# Demystifying
# Data Science

**Venkatesh Vinayakarao**

venkateshv@cmi.ac.in

http://vvtesh.co.in

SSN School of Advanced Career Education

The world's most valuable resource is no longer oil, but data. – **Economist Report, 2017.**

# What Comes Next?

byte

kilobyte
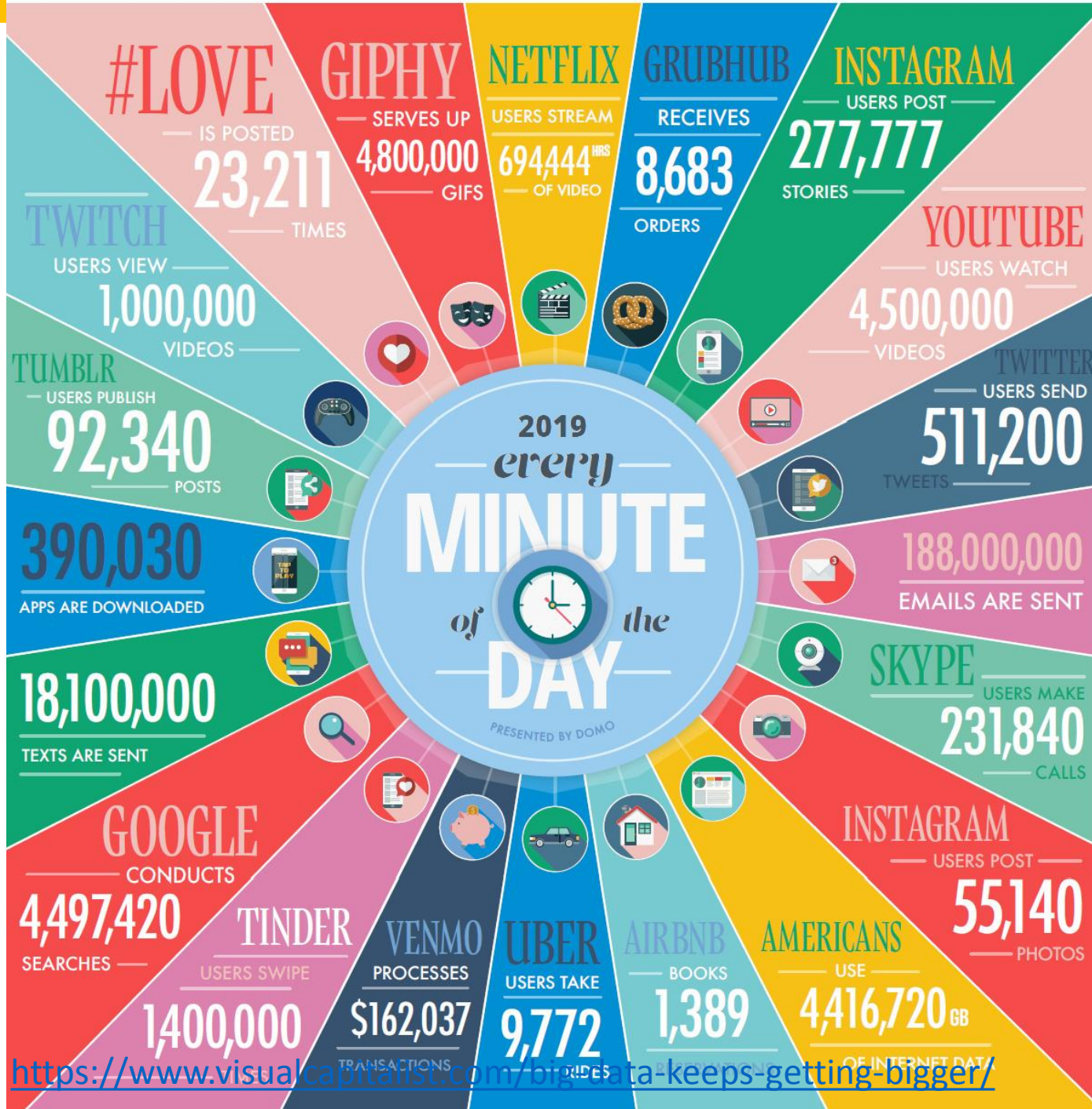
megabyte

gigabyte

??

???

????

?????

