06.02.2019 08:27	Data Scientist	Seattle	WA	Usa	Alexandria	VA	Alexandria	05.02.2019	to hire roles with a leading client in Seattle.	Full Time	careerbuilder
06.02.2019 08:28	Data Scientist	Grant Thornton LLP	Tampa	FL	Usa	Tampa	Tampa	05.02.2019	Data Scientist Tampa, FL \$110-130K Job Description - A highly respected client is seeking a skilled Data Scientist.	Full Time	careerbuilder
06.02.2019 08:31	Data Scientist Tampa, FL \$110-130K	Experis	Seattle	WA	Usa	Seattle	Seattle	01.02.2019	The Team The AIML (Artificial Intelligence and Machine Learning) team contributes to the vision and delivery of AI solutions.	Full Time	careerbuilder
06.02.2019 07:48	Data Scientist	Jefferson Frank	Tampa	FL	Usa	Tampa	Tampa	05.02.2019	POSITION SUMMARY: The Contract Clinical Data Scientist in CDS is responsible for supporting the clinical needs of the organization.	Full Time	careerbuilder
06.02.2019 08:29	Clinical Data Scientist	Staffing Technologies	Boston	MA	Usa	Boston	Boston	06.02.2019	Read what people are saying about working here. Our CompanyThe Kraft Heinz Company is a global food and beverage company with operations in over 100 countries.	Full Time	indeed
06.02.2019 09:37	Sr. Mgr, Data Scientist	Aerotek	Waltham	MA	Usa	Waltham	Waltham	05.02.2019	POSITION SUMMARY: The Contract Clinical Data Scientist in CDS is responsible for supporting the clinical needs of the organization.	Full Time	careerbuilder
06.02.2019 09:37	Data Scientist Intern, Engineering - Software Development	Kraft Heinz Company	San Francisco	CA	Usa	San francisco	San francisco	06.02.2019	Read what people are saying about working here. InternshipWho we areCriteo, we are a leading provider of AI-powered solutions for the retail and consumer goods industries.	Full Time	indeed
06.02.2019 09:49	Data Scientist	Criteo	Palo Alto	CA	Usa	Palo alto	Palo alto	06.02.2019	Read what people are saying about working here. InternshipWho we areCriteo, we are a leading provider of AI-powered solutions for the retail and consumer goods industries.	Full Time	indeed
06.02.2019 09:49	Data Scientist	Comcast	Philadelphia	PA	Usa	Philadelphia	Philadelphia	06.02.2019	Read what people are saying about working here. Comcast brings together the best in media and technology to create an exceptional customer experience.	Full Time	indeed
06.02.2019 09:49	Principal Data Scientist (REMOTE)	Fiserv	Boston	MA	Usa	Boston	Boston	05.02.2019	Read what people are saying about working here. Job DescriptionIf you dream of a job where you can make a difference, this is the place for you.	Full Time	indeed

Target: category

Given: job description

Step 1: Feature Extraction

job_description
Read what people are saying about working here. We are Farmers!Join a team of diverse p
We have an immediate opening for a Sharp Data Scientist with a strong Mathematical/Stat
Candidates should have the following background, skills and characteristics: Â Experience
Read what people are saying about working here. About BlackRockBlackRock helps investo
We are seeking an extraordinary Data Scientist in Charlotte to join our fast growing healthc
Read what people are saying about working here. OpportunityThe opportunity is to join ou
Read what people are saying about working here. At Spotify our mission is to provide the v
Job Title: - Sr. Data Science Consultant Duration: 1+ Yrs (will get extend) Mode of interview
Read what people are saying about working here. About BlackRockBlackRock helps investo
â€¢ 3+ years related a professional experienceÂ ¢ Proven achievements resulting from p
Read what people are saying about working here. At Northrop Grumman, innovation isn't j
Data Scientist EnvoyIT is looking for a Data Scientist for a Full-time position with one of our
Read what people are saying about working here. RESPONSIBILITIESThe Citizen Services Mi
Read what people are saying about working here. Job DetailsDescriptionSchedule:Monday
Every day at Perspecta, we enable hundreds of thousands of people to take on our nation's
Immediate need for Senior Data Scientist with experience in the IT Industry. This is 12 Mo
Work Authorization: Only US Citizen's & Green CardÂ ¢ Job description:Â ¢ Business
Hi, Hope you are doing great. We are having requirement forÂ Data Scientist with one of o
Read what people are saying about working here. Job Description:Position SummaryData S
Read what people are saying about working here. At Operations strategy and Analytics, we
Read what people are saying about working here. Read something interesting online today
Read what people are saying about working here. About UsVideoAmpâ€™s mission is to br
Read what people are saying about working here. SummaryPosted: Feb 5, 2019Weekly Hou
Data Scientist 2 years Qualifications and Skills: Â·Â·Â·Â·Â· 3+ years work experience in a

How do we get from plain text to a feature vector?

Term-document count matrices

- Consider the number of occurrences of a term in a document:
 - Each document is a **count vector** in \mathbb{N}^v : a column below

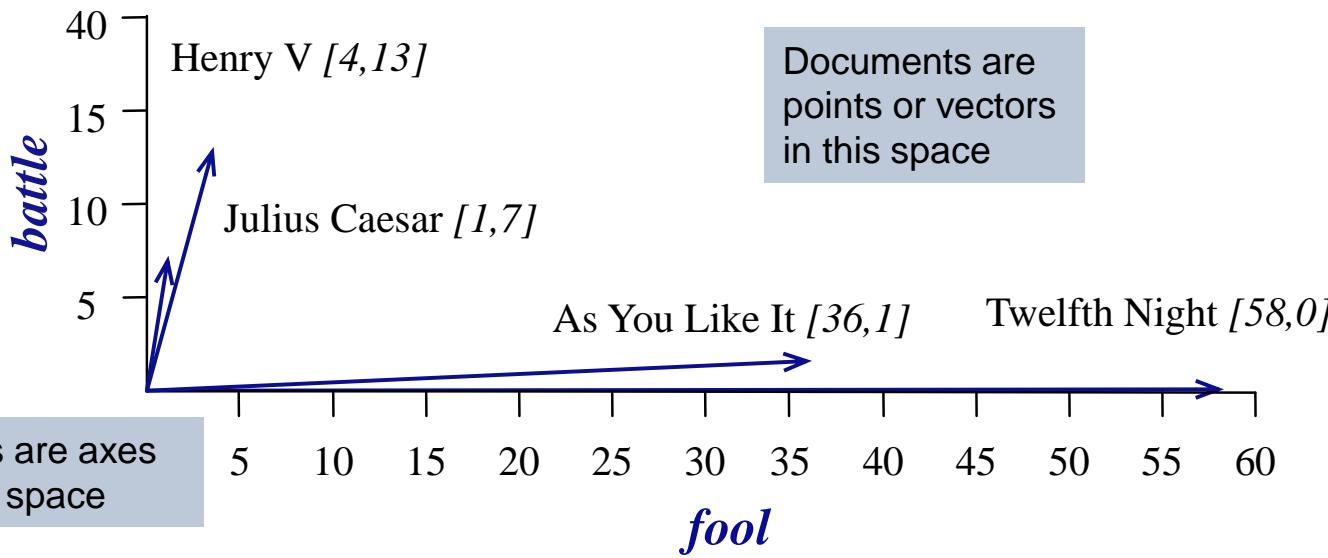
	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	157	73	0	0	0	0
Brutus	4	157	0	1	0	0
Caesar	232	227	0	2	1	1
Calpurnia	0	10	0	0	0	0
Cleopatra	57	0	0	0	0	0
mercy	2	0	3	5	5	1
worser	2	0	1	1	1	0

Bag of words model

- Vector representation doesn't consider the ordering of words in a document
- *John is quicker than Mary and Mary is quicker than John have the same vectors*
- This is called the bag of words model.

Document vectors visualized

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	1	0	7	13
good	114	80	62	89
fool	36	58	1	4
wit	20	15	2	3

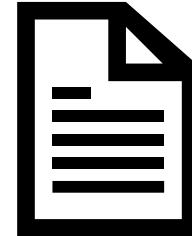


Term frequency tf

- The term frequency $tf_{t,d}$ of term t in document d is defined as the number of times that t occurs in d .
- Compute query-document match score



search term
count = 1



search term
count = 10

**More relevant
10x more
relevant?**

Rare terms are more informative



Term frequency (arachnophobia) = 1

Term frequency (actor) = 1

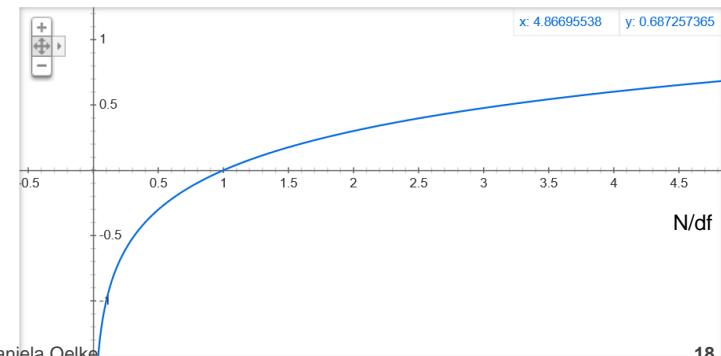
- Search term: *arachnophobia* → document most likely relevant
- Search term: *actor* → document not necessarily relevant

idf weight

- N is the number of documents in the collection
- df_t is the document frequency of t : the number of documents that contain t
 - df_t is an inverse measure of the informativeness of t
 - $df_t \leq N$
- We define the idf (inverse document frequency) of t by

$$\text{idf}_t = \log_{10} (N/\text{df}_t)$$

- We use $\log (N/\text{df}_t)$ instead of N/df_t to “dampen” the effect of idf.



idf example, suppose $N = 1$ million

term	df _t	idf _t
calpurnia	1	6
animal	100	4
sunday	1,000	3
fly	10,000	2
under	100,000	1
the	1,000,000	0

$$\text{idf}_t = \log_{10} (N/\text{df}_t)$$

There is one idf value for each term t in a collection.

TF-IDF weighting

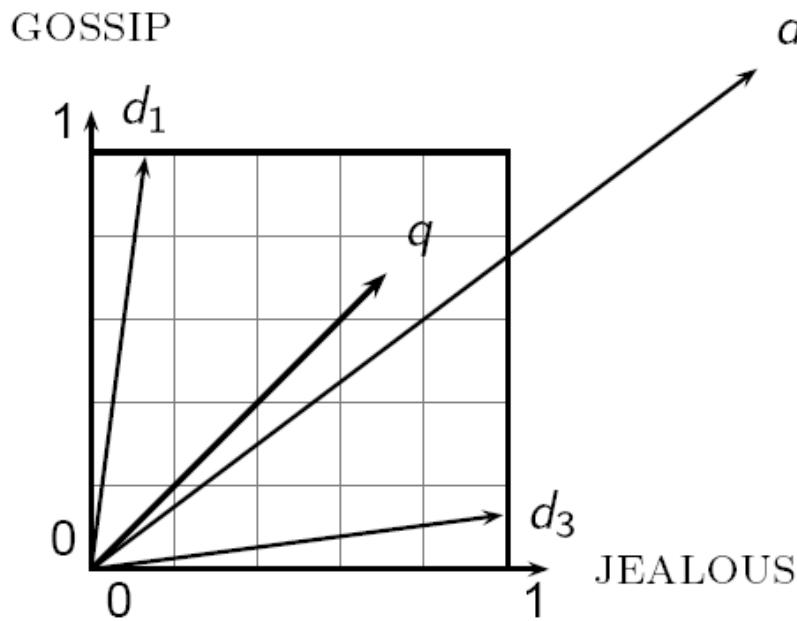
$$w_{t,d} = tf_{t,d} \times \log_{10}\left(\frac{N}{df_t}\right)$$

When does the tf-idf value of a term increase?

- with the number of occurrences of the term within a document
- with the rarity of the term in the collection

Variants of how TF-IDF values are calculated exist

Calculating the similarity between two document vectors



How would you calculate the similarity between these document vectors?

Cosine similarity

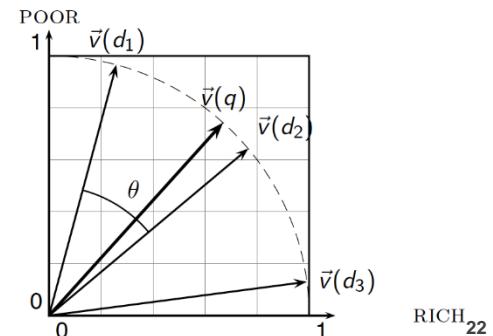
$$\text{cosine}(\vec{v}, \vec{w}) = \frac{\vec{v} \cdot \vec{w}}{|\vec{v}| |\vec{w}|} = \frac{\sum_{i=1}^N v_i w_i}{\sqrt{\sum_{i=1}^N v_i^2} \sqrt{\sum_{i=1}^N w_i^2}}$$

v_i is the tf-idf value for word i in document v
 w_i is the tf-idf value for word i in document w

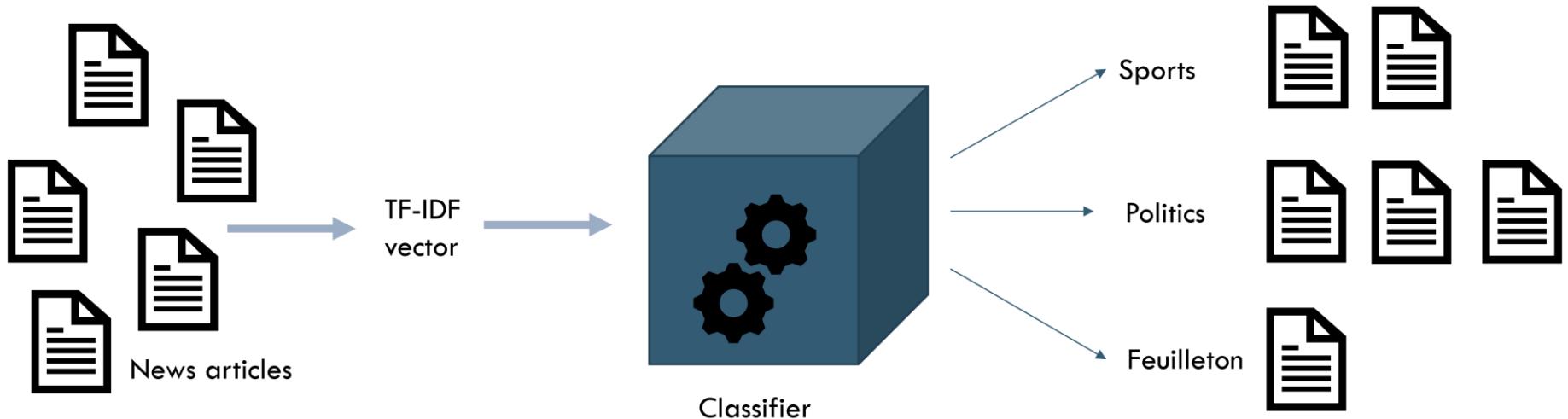
$\text{Cos}(v, w)$ is the cosine similarity of v and w

