data
analytics
group

# STATISTICAL MODELING VS. MACHINE LEARNING – SIMILARITIES AND DIFFERENCES

## WK RECH 2019, Frankfurt School of Finance & Management

# ILLUSTRATIVE EXAMPLES: GOOGLE BOOKS AND GOOGLE FLU TRENDS

**Paper**

## Data and Methods



Source: Michel et al. (2011)

## Selected Findings



Source: Michel et al. (2011)

## Selected Findings



Source: Michel et al. (2011)

**Try It Yourself at https://books.google.com/ngrams**

## Google Search Terms

## Original Paper

# LETTERS

## Detecting influenza epidemics using search engine query data

Jeremy Ginsberg[1], Matthew H. Mohebbi[1], Rajan S. Patel[1], Lynnette Brammer[2], Mark S. Smolinski[1] & Larry Brilliant[1]
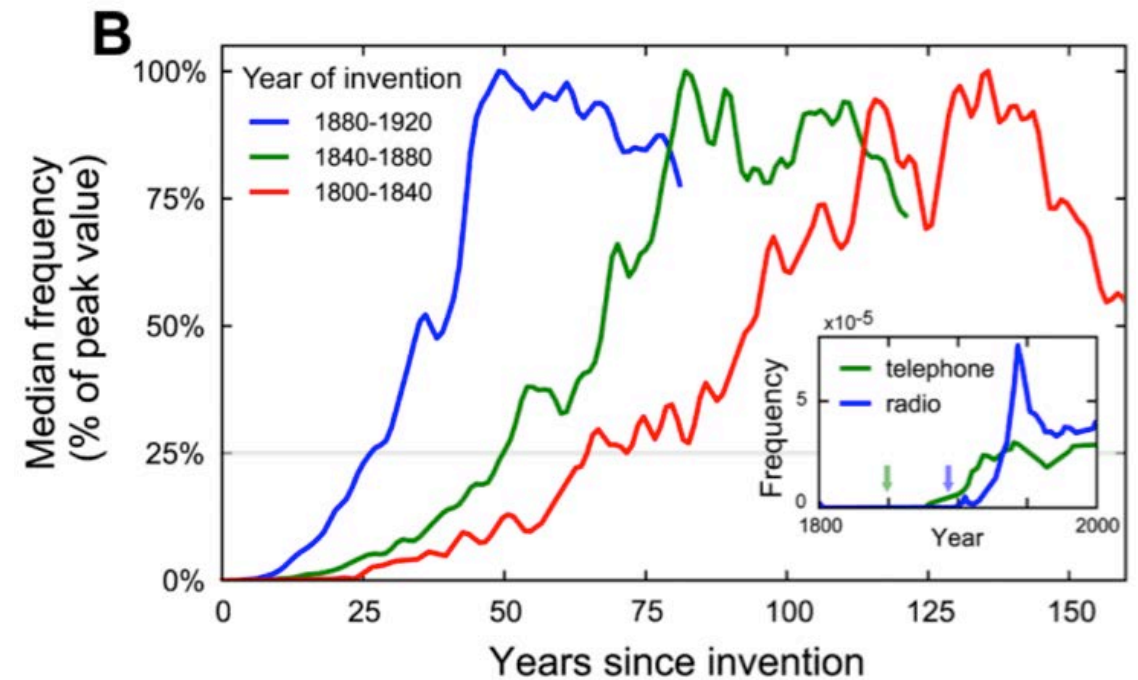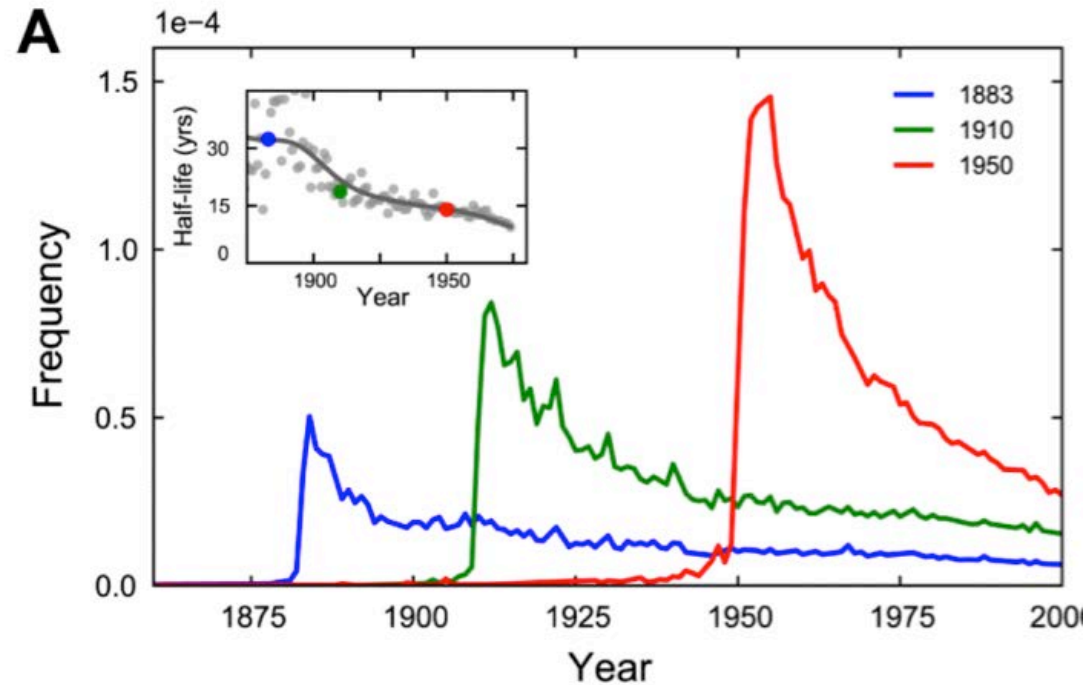
Seasonal influenza epidemics are a major public health concern, causing tens of millions of respiratory illnesses and 250,000 to 500,000 deaths worldwide each year[1]. In addition to seasonal influenza, a new strain of influenza virus against which no previous immunity exists and that demonstrates human-to-human transmission could result in a pandemic with millions of fatalities[2]. 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 lines[5] and over-the-counter drug sales[6]. About 90 million American adults are believed to search online for information about specific diseases or medical problems each year[7], 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 website[8], and user clicks on a search keyword advertisement in Canada[9]. A set of Yahoo search queries containing the words 'flu' or 'influenza' were found to correlate with virological and mortality surveillance data over multiple years[10].

