



UNIVERSIDAD
DE GRANADA



Sistemas Inteligentes para la Gestión de la Empresa

E.T.S. de Ingenierías Informática y de Telecomunicación
Universidad de Granada

Juan Gómez Romero
jgomez@decsai.ugr.es

Departamento de Ciencias de la
Computación e Inteligencia Artificial
<http://decsai.ugr.es>

nature

LETTERS

Vol 457 | 9 February 2009 | doi:10.1038/nature07634

Detecting influenza epidemics using search engine query data

Jeremy Ginsberg¹, Matthew H. Mohebbi¹, Rajan S. Patel¹, Lynnette Brammer², Mark S. Smolinski¹ & Larry Brilliant¹

Seasonal influenza epidemics are a major public health concern, causing tens of millions of respiratory illnesses and 250,000 to 500,000 deaths worldwide each year¹. In addition to seasonal influenza, a new strain of influenza virus against which no previous immunity exists and that demonstrates human-to-human transmission could result in a pandemic with millions of fatalities². Early detection of disease activity, when followed by a rapid response, can reduce the impact of both seasonal and pandemic influenza^{3,4}. One way to improve early detection is to monitor health-seeking behaviour in the form of queries to online search engines, which are submitted by millions of users around the world each day. Here we present a method of analysing large numbers of Google search queries to track influenza-like illness in a population. Because the relative frequency of certain queries is highly correlated with the percentage of physician visits in which a patient presents with influenza-like symptoms, we can accurately estimate the current level of weekly influenza activity in each region of the United States, with a reporting lag of about one day. This approach may make it possible to use search queries to detect influenza epidemics in areas with a large population of web search users.

Traditional surveillance systems, including those used by the US Centers for Disease Control and Prevention (CDC) and the European Influenza Surveillance Scheme (EISS), rely on both virological and clinical data, including influenza-like illness (ILI) physician visits. The CDC publishes national and regional data from these surveillance systems on a weekly basis, typically with a 1–2-week reporting lag.

In an attempt to provide faster detection, innovative surveillance systems have been created to monitor indirect signals of influenza activity, such as call volume to telephone triage advice lines⁵ and over-the-counter drug sales⁶. About 90 million American adults are believed to search online for information about specific diseases or medical problems each year⁷, making web search queries a uniquely valuable source of information about health trends. Previous attempts at using online activity for influenza surveillance have counted search queries submitted to a Swedish medical website (A. Hulth, G. Rydevik and A. Linde, manuscript in preparation), visitors to certain pages on a US health website⁸, and user clicks on a search keyword advertisement in Canada⁹. A set of Yahoo! search queries containing the words ‘flu’ or ‘influenza’ were found to correlate with virological and mortality surveillance data over multiple years¹⁰.

Our proposed system builds on this earlier work by using an automated method of discovering influenza-related search queries. By processing hundreds of billions of individual searches from 5 years of Google web search logs, our system generates more comprehensive models for use in influenza surveillance, with regional and state-level estimates of ILI activity in the United States. Widespread global usage of online search engines may eventually enable models to be developed in international settings.

By aggregating historical logs of online web search queries submitted between 2003 and 2008, we computed a time series of weekly counts for 50 million of the most common search queries in the United States. Separate aggregate weekly counts were kept for every query in each state. No information about the identity of any user was retained. Each time series was normalized by dividing the count for each query in a particular week by the total number of online search queries submitted in that location during the week, resulting in a query fraction (Supplementary Fig. 1).

We sought to develop a simple model that estimates the probability that a random physician visit in a particular region is related to an ILI; this is equivalent to the percentage of ILI-related physician visits. A single explanatory variable was used: the probability that a random search query submitted from the same region is ILI-related, as determined by an automated method described below. We fit a linear model using the log-odds of an ILI physician visit and the log-odds of an ILI-related search query: $\text{logit}(I(t)) = \text{zlogit}(Q(t)) + z$, where $I(t)$ is the percentage of ILI physician visits, $Q(t)$ is the ILI-related query fraction at time t , z is the multiplicative coefficient, and z is the error term. $\text{logit}(p)$ is simply $\ln(p/(1-p))$.

Publily available historical data from the CDC's US Influenza Sentinel Provider Surveillance Network (<http://www.cdc.gov/flu/weekly>) was used to help build our models. For each of the nine surveillance regions of the United States, the CDC reported the average percentage of all outpatient visits to sentinel providers that were ILI-related on a weekly basis. No data were provided for weeks outside of the annual influenza season, and we excluded such dates from model fitting, although our model was used to generate unvalidated ILI estimates for these weeks.

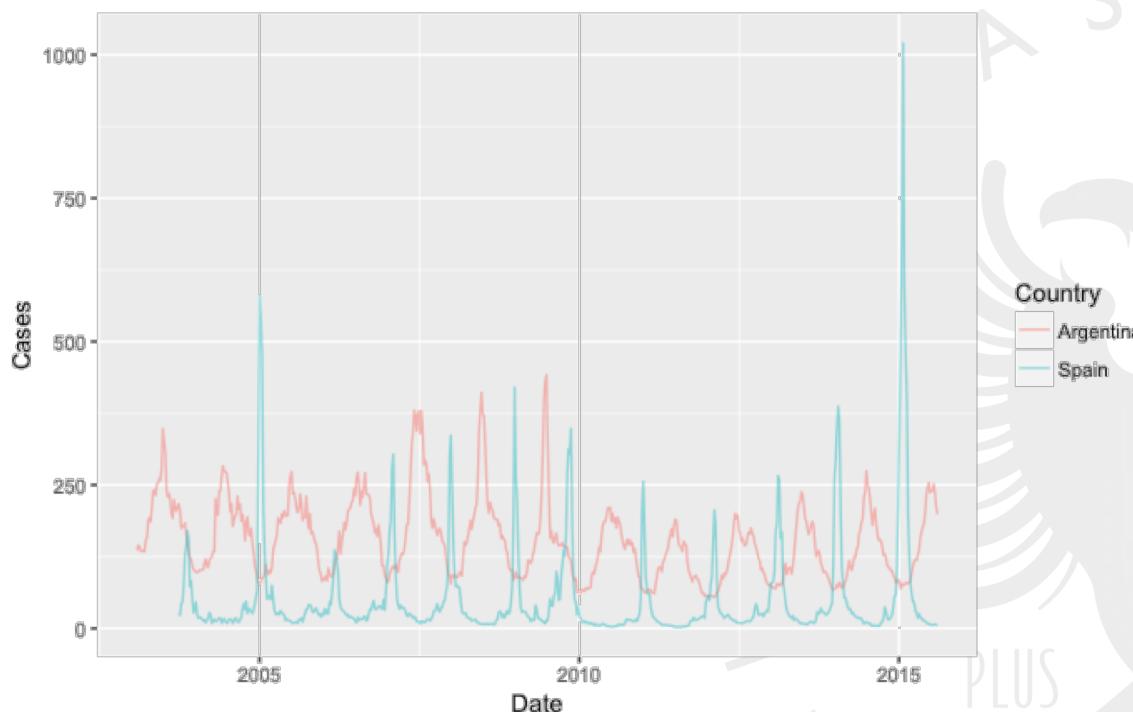
We designed an automated method of selecting ILI-related search queries, requiring no previous knowledge about influenza. We measured how effectively our model would fit the CDC ILI data in each region if we used only a single query as the explanatory variable, $Q(t)$. Each of the 50 million candidate queries in our database was separately tested in this manner, to identify the search queries which could most accurately model the CDC ILI visit percentage in each region. Our approach rewarded queries that showed regional variations similar to the regional variations in CDC ILI data: the chance that a random search query can fit the ILI percentage in all nine regions is considerably less than the chance that a random search query can fit a single location (Supplementary Fig. 2).

The automated query selection process produced a list of the highest scoring search queries, sorted by mean Z-transformed correlation across the nine regions. To decide which queries would be included in the ILI-related query fraction, $Q(t)$, we considered different sets of n top-scoring queries. We measured the performance of these models based on the sum of the queries in each set, and picked n such that we obtained the best fit against out-of-sample ILI data across the nine regions (Fig. 1).

¹Google Inc., 1600 Amphitheatre Parkway, Mountain View, California 94043, USA. ²Centers for Disease Control and Prevention, 1600 Clifton Road, NE, Atlanta, Georgia 30333, USA.

J. Ginsberg, M.H. Mohebbi,
R.S. Patel, L. Brammer, M.S.
Smolinski, L. Brilliant (2009)
Detecting influenza epidemics
using search engine query data.
Nature 457, pp. 1012–1015.

```
library(tidyr)
plotdata <-
  data %>%
  gather(key = Country, value = Cases, Argentina:Uruguay, na.rm = TRUE) %>%
  filter(Country == 'Spain' | Country == "Argentina")
library(ggplot2)
ggplot(plotdata, aes(x = Date, y = Cases, color = Country)) +
  geom_line(alpha = 0.5)
```