If the vectors are length-normalized:

$$\text{cosine}(\vec{v}, \vec{w}) = \sum_{i=1}^N v_i w_i$$



Application: TF-IDF for classification



Exercise 1: Text Classification

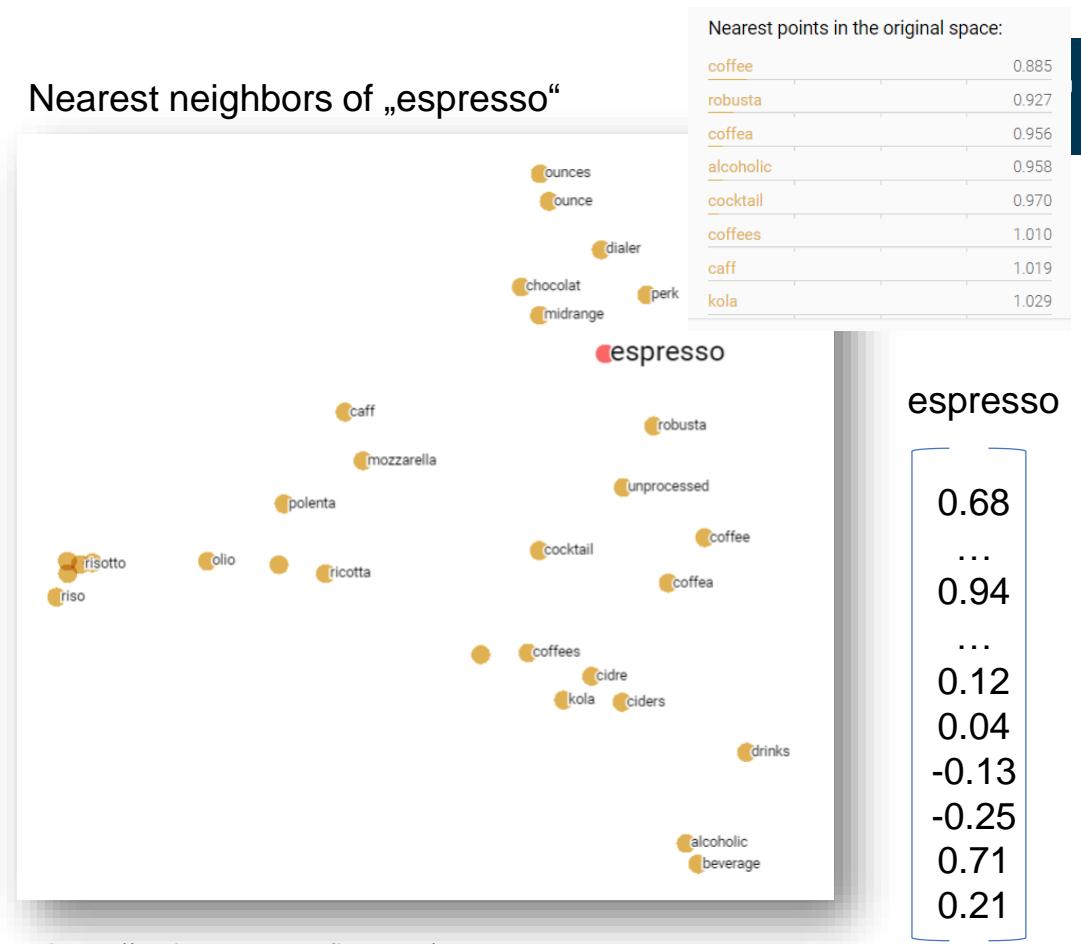
What's next? Word Embeddings & Co.

Word2Vec

- Word Embedding = Representation of words as vectors in an n-dimensional space
- Similar Words are located close to each other in the feature space

Caveat:
Projection from
200D → 3D!

Nearest neighbors of „espresso“

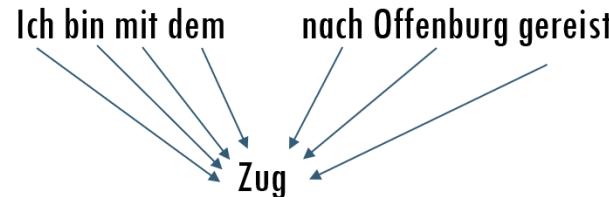


How do we get the vectors?

Method 1:

Continuous Bag of Words (CBOW)

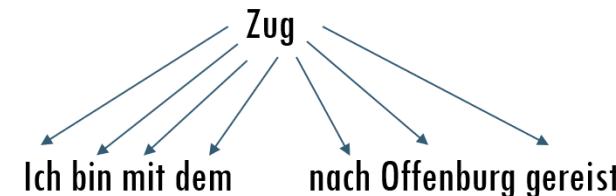
Task: Using the context, predict the missing word



Method 2:

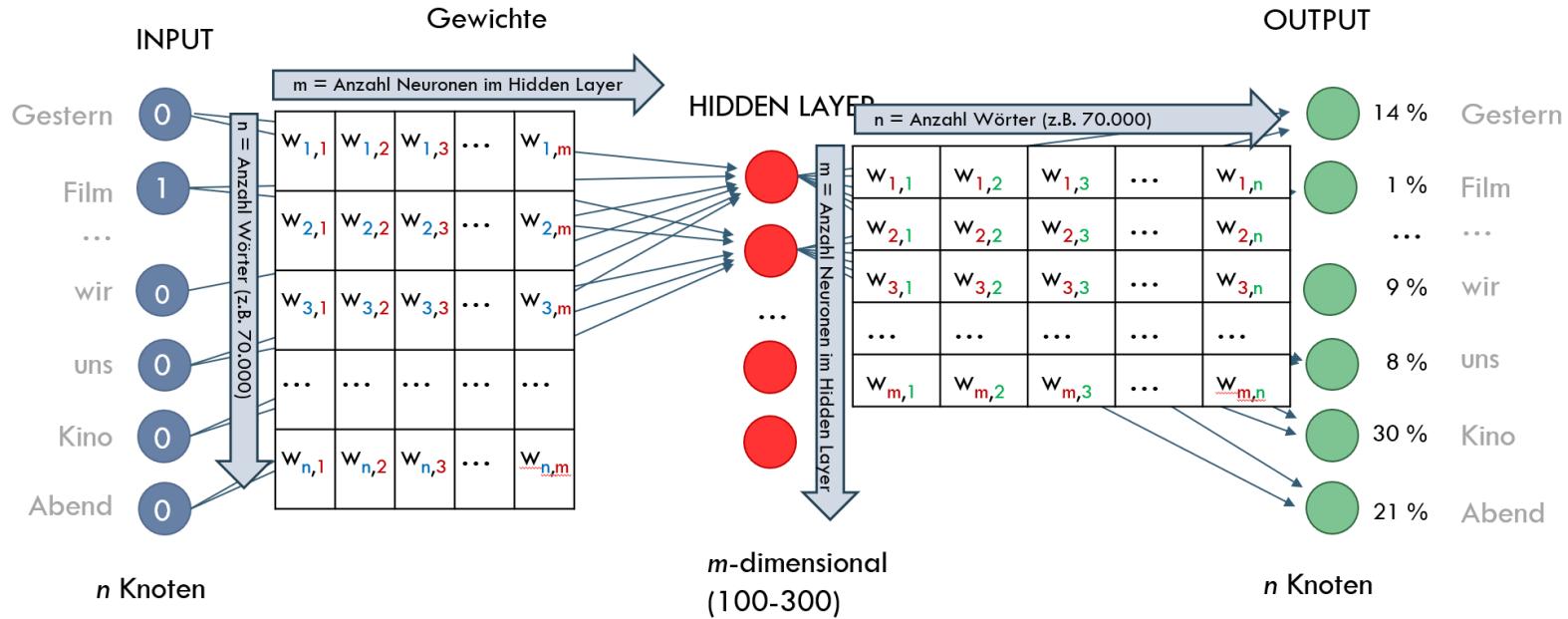
Skip-gram

Task: Given a word, predicts its context words



„You shall know a word
by the company it keeps“

Learning the word vectors



Disadvantages of Word2Vec

- Using Word2Vec, any given word in a vocabulary (e.g., “go”) has its own word vector, and those vectors are effectively stored in a lookup table or dictionary
- Word2Vec **does not address polysemy**, or the co-existence of many possible meanings for a given word or phrase
 - “go” is a verb and it is also a board game
- The **meaning** of a given word type **varies according to its context**
- In recent years, several models have been presented that learn **context-dependent** word vectors (e.g., ELMo, BERT)

- **Context independent models** (e.g., Word2vec, Glove)
 - these models output just one vector (embedding) for each word regardless of where the words occur in a sentence and regardless of the different meanings they may have
- **Context dependent models** (e.g., ELMo, BERT)
 - these models can generate different word embeddings for a word that **captures the context** of a word
 - they take the entire input sentence into equation for calculating the word embeddings

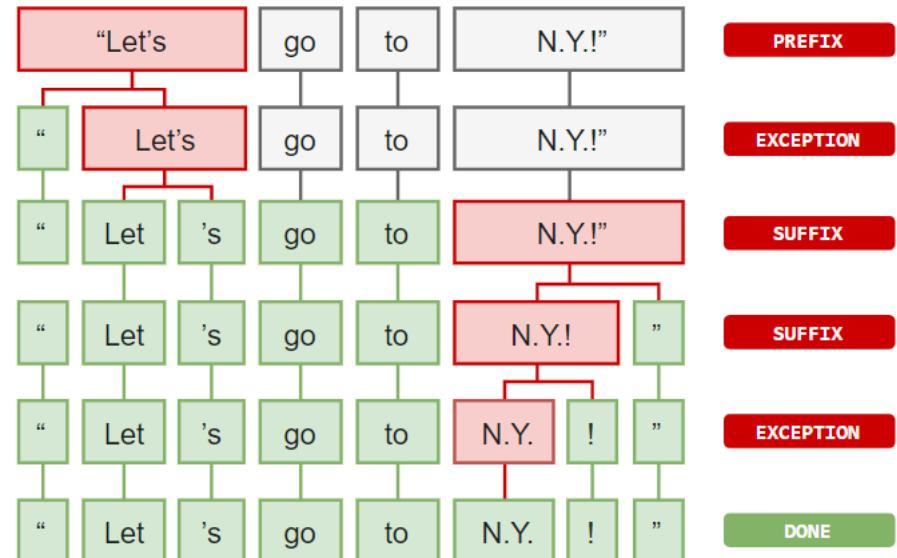
Sentence Embeddings

- Simple baseline: Averaging Word Embeddings
 - The simplest way to calculate sentence embeddings is to calculate the average word embeddings among all words in a sentence
 - Use a *simple* model (e.g., Word2Vec, Glove) to get word embeddings
 - Even though the approach is simple, the results are satisfying
- Deep Learning Approaches

Normalization of Natural Language Texts

Tokenization

- Tokenization: The task of converting a text from a single string to a list of tokens.
- It is harder than it seems:
 - I'll see you in New York.
 - The aluminum-export ban.
- Approaches
 - separate words using whitespace
 - use regular expressions to specify which substrings are valid words



Source: Spacy.io

Stopword Removal

- Stopwords: very common words, „without meaning“ (in the specific context) or of little value for the specific task
- Benefit of removing:
 - dataset size decreases
 - can increase accuracy, because fewer and only meaningful tokens are left
 - can increase speed
- Avoid stopword removal for
 - machine translation
 - language modeling
 - text summarization
 - question-answering problems

Stopword Removal

i	they	having	until	off	nor
me	them	do	while	over	not
my	their	does	of	under	only
myself	theirs	did	at	again	own
we	themselves	doing	by	further	same
our	what	a	for	then	so
ours	which	an	with	once	than
ourselves	who	the	about	here	too
you	whom	and	against	there	very
your	this	but	between	when	s
yours	that	if	into	where	t
yourself	these	or	through	why	can
yourselves	those	because	during	how	will
he	am	as	before	all	just
him	is	until	after	any	don
his	are	while	above	both	should
himself	was	of	below	each	now
she	were	at	to	few	
her	be	by	from	more	
hers	been	for	up	most	
herself	being	with	down	other	
it	have	about	in	some	
its	has	against	out	such	
itself	had	between	on	no	

Customization
of stopword
lists may be
important for
some
applications.

English stopwords in NLTK

<https://gist.github.com/sebleier/554280>

Word normalization

- Putting words / tokens in a standard format

e.g.

Usa, US, USA → USA

uh-huh, uhhuh → uhhuh

Case folding

- Mapping everything to lower case
- Helpful in tasks like information retrieval or speech recognition
- Not recommended for tasks like text classification, information extraction or machine translation (case = important information!)

Stemming

- Naive version of morphological analysis
- Chops off word-final affixes
- Example:

Pierre Vinken, 61 years old, will join the board as a nonexecutive director Nov. 29. Mr. Vinken is chairman of Elsevier N.V., the Dutch publishing group. Rudolph Agnew, 55 years old and former chairman of Consolidated Gold Fields PLC, was named a nonexecutive director of this British industrial conglomerate. A form of asbestos once used to make Kent cigarette filters has caused a high percentage of cancer deaths among a group of workers exposed to it more than 30 years ago, researchers reported.

is turned into:

Pierr Vinken, 61 year old, will join the board as a nonexecut director Nov. 29. Mr. Vinkén is chairman of Elsevi N.V., the Dutch publish group. Rúdolph Agnew, 55 year old and former cháirman of Consolid Gold Field PLC, wa name a nonexecut director of thi British industri conglomérat. A fórm of asbestos onc use to make Kent cigarett filter ha caus a high percentag of cancer death among a group of worker expos to it more than 30 year ago, research report .

