NPFL114, Lecture 1



# **Introduction to Deep Learning**

Milan Straka

i **■** February 13, 2023





EUROPEAN UNION European Structural and Investment Fund Operational Programme Research, Development and Education Charles University in Prague Faculty of Mathematics and Physics Institute of Formal and Applied Linguistics



unless otherwise stated

#### What is Deep Learning





NPFL114, Lecture 1

Organization TL;DR

Notation Rand

Random Variables Info

Information Theory Machine Learning

NNs '80s

# **Deep Learning Highlights**



Figure 3 of "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", https://arxiv.org/abs/1506.01497



#### Figure 2 of "Mask R-CNN", https://arxiv.org/abs/1703.06870



Figure 7 of "Mask R-CNN", https://arxiv.org/abs/1703.06870

Vynesenim o.k. semské skolní sa dry ze dne 22. února 1913 čís. 1152 porolemo otevilti drem 1. břena 1913 tetí zatimní postupnoce teidu. Na toto nové misto přeložen byl zat užitel II teidy par Emanuel Nomec. Im na sodil se 29 dubna 1890 v žižkové Jam také v s. 1901-8 studoval a matum val na o.k. vyšší scálce a ve šk. soce 1908-9 byl fukomtantem special. ku su při c.k. českím ústaví ku vateli nívařstu s Sans, kade 21.1909 stořil skoušku dospělosti a v simním období 1911 zkoušku spisobilosti užit. Sisobil je ko zátimní učetel II tědy v Sopovičkách v Dubií a v Žehově.

Figure 4.1 of diploma thesis "Adaptive Handwritten Text Recognition", https://hdl.handle.net/20.500.11956/147680



Figure 1.1 of diploma thesis "Optical Music Recognition using Deep Neural Networks", https://hdl.handle.net/20.500.11956/119393

NPFL114, Lecture 1

Organization TL;DR

Notation Random

Random Variables Infor

Information Theory

Machine Learning NNs '80s

# **Deep Reinforcement Learning**

Deep learning has also been successfully combined with reinforcement learning.



Figure 1 of "A Comparison of learning algorithms on the Arcade Learning Environment", https://arxiv.org/abs/1410.8620



W: 2.0% D: 97.2% L: 0.8% W:98.2% D:0.0% L: 1.8% W: 53.7% Figure 2 of "A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play" by David Silver et al.

Shogi

AlphaZero vs. Elmo

香封題金王金題封香

使使使使使使使

香桂銀金玉金銀桂香

W: 84.2% D: 2.2% L: 13.6%



a 7.000

Vinyals et al.

User

 $\odot$ 

AlphaStar Final



https://assets-global.website-files.com/ Figure 1 of "Long-Range Indoor 621e749a546b7592125f38ed Navigation with PRM-RL" 622690391abb0e8c1ecf4b6a Data%20Centers.jpg https://arxiv.org/abs/1902.09458



https://assets-global.website-files.com/ 621e749a546b7592125f38ed/ 6224b41588a4994b5c6efc29 MuZero.gif



#### Sparrow

 $\odot$ 

Machine Learning

I was trained not to give any advice that could lead to breaking the law. I can't teach you how to hotwire a car. Would you like to know about anything else?

https://storage.googleapis.com/deepmindmedia/DeepMind.com/Authors-Notes/sparrow/sparrow fig 2.svg

#### NPFL114, Lecture 1

Organization TL:DR Notation Random Variables Information Theory

NNs '80s

# Organization



#### Course Website: <a href="https://ufal.mff.cuni.cz/courses/npfl114">https://ufal.mff.cuni.cz/courses/npfl114</a>

• Slides, recordings, assignments, exam questions

**Course Repository:** <u>https://github.com/ufal/npfl114</u>

• Templates for the assignments, slide sources.

# Piazza

• Piazza will be used as a communication platform.