<https://www.google.org/flutrends/about/>

BUSINESS
INSIDER

TECH

FINANCE

POLITICS

STRATEGY

L

The price of bitcoin has a 91% correlation with Google searches for bitcoin



Jim Edwards



Sep. 19, 2017, 5:08 AM

42,207



FACEBOOK



LINKEDIN



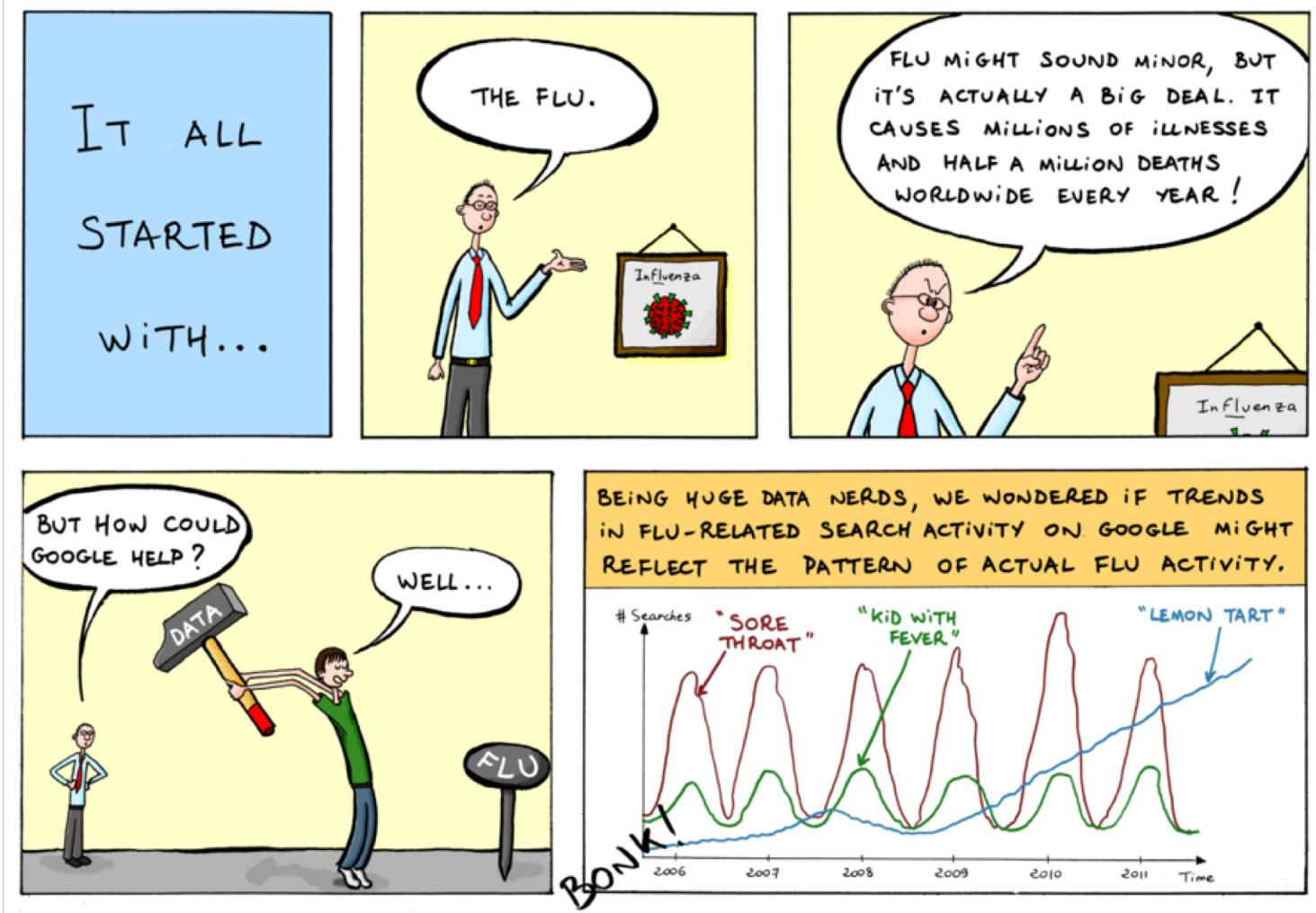
TWITTER



The current price of bitcoin has a 91% correlation with the volume of Google search requests for bitcoin-related terms, according to a study by SEMrush, a search engine marketing agency.

The study drew from a database of 120 million US keyword searches linked to the cryptocurrency. The overall search volume of bitcoin-related keywords is estimated to be 51.4 million requests over a period of a year. It showed that the price of bitcoin in US dollars rose and fell largely in tandem with the number of search requests for terms like "bitcoin," "bitcoin price," and "bitcoin value."

At one level, the study merely confirms the obvious: As bitcoin becomes more expensive, and thus more exciting, more people search online to find out how it is doing. On the other hand, it is nice to see statistics confirm your hunches. While the study looked at the correlation between searches and prices, it did not — sadly — say whether searches *predicted* or *trailed* the bitcoin/dollar exchange rate.



<https://www.google.com/trends/correlate/comic?p=1>

=> Creación de modelos de predicción

Modelo formal que permite obtener la salida de un sistema ante unos estímulos previamente no conocidos

=> Soporte a la toma de decisiones

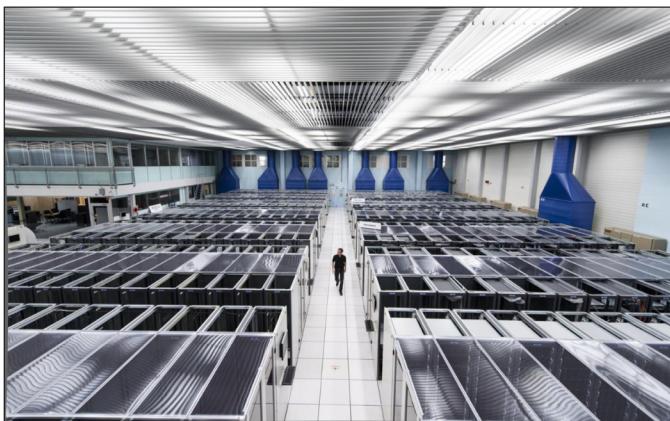
Plantear escenarios alternativos y valorar los resultados esperados

Big Data

Grandes conjuntos de datos permiten abordar problemas minimizando o incluso eliminando la necesidad de realizar muestreos aleatorios

Computing

Experiments at CERN generate colossal amounts of data. The Data Centre stores it, and sends it around the world for analysis



Approximately 600 million times per second, particles collide within the [Large Hadron Collider](#) (LHC). Each collision generates particles that often decay in complex ways into even more particles. Electronic circuits record the passage of each particle through a detector as a series of electronic signals, and send the data to the CERN Data Centre (DC) for digital reconstruction. The digitized summary is recorded as a "collision event". Physicists must sift through the 30 petabytes or so of data produced annually to determine if the collisions have thrown up any interesting physics.

CERN does not have the computing or financial resources to crunch all of the data on site, so in 2002 it turned to grid computing to share the burden with computer centres around the world. The [Worldwide LHC Computing Grid](#) (WLCG) – a distributed computing infrastructure [arranged in tiers](#) – gives a community of over 8000 physicists near real-time access to LHC data. The Grid builds on the technology of the World Wide Web, [which was invented at CERN in 1989](#).

<https://home.cern/about/computing>

YouTube

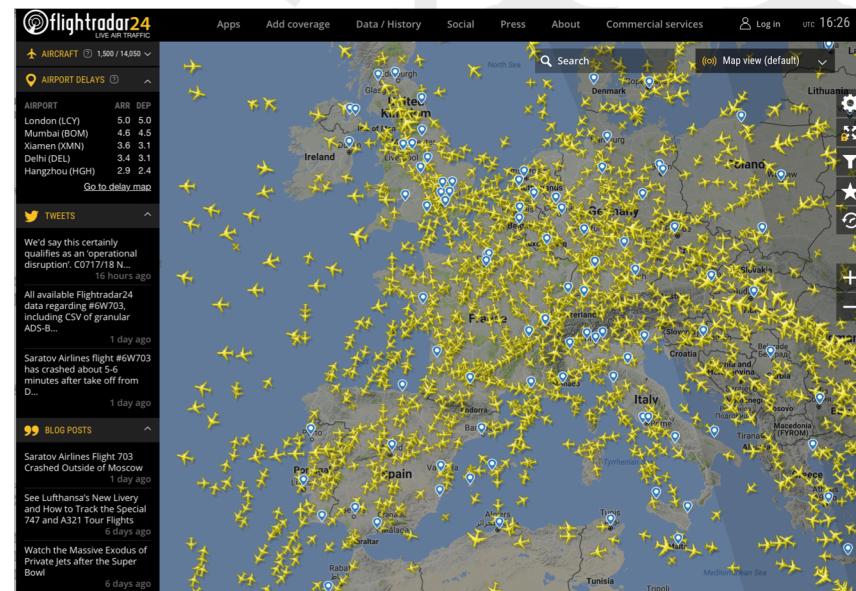
YouTube in numbers

One billion hours watched daily

This is the number of hours of video watched on YouTube every day, generating billions of views.

← →

Close



¿Quién genera datos?

Ciencia

Astronomía, genómica, física experimental, medioambiente

Ciencias Sociales y Humanidades

Libros escaneados, documentos históricos, datos sociales

Negocio y Comercio

Transacciones de mercados, logística, tráfico (aéreo, marítimo, carretera)

Entretenimiento y Ocio

Recursos multimedia en internet

Medicina

Bases de datos de pacientes, imágenes médicas

Industria

Sensores (redes eléctricas inteligentes, control de edificios, etc.)