Example rules of
Porter Stemmer:

sses → ss

ss → ss

s → -

ies → i

y → i

Lemmatization

- Determine if two words have the same root
- Examples:
 - am, are, is → be
 - dinner, dinners → dinner
 - He is reading detective stories → He be read detective story
- Lemmatization is done by morphological parsing
- Lemmatization algorithms can be complex!

Example

Word	Inflection	Stem	Morphological information	Lemma
study	-y	stud	Infinitive of the verb “study”	study
studies	-ies	stud	Third person, singular, Present Simple of the verb “study”	study
studying	-ing	stud	Gerund of the verb “study”	study

<https://chatbotsmagazine.com/how-to-use-nlp-for-building-a-chatbot-fac05476a58e?source=-----9-----&gi=12a9d7ac3f50>

Sentence Segmentation

- Cues are punctuation, like periods, question marks, and exclamation marks
 - But: Periods are ambiguous!

The period character is ambiguous. It can be a sentence boundary marker but also a marker of abbreviations like Mr. or Inc.

 Abbreviation  Sentence boundary marker & Abbreviation  Sentence boundary marker

- How solved?
 - Use rules or ML to decide whether a period is part of a word or a sentence-boundary marker.
 - Use abbreviation dictionaries

spaCy uses the dependency parser for sentence segmentation

Feature Extraction from Natural Language Text

Parts of speech (POS)

aka word classes, syntactic categories

- Noun
 - Verb
 - Adjective / Adverb
 - Pronoun
 - Preposition
 - Conjunction
 - Participle
 - Article
-
- The diagram illustrates the classification of parts of speech. On the left, a vertical list of nine items is shown, each preceded by a blue square bullet point. To the right of this list, two blue curly braces group the items into two categories: 'open class types' and 'closed class types'. The first brace groups the first three items (Noun, Verb, Adjective / Adverb). The second brace groups the remaining six items (Pronoun, Preposition, Conjunction, Participle, Article).

What does Tagging do?

1. Collapses Distinctions

- Lexical identity may be discarded
- e.g., all personal pronouns tagged with PRP

2. Introduces Distinctions

- Ambiguities may be resolved
- e.g. *deal* tagged with NN or VB

3. Helps in classification and prediction

Terminology

- **Tagging**
 - The process of associating labels with each token in a text
- **Tags**
 - The labels
- **Tag Set**
 - The collection of tags used for a particular task
 - Tag Sets are defined by linguists and not unambiguous. Different languages require different tag sets -> difficult to compare tags of different languages.

Examples

The/DT grand/JJ jury/NN commented/VBD on/IN a/DT number/NN of/IN other/JJ topics/NNS ./.

There/EX are/VBP 70/CD children/NNS there/RB

Preliminary/JJ findings/NNS were/VBD reported/VBN in/IN today/NN 's/POS New/NNP England/NNP Journal/NNP of/IN Medicine/NNP ./.

Tag	Description	Example	Tag	Description	Example	Tag	Description	Example
CC	coordinating conjunction	<i>and, but, or</i>	PDT	predeterminer	<i>all, both</i>	VBP	verb non-3sg present	<i>eat</i>
CD	cardinal number	<i>one, two</i>	POS	possessive ending	's	VBZ	verb 3sg pres	<i>eats</i>
DT	determiner	<i>a, the</i>	PRP	personal pronoun	<i>I, you, he</i>	WDT	wh-determ.	<i>which, that</i>
EX	existential 'there'	<i>there</i>	PRP\$	possess. pronoun	<i>your; one's</i>	WP	wh-pronoun	<i>what, who</i>
FW	foreign word	<i>mea culpa</i>	RB	adverb	<i>quickly</i>	WPS	wh-possess.	<i>whose</i>
IN	preposition/ subordin-conj	<i>of, in, by</i>	RBR	comparative	<i>faster</i>	WRB	wh-adverb	<i>how, where</i>
JJ	adjective	<i>yellow</i>	RBS	superlatv. adverb	<i>fastest</i>	\$	dollar sign	\$
JJR	comparative adj	<i>bigger</i>	RP	particle	<i>up, off</i>	#	pound sign	#
JJS	superlative adj	<i>wildest</i>	SYM	symbol	<i>+, %, &</i>	"	left quote	' or "
LS	list item marker	<i>1, 2, One</i>	TO	"to"	<i>to</i>	"	right quote	' or "
MD	modal	<i>can, should</i>	UH	interjection	<i>ah, oops</i>	(left paren	[, (, {, <
NN	sing or mass noun	<i>llama</i>	VB	verb base form	<i>eat</i>)	right paren	I,), }, >
NNS	noun, plural	<i>llamas</i>	VBD	verb past tense	<i>ate</i>	,	comma	,
NNP	proper noun, sing.	<i>IBM</i>	VBG	verb gerund	<i>eating</i>	.	sent-end punc	. ! ?
NNPS	proper noun, plu.	<i>Carolinas</i>	VBN	verb past part.	<i>eaten</i>	:	sent-mid punc	: ; ... - -

Figure 8.1 Penn Treebank part-of-speech tags (including punctuation).

Some of the best-known Tagsets

- Brown corpus: 87 tags
 - (more when tags are combined)
- Penn Treebank: 45 tags
- Lancaster UCREL C5 (used to tag the BNC): 61 tags
- Lancaster C7: 145 tags

For the German language:

- Stuttgart-Tübingen Tagset (STTS): 54 tags

The Tagging Process

Training data
(Annotated text)

<i>This</i>	<i>sentence</i>	<i>serves</i>	<i>as</i>	<i>an</i>	<i>example</i>	<i>of</i>	<i>annotated</i>	<i>text...</i>
Det	N	V1	P	Det	N	P	V2	N

"This is a new sentence."

POS Tagger

This is a new sentence.

Det	Aux	Det	Adj	N
-----	-----	-----	-----	---

Pick the **most likely** tag sequence.

$$p(w_1, \dots, w_k, t_1, \dots, t_k) = \frac{p(t_1 | w_1) \dots p(t_k | w_k) p(w_1) \dots p(w_k)}{\prod_{i=1}^k p(w_i | t_i) p(t_i | t_{i-1})}$$

Independent assignment
Most common tag

Partial dependency
(HMM)

Training and Testing of Learning Algorithms

- Algorithms that “learn” from data see a set of examples and try to generalize from them.
- Training set:
 - Examples trained on
- Test set:
 - Also called held-out data and unseen data
 - Use this for evaluating your algorithm
 - Must be separate from the training set
 - Otherwise, you cheated!
- “Gold” standard
 - A test set that a community has agreed on and uses as a common benchmark.

Three main corpora tagged with POS for English: Brown (different genres), WSJ (Wall Street Journal), Switchboard (phone conversations)

Syntactic parsing

Syntactic parsing is the task of recognizing a sentence and assigning a syntactic structure to it.

Comparison of POS tagging, chunking and full parsing

He reckons the current account deficit will narrow to only 2 billion in September.

Output of POS tagger:

He\PRP reckons\VBZ the\DT current\JJ account\NN deficit\NN will\MD narrow\VB to\TO only\RB 2\CD billion\CD in\IN September\NNP .\.

Comparison of POS tagging, chunking and full parsing

He reckons the current account deficit will narrow to only 2 billion in September.

Output of chunker:

[NP He] [VP reckons] [NP the current account deficit] [VP will narrow] [PP to]
[NP only 2 billion] [PP in] [NP September].

Constituency

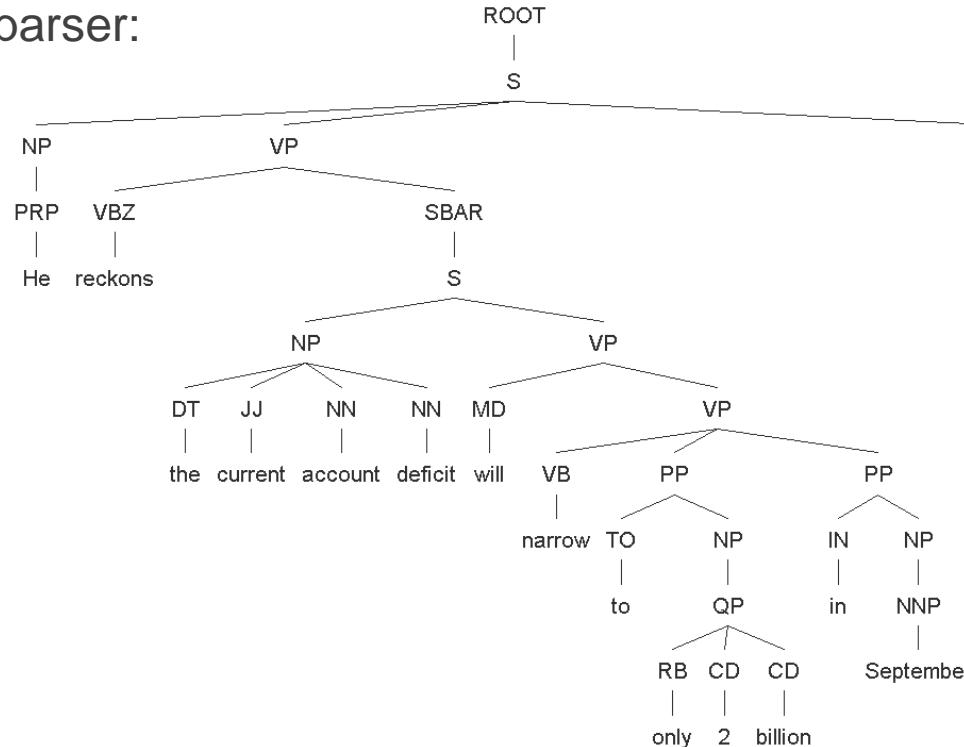
- Constituents: groups of words behaving as single units
- Example:
 - *On September seventeenth*, I'd like to fly from Atlanta to Denver
 - I'd like to fly *on September seventeenth* from Atlanta to Denver
 - I'd like to fly from Atlanta to Denver *on September seventeenth*
- What about
 - On September, I'd like to fly seventeenth from Atlanta to Denver
 - On I'd like to fly September seventeenth from Atlanta to Denver
 - I'd like to fly on September from Altanta to Denver seventeenth

What about German?
Determine the constituents *in the following sentence*:

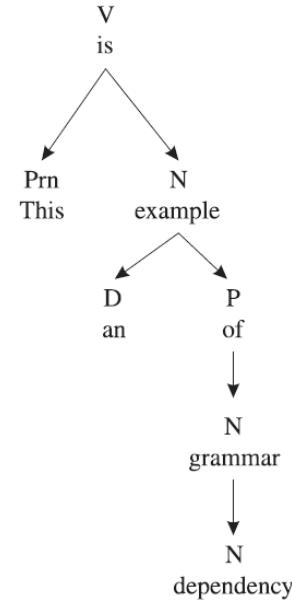
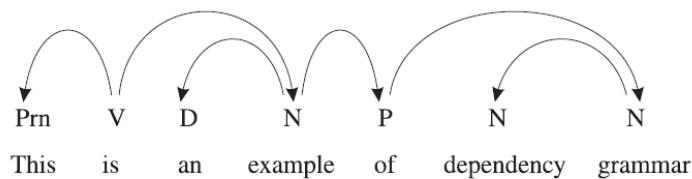
„Ich sah gestern Abend einen schönen Film im neuen Kino unserer Stadt.“

Comparison of POS tagging, chunking and full parsing

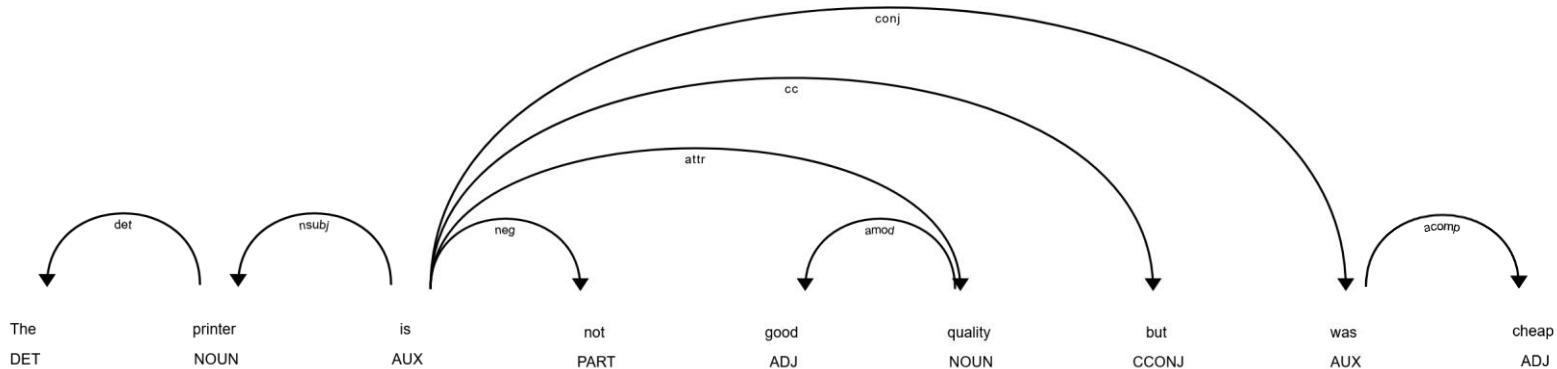
Output of full parser:



Dependency Parsing



Example with negation



Application 1: Grammar checking

- Check if sentence can be parsed (full parse tree)
- Yes: likely grammatical
- No: likely ungrammatical or at least difficult to read

Application 2: Question answering

- Example question:
„What books were written by British women authors before 1800?“
- Parsing helps to determine if the user wants a list of books or authors

Application 3: Information extraction systems

- With chunking, only segments in a text that contain valuable information are identified and classified.

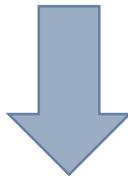
Application 4: Information Retrieval

- Using chunking, Information Retrieval systems may index texts according to a subset of the constituents found in them

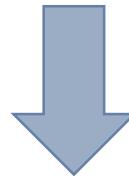
Information Extraction

Information Extraction

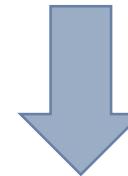
Information Extraction (IE) is the **process of extracting structured information** (e.g., database tables) from unstructured machine-readable documents (e.g., Web documents).



