

Any reproduction or distribution of this document, in whole or in part, is prohibited unless permission is granted by the authors.

EE-559

Deep Learning

What's on today?

- **Loss function**: training signal to optimize the parameters
- **Data**: how you can get the most from them
- **About the course**: the learning journey
- **Assessment**: earning your grade
- **Group mini-project**: where principles meet deep learning
- **Exercises**: hands-on practice on activation and loss functions

Practice exercises – students' questions

EE-559 Deep learning – Practice 1, Students' questions

A question is denoted by **Q**; the corresponding answer is denoted by **A**. The questions-related exercises are marked by their numbers in Practice_1.pdf and Practice_1.ipynb documents.

Environment

Q: How to create an environment in Noto?

A: Instructions for environment creation are mentioned in the first cell of Practice_1.ipynb. To run the commands, you need to open the bash terminal in Noto, which you can do by going to Launcher → Other → Terminal.

Q: Do you provide an environment for local installation ?

A: No, we do not provide an environment for local installation. We recommend using [Noto](#) platform for week 1, and [Gnoto](#) platform for the following weeks, as the exercises have been developed and tested on these platforms.

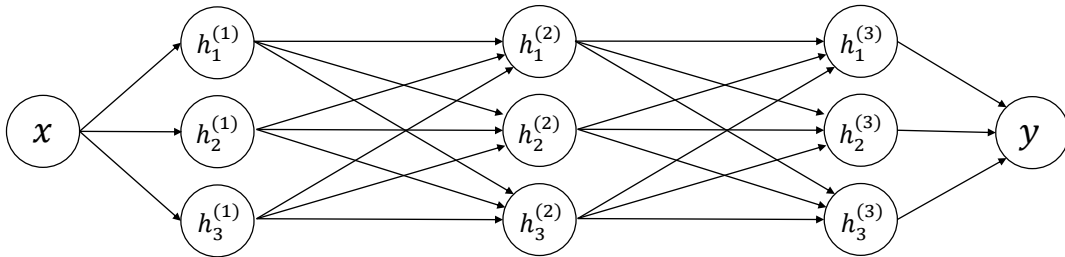
Q: When we call `model(input)` (which is equivalent to `model.forward(input)`) how does the system know that it has to go to the `forward()` definition inside the main class? If a new function is declared in the main class (say `forward2`) why does it not go there?

A: Every model class should be inherited from `nn.Module`. Let's say you want to create a class `MyModelClass(nn.Module)` and make an object `model=MyModelClass()`. The parent class `nn.Module` has two main functions that you have to redefine for your `MyModelClass`: `__init__` and `forward`. Calling `model(input)` is equivalent to calling `model.forward(input)` or `model.__call__(input)`. Additionally, parent `nn.Module` class has many functions that compute backwards pass, output parameters (`.parameters()`), load weights (`.load_state_dict()`) and so on, which means that the child class inherits those functions too. You do not have to implement any of those additional functions as they are implemented in the parent class and will automatically work correctly if you define `__init__` and `forward` functions for your

<https://go.epfl.ch/EE-559>

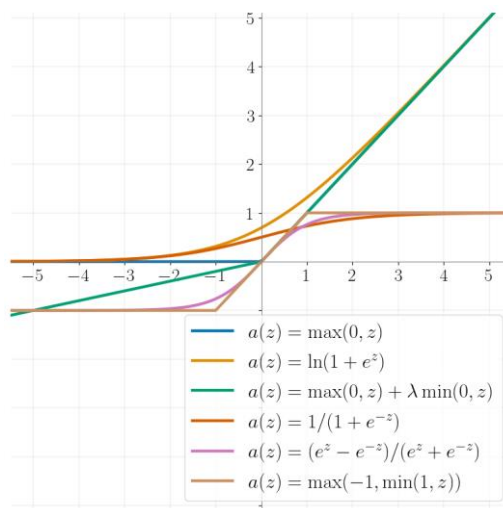
Recap

Network diagram & matrix notation



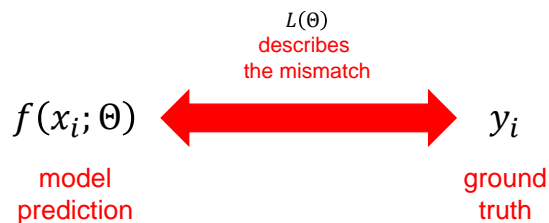
$$y = \theta_0^{(3)} + \boldsymbol{\theta}^{(3)} a \left[\theta_0^{(2)} + \boldsymbol{\theta}^{(2)} a \left[\theta_0^{(1)} + \boldsymbol{\theta}^{(1)} a \left[\theta_0^{(0)} + \boldsymbol{\theta}^{(0)} x \right] \right] \right]$$

Activation functions



Loss function

$L(\Theta)$: training signal



Recall:

$$L(\Theta) = \sum_{i=1}^N (f(x_i; \Theta) - y_i)^2$$

Least square loss
(for univariate regression)

Concepts:
Labels, annotation

Model parameters

$$y = f(x; \Theta)$$

(family of) family of possible relationships between x and y
 Θ : model parameters

$$\{x_i, y_i\}_{i=1}^N$$

training dataset of N input-output pairs

$$\Theta^* = \arg \min_{\Theta} L(\Theta)$$

minimization of the loss to determine the model parameters Θ

Prediction of y from x



Computing the parameters Θ of $P(y|\Theta)$ over the output space

Concept:

Conditional probability distribution

How do we build a loss function?

