

# NLP CoE

Kristin Chen, Aishwarya Bhangale, Raj Desai, Meet Paradia, Wenxi Li

# Purpose

- Find dataset and research problems aligns with business needs and industrial trend

# Problem 1: Drug Review

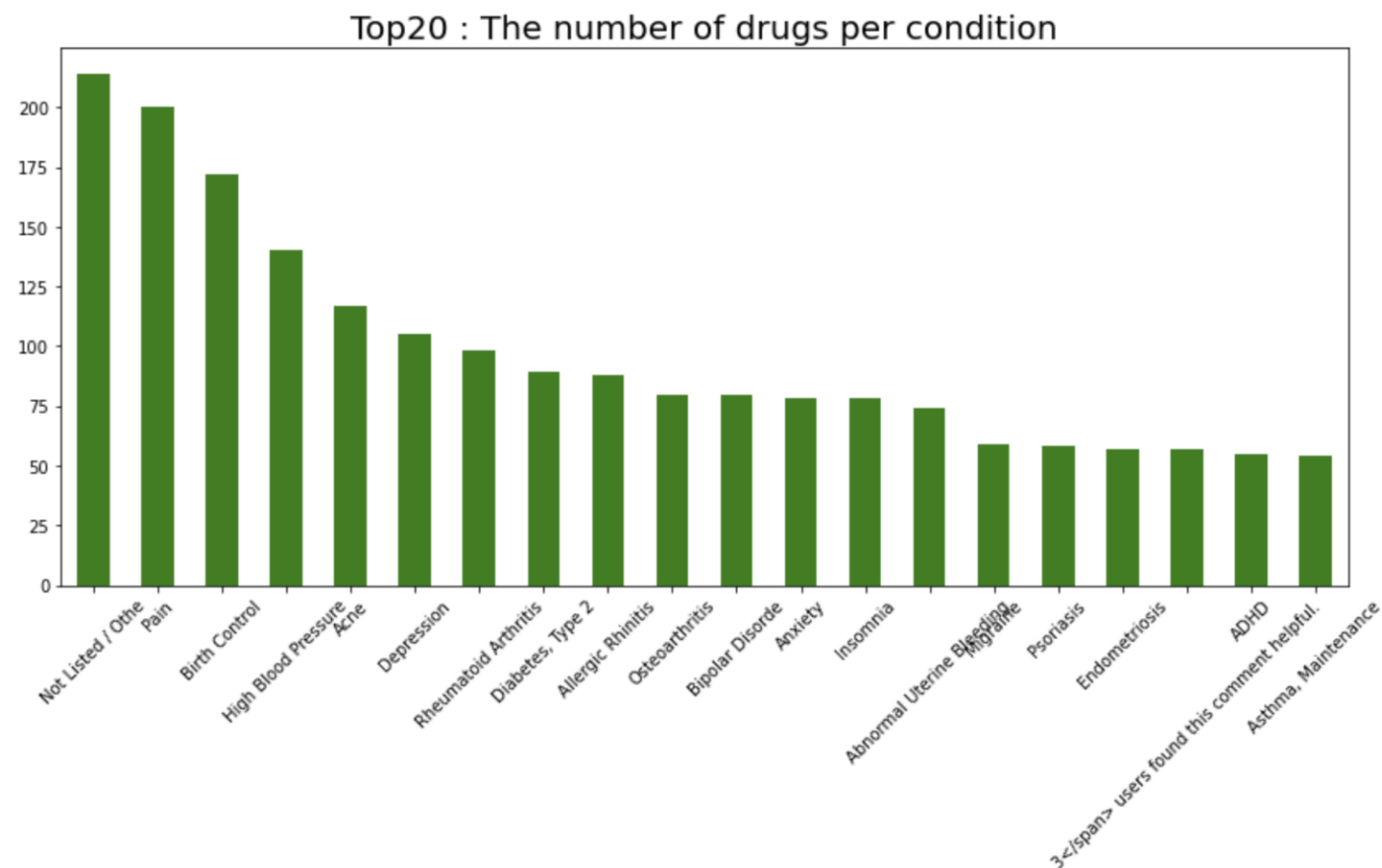
- Data: 161,297 unique reviews for train dataset about drugs related to 884 different health conditions (rating from 1 to 10), and the number of users who found review useful (ranging from 0 to 1291)
- Problem:
  - Text classification: classify health conditions and rating based on the review
  - Regression: predict the rating of the drug based on the review
  - Sentiment analysis
- Potential business application: drug reviews for healthcare/pharma industries

