

Introduction to Deep Learning

Chapter 5: MLP – Matrix Notation

Bernd Bischl

Department of Statistics – LMU Munich

Winter term 2021

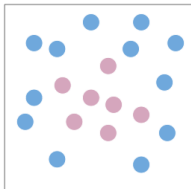


MOTIVATION

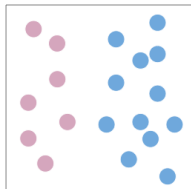
- The graphical way of representing simple functions/models, like logistic regression. Why is that useful?
- Because individual neurons can be used as building blocks of more complicated functions.
- Networks of neurons can represent extremely complex hypothesis spaces.
- Most importantly, it allows us to define the “right” kinds of hypothesis spaces to learn functions that are more common in our universe in a data-efficient way (see Lin, Tegmark et al. 2016).

MOTIVATION

- As a single neuron is restricted to learning only linear decision boundaries, its performance on the following task is quite poor:

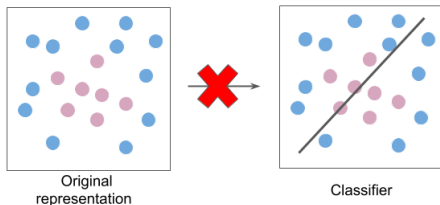


- However, the neuron can easily separate the classes if the original features are transformed (e.g., from Cartesian to polar coordinate):



MOTIVATION

- Instead of classifying the data in the original representation,

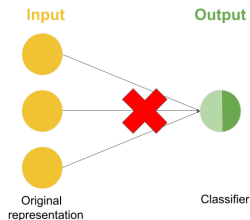


- we classify it in a new feature space.

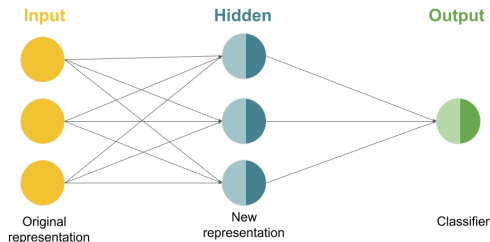


MOTIVATION

- Analogously, instead of a single neuron,



- we use more complex networks.



REPRESENTATION LEARNING

- It is therefore *very* critical to feed a classifier the “right” features in order for it to perform well.
- Before deep learning took off, features for tasks like machine vision and speech recognition were “hand-designed” by domain experts. This step of the machine learning pipeline is called **feature engineering**.
- The single biggest reason DL is so important is that it automates feature engineering. This is called **representation learning**.

SINGLE HIDDEN LAYER NETWORKS

Recall: a single neuron performs a 2-step computation:

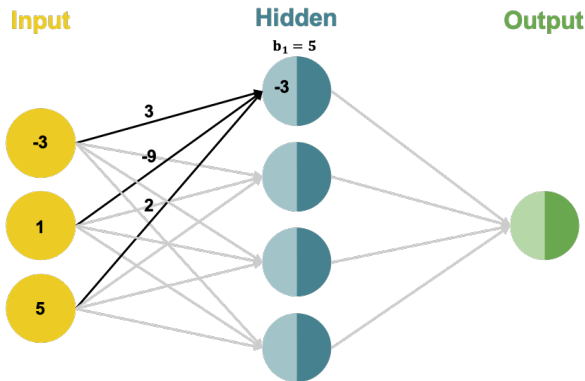
- ➊ **Affine Transformation:** weighted sum of inputs plus bias.
- ➋ **Non-linear Activation:** a non-linear transformation applied to the weighted sum.
- A single hidden layer network is composed of a **hidden layer** having a set of neuron.
- Each neuron located in the hidden layer performs a 2-step computation.
- It is possible to feed multiple inputs to a neural network simultaneously.
- The inputs $\mathbf{x}^{(i)}$, for $i \in \{1, \dots, n\}$, are arranged as rows in the **design matrix \mathbf{X}** .
- The weighted sum in the hidden layer is now computed as $\mathbf{XW} + \mathbf{B}$, where,

SINGLE HIDDEN LAYER NETWORKS

- The *matrix* of hidden activations $\mathbf{Z} = \sigma(\mathbf{XW} + \mathbf{B})$
- The final output of the network, which contains a prediction for each input, is $\tau(\mathbf{Zu} + \mathbf{C})$, where

SINGLE HIDDEN LAYER NETWORKS: EXAMPLE

- Weights (and biases) of the network.

