**Part 1/2**
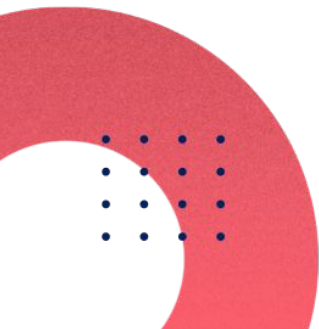
# No Mercy for Manual Entry

29/Sep/2021
Workshop @ AMLD2021
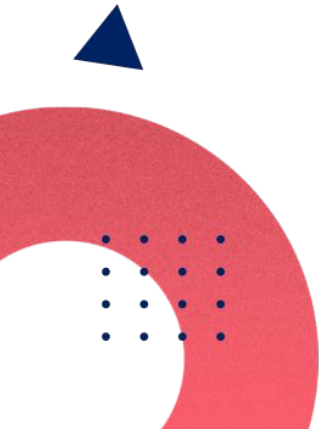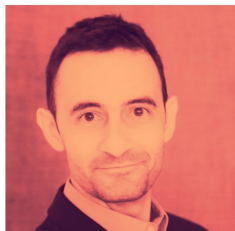
# Wifi Info

Network:     Free_STCC
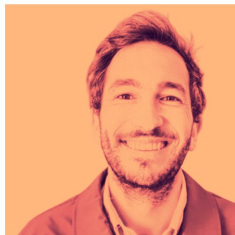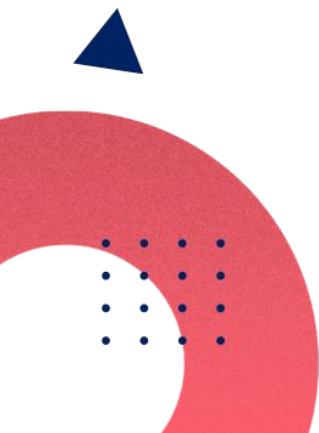User ID :    2317791957
Password:    5197

# Authors

**Valerio Rossetti, PhD**
Co-founder of SamurAI
Senior Data Scientist

**Giulio Grossi, PhD**
Senior Quantitative Portfolio Manager
at ONE swiss bank

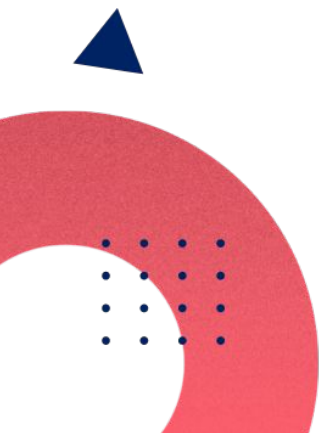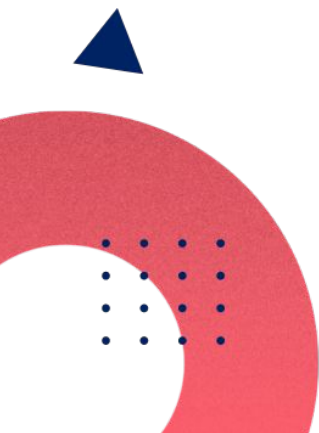# Get to know the audience

- Knowledge in Supervised Learning and Computer Vision

- Your coding skills in Python

- Do you have practical applications of the techniques in this workshop?

# Manual Entry

- A lot of information in paper/pdf documents:
  invoices, contracts, personal information, surveys …

SamurAI

AMLD EPFL

# Manual Entry

- A lot of information in paper/pdf documents: invoices, contracts, personal information, surveys ...
- These documents are then treated manually:
  - costly,
  - time-consuming and
  - error-prone

Manual image-to-data process

# Manual Entry

- A lot of information in paper/pdf documents: invoices, contracts, personal information, surveys …
- These documents are then treated manually:
  - costly,
  - time-consuming and
  - error-prone
- Main solution: go fully digital and eliminate the paper document!

