Basic data literacy for AI reporters

What we will cover in this module

- A brief and incomplete history of how data evolved
- Data in the realm of AI
- Three examples:
 - Structured data feeding decision-making algorithms
 - Longitudinal data fueling your social media timelines
 - Text data for large language models

Let's start with an exercise

What kind of data is out there?

Let's break down a tweet.

- Read this.
- Laugh a little.
- Then tell me what kind of data this tweet may yield.





Purebread dogs vs. inbread dogs



RETWEETS 1,446

LIKES 2,758 🚵 👔 🚮 🔯 🌘 B. 📆 🐧 🕥













12:25 PM - 22 Feb 2017









Person who posted it: Picture link Display name Twitter handle

Text

Media

Engagement: Retweets Likes Replies (further down) Timestamp





Purebread dogs vs. inbread dogs



RETWEETS 1,446

LIKES

2,758















12:25 PM - 22 Feb 2017







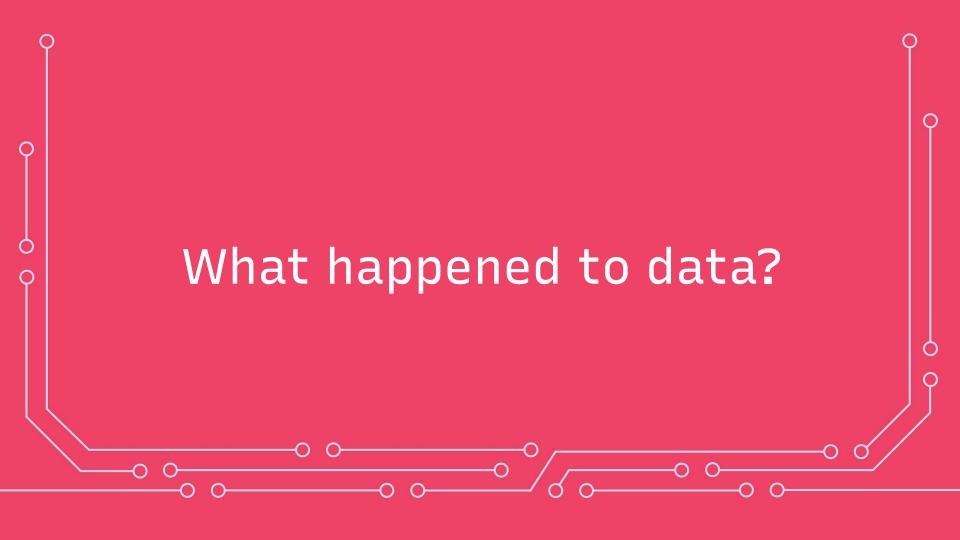




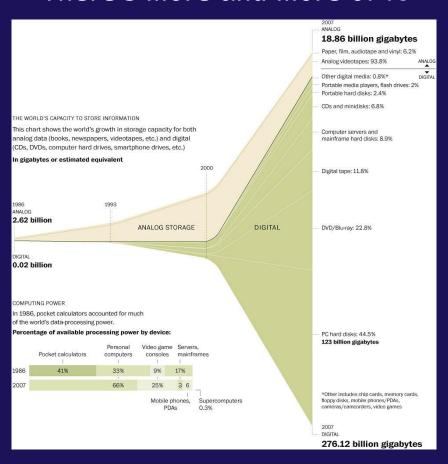
		"contributors": null,
		"truncated": false,
		"text": "Purebread dogs vs. inbread dogs <a "="" href="https://t.co/e1nC0d7EkR">https://t.co/e1nC0d7EkR ",
		"is_quote_status": false,
		"in_reply_to_status_id": null,
		"id": 834454023260532737,
		"favorite_count": 2760,
Camo twoot		"source": " Twitter for
Same tweet,		Android",
different format! The		"retweeted": false,
entire X.com site is		"coordinates": null,
		"entities": {
fueled by data		"symbols": [],
streams like this.	14	"user_mentions": [],
		"hashtags": [],
		"urls": [],
APIs or Application		"media": [{
Programming Interface*		"expanded_url": "https://twitter.com/elle91/status/834454023260532737/photo/1",
Programming interface		"display_url": "pic.twitter.com/e1nC0d7EkR",
		"url": "https://t.co/e1nC0d7EkR",
		"media_url_https": "https://pbs.twimg.com/media/C5STfk1WQAEzRet.jpg",
*fancy word for a stream		"id_str": "834454001643044865",
of data in a format that		"sizes": {
		"large": {
robots/computers can		"h": 413, "resize": "fit", "w": 550
understand 🤖		}
		"small": {
		"h": 413, "resize": "fit", "w": 550
		}
		"medium": {
		"h": 413, "resize": "fit", "w": 550

Data is everywhere!

