

# HANDS-ON AI I

## Your First Neural Networks



Sohvi Luukkonen

**Institute for Machine Learning**

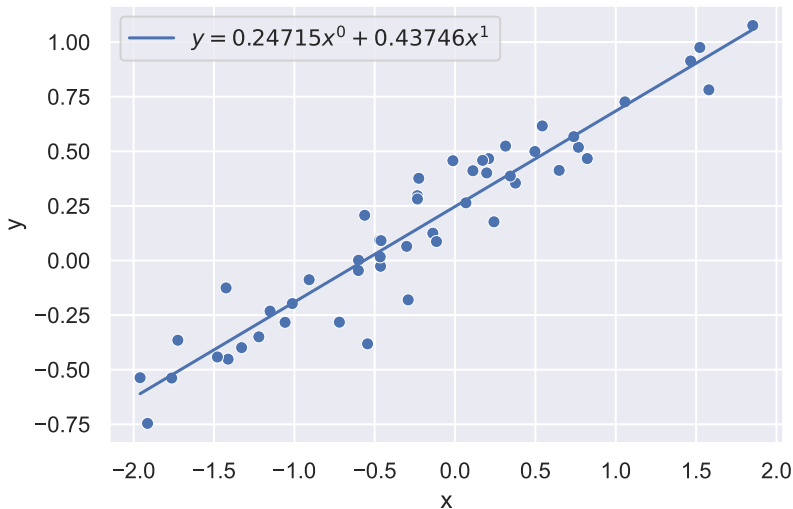
## Copyright Statement

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.

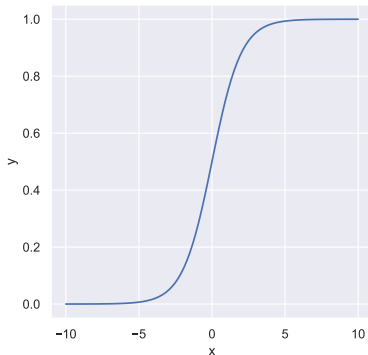
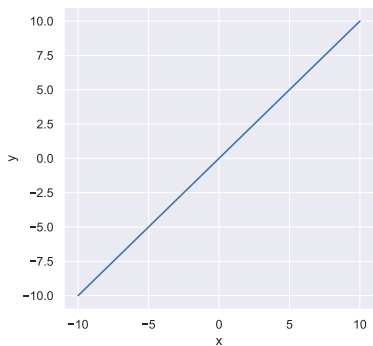
## Content of Unit 5

- Let's put everything together and build our first neural networks.
- We elaborate a bit on:
  - Loss functions
  - Optimization of loss functions
  - From Logistic Regression towards Neural Networks
  - Logistic Regression and Neural Networks coded in PyTorch

