

# Biotechnology Big Data Artificial Intelligence

Promises, Dangers, Worries, Gotchas

Slides available at [https://webbedfeet.github.io/GMU\\_2019](https://webbedfeet.github.io/GMU_2019)

Abhijit Dasgupta, PhD  
Chief Data Scientist  
Co-Founder





**Biotechnology**

also includes techniques called inserting genes into plants to produce desired products. This is called genetic engineering. It involves manipulating the DNA of an organism to produce a specific product or trait. The process can be used to create new varieties of crops that are resistant to pests or diseases, or to produce pharmaceuticals like insulin or interferon. Biotechnology has also been used to develop new treatments for diseases like cancer and AIDS.

Biotechnology is a rapidly growing field with many applications in agriculture, medicine, and industry. It has the potential to improve our lives in many ways, but it also raises important ethical and social issues that need to be addressed.

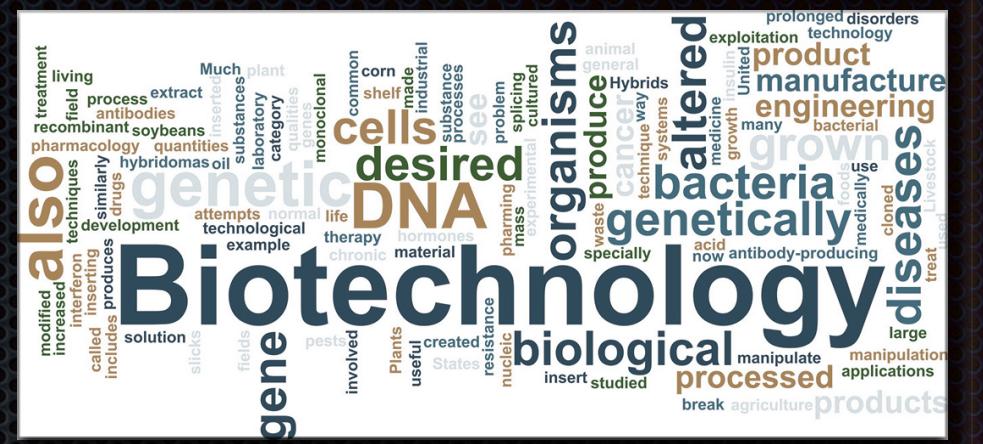
- Evidence of effectiveness
- Reproducibility
- Manufacturability
- Quality

# Context (why we need it)

- Market research
- Competition research
- Focus groups
- Gaps in the market
- Market potential

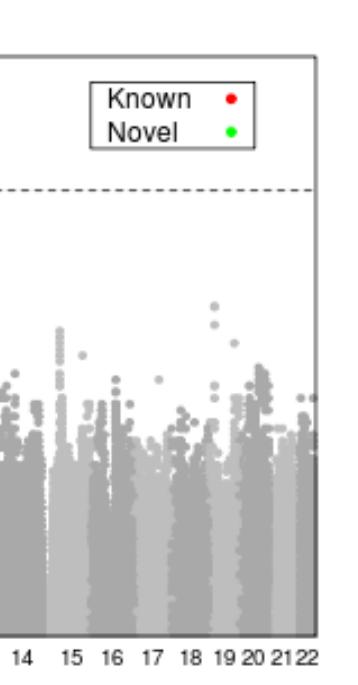
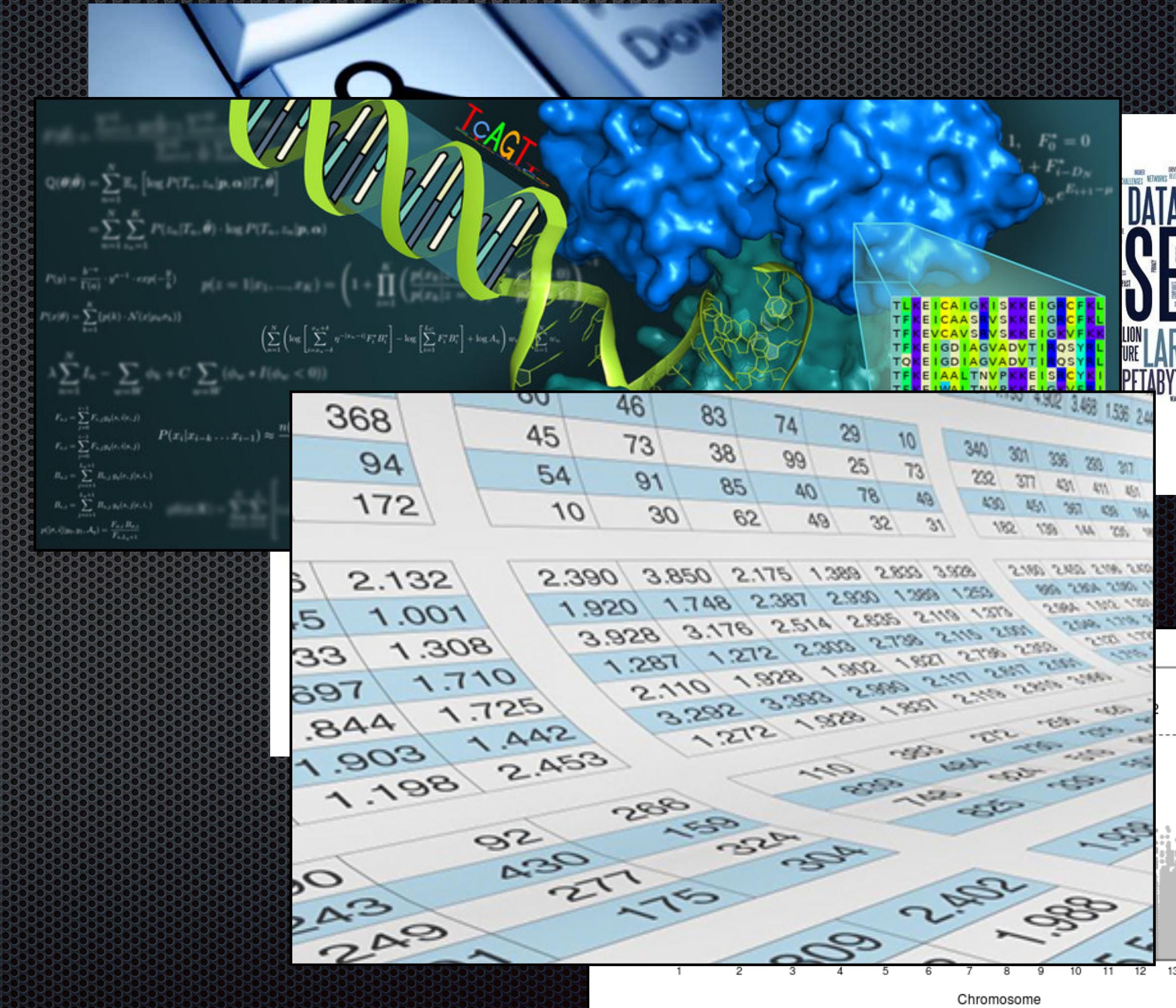


368	60	46	83	74	29	10	340	301	326	289	317	130	4,902	3,468	1,536	2,400
94	45	73	38	99	25	73	232	377	431	411	451	130	232	377	431	411
172	54	91	85	40	78	49	430	451	397	438	394	130	430	451	397	438
	10	30	62	48	32	31	182	139	144	128	140	130	182	139	144	128
6	2.132						2,390	3,850	2,175	1,389	2,833	3,928	2,180	2,820	2,190	2,820
45	1.001						1,920	1,748	2,387	2,930	1,389	1,250	190	2,384	1,250	190
33	1.308						3,928	3,176	2,514	2,835	2,110	1,273	2,944	1,942	1,273	2,944
697	1.710						1,287	1,272	2,303	2,738	2,115	2,381	2,046	1,287	2,303	2,046
.844	1.725						2,110	1,928	1,982	1,821	2,739	2,117	2,381	2,046	1,928	1,982
1.903	1.442						3,292	3,393	2,966	2,117	2,937	2,118	2,381	2,046	3,393	2,117
	1.198	2.453					1,272	1,928	1,857	2,117	2,937	2,118	2,381	2,046	1,928	1,857
							110	329	232	276	153	424	110	329	232	276
							809	424	504	524	539	424	809	424	504	524
							743	277	175	324	304	277	743	277	175	324
							209	2,402	1,988	1,928	1,928	1,928	209	2,402	1,988	1,928
							249	243	241	241	241	241	249	243	241	241



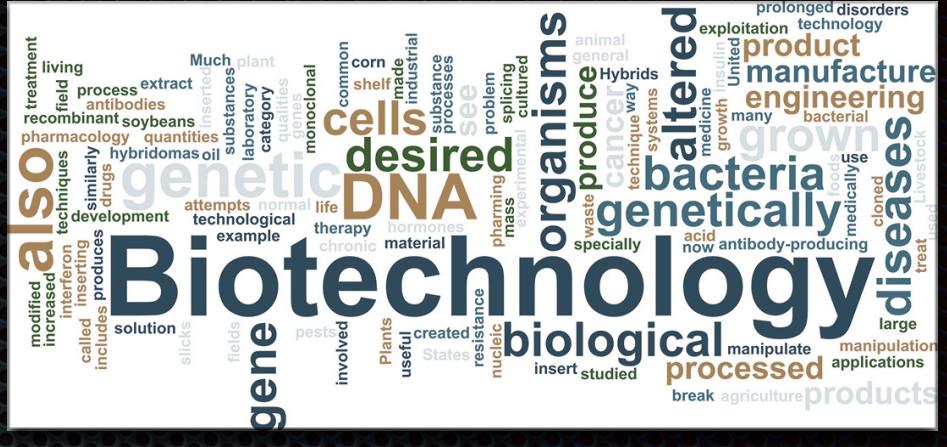
# Discovery (finding it)

