

Theory of Statistical Learning

Part I

Damien Garreau

Université Côte d'Azur

2022

1. General presentation

Who am I?

- ▶ maître de conférence (= assistant professor) in LJAD (Laboratoire Jean Dieudonné)
- ▶ before that: postdoctoral researcher (Max Planck Institute, Tübingen, Germany)
- ▶ even before: PhD in Inria Paris
- ▶ teaching (≈ 200 hours per year)
- ▶ research (rest of the time):
 - ▶ interpretability of machine learning algorithms
 - ▶ using statistical tools to understand methods in computer vision

Who are you?

- ▶ online teaching is suboptimal :(
- ▶ speaking to a black screen is weird
- ▶ **Please introduce yourselves!** (with camera on if possible)
- ▶ I will call your name then you can briefly introduce yourself

Goal of the course

- ▶ **Goal i):** understand the *maths* behind the algorithms that you learn
- ▶ this can save a lot of time!
- ▶ one often works with limited resources
- ▶ **Goal ii):** learn about theoretical guarantees on existing algorithms
- ▶ a way to be reassured: under some assumptions, my method works
- ▶ see more clearly the *limitations* of the methods: if some assumption is not satisfied, we can prove that it will fail

Organization of the course

- ▶ all the information, documents → Slack
- ▶ (provisional) calendar:
 1. January 12 (today), 9am-12am
 2. January 21, 9am-12am
 3. January 28, 9am-12am
 4. February 4, 9am-12am
 5. February 11, 9am-12am (midterm)
 6. February 25, 9am-12am
 7. March 11, 9am-12am
 8. March 18, 9am-12am
 9. March 22, 9am-12am
 10. March 25, 9am-12am (exam)
- ▶ final grade = (midterm + final)/2

Outline

1. General presentation
2. Introduction to Statistical Learning
 - General introduction
 - First concepts
 - Empirical risk minimization
 - Overfitting
 - PAC learning
 - Uniform convergence
3. A closer look at test error
4. VC dimension
 - Infinite classes can be PAC learnable
 - Examples
 - The fundamental theorem of learning

Requirements

- ▶ **Elementary real analysis:** functions of a real variable, usual functions, continuity, Lipschitz continuity
- ▶ **Calculus:** derivative, partial derivatives, gradient, Taylor series
- ▶ **Basic probability theory:** measurable space, probability measure, random variable, expectation, conditional expectation, probability density function, cumulative density functions
- ▶ **Limit theorems:** law of large numbers, central limit theorem
- ▶ **Linear algebra:** vector space, matrix, norms, diagonalization of a matrix, singular value decomposition

If you feel like you are not up to date on one of these points, write me and I will point you towards some resources.

Useful resources

- ▶ **Main reference:** Shalev-Schwartz, Ben-David, *Understanding Machine Learning: from Theory to Algorithms*, Cambridge University Press, 2014
- ▶ **Also a good read:** Hastie, Tibshirani, Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, Springer Series in Statistics, 2001 (second edition: 2009)
- ▶ **Wikipedia:** as good as ever.
- ▶ **Wolfram alpha:** if you have computations to make and you do not know want to use a proper language:
<https://www.wolframalpha.com/>
- ▶ **Google scholar:** use it!

2. Introduction to Statistical Learning

2.1. General introduction

The goal of statistical learning

- ▶ **Fundamental example:** image classification
- ▶ **Goal:** given any image x , we want to predict which object / animal y is in the image



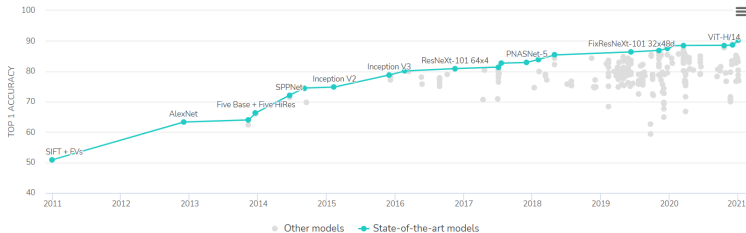
\mapsto “lion”

- ▶ **Main idea:** instead of defining the function f ourselves, we are going to *learn it* from data
- ▶ **Why?** no clear definition of a “lion”
- ▶ **Motivation:** industry (advertisement), healthcare (automated patient triage), military (automated defense systems)

How is this even possible?

