



BỘ CÔNG THƯƠNG  
TRƯỜNG ĐẠI HỌC CÔNG THƯƠNG THÀNH PHỐ HỒ CHÍ MINH

# MACHINE LEARNING & DATA MINING

**GV. Vũ Đức Thịnh**



# Nội dung

## Chương 1. Giới thiệu chung

Chương 2. Thu thập và tiền xử lý dữ liệu

Chương 3. Hồi quy

Chương 4. Phân cụm

Chương 5. Phân cụm phân cấp

Chương 6. Học dựa trên láng giềng (KNN)

Chương 7. Cây quyết định và rừng ngẫu nhiên

Chương 8. Máy vectơ hỗ trợ (SVM)

Chương 9. Đánh giá hiệu quả của mô hình

Chương 10. Mạng nơron nhân tạo

Chương 11. Mô hình xác suất

# Tài liệu tham khảo

1. Tom Mitchell. Machine Learning. McGraw-Hill, 1997.
2. Ian Goodfellow, Yoshua Bengio, and Aaron Courville. Deep Learning. MIT press, 2016
3. Jiawei Han, Micheline Kamber, Jian Pei. Data Mining: Concepts and Techniques (3rd Edition). Morgan Kaufmann, 2011.
4. Trevor Hastie, Robert Tibshirani, Jerome Friedman. The Elements of Statistical Learning: Data Mining, Inference, and Prediction (12th Edition). Springer, 2017.

# Chương 1. Giới thiệu chung

## Machine Learning vs Data Mining

- Machine Learning  
(ML - Học máy)

To build computer systems  
that can improve themselves  
by learning from data.

(Xây dựng những hệ thống mà  
có khả năng tự cải thiện bản  
thân bằng cách học từ dữ liệu.)

- Some venues: NeurIPS,  
ICML, IJCAI, AAAI, ICLR,  
ACML, ECML

- Data Mining  
(DM - Khai phá dữ liệu)

To find new and useful  
knowledge from datasets.

(Tìm ra/Khai phá những tri thức  
mới và hữu dụng từ các tập dữ  
liệu lớn.)

- Some venues: KDD, PKDD,  
PAKDD, ICDM, CIKM



# Data

## Structured – relational (table-like)

	A	B	C	D	E	F	G
1	Country	Region	Population	Under15	Over60	Fertil	LifeExp
2	Zimbabwe	Africa	13724	40.24	5.68	3.64	54
3	Zambia	Africa	14075	46.73	3.95	5.77	55
4	Yemen	Eastern M	23852	40.72	4.54	4.35	64
5	Viet Nam	Western P	90796	22.87	9.32	1.79	75
6	Venezuela (Bo	Americas	29955	28.84	9.17	2.44	75
7	Vanuatu	Western P	247	37.37	6.02	3.46	72
8	Uzbekistan	Europe	28541	28.9	6.38	2.38	68
9	Uruguay	Americas	3395	22.05	18.59	2.07	77

## Un-structured

```
{
  "code": "1473a6fd39d1d8fa48654aac9d8cc2754232",
  "title": "[Updating] Câu chuyện xuyên mưa về :  
",
  "url": "http://techtalk.vn/updating-cau-chuye",
  "labels": "techtalk/Cong nghe",
  "content": "Vào chiều tối ngày 09/12/2016 vừa  
",
  "image_url": "",
  "date": "2016-12-10T03:51:10Z"
}
```

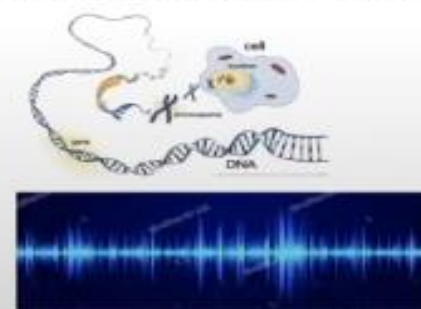
texts in websites, emails, articles, tweets



