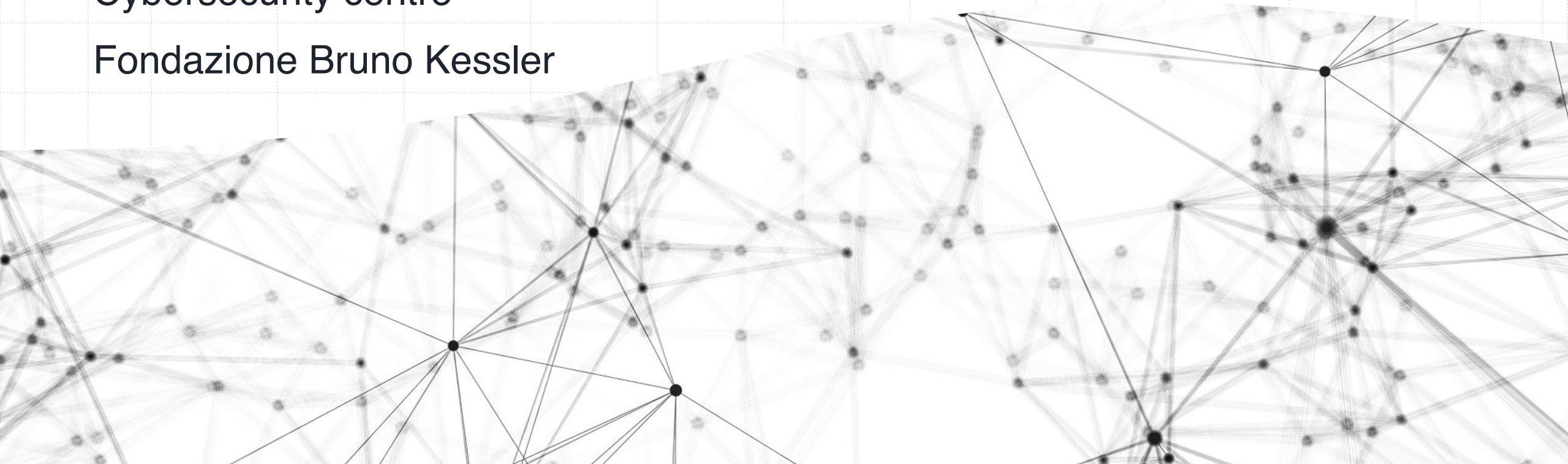# Model training

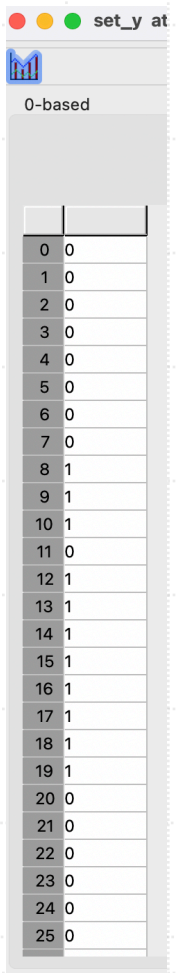