

# Introduction to Deep Learning

Milan Straka

 March 01, 2021

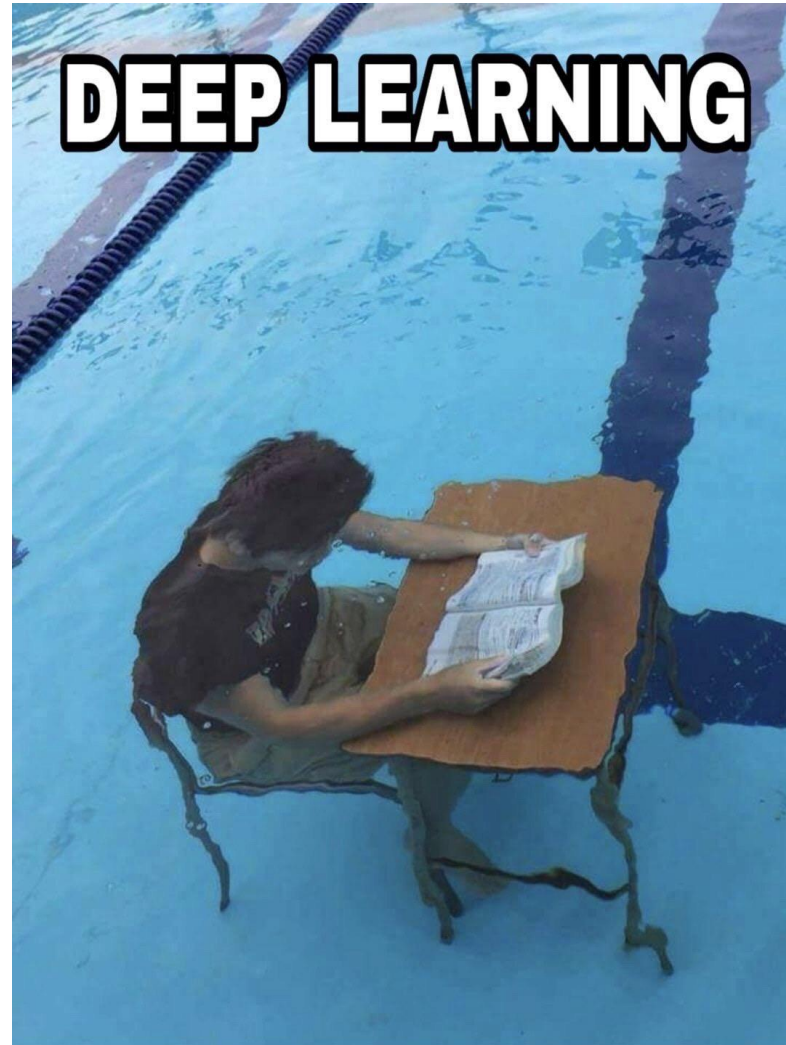


EUROPEAN UNION  
European Structural and Investment Fund  
Operational Programme Research,  
Development and Education

Charles University in Prague  
Faculty of Mathematics and Physics  
Institute of Formal and Applied Linguistics



unless otherwise stated



<https://i.redd.it/t87gswsbmq41.jpg>

- Image recognition
- Object detection
- Image segmentation,
- Human pose estimation
- Image labeling
- Visual question answering
- Speech recognition and generation
- Lip reading
- Machine translation
- Machine translation without parallel data
- Chess, Go and Shogi
- Multiplayer Capture the flag

**Course Website** <https://ufal.mff.cuni.cz/courses/npfl114>

- Recordings of lectures and practicals, slides, assignments, exam questions.

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

- Templates for the assignments, slide sources.

## Piazza

- Piazza will be used as a communication platform.

You can post questions or notes,

- privately to the instructors, or
- to everyone (signed or anonymously).

Students can answer other student's questions too, which allows you to get faster response. Please do not send complete solutions to other students, only excerpts of the source files.

- Please use Piazza for **all communication** with the instructors.
- You will get the invite link after the first lecture.

<https://recodex.mff.cuni.cz>

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

## Practicals

- There will be 2-3 assignments a week, each with 2-week deadline.
  - Deadlines can be extended, but you need to write **before** the deadline.
- 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 obtain 50 bonus points.

## 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; 75 points for PhD students.

# Notation

- $a$ ,  $\mathbf{a}$ ,  $\mathbf{A}$ ,  $\mathbf{A}$ : scalar (integer or real), vector, matrix, tensor
  - all vectors are always **column** vectors
  - transposition changes a column vector into a row vector, so  $\mathbf{a}^T$  is a row vector
  - we denote **scalar product** between vectors  $\mathbf{a}$  and  $\mathbf{b}$  as  $\mathbf{a}^T \mathbf{b}$ 
    - we understand it as matrix multiplication
- $a$ ,  $\mathbf{a}$ ,  $\mathbf{A}$ : 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  $x$
- $\nabla_{\mathbf{x}} f(\mathbf{x})$ : gradient of  $f$  with respect to  $\mathbf{x}$ , i.e.,  $\left( \frac{\partial f(\mathbf{x})}{\partial x_1}, \frac{\partial f(\mathbf{x})}{\partial x_2}, \dots, \frac{\partial f(\mathbf{x})}{\partial x_n} \right)$

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

## Probability Distribution

A probability distribution describes how likely are individual values 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  $\mathbf{x}$  takes a value  $x$  is denoted as  $P(x)$  or explicitly as  $P(\mathbf{x} = x)$ . All probabilities are non-negative and sum of probabilities of all possible values of  $\mathbf{x}$  is  $\sum_x P(\mathbf{x} = x) = 1$ .

For continuous variables, the probability that the value of  $\mathbf{x}$  lies in the interval  $[a, b]$  is given by  $\int_a^b p(x) dx$ .



## Expectation

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

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

For continuous variables it is computed as:

$$\mathbb{E}_{x \sim p}[f(x)] \stackrel{\text{def}}{=} \int_x p(x) f(x) dx$$

If the random variable is obvious from context, we can write only  $\mathbb{E}_P[x]$  or even  $\mathbb{E}[x]$ .

Expectation is linear, i.e.,

$$\mathbb{E}_x[\alpha f(x) + \beta g(x)] = \alpha \mathbb{E}_x[f(x)] + \beta \mathbb{E}_x[g(x)]$$

## Variance

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

$$\text{Var}(x) \stackrel{\text{def}}{=} \mathbb{E} \left[ (x - \mathbb{E}[x])^2 \right], \text{ or more generally}$$

$$\text{Var}(f(x)) \stackrel{\text{def}}{=} \mathbb{E} \left[ (f(x) - \mathbb{E}[f(x)])^2 \right]$$

It is easy to see that

$$\text{Var}(x) = \mathbb{E} \left[ x^2 - 2x\mathbb{E}[x] + (\mathbb{E}[x])^2 \right] = \mathbb{E} [x^2] - (\mathbb{E}[x])^2.$$