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: $\mathrm{logit}(I(t)) = \alpha\,\mathrm{logit}(Q(t)) + \varepsilon$, where $I(t)$ is the percentage of ILI physician visits, $Q(t)$ is the ILI-related query fraction at time $t$, $\alpha$ is the multiplicative coefficient, and $\varepsilon$ is the error term. $\mathrm{logit}(p)$ is simply $\ln(p/(1 - p))$.

Publicly 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).

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

1012

## Google Flu Study



Data available as of 4 February 2008

Data available as of 3 March 2008

Data available as of 31 March 2008

Data available as of 12 May 2008

"The final model was validated on 42 points per region of previously untested data from 2007 to 2008, which were excluded from all previous steps. Estimates generated for these 42 points obtained a **mean correlation of 0.97** (min: 0.92, max: 0.99, n: 9 regions) with the CDC-observed ILI percentages.

Harnessing the collective intelligence of millions of users, **Google web search logs can provide one of the most timely, broad-reaching influenza monitoring systems available today**.

Source: Ginsberg et al. (2009)

# GOOGLE FLU TRENDS STUDY

## How to Forecast the Flu with Google Search Terms?

X                                                                Y

| | Date | Location | "coughing" | „soar throat" | "cold" | … | ILI level |
|---|---|---|---|---|---|---|---|
| **Observation 1** | 01.01.2019 | Frankfurt | 2321 | 3441 | 5513 | ... | 0.020 |
| **Observation 2** | 01.01.2019 | Berlin | 1968 | 3201 | 4236 | ... | 0.008 |
| **Observation 3** | 02.01.2019 | Frankfurt | 2331 | 3446 | 5657 | ... | 0.021 |
| **…** | ... | ... | ... | ... | ... | ... | |

$$Y = f(X) + \varepsilon$$

**Response**

## Traps in Big Data Analysis

**TWO PARADIGMS:
STATISTICAL MODELING VS.
MACHINE LEARNING**

## What is Statistics?



Source: https://lagunita.stanford.edu/courses/HumanitiesSciences/StatLearning/Winter2016/about

" "Vast amounts of data are being generated in many fields, and the statistician's job is to make sense of it all: to extract important patterns and trends, and to understand "what the data says"." (Friedman et al., 2001)

" "A branch of mathematics dealing with the collection, analysis, interpretation, and presentation of masses of numerical data" (Merriam-Webster)

## The Scientific Method: Hypothetico-deductive



Source: https://www.sciencebuddies.org

# THE STATISTICAL MODELING PARADIGM

**Example: Lady Tasting Tea**
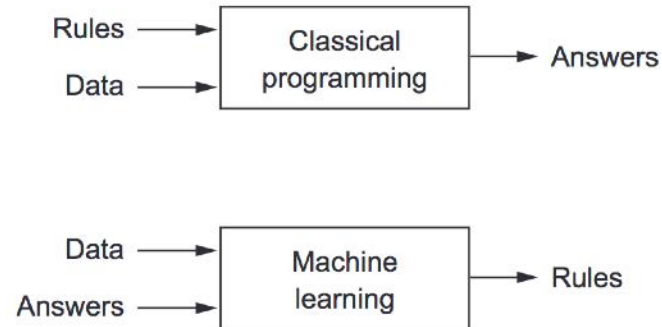


Source: https://www.youtube.com/watch?v=lgs7d5saFFc
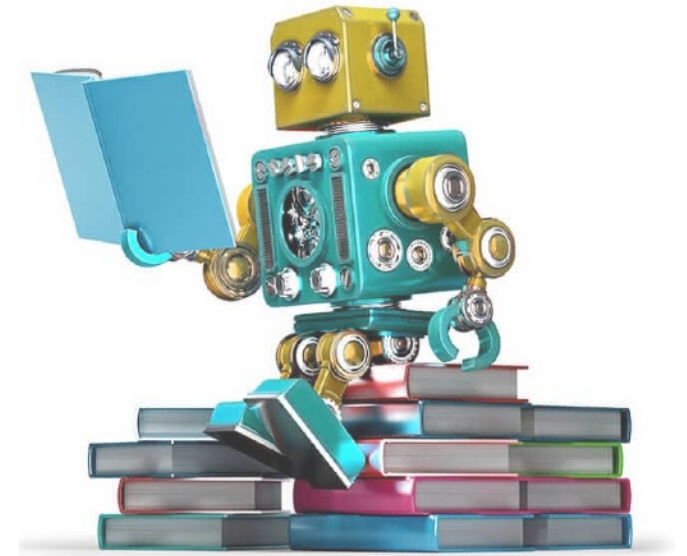
## What is Machine Learning?