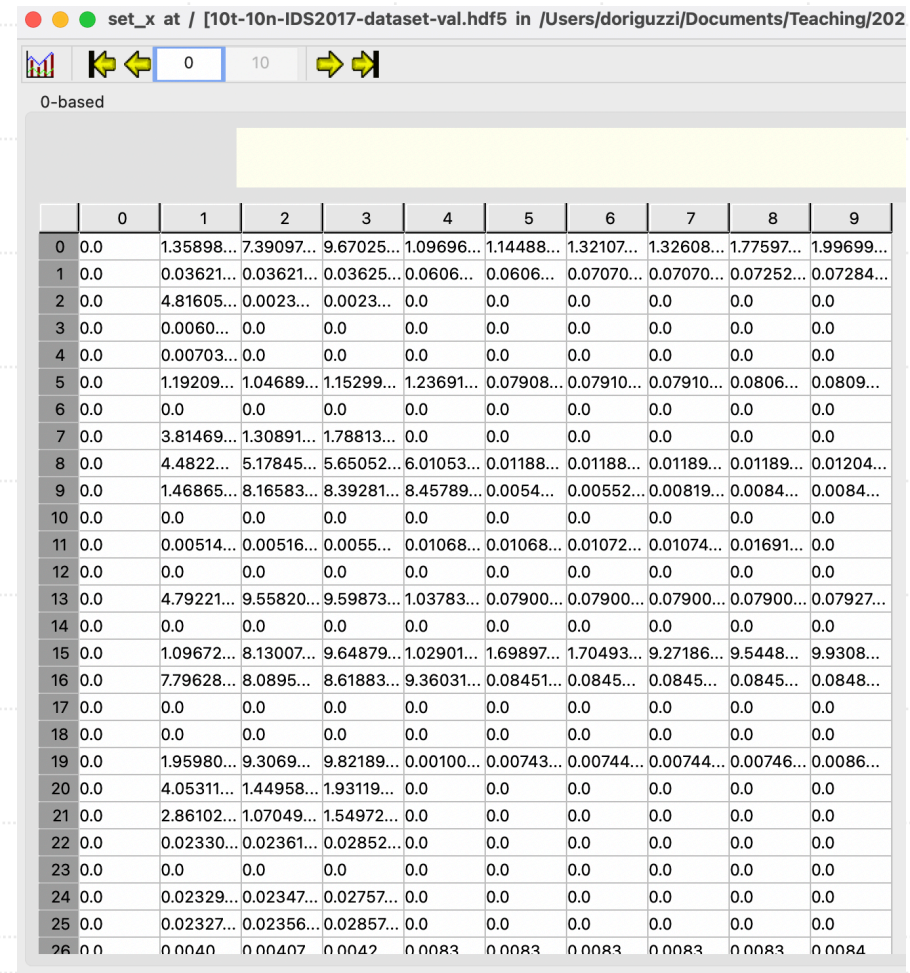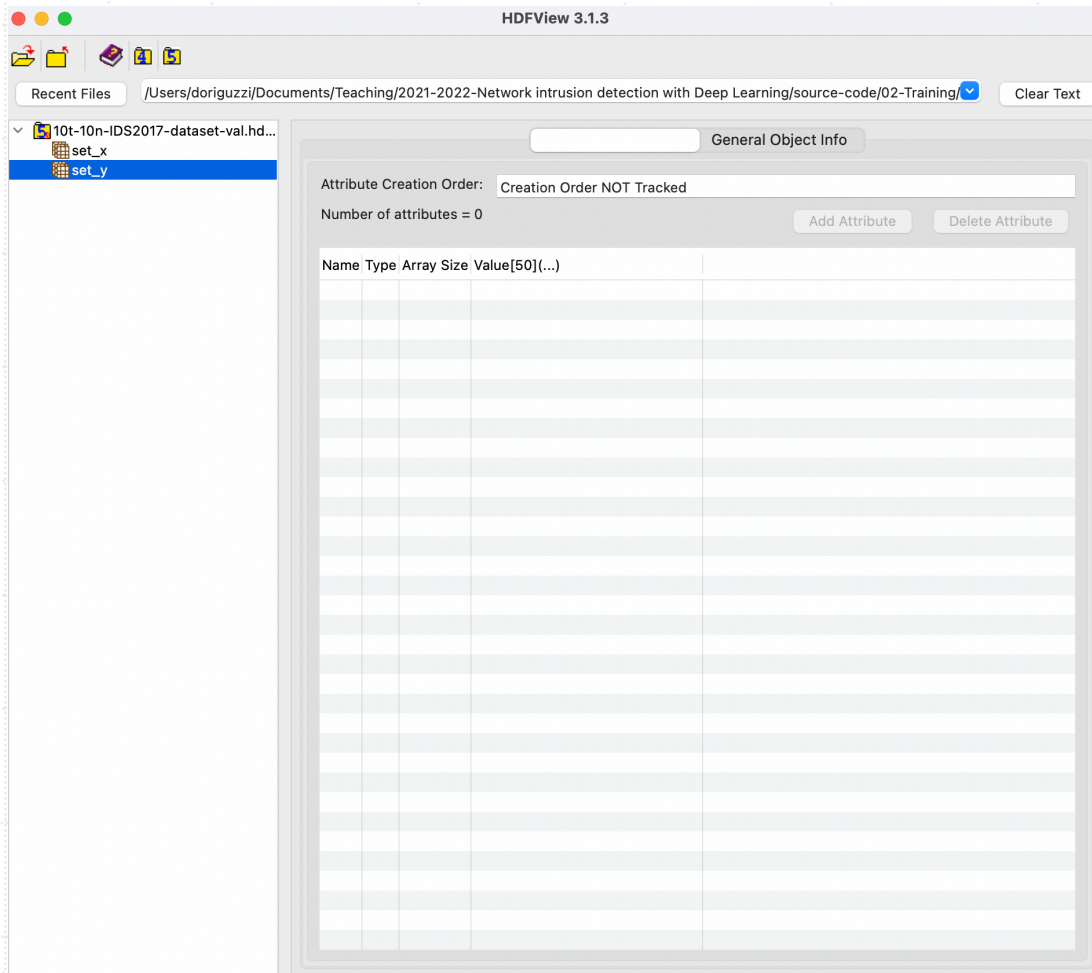**Roberto Doriguzzi Corin**

