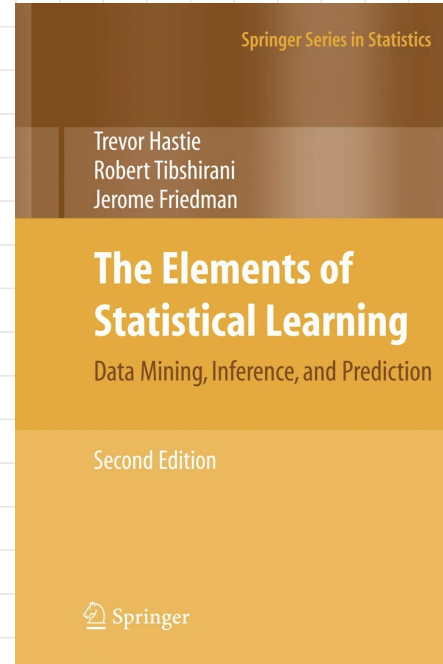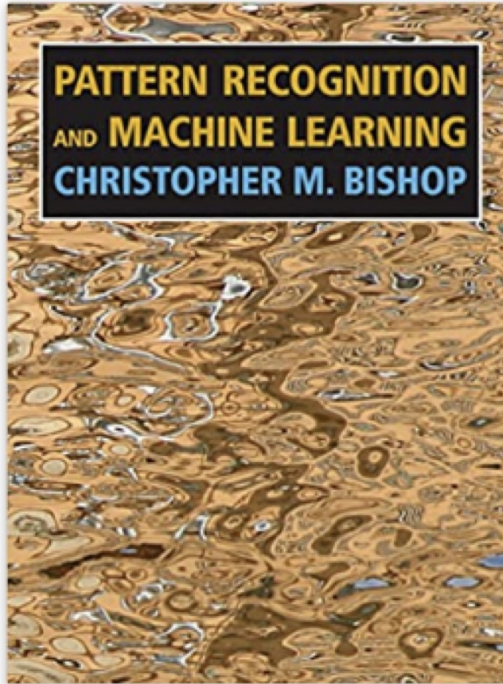ULB, Oct/2023

# Introduction to Machine Learning Methods

Bryan Zaldívar (IFIC, Valencia)

# OUTLINE OF THE LECTURES

1. Overview of Machine Learning ⎫ today
2. Summary of statistics
3. Regression & overfitting control
4. Bayesian learning
5. Classification methods
6. Neural networks
7. ... you decide