"The field of study that gives computers the ability to learn without being explicitly programmed." (Samuel, 1959)



Source: Chollet (2018)
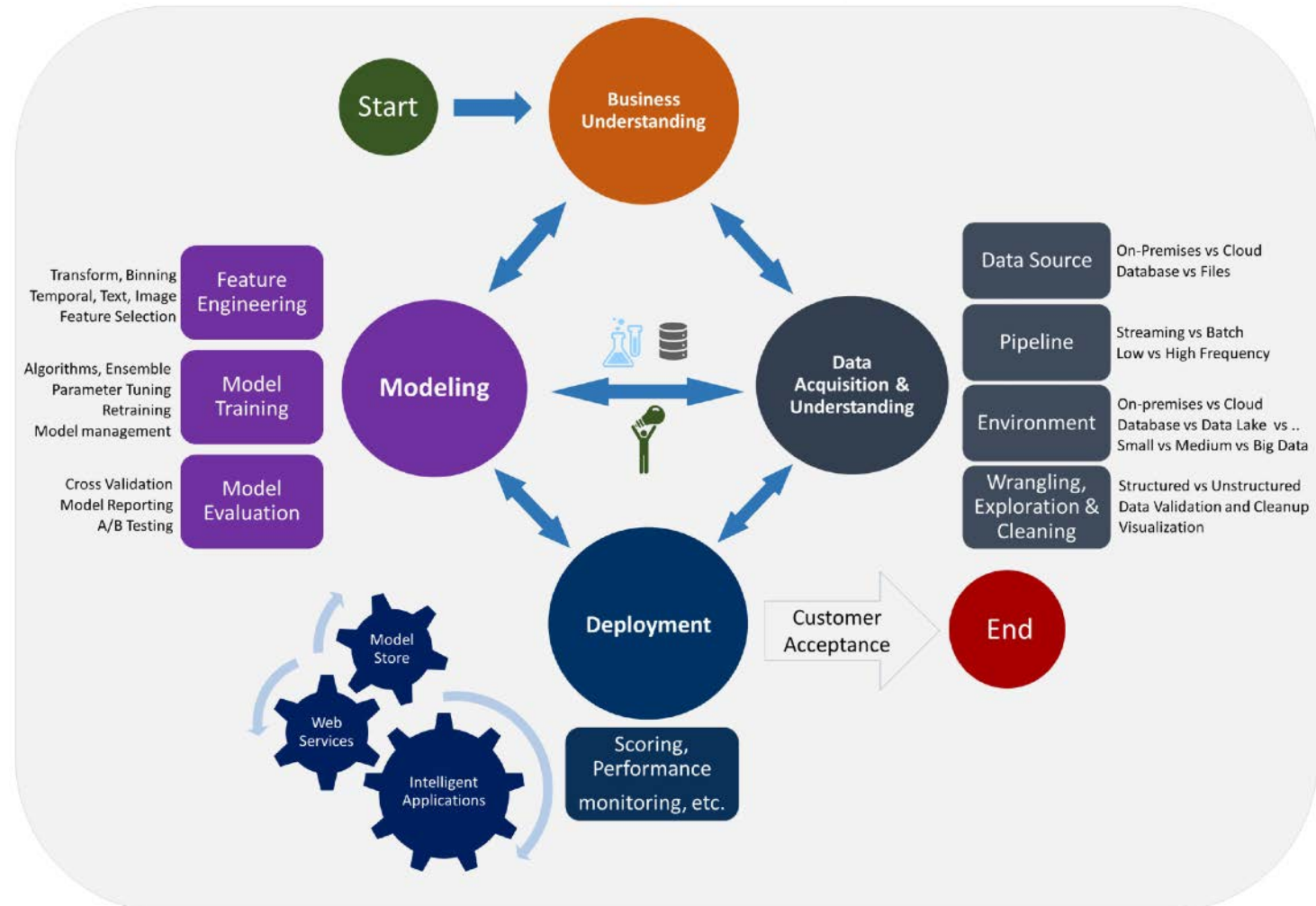
"A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if its performance at tasks in T, as measured by P, improves with experience E." (Mitchell, 1997)
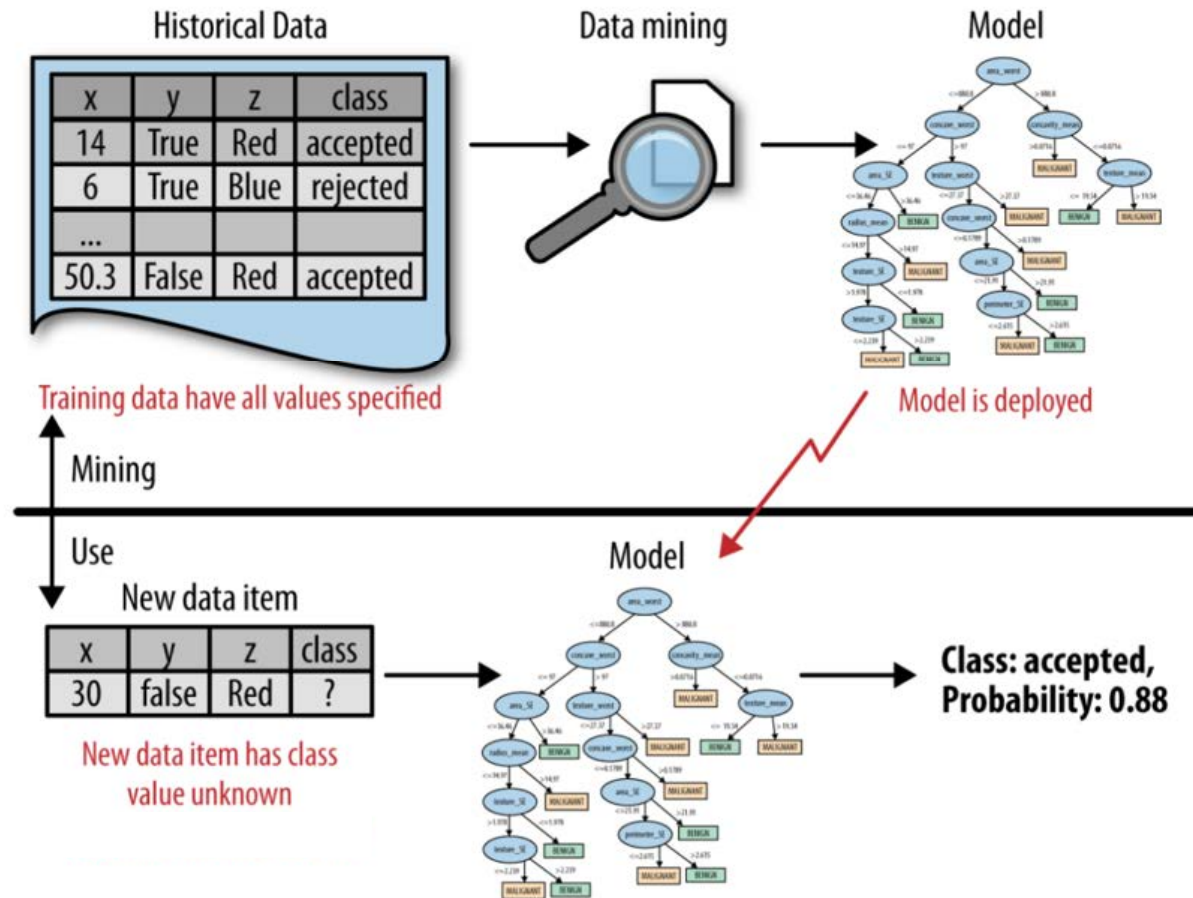
## The Data Science Lifecycle: Data-driven, Inductive, Iterative



