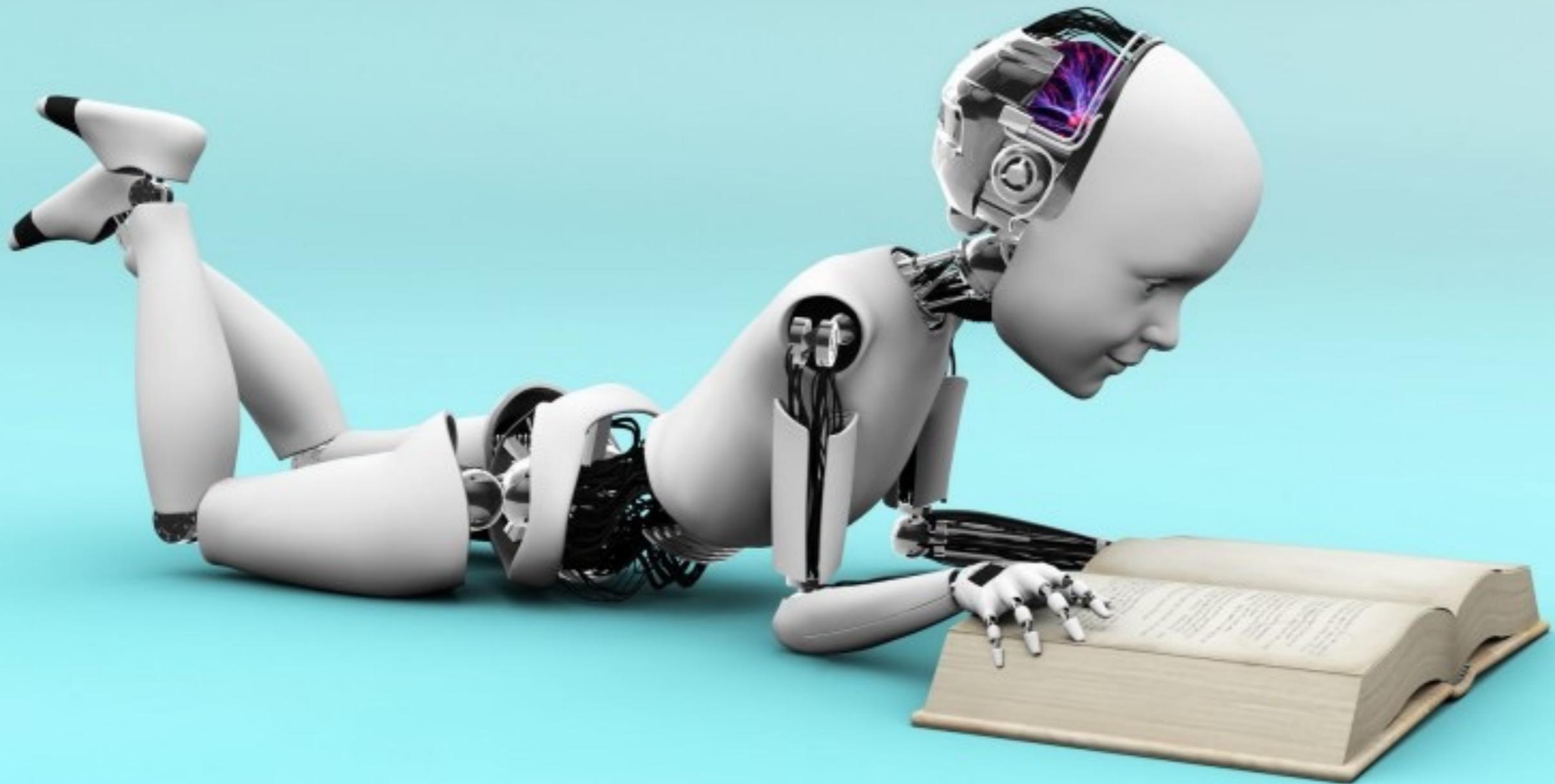


# Crash Course on Machine Learning

August 22–24, 2016



**simula**  
**education**

**simula**

# Instructors



**Valeriya Naumova**

Senior Research Scientist

[valeriya@simula.no](mailto:valeriya@simula.no)

**Interests:** applied mathematics, learning theory, inverse problems



**Arnaud Gotlieb**

Chief Research Scientist

[arnaud@simula.no](mailto:arnaud@simula.no)

**Interests:** software engineering & testing

# Lab Instructors



**Timo Klock**

PhD Student,  
Biomedical Computing  
Department

[timo@simula.no](mailto:timo@simula.no)

# Organisers



**Are Magnus Bruaset**

Head of SSRI,  
Section Director  
Computing and Software



**Carlo Leva**

Senior Research Engineer,  
Software Engineering  
Department

[carlo@simula.no](mailto:carlo@simula.no)



**Elin Backe**

**Christophersen**  
SSRI Advisor



**Sigurd Lekve**

Summer Intern,  
NTNU Student

[sigurd.lekve@gmail.com](mailto:sigurd.lekve@gmail.com)



**Teodora Tufa**

Executive Officer

# Syllabus

<b>Mon 22</b>	<b>Morning</b>	Introduction to ML. Local Methods. Model Selection. Regularisation I.
	<b>Afternoon</b>	Machine Learning in Practice. Lab Work.
<b>Tue 23</b>	<b>Morning</b>	Regularisation II. Dimensionality Reduction, Variable Selection.
	<b>Afternoon</b>	Machine Learning Applications and Trends. Lab Work.
<b>Wed 24</b>	<b>Morning</b>	Basics in Logic. Programming in Logic. Decision Tree Induction.
	<b>Afternoon</b>	Version Space Search. Lab Work.
<b>Thu 25</b>	<b>Morning</b>	Open Problems Session and Discussion.
<b>Sept 16 - May 17</b>		Seminar and Mini-tutorial Series by Invited Guests (TBA).

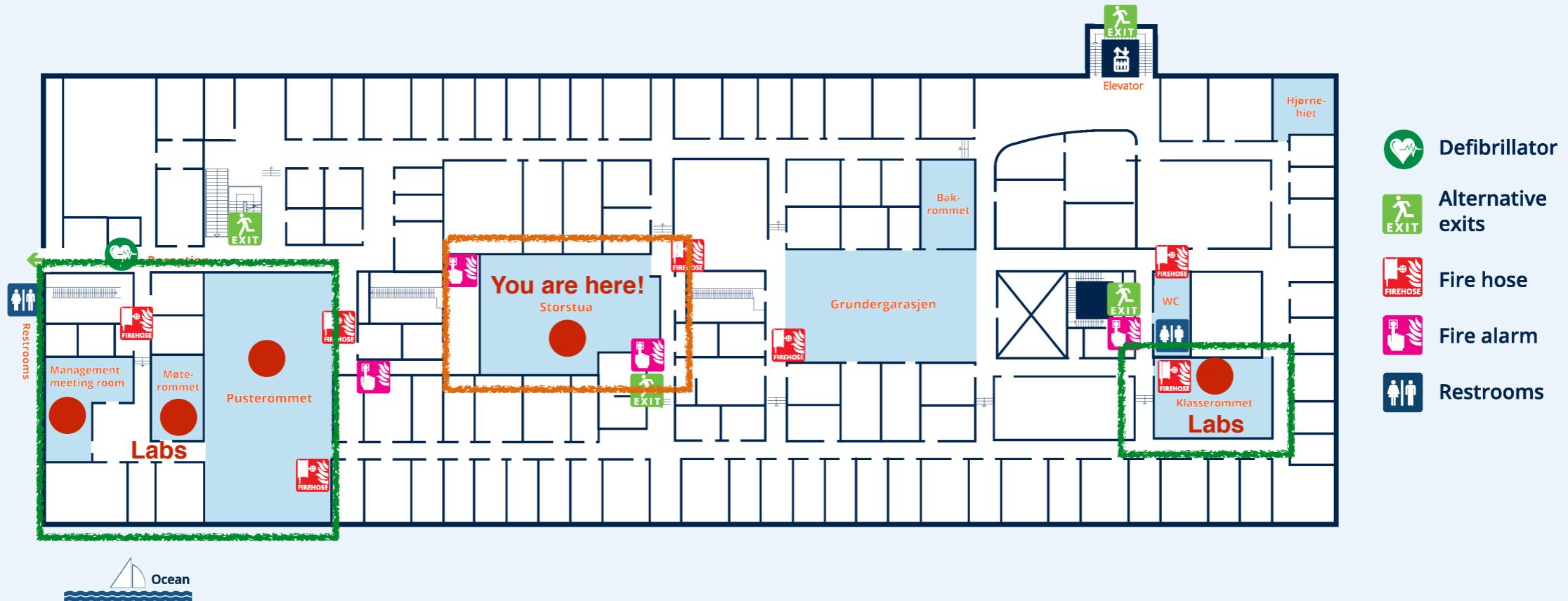
An introduction to *essential Machine Learning*:

- Concepts
- Algorithms

11:00-11:30 - Break  
13:00-14:00 - Lunch  
18:00-20:00 (Monday) - Social event

# Practical Information

## Important information in case of an emergency



### Important information in case of an emergency

Below you will find the central emergency contact numbers to the fire brigade, police and ambulance, as well as the link to the official Norwegian website providing information to the general public before, during and after a crisis. In addition, you will find information on how any emergency event affecting Simula's premises, employees, students or guests should be notified internally.

**Fire:** 110  
**Police:** 112  
**Ambulance:** 113

Crisis information from the authorities: [www.kriseinfo.no](http://www.kriseinfo.no)

Security at Technopolis (24h): 67 82 70 00

### Emergency notification internally at Simula

If an emergency occurs that involves Simula's employees, students, guests or premises, after having alerted the appropriate emergency contact number, you should report directly to either of the contact persons below:

Contact	Role	Phone number
• Aslak Tveito	Managing Director	(+47) 906 87 348
• Kyrre Lekve	Deputy Managing Director	(+47) 934 24 311
• Marianne M. Sundet	Director of Administration	(+47) 900 18 483
• Halvard Moe	Technical Director	(+47) 928 10 236
• Maria Benterud	Administrative Manager	(+47) 954 95 616

**simula**

# Course Material

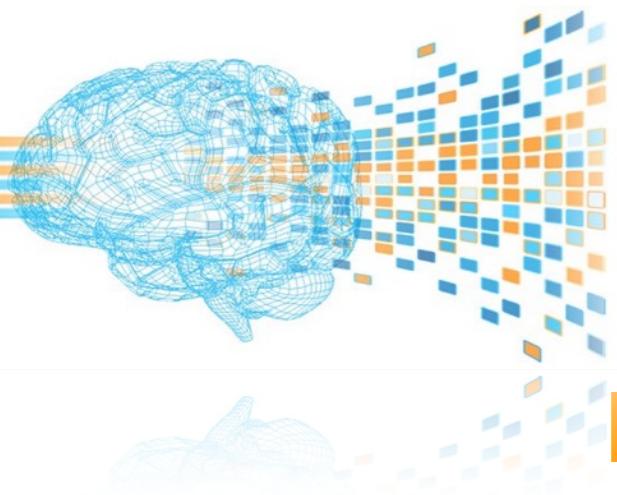
**Lectures and materials for lab sessions are available via**

- ▶ Github vnaumova / ML\_Course2016
- ▶ Dropbox

**Software requirements for lab sessions**

- ▶ Python with Matplotlib and Scipy packages (for Monday-Tuesday)
- ▶ SICStus Prolog 4.3 (for Wednesday)





# Introduction to Machine Learning

## Basic Concepts. Regularisation I

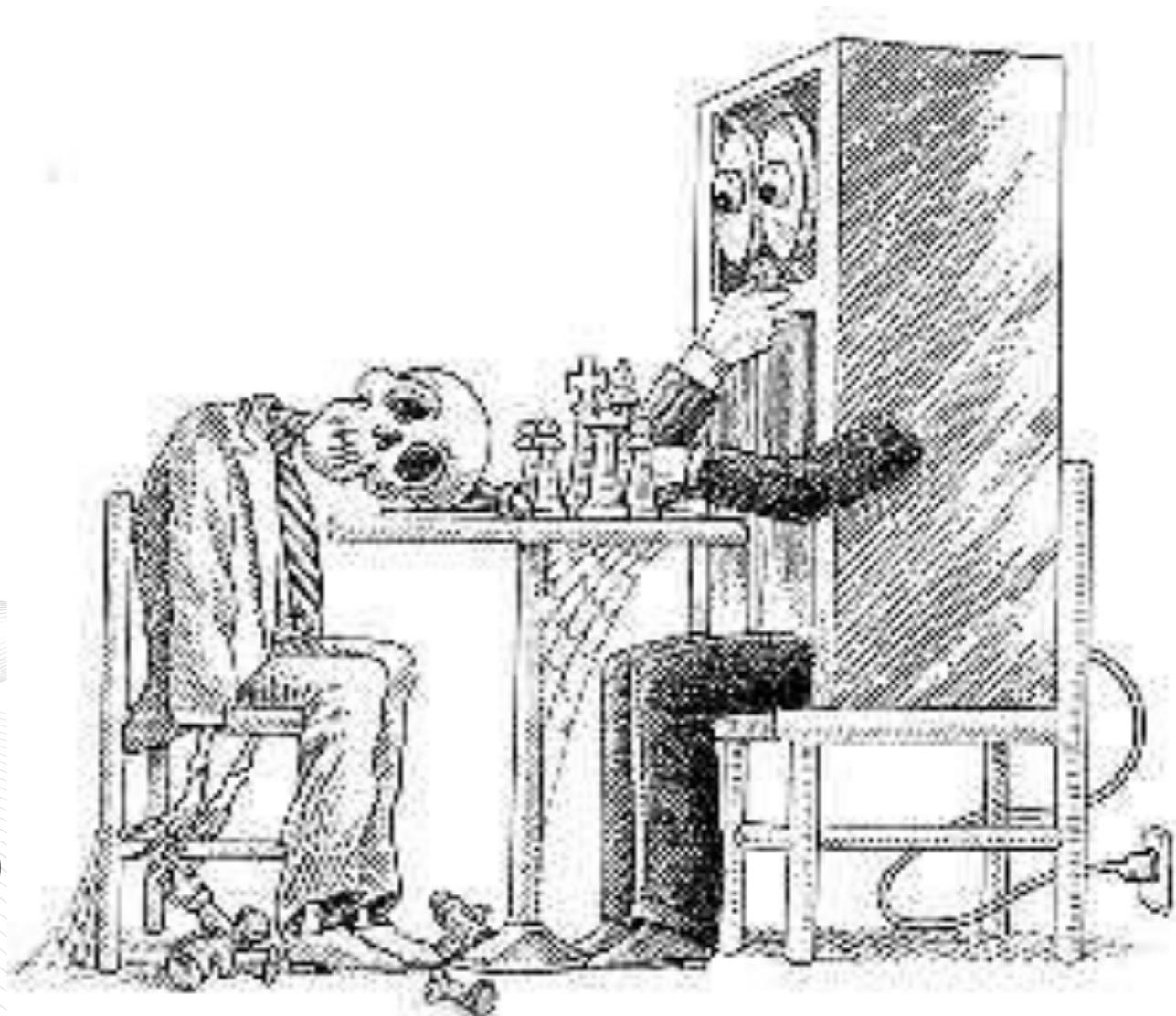
### Crash Course on Machine Learning

August 22-24, 2016

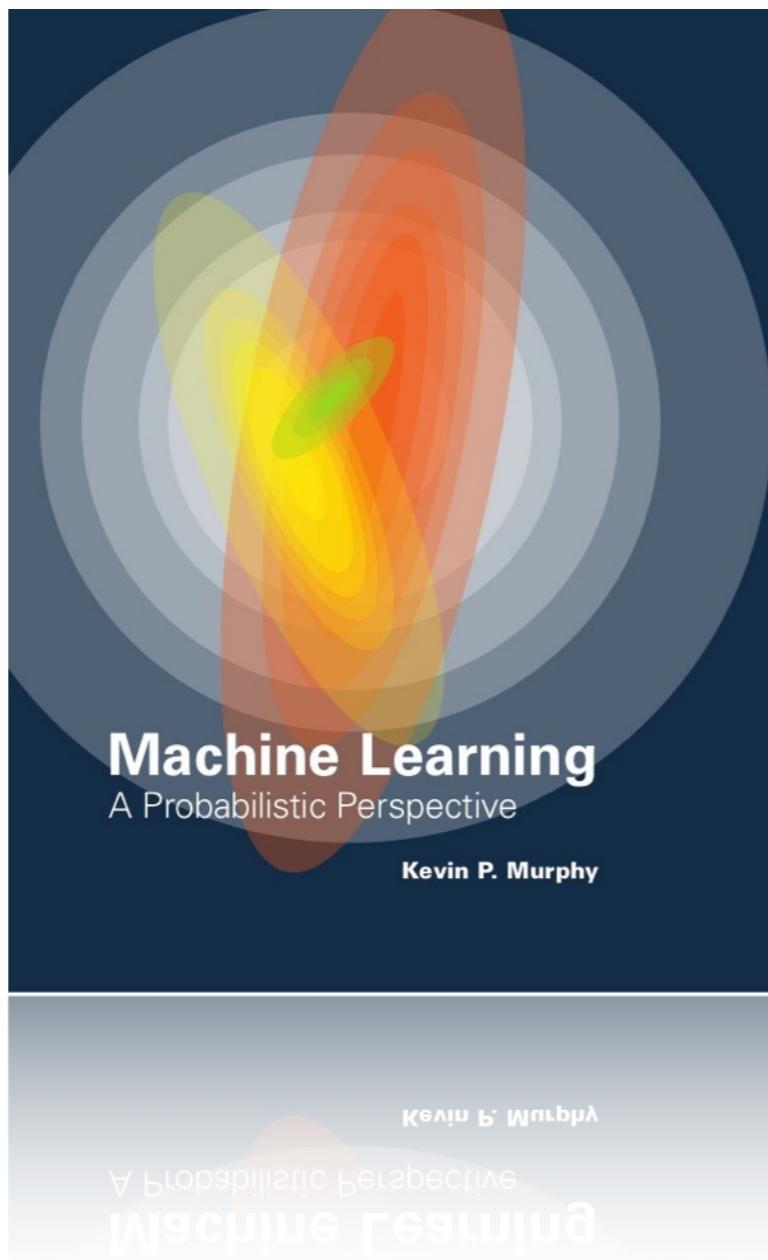
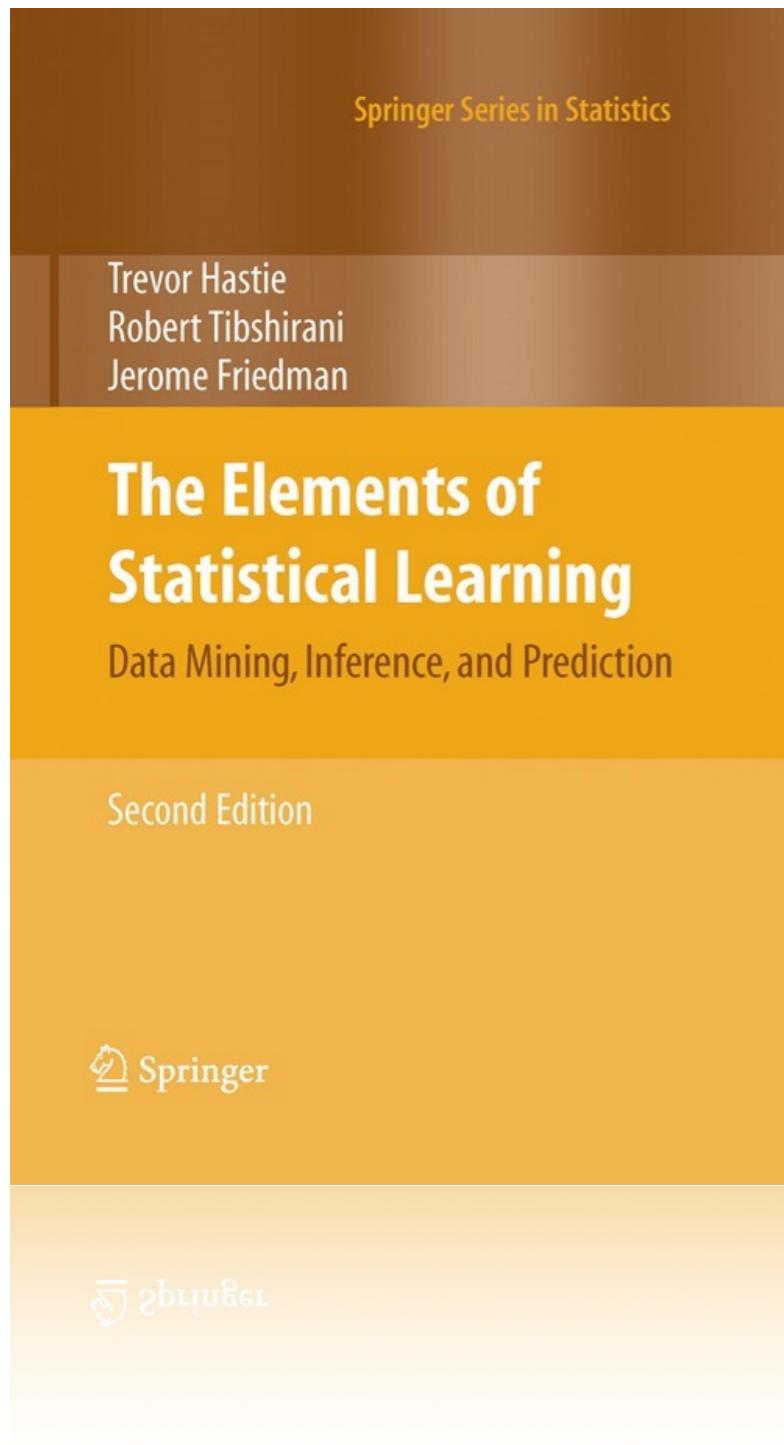
**Valeriya Naumova and Arnaud Gotlieb**

Simula Research Laboratory AS

Simula School of Research and Innovation (SSRI)



# Literature



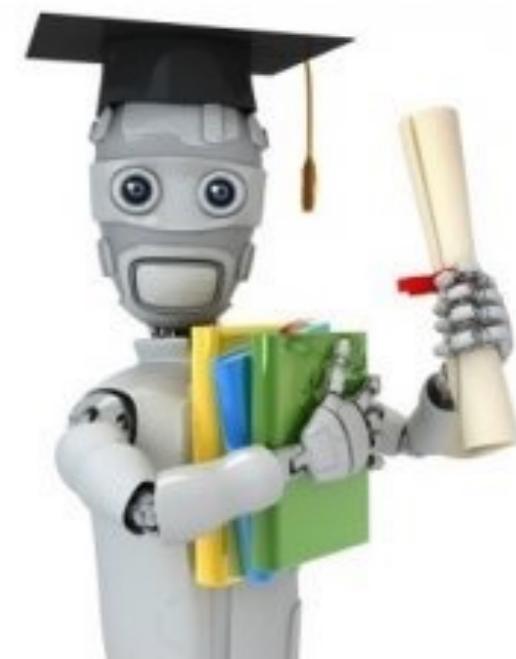
We are drowning in information and  
starving for knowledge.