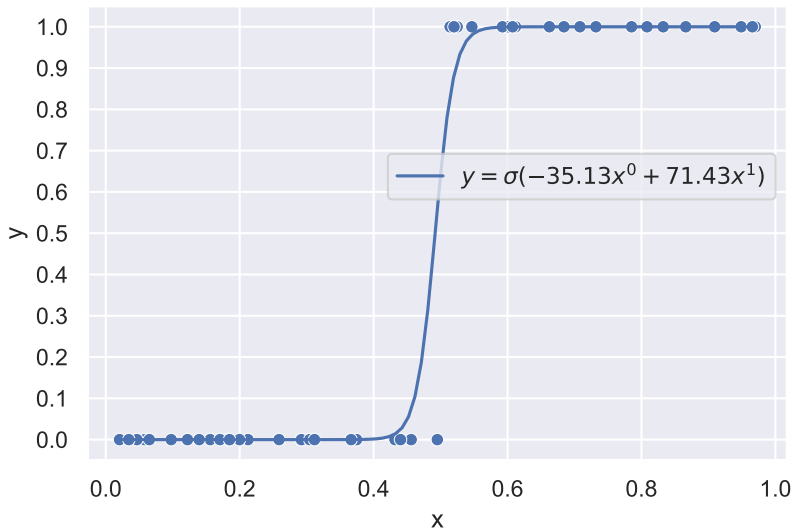
# Recap: Linear Regression



# Recap: Linear Function $\Rightarrow$ Logistic Function



## Recap: Logistic Regression

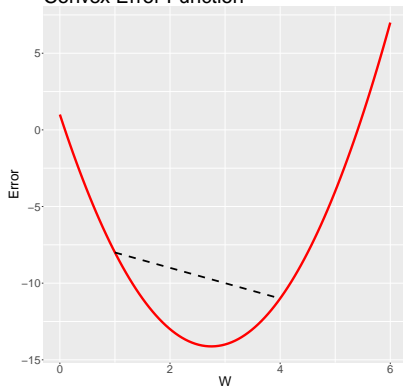


## Recap: Loss Function

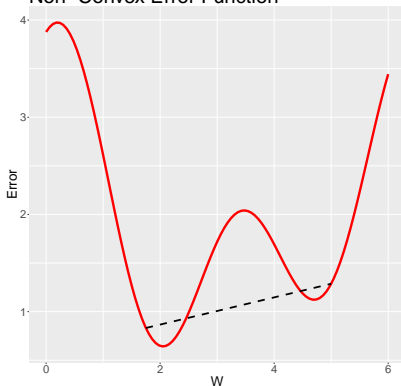
- The loss function tells **how good the prediction is**.
- The loss function gives the **error on the predictions**.
- We can optimize our machine learning algorithm knowing this error → we want to **minimize** the error/loss.
- Common loss functions:
  - Regression: **Mean-Squared Error** (MSE)
  - Classification: (Binary) **Cross Entropy** ((B)CE)

# Recap: Convex vs. Non-convex Functions

Convex Error Function



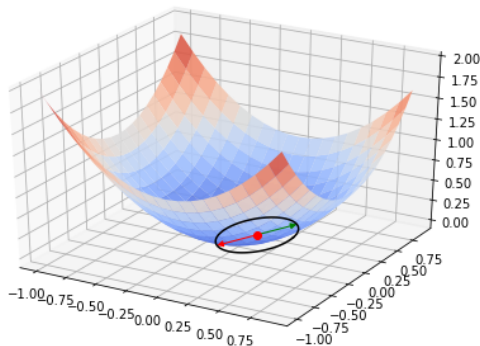
Non-Convex Error Function





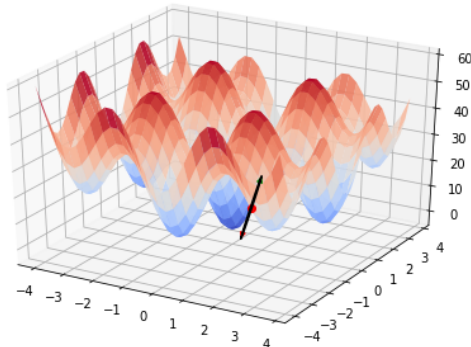
# Convex Loss Functions

- Convex loss functions only have one **minimum** (global).
- Often (but not always!) easy to optimize (no iterative algorithms required).



# Non-Convex Loss Functions

- Non-convex loss functions are harder to handle since they have different **local** and **global** minima.
- Difficult to optimize (in the majority of the cases, **iterative methods needed**).
- Pretty common in machine learning.

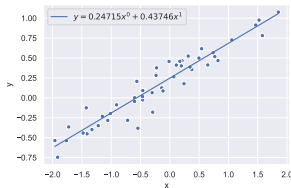


## One Subtlety

- A non-convex loss function **often** has a non-closed-form solution, i.e., iterative methods are needed, there is no analytical solution.
- A convex loss function has **sometimes** a closed-form solution (linear regression), and sometimes not (logistic regression).
- While a loss function itself can be convex with respect to its own parameters (e.g., MSE), this can (and typically will) change when applied to a neural network with respect to its parameters.

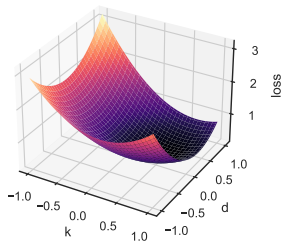
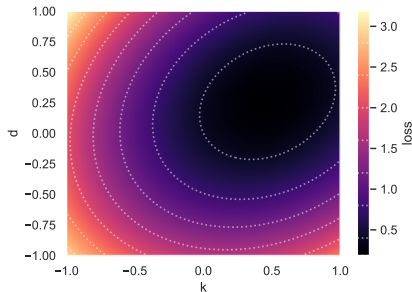
# Mean-Squared Error for Linear Regression

$$\begin{aligned}\text{MSE} &= \frac{1}{n} \sum_{i=1}^n (y_i - g(x_i; d, k))^2 && \text{with } g(x_i; d, k) = d + k \cdot x_i \\&= \frac{1}{n} \sum_{i=1}^n (y_i - d - k \cdot x_i)^2 \\&= \frac{1}{n} \sum_{i=1}^n y_i^2 - 2y_i d - 2y_i k x_i - d^2 - 2dk x_i - k^2 x_i^2 \\&= -1 \cdot d^2 - \left( \frac{1}{n} \sum_{i=1}^n x_i^2 \right) \cdot k^2 - \left( \frac{1}{n} \sum_{i=1}^n 2x_i \right) \cdot d \cdot k - \left( \frac{1}{n} \sum_{i=1}^n 2y_i \right) \cdot d \\&\quad - \left( \frac{1}{n} \sum_{i=1}^n 2y_i x_i \right) \cdot k - \left( \frac{1}{n} \sum_{i=1}^n y_i^2 \right) \\&= a_1 \cdot d^2 + a_2 \cdot k^2 + a_3 \cdot d \cdot k + a_4 \cdot d + a_5 \cdot k + a_6\end{aligned}$$