- Identifying potential
- Searching the literature
- Searching databases
- Analyzing relationships
- Culling candidates



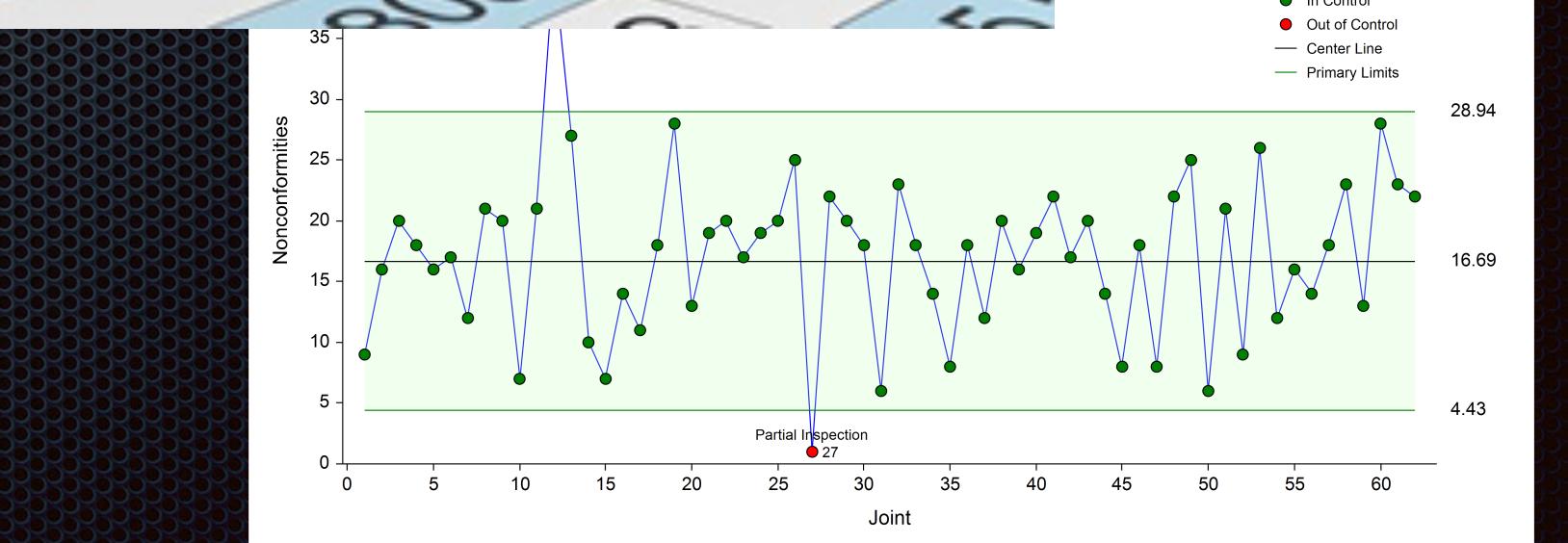
# Evidence of effectiveness

- Laboratory experiments
- Field experiments
- Observational studies
- Randomized trials



# Manufacturability & Quality

- Translation to a process
  - Establishing standards
  - Reproducibility
  - Quality control / Six Sigma







# Data drives

# Research

# Development

# Production

# Marketing

# Sales

# Revisions



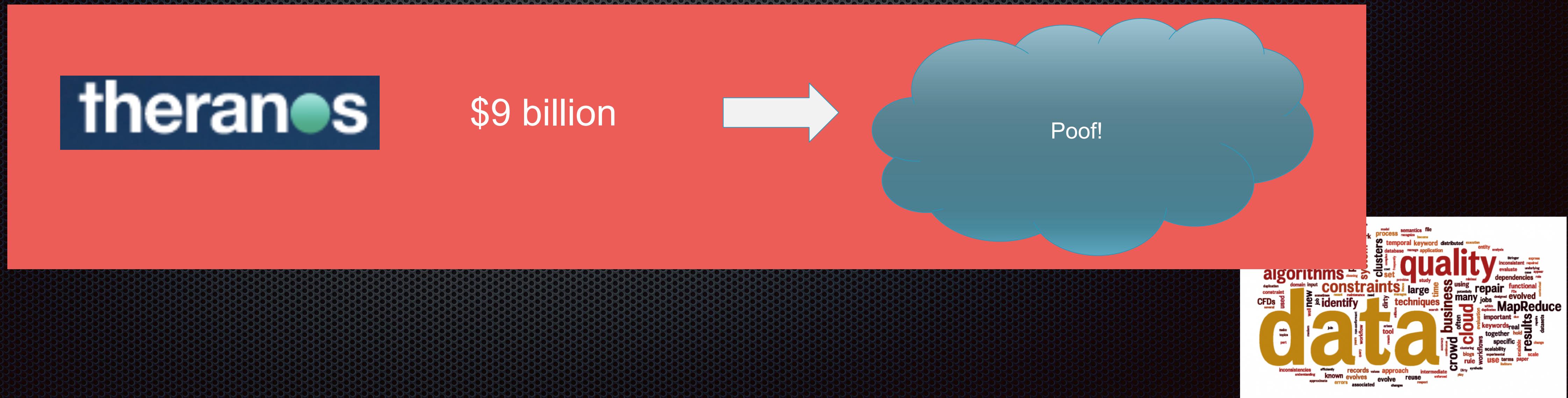
# Central issue for data-driven business

- The data must be reliable and valid
    - Repeating the experiment must yield similar results (reliability)
    - We're measuring what we think we're measuring (internal validity)
    - The results can be generalized to larger populations (external validity)



# Central issue for data-driven business

- This is a matter of trust
  - If the foundational data on which the business is built isn't trustworthy, you can't sell anything



# Central issue for data-driven business

- It doesn't matter how much data you can get
  - It matters how you got it
    - Experimental design
    - Ethics
    - Provenance



# Central issue for data-driven business

- It doesn't matter how much data you can get
  - It matters how you use it

# Analytics

- Translation to business practice
  - Follow-up



system automation engineering brain interface mechanical communication

internet relations business machine robotic

cybernetics cyber modern scifi cyborg human

technology science design

artificial robot concept tech

intelligence

information future computer virtual digital

cooperation futuristic people

mind man connection

internet business line  
cyber modern cyborg human  
scifi electronic creation  
**technology science design**

system automation  
engineering brain in  
mechanical commun  
ar  
in.  
  
information f  
coo

The word cloud is centered around the word 'learning'. Other prominent words include 'data', 'analysis', 'research', 'processing', 'mining', 'statistics', 'machine', 'intelligence', 'systems', 'modelling', 'deep-learning', 'bigdata', and 'data'. The words are in various sizes and colors (black, red, orange, brown), and smaller text provides additional context or related terms.

learning

data

analysis

research

processing

mining

statistics

intelligence

systems

modelling

deep-learning

bigdata

data

intelligence

systems

modelling

deep-learning

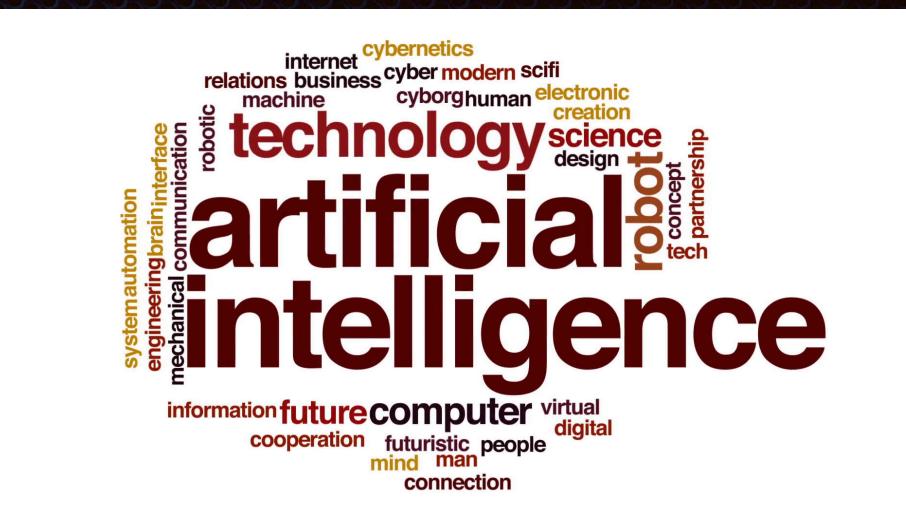
bigdata

data

**data** collection interpretation  
probability prediction experiments statistic  
population science theory inference  
science experiments theory probability  
mathematical statistics  
statistical inference