Named Entity
Recognition



Relation extraction



Event extraction

Named Entities

- Roughly speaking anything that can be referred to with **a proper name**
- E.g., **person, location, organization**
- Commonly also **dates, times, numerical expressions, ...**
- Can also include **domain-specific entities** like protein names or commercial products

Example

Citing high fuel prices, [ORG United Airlines] said [TIME Friday] it has increased fares by [MONEY \$6] per round trip on flights to some cities also served by lower-cost carriers. [ORG American Airlines], a unit of [ORG AMR Corp.], immediately matched the move, spokesman [PER Tim Wagner] said. [ORG United], a unit of [ORG UAL Corp.], said the increase took effect [TIME Thursday] and applies to most routes where it competes against discount carriers, such as [LOC Chicago] to [LOC Dallas] and [LOC Denver] to [LOC San Francisco].

A list of generic named entity types

Type	Tag	Sample Categories	Example sentences
People	PER	people, characters	Turing is a giant of computer science.
Organization	ORG	companies, sports teams	The IPCC warned about the cyclone.
Location	LOC	regions, mountains, seas	The Mt. Sanitas loop is in Sunshine Canyon.
Geo-Political Entity	GPE	countries, states, provinces	Palo Alto is raising the fees for parking.
Facility	FAC	bridges, buildings, airports	Consider the Golden Gate Bridge.
Vehicles	VEH	planes, trains, automobiles	It was a classic Ford Falcon.

Families of algorithms for NER tagging

- Feature based ML systems
- Neural (bi-LSTM)
- Rule-based

Input for feature based NER systems

Typical features for a feature-based NER system

- identity of w_i , identity of neighboring words
- embeddings for w_i , embeddings for neighboring words
- part of speech of w_i , part of speech of neighboring words
- base-phrase syntactic chunk label of w_i and neighboring words
- presence of w_i in a **gazetteer**
- w_i contains a particular prefix (from all prefixes of length ≤ 4)
- w_i contains a particular suffix (from all suffixes of length ≤ 4)
- w_i is all upper case
- word shape of w_i , word shape of neighboring words
- short word shape of w_i , short word shape of neighboring words
- presence of hyphen

Deep dive: Word shape features

- Abstract letter patterns of a word
- Mapping
 - lower-case letters → x
 - upper-case letters → X
 - numbers → d
 - punctuation → retained
- Example:
 - DC10-30 → XXdd-dd
- Short version: remove consecutive character types
 - DC10-30 → XXdd-dd → Xd-d

Deep-dive: Gazetteer

- List of place names
- Related: name-lists / corporations / commercial products...
- Typically implemented as binary feature (in / not in gazetteer)
- Problems:
 - can be difficult to create and maintain
 - usefulness varies considerably

Which features are best?

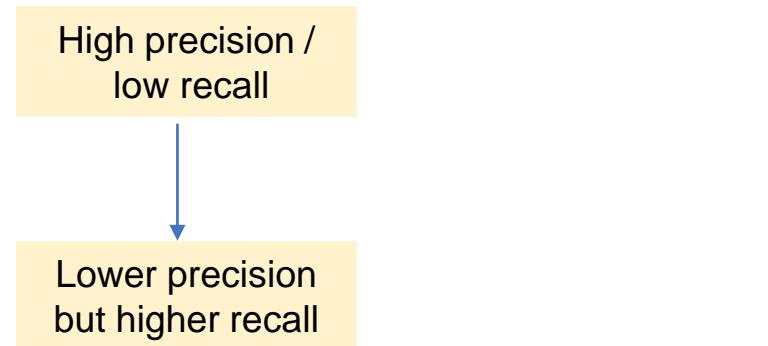
- Which features are best?
This also depends on the application, genre, media, and language

identity of w_i , identity of neighboring words
embeddings for w_i , embeddings for neighboring words
part of speech of w_i , part of speech of neighboring words
base-phrase syntactic chunk label of w_i and neighboring words
presence of w_i in a **gazetteer**
 w_i contains a particular prefix (from all prefixes of length ≤ 4)
 w_i contains a particular suffix (from all suffixes of length ≤ 4)
 w_i is all upper case
word shape of w_i , word shape of neighboring words
short word shape of w_i , short word shape of neighboring words
presence of hyphen

Which features would you use when applying NER to Twitter posts?

Rule-based NER

- Commercial approaches often are based on pragmatic combinations of lists and rules (with smaller amount of ML)
- Often repeated rule-based passes over a text
→ Output of one pass becomes input of next one



Natural Language Processing with spaCy

Acknowledgements

The following slides are based on material provided by

- spacy.io

What's spaCy?

- a free, open-source library for NLP in Python
- can be used for
 - text preprocessing
 - information extraction
 - natural language understanding systems

spaCy vs. nltk

nltk

- created with focus on teaching and research
- user can choose between multiple algorithms with equivalent functionality

spaCy

- created with focus on production-use
- choice has been made by spaCy-developers (easier for developers + performance-optimization)

Features

NAME	DESCRIPTION	NAME	DESCRIPTION
Tokenization	Segmenting text into words, punctuations marks etc.	Entity Linking (EL)	Disambiguating textual entities to unique identifiers in a Knowledge Base.
Part-of-speech (POS) Tagging	Assigning word types to tokens, like verb or noun.	Similarity	Comparing words, text spans and documents and how similar they are to each other.
Dependency Parsing	Assigning syntactic dependency labels, describing the relations between individual tokens, like subject or object.	Text Classification	Assigning categories or labels to a whole document, or parts of a document.
Lemmatization	Assigning the base forms of words. For example, the lemma of "was" is "be", and the lemma of "rats" is "rat".	Rule-based Matching	Finding sequences of tokens based on their texts and linguistic annotations, similar to regular expressions.
Sentence Boundary Detection (SBD)	Finding and segmenting individual sentences.	Training	Updating and improving a statistical model's predictions.
Named Entity Recognition (NER)	Labelling named "real-world" objects, like persons, companies or locations.	Serialization	Saving objects to files or byte strings.

Statistical (ML) models

- Some features require a **statistical model** to be loaded
→ **prediction** of linguistic annotations

Which statistical models are available?

- Different languages (as of 9/2020: 16)
- Trained on different corpora (web, news, ...)
- Different sizes, speed, memory usage and accuracy

→ You have to choose and load the right model!

<https://spacy.io/usage/models>

How do I use models?

1. Download and install the required model
(It's a Python package)
<https://spacy.io/usage/models>

```
# Download best-matching version of specific model for your spaCy installation
python -m spacy download en_core_web_sm

# Out-of-the-box: download best-matching default model and create shortcut
python -m spacy download en

# Download exact model version (doesn't create shortcut link)
python -m spacy download en_core_web_sm-2.2.0 --direct
```

How do I use models? (2)

2. Loading models

```
import spacy
```

```
# load model package "en_core_web_sm"  
nlp = spacy.load("en_core_web_sm")
```

```
# load package from a directory  
nlp = spacy.load("/path/to/en_core_web_sm")
```

```
# load model with shortcut link "en"  
nlp = spacy.load("en")
```

How do I use models? (3)

3. Using the model

```
nlp = spacy.load(„...“)
```

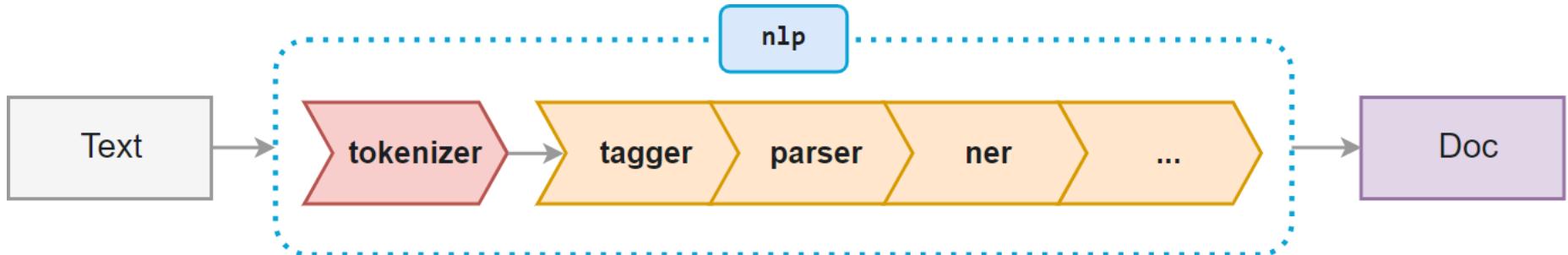
Returns a Language
object, usually called „nlp“

```
doc = nlp("This is a sentence.")
```

Returns a
Doc object

Can be called
with a string.

Processing pipeline



Lemmatization

```
import spacy

nlp = spacy.load("en_core_web_sm")
doc = nlp("Apple is looking at buying U.K. startup for $1 billion")

for token in doc:
    print(token.text, token.lemma_)
```

TEXT LEMMA

Apple	apple
is	be
looking	look
at	at
buying	buy
U.K.	u.k.
startup	startup
for	for
\$	\$
1	1
billion	billion

Text: The original word text.

Lemma: The base form of the word.

Part-of-speech (POS) tagging

```
import spacy

nlp = spacy.load("en_core_web_sm")
doc = nlp("Apple is looking at buying U.K. startup for $1 billion")

for token in doc:
    print(token.text, token.lemma_, token.pos_, token.tag_)
```

TEXT	LEMMA	POS	TAG
Apple	apple	PROPN	NNP
is	be	AUX	VBZ
looking	look	VERB	VBG
at	at	ADP	IN
buying	buy	VERB	VBG
U.K.	u.k.	PROPN	NNP
startup	startup	NOUN	NN
for	for	ADP	IN
\$	\$	SYM	\$
1	1	NUM	CD
billion	billion	NUM	CD

POS: The simple UPOS part-of-speech tag.

Tag: The detailed part-of-speech tag.

Dependency Parsing

Dep: Syntactic dependency, i.e. the relation between tokens.

```
import spacy

nlp = spacy.load("en_core_web_sm")

doc = nlp("Apple is looking at buying U.K. startup for $1 billion")

for token in doc:
    print(token.text, token.lemma_, token.pos_, token.tag_, token.dep_)
```

TEXT	LEMMA	POS	TAG	DEP
Apple	apple	PROPN	NNP	nsubj
is	be	AUX	VBZ	aux
looking	look	VERB	VBG	ROOT
at	at	ADP	IN	prep
buying	buy	VERB	VBG	pcomp
U.K.	u.k.	PROPN	NNP	compound
startup	startup	NOUN	NN	dobj
for	ADP	IN	IN	prep
\$	SYM	SYM	\$	quantmod
1	NUM	NUM	CD	compound
billion	billion	NUM	CD	pobj

```
from spacy import displacy
displacy.render(doc, style = "dep")
```

More features

```
import spacy

nlp = spacy.load("en_core_web_sm")
doc = nlp("Apple is looking at buying U.K. startup for $1 billion")

for token in doc:
    print(token.text, token.lemma_, token.pos_, token.tag_, token.dep_,
          token.shape_, token.is_alpha, token.is_stop)
```

Shape: The word shape – capitalization, punctuation, digits.

is alpha: Is the token an alpha character?

is stop: Is the token part of a stop list, i.e. the most common words of the language?

TEXT	LEMMA	POS	TAG	DEP	SHAPE	ALPHA	STOP
Apple	apple	PROPN	NNP	nsubj	Xxxxx	True	False
is	be	AUX	VBZ	aux	xx	True	True
looking	look	VERB	VBG	ROOT	xxxx	True	False
at	at	ADP	IN	prep	xx	True	True
buying	buy	VERB	VBG	pcomp	xxxx	True	False
U.K.	u.k.	PROPN	NNP	compound	X.X.	False	False
startup	startup	NOUN	NN	dobj	xxxx	True	False
for	for	ADP	IN	prep	xxx	True	True
\$	\$	SYM	\$	quantmod	\$	False	False
1	1	NUM	CD	compound	d	False	False
billion	billion	NUM	CD	pobj	xxxx	True	False

Named Entities

Named Entities = „real-world objects“ that have been assigned a name (i.e. a person, a country, a product, a book title...)

```
import spacy

nlp = spacy.load("en_core_web_sm")
doc = nlp("Apple is looking at buying U.K. startup for $1 billion")

for ent in doc.ents:
    print(ent.text, ent.start_char, ent.end_char, ent.label_)
```

TEXT	START	END	LABEL	DESCRIPTION
Apple	0	5	ORG	Companies, agencies, institutions.
U.K.	27	31	GPE	Geopolitical entity, i.e. countries, cities, states.
\$1 billion	44	54	MONEY	Monetary values, including unit.

```
from spacy import displacy
displacy.render(doc, style = "ent")
```



Apple **ORG** is looking at buying **U.K. GPE** startup for **\$1 billion MONEY**

Result visualized with displaCy visualizer

Word Vectors / Word embeddings

- multi-dimensional meaning representation of a word
- e.g. generated with word2vec
- only available in larger models!
- available on the level of Tokens, Documents and Spans
- out-of vocabulary words do not have a word vector

BANANA.VECTOR

```
array([ 2.02280000e-01, -7.66180009e-02,  3.70319992e-01,
       3.28450017e-02, -4.19569999e-01,  7.20689967e-02,
      -3.74760002e-01,  5.74599989e-02, -1.24009997e-02,
       5.29489994e-01, -5.23800015e-01, -1.97710007e-01,
      -3.41470003e-01,  5.33169985e-01, -2.53309999e-02,
       1.73800007e-01,  1.67720005e-01,  8.39839995e-01,
       5.51070012e-02,  1.05470002e-01,  3.78719985e-01,
       2.42750004e-01,  1.47449998e-02,  5.59509993e-01,
       1.25210002e-01, -6.75960004e-01,  3.58420014e-01,
       # ... and so on ...
       3.66849989e-01,  2.52470002e-03, -6.40089989e-01,
      -2.97650009e-01,  7.89430022e-01,  3.31680000e-01,
      -1.19659996e+00, -4.71559986e-02,  5.31750023e-01], dtype=float32)
```

Word Similarity

- Doc, Span and Token come with a method `.similarity()`

```
import spacy

nlp = spacy.load("en_core_web_md") # make sure to use larger model!
tokens = nlp("dog cat banana")

for token1 in tokens:
    for token2 in tokens:
        print(token1.text, token2.text, token1.similarity(token2))
```

```
dog dog 1.0
dog cat 0.80168545
dog banana 0.24327643
cat dog 0.80168545
cat cat 1.0
cat banana 0.28154364
banana dog 0.24327643
banana cat 0.28154364
banana banana 1.0
```

Exercise 2: spaCy and NER

Named Entities

The EIB ORG will finance around 40 CARDINAL projects in Barcelona GPE that aim to support climate change mitigation and adaptation in the city. To this end, the EU ORG bank will provide €95 million MONEY to promote urban regeneration, with a focus on the environment but also on social inclusion and job creation to boost the economic recovery in the wake of the COVID-19 crisis.

