

AI for Health

Erwan Scornet (Maître de conférences, Ecole Polytechnique)

1 Introduction to AI

- No data project without data
- Data science, Machine Learning, Artificial Intelligence... Which term should we use?
- Different applications of AI

2 How does Machine Learning work?

3 Application of AI in health

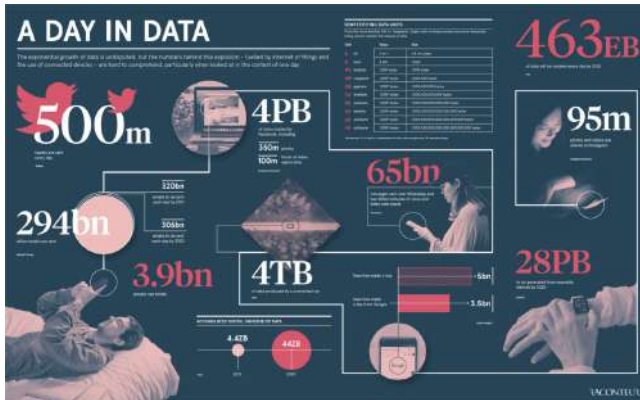
- Radiology
 - Chest X-ray
 - Liver lesion segmentation
- Genomics
 - Gene network
 - Toxicogenetics
- Medical

4 Perspective and issues

- Limitations of data projects
- Data Project Organization
- Unveiling the mystery of Deep Learning

Where are the data?

Everywhere!



Good point: impossible to do a data project without data

But how much data exactly?

Wooclap: *What is the average quantity of data created per person per day in 2020?*

Some scale (on average):

- An office document (Word, PowerPoint...): 321 KB
- 1 picture with a smartphone : 10 MB
- Recording a 3 hour zoom meeting: 1 GB

But how much data exactly?

Wooclap: *What is the average quantity of data created per person per day in 2020?*

Some scale (on average):

- An office document (Word, PowerPoint...): 321 KB
- 1 picture with a smartphone : 10 MB
- Recording a 3 hour zoom meeting: 1 GB

The quantity of data produced each day per person in 2020 is 2GB which is equivalent to

6h of zoom meeting recording or 200 pictures or 6.000 Office documents

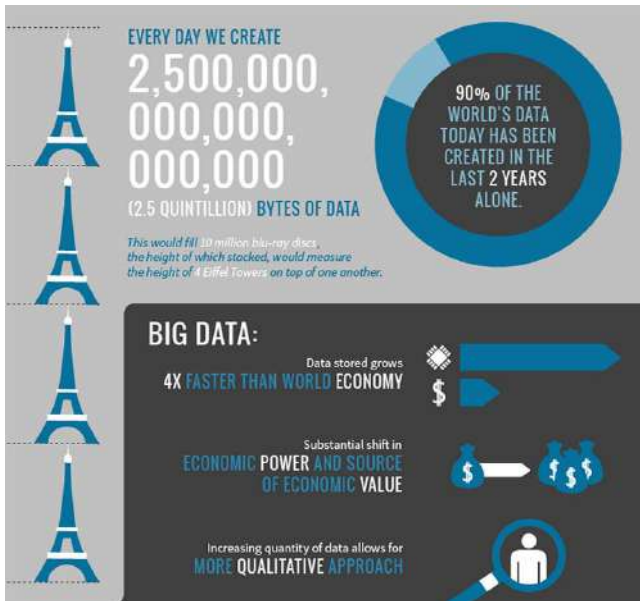
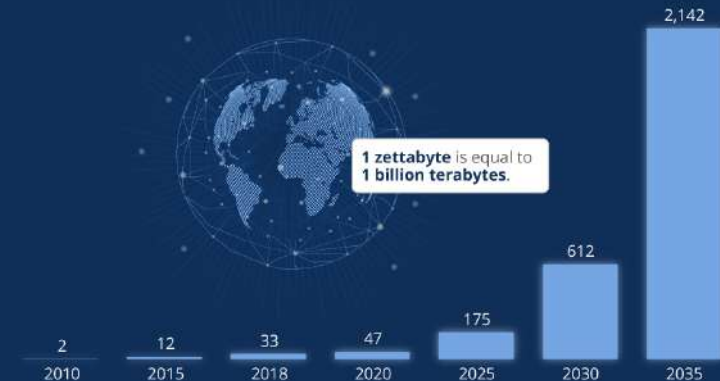


Figure: OECD, 2019

Big picture on data growth

Global Data Creation is About to Explode

Actual and forecast amount of data created worldwide 2010-2035 (in zettabytes)

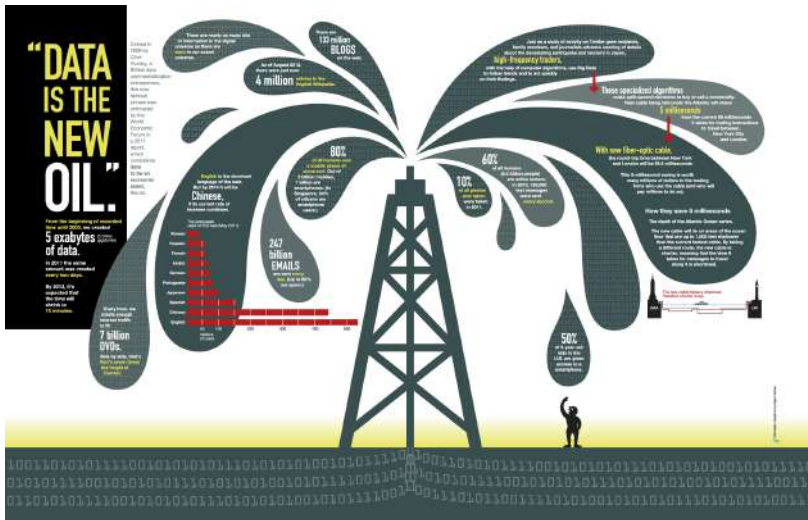


@StatistaCharts

Source: Statista Digital Economy Compass 2019

statista

Data is power



1 Introduction to AI

- No data project without data
- Data science, Machine Learning, Artificial Intelligence... Which term should we use?
- Different applications of AI

2 How does Machine Learning work?

3 Application of AI in health

- Radiology
 - Chest X-ray
 - Liver lesion segmentation
- Genomics
 - Gene network
 - Toxicogenetics
- Medical

4 Perspective and issues

- Limitations of data projects
- Data Project Organization
- Unveiling the mystery of Deep Learning

Data Terminology



When you're fundraising, it's AI / When you're hiring, it's ML / When you're implementing, it's linear regression / When you're debugging, it's printf()

Baron Schwartz, Twitter, Nov 2017

Wooclap: *Assign to each term its corresponding definition.*

- Data Management
 - Business Intelligence
 - Statistics
 - Data science
 - Big data
 - Machine learning
 - Artificial Intelligence
 - Deep Learning
- is the study of the collection, analysis, interpretation, presentation and organization of data.
 - comprises the strategies and technologies used by enterprises for the data analysis of business information.
 - is part of a broader family of machine learning methods based on learning data representations, as opposed to task-specific algorithms.
 - is the study of the generalizable extraction of knowledge from data.
 - is an all-encompassing term for any collection of data sets so large and complex that it becomes difficult to process using traditional data processing applications.
 - comprises all disciplines related to handling data as a valuable resource.
 - is the subfield of computer science that gives computers the ability to learn without being explicitly programmed.
 - aims at designing and studying *devices* that perceive its environment and take actions that maximize its chance of success at some goal.

Data terminology

- **Data Management** comprises all disciplines related to handling data as a valuable resource.
- **Business Intelligence** comprises the strategies and technologies used by enterprises for the data analysis of business information.
- **Statistics** is the study of the collection, analysis, interpretation, presentation and organization of data.
- **Data science** is the study of the generalizable extraction of knowledge from data.
- **Big data** is an all-encompassing term for any collection of data sets so large and complex that it becomes difficult to process using traditional data processing applications.
- **Machine learning** is the subfield of computer science that gives computers the ability to learn without being explicitly programmed.
- **Artificial Intelligence** aims at designing and studying *devices* that perceive its environment and take actions that maximize its chance of success at some goal.
- **Deep Learning** is part of a broader family of machine learning methods based on learning data representations, as opposed to task-specific algorithms.

Artificial Intelligence - an old buzzword (Dartmouth conference)

On September 1955, a project was proposed by McCarthy, Marvin Minsky, Nathaniel Rochester and Claude Shannon introducing formally for the first time the term "Artificial Intelligence".

The study is to proceed on the basis of the conjecture that every aspect of learning or any other feature of intelligence can in principle be so precisely described that a machine can be made to simulate it. An attempt will be made to find how to make machines use language, form abstractions and concepts, solve kinds of problems now reserved for humans, and improve themselves.