Source: https://docs.microsoft.com/en-us/azure/machine-learning/team-data-science-process/lifecycle

## Example: Supervised Machine Learning for Credit Risk Scoring



Source: Provost & Fawcett (2013)

# BANKRUPTCY PREDICTION: THE SM VS. ML WAY

## Governing Machine Learning in Governments

Per Rådberg Nagbøl
Phd Student
ITU Copenhagen

EARLY WARNING EUROPE

EARLY WARNING EUROPE

**Early Warning Europe provides free, impartial and confidential counselling to companies in distress**

## Research Design

$$Z = .012X_1 + .014X_2 + .033X_3 + .006X_4 + .999X_5$$

where
- $X_1 =$ Working capital/Total assets
- $X_2 =$ Retained Earnings/Total assets
- $X_3 =$ Earnings before interest and taxes/Total assets
- $X_4 =$ Market value equity/Book value of total debt
- $X_5 =$ Sales/Total assets
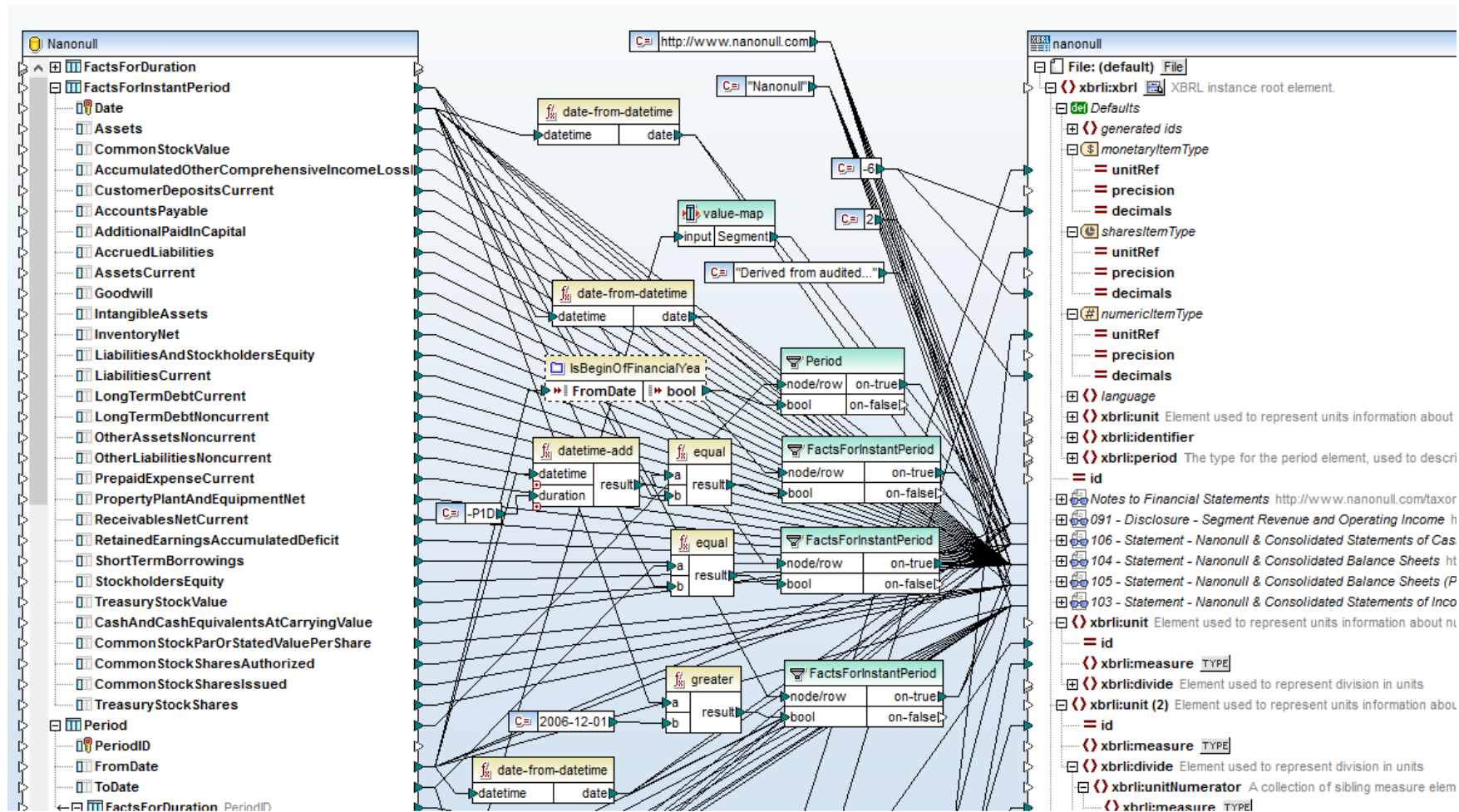- $Z =$ Overall Index