Reference: <https://www.kaggle.com/jessicali9530/kuc-hackathon-winter-2018>

# Problem 1: Drug Review

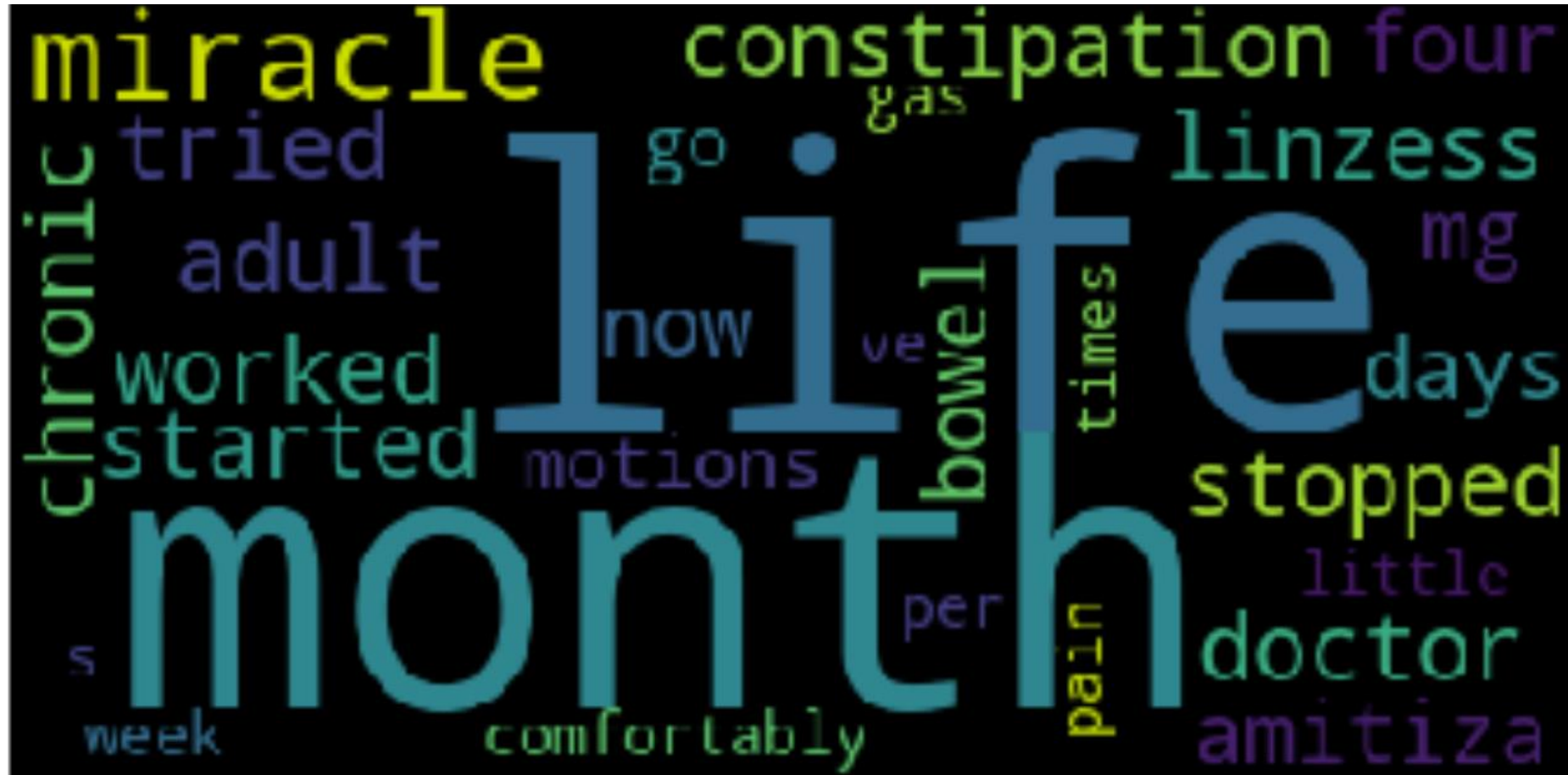
- There are 885 unique conditions , 3436 unique drugs and 389 unique usefulCounts.
- The reviews are given in the timeframe between 1-Apr-08 to 9-Sep-17.