$y = f(x; \Theta)$ predicts y from x



compute $P(y|x)$ over the output space



encourage $P(y_i|x_i)$ to represent y_i with high probability



select a parametric distribution $P(y|\varphi)$ defined on the output space



use the network $y = f(x; \Theta)$ to determine the parameter(s) φ of the distribution

How do we build a loss function?

$\varphi_i = f(x_i; \theta)$ parameter of the distribution corresponding to training input x_i



each training output y_i has to have a high probability under $P(y_i|\varphi_i)$

Determining the parameters

$$\begin{aligned}
 \theta^* &= \arg \max_{\theta} \prod_{i=1}^N P(y_i|x_i) = \arg \max_{\theta} \prod_{i=1}^N P(y_i|\varphi_i) \\
 &= \arg \max_{\theta} \prod_{i=1}^N P(y_i|f(x_i; \theta)) && \text{maximum likelihood criterion} \\
 &= \arg \max_{\theta} \log \left[\prod_{i=1}^N P(y_i|f(x_i; \theta)) \right] \\
 &= \arg \max_{\theta} \sum_{i=1}^N \log[P(y_i|f(x_i; \theta))] && \text{log-likelihood criterion}
 \end{aligned}$$

Determining the parameters

$$\begin{aligned}
 \Theta^* &= \arg \max_{\Theta} \sum_{i=1}^N \log[P(y_i | f(x_i; \Theta))] \\
 &= \arg \min_{\Theta} \left[- \sum_{i=1}^N \log[P(y_i | f(x_i; \Theta))] \right] \quad \text{negative log-likelihood criterion} \\
 &= \arg \min_{\Theta} L(\Theta)
 \end{aligned}$$

Univariate regression

$$y \in \mathbb{R} \quad y = f(x; \Theta) \quad \Theta^* = \arg \max_{\Theta} \prod_{i=1}^N P(y_i | f(x_i; \Theta))$$

$$P(y | \varphi) = P(y | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\mu)^2}{2\sigma^2}} \quad \text{univariate normal distribution}$$

$$P(y | f(x; \Theta), \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-f(x; \Theta))^2}{2\sigma^2}} \quad f \text{ to compute } \mu$$

$$L(\Theta) = - \sum_{i=1}^N \log[P(y_i | f(x_i; \Theta), \sigma^2)] = - \sum_{i=1}^N \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - f(x_i; \Theta))^2}{2\sigma^2}} \right]$$

Loss minimization

$$\begin{aligned}
 \Theta^* &= \arg \min_{\Theta} \left[- \sum_{i=1}^N \log[P(y_i | f(x_i; \Theta))] \right] \\
 &= \arg \min_{\Theta} \left[- \sum_{i=1}^N \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - f(x_i; \Theta))^2}{2\sigma^2}} \right] \right] \\
 &= \arg \min_{\Theta} \left[- \sum_{i=1}^N \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \right] - \frac{(y_i - f(x_i; \Theta))^2}{2\sigma^2} \right] \\
 &= \arg \min_{\Theta} \left[\sum_{i=1}^N (y_i - f(x_i; \Theta))^2 \right] \quad \text{Least square loss}
 \end{aligned}$$

Inference

$$y \in \mathbb{R} \quad y = f(x; \Theta)$$

$P(y_i | f(x_i; \Theta))$ point estimate from the distribution

$$\hat{y} = \arg \max_y P(y | f(x; \Theta^*)) \quad \text{maximum is determined by } \mu \text{ of the normal distribution}$$

$$\hat{y} = f(x; \Theta^*)$$

The role of data

On the dimensionality of x

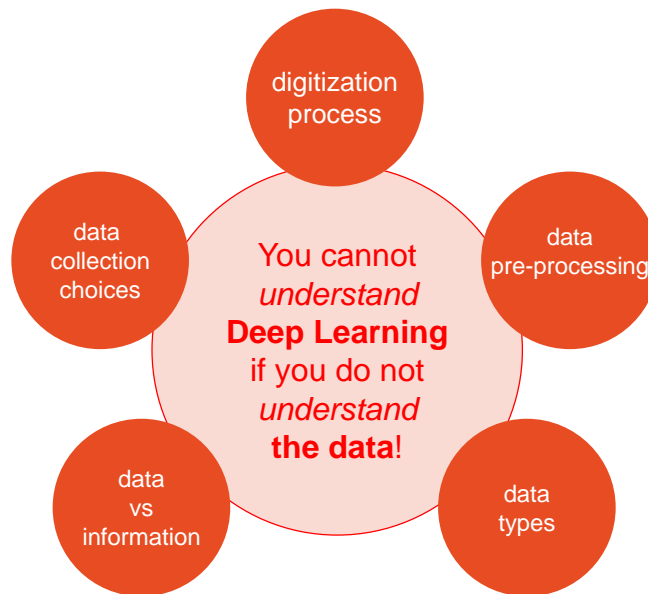
- Picture taken by an iPhone: e.g. 24,000,000 pixels (x 3 channels)
 - *do we need all those numbers?*
- Toy example: object recognition