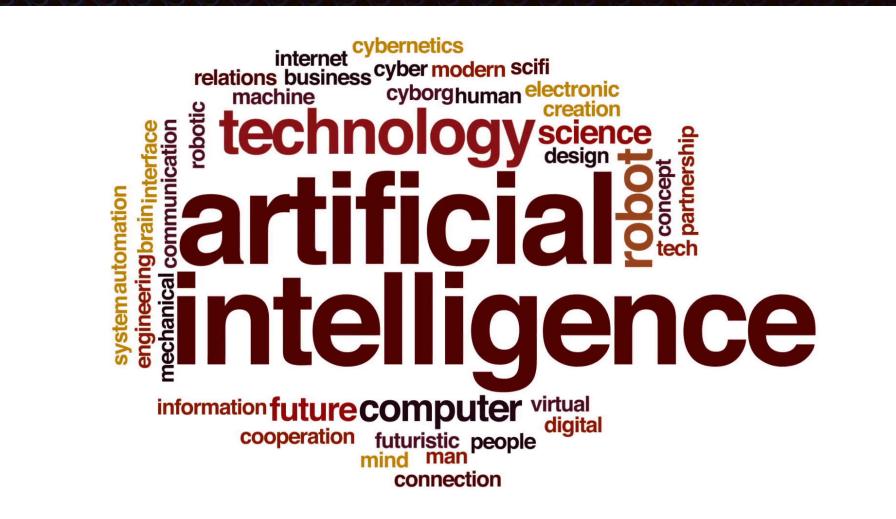
# What is AI?

- The ability of machines (computers) to learn, reason and act on their own
- The ability of machines to demonstrate intelligence
- The ability of machines to be autonomous, independent, and make decisions based on their environment



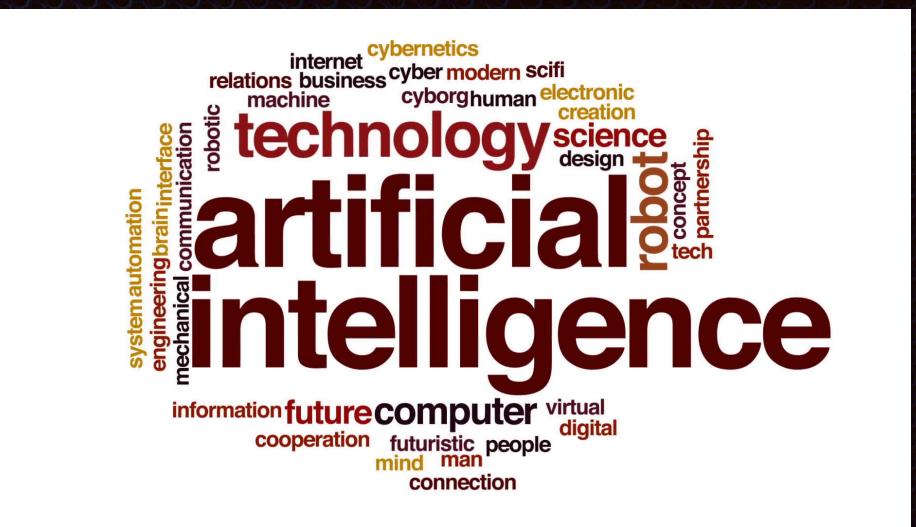
# AI in Biotech

- Drug discovery
  - Culling through millions of candidates to identify the next molecule for drug development
- Patient management and decision support
  - Fast analytics to provide likely causes, drug effects and interactions, risks and side-effects based on patient history, EMR mining and literature mining



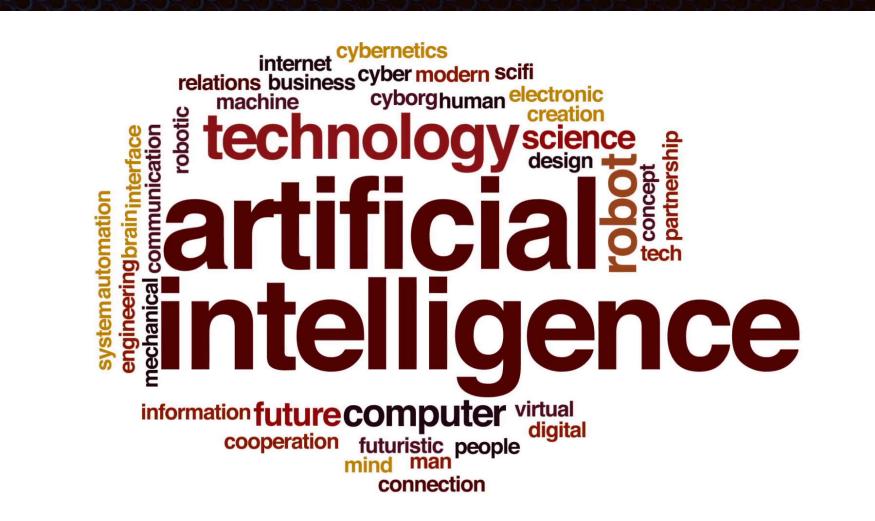
# AI in Biotech

- ❖ Prosthetics/exoskeletons and predictive movement
- ❖ Environmental assessment and closed-loop therapeutic deployment
- ❖ Performance optimization through sensing, adaptation and intelligent coaching
- ❖ Precision medicine
  - ❖ Identifying therapies better targeted at the individual