Proposal for Dartmouth conference on AI (1956)

Misconception of AI

AI is about electronic device able to **mimic** human thinking:

- Artificial **intelligence**
- One famous class of AI algorithms are called **neural networks**.
- **Android** are close to humans in shape so they must think like humans.

Most AI algorithms do **not** aim at **reproducing human reasoning**.

Artificial intelligence is the science of making machines do things that would require intelligence if done by men

Marvin Minsky (1968)



2001: A Space Odyssey

Artificial Intelligence is not human intelligence

What often happens is that an engineer has an idea of how the brain works (in his opinion) and then designs a machine that behaves that way. This new machine may in fact work very well. But, I must warn you that that does not tell us anything about how the brain actually works, nor is it necessary to ever really know that, in order to make a computer very capable. It is not necessary to understand the way birds flap their wings and how the feathers are designed in order to make a flying machine [...] **It is therefore not necessary to imitate the behavior of Nature in detail in order to engineer a device which can in many respects surpass Nature's abilities.**

Richard Feynman (1999)

AI technology - Autonomous cars

- Originates from 1920 (NY)
- First use of neural networks to control autonomous cars (1989)
- Four US states allow self-driving cars (2013)
- First known fatal accident (May 2016)
- Singapore launched the first self-driving taxi service (Aug. 2016)
- A Arizona pedestrian was killed by an Uber self-driving car (March 2018).



AI technology - virtual assistant / chatbot

- Voice recognition tool "Harpy" masters about 1000 words (1970s, CMU, US Defense).
- System capable of analyzing entire word sequences (1980).
- Siri was the first modern digital virtual assistant installed on a smartphone (2011).
- Watson won the TV show Jeopardy! (2011)



Different uses of AI



Different uses of AI



Different uses of AI

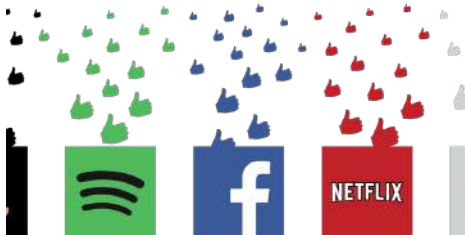


The study found that Google Home performed the best, recognizing 98 per cent of topics accurately and providing advice that matched with Red Cross first aid guidelines 56 per cent of the time.

Alexa recognized 92 per cent of topics, and gave appropriate advice 19 per cent of the time.

The responses from Siri and Cortana were so low that researchers determined that they couldn't analyze them.

Different uses of AI



Different uses of AI



NEWS • 10 OCTOBER 2019

Google AI beats top human players at strategy game *StarCraft II*

DeepMind's AlphaStar beat all but the very best humans at the fast-paced sci-fi video game.

Different uses of AI



AI artwork sells for \$432,500 — nearly 45 times its high estimate — as Christie's becomes the first auction house to offer a work of art created by an algorithm

Different uses of AI



824 views | Jan 18, 2020, 00:00am

A Look Inside Augmented Analytics And Its Business Value In 2020

Different uses of AI



1 Introduction to AI

- No data project without data
- Data science, Machine Learning, Artificial Intelligence... Which term should we use?
- Different applications of AI

2 How does Machine Learning work?

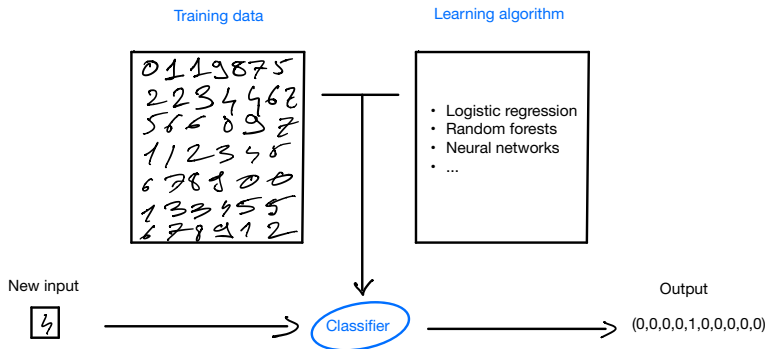
3 Application of AI in health

- Radiology
 - Chest X-ray
 - Liver lesion segmentation
- Genomics
 - Gene network
 - Toxicogenetics
- Medical

4 Perspective and issues

- Limitations of data projects
- Data Project Organization
- Unveiling the mystery of Deep Learning

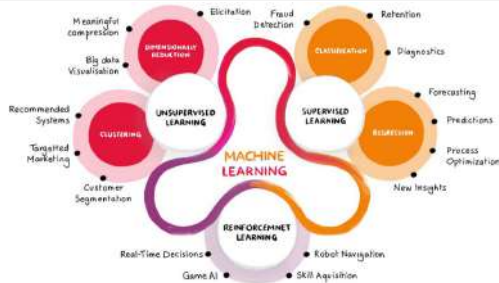
Supervised learning



A definition by Tom Mitchell (<http://www.cs.cmu.edu/~tom/>)

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.

Three Kinds of Learning



Unsupervised Learning

- **Task:** Clustering/DR
- **Performance:** Quality
- **Experience:** Raw dataset (No Ground Truth)

Supervised Learning

- **Task:** Prediction
- **Performance:** Average error
- **Experience:** Predictions (Ground Truth)

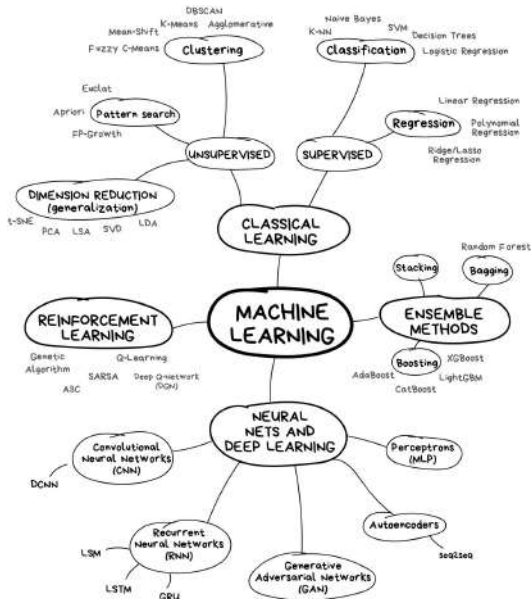
Reinforcement Learning

- **Task:** Action
- **Performance:** Total reward
- **Experience:** Reward from env. (Interact. with env.)

- **Timing:** Offline/Batch (learning from past data) vs Online (continuous learning)

Figure Source: BCG

Algorithms



Difficulties related to (Big) data

- The prediction must be **accurate**: difficult for some tasks like image classification, video captioning...
- Predictions must be **fast**: online recommendation should not take minutes.
- Data must be **stored** and **easily accessible**.
- It may be difficult to **access all data simultaneously**. Data may come sequentially.
- Data must be **clean**.
- Data should be **relevant**.



Wooclap: *How would you evaluate the performance of an ML algorithm aiming at diagnosing a patient?*

Various applications of AI in health

Wooclap: *Can you think of one application of AI in health/medicine?*

Wooclap: *Can you think of one application of AI in health/medicine?*

- **Diagnosis**

- From 1970s, an **expert system** to identify blood infections.
- Based on clinical data, omics data, medical imaging...

Various applications of AI in health

Wooclap: *Can you think of one application of AI in health/medicine?*

- **Diagnosis**

- From 1970s, an **expert system** to identify blood infections.
- Based on clinical data, omics data, medical imaging...

- **Predicting / Forecasting**

- Patient journey (predict missing appointments, ER waiting time)
- **Forecasting glycemic concentration** (D2P project, DiabeLoop)
- Classifying signals: Brainwaves / ECG (Cardiologists)

Various applications of AI in health

Wooclap: *Can you think of one application of AI in health/medicine?*

- **Diagnosis**

- From 1970s, an **expert system** to identify blood infections.
- Based on clinical data, omics data, medical imaging...

- **Predicting / Forecasting**

- Patient journey (predict missing appointments, ER waiting time)
- **Forecasting glycemic concentration** (D2P project, DiabeLoop)
- Classifying signals: Brainwaves / ECG (Cardiologists)

- **Drugs**

- Determining **protein 3D shapes** (AlphaFold, 2020)
- **Finding new drugs** (hit identification and de novo molecular design)

Various applications of AI in health

Wooclap: *Can you think of one application of AI in health/medicine?*

- **Diagnosis**

- From 1970s, an **expert system** to identify blood infections.
- Based on clinical data, omics data, medical imaging...