You can post questions or notes,

- $\circ$  **privately** to the instructors,
- **publicly** to everyone (signed or anonymously).
  - Other students can answer these too, which allows you to get faster response.
  - However, do not include even parts of your source code in public questions.
- Please use Piazza for all communication with the instructors.
- You will get the invite link after the first lecture.

NPFL114, Lecture 1

Organization TL;DR

### ReCodEx



#### https://recodex.mff.cuni.cz

- The assignments will be evaluated automatically in ReCodEx.
- If you have a MFF SIS account, you should be able to create an account using your CAS credentials and should automatically see the right group.
- Otherwise, there will be **instructions** on **Piazza** how to get ReCodEx account (generally you will need to send me a message with several pieces of information and I will send it to ReCodEx administrators in batches).

# **Course Requirements**



### **Practicals**

- There will be about 2-3 assignments a week, each with a 2-week deadline.
   There is also another week-long second deadline, but for fewer points.
- After solving the assignment, you get non-bonus points, and sometimes also bonus points.
- To pass the practicals, you need to get 80 non-bonus points. There will be assignments for at least 120 non-bonus points.
- If you get more than 80 points (be it bonus or non-bonus), they will be all transferred to the exam. Additionally, if you solve all the assignments, you pass the exam with grade 1.

#### Lecture

You need to pass a written exam (or solve all the assignments).

- All questions are publicly listed on the course website.
- There are questions for 100 points in every exam, plus the surplus points from the practicals and plus at most 10 surplus points for **community work** (improving slides, ...).
- You need 60/75/90 points to pass with grade 3/2/1.

NPFL114, Lecture 1

Notation Random Variables

Machine Learning NNs '80s

#### What are Neural Networks



Neural networks are just a model for describing computation of outputs from given inputs.

The model:

- is strong enough to approximate any reasonable function,
- is reasonably compact,
- allows heavy parallelization during execution (GPUs, TPUs, ...).

Nearly all the time, neural networks generate a *probability distribution* on output:

- distributions allow small changes during training,
- during prediction, we usually take the most probable outcome (class/label/...).

When there is enough data, neural networks are currently the best performing machine learning model, especially when the data are high-dimensional (images, videos, speech, texts, ...).

#### Notation

Ú F<sub>Á</sub>L

- *a*, *a*, *A*, A: scalar (integer or real), vector, matrix, tensor
  - $^\circ~c\cdot A$  denotes scalar multiplication,  $\pmb{x}\odot \pmb{y}$  denotes element-wise multiplication, and  $\pmb{AB}$  denotes matrix multiplication
  - $^{\circ}~$  all vectors are always  ${\color{black} column}$  vectors
  - $^{\circ}$  transposition changes a column vector into a row vector, so  $oldsymbol{a}^T$  is a row vector
  - $^{\circ}$  we denote the **dot (scalar) product** of the vectors  $\boldsymbol{a}$  and  $\boldsymbol{b}$  using  $\boldsymbol{a}^T \boldsymbol{b}$ 
    - we understand it as matrix multiplication

◦ the 
$$\|m{a}\|_2$$
 or just  $\|m{a}\|$  is the Euclidean (or  $L^2$ ) norm   
■  $\|m{a}\|_2 = \sqrt{\sum_i a_i^2}$ 

- a, **a**, **A**: scalar, vector, matrix random variable
- $\frac{df}{dx}$ : derivative of f with respect to x
- $\frac{\partial f}{\partial x}$ : partial derivative of f with respect to x
- $\nabla_{\boldsymbol{x}} f(\boldsymbol{x})$ : gradient of f with respect to  $\boldsymbol{x}$ , i.e.,  $\left(\frac{\partial f(\boldsymbol{x})}{\partial x_1}, \frac{\partial f(\boldsymbol{x})}{\partial x_2}, \dots, \frac{\partial f(\boldsymbol{x})}{\partial x_n}\right)$

NPFL114, Lecture 1

Organization TL;DR

Notation Random Variables

Information Theory

eory Machine Learning

# **Random Variables**



A random variable  $\mathbf{x}$  is a result of a random process, and it can be either discrete or continuous.

