DATA IN (ASTRO) PARTICLE PHYSICS and COSMOLOGY: the BRIDGE to INDUSTRY **SCIENCE**

Machine learning: Neural Networks

Márcio Ferreira

Centro de Física da Universidade de Coimbra

27 - 30 June (2022)

Contact: marcio.ferreira@uc.pt

Márcio Ferreira (CFisUC)

Machine learning: Neural Networks

Contents

- Perceptron
- Neural Networks
- Deep Neural Networks
- How to train DNN
 - Backpropagation algorithm
- Classification task with DNN

Biological Inspiration

• Biological neural network



Artificial neuron

- Neural network as computing machines (McCulloch-Pitts model)
- Perceptrons are based on Linear Threshold Units (LTU)



• A single LTU can be used for simple linear binary classification

Perceptron

- A Perteptron is made of a single layer of LTUs
- The first type of artificial neural network (ANN)



• Multi-class **linear** classifier: the decision boundary of each neuron is linear

Márcio Ferreira (CFisUC)

Perceptron

• Single-layer perceptrons only learn linearly separable patterns

Binary classification: Multi-class classification: x_2 x_2 x_1 x_1 x_1

Feed-forward Deep Neural Networks

• Multiple hidden layers increase the representational power of NN



The math of Neural Networks



[https://doi.org/10.1190/tle37080616.1]

Márcio Ferreira (CFisUC)	Machine learning: Neural Networks	27 - 30 June (2022)	8 / 38

The math of Neural Networks

Hidden layer:
$$\mathbf{a}_1^{(i)} = \sigma\left(\mathbf{\Sigma}_1^{(i)}\right)$$
, where $\mathbf{\Sigma}_1^{(i)} = \mathbf{W}\mathbf{x}^{(i)} + \mathbf{b}_1$

$$\begin{bmatrix} a_{11}^{(i)} \\ a_{12}^{(i)} \\ a_{13}^{(i)} \\ a_{14}^{(i)} \\ a_{15}^{(i)} \end{bmatrix} = \boldsymbol{\sigma} \left(\begin{bmatrix} \Sigma_{11}^{(i)} \\ \Sigma_{12}^{(i)} \\ \Sigma_{13}^{(i)} \\ \Sigma_{14}^{(i)} \\ \Sigma_{15}^{(i)} \end{bmatrix} \right) = \boldsymbol{\sigma} \left(\begin{bmatrix} W_{11} & W_{12} & W_{13} \\ W_{22} & W_{22} & W_{23} \\ W_{33} & W_{32} & W_{33} \\ W_{44} & W_{42} & W_{43} \\ W_{55} & W_{52} & W_{53} \end{bmatrix} \begin{bmatrix} x_{1}^{(i)} \\ x_{2}^{(i)} \\ x_{3}^{(i)} \end{bmatrix} + \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \\ b_{14} \\ b_{15} \end{bmatrix} \right)$$

Output layer:
$$\mathbf{\hat{y}}^{(i)} = \phi\left(\Sigma_2^{(i)}\right)$$
, where $\Sigma_2^{(i)} = \mathbf{Da}_1^{(i)} + \mathbf{b}_2$

$$\begin{bmatrix} \hat{y}_{1}^{(i)} \end{bmatrix} = \phi \left(\begin{bmatrix} \Sigma_{21}^{(i)} \end{bmatrix} \right) = \phi \left(\begin{bmatrix} D_{11} & D_{12} & D_{13} & D_{15} & D_{15} \end{bmatrix} \begin{bmatrix} a_{11}^{(i)} \\ a_{12}^{(i)} \\ a_{13}^{(i)} \\ a_{14}^{(i)} \\ a_{15}^{(i)} \end{bmatrix} + \begin{bmatrix} b_{21} \end{bmatrix} \right)$$