# Sizes

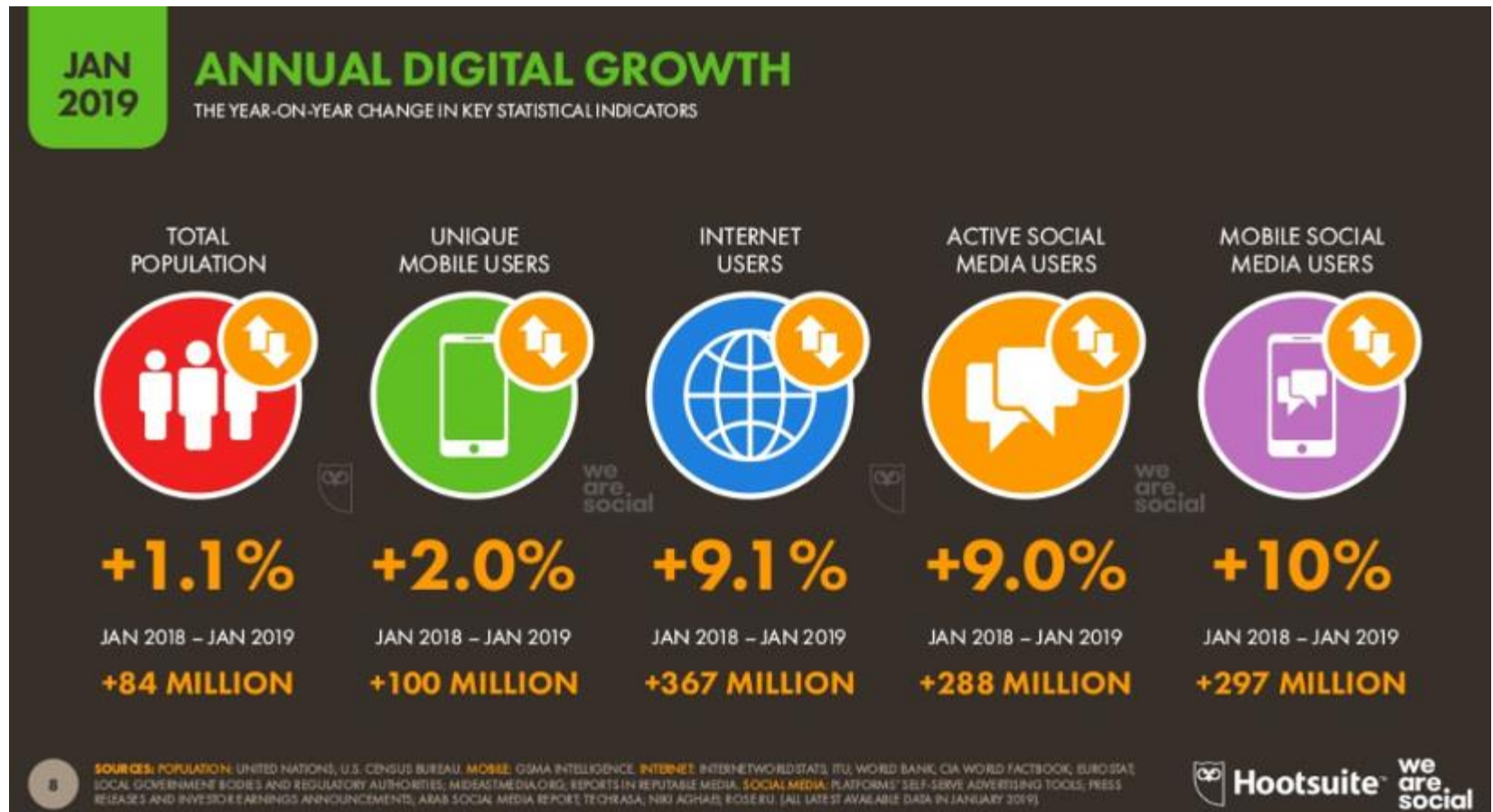| Name | Size |
|------|------|
| **Byte** | 8 bits |
| **Kilobyte** | 1024 bytes |
| **Megabyte** | 1024 kilobytes |
| **Gigabyte** | 1024 megabytes |
| **Terabyte** | 1024 gigabytes |
| **Petabyte** | 1024 terabytes |
| **Exabyte** | 1024 petabytes |
| **Zettabyte** | 1024 exabytes |
| **Yottabyte** | 1024 zettabytes |