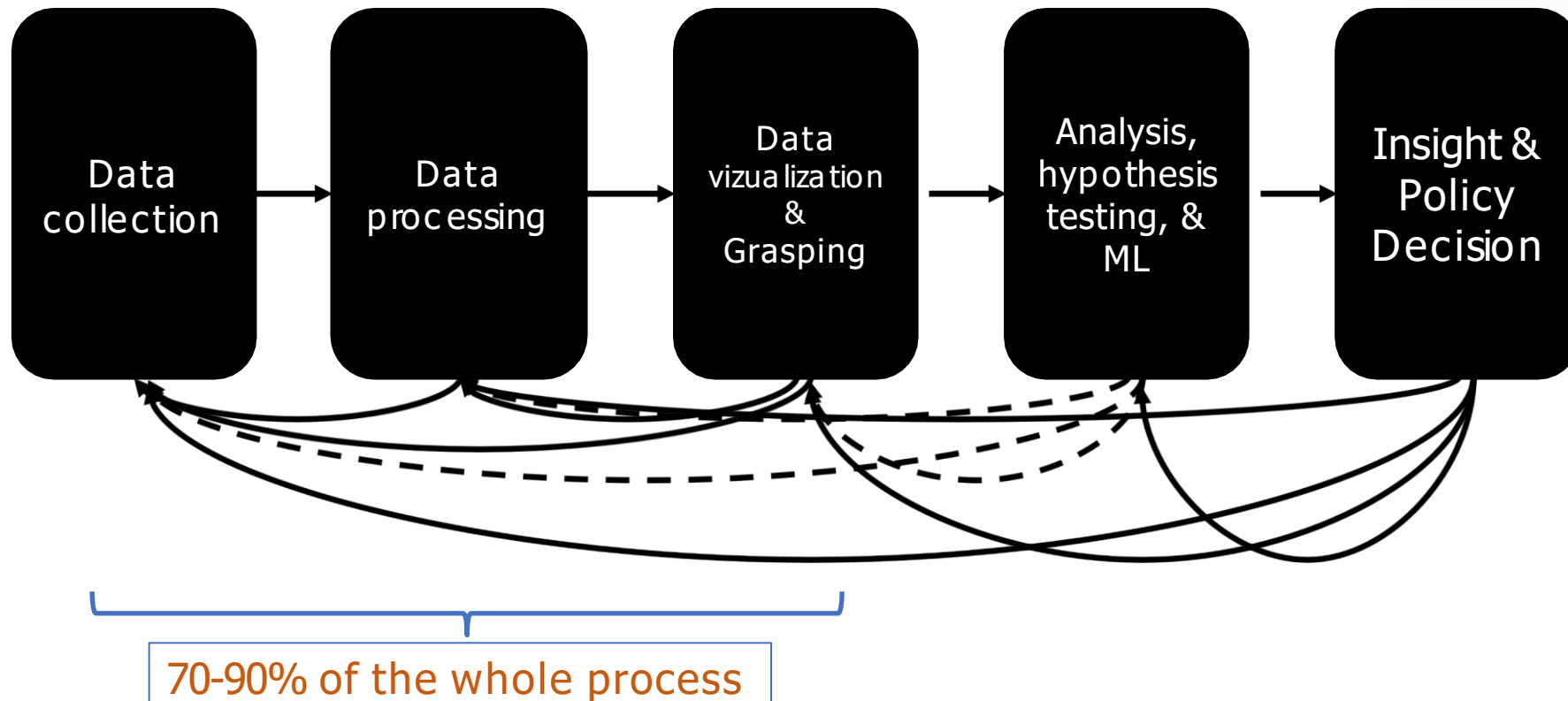
2D/3D images, videos + meta



spectrograms, DNAs, ...