Explosión de información

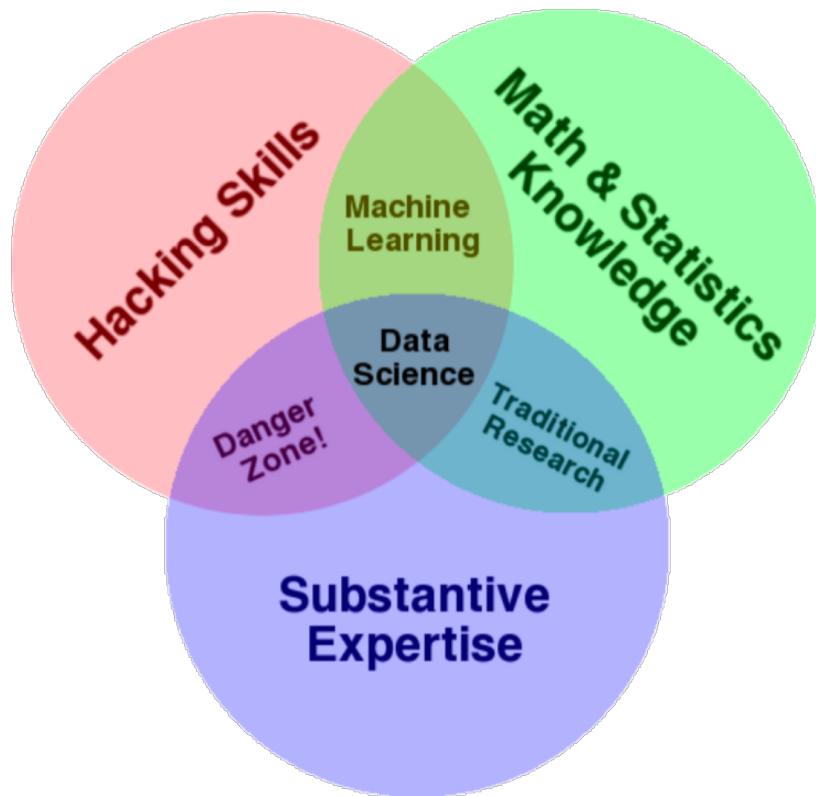
Nuevas herramientas para manejar datos

- Bajo precio de hardware
- Tecnologías de bases de datos
- Software para procesamiento distribuido

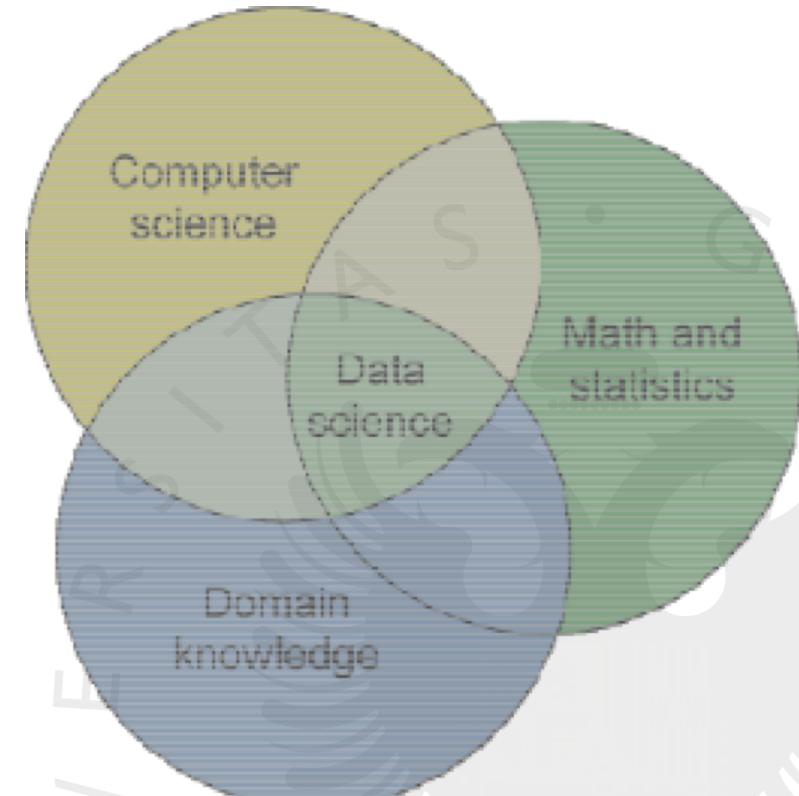
Podemos manejar grandes volúmenes de datos, pero resulta difícil extraer conocimiento

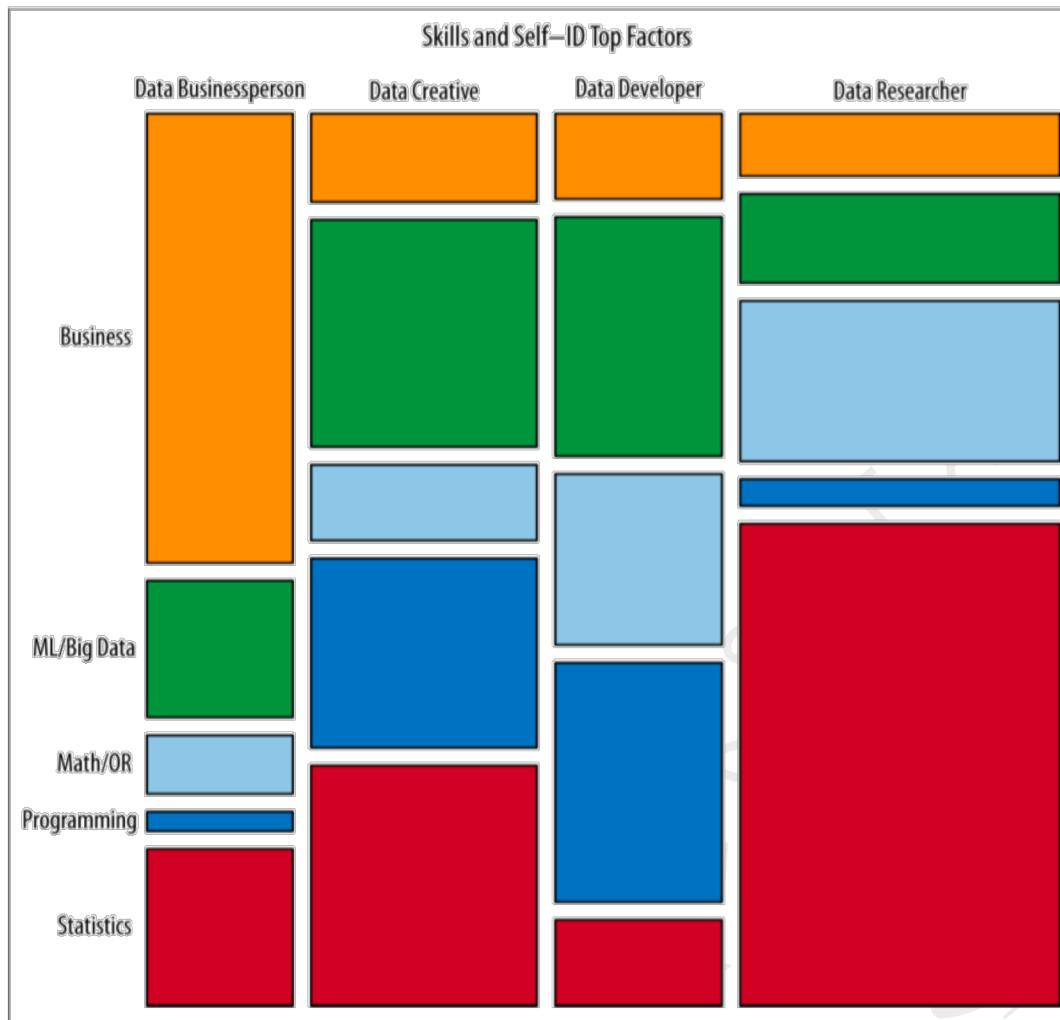
El progreso y la innovación ya no se ven obstaculizados por la capacidad de recopilar datos, sino por la capacidad de gestionar, analizar, sintetizar, visualizar, y descubrir el conocimiento de los datos recopilados de manera oportuna y en una forma escalable

Ciencia de Datos (*Data Science*) incorpora diferentes elementos y se basa en las técnicas y teorías de muchos campos, incluyendo las matemáticas, estadística, ingeniería de datos, reconocimiento de patrones y aprendizaje, computación avanzada, visualización, modelado de la incertidumbre, almacenamiento de datos y la informática de alto rendimiento **con el objetivo de extraer el significado de datos y la creación de productos de datos**



<http://drewconway.com>





M. Vaismann, H. Harris, S. Murphy (2013) Analyzing the analyzers: an introspective survey of Data Scientists and their work. O'Reilly

BUSINESS ANALYTICS VS DATA SCIENCE

In the age of big data, parsing unwieldy amounts of information can lead to world-changing innovation.

To understand this data, companies are hiring a variety of specialists, including business analysts and data scientists.

Who Are They?

BUSINESS ANALYSTS



Research and extract valuable information from structured and unstructured sources to explain historical, current, and future business performance; determine the best analytical models and approaches to present and explain solutions to business users.

DATA SCIENTISTS

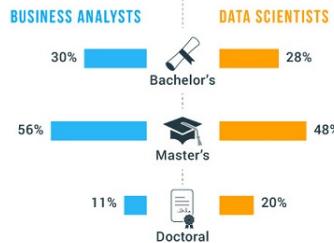


Design, develop, and deploy algorithms through statistical programming that support business decision-making tools; manage large amounts of data; create visualizations to aid in understanding.

What Education Do They Have?

The majority of business analysts come from a variety of backgrounds including business and humanities, while data scientists tend to come from computer science, mathematics, and technology backgrounds.