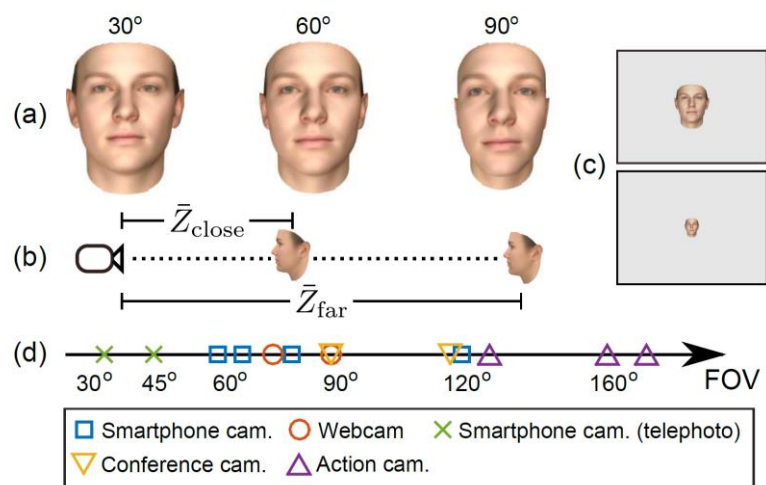
Concepts:

High-dimensional space, data capturing (digitization), data manifold

Data



Data collection process



Sariyanidi et al., *Can Facial Pose and Expression Be Separated With Weak Perspective Camera?* CVPR 2020

Data vs information

- How many people are here?
- Who is here?
- Which people are in a group?
- What is the building?
- What is the name of the peak?
- What is the weather like?
- Is there an exit anywhere?
- ...



Data containers

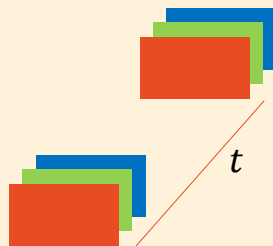
Tensors: higher dimensional arrays

- rank 0: scalar
- rank 1: vector
- rank 2: matrix

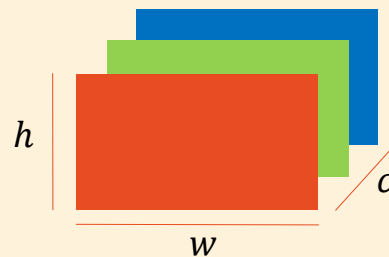
Rank 3: examples

video
video frame

$[x_t]_{t=1}^N$



x_t



Types of data

- Text
- Molecular genetic data
- Point clouds
- Social interaction (graph-structured data)
- Images (including multi-spectral)
- Videos
- Time series, ...

Time series: data points indexed over time

- Raw audio
- Bio-signals
 - heart rate monitoring (ECG)
 - brain monitoring (EEG)
- Weather data
 - rainfall measurements
 - temperature readings
- Financial data
 - quarterly sales
 - stock prices
- Control engineering
- Communication engineering
- ...

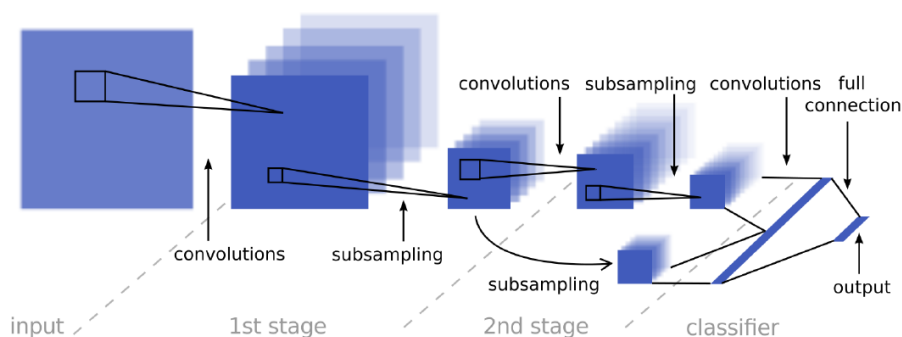
Concepts:

Time vs frequency domain, seasonality

Know your data

Convolutional Neural Networks (CNNs)

Learn image representations (spatial features) from pixel values



[arXiv:1204.3968](https://arxiv.org/abs/1204.3968)

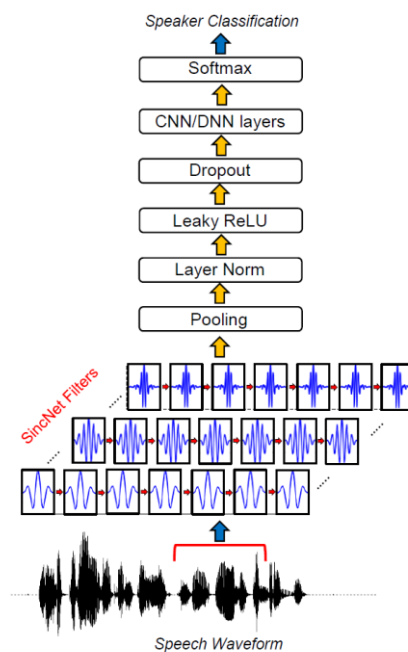
Know your data

Learn low-level speech representations from waveforms

Only **low and high cut-off frequencies** are directly learned from data (unlike standard CNNs that learn **all** elements of each filter)

Capture important narrow-band speaker characteristics (e.g. pitch and formants)

Customized filter bank: parametrized sinc functions, which implement band-pass filters



[arXiv:1808.00158](https://arxiv.org/abs/1808.00158)

Secondary information in time-series data

- Audio data (voice)