- ▶ four ingredients made statistical learning a viable paradigm:
- ▶ **Ingredient (i):** data to feed to the models
- ▶ previous example from ImageNet¹: roughly 1 million images for training (150GB of data)

Image Classification on ImageNet



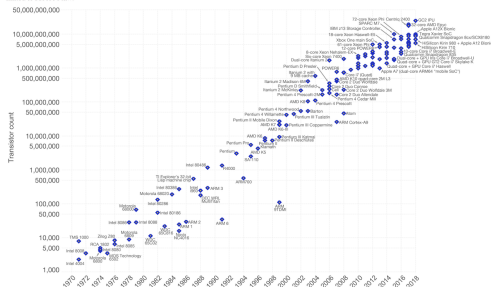
¹Deng et al., *ImageNet: A large hierarchical image database*, CVPR, 2009

How is this even possible?

- **Ingredient (ii):** computing power
- we have the processing power to deal with these data

Moore's Law – The number of transistors on integrated circuit chips (1971-2018)

Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years. This advancement is important as other aspects of technological progress – such as processing speed or the price of electronic products – are linked to Moore's law.

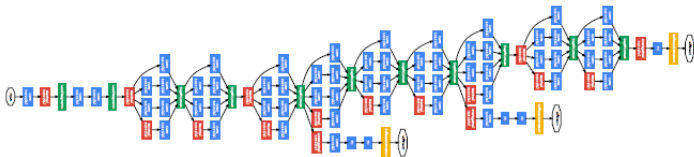


Data source: Wikipedia (https://en.wikipedia.org/wiki/Transistor_count)
The data visualization is available at OurWorldInData.org. There you find more visualizations and research on this topic.

Licensed under CC-BY-SA by the author Max Roser.

How is this even possible?

- ▶ **Ingredient (iii):** models that are complex enough
- ▶ state-of-the-art today: (deep) neural networks (originating from much earlier research²)
- ▶ Inception:³ 24M parameters, GTP-3:⁴ 175B



²Rosenblatt, *The perceptron, a perceiving and recognizing automaton*, tech report, 1957

³Szegedy et al., *Going deeper with convolutions*, CVPR, 2015

⁴Brown et al., *Language Models are Few-Shot Learners*, tech report, 2020

How is this even possible?

SWITCH TRANSFORMERS: SCALING TO TRILLION PARAMETER MODELS WITH SIMPLE AND EFFICIENT SPARSITY

William Fedus*
Google Brain

liamfedus@google.com

Barret Zoph*
Google Brain

barretzoph@google.com

Noam Shazeer
Google Brain

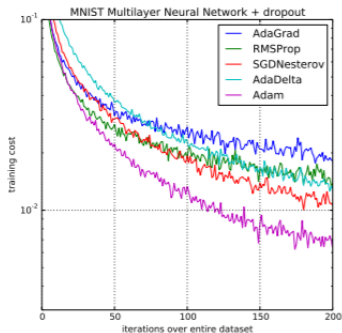
noam@google.com

ABSTRACT

In deep learning, models typically reuse the same parameters for all inputs. Mixture of Experts (MoE) models defy this and instead select *different* parameters for each incoming example. The result is a sparsely-activated model – with an outrageous number of parameters – but a constant computational cost. However, despite several notable successes of MoE, widespread adoption has been hindered by complexity, communication costs, and training instability. We address these with the Switch Transformer. We simplify the MoE routing algorithm and design intuitive improved models with reduced communication and computational costs. Our proposed training techniques mitigate the instabilities, and we show large sparse models may be trained, for the first time, with lower precision (bfloat16) formats. We design models based off T5-Base and T5-Large (Raffel et al., 2019)

How is this even possible?

- ▶ **Ingredient (iv):** efficient algorithms to train the models
- ▶ gradient descent on steroids⁵
- ▶ efficient gradient computations⁶



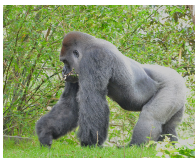
⁵Kingma, Ba, *ADAM: A method for stochastic optimization*, ICLR, 2015

⁶Rumelhart et al., *Learning representations by back-propagating errors*, Nature, 1986