$$\log \left( \frac{p(X)}{1-p(X)} \right) = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p$$



**Ask a Question**

**Do Background Research**

**Construct a Hypothesis**

**Test with an Experiment**

**Procedure Working?**

**Troubleshoot procedure. Carefully check all steps and set-up.**

**No** **Yes**

**Experimental data becomes background research for new/future project. Ask new question, form new hypothesis, experiment again!**

**Analyze Data and Draw Conclusions**

**Results Align with Hypothesis**

**Results Align Partially or Not at All with Hypothesis**

**Communicate Results**

## Dataset: 62.000+ Annual Reports in XBRL from 2014

## Results of Logistic Regression

```
=============================================
                    Dependent variable:
                    -----------------------------
                    bankruptcy
---------------------------------------------
Constant            -4.027*** (0.036)
X1                  -1.501*** (0.067)
X2                  -0.00000 (0.00004)
X3                   0.00001 (0.00005)
X4                  -0.00000 (0.00000)
---------------------------------------------
Observations             50,377
Log Likelihood        -4,432.134
Akaike Inf. Crit.      8,874.268
=============================================
Note:           *p<0.1; **p<0.05; ***p<0.01
```

$X_1$ = Working capital/Total assets
$X_2$ = Retained Earnings/Total assets
$X_3$ = Earnings before interest and taxes/Total assets
$X_4$ = Market value equity/Book value of total debt
$X_5$ = Sales/Total assets

## Training and Test Sets



Dataset

Training set

Test set

Source: James et al. (2013)

Wirtschaftsinformatik, insb. Data Analytics | Prof. Dr. Oliver Müller

30

## Predictive Accuracy of Logistic Regression on Test Set

## Going Beyond Financial Ratios (i.e., reading 62.000 annual reports)

## The Bag-Of-Words (BOW) Model

## The Bag-Of-Words (BOW) Model

- Treat every document as a unordered set of words
- Ignore word order, sentence structure, and punctuation

- Tidy data frame:
  - Every document is an observation (row)
  - Every word is a variable (column)
  - The presence of a word in a document (aka. token) is represented by the cell values

## The Bag-Of-Words (BOW) Model

| | "word 1" | "word 2" | "word 3" | "word 4" | "word 5" | … | Bankruptcy? |
|---|---|---|---|---|---|---|---|
| **Lego** | 0 | 0 | 0 | 2 | 1 | … | No |
| **Maersk** | 0 | 0 | 0 | 2 | 0 | … | No |
| **Jysk Fragt** | 0 | 1 | 1 | 1 | 0 | … | Yes |
| **…** | … | … | … | … | … | … | … |

$$Y = f(X) + \varepsilon$$

## Tree-based Classification Algorithms

## Example: Loan Default

## Example: Loan Default



Source: Provost & Fawcett (2013)

## Example: Loan Default



| Claudio | $115,000 | 40 | no |
|---------|----------|-----|-----|

Source: Provost & Fawcett (2013)

## The CART Algorithm: Top-down, Greedy Search

- It is computationally infeasible to consider all possible sequences and combinations of splits.

- Instead, do **recursive binary partitioning**
  - **Top-down**: Start with zero splits and successively partition the feature space into two parts.
  - **Greedy**: At each step, make the best possible split at that particular step (i.e., the split with the highest information gain, i.e., reduction in entropy).
  - Stop when some condition (e.g., minimal number of observations in one leaf) is met.

