

**JOHANNES KEPLER
UNIVERSITY LINZ**

Deep Learning and Neural Networks I

Notation

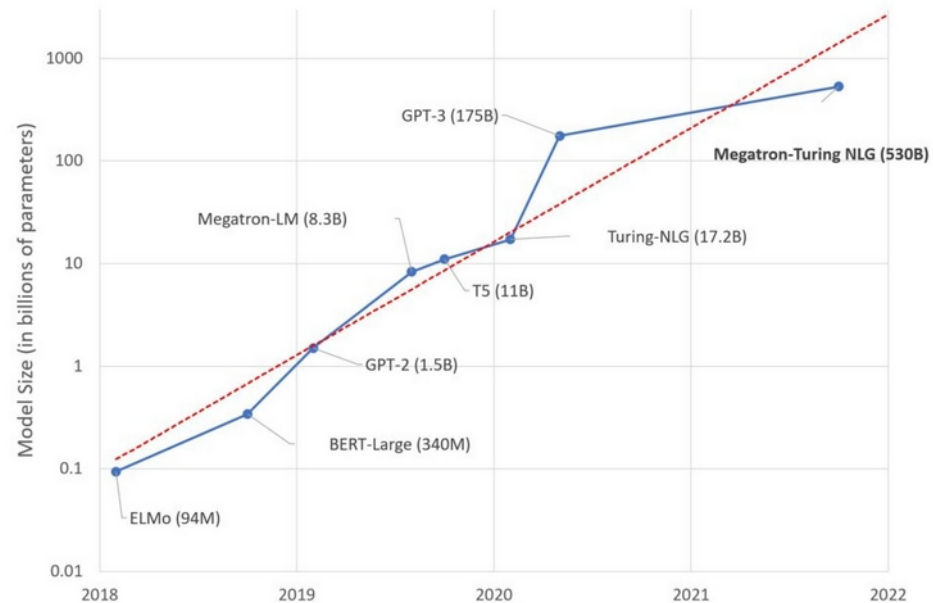


Günter Klambauer
LIT AI Lab & Institute for Machine Learning

This material, no matter whether in printed or electronic form, may be used for personal and non-commercial educational use only. Any reproduction of this material, no matter whether as a whole or in parts, no matter whether in printed or in electronic form, requires explicit prior acceptance of the authors.

Current topics

- Microsoft and NVIDIA have just announced a deep language model with 530B parameters!
 - Parallelized across thousands of GPUs
- Kunihiro Fukushima just won the 2021 Bower Award for Achievement in Science
 - ML community discusses:
[https://www.reddit.com/r/MachineLearning/comments/q76js4/schmidhuber_pays_tribute_to_kunihiro_fukushima/](https://www.reddit.com/r/MachineLearning/comments/q76js4/schmidhuber_paystribute_to_kunihiro_fukushima/)
- We are looking for a technician for a Deep Learning cluster
 - Please contact: Doris Kaiserrainer
kaiserreiner@ai-lab.jku.at if you are interested!
 - Email subject “Technician”
 - Details on Moodle



Source: <https://www.microsoft.com/en-us/research/blog/using-deepspeed-and-megatron-to-train-megatron-turing-nlg-530b-the-worlds-largest-and-most-powerful-generative-language-model/>

Notation example: Supervised learning

X			y
x^1	0.13	0.03	0.35
x^2	0.15	0.14	0.57
	0.79	0.19	0.87
	0.48	0.28	1.21
	0.93	0.43	0.48
x^n	0.43	0.41	0.70
	0.62	0.63	-0.44
	0.69	0.69	-1.05
	0.19	0.79	-1.29
	0.23	0.85	-1.11
x^N	0.11	0.99	0.32

- Our supervised data set is:
- Samples are rows in the data matrix \mathbf{X} :

$$\mathbf{X} = (\mathbf{x}^1, \dots, \mathbf{x}^N)$$

- Features are columns of \mathbf{X} :
- Single features, e.g.: $x_{21} = 0.15$
- Feature vector, e.g.: $\mathbf{x}_{.1}$

- Scalar label for each data point:

$$\mathbf{y} = (y^1, \dots, y^N)^T$$

Data objects

- $\mathbf{x}^n \in \mathbb{R}^D$ or $\mathbf{x} \in \mathbb{R}^D$: the n -th input data point or a general data point, respectively. A column vector. Sometimes also used for the inputs of a particular neural network layer, then the dimensions could be $\mathbf{x} \in \mathbb{R}^J$.
- $y^n \in \mathbb{R}$ or $y \in \mathbb{R}$: a scalar label for the n -th data point or a general label y , respectively.
- $\mathbf{y}^n \in \mathbb{R}^K$ or $\mathbf{y} \in \mathbb{R}^K$: For multi-class or multi-task problems, this is a *label vector* for the n -th data point or a general label vector, respectively. A column vector.
- $\hat{y} \in \mathbb{R}$ or $\hat{\mathbf{y}} \in \mathbb{R}^K$: the predicted scalar label or – for a multi-class problem – the predicted label vector, respectively.
- $p \in \mathbb{R}$ or $\mathbf{p} \in \mathbb{R}^K$: same as above if outputs can be interpreted probabilistically, the predicted scalar label or – for a multi-class problem – the predicted label vector.
- $\mathbf{a} \in \mathbb{R}^I$: *activations* of a neural network.
- $\mathbf{s} \in \mathbb{R}^I$: *pre-activations*, also called *netI*, of a neural network.