*Rutherford D. Roger, American Librarian*

coursera

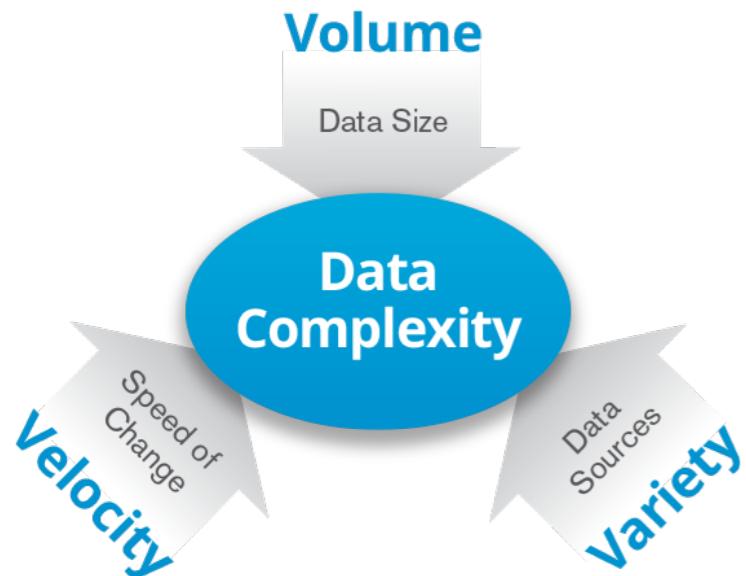


Andrew Ng



# Table of Contents

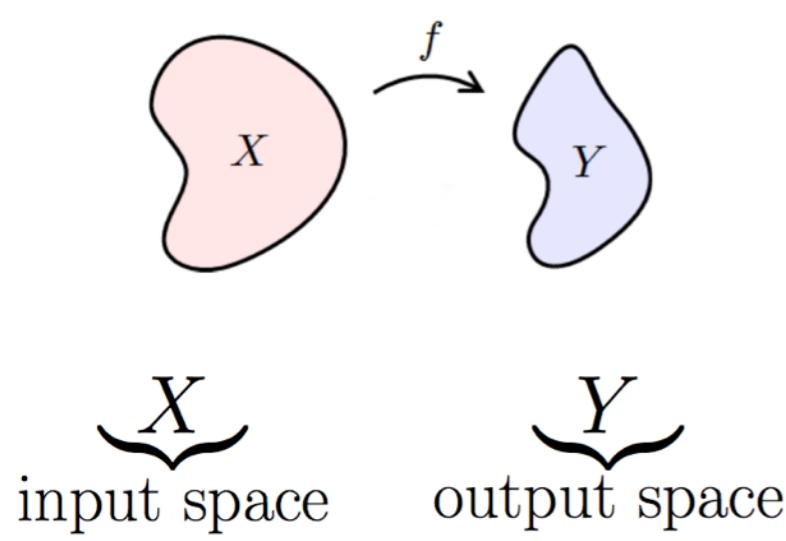
## What is Big Data?



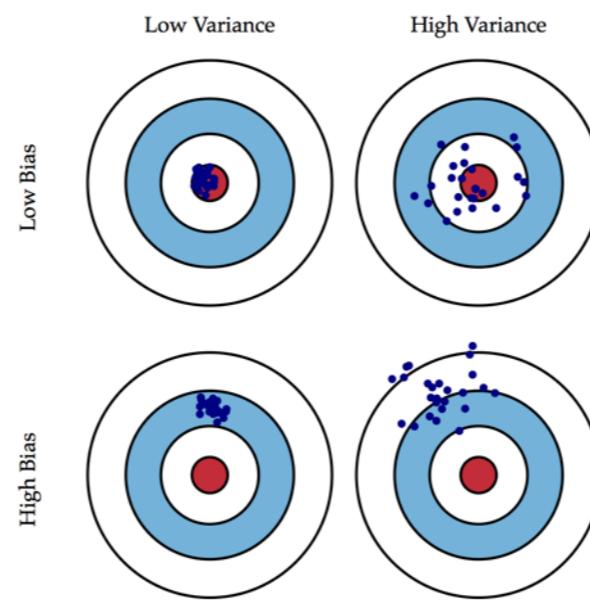
## ML: History, Types, Applications



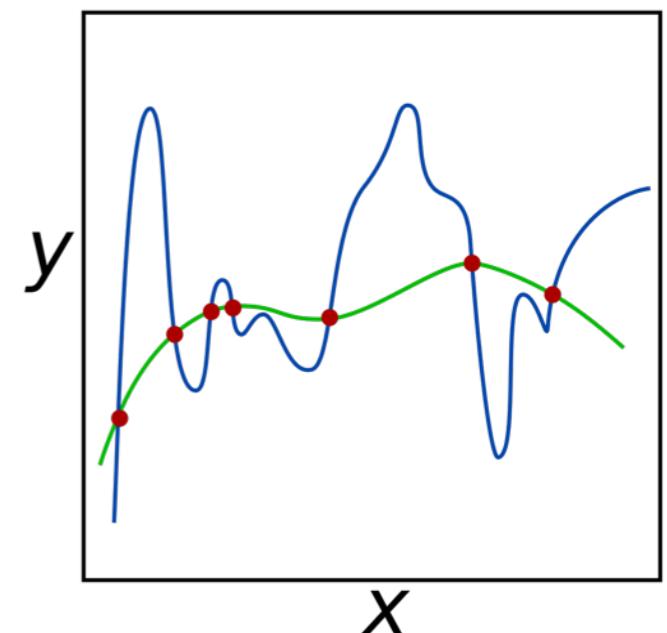
## Basic Statistical Learning Theory



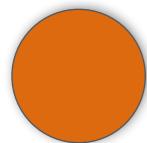
## Local Methods & Bias-Variance



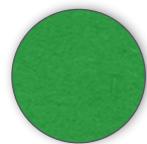
## Regularisation I



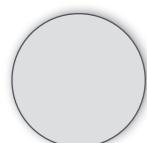
# Colour Coding / Feedback



**Really important stuff**



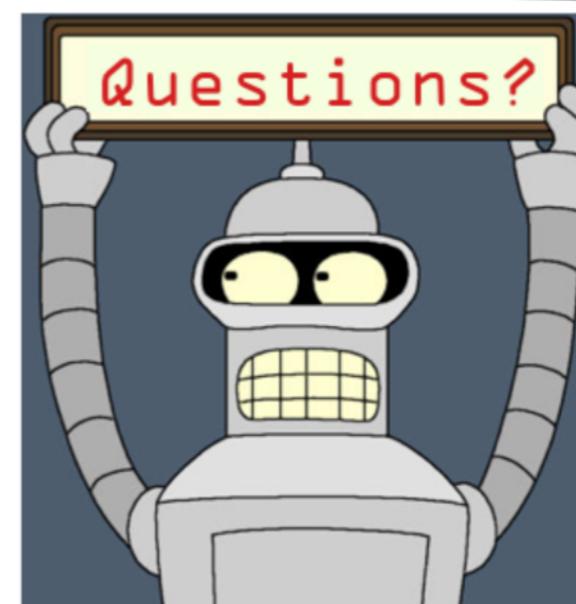
**Important stuff**



**Regular stuff**

**Let us know if you have  
comments, concerns,  
suggestions!**

The course contains many ideas and  
(quite) a bit of math, questions help  
prevent sleeping...



# Harvard Business Review

GETTING  
CONTROL  
OF



How vast new streams of  
information are changing  
the art of management  
**PAGE 59**

OCTOBER 2012

46 The Big Idea  
The True Measures  
Of Success  
Michael J. Mauboussin

84 International Business  
10 Rules for Managing  
Global Innovation  
Keeley Wilson and Yves L. Duzé

93 Leadership  
What Ever Happened  
To Accountability?  
Thomas E. Ricks

Data scientist ... 'sexiest job of the 21st Century'.

Harvard Business Review

# The Economist

FEBRUARY 22ND-MARCH 8TH 2014

Economist.com

Obama the warrior  
Misgoverning Argentina  
The economic shift from West to East  
Genetically modified crops blossom  
The right to eat cats and dogs

## The data deluge

AND HOW TO HANDLE IT: A 14-PAGE SPECIAL REPORT



In God we trust, all others bring data.

William Deming, Engineer & Statistician

**five**

# The four dimensions (V's) of Big Data



- ▶ A collection of large and complex data sets which are difficult to process using common database management tools or traditional data processing applications.
- ▶ “Big data refers to the tools, processes and procedures allowing an organisation to create, manipulate, and manage very large data sets and storage facilities”.

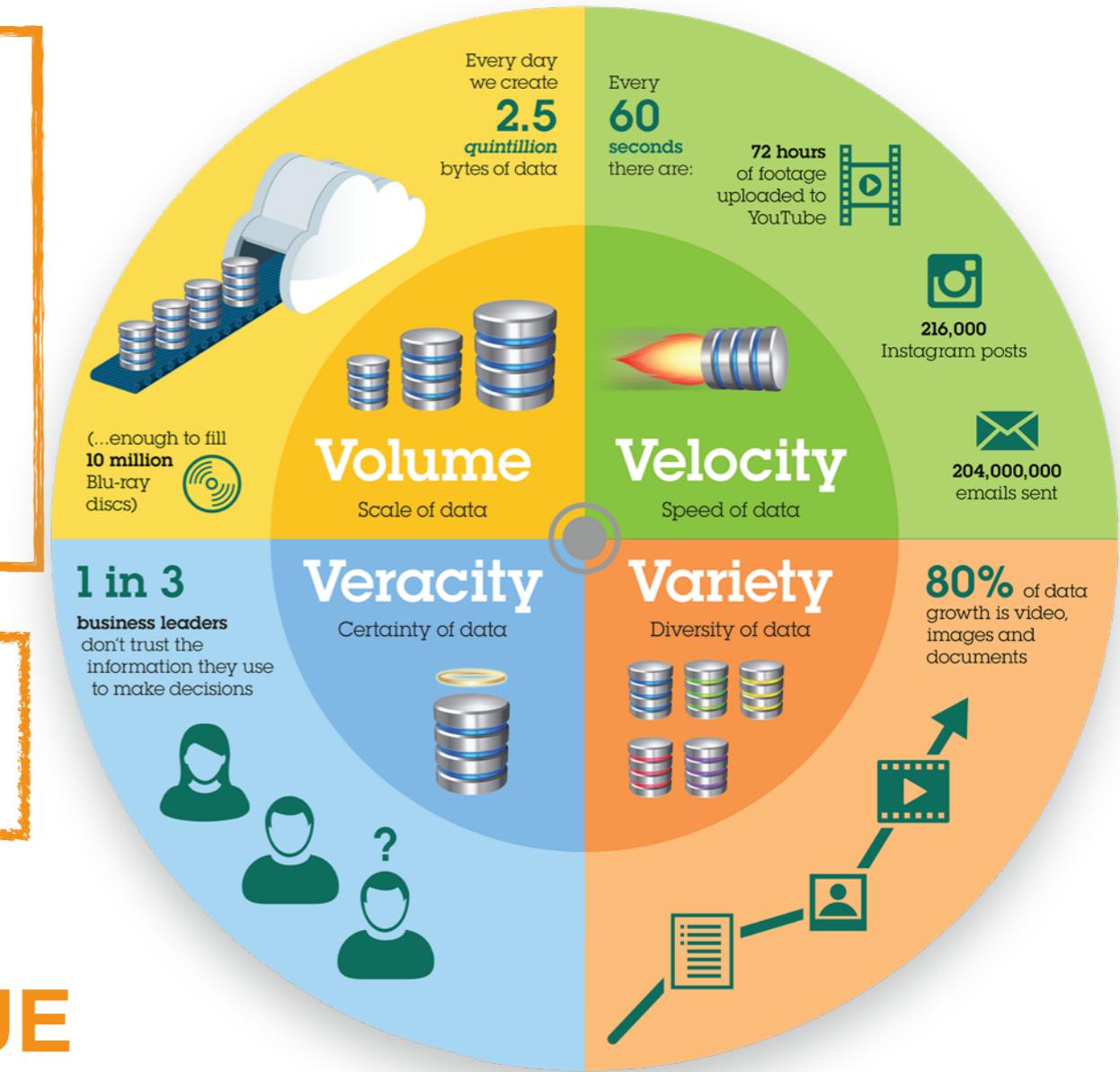
according to [zdnet.com](http://zdnet.com)

## Big data is not just about size.

- ▶ Finds insights from complex, noisy, heterogeneous, longitudinal, and voluminous data.
- ▶ It aims to answer questions that were previously unanswered.

The challenges include capturing, storing, searching, sharing & analysing.

+ VALUE

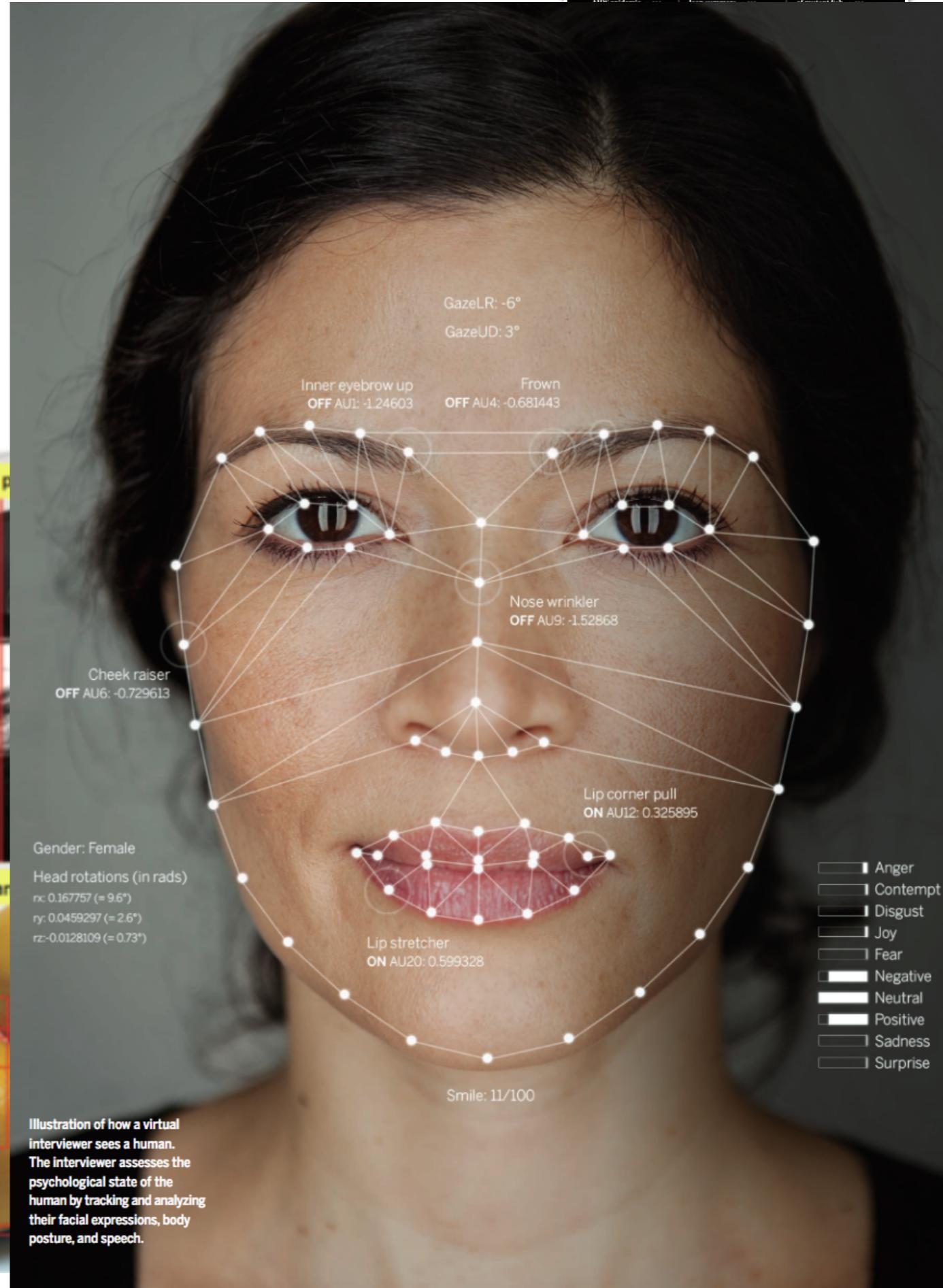
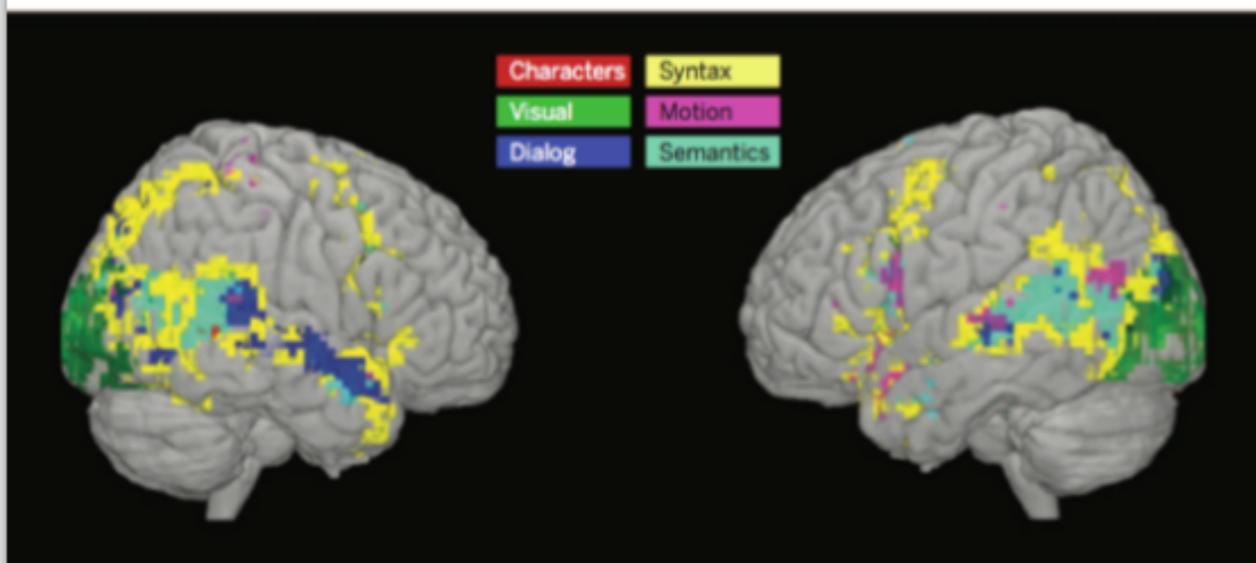
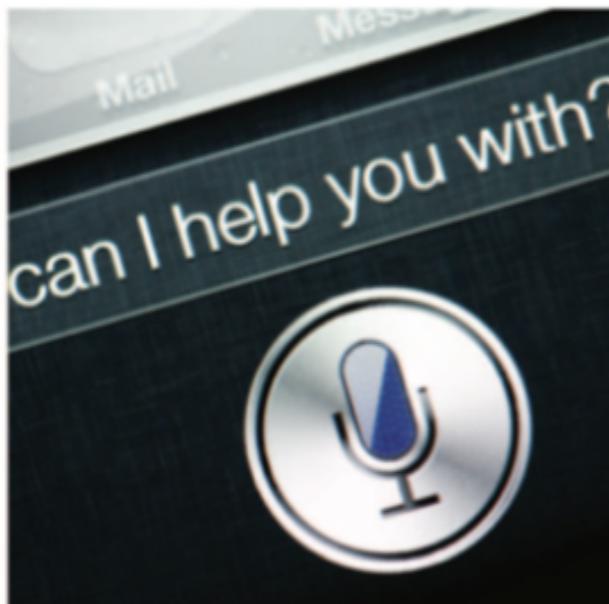


# Big Data Analysis - How?

→ machine / statistical learning

A breakthrough in machine learning  
would be worth ten Microsofts.