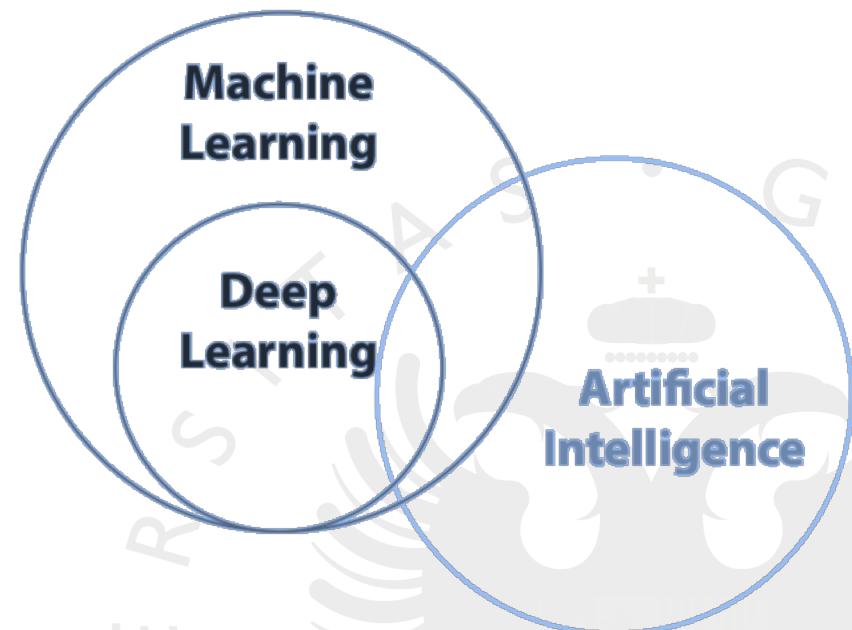
BUSINESS ANALYSTS



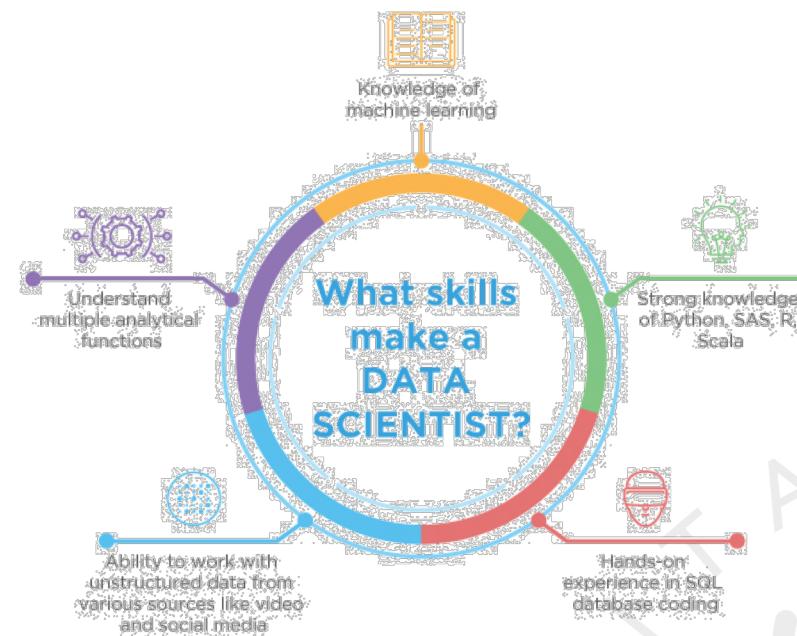
Sample Undergraduate Majors

- Business Administration
 - Information Technology
 - Finance
 - Political Science
 - Anthropology
 - Economics
 - History
 - Psychology
- Mathematics
 - Computer Science
 - Information Science
 - Computer Software Engineering
 - Technical Communications
 - Computer Information Systems
 - Statistics

[https://www.kdnuggets.com/2015/10/infographic-data-scientist-business-analyst-difference.html](https://www.kdnuggets.com/2015/10/infoographic-data-scientist-business-analyst-difference.html)



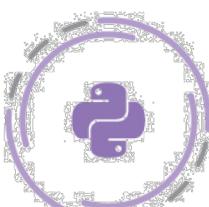
Andrew W. Trask (2018) Grokking Deep Learning. Manning.



What are the skills required to become a data analyst?



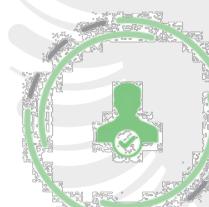
Knowledge of mathematical statistics



Fluent understanding of R and Python



Data wrangling



Understand PIG/HIVE

<https://www.simplilearn.com/data-science-vs-data-analytics-vs-machine-learning-article>



Cerveza => Pañales **60%**

Pañales => Cerveza **100%**

https://www.theregister.co.uk/2006/08/15/beer_diapers/

Tarjetas de crédito

*"A division called MasterCard Advisors aggregates and analyzes 65 billion transactions from 1.5 billion cardholders in 210 countries in order to divine business and consumer trends. Then it sells that information to others. It discovered, among other things, that if people fill up their gas tanks at around four o'clock in the afternoon, they're quite likely to spend between \$35 and \$50 in the next hour at a grocery store or restaurant. **A marketer might use that insight to print out coupons for a nearby supermarket on the back of gas-station receipts around that time of day.**"*

Cesta de la compra

*"The analytics team reviewed the shopping histories of women who signed up for its baby gift-registry. They noticed that these women bought lots of unscented lotion at around the third month of pregnancy, and that a few weeks later they tended to purchase supplements like magnesium, calcium, and zinc. **The team ultimately uncovered around two dozen products that, used as proxies, enabled the company to calculate a "pregnancy prediction" score for every customer who paid with a credit card or used a loyalty card or mailed coupons.** The correlations even let the retailer estimate the due date within a narrow range, so it could send relevant coupons for each stage of the pregnancy."*

V. Mayer-Schonberger, K. Cukier (2013) Big Data, A revolution that will transform how we live, work and think. Harcourt Publishing.

ARTICLE

doi:10.1038/nature16961

Mastering the game of Go with deep neural networks and tree search

David Silver^{1*}, Aja Huang^{1*}, Chris J. Maddison¹, Arthur Guez¹, Laurent Sifre¹, George van den Driessche¹, Julian Schrittwieser¹, Ioannis Antonoglou¹, Veda Pannenberghelvam¹, Marc Lanctot¹, Sander Dieleman¹, Dominik Grewe¹, John Nham¹, Nal Kalchbrenner¹, Ilya Sutskever¹, Timothy Lillicrap¹, Madeleine Leach¹, Koray Kavukcuoglu¹, Thore Graepel¹ & Demis Hassabis¹

The game of Go has long been viewed as the most challenging of classic games for artificial intelligence owing to its enormous search space and the difficulty of evaluating board positions and moves. Here we introduce a new approach to computer Go that uses ‘value networks’ to evaluate board positions and ‘policy networks’ to select moves. These deep neural networks are trained by a novel combination of supervised learning from human expert games, and reinforcement learning from games of self-play. Without any lookahead search, the neural networks play Go at the level of state-of-the-art Monte Carlo tree search programs that simulate thousands of random games of self-play. We also introduce a new search algorithm that combines Monte Carlo simulation with value and policy networks. Using this search algorithm, our program AlphaGo achieved a 99.8% winning rate against other Go programs, and defeated the human European Go champion by 5 games to 0. This is the first time that a computer program has defeated a human professional player in the full-sized game of Go, afeat previously thought to be at least a decade away.

All games of perfect information have an optimal value function, $v^*(s)$, which determines the outcome of the game, from every board position or state s , under perfect play by all players. These games may be solved by recursively computing the optimal value function in a search tree containing approximately b^d possible sequences of moves, where b is the game’s breadth (number of legal moves per position) and d is its depth (game length). In large games, such as chess ($b=35$, $d \approx 80$)¹ and especially Go ($b \approx 250$, $d \approx 150$)², exhaustive search is infeasible^{3,4}, but the effective search space can be reduced by two general principles. First, the depth of the search may be reduced by position evaluation: truncating the search tree at state s and replacing the subtree below s by an approximate value function $v(s) \approx v^*(s)$ that predicts the outcome from state s . This approach has led to superhuman performance in chess⁵, checkers⁶ and othello⁷, but it was believed to be intractable in Go due to the complexity of the game⁸. Second, the breadth of the search may be reduced by sampling actions from a policy $p(a|s)$ that is a probability distribution over possible moves a in position s . For example, Monte Carlo rollouts⁹ search to maximum depth without branching at all, by sampling long sequences of actions for both players from a policy p . Averaging over such rollouts can provide an effective position evaluation, achieving superhuman performance in backgammon⁸ and Scrabble⁹, and weak amateur level play in Go¹⁰.

Monte Carlo tree search (MCTS)^{11,12} uses Monte Carlo rollouts to estimate the value of each state in a search tree. As more simulations are executed, the search tree grows larger and the relevant values become more accurate. The policy used to select actions during search is also improved over time, by selecting children with higher values. Asymptotically, this policy converges to optimal play, and the evaluations converge to the optimal value function¹². The strongest current Go programs are based on MCTS, induced by policies that are trained to predict human expert moves¹³. These policies are used to narrow the search to a beam of high-probability actions, and to sample actions during rollouts. This approach has achieved strong amateur play^{13–15}. However, prior work has been limited to shallow

policies^{13–15} or value functions¹⁶ based on a linear combination of input features.

Recently, deep convolutional neural networks have achieved unprecedented performance in visual domains: for example, image classification¹⁷, face recognition¹⁸, and playing Atari games¹⁹. They use many layers of neurons, each arranged in overlapping tiles, to construct increasingly abstract, localized representations of an image²⁰. We employ a similar architecture for the game of Go. We pass in the board position as a 19×19 image and use convolutional layers to construct a representation of the position. We use these neural networks to reduce the effective depth and breadth of the search tree: evaluating positions using a value network, and sampling actions using a policy network.

We train the neural networks using a pipeline consisting of several stages of machine learning (Fig. 1). We begin by training a supervised learning (SL) policy network p_s , directly from expert human moves. This provides fast, efficient learning updates with immediate feedback and high-quality gradients. Similar to prior work^{13,15}, we also train a fast policy p_f , that can rapidly sample actions during rollouts. Next, we train a reinforcement learning (RL) policy network p_r , that improves the SL policy network by optimizing the final outcome of games of self-play. This adjusts the policy towards the correct goal of winning games, rather than maximizing predictive accuracy. Finally, we train a value network v that predicts the winner of games played by the RL policy network against itself. Our program AlphaGo efficiently combines the policy and value networks with MCTS.

Supervised learning of policy networks

For the first stage of the training pipeline, we build on prior work on predicting expert moves in the game of Go using supervised learning^{13,21–23}. The SL policy network $p_s(a|s)$ alternates between convolutional layers with weights σ , and rectifier nonlinearities. A final softmax layer outputs a probability distribution over all legal moves a . The input s to the policy network is a simple representation of the board state (see Extended Data Table 2). The policy network is trained on randomly

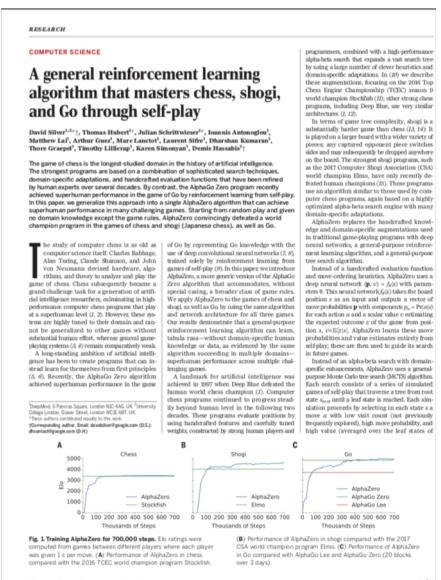
¹Google DeepMind, 5 New Street Square, London EC4A 3TW, UK. ²Google, 1600 Amphitheatre Parkway, Mountain View, California 94043, USA.