Sizes and dimensions

- N : number of training examples, i.e. objects or samples, in the *training data set*. Running index: n .
- D : number of input units which is also the number of features that a sample has. Running index: d
- M : number of test or validation examples, i.e. objects or samples, in the *test data set*. Running index: m with $N + 1 \leq m \leq N + M$.
- K : number of output units. Running index: k .
- L : number of layers in a network without counting the input layer. Running index: l .
- J : input dimension of a general neural network layer. Running index: j .
- I : output dimension of a general neural network layer. Running index: i .

Stacked data and labels

- $\mathbf{X} \in \mathbb{R}^{N \times D}$: input data matrix. The rows represent objects, i.e. samples. We assume that objects, i.e. samples, are represented or described by *feature vectors* \mathbf{x}^n .
- $\mathbf{y} \in \mathbb{R}^N$: the scalar labels of all samples stacked to a column vector.
- $\mathbf{Y} \in \mathbb{R}^{N \times K}$: for multi-class or multi-task problems, stacked labels yield a *label matrix*.

Parameter objects

- $\mathbf{w} \in \mathbb{R}^D$: a *weight* or *parameter* vector of a simple machine learning method, such as linear regression. A column vector.
- $\mathbf{W} \in \mathbb{R}^{I \times J}$: a *weight* or *parameter* matrix of a learning method mapping from an input space with dimension J to an output space with dimension I .
- $\mathbf{b} \in \mathbb{R}^I$: a bias vector.
- θ : a set of parameters of a probabilistic model.

Multiple layers

- $n_h^{[l]}$: number of hidden units of the l -th layer. Thus, $D = n_h^{[0]}$ and $K = n_h^{[L+1]}$.
- $\mathbf{a}^{[l]} \in \mathbb{R}^{n_h^{[l]}}$: *activations* of a neural network in the l -th layer.
- $\mathbf{s}^{[l]} \in \mathbb{R}^{n_h^{[l]}}$: *pre-activations*, also called *netI*, of a neural network in the l -th layer.
- $\mathbf{W}^{[l]} \in \mathbb{R}^{n_h^{[l-1]} \times n_h^{[l]}}$: the weight matrix connecting the $(l - 1)$ -th layer with the l -th layer.
- $\mathbf{b}^{[l]} \in \mathbb{R}^{n_h^{[l]}}$: the bias vector in the l -th layer.

Common transformations

$$\mathbf{w}^T \mathbf{x} + b = w_1 x_1 + \dots w_D x_D + b. \quad (1)$$

$$\mathbf{s} = \mathbf{W} \mathbf{x} + \mathbf{b}. \quad (2)$$

$$\mathbf{a} = f(\mathbf{W} \mathbf{x} + \mathbf{b}). \quad (3)$$

$$\mathbf{w}^T \mathbf{w} = \sum_{i=1}^I w_i^2 = \|\mathbf{w}\|^2, \quad (4)$$

Functions

- $f : \mathbb{R} \rightarrow \mathbb{R}$: a non-linear *activation function* that is applied element-wise to a vector or matrix. Sometimes also denoted as ϕ .
- $g(\mathbf{x}; \mathbf{w})$: a machine learning model with input \mathbf{x} and parameters \mathbf{w} .
- $p(\mathbf{x}; \boldsymbol{\theta})$: a probabilistic model with data point \mathbf{x} and set of parameters $\boldsymbol{\theta}$.
- $L(y, g(\mathbf{x}; \mathbf{w}))$ or $L(y, \hat{y})$: a loss function L .
- $R_{\text{emp}}(\mathbf{y}, \mathbf{X}, \mathbf{w})$: an empirical error or risk function. Typically, the empirical error is an average of the loss function for a single training data point. For neural networks, this serves as a *cost function* that is minimized.
- $R_{\text{emp}}(\mathbf{w})$: an empirical error or risk function. The same as above only that occasionally the dependency on \mathbf{X} and \mathbf{y} is dropped to keep the notation uncluttered.

Derivatives and matrix layout

We use the so-called *numerator layout* of matrix calculus:

- $\frac{\partial R(\mathbf{w})}{\partial \mathbf{w}}$: a row vector of length W according to the definition of the Jacobian.
- $\frac{\partial \mathbf{a}}{\partial x}$: column vector
- $\frac{\partial a}{\partial \mathbf{x}}$: row vector
- $\frac{\partial \mathbf{a}}{\partial \mathbf{x}}$: matrix with as many rows as dimension of \mathbf{a} , as many columns as dimension of \mathbf{x} .
- $\frac{\partial a}{\partial \mathbf{X}}$: dimension of transposed \mathbf{X}

With the *Nabla operator* ∇ , we can switch to *denominator layout*:

- $\nabla R(\mathbf{w})$ or $\nabla_{\mathbf{w}} R(\mathbf{w})$: a column vector of length W . $\nabla_{\mathbf{w}} R(\mathbf{w}) = \left(\frac{\partial R(\mathbf{w})}{\partial \mathbf{w}} \right)^T$

Distributions

- $\mathcal{U}(a, b)$: a uniform distribution in the interval $[a, b]$. This distribution has mean $\frac{1}{2}(a + b)$ and variance $\frac{1}{12}(b - a)^2$.
- $\mathcal{N}(\mu, \sigma^2)$: a normal distribution with mean μ and variance σ^2 . Also commonly referred to as the Gaussian distribution.
- $\mathcal{B}(n, p)$: a binomial distribution with size parameter n and probability parameter p . Sometimes used as $\mathcal{B}(1, p)$ to denote Bernoulli distributions.