

Social Media Data and its availability in Research

Debjyoti Paul
University of Utah
deb@cs.utah.edu

OVERVIEW

The term *social media* gained attention with the advent of Web 2.0 in the first decade of 20th century [9]. Web 2.0 is also known as *Participative* or *Social Web* that emphasize on user interaction and user generated content encouraging participatory culture. Before we jump into more details of social media it would be wiser to define it. Though ever evolving social media services makes it hard to define them, most of the research work define it as follows.

Definition 1 (Social Media). Social media are interactive computer-mediated technologies that facilitate the creation and sharing of information, ideas, career interests and other forms of expression via virtual communities and networks [10].

In contrast to the *traditional media* which operates under a monologic transmission model i.e. one source to many receivers, such as a television, newspaper or a radio station which broadcasts the same programs to an entire city; *social media* are dialogic transmission system which brings interaction, usability and a notion of individual entity in digital world.

Marketing and social media experts broadly agrees to classify social media with respect to media type and its usage i.e *blogs*, *social networks*, *microblogs*, *photo sharing*, *video sharing*, *business networks*, *enterprise social networks*, *forums*, *products/services review*, *social bookmarking*, *social gaming*, *collaborative projects* and *virtual worlds*.

1 PART A

Social Media Data in Numbers

Marketing and social media experts broadly agrees to classify social media with respect to media type and its usage i.e *blogs*, *social networks*, *private messaging*, *microblogs*, *photo sharing*, *video sharing*, *professional networks*, *enterprise social networks*, *forums*, *products/services review*, *social bookmarking*, *social gaming*, *collaborative projects* and *virtual worlds* [1]. We now present a list of relevant social media according to the classification stated in Table 1.

The popularity of a social media site is primarily determined by the total number of users or monthly active users. Table 2 presents facts about social media sites user base which gives some sense of its popularity [4, 5, 11, 13]. The attribute *type* with values (a) *Total* (b) *Active* represents whether the statistic is of total users or active monthly users respectively.

Other than the social media sites mentioned in Table 2 there are some significant sites where only the current user statistics are available. For example Flickr, the photo sharing platform has 90 million users. Quora, a question answer social platform has 300 million users worldwide. Reddit, a social forum has 330 million active users.

Table 1: List of Relevant Social Media

Category	Social media sites with link
Social Networks	Facebook, Snapchat, WeChat, Quora
Private Messaging	Messenger, Whatsapp, QQ, WeChat, Skype
Microblogs	Twitter, Sina Weibo, Tumblr
Photo Sharing	Instagram, Photobucket, Flickr
Video Sharing	Youtube, Vimeo, Dailymotion
Professional Networks	LinkedIn, AngelList, Meetup
Enterprise Social Networks	Workday
Blogs	Wordpress, Medium, Buffer Blog
Forums	Reddit, Hacker News, Quora
Products/Services Review	Yelp, Foursquare, Google Places
Social Bookmarking	Pinterest, Digg, Stumble Upon Mix
Social Gaming	Pokemon Go, IGN, Gamespot [12]
Collaborative Projects	Slack, Invision, Trello, Github, Bitbucket

Table 2: Social media sites and number of users (in millions).

Category	Site	Years						Type
		2013	2014	2015	2016	2017	2018	
Social Networks	Facebook	1228	1393	1591	1860	2129	2271	Total
	WeChat	355	500	697	889	989	1082	Total
Microblogs	Twitter	241	284	305	318	330	332	Active
	Weibo	140	175	237	310	340	392	Active
	Tumblr	175	—	—	—	460	550	Total
Photo Sharing	Instagram	150	300	460	600	870	1000	Active
	Snapchat	33	100	180	301	—	400	Total
Video	Youtube	700	1100	1431	1618	1767	1900	Active
Professional	LinkedIn	277	347	414	467	530	576	Total
Services	Yelp	96	135	150	158	170	178	Active
	Foursquare	33	30	50	—	—	55	Total
	Ridesharing	—	—	208	272	338	400	Active
Bookmarking	Pinterest	—	—	110	160	220	250	Total

Table 3: Social media sites and media units created per day (in millions).