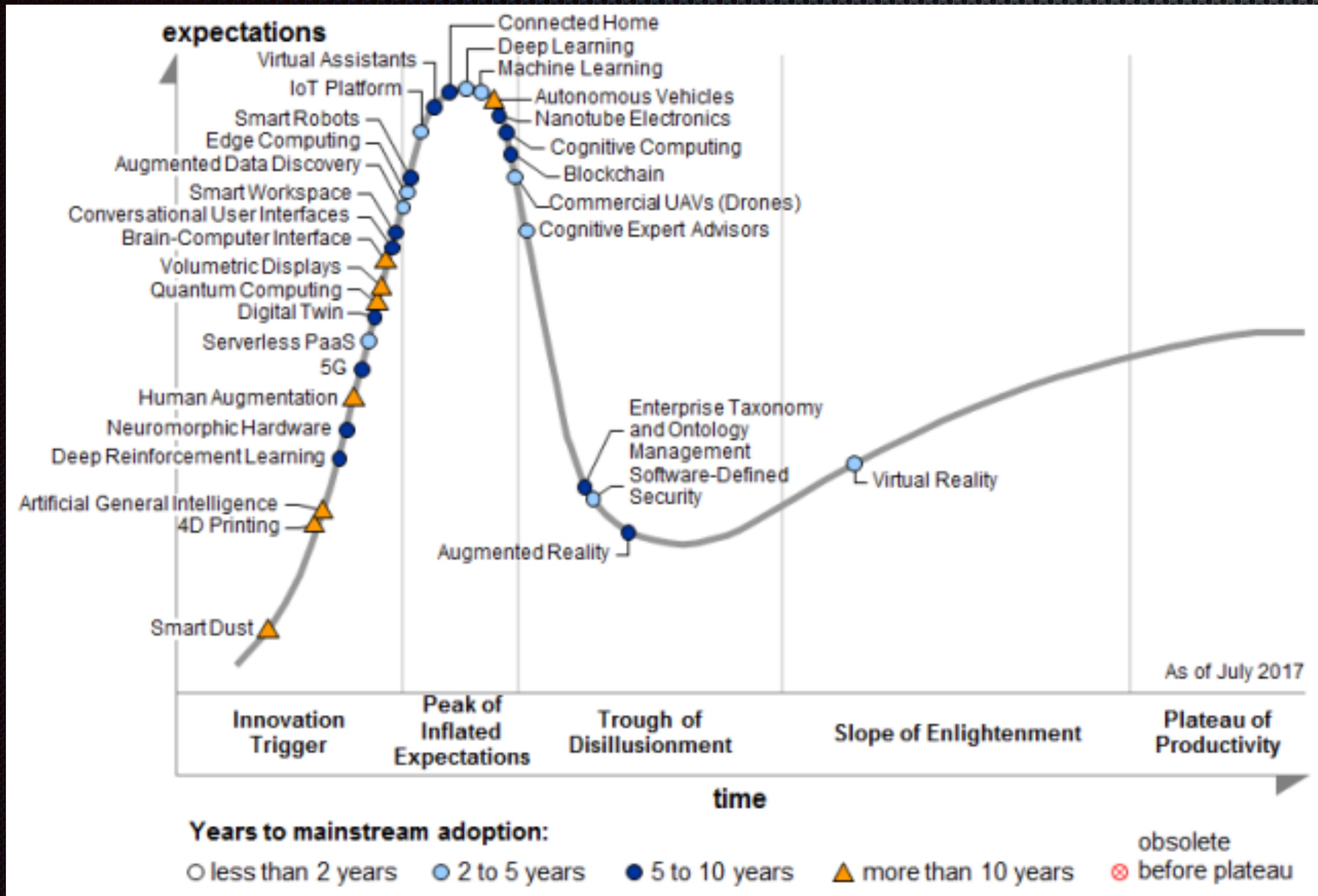


# Are we there yet?

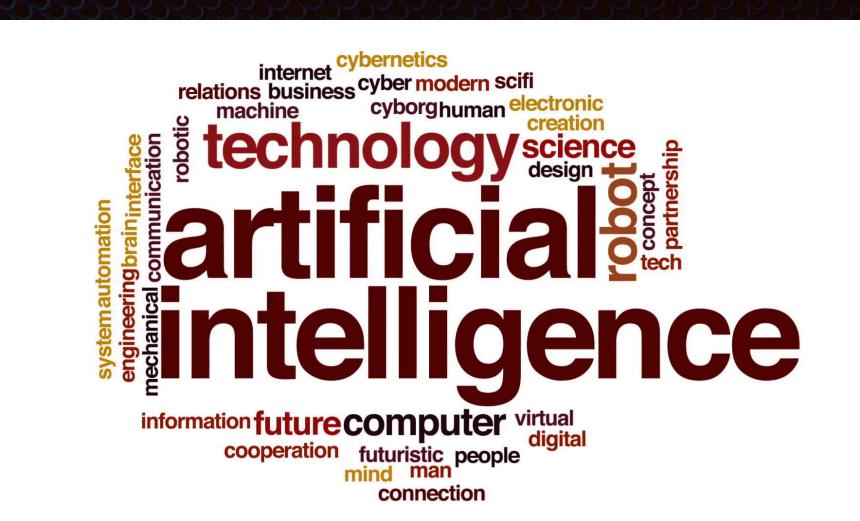
- Not really
  - Self-driving cars and planes on auto-pilot
  - Experimental situations
  - Some medical decision support



# Are we there yet?



## Gartner Hype Cycle 2017



# So where are we?

system automation  
engineering brain interface  
mechanical communication

# artificial intelligence

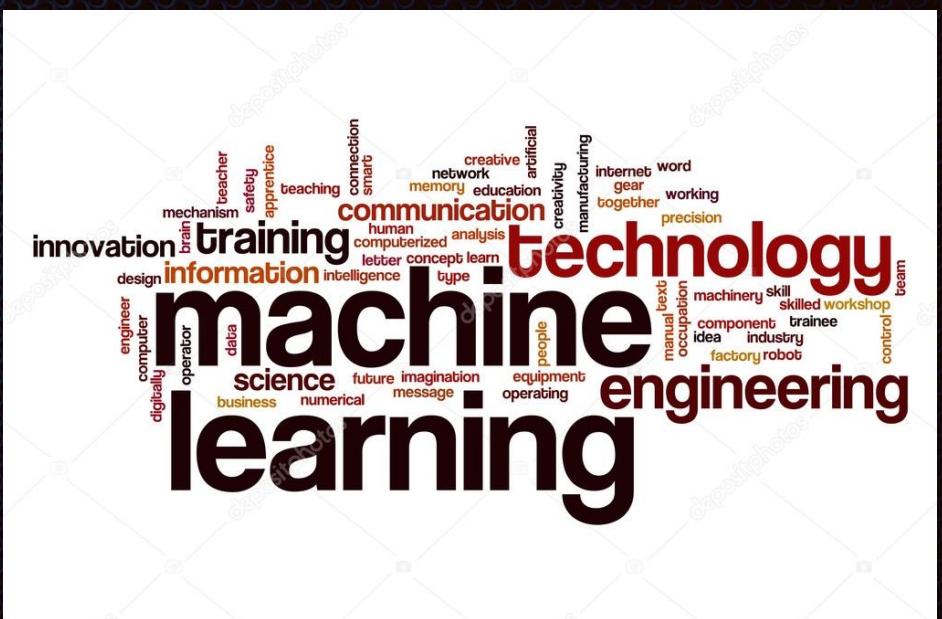
information future computer virtual  
cooperation futuristic people digital  
mind man connection

internet cybernetics  
relations business cyber modern scifi  
machine cyborg human electronic creation  
robotic science design  
robot concept tech partnership

# Common use

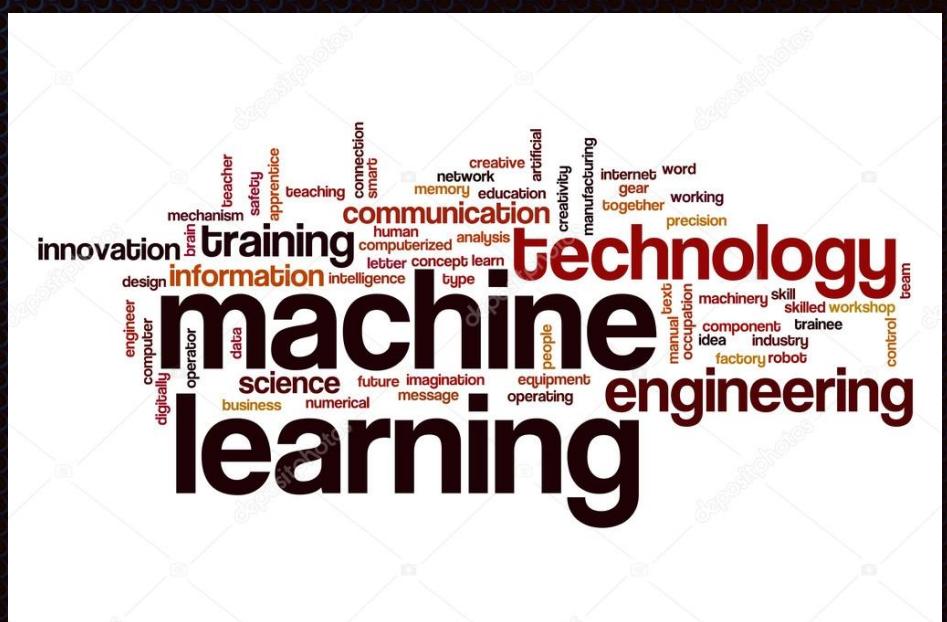
# Machine learning

- Artificial intelligence is about learning, reasoning and acting
  - We are still trying to get the learning part down well
    - Learning innate patterns in data
    - Learning how an outcome is influenced by different factors
    - Learning to predict the future by observing the past
    - Learning to discern what is real and what is artifact



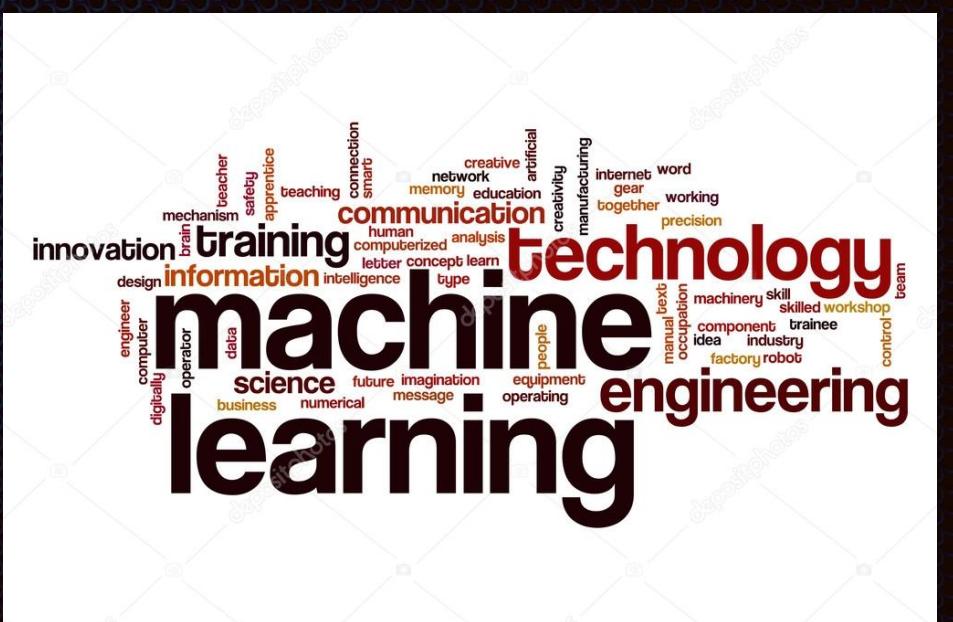
# Machine learning success stories

- Breast cancer therapeutics
    - Bioinformatics used to understand innate differences in different types of BrCa, leading to very successful targeted therapeutics
    - Tamoxifen, Herceptin, Fluoracil, Paclitaxel