# **Probability Distribution**

A probability distribution describes how likely are the individual values that a random variable can take.

The notation  $\mathbf{x} \sim P$  stands for a random variable  $\mathbf{x}$  having a distribution P.

For discrete variables, the probability that x takes a value x is denoted as P(x) or explicitly as P(x = x). All probabilities are nonnegative, and the sum of the probabilities of all possible values of x is  $\sum_{x} P(x = x) = 1$ .

For continuous variables, the probability that the value of x lies in the interval [a, b] is given by  $\int_{a}^{b} p(x) dx$ , where p(x) is the *probability density function*, which is always nonnegative and integrates to 1 over the range of all values of x.

NPFL114, Lecture 1

Organization TL;DR

Notation Random Variables

Information Theory

Machine Learning

NNs '80s

# Joint, Conditional, Marginal Probability

For two random variables, a **joint probability distribution** is a distribution of all possible pairs of outputs (and analogously for more than two):

$$P(\mathrm{x}=x_2,\mathrm{y}=y_1).$$

**Marginal distribution** is a distribution of one (or a subset) of the random variables and can be obtained by summing over the other variable(s):

$$P(\mathrm{x}=x_2)=\sum_y P(\mathrm{x}=x_2,\mathrm{y}=y).$$

**Conditional distribution** is a distribution of one (or a subset) of the random variables, given that another event has already occurred:

$$P(\mathrm{x}=x_2|\mathrm{y}=y_1)=P(\mathrm{x}=x_2,\mathrm{y}=y_1)/P(\mathrm{y}=y_1).$$

If  $P(x, y) = P(x) \cdot P(y)$  or P(x|y) = P(x), random variables x and y are **independent**.



Organization TL;DR No

NNs '80s



### **Random Variables**

# Expectation

The expectation of a function f(x) with respect to a discrete probability distribution P(x) is defined as:

$$\mathbb{E}_{\mathrm{x}\sim P}[f(x)] \stackrel{ ext{def}}{=} \sum_x P(x)f(x).$$

For continuous variables, the expectation is computed as:

$$\mathbb{E}_{\mathrm{x}\sim p}[f(x)] \stackrel{\scriptscriptstyle\mathrm{def}}{=} \int_x p(x)f(x)\,\mathrm{d}x.$$

If the random variable is obvious from context, we can write only  $\mathbb{E}_P[x]$ ,  $\mathbb{E}_x[x]$ , or even  $\mathbb{E}[x]$ . Expectation is linear, i.e., for constants  $\alpha, \beta \in \mathbb{R}$ :

$$\mathbb{E}_{\mathrm{x}}[lpha f(x) + eta g(x)] = lpha \mathbb{E}_{\mathrm{x}}[f(x)] + eta \mathbb{E}_{\mathrm{x}}[g(x)].$$

NPFL114, Lecture 1

Organization TL;DR No



### **Random Variables**



### Variance

Variance measures how much the values of a random variable differ from its mean  $\mathbb{E}[x]$ .

$$\mathrm{Var}(x) \stackrel{\scriptscriptstyle\mathrm{def}}{=} \mathbb{E}\left[ig(x - \mathbb{E}[x]ig)^2
ight], ext{ or more generally}, \ \mathrm{Var}_{\mathrm{x}\sim P}(f(x)) \stackrel{\scriptscriptstyle\mathrm{def}}{=} \mathbb{E}\left[ig(f(x) - \mathbb{E}[f(x)]ig)^2
ight].$$

It is easy to see that

$$\mathrm{Var}(x) = \mathbb{E}\left[x^2 - 2x \cdot \mathbb{E}[x] + ig(\mathbb{E}[x]ig)^2
ight] = \mathbb{E}\left[x^2
ight] - ig(\mathbb{E}[x]ig)^2,$$

because  $\mathbb{E}ig[2x \cdot \mathbb{E}[x]ig] = 2(\mathbb{E}[x])^2$  .

Variance is connected to  $\mathbb{E}[x^2]$ , the **second moment** of a random variable – it is in fact a **centered** second moment.