The world wide web is fueled by data

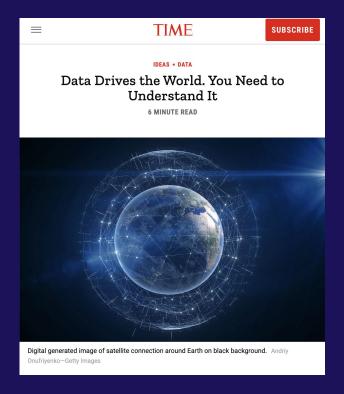


There's more and more of it



Source: Washington Post (2011)

There's a lot more than just government data now





Source: <u>IBM</u> via <u>Time</u> Source: <u>CNET</u>

Digital footprint example: WSJ



How Pizza Night Can Cost More in Data Than Dollars

Even a low-key evening at home can mean handing over a trove of personal information to high-tech companies

By Stephanie Stamm, Tripp Mickle and Jessica Kuronen
Published April 10, 2018 at 5:30 a.m. ET

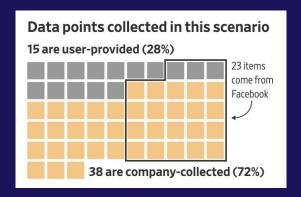
The smartphones, Facebook accounts and other technology products deeply embedded in modern life help people get more things done every day. They also gather more information about us than we often realize.

But Facebook's crisis over how it handles and protects user data has led some to ask: What data am I giving up?

Imagine "Sally" sets up a pizza-and-movie night with her friend "Kristen." The Wall Street Journal reviewed privacy statements to assess just how much data could be unknowingly shared on top of the price of that pepperoni pie.

The Cost in Data

Sally and Kristen potentially gave up at least 53 pieces of information together. The data detailed in the scenario reflect information the companies could collect according to their privacy statements, terms of service and related documents.



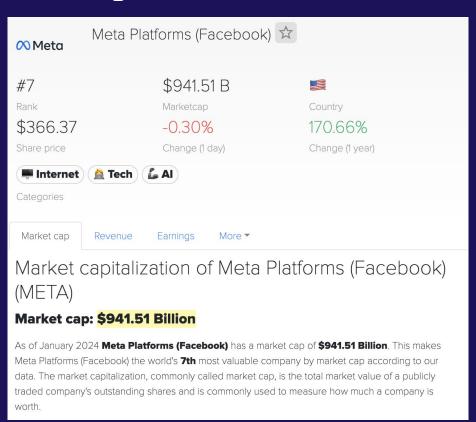
Source: WSJ

Digital footprint example: Meta data servers



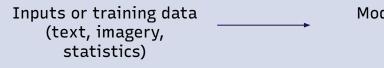
Big Data

It's making a lot of people a lot of money



Data in the realm of AI

Lifecycle of AI Where in the life cycle is training data?



Model/Algorithm

Outputs
('predictions,'
evaluative scores, new
images, text)





AI

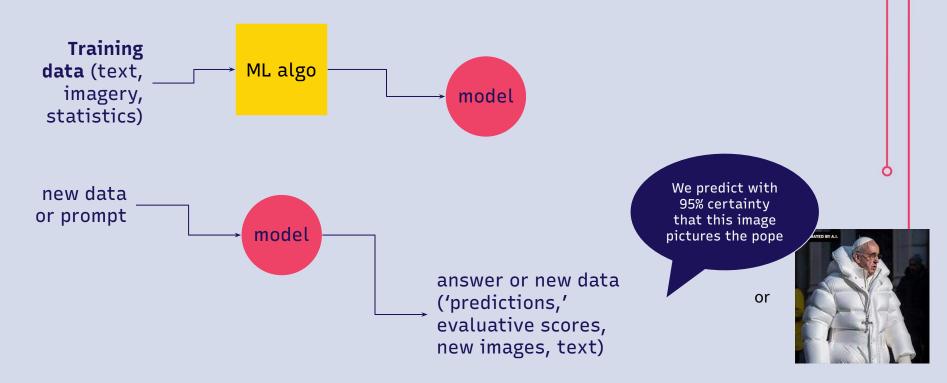
We predict with 95% certainty that this image pictures the pope

or

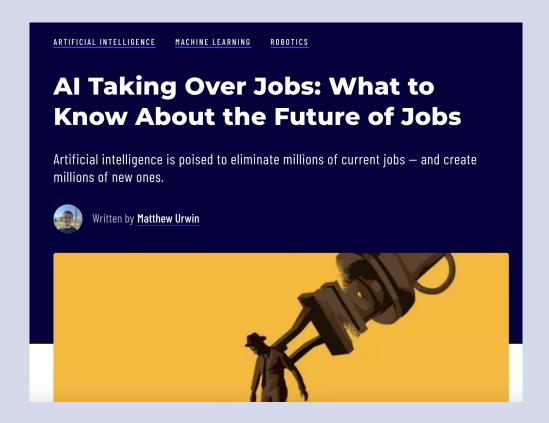


Lifecycle of AI

Where in the life cycle is training data?