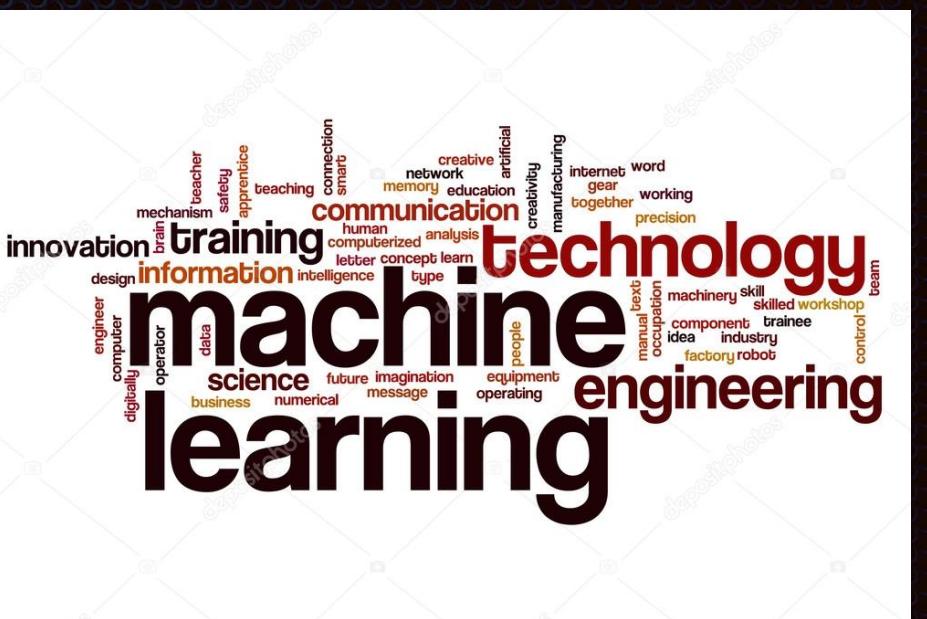
# Machine learning success stories

- Identification of subgroups in various diseases, leading to more targeted interventions
  - EGFR & Lung Cancer
  - APoE & Alzheimer's prevention
  - BRCA & Breast cancer prophylaxis
  - Immunotherapies
  - Gene editing (CRISP-R)



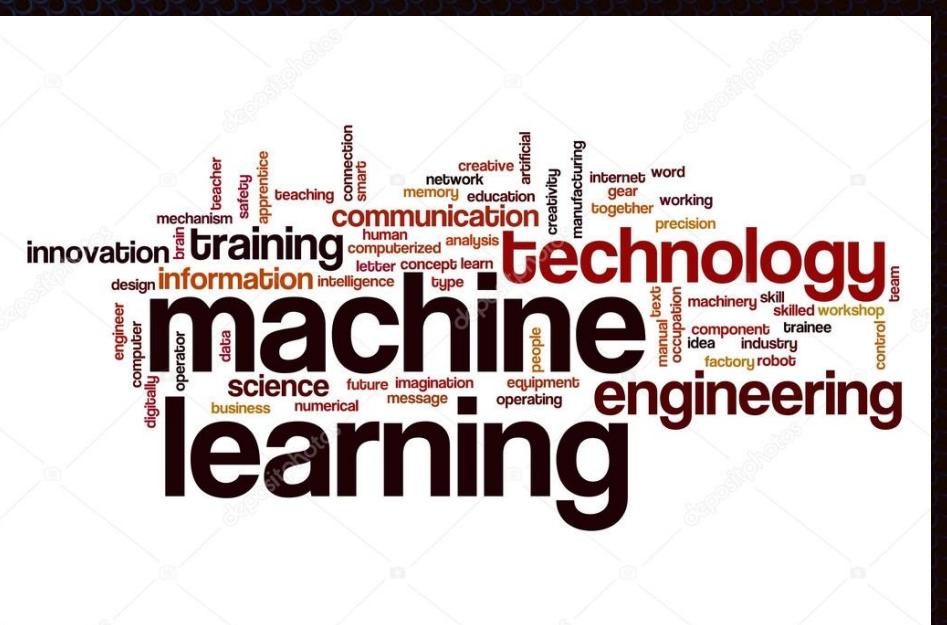
# Machine learning success stories

- Risk assessment, epidemiology and causality
  - Cancer risk prediction
    - Gail model for Breast Cancer
  - Impact of inflammation on health
  - Identification of genetic and proteomic risk factors in complex disease



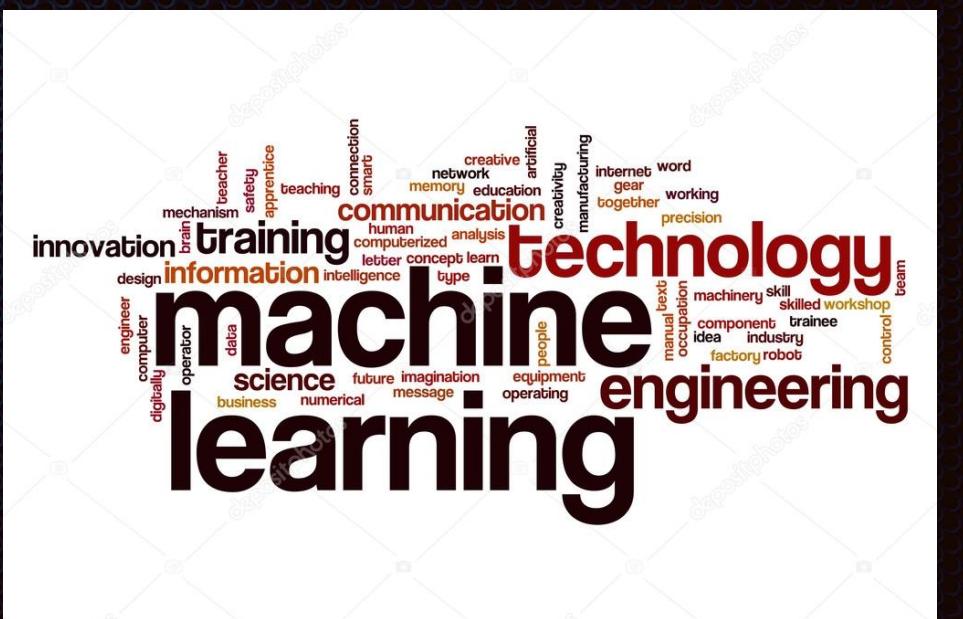
# What is ML?

- Algorithms that ingest data and adapt to mimic and express inherent data patterns
  - Learns from data
  - Can provide predictions
- Algorithms can be
  - Simple (linear regression)
  - Complex (decision trees)
- Indecipherable or magical (neural networks / deep learning)

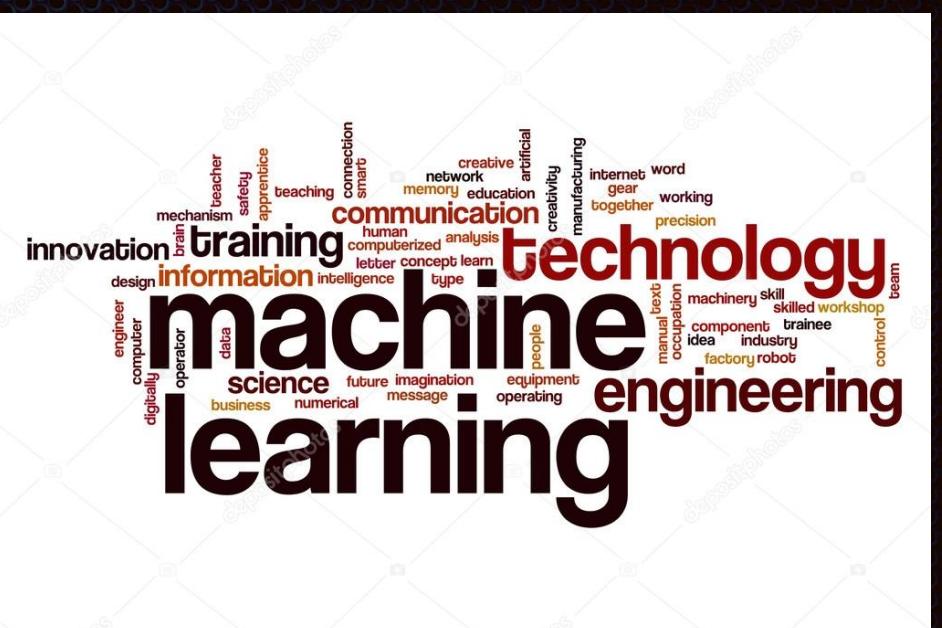


# What is ML?

- Algorithms that are mainly meant for prediction
  - Supervised learning
  - They cannot often express why the predictions happen
    - What factors influenced the prediction in what way?
    - It's a function of the complexity of the model and our inability to express or comprehend that complexity easily

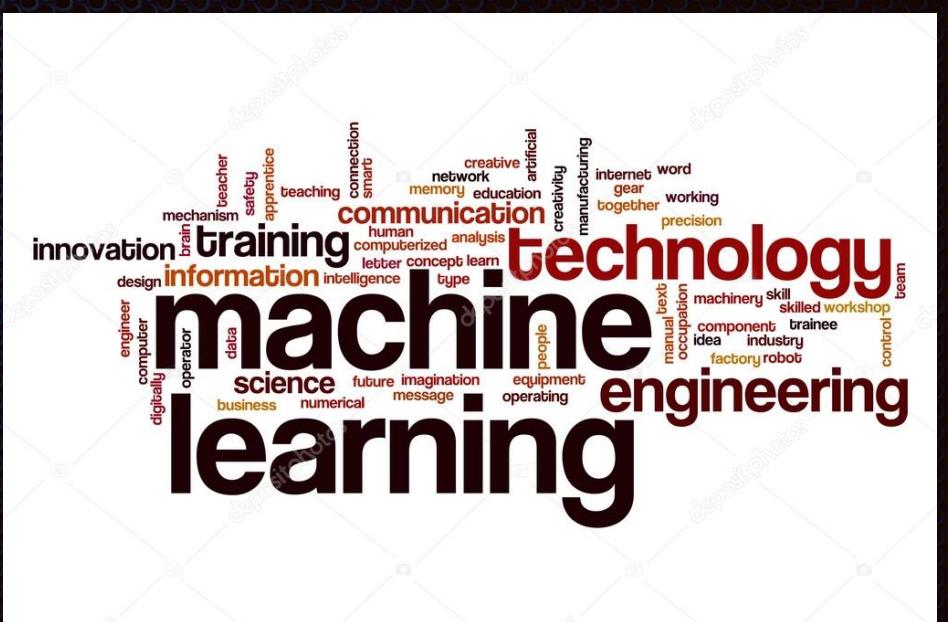


# ML and Big Data



# ML and Big Data

- Complex models require a lot of data to work well
    - Deep learning works for Google and Facebook because of their data stores
    - So we have been collecting lots of data in the hope (and prayer) that we can learn from it
  - And, by lots, I mean....