- **Predicting / Forecasting**

- Patient journey (predict missing appointments, ER waiting time)
- **Forecasting glycemic concentration** (D2P project, DiabeLoop)
- Classifying signals: Brainwaves / ECG (Cardiologists)

- **Drugs**

- Determining **protein 3D shapes** (AlphaFold, 2020)
- **Finding new drugs** (hit identification and de novo molecular design)

- **Robots**

- **Surgical** robots (robotic laparoscopes, sutures, cuts...)
- Robots for **rehabilitation, physical therapy**
- **Prosthetic arm control**

Various applications of AI in health

Wooclap: *Can you think of one application of AI in health/medicine?*

- **Diagnosis**

- From 1970s, an **expert system** to identify blood infections.
- Based on clinical data, omics data, medical imaging...

- **Predicting / Forecasting**

- Patient journey (predict missing appointments, ER waiting time)
- **Forecasting glycemic concentration** (D2P project, DiabeLoop)
- Classifying signals: Brainwaves / ECG (Cardiologists)

- **Drugs**

- Determining **protein 3D shapes** (AlphaFold, 2020)
- **Finding new drugs** (hit identification and de novo molecular design)

- **Robots**

- **Surgical** robots (robotic laparoscopes, sutures, cuts...)
- Robots for **rehabilitation, physical therapy**
- **Prosthetic arm control**

- **Automation**

Various applications of AI in health

Wooclap: *Can you think of one application of AI in health/medicine?*

- **Diagnosis**

- From 1970s, an **expert system** to identify blood infections.
- Based on clinical data, omics data, medical imaging...

- **Predicting / Forecasting**

- Patient journey (predict missing appointments, ER waiting time)
- **Forecasting glycemic concentration** (D2P project, DiabeLoop)
- Classifying signals: Brainwaves / ECG (Cardiologists)

- **Drugs**

- Determining **protein 3D shapes** (AlphaFold, 2020)
- **Finding new drugs** (hit identification and de novo molecular design)

- **Robots**

- **Surgical** robots (robotic laparoscopes, sutures, cuts...)
- Robots for **rehabilitation, physical therapy**
- **Prosthetic arm control**

- **Automation**

- **Treatment efficiency** (causality)

1 Introduction to AI

- No data project without data
- Data science, Machine Learning, Artificial Intelligence... Which term should we use?
- Different applications of AI

2 How does Machine Learning work?

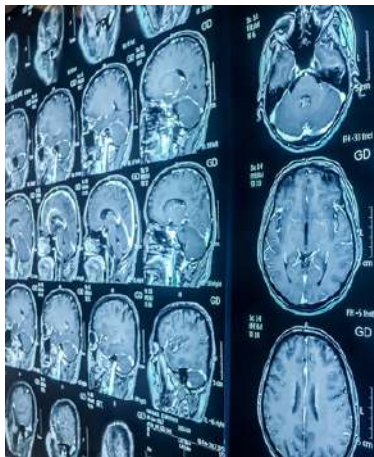
3 Application of AI in health

- Radiology
 - Chest X-ray
 - Liver lesion segmentation
- Genomics
 - Gene network
 - Toxicogenetics
- Medical

4 Perspective and issues

- Limitations of data projects
- Data Project Organization
- Unveiling the mystery of Deep Learning

Can Deep Learning be useful in radiology?



Deep Learning requires large annotated data sets (typically several millions of images, as in ImageNet)

Chest pathologies [1]

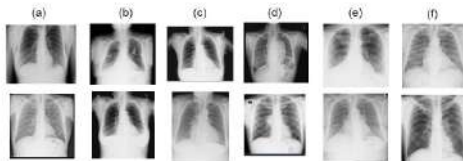


Fig. 1. Chest X-rays categories examples. (a) healthy, (b) left or right effusion, (c) enlarged heart (cardiomegaly), (d) enlarged mediastinum, (e) left or right consolidation, (f) multiple pathologies: enlarged heart, mediastinum, left and right effusion and left or right consolidation (Source: Diagnostic Imaging Department, Sheba Medical Center, Tel Hashomer, Israel)

637 X-ray images (6 chest pathologies):

- Right Pleural Effusion (73 images)
- Left Pleural Effusion (74 images)
- Right Consolidation (58 images)
- Left Consolidation (45 images)
- Cardiomegaly (154 images)
- Abnormal Me-diastinum (145 images)

Algorithm

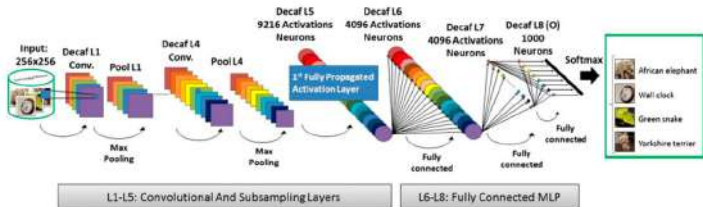


Figure 2. A schematic illustration of Donahue et al. (2013) CNN architecture and training process.

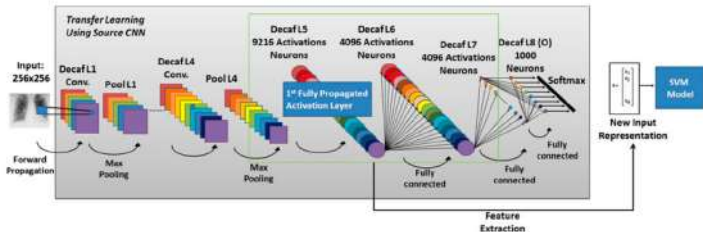
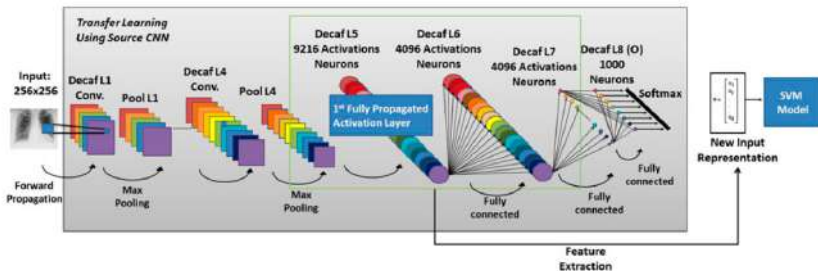


Figure 3. Feature extraction from Donahue et al. (2013) pre-trained CNN.

Convolutional Neural Network



Use SVM to predict a class based on the following features:

- Deep Learning features of different layers,
- Bag-of-Visual-Words (BoVW), particularly useful to categorize X-rays on organ level (ImageClefcompetitions, <http://www.imageclef.org>)
- GIST descriptors [4]

Which features are relevant?

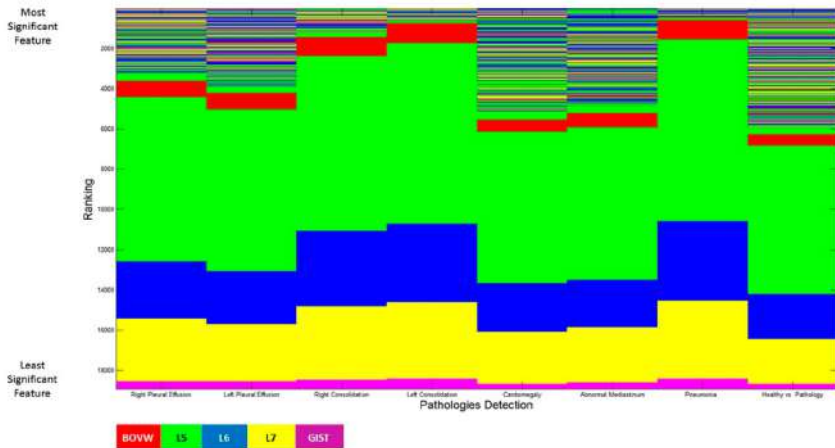


Figure 4. Features ranking, evaluated on the testing set, applied on the combined feature vector [BoVW, Deep features (*Decaf5*, *Decaf6* and *Decaf7*), Gist]. Ranking is performed on all examined identification cases. We use the following abbreviation: Ln for *Decafn*.

Table 1. AUC accuracy metric classification performance.