# MAIN BIBLIOGRAPHY

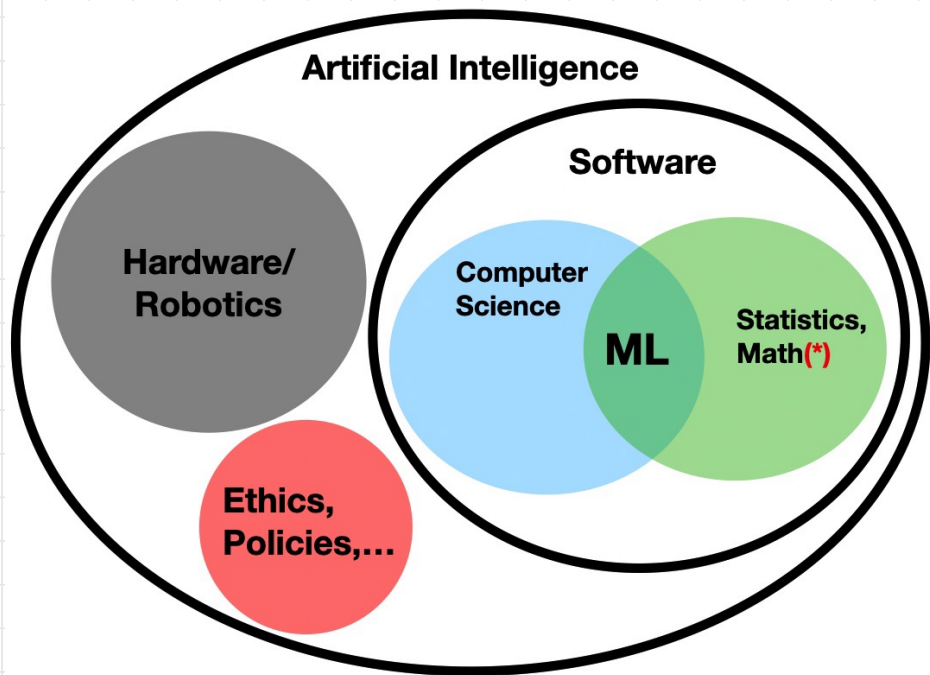PATTERN RECOGNITION
AND MACHINE LEARNING
CHRISTOPHER M. BISHOP

**Springer Series in Statistics**

Trevor Hastie
Robert Tibshirani
Jerome Friedman

**The Elements of
Statistical Learning**

Data Mining, Inference, and Prediction

**Second Edition**

Springer

# Lecture #1: Overview of ML

you are not supposed to understand these slides in detail...

just watch them in a relaxed mood...

# DEFINITION



**Artificial Intelligence**

**Hardware/ Robotics**

**Software**

**Computer Science**

**ML**

**Statistics, Math(*)**

**Ethics, Policies,...**

(*) should be tailored to the domain expertise (e.g. physics)!

Very broadly... ML is about implementing stats on a computer

**ML**

Stats methods in "Numerical Recipes"

Multivariate analyses (e.g. LEP in the 90's)

Looking for best polynomial to fit your data

# DEFINITION

"Machine Learning" coined by Arthur Samuel
(IBM, 1959)

Quote by Tom M. Mitchell (prof. US, 1997):

"A computer program is said to learn from experience 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"

# DEFINITION

**Historically**, ML associated to:

- Games (chess, checkers, Go, ...)
- Image/sound/video recognition
- Natural language processing

**Nowadays**, ML synonym of:

- Data mining
- Big Data
- Deep Learning

# SOME HISTORICAL MILESTONES

1950: Alan Turing's learning machine (based on primitive form of genetic algorithm)

1951: The 1st neural network is created (founded by Air Force Office)

1952: Arthur L. Samuel designed a computer program able to play checkers (IBM)

1967: Nearest Neighbor algorithm is created. The algorithm was used to map routes

1970: Back-propagation was invented

1985: NetTalk: a program that learns to pronounce written words in English

1997: IBM's Deep Blue beats Kasparov

2016: Google's "Alpha Go" beats an unhandicapped Go player

2020: Google's "Alpha Fold 2" is able to predict how proteins fold from amino acid sequence

2022: Open AI launches "ChatGPT 3", marking a breakthrough in Language Processing

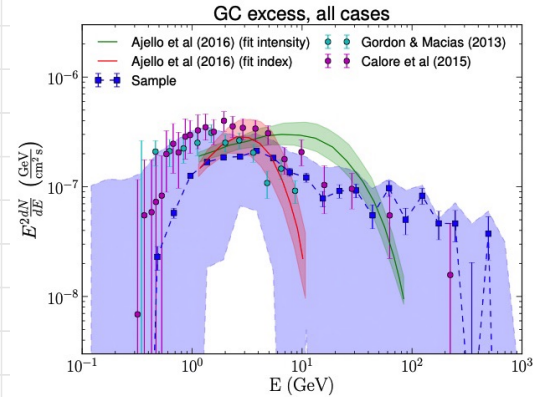WHAT DOES ML BRINGS TO PHYSICS ?