## Methodology: *insight-driven*

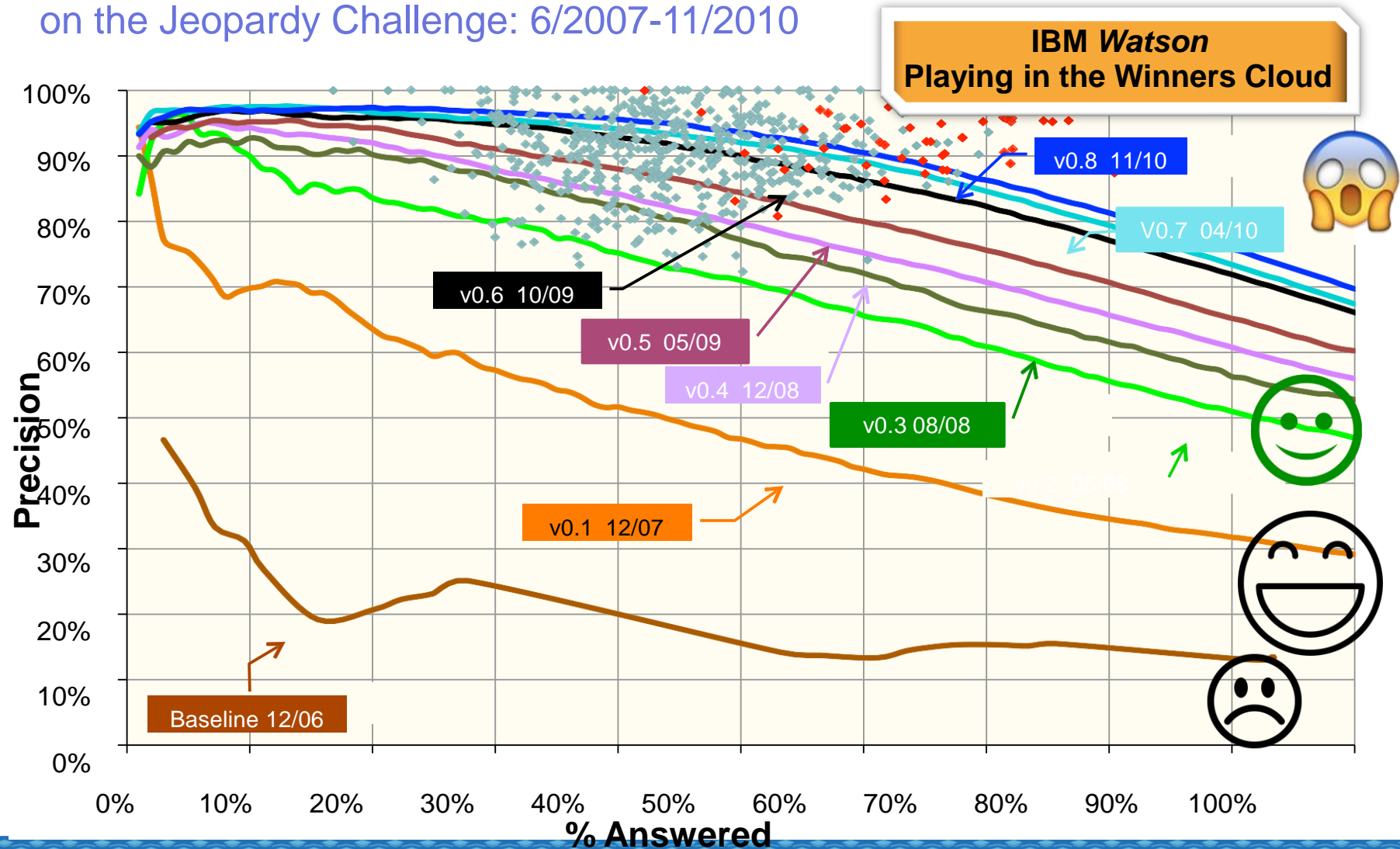


(John Dickerson, University of Maryland)



# Product development: experience

DeepQA: Incremental Progress in Answering Precision  
on the Jeopardy Challenge: 6/2007-11/2010





# What is Machine Learning?

- Machine Learning (ML) is an active subfield of Artificial Intelligence.
- ML seeks to answer the question [Mitchell, 2006]
  - *How can we build computer systems that automatically improve with experience, and what are the fundamental laws that govern all learning processes?*
- Some other views on ML:
  - Build systems that automatically improve their performance [Simon, 1983].
  - Program computers to optimize a performance objective at some task, based on data and past experience [Alpaydin, 2020]



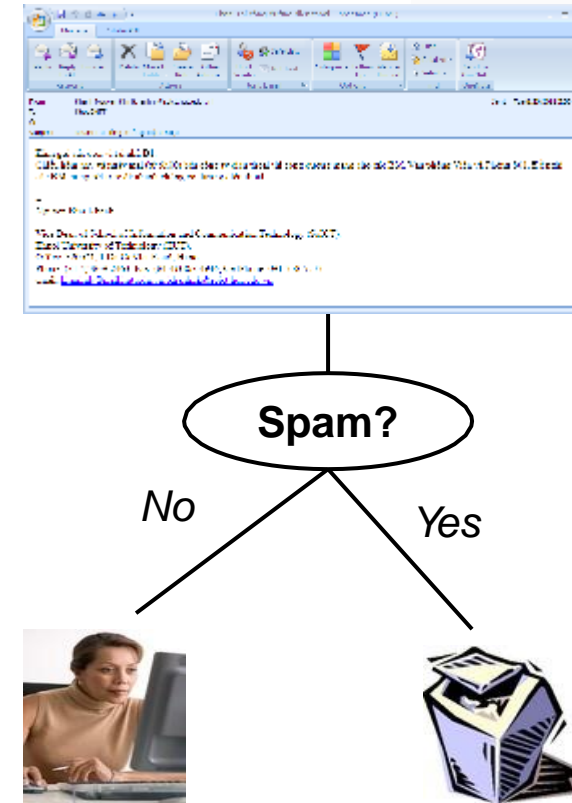
# A learning machine

- We say that a machine *learns* if the system reliably improves its performance **P** at task **T**, following experience **E**.
- A *learning problem* can be described as a triple (**P**, **T**, **E**).
- ML is close to and intersects with many areas.
  - Computer Science,
  - Statistics, Probability,
  - Optimization,
  - Psychology, Neuroscience,
  - Computer Vision,
  - Economics, Biology, Bioinformatics, ...

# Some real examples (1)

## ■ Spam filtering for emails

- **T**: filter/predict all emails that are spam.
- **P**: the accuracy of prediction, that is the percentage of emails that are correctly classified into normal/spam.
- **E**: set of old emails, each with a label of spam/normal.



# Some real examples (2)

## ■ Image tagging

- **T**: give some words that describe the meaning of a picture.
- **P**: ?
- **E**: set of pictures, each has been labelled with a set of words.



