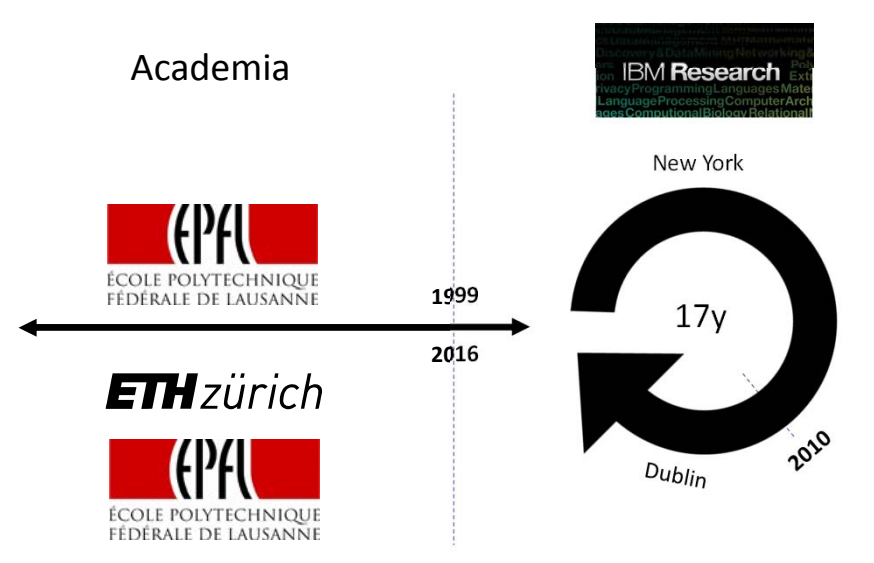ETH-RAT

SDSC

# REAL-WORLD DATA SCIENCE

Olivier Verscheure, PhD

Swiss Data Science Center

EPFL & ETH Zurich

RH Vaud, Lausanne – January 16, 2018

## About me



Academia

IBM Research

New York

1999
2016

17y

2010

Dublin

ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

**ETH**zürich

ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

SDSC

# What Do You See?



# Dublin City Data Hub

## Data *is the New Oil*



*The Economist, May 2017*

## Big Data, Bad Data

## A Fragmented Ecosystem



Algorith... ...cs

Data management

...ne learning

**GAP**

Visual ana... ...esearch

What is the hyperplane that best separates two classes of points in multidimensional space?

How can I best match the right drug with the right dosage to the right patient at the right time?

## Machine Learning in a nutshell



$$f : \quad \longrightarrow \text{Cat}$$

SDSC

## Machine Learning in a nutshell



$$\min_{f \in \mathcal{F}} \sum_{(\mathbf{x}, y) \in D} \ell(y, f(\mathbf{x}))$$

## A Sobering View on Machine Learning



*THIS* IS YOUR MACHINE LEARNING SYSTEM?

YUP! YOU POUR THE DATA INTO THIS BIG PILE OF LINEAR ALGEBRA, THEN COLLECT THE ANSWERS ON THE OTHER SIDE.

WHAT IF THE ANSWERS ARE WRONG?

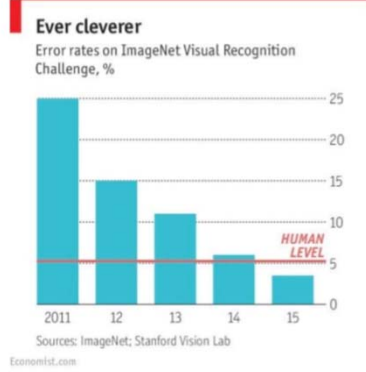JUST STIR THE PILE UNTIL THEY START LOOKING RIGHT.

© XKDC

## Recent advances in Machine Learning

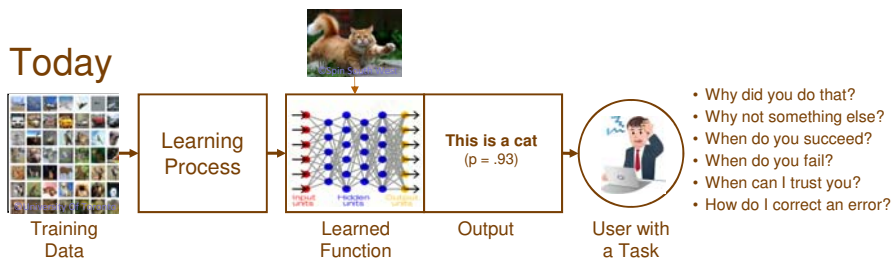— Recognise an object from a photo

**ImageNet Challenge**

IMAGENET

- 1,000 object classes (categories).
- Images:
  - 1.2 M train
  - 100k test.