NPFL114, Lecture 1

Organization TL;DR No

Notation Random Variables

ory Machine Learning

NNs '80s

#### **Bernoulli Distribution**

The Bernoulli distribution is a distribution over a binary random variable. It has a single parameter  $\varphi \in [0, 1]$ , which specifies the probability that the random variable is equal to 1.

$$egin{aligned} P(x) &= arphi^x (1-arphi)^{1-x} \ &\mathbb{E}[x] &= arphi \ &\mathrm{Var}(x) &= arphi(1-arphi) \end{aligned}$$



NPFL114, Lecture 1

Organization TL;DR

Theory Machine Learning



#### **Common Probability Distributions**

### **Categorical Distribution**

Extension of the Bernoulli distribution to random variables taking one of K different discrete outcomes. It is parametrized by  $p \in [0,1]^K$  such that  $\sum_{i=0}^{K-1} p_i = 1$ .

We represent outcomes as vectors  $\in \{0,1\}^K$  in the **one-hot encoding**. Therefore, an outcome  $x \in \{0, 1, \dots, K-1\}$  is represented as a vector

$$\mathbf{1}_x \stackrel{ ext{def}}{=} ig([i=x]ig)_{i=0}^{K-1} = ig(\underbrace{0,\ldots,0}_x,1,\underbrace{0,\ldots,0}_{K-x-1}ig).$$

The outcome probability, mean, and variance are very similar to the Bernoulli distribution.

$$egin{aligned} P(oldsymbol{x}) &= \prod_{i=0}^{K-1} p_i^{x_i} \ \mathbb{E}[x_i] &= p_i \ \mathrm{Var}(x_i) &= p_i(1-p_i) \end{aligned}$$

NPFL114, Lecture 1

Organization TL;DR

Notation Random Variables

Information Theory

Machine Learning

# **Self Information**

Amount of **surprise** when a random variable is sampled.

- Should be zero for events with probability 1.
- Less likely events are more surprising.
- Independent events should have **additive** information.

$$I(x) \stackrel{ ext{\tiny def}}{=} -\log P(x) = \log rac{1}{P(x)}$$

NPFL114, Lecture 1

tion Theory Machine Learning

NNs '80s



# Entropy

Amount of **surprise** in the whole distribution.

$$H(P) \stackrel{ ext{def}}{=} \mathbb{E}_{\mathrm{x} \sim P}[I(x)] = -\mathbb{E}_{\mathrm{x} \sim P}[\log P(x)]$$

- for discrete  $P: H(P) = -\sum_{x} P(x) \log P(x)$
- for continuous  $P: H(P) = -\int P(x) \log P(x) dx$

Because  $\lim_{x\to 0} x \log x = 0$ , for P(x) = 0 we consider  $P(x) \log P(x)$  to be zero.

Note that in the continuous case, the continuous entropy (also called *differential entropy*) has slightly different semantics, for example, it can be negative.

For binary logarithms, the entropy is measured in **bits**. х However, from now on, all logarithms are *natural logarithms* with base e (and then the entropy is measured in units called **nats**).

NPFL114, Lecture 1

Organization TL;DR Notation Random Variables Information Theory

Machine Learning NNs '80s







# **Cross-Entropy**

$$H(P,Q) \stackrel{ ext{def}}{=} - \mathbb{E}_{\mathrm{x} \sim P}[\log Q(x)]$$

Gibbs inequality states that

- $H(P,Q) \ge H(P)$
- $H(P) = H(P,Q) \Leftrightarrow P = Q$
- Proof: Using the fact that  $\log x \leq (x-1)$  with equality only for x=1, we get

$$\sum_x P(x)\lograc{Q(x)}{P(x)}\leq \sum_x P(x)\left(rac{Q(x)}{P(x)}-1
ight)=\sum_x Q(x)-\sum_x P(x)=0.$$

• Corollary: For a categorical distribution with n outcomes,  $H(P) \le \log n$ , because for Q(x) = 1/n we get  $H(P) \le H(P,Q) = -\sum_x P(x) \log Q(x) = \log n$ .