Descriptor	Right pleural effusion	Left pleural effusion	Right consolidation	Left consolidation
GIST	0.85	0.79	0.77	0.41
BoVW	0.89	0.87	0.78	0.65
L5	0.91	0.81	0.80	0.75
L6	0.91	0.82	0.85	0.76
L7	0.90	0.79	0.75	0.79
L5 + L6 + GIST	0.92	0.82	0.83	0.68
L5+L6+L7	0.92	0.82	0.83	0.78
FS (5000)	0.93	0.82	0.84	0.78

Cardiomegaly	Abnormal mediastinum	Healthy vs. pathology
0.96	0.73	0.88
0.94	0.74	0.85
0.95	0.79	0.90
0.94	0.80	0.90
0.93	0.79	0.89
0.96	0.79	0.91
0.94	0.80	0.91
0.95	0.80	0.92

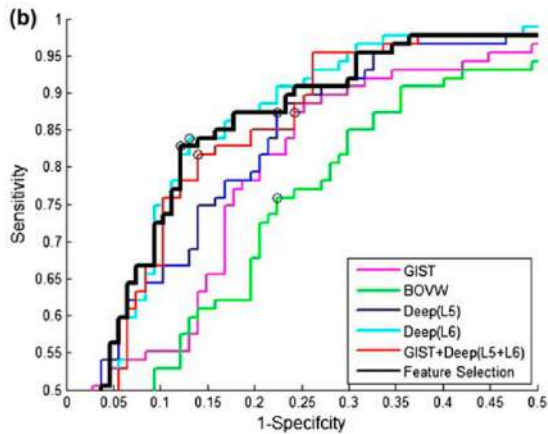


Figure: Healthy vs. Pathology ROC;

1 Introduction to AI

- No data project without data
- Data science, Machine Learning, Artificial Intelligence... Which term should we use?
- Different applications of AI

2 How does Machine Learning work?

3 Application of AI in health

- Radiology
 - Chest X-ray
 - Liver lesion segmentation
- Genomics
 - Gene network
 - Toxicogenetics
- Medical

4 Perspective and issues

- Limitations of data projects
- Data Project Organization
- Unveiling the mystery of Deep Learning

Liver lesion segmentation [2]

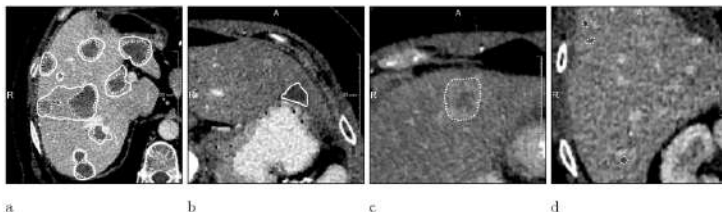


Figure 3. MTRA (*dashed*) vs. LiTS (*solid*) annotations. (a) Case with low dice/correspondence (b) Case where a LiTS reference tumor was missed (c) Case where MTRA found a lesion in a case with no tumors according to LiTS reference (d) Case where small additional tumors were found by the MTRA.

Data:

- 131 abdominal CT scans
- The CT scans come with reference annotations of the liver and tumors done by trained radiologists (LiTS reference)
- Another set of annotations given by a medical-technical radiology assistant (MTRA)

Fully Convolutional Neural Networks

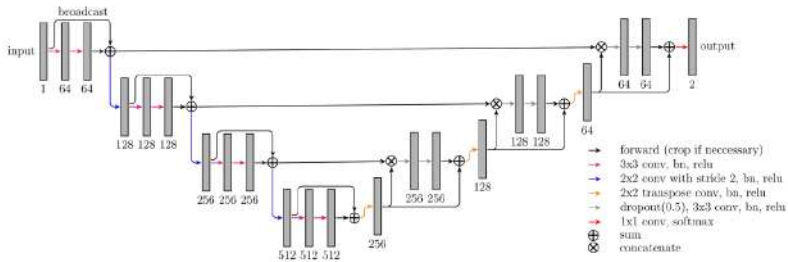


Figure 1. Overview of the neural network architecture. The numbers denote the feature map count.

Based on U-net [6] designed specifically for Biomedical Image Segmentation.

Neural Network Output compared to annotated data

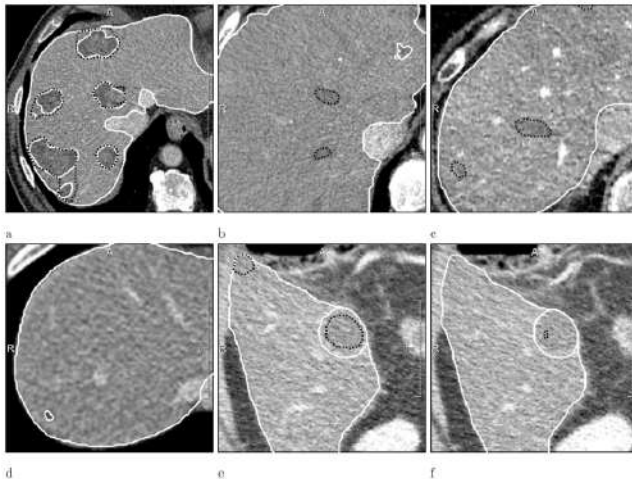


Figure 5. Neural network (*black*) compared with the LiTS (*white*) annotations. (a) Case with 0.85 dice/ case (b,c) Cases with 19 and 16 FPs (d) Case where a small tumor was not detected (e,f) Case where tumor segmentation strongly differed on consecutive slices.

Comparison between Human and Computer recognition performances

	Recall	Recall ≥ 10 mm	FP per case	Dice per case	Dice per correspondence	Merge error	Split error
Human vs. Human							
MTRA (LiTS)	0.92	0.94	2.6	0.70 ± 0.27	0.72 ± 0.11	11	5
LiTS (MTRA)	0.62	0.85	0.3	0.70 ± 0.27	0.72 ± 0.11	5	12
Computer vs. Human							
FCN (MTRA)	0.47	0.75	4.7	0.53 ± 0.37	0.72 ± 0.11	7	13
FCN (LiTS)	0.72	0.86	4.6	0.51 ± 0.37	0.65 ± 0.16	12	14
FCN + RF (LiTS)	0.63	0.77	0.7	0.58 ± 0.36	0.69 ± 0.18	11	10

Table 1. Mean metric values for human vs. human and computer vs. human comparisons. The parentheses denote the dataset used as a reference for the computation of evaluation metrics.

1 Introduction to AI

- No data project without data
- Data science, Machine Learning, Artificial Intelligence... Which term should we use?
- Different applications of AI

2 How does Machine Learning work?

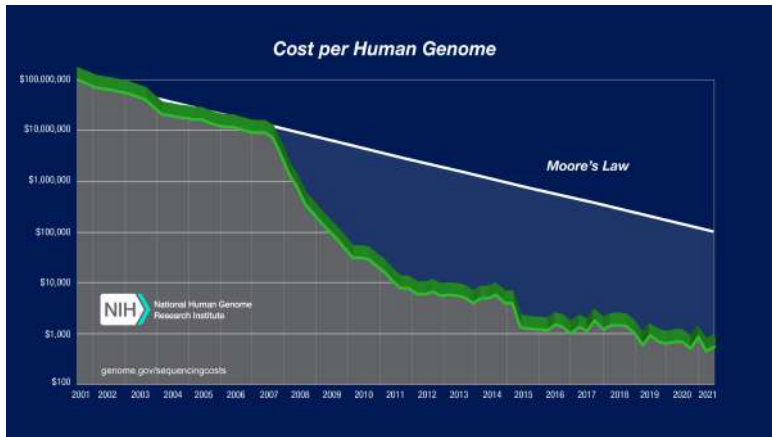
3 Application of AI in health

- Radiology
 - Chest X-ray
 - Liver lesion segmentation
- Genomics
 - Gene network
 - Toxicogenetics
- Medical

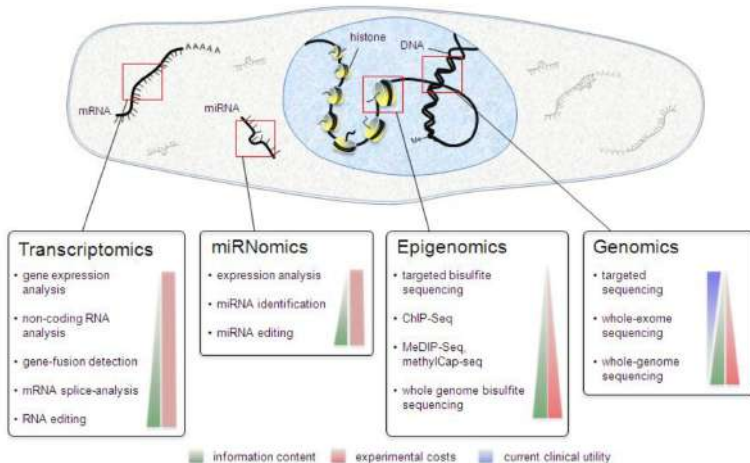
4 Perspective and issues

- Limitations of data projects
- Data Project Organization
- Unveiling the mystery of Deep Learning