*These authors contributed equally to this work.

D. Silver et al. (2016)
Mastering the game of Go with
deep neural networks and tree
search. *Nature* 529, pp. 484–
489.



<http://www.alphagomovie.com/>



arXiv:1810.04805v1 [cs.CL] 11 Oct 2018

ARTICLE OPEN

Scalable and accurate deep learning with electronic health records

Alvin Rajkumar,^{1,2} Eyal Ofen,¹ Kai Chen,¹ Andrew M. Dahl,¹ Nisan Hsiao,¹ Michaela Hanif,¹ Peter J. Liu,¹ Xiangbin Liu,¹ Jake Marcus,¹ Mimi Sun,¹ Patrick Sweeney,¹ Helen Yu,¹ Kun Zhang,¹ Yi Zhang,¹ Gerardo Flores,¹ Gavin E. Duggan,¹ James Viik,¹ Quoc Le,¹ Hart K. Lee,¹ Daniel L. Rabinowitz,¹ Austin Tawarayama,¹ Da-Heng Jenny Jiang,¹ Michael J. Kupersmith,¹ Daniel S. Markin,¹ Samir S. Umeshwar,¹ Katherine Choi,¹ Michael Pearson,¹ Stevenwas Machekon,¹ Rajeev M. Shah,¹ Andi J. Butte,¹ Michael D. Howell,¹ Clare Cul,¹ Greg S. Corrado,¹ and Jeffrey Dean¹

In terms of game tree complexity, chess is a challenging domain for reinforcement learning (RL). In AlphaZero, we have developed a general RL algorithm to play a large board with a wider variety of pieces than chess. We have also developed a general RL architecture that can be applied to many other domains and may subsequently be dropped anywhere as the need arises. We present our work at the 2017 Computer Shogi Association (CSA) world championship. Elmo, have only recently demonstrated superhuman performance in the game of Go by reinforcement learning from self-play. AlphaZero, however, has achieved superhuman performance in many challenging games. Starting from random play and given no domain knowledge except the game rules, AlphaZero convincingly defeated a world champion in chess, shogi, and Go.

The study of computer chess is old as computer science itself. Charles Babbage, Alan Turing, Claude Shannon, and John von Neumann all contributed to the field. Games like chess and shogi have been studied for centuries, and the search for a general algorithm to play them has been a grand challenge for a generation of artificial intelligence researchers. AlphaZero, a general-purpose computer chess program that plays chess, shogi, and Go, has achieved superhuman performance in all three games. Our results demonstrate that a general-purpose reinforcement learning algorithm can learn to play multiple games without being explicitly told what knowledge or data, as evidenced by the same algorithm achieving superhuman performance across multiple chess variants.

A landmark for artificial intelligence was the 2016 DeepMind Go program's victory over the world chess champion.¹ Computer Go is a well-studied problem, and the search for a general algorithm to play beyond human level has been following two paths. These programs evaluate position by selecting in each state a move with a small cost (not necessarily minimum), and then select the move with the highest value (averaged over the leaf states of

the search tree). AlphaZero follows these two paths simultaneously, but uses a much more limited set of moves to guide its search.

Instead of an alphabeta search with domain-specific knowledge, AlphaZero uses a general search function that can search any game variant. It uses Monte Carlo tree search (MCTS) algorithms to search for the best move of each game that traverses from root node to leaf node. AlphaZero uses a general search procedure by selecting in each state a move with a small cost (not necessarily minimum), and then select the move with the highest value (averaged over the leaf states of the search tree).

In addition to the development of available data, scaling the development of predictive models is difficult because, for medical applications, the data required for training a deep learning model requires the creation of a custom dataset with specific samples and labels. This is time-consuming and costly, leading to poor clinical outcomes.^{1,2,3} Incorporating the scalability of deep learning into the workflow of the clinic is a challenge that is the number of potential predictor variables and the volume of data that can be used to overcome these shortcomings, but is unlikely to meet for most predictive modeling techniques.

Recent work in deep learning and artificial neural networks may allow us to address some of these challenges and overcome the limitations of current approaches. Preferred machine learning approaches in machine prediction are deep learning models, which have more recently been shown to provide many benefits not only for safety and quality but also in reducing healthcare costs.

In this work, we propose a novel approach to the problem of overcoming these shortcomings, but is unlikely to meet for most predictive modeling techniques.

Recent work in deep learning and artificial neural networks may allow us to address some of these challenges and overcome the limitations of current approaches. Preferred machine learning approaches in machine prediction are deep learning models, which have more recently been shown to provide many benefits not only for safety and quality but also in reducing healthcare costs.

Previous work in dermatology called *alpha*-diagnostic^{4,5 has highlighted the diagnostic capability of medical practitioners to identify difficult data and on standardized tasks such as skin lesion classification. However, this work has more recently been used in natural language processing, sequence prediction, and mixed modal data fusion. In this work, we propose a novel approach to overcome these shortcomings, but is unlikely to meet for most predictive modeling techniques.}

To take advantage of fine-grained information contained within the electronic health record (EHR) data, we propose a novel approach to partition diseases into fine-grained training classes (for example, skin cancer types) and then use a deep learning model to make inferences. The CNN outputs a probability distribution over these fine classes. To recover the probability for coarse-level classes of interest (such as melanoma), we use a softmax layer. This work is critical, as the estimated 5-year survival risk for melanoma drops from over 99% detected in its earliest stages to about 11% detected in its latest stages.⁶ This work is critical to the success of different clinical curated, open-access online repositories, as well as from the medical community. This work is critical to the success of a subset of the full taxonomy, which has been organized clinically and visually by medical experts. We split our data into 127,463 training cases and 12,463 testing cases.

To take advantage of fine-grained information contained within the electronic health record (EHR) data, we propose a novel approach to partition diseases into fine-grained training classes (for example, skin cancer types) and then use a deep learning model to make inferences. The CNN outputs a probability distribution over these fine classes. To recover the probability for coarse-level classes of interest (such as melanoma), we use a softmax layer. This work is critical, as the estimated 5-year survival risk for melanoma drops from over 99% detected in its earliest stages to about 11% detected in its latest stages.⁶ This work is critical to the success of different clinical curated, open-access online repositories, as well as from the medical community. This work is critical to the success of a subset of the full taxonomy, which has been organized clinically and visually by medical experts. We split our data into 127,463 training cases and 12,463 testing cases.

Stanford University, Stanford, California, USA; Dermatology Service, Veterans Affairs Palo Alto Health Care System, Palo Alto, California, USA; Center Laboratory for Stem Cell Biology, Department of Pathology, Stanford University, Stanford, California, USA; Department of Dermatology, Stanford University, Stanford, California, USA; Department of Oncologic Surgery, Stanford University, Stanford, California, USA

*These authors contributed equally to this work.

Received: 26 January 2018; Revised: 14 March 2018; Accepted: 26 March 2018

Published online: 08 May 2018

Published in partnership with the Scopus Translational Science Institute

npj partner journals

© 2018 The Author(s). *npj Digital Medicine* published by Springer Nature Limited on behalf of Scopus Translational Science Institute.

This article is an open access publication

2 FEBRUARY 2019 | VOL 542 | NATURE | 113

Automatic Handgun Detection Alarm in Videos Using Deep Learning

Roberto Olmos¹, Siham Tabik¹, and Francisco Herrera^{1,2}

¹Soft Computing and Intelligent Information Systems research group, University of Granada, Spain; ²Department of Computer Science and Artificial Intelligence, University of Granada, Spain. emails: siham.olmos@ugr.es, francesco.herrera@decsai.ugr.es

February 20, 2017

Abstract

Current surveillance and control systems still require human supervision and intervention. This work presents a novel automatic handgun detection system in video applications, monitoring surveillance control systems. We implement this detection process by solving the problem of minimizing false positives and false negatives by building the key training dataset guided by the results of a deep Convolutional Neural Network (CNN) trained by the authors. The proposed system is based under two approaches, the sliding window approach and region proposal approach. The most promising results are obtained by Faster R-CNN based model trained on the new dataset. The proposed system has a high performance in low quality youtube videos and provides satisfactory results as automatic alarm system. Among 30 scenes, it successfully activates the alarm after 12.2 seconds and activates the alarm after 12.2 seconds. This work also defines a new metric, Alarm Activation per Interval (AApi), to assess the performance of a detection model as an automatic detection system.

Index terms — Classification, Detection, Deep learning, Convolutional Neural Networks (CNNs), Faster R-CNN, VGG-16, Alarm Activation per Interval

1 Introduction

Language model pre-training has shown to be effective for improving many natural language processing tasks (Dai et al., 2015; Peters et al., 2018; Radford et al., 2018; Vaswani et al., 2018; Ruder, 2018). These tasks include sentence-level tasks such as natural language inference (Bowman et al., 2015; Williams et al., 2018) and question answering tasks such as SQuAD question answering (Rajpurkar et al., 2016; Williams et al., 2018) and reading comprehension tasks such as CoQA (Rajpurkar et al., 2016). In this paper, we propose a novel pre-training objective, the “masked language

modeling”, conditioned with a high performancealphabet that regards a word as unknown by using a large number of clever heuristics and thus avoiding the need to manually tune these heuristics. Focusing on the 2018 Top Chess Engine Championship (TCEC) round 9 tournament, we present the first results of this program, including Deep Blue, use very similar architecture.

In terms of game tree complexity, chess is a challenging domain for reinforcement learning (RL). In AlphaZero, we have developed a general RL algorithm to play a large board with a wider variety of pieces than chess, shogi, and Go. The search space and the number of possible moves are exponentially larger than in chess and Go, and may subsequently be dropped anywhere as the need arises. We present our work at the 2017 Computer Shogi Association (CSA) world championship. Elmo, have only recently demonstrated superhuman performance in the game of Go by reinforcement learning from self-play. AlphaZero, however, has achieved superhuman performance in many challenging games. Starting from random play and given no domain knowledge except the game rules, AlphaZero convincingly defeated a world champion in chess, shogi, and Go.

