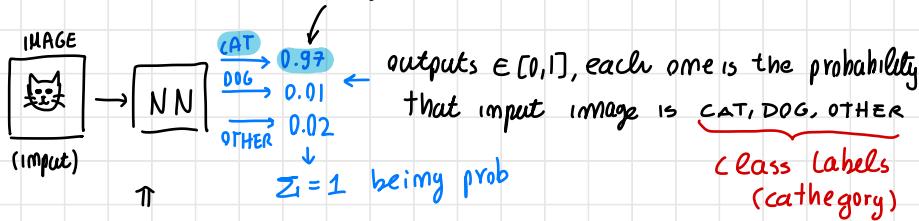


## GETTING STARTED WITH NN

→ **NN: What they are, and why they matter**

Simplify details to grasp most basic concept

- **NN as a black-box** highest PROB is class CAT → "The NET predict the image as CAT"



this is an image classification problem

where input output

image → numerical  
value & class  
(NOT labels itself)

the chosen predicted label is the one with highest probability

even if  $P_{CAT} = 0.51 \rightarrow$  still CAT,

$P_{DOG} = 0.48$  but less confident about prediction

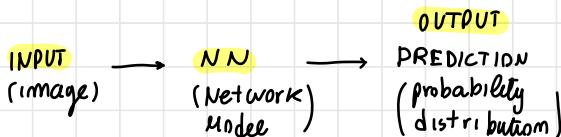
$P_{OTHER} = 0.01$

a perfect NN output  $P_{CAT} = 1$ ,  $P_{DOG} = 0$ ,  $P_{OTHER} = 0$  when INPUT = cat image

↳ in reality even well

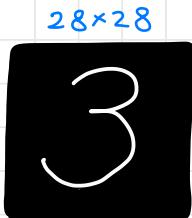
trained NETWORK doesn't give such result!

how all these elements on the predicting steps are represented?



**INPUT** ↗ How it can be represented as a **numeric value?**

**GRAY SCALE images** are represented as an array of pixel values where each pixel has intensity  $0 \div 255$   
(black  $\div$  white)



Gray scale  
image

→

$$\begin{bmatrix} [0\ 0\ 0\dots\ 1\ 12\ 0\ 0\dots 0], \\ [0\ \dots\ 50\ 255\ 255\ 0\dots 0], \\ \vdots \\ [0\dots 0\dots \dots \dots 0] \end{bmatrix}$$

$\uparrow$   
 $28 \times 28$  matrix  $M_{ij} \in [0, 255]$

Similarly color images has 3 components  $28 \times 28 \times 3$  matrix of RGB intensity levels!

$28 \times 28 \times 3$  → numbers represent a colored  $28 \times 28$  pixel image!

total amount  
of numbers as  
Input to the  
Network

HOW IT IS REPRESENTED?

⇒ It depends on the network..

- ✓  $\left\{ \begin{array}{l} \bullet 1 \text{ dimensional Vector } \text{input} \in \mathbb{R}^{28 \cdot 28 \cdot 3} \\ \bullet \text{three 2D arrays each } 28 \times 28 \text{ input} \in \mathbb{R}^{28 \cdot 28 \times 3} \end{array} \right.$

the Network design expect a fixed size and shape for data

When the image is NOT coherent with the expected input...

↪ we resize/crop to the expected size

IF GRAYSCALE we can replicate it for all 3 channels (when RGB expected)  
or opposite when grayscale is expected RGB is converted

**PRE - PROCESSING step to create suitable input**

NN are designed to accept a certain shape / size INPUT

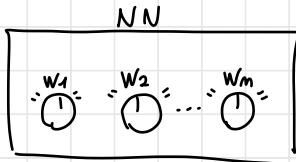


different image classification networks use different INPUT based on the application they are designed for!

"TRAIN a NEURAL NETWORK"



$\boxed{\text{NN}}$



contains many tunable

parameters ~ KNOBS inside NN

Weights



← Network weights w<sub>1</sub>, w<sub>2</sub>, ..., w<sub>m</sub>

essential to use the NN...

you find the correct Weights by TRAINING the NN



showing thousands of examples of classes you want to learn.