# Closed-Form Solution

- The coefficients ( $a_i$ ) of the function can **directly be calculated from the data set**.
- For quadratic functions like the MSE in combination with linear regression, the **minimum can be computed analytically**, i.e., there are equations we can plug in our coefficients to obtain the coordinates of the minimum.

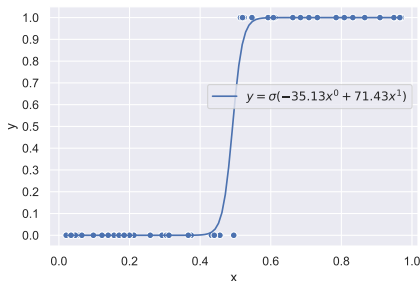


# Binary Cross Entropy for Logistic Regression

$$\text{BCE} = \frac{1}{n} \sum_{i=1}^n (-y_i \log(g(x_i; d, k)) - (1 - y_i) \log(1 - g(x_i; d, k)))$$

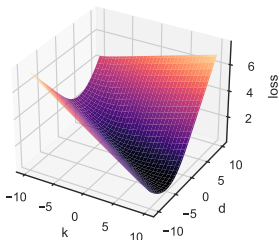
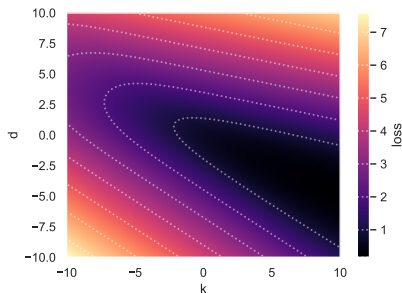
with  $g(x_i; d, k) = \sigma(d + k \cdot x_i)$

- Again, BCE can be regarded as a function that depends on the data set and the model parameters  $d, k$ .



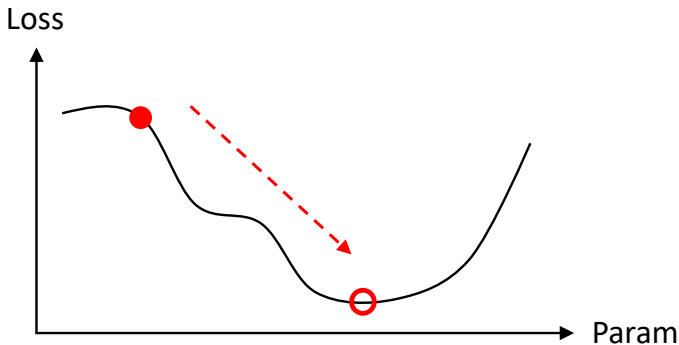
# No Closed-Form Solution

- There is **no closed-form solution** for the minimum that we can compute from the data set.
- We must use **iterative methods** such as **gradient descent** to solve this problem.



## Recap: Gradient Descent

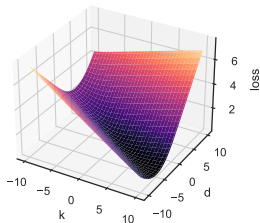
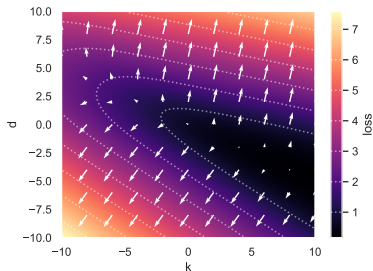
- Given: a function  $f(x)$
- Task: find  $x$  that maximizes (or minimizes)  $f(x)$
- Idea: start at some value  $x_0$ , and take a small step  $\eta$  in the direction in which the function decreases strongest





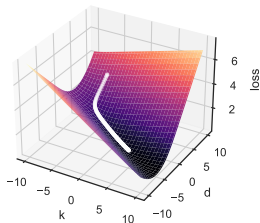
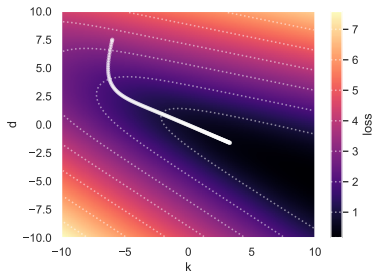
# Gradient Descent for Logistic Regression

- For iterative optimization, we need the **gradient of the loss function**.
- This gradient is another function of  $d$  and  $k$  which tells us the **slope of the loss function** at a particular position.
- At each position  $(d, k)$ , the gradient gives the **direction of steepest ascent**.