*Bill Gates, Microsoft*



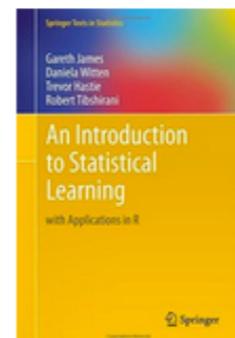
# Big Data Analysis : Amazon / Netflix



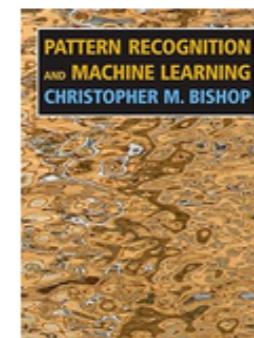
Netflix  
recommendations

Amazon  
Books

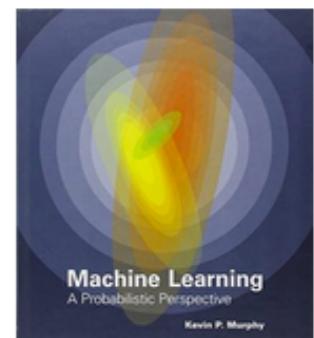
Customers Who Bought This Item Also Bought



An Introduction to  
Statistical Learning: with  
Applications in R...  
› Gareth James  
★★★★★ 3  
Gebundene Ausgabe  
EUR 55,99 ✓Prime



Pattern Recognition and  
Machine Learning  
(Information Science and...  
Christopher Bishop  
★★★★★ 13  
Gebundene Ausgabe  
EUR 62,45 ✓Prime

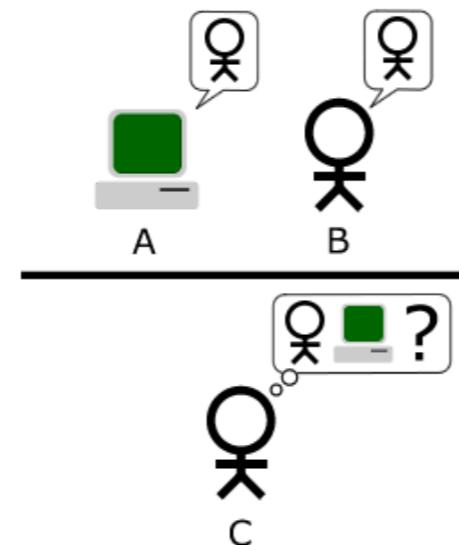


Machine Learning: A  
Probabilistic Perspective  
(Adaptive computation...  
› Kevin P Murphy  
★★★★★ 1  
Gebundene Ausgabe  
EUR 79,53 ✓Prime

Machine learning is the next Internet.

Tony Tether, DARPA

# The Quest for AI. Birth of Dream



Alan Turing (1912–1954) proposed a test that calls for a panel of judges to review typed answers to any question that has been addressed to both a computer and a human. If the judges can make no distinctions between the two answers, the machine may be considered **intelligent**.

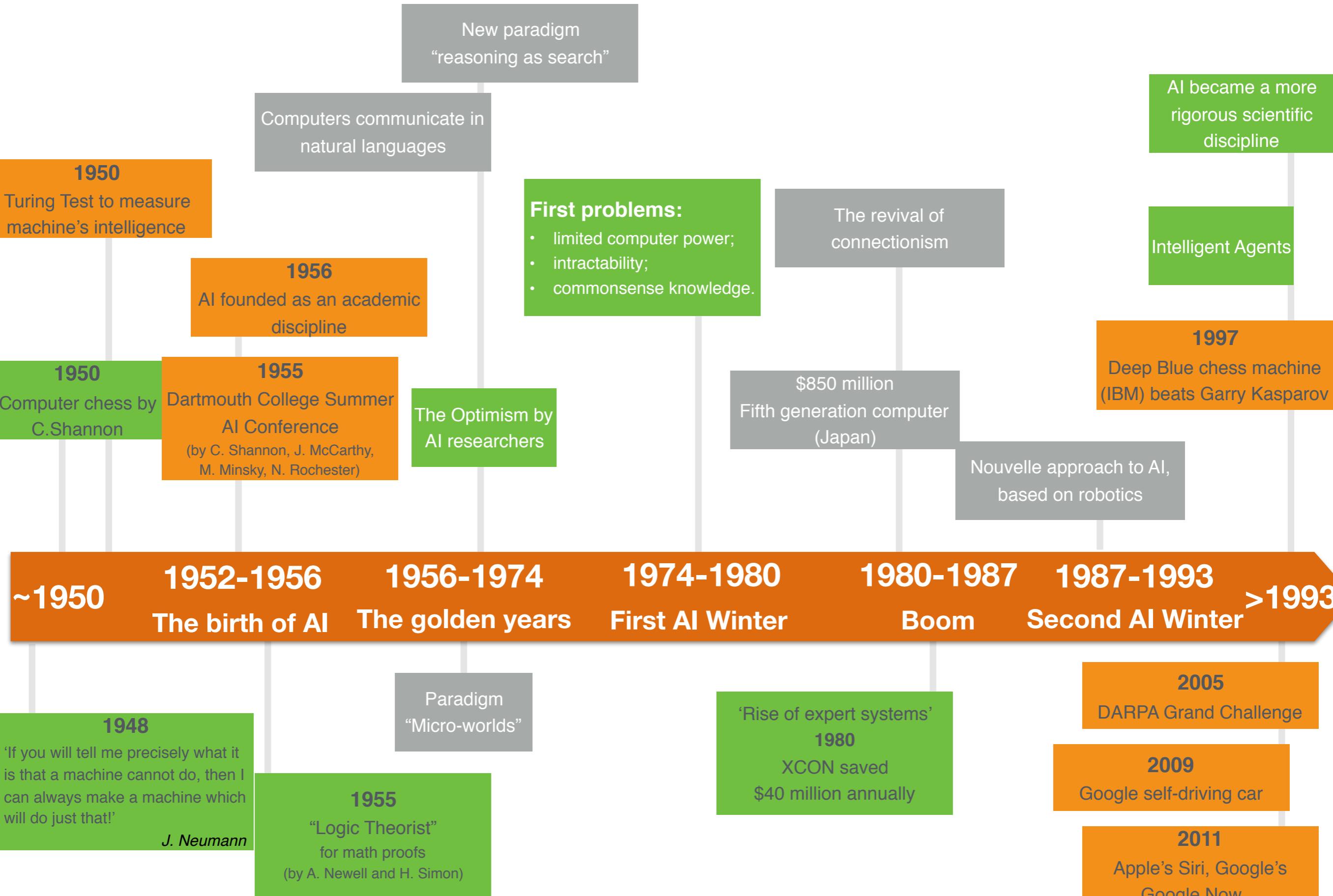
## Intelligence: a Working Definition

- *abstract reasoning, knowledge acquisition, decision making.*
- *knowledge acquisition: memorisation vs learning.*

## Ingredients of AI

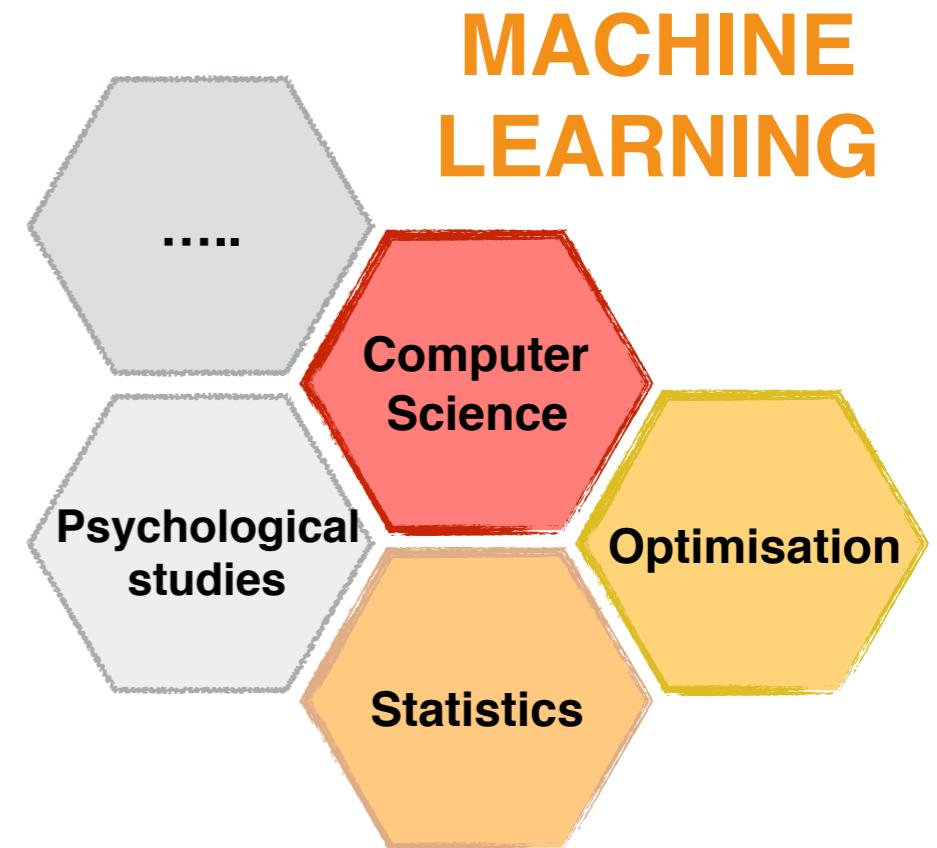
- natural language processing;
- knowledge representation;
- automated reasoning;
- machine learning;
- computer vision;
- robotics to manipulate.

# The Quest for AI. Birth of Dream



# So What is Machine Learning?

- ▶ Automating automation
- ▶ Getting computers to program themselves
- ▶ Writing software is the bottleneck
- ▶ Let the data do the work instead!



## Traditional Programming



## Machine Learning



# Is Machine Learning Magic?

No, it is more like gardening...

▶ **Seeds** = Algorithms

▶ **Nutrients** = Data

▶ **Gardener** = You

▶ **Plants** = Programs

Machine learning is not a magic; it can't get something from nothing. What it does it gets more from less.

*Pedro Domingos, UC Washington*



# Sample Applications: Prediction

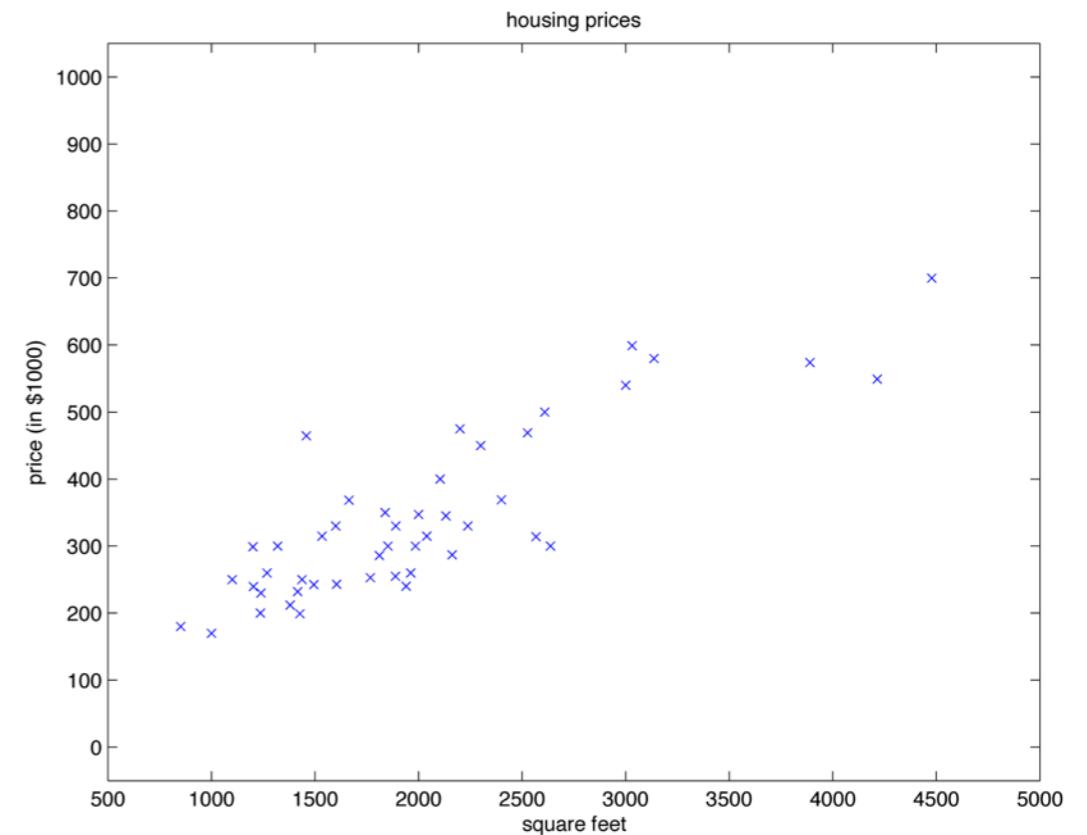
Living area (feet <sup>2</sup> )	Price (1000\$s)
2104	400
1600	330
2400	369
1416	232
3000	540
:	:

# Sample Applications: Prediction

Living area (feet <sup>2</sup> )	Price (1000\$s)
2104	400
1600	330
2400	369
1416	232
3000	540
$\vdots$	$\vdots$
$S = \{(x_1, y_1), \dots, (x_n, y_n)\}.$	

# Sample Applications: Prediction

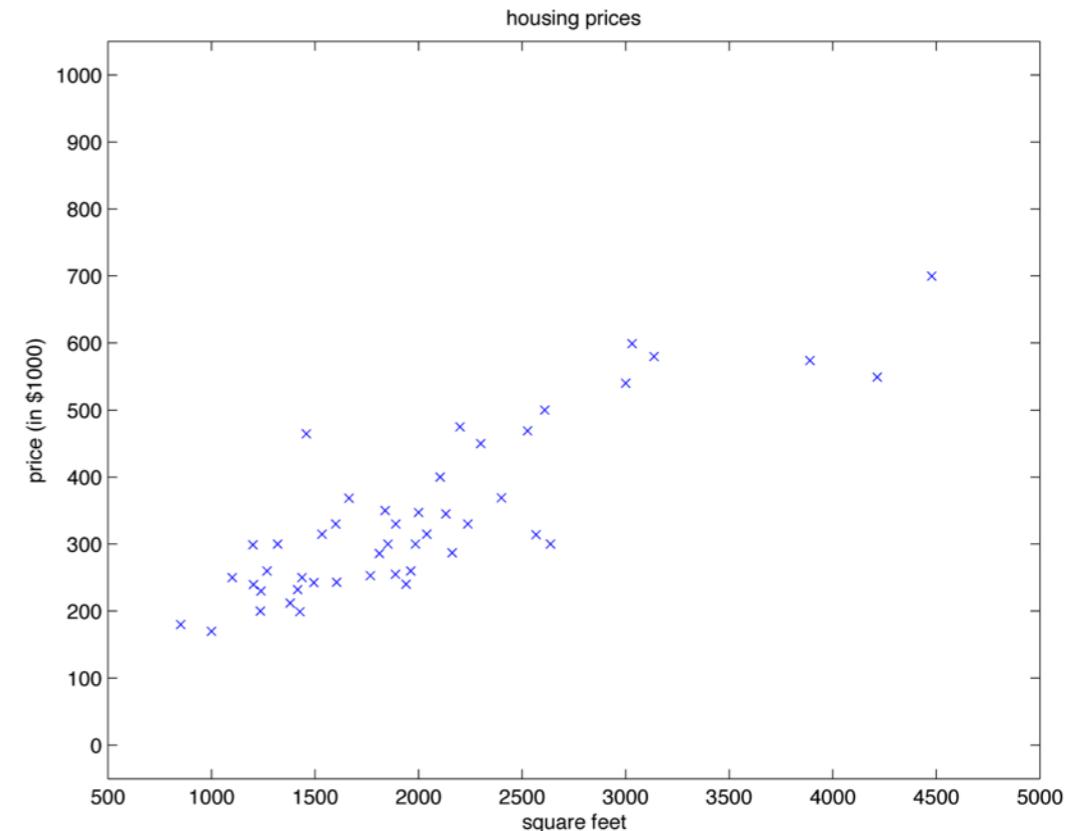
Living area (feet <sup>2</sup> )	Price (1000\$s)
2104	400
1600	330
2400	369
1416	232
3000	540
$\vdots$	$\vdots$
$S = \{(x_1, y_1), \dots, (x_n, y_n)\}$ .	



# Sample Applications: Prediction

Living area (feet <sup>2</sup> )	Price (1000\$)
2104	400
1600	330
2400	369
1416	232
3000	540
$\vdots$	$\vdots$
$S = \{(x_1, y_1), \dots, (x_n, y_n)\}$ .	

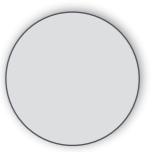
Living area (feet <sup>2</sup> )	#bedrooms	Price (1000\$)
2104	3	400
1600	3	330
2400	3	369
1416	2	232
3000	4	540
$\vdots$	$\vdots$	$\vdots$



$$y_i = f(x_i) + \sigma \varepsilon_i, \quad \sigma > 0$$

e.g.  $f(x) = w^T x, \quad \varepsilon_i \sim N(0, 1)$

# Sample Applications: Spam Filtering



True loneliness is when you don't even receive spam

[ simula ]

Mail ▾

COMPOSE

Inbox (2)

Starred

Important

Sent Mail

Drafts (233)

Administration/Apa...