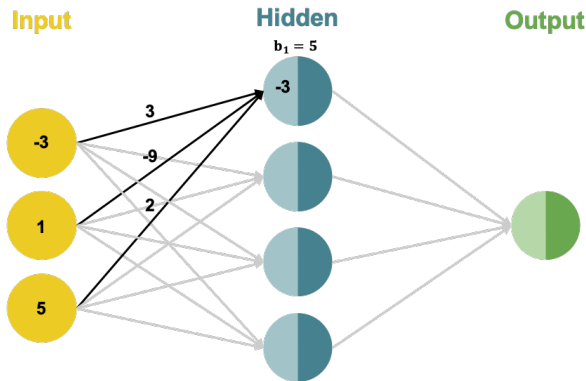


$$z_{\text{in}}^{(1)} = w_{11}x^{(1)} + w_{21}x^{(2)} + w_{31}x^{(3)} + b_1$$

$$z_{\text{in}}^{(1)} = 3 * (-3) + (-9) * 1 + 2 * 5 + 5 = -3$$

SINGLE HIDDEN LAYER NETWORKS: EXAMPLE

- Weights (and biases) of the network.

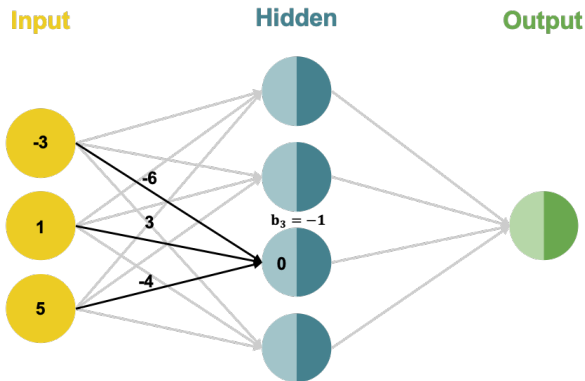


$$z_{in}^{(1)} = w_{11}x^{(1)} + w_{21}x^{(2)} + w_{31}x^{(3)} + b_1$$

$$z_{in}^{(1)} = 3 * (-3) + (-9) * 1 + 2 * 5 + 5 = -3$$

SINGLE HIDDEN LAYER NETWORKS: EXAMPLE

- Weights (and biases) of the network.

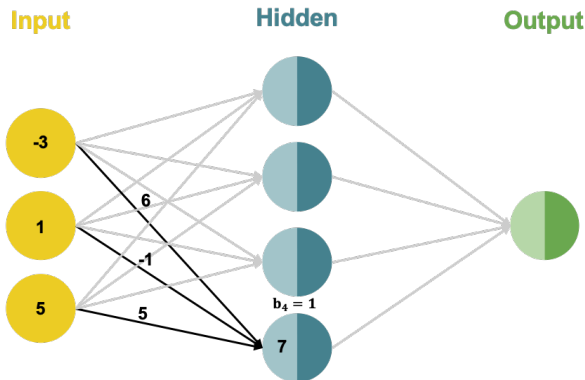


$$z_{\text{in}}^{(3)} = w_{13}x^{(1)} + w_{23}x^{(2)} + w_{33}x^{(3)} + b_3$$

$$z_{\text{in}}^{(3)} = (-6) * (-3) + 3 * 1 + (-4) * 5 - 1 = 0$$

SINGLE HIDDEN LAYER NETWORKS: EXAMPLE

- Weights (and biases) of the network.