# Gradient Descent for Logistic Regression

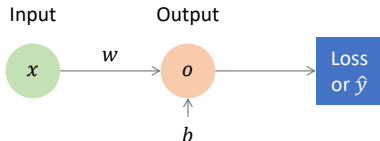
- **Gradient Descent** exploits exactly this principle: It starts at a random position and repeatedly takes **small steps in the direction opposing the gradient**.



# A Simple Logistic Regression Example

- Given binary labeled 1D data  $((x_1, y_1), \dots (x_n, y_n))$ , where  $x \in \mathbb{R}^1$  and  $y \in \{0, 1\}$  (i.e., two classes 0 and 1).
- A matching logistic regression model would look like  $o = g(x; \{d, k\}) = \sigma(z) = \sigma(d + k \cdot x)$ , where  $z$  is the “raw” model output and  $o$  the final model output (a probability) for some given input  $x$  and the model parameters  $d$  and  $k$ .
- As preparation for the following steps, let's substitute the model parameters as follows:  $d \rightarrow b, k \rightarrow w$ , so our model now looks like  $g(x; \{b, w\}) = \sigma(b + w \cdot x)$ . Also, for brevity, we can optionally write just  $g(x)$ , i.e., we do not explicitly include the parameters for a more compact notation.
- We will call  $b$  the **bias** and  $w$  the **weight**.

# Visualization and Output Calculation



- The above is a simplified visualization of the data flow through our model:  $o = \sigma(z) = \sigma(b + w \cdot x)$ .
- The output  $o$  can then be used for two things:
  - Calculating the **loss** during gradient descent to repeatedly update our model parameters (bias  $b$  and weight  $w$ ) to decrease the loss.
  - Determining the class **prediction**  $\hat{y}$  (recall: we have to apply some decision rule, e.g., via a threshold).

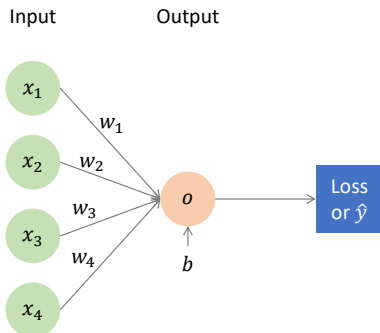
# Multi-Dimensional Example

- Given binary labeled  $M$ -dimensional data  $((\mathbf{x}_1, y_1), \dots (\mathbf{x}_n, y_n))$ , where  $\mathbf{x} \in \mathbb{R}^M$  and  $y \in \{0, 1\}$  (i.e., two classes 0 and 1).
- A matching logistic regression model would look like  $o = g(\mathbf{x}; \{b, \mathbf{w}\}) = \sigma(z) = \sigma(b + \mathbf{w} \cdot \mathbf{x})$ , where  $z$  is the “raw” model output and  $o$  the final model output (a probability) for some given input  $\mathbf{x}$  and the model parameters  $b$  and  $\mathbf{w}$ .
- Note that the parameter  $\mathbf{w}$  is now a vector,<sup>1</sup> and  $\mathbf{w} \cdot \mathbf{x}$  is the scalar/dot product, i.e.,  $\sum_{i=1}^M w_i x_i$ .

---

<sup>1</sup>Strictly speaking,  $b$  is internally also a vector since  $b + \mathbf{w} \cdot \mathbf{x}$  is actually  $b \cdot \mathbf{x}^0 + \mathbf{w} \cdot \mathbf{x}^1$ . However, because the computation does not depend on the input  $\mathbf{x}$ , i.e.,  $b \cdot \mathbf{x}^0 = b \cdot \mathbf{1} = \sum_{i=1}^M b_i \cdot 1 = \sum_{i=1}^M b_i$ , we simplify it to just  $b$ .

# Visualization and Output Calculation

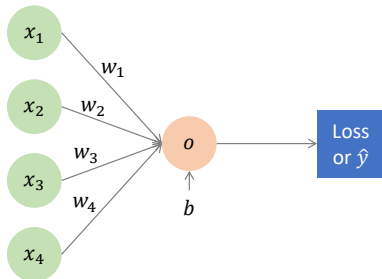


- The above is a simplified visualization of the data flow through our model with an  $M = 4$ -dimensional input size. See the next slide for how the final model output is calculated.

