

Lecture I: Machine Learning

Outline:

Lecture I: → Motivation

→ Machine Learning basics

→ Example of polynomial regression

Lecture II: → Training procedure

→ Neural networks

Lecture III: → ML in high energy physics

Motivation: Why using ML in physics?

→ Deep Learning revolution in many fields outside of physics, such as:

↳ Computer vision:

• image ~~recognition~~ classification

• image generation

• image segmentation: locate objects and boundaries

↳ Speech recognition

⋮

→ We should take advantage of this and apply it to physics!

What do we need?

→ large data sets

(→ truth labels)

example:
LHC physics

✓

(via simulations) ✓

ML basics

What is ML?

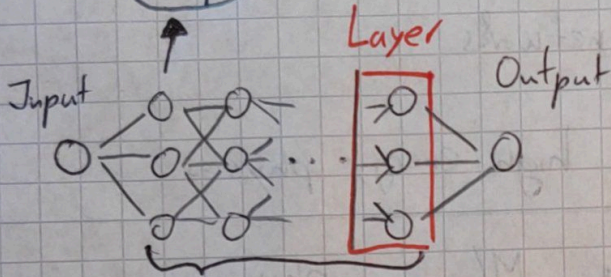
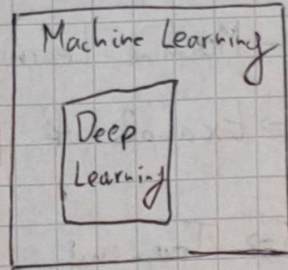
Study of algorithms which are learnt through the use of data

vs. handcrafted algorithms

What is Deep Learning?

↳ subfield of ML

↳ use of deep neural networks



→ Lecture II

deep $\hat{=}$ many layers \sim many parameters / trainable weights

What is the advantage of deep networks?

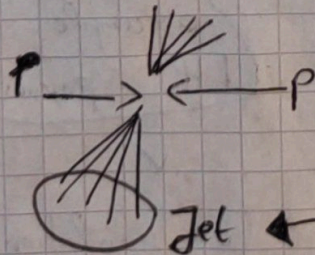
→ allows to learn features of high dimensional data

⇔

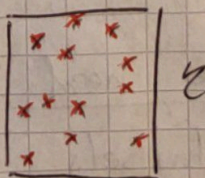
→ "non deep" learning:

- require more human work
- ~~can~~ reduce dimensionality by extracting features by hand!

Example: LHC physics



Jet image



high dimensional data

Deep network

features

vs

$(p_T^J, \psi_{01}^J, \dots, N\text{-subjettiness})$

BDT

Categories of Machine Learning

Supervised: labeled data!

Example:

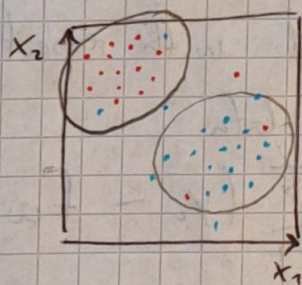
(i) Classify images of cats and dogs. Each image is either labeled "cat" or "dog".

(ii) top tagging: hadronically decaying top vs. QCD jet.
For each jet of the training data set we know that it originated from a top or light quark / gluon!
→ label: "top", "QCD"

Unsupervised: NO labels! Learn patterns in data, detect anomalies

Example:

(i) clustering



(ii) data set with mostly dogs but a few cats.

However, no labels! We don't know which image is a cat or a dog!
↑ anomaly

reinforcement learning: Similar to supervised but no simple data set. Instead an "agent" interacts with an environment and has to adopt behaviour accordingly.

↳ ~~Sequential~~ Sequential decisions: trial and error

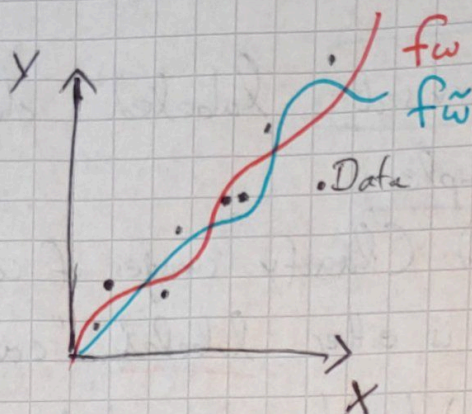
Example:

(i) games where agent has to adapt constantly to new environment

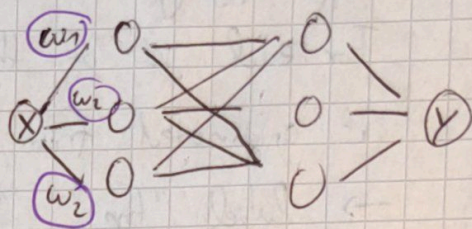
(ii) robot which has to adapt to environment

Basic ingredients of ML (supervised)

→ Dataset $D = (X, Y)$
Input ↑ Label



→ Model: $f_w(x)$
↑ w : trainable parameters



→ Loss function: Defines the task!