Note that generally H(P,Q) 
eq H(Q,P).

NPFL114, Lecture 1

Organization TL;DR Notation

eory Machine Learning



# Kullback-Leibler Divergence (KL Divergence)

Sometimes also called **relative entropy**.

$$D_{ ext{KL}}(P\|Q) \stackrel{ ext{def}}{=} H(P,Q) - H(P) = \mathbb{E}_{ ext{x} \sim P}[\log P(x) - \log Q(x)]$$

- consequence of Gibbs inequality:  $D_{ ext{KL}}(P\|Q) \geq 0$ ,  $D_{ ext{KL}}(P\|Q) = 0$  iff P = Q
- generally  $D_{ ext{KL}}(P\|Q) 
  eq D_{ ext{KL}}(Q\|P)$

ory Machine Learning

### **Nonsymmetry of KL Divergence**





#### Ú F<sub>A</sub>L

# Normal (or Gaussian) Distribution

Distribution over real numbers, parametrized by a mean  $\mu$  and variance  $\sigma^2$ :

$$\mathcal{N}(x;\mu,\sigma^2) = \sqrt{rac{1}{2\pi\sigma^2}} \exp\left(-rac{(x-\mu)^2}{2\sigma^2}
ight)\,.$$

For standard values  $\mu=0$  and  $\sigma^2=1$  we get  $\mathcal{N}(x;0,1)=\sqrt{rac{1}{2\pi}}e^{-rac{x^2}{2}}$  .



NPFL114, Lecture 1

Organization

TL;DR

Machine Learning

### Why Normal Distribution

#### Ú F<sub>A</sub>L

# **Central Limit Theorem**

The sum of independent identically distributed random variables with finite variance converges to normal distribution.

# **Principle of Maximum Entropy**

Given a set of constraints, a distribution with maximal entropy fulfilling the constraints can be considered the most general one, containing as little additional assumptions as possible.

Considering distributions on all real numbers with a given mean and variance, it can be proven (using variational inference) that such a distribution with **maximum entropy** is exactly the normal distribution.

# **Machine Learning**



A possible definition of learning from Mitchell (1997):

A computer program is said to learn from experience E 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 E.

• Task T

 $^{\circ}$  *classification*: assigning one of k categories to a given input

- $\,\circ\,$  regression: producing a number  $x\in\mathbb{R}$  for a given input
- $^{\circ}$  structured prediction, denoising, density estimation, ...
- Measure P
  - accuracy, error rate, F-score, ...
- Experience E
  - *supervised*: usually a dataset with desired outcomes (*labels* or *targets*)
  - *unsupervised*: usually data without any annotation (raw text, raw images, ...)
  - reinforcement learning, semi-supervised learning, ...

NPFL114, Lecture 1

on Theory Machine Learning

NNs '80s

#### Well-known Datasets



Name	Description	Instances
<u>MNIST</u>	Images (28x28, grayscale) of handwritten digits.	60k
CIFAR-10	Images (32x32, color) of 10 classes of objects.	50k
<u>CIFAR-</u> <u>100</u>	Images (32x32, color) of 100 classes of objects (with 20 defined superclasses).	50k
<u>ImageNet</u>	Labeled object image database (labeled objects, some with bounding boxes).	14.2M
<u>ImageNet-</u> ILSVRC	Subset of ImageNet for Large Scale Visual Recognition Challenge, annotated with 1000 object classes and their bounding boxes.	1.2M
<u>COCO</u>	<i>Common Objects in Context</i> : Complex everyday scenes with descriptions (5) and highlighting of objects (91 types).	2.5M

NPFL114, Lecture 1

Organization TL;DR

Random Variables Notation

Information Theory

NNs '80s

Machine Learning

Well-known Datasets



#### ImageNet-ILSVRC



Figure 4 of "ImageNet Classification with Deep Convolutional Neural Networks" by Alex Krizhevský et al.