Administration/Boo...

Administration/Har...

Administration/Lan...

Administration/VN

Administration/Slm...

Administration Slm...

Valeriya

John Vass

SEARCH

Ham

1–50 of 2,489 < > ⚙️

From	Subject	Date
Kristin McLeod	Hiking/cabin weekend August 20-21 - Hey gang Hermenegild, Johannes, and I were thinking of going for a weekend hiking/c	10:21
Airbnb	Reservation reminder - August 13, 2016 - Airbnb Pack your bags! It's almost time for your trip to \ <a href="#">Modify reservation</a>	10:18
H2020 - IKT	Don't miss this autumn's interesting events - H2020 ICT 10/2016 - H2020 ICT - National Contact Points - NCP Hvis du ik	5 Aug
Andy Edwards	Fwd: Your reservation: Aug 7 through 12 - VRBO.com #857936 - Google maps, what a tremendous piece of software. Forwa	4 Aug
Sigurd, me (2)	Report: Wednesday and Thursday (08.03 and 08.04) - Hi Sigurd, Thanks for the update and your hard work (so late)! It look	4 Aug
Viviane, me (2)	Awesome Summer School Pool Party! - Hi Viviane, Thanks a lot for the invitation, which I gladly accept :) Looking forward to	4 Aug
Massimo, me (6)	Funny. - I fully agree :) I think I need at least one year to decide for myself, finish current work, see how it works with the gro	4 Aug
me, Kereta (6)	Projects/FRIPRO/VN Reimbursement - Dear Zeljko, Great, many thanks for the response and your wishes! I am very happy \	4 Aug
Inger, me (2)	EU news item - Thanks Karoline for your kind email and changing the news! Best Valeriya On 04 Aug 2016, at 15:22, Inger h	4 Aug
Jan, me, Julio (7)	[Reminder] CanPathPro - WP5 Skype meeting - Dear all, this is a brief reminder regarding our regular Skype meeting. As \	4 Aug
Lars, me (2)	Crash Course on Machine Learning - Dear Lars, Yes for sure! You are more than welcome to join. Best Valeriya Dr. Valeriya	4 Aug

[ simula ]

in:spam

Mail ▾

COMPOSE

Inbox

Starred

Important

Sent Mail

Drafts (233)

Administration/Apa...

Administration/Boo...

Administration/Har...

Administration/Lan...

Administration/VN

Administration/Slm...

Administration Slm...

Valeriya

John Vass

SEARCH

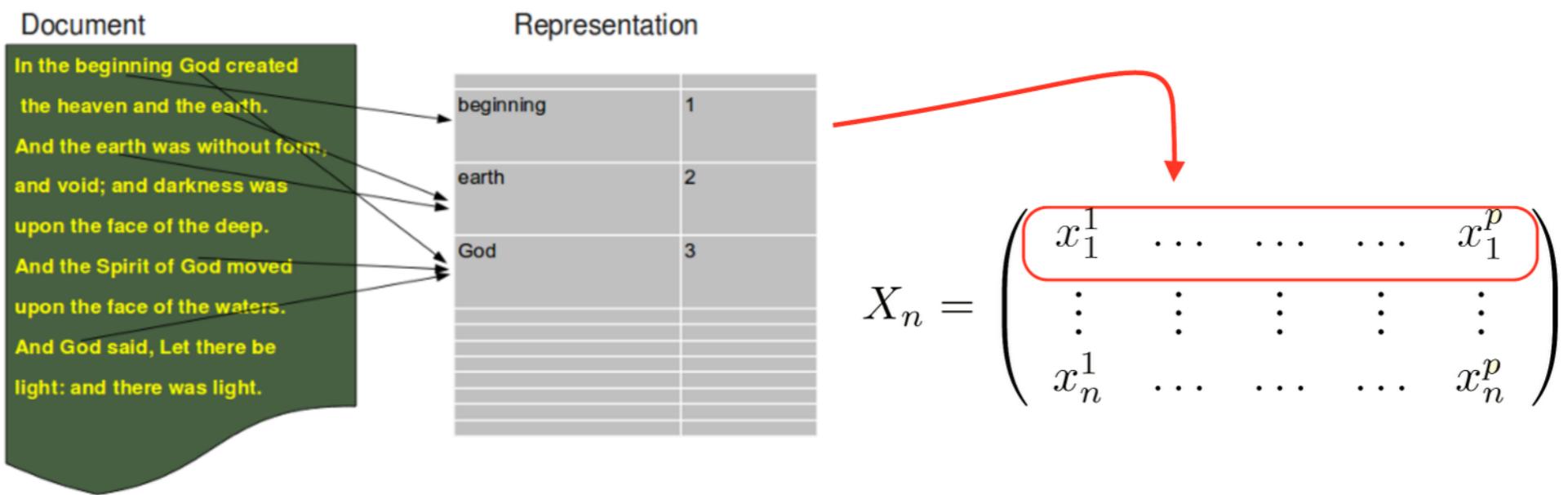
Spam

1–12 of 12 < > ⚙️

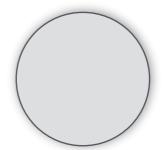
Delete all spam messages now (messages that have been in Spam more than 30 days will be automatically deleted)

From	Subject	Date
Markus Stahl	Mehr als 448 Euro pro Tag verdienen? HEUTE noch möglich - Sehr geehrter Interessent, wir möchten Sie mit dieser Em	03:39
Laura Schubert (4)	RE: Kundendienst schuldet dir Geld - Hallo valeriya@simula.no! Das Customer Service Team sendete mir deine Details \	5 Aug
Dr. Matthias Reich	Der unerwartete Geldsegen - Guten Tag, Unverhofft kommt oft? Naja, in diesem Fall stimmt es zumindest. Klicken Sie hier	5 Aug
Google Trading Inc. (5)	TAG 1 - Wettbewerb einladen for valeriya@simula.no! - Hallo valeriya@simula.no, TAG 1 competition.jpeg Ich habe eine	5 Aug
Linnea	Symposium 2016 - Invitation Letter from China - Symposium 2016 - Invitation Letter from China Dear Professor/Researc	5 Aug
Mr.Brian William.	Attention. - BG Group PLC Thames Valley Park, Reading, RG6 1PT, Berkshire, United Kingdom. Attention. I am Mr. Brian V	5 Aug
Sven Lindholm	Recently posted academic job vacancies at Educaloxy - Dear Colleague, We are pleased to present you our specialised	5 Aug
Google Benachrichtigung	Wettbewerb einladen for valeriya@simula.no! - IEMB16051 Hallo valeriya@simula.no, Ich habe eine Woche gewartet oh	4 Aug
DNB Bank ASA	Varsling! - Kjære kunde, Du har mottatt en intern melding. Klikk her Vennlig hilsen DNB ASA	4 Aug
giuseppelunardi	RE: valeriya	2 Aug

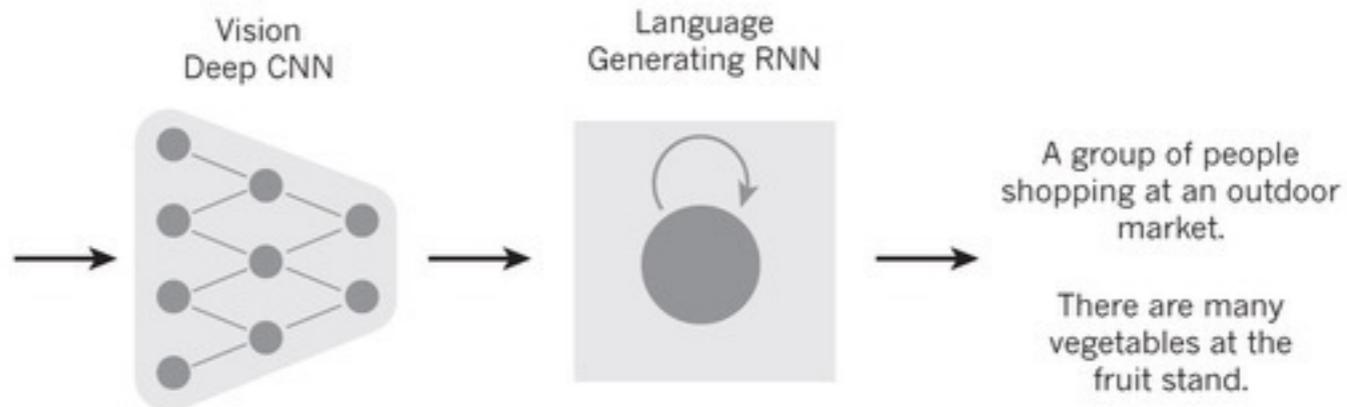
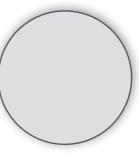
# Sample Applications: Text Classification



# Sample Applications: Object Recognition



# Sample Applications: from Image to Text



A group of people shopping at an outdoor market.

There are many vegetables at the fruit stand.



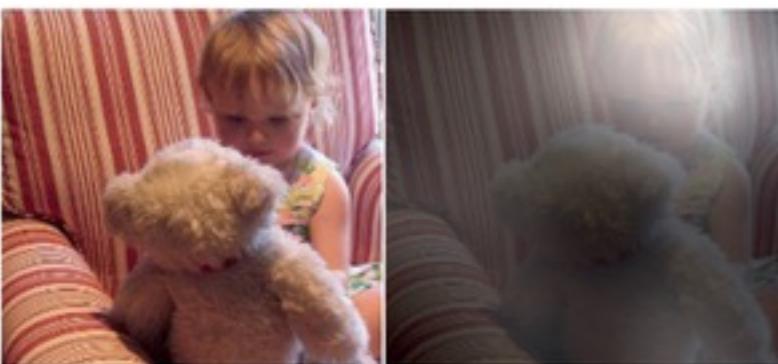
A woman is throwing a **frisbee** in a park.



A **dog** is standing on a hardwood floor.



A **stop** sign is on a road with a mountain in the background



A little girl sitting on a bed with a teddy bear.

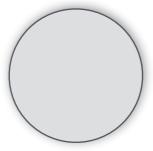


A group of **people** sitting on a boat in the water.



A giraffe standing in a forest with **trees** in the background.

# Sample Applications



Web ranking

Social Networks

Debugging

Finance

Computational Biology

Mobile technology

E-Commerce

Robotics

Space exploration

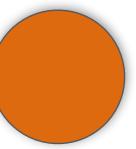
Navigation

...[Your favourite area]

Tens of thousands of machine learning algorithms

Hundreds new every year

# Programming with Data / Data Analysis



★ Want adaptive robust and fault tolerant systems

★ Rule-based implementation is (often)

- ▶ difficult (for the programmer)
- ▶ brittle (can miss many edge-cases)
- ▶ becomes a nightmare to maintain explicitly
- ▶ often doesn't work too well (e.g. OCR)

We say that a program for performing a task has been acquired by learning if it has been acquired by any means other than explicit programming.

*Valiant, 1984*

★ Usually easy to obtain examples of what we want **IF x THEN DO y**

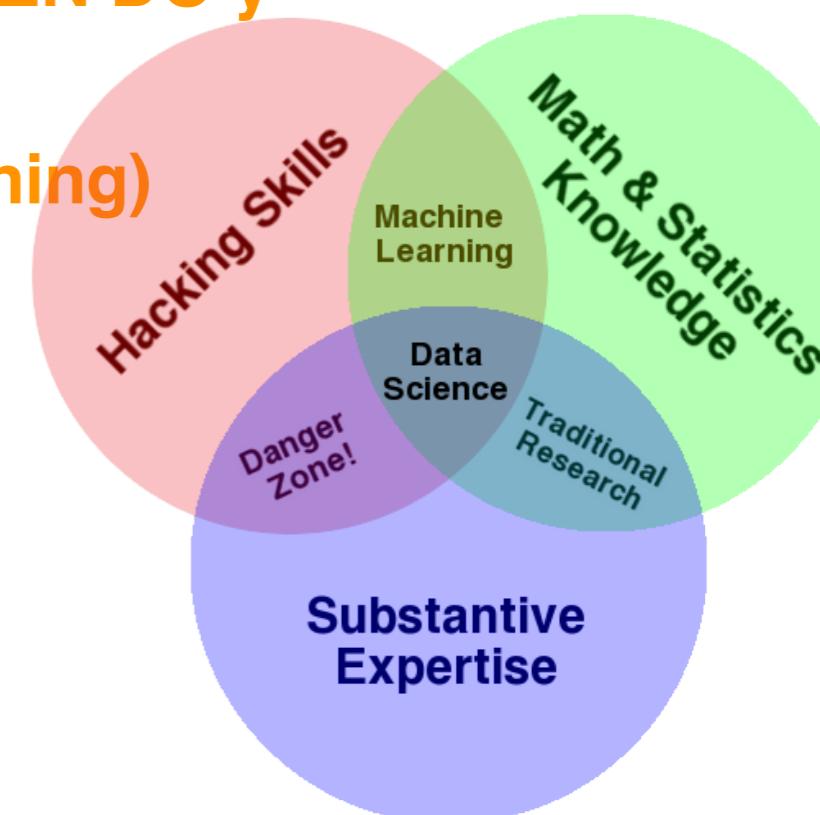
★ Collect many pairs  $\{(x_i, y_i)\}_{i=1}^N$

★ Estimate function such that  $f(x_i) = y_i$  (**supervised learning**)

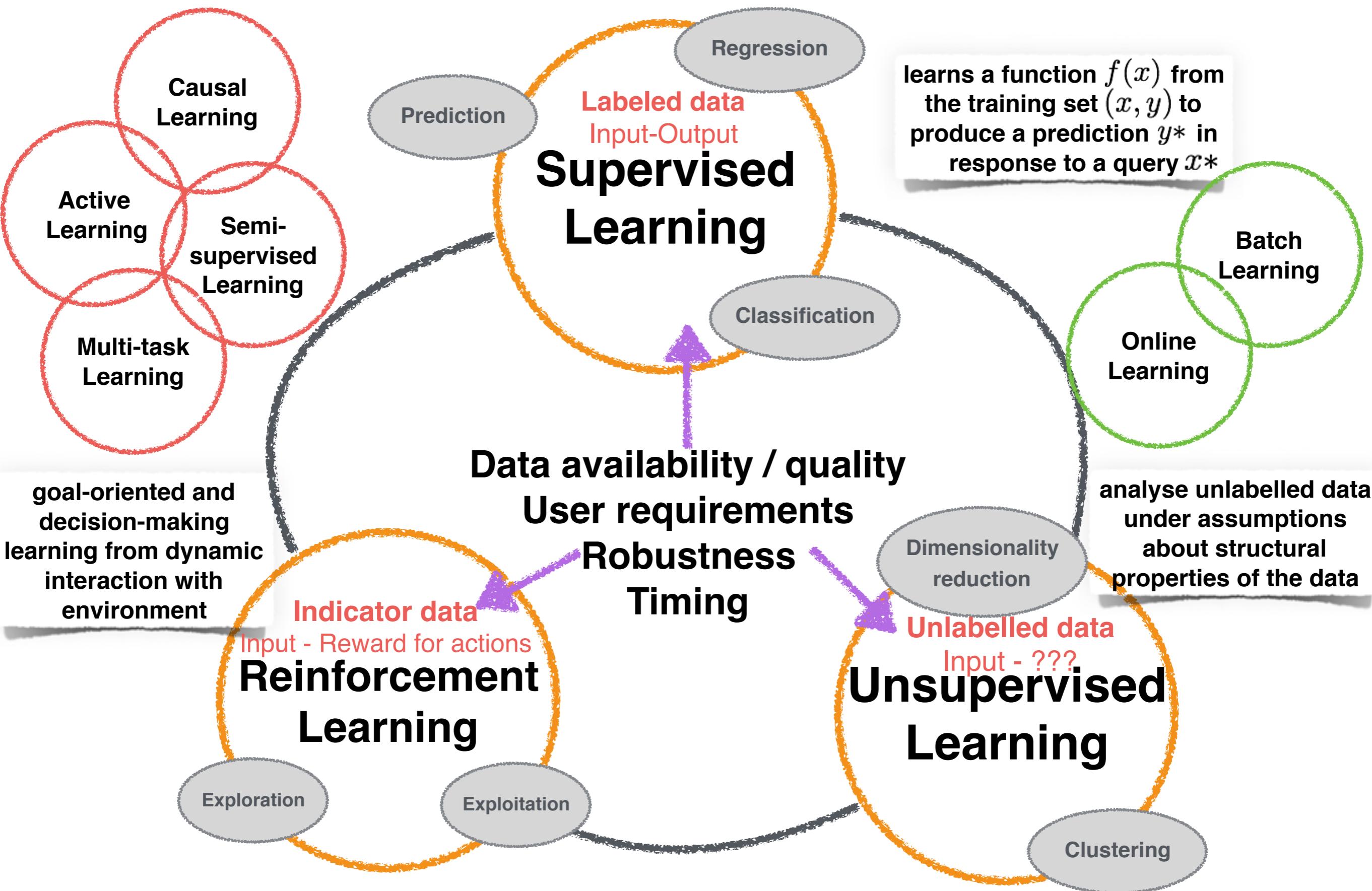
★ Detect patterns in data (**unsupervised learning**)

... learning from examples, refers to systems that are trained instead of programmed with a set of examples, that is, a set of input/output pairs...