# Output Calculation

- Model parameters  $b$  and  $w$ :

$$b, w = \begin{bmatrix} w_1 & w_2 & w_3 & w_4 \end{bmatrix}$$



- “Raw”/linear model output  $z$  given  $M = 4$ -dimensional input  $x$  (with  $w \cdot x$  being the scalar product):

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}, z = b + \mathbf{w} \cdot \mathbf{x} = b + \sum_{i=1}^4 w_i x_i$$
$$= b + w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4$$

- Final model output  $o$  (class probability):

$$o = \sigma(z)$$

## Multi-Class Example

- Given  $K$ -multi-class labeled 1D data  $((x_1, y_1), \dots, (x_n, y_n))$ , where  $x \in \mathbb{R}^1$  and  $y \in \{1, \dots, K\}$ .
- Recall that the softmax function is suitable for multi-class classification, so we replace our logistic function  $\sigma(z)$  with the softmax function  $\sigma(z)_i$ , where  $i$  means the  $i$ -th class and  $z = (z_1, \dots, z_K)^\top$  is the vector of “raw” model outputs.



## Recap: Softmax

- Generalization of the sigmoid function.
- Suitable for **multi-class** classification.
- For  $K$  classes with  $y \in \{1, \dots, K\}$  the probability of  $x$  belonging to class  $i$  is:

$$p(y = i \mid \mathbf{x}) = \sigma(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

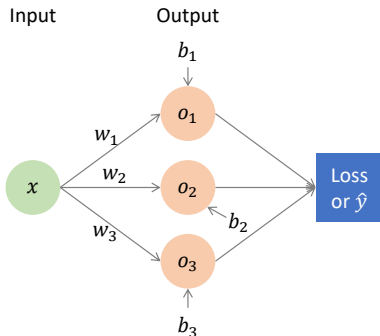
where  $\mathbf{z} = (z_1, \dots, z_K)$  is the vector of “raw” model outputs, i.e., there is an output for each class (see example later).

- While not necessary, as a regular sigmoid function would suffice, this also works for binary classification, i.e.,  $K = 2$ .

## Multi-Class Example

- Given  $K$ -multi-class labeled 1D data  $((x_1, y_1), \dots, (x_n, y_n))$ , where  $x \in \mathbb{R}^1$  and  $y \in \{1, \dots, K\}$ .
- Recall that the softmax function is suitable for multi-class classification, so we replace our logistic function  $\sigma(z)$  with the softmax function  $\sigma(z)_i$ , where  $i$  means the  $i$ -th class and  $z = (z_1, \dots, z_K)^\top$  is the vector of “raw” model outputs.
- A matching logistic regression model would look like  $\mathbf{o} = (o_1, \dots, o_K)^\top = g(x; \{\mathbf{b}, \mathbf{w}\}) = (\sigma(z)_1, \dots, \sigma(z)_K)^\top$ , where  $\mathbf{o}$  is the final vector of  $K$  class probabilities for some given input  $x$  and the model parameters  $\mathbf{b}$  and  $\mathbf{w}$ .
- Note that the parameters  $\mathbf{b}$  and  $\mathbf{w}$  are now vectors.

# Visualization

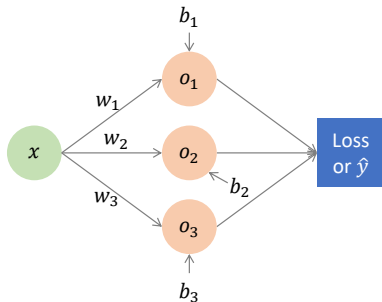


- The above is a simplified visualization of the data flow through our model with 1D input size and  $K = 3$  classes. See the next slide for how the final model output is calculated.

# Output Calculation

- Model parameters  $\mathbf{b}$  and  $\mathbf{w}$ :

$$\mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}, \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}$$



- “Raw”/linear model output  $\mathbf{z}$  given input  $x$ :

$$\mathbf{z} = \begin{bmatrix} z_1 \\ z_2 \\ z_3 \end{bmatrix} = \mathbf{b} + \mathbf{w} \cdot x = \begin{bmatrix} b_1 + w_1 x \\ b_2 + w_2 x \\ b_3 + w_3 x \end{bmatrix}$$