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

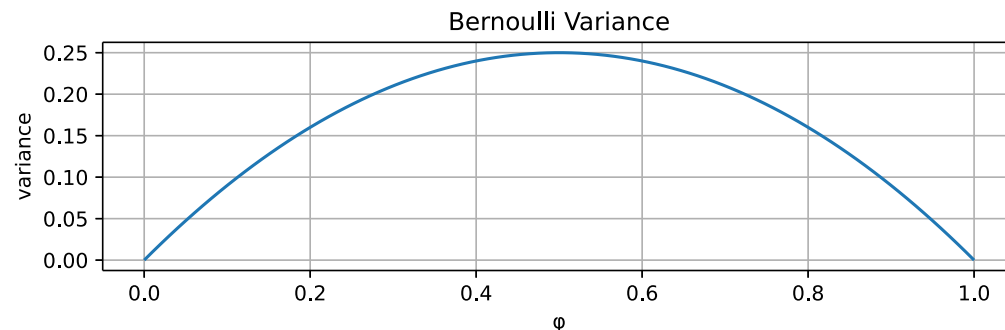
## Bernoulli Distribution

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

$$P(x) = \varphi^x (1 - \varphi)^{1-x}$$

$$\mathbb{E}[x] = \varphi$$

$$\text{Var}(x) = \varphi(1 - \varphi)$$



## Categorical Distribution

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

$$P(\mathbf{x}) = \prod_i^k \mathbf{p}_i^{\mathbf{x}_i}$$

$$\mathbb{E}[\mathbf{x}_i] = \mathbf{p}_i, \text{Var}(\mathbf{x}_i) = \mathbf{p}_i(1 - \mathbf{p}_i)$$

## 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.

$$I(x) \stackrel{\text{def}}{=} -\log P(x) = \log \frac{1}{P(x)}$$

## Entropy

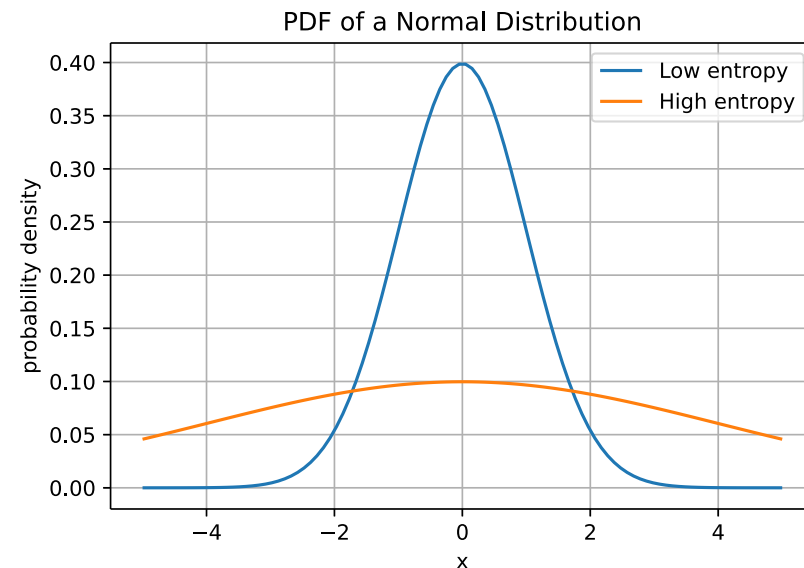
Amount of **surprise** in the whole distribution.

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

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

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

From now on, all logarithms are *natural logarithms* with base  $e$ .



## Cross-Entropy

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

- Gibbs inequality
  - $H(P, Q) \geq H(P)$
  - $H(P) = H(P, Q) \Leftrightarrow P = Q$
  - Proof: Using Jensen's inequality, we get

$$\sum_x P(x) \log \frac{Q(x)}{P(x)} \leq \log \sum_x P(x) \frac{Q(x)}{P(x)} = \log \sum_x Q(x) = 0.$$

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

## Kullback-Leibler Divergence (KL Divergence)

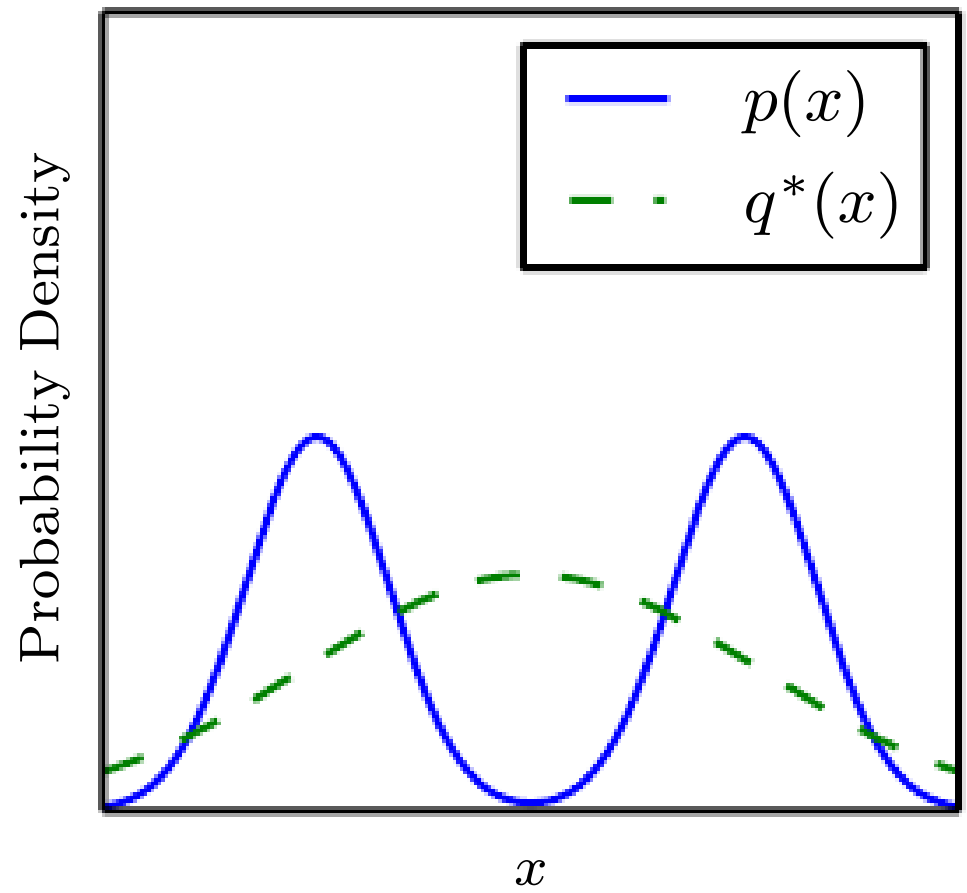
Sometimes also called **relative entropy**.