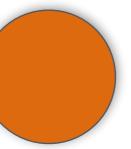
*Poggio and Smale, 2003*



# Three paradigms of learning



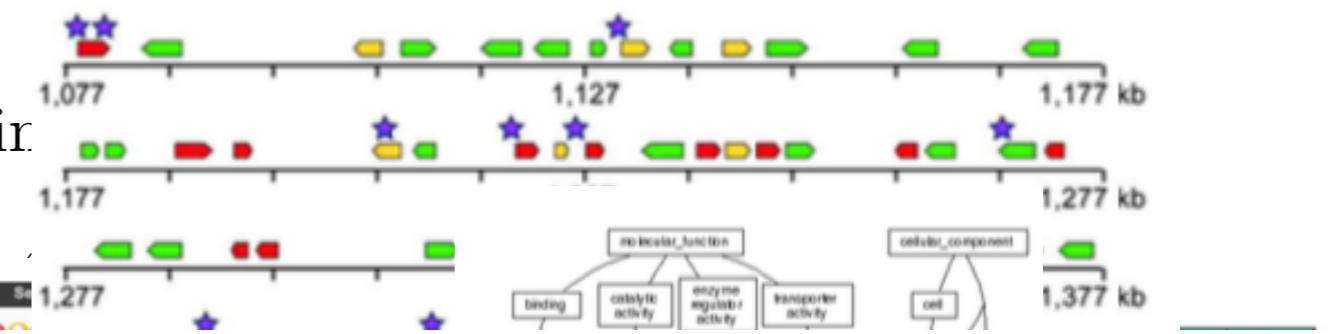
# Problem Prototypes: Supervised Learning $y = f(x)$



We are given a set of input-output pairs (training set)  $\{(x_i, y_i)\}_{i=1}^n$

- ▶ **Binary classification:** Given  $x$  find  $y$  in  $\{-1, 1\}$
- ▶ **Multi-category (Multi-class) classification:** Given  $x$  find  $y$  in  $\{1, \dots, K\}$
- ▶ **Regression:** Given  $x$  find  $y$  in  $\mathbb{R}$  (or  $\mathbb{R}^d$ )
- ▶ **Sequence annotation:** Given sequence  $x_1, \dots, x_k$  find  $y_1, \dots, y_l$
- ▶ **Hierarchical categorisation (Ontology):**

Given  $x$  find a point in



- ▶ **Prediction:** Given  $x_t$  and  $y_{t-1}, \dots, y_1$  find  $y_t$

Forms of mapping  $f$ :

- ▶ decision trees / forests
- ▶ logistic regression,
- ▶ support vector machine
- ▶ neural networks,
- ▶ kernel machines,
- ▶ Bayesian classifiers.

given sequence

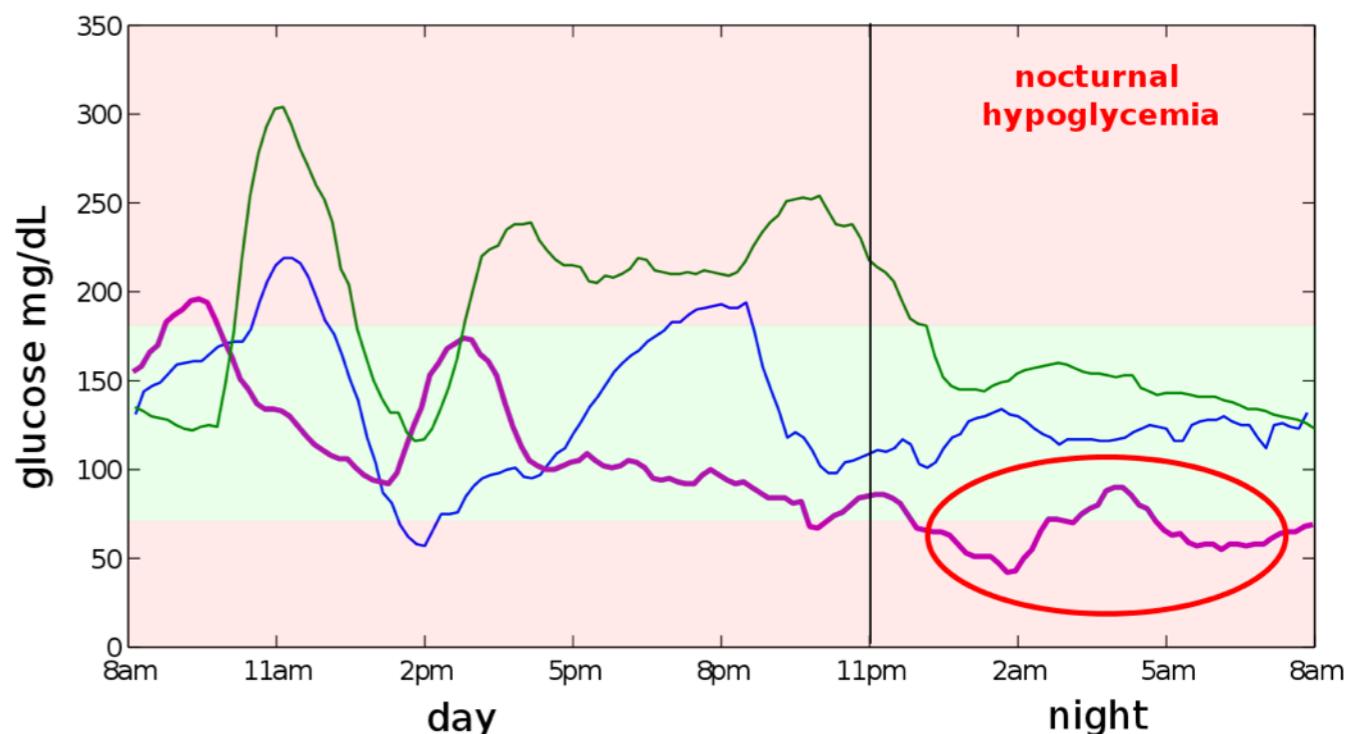
gene finding

Map image to digitization

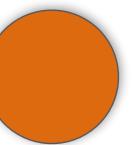
activity segmentation

Often with loss function

$$V(y, f(x))$$

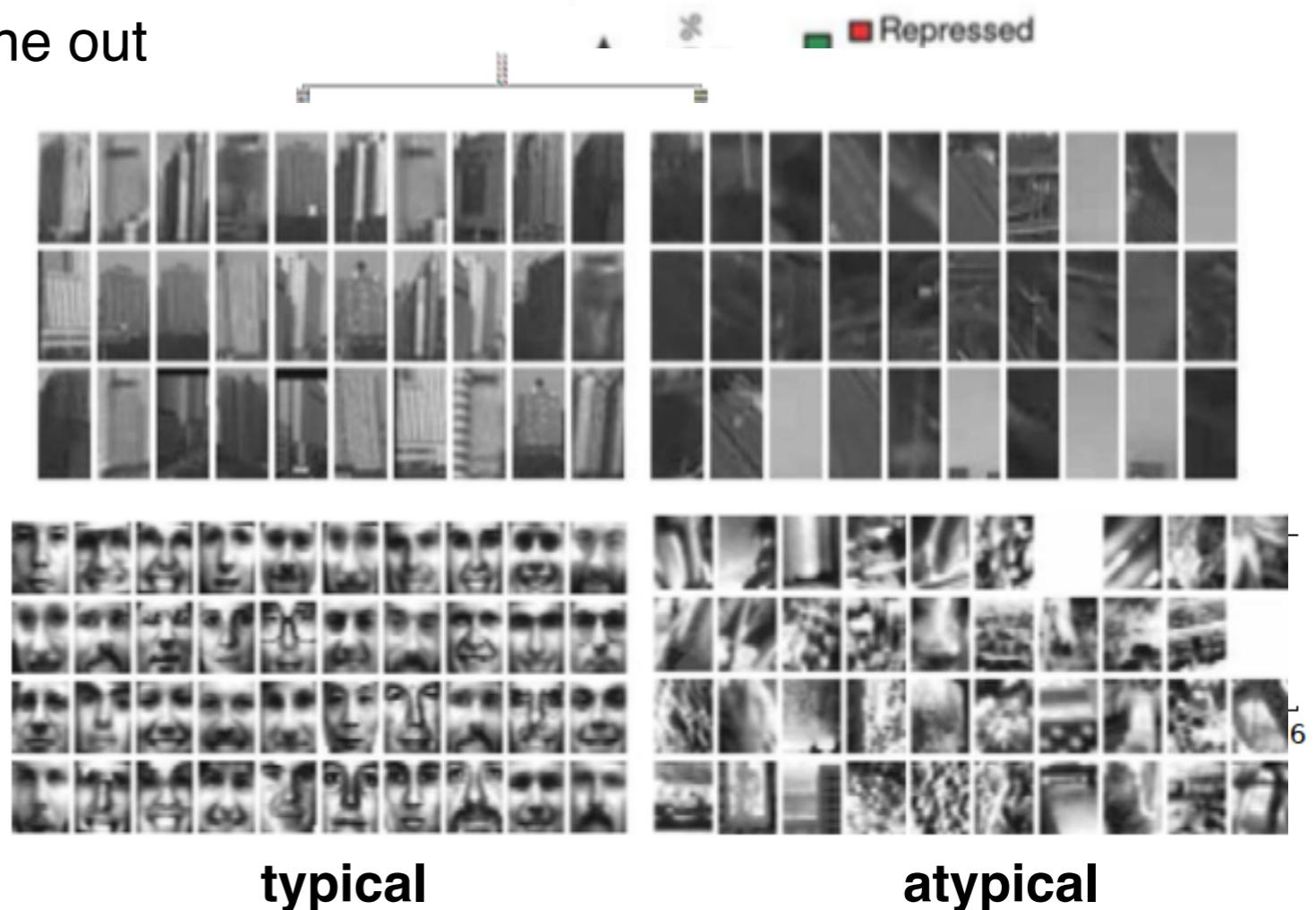


# Problem Prototypes: Unsupervised Learning

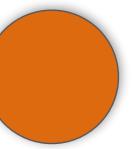


Given data  $x$ , ask a good question.... about  $x$  or about model for  $x$

- ▶ **Clustering:** Find a set of prototypes representing the data
- ▶ **Principal components:** Find a subspace representing the data
- ▶ **Sequence analysis:** Find a latent causal sequence for observations  
(Kalman Filter, Hidden Markov Model)
- ▶ **Hierarchical representation**
- ▶ **Dictionary learning:** Find (small) set of factors for observation
- ▶ **Novelty detection:** Find the odd one out



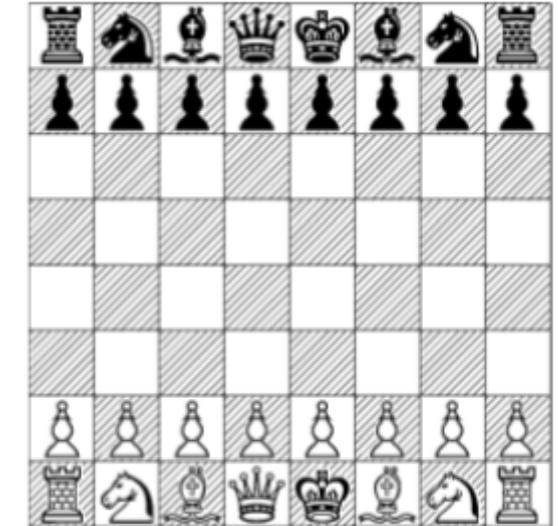
# Problem Prototypes: Reinforcement Learning



Examples of desired behaviour are not available but it is possible to score examples of behaviour according to some performance measure

## Sequential decision tasks:

- ▶ Take action
- ▶ Environment responds
- ▶ Observe stuff
- ▶ Update model
- ▶ Repeat



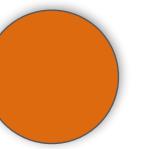
## Bandits: Pick arm, get reward, pick new arm

- ▶ Choose an option
- ▶ See what happens (get reward)
- ▶ Update model
- ▶ Choose next option

Reinforcement learning  
combines search and  
long-term memory

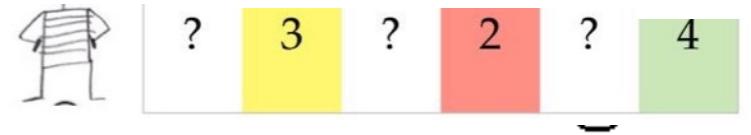
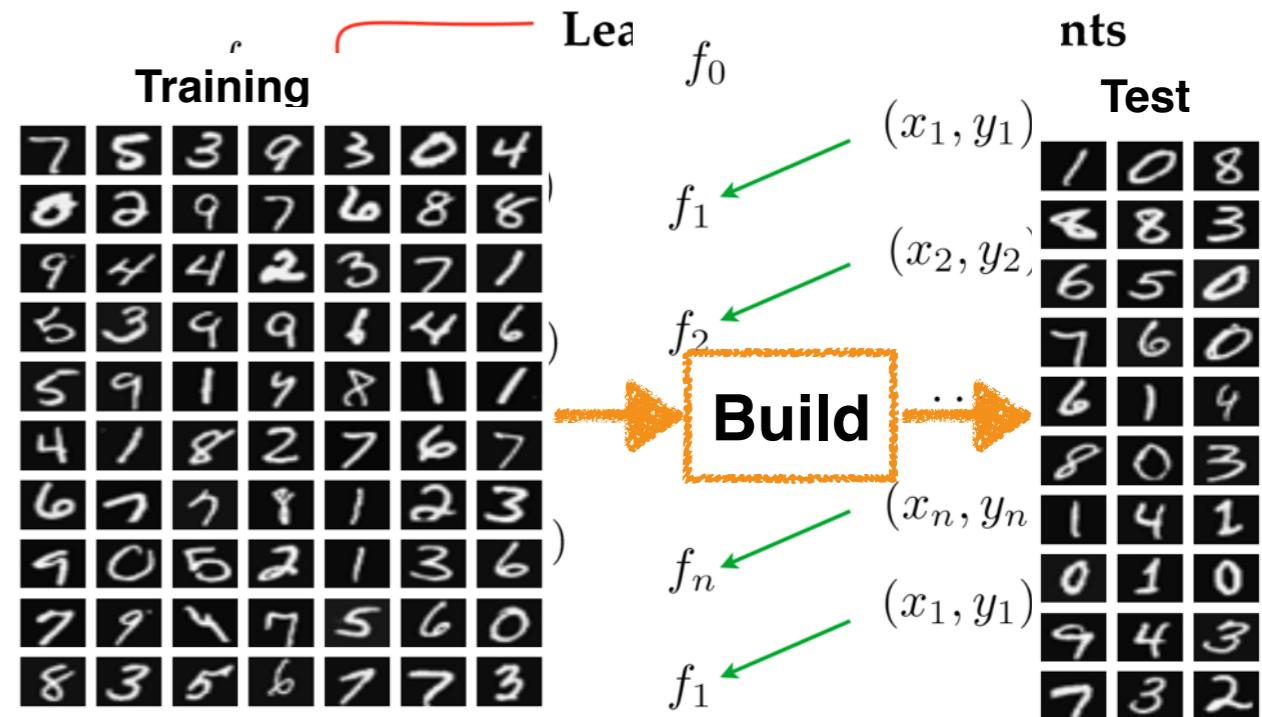
Differences to  
Supervised learning?



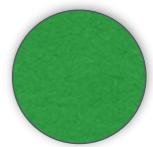


# Problem Prototypes: Other types

- ▶ **Induction:** Reasoning from observed training cases to general rules, which are then applied to the test cases
- ▶ **Transductive inference:** Reasoning from observed, specific (training) cases to specific (test) cases
- ▶ **Multi-task learning:** Learns a problem together with other related problems at the same time, using a shared representation.
- ▶ **Active learning:** Interactive querying of the user to obtain the desired outputs at anew data points
- ▶ **Online learning:** Analyse each training example as it is presented
- ▶ **Batch learning:** Collect training examples, analyse them, output an hypothesis



# Key Issues in Machine Learning

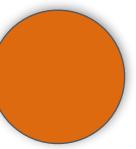


- ▶ **What are good hypothesis spaces?** Which spaces have been useful in practical applications and why?
- ▶ **What algorithms can work with these spaces?** Are there general design principles for machine learning algorithms?
- ▶ **How can we optimise accuracy on future data points?** This is sometimes referred to “generalisation” ability or problem of overfitting.
- ▶ **How can we have confidence in the results?** (*the statistical question*)
- ▶ **Are some learning problems computationally intractable?**  
*(the computational question)*
- ▶ **How can we formulate application problems as machine learning problems?** (*the engineering question*)

**Hypothesis:** a function produced by a learning algorithm, which is believed to be close to the true function.

**Hypotheses space:** the space of all hypotheses that can be an output by a learning algorithm.

# Part 0. Settings for Supervised Learning



The goal of **supervised learning** is to find an underlying input-output relation

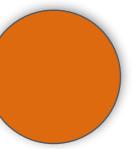
$$f(x_{new}) \sim y$$

given data.

The data, called **training set**, is a set of  $n$  input-output pairs,

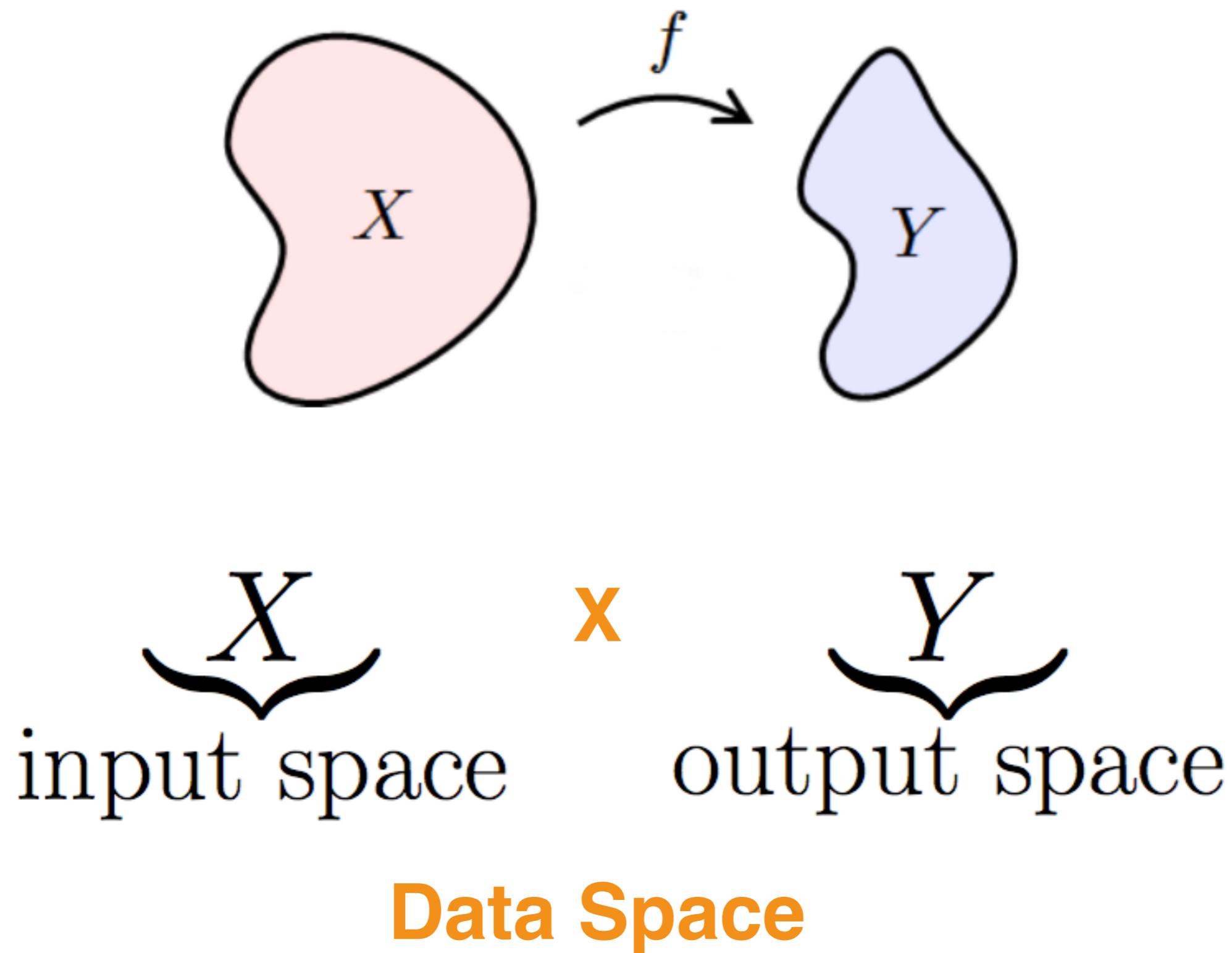
$$S = \{(x_1, y_1), \dots, (x_n, y_n)\}.$$

# We Need a Model to Learn

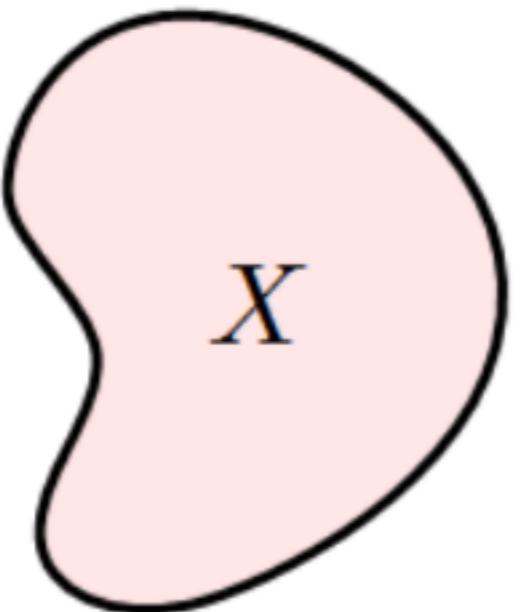


- ▶ We consider the approach to machine learning based on the *learning from examples* paradigm.
- ▶ **Goal:** Given the training set, learn a corresponding I/O relation.
- ▶ We have to postulate the existence of a model for the data.
- ▶ The model should take into account the possible *uncertainty* in the task and in the data.