- height & weight
- emotional state
- health conditions

Krauss et al. "Inferring speakers' physical attributes from their voices"

Trigeorgis et al. "Adieu features? End-to-end speech emotion recognition using a deep convolutional recurrent network"

Schuller et al. "The INTERSPEECH 2013 computational paralinguistics challenge: Social signals conflict emotion autism"

- Motion data (wearables)

- height & weight
- level of activity
- changes in behaviors

Masuda & Maekawa "Estimating physical characteristics with body-worn accelerometers based on activity similarities"

Zainudin et al. "Monitoring daily fitness activity using accelerometer sensor fusion"

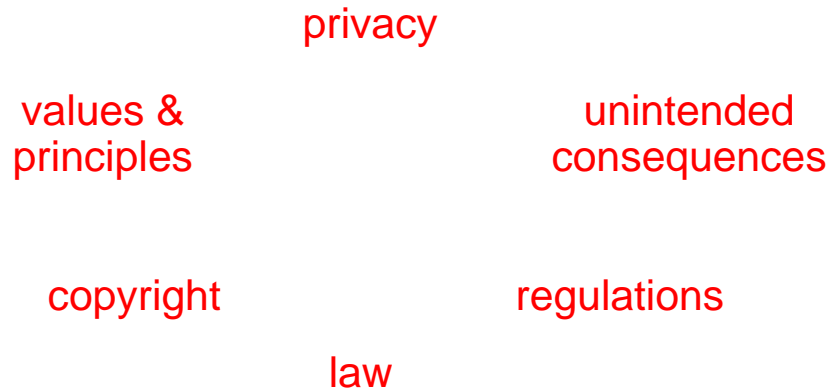
Gruenerbl et al. "Using smartphone mobility traces for the diagnosis of depressive and manic episodes in bipolar patients"

Art.4 GDPR

'personal data' means any information relating to an identified or identifiable natural person ('data subject');

an identifiable natural person is one who can be identified, directly or indirectly, in particular by reference to an identifier such as a name, an identification number, location data, an online identifier or to one or more factors specific to the physical, physiological, genetic, mental, economic, cultural or social identity of that natural person

Data: asset or liability?



WIRED

KATE KNIBBS BUSINESS JAN 9, 2025 5:33 PM

Meta Secretly Trained Its AI on a Notorious Piracy Database, Newly Unredacted Court Docs Reveal

One of the most important AI copyright legal battles just took a major turn.



<https://www.wired.com/story/new-documents-unredacted-meta-copyright-ai-lawsuit/>

The New York Times

The Times Sues OpenAI and Microsoft Over A.I. Use of Copyrighted Work

Millions of articles from The New York Times were used to train chatbots that now compete with it, the lawsuit said.

<https://www.nytimes.com/2023/12/27/business/media/new-york-times-open-ai-microsoft-lawsuit.html>

MIT
Technology
Review

Featured Topics Newsletters

In 2016, hoping to spur advancements in facial recognition, Microsoft released the largest face database in the world. Called MS-Celeb-1M, it contained 10 million images of 100,000 celebrities' faces. "Celebrity" was loosely defined, though.

Three years later, researchers Adam Harvey and Jules LaPlace scoured the data set and found many ordinary individuals, like journalists, artists, activists, and academics, who maintain an online presence for their professional lives. None had given consent to be included, and yet their faces had found their way into the database and beyond; research using the collection of faces was conducted by companies including Facebook, IBM, Baidu, and SenseTime, one of China's largest facial recognition giants, which sells its technology to the Chinese police.

Shortly after Harvey and LaPlace's investigation, and after receiving criticism from journalists, Microsoft removed the data set, stating simply: "The research challenge is over." But the privacy concerns it created linger in an internet forever-land. And this case is hardly the only one.

<https://www.technologyreview.com/2021/08/13/1031836/ai-ethics-responsible-data-stewardship/>

TECH / ARTIFICIAL INTELLIGENCE

AI image training dataset found to include child sexual abuse imagery



/ Stanford researchers discovered LAION-5B, used by Stable Diffusion, included thousands of links to CSAM.

A popular training dataset for AI image generation contained links to child abuse imagery, [Stanford's Internet Observatory found](#), potentially allowing AI models to create harmful content.

LAION-5B, a dataset used by Stable Diffusion creator Stability AI, included at least 1,679 illegal images scraped from social media posts and popular adult websites.

<https://www.theverge.com/2023/12/20/24009418/generative-ai-image-laion-csam-google-stability-stanford>

The **Register**



Google says public data is fair game for training its AIs

Hey, we're just being honest, says web giant

[Katyanna Quach](#)

Thu 6 Jul 2023 // 01:50 UTC

Google has updated its privacy policy to confirm it scrapes public data from the internet to train its AI models and services – including its chatbot Bard and its cloud-hosted products.

The [fine print](#) under research and development now reads: "Google uses information to improve our services and to develop new products, features and technologies that benefit our users and the public. For example, we use publicly available information to help train Google's AI models and build products and features like Google Translate, Bard and Cloud AI capabilities."

https://www.theregister.com/2023/07/06/google_ai_models_internet_scraping

TECH Help Desk Artificial Intelligence Internet Culture Space Tech Policy

Google takes down Gemini AI image generator. Here's what you need to know.