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

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

# Nonsymmetry of KL Divergence

$$q^* = \operatorname{argmin}_q D_{\text{KL}}(p||q)$$



$$q^* = \operatorname{argmin}_q D_{\text{KL}}(q||p)$$

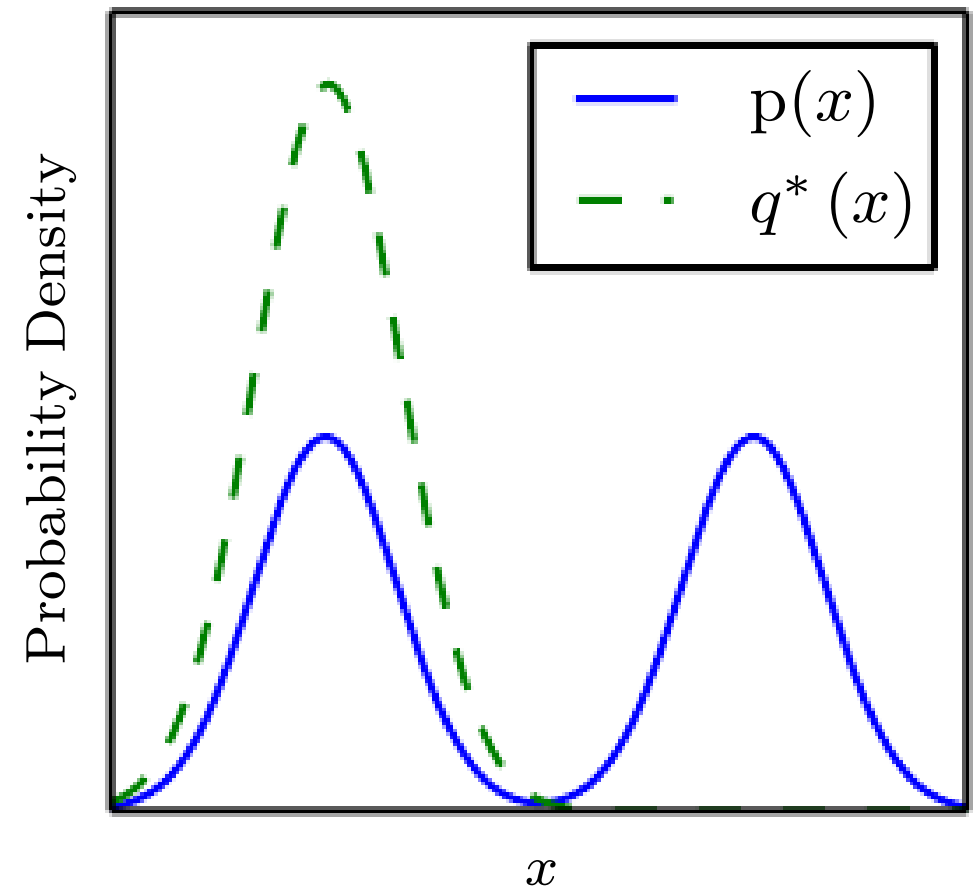


Figure 3.6, page 76 of Deep Learning Book, <http://deeplearningbook.org>



## Normal (or Gaussian) Distribution

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

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

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

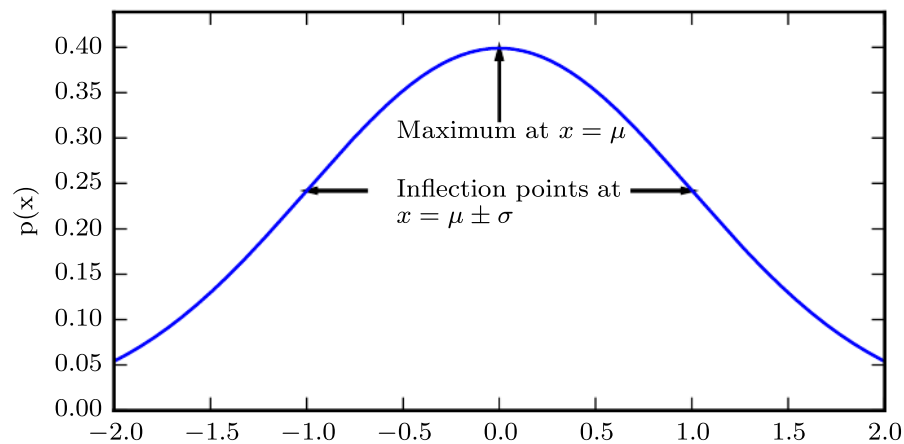


Figure 3.1, page 64 of Deep Learning Book, <http://deeplearningbook.org>.

## 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.

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  $T$ 
  - *classification*: assigning one of  $k$  categories to a given input
  - *regression*: producing a number  $x \in \mathbb{R}$  for a given input
  - *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, ...*

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

## ImageNet-ILSVRC

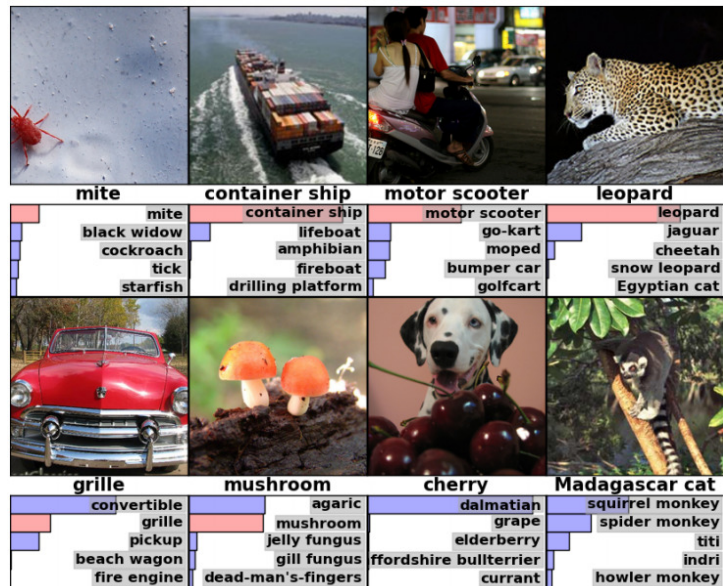


Image from "ImageNet Classification with Deep Convolutional Neural Networks" paper by Alex Krizhevsky et al.

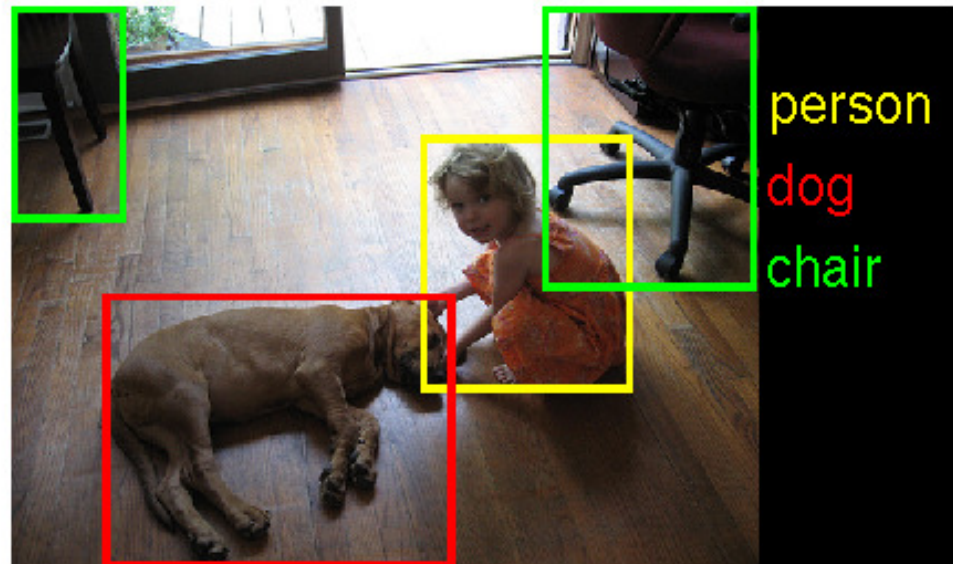


Image from <http://image-net.org/challenges/LSVRC/2014/>.