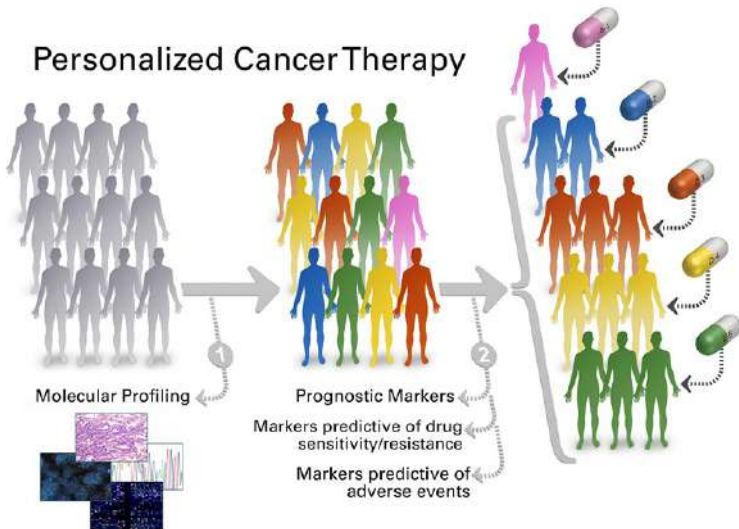
Cost per genome



Different fields in Omics [3]



Personalized Cancer Therapy

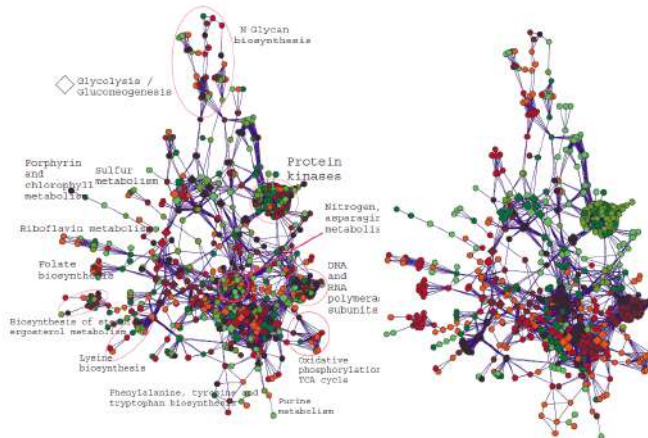


<https://pct.mdanderson.org>

Need for robustness - incorporating knowledge in the model

[5] Effect of low irradiation on *Saccharomyces cerevisiae* strains: 6 irradiated groups vs 11 non-irradiated groups.

Aim: Detect the irradiation looking only at the transcriptional changes.



- 1 Introduction to AI
 - No data project without data
 - Data science, Machine Learning, Artificial Intelligence... Which term should we use?
 - Different applications of AI
- 2 How does Machine Learning work?
- 3 Application of AI in health
 - Radiology
 - Chest X-ray
 - Liver lesion segmentation
 - Genomics
 - Gene network
 - Toxicogenetics
 - Medical
- 4 Perspective and issues
 - Limitations of data projects
 - Data Project Organization
 - Unveiling the mystery of Deep Learning

What is toxicogenetics ?



- Different responses to drugs or environmental chemicals according to genotypes
- Aim : provide a personalized treatment for patients

The Dream Project |

www.the-dream-project.org

Home Challenges Team Ranking Conferences Discussion Literature Reverse Engineering Code News Contact us [Login / Register](#)

DREAM is a Dialogue for Reverse Engineering Assessments and Methods. The main objective is to catalyze the interaction between experiment and theory in the area of cellular network inference and quantitative model building in systems biology.

Dialogue for Reverse Engineering Assessments and Methods

DREAM8 Challenges Are Open

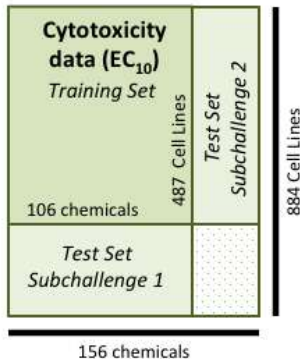
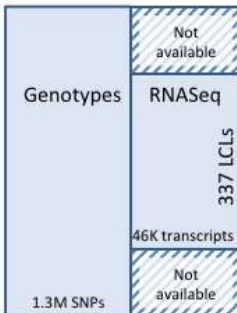
We are pleased to announce that the DREAM8 Challenges are now open for participation. During the "Challenge season" spanning from June 10 to September 15, 2013, Sage Bionetworks and DREAM will run the following three Challenges:

1. [HPN-DREAM Breast Cancer Network Inference Challenge](#) – Infer the signaling networks in breast cancer cell lines
2. [NIEHS-NCATS-UNC DREAM Toxicogenetics Challenge](#) – Predict individual response to environmental and pharmaceutical chemicals
3. [The Whole-Cell Parameter Estimation DREAM Challenge](#) – Infer the kinetic parameters underlying biological processes in whole cell models

To sign up for a Challenge, and access the data sets and descriptions of the DREAM8 Challenges, please go to <https://www.synapse.org/#/Challenges/DREAM8>.

DREAM8 Toxicogenetics challenge

Toxicogenetics Challenge Data



The two G projects

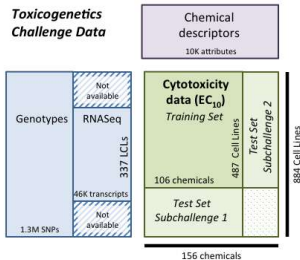
The 1000 Genomes Project

- Decreasing of sequencing cost allow to sequence a large number of people
- Find most genetic variants that have frequencies of 1%.
- Data are publicly available at <http://www.1000genomes.org/data>

The Geuvadis Project

- European project in high-throughput sequencing
- RNA sequencing of 465 cell lines belonging to the 1000 Genomes Project
- So RNA data for non-european people are missing in the challenge
- Data publicly available at <http://www.geuvadis.org>

Our strategy



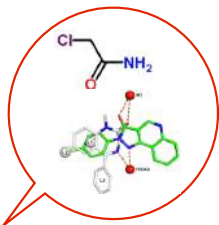
- Need to integrate large-scale, heterogeneous data with missing information → **kernels**
- Share information across chemicals → **multitask learning**

So we developed a **joint** model on both **patients** and **drugs**.