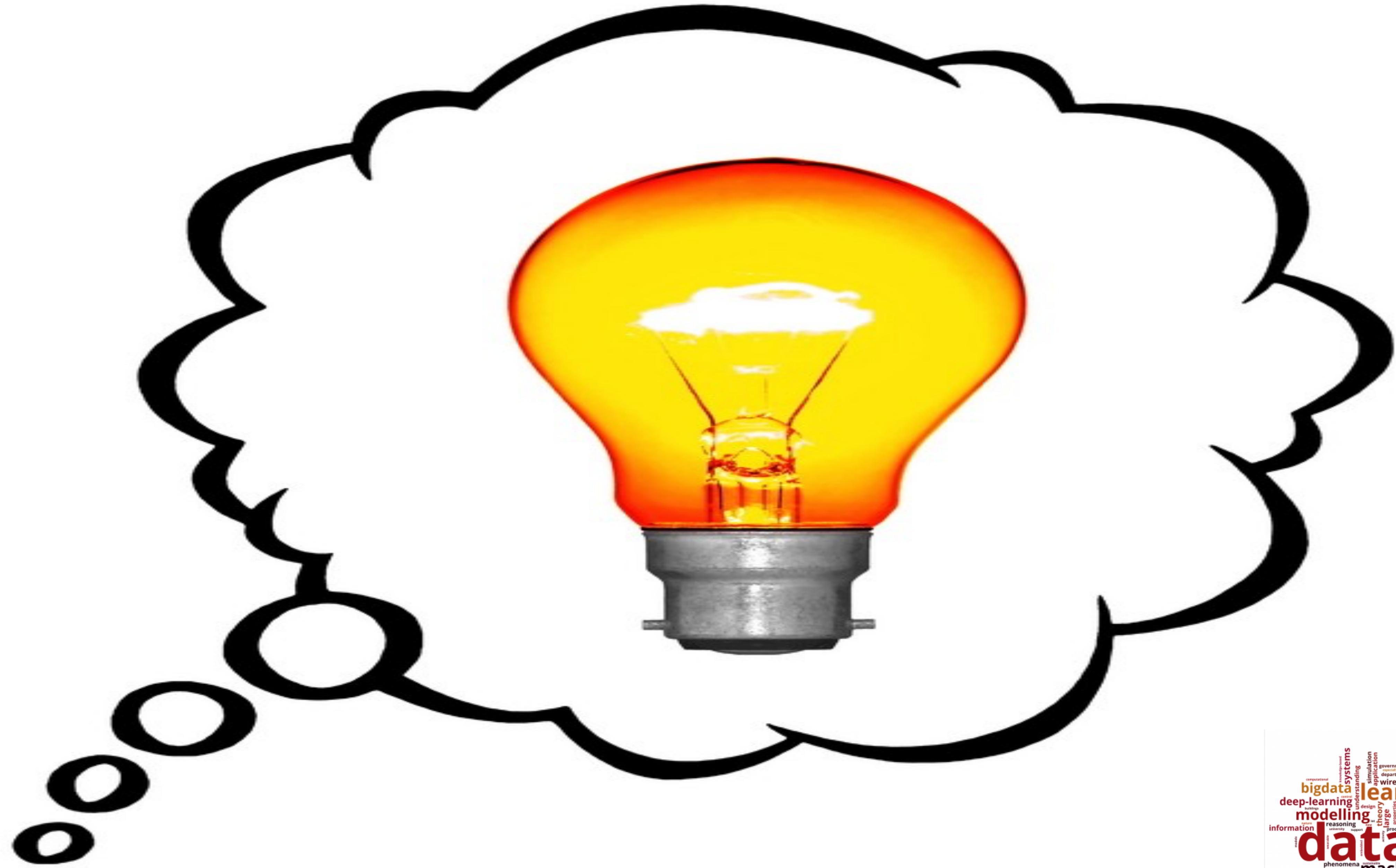




# Big Data in Drug Discovery

- Pharma is losing patent control
- Getting a drug to market requires many years and billions of dollars
- However, we have plenty of clinical, bioinformatic, molecular and trial data





# A bright idea

- Let's mine this data for good candidate molecules
  - No assays, just data
  - To confirm, outsource to do targeted assays



Grab public data.  
Start a company.  
Sell the company.  
Repeat

*Dr. Atul Butte, UCSF, 2015*

Data sitting still is useless  
Data needs to stay in motion  
Data needs to get us someplace



# Where's the data



# Where's the data

- Our bodies
    - Genetic/proteomic/molecular code (1.5 Gb per cell)
    - Movement / behavior
      - Sleep
      - Gait
      - Sounds



# Where's the data

- Interactions with the world
    - Health records and Insurance
    - Government data
    - Telecommunication
    - Social media
    - Economic behavior



# Where's the data

- Data we voluntarily give, maybe inadvertently
    - Sensors we wear and carry
    - Movement through GPS (Google Maps, anyone)
    - Electronic communication
    - Economic activity



# Where's the data

- Data we voluntarily give, by choice
    - Medical and health data
    - Surveys



# All told, 2 Zb\* per year

\*A Zb is  $10^{21}$  bytes



# Can we trust scientific discoveries made using machine learning?

Rice U. expert: Key is creating ML systems that question their own predictions

RICE UNIVERSITY

# **Statistician: Machine Learning Is Causing A “Crisis in Science”**

Many researchers now use machine learning to analyze data.  
There's just one glaring problem.

Jon Christian

February 18th 2019

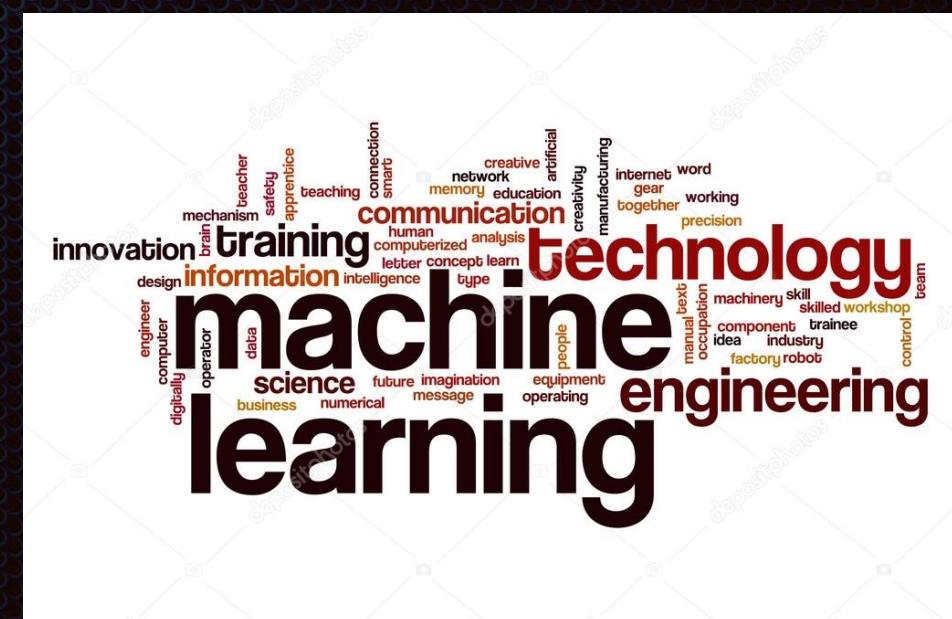
# February 16, 2019

Science & Environment

# AAAS: Machine learning 'causing science crisis'

By Pallab Ghosh

Science correspondent, BBC News, Washington



*open access, freely available online*

Essay

# Why Most Published Research Findings Are False

John P. A. Ioannidis

## Summary

There is increasing concern that most current published research findings are false. The probability that a research claim is true may depend on study power and bias, the number of other studies on the same question, and, importantly, the ratio of true to no relationships among the relationships probed in each scientific field. In this framework, a research finding is less likely to be true when the studies conducted in a field are smaller; when effect sizes are smaller; when there is a greater number and lesser preselection of tested relationships; where there is greater flexibility in designs, definitions, outcomes, and analytical modes; when there is greater financial and other interest and prejudice; and when more teams are involved in a scientific field in chase of statistical significance.