# Big Data is Ubiquitous

- Facebook Statistics
  - 1.5 billion people are active on Facebook **daily!**
  - **Every minute** there are 510,000 comments posted and 293,000 statuses updated!
  - More than 300 million photos get uploaded **per day**!
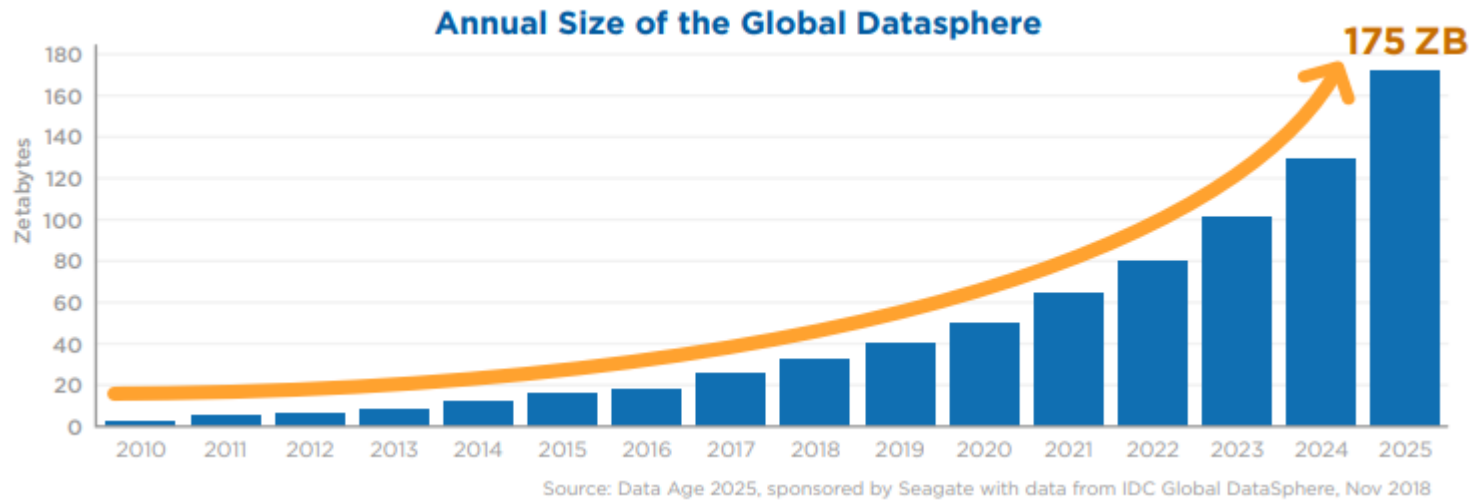  - Totally, more than 2.5 Trillion posts!

Source: Forbes

2019
every
MINUTE
of the
DAY
PRESENTED BY DOMO

#LOVE IS POSTED 23,211 TIMES

GIPHY SERVES UP 4,800,000 GIFS

NETFLIX USERS STREAM 694,444 HRS OF VIDEO

GRUBHUB RECEIVES 8,683 ORDERS

INSTAGRAM USERS POST 277,777 STORIES

TWITCH USERS VIEW 1,000,000 VIDEOS

YOUTUBE USERS WATCH 4,500,000 VIDEOS

TUMBLR USERS PUBLISH 92,340 POSTS

TWITTER USERS SEND 511,200 TWEETS

390,030 APPS ARE DOWNLOADED

188,000,000 EMAILS ARE SENT

18,100,000 TEXTS ARE SENT

SKYPE USERS MAKE 231,840 CALLS

GOOGLE CONDUCTS 4,497,420 SEARCHES

INSTAGRAM USERS POST 55,140 PHOTOS

TINDER USERS SWIPE 1,400,000

VENMO PROCESSES $162,037 TRANSACTIONS

UBER USERS TAKE 9,772 RIDES

AIRBNB BOOKS 1,389

AMERICANS USE 4,416,720 GB OF INTERNET DATA

# And, It is Growing!

# Data Growth

**Annual Size of the Global Datasphere**



Source: Data Age 2025, sponsored by Seagate with data from IDC Global DataSphere, Nov 2018

Mankind's quest to digitize the world!
33 ZB (2018) → 175 ZB (2025)
size of global datasphere*

*Source: https://www.seagate.com/files/www-content/our-story/trends/files/idc-seagate-dataage-whitepaper.pdf