SUPERVISED LEARNING

↳ you provide the NN with image of a class and you explicitly tell the result



IF the network model makes a wrong prediction, we compute the error related to the bad prediction, and use it to adjust network weights → so that subsequent predictions accuracy are improved

↳ details on the training step are discussed next..

## FUNDAMENTALS OF TRAINING A NEURAL NETWORK

most essential elements required for TRAINING a NN

example specific for an IMAGE CLASSIFICATION

- Label training data
- Loss function
- Updating Weights



the TRAINING

requires

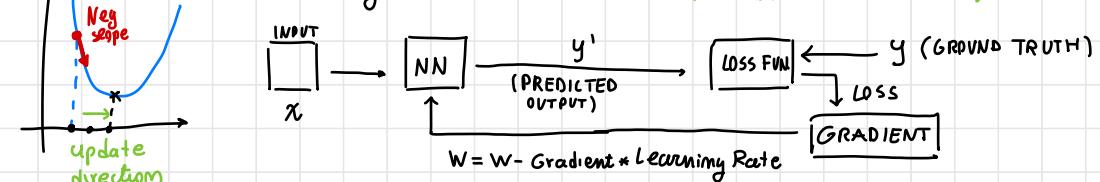
### NETWORK TRAINING

1) Labeled training data = thousands of images  $\forall$  class with expected output (Label)

2) Loss Function  $\approx$  cost Function: numerical computation that quantify error between network output and expected result

here we rely on MEAN SQUARED ERROR

3) Process for tuning the NN's weights based on the value of the LOSS FUNCTION (GRADIENT DESCENT)



1 and 2: Representation of training data and loss function computation

### 1. TRAINING DATA:

Input image ( $x$ ) (numerical matrix)



"One-Hot Encoding"  $\rightarrow$

is recommended

when classes are

NOT related..

representing categorical labels as BINARY VECTORS

$$y = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad \text{class labels (numerically)}$$

"CAT" label

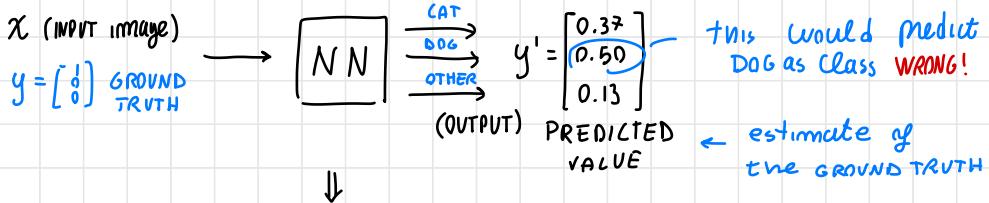
$$\text{"DOG" would be } y = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

One-Hot Encoding

$\leftarrow$  Used to represent class labels

$\nwarrow$  this known labels are called

"GROUND TRUTH" := y



to quantify the ERROR between the Network output (PREDICTION) and the RESULT EXPECTED is to compute the sum of squared error

$$y \in \mathbb{R}^N, y' \in \mathbb{R}^N \quad SSE = \sum_{j=1}^n (y[j] - y'[j])^2$$

error for a single training sample using prediction and ground truth

When training a NN MANY (multiple) images are used to compute the loss before updating Network weights

we use Mean Squared Error for multiple ( $m$ ) training images referred as BATCH SIZE (32 typical)

$$\underline{MSE} = \text{mean}\{SSE\} = \frac{1}{m} \sum_{i=1}^m (y_i - y'_i)^2$$

each batch of image is defined an ITERATION

other loss functions can be used requiring a numerical representation for input class and prediction ( $y$  and  $y'$ )

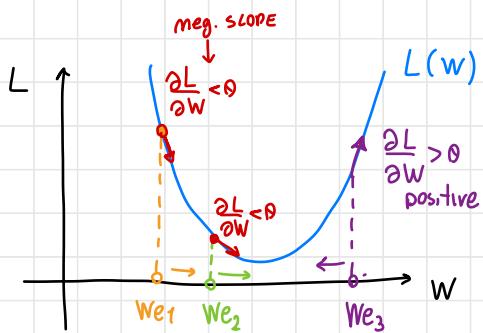
representation of error between expected output and input to network



3. How can we use this loss information to train the network?  
 this TUNING of weights of NN can be achieved by GRADIENT DESCENT Algorithm

## GRADIENT DESCENT (Gradient ~ slope)

example of this Algorithm with only  $w$  (one parameter) and loss  $L$  convex



⇒ the GRADIENT sign  $\frac{\partial L}{\partial w}$  is what matters!