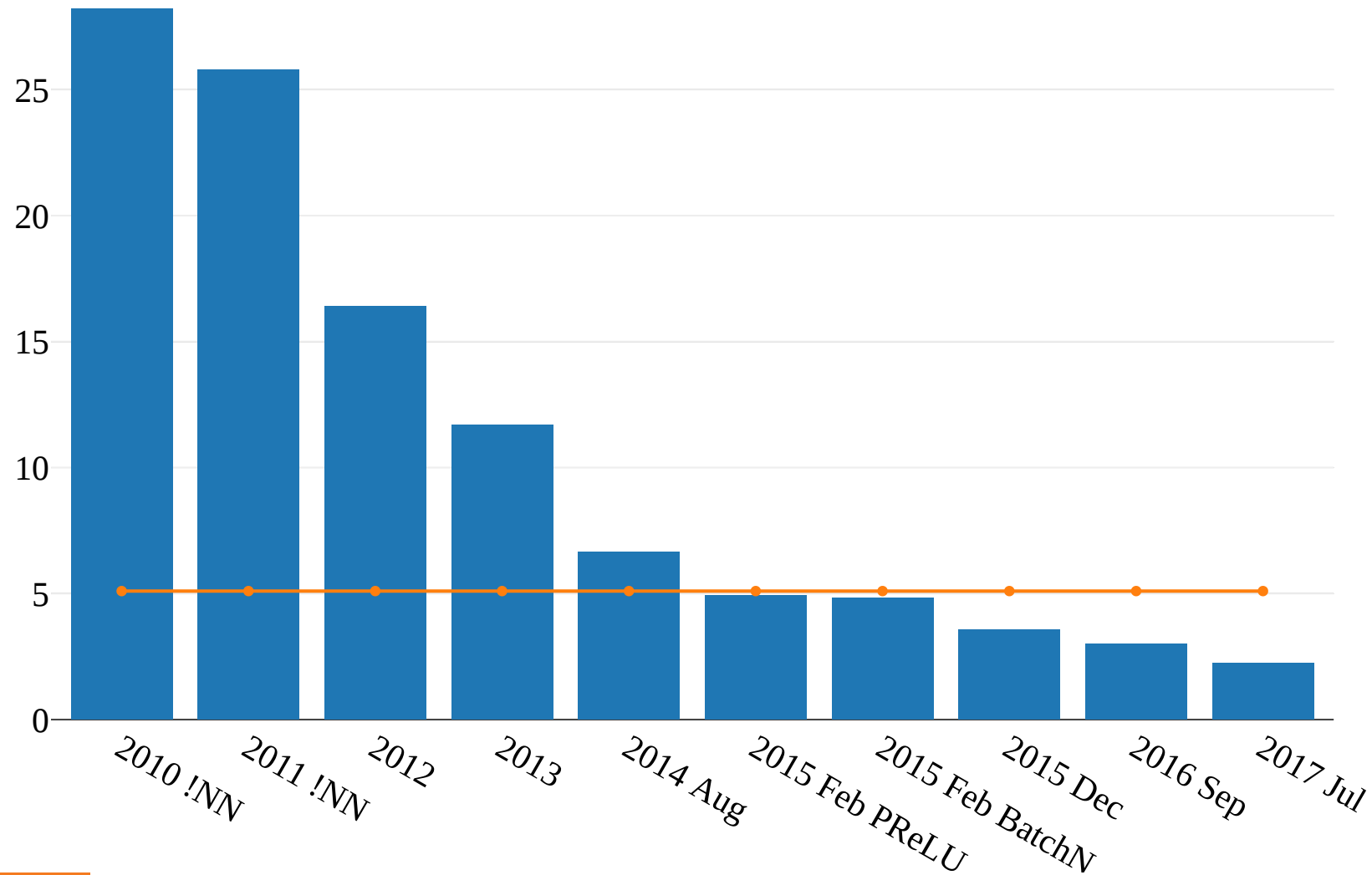
## COCO



Image from <http://mscoco.org/dataset/\#detections-challenge2016>.

Name	Description	Instances
<a href="#">IAM-OnDB</a>	Pen tip movements of handwritten English from 221 writers.	86k words
<a href="#">TIMIT</a>	Recordings of 630 speakers of 8 dialects of American English.	6.3k sents
<a href="#">CommonVoice</a>	400k recordings from 20k people, around 500 hours of speech.	400k
<a href="#">PTB</a>	<i>Penn Treebank</i> : 2500 stories from Wall Street Journal, with POS tags and parsed into trees.	1M words
<a href="#">PDT</a>	<i>Prague Dependency Treebank</i> : Czech sentences annotated on 4 layers (word, morphological, analytical, tectogrammatical).	1.9M words
<a href="#">UD</a>	<i>Universal Dependencies</i> : Treebanks of 104 languages with consistent annotation of lemmas, POS tags, morphology, syntax.	183 treebanks
<a href="#">WMT</a>	Aligned parallel sentences for machine translation.	gigawords

# ILSVRC Image Recognition Error Rates





# ILSVRC Image Recognition Error Rates

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

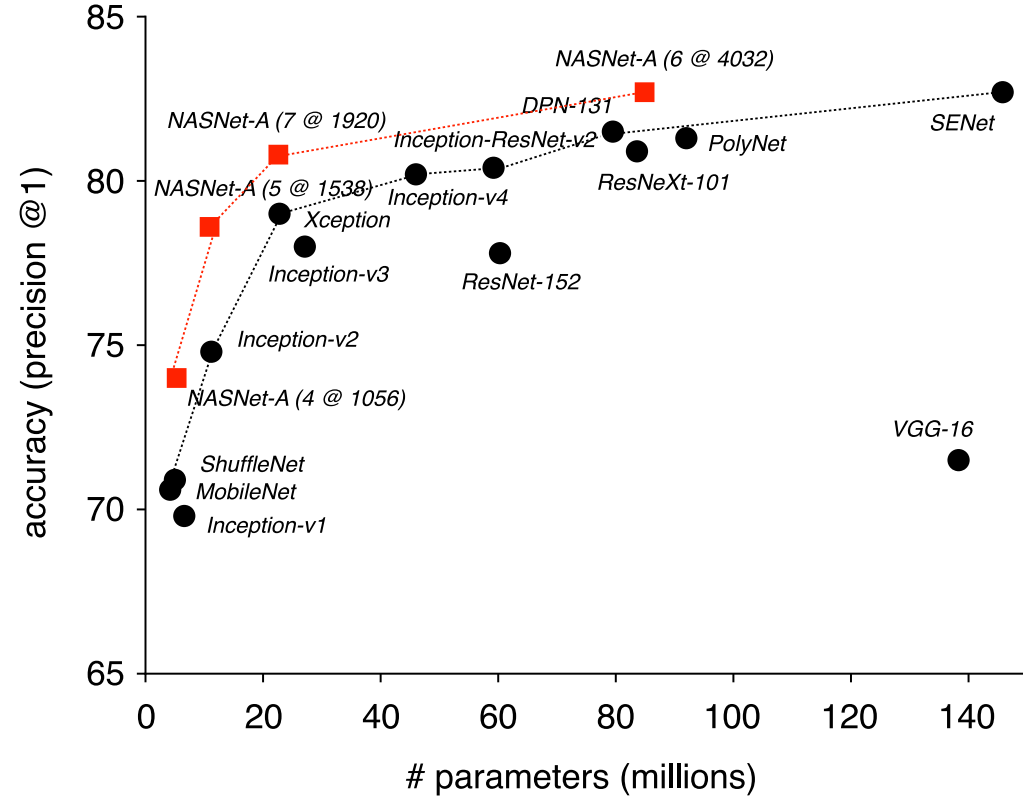
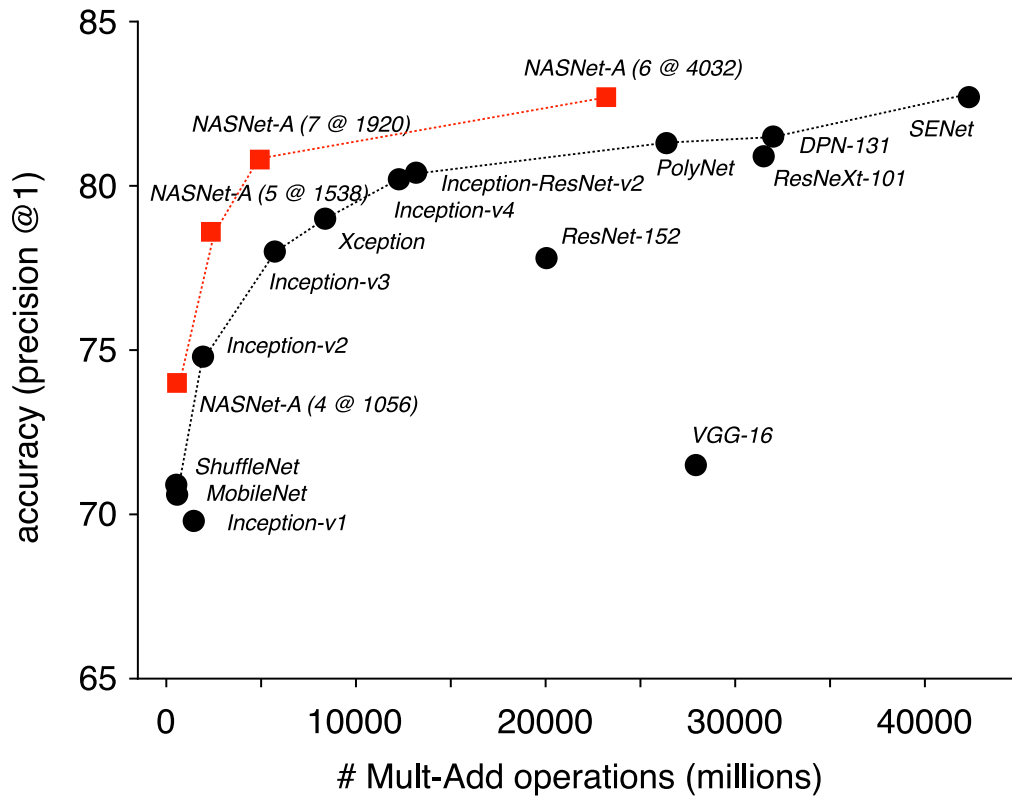