Critics said the company's tool created images of a woman pope and Black founding father

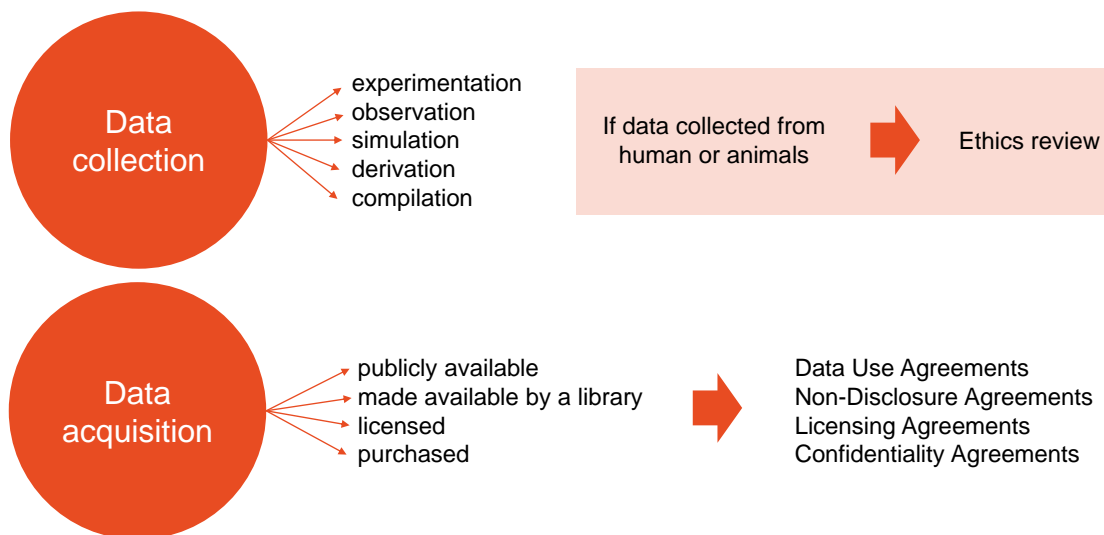
By [Gerrit De Vynck](#) and [Nitasha Tiku](#)

Updated February 23, 2024 at 3:27 p.m. EST | Published February 22, 2024 at 10:54 p.m. EST



<https://www.washingtonpost.com/technology/2024/02/22/google-gemini-ai-image-generation-pause>

Data collection and acquisition



<https://researchdatamanagement.harvard.edu/acquisition-agreements>

Compliance

Human Research Ethics Committee (HREC)

Animal Research Ethics Committee (AREC)

Research Integrity

Research Misconduct

Human Research Ethics Committee (HREC)

EPFL is committed to respect for the basic principles of research ethics. The EPFL Human Research Ethics Committee reviews projects conducted by EPFL researchers that involve human participants and/or personal data. These projects do not fall within the scope of the Federal Act on Research with Human Beings (HRA).



HREC at a glance



Review process

<https://www.epfl.ch/research/ethic-statement/human-research-ethics-committee/>

PUBLISHING SUPPORT

Open Access

Financial support for Open Access

Research and copyright

Infoscience

ORCID

Text and Data Mining

RESEARCH DATA MANAGEMENT

Data planning and guidelines

Active data management

Data publication

ACQUA: Long-term preservation

Data policies

Data services, expertise, tools and training

OPEN SCIENCE

Search in the BEAST catalog

Registration

My account

New acquisitions

Data planning and guidelines

Data planning and guidelines

Active data management

Data publication

ACQUA: Long-term preservation

Data policies

Data services, expertise, tools and training

At EPFL, the Library Research Data Management team is at your disposal to provide you with expert advice, support, and solutions.

<https://www.epfl.ch/campus/library/services-researchers/data-planning-guidelines/>

The Data Provenance Initiative

A Large Scale Audit of Dataset Licensing & Attribution in AI



Data Provenance Explorer

The Data Provenance Initiative is a large-scale audit of AI datasets used to train large language models. As a first step, we've traced 1800+ popular, text-to-text finetuning datasets from origin to creation, cataloging their data sources, licenses, creators, and other metadata, for researchers to explore using this tool. The purpose of this work is to improve transparency, documentation, and informed use of datasets in AI.

You can download this data (with filters) directly from the [Data Provenance Collection](#).

[arXiv:2310.16787](#)

Concepts:

Data transparency, attribution, dataset documentation, ethical and legal risks.

About the course

Instructor and TAs



Andrea
Cavallaro



Egor
Rumiantsev



Qin
Liu



Corentin
Genton



Ti
Wang



Olena
Hrynenko



Darya
Baranouskaya

EE-559: Deep Learning

This course explores how to design reliable **discriminative** and **generative** neural networks, the ethics of **data** acquisition and **model** deployment, as well as modern **multi-modal** models.