- That is, we consider all predictors $X_1, . . ., X_p$, and all possible split points $s$ for each of the predictors, and then choose the predictor and split point with the highest information gain at each step.
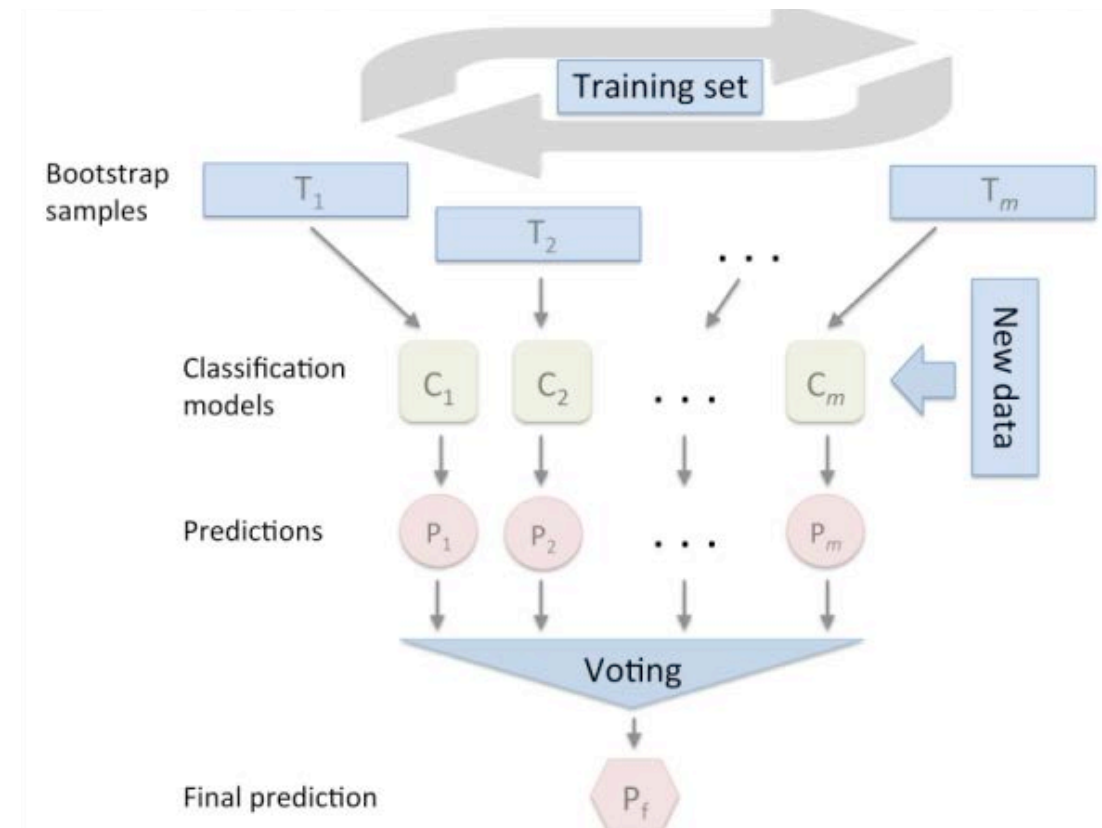
Source: Provost & Fawcett (2013)

**Random Forests**

## Bootstrap Aggregation (Bagging)

- A way to reduce overfitting of a machine learning algorithm is to take many training sets from the population, build a separate model on each training set, and average the resulting predictions.
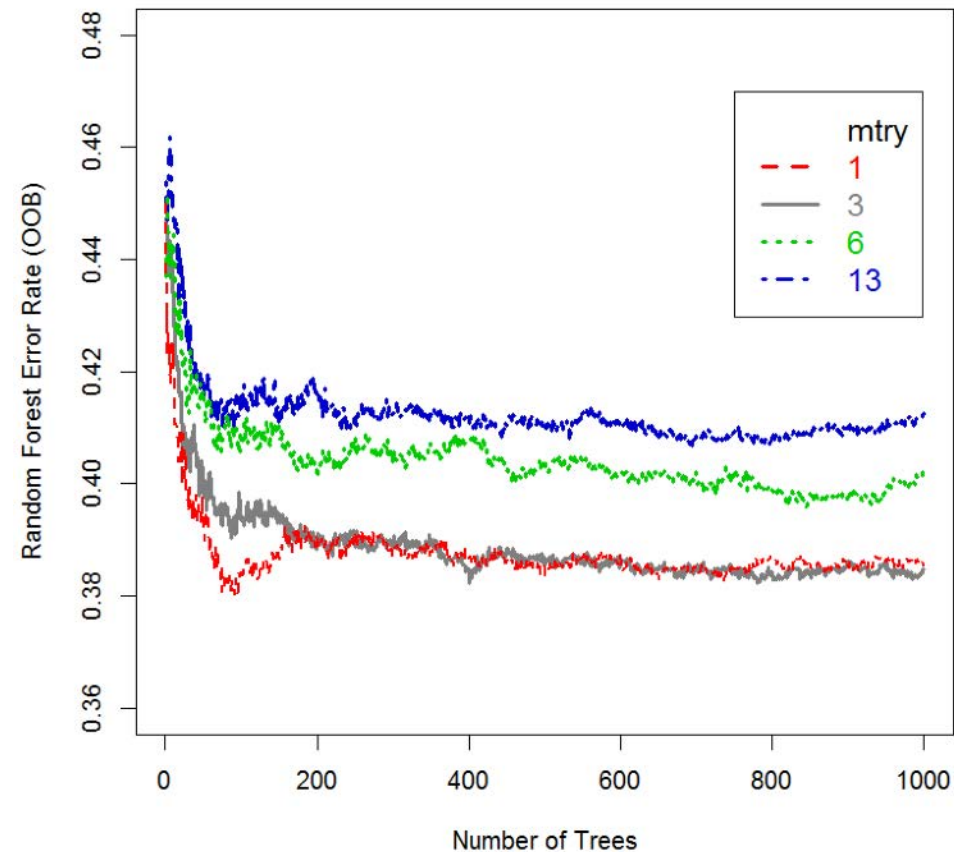