The study of computer chess is old as computer science itself. Charles Babbage, Alan Turing, Claude Shannon, and John von Neumann all contributed to the field. Games like chess and shogi have been studied for centuries, and the search for a general algorithm to play them has been a grand challenge for a generation of artificial intelligence researchers. AlphaZero, a general-purpose computer chess program that plays chess, shogi, and Go, has achieved superhuman performance in all three games. Our results demonstrate that a general-purpose reinforcement learning algorithm can learn to play multiple games without being explicitly told what knowledge or data, as evidenced by the same algorithm achieving superhuman performance across multiple chess variants.

A landmark for artificial intelligence was the 2016 DeepMind Go program's victory over the world chess champion.¹ Computer Go is a well-studied problem, and the search for a general algorithm to play beyond human level has been following two paths. These programs evaluate position by selecting in each state a move with a small cost (not necessarily minimum), and then select the move with the highest value (averaged over the leaf states of the search tree).

In addition to the development of available data, scaling the development of predictive models is difficult because, for medical applications, the data required for training a deep learning model requires the creation of a custom dataset with specific samples and labels. This is time-consuming and costly, leading to poor clinical outcomes.^{1,2,3} Incorporating the scalability of deep learning into the workflow of the clinic is a challenge that is the number of potential predictor variables and the volume of data that can be used to overcome these shortcomings, but is unlikely to meet for most predictive modeling techniques.

Previous work in dermatology called *alpha*-diagnostic^{4,5 has highlighted the diagnostic capability of medical practitioners to identify difficult data and on standardized tasks such as skin lesion classification. However, this work has more recently been used in natural language processing, sequence prediction, and mixed modal data fusion. In this work, we propose a novel approach to overcome these shortcomings, but is unlikely to meet for most predictive modeling techniques.}

To take advantage of fine-grained information contained within the electronic health record (EHR) data, we propose a novel approach to partition diseases into fine-grained training classes (for example, skin cancer types) and then use a deep learning model to make inferences. The CNN outputs a probability distribution over these fine classes. To recover the probability for coarse-level classes of interest (such as melanoma), we use a softmax layer. This work is critical, as the estimated 5-year survival risk for melanoma drops from over 99% detected in its earliest stages to about 11% detected in its latest stages.⁶ This work is critical to the success of different clinical curated, open-access online repositories, as well as from the medical community. This work is critical to the success of a subset of the full taxonomy, which has been organized clinically and visually by medical experts. We split our data into 127,463 training cases and 12,463 testing cases.

To take advantage of fine-grained information contained within the electronic health record (EHR) data, we propose a novel approach to partition diseases into fine-grained training classes (for example, skin cancer types) and then use a deep learning model to make inferences. The CNN outputs a probability distribution over these fine classes. To recover the probability for coarse-level classes of interest (such as melanoma), we use a softmax layer. This work is critical, as the estimated 5-year survival risk for melanoma drops from over 99% detected in its earliest stages to about 11% detected in its latest stages.⁶ This work is critical to the success of different clinical curated, open-access online repositories, as well as from the medical community. This work is critical to the success of a subset of the full taxonomy, which has been organized clinically and visually by medical experts. We split our data into 127,463 training cases and 12,463 testing cases.

Stanford University, Stanford, California, USA; Dermatology Service, Veterans Affairs Palo Alto Health Care System, Palo Alto, California, USA; Center Laboratory for Stem Cell Biology, Department of Pathology, Stanford University, Stanford, California, USA; Department of Dermatology, Stanford University, Stanford, California, USA; Department of Oncologic Surgery, Stanford University, Stanford, California, USA

*These authors contributed equally to this work.

Received: 26 January 2018; Revised: 14 March 2018; Accepted: 26 March 2018

Published online: 08 May 2018

Published in partnership with the Scopus Translational Science Institute

npj partner journals

© 2018 The Author(s). *npj Digital Medicine* published by Springer Nature Limited on behalf of Scopus Translational Science Institute.

This article is an open access publication

2 FEBRUARY 2019 | VOL 542 | NATURE | 113

LETTER

doi:10.1038/nature21096

Dermatologist-level classification of skin cancer with deep neural networks

Andre Esteva¹, Brett Kuprel², Robert A. Novak³, Justin M. Ko⁴, Susan M. Swetter⁵, Helen M. Blau⁶ & Sebastian Thrun¹

Skin cancer, the most common human malignancy^{1–3}, is primarily diagnosed visually, beginning with initial clinical screening and followed potentially by dermoscopic analysis, a biopsy and histopathological examination. Skin cancer classification is a challenging task owing to the great variability in both photographic and dermoscopic images. Previous work has made classification robust to photographic variability. Many previous techniques require extensive preprocessing, lesion segmentation and highly variable tasks across many fine-grained object categories^{4–6}. However, these methods are slow and computationally expensive. In contrast, our system requires no hand-coded features; it is trained on raw images directly, using only pixels and simple geometric features to end from images directly, using only pixels and simple geometric features to end from images directly. We demonstrate that deep learning with electronic health records (EHRs) is anticipated to drive personalized medicine and improve healthcare quality. Constructing predictive statistical models typically requires extraction of causal predictor variables from normalized EHR data. In this work, we demonstrate that deep learning representations are capable of accurately predicting multiple clinical outcomes from EHRs. We use the EHRs of 216,221 adult patients hospitalized for at least 24 h. In the sequential format we use, the primary outcome is the total of 86,447 melanomas across sites (9.5% of all hospitalizations). We use a deep learning model to predict the probability of a patient having melanoma versus benign lesion. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. The study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

Predictive modeling with electronic health record (EHR) data is a challenging task, especially when dealing with visual data. We overcome this challenge by using a data-driven approach to predict skin cancer using only pixels and simple geometric features to end from images directly, using only pixels and simple geometric features to end from images directly. We demonstrate that deep learning with electronic health records (EHRs) is anticipated to drive personalized medicine and improve healthcare quality. Constructing predictive statistical models typically requires extraction of causal predictor variables from normalized EHR data. In this work, we demonstrate that deep learning representations are capable of accurately predicting multiple clinical outcomes from EHRs. We use the EHRs of 216,221 adult patients hospitalized for at least 24 h. In the sequential format we use, the primary outcome is the total of 86,447 melanomas across sites (9.5% of all hospitalizations). We use a deep learning model to predict the probability of a patient having melanoma versus benign lesion. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. The study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

Previous work^{1–6} has shown that deep learning can be used to predict skin cancer with accuracy comparable to dermatologists. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. The study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

Recent work^{7–10} has shown that deep learning can be used to predict skin cancer with accuracy comparable to dermatologists. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. The study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

Current work^{11–13} has shown that deep learning can be used to predict skin cancer with accuracy comparable to dermatologists. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. The study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

Most recently, work^{14–16} has shown that deep learning can be used to predict skin cancer with accuracy comparable to dermatologists. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. The study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

Electronic health records (EHRs) contain a wealth of information that can be used to predict skin cancer with accuracy comparable to dermatologists. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. The study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

Dermatologists are the most common human malignancy^{1–3}, as it is primarily diagnosed visually, beginning with initial clinical screening and followed potentially by dermoscopic analysis, a biopsy and histopathological examination. Skin cancer classification is a challenging task owing to the great variability in both photographic and dermoscopic images. Previous work has made classification robust to photographic variability. Many previous techniques require extensive preprocessing, lesion segmentation and highly variable tasks across many fine-grained object categories^{4–6}. However, these methods are slow and computationally expensive. In contrast, our system requires no hand-coded features; it is trained on raw images directly, using only pixels and simple geometric features to end from images directly. We demonstrate that deep learning with electronic health records (EHRs) is anticipated to drive personalized medicine and improve healthcare quality. Constructing predictive statistical models typically requires extraction of causal predictor variables from normalized EHR data. In this work, we demonstrate that deep learning representations are capable of accurately predicting multiple clinical outcomes from EHRs. We use the EHRs of 216,221 adult patients hospitalized for at least 24 h. In the sequential format we use, the primary outcome is the total of 86,447 melanomas across sites (9.5% of all hospitalizations). We use a deep learning model to predict the probability of a patient having melanoma versus benign lesion. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. The study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

Previous work^{1–6} has shown that deep learning can be used to predict skin cancer with accuracy comparable to dermatologists. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. The study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

Recent work^{7–10} has shown that deep learning can be used to predict skin cancer with accuracy comparable to dermatologists. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. The study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

Current work^{11–13} has shown that deep learning can be used to predict skin cancer with accuracy comparable to dermatologists. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. The study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

Most recently, work^{14–16} has shown that deep learning can be used to predict skin cancer with accuracy comparable to dermatologists. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. The study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

Electronic health records (EHRs) contain a wealth of information that can be used to predict skin cancer with accuracy comparable to dermatologists. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. The study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

Dermatologists are the most common human malignancy^{1–3}, as it is primarily diagnosed visually, beginning with initial clinical screening and followed potentially by dermoscopic analysis, a biopsy and histopathological examination. Skin cancer classification is a challenging task owing to the great variability in both photographic and dermoscopic images. Previous work has made classification robust to photographic variability. Many previous techniques require extensive preprocessing, lesion segmentation and highly variable tasks across many fine-grained object categories^{4–6}. However, these methods are slow and computationally expensive. In contrast, our system requires no hand-coded features; it is trained on raw images directly, using only pixels and simple geometric features to end from images directly. We demonstrate that deep learning with electronic health records (EHRs) is anticipated to drive personalized medicine and improve healthcare quality. Constructing predictive statistical models typically requires extraction of causal predictor variables from normalized EHR data. In this work, we demonstrate that deep learning representations are capable of accurately predicting multiple clinical outcomes from EHRs. We use the EHRs of 216,221 adult patients hospitalized for at least 24 h. In the sequential format we use, the primary outcome is the total of 86,447 melanomas across sites (9.5% of all hospitalizations). We use a deep learning model to predict the probability of a patient having melanoma versus benign lesion. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. The study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