# Problem 1: Drug Review



Condition	
Not Listed / Other	214
Pain	200
Birth Control	172
High Blood Pressure	140
Acne	117
Depression	105
Rheumatoid Arthritis	98
Diabetes, Type 2	89
Allergic Rhinitis	88
Osteoarthritis	80
Bipolar Disorder	80
Anxiety	78
Insomnia	78
Abnormal Uterine Bleeding	74
Migraine	59
Psoriasis	58
Endometriosis	57
3</span> users found this comment helpful.	57
ADHD	55
Asthma, Maintenance	54

# WordCloud



# Problem 1: Drug Review

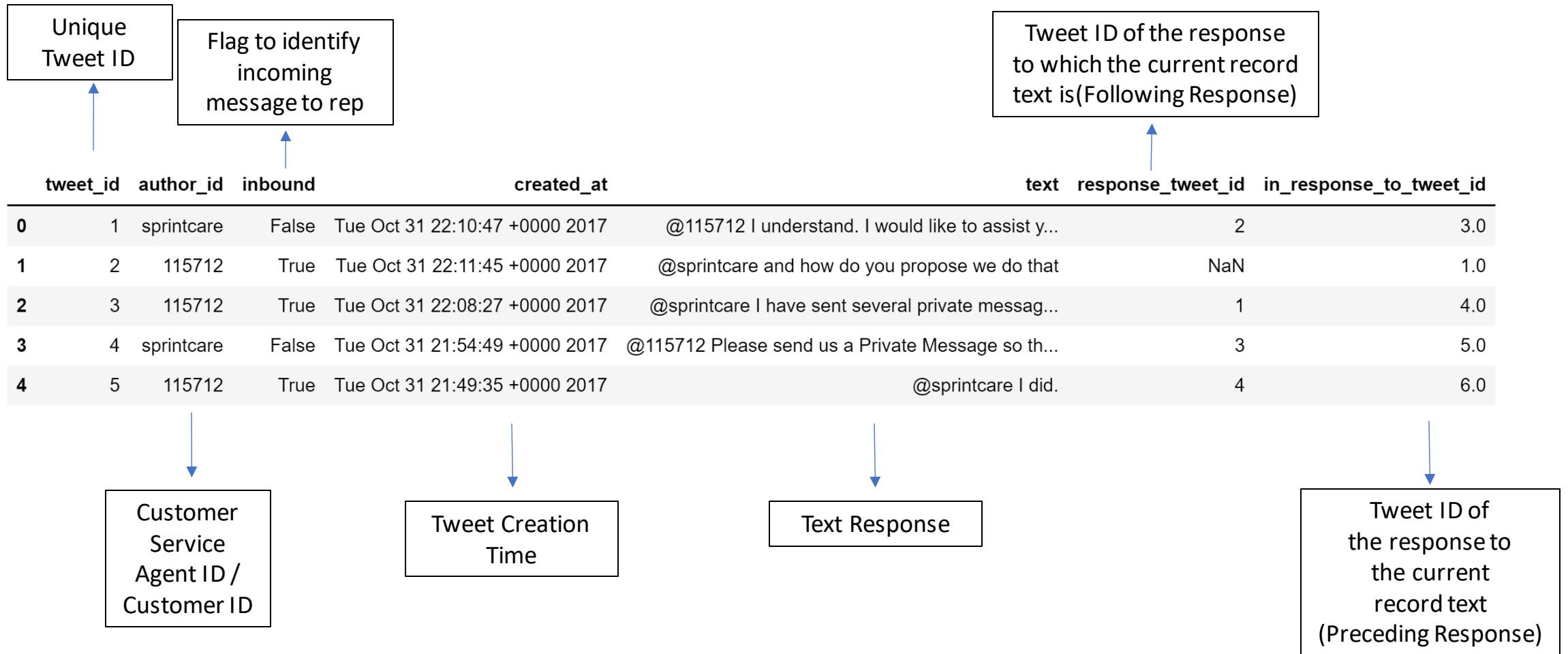
- the number of users who found review useful distribution