Manual image-to-data process

# Manual Entry: a zombie among us

- In some cases we don't manage to eliminate the paper/pdf document

- Examples:
  - Traders send emails with pdf documents treated by the back-office
  - Commodity trading dealing with letters of credit
  - Banks dealing with client documents (passport, ID, scanned contract, …)
  - Medical research dealing with huge volumes of paper documents from patients
  - Auditors oftentimes compare tables in pdf documents to Excel spreadsheets



ONE swiss bank

SamurAI

AMLD EPFL

# This workshop: the goal

- ML can help you to automate processes with scanned paper documents
- This workshop presents few selected techniques to:
  - Classify documents
  - Extract information from documents
- Disclaimer:
  - This workshop doesn't present the fanciest / more powerful techniques to do these tasks
  - This workshop shows you few techniques that are easy to understand, implement and deploy as a beginner data scientist

Automated image-to-data process



ONE swiss bank

SamurAI
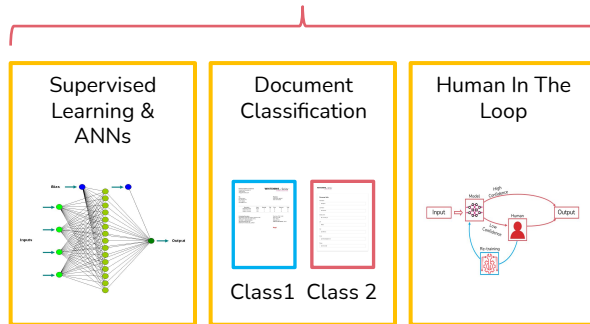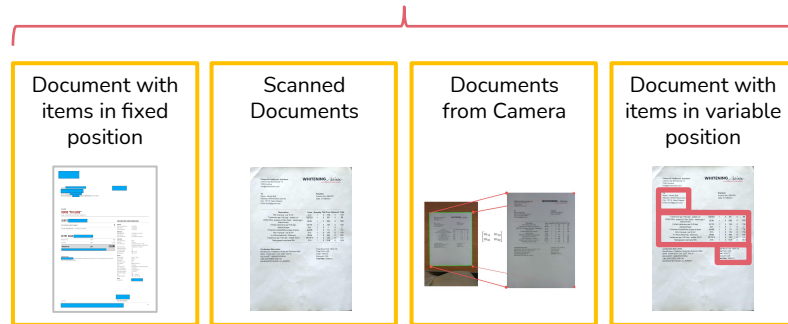
AMLD EPFL

# This workshop

- For beginners
- Get your hands dirty in code
- Adding levels of complexity

Part 1
(~1h30 )

Part 2
(~2h00)



| Supervised Learning & ANNs | Document Classification | Human In The Loop | Document with items in fixed position | Scanned Documents | Documents from Camera | Document with items in variable position |

Class1  Class 2

# Workshop: Schedule

| Session | Duration | Start - End | Subjects |
|---------|----------|-------------|----------|
| Part 1 | ~1:30 | 13:30 – 15:00 | Intro + Document Classification + HITL |
| Break | 0:30 | 15:00 – 15:30 | |
| Part 2 | ~2:00 | 15:30 – 17:30 | Information Extraction |

ONE
swiss bank

SamurΛI

AMLD EPFL

# Your questions are welcomed!

Your 1st question



=

Swiss chocolate bar



«I don't get the big picture»