# MOTIVATION

* Physical knowledge of
  background is limited
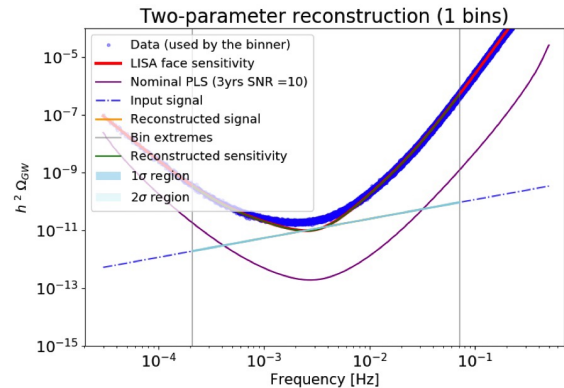  ( common problem in astrophysics)

# MOTIVATION

\* Physical knowledge of background is limited
(common problem in astrophysics)

e.g. • $\gamma$-rays
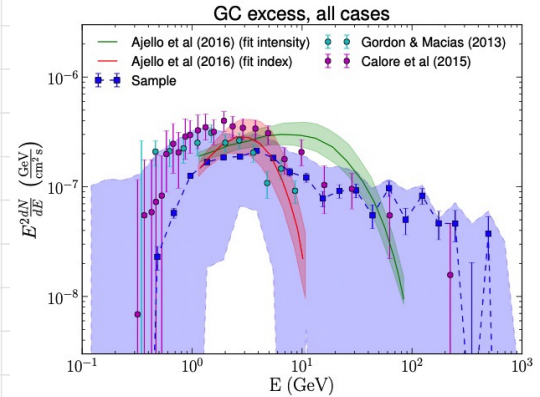• GW's

Fermi-LAT, 1704.03910



Caprini et al, 1906.09244

# MOTIVATION

**\* Physical knowledge of background is limited**
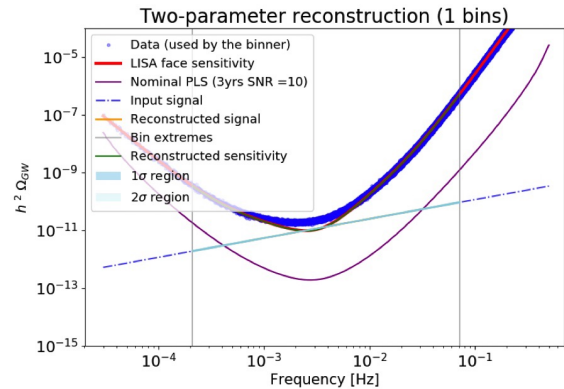
(common problem in astrophysics)

e.g. • $\gamma$-rays
• GW's

<span style="color:blue">Signal</span> + <span style="color:orange">Background</span> ⟺ Data

<span style="color:blue">( ↳ interest
(Physical model)</span>

<span style="color:orange">↳ nuisance
(data-driven model)</span>

Fermi-LAT, 1704.03910



GC excess, all cases

Caprini et al, 1906.09244



Two-parameter reconstruction (1 bins)

# MOTIVATION

* Physics well known, but observables very complicated to compute

(complex topology, particle misidentification, etc)

# MOTIVATION

* Physics well known, but observables very complicated to compute
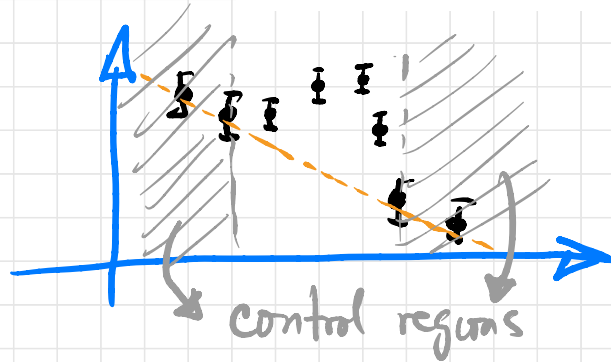
(complex topology, particle misidentification, etc)
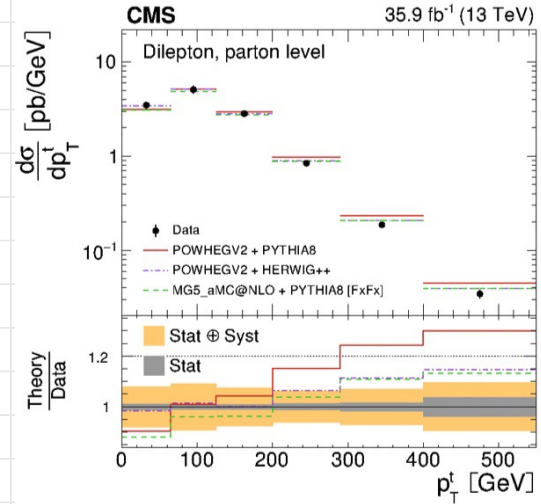
# MOTIVATION