# Data Space



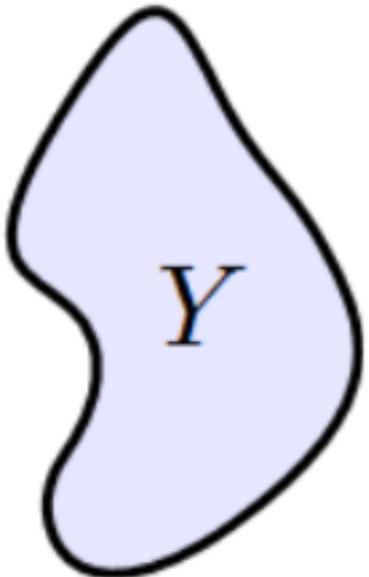
# Examples of Input Space



$X$   
input space

- ▶ linear spaces, e. g.
  - vectors,
  - functions,
  - matrices/operators
- ▶ “structured” spaces, e. g.
  - strings,
  - probability distributions,
  - graphs

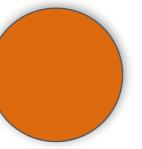
# Examples of Output Space



$Y$   
output space

- ▶ linear spaces, e. g.
  - $Y = \mathbb{R}$ , regression,
  - $Y = \{+1, -1\}$ , classification
  - $Y = \mathbb{R}^T$ , multi-task regression
  - $Y = \{1, \dots, T\}$ , multi-label classification
  
- ▶ “structured” spaces, e. g.
  - strings,
  - probability distributions,
  - graphs

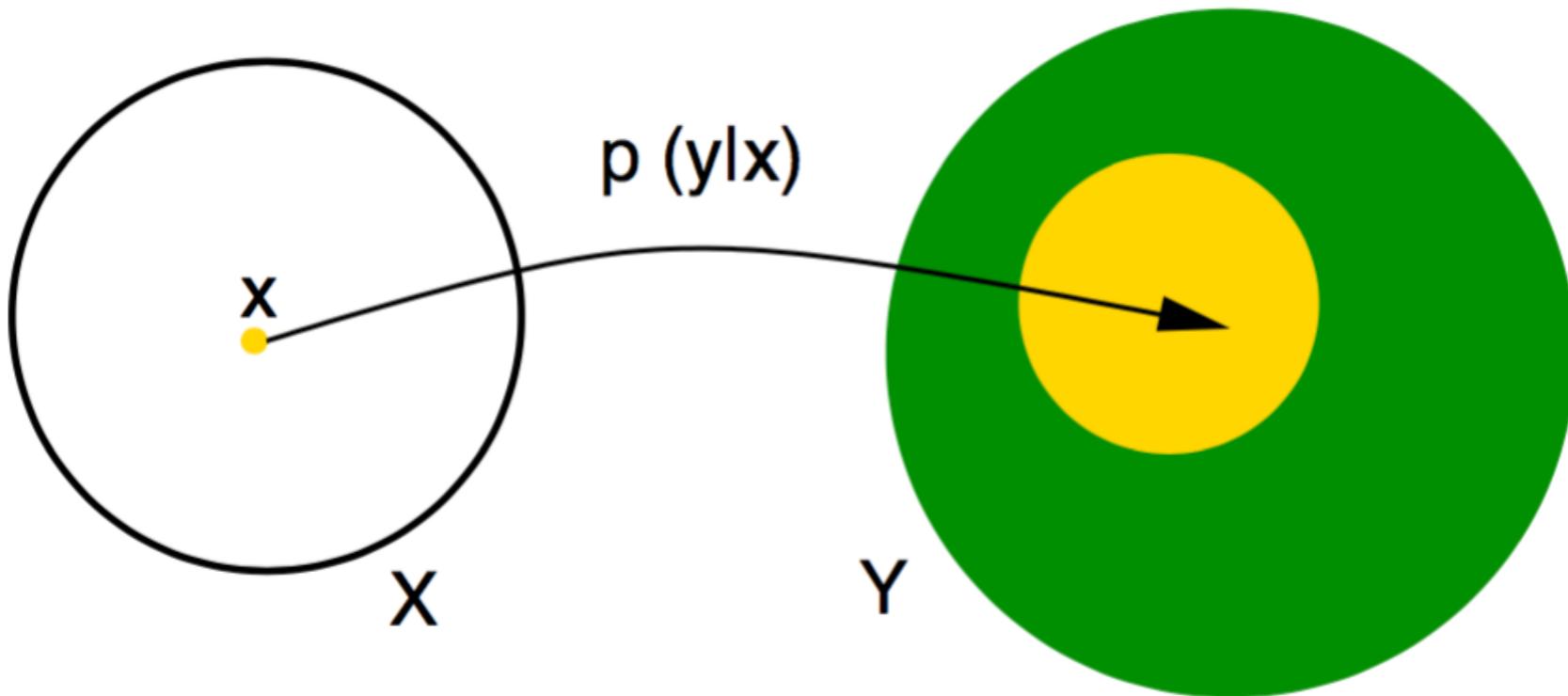
# Modelling Uncertainty in the Data Space



- ▶ **Assumption:** there exists a fixed unknown distribution  $\rho(x, y)$ , according to which the data are identically and independently sampled
- ▶ The distribution  $\rho$  models different sources of uncertainty
- ▶ **Assumption:**  $\rho$  factorizes as  $\rho(x, y) = \rho_X(x)\rho(y|x)$

# Marginal and Conditional

$\rho(y|x)$  can be seen as a form of noise in the output



For each input  $x$  there is a distribution of possible outputs  $\rho(y|x)$

The marginal distribution  $\rho_X(x)$  models uncertainty in the sampling of the input samples

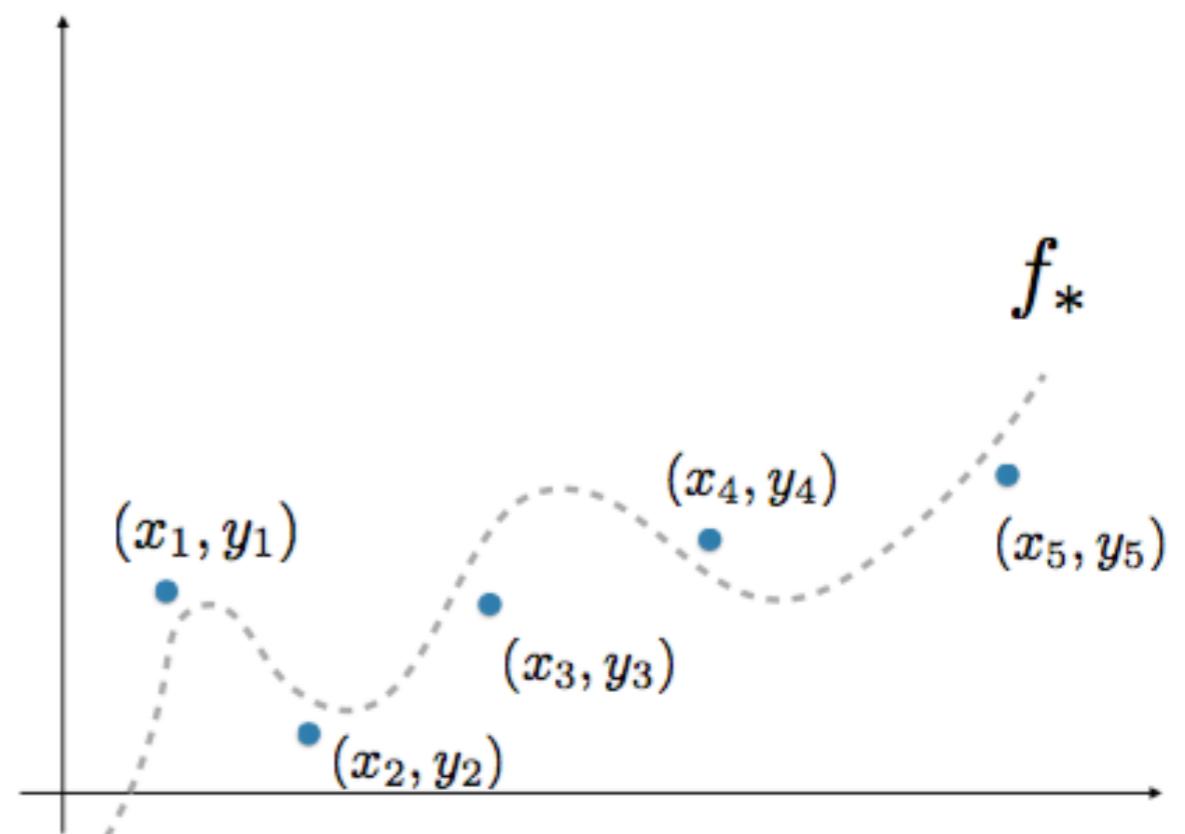
# Data Models

In **regression**, the following model is often considered

$$y = f^*(x) + \epsilon$$

where

- $f^*(x)$  fixed unknown (regression) function
- $\epsilon$  random noise, e.g., standard Gaussian



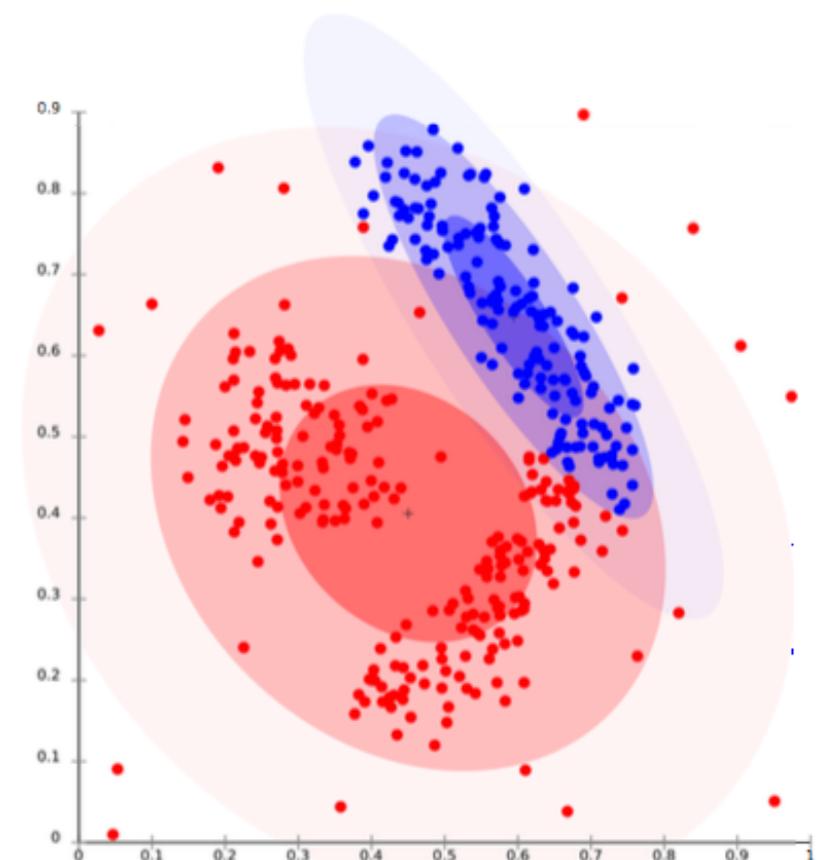
# Data Models

In **classification**,

$$\rho(1|x) = 1 - \rho(-1|x) \quad \forall x$$

Noiseless classification

$$\rho(1|x) = \{1, 0\} \quad \forall x$$



# Loss Function

**Goal of learning:** Estimate “best” I/O relation, not the whole  $\rho(x, y)$

We need to fix a measurable loss function

$$V : Y \times Y \rightarrow [0, \infty)$$

It is a point-wise error measure: cost of when predicting  $f(x)$  in place of  $y$

# Expected Risk and Target Function

The *expected loss* (or *expected risk*)

$$\mathcal{E}(f) = \mathbb{E}[V(y, f(x))] = \int \rho(x, y)V(y, f(x))dxdy$$

can be seen as a measure of the error on past as well as future data.

Given the loss function and a distribution, the ‘best’ I/O relation is the *target function*

$$f^* : X \rightarrow Y$$

that minimises the expected risk.

# Learning from Data. Learning Algorithms

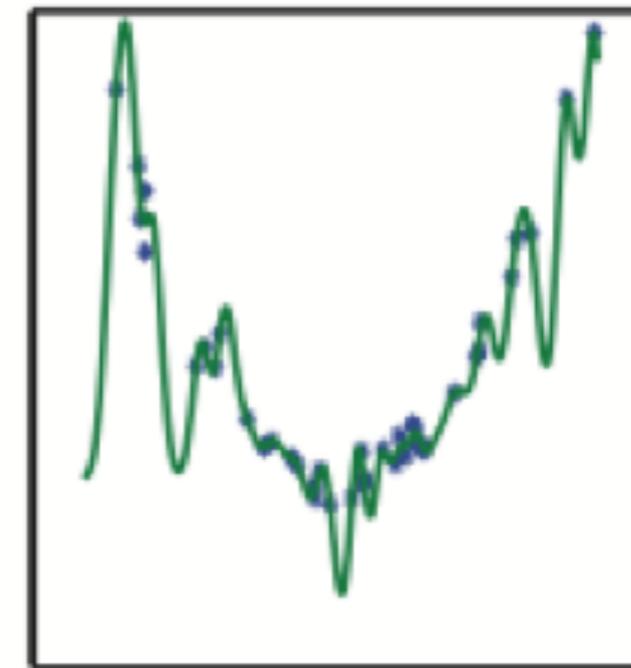
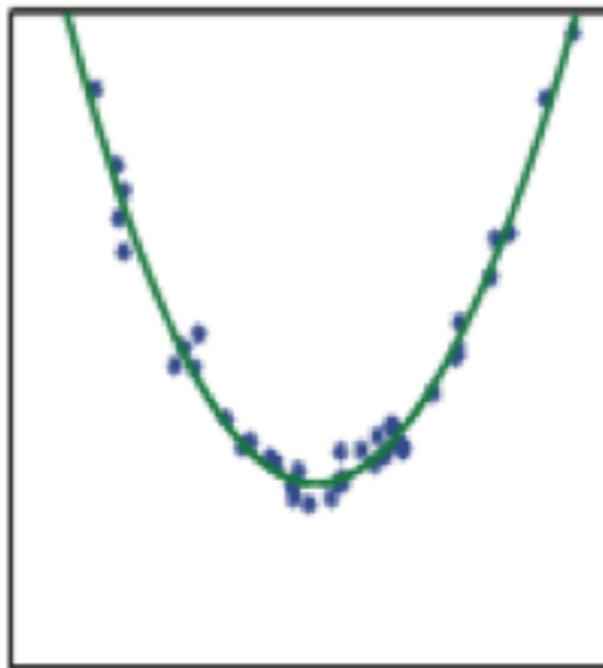
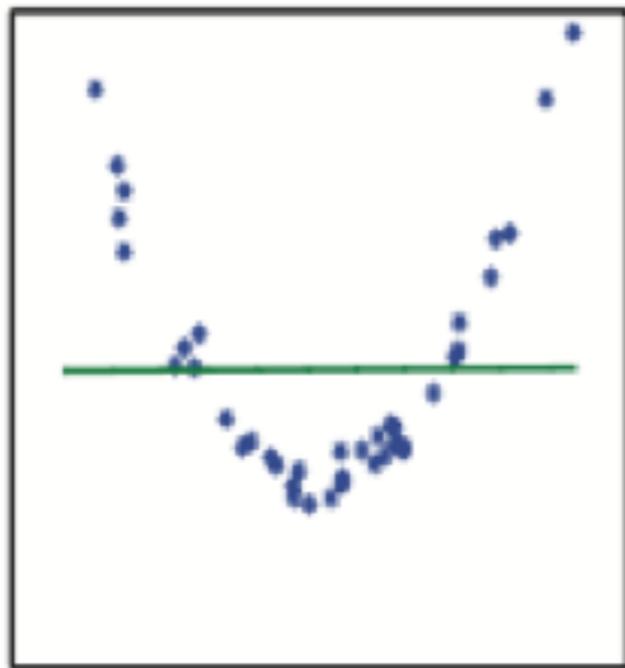
- ▶ The target function cannot be computed, since the probability is unknown.
- ▶ The goal of learning is to find an estimator of the target function from data.
- ▶ A learning algorithm is a procedure that given a training set  $S$  computes an estimator  $f_S$ .
- ▶ An estimator should mimic the target function, in which case we say that it **generalises**.

# Generalisation: Fitting and Stability

How to design a good algorithm?

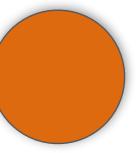
Two concepts are key:

- ▶ **Fitting**: an estimator should *fit* data well.
- ▶ **Stability**: an estimator should be stable, it should not change much if data change slightly.



- ▶ We say that an algorithms **overfits**, if it fits the data while being unstable.
- ▶ We say that an algorithms **oversmooth**, if it is stable while disregarding the data.

# Part I. Local Methods / Bias-Variance

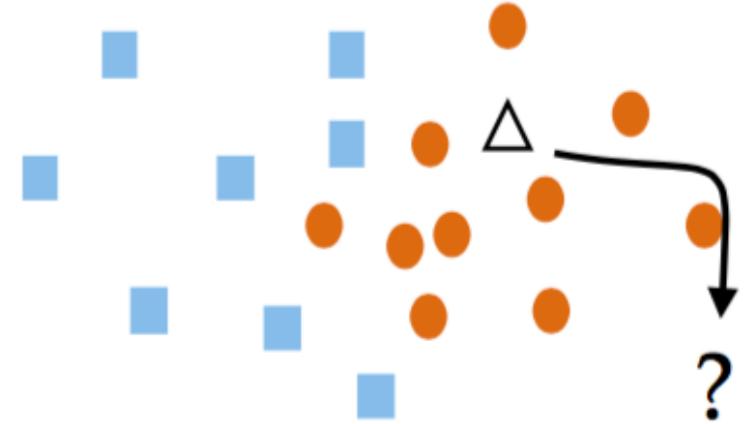


## GOAL:

Discuss the fundamental concept of bias-variance trade-off starting from simple machine learning approaches to understand parameter tuning

# Local Methods: Nearest Neighbour

Nearby points have similar labels



Given a training set

$$S = \{(x_1, y_1), \dots, (x_n, y_n)\} \text{ with } x_i \in \mathbb{R}^D, y_i \in \mathbb{R}.$$

and a new input  $\bar{x}$ , let

$$i' = \arg \min_{i=1, \dots, n} \|\bar{x} - x_i\|^2$$

In general,

$$i' = \arg \min_{i=1, \dots, n} d(\bar{x}, x_i)$$

and define the nearest neighbour (NN) estimator as

$$\hat{f}(\bar{x}) = y_{i'}$$

Can we do any better?