↓  
we want  $L$  to decrease  
so based on the fact that

$$\text{Slope} = \frac{\text{rise}}{\text{run}} < 0 \quad \begin{matrix} \text{When} \\ \text{moving} \\ \text{towards} \\ \text{minimum!} \end{matrix}$$

↓  
we need to adjust  
the weights in the direction  
toward decreasing loss

IF  $\frac{\partial L}{\partial w} > 0$  we move left ←

IF  $\frac{\partial L}{\partial w} < 0$  we move right →

↓

We need to adjust the weights in the direction opposite to the sign of the gradient

$$We_{t+1} = We_t - \underbrace{\text{GRADIENT}}_{\text{the sign of gradient}} \times \underbrace{\text{LEARNING Rate}}_{\text{determine the direction of motion}}$$

the amount of motion (step size)  
depends on the Learning Rate

parameter  $< 1$  it is a parameter we  
initialize (NOT determined by the net)  
that determines step  
this is an **HYPER PARAMETER**

↑  
because it is different from trainable  
parameters such as Net Weights

equation to update  
← the weight in the  
proper direction  
regardless of current  
W relative to optimal W\*

In practice  $L(w_1 \dots w_n)$  has many dimensions, NOT convex, many pick/valley

Its slope is called GRADIENT and it is function of all Network Weights

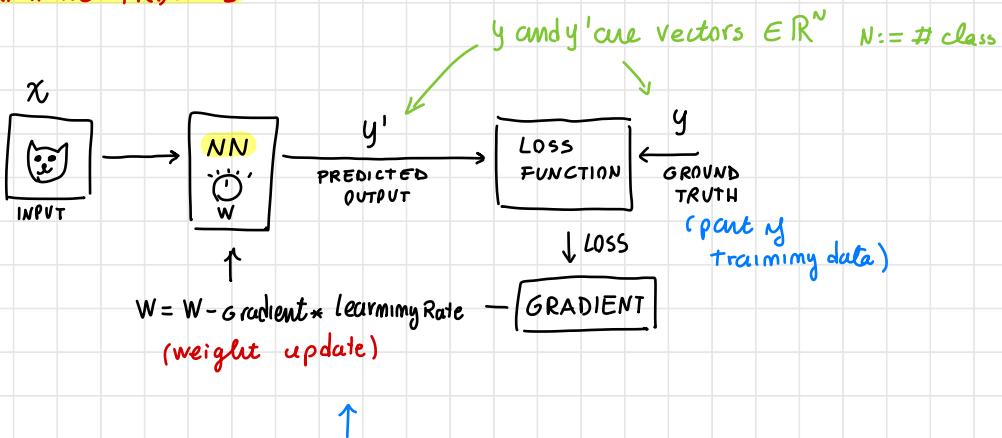
↳ so GRADIENT is a vector in multi dimensional space...

(but the previous approach is used)

this GRADIENT of LOSS FUNCTION with respect to network weights is computed with an Algorithm BACKPROPAGATION

(built in in DL Frameworks as TensorFlow OR Pytorch)

## NN TRAINING PROCESS

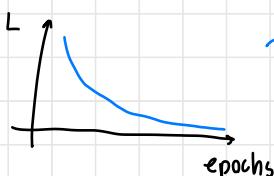


Training a NN is an ITERATIVE PROCESS, often requires passing entire TRAINING set through the NET Multiple times!

↳ each time the TRAINING SET is passed through the Net

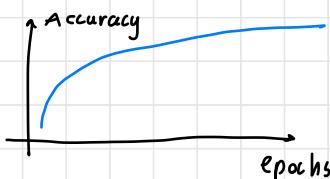
↳ Is define as TRAINING EPOCH

Training a NN often Requires many training epochs until Loss decrease satisfactory



~ training progresses until loss Function stabilize and stop decreasing

Similarly ACCURACY is evaluated



⇒ many advanced topics on NN training are neglected here... such as DATA SPLITTING

What you do with your TRAINED NN?