Simulations show that for most study designs and settings, it is more likely for a research claim to be false than true.

Moreover, for many current scientific fields, claimed research findings may often be simply accurate measures of the prevailing bias. In this essay, I discuss the implications of these problems for the conduct and interpretation of research.

**P**ublished research findings are sometimes refuted by subsequent evidence, with ensuing confusion

factors that influence this problem and some corollaries thereof.

## Modeling the Framework for False Positive Findings

Several methodologists have pointed out [9–11] that the high rate of nonreplication (lack of confirmation) of research discoveries is a consequence of the convenient, yet ill-founded strategy of claiming conclusive research findings solely on the basis of a single study assessed by formal statistical significance, typically for a  $p$ -value less than 0.05. Research is not most appropriately represented and summarized by  $p$ -values, but, unfortunately, there is a widespread notion that medical research articles

**It can be proven that most claimed research findings are false.**

should be interpreted based only on  $p$ -values. Research findings are defined here as any relationship reaching formal statistical significance, e.g., effective interventions, informative predictors, risk factors, or associations. “Negative” research is also very useful. “Negative” is actually a misnomer, and the misinterpretation is widespread. However, here we will target relationships that investigators claim

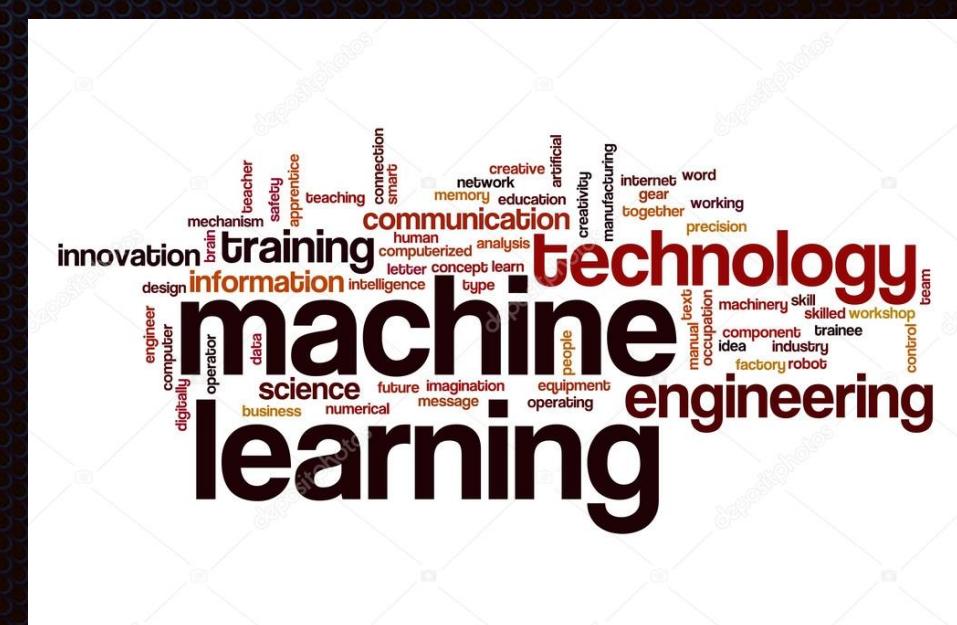
characteristic of the field and can vary a lot depending on whether the field targets highly likely relationships or searches for only one or a few true relationships among thousands and millions of hypotheses that may be postulated. Let us also consider, for computational simplicity, circumscribed fields where either there is only one true relationship (among many that can be hypothesized) or the power is similar to find any of the several existing true relationships. The pre-study probability of a relationship being true is  $R/(R + 1)$ . The probability of a study finding a true relationship reflects the power  $1 - \beta$  (one minus the Type II error rate). The probability of claiming a relationship when none truly exists reflects the Type I error rate,  $\alpha$ . Assuming that  $c$  relationships are being probed in the field, the expected values of the  $2 \times 2$  table are given in Table 1. After a research finding has been claimed based on achieving formal statistical significance, the post-study probability that it is true is the positive predictive value, PPV. The PPV is also the complementary probability of what Wacholder et al. have called the false positive report probability [10]. According to the  $2 \times 2$  table, one gets  $PPV = (1 - \beta)R/(R + \beta R + \alpha)$ . A research finding is thus

Starkey, Leannick, & P. (2005) Why most published

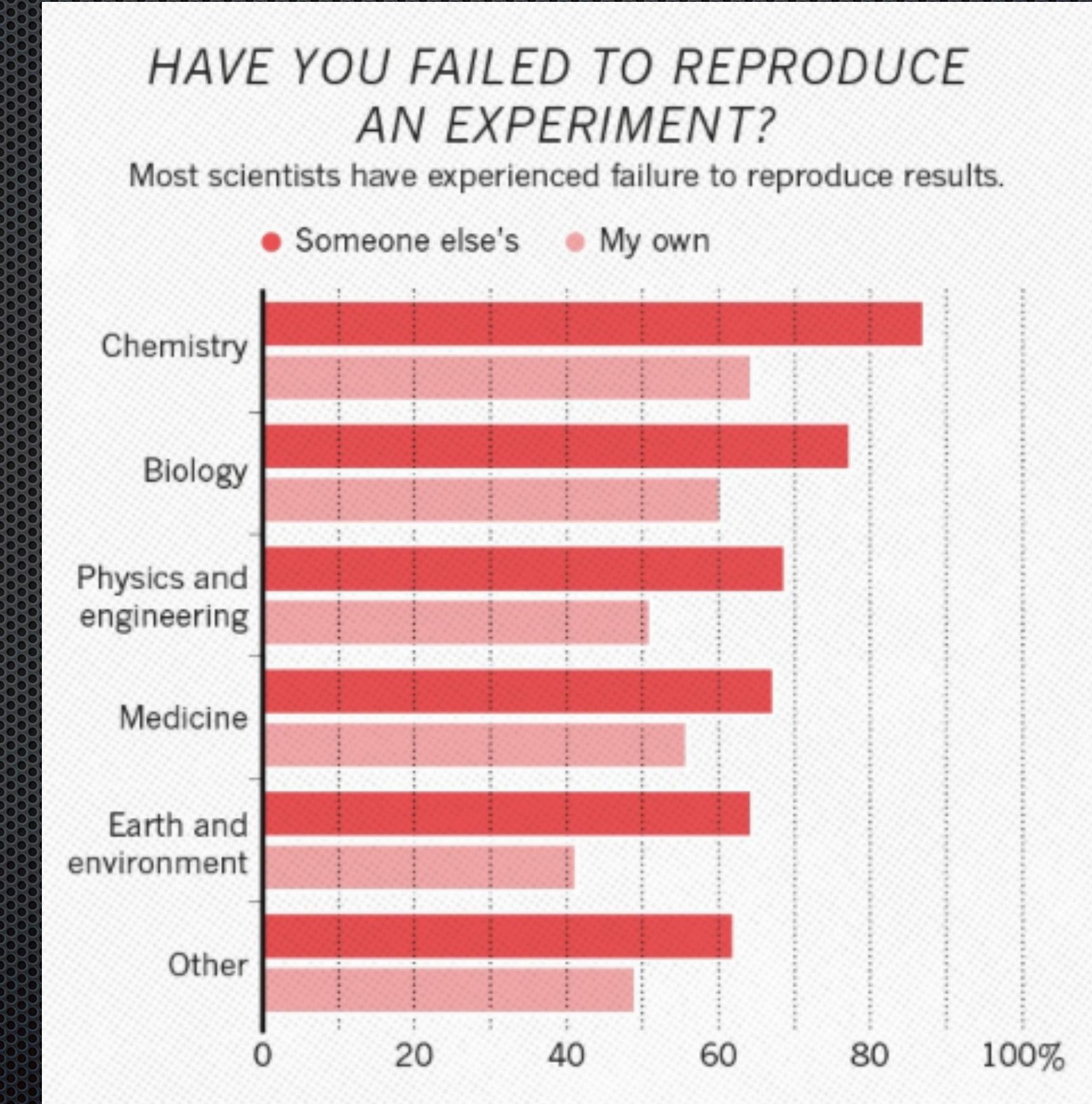
# Pharma report:

89% of research findings  
are not reproducible

This has huge implications  
for biotech and pharma  
development, right?!!



# IS THERE A REPRODUCIBILITY CRISIS?



# Issues with Big Data

- Variability – Inconsistencies in the data
  - Veracity – Poor data quality
  - Complexity – Data management challenges

These are really the more pressing challenges

This is behind most of the problems!