The Italian NORP healthcare system is also being reinforced to tackle the emergency situation caused by the COVID-19 pandemic. This is being conducted with the backing of the EU ORG bank, the EIB ORG, which is providing the Italian NORP government with a €2 billion MONEY loan covering around two-thirds CARDINAL of the resources needed for the operations contained in the Decree GPE for revival of the healthcare system

The EIB Group ORG and Banco Santander Consumer Portugal ORG (BSCP) are joining forces to support Portuguese NORP small and medium-sized enterprises (SMEs) and mid-caps affected by the COVID-19 crisis. The EU ORG bank and BSCP ORG have signed two CARDINAL agreements to provide EUR CPE 587 million CARDINAL to inject liquidity and finance investments at a critical time.

Amadou Hott PERSON, Senegalese NORP Minister of the Economy, Planning and Cooperation ORG, and Ambroise Fayolle PERSON, EIB ORG Vice-President responsible for Africa LOC, today DATE formally agreed a CFAF 49 billion CARDINAL concessional loan to the Republic of Senegal GPE. It comes on top of a CFAF 200 billion CARDINAL financing mechanism established by Macky Sall PERSON, President of the Republic of Senegal GPE, as part of the country's Economic and Social Resilience Programme.

The health, economic and social impact of COVID-19 across Africa LOC continues to evolve with more than 12,000 CARDINAL deaths, in excess of 540,000 CARDINAL confirmed cases and the livelihoods of millions CARDINAL of Africans NORP damaged.

This has been one of the Italian NORP government's top priorities for at least three years DATE. The government has decidedly sped up the process following the restrictions caused by the COVID-19 pandemic, and this has now become a key pillar to aid the country's recovery. Digitalisation of public administrations (PAs) and their relations with the country's citizens and businesses have become the focus of the digital innovation measures contained in the "Simplifications" Decree Law approved by the Prime Minister's Office ORG on 7 July DATE. It is against this backdrop that the four-year

NER on the German Press releases

Die EIB-Gruppe **ORG** und die Banco Santander Consumer Portugal **ORG** (BSCP **ORG**) unterstützen gemeinsam von der Coronakrise betroffene kleine und mittlere Unternehmen (KMU) und Midcap-Unternehmen in **Portugal LOC** .

Die Bank der EU **ORG** und die BSCP **ORG** haben zwei Vereinbarungen unterzeichnet, um 587 Millionen Euro als Liquiditätshilfe und zur Finanzierung von Investitionen in einer kritischen **Zeit MISC** bereitzustellen.

Ganz **Afrika LOC** leidet zunehmend unter den gesundheitlichen, wirtschaftlichen und sozialen Folgen von **Covid-19 ORG**

Inzwischen sind auf dem Kontinent über 12&nbsnbsp;000&nbsnbsp;Menschen an dem **Virus MISC** gestorben, mehr als 540&nbsnbsp;000 haben sich infiziert, und für Millionen sind die Existenzgrundlagen gefährdet.

Seit mindestens drei Jahren steht die Digitalisierung der **öffentlichen Verwaltung ORG** ganz oben auf der Tagesordnung der **italienischen Regierung ORG** .

Aufgrund der **Covid-19-Beschränkungen MISC** wurde sie stark beschleunigt und bildet heute eine der Säulen für den Aufbau des **Landes LOC** .

Das Gesetzesdekrekt **Nr.&nbsnbsp;76 MISC** vom 16.&nbsnbsp;Juli 2020 über dringende Ma&srlig;nahmen zur Vereinfachung und zur digitalen Innovation zielt ebenfalls auf die Digitalisierung der öffentlichen Verwaltung und deren Beziehungen zu Bürgern und Unternehmen ab.

In diesem Rahmen ist die Tätigkeit des staatlichen Unternehmens **PagoPA&nbsnbsp;S.p LOC** . A. **PER** zu sehen, dessen Vier-Jahres-Investitionsplan die **EIB MISC** unterstützt.

Die **EIB MISC** unterstützt in der Autonomen Provinz **Trient LOC** die Konjunkturerholung nach der Coronapandemie und nachhaltige Projekte des öffentlichen **Sektors LOC** .

Heute wurde die entsprechende Vereinbarung unterzeichnet – die **EIB MISC** stellt 300 Millionen Euro für die norditalienische Provinz bereit, wovon eine erste Tranche über 160 Millionen Euro bereits abgeschlossen wurde.

Die **EIB MISC** hat heute neue Finanzierungen im Umfang von 16,6 Milliarden Euro für Projekte in **Europa LOC** und anderen Regionen der Welt genehmigt, davon mehr als zehn Milliarden Euro für Investitionen im Zusammenhang mit **Covid-19**.

NER on job posts

Extract data from a variety of sources (databases, web, text files) and in a variety of formats (structured, unstructured, JSON, Parquet GPE, etc.) and prepare the data for analysis. Describe data using qualitative and quantitative aggregation/summarization and build tools (static reports, interactive dashboards, web applications) that allow clients to easily view and understand important patterns and trends in the data. Provide quantitative comparisons of data points using advanced statistical analysis including data reduction techniques (PCA ORG, clustering), resampling (permutation tests, Monte Carlo PERSON testing) as well as traditional statistical hypothesis testing Construct ORG and validate predictive models using a variety of techniques.

Physical Actions ORG Physical Environment

Education Requirements

Master's degree in scientific or quantitative field and documented academic excellence. PhD in scientific field preferred

Experience Requirements

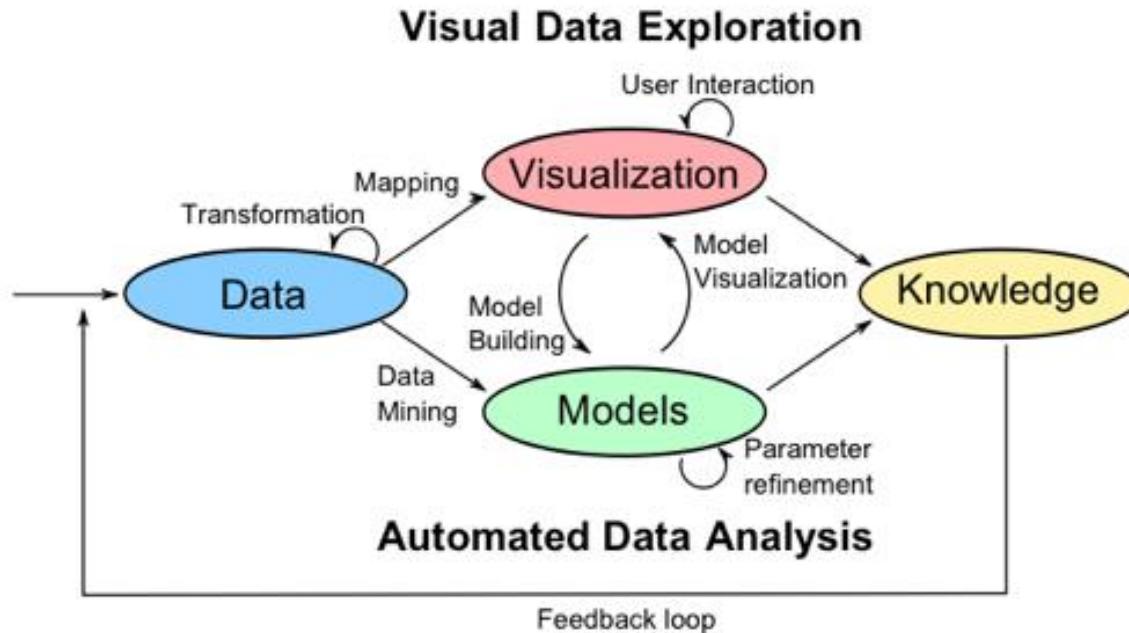
At least year DATE in a data analytics role (business intelligence analyst, data analyst, etc.)

Special Skill Requirement

Proficiency in at least one CARDINAL of the following is absolutely required: R, Python Basic ORG proficiency in at least one CARDINAL of the following: PowerBI, Tableau ORG, QlikView ORG, D3.js Basic proficiency in SQL Able ORG to implement and utilize advanced statistical and machine learning techniques (GLMs, PCA ORG, cluster analyses, ARIMA ORG, ETS ORG, decision trees, SVM ORG, neural networks, ensemble methods, etc.) Basic familiarity with project management tools (Asana PERSON, Basecamp PERSON, LiquidPlanner ORG)

Visual Analytics

Visual Analytics

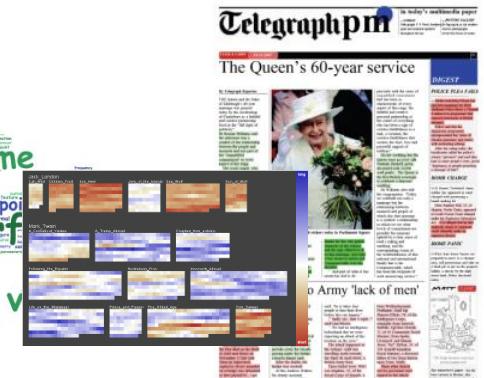
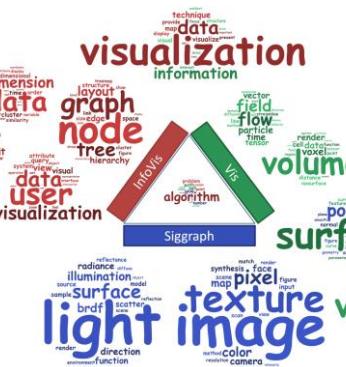
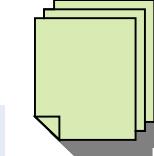
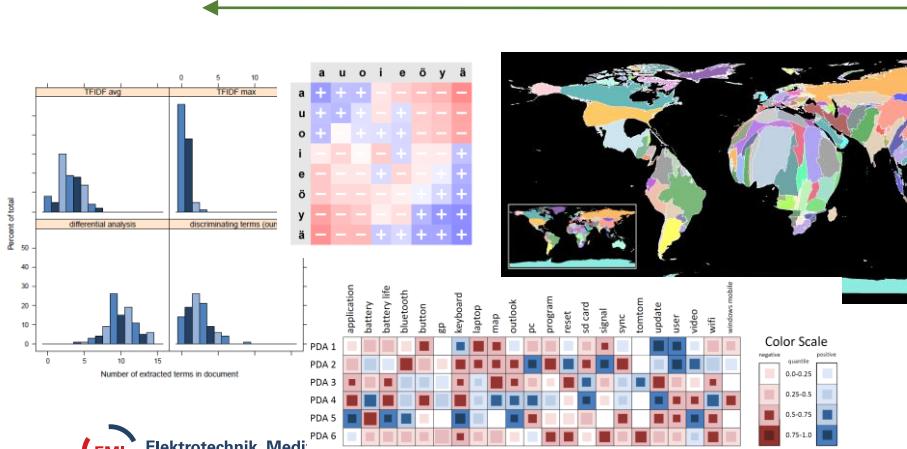


Documents

automatic
(pre)processing

completely structured

text properties are important

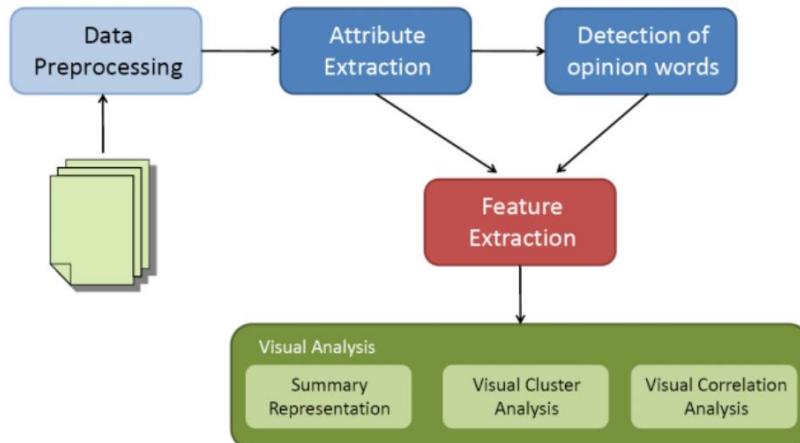


Analysis of customer reviews

Example review:

„I feel obligated to counter the bad reviews. This printer is just fine. I don't know what people are complaining about regarding the software, but it installed seamlessly and is intuitive in its operation. Even though the paper tray jams sometimes altogether I am happy that I bought this wonderful printer.“

Preprocessing of reviews



Step 1: Identification of attributes and sentiment

I feel obligated to counter the **bad** reviews. This **printer** is just **fine**. I don't know what people are **complaining** about regarding the **software**, but it installed **seamlessly** and is **intuitive** in its operation. Even though the **paper tray jams** sometimes altogether I am **happy** that I bought this **wonderful** **printer**.

Step 2: Mapping between attributes and sentiment

[...] Even though the **paper tray jams** sometimes altogether
I am **happy** that I bought this **wonderful** **printer**.