**Ever cleverer**
Error rates on ImageNet Visual Recognition Challenge, %

HUMAN LEVEL

2011   12   13   14   15

Sources: ImageNet; Stanford Vision Lab
Economist.com

SDSC

*The Economist, May 2017*

---

## Explainable AI – What Are We Trying To Do?

DARPA

**Today**

Training Data → Learning Process → Learned Function → Output **This is a cat** (p = .93) → User with a Task

Input Units   Hidden Units   Output Units

- Why did you do that?
- Why not something else?
- When do you succeed?
- When do you fail?
- When can I trust you?
- How do I correct an error?

© David Gunning, DARPA/I20

Distribution Statement "A" (Approved for Public Release, Distribution Unlimited)

12

6

## Fooling deep neural net classifiers

**Title:** Universal adversarial perturbations
**Authors:** Moosavi-Dezfooli, Seyed-Mohsen; Fawzi, Alhussein; Fawzi, Omar; Frossard, Pascal
**Publication:** eprint arXiv:1610.08401
**Publication Date:** 10/2016



*This is not a sock*

- It's an Indian elephant!
- At least after adding a universal noise to the image
- Deep learning models do not mimic brain activity
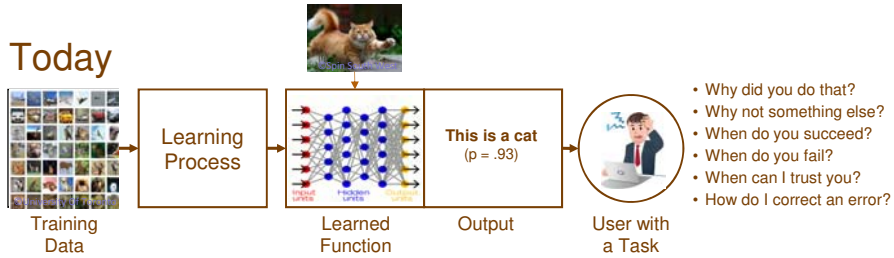
## Fooling deep neural net classifiers

- Autonomous driving!



- Turning a **STOP** sign into a **45 mph speed limit**

SDSC

Explainable AI – What Are We Trying To Do?

Today

Training Data → Learning Process → Learned Function → Output: This is a cat (p = .93) → User with a Task

- Why did you do that?
- Why not something else?
- When do you succeed?
- When do you fail?
- When can I trust you?
- How do I correct an error?

Tomorrow

Training Data → New Learning Process → Explainable Model → Explanation Interface → User with a Task

This is a cat:
- It has fur, whiskers, and claws.
- It has this feature:

- I understand why
- I understand why not
- I know when you'll succeed
- I know when you'll fail
- I know when to trust you
- I know why you erred

© David Gunning , DARPA/I20

Distribution Statement "A" (Approved for Public Release, Distribution Unlimited)

15

# Anecdotal digression

- Forecasting demand in electricity (France)