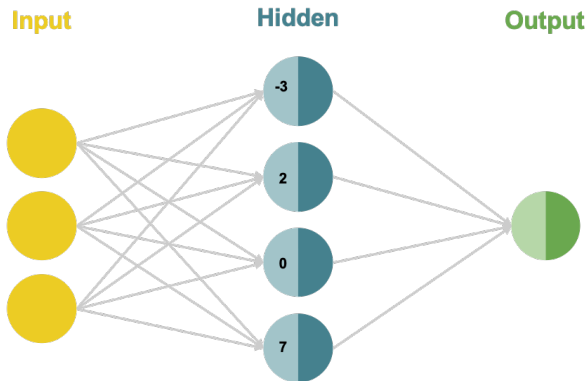


$$z_{in}^{(4)} = w_{14}x^{(1)} + w_{24}x^{(2)} + w_{34}x^{(3)} + b_4$$

$$z_{in}^{(4)} = 6 * (-3) + (-1) * 1 + 5 * 5 + 1 = 7$$

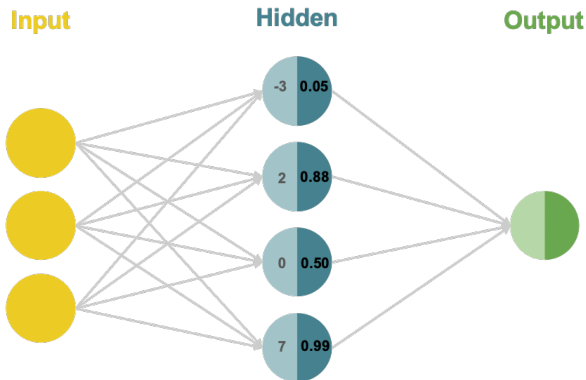
SINGLE HIDDEN LAYER NETWORKS: EXAMPLE

- Weights (and biases) of the network.



SINGLE HIDDEN LAYER NETWORKS: EXAMPLE

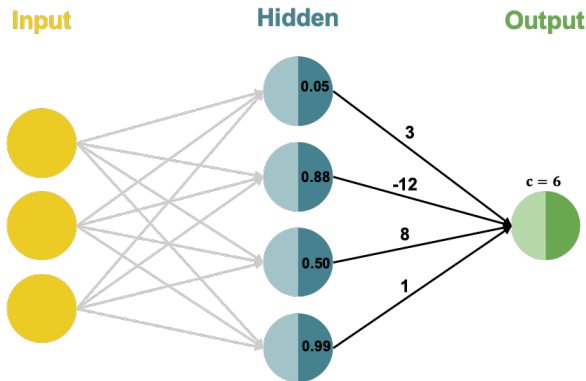
- Weights (and biases) of the network.



$$\mathbf{z}_{\text{out}}^{(i)} = \sigma \left(\mathbf{z}_{\text{in}}^{(i)} \right) = \frac{1}{1 + e^{\mathbf{z}_{\text{in}}^{(i)}}}$$

SINGLE HIDDEN LAYER NETWORKS: EXAMPLE

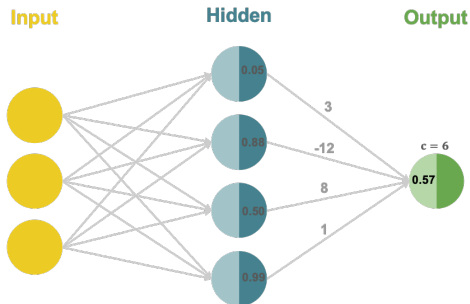
Forward pass through the shallow neural network.



$$\mathbf{f}_{\text{in}} = \mathbf{u}_1 \mathbf{z}_{\text{out}}^{(1)} + \mathbf{u}_2 \mathbf{z}_{\text{out}}^{(2)} + \mathbf{u}_3 \mathbf{z}_{\text{out}}^{(3)} + \mathbf{u}_4 \mathbf{z}_{\text{out}}^{(4)} + c$$

SINGLE HIDDEN LAYER NETWORKS: EXAMPLE

Forward pass through the shallow neural network.