You can then feed the NN with new images  $x$  and predict  $y'$  on NON CLASSIFIED images

using the NETWORK to perform inference  
without a label training data

- Remember that after INPUT must be PRE-PROCESSED to be coherent with expected NN input
- NN weights are initialized to small random values
- BATCH, ITERATIONS and EPOCHS

TRAINING SET

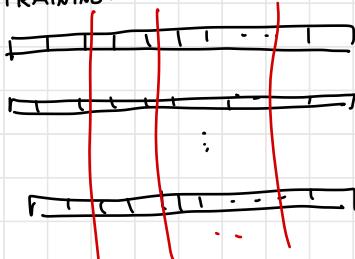


$M$  Labeled Training data

$B :=$  batch size → ITERATION

Iterations & epoch e

↳ TRAINING:



$$\# \text{ Iterations} = I = M / B$$

# LINEAR REGRESSION (LR)

how simple NN can be used to solve LR problem using Tensorflow and Keras

↓ it make easy to load the dataset, define NN and train it

LR, fit a straight line through a set of data points

the power of NN comes from the ability to model highly non linear functions → here we use `cmean` function to define basic terms

→ set up a Random Number SEED so that results are repeatable

1) Load dataset: boston housing dataset available in Keras (simple example)

↓

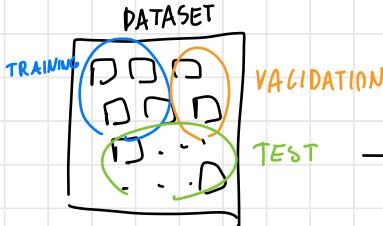
we get training and test dataset with input features and target variables  
(se)

$\begin{cases} x: \text{input} \\ y: \text{output target} \leftarrow \text{here it is a number} \end{cases}$

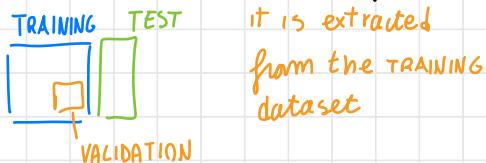
this is the medium price of the house for a given features

→ explore the dataset  
.shape → (404, 13)  
↑  
Training samples ↗ attributes  
" FEATURES "

(\*) DATA SPLITTING Data set  $\begin{cases} \text{training} \\ \text{validation} \\ \text{test} \end{cases}$  ⇒ help prevent overfitting  
(when model doesn't generalize unseen data)



→ IF Validation dataset NOT defined,



it is extracted  
from the TRAINING  
dataset

DATA SPLITTING IS fundamental in ML!

HERE we use single features to keep simple + plot the results

BUT it is important in general that many features are used to predict a single target value

→ 2) EXTRACT the 6<sup>th</sup> feature only from dataset  
(here the highly correlated)



good practice to visualize dataset, plot a portion of it

+ to check if a reasonable feature has been chosen for prediction

Inspect data to be aware of Anomalies



we try to obtain  
the model

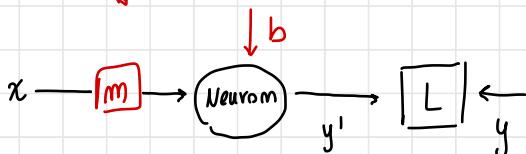
- y-intercept:  $b$
- slope :  $m$

↑  
best values of  $b, m$   
that defines best FIT lines  
through this data

How to use NN to create such a model

In this problem

$(m, b)$  are the weights of NN model

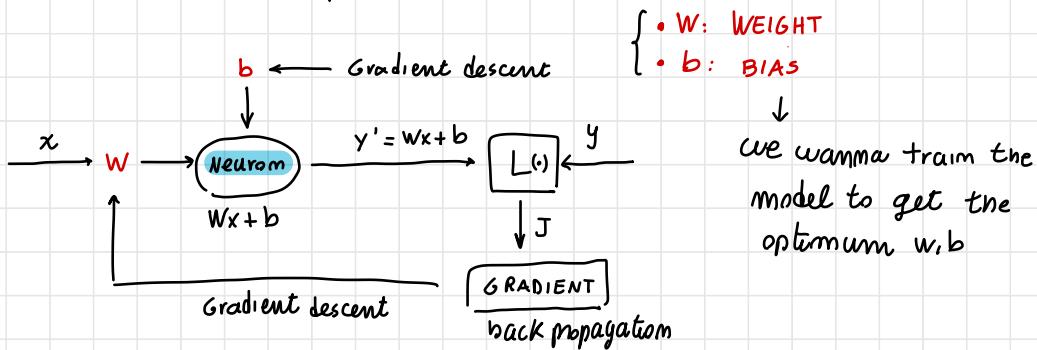


In general a NN has  
many Neurons arranged  
in Layers...