<https://go.epfl.ch/EE-559>

Course content

- Loss functions, data and labels, data provenance
- Training models: gradients and initialization
- Generalization and performance
- Graph neural networks and transformers
- Multi-modal models
- Generative adversarial networks
- Variational autoencoders
- Diffusion models
- Interpretability, explanations, bias and fairness
- Principles and regulations

Applications:

Natural language processing, audio processing, computer vision, robotics, biology, science

Learning outcomes

- By the end of the course, you will be able to:
 - **Justify** the choices for training and testing a deep learning model
 - **Interpret** the performance of a deep learning model
 - **Analyze** the limitations of a deep learning model
 - **Propose new solutions** for a given application

Prerequisites for this course

- Linear algebra
- Differential calculus
- Python programming
- Basics in
 - probabilities and statistics
 - optimization
 - algorithmic
 - signal processing

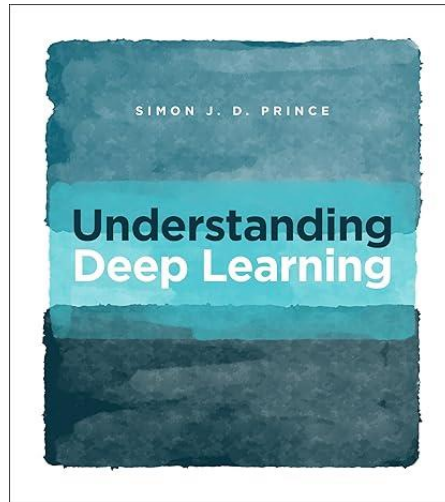
Concepts to start the course:

Discrete and continuous distributions, normal density, law of large numbers, conditional probabilities, Bayes, PCA, vector, matrix operations, Euclidean spaces, Jacobian, Hessian, chain rule, notion of minima, gradient descent, computational costs, Fourier transform, convolution.

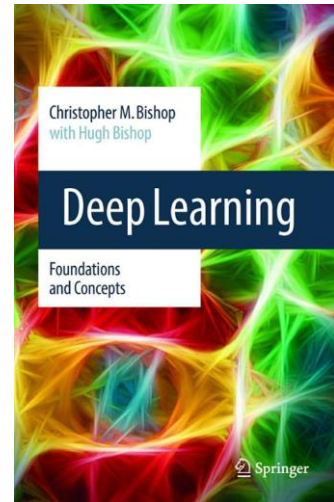
Activities

- Attending **lectures**
- Participating in class discussion (using slido)
- Completing **exercises** (using Python, *please come with your laptop*)
 - *practice exercises*
 - *marked exercises*
- Choosing, completing and presenting a **group project**
- Reading written material: books and scientific papers

Books

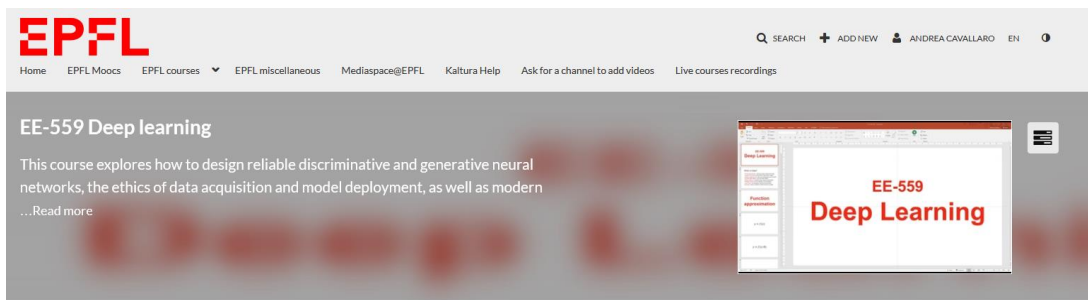


<https://udlbook.github.io/udlbook/>



<https://www.bishopbook.com/>

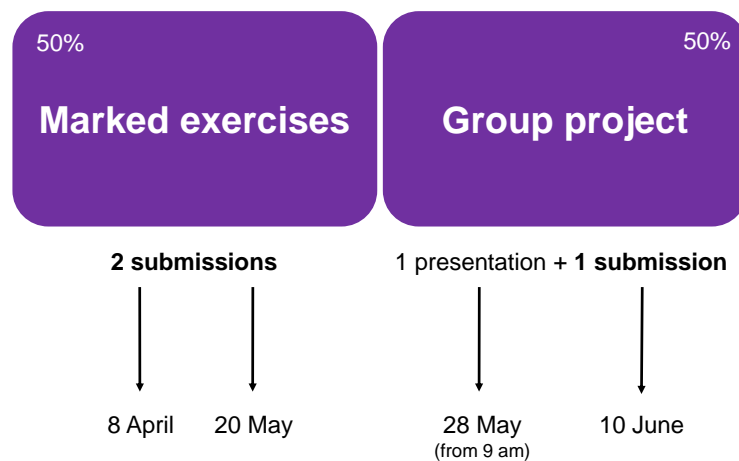
Recording of the lectures



<https://mediaspace.epfl.ch/channel/EE-559+Deep+learning>

Assessment

Assessment



Overall mark

- The final EE-559 mark M_{EE-559} is a weighted average of the marks of two *marked exercises* submissions (S_1 , S_2) and the *group mini-project* (P):