**Computational cost:**  $O(nD)$

We compute  $n$  times the distance  $\|\bar{x} - x_i\|$  which costs  $O(D)$

# Local Methods: K-Nearest Neighbours

A new point is the mean values of the K closest points

Consider

$$d_{\bar{x}} = (\|\bar{x} - x_i\|^2)_{i=1}^n$$

the array of distances of a new point  $\bar{x}$  to the input points in the training set. Let

$$s_{\bar{x}}$$

be the above array sorted in increasing order and

$$I_{\bar{x}}$$

the corresponding vector of indices,

$$K_{\bar{x}} = \{I_{\bar{x}}^1, \dots, I_{\bar{x}}^K\}$$

be the array of the first K entries of  $I_{\bar{x}}$ . The K-nearest neighbours estimator (KNN) is defined as,

$$\hat{f}(\bar{x}) = \sum_{i' \in K_{\bar{x}}} y_{i'}$$

**Computational cost:**  $O(nD + n \log n)$

We compute the n distance for each point and order them

# Local Methods. Remarks

*Generalisation I:* closer points should count more (*Parzen windows*)

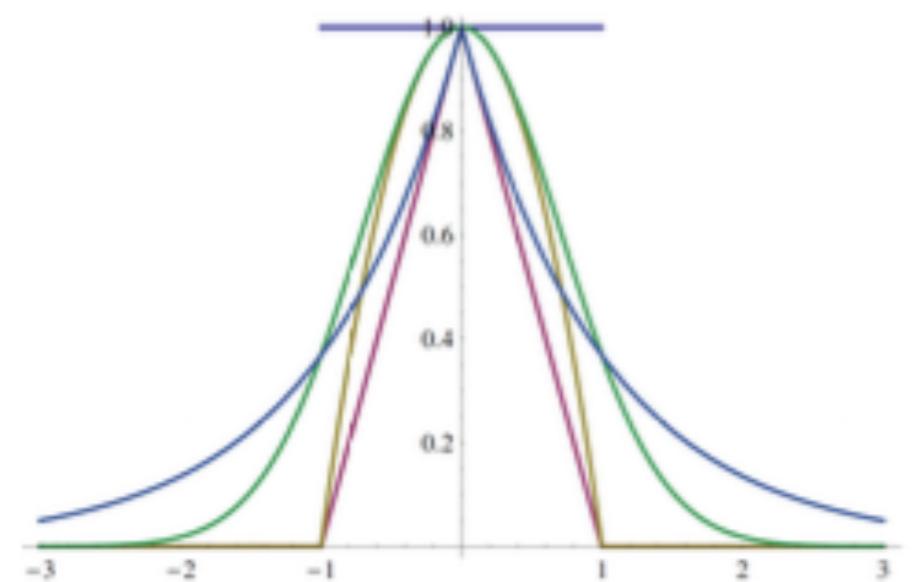
$$\hat{f}(\bar{x}) = \frac{\sum_{i=1}^n y_i k(\bar{x}, x_i)}{\sum_{i=1}^n k(\bar{x}, x_i)}$$

where  $k$  is a similarity function

$$k(x, \bar{x}) \geq 0 \text{ for all } x, \bar{x} \in \mathbb{R}^D;$$

$$k(x, \bar{x}) \rightarrow 1 \text{ when } \|x - \bar{x}\| \rightarrow 0;$$

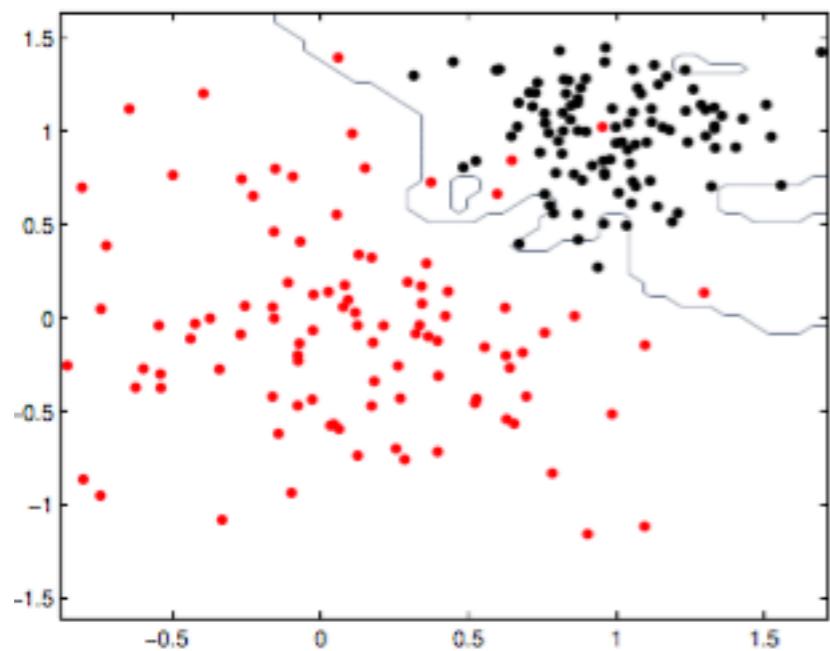
$$k(x, \bar{x}) \rightarrow 0 \text{ when } \|x - \bar{x}\| \rightarrow \infty.$$



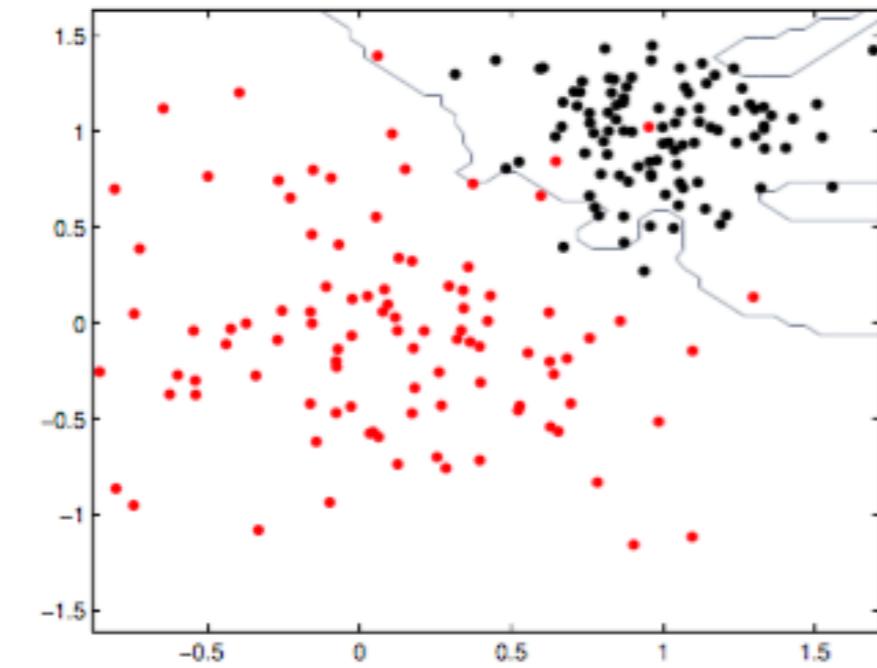
*Generalisation II:* other metric / similarities

$$X = \{0, 1\}^D \quad d_H(x, \bar{x}) = \frac{1}{D} \sum_{j=1}^D \mathbf{1}_{[x^j \neq \bar{x}^j]}$$

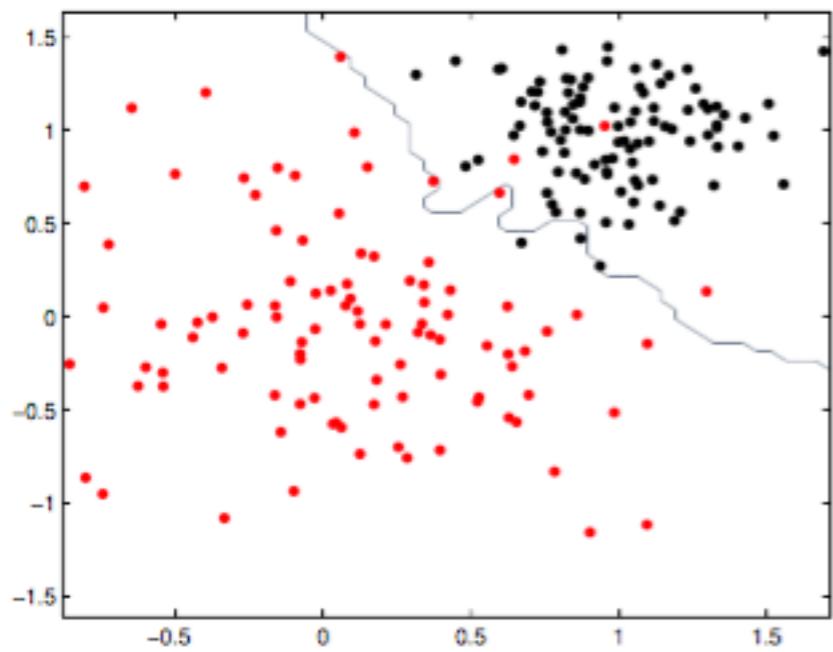
# Local Methods



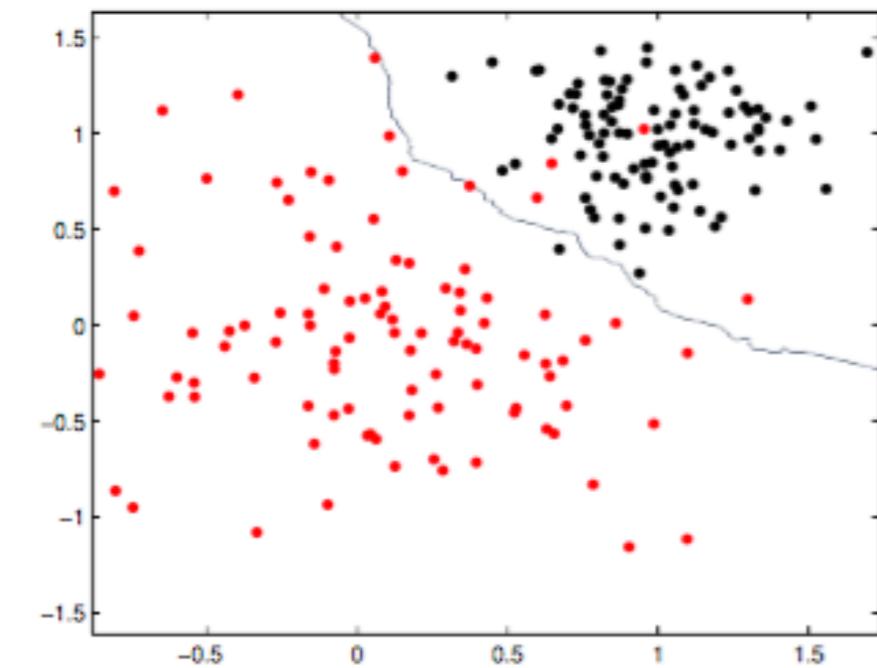
1-Nearest Neighbour



2-Nearest Neighbours



5-Nearest Neighbours



9-Nearest Neighbours

**How can we choose it?**

**Is there an optimal value?**

**Can we compute it?**

# Model Selection Bias-Variance Trade-Off

Is there an optimal parameter?

Ideally we would like to choose  $K$  that minimises the expected error

$$\mathcal{E}_K = \mathbf{E}_S \mathbf{E}_{x,y} (y - \hat{f}_K(x))^2.$$

Optimal hyperparameter  $K^*$  should minimise  $\mathcal{E}_K$

$$K^* = \arg \min_{K \in \mathbb{K}} \mathcal{E}_K$$

Ideally! (In practice we don't have access to the distribution)

- ▶ We can still try to understand the above minimisation problem: **Does a solution exists? What does it depend on?**
- ▶ Yet, ultimately, we need something we can compute!

Next: Characterise corresponding minimisation problem to uncover one of the **most fundamental aspect of machine learning**.

# Example: Regression

Regression model  $y_i = f^*(x_i) + \delta_i$

$$\mathbf{E}\delta = 0, \mathbf{E}\delta_i^2 = \sigma^2, i = 1, \dots, n$$

**Fix  $\mathbf{x}$**

Now  $\mathcal{E}_K = \mathbf{E}_x \mathbf{E}_S \mathbf{E}_{y|x} (y - \hat{f}_K(x))^2 \rightarrow \mathbf{E}_S \mathbf{E}_{y|x} (f^*(x) + \delta - \hat{f}_K(x))^2$   
that is

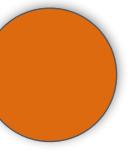
$$\mathcal{E}_K(x) = \mathbf{E}_S (f^*(x) - \hat{f}_K(x))^2 + \sigma^2$$

...

# Bias-Variance Trade-Off for K-NN

Define the noiseless K-NN (it is ideal!)

$$\mathbf{E}_{y|x} \hat{f}_K(x) = \frac{1}{K} \sum_{l \in K_x} f^*(x_l).$$



# Bias-Variance Trade-Off for K-NN

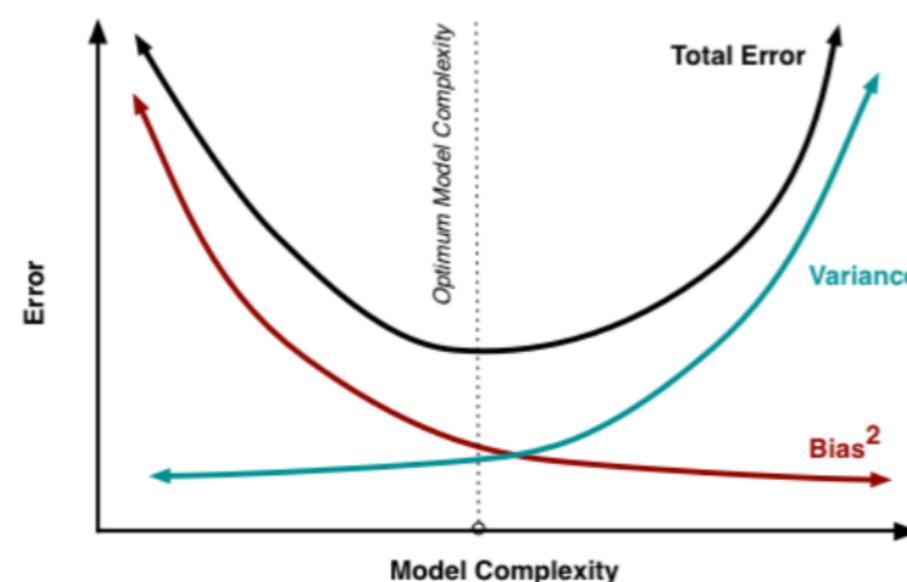
Define the noiseless K-NN (it is ideal!)

$$\mathbf{E}_{y|x} \hat{f}_K(x) = \frac{1}{K} \sum_{l \in K_x} f^*(x_l).$$

Consider

$$\begin{aligned} \mathcal{E}_K(x) &= \mathbf{E}_S \mathbf{E}_{y|x} (f^*(x) - \hat{f}_K(x))^2 = \\ &= \underbrace{(f^*(x) - \mathbf{E}_S \mathbf{E}_{y|x} \hat{f}_K(x))^2}_{bias} + \underbrace{\mathbf{E}_S \mathbf{E}_{y|x} (\mathbf{E}_{y|x} \hat{f}_K(x) - \hat{f}_K(x))^2}_{variance} + \sigma^2 \\ &\quad \text{+ } \frac{\sigma^2}{K} \end{aligned}$$

Is there an  
optimal value?  
Yes



Can we compute it?

# Model Selection Bias-Variance Trade-Off

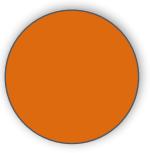
**Bias-Variance Trade-Off** is theoretical but shows that:

- ▶ an optimal parameter exists and
- ▶ it depends on the noise and the unknown target function.

How to choose  $K$  in practice? **Cross Validation!!!**

- ▶ **Idea:** train on some data and validate the parameter on new unseen data as proxy for the ideal case.

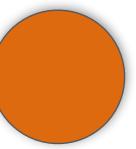
# Cross-Validation Flavours



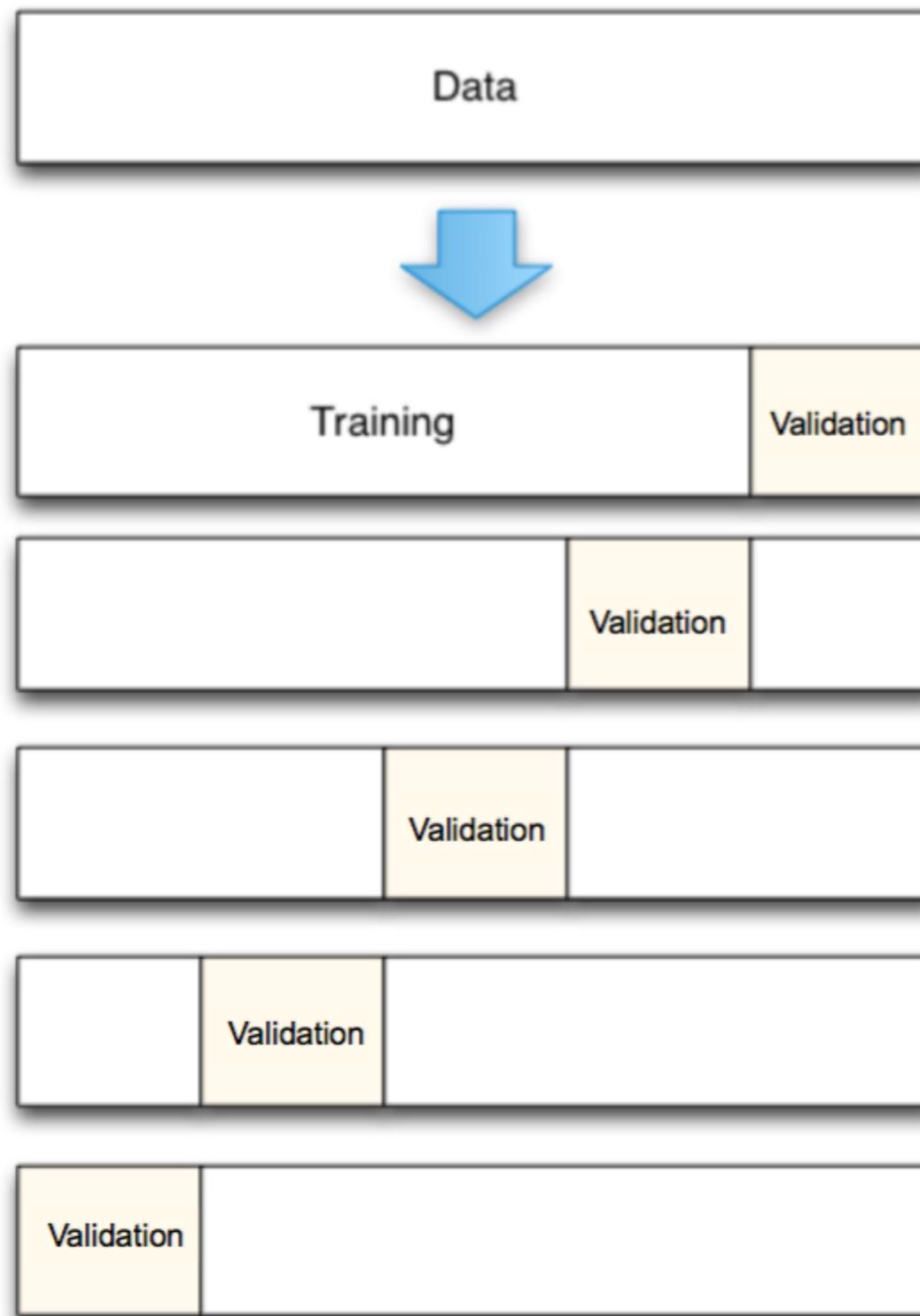
## Hold-out



# Cross-Validation Flavours



## K-Fold



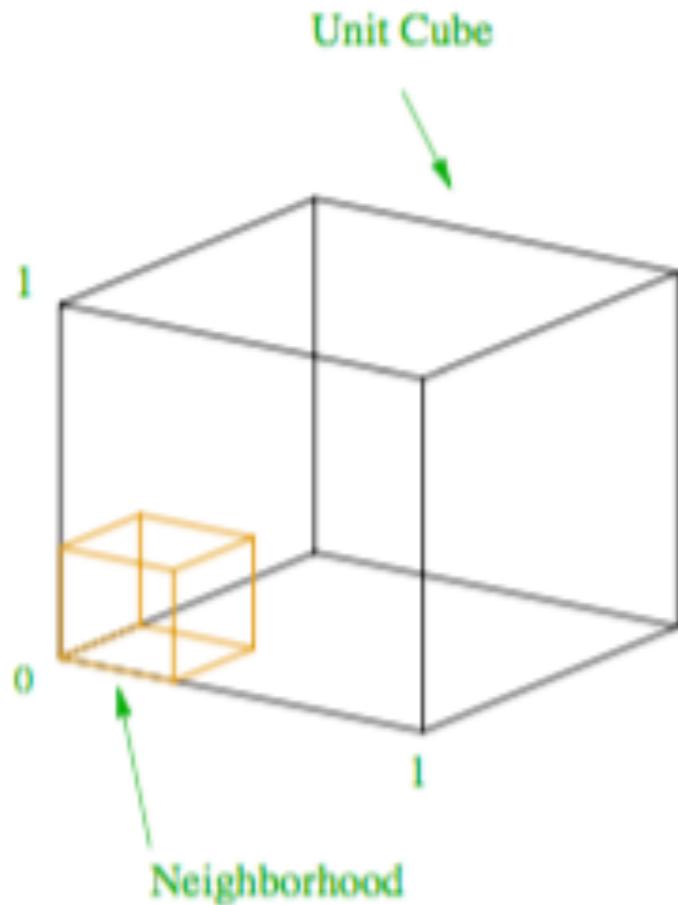
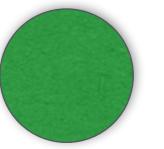
# End of Part I

We made our first encounter with learning algorithms (local methods) and the problem of tuning their parameters (via bias-variance trade-off and cross-validation) to avoid overfitting and achieve generalisation.

**Stability - Overfitting - Bias/Variance - Cross-Validation**

**End of story?**

# Local Methods in High Dimensions



What is the length of the edge of  
a cube containing 1% of the  
volume of a cube with edge 1?

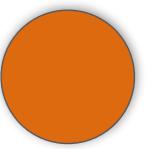
## Part II. Regularisation

### GOAL:

Introduce the basic (global) regularisation methods and study their computational aspects

**Going Global + Impose Smoothness**

# Regularisation. Expected Risk Minimisation



Of all the principles which can be proposed for that purpose, I think there is none more general, more exact, and more easy of application, that of which we made use in the preceding researches and which consists of rendering the **sum of square of the errors** a minimum.