# Solitary Confinement is Cruel
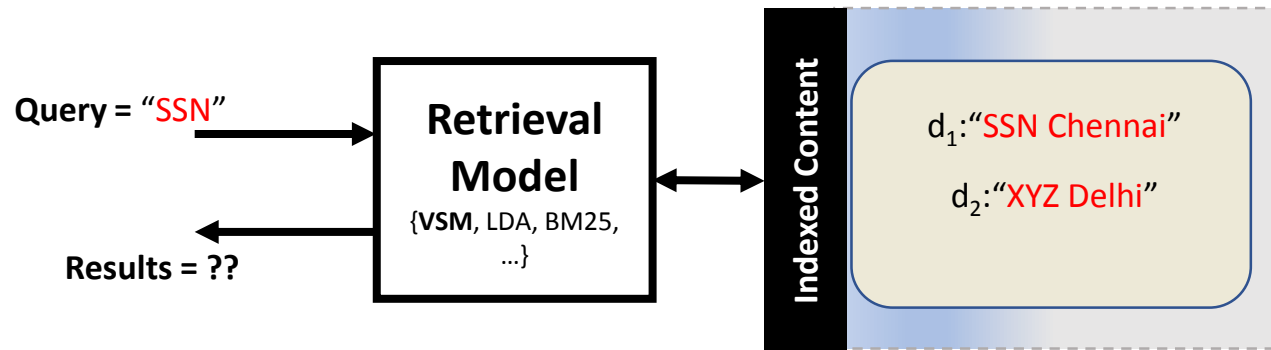
# DATA & AI LANDSCAPE 2019

**World needs data scientists!**

Jun 27, 2019 — © Matt Turck (@mattturck), Lisa Xu (@lisaxu92), & FirstMark (@firstmarkcap) — mattturck.com/bigdata2019

FIRSTMARK

9

# Data Science

**Loads of (structured and unstructured) data available.**

**Need scientifically sound methods to capture, maintain, process, communicate and analyze data.**
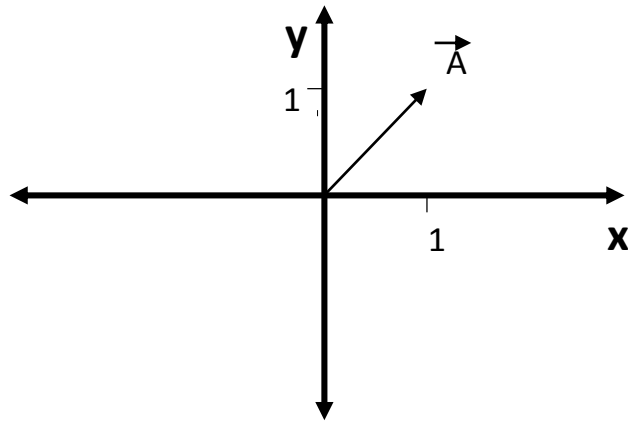
# Modern Text Processing

Vector Space Model
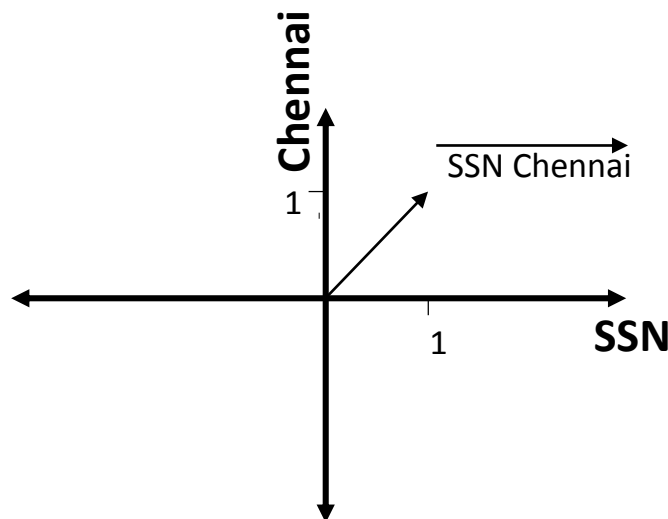
# Which Document to Retrieve?

**Query = "SSN"**

**Retrieval Model**
{**VSM**, LDA, BM25, …}

**Results = ??**

**Indexed Content**

$d_1$:"SSN Chennai"

$d_2$:"XYZ Delhi"