FISH WATER OCEAN  
TREE CORAL



PEOPLE MARKET PATTERN  
TEXTILE DISPLAY



BIRDS NEST TREE  
BRANCH LEAVES

# What does a machine learn?

---

- A **mapping** (function):

$$f : x \mapsto y$$

- $x$ : observations (data), past experience
  - $y$ : prediction, new knowledge, new experience,...
- A **model** (mô hình)
    - Data are often supposed to follow or be generated from an unknown model.  
(Ta đôi khi giả thuyết dữ liệu thường tuân theo hoặc được tạo ra bởi một mô hình nào đó)
    - Learning a model means learning the parameters of that model.  
(Học một mô hình có nghĩa là học/tìm những tham số của mô hình đó)



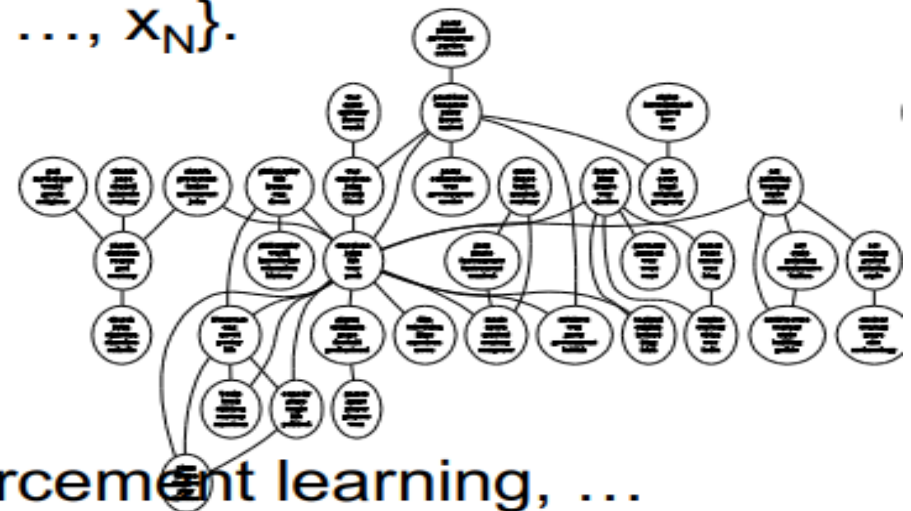
# Where does a machine learn from?

- Learn from a set of training examples (**training set**, tập học, tập huấn luyện)  $\{ \{x_1, x_2, \dots, x_N\}; \{y_1, y_2, \dots, y_M\} \}$ 
  - $x_i$  is an observation (quan sát, mẫu, điểm dữ liệu) of  $x$  in the past.
  - $y_j$  is an observation of  $y$  in the past, often called *label* (nhãn) or *response* (phản hồi) or *output* (đầu ra).
- After learning:
  - We obtain a model, new knowledge, or new experience ( $f$ ).
  - We can use that model/function to do **prediction** or **inference** for future observations, e.g.,

$$y = f(x)$$

# Two basic learning problems

- **Supervised learning (học có giám sát):** learn a function  $y = f(x)$  from a given training set  $\{x_1, x_2, \dots, x_N, y_1, y_2, \dots, y_N\}$  so that  $y_i \cong f(x_i)$  for every  $i$ .
  - **Classification** (categorization, phân loại, phân lớp): if  $y$  only belongs to a discrete set, for example {spam, normal}
  - **Regression** (hồi quy): if  $y$  is a real number
- **Unsupervised learning (học không giám sát):** learn a function  $y = f(x)$  from a given training set  $\{x_1, x_2, \dots, x_N\}$ .
  - $y$  can be a data cluster
  - $y$  can be a hidden structure
  - $y$  can be a trend
- Other: semi-supervised learning, reinforcement learning, ...



## Supervised learning: classification

- **Multiclass** classification (*phân loại nhiều lớp*):  
when the output  $y$  is one of the pre-defined labels  $\{c_1, c_2, \dots, c_L\}$   
(mỗi đầu ra chỉ thuộc 1 lớp, mỗi quan sát  $x$  chỉ có 1 nhãn)
  - Spam filtering:  $y$  in {spam, normal}
  - Financial risk estimation:  $y$  in {high, normal, no}
  - Discovery of network attacks: ?
- **Multilabel** classification (*phân loại đa nhãn*):  
when the output  $y$  is a subset of labels  
(mỗi đầu ra là một tập nhỏ các lớp;  
mỗi quan sát  $x$  có thể có nhiều nhãn)
  - Image tagging:  $y = \{\text{birds, nest, tree}\}$
  - sentiment analysis