Figure 5 of paper "Learning Transferable Architectures for Scalable Image Recognition", <https://arxiv.org/abs/1707.07012>.

# ILSVRC Image Recognition Error Rates

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

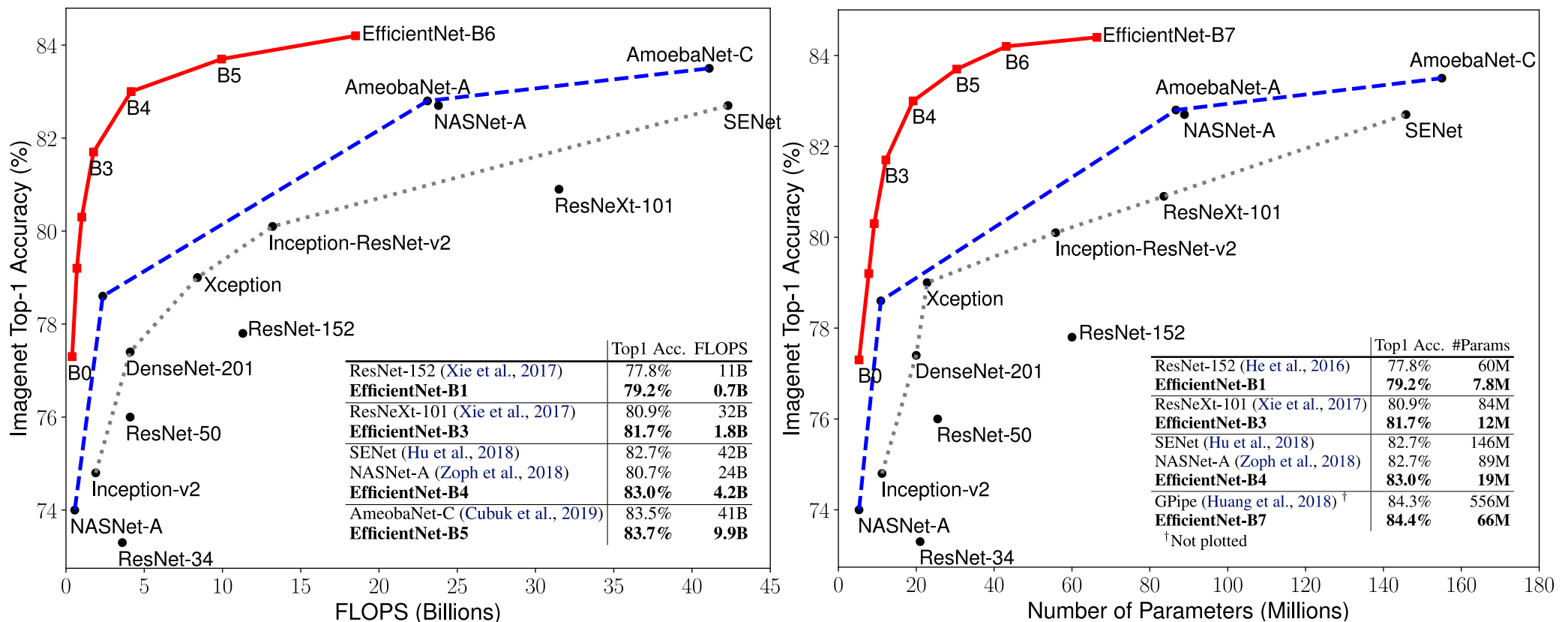
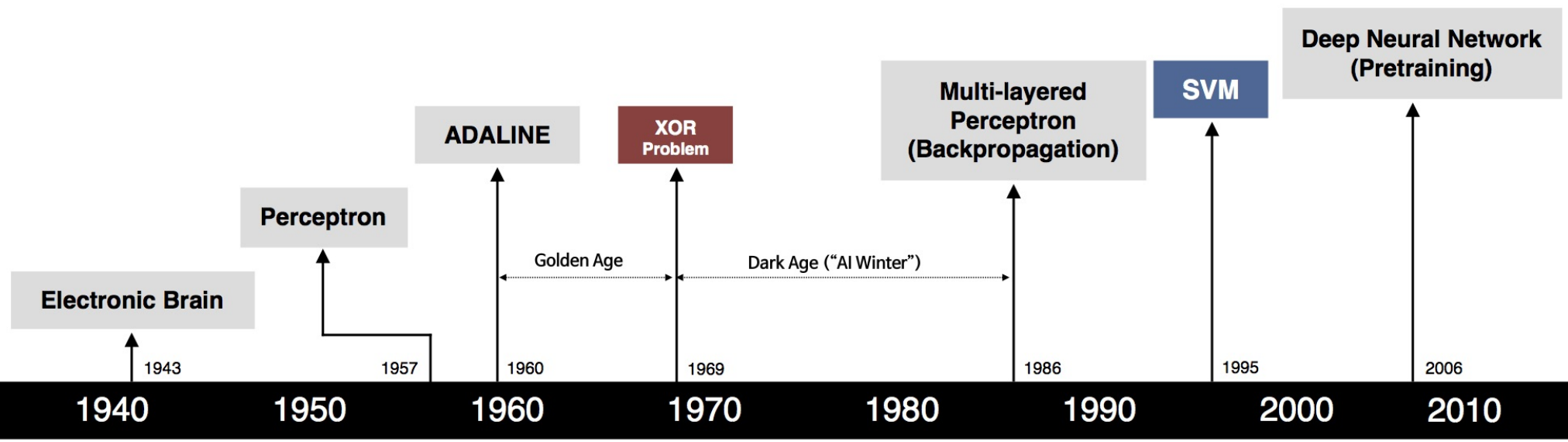
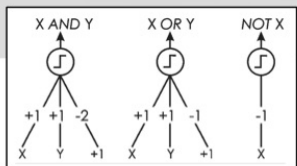


Figure 5 of paper "EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks", Figure 1 of paper "EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks", <https://arxiv.org/abs/1905.11946>.