Category	Site	Years						Unit
		2013	2014	2015	2016	2017	2018	
Social	Facebook	3600	—	—	4320	—	—	posts
Microblogs	Twitter	144	399	500	—	657	682	tweets
	Tumblr	120 ^a	205 ^a	270 ^a	315 ^a	380 ^a	448 ^a	total blogs ^a
Photo Sharing	Instagram	5	31	—	—	—	67	photos
	Snapchat	—	—	—	—	760	3000	photos
Video	Youtube	69,120 ¹	103,680 ¹	432,000 ¹	576,000 ¹	720,000 ¹	—	hours video ¹
	Yelp	40 ²	55 ²	75 ²	95 ²	135 ²	171 ²	total reviews ²
Services	Foursquare	33	—	—	—	—	12000 ³	total checkins ³
Bookmarking	Pinterest	—	5	13	—	—	—	pins

Number of users is not just important to measure the popularity of a social media site but also to estimate the amount of data storage it maintains. Another feature that will help us to estimate data storage is the amount of media units (e.g. posts, photos, microblogs, videos etc.) ingested per day. Table 3 presents all the statistics of relevant social media from open internet [4, 5, 13, 15]. The statistics for social media sites missing in Table 3 but mentioned in Table 2 are almost impossible to find in open internet.

Social Media Storage Estimate:

Social media sites seldom reveals the amount of data they store or ingest on daily basis. Also the ever growing social media makes it hard to estimate their storage capacity. I present few methods in the following section to estimate social media storage.

1. Storage space estimate from media units: This method works for all the social media sites mentioned in Table 3 where the approximate storage space required by media unit is known.

Youtube: Lets take an example of Youtube video data. From the Table 3 we find by year 2017 users upload 720,000 hours of video in Youtube. First, assuming the fact that Youtube pretty much stores almost video in 1080p and it stores video in multiple resolution such as 240p, 360p, 720p, 1080p and format e.g. Webm, flv, mp4, 3gp, mp3. We can determine the amount of storage space needed for a 1 minute video [8].

$$\begin{aligned} & 27.71 \text{ MB (Webm)} + 17.00 \text{ MB (flv)} + 554.43 \text{ KB (3gp)} \\ & + 45.80 \text{ MB (mp4)} + 2.81 \text{ MB (mp3)} \\ & = 93.8614355 \text{ MB} \end{aligned} \quad (1)$$

From the above we find that $720,000 \times 60 \times 93.8614355 \approx 4.055$ petabytes (PB) of storage space is required by Youtube everyday. We can also calculate the total amount of storage space ingested during the period of 2013 to 2017 from Table 3 by utilizing area under the curve method with interpolation. The above method reckons 3096.17 PB or 3.096 exabytes (EB) of storage. Considering videos before 2013 and new 4K video which takes more space it can be easily assumed that Youtube use 10-15 EB storage space.

Twitter: Similar to the method above we can find the space required to store a tweet. A tweet is stored in Twitter as UTF-8 format. This takes 140 characters tweets atmost 560 bytes of space. However the metadata attached with a tweet is much more than the tweet itself. I personally did a random sample experiment of 100K tweets stored in our databases to find the average storage space for tweet json object obtained from streaming api. I find one json tweet object takes 3247 bytes of space in average. 682 million tweets per day will require around 2.2145 terabytes of data per day. Using the interpolation method for area under the curve we can find that Twitter use 3.13 petabyte of space for storing the tweet alone. It is also worth noting that 42% of tweets contains images [14]. If we assume the average image size be 100 KB then we will see $(100 * 1024)/3247 * 42\% \approx 13.2$ times increase in storage space requirement.