Source: James et al. (2013)
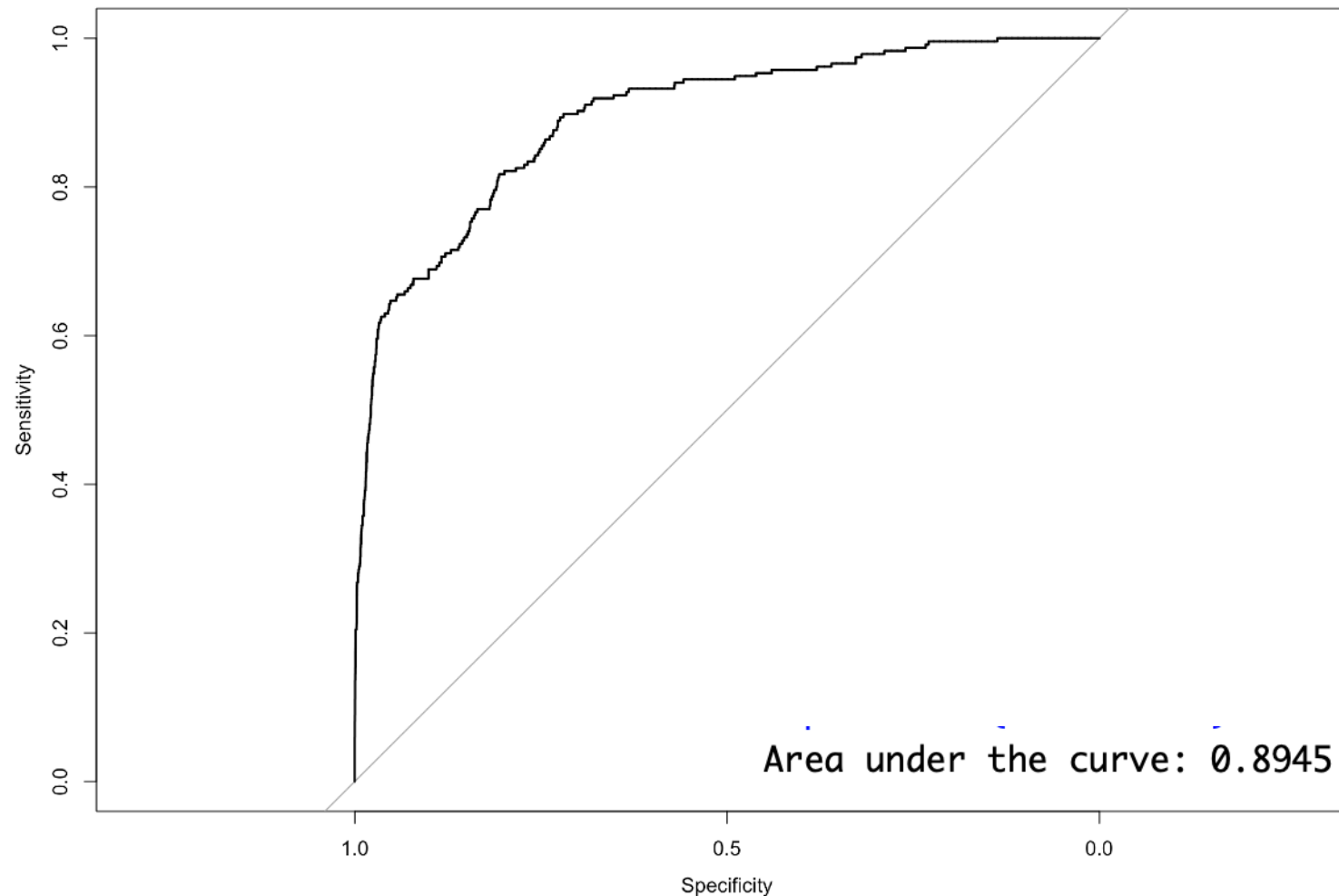
## From Bagging to Random Forest

- **Following the idea of bagging**, we draw multiple random samples (bootstrap samples) from the training data and create a decision tree on each sample
  - Typically, 2/3 of the rows in the training set

- However, in Random Forests we **allow only a random subset ($m$) of all the predictors ($p$) to be used at each split of the decision tree**
  - Typically, $m = SQRT(p)$

- **Why** does this work?
  - In bagging, if there is one strong predictor, all the trees will use this predictor in the top split
    - ➔ all of the trees will look quite similar to each other
    - ➔ their predictions will be highly correlated
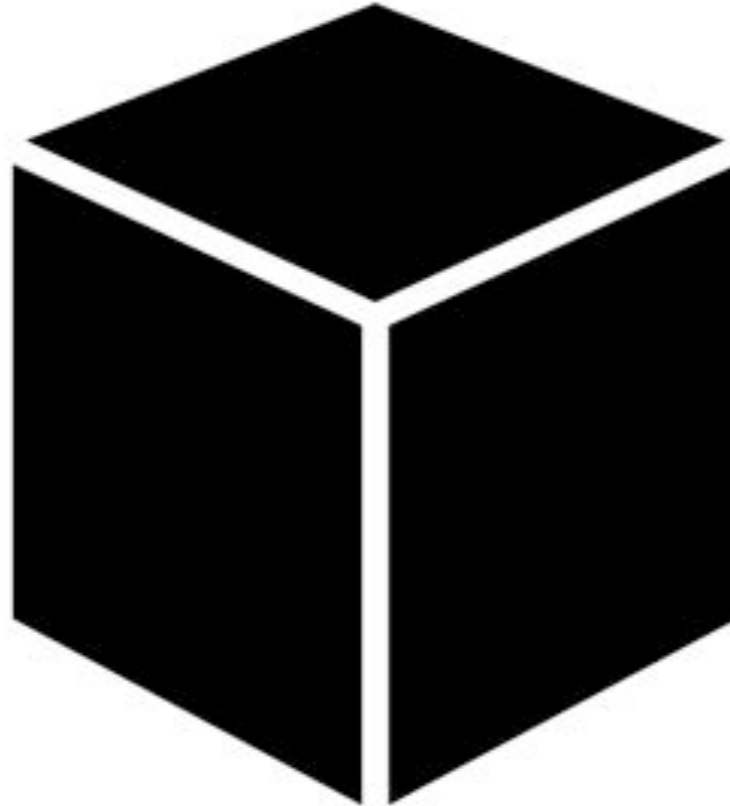    - ➔ only a little bit of variance will be removed

