

Stochastic Machine Learning

Chapter 01 - Introduction

Thorsten Schmidt

Abteilung für Mathematische Stochastik

www.stochastik.uni-freiburg.de
thorsten.schmidt@stochastik.uni-freiburg.de

SS 2024

Motivation

- ▶ **Machine Learning** is nowadays used at many places (Google, Amazon, etc.) with a great variety of applications.
- ▶ It is a great job opportunity ! It needs maths and probability !
- ▶ Many applications are surprisingly successful (speech / face recognition / robotic / autonomous cars / medicine / chemistry / chat GPT) and currently people are seeking further applications
- ▶ Here we want to learn about the foundations, discuss implications and what can be done by ML and what not.

Topics include:

- ▶ Foundations and deeper understanding
- ▶ Uncertainty Quantification
- ▶ Optimization under Uncertainty
- ▶ Regulation and Fairness
- ▶ Risk quantification (for example in transfer learning)
- ▶ Advancement of fundamental theories (LEAN)

Organization

Slides and code will be available on github. But - proofs and theorems will mostly remain on blackboard only.

- ▶ The lectures will mix python implementation with theory - so it is now a good time for you to start learning python. Every 4th lecture will focus on implementation and projects.
- ▶ Homework in the 2nd part will be done by projects. You can choose a topic which interests you, and we will provide topics. Groups up to 6 people work on a project, more than one group can also work on the same project.
- ▶ We provide a shared repository for all projects, such that you share your current work and can profit from the work from others.
- ▶ **Moritz Ritter** is organizing the projects - please contact him for questions.
- ▶ You need 50% of the points for the oral exam at the end. The oral exam will cover all **theoretical** parts of the lecture.
- ▶ Code and slides are on https://github.com/tschmidtfreiburg/ML_lecture_2024 (see homepage).

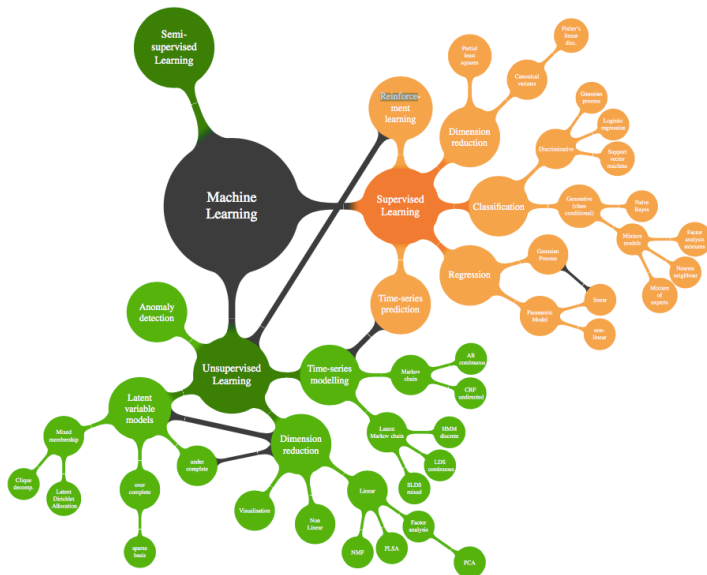
Overview¹

- ▶ Artificial intelligence is the field where computers solve problems.
- ▶ It is easy for a computer to solve tasks which can be described formally (Chess, Tic-Tac-Toe). The challenge is to solve a tasks which are hard to describe formally (but are easy for humans: walk, drive a car, speak, recognize people ...)
- ▶ The solution is to allow computers to learn from experience and to understand the world by a hierarchy of concepts, each concept defined in terms of its relation to simpler concepts.
- ▶ A fixed knowledge-base would be somehow limiting such that we are interested in such attempts where the systems acquire their own knowledge, which we call **Machine Learning**.

¹This introduction follows somehow Goodfellow et.al. (2016) and my previous lectures.

Artificial Intelligence

- ▶ First examples of machine learning are **logistic regression** or **naive Bayes** → standard statistical procedures (E.g. the recognition of spam, more examples to follow)
- ▶ Problems become simpler with a nice representation. Of course it would be nice if the system itself could find such a representation, which we call **representation learning**.
- ▶ An example is the so-called **auto-encoder**. This is a combination of an encoder and a decoder. The encoder converts the input to a certain representation and the decoder converts it back again, such that the result has nice properties.
- ▶ Speech for example might be influenced by many factors of variation (age, sex, origin, ...) and it needs nearly human understanding to disentangle the variation from the content we are interested in.
- ▶ **Deep Learning** solves this problem by introducing hierarchical representations.
- ▶ This leads to the following hierarchy:
- ▶ AI → machine learning → representation learning → deep learning.