↓  
so DNN has many layers  
and many Neurons per layer

here we have  
one Neuron  
in a single layer

here we have two tunable parameters



In complex networks

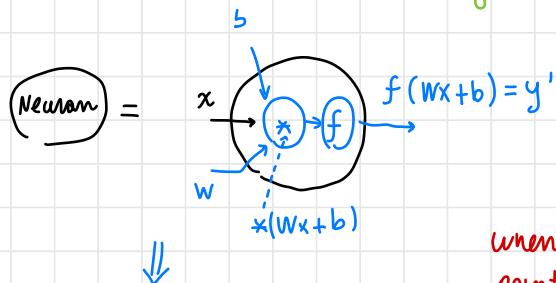
with many  $w$  and  $b$ , the weight  $w$ s are  $\sim$  function slope  
 $b$ s are  $\sim$  bias offset

• Neuron

"

it's a computation unit (PERCEPTRON)  $\rightarrow$  computing  $Wx + b$

then  $Wx + b$   $\leadsto$  is passed to an ACTIVATION FUNCTION  $f$  to produce the output of the neuron



in a LR We don't need an activation function  
 $f(\tilde{x}) = \tilde{x}$

when  $f(\tilde{x})$  is specified,  $\tilde{x}$  is floating point number and  $y' = f(\tilde{x})$  is also a floating point number

Network like this are

SINGLE-LAYER PERCEPTRON (SLP)

(here even special case, because it has just one Neuron)

Network with multiple layers are MULTI-LAYER PERCEPTRON (MLP)

Input layer      Hidden layer      Output layer



In our special SLP, we use  $y'$  and  $y$  to compute MSE (loss)  
 then GRADIENT to update weights ← (using GRADIENT DESCENT)  
 ↑  
 (managed by KERAS with BACKPROPAGATION) ← ↑ hidden in keras

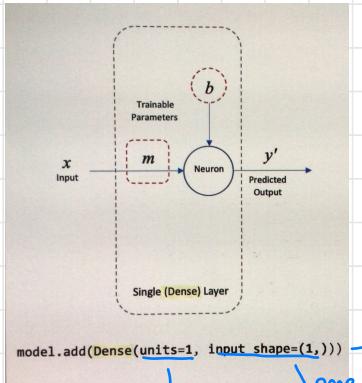
$$J = \frac{1}{m} \sum_{i=1}^m (y'_i - y_i)^2$$

- code implementation:

- { 1) Define and compile the model
- 2) Train model with fit
- 3) use predict to make predictions

(1) **Instantiate** the model = `Sequential()`

then define it by `model.add(...)`  
 ↑ with a single neuron



Dense (Layer): every input is connected to every output  
 (here one neuron)

the output is  
 ↗ by definition  
 a single number

then `model.summary()`  
 summarize the net architecture

- compile** the model : prepare it for training ...

`model.compile(...)`

↑ specifies the type of optimizer used to perform gradient descent  
 and the type of loss function

② TRAIN the model with `fit(.)` method

It execute the training process := iterative cycle that updates the Network weights

```
fit(x_train, y_train, batch_size, epochs, validation_split)
```

↑  
↑  
features  
training  
data      target  
value  
(ground  
truth)

↑

we use a portion of training data to update ...

+ number of epochs we want to train model...

validation  
dataset, instruct  
to hold a  
percentage  
of data  
used for  
validation

batch size is  $cm$   $\leftarrow$  before update

HYPERTPARAMETER because  $[16 \div 256]$  depend  
it is melaninarily specified, but it is  
important and can affect the model  
accuracy

11

this output the TRAINING EPOCHS for each line...

each epoch has # iterations = train # / batch size if update  
also train and validation

loss are returned by `fit()` method

history = model.fit(...)