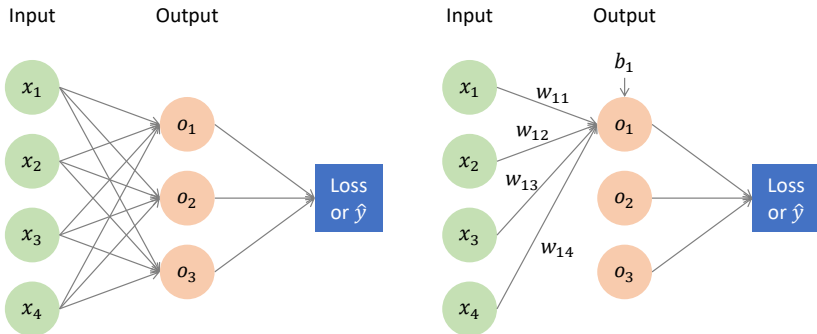
- Final model output  $\mathbf{o}$  (class probability vector):

$$\mathbf{o} = \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} = \begin{bmatrix} \sigma(\mathbf{z})_1 \\ \sigma(\mathbf{z})_2 \\ \sigma(\mathbf{z})_3 \end{bmatrix}$$

## Multi-Dimensional and Multi-Class Example

- Given  $K$ -multi-class labeled  $M$ -dimensional data  $((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n))$ , where  $\mathbf{x} \in \mathbb{R}^M$  and  $y \in \{1, \dots, K\}$ .
- A matching logistic regression model would look like  $\mathbf{o} = (o_1, \dots, o_K)^\top = g(\mathbf{x}; \{\mathbf{b}, \mathbf{W}\}) = (\sigma(\mathbf{z})_1, \dots, \sigma(\mathbf{z})_K)^\top$ , where  $\mathbf{o}$  is the final vector of  $K$  class probabilities for some given input  $\mathbf{x}$  and the model parameters  $\mathbf{b}$  and  $\mathbf{W}$ .
- Note that the parameter  $\mathbf{b}$  is now a vector and the parameter  $\mathbf{W}$  a matrix.

# Visualization

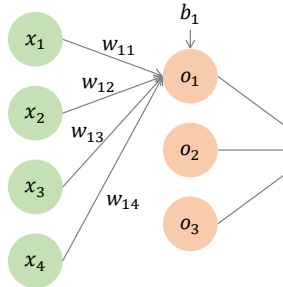


- The above is a simplified visualization of the data flow through our model with  $M = 4$ -dimensional input size and  $K = 3$  classes. See the next slide for how the final model output is calculated.

# Output Calculation

- Model parameters  $\mathbf{b}$  and  $\mathbf{W}$ :

$$\mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}, \mathbf{W} = \begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{14} \\ w_{21} & w_{22} & w_{23} & w_{24} \\ w_{31} & w_{32} & w_{33} & w_{34} \end{bmatrix}$$



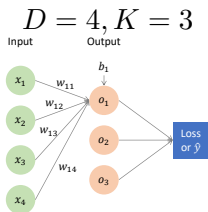
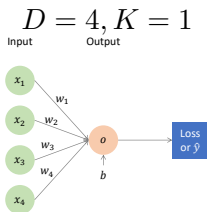
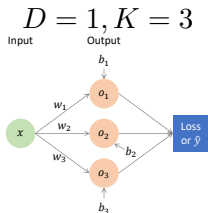
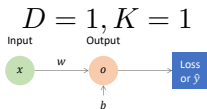
- “Raw”/linear model output  $\mathbf{z}$  given  $M = 4$ -dimensional input  $\mathbf{x}$  (with  $\mathbf{W}\mathbf{x}$  being the matrix product):

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}, \mathbf{z} = \begin{bmatrix} z_1 \\ z_2 \\ z_3 \end{bmatrix} = \mathbf{b} + \mathbf{W}\mathbf{x} = \begin{bmatrix} b_1 + \sum_{i=1}^4 w_{1i}x_i \\ b_2 + \sum_{i=1}^4 w_{2i}x_i \\ b_3 + \sum_{i=1}^4 w_{3i}x_i \end{bmatrix}$$

- Final model output  $\mathbf{o}$  (class probability vector):

$$\mathbf{o} = \begin{bmatrix} o_1 \\ o_2 \\ o_3 \end{bmatrix} = \begin{bmatrix} \sigma(\mathbf{z})_1 \\ \sigma(\mathbf{z})_2 \\ \sigma(\mathbf{z})_3 \end{bmatrix}$$

# Recap - Logistic Regression Calculation



$$\mathbf{o}_K^\top = \sigma \left( \mathbf{b}_K^\top + \mathbf{W}_{K \times D} \mathbf{x}_D \right)$$