# Introduction to Deep Learning History



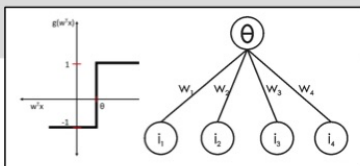
S. McCulloch - W. Pitts



- Adjustable Weights
- Weights are not Learned



F. Rosenblatt



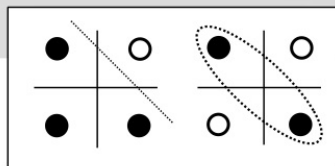
- Learnable Weights and Threshold



B. Widrow - M. Hoff



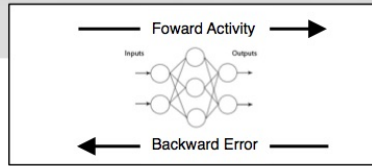
M. Minsky - S. Papert



- XOR Problem



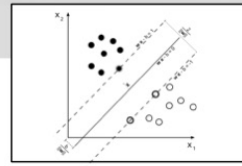
D. Rumelhart - G. Hinton - R. Williams



- Solution to nonlinearly separable problems
- Big computation, local optima and overfitting



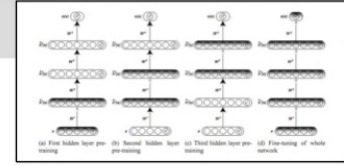
V. Vapnik - C. Cortes



- Limitations of learning prior knowledge
- Kernel function: Human Intervention



G. Hinton - S. Ruslan



- Hierarchical feature Learning

<https://www.slideshare.net/devview/251-implementing-deep-learning-using-cu-dnn/4>

# 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
  - Stupid errors



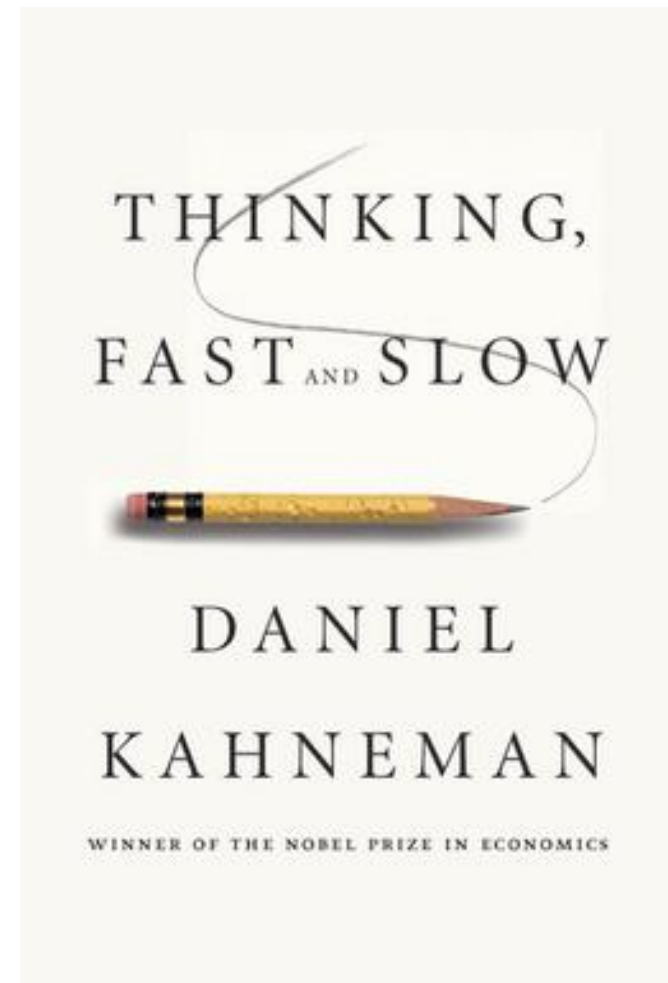
[https://intl.startrek.com/sites/default/files/styles/content\\_full/public/images/2019-07/c8ffe9a587b126f152ed3d89a146b445.jpg](https://intl.startrek.com/sites/default/files/styles/content_full/public/images/2019-07/c8ffe9a587b126f152ed3d89a146b445.jpg)

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

Current DL

- System 2
  - logical
  - slow
  - effortful
  - infrequent
  - conscious

Future DL



[https://en.wikipedia.org/wiki/File:Thinking,\\_Fast\\_and\\_Slow.jpg](https://en.wikipedia.org/wiki/File:Thinking,_Fast_and_Slow.jpg)

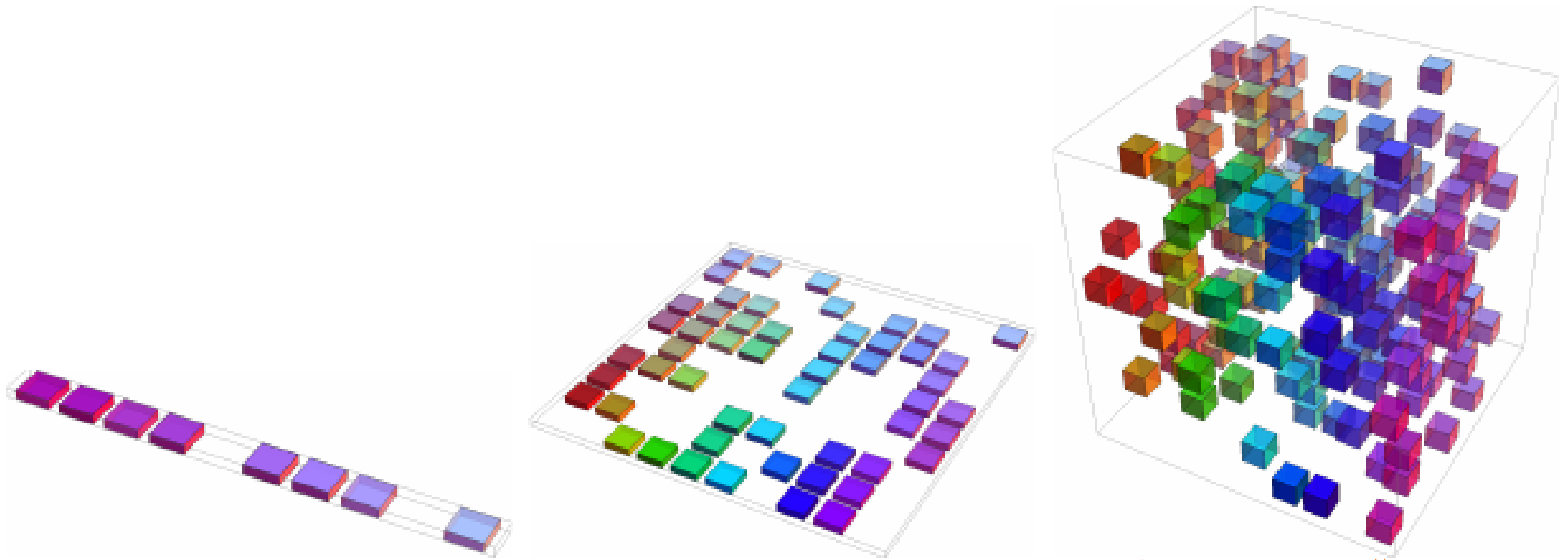


Figure 5.9, page 156 of Deep Learning Book, <http://deeplearningbook.org>.