AI: Data is changing our society

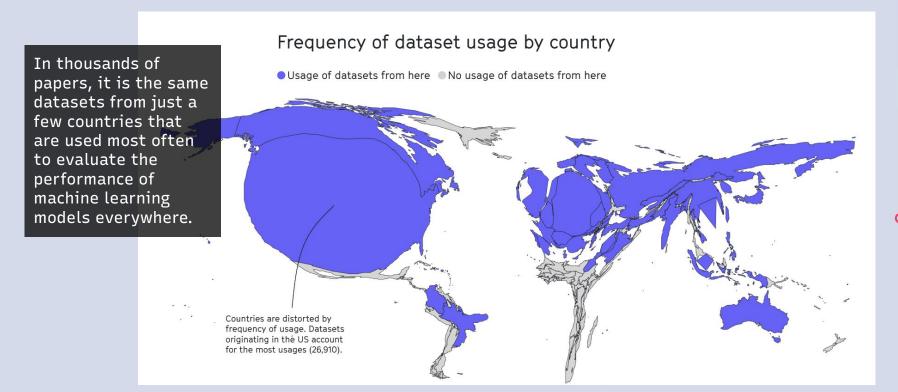


AI: Data is changing our society



Source: <u>Lighthouse News</u>

A lot of the data is reproducing bias



Source: Mozilla

Structured data feeding

decision-making algorithms

Under Haryana's Old Age Samman Allowance scheme, people aged 60 years and above, whose income together with that of their spouse doesn't exceed 300,000 rupees (\$3,600) per annum, are eligible for a monthly pension of 2,750 rupees (\$33).

In June 2020, the state started using a newly built algorithmic system — the Family Identity Data Repository or the Parivar Pehchan Patra (PPP) database — to determine the eligibility of welfare claimants.



Dhuli Chand led a wedding procession to prove to government officials he was alive and should be paid his pension. Image courtesy of The Reporters' Collective.

Algorithms and welfare Source: Pulitzer Center

The PPP is an eight-digit unique ID provided to each family in the state and has details of birth and death, marriage, employment, property, and income tax, among other data, of the family members. It maps every family's demographic and socioeconomic information by linking several government databases to check their eligibility for welfare schemes.

The state said that the PPP created "authentic, verified and reliable data of all families", and made it mandatory for citizens to access all welfare schemes.

But in practice, the PPP wrongly marked Chand as "dead", denying him his pension for several months. Worse, the authorities did not change his "dead" status even when he repeatedly met them in person.



Dhuli Chand led a wedding procession to prove to government officials he was alive and should be paid his pension. Image courtesy of The Reporters' Collective.

Algorithms and welfare Source: Pulitzer Center

Longitudinal data fueling your

social media timelines

Data used by Facebook

Folders that are included in a downloadable Facebook data archive

The Facebook algorithm is a ranking system that uses machine learning to arrange content in users' feeds. Data used:

- Data you put into Facebook
- Data that you produce by using Facebook
- Data that is tracked of you by pages with Facebook buttons

To get an idea of how deep this data goes users can download their data following Facebook's instructions:

https://www.facebook.com/help/1701730696756992?helpref=hc global nav

More on the subject here:

How Facebook watches you online (<u>BuzzFeed</u> News)