Kernel Trick

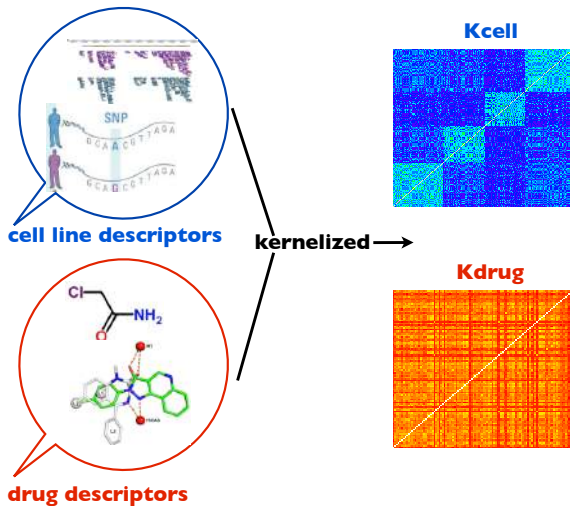


cell line descriptors

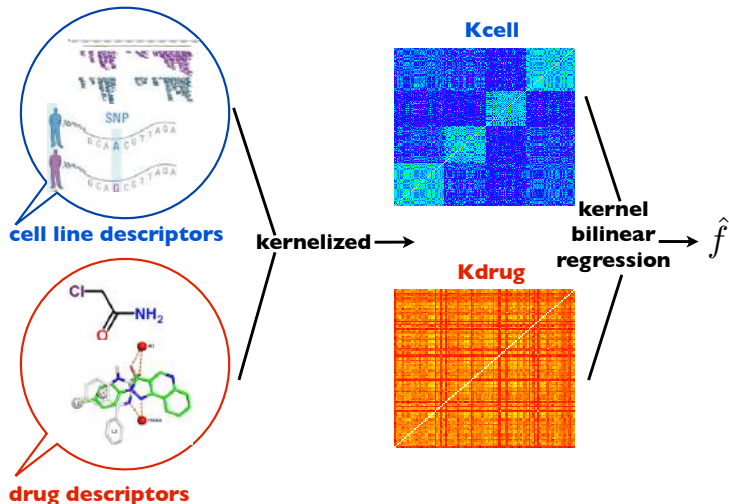


drug descriptors

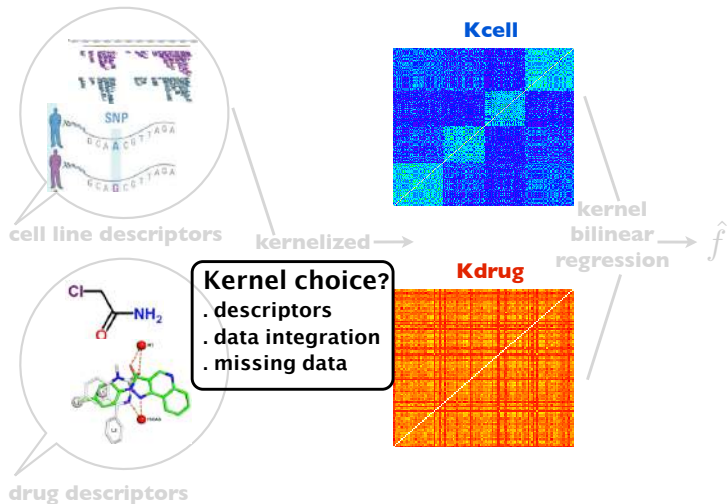
Kernel Trick



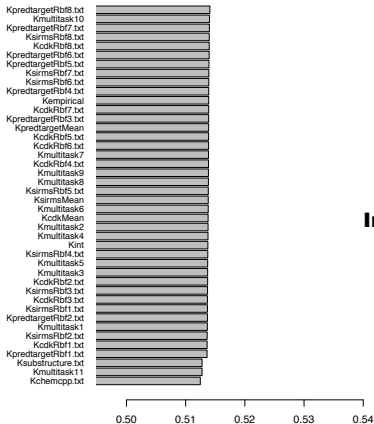
Kernel Trick



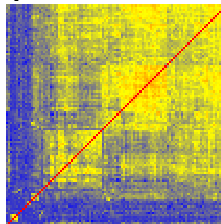
Kernel Trick



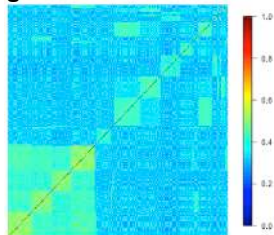
Mean CI for chemicals kernels



Empirical kernel on drugs



Integrated kernel on cell lines



1 Introduction to AI

- No data project without data
- Data science, Machine Learning, Artificial Intelligence... Which term should we use?
- Different applications of AI

2 How does Machine Learning work?

3 Application of AI in health

- Radiology
 - Chest X-ray
 - Liver lesion segmentation
- Genomics
 - Gene network
 - Toxicogenetics
- Medical

4 Perspective and issues

- Limitations of data projects
- Data Project Organization
- Unveiling the mystery of Deep Learning

15000 patients / 250 variables / 11 hospitals, from 2011

- 4000 new patients / year

-----	Center	Accident	Age	Sex	Weight	Height	BMI	BP	SBP
1	Beaujon	Fall	54	m	85	NR	NR	180	110
2	Lille	Other	33	m	80	1.8	24.69	130	62
3	Pitie Salpetriere	Gun	26	m	NR	NR	NR	131	62
4	Beaujon	AVP moto	63	m	80	1.8	24.69	145	89
6	Pitie Salpetriere	AVP bicycle	33	m	75	NR	NR	104	86
7	Pitie Salpetriere	AVP pedestrian	30	w	NR	NR	NR	107	66
9	HEGP	White weapon	16	m	98	1.92	26.58	118	54
10	Toulon	White weapon	20	m	NR	NR	NR	124	73

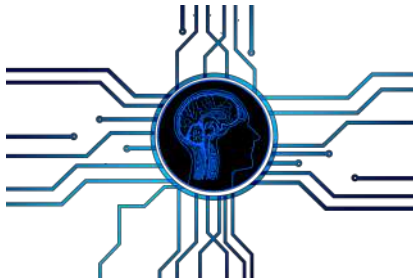
---	SpO2	Temperature	Lactates	Hb	Glasgow	Transfusion
1	97	35.6	<NA>	12.7	12	yes
2	100	36.5	4.8	11.1	15	no
3	100	36	3.9	11.4	3	no
4	100	36.7	1.66	13	15	yes
6	100	36	NM	14.4	15	no
7	100	36.6	NM	14.3	15	yes
9	100	37.5	13	15.9	15	yes
10	100	36.9	NM	13.7	15	no

http://www.traumabase.eu/fr_FR

Major trauma: any injury that endangers the life or the functional integrity of a person. Road traffic accidents, interpersonal violence, self-harm, falls, etc

→ hemorrhage and traumatic brain injury.

Patient prognosis can be improved: **standardized and reproducible procedures** but **personalized** for the patient and the trauma system.



⇒ Can AI help?

-----	Center	Accident	Age	Sex	Weight	Height	BMI	BP	SBP
1	Beaujon	Fall	54	m	85	NR	NR	180	110
2	Lille	Other	33	m	80	1.8	24.69	130	62
3	Pitie Salpetriere	Gun	26	m	NR	NR	NR	131	62
4	Beaujon	AVP moto	63	m	80	1.8	24.69	145	89
6	Pitie Salpetriere	AVP bicycle	33	m	75	NR	NR	104	86
7	Pitie Salpetriere	AVP pedestrian	30	w	NR	NR	NR	107	66
9	HEGP	White weapon	16	m	98	1.92	26.58	118	54
10	Toulon	White weapon	20	m	NR	NR	NR	124	73

---	SpO2	Temperature	Lactates	Hb	Glasgow	Transfusion
1	97	35.6	<NA>	12.7	12	yes
2	100	36.5	4.8	11.1	15	no
3	100	36	3.9	11.4	3	no
4	100	36.7	1.66	13	15	yes
6	100	36	NM	14.4	15	no
7	100	36.6	NM	14.3	15	yes
9	100	37.5	13	15.9	15	yes
10	100	36.9	NM	13.7	15	no

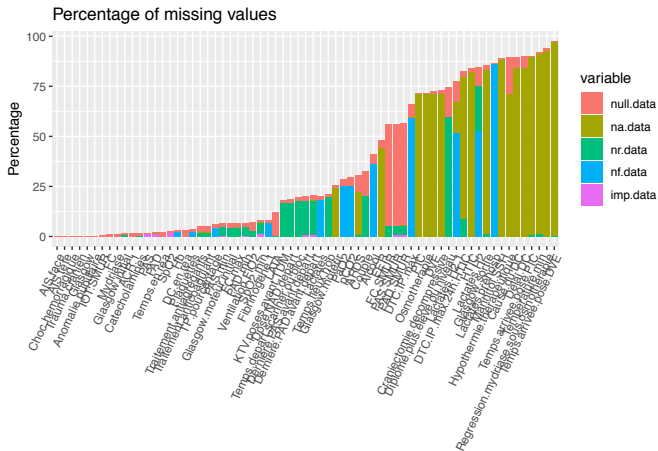
⇒ **Predict** whether to start a blood transfusion, to administer fresh frozen plasma, etc...

-----	Center	Accident	Age	Sex	Weight	Height	BMI	BP	SBP
1	Beaujon	Fall	54	m	85	NR	NR	180	110
2	Lille	Other	33	m	80	1.8	24.69	130	62
3	Pitie Salpetriere	Gun	26	m	NR	NR	NR	131	62
4	Beaujon	AVP moto	63	m	80	1.8	24.69	145	89
6	Pitie Salpetriere	AVP bicycle	33	m	75	NR	NR	104	86
7	Pitie Salpetriere	AVP pedestrian	30	w	NR	NR	NR	107	66
9	HEGP	White weapon	16	m	98	1.92	26.58	118	54
10	Toulon	White weapon	20	m	NR	NR	NR	124	73

---	SpO2	Temperature	Lactates	Hb	Glasgow	Transfusion
1	97	35.6	<NA>	12.7	12	yes
2	100	36.5	4.8	11.1	15	no
3	100	36	3.9	11.4	3	no
4	100	36.7	1.66	13	15	yes
6	100	36	NM	14.4	15	no
7	100	36.6	NM	14.3	15	yes
9	100	37.5	13	15.9	15	yes
10	100	36.9	NM	13.7	15	no

⇒ **Presence of missing data:** Not Recorded, Made, Applicable, etc.

Missing values



How to solve this problem?

- 1 Delete all missing values → bad idea!

How to solve this problem?

- ❶ Delete all missing values → bad idea!
- ❷ Impute data with your favorite imputation method
 - Replace all missing values by the mean/median/mode of the corresponding variable.
 - Good point: the mean/median/mode is unchanged!
 - Bad point 1: the variance of the imputed data is lower than reality
 - Bad point 2: structure of dependence between variable is destroyed.
 - Multiple imputation.

How to solve this problem?

- ❶ Delete all missing values → bad idea!
- ❷ Impute data with your favorite imputation method
 - Replace all missing values by the mean/median/mode of the corresponding variable.
 - Good point: the mean/median/mode is unchanged!
 - Bad point 1: the variance of the imputed data is lower than reality
 - Bad point 2: structure of dependence between variable is destroyed.
 - Multiple imputation.
- ❸ Design methods that handle missing values.