count	161297	10	31.61%
mean	28.00476	9	17.07%
std	36.40374	1	13.40%
min	0	8	11.71%
25%	6	7	5.86%
50%	16	5	4.97%
75%	36	2	4.30%
max	1291	3	4.04%
		6	3.93%
		4	3.11%

- Rating distribution
- Top 10 health conditions

Birth Control	17.97%
Depression	5.69%
Pain	3.86%
Anxiety	3.65%
Acne	3.48%
Bipolar Disorder	2.62%
Insomnia	2.29%
Weight Loss	2.27%
Obesity	2.22%
ADHD	2.11%

# Problem 2: Twitter Conversational Bot





# Twitter Agent Description

With 1.5MM total responses sent by Twitter agents and 1.27MM total responses received, following is a table showing agents classified into various businesses.

Business Name	Agents
Airline	AirAsia Support, Delta, AmericanAir, SouthwestAir, VirginAtlantic, AlaskaAir, VirginAmerica, JetBlue
Software	AdobeCare, YahooCare, SpotifyCare, AskeBay, DropBoxSupport, AzureSupport, GoDaddyHelp, asksalesforce, GooglePlayMusic, OfficeSupport, TwitterSupport, NortonSupport, SCsupport, AWSSupport, mediatemplehelp, AskTigogh, PandoraSupport, AmazonHelp
Super Market	Morrisons, Tesco, sainsburys, ArgosHelpers, AskTarget, Walmart, AldiUK
Transport	Nationalrailenq, AskLyft, UPSHelp, VirginTrains, GWRHelp, TfL, LondonMidland

*\*Additional business include Beauty, Education, Clothing, Electronics, Finance, Food, Gaming, Hotel and Rentals*

# Most Frequently Used Words By Customers

Agent	Responses Sent	Responses Received	Most Common Words Used By Customer				
			1st	2nd	3rd	4th	5th
Sprintcare	22,381	13,876	phone	store	bill	account	upgrade
ChipotleTweets	18,749	21,593	burrito	order	queso	bowl	chip
XboxSupport	24,557	28,083	update	account	live	buy	console
JetBlue	8,020	9,475	delay	seat	book	bag	cancel
AskPayPal	11,298	10,164	account	money	bank	card	hold