«I didn't understand what this code does»
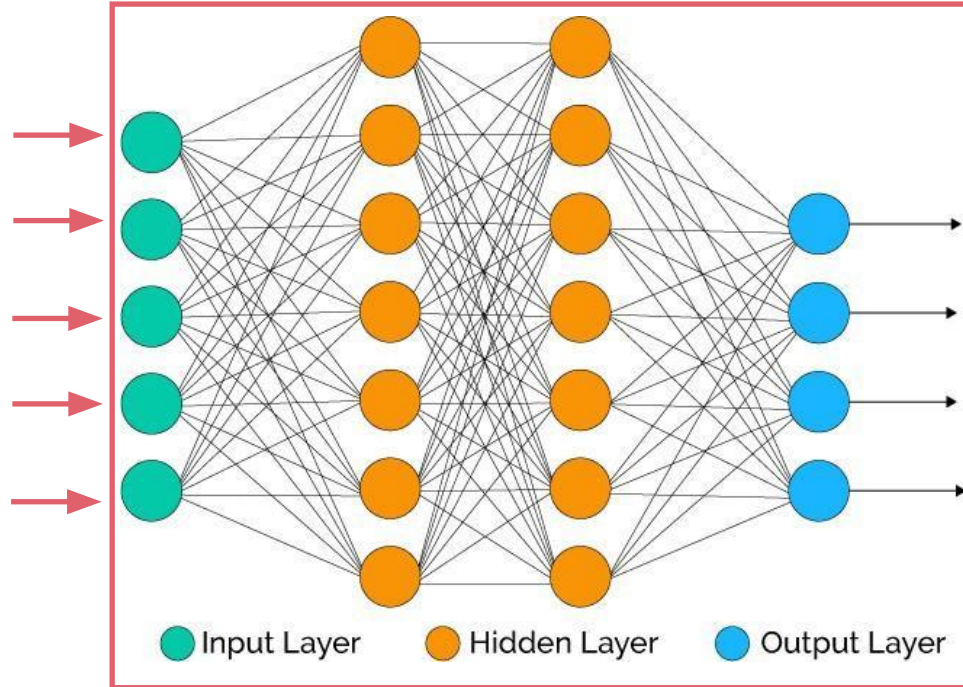
«Why are we doing this?»

ONE swiss bank

SamurAI

AMLD EPFL

# Supervised Learning



Attributes/ Features

Predictions

It's a document Invoice!

Input Layer   Hidden Layer   Output Layer

# Supervised Learning

# Supervised Learning



Attributes/ Features

Predictions   Labels

| invoice | 🟢 | invoice |
| registration | 🟢 | registration |
| registration | 🔴 | invoice |
| invoice | 🔴 | registration |
| registration | 🟢 | registration |
| invoice | 🟢 | invoice |
| invoice | 🔴 | registration |
| invoice | 🟢 | invoice |
| invoice | 🟢 | invoice |
| registration | 🟢 | registration |

🟢 Input Layer   🟠 Hidden Layer   🔵 Output Layer

Loss/CostFunction:    600
Accuracy:             70%

Values shown on network: 2.52, 0.32, -0.11, -0.51, 1.03, -1.53, -0.45, -0.38, 0.76, 0.50, -0.26, 2.18, 1.21, -0.51, -0.38, -0.12, 0.08, 0.17, -0.29

# Supervised Learning



Attributes/
Features

Predictions

Labels

| Predictions | | Labels |
|---|---|---|
| invoice | 🟢 | invoice |
| registration | 🟢 | registration |
| invoice | 🟢 | invoice |
| registration | 🟢 | registration |
| registration | 🟢 | registration |
| invoice | 🟢 | invoice |
| invoice | 🔴 | registration |
| invoice | 🟢 | invoice |
| invoice | 🟢 | invoice |
| registration | 🟢 | registration |

Loss/CostFunction:      300
Accuracy:                      90%

🟢 Input Layer    🟠 Hidden Layer    🔵 Output Layer

# Training as a loss minimization

- The loss quantifies the spread between labels and predictions
- The optimizers are algorithms that find the (possibly absolute) minimum of the loss

During the training the optimizer finds the path to minimize the loss, like a river flowing downhill

Loss

θ2

θ1

Batch gradient descent
Mini-batch gradient Descent
Stochastic gradient descent

Different optimizers can be quicker or slower

# Learning is hard

# ANN: Dense Layers

- Multi-Layer Perceptrons are the simplest ANN
- Every node of a layer is connected with all the nodes in the previous and in the following layer
- These layers are called Fully connected or Dense



- Image of 500 x 500 x 3 pixels
- First layer 100 nodes
- Already for the first layer we need ~750k parameters

The MLP is great, but too many parameters!