Source: James et al. (2013)

## A Random Forest has many Trees

**Predictive Accuracy of Random Forest with BOW on Test Set**

## How Does the Model Look Like?
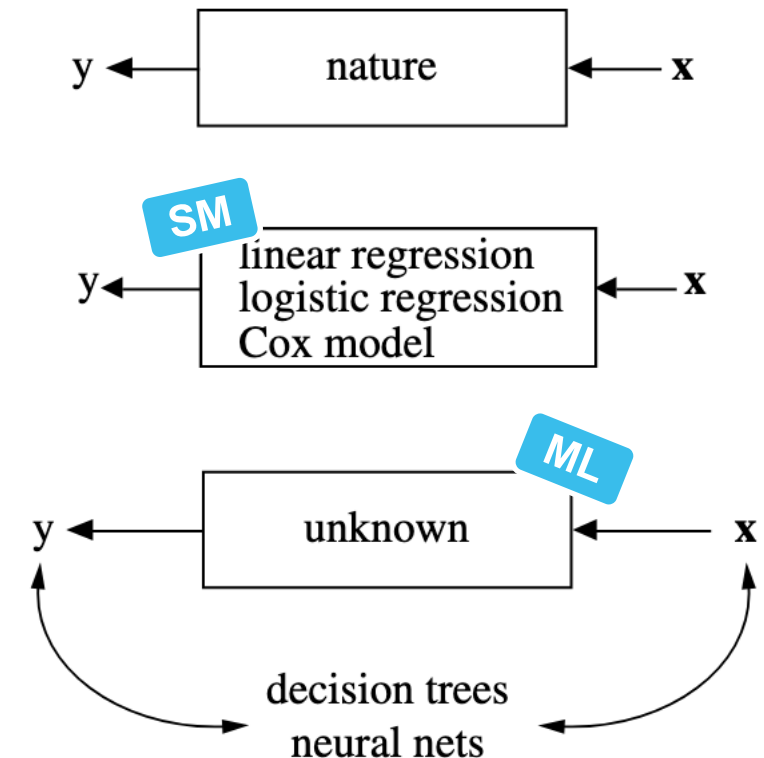
**Black Box**

## Next Steps

- Include more text from annual reports (e.g., management review, management's statement, CSR)

- Quantify the informativeness of different sections of annual reports

- Use artificial neural networks to better capture syntax and semantics of text

- Try to open the black box of machine learning algorithms

# REFLECTIONS

## Similarities and Differences of Statistical Modeling and Machine Learning

- **Data**
  - Both work with almost the same data structures
  - ML wrangles messy data until it fits into rows and columns

- **Methods**
  - Both use regression and classification techniques
  - SM applies mainly additive linear models
  - ML uses on non-linear methods that work on high-dimensional data
  - ML makes use of unsupervised techniques for data preparation

- **Process**
  - SM is theory/hypothesis-driven (no fishing for correlations!)
  - ML is mainly data-driven

- **Outputs**
  - SM: focus on causal explanations
  - ML: focus on predictive accuracy (on unseen test data!)

Source: Breiman (2001)

## Paper

**Prof. Dr. Oliver Müller**
Lehrstuhl für Wirtschaftsinformatik, insb. Data Analytics
Universität Paderborn
Warburger Str. 100, 33098 Paderborn
R: Q2.457
E: oliver.mueller@uni-paderborn.de
T: +49-5251-605245
W: https://wiwi.uni-paderborn.de/dep3/mueller/

Mat Velloso
@matvelloso

Difference between machine learning and AI:

If it is written in Python, it's probably machine learning

If it is written in PowerPoint, it's probably AI

Tweet übersetzen

02:25 · 23.11.18 · Twitter Web Client

8.083 Retweets   22,5K „Gefällt mir"-Angaben

- Michel, J. B., Shen, Y. K., Aiden, A. P., ... & Pinker, S. (2011). Quantitative analysis of culture using millions of digitized books. Science, 331(6014), 176-182.
- Ginsberg, J., Mohebbi, M. H., Patel, R. S., Brammer, L., Smolinski, M. S., & Brilliant, L. (2009). Detecting influenza epidemics using search engine query data. Nature, 457(7232), 1012-1014.
- Lazer, D., Kennedy, R., King, G., & Vespignani, A. (2014). The parable of Google Flu: traps in big data analysis. Science, 343(6176), 1203-1205.
- Samuel, A. L. (1959). Some studies in machine learning using the game of checkers. IBM Journal of Research and Development, 3(3), 210-229.
- Mitchell, T. M. (1997). Machine learning. McGraw Hill.
- Friedman, J., Hastie, T., & Tibshirani, R. (2001).The elements of statistical learning. Springer.
- Provost, F., & Fawcett, T. (2013). Data Science for Business: What you need to know about data mining and data-analytic thinking, O'Reilly Media.
- James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). An Introduction to Statistical Learning. Springer.
- Müller, O., Junglas, I., Brocke, J. V., & Debortoli, S. (2016). Utilizing big data analytics for information systems research: challenges, promises and guidelines. European Journal of Information Systems, 25(4), 289-302.
- Breiman, L. (2001). Statistical modeling: The two cultures (with comments and a rejoinder by the author). Statistical science, 16(3), 199-231.

- Icons: Depb Dew, Jemis Mali from the Noun Project