$$
\begin{aligned}
y_k \;=\; & \beta^{\text{Intercept}} + f^{\text{Trend}}(k) + f^{\text{LagLoad}}(y_{k-48}) + \sum_{l=1}^{6} \mathbf{1}(x_k^{\text{DayType}} = l)(\beta_l^{\text{DayType}} + f_l^{\text{TimeOfDay}}(x_k)) \\
& + f^{\text{CloudCover}}(x_k) + f^{\text{Temperature/TimeOfDay}}(x_k) + f^{\text{LagTemperature}}(x_{k-48}) \\
& + f^{\text{TimeOfYear}}(x_k) + x_k^{\text{LoadDecrease}} f^{\text{LoadDecrease}}(x_k) + \epsilon_k.
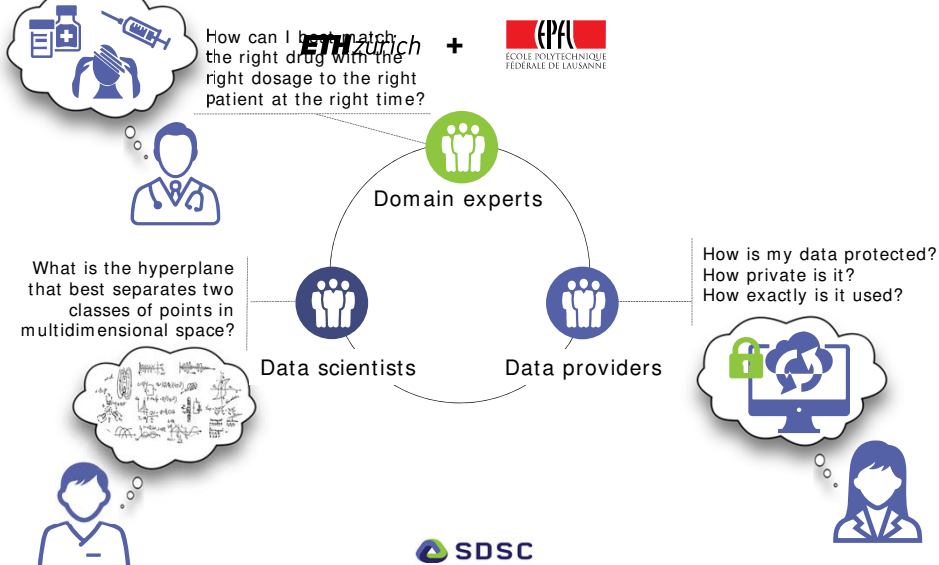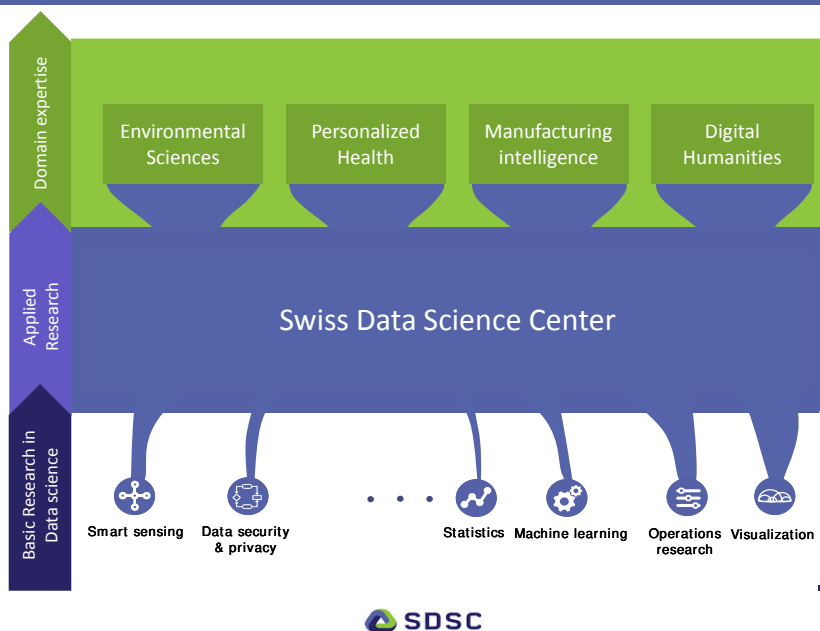\end{aligned}
$$

**Transfer functions learned from data:**



Time of Day
Temperature (°C)

Jan 1st                    Dec 31st
Time of Year

© IBM Research

8

# Swiss Data Science Center (SDSC)

Multi-disciplinary teams of 40 cell both in academia and dlatastry scientists, and domain experts



How can I best match the right drug with the right dosage to the right patient at the right time?

**ETH** zürich + EPFL ÉCOLE POLYTECHNIQUE FÉDÉRALE DE LAUSANNE

Domain experts

What is the hyperplane that best separates two classes of points in multidimensional space?

Data scientists          Data providers

How is my data protected?
How private is it?
How exactly is it used?