Step 3: Determining the overall sentiment of an attribute

[...] This **printer** is just **fine**. [...] Even though the paper tray jams sometimes altogether I am happy that I bought this wonderful **printer**.
→ Overall sentiment for **printer**: positive

Resulting Feature Vector

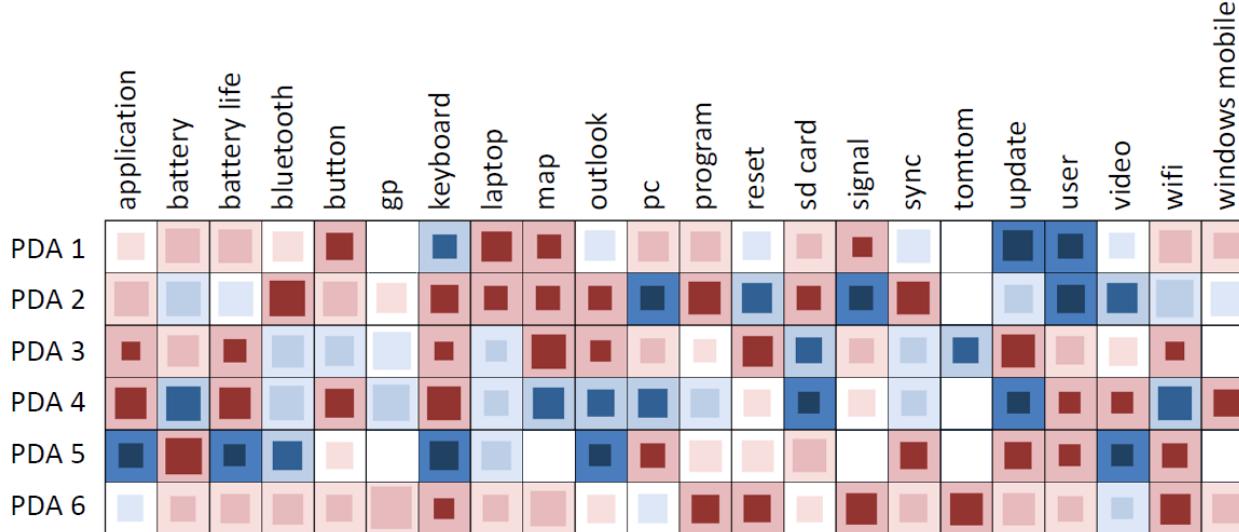
Printer	Ink	Software	Paper tray	Price
1	0	1	-1	0

Evaluation of the algorithm

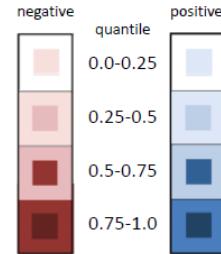
- Overall accuracy: 72 %
(only subjective sentences, the attributes were given)

Main Error Sources	In % of all Errors
1. Domain-dependency of opinion words	23.5%
2. Fixed expressions and phrases	21.6%
3. Errors in opinion word list	19.6%
4. Opinion word combinations	13.7%
5. But-clauses	7.8%
6. Attribute-dependency of opinion words	5.9%
7. Implicit attributes	2.0%

Review Analysis (example: printer reviews)



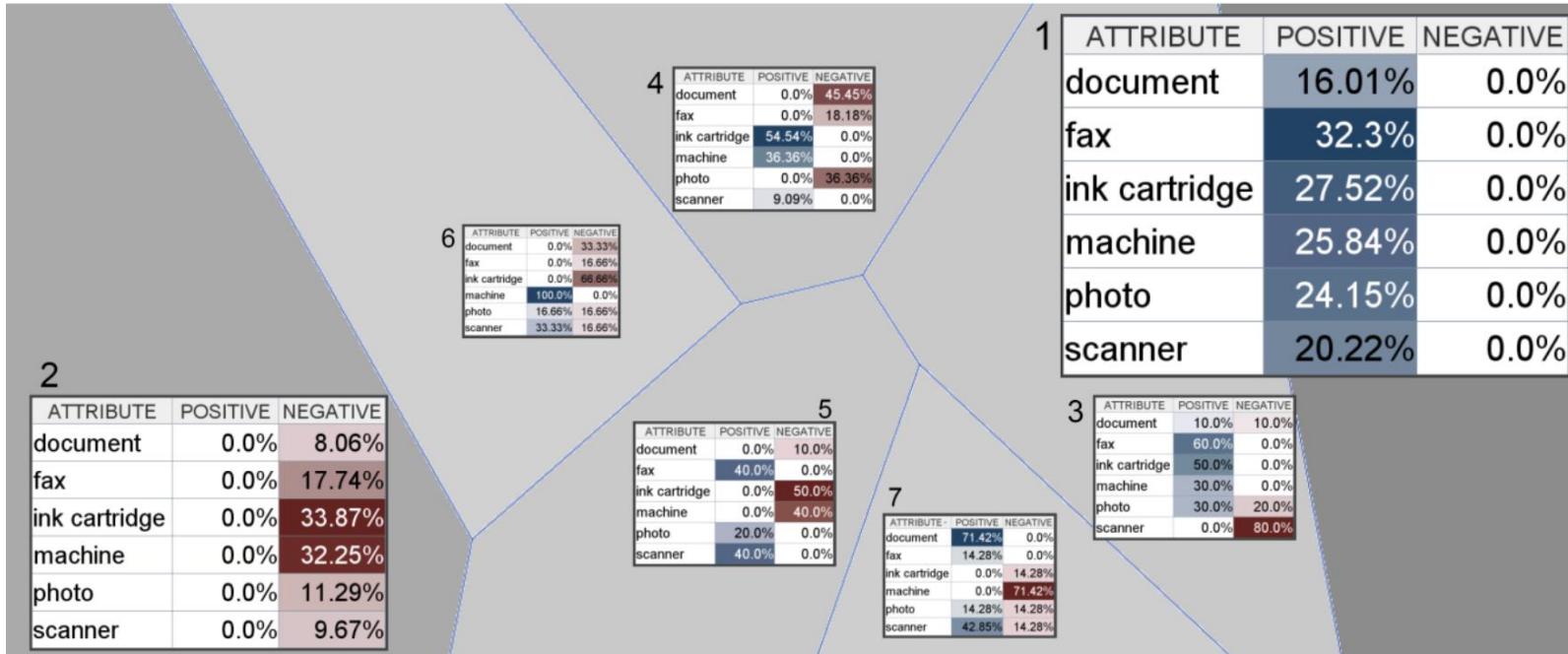
Color Scale



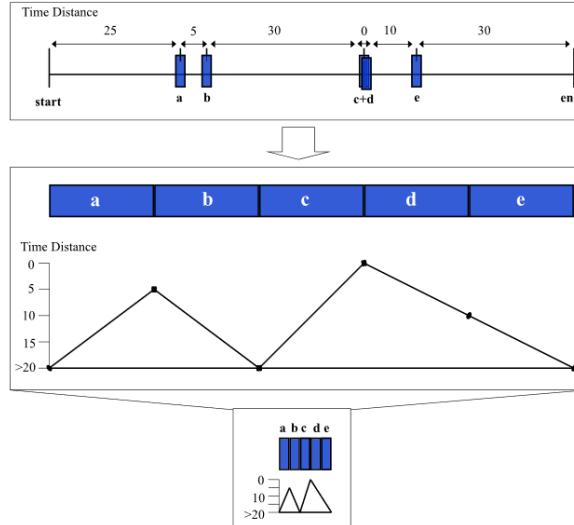
- **size inner rectangle:** amount of customers that commented on the attribute
- **blue color:** positive opinions
- **red color:** negative opinions
- **brightness of color:** degree of positiveness / negativeness of comments

D. Oelke, M. Hao, C. Rohrdantz, U. Dayal, L. Haug, H. Janetzko, "Visual opinion analysis of customer feedback data," *IEEE Symposium on Visual Analytics Science and Technology*, Atlantic City, NJ, 2009, pp. 187-194.

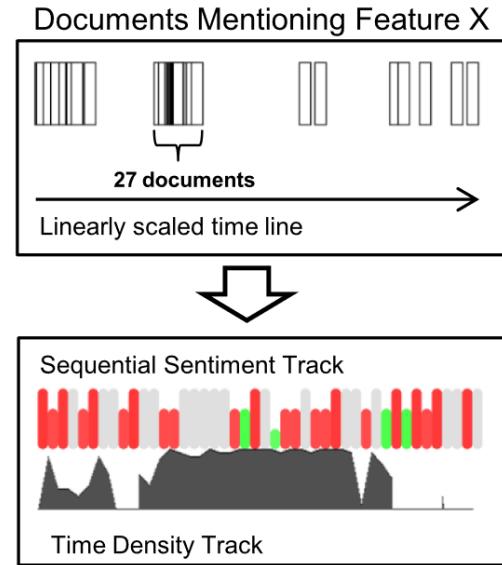
Opinion clusters



Review Analysis over time (Time Density Plots)



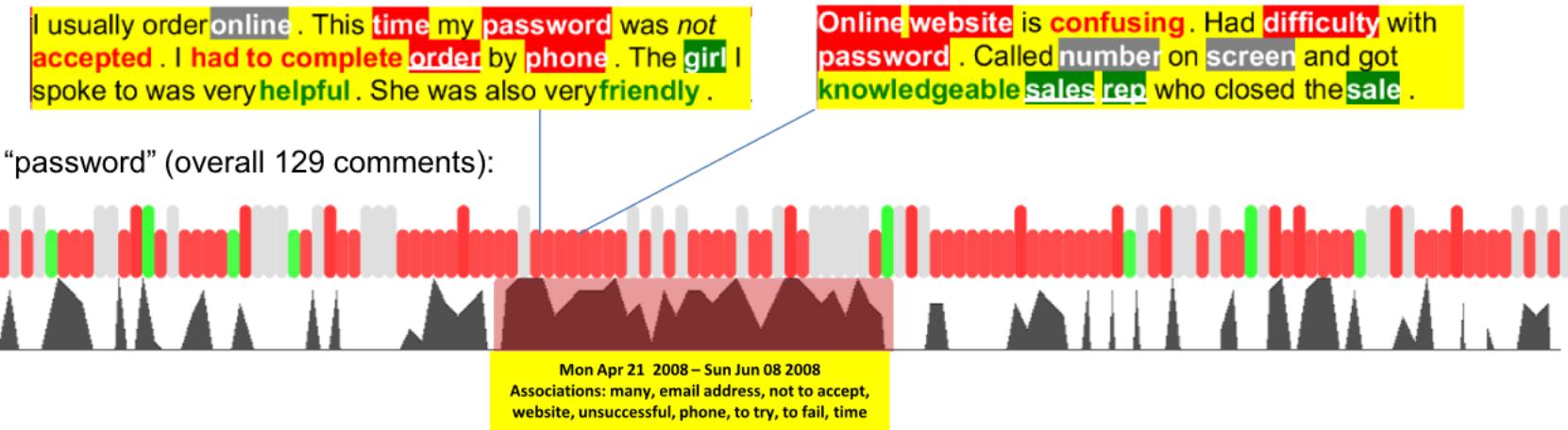
(a) details about the construction of time density curves



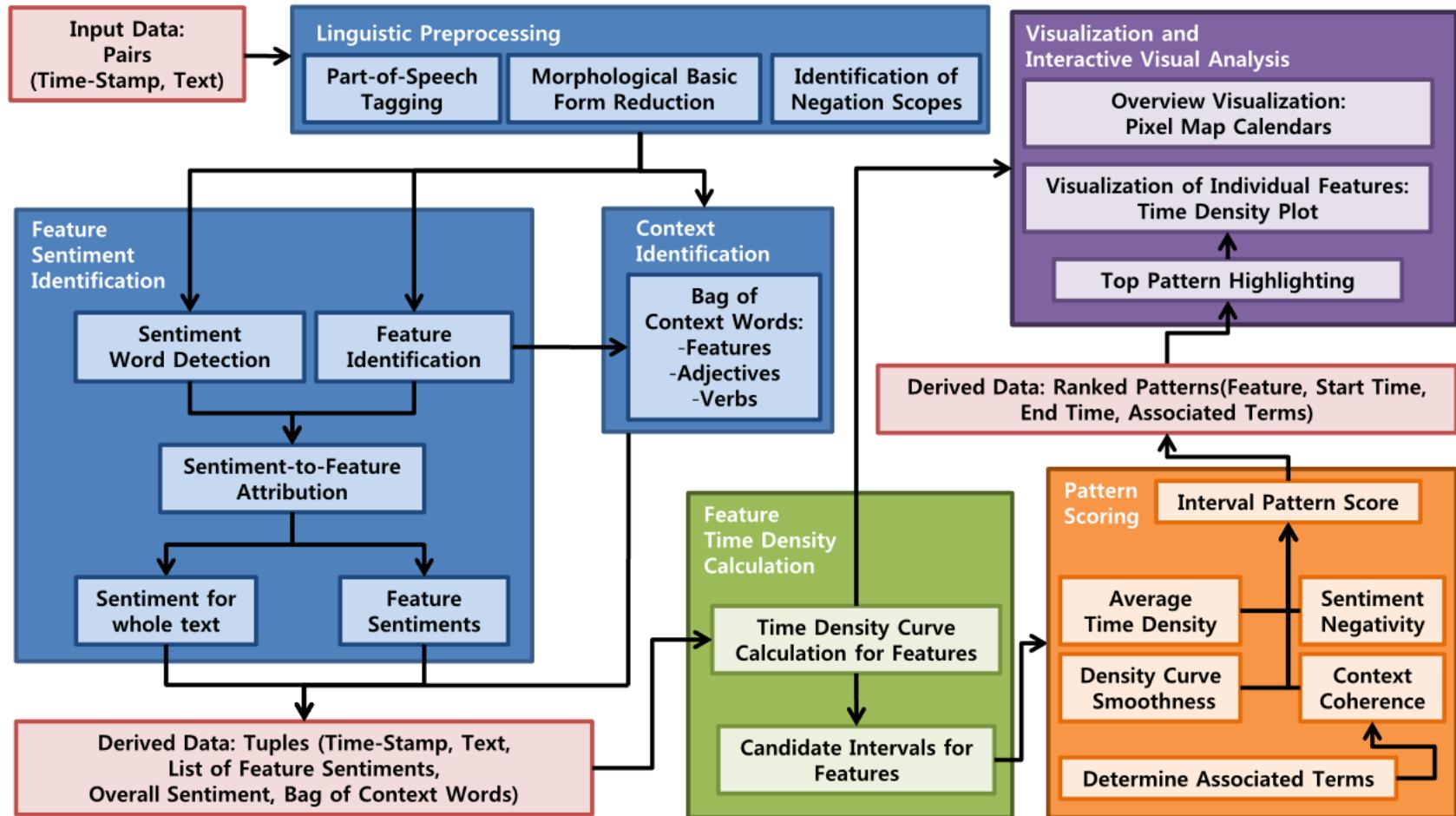
(b) example for a beneficial application of *time density plots*