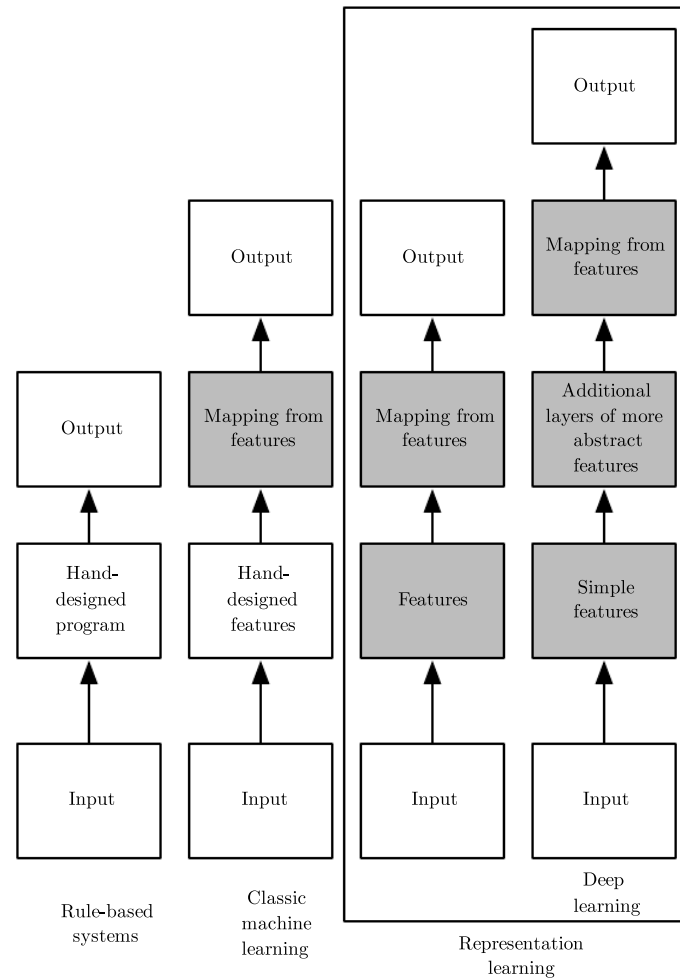
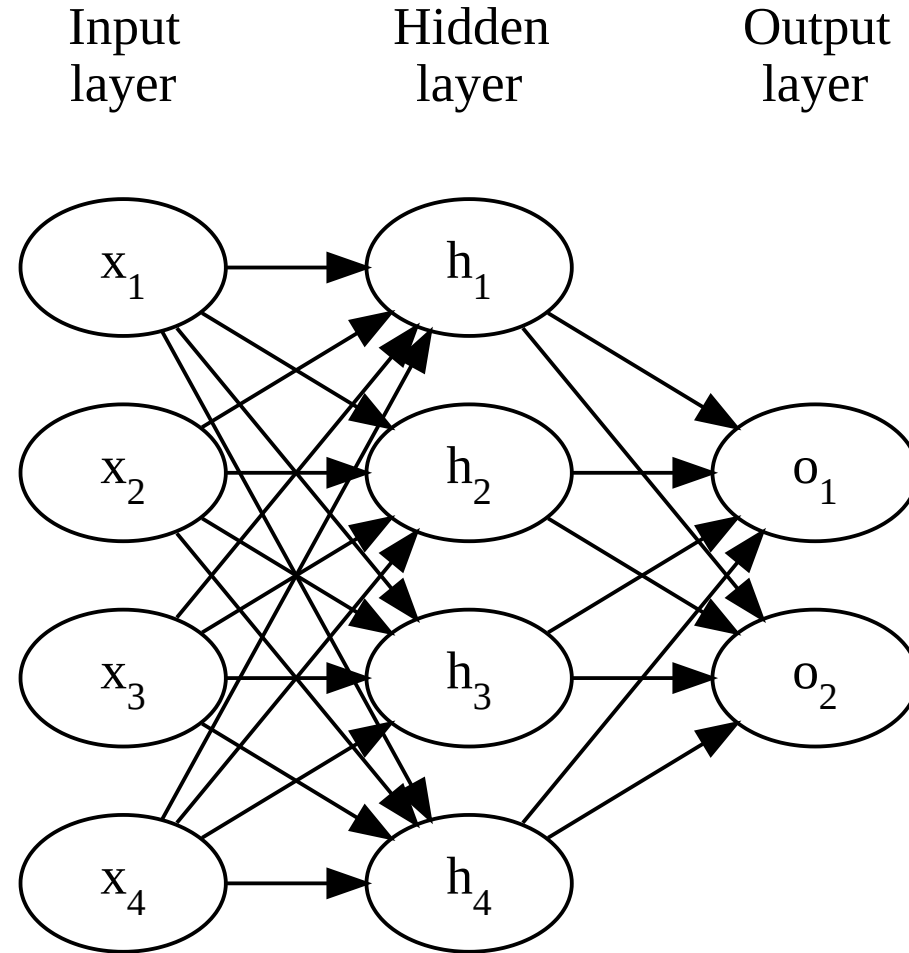


Figure 1.5, page 10 of Deep Learning Book, <http://deeplearningbook.org>.





There is a weight on each edge, and an activation function  $f$  is performed on the hidden layers, and optionally also on the output layer.

$$h_i = f \left( \sum_j w_{i,j} x_j + b_i \right)$$

If the network is composed of layers, we can use matrix notation and write

$$\mathbf{h} = f(\mathbf{W}\mathbf{x} + \mathbf{b}),$$

where  $\mathbf{W} \in \mathbb{R}^{|\text{hidden neurons}| \times |\text{input neurons}|}$  is a matrix of weights and  $\mathbf{b} \in \mathbb{R}^{|\text{hidden neurons}|}$  is a vector of biases.

## Output Layers

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

$$\sigma(x) \stackrel{\text{def}}{=} \frac{1}{1 + e^{-x}}$$

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

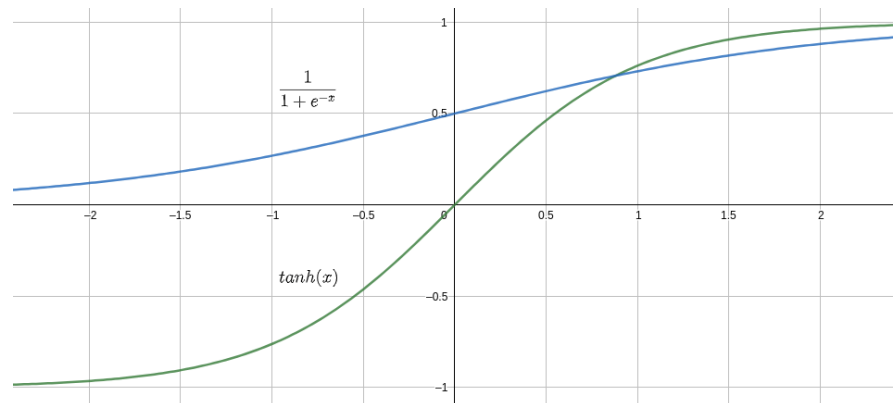
- softmax (maximum entropy model if there are no hidden layers)

$$\begin{aligned} \text{softmax}(\mathbf{x}) &\propto e^{\mathbf{x}} \\ \text{softmax}(\mathbf{x})_i &\stackrel{\text{def}}{=} \frac{e^{x_i}}{\sum_j e^{x_j}} \end{aligned}$$

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

## Hidden Layers

- none: does not help, composition of linear mapping is a linear mapping
- $\sigma$ : however, it works badly – nonsymmetrical, repeated application converges to the fixed point  $x = \sigma(x) \approx 0.659$ , and  $\frac{d\sigma}{dx}(0) = 1/4$
- tanh
  - result of making  $\sigma$  symmetrical and making the derivative in zero 1
  - $\tanh(x) = 2\sigma(2x) - 1$



- ReLU:  $\max(0, x)$

Let  $\varphi : \mathbb{R} \rightarrow \mathbb{R}$  be a nonconstant, bounded and nondecreasing continuous function.

(Later a proof was given also for  $\varphi = \text{ReLU}$  and even for any nonpolynomial function.)