Cybersecurity centre

Fondazione Bruno Kessler

# Outline

- Hdf5 file format

- Optimizers

- Training a binary model for DDoS attack detection

- Laboratory: implement a multi-class model for DDoS attack detection
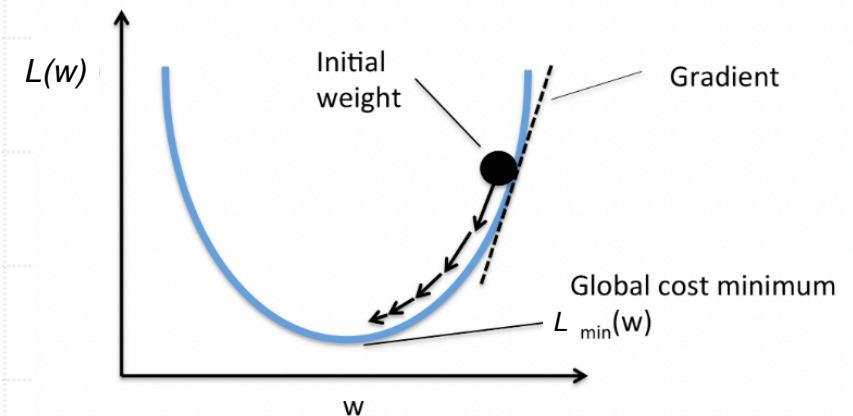
# Hdf5 file format

# Optimizers

**Optimizers** *are algorithms or methods used to change the parameters of your neural network (e.g., weights and learning rate) in order to reduce the loss.*

**Learning Rate** *defines how big/small the steps are gradient descent takes into the direction of the local minimum are determined by the learning rate.*

# Gradient Descent

- **Batch gradient descent:** all the data used into a single step per training epoch. We take the average of the gradients of all the training samples. Vectorization can be used to leverage parallel computation.

- **Stochastic gradient descent:** one sample for each step. The parameters are updated more frequently.

- **Mini-batch gradient descent:** steps with mini-batches of samples. Multiple steps each epoch, with the gain in speed thanks to vectorization.
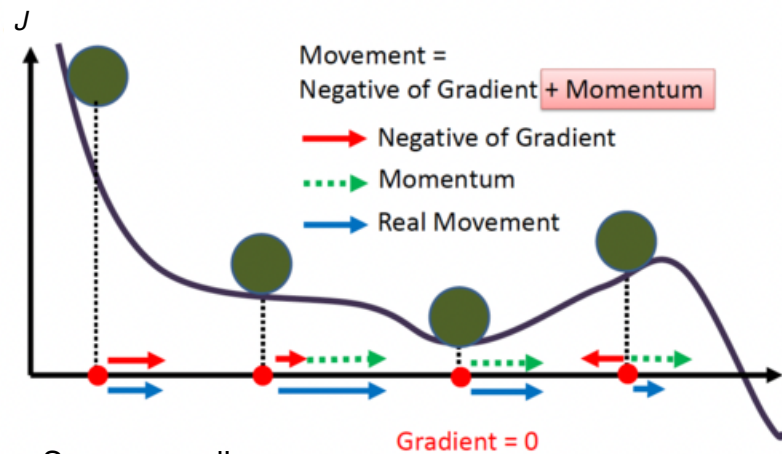


Source: sebastianraschka.com

$$w_t = w_{t-1} - \eta \frac{\partial L}{\partial w_{t-I}}$$

# Gradient Descent with Momentum

- Momentum is a hyper-parameter (like the learning rate) that has been introduced to:

  - Faster escape from plateaus
  - Escape from a local mimima
  - Tune the updates based on past gradients
  - Momentum $\beta = 0.9$ works well in practice