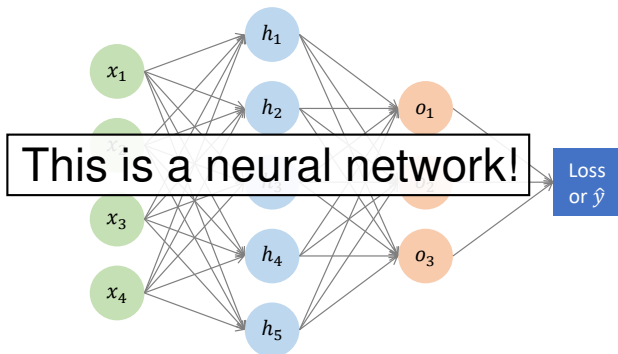


# Logistic Regression $\Rightarrow$ Neural Networks

- Again as preparation, let's rename *Input* and *Output*:
- What if we made our model *g* **more complex/powerful**?
- We could just reuse the idea of our bias-weight calculation by **adding** an intermediate/hidden **layer** *h* with output

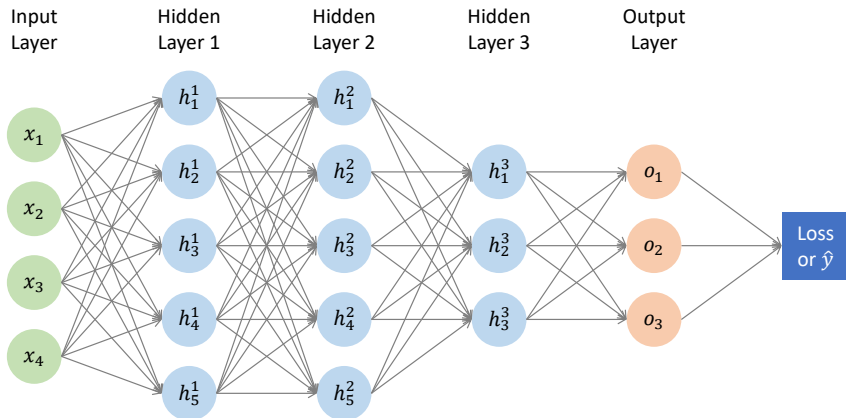
$$h = h(x; \{b, w\})$$

Input Layer
Hidden Layer
Output Layer



# Neural Networks

- We can add even more such layers, and the number of **nodes** (a.k.a. **units** or **neurons**) can also vary.
- E.g., 3 layers with a final output of  $o = g(h^3(h^2(h^1(x))))$ :



# Stacking Linear Layers

- There is a **caveat** when naively stacking such linear layers.
- If we perform **multiple linear transformations** in a row, they can be **collapsed into** a **single** linear transformation.

For  $n$  hidden layers  $h^i$  with output  $h^i = h^i(x; \{b^i, W^i\}) = b^i + W^i x$  or simply  $h^i(x)$  for short, we get:

$$\begin{aligned} h^n &= h^n(h^{n-1}(\dots h^1(x))) \\ &= b^n + W^n(b^{n-1} + W^{n-1}(\dots b^1 + W^1 x)) \\ &= b^n + W^n b^{n-1} + W^n W^{n-1} b^{n-2} + \dots + W^n W^{n-1} \dots W^1 x \\ &= b' + W' x \end{aligned}$$

- This means that simply stacking such layers **does not increase the complexity/power of our model** (it is just a single linear transformation again).

# Non-linearity

- To get rid of the linear transformation problem, we thus add a **non-linear function**  $f$  after each layer.

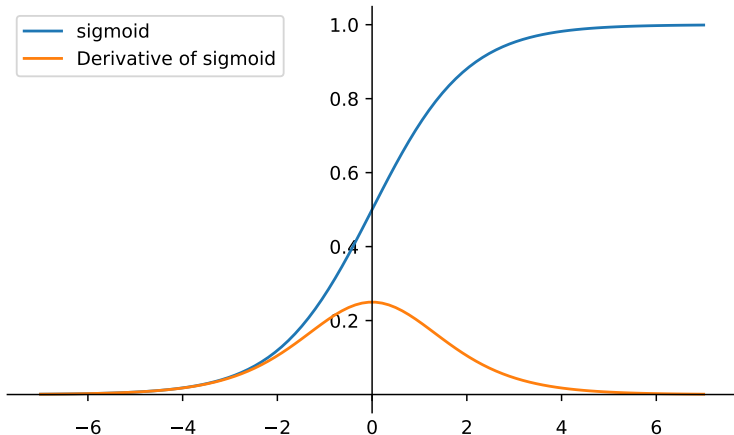
$$\text{Before: } h(x) = \mathbf{b} + \mathbf{W}x$$

$$\text{After: } h(x) = f(\mathbf{b} + \mathbf{W}x)$$

- Looks familiar? Indeed it does, this is the same idea we had when going from linear regression to logistic regression with  $f = \sigma$ .
- The logistic/sigmoid function is actually already such a function that we can use to introduce non-linearity.
- These functions are also called **activation functions**.