# Problem #3: resume parsing

name	phone_num	email	
Karthik	(410)-292-1151	karthikr2194@gmail.com	{'EDUCATIO
Robert		hpundir@umd.edu	
Data Analysis	469-370-9437	tirth2410@gmail.com	{'University', 'Ha
Mythri	(281) 725-7080	mythripartha8@gmail.com	{
Qizhe		ziyaotingyu@gmail.com	{'EDUCATION Univer
Stamford	(475) 685 0166	rachan_vamsi.bhooshi@uconn.edu	{'Un
Github	240-713-8296	rmadireddy1@student.gsu.edu	
Data Science	443-833-6344	saidam1@umbc.edu	
Akash Patel	+1 4845387112	adp178@scarletmail.rutgers.edu	{'Sti
Analytics		akshara@ou.edu	{'EDUCA
Alexandra	630.818.6275	manetas.alexandra@gmail.com	{'DePaul University', 'University Research Institutes'}
Pranav Premdas Gulghane	-4695140739	pranavpremdasgulghane@gmail.com	{'University'}
Lu Berkeley	(559) 387-0880	luwinnie12@gmail.com	{'University'}
Khoury College	(617) 818-4953	nagaraj.m@northeastern.edu	{'National University', 'D Institute', 'REVA University'}
Data Management	18572077337	mohitmanjaria55333@gmail.com	{'UNIVERSITY Master', 'UNIVERSITY Masters'}
Yuchen	949-413-2863	yuchen724@ucla.edu	{'EDUCATION University', 'University'}
Dimple		dimple8997@gmail.com	set()
Kompi S		indupriyakompi@gmail.com	{'GPA Northeastern University'}
James	(706) 305-6369	james.domingo@gatech.edu	{'University', 'EDUCATION Georgia Institute'}
Machine Learning	805 453 1532	anyampatel@gmail.com	set()
Ryan Martin Goodwin LinkedIn Github	(910) 547-7027	rgoodwin1997@gmail.com	{'Science University'}
Andrew	(443)-742-3540	eckan01@gettysburg.edu	set()
Git	(202) 212 -9607	jiatingchen0107@gmail.com	{'UNIVERSITY Master'}
Data Scientist	774-994-4106	manishap2690@gmail.com	{'State University'}
Python	12678819188	sgheereddy@gmail.com	{'EDUCATION University', 'Jawaharlal Nehru Technological University'}
Senthil	-8553247568	msnathan55@yahoo.com	{'University', 'Toyohashi University', 'PSG College'}

Name	Date modified	Type
721091408_phonescreening	10/19/2021 1:22 PM	Adobe Acrobat ...
821101104_phonescreening	10/19/2021 1:23 PM	Adobe Acrobat ...
821101136_phonescreening	10/20/2021 4:00 PM	Adobe Acrobat ...
821110440_Mythri	11/9/2021 11:47 PM	Adobe Acrobat ...
821110445_Qizhe	11/30/2021 11:13 AM	Adobe Acrobat ...
821110447_rachan	11/30/2021 11:16 AM	Adobe Acrobat ...
821110448_RAGHUVI	11/30/2021 11:17 AM	Adobe Acrobat ...
821110464_Saida	11/9/2021 11:51 PM	Adobe Acrobat ...
821110509_Akash	11/18/2021 1:21 PM	Adobe Acrobat ...
821110510_AKSHARA	11/18/2021 1:24 PM	Adobe Acrobat ...
821110511_Alexandra	11/18/2021 1:24 PM	Adobe Acrobat ...
821110527_	11/18/2021 1:25 PM	Adobe Acrobat ...
821110528_Wanting	11/18/2021 1:26 PM	Adobe Acrobat ...
821111539_manaswini	12/13/2021 11:19 AM	Adobe Acrobat ...
821112919_mohit	12/13/2021 11:19 AM	Adobe Acrobat ...
821112928_yechen	12/13/2021 11:18 AM	Adobe Acrobat ...
821122313_Dimple	1/10/2022 9:08 AM	Adobe Acrobat ...
821122315_Indupriya	1/10/2022 9:11 AM	Adobe Acrobat ...
821122316_James	1/10/2022 9:13 AM	Adobe Acrobat ...
822010605_anya	1/21/2022 11:43 AM	Adobe Acrobat ...
822010628_ryan	1/21/2022 11:50 AM	Adobe Acrobat ...
C01-21111807_Andrew_Decker	12/8/2021 2:12 PM	Adobe Acrobat ...
KRISTINJIATINGCHEN_RESUME_06082...	12/13/2021 2:07 PM	Adobe Acrobat ...
Manisha Patel Resume	12/28/2021 2:04 PM	Adobe Acrobat ...
SAIKUMAR	12/29/2021 1:46 PM	Adobe Acrobat ...