## BIRDS NEST TREE



# Supervised learning: Regression

- Prediction of stock indices

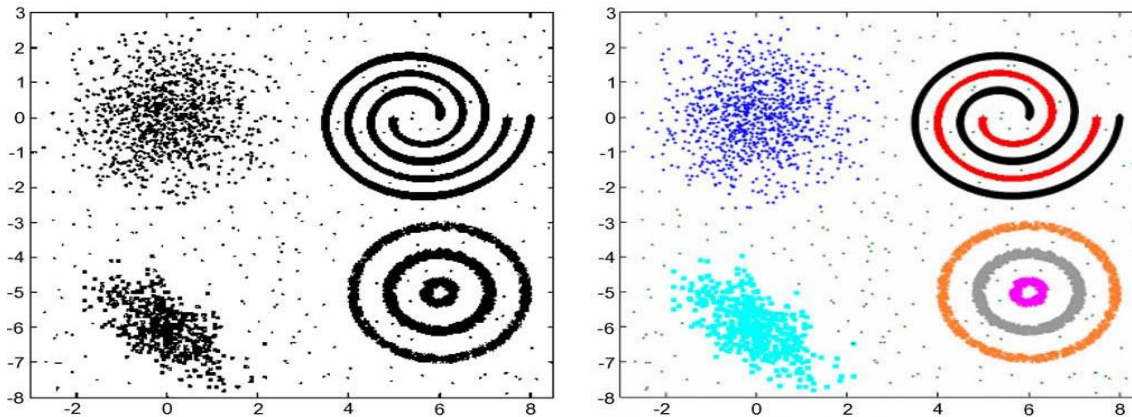


12.86	46.34	6	12.26	12.25	12.45	-4.25	-0.45
34.49	88.90	12	435.86	435.63	120.58	+6.63	+3.5
35.63	34.75	1	54.23	54.33	54.10	-0.33	-2
21.87	75.33	7	46.32	46.34	23.64	+1.34	+6
89.12	12.25	45	88.54	88.90	64.15	+2.90	+
34.3	35.63	6	43.45	43.66	43.62	-1.66	
25	21.87	45	12.23	12.86	75.21	+4.86	
96	89.12	7	434.64	434.49	632.55	-7.49	
7	23.43	34	32.21	32.00	12.21	-3.00	
	65.25	5	65.75	65.22	23.46	+0.7	
	42.96	12	123.74	123.76	121.51	-9	