$$M_{EE-559} = 0.25S_1 + 0.25S_2 + 0.5P$$

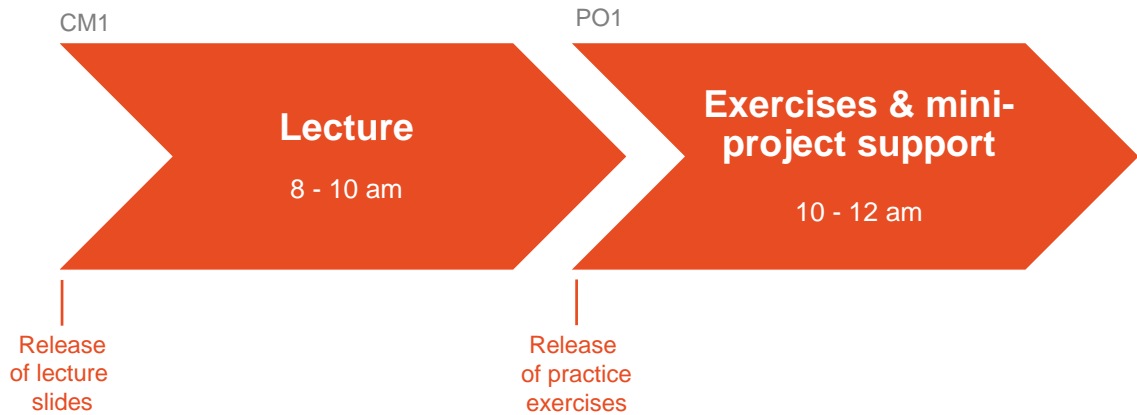
- The highest mark for each submission (S_1 , S_2) and the mini-project (P) is 60
- The outcomes of the *group mini-project* will be presented at an open poster session

Submissions of marked exercises

- Practice exercises*
 - to prepare for submissions S_1 and S_2
 - questions released on Wednesdays at 10 am
 - solutions released on Fridays at 3 pm
- Marked exercises*
 - the content of submission S_1 and S_2 spans several lectures

Submission	Weight	Deadline
S_1	25%	8 Apr, 11:59 pm
S_2	25%	20 May, 11:59 pm

Wednesday's activities



Group Mini-Project

Group mini-project

- Working on a deep learning problem in a group of three students
 - group registration to open next week
- Deliverables
 - the code + a screencast of the code running
 - a 3-page paper (*template provided*)
 - a poster (*template provided*)

Weight	Group formed by	Presentation and Q&A	Deliverables deadline
50%	18 March, 11.59 pm	28 May, from 9 am	10 June, 11:59 pm

Group mini-project

- **Assessment criteria**
 - Literature review and methodology (max 20 marks)
 - Evaluation, testing, and analysis (max 20 marks)
 - Communication of the findings (max 20 marks)

Group mini-project: assessment criteria (1/3)

- **Literature review and methodology**

- **thoroughness** of literature review
- **clarity** of the objectives and problem definition
- evidence of **creativity and novelty** of the adopted methodology

Group mini-project: assessment criteria (2/3)

- **Evaluation, testing, and analysis**

- quality of the **results** of the proposed solution and their **analysis**
- **discussion** of the limitations of the proposed solution
- **justification** of the choices for the experiments
- evidence of **critical thinking**

Group mini-project: assessment criteria (3/3)

- **Communication of the findings**

- ability to clearly communicate the findings
 - in the poster
 - during the oral presentation of the poster
 - in the written report

Group Mini-Project Title

Name Surname 1¹ Name Surname 2¹ Name Surname 3¹

Abstract

Add your abstract here. Mention the issue(s) you have addressed, why they are important, and describe your proposed solution. Do not edit the style of this document (e.g., font size, margins) and do not to exceed the 3-page limit.

Keywords: add your keywords here.

1. Introduction

Add your introduction here. Select a current limitation or a relevant new problem; identify a specific problem you want to address. Clarify the objective(s) and the problem definition. State any hypotheses you made and reference the sources you used. Some common L^AT_EX commands are listed below.

Section: you can refer to a section as Sec. 2.

Equation:

$$x + y = 0. \tag{1}$$

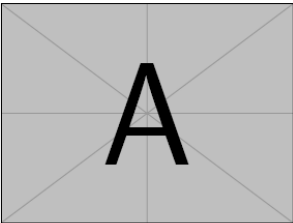


Figure 1. Insert your caption here.

DATA SET	NAIVE	FLEXIBLE	BETTER?
CLEVELAND	83.3± 0.6	80.0± 0.6	×
GLASS2	61.9± 1.4	83.8± 0.7	✓
CREDIT	74.8± 0.5	78.3± 0.6	

Table 1. INSERT YOUR CAPTION HERE.

You can refer to an equation as Equation 1.

Figure: you can refer to a figure as Fig. 1.

Table: you can refer to a table as Tab. 1 in your text. <https://www.tablesgenerator.com/> is useful for creating custom tables.

Reference: you can cite a source using command `\cite`, e.g. [1]. To add a reference, you can find your source on Google Scholar, click "Cite" and select BibTeX. Then copy the reference to `main.bib`.

2. Related Work

Add your literature review here. **Discuss the limitations** of the literature.

3. Method

Use what you have learnt in EE-559 to address the limitations you identified. Describe and **motivate** your methodol-

¹Group GroupNumber.

ogy.

4. Validation


Implement your ideas and test them. Add your evaluation, testing and **analysis** here. Justify the **choices** for the experiments. Analyse the results and the performance: why does your hypothesis work / doesn't work? **Compare** with alternative ideas / hypotheses. Discuss the **limitations** of the proposed solution.

5. Conclusion


Add your conclusion here.

References

- [1] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.



Group Mini-Project Title
Name Surname 1, Name Surname 2, Name Surname 3
Group GroupNumber



Problem definition

- Add your Problem Definition here.
- Note: do not copy-paste text from your report.

Validation

- Add your Evaluation, Testing and Analysis here.
- You can add plots, tables with results to showcase your method.