- Object which is minimized with respect to w to best describe the data

Example: $MSE = \sum_{i=1}^N (y_i - f_w(x_i))^2$, N : Number of data points
~ Error

(neg log of Gaussian)

→ Training: Minimization of loss function → Lecture 2

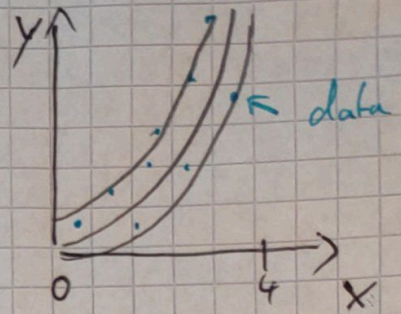
Example of polynomial regression

Motivation: • get some intuition about ML

- key points extrapolate to Deep Learning

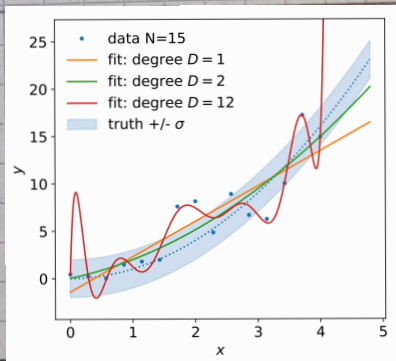
(This is stolen from the great review: arXiv: 1803.08823)

Data set: $X \in [0, 4]$, $y = X^2 + \text{Noise}$
 \rightarrow true dependence!
 \rightarrow we don't know!

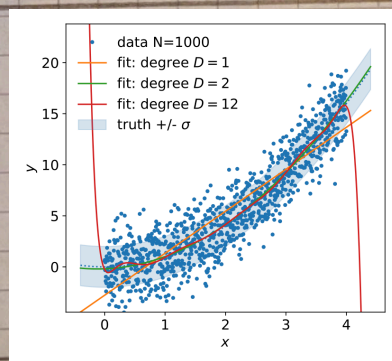
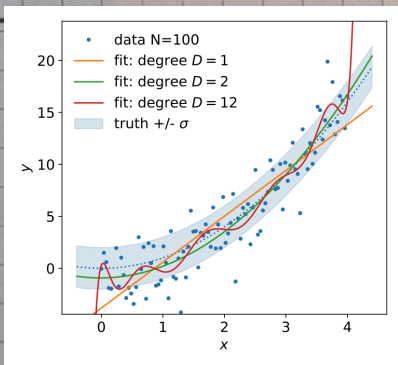


Model: $f_w(x) = w_0 + w_1x + w_2x^2 + \dots + w_Dx^D$
 Polynomial of degree D
 $\equiv f_D(x)$

Loss function: MSE



$\rightarrow f_1$: not "complex" enough
 $\rightarrow f_{12}$: fits all statistical fluctuations! \equiv Overfitting
 How to fix this?



\rightarrow Complex models need large data sets!
 Deep Learning requires large data sets!

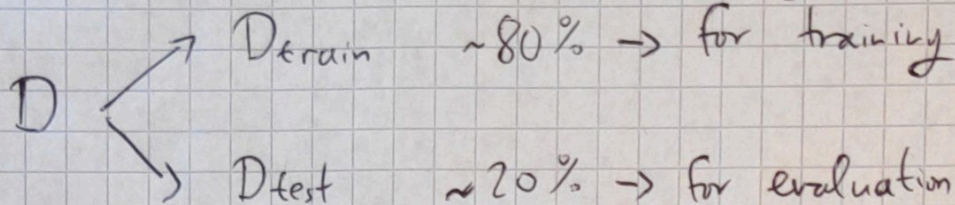
(Example for state of the art deep networks: GPT-3 \sim 172 B trainable weights)

How to quantify performance?

MSE \sim Error

But MSE is smallest for $f_{12}(x)$ ~~*~~. Even for $N=15$!

\Rightarrow Have to split data set before training!

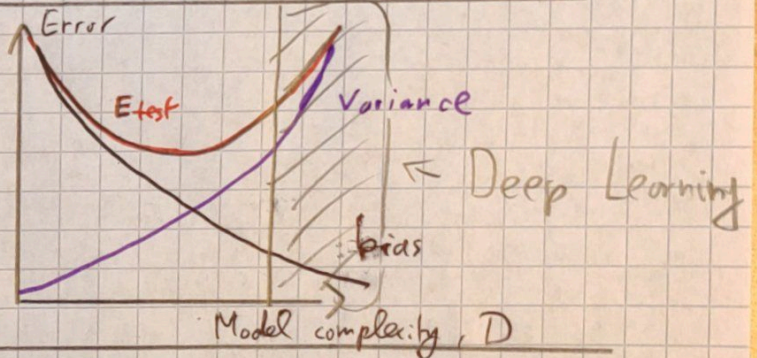
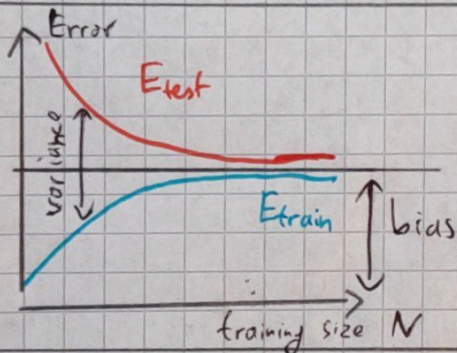
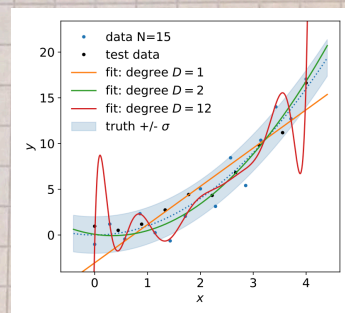


Compare MSE =: E

$$E_{\text{train}} = \text{MSE}(D^{\text{train}})$$

$$E_{\text{test}} = \text{MSE}(D^{\text{test}})$$

vs.



- \rightarrow Need of separate data set for evaluation!
- \rightarrow For Deep Learning: need of regularisation methods ~~to~~ large data sets and potentially regularisation methods

\rightarrow lecture 2 or 3?