# Data is NOT information

- Data is just that, data
  - It only contributes to information when given a context (your question)
    - In that context, the data-derived information can be translated into action
      - Products
      - Services



# The question drives the learning

# Not the other way around



# My personal experience

- You've got data
- I need data and a question
- 80% of my time will be to help you figure out the question
- 80% of the remaining time will be spent getting your data into a usable state for your question

Need to find the right Big Data  
for your question

Need to see if your question is  
answerable given available data

# Biases and ethics

- Our learning is only as good as the training data that is used
  - Training data is often biased
    - Funding biases
    - Historical biases
    - Incompatibilities



# Biases and ethics

- Genomic and pathway databases ingest the scientific literature
  - The relative volume of literature in various fields is driven by funding
  - Databases are biased to cancer, heart disease and diabetes
  - Good luck if you have a rarer disease to treat



# Biases and ethics

- Historically, minorities had lower chances of getting loans, mortgages, etc.
  - This is reflected in the historical data
  - This data is used to train ML models
  - Now the historical racial bias is codified, since the models will make predictions based on patterns it saw in the training data
  - Similar issue with clinical trials
    - Most studies were done on middle aged white men
    - Do these data apply to women, children, other ethnicities?
  - Huge issue recognized in data ethics



# Biases and ethics

- aka, we like shiny new things
  - Historical bioinformatic data incompatible with currently produced data
  - We have lost the opportunity to use this large corpus of data
    - Billions of dollars in international funding
  - We can't tell how these data relate to current findings easily

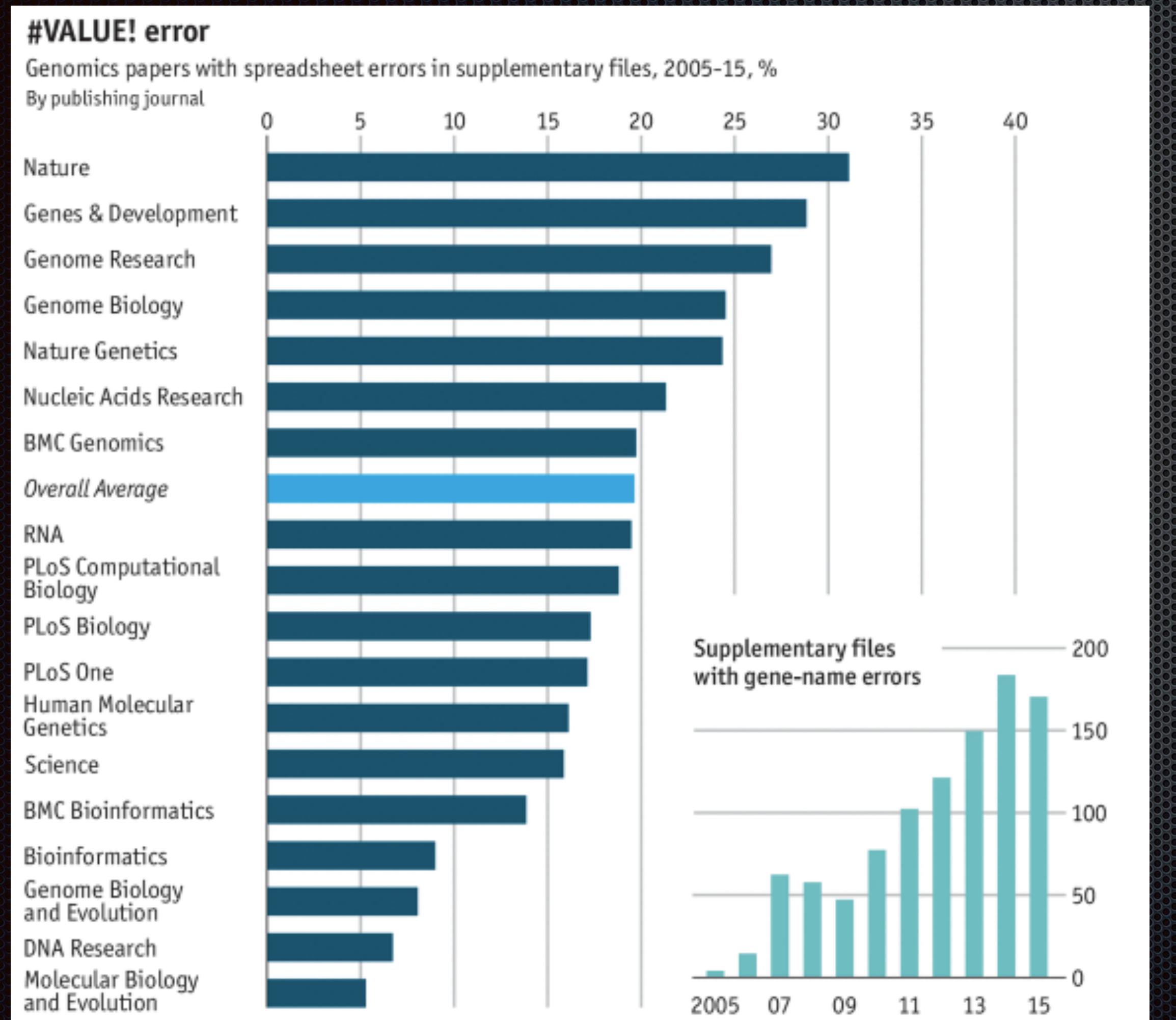


# Garbage in, garbage out

- Make sure the data you're leveraging for your business and your career is of the highest quality
  - Provenance
  - How it was collected
  - How it is stored



# Garbage in, garbage out



# Garbage in, garbage out

- Make sure you use the data
    - Appropriately
    - Ethically
  - Really think through and understand your data, your question, and how they relate
  - Don't just dump things into AI just because it's the cool thing



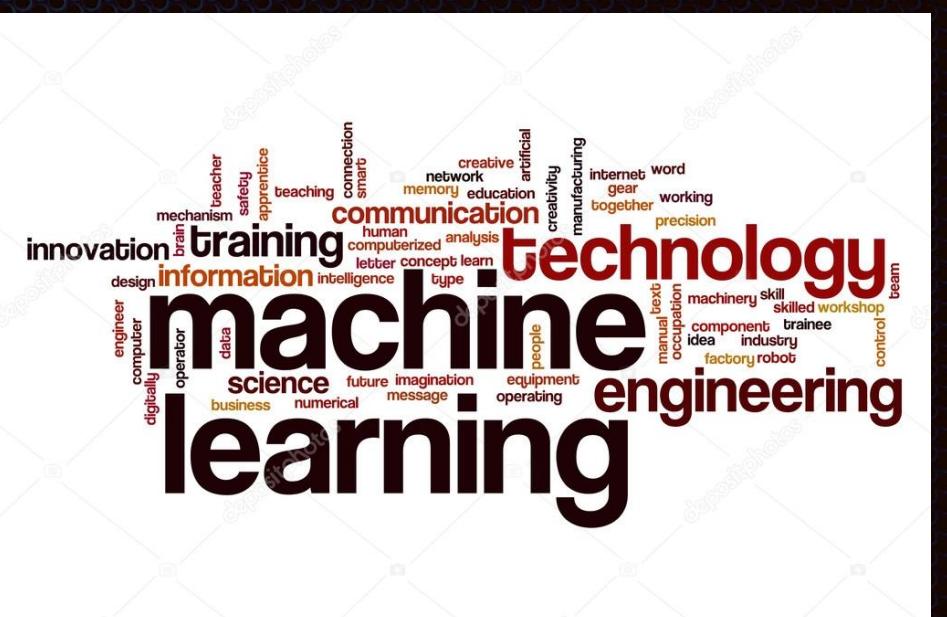
# Big Data + ML/AI ≠ Sorceror's Stone

- Big Data can mean Big Noise
  - No signal to be had
  - ML will give spurious, irreproducible answers
  - Dietary studies are notorious for this, as are psych studies



# Big Data + ML/AI ≠ Sorceror's Stone

- Finding and assessing rare events is VERY HARD
  - Suppose an event occurs in 1 out of 10,000 people
  - You fit a ML model to this huge database (think Medicare)
  - I use the amazing “Ostrich Method”: Just say no
    - I’ll have an accuracy of 99.99% in this scenario
    - Beat that, you fake intelligence.



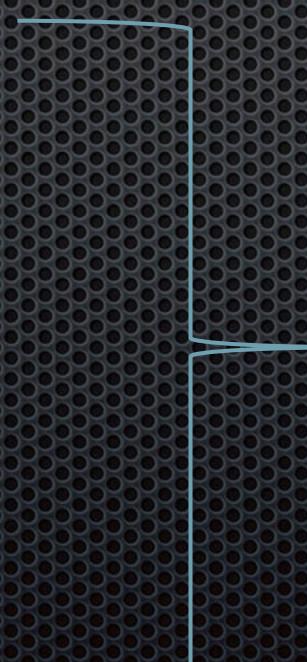
# The crystal ball

- Storage costs are a huge issue
  - Analytic methods need progress
    - Still leveraging methods from 20 or more years ago
    - A fast moving area



# The crystal ball

- Lots of unusable data
  - Lack of compatibility
  - Lack of documentation
- More efficiency and focus needed
  - Data generation
  - Data analysis
  - Data products



Symbiosis



“Begin with the end in mind”

- Find a problem that needs solving
  - Find the tools, data and methods that will help you figure it out
  - Create a diverse team (in skills, talent and background)
  - Work as a team to develop a product that solves it
  - Use evidence-based best practices to convince the market that you've solved it



# Contact info

- [abhijit@zansors.com](mailto:abhijit@zansors.com)
- [@webbedfeet / @Zansors](https://twitter.com/webbedfeet)
- <https://webbedfeet.netlify.com>



**Abhijit Dasgupta**

Director/Chief Data Scientist

ARAASTAT

Zansors