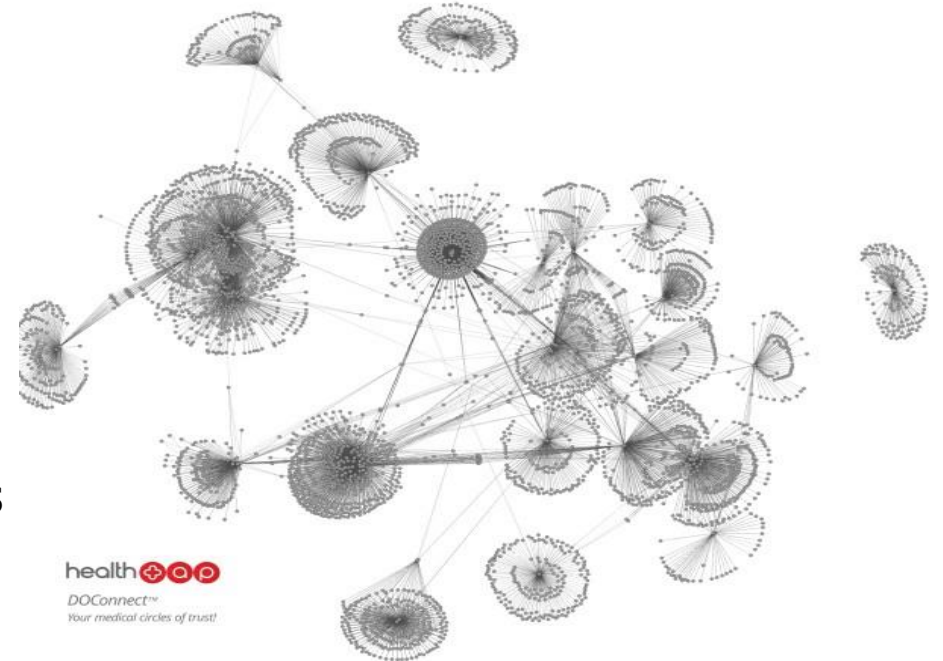


# Unsupervised learning: examples (1)

- Clustering data into clusters
  - Discover the data groups/clusters



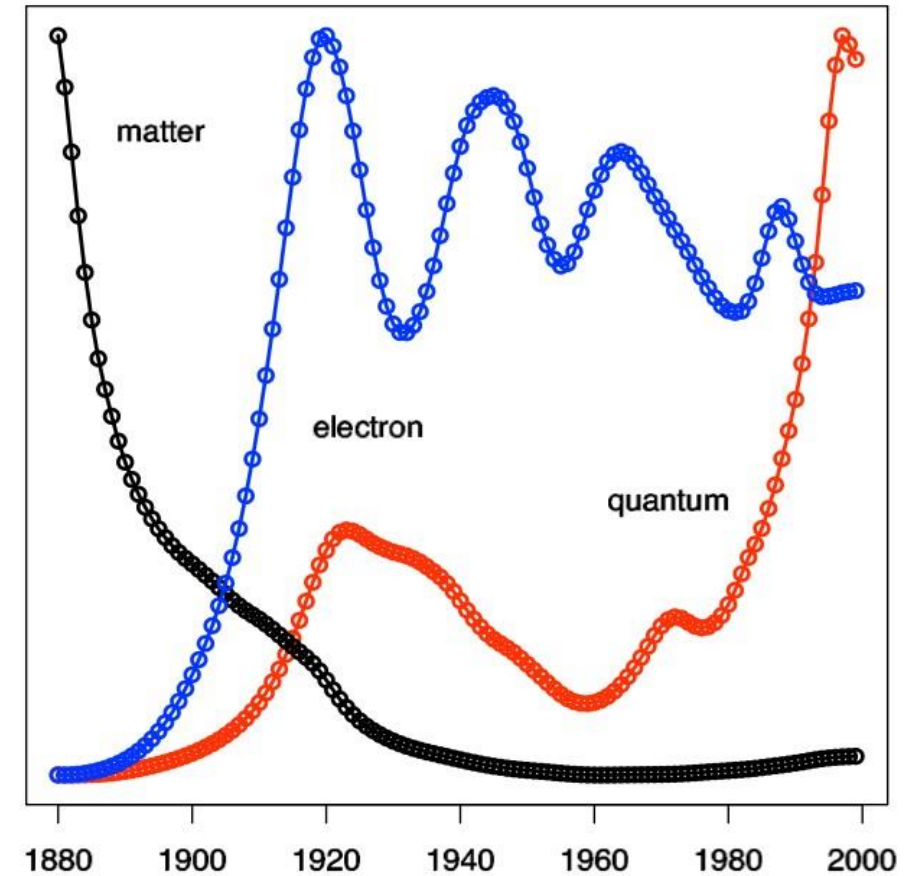
- Community detection
  - Detect communities in online social networks





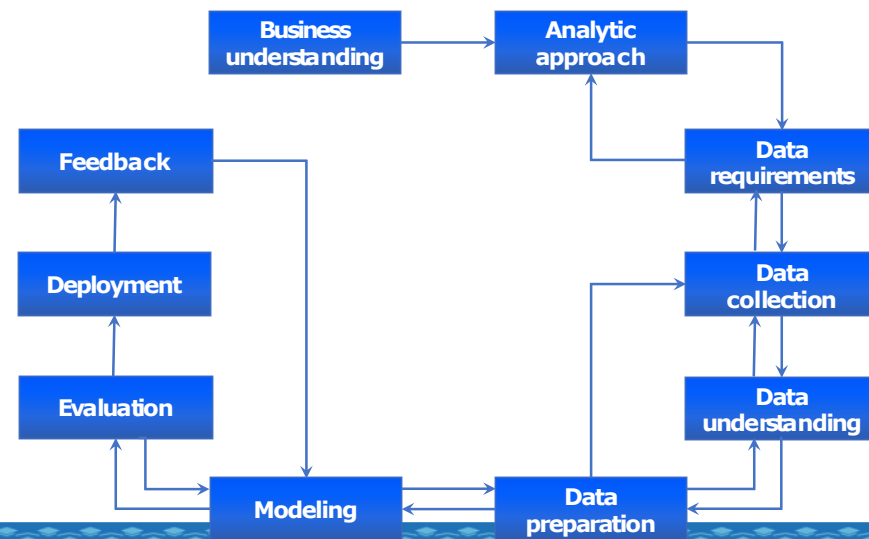
# Unsupervised learning: examples (2)

- Trends detection
  - Discover the trends, demands, future needs of online users



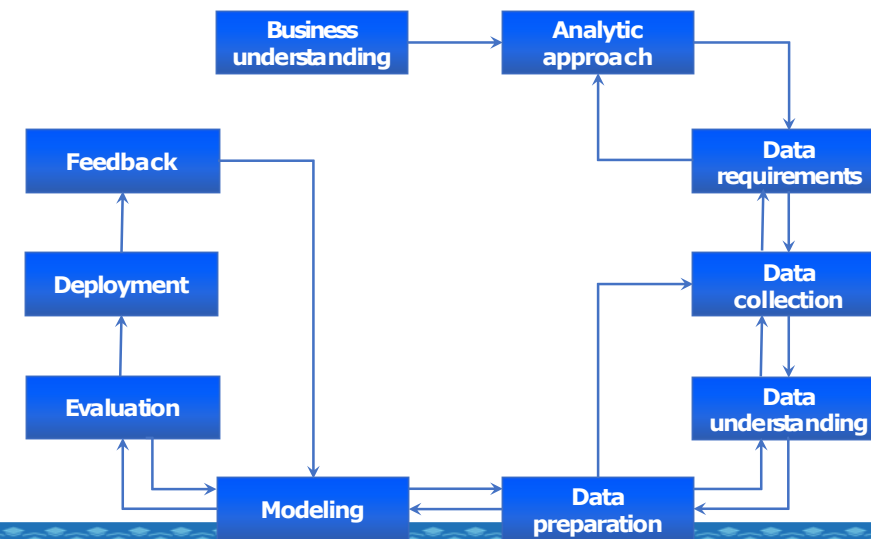
# Design a learning system (1)