2.2. First concepts

Input space

- ▶ **Input space:** measurable space \mathcal{X} containing all the objects that we want to label
- ▶ also called *domain*, or *domain set*
- ▶ elements $x \in \mathcal{X}$ are usually described as vectors ($\mathcal{X} = \mathbb{R}^d$)
- ▶ coordinates of the vector = *features*
- ▶ **Example:** ImageNet images: RGB images \rightarrow 3 8-bits channels



$$\in \llbracket 0, 255 \rrbracket^{299 \times 299 \times 3}$$

- ▶ \mathcal{X} can be very high-dimensional in modern applications (here $299 \times 299 \times 3 = 268,203$)

Labels and training data

- ▶ **Output:** labels / responses belong to a set \mathcal{Y}
- ▶ **Example:** \mathcal{Y} is the set of names of object and animals of the dataset

```
1  {0: 'tench, Tinca tinca',  
2    1: 'goldfish, Carassius auratus',  
3    2: 'great white shark, white shark, man-eater, man-eating shark, Carcharodon carcharias',  
4    3: 'tiger shark, Galeocerdo cuvieri',  
5    4: 'hammerhead, hammerhead shark',  
6    5: 'electric ray, crampfish, numbfish, torpedo',  
7    6: 'stingray',  
8    7: 'cock',  
9    8: 'hen',  
10   9: 'ostrich, Struthio camelus',
```

- ▶ we restrict ourselves to $\mathcal{Y} = \{0, 1\}$ for the time being, but can be much larger in modern applications (1,000 for ImageNet)
- ▶ **Training data:** $S = ((x_1, y_1), \dots, (x_n, y_n))$ *finite* sequence of points of $\mathcal{X} \times \mathcal{Y}$
- ▶ also called *training set*

Hypothesis class

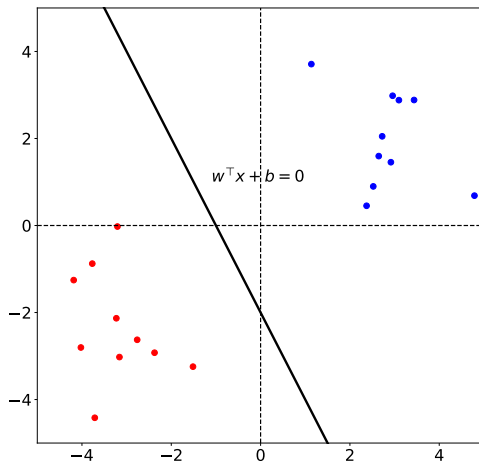
- ▶ **Hypothesis:** $h : \mathcal{X} \rightarrow \mathcal{Y}$ a prediction rule. also called *predictor*, *classifier* (in the context of classification)
- ▶ we are looking for a good h
- ▶ **Hypothesis class:** \mathcal{H} some space of functions. if no restrictions, set of all measurable functions
- ▶ **Example:** linear classifiers:

$$\mathcal{H} = \{h : x \mapsto \text{sign}(w^\top x + b), w \in \mathbb{R}^d, b \in \mathbb{R}\},$$

where $w^\top x$ denotes the scalar product between w and x

- ▶ given an algorithm A and a dataset S , we will write $h = A(S)$ the output of our algorithm on S

Example: linear classifier



Data generation

- ▶ **Data generation:** for now, we assume that there is a true distribution \mathcal{D} of the data on \mathcal{X}
- ▶ the training examples are i.i.d. samples from \mathcal{D}
- ▶ i.i.d.: *independent identically distributed*
- ▶ **Example:** sample images uniformly at random from a larger set (all the images on the internet)
- ▶ hard to satisfy: there is always a bias in the way your dataset is constructed

Assumption (noiseless setting): there exists a function $f : \mathcal{X} \rightarrow \mathcal{Y}$ such that $y = f(x)$ for any $x \in \mathcal{X}$.

- ▶ **Important:** we know neither \mathcal{D} nor f ! we only have access to S

Measure of success

- ▶ **Risk of a classifier:** probability that h does not return the correct label on a (new) random sample:

$$\mathcal{R}_{\mathcal{D},f}(h) = \mathbb{P}_{x \sim \mathcal{D}} (h(x) \neq f(x)) .$$

- ▶ **Intuition:** we want to be good, *on average*, for new samples of the same distribution
- ▶ subscript often omitted when clear from context
- ▶ also called *generalization error*, *true error* (notation L or \mathcal{E})
- ▶ **Important:** we want to find h with small generalization error. Ideally,

$$\mathcal{R}_{\mathcal{D},f}(h) = 0 .$$

- ▶ **Question:** how to do this?