# Vectors

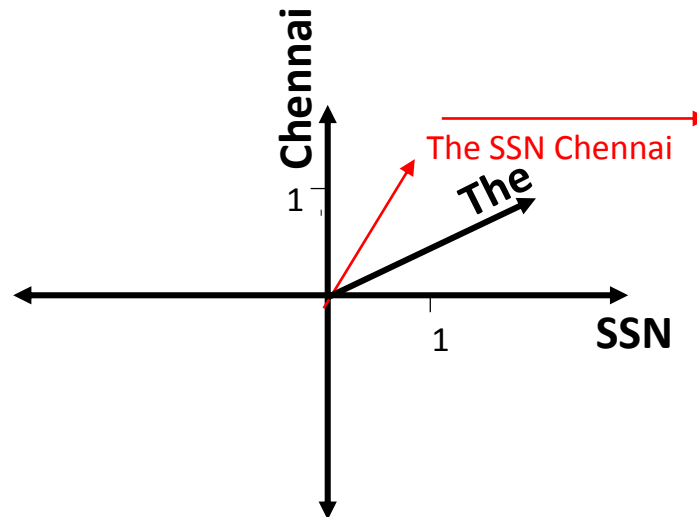- Geometric entity which has magnitude and direction

# Sentences are vectors
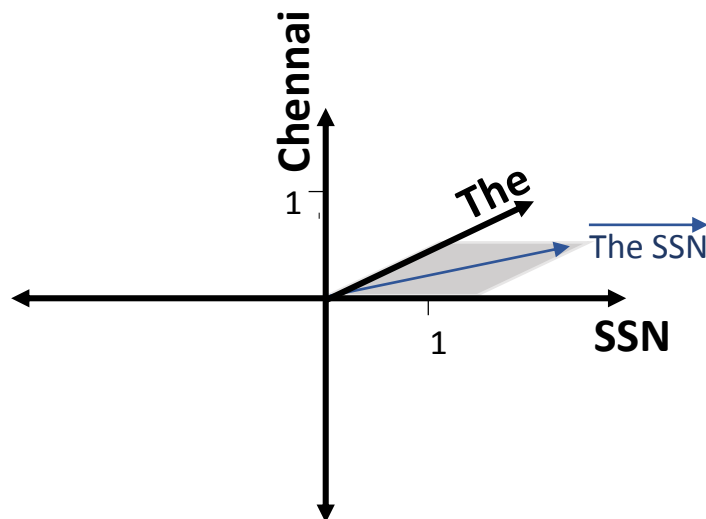
- "SSN Chennai" as a vector

# Sentences are vectors

- "The SSN Chennai" is a 3-dimensional vector
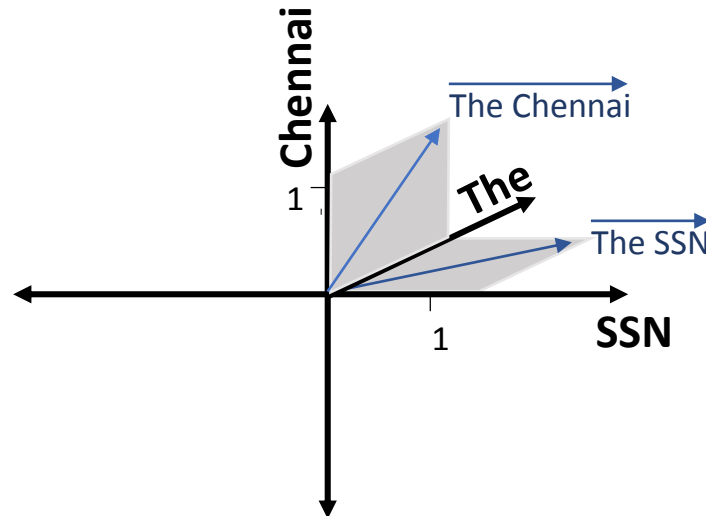
# Sentences are vectors

- On this 3D space, "The SSN" vector will lie on the x (The) and z (SSN) plane.

# Comparing Sentences

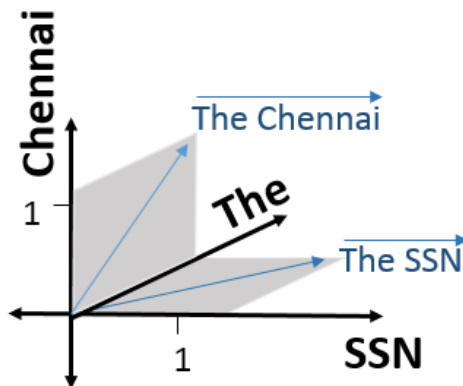- We can compare sentences using the angle between vectors

# Angle between two vectors

- What is the angle between $\overrightarrow{\text{The}}$ and $\overrightarrow{\text{SSN}}$ vectors?
- What is the angle between $\overrightarrow{\text{SSN}}$ and $\overrightarrow{\text{Chennai}}$ vectors?
- What is the angle between $\overrightarrow{\text{The SSN}}$ and $\overrightarrow{\text{The SSN}}$ vectors?