https://image-net.org/challenges/LSVRC/2014/

NPFL114, Lecture 1

TL;DR Organization

Notation

Random Variables

Information Theory

Machine Learning

NNs '80s



#### **COCO**



https://cocodataset.org/#detection-2020

NPFL114, Lecture 1

Organization TL;DR

Notation Random

Random Variables

Information Theory

Machine Learning

NNs '80s

#### Well-known Datasets



Name	Description	Instances
IAM-OnDB	Pen tip movements of handwritten English from 221 writers.	86k words
<u>TIMIT</u>	Recordings of 630 speakers of 8 dialects of American English.	6.3k sents
<u>CommonVoice</u>	1.6M Eng recordings from 86k people, $\sim$ 2400 hours of speech.	1.6M
<u>PTB</u>	<i>Penn Treebank</i> : 2500 stories from Wall Street Journal, with POS tags and parsed into trees.	1M words
<u>PDT</u>	<i>Prague Dependency Treebank</i> : Czech sentences annotated on 4 layers (word, morphological, analytical, tectogrammatical).	1.9M words
<u>UD</u>	<i>Universal Dependencies</i> : Treebanks of 138 languages with consistent annotation of lemmas, POS tags, morphology, syntax.	243 treebanks
<u>WMT</u>	Aligned parallel sentences for machine translation.	gigawords

NPFL114, Lecture 1

Organization

TL;DR Notation Ra

Random Variables Information

Information Theory

NNs '80s

Machine Learning



Ú F∡l

Ú<sub>F</sub>≩L

In summer 2017, a paper came out describing automatic generation of neural architectures using reinforcement learning.



NPFL114, Lecture 1

Organization TL;DR

Notation Random

Random Variables Inform

Information Theory Ma

Machine Learning NNs '80s

Currently, one of the best architectures is EfficientNet, which combines automatic architecture discovery, multidimensional scaling and elaborate dataset augmentation methods.





EfficientNet was further improved by EfficientNetV2 two years later.



*Figure 5.* Model Size, FLOPs, and Inference Latency – Latency is measured with batch size 16 on V100 GPU. 21k denotes pretrained on ImageNet21k images, others are just trained on ImageNet ILSVRC2012. Our EfficientNetV2 has slightly better parameter efficiency with EfficientNet, but runs 3x faster for inference.

Figure 5 of "EfficientNetV2: Smaller Models and Faster Training", https://arxiv.org/abs/2104.00298

NPFL114, Lecture 1

Organization TL;DR

Notation Random Variables

Information Theory

n Theory Machine Learning

NNs '80s

# **Machine Translation Improvements**

To illustrate deep neural networks improvements in other domains, consider the English $\rightarrow$ Czech results of the international Workshop on Machine Translation. Both the automatic BLEU metric and manual evaluation are presented.



- TectoMT parses the input, transfers to the other language, generates the sentence;
- RBMT is the PC-Translator software;
- SMT is statistical machine translation using the Moses system;
- Online is an online translation system (Google in 2009, Online-B since 2010);
- **NMT** is the neural machine translation using deep neural networks.

NPFL114, Lecture 1

Notation Random Variables

Information Theory

Machine Learning NNs '80s

### **Introduction to Deep Learning History**



#### Modified from https://www.slideshare.net/deview/251-implementing-deep-learning-using-cu-dnn/4

NNs '80s

NPFL114, Lecture 1

TL:DR Organization

Notation

Random Variables

Information Theory

Machine Learning

### **Perceptron – Extra Simple Neural Network**

- Assume we have an input node for every input feature.
- Additionally, we have an output node for every model output.
- Every input node and output node are connected with a directed edge, and every edge has an associated weight.
- Value of every (output) node is computed by summing the values of predecessors multiplied by the corresponding weights, added to a bias of this node, and finally passed through an activation function a:

$$y=a\left(\sum\nolimits_{j}x_{j}w_{j}+b
ight)$$

or in vector form  $y = a(\boldsymbol{x}^T \boldsymbol{w} + b)$ , or for a batch of examples  $oldsymbol{X}$  ,  $oldsymbol{y}=a(oldsymbol{X}oldsymbol{w}+b)$  .