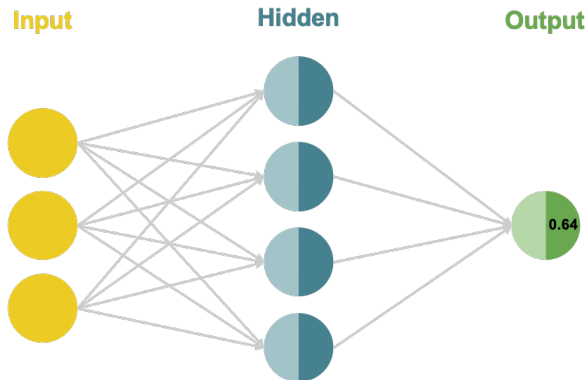


$$\mathbf{f}_{\text{in}} = \mathbf{u}_1 \mathbf{z}_{\text{out}}^{(1)} + \mathbf{u}_2 \mathbf{z}_{\text{out}}^{(2)} + \mathbf{u}_3 \mathbf{z}_{\text{out}}^{(3)} + \mathbf{u}_4 \mathbf{z}_{\text{out}}^{(4)} + c$$

$$\mathbf{f}_{\text{in}} = 3 * 0.05 + (-12) * 0.88 + 8 * 0.50 + 1 * 0.99 + 6 = 0.57$$

SINGLE HIDDEN LAYER NETWORKS: EXAMPLE

Forward pass through the shallow neural network.



$$\mathbf{f}_{\text{out}} = \sigma(\mathbf{f}_{\text{in}}) = \frac{1}{1 + e^{-\mathbf{f}_{\text{in}}}}$$
$$\mathbf{f}_{\text{out}} = \frac{1}{1 + e^{0.57}} = 0.64$$

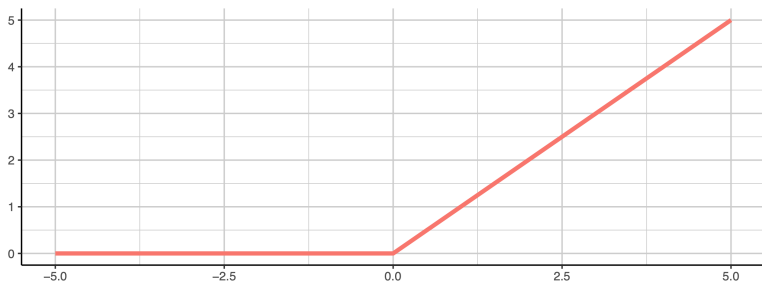
HIDDEN LAYER: ACTIVATION FUNCTION

Note: if the hidden layer does not have a non-linear activation, the network can only learn linear decision boundaries.

ReLU Activation:

- Currently the most popular choice is the ReLU (rectified linear unit):

$$\sigma(v) = \max(0, v)$$

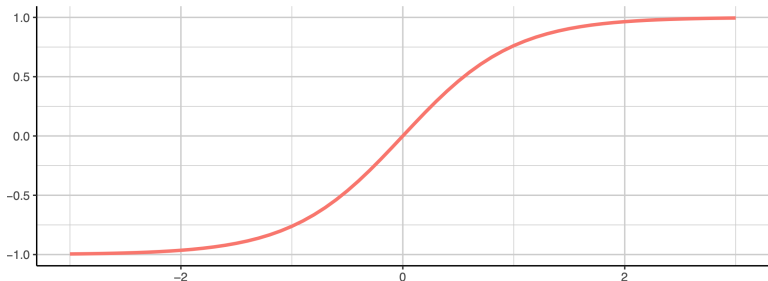


HIDDEN LAYER: ACTIVATION FUNCTION

Hyperbolic Tangent Activation:

- Another choice might be the hyperbolic tangent function:

$$\sigma(v) = \tanh(v) = \frac{\sinh(v)}{\cosh(v)} = 1 - \frac{2}{\exp(2v) + 1}$$