SDSC

# Where does SDSC fit?



| Domain expertise | Environmental Sciences | Personalized Health | Manufacturing intelligence | Digital Humanities |

Applied Research

Swiss Data Science Center

Basic Research in Data science

Smart sensing    Data security & privacy    . . .    Statistics    Machine learning    Operations research    Visualization

SDSC

9

# What will the SDSC offer?

### Embedded R&D collaboration

We engage in academic and industrial collaborations requiring large-scale distributed data processing (Big & Fast Data) and/or advanced analytics (machine learning & statistics) combined with an in-depth knowledge in select domains

### Domain-specific Insights as a Service

We provide secure access to our cloud-hosted analytics platform - **RENGA**, a highly scalable open software platform offering a one-stop-shop for hosting and exploring curated, calibrated and possibly anonymized data at scale, at-rest or in-motion.

### Open (Data) Science

RENGA offers user-friendly tooling and services to help with the adoption of Open Science, fostering research productivity and excellence.

*SDSC Analytics Platform*

SDSC

# Status quo in Data Science



© Oxford Creativity 2012

credit: oxford creativity, https://www.triz.co.uk/

## Facilitate communication to foster innovation



credit: oxford creativity, https://www.triz.co.uk/

## Foster multidisciplinary collaborations



credit: oxford creativity, https://www.triz.co.uk/

## Available as Open Source (Apache v2)



http://get-renga.io

---



## SIMPLE USE CASE
### PREDICTING WHEN PEOPLE QUIT THEIR JOBS

## Objectives

1. Understand why some of the company's **best** and **most experienced** employees are leaving prematurely
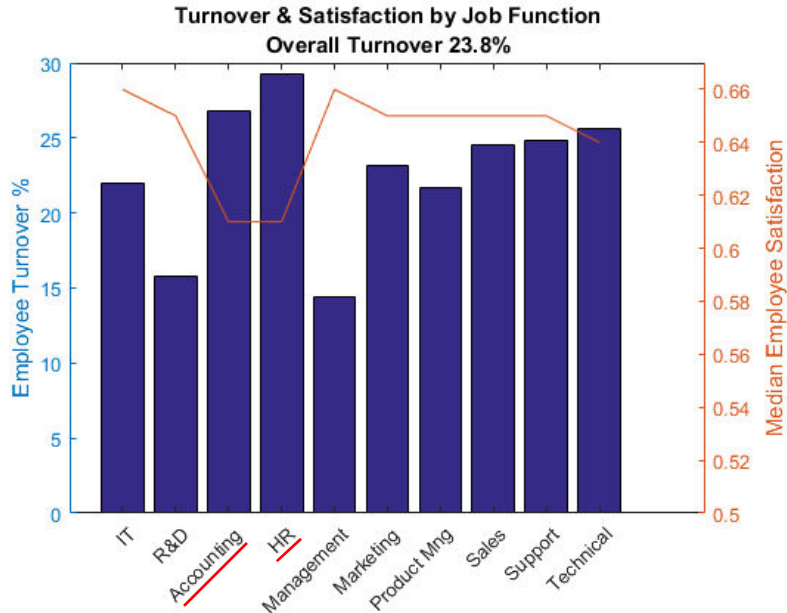
2. Predict which valuable employees will leave next

SDSC

## Dataset published by IBM

- Fictitious large company
- 14,999 employees
- 10 data fields include:
  - Employee satisfaction level, scaling 0 to 1
  - Last evaluation, scaling 0 to 1
  - Number of projects
  - Average monthly hours
  - Time spent at the company in years
  - Whether they have had a work accident
  - Whether they have had a promotion in the last 5 years
  - Sales (which actually means job function)
  - Salary - low, medium or high
  - Whether the employee has left

SDSC

## How Bad Is Turnover at This Company?



Turnover & Satisfaction by Job Function
Overall Turnover 23.8%

## Defining Who Are the "Best"



Distribution of Last Evaluation

## Defining Who Are the "Most Experienced"



Time Spent @ Company Among High Performers

## Job Satisfaction Among High Risk Group



Satisfaction Level of the Best and Most Experienced

# Was It For Money?



# Building a Predictive Model

# Explaining the Model



# Operationalizing Action

# THANK YOU!

http://www.datascience.ch

**Twitter**: @SDSCdatascience