## Validationism

this history object allows to plot training results

history.history['...'] dictionaries containing training/validation loss

reasonably  $J_{TRAIN} < J_{VAL}$  because validation is done on data never seen before

⑨ we can use the model .predict(x) to make prediction over a list of x values



We can plot the best fit line overlayed on training data to check accuracy of our LR model



how to evaluate the accuracy? → Compute meaningful metrics on test data to determine if predictions is within acceptable limit for example the

When evaluating a model, use a metric easy to interpret!



RMS (Root Mean Squared)

MAE (Mean Absolute Error)

from best fit on test dataset  
(MSE is harder to interpret because it involve squared of price...)

the LOSS FUNCTION used for optimization is NOT necessarily the same of performance metric!



this metric are used for different purpose!

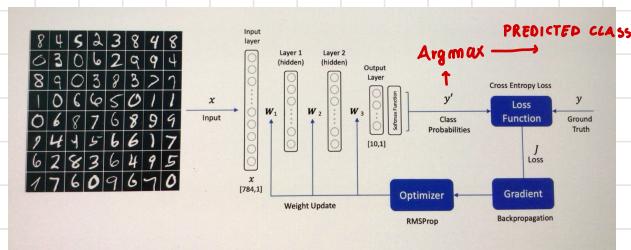
from the Validation - prediction plot we can see big errors ... than maybe we should consider using additional features to improve predictions

With 2 input features the model became a plane while hyperplane with more than 2 → still LINEAR MODEL

otherwise, we can consider to use non linear model!

# MNIST DIGIT CLASSIFICATION USING MLP

How to implement a Feed Forward NN in Keras to classify hand written digits from MNIST Dataset



concept in the context of  
image classification involving  
more than two classes

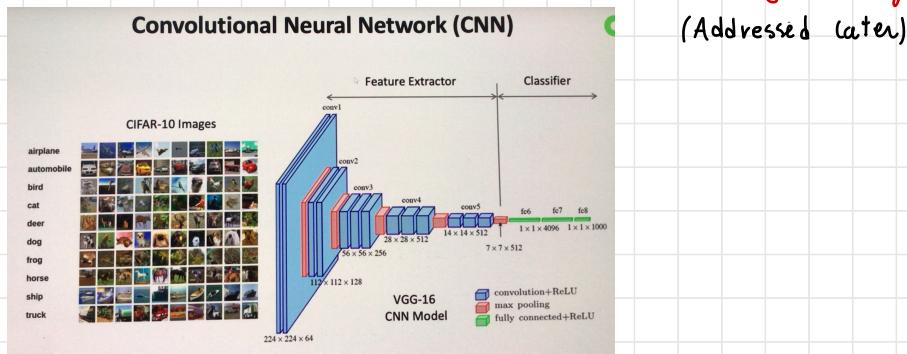


multinomial Regression  
(soft max Regression) → more  
than 2 classes

In general, MLP are NOT the best way to process

image data BUT it is an example before covering advanced Network

CONVOLUTIONAL NN (CNN) better suited for image processing



MNIST dataset is included in Keras → easy to load



70000 images: partitioned { 60000 training → We reserve a portion for validation (10000)  
10000 testing }



28x28 pixels

Images (GRAYSCALE)

{ 50000 training  
10000 validation

1) LOAD and SPLIT Dataset ↑

2) PREPARE dataset (PRE-PROCESS) for processing through Network  
by transforming

↳ here, the single input features (IMAGES)

logically it is represented by pixel intensity as features

↓

flatten the 2D image ( $28 \times 28$ ) data to 1D vector

Reshaping sample images on each set

+ Normalize pixel intensity by 255 to be  $I_{xy} \in [0, 1]$

(common when working with img data → help model to train efficiently)

How Image Labels are represented? In categorical Data the target labels (STRING) must be converted to numerical values

STRING → Numerical : "Label Encoding"  
values

{ Option 1: ORDINAL INTEGER ENCODING  
(number int & class)  
Option 2: ONE-HOT ENCODING  
(separate binary vector  
to encode each class label)

Option 1: also mapping from  
class name to int representation is  
reported in a config file of dataset  
→ unique integer are used to encode class labels  
when class labels are NOT related one another,