Source: Barber (2012).

Examples of Machine Learning

Some of the most prominent examples:

- ▶ LeCun et.al.² recognition of handwritten digits. The MNIST Database³ provides 60.000 samples for testing algorithms. The NIST database is of increased size⁴
- ▶ The Viola & Jones face recognition,⁵. This path-breaking work proposed a procedure to combine existing tools with machine-learning algorithms. One key is the use of approx. 5000 learning pictures to train the routine. We will revisit this procedure shortly.
- ▶ Imagenet is an image database containing many images classified (cats, cars, etc.)⁶
- ▶ Various twitter datasets are available, for example for learning to detect hate speech.
- ▶ Kaggle⁷ is a platform where computational competitions are hosted. It also provides many many data examples with it.
- ▶ Datasets for machine-learning research on Wikipedia⁸.

²Y. LeCun et al. (1998). „Gradient-based learning applied to document recognition“. In: *Proceedings of the IEEE* 86.11, pp. 2278–2324.

³<http://yann.lecun.com/exdb/mnist/>

⁴<https://www.nist.gov/srd/nist-special-database-19>

⁵P. Viola and M. Jones (2001). „Robust Real-time Object Detection“. In: *International Journal of Computer Vision*. Vol. 4. 34–47.

⁶<http://image-net.org>

⁷www.kaggle.com

⁸https://en.wikipedia.org/wiki/List_of_datasets_for_machine-learning_research

MNIST

We will see this data set in more detail.



As an extension we also have the EMNIST data set - with letters and more challenging.



- ▶ Speech recognition has long been a difficult problem for computers (first works date to the 50's) and only recently been solved with high computer power. It may seem surprising, that mathematical tools are at the core of these solutions. Let us quote Hinton et.al.⁹

Most current speech recognition systems use hidden Markov models (HMMs) to deal with the temporal variability of speech and Gaussian mixture models (GMMs) to determine how well each state of each HMM fits a frame or a short window of frames of coefficients that represents the acoustic input. (...)

Deep neural networks (DNNs) that have many hidden layers and are trained using new methods have been shown to outperform GMMs on a variety of speech recognition benchmarks, sometimes by a large margin

So, one of our tasks will be to develop a little bit of mathematical tools which we will need later. Most notably, some of the mathematical parts can be replaced by deep learning, which will be of high interest to us.

⁹Geoffrey Hinton et al. (2012). „Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups“. In: **IEEE Signal Processing Magazine** 29.6, pp. 82–97.

NEWS | 15 April 2024

AI now beats humans at basic tasks — new benchmarks are needed, says major report

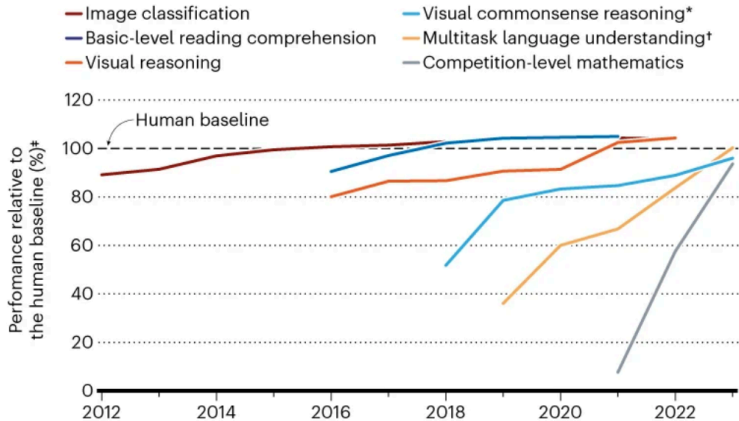
Stanford University's 2024 AI Index charts the meteoric rise of artificial-intelligence tools.

By [Nicola Jones](#)



SPEEDY ADVANCES

In the past several years, some AI systems have surpassed human performance on certain benchmark tests, and others have made rapid progress.



*Requires an AI system to answer questions about an image and provide a rationale for why its answers are true.