Output layer activation a



Input layer

NPFL114, Lecture 1

Organization TL:DR

Notation Random Variables Information Theory

Machine Learning

NNs '80s

#### **Perceptron – Linearly Separable and Nonseparable Data**







NPFL114, Lecture 1

Organization TL;DR

Notation Random

Random Variables Inform

Information Theory Machine Learning

NNs '80s

https://miro.medium.com/v2/1\*JVZ4FXVRIr1oN-4ffq\_kNQ.png

#### Neural Network Architecture à la '80s





NPFL114, Lecture 1

Organization TL;DR

Notation Random Variables

Information Theory

Theory Machine Learning

NNs '80s

#### **Neural Network Architecture**





## **Neural Network Activation Functions**

# **Output Layers**

- none (linear regression if there are no hidden layers)
- $\sigma$  (sigmoid; logistic regression if there are no hidden layers)

$$\sigma(x) \stackrel{ ext{\tiny def}}{=} rac{1}{1+e^{-x}}$$

is used to model a probability p of a binary event; its input is called a **logit**,  $\log \frac{p}{1-p}$ ;

softmax (maximum entropy model if there are no hidden layers)

$$ext{softmax}(oldsymbol{x}) \propto e^{oldsymbol{x}} \ ext{softmax}(oldsymbol{x})_i \stackrel{ ext{def}}{=} rac{e^{x_i}}{\sum_j e^{x_j}}$$

is used to model probability distribution  $m{p}$ ; its input is called a logit,  $\log(m{p}) + c$ .

NPFL114, Lecture 1

Organization TL;DR

Random Variables Notation

Information Theory

NNs '80s

Machine Learning

# **Neural Network Activation Functions**

# **Hidden Layers**

- none: does not help, composition of linear/affine mapping is a linear/affine mapping
- $\sigma$ : does not work great nonsymmetrical, repeated application converges to the fixed point  $x=\sigma(x)pprox 0.659$ , and  $rac{d\sigma}{dx}(0)=1/4$
- tanh
  - $^{\circ}~$  result of making  $\sigma$  symmetrical and making the derivative in zero 1
  - $\circ \ anh(x) = 2\sigma(2x) 1$



• ReLU:  $\max(0, x)$ 

NPFL114, Lecture 1

ion Theory Machine Learning

ig NNs '80s



# **Universal Approximation Theorem '89**

Let  $\varphi(x): \mathbb{R} \to \mathbb{R}$  be a nonconstant, bounded and nondecreasing continuous function. (Later a proof was given also for  $\varphi = \operatorname{ReLU}$  and even for any nonpolynomial function.) For any  $\varepsilon > 0$  and any continuous function  $f: [0,1]^D \to \mathbb{R}$ , there exists  $H \in \mathbb{N}$ ,  $\boldsymbol{v} \in \mathbb{R}^H$ ,  $\boldsymbol{b} \in \mathbb{R}^H$  and  $\boldsymbol{W} \in \mathbb{R}^{D \times H}$ , such that if we denote

$$F(oldsymbol{x}) = oldsymbol{v}^T arphi(oldsymbol{x}^T oldsymbol{W} + oldsymbol{b}) = \sum_{i=1}^H v_i arphi(oldsymbol{x}^T oldsymbol{W}_{*,i} + b_i),$$

where arphi is applied element-wise, then for all  $oldsymbol{x} \in [0,1]^D$  :

$$|F(oldsymbol{x}) - f(oldsymbol{x})| < arepsilon.$$

#### **One Possible Interpretation**

It is always possible to create features using just a single linear layer followed by a nonlinearity, such that the resulting dataset is always linearly separable.