$J$

Movement =
Negative of Gradient + Momentum

→ Negative of Gradient

···▶ Momentum

→ Real Movement

Gradient = 0

Learning rate

$$w_t = w_{t-1} - \alpha \cdot m_t$$

$$m_t = \beta m_{t-1} + (1 - \beta)\frac{\partial J}{\partial w_{t-1}}$$

Exponential moving average
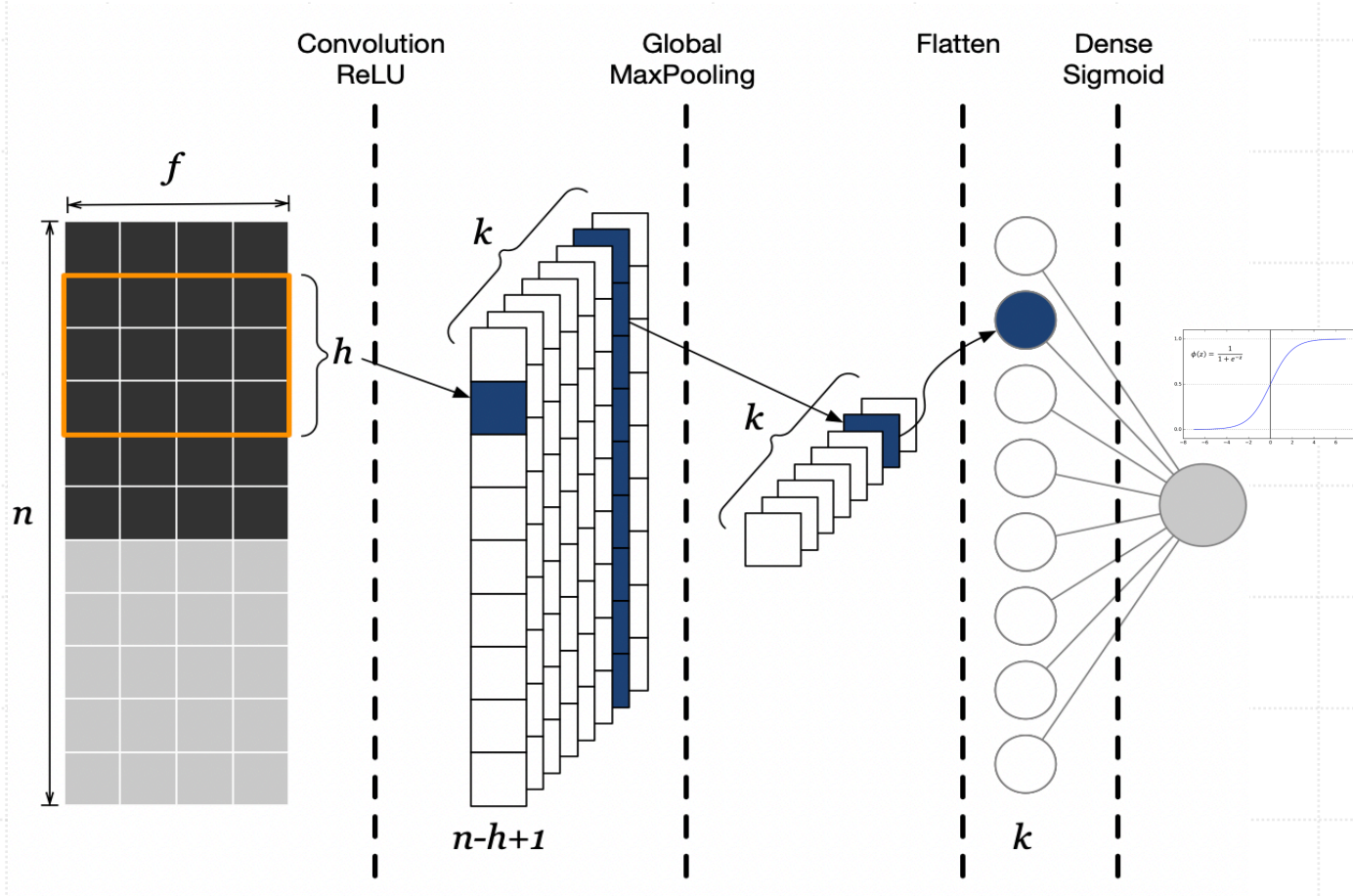
6

# Adam (adaptive momentum estimation)

- Adam is an optimisation algorithm that can be used instead of the classical stochastic gradient descent

- The algorithm calculates an **exponential moving average** of the gradient and the squared gradient, and the parameters $\beta_1$ and $\beta_2$ control the decay rates of these moving averages (usually set to $\beta_1 = 0.9$ and $\beta_2 = 0.999$).