# Problem #3: resume parsing (variants in phone numbers)

name	phone_num	email	education
Robert		hpundir@umd.edu	{'University'}
Dimple		dimple8997@gmail.com	set()
Senthil	-8553247568	msnathan55@yahoo.com	{'University', 'Toyohashi University', 'PSG College'}

## HARSH PUNDIR

3425 Tulane Dr. Hyattsville, MD 20783 240.423.5453

[hpundir@umd.edu](mailto:hpundir@umd.edu) | [www.linkedin.com/in/hpundir](https://www.linkedin.com/in/hpundir) | <https://github.com/HARSHPUNDIR>


## EDUCATION

**Robert H. Smith School of Business, University of Maryland**, College Park, MD

August 2020 - Present

**Master of Science in Business Analytics**, Focus Area: Data Science GPA: 4.0

■ Relevant Coursework – DBMS, Python, Data Models & Decisions, Data Mining, Big Data & AI, Data Visualization.

**Dimple Mehra** 

San Jose, CA 95035

Email : [dimple8997@gmail.com](mailto:dimple8997@gmail.com)

Mobile : +1-(312-358-1359)

## Senthil Nathan M

Bangalore, Karnataka, India • +91-8553247659 • [msnathan55@yahoo.com](mailto:msnathan55@yahoo.com) • <https://www.linkedin.com/in/senthil-nathan-murugappan-585a4967/>

(240) 917-4861 • • [jia Yue.fei@marylandsmith.umd.edu](mailto:jia Yue.fei@marylandsmith.umd.edu) • [www.linkedin.com/in/jifei](http://www.linkedin.com/in/jifei)