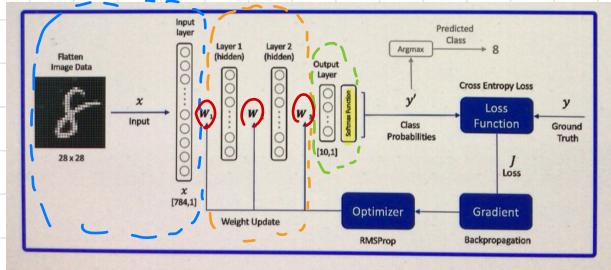
↓

Option 2 pre-processing, each labels → converted in N dimensional  
vector (with N as # classes)

↓

Transform easily to one-hot encoding by to\_categorical()

### 3) DEFINE Network architecture



#### 3.1 Deep Neural Network Architecture

The network architecture shown below has multiple layers. An input layer, two hidden layers, and an output layer. There are several things to note about this architecture.

1. Input Data: The image input data is pre-processed (flattened) from a 2-Dimensional array [28x28] to 1-Dimensional vector of length [784x1] where the elements in this input vector are the normalized pixel intensities. The input to the network is sometimes referred to as the input "layer", but it's not technically a layer in the network because there are no trainable parameters associated with it.
2. Hidden Layers: We have two hidden layers that contain some number of neurons (that we need to specify). Each of the neurons in these layers has a non-linear activation function (e.g., ReLU, Sigmoid, etc.).
3. Output Layer: We now have ten neurons in the output layer to represent the ten different classes (digits: 0 to 9), instead of a single neuron as in the regression example.
4. Dense Layers: All the layers in the network are fully connected, meaning that each neuron in a given layer is fully connected (or dense) to each of the neurons in the previous layer. The weights associated with each layer are represented in bold to indicate that these are matrices that contain each of the weights for all the connections between adjacent layers in the network.
5. Softmax Function: The values from each of the neurons in the output layer are passed through a softmax function to produce a probability score for each of the ten digits in the dataset.
6. Network Output: The network output ( $y'$ ) is a vector of length ten, that contains the probabilities of each output neuron. Predicting the class label simply requires passing ( $y'$ ) through the argmax function to determine the index of the predicted label.
7. Loss Function: The loss function used is Cross Entropy Loss, which is generally the preferred loss function for classification problems. It is computed from the ground truth labels ( $y$ ) and the output probabilities of the network ( $y'$ ). Note that  $y$  and  $y'$  are both vectors whose length is equal to the number of classes.

Although the diagram looks quite a bit different from the single-layer perceptron in the linear regression example, it is fundamentally very similar in terms of the processing that takes place during training and prediction. We still compute a loss based on the predicted output of the network and the ground truth label of the inputs. Backpropagation is used to compute the gradient of the loss with respect to the weights in the network. An optimizer (which implements gradient descent) is used to update the weights in the neural network.

the Hidden layers contains neurons to specify:



each neuron has a Non Linear Activation function such as ReLU / sigmoid

the Output layer has N neurons (as # classes)

↑  
Here all layers are fully connected (DENSE) because each Neuron in next layer is fully connected to previous layer

the weights  $W_i$  are matrix ( $\forall$  layer) containing the weight for all connection between adjacent layers

Softmax function is used to the values of all output layers neuron ...

↑  
to produce a probability score for each of the 10 digits in dataset  
≈ it is a way of normalizing a set of inputs

Multiple layers:

- 1 INPUT LAYER
- 2 HIDDEN LAYERS
- 1 OUTPUT LAYER

- the INPUT image is transformed from 28x28 2D vector to 1D 784 Vector with pixel intensity
- ↑  
INPUT LAYER, even y  
there are NO trainable parameters associated  
With that it is NOT a real NETWORK (layer)

Softmax: Output layer pass through softmax function to normalize them into probabilities

then  $y' = \begin{bmatrix} y'_1 \\ y'_2 \\ \vdots \\ y'_n \end{bmatrix}$  each  $y'_j$  contain the probability of each output neuron (class j)

then the argmax is used to predict the label

As LOSS FUNCTION a CROSS-ENTROPY loss is used for classification problems  $\rightarrow$  computed from  $y, y'$  (vectors of length = # classes)

↑

even if different to SLP, in MLP the processing to update the weights is the same!