- Converges faster than SGD on large problems in terms of data and parameters

$$m_t = \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot \frac{\partial J}{\partial w_{t-1}} \qquad v_t = \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot \left[\frac{\partial J}{\partial w_{t-1}}\right]^2$$

# Training a binary classifier

- CNN model

- Pre-processed balanced dataset of benign and DDoS attack traffic (HTTP attack from the CIC-IDS2017 dataset, generated with the LOIC tool https://www.imperva.com/learn/ddos/low-orbit-ion-cannon/)

- Make it run and play with the hyper-parameters

- Change the optimizer

- Change the model (e.g., replace the CNN with a MLP)
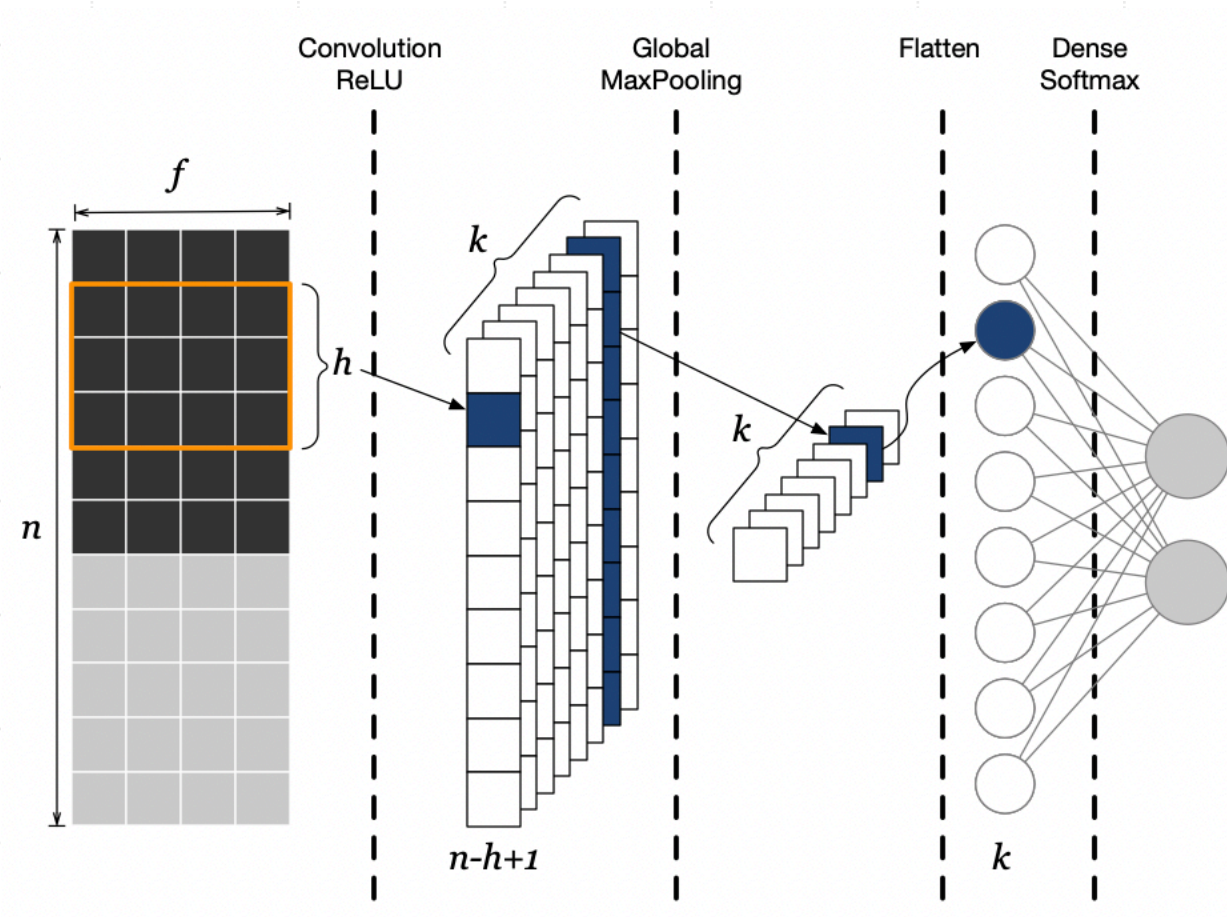
# Architecture of the CNN model



Hyperparameters
- $f$=11
- $n$=10
- $k$=32
- $h$=3
- Learning rate = 0.01

# Lab: implement a 2-class classifier

- Same CNN model with different output

- Replace the binary classifies with a 2-class classifier

- Pay attention to the format of the labels

# Architecture of the multi-class CNN model



Hyperparameters
- $f$=11
- $n$=10
- $k$=32
- $h$=3
- Learning rate = 0.01