2.3. Empirical risk minimization

Empirical risk minimization

- ▶ as we have seen, what we would *like* to do is find

$$h \in \arg \min_{h \in \mathcal{H}} \mathcal{R}_{\mathcal{D},f}(h) = \arg \min_{h \in \mathcal{H}} \mathbb{P}_{x \sim \mathcal{D}} (h(x) \neq f(x)) .$$

- ▶ **Problem:** we know neither \mathcal{D} nor f ...
- ▶ ...and even if we did it would still be a very difficult problem (there are *a lot* of measurable functions!)
- ▶ **Idea:** replace $\mathcal{R}_{\mathcal{D},f}$ by an *empirical* version
- ▶ empirical risk (or training error):

$$\hat{\mathcal{R}}_S(h) = \frac{1}{n} |\{i \in \{1, \dots, n\} \text{ s.t. } h(x_i) \neq y_i\}| ,$$

where $|E|$ denotes the cardinality of (finite) set E

- ▶ minimizing the empirical risk = empirical risk minimization⁷ (ERM)

⁷Vapnik, *Principles of risk minimization for learning theory*, NIPS, 1992

Exercise

Exercise: set $h \in \mathcal{H}$. Let n be a fixed integer.

1. Show that

$$\mathbb{E}_S \left[\hat{\mathcal{R}}_S(h) \right] = \mathcal{R}_{\mathcal{D},f}(h),$$

where the expectation is taken with respect to all i.i.d. draws of S .

2. Show that $\hat{\mathcal{R}}_S(h) \xrightarrow{\mathbb{P}} \mathcal{R}_{\mathcal{D},f}(h)$ when $n \rightarrow +\infty$.

Solution

1. First, we see that

$$|\{i \in \{1, \dots, n\} \text{ s.t. } h(x_i) \neq y_i\}| = \sum_{i=1}^n \mathbb{1}_{h(x_i) \neq y_i}.$$

Then we write

$$\begin{aligned} \mathbb{E} [\hat{\mathcal{R}}_S(h)] &= \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{h(x_i) \neq y_i} \right] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\mathbb{1}_{h(x_i) \neq y_i}] && \text{(linearity)} \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{P}(h(x_i) \neq y_i) && (\mathbb{E} [\mathbb{1}_A] = \mathbb{P}(A)) \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{P}(h(x_i) \neq f(x_i)) && \text{(noiseless assumption)} \end{aligned}$$

Solution, ctd.

Further, since the x_i are i.i.d., for any $1 \leq i \leq n$,

$$\mathbb{P}(h(x_i) \neq f(x_i)) = \mathbb{P}(h(x) \neq f(x)) .$$

We recognize the definition of the true risk. Therefore,

$$\begin{aligned}\mathbb{E} \left[\hat{\mathcal{R}}_S(h) \right] &= \frac{1}{n} \sum_{i=1}^n \mathcal{R}_{\mathcal{D},f}(h) \\ &= \mathcal{R}_{\mathcal{D},f}(h) \quad \quad \quad (\text{does not depend on } i)\end{aligned}$$

2. Since the x_i are i.i.d. random variables, so are the Z_i defined by

$$Z_i = \mathbb{1}_{h(x_i) \neq f(x_i)} .$$

Solution, ctd.

Moreover, the Z_i s are bounded almost surely (by 1). In particular, they are integrable. Therefore, we can use the law of large numbers and write

$$\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{h(x_i) \neq y_i} \xrightarrow{\mathbb{P}} \mathbb{E} [\mathbb{1}_{h(x_1) \neq f(x_1)}] .$$

From question 1., we deduce that

$$\hat{\mathcal{R}}_S(h) \xrightarrow{\mathbb{P}} \mathcal{R}_{\mathcal{D},f}(h) .$$



Reminder: the law of large numbers

Theorem (Weak Law of Large Numbers = WLLN): Let Z_1, Z_2, \dots be a sequence of i.i.d. random variables. Assume that $\mathbb{E}[|Z_1|] < +\infty$ and set $\mu := \mathbb{E}[Z_1]$. Then

$$\frac{Z_1 + \dots + Z_n}{n} \xrightarrow{\mathbb{P}} \mu.$$

- ▶ **Intuition:** average of measurements converges towards the true value
- ▶ stronger statement is true, *strong* law of large numbers, with almost sure convergence instead of in probability
- ▶ multivariate extension: coordinate-wise