## EDUCATION

College Park, MD, USA

May 2022

- Nanjing, China

Jun 2019

- ## TECHNICAL SKILLS

- ## PROJECT EXPERIENCE

### International market segmentation using Normal Mixture Regression Models

- ## WORK EXPERIENCE

Shanghai, China

Apr 2021 –Jul2021

- DIDI Global**

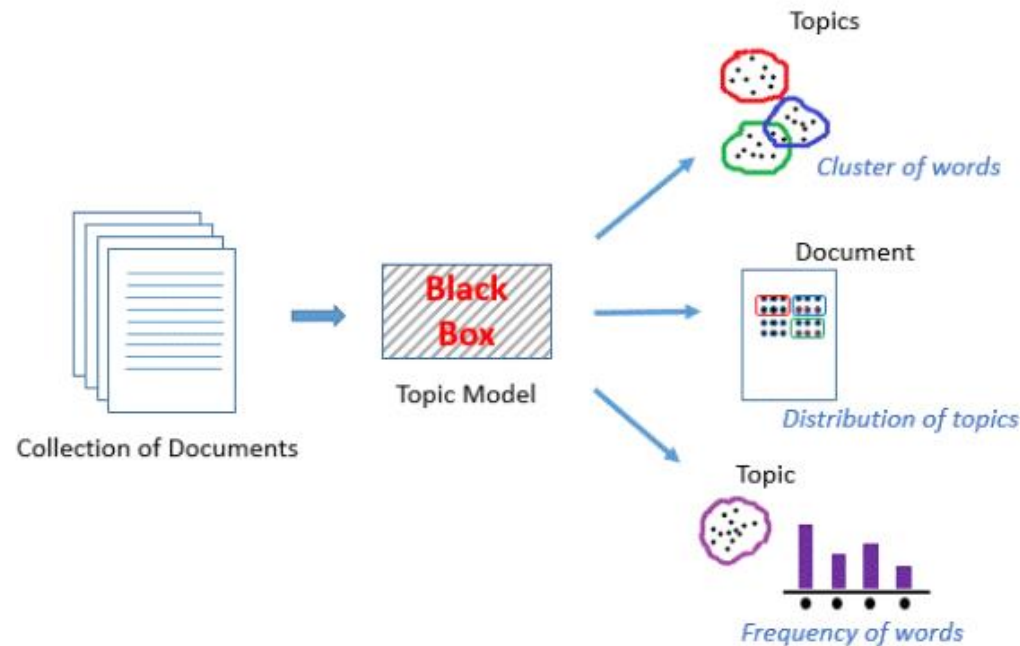
Nanjing, China

Jan 2021 – Mar 2021

- ### LEADERSHIP AND VOLUNTEER EXPERIENCE

# Topic Modelling

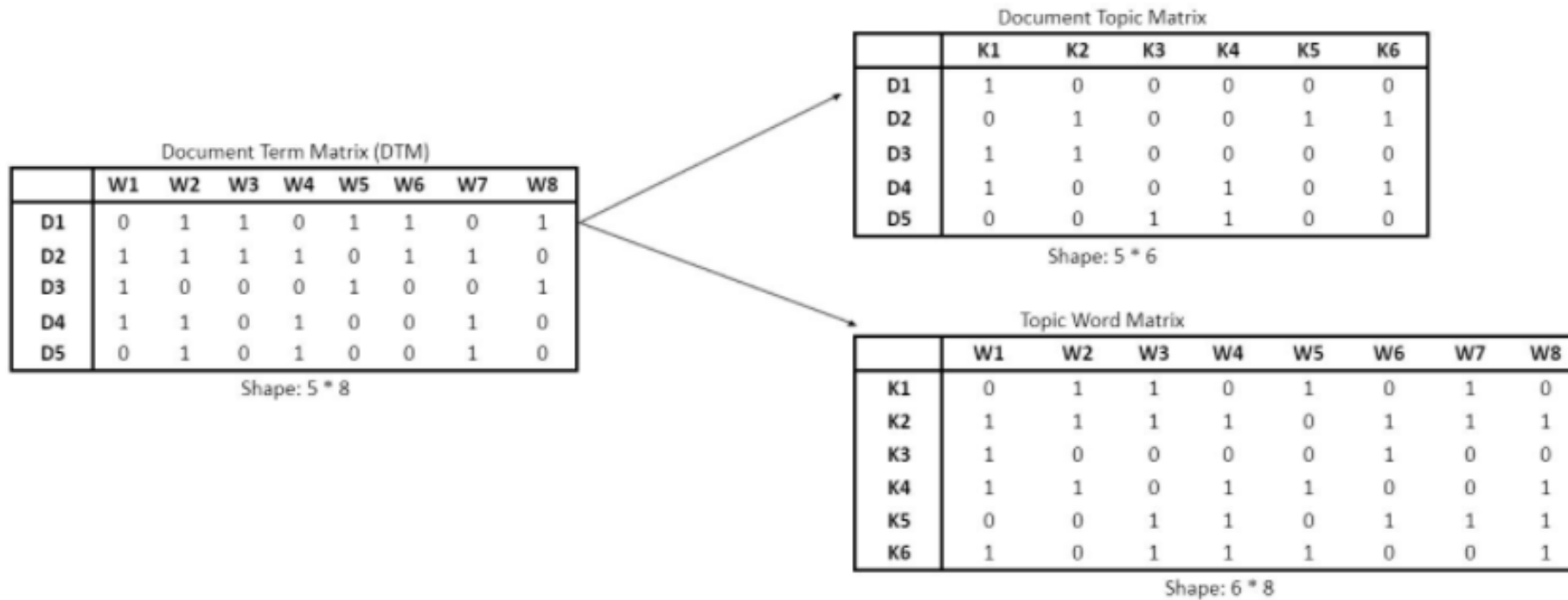
- Unsupervised ML technique to discover a set of topics that best segregate a corpus of documents/ bag of words
- Tries to figure out which topics are present in the documents of the corpus and how strong is that presence
- Eg : Latent Dirichlet Allocation (*Gensim Library in Python*)



# Latent Dirichlet Allocation

Two key assumptions:

- Documents are a distribution of topics
- Topics are a distribution of words





# How the Algorithm functions?

- For each document in the corpus, a topic word distribution is generated
- Each document will be randomly assigned to topics in the first iteration
- Post the first iteration, LDA provides per document topic distribution and per topic word distribution
- End goal is to optimize these two output matrices by updating the topic for each word in each document keeping the topics assigned to other words in the document constant
- This update is done by calculating the following two probabilities:
  - 1) **Proportion (Topic  $k$  / Document  $D$ )**
  - 2) **Proportion (word  $w$  / Topic  $k$ )**
- Based on the product of these two probabilities, LDA assigns a new topic to the word