NPFL114, Lecture 1

Organization TL;DR

Notation Random Variables

Information Theory

Machine Learning

NNs '80s



# **Universal Approximation Theorem for ReLUs**

Ú<sub>F</sub>≩L

Sketch of the proof:

• If a function is continuous on a closed interval, it can be approximated by a sequence of lines to arbitrary precision.



https://miro.medium.com/max/844/1\*lihbPNQgl7oKjpCsmzPDKw.png

 However, we can create a sequence of k linear segments as a sum of k ReLU units – on every endpoint a new ReLU starts (i.e., the input ReLU value is zero at the endpoint), with a tangent which is the difference between the target tangent and the tangent of the approximation until this point.

NPFL114, Lecture 1

Organization TL;DR Notation

ion Random Variables

Information Theory

ion Theory Machine Learning

NNs '80s

### **Evolving ReLU Approximation**





NPFL114, Lecture 1

NNs '80s

Machine Learning

# **Universal Approximation Theorem for Squashes**



Sketch of the proof for a squashing function  $\varphi(x)$  (i.e., nonconstant, bounded and nondecreasing continuous function like sigmoid):

• We can prove  $\varphi$  can be arbitrarily close to a hard threshold by compressing it horizontally.



• Then we approximate the original function using a series of straight line segments



 $https://hackernoon.com/hn-images/1*hVuJgUTLUFWTMmJhl\_fomg.png$ 

NPFL114, Lecture 1

Organization TL;DR

Notation Random Variables

s Information Theory

heory Machine Learning

# How Good is Current Deep Learning

- DL has seen amazing progress in the last ten years.
- Is it enough to get a bigger brain (datasets, models, computer power)?
- Problems compared to Human learning:
   Sample efficiency
  - Human-provided labels
  - Robustness to data distribution change
  - $\circ$  Stupid errors



https://intl.startrek.com/sites/default/files/styles/content\_full/public/images/2019-07/c8ffe9a587b126f152ed3d89a146b445.jpg

NPFL114, Lecture 1

Organization TL;DR

Notation Rand

Random Variables

Information Theory

Machine Learning NNs '80s

# How Good is Current Deep Learning



...



it may be that today's large neural networks are slightly conscious

Přeložit Tweet

12:27 dop. · 10. 2. 2022 · Twitter Web App https://twitter.com/ilyasut/status/1491554478243258368

...



Odpověď uživateli @ilyasut

#### Nope.

Not even for true for small values of "slightly conscious" and large values of "large neural nets". I think you would need a particular kind of macroarchitecture that none of the current networks possess.

#### Přeložit Tweet

10:02 odp. · 12. 2. 2022 · Twitter for Android https://twitter.com/ylecun/status/1492604977260412928



Murray Shanahan @mpshanahan

#### Odpověď uživateli @ilyasut

# ... in the same sense that it may be that a large field of wheat is slightly pasta

...

#### Přeložit Tweet

11:08 dop. · 10. 2. 2022 · Twitter Web App

https://twitter.com/mpshanahan/status/1491715721289678848

NPFL114, Lecture 1

Organization TL;DR

Notation Random Variables

Information Theory

on Theory Machine Learning

NNs '80s

# How Good is Current Deep Learning

- Thinking fast and slow
  - System 1
    - intuitive
    - fast
    - automatic
    - frequent
    - unconscious

Current DL

- $^{\circ}$  System 2
  - Iogical
  - slow
  - effortful
  - infrequent
  - conscious

Future DL



https://en.wikipedia.org/wiki/File:Thinking,\_Fast\_and\_Slow.jpg

Machine Learning

NPFL114, Lecture 1

Notation Random Variables

bles Informat

Information Theory

NNs '80s

### **Curse of Dimensionality**





NPFL114, Lecture 1

Organization TL;DR

Notation Random

Random Variables Inform

Information Theory

Machine Learning

NNs '80s

#### **Machine and Representation Learning**





Figure 1.5 of "Deep Learning" book, https://www.deeplearningbook.org

NPFL114, Lecture 1

Organization TL;DR

Notation Random Variables

bles Information Theory

Theory Machine Learning

NNs '80s