* Physics well known, but observables very complicated to compute

(complex topology, particle misidentification, etc)

Data-driven methods to infer the background from inter/extrapolations

control regions

# MOTIVATION

* Statistical bottleneck

# MOTIVATION

* Statistical bottleneck

- More complex datasets
  ↓
More complex physical modeling
  ↓
More complex simulators

# MOTIVATION

<span style="color:red">\* Statistical bottleneck</span>

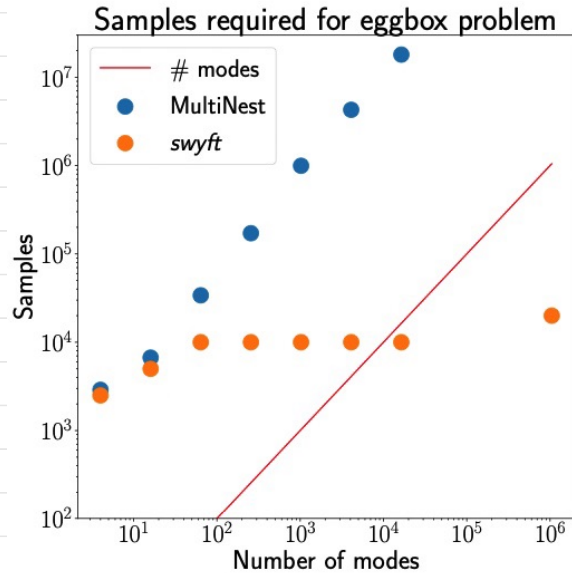- More complex datasets
  ↓
  More complex physical modeling
  ↓
  More complex simulators
  ↓

<span style="color:blue">Better statistical treatment !</span>

<span style="color:gray">Miller et al, 2011.13951</span>



Samples required for eggbox problem

<span style="color:orange">
- better scaling
- more descriptive
- Higher statistical power
</span>

# MOTIVATION

* Hints about the underlying physics

# MOTIVATION

* Hints about the underlying physics

  → # physical variables

  • The intrinsic dimension of
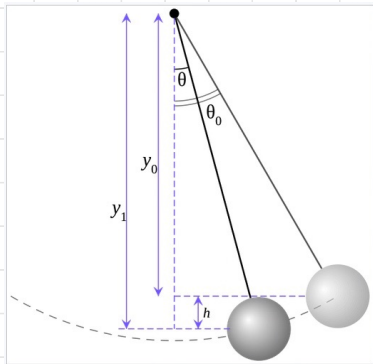    - Single pendulum : 2
    - Lava lamp : ?

# MOTIVATION

* **Hints about the underlying physics**

→ # physical variables

- The intrinsic dimension of
  - Single pendulum : 2
  - Lava lamp : **?**



ID = 7 − 8 ⟸ | Statistical model (ML) |



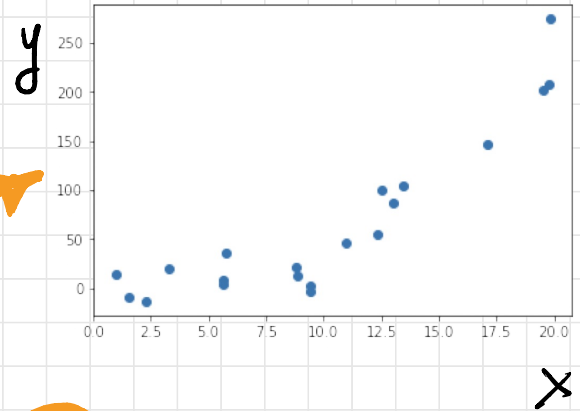(see 2112.10755)