(:) =

The math of Neural Networks

Normally, instead of single points, mini-batches of N points are fed into NN

$$\begin{array}{cccc} \mathbf{Mini-batch:} & [\mathbf{x}^{(1)} & \dots & \mathbf{x}^{(N)}] = \begin{bmatrix} \mathbf{x}_{1}^{(1)} & \dots & \mathbf{x}_{1}^{(N)} \\ \mathbf{x}_{2}^{(1)} & \dots & \mathbf{x}_{2}^{(N)} \\ \mathbf{x}_{3}^{(1)} & \dots & \mathbf{x}_{3}^{(N)} \end{bmatrix} \\ \begin{bmatrix} \mathbf{a}_{11}^{(1)} & \dots & \mathbf{a}_{12}^{(N)} \\ \mathbf{a}_{12}^{(1)} & \dots & \mathbf{a}_{13}^{(N)} \\ \mathbf{a}_{13}^{(1)} & \dots & \mathbf{a}_{14}^{(N)} \\ \mathbf{a}_{15}^{(1)} & \dots & \mathbf{a}_{15}^{(N)} \end{bmatrix} = \sigma \begin{pmatrix} \begin{bmatrix} W_{11} & W_{12} & W_{13} \\ W_{22} & W_{22} & W_{23} \\ W_{33} & W_{32} & W_{33} \\ W_{44} & W_{42} & W_{43} \\ W_{55} & W_{52} & W_{53} \end{bmatrix} \begin{bmatrix} \mathbf{x}_{1}^{(1)} & \dots & \mathbf{x}_{1}^{(N)} \\ \mathbf{x}_{2}^{(1)} & \dots & \mathbf{x}_{3}^{(N)} \end{bmatrix} + \begin{bmatrix} b_{11} & \dots & b_{11} \\ b_{12} & \dots & b_{12} \\ b_{13} & \dots & b_{13} \\ b_{14} & \dots & b_{14} \\ b_{15} & \dots & b_{15} \end{bmatrix} \end{pmatrix} \\ \hat{\mathbf{y}}^{(1)} & \dots & \hat{\mathbf{y}}^{(N)} \end{bmatrix} = \phi \begin{pmatrix} \begin{bmatrix} D_{11} & D_{12} & D_{13} & D_{15} & D_{15} \end{bmatrix} \begin{bmatrix} \mathbf{a}_{11}^{(1)} & \dots & \mathbf{a}_{16}^{(N)} \\ \mathbf{a}_{12}^{(1)} & \dots & \mathbf{a}_{16}^{(N)} \\ \mathbf{a}_{13}^{(1)} & \dots & \mathbf{a}_{16}^{(N)} \\ \mathbf{a}_{15}^{(1)} & \dots & \mathbf{a}_{16}^{(N)} \end{bmatrix} + \begin{bmatrix} b_{21} & \dots & b_{21} \end{bmatrix} \end{pmatrix} \end{pmatrix}$$

The output of a mini-batch is obtained by matrix operations

Training a NN

• Training the model consists in adjusting the network' weights to minimize a cost function

$$\mathbf{\hat{w}} = \operatorname*{arg\,min}_{\mathbf{w}} \mathcal{C}(\mathbf{x}, \mathbf{y}, \mathbf{w})$$

• Moving in the direction of negative and large $\nabla C(\mathbf{x}, \mathbf{y}, \mathbf{w})$ (gradient descent): initialize \mathbf{w} to some value \mathbf{w}_0 and update its value according to

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \nabla_{\mathbf{w}} \mathcal{C}(\mathbf{w}_t)$$

• To speed up calculations and add stochasticity (low probability of getting stuck in local minimum), the gradients are approximated on a subset of data (mini-batch)

$$abla_{\mathbf{w}} \mathcal{C}(\mathbf{w}_t)
ightarrow
abla_{\mathbf{w}} \mathcal{C}^{\mathsf{MB}}(\mathbf{w}) = \sum_{i \in B_k}
abla_{\mathbf{w}} \mathcal{C}(\mathbf{