†Tests an AI model's knowledge and problem-solving ability with regard to 57 subjects, including broader topics such as mathematics and history, and narrower areas such as law and ethics.

*Data indicate the best performance of an AI model that year.

Cost of business

As performance is skyrocketing, so are costs. GPT-4 – the LLM that powers ChatGPT and that was released in March 2023 by San Francisco-based firm OpenAI – reportedly cost US\$78 million to train. Google's chatbot Gemini Ultra, launched in December, cost \$191 million. Many people are concerned about the energy use of these systems, as well as the amount of water needed to cool the data centres that help to run them². "These systems are impressive, but they're also very inefficient," Maslej says.

- What kind of consequences do we draw from this ?

If we want to give this a little bit more context we could compare these results.

- ▶ The human brain has approximately 86 billion neurons
- ▶ GPT4 is announced to have approx. 175 trillion parameters
- ▶ The cost of training the human brain is (compare this to 80 Million dollar to train GPT4). But ...
- ▶ How many parameters does a neural network with n neurons have? Think of a simple fully connected NN with L layers and n_i neurons per layer.
- ▶ For each neuron in layer i we connect to n_{i+1} neurons, which makes $n_i \cdot n_{i+1}$ parameters. Additionally we include one parameter for a shift in mean which makes

$$\sum_{i=1}^{L-1} (n_i \cdot n_{i+1} + n_{i+1})$$

- ▶ So for a single-layer NN with 100 input nodes and 100 output nodes we have $100 \cdot 100 + 100 \cdot 100 + 100 = 20.100$ parameters.

The AI Index Report 2024

1. AI beats humans on some tasks, but not on all.

AI has surpassed human performance on several benchmarks, including some in image classification, visual reasoning, and English understanding. Yet it trails behind on more complex tasks like competition-level mathematics, visual commonsense reasoning and planning.

2. Industry continues to dominate frontier AI research.

In 2023, industry produced 51 notable machine learning models, while academia contributed only 15. There were also 21 notable models resulting from industry-academia collaborations in 2023, a new high.

3. Frontier models get way more expensive.

According to AI Index estimates, the training costs of state-of-the-art AI models have reached unprecedented levels. For example, OpenAI's GPT-4 used an estimated \$78 million worth of compute to train, while Google's Gemini Ultra cost \$191 million for compute.

4. The United States leads China, the EU, and the U.K. as the leading source of top AI models.

In 2023, 61 notable AI models originated from U.S.-based institutions, far outpacing the European Union's 21 and China's 15.

5. Robust and standardized evaluations for LLM responsibility are seriously lacking.

New research from the AI Index reveals a significant lack of standardization in responsible AI reporting. Leading developers, including OpenAI, Google, and Anthropic, primarily test their models against different responsible AI benchmarks. This practice complicates efforts to systematically compare the risks and limitations of top AI models.

6. Generative AI investment skyrockets.

Despite a decline in overall AI private investment last year, funding for generative AI surged, nearly octupling from 2022 to reach \$25.2 billion. Major players in the generative AI space, including OpenAI, Anthropic, Hugging Face, and Inflection, reported substantial fundraising rounds.

7. The data is in: AI makes workers more productive and leads to higher quality work.

In 2023, several studies assessed AI's impact on labor, suggesting that AI enables workers to complete tasks more quickly and to improve the quality of their output. These studies also demonstrated AI's potential to bridge the skill gap between low- and high-skilled

8. Scientific progress accelerates even further, thanks to AI.

In 2022, AI began to advance scientific discovery. 2023, however, saw the launch of even more significant science-related AI applications—from AlphaDev, which makes algorithmic sorting more efficient, to GNoME, which facilitates the process of materials discovery.

9. The number of AI regulations in the United States sharply increases.

The number of AI-related regulations in the U.S. has risen significantly in the past year and over the last five years. In 2023, there were 25 AI-related regulations, up from just one in 2016. Last year alone, the total number of AI-related regulations grew by 56.3%.

Questions

We repeatedly state questions after some slides. These allow you to reflect on the content and also invite you to research / experiment with some topics yourself.

- ▶ Was is artificial intelligence ?
- ▶ Was is machine learning ?
- ▶ Do you know what a neural network is (look for the history in the internet)?
- ▶ What are shallow / deep networks ?
- ▶ What are the applications which you find most exciting ?
- ▶ What are the applications that you think will have the largest impact on our future?
- ▶ Research a bit yourself: look for datasets, look for latest applications etc.