1 Introduction to AI

- No data project without data
- Data science, Machine Learning, Artificial Intelligence... Which term should we use?
- Different applications of AI

2 How does Machine Learning work?

3 Application of AI in health

- Radiology
 - Chest X-ray
 - Liver lesion segmentation
- Genomics
 - Gene network
 - Toxicogenetics
- Medical

4 Perspective and issues

- Limitations of data projects
- Data Project Organization
- Unveiling the mystery of Deep Learning

Cost of storing data

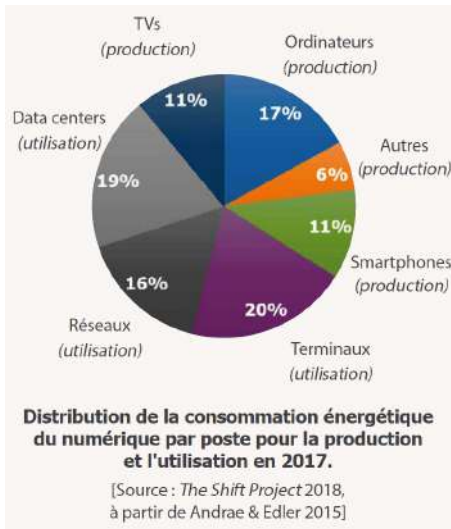
Cost of storing data

- Money: 300.000 US dollars in Google Cloud to store 1 petabyte during one year.

Cost of storing data

- Money: 300.000 US dollars in Google Cloud to store 1 petabyte during one year.
- Data centers and environment
 - 2% of the total electricity consumption in the US.
 - 626 billion liters of water.
 - 2% of total global greenhouse emissions.





The Shift Project <https://theshiftproject.org/lean-ict/>

Car accidents

Female drivers and right front passengers are approximately

17 percent more likely

to be killed

in a car crash than a male occupant of the same age.

Any seatbelt-wearing female vehicle occupant has

73 percent greater odds of being

seriously injured

in a frontal car crash than the odds of a seatbelt-wearing male occupant being injured in the same kind and severity of crash.

Sources: NHTSA and the Journal Traffic Injury Prevention

Analysis of crash and injury data compiled from the National Automotive Sampling System Crashworthiness Data System for the years 1998 to 2015.

Bias in crash test



<https://www.consumerreports.org/car-safety/crash-test-bias-how-male-focused-testing-puts-female-drivers-at-risk/>

What is COMPAS?

Correctional Offender Management Profiling for Alternative Sanction used in US justice courts to predict the reoffending probability.



What is COMPAS?

Correctional Offender Management Profiling for Alternative Sanction used in US justice courts to predict the reoffending probability.

Assessing the fairness of COMPAS

(A) Calibration

Given a score, the percentage of black people who reoffend is the same as the percentage of white people who reoffend.

(B) Parity - False Positive rate

The false positive rates (probability of being classified at risk while being not at risk) are the same for the group of black people and white people.

(C) Parity - False Negative rate

The false negative rates (probability of being classified not at risk while being at risk) are the same for the group of black people and white people.

- (A) According to Northpoint, **COMPAS is calibrated**.
Among defendants who scored a seven on the COMPAS scale, 60 percent of white defendants reoffended, which is nearly identical to the 61 percent of black defendants who reoffended.
- (B) According to ProPublica, **COMPAS does not satisfy parity** for false Positive rate.
Among defendants who ultimately did not reoffend, blacks were more than twice as likely as whites to be classified as medium or high risk (42 percent vs. 22 percent).
- (C) Parity - False Negative rate

Theorem: assume that reoffending cannot be **exactly** predicted via the input features (life is always a bit random), then there is no algorithm that satisfies (A), (B), (C).

Washington Post : A computer program used for bail and sentencing decisions was labeled biased against blacks. It's actually not that clear.

- Car crash tests lead to unfair vehicles.
→ **Debias the data** : collect more, better quality, better representativity.
- Correctional Offender Management Profiling for Alternative Sanctions (Compas) used in the US.
→ **Debias the algorithm** : twist predictions to annihilate one bias.
- Social Credit System / DeepNude
→ **Impact on society** : do we want these algorithms in our life?

1 Introduction to AI

- No data project without data
- Data science, Machine Learning, Artificial Intelligence... Which term should we use?
- Different applications of AI

2 How does Machine Learning work?

3 Application of AI in health

- Radiology
 - Chest X-ray
 - Liver lesion segmentation
- Genomics
 - Gene network
 - Toxicogenetics
- Medical

4 Perspective and issues

- Limitations of data projects
- Data Project Organization
- Unveiling the mystery of Deep Learning

Steps of a data project

Wooclap: *Order the following different steps of a data project.*

- Data collection
- Model evaluation
- Predictive modeling
- Continuous Optimization
- Solution Deployment
- Data Wrangling (gathering data in a usable format)
- Business understanding
- Testing / validation

You can also mention how each step interacts with the others.

Data Project Framework



Predictive Analytical Model Process Flow

Figure: <http://www.anovaanalytics.com/data-science-consulting/>

Data Science in 1 Slide

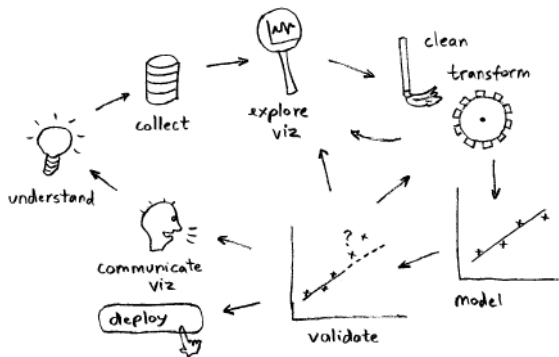
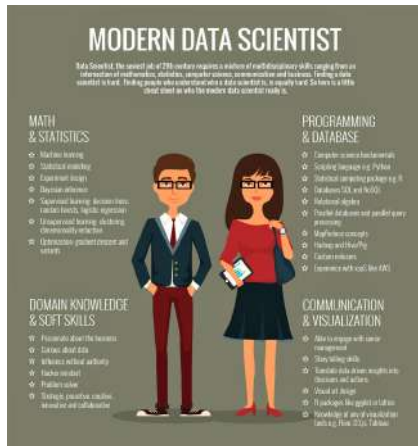


Figure: Source: Sz. Pafka

CRISP-DM

- CCross-Industry Standard Process for Data Mining (1999)
- Adapted by Szilard Pafka.

Data Scientists and Challenges



Data Scientist

- **Mix** of various skills.
- **Hard** to be an expert of everything!

Different occupations in a data project

Wooclap: Assign the following job names to the job descriptions below.

Data engineer, business analyst, statistician, data and analytics manager, data scientist, data architect, data analyst.



AS RARE AS UNICORNS

Role
Churn, manages and ingests Big Data

Languages
R, SAS, Python, Matlab, SQL, Hive, Pig, Spark

Mindset
Fastest data scientist

Skills & Talents
• Data kernel computing
• Predictive modeling
• Data mining and visualization
• Math, Stats, Machine Learning



HISTORIC LEADERS OF DATA

Role
Collects, analyzes and interprets qualitative as well as quantitative data with statistical theories and methods

Languages
R, SAS, SPSS, Matlab, Stata, Python, Perl, Hive, Pig, Spark, SQL

Mindset
Logical and methodical data genius

Skills & Talents
• Statistical theories & methodology
• Data mining & machine learning
• Distributed Computing (Hadoop)
• Database systems (SQL and NO SQL, Hetero)
• Cloud tools



DATA DETECTIVE

Role
Collects, processes and performs statistical data analysis

Languages
R, Python, HIVE, Java/Scala, C++/H, SAS

Mindset
Visual data junkie with high "figural-verbal" quotient

Skills & Talents
• Spreadsheet tools (e.g., Excel)
• Database systems (SQL and NO SQL, based)
• Data visualization & methodology
• Math, Stats, Machine Learning



THE CONTEMPORARY DATA MODELLER

Role
Creates blueprints for data management systems that integrate, combine, connect and maintain data sources

Languages
SQL, HIVE, Hive, Pig, Spark

Mindset
Inspiring role with a flair for data architecture design patterns

Skills & Talents
• Data warehousing solutions
• In-depth knowledge of database architecture
• Extraction Transformation and Load (ETL) spreadsheet and BI tools
• Data modeling
• Systems development



SOFTWARE ENGINEERS BY TRADE

Role
Develops, constructs, tests and integrates architectures (such as databases and large-scale processing systems)

Languages
SQL, Hive, Pig, R, Matlab, SAS, SPSS, Python, Java, Ruby, C#, Perl

Mindset
All-purpose engineer

Skills & Talents
• Database systems (SQL and NO SQL, based)
• Data modeling & ETL tools
• Data APIs
• Data warehousing solutions



CHANGE AGENT

Role
Implements business processes as interactivity between business and IT

Languages
SQL

Mindset
Realizes project goals

Skills & Talents
• Business (e.g., MS Office)
• Data visualization tools (e.g., Tableau)
• Statistical learning and data mining
• Business intelligence understanding
• Data modeling



DATA SCIENCE TEAM LEADER

Role
Manages a team of analysts and data scientists

Languages
SQL, R, SAS, Python, Matlab, Java

Mindset
Data Whiz's Chief leader

Skills & Talents
• Database systems (SQL and NO SQL, based)
• Leadership & project management
• Interpersonal communication
• Data mining & predictive modeling