Dataset	Naïve	Flexible	Better?
CLEVELAND	83.3 ± 0.5	80.0 ± 0.5	+
GLASS2	61.9 ± 1.4	63.6 ± 0.7	✓
CREDIT	74.8 ± 0.5	76.3 ± 0.6	

Key Related Works

- Add your Related Works here.

Method

- Add your Method here.
- You can add diagrams and/or formulas to explain your method.



Dataset(s)

- Add Datasets that you use here.

Limitations

- Add your Limitations here.

Conclusion

- Add Conclusions here.

References

Add your References here using IEEE referencing style. For example:
 [1] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.

EE-559: Deep Learning, 2025

Mini-project: timeline



Mini-project Q&A session in the lab every Wednesday.
Drop by when you need help!

Deep learning to foster safer online spaces

The group mini-project aims to support a safer online environment by tackling hate speech in various forms, ranging from **text** and **images** to **memes**, **videos**, and **audio** content.

Deep learning to foster safer online spaces

The objective of the mini-project is to develop deep learning models that help foster healthier online interactions by **automatically identifying hate speech across diverse content formats**.

These deep learning models shall be carefully designed to prioritize **accuracy** and **context comprehension**, ensuring they differentiate between harmful hate speech and **legitimate critical discourse or satire**.

Developing deep learning models that help prevent the surfacing of hateful rhetoric will lead to a **more respectful online environment** where diverse voices can coexist and thrive.

Deep learning to foster safer online spaces

What is hate speech?

In common language, "hate speech" refers to offensive discourse targeting a group or an individual based on inherent characteristics (such as race, religion or gender) and that may threaten social peace.

To provide a unified framework for the United Nations to address the issue globally, the UN Strategy and Plan of Action on Hate Speech defines hate speech as...*"any kind of communication in speech, writing or behaviour, that attacks or uses pejorative or discriminatory language with reference to a person or a group on the basis of who they are, in other words, based on their religion, ethnicity, nationality, race, colour, descent, gender or other identity factor."*

However, to date there is no universal definition of hate speech under international human rights law. The concept is still under discussion, especially in relation to freedom of opinion and expression, non-discrimination and equality.



<https://www.un.org/en/hate-speech/understanding-hate-speech/what-is-hate-speech>

Deep learning to foster safer online spaces

1

Hate speech can be conveyed through any form of expression, including **images, cartoons, memes, objects, gestures and symbols** and it can be disseminated offline or online.

2

Hate speech is “**discriminatory**” (biased, bigoted or intolerant) or “**pejorative**” (prejudiced, contemptuous or demeaning) of an individual or group.

3

Hate speech calls out real or perceived “**identity factors**” of an individual or a group, including: “**religion, ethnicity, nationality, race, colour, descent, gender,**” but also characteristics such as language, economic or social origin, disability, health status, or sexual orientation, among many others.

It's important to note that hate speech can only be directed at individuals or groups of individuals. It does not include communication about States and their offices, symbols or public officials, nor about religious leaders or tenets of faith.



<https://www.un.org/en/hate-speech/understanding-hate-speech/what-is-hate-speech>

Your mini-project choices

To be aligned with the course learning objectives and your broader transversal skill development

- **Learning objectives**
 - interpret the performance of a deep learning model
 - analyze the limitations of a deep learning model
 - justify the choices for training and testing a deep learning model
 - propose new solutions for a given application
- **Transversal skill development**
 - respect relevant legal guidelines and ethical codes
 - take account of the social and human dimensions
 - ...

Framing your mini-project

- **Select a current limitation or a relevant new problem**
 - identify a specific **problem** you want to address
 - discuss in the report the **limitations** of the literature
 - use what you have learnt in EE-559 to address the limitations you identified
- **Implement your ideas**
- **Test them**
- **Analyse the performance**
 - why does your hypothesis work / doesn't work?
 - **compare** with alternative ideas / hypotheses

Do not forget to ...

- State in the report any **hypotheses** you made
- Reference the **sources** you use
- **Comment** your code

Data

Objective: high-quality, well-labeled data

Get familiar with the benchmarks on the task you are addressing

Do you really need to train a model from scratch?

Can you apply fine-tuning?

Check if your model overfits

Discuss your choices in the report

The amount of data required for training a neural network depends on several factors, such as problem complexity and model size

Model

- **Multi-modal models are cool!**

however check if they are necessary for your problem

- **You can use a pre-trained model as a basis of your project**

justify in the report why you are using a pre-trained model

- **Having the best-performing model in its class is highly gratifying!**

however this is not the main goal of the mini-project

it depends on the other contributions/innovations with respect to the learning outcomes (performance can be slightly lower if compensated by innovation/creativity elsewhere)

After the group mini-project ...

Based on the assessment,
qualifying groups will be offered one of two
mentorship opportunities:

- 1) support to develop further your project for the
submission of a scientific paper
- 2) guidance and sponsorship to
build & launch your start-up!

Today's practice exercises

You will

- investigate **data-induced bias** within and across text encoders
- familiarize yourself with
 - **loading** image datasets
 - creating dataset loaders
 - image **augmentation**
- analyse the **licenses** of images

Today's marked exercises

You will

- implement **activation functions**
 - sigmoid
 - leaky ReLU
- implement **loss functions**
 - mean squared error (MSE)
 - contrastive

What did we learn today?

- Loss function
- Data
- About the course
- Assessment
- Group mini-project
- Exercises

EE-559

Deep Learning

andrea.cavallaro@epfl.ch