2. Storage space estimate from data center power usage: This section presents an approximate method to estimate space capacity of large social media companies like Facebook and Google. A typical breakdown of energy consumption by data center given in Figure ?? . The largest energy consuming component is cooling infrastructure with 50% of total energy. Rest of the energy is used by power conversion, lighting, network and server components [3, 7]. Facebook data centers use efficient data center architecture and hardware tweaks saves 8-12% of energy spent in cooling, 13-25% in power conversion, 10% in motherboard [6]. That implies atmost 11% more efficient than typical data centers. Hence, it can

be claimed that Facebook servers use 37-38% of energy. Considering Facebook's 138 MW Altoona data center equipped with 200 Watts servers each with 6×4 TB of HDD as used in their experiment for [6]. Assuming the datacenter is running at peak energy $(138 \times 0.37 \times 24)/200 = 6127200 \text{ TB} = 6.1272 \text{ exabytes (EB)}$. Taking all the data centers in consideration and diving them with replication factor we can estimate the storage capacity of Facebook. The analysis provided above supports news *Facebook Builds Exabyte Data Centers for Cold Storage* in 2013 [2].

Social Media Data for Researchers:

2 PART B.

Scraping Social Media:

Difficulties in Scraping:

Non availability of good free APIs:

API call limit:

Scraping Technology:

Legal Terms and Conditions:

IP block while scraping:

Missing link metadata:

REFERENCES

- [1] T. Aichner and F. Jacob. Measuring the degree of corporate social media use. *International Journal of Market Research*, 57(2):257–276, 2015.
- [2] DataCenterKnowledge. Data center knowledge, facebook builds new exabyte data center, 2018. <https://www.datacenterknowledge.com/archives/2013/01/18/facebook-builds-new-data-centers-for-cold-storage>, Online; accessed 27-Nov-2018.
- [3] M. Dayarathna, Y. Wen, and R. Fan. Data center energy consumption modeling: A survey. *IEEE Communications Surveys & Tutorials*, 18(1):732–794, 2016.
- [4] Domo. Domo, 2018. [Online; accessed 27-Nov-2018].
- [5] ExpandedRamblings. Expanded ramblings, 2018. <https://expandedramblings.com>, [Online; accessed 27-Nov-2018].
- [6] E. Frachtenberg, A. Heydari, H. Li, A. Michael, J. Na, A. Nisbet, and P. Sarti. High-efficiency server design. In *Proceedings of 2011 International Conference for High Performance Computing, Networking, Storage and Analysis*, page 27. ACM, 2011.
- [7] I.-T. R. Group et al. Top 10 energy-saving tips for a greener data center. *Info-Tech Research Group. fi...*, 2007.
- [8] jdownloader. Youtube jdownloader, 2018. <http://i.imgur.com/CgX5t.png>, Online; accessed 27-Nov-2018.
- [9] A. M. Kaplan and M. Haenlein. Users of the world, unite! the challenges and opportunities of social media. *Business horizons*, 53(1):59–68, 2010.
- [10] J. H. Kietzmann, K. Hermkens, I. P. McCarthy, and B. S. Silvestre. Social media? get serious! understanding the functional building blocks of social media. *Business horizons*, 54(3):241–251, 2011.
- [11] NetImperative. Sina weibo overtakes twitter user numbers, 2018. [Online; accessed 27-Nov-2018].
- [12] Statista. Most popular games in us, 2018. <https://www.statista.com/statistics/580150/most-popular-us-gaming-apps-ranked-by-audience/>, [Online; accessed 27-Nov-2018].
- [13] Statista. Statista, 2018. <https://www.statista.com>, [Online; accessed 27-Nov-2018].
- [14] Thenextweb. Next web, 2018. [Online; accessed 27-Nov-2018].
- [15] Zephoria. Facebook statistics, 2018. [Online; accessed 27-Nov-2018].