18

# Mathematical Notation

- We represent vectors as follows:
  - Vector = (dimension1, dimension2, dimension3, …)
    - First, define the dimensions
    - Next, put "1" if the word is present in the sentence, else "0"
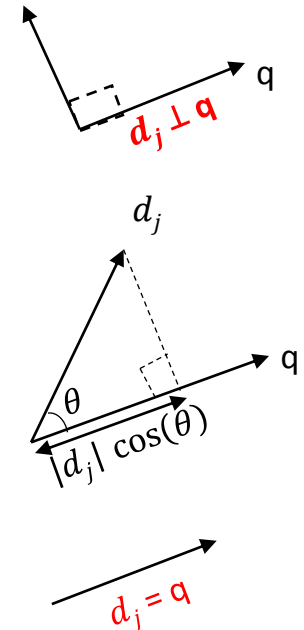
- Example

In our example,
 vector = (The, SSN, Chennai)

So,

$\overrightarrow{\text{The Chennai}}$ = (1,0,1)
$\overrightarrow{\text{The SSN}}$ = (1,1,0)

# Converting from "$0 - 90$" to "$0 - 1$"

- For convenience, We convert the angles 0 – 90 to values 0 – 1
  - When vectors are same, we want to output 1.
  - When vectors are perpendicular, we want to output 0.

|  | 0° | 30° | 45° | 60° | 90° |
|---|---|---|---|---|---|
| sin $\theta$ | 0 | $\frac{1}{2}$ | $\frac{1}{\sqrt{2}}$ | $\frac{\sqrt{3}}{2}$ | 1 |
| cos $\theta$ | 1 | $\frac{\sqrt{3}}{2}$ | $\frac{1}{\sqrt{2}}$ | $\frac{1}{2}$ | 0 |
| tan $\theta$ | 0 | $\frac{1}{\sqrt{3}}$ | 1 | $\sqrt{3}$ | Not defined |

$d_j \perp q$

$d_j$

$\theta$

q

$|d_j| \cos(\theta)$

$d_j = q$

20

# A Way to Calculate cosθ

- $\cos(\theta) = \dfrac{x.y}{||x||\,||y||}$

- Here,
  - x.y  is the "dot product" of x and y vectors.

- So, similarity between "The SSN" and "SSN Chennai"

  $= \dfrac{1.0 + 1.1 + 0.1}{\sqrt{1^2 + 1^2 + 0^2}\sqrt{0^2 + 1^2 + 1^2}} = \dfrac{1}{\sqrt{2}\sqrt{2}} = 0.5$

# Which Document to Retrieve?

**Query = "SSN"**

**Results = ??**

**Retrieval Model**
{**VSM**, LDA, BM25, …}

**Indexed Content**

$d_1$:"SSN Chennai"

$d_2$:"XYZ Delhi"

# Example

Let query q = "SSN".

Let document, $d_1$ = "SSN Chennai" and $d_2$ = "XYZ Delhi".

|     | SSN | Chennai | XYZ | Delhi |
| --- | --- | --- | --- | --- |
| q | 1 | 0 | 0 | 0 |
| $d_1$ | 1 | 1 | 0 | 0 |
| $d_2$ | 0 | 0 | 1 | 1 |

In our VSM, q = (1,0,0,0), $d_1$= (1,1,0,0) and $d_2$ = (0,0,1,1)

**similarity(d$_1$, q)** $= \dfrac{d_1.q}{||d_1||\ ||q||} = \dfrac{1.1+1.0+0.0+0.0}{\sqrt{1^2+1^2}\sqrt{1^2}} = \dfrac{1}{\sqrt{2}} = 0.71$

**similarity(d$_2$, q)** $= \dfrac{d_2.q}{||d_2||\ ||q||} = \dfrac{1.0+0.0+0.1+0.1}{\sqrt{1^2+1^2}\sqrt{1^2}} = 0.$

# Which Document to Retrieve?



**Query** = "SSN"

**Retrieval Model**
{**VSM**, LDA, BM25, ...}

**Results = ??**

**Indexed Content**

$d_1$:"SSN Chennai"

$d_2$:"XYZ Delhi"

# Summary

- Data is Ubiquitous
  - and it is growing too!

- Modern Text Processing
  - Vector Space Model

- Remember
  - Data processing goes beyond common sense… we need techniques and tools.
  - Products are good to learn. Principles are even more important. Don't ignore them.

# Memories