#### 4) MODEL IMPLEMENTATION (in Keras)

↓

Instantiate as Sequential() then use .add(...) to build the architecture by 2 hidden layers + 1 output layers

'relu'

↑

the first hidden layer require the input shape + specify # Neurons as design choice  
each neuron has its own activation function!)

'softmax'

↓

here in output layers, the activation function Softmax is used to pass all the outputs through it!

~ 118.282 trainable parameters SMALL NET!  
usually millions of parameters

## 5) COMPILE MODEL

specify {

- type of optimizer (rmsprop)
- loss function (categorical-crossentropy)
- additional metrics to track during training process [accuracy] training/validation loss automatically saved

specify the label encoding  
if integer encoded you use sparse-crossentropy

$L := \text{cross entropy}$  → simple computation that allow to quantify error between output probabilities and one-hot encoded vector & class

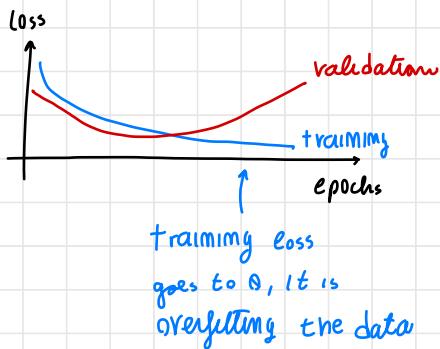
## 6) TRAIN the MODEL by `fit(...)` method

train dataset, epochs, batch-size, validation-data  
return training-results history of training data

specify explicitly with separate validation dataset

## 7) ANALYZE training Results

after first epochs... while training loss decrease, validation loss can increase when OVERTFITTING is occurring!



↑  
Common problem in machine learning, several techniques can mitigate overfitting issues...

↓  
still accuracy ~ 98% is still very good... on average 1 mistake every 50 samples!

## 8) MODEL EVALUATION

1. We can use the trained model to obtain prediction for any number of samples



containing the probability for each of the  $N$  (# class) neurons in the output layer

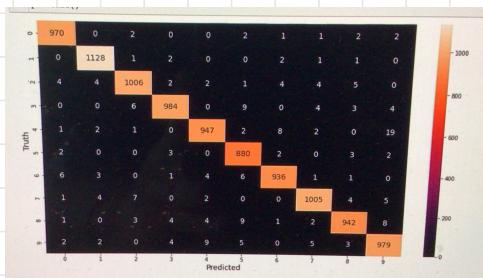
we can apply predict to all the test data



an efficient way (INFORMATIVE) to evaluate the model on the entire Test set → is by CONFUSION MATRIX



common metric used to summarize results of a classification problem



↳ this is presented as a table (matrix) where one axis is ground truth of each class while other axis is the prediction from model

↳ total # of entries are the # of instances from experiment

to evaluate it,

we extract for

all test dataset the argmax prediction

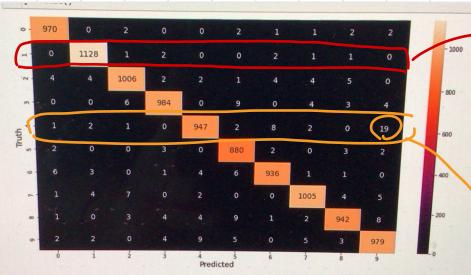
since labels are one-hot encoded → we convert the prediction back to integer

then the confusion matrix is generated easily passing list of ground truth labels and predicted



then we can plot this as heatmap

(or also provided as simple table of numbers or percentage)



the row 1 is the digit 1 ground truth..  
here we get that we have  
1128 correct prediction  
and 7 samples wrong

WORST classification  
because 4 ~ 9 similar shape

↓  
Overall accuracy of 98% on validation...

from confusion matrix we can learn which  
class is more difficult to predict → INSIGHT!



NN involves millions parameters in non-visualizable dimensions...  
modeling complex non linear relationship



focus on understanding the high level components  
of a Neural Network!

- {
- GRADIENT DESCENT
  - LOSS FUNCTION
  - BACKPROPAGATION
  - NN ARCHITECTURE
- ....

well established  
algorithms  
embedded in  
all modern DNN  
frameworks