3 Simple Ways We Give Up A Ton Of Very Personal **Information To Facebook And Random Apps** (BuzzFeed News)

about_you ads apps_and_websites calls_and_messages comments events

> following_and_followers friends groups

https/ likes and reactions location_history

marketplace messages

other_activity

pages payment_history

photos and videos posts

profile information

saved items search history

security_an...information your_places

止の奇色

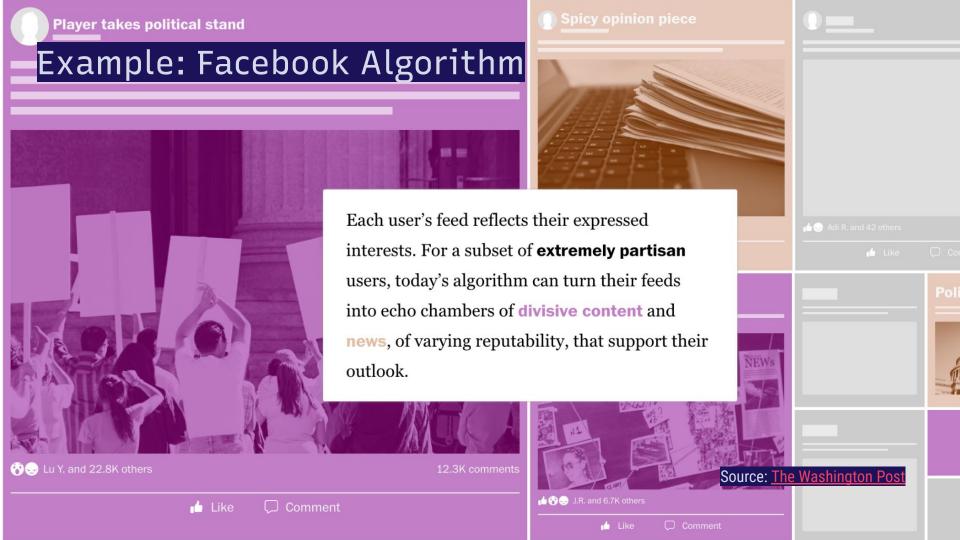
Example: Facebook Algorithm

The top post on a Facebook user's news feed, shown as the biggest box, is a prized position based on thousands of data points related to the user and post itself, such as the poster, reactions and comments.

#000

Source: The Washington Post





Text data for natural

language processing

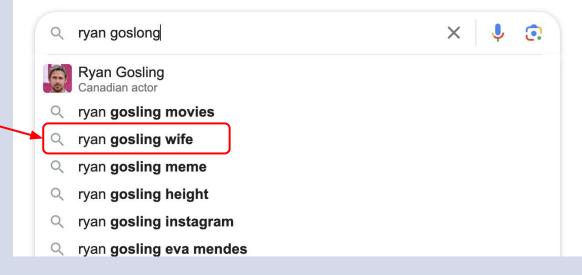
Example: Google

Predictive text in Google search is based on countless searches.

How many broken-hearted people googled Ryan Gosling's name to find out whether he was married?! (Plus a lot of people likely also googled a male celebrity name + wife)

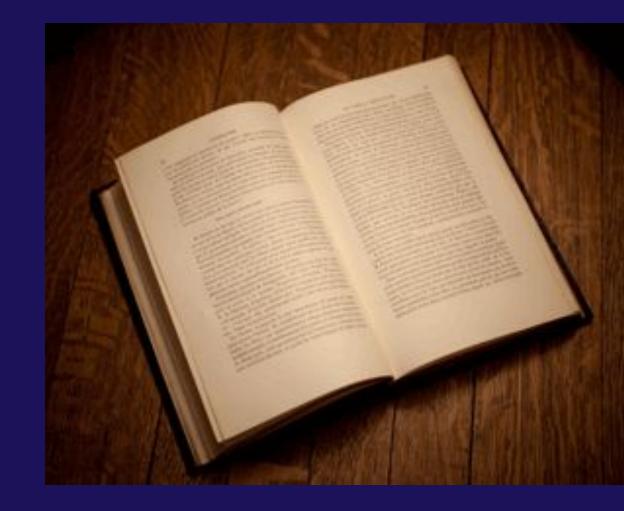


Google



- 1. Corpus
- 2. Tokenizing
- 3. Stemming or lemmatizing

- 1. Corpus
- 2. Tokenizing
- 3. Stemming or lemmatizing



- 1. Corpus
- 2. Tokenizing
- 3. Stemming or lemmatizing



- 1. Corpus
- 2. Tokenizing
- 3. Stemming or lemmatizing

The goal of both stemming and lemmatization is to reduce inflectional forms and sometimes derivationally related forms of a word to a common base form.

For instance:

am, are, is → be car, cars, car's, cars' → car

The result of this mapping of text will be something like:

the boy's cars are different colors

→ the boy car be differ color

- 1. Corpus
- 2. Tokenizing
- Stemming or lemmatizing

Stemming vs. lemmatization

Example: "saw"

(As in "I saw a bird, or was it a plane? I think it may have been Supergirl.")

Stem:

saw → s

Lemma:

saw → see*

* if the token is a verb!

Analyzing text

Once you've prepared your corpus into tokens (either as stems or lemmas) you can then analyze the text. This could include:

- Creating lists of words that are used the most
- Predicting which words are most likely to follow one another (just like in the good search!)
- Detecting sentiment in a text