Previous work^{1–6} has shown that deep learning can be used to predict skin cancer with accuracy comparable to dermatologists. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. The study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

Recent work^{7–10} has shown that deep learning can be used to predict skin cancer with accuracy comparable to dermatologists. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. The study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

Current work^{11–13} has shown that deep learning can be used to predict skin cancer with accuracy comparable to dermatologists. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. The study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

Most recently, work^{14–16} has shown that deep learning can be used to predict skin cancer with accuracy comparable to dermatologists. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. The study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

Electronic health records (EHRs) contain a wealth of information that can be used to predict skin cancer with accuracy comparable to dermatologists. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. The study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

Dermatologists are the most common human malignancy^{1–3}, as it is primarily diagnosed visually, beginning with initial clinical screening and followed potentially by dermoscopic analysis, a biopsy and histopathological examination. Skin cancer classification is a challenging task owing to the great variability in both photographic and dermoscopic images. Previous work has made classification robust to photographic variability. Many previous techniques require extensive preprocessing, lesion segmentation and highly variable tasks across many fine-grained object categories^{4–6}. However, these methods are slow and computationally expensive. In contrast, our system requires no hand-coded features; it is trained on raw images directly, using only pixels and simple geometric features to end from images directly. We demonstrate that deep learning with electronic health records (EHRs) is anticipated to drive personalized medicine and improve healthcare quality. Constructing predictive statistical models typically requires extraction of causal predictor variables from normalized EHR data. In this work, we demonstrate that deep learning representations are capable of accurately predicting multiple clinical outcomes from EHRs. We use the EHRs of 216,221 adult patients hospitalized for at least 24 h. In the sequential format we use, the primary outcome is the total of 86,447 melanomas across sites (9.5% of all hospitalizations). We use a deep learning model to predict the probability of a patient having melanoma versus benign lesion. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. The study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

Previous work^{1–6} has shown that deep learning can be used to predict skin cancer with accuracy comparable to dermatologists. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. The study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

Recent work^{7–10} has shown that deep learning can be used to predict skin cancer with accuracy comparable to dermatologists. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. The study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

Current work^{11–13} has shown that deep learning can be used to predict skin cancer with accuracy comparable to dermatologists. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. The study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

Most recently, work^{14–16} has shown that deep learning can be used to predict skin cancer with accuracy comparable to dermatologists. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. The study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

Electronic health records (EHRs) contain a wealth of information that can be used to predict skin cancer with accuracy comparable to dermatologists. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. The study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

Dermatologists are the most common human malignancy^{1–3}, as it is primarily diagnosed visually, beginning with initial clinical screening and followed potentially by dermoscopic analysis, a biopsy and histopathological examination. Skin cancer classification is a challenging task owing to the great variability in both photographic and dermoscopic images. Previous work has made classification robust to photographic variability. Many previous techniques require extensive preprocessing, lesion segmentation and highly variable tasks across many fine-grained object categories^{4–6}. However, these methods are slow and computationally expensive. In contrast, our system requires no hand-coded features; it is trained on raw images directly, using only pixels and simple geometric features to end from images directly. We demonstrate that deep learning with electronic health records (EHRs) is anticipated to drive personalized medicine and improve healthcare quality. Constructing predictive statistical models typically requires extraction of causal predictor variables from normalized EHR data. In this work, we demonstrate that deep learning representations are capable of accurately predicting multiple clinical outcomes from EHRs. We use the EHRs of 216,221 adult patients hospitalized for at least 24 h. In the sequential format we use, the primary outcome is the total of 86,447 melanomas across sites (9.5% of all hospitalizations). We use a deep learning model to predict the probability of a patient having melanoma versus benign lesion. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. The study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

Previous work^{1–6} has shown that deep learning can be used to predict skin cancer with accuracy comparable to dermatologists. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. The study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

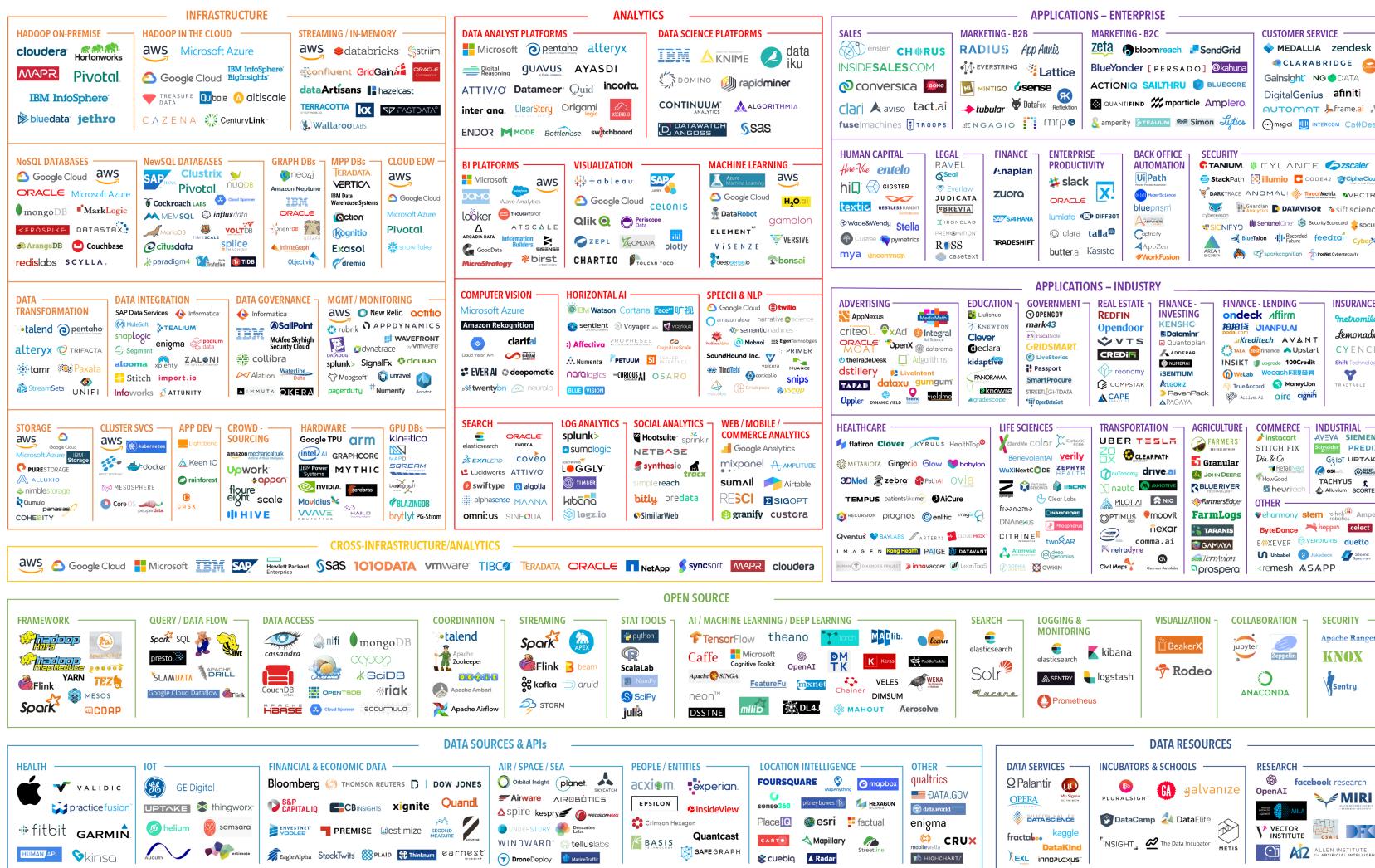
Recent work^{7–10} has shown that deep learning can be used to predict skin cancer with accuracy comparable to dermatologists. We believe that this approach can be used to create accurate and scalable predictions for a variety of clinical scenarios. The study of a particular prediction, we demonstrate that neural networks can be used to identify relevant information from the patient's chart.

Current work^{11–13} has shown that deep learning can be used to

Introducción a la Ciencia de Datos

Tecnologías

BIG DATA & AI LANDSCAPE 2018



Final 2018 version, updated 07/15/2018

© Matt Turck (@mattturck), Demi Obavomi (@demi_ obavomi), & FirstMark (@firstmarkcap)

mattturck.com/bigdata2018

FIRSTMARK
EARLY STAGE VENTURE CAPITAL

Herramientas

KNIME

WEKA

KEEL

Tableau

Lenguajes

Java, Scala

R

Python

Bibliotecas

scikit-learn

tensorflow

caret

PyTorch

h2o

Frameworks

Spark

Recursos

MOOC (Coursera,
Udacity, edX, etc.)

Kaggle

Cloud

Microsoft Azure

Google Cloud

Amazon Web Services



Unique in the Crowd: The privacy bounds of human mobility

SUBJECT AREAS:
APPLIED PHYSICS

APPLIED MATHEMATICS
STATISTICS
COMPUTATIONAL SCIENCE

Received
1 October 2012

Accepted
4 February 2013

Published
25 March 2013

Correspondence and
requests for materials
should be addressed to
Y.A. de M. (yao@mit.
edu)

Yves-Alexandre de Montjoye^{1,2}, César A. Hidalgo^{1,3,4}, Michel Verleysen² & Vincent D. Blondel^{1,5}

¹Massachusetts Institute of Technology, Media Lab, 20 Ames Street, Cambridge, MA 02139 USA, ²Université catholique de Louvain, Institute for Information and Communication Technologies, Electronics and Applied Mathematics, Avenue Georges Lemaitre 4, B-1348 Louvain-la-Neuve, Belgium, ³Harvard University, Center for International Development, 79 JFK Street, Cambridge, MA 02138, USA, ⁴Instituto de Sistemas Complejos de Valparaíso, Paseo 21 de Mayo, Valparaíso, Chile, ⁵Massachusetts Institute of Technology, Laboratory for Information and Decision Systems, 77 Massachusetts Avenue, Cambridge, MA 02139, USA.