# MAIN ML PARADIGMS

# SUPERVISED LEARNING

Data $\{ \mathbf{X}, \vec{y} \}$

| $x_1$ | $x_2$ | $x_3$ | $y$ |
|---|---|---|---|
| | | | |
| | | | |

$\vdots$

Regression
$y \in \mathbb{R}$



Classification
$y \in \{ l_1, l_2, l_3, \ldots \}$



Tasks:

1 - Statistical inference of the data's probability distribution

2 - Predict the output for a new point

# SUPERVISED LEARNING

Data $\{ \mathbf{X}, \vec{y} \}$

| $x_1$ | $x_2$ | $x_3$ | $y$ |
|---|---|---|---|
| | | | |
| | | | |

$\vdots$

Regression
$y \in \mathbb{R}$



Classification
$y \in \{ l_1, l_2, l_3, \ldots \}$



Tasks:

1 - Statistical inference of
the data's probability distribution

2 - Predict the output for a new point
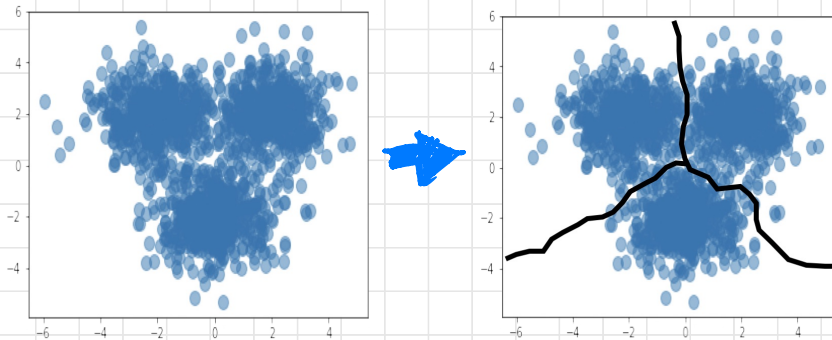
# SUPERVISED LEARNING

How is this different from fitting a function to some data?

- Optimization algorithms, adapted to arbitrarily complex fitting functions

- Procedures to avoid overfitting (beyond e.g. $x^2/_{D.O.F}$)

- Choose the best function from a catalog

- Types of functions typically used (Neural networks, Decision Trees, Kernels,...)
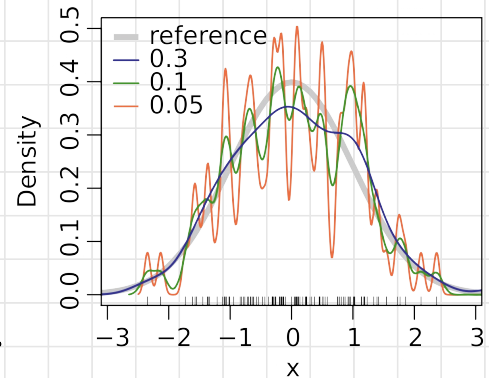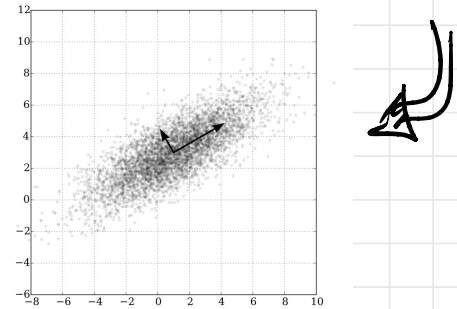
# UNSUPERVISED LEARNING

Data $\{X\}$ [no output/labels]

1) Clusterize the data



2) Dimensionality reduction



3) Probability density estimation

# REINFORCEMENT LEARNING

(similar in spirit to the way humans learn from their environment)

computer program

"Environment"

System in which the program operates



state $S_t$

reward $R_t$

Agent

$R_{t+1}$

$S_{t+1}$

Environment

action $A_t$

"Policy"

Method to map the program's state to an action

"Action"

Done by the program to move to a new state

"Reward"

Feedback from the environment

RL commonly used for learning to play games (chess, Go, ...)

Nowadays also used in science; e.g. physics. (quantum systems & computing)

MCMC improvements, ...

BTW ChatGPT uses a form of RL!