For any  $\varepsilon > 0$  and any continuous function  $f : [0, 1]^D \rightarrow \mathbb{R}$ , there exists  $N \in \mathbb{N}$ ,  $\mathbf{v} \in \mathbb{R}^N$ ,  $\mathbf{b} \in \mathbb{R}^N$  and  $\mathbf{W} \in \mathbb{R}^{N \times D}$ , such that if we denote

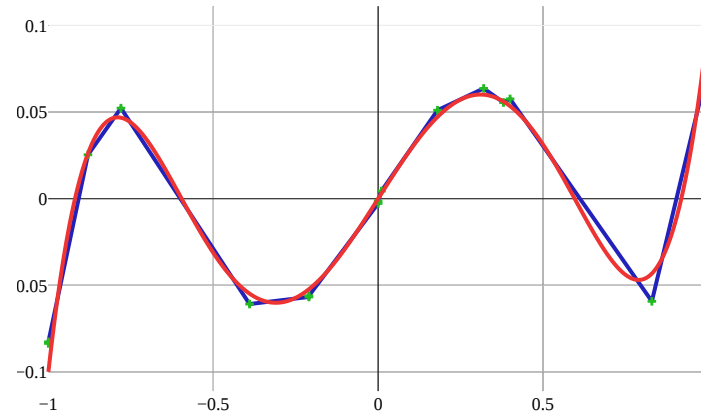
$$F(\mathbf{x}) = \mathbf{v}^T \varphi(\mathbf{W}\mathbf{x} + \mathbf{b}),$$

where  $\varphi$  is applied elementwise, then for all  $\mathbf{x} \in [0, 1]^D$ :

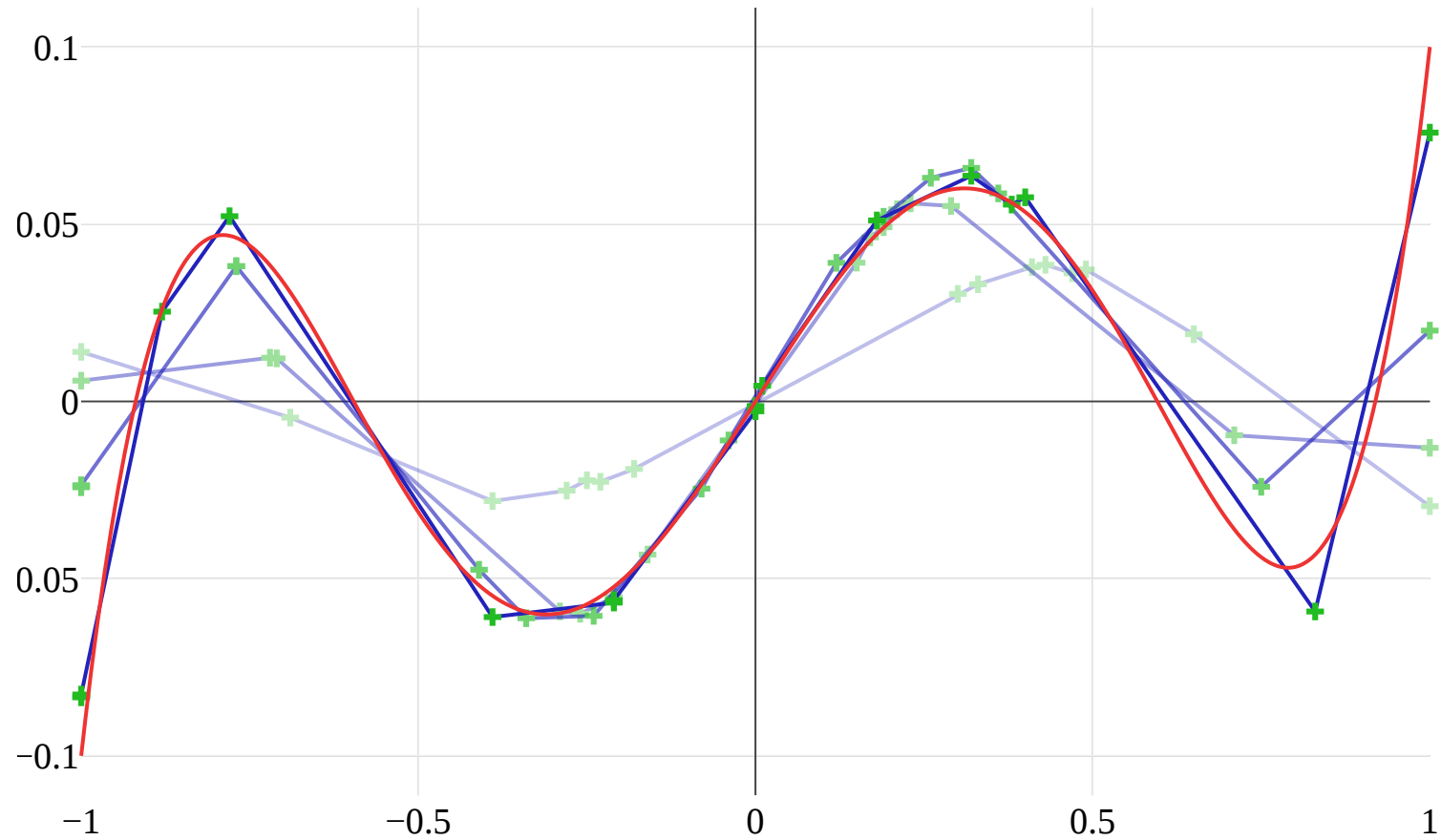
$$|F(\mathbf{x}) - f(\mathbf{x})| < \varepsilon.$$

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.



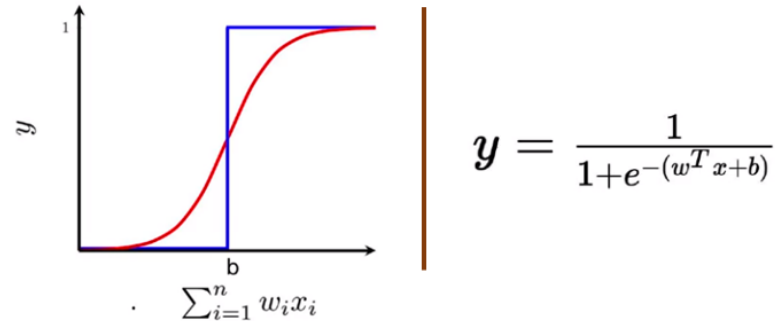
- However, we can create a sequence of  $k$  linear segments as a sum of  $k$  ReLU units – on every endpoint a new ReLU starts (i.e., the input ReLU value is zero at the endpoint), with a tangent which is the difference between the target target and the tangent of the approximation until this point.



# Universal Approximation Theorem for Squashes

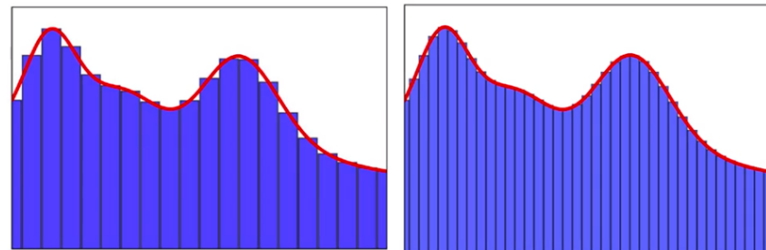
Sketch of the proof for a squashing function  $\varphi(x)$  (i.e., nonconstant, bounded and nondecreasing continuous function like sigmoid):

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



[https://hackernoon.com/hn-images/1\\*N7dfPwbiXC-Kk4TCbfRerA.png](https://hackernoon.com/hn-images/1*N7dfPwbiXC-Kk4TCbfRerA.png)

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



[https://hackernoon.com/hn-images/1\\*hVuJgUTLUFWTMmJhl\\_fomg.png](https://hackernoon.com/hn-images/1*hVuJgUTLUFWTMmJhl_fomg.png)