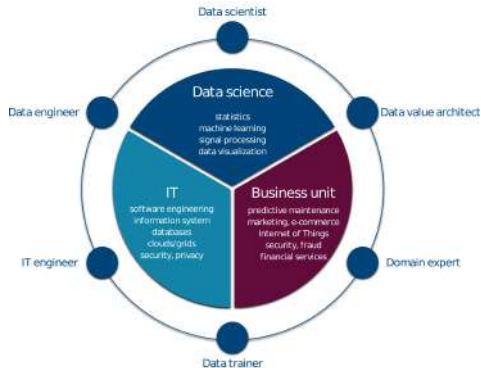
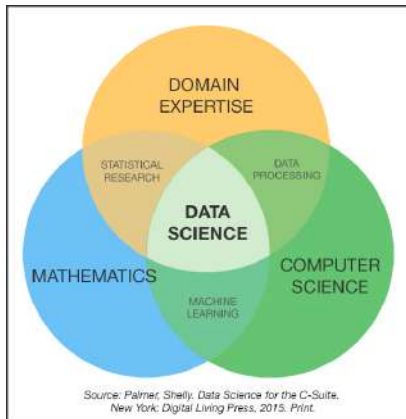
Different occupations in a data project



Several Profiles

- Several kind of problem / several kind of tools
- Much more variety than this...
- Importance of balanced **teams**.

Different training and different occupations



1 Introduction to AI

- No data project without data
- Data science, Machine Learning, Artificial Intelligence... Which term should we use?
- Different applications of AI

2 How does Machine Learning work?

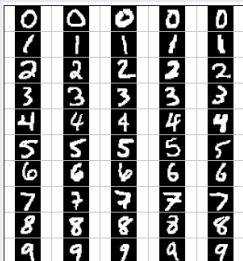
3 Application of AI in health

- Radiology
 - Chest X-ray
 - Liver lesion segmentation
- Genomics
 - Gene network
 - Toxicogenetics
- Medical

4 Perspective and issues

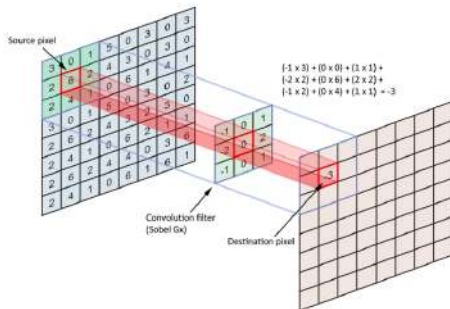
- Limitations of data projects
- Data Project Organization
- Unveiling the mystery of Deep Learning

Number Recognition



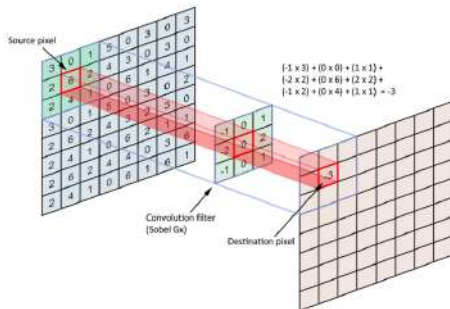
- Data: Annotated database of images (each image is represented by a vector of $28 \times 28 = 784$ pixel intensities)
- Input: Image
- Output: Corresponding number

Fundamental elements

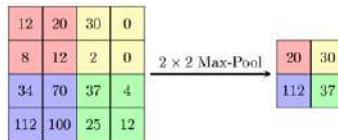


Convolution

Fundamental elements

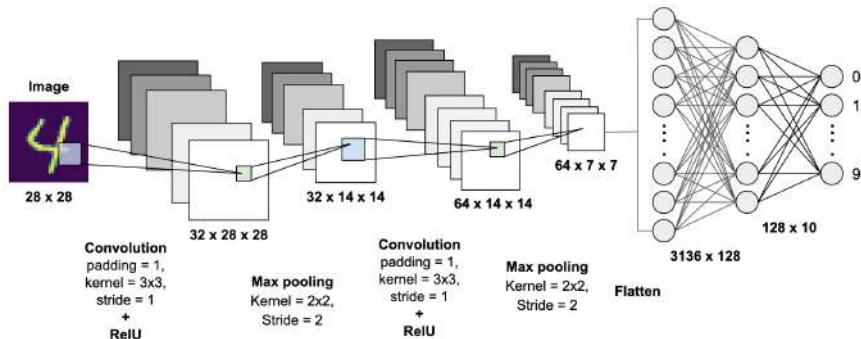


Convolution



Max-Pooling

Convolutional neural network



Results



The 82 patterns misclassified by LeNet5. Below each image is displayed the correct answer (left) and the prediction (right). These errors are mostly caused by genuinely ambiguous patterns, or by digits written in a style that are under represented in the training set.

Other generic applications of CNN



[Krizhevsky 2012]



[Ciresan et al. 2013]



[Faster R-CNN - Ren 2015]



[NVIDIA dev blog]

Far from terminator

- Stephen Hawking BBC, Dec 2 2014

The development of full artificial intelligence could spell the end of the human race. We cannot quite know what will happen if a machine exceeds our own intelligence, so we can't know if we'll be infinitely helped by it, or ignored by it and sidelined, or conceivably destroyed by it.



Take-home messages

- No data projects without data
 - A lot of data are available in the world
 - Difficulty of gathering the relevant ones and cleaning them (70% of the data project)
 - Environmental/Technical point of view: collecting/creating data is expensive for the planet (and the company)
- Different terms used in a data project
 - Keep in mind that AI has nothing to do with intelligence.
 - AI does not mimic human reasoning

Take-home messages

- No data projects without data
 - A lot of data are available in the world
 - Difficulty of gathering the relevant ones and cleaning them (70% of the data project)
 - Environmental/Technical point of view: collecting/creating data is expensive for the planet (and the company)
- Different terms used in a data project
 - Keep in mind that AI has nothing to do with intelligence.
 - AI does not mimic human reasoning
- How does Machine learning works?
 - Machine learning requires data to detect and learn patterns in the data.
 - Different tasks can be solved depending on the data (supervised, unsupervised, images, texts...)
 - Different tasks cannot be solved with ML notably if relevant information are not inside the collected data
 - Specific questions require specific data

Take-home message II

- Limitations of ML
 - **Data may be biased** because our world is, and data are nothing but a reflection of it.
 - Detecting and removing these biases is tricky but very important if individuals are impacted by the ML solution.
 - **ML may have trouble to adapt in temporal changes.** ML has a tendency to reproduce the past.

Take-home message II

- Limitations of ML
 - **Data may be biased** because our world is, and data are nothing but a reflection of it.
 - Detecting and removing these biases is tricky but very important if individuals are impacted by the ML solution.
 - **ML may have trouble to adapt in temporal changes.** ML has a tendency to reproduce the past.
- Data cycle and Data jobs
 - A data cycle is composed of iterations, **nothing is ever over.**
 - Business analysis is very important through the cycle
 - Many different actors are involved in a data project
 - **Good communication is required!**

Take-home message II

- Limitations of ML
 - **Data may be biased** because our world is, and data are nothing but a reflection of it.
 - Detecting and removing these biases is tricky but very important if individuals are impacted by the ML solution.
 - **ML may have trouble to adapt in temporal changes.** ML has a tendency to reproduce the past.
- Data cycle and Data jobs
 - A data cycle is composed of iterations, **nothing is ever over.**
 - Business analysis is very important through the cycle
 - Many different actors are involved in a data project
 - **Good communication is required!**
- Data Science is evolving constantly.
 - New opportunities appear
 - New challenges are detected
 - Need for adaptability



Thank you!

- [1] Y. Bar, I. Diamant, L. Wolf, S. Lieberman, E. Konen, and H. Greenspan. Chest pathology identification using deep feature selection with non-medical training. *Computer Methods in Biomechanics and Biomedical Engineering: Imaging & Visualization*, 6(3):259–263, 2018.
- [2] G. Chlebus, A. Schenk, J. H. Moltz, B. van Ginneken, H. K. Hahn, and H. Meine. Automatic liver tumor segmentation in ct with fully convolutional neural networks and object-based postprocessing. *Scientific reports*, 8(1):15497, 2018.
- [3] K. Frese, H. Katus, and B. Meder. Next-generation sequencing: from understanding biology to personalized medicine. *Biology*, 2(1):378–398, 2013.
- [4] A. Oliva and A. Torralba. Modeling the shape of the scene: A holistic representation of the spatial envelope. *International journal of computer vision*, 42(3):145–175, 2001.
- [5] F. Rapaport, A. Zinovyev, M. Dutreix, E. Barillot, and J.-P. Vert. Classification of microarray data using gene networks. *BMC bioinformatics*, 8(1):35, 2007.
- [6] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015.