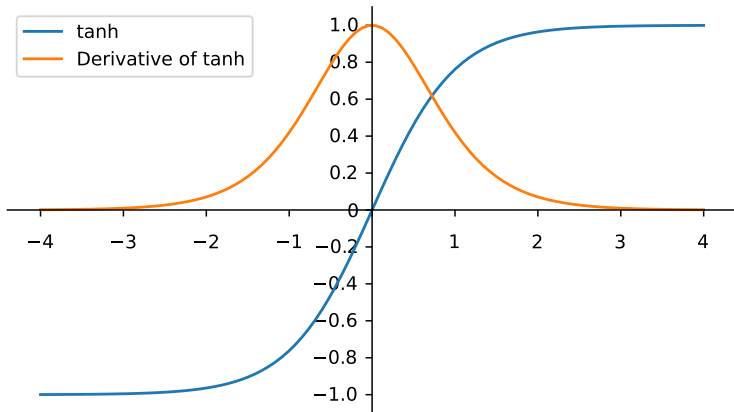
# Common Activation Functions

- There are various activation functions, each with their own benefits and disadvantages.
- Here are some often used activation functions:



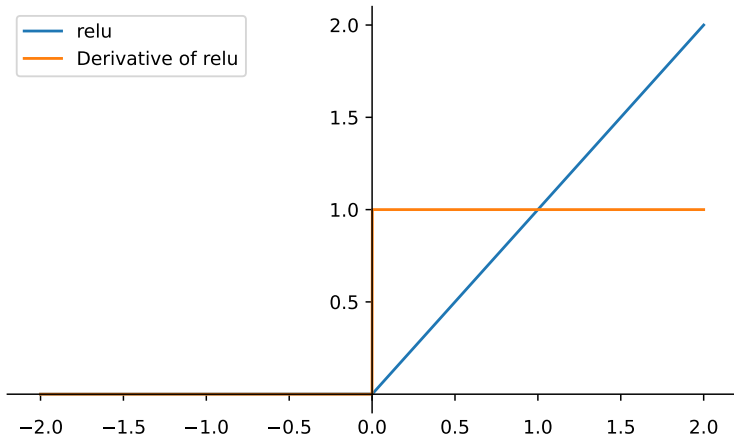
# Common Activation Functions

- There are various activation functions, each with their own benefits and disadvantages.
- Here are some often used activation functions:



# Common Activation Functions

- There are various activation functions, each with their own benefits and disadvantages.
- Here are some often used activation functions:

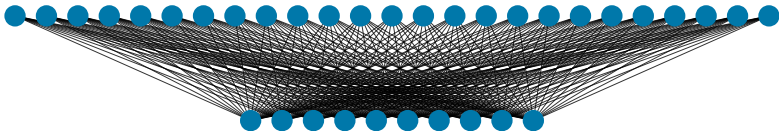


# Logistic Regression in PyTorch

- The following is just a brief summary. For more details, see the exercise and the accompanying materials.
- Define the model, e.g, for  $M = 25$ -dimensional input and  $K = 10$  classes, and define the loss:

```
import torch
```

```
model = torch.nn.Sequential(  
    torch.nn.Flatten(), # Ensure flat vector  
    torch.nn.Linear(25, 10)  
    #torch.nn.Softmax() # Not needed (incl. in loss)  
)  
loss = torch.nn.functional.cross_entropy
```

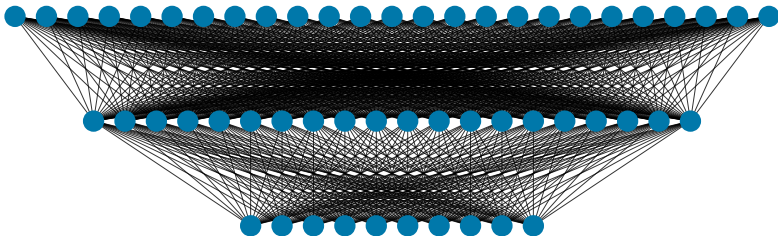




## Adding Another Layer

- Neural network model with one hidden layer and 20 nodes in PyTorch:

```
model = torch.nn.Sequential(  
    torch.nn.Flatten(),  
    torch.nn.Linear(25, 20),  
  
    torch.nn.Linear(20, 10)  
)
```



## Adding Non-linearity

- Neural network model with one hidden layer and 20 nodes in PyTorch including appropriate non-linearity:

```
model = torch.nn.Sequential(  
    torch.nn.Flatten(),  
    torch.nn.Linear(25, 20),  
    torch.nn.Sigmoid(), # or: Tanh(), ReLU()  
    torch.nn.Linear(20, 10)  
)
```