- Some issues should be carefully considered when designing a learning system.
  - **Select a training set:**
    - The training set plays the key role in the effectiveness of the system.
    - Do the observations have any label?
    - The training observations should characterize the whole data space → good for future predictions.
  - **Determine the type of the function to be learned**
    - $F: X \rightarrow \{0,1\}$
    - $F: X \rightarrow \text{set of labels/tags}$
    - $F: X \rightarrow \mathbb{R}$
- 
- ```
graph TD; BU[Business understanding] --> AA[Analytic approach]; AA --> DR[Data requirements]; DR --> DC[Data collection]; DC --> DU[Data understanding]; DU --> E[Evaluation]; E --> D[Deployment]; D --> F[Feedback]; F --> AA;
```



# Design a learning system (2)

- Select a representation for the function: (model)
  - Linear?
  - A neural network?
  - A decision tree? ...
- Select a good algorithm to learn the function:
  - Ordinary least square? Ridge regression?
  - Back-propagation?
  - ID3?



# ML: some issues (1)

## ■ Learning algorithm

- Under what conditions the chosen algorithm will (asymptotically) converge?
- For a given application/domain and a given objective function, what algorithm performs best?

- ***No-free-lunch theorem*** [Wolpert and Macready, 1997]:  
if an algorithm performs well on a certain class of problems then it necessarily pays for that with degraded performance on the set of all remaining problems.

- *No algorithm can beat another on all domains.*  
(không có thuật toán nào luôn hiệu quả nhất trên mọi miền ứng dụng)



# ML: some issues (2)

---

## ■ Training data

- *How many observations* are enough for learning?
- Whether or not does the *size of the training set* affect performance of an ML system?
- What is the effect of the *disrupted* or *noisy* observations?





# ML: some issues (3)

---

## ■ Learnability:

- The goodness/limit of the learning algorithm?
- What is the **generalization** (tổng quát hoá) of the system?
  - ✧ Predict well new observations, not only the training data.
  - ✧ Avoid overfitting.

# Overfitting (quá khớp, quá khít)

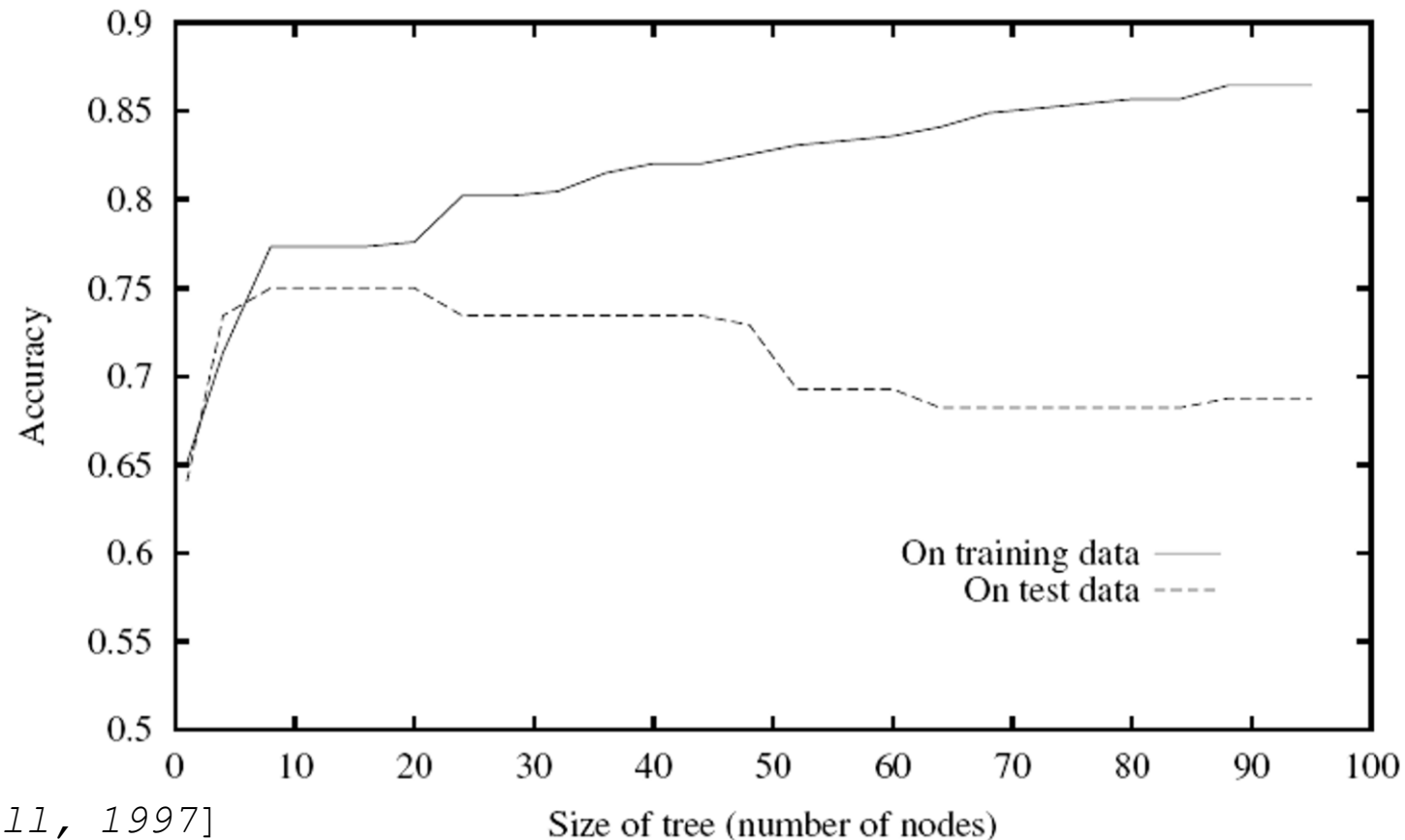
- Function  $h$  is called *overfitting* [Mitchell, 1997] if there exists another function  $g$  such that:
  - $g$  might be worse than  $h$  for the training data, but
  - $g$  is better than  $h$  for future data.
- A learning algorithm is said to overfit relative to another one if it is *more accurate in fitting* known data, but *less accurate in predicting* unseen data.
- Overfitting is caused by many factors:
  - The trained function/model is **too complex** or have too much parameters.
  - **Noises or errors** are present in the training data.
  - The training size is **too small**, not characterizing the whole data space.