Legendre, 1805

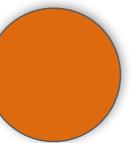
- ▶ **Empirical Risk Minimisation (ERM):** probably the most popular approach to design learning algorithms.
- ▶ General idea: considering the empirical error

$$\hat{\mathcal{E}}(f) = \frac{1}{n} \sum_{i=1}^n V(y_i, f(x_i))$$

as a proxy for the expected risk

$$\mathcal{E}(f) = \mathbf{E}[V(y, f(x))] = \int p(x, y) V(y, f(x)) dx dy$$

# Regularisation. Expected Risk Minimisation



- ▶  $V$  measures the price we pay predicting  $f(x)$  when the true label is  $y$
- ▶  $\mathcal{E}(f)$  cannot be directly computed, since  $p(x, y)$  is unknown

To turn the above idea into an actual algorithm, we need to

- ▶ Fix a suitable hypothesis space  $H$
- ▶ Minimise  $\hat{\mathcal{E}}(f)$  over  $H$

The space should allow *feasible computations* and be *rich*, since the complexity of the problem is not known a priori.

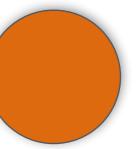
**Simplest example:** Space of linear functions

$$H = \{f : \mathbb{R}^D \rightarrow \mathbb{R} : \exists \omega \in \mathbb{R} \text{ such that } f(x) = x^T \omega, \forall x \in \mathbb{R}^D\}$$

- ▶ Each function  $f$  is defined by the vector  $\omega$ :  $f(x) = x^T \omega$

If  $H$  is rich enough, solving ERM may cause overfitting  
(solutions highly dependent on the data)

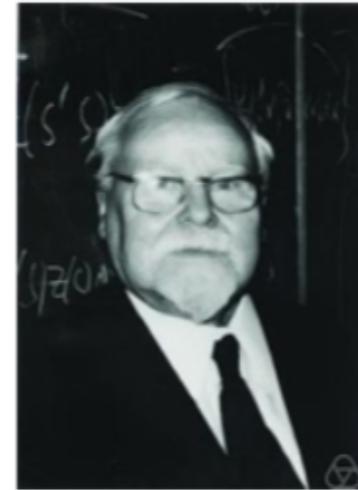
# Regularisation. Tikhonov regularisation



Regularisation techniques restore stability and ensure generalisation

Consider the Tikhonov regularisation scheme,

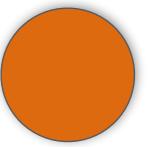
$$\min_{\omega \in \mathbb{R}^d} \hat{\mathcal{E}}(f) + \lambda \|\omega\|^2$$



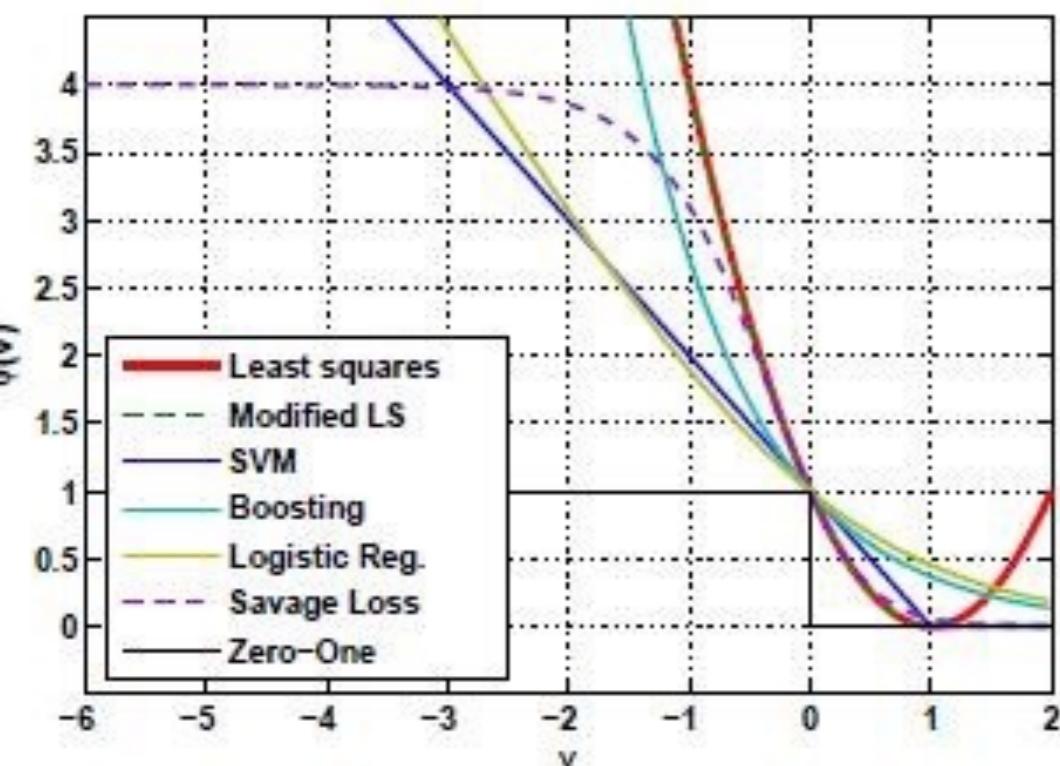
Tikhonov '62   Phillips '62

It describes a large class of methods sometimes called **Regularisation Networks**.

- ▶  $\omega$  is called the regulariser, which controls the stability of the solution and prevents overfitting.
- ▶  $\lambda$  balances the error term and the regulariser.

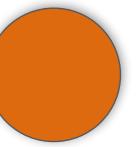


# Regularisation. Loss functions



- ▶ Different loss functions  $V$  induce different classes of methods.
- ▶ We will see common aspects and differences in considering different loss function.
- ▶ There exists no general computational scheme to solve Tikhonov Regularisation.
- ▶ The solution depends on the considered loss function.

# The Regularised Least Squares Algorithm



**Regularised Least Squares (RLS):** Tikhonov regularisation

$$\min_{\omega \in \mathbb{R}^d} \hat{\mathcal{E}}(f) + \lambda \|\omega\|^2, \quad \hat{\mathcal{E}}(f) = \frac{1}{n} \sum_{i=1}^n V(y_i, f(x_i))$$

*Square loss function:*

$$V(y, f(x)) = (y - f(x))^2$$

We then obtain the RLS optimisation problem (linear model):

$$\min_{\omega \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n (y_i - \omega^T x_i)^2 + \lambda \omega^T \omega, \quad \lambda \geq 0$$

**Computations?**

**Statistics?**

# The Regularised Least Squares Algorithm

$$\min_{\omega \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n (y_i - \omega^T x_i)^2 + \lambda \omega^T \omega, \quad \lambda \geq 0$$

## Computations?

### Notation

$$\frac{1}{n} \sum_{i=1}^n (y_i - \omega^T x_i)^2 = \frac{1}{n} \|Y_n - X_n \omega\|^2$$

**Setting gradients....**  $-\frac{2}{n} X_n^T (Y_n - X_n \omega), \quad \text{and,} \quad 2\omega$

**... to zero**  $(X_n^T X_n + \lambda n I) \omega = X_n^T Y_n$

**What does regularisation doing???**

$\lambda$  controls invertibility of  $(X_n^T X_n + \lambda n I)$

# The Regularised Least Squares Algorithm

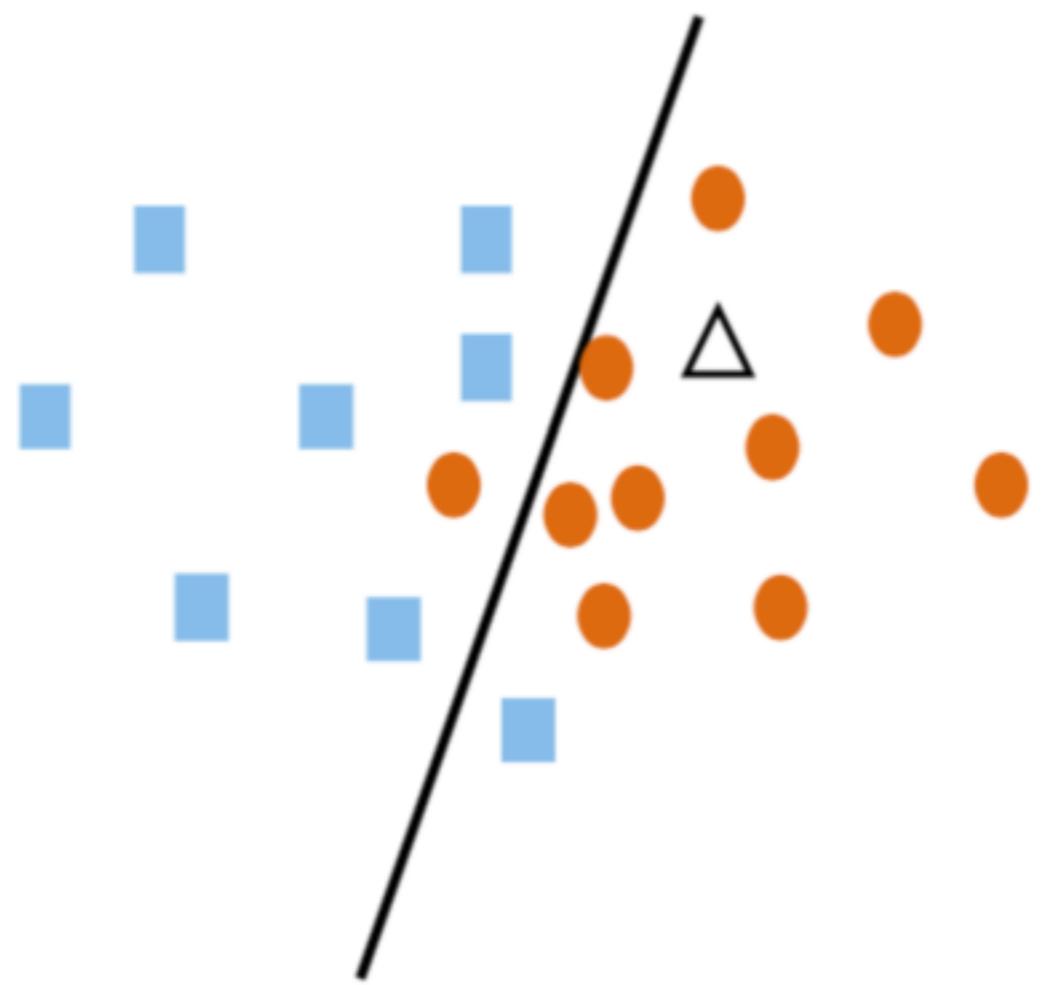
$$\min_{\omega \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n (y_i - \omega^T x_i)^2 + \lambda \omega^T \omega, \quad \lambda \geq 0$$

**Statistics?**

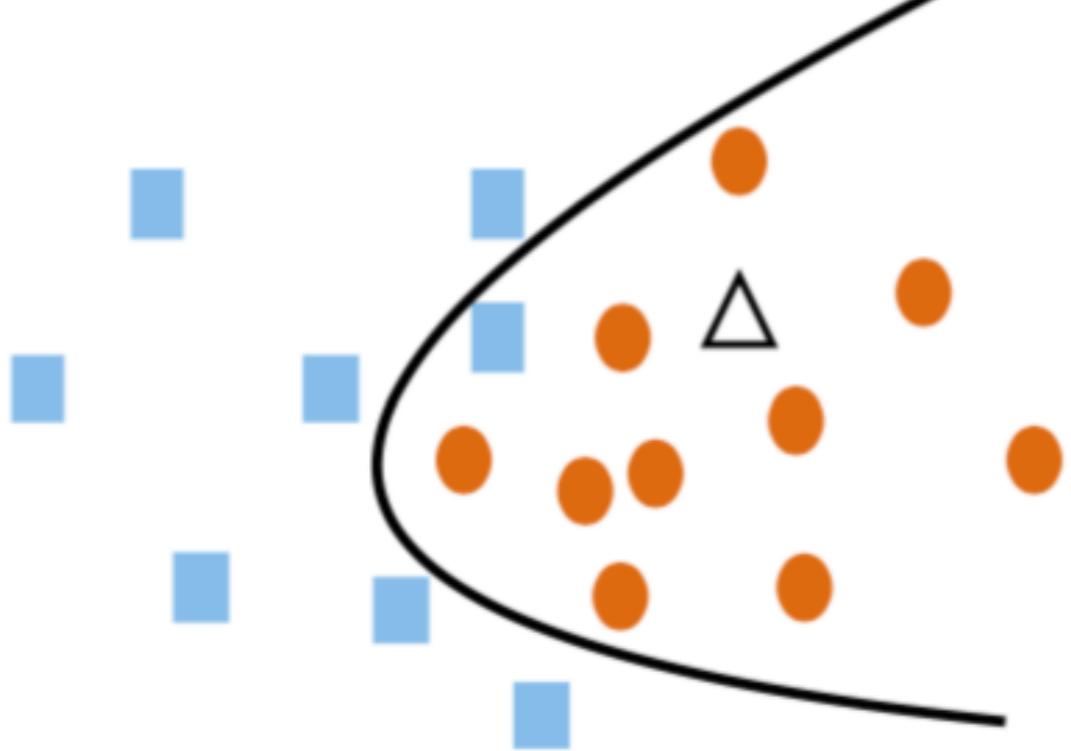
$$(X_n^T X_n + \lambda n I) \omega = X_n^T Y_n$$

***another story that shall be told another time***  
(Stein 1956, James and Stein 1961)

**Shrinkage - Stein Effect- Admissible Estimator**



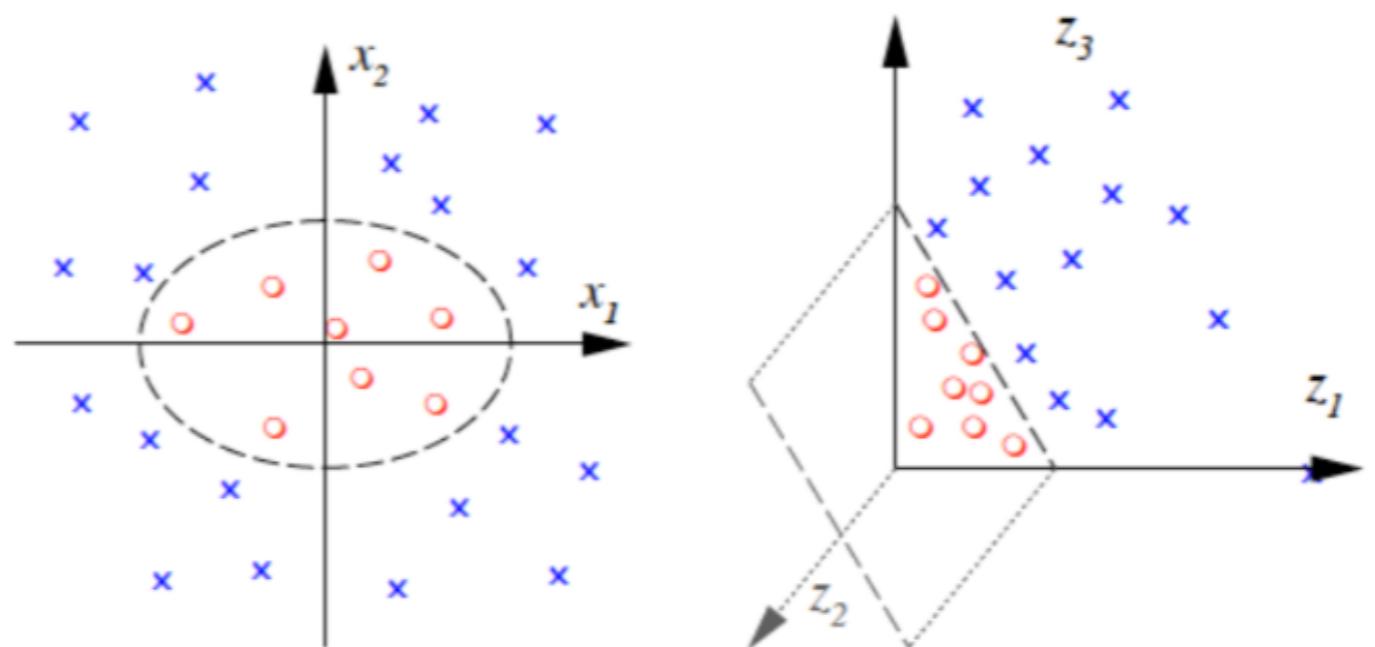
**Why a linear decision rule?**



## Dictionaries

$$x \mapsto \tilde{x} = (\phi_1(x), \dots, \phi_p(x)) \in \mathbb{R}^p$$

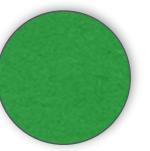
$$f(x) = \omega^T \tilde{x} = \sum_{j=1}^p \phi_j(x) \omega^j$$



$$\Phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$$

$$(x_1, x_2) \mapsto (z_1, z_2, z_3) := (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$

# Computational Complexity



$$(X_n^T X_n + \lambda n I) \omega = X_n^T Y_n \quad \longrightarrow \quad (\tilde{X}_n^T \tilde{X}_n + \lambda n I) \omega = \tilde{X}_n^T Y_n$$

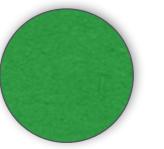
**What about computational complexity?**

## Complexity vademecum

$M \in \mathbb{R}^{n \times p}$  is a matrix and  $v, v' \in \mathbb{R}^p$  are  $p$ -dimensional vectors

- $v^T v' \mapsto O(p)$
- $Mv' \mapsto O(np)$
- $MM^T \mapsto O(n^2p)$
- $(MM^T)^{-1} \mapsto O(n^3)$

# Computational Complexity



$$(X_n^T X_n + \lambda n I) \omega = X_n^T Y_n \quad \longrightarrow \quad (\tilde{X}_n^T \tilde{X}_n + \lambda n I) \omega = \tilde{X}_n^T Y_n$$

**What about computational complexity?**

$$O(p^3) + O(p^2 n)$$

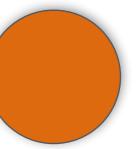
**What if  $p$  is much larger than  $n$ ?**

$$(X_n^T X_n + \lambda n I)^{-1} X_n^T = X_n^T (X_n X_n^T + \lambda n I)^{-1}$$

$$\omega = X_n^T \underbrace{(X_n X_n^T + \lambda n I)^{-1} Y_n}_{c} = \sum_{i=1}^n x_i^T c_i$$

**Computational complexity:**  $O(n^3) + O(pn^2)$

# Kernels



$$(X_n^T X_n + \lambda n I)^{-1} X_n^T = X_n^T (X_n X_n^T + \lambda n I)^{-1}$$

$$\omega = X_n^T \underbrace{(X_n X_n^T + \lambda n I)^{-1} Y_n}_c = \sum_{i=1}^n x_i^T c_i$$



$$f(x) = x^T \omega = \sum_{i=1}^n \underbrace{x^T x_i}_{K(x, x_i)} c_i$$

$$c = (K_n + \lambda n I)^{-1} Y_n, \quad (K_n)_{i,j} = K(x_i, x_j)$$

- linear kernel  $K(x, x') = x^T x'$
- polynomial kernel  $K(x, x') = (x^T x' + 1)^d$
- Gaussian kernel  $K(x, x') = e^{-\|x-x'\|^2/2\sigma^2}$

# **Things I won't tell you about this time....**

$$\hat{f}(x) = \sum_{i=1}^n K(x_i, x)c_i$$

- ▶ Reproducing Kernel Hilbert Spaces
- ▶ Gaussian Processes
- ▶ Integral Equations
- ▶ Sampling Theory / Inverse Problems
  
- ▶ Multi-task, labels, outputs, classes

# **Things I will tell you about tomorrow ....**

- ▶ Loss functions: SVM, Logistic Regression...

## End of sub-Part II

We considered the regularised empirical risk minimisation with squared loss function, leading to regularised least squares algorithm. We discussed how to deal with non linear problems via kernels and dictionaries.

### **Regularised Least Squares - Dictionaries - Kernels**

Today's Afternoon Lecture:  
**Machine Learning in Practice**