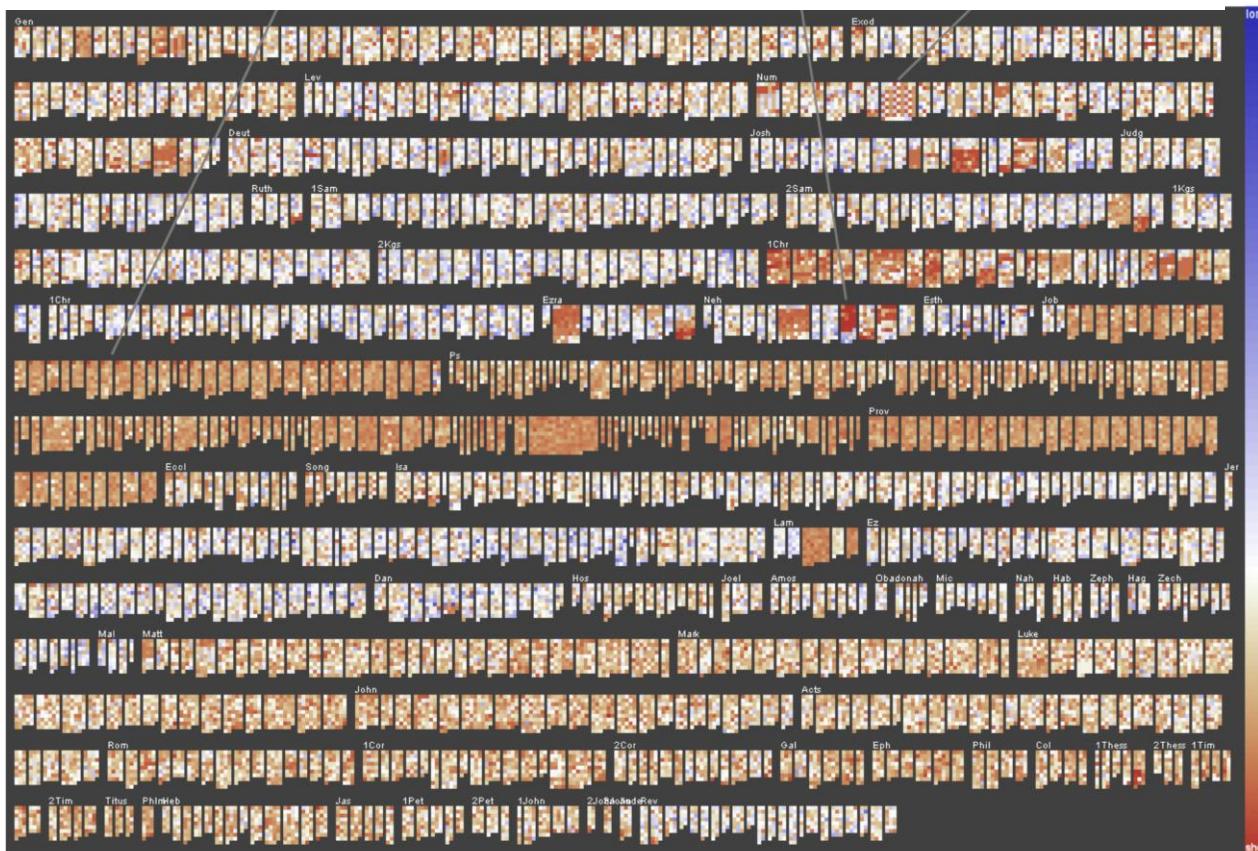
Christian Rohrdantz, Ming C. Hao, Umeshwar Dayal, Lars-Erik Haug, and Daniel A. Keim. *Feature-Based Visual Sentiment Analysis of Text Document Streams*. ACM Transactions on Intelligent Systems and Technology (TIST), vol. 3, no. 2, pp. 26:1-26:25, 2012.

Review Analysis over time (Time Density Plots)



Christian Rohrdantz, Ming C. Hao, Umeshwar Dayal, Lars-Erik Haug, and Daniel A. Keim. Feature-Based Visual Sentiment Analysis of Text Document Streams. ACM Transactions on Intelligent Systems and Technology (TIST), vol. 3, no. 2, pp. 26:1-26:25, 2012.





D. A. Keim and D. Oelke,
"Literature Fingerprinting: A New
Method for Visual Literary
Analysis," 2007 IEEE
*Symposium on Visual Analytics
Science and Technology*,
Sacramento, CA, 2007, pp. 115-
122, doi:
10.1109/VAST.2007.4389004.

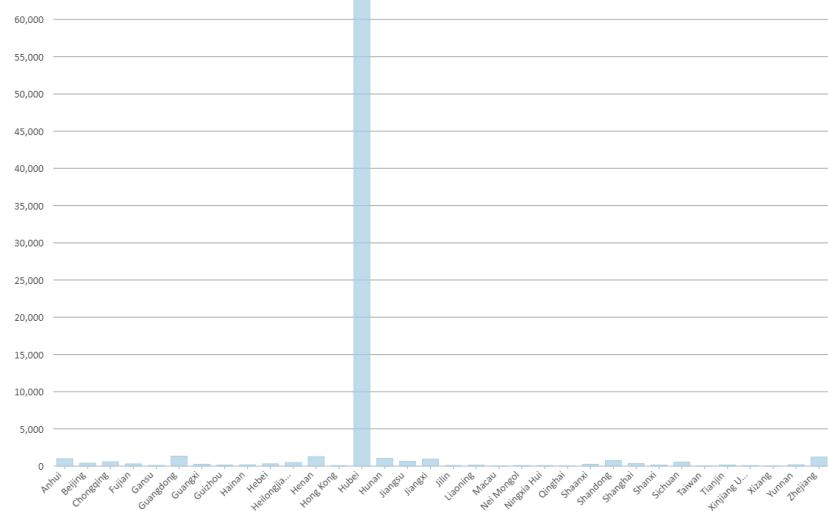
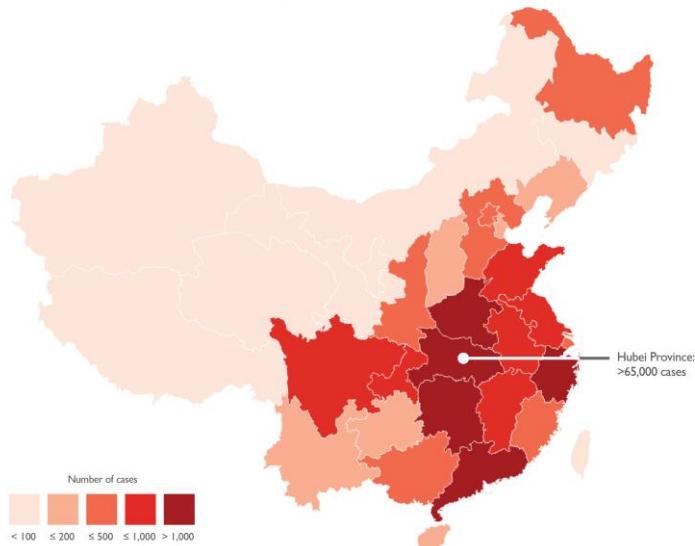
Visual Analytics of Twitter News for Crisis Management



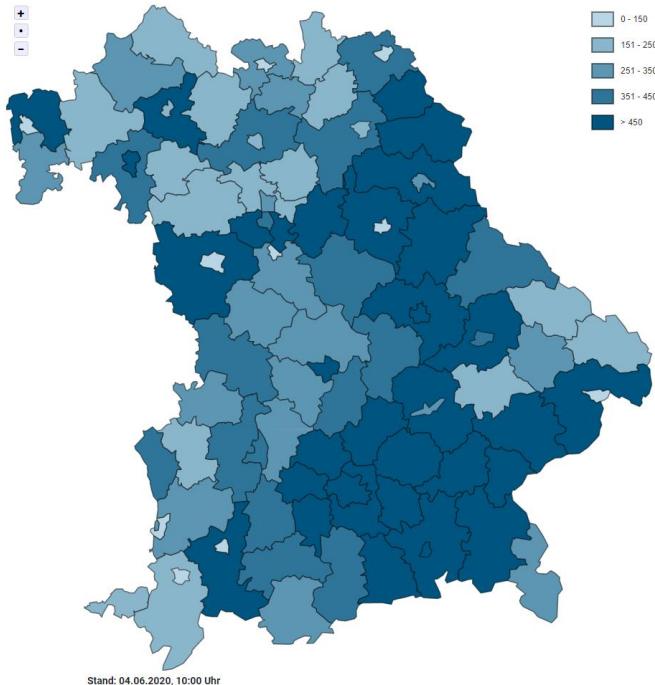
Learning from Covid-19 Visualizations

Wie unterschiedlich stark waren die einzelnen chinesischen Provinzen am 24.2.20 vom Corona-Virus betroffen?

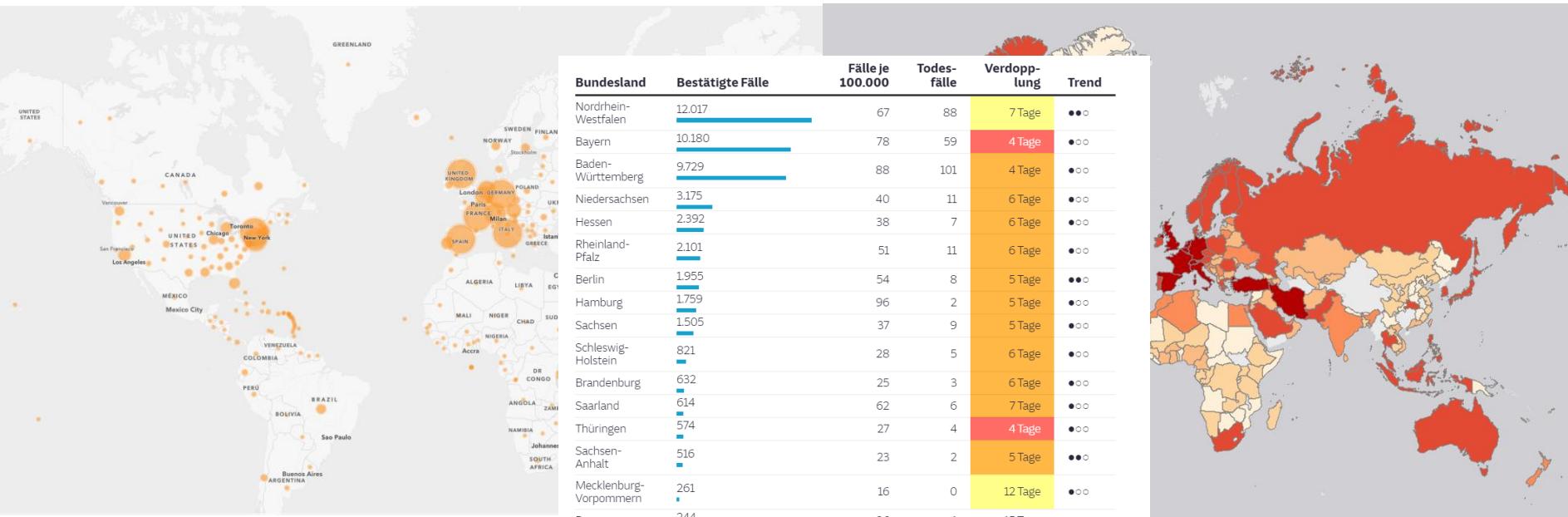
Coronavirus in China: 24th February 2020



Welche bayrischen Landkreise waren am 4.6.20 besonders stark von Corona betroffen?



Welche Regionen der Erde waren zum Zeitpunkt des Screenshots am stärksten von Covid-19 betroffen?

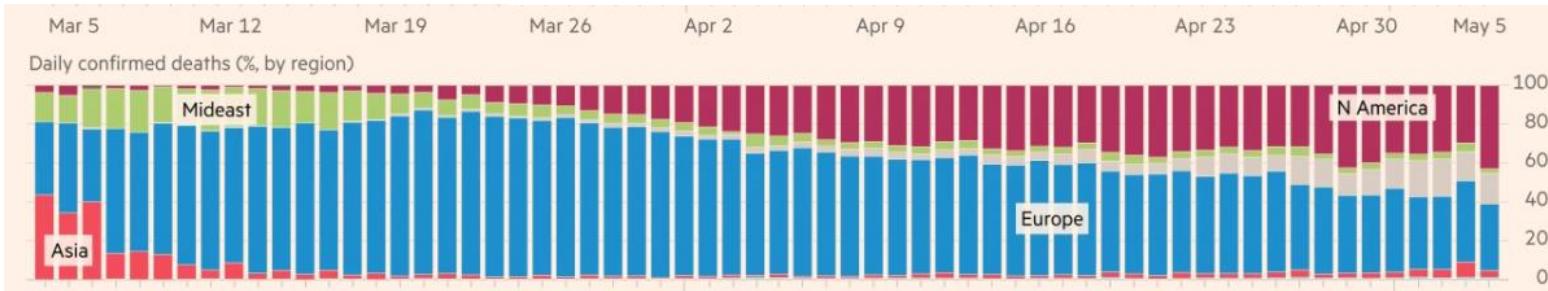
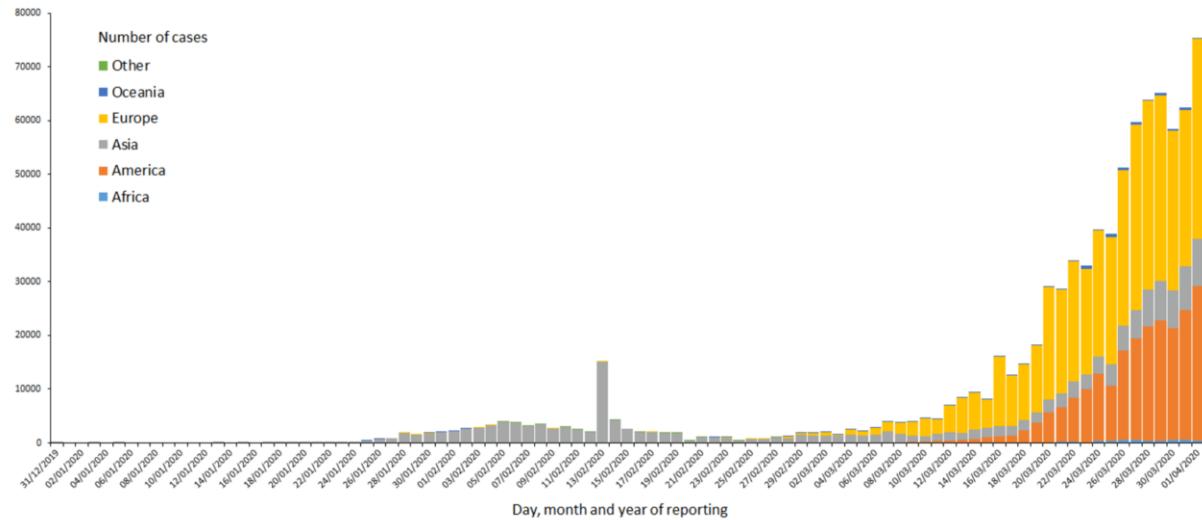