# Overfitting



# Overfitting: *example*

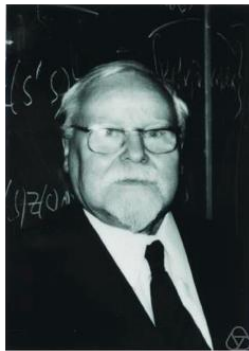
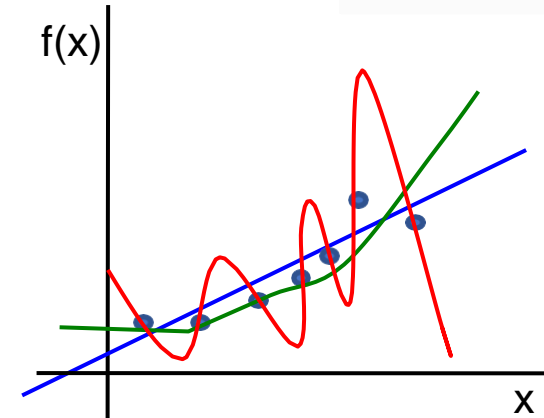
- Increasing the size of a decision tree can degrade prediction on unseen data, even though increasing the accuracy for the training data.



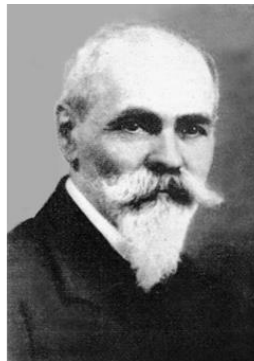
[Mitchell, 1997]

# Overfitting: Regularization

- Among many functions, which one can generalize best from the given training data?
  - *Generalization is the main target of ML.*
  - Predict unseen data well.
- **Regularization:** a popular choice (Hiệu chỉnh)



Tikhonov,  
smoothing an ill-  
posed problem



Zaremba, model  
complexity  
minimization



Bayes: priors  
over parameters



Andrew Ng: need no  
maths, but it prevents  
overfitting!



# References

---

- Alpaydin E. (2020). Introduction to Machine Learning. The MIT Press.
- Mitchell, T. M. (1997). Machine learning. *McGraw Hill*.
- Mitchell, T. M. (2006). *The discipline of machine learning*. Carnegie Mellon University, School of Computer Science, Machine Learning Department.
- Simon H.A. (1983). Why Should Machines Learn? In R. S. Michalski, J. Carbonell, and T.M. Mitchell (Eds.): Machine learning: An artificial intelligence approach, chapter 2, pp. 25-38. Morgan Kaufmann.
- Valiant, L. G. (1984). A theory of the learnable. *Communications of the ACM*, 27(11), 1134-1142.
- Wolpert, D.H., Macready, W.G. (1997), "No Free Lunch Theorems for Optimization", *IEEE Transactions on Evolutionary Computation* **1**, 67.