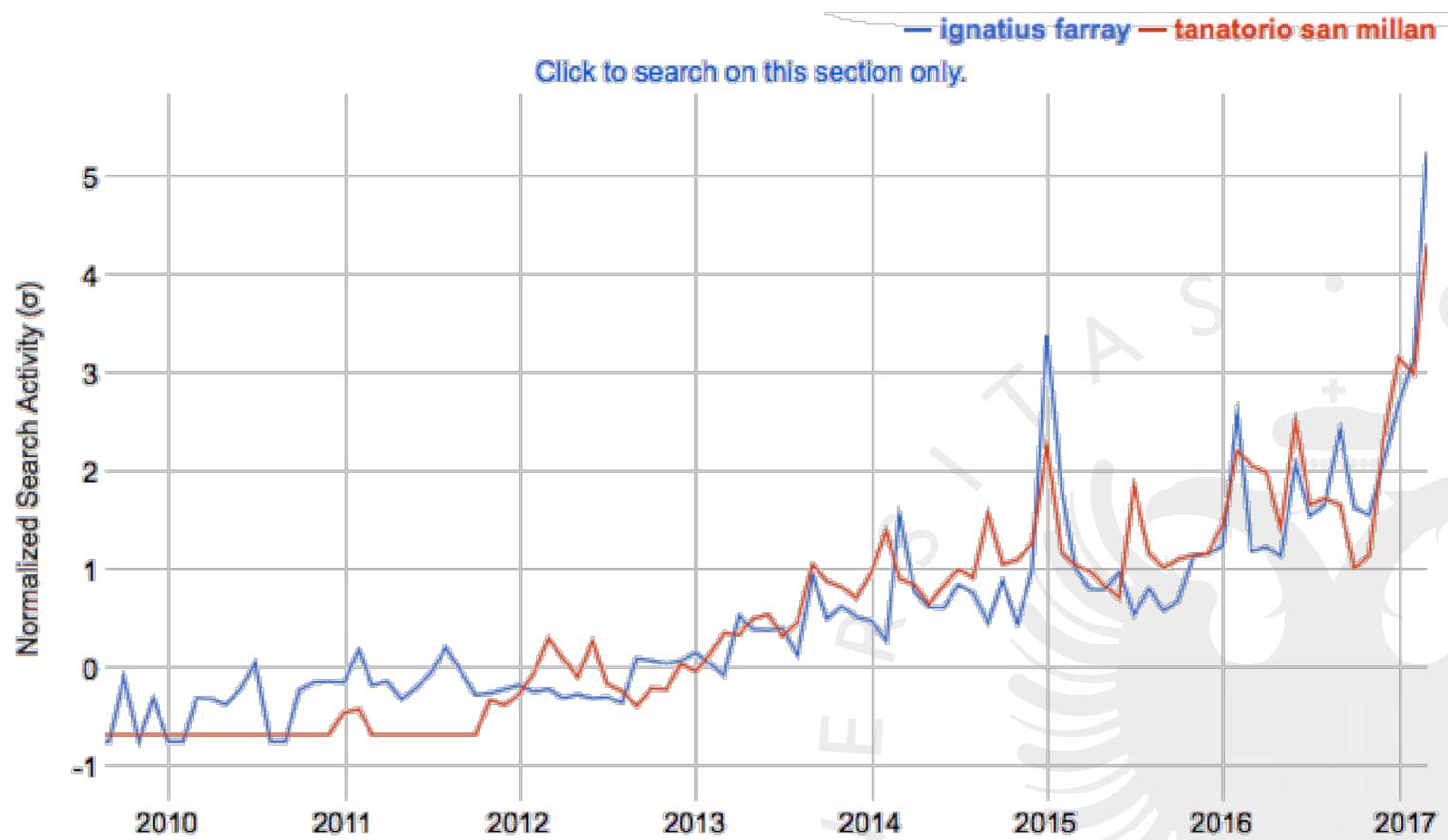
We study fifteen months of human mobility data for one and a half million individuals and find that human mobility traces are highly unique. In fact, in a dataset where the location of an individual is specified hourly, and with a spatial resolution equal to that given by the carrier's antennas, four spatio-temporal points are enough to uniquely identify 95% of the individuals. We coarsen the data spatially and temporally to find a formula for the uniqueness of human mobility traces given their resolution and the available outside information. This formula shows that the uniqueness of mobility traces decays approximately as the $1/10$ power of their resolution. Hence, even coarse datasets provide little anonymity. These findings represent fundamental constraints to an individual's privacy and have important implications for the design of frameworks and institutions dedicated to protect the privacy of individuals.

Derived from the Latin *privatus*, meaning "withdraw from public life," the notion of privacy has been foundational to the development of our diverse societies, forming the basis for individuals' rights such as free speech and religious freedom¹. Despite its importance, privacy has mainly relied on informal protection mechanisms. For instance, tracking individuals' movements has been historically difficult, making them de-facto private. For centuries, information technologies have challenged these informal protection mechanisms. In 1086, William I of England commissioned the creation of the *Domesday book*, a written record of major property holdings in England containing individual information collected for tax and draft purposes². In the late 19th century, de-facto privacy was similarly threatened by photographs and yellow journalism. This resulted in one of the first publications advocating privacy in the U.S. in which Samuel Warren and Louis Brandeis argued that privacy law must evolve in response to technological changes³.

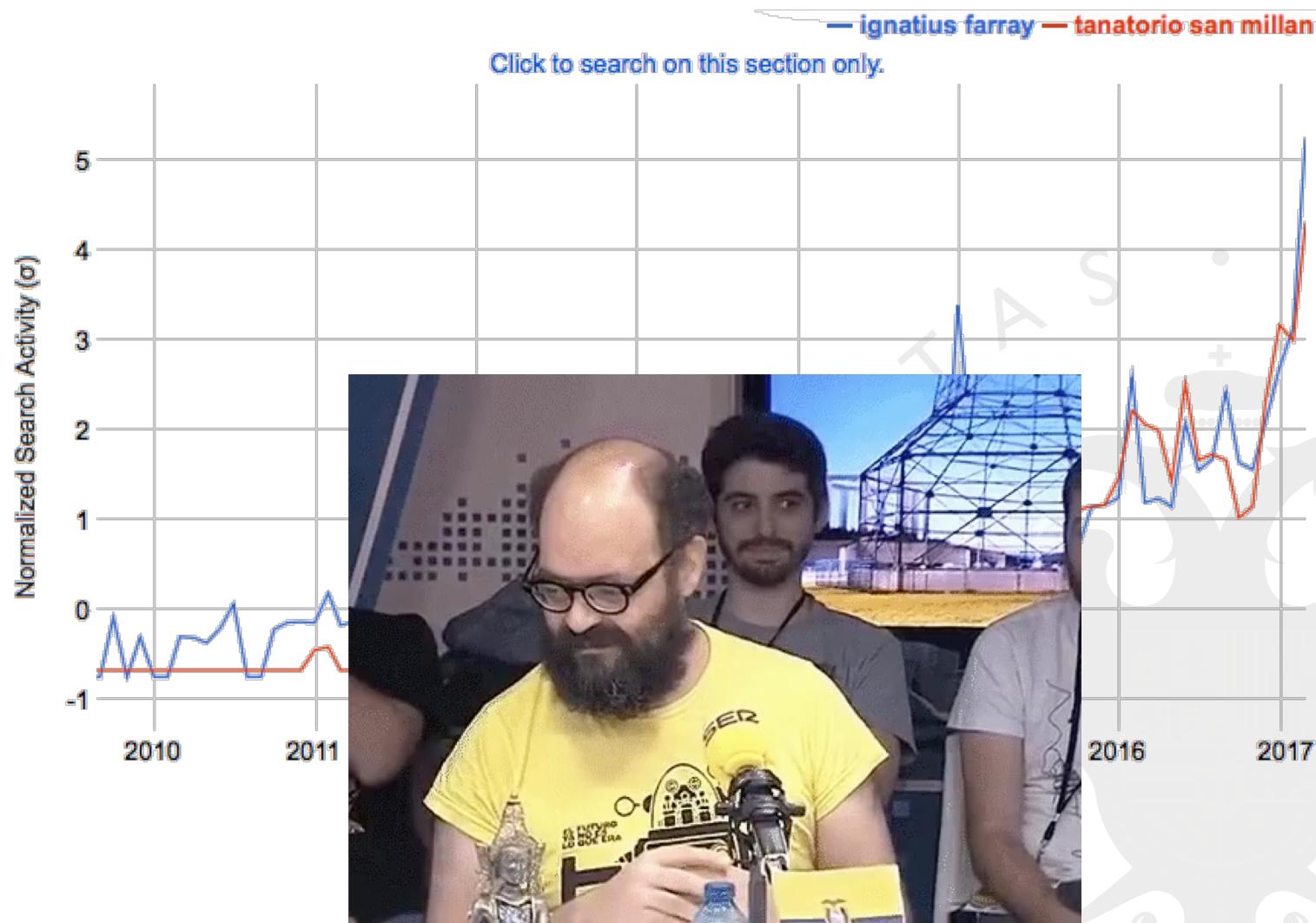
Y.A. de Montjoye, C.A. Hidalgo, M. Verleysen, V.D. Blondel (2013) Unique in the crowd: the privacy bounds of human mobility. *Nature Scientific Reports* 3, no. 1376.



High-Level Expert Group on Artificial Intelligence (2018)
Draft – Ethics Guidelines for Trustworthy AI.



<https://www.google.com/trends/correlate/search?e=Ignatius+Farray&e=tanatorio+san+millan&t=monthly&p=es&shift=2#default,20>



V. Mayer-Schonberger, K. Cukier (2013) *Big Data, A revolution that will transform how we live, work and think.* Harcourt Publishing.

C. Rudder (2014) *Dataclysm*. 4th State.

C. Hidalgo (2015) *Why information grows*. Penguin Books.