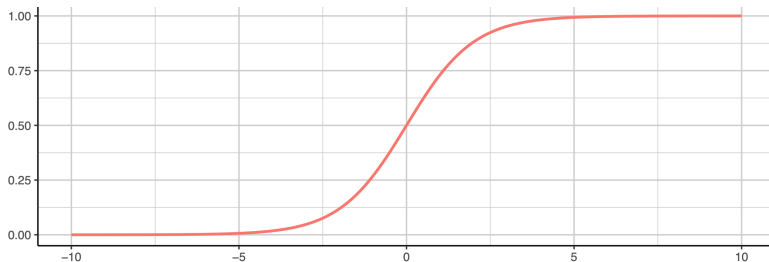


HIDDEN LAYER: ACTIVATION FUNCTION

Sigmoid Activation Function:

- The sigmoid function can be used even in the hidden layer:

$$\sigma(v) = \frac{1}{1 + \exp(-v)}$$



DEEP FEEDFORWARD NETWORK

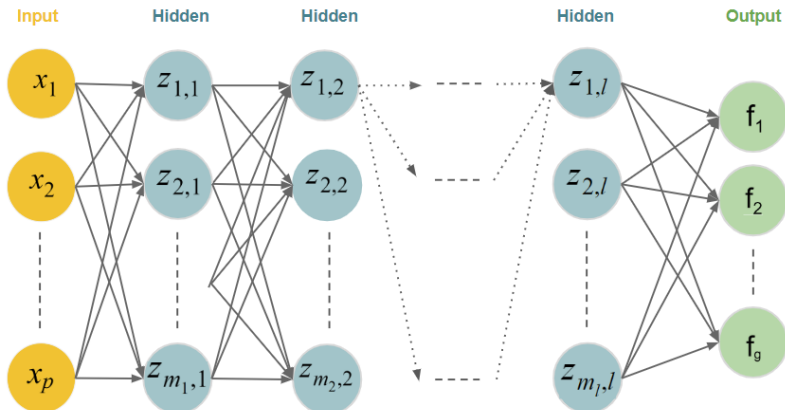


Figure: Structure of a deep neural network with l hidden layers.

WHY ADD MORE LAYERS?

- Multiple layers allow for the extraction of more and more abstract representations.
- Each layer in a feed-forward neural network adds its own degree of non-linearity to the model.

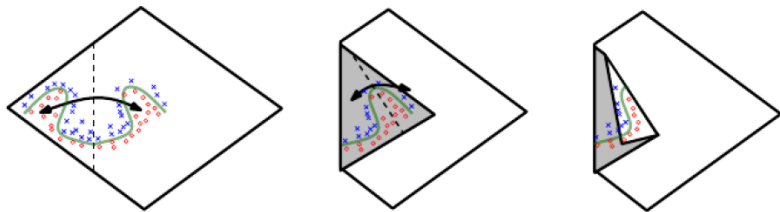


Figure: An intuitive, geometric explanation of the exponential advantage of deeper networks formally (Montúfar et al. (2014)).

DEEP NEURAL NETWORKS

Neural networks today can have dozens or even hundreds of hidden layers. The greater the number of layers, the "deeper" the network. Historically, however, deep neural networks were very challenging to train for several reasons:

- ❶ For one thing, the use of sigmoid activations (e.g., logistic sigmoid and tanh) significantly slowed down training due to a phenomenon known as “vanishing gradients”. The introduction of the ReLU activation largely solved this problem.
- ❷ Training deep neural networks on CPUs was too slow to be practical. Switching over to GPUs cut down training time by more than an order of magnitude.
- ❸ Another reason neural networks were not popular until the late '00s is that when dataset sizes are small, other models (such as SVMs) and techniques (such as feature engineering) outperform them.

DEEP NEURAL NETWORKS

- The availability of large datasets and novel architectures that are capable to handle even complex tensor-shaped data (e.g. CNNs for image data), faster hardware, and better optimization and regularization methods made it feasible to successfully implement deep neural networks in the last decade.
- An increase in depth often translates to an increase in performance on a given task. State-of-the-art neural networks, however, are much more sophisticated than the simple architectures we have encountered so far.
- The term "**deep learning**" encompasses all of these developments and refers to the field as a whole.