<https://storymaps.arcgis.com/stories/4fdc0d03d3a34aa485de1fb>
 **EMI** Elektrotechnik, Medizintechnik und Informatik
<https://nssabibl.virginia.edu/covid-19/dashboard/>

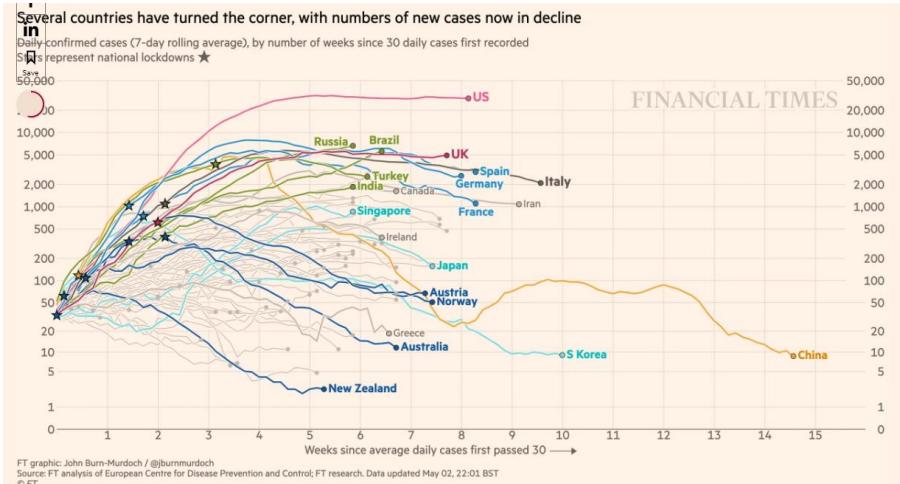
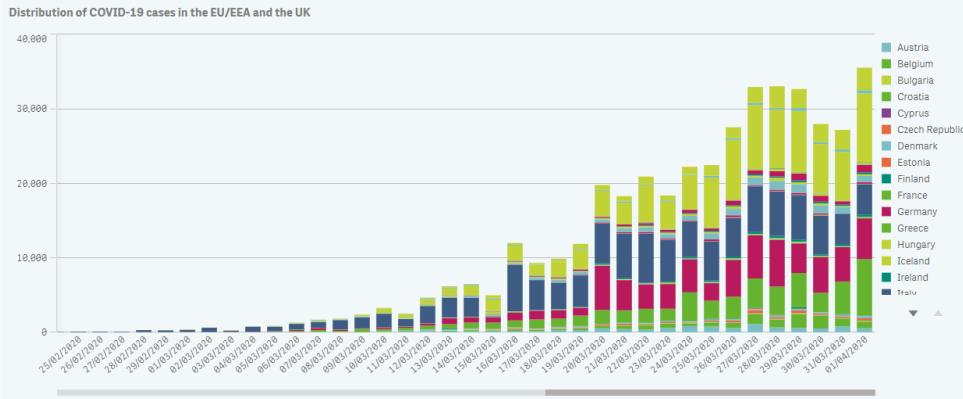
Die Verdopplungszeit gibt an, wie schnell sich die Epidemie ausbreitet. Der Trend zeigt an, wie sich dieses Tempo verändert: wird langsamer ●○○, bleibt gleich ●○○, wird schneller ●●●. Letzter Stand der Daten: 27.03.2020 20:30 Uhr

Summerschool 2020, Prof. Dr. Daniela Oelke

In welcher Region der Erde ist im betrachteten Zeitraum die stärkste Zunahme zu verzeichnen?



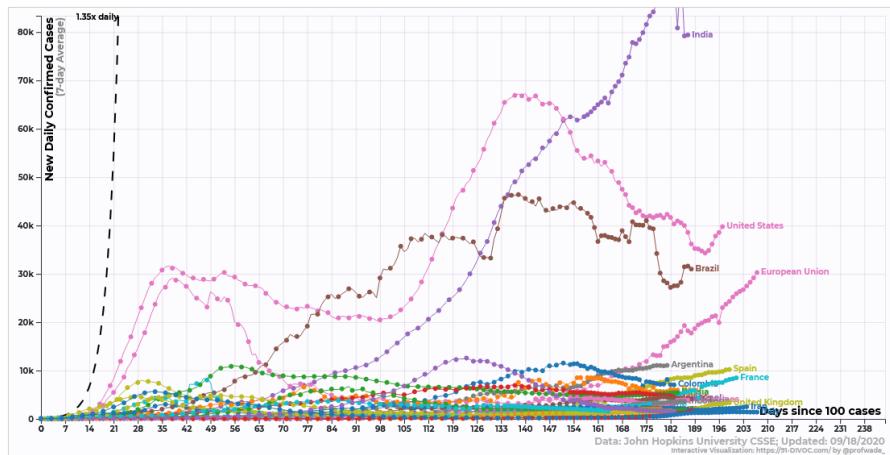
Gab es in Griechenland im betrachteten Zeitraum eine konstante Zunahme der Infektionszahlen oder zwischendrin Einbrüche?



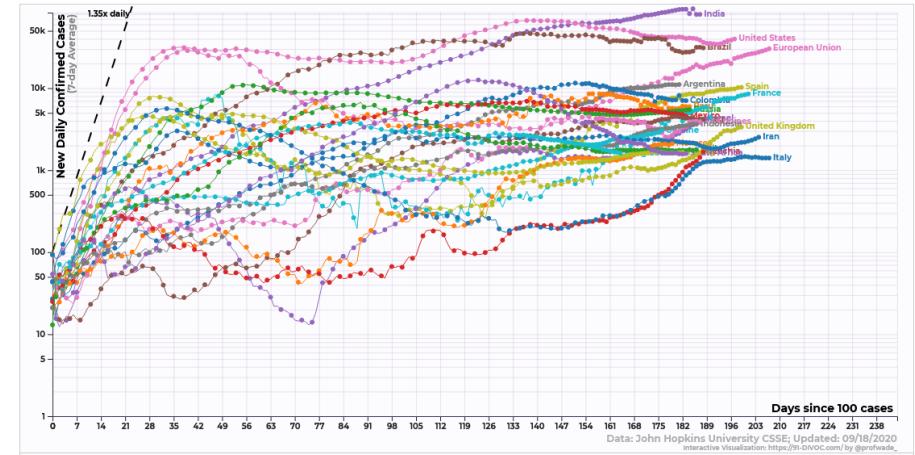
<https://qap.ecdc.europa.eu/public/extensions/COVID-19/COVID-19.html>, <https://www.ft.com/content/a26fbf7e-48f8-11ea-aeb3-955839e06441>

Um welchen Faktor sind die Fallzahlen in der EU zum Zeitpunkt des Erstellung der Graphik kleiner als in den USA?

New Confirmed COVID-19 Cases per Day



New Confirmed COVID-19 Cases per Day



<http://91-divoc.com/pages/covid-visualization/>

Death rates have climbed far above historical averages in many countries that have faced Covid-19 outbreaks



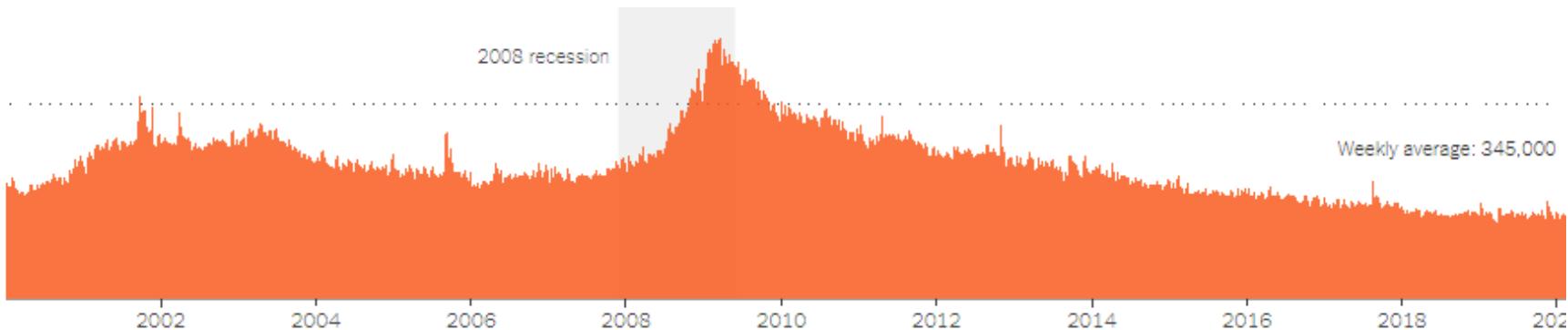
*Italian data are a representative sample of 86% of the country

Source: FT analysis of mortality data. Data updated May 06

FT graphic: John Burn-Murdoch / @burnmurdoch

© FT

Arbeitslosenzahlen USA



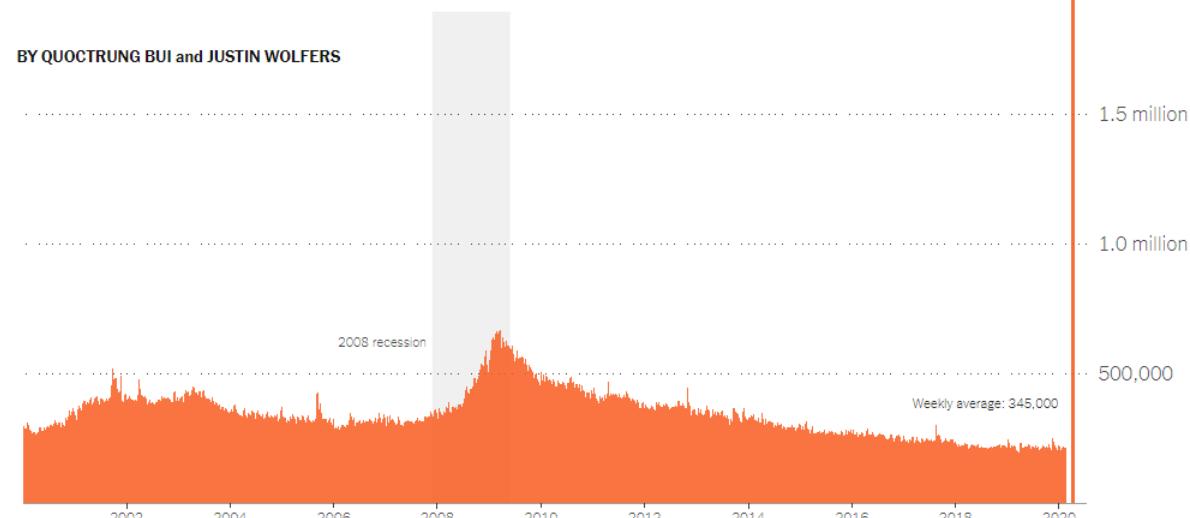
Note: Official figures are seasonally adjusted. Source: Department of Labor

<https://www.nytimes.com/interactive/2020/03/26/upshot/coronavirus-millions-unemployment-claims.html>

More Than 3 Million Americans Lost Their Jobs Last Week. See Your State.

Official statistics have revealed how severely coronavirus has hurt the job market. But it may take several months before we know whether this economic disaster will resemble a storm or a long winter.

BY QUOCTRUNG BUI and JUSTIN WOLFERS



3,283,000
claims filed last week

Visualization Tools & Libraries

Visualization in Python

- **matplotlib**
- **seaborn**
- **plotly / plotly express** → interactivity
- **bokeh** → for building dashboards
- **dash** → for building dashboards

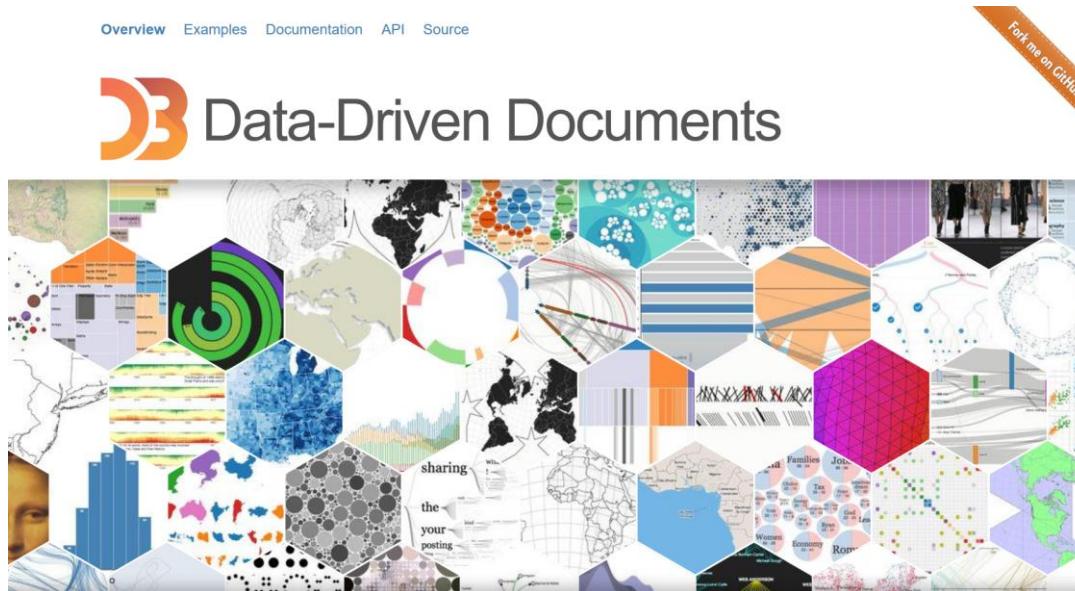
Interactive Dashboards

- Tableau
- Microsoft PowerBI
- QlikView
- SAS Visual Analytics
- ...

Comparison: <https://commercialtools.dbvis.de/>

Visualizations in JavaScript

- Many libraries for „standard“ visualizations exist
- Use the library D3 if you need a maximum of flexibility



Credits & Recommended reading

- Speech and Language Processing. An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition. Daniel Jurafsky, James H. Martin, 2019.
https://web.stanford.edu/~jurafsky/slp3/edbook_oct162019.pdf
- Introduction to Information Retrieval. Manning, Raghavan, Schütze. Cambridge University Press, 2008.
<https://nlp.stanford.edu/IR-book/information-retrieval-book.html>
- spaCy, <http://spacy.io>