# ANN: Convolutional Layers

- We want to reduce the model parameters
- We introduce the concept of filter



Inputs     Filter     Outputs

# ANN: Convolutional Layers

- We want to reduce the model parameters
- We introduce the concept of filter
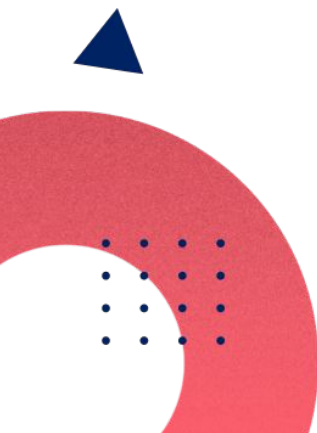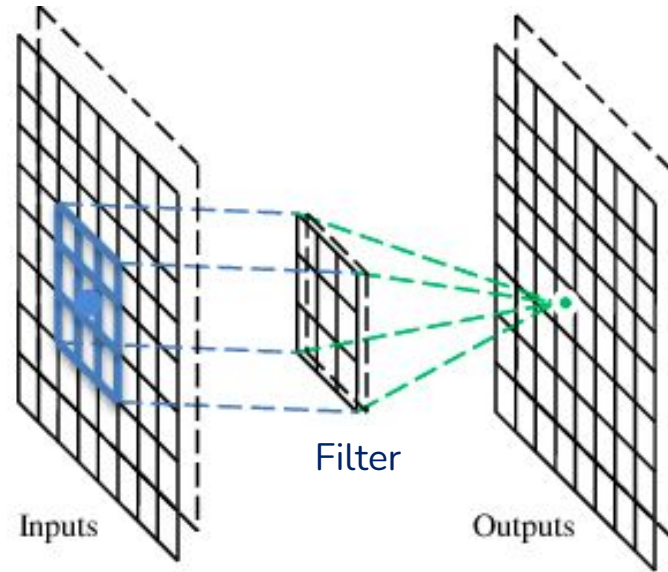
| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

Dot product

Filter

| 1 | -1 | -1 |
|---|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

=

3

Conv. layer parameters

# ANN: Convolutional Layers

- We want to reduce the model parameters
- We introduce the concept of filter

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

Dot product →

### Filter

| 1 | -1 | -1 |
|---|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

=

3    -1

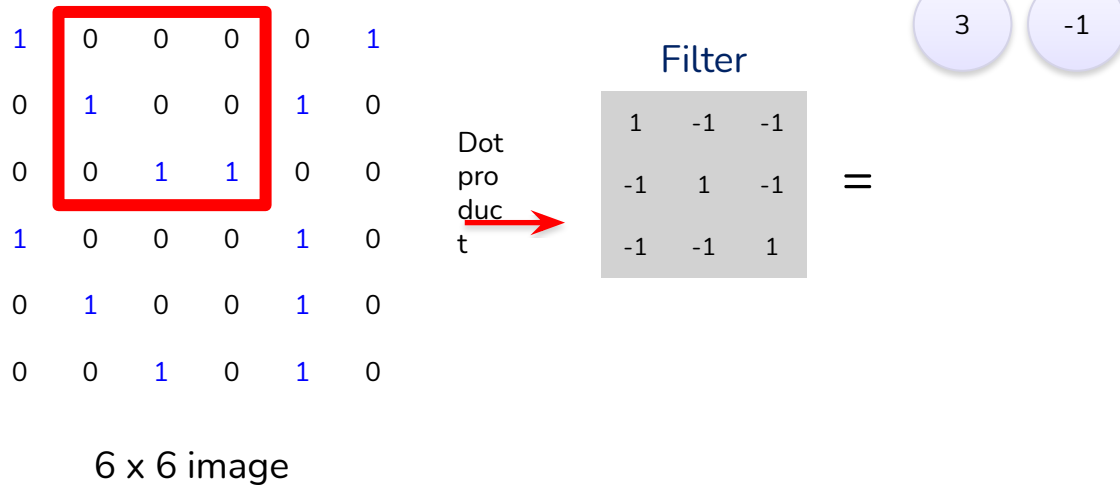# ANN: Convolutional Layers

- We want to reduce the model parameters
- We introduce the concept of filter

| 1 | 0 | 0 | 0 | 0 | 1 |
|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

Dot product

Filter

| 1 | -1 | -1 |
|---|----|----|
| -1 | 1 | -1 |
| -1 | -1 | 1 |

=

| 3 | -1 | -3 | -1 |
|---|----|----|----|
| -3 | 1 | 0 | -3 |
| -3 | -3 | 0 | 1 |
| 3 | -2 | -2 | -1 |

# ANN: Convolutional Layers

- We want to reduce the model parameters
- We introduce the concept of filter

|   |   |   |   |   |   |
|---|---|---|---|---|---|
| 1 | 0 | 0 | 0 | 0 | 1 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 |

6 x 6 image

Dot product

Filter

|   |   |   |
|---|---|---|
| 1 | -1 | -1 |
| -1 | 1 | -1 |
| -1 | -1 | 1 |

=

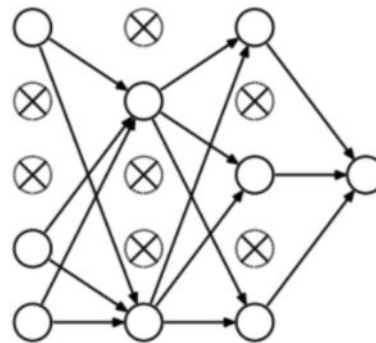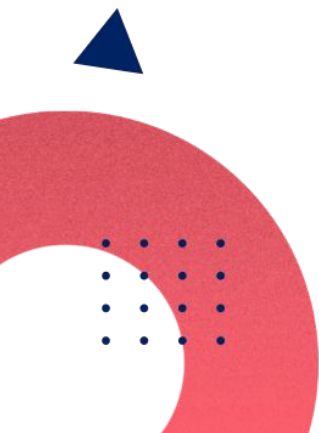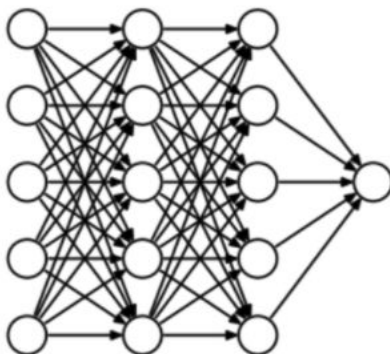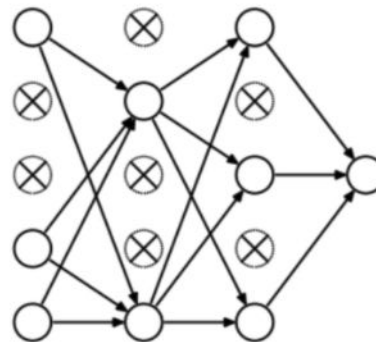|   |   |   |   |
|---|---|---|---|
| 3 | -1 | -3 | -1 |
| -3 | 1 | 0 | -3 |
| -3 | -3 | 0 | 1 |
| 3 | -2 | -2 | -1 |

# Dropout

- Regularization technique
- A fraction of the nodes are not considered in a training step
- This forces the network to have several "routes" in the nodes to ensure good performances

Standard Training

Training with Droupout

# Dropout

- Regularization technique
- A fraction of the nodes are not considered in a training step
- This forces the network to have several "routes" in the nodes to ensure good performances

```python
# Convolutional Layer 1
headmodel.add(Conv2D(8, (5, 5), padding='same', activation='relu'))
headmodel.add(Dropout(0.2))
```
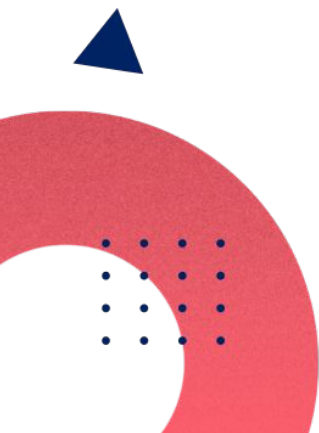
This means "drop 20% of the nodes in the previous layer while training
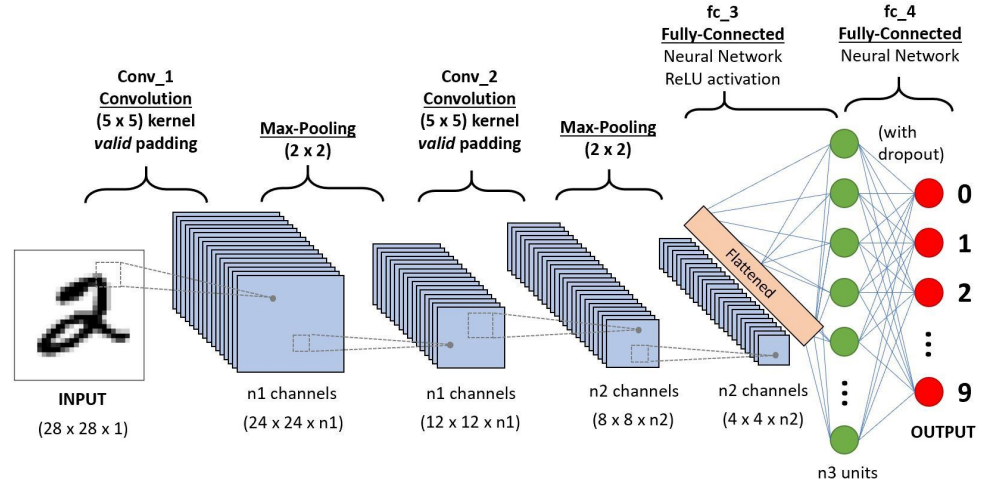


Standard Training

Training with Droupout

# Typical CNN structure

- Typical CNNs are made of :
  - A series of convolutional layer + dropout + max-pooling
  - Fully-connected layers at the end
- Many parameters are subject to tuning:
  - Number and type of layer
  - Number of filters
  - Size of filters



ONE swiss bank

SamurAI

AMLD EPFL

# CNN in Keras

```python
def define_model(num_classes,epochs):
    # Create the model
    model = Sequential()

    # Layer 1 (Convolutional)
    model.add(Conv2D(4, (5, 5), input_shape=(X.shape[1], X.shape[2], 1), padding='same', activation='relu',
    model.add(Dropout(0.2))
    model.add(MaxPooling2D(pool_size=(2, 2)))

    # Layer 2 (Convolutional)
    model.add(Conv2D(4, (3, 3), activation='relu', padding='same', kernel_constraint=maxnorm(3)))
    model.add(Dropout(0.2))
    model.add(MaxPooling2D(pool_size=(2, 2)))

    # Layer 3 (Convolutional)
    #model.add(Conv2D(4, (3, 3), activation='relu', padding='same', kernel_constraint=maxnorm(3)))
    #model.add(Dropout(0.2))
    #model.add(MaxPooling2D(pool_size=(2, 2)))

    # Additional Convolutional layers
    # ...

    # Additional Dense Layers
    model.add(Flatten())
    # model.add(Dense(6, activation='relu', kernel_constraint=maxnorm(3)))
    model.add(Dense(num_classes, activation='softmax'))
```
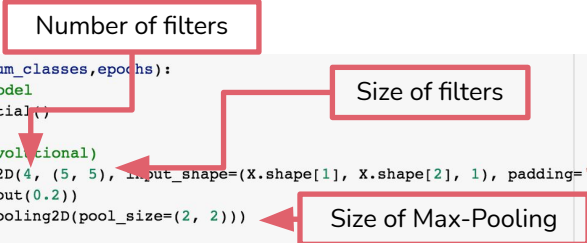
Conv. Layer +
Dropout +
Max-Pooling

Conv. Layer +
Dropout +
Max-Pooling

Conv. Layer +
Dropout +
Max-Pooling

Fully Connected Layers

Output

# CNN in Keras

Number of filters

Size of filters

Size of Max-Pooling

```python
def define_model(num_classes,epochs):
    # Create the model
    model = Sequential()

    # Layer 1 (Convolutional)
    model.add(Conv2D(4, (5, 5), input_shape=(X.shape[1], X.shape[2], 1), padding='same', activation='relu',
    model.add(Dropout(0.2))
    model.add(MaxPooling2D(pool_size=(2, 2)))

    # Layer 2 (Convolutional)
    model.add(Conv2D(4, (3, 3), activation='relu', padding='same', kernel_constraint=maxnorm(3)))
    model.add(Dropout(0.2))
    model.add(MaxPooling2D(pool_size=(2, 2)))

    # Layer 3 (Convolutional)
    #model.add(Conv2D(4, (3, 3), activation='relu', padding='same', kernel_constraint=maxnorm(3)))
    #model.add(Dropout(0.2))
    #model.add(MaxPooling2D(pool_size=(2, 2)))

    # Additional Convolutional layers
    # ...

    # Additional Dense Layers
    model.add(Flatten())
    # model.add(Dense(6, activation='relu', kernel_constraint=maxnorm(3)))
    model.add(Dense(num_classes, activation='softmax'))
```

# Document Classification

- We are going to train an algorithm to discriminate between different types of documents, <u>only using their images</u>
- The dataset:
  - ~400 invoices
  - ~400 registration
  - ~600 other
- We want the model to be robust enough to use images taken from a phone
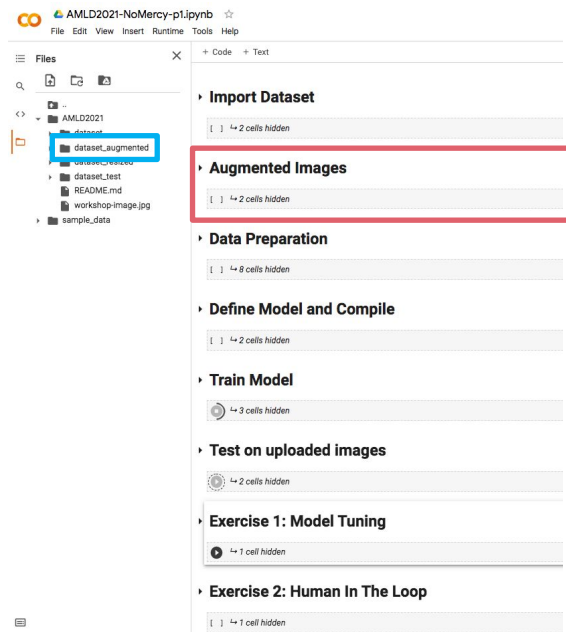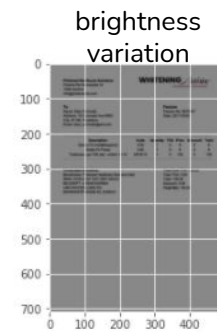
invoice

registration

other

# Notebook structure



- Import images from github
- On the left, you will see the folder AMLD2021 appear
- The main dataset used for this notebook is in AMLD2021/dataset_resized/
- ~1600images of 708x500 pixels

# Notebook structure



- More augmented images are created with distortions, brightness variations, tilt, …



(un)zoom      rotation      brightness variation

- Augmented images are moved to AMLD2021/dataset_augmented/

# Notebook structure



- Images are resized to a smaller size
- Dataset is created
- Train-test split

# Notebook structure



- Model is defined as a Keras ANN with Conv and Dense layers
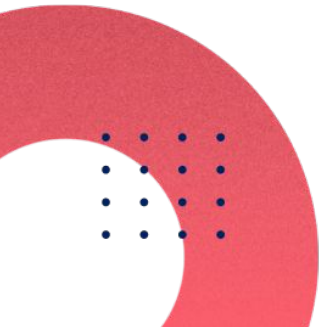
# Notebook structure



- Training and evaluation of performance

# Notebook structure



- You can take a picture of a document and test the model on it
- Upload your picture to Colab and in the folder AMLD2021/dataset_test/
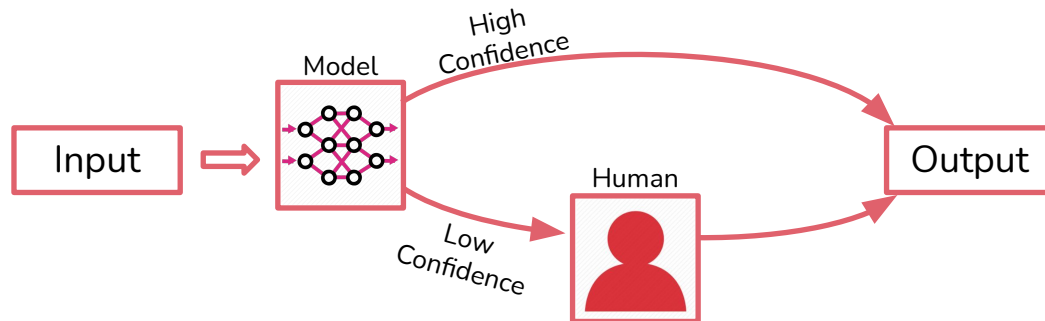- The code will make predictions on all pictures in this folder

# Exercise 1: Tuning the model

- Hyper-parameter tuning and feature engineering:
  - CNN layers, filters, dropout
  - Dense layers, number of nodes, dropout
  - Optimizer, Batch size, Num. of Epochs
  - Resize of the images (up in the data preparation)
- Goal: reach an accuracy of > 99% on the training and validation sets

# Human in the loop (HITL)

- AI systems are typically <100% accurate
- We can keep a HITL to mitigate for mistakes
  - Manually do the tasks for which the model has low confidence (low h)
  -

# Human in the loop (HITL)

- AI systems are typically <100% accurate
- We can keep a HITL to mitigate for mistakes
  - Manually do the tasks for which the model has low confidence (low h)
  - Prepare data to retrain the model

# Exercise 2: Human in the loop

- Make photos of several documents (~5-10 per class)
  - Jpeg format is perfect
  - Advice: name the images according to their class (ex: my_invoice_1.jpg). This will help for the rest of the exercise
- Upload them on Colab and copy them to the right folder for re-training (for example AMLD2021/dataset_resize/invoice/ for invoices)
  - example: !cp my_invoice_* AMLD2021/dataset_resize/invoice/
- Re-build augmented images, retrain the model, and re-evaluate performance

ONE
swiss bank

SamurAI

AMLD EPFL

# Possible Extensions: OCR

- In our classification we used only the images of our documents. We didn't use the text within the document

- Here an example of how to do document classification using only the text extracted with OCR: link

    - Using an LDA model: https://www.jmlr.org/papers/volume3/blei03a/blei03a.pdf



Image

OCR

Text

LDA model

invoice

registration

# Possible Extensions: OCR + Doc2Vec

- Document classification using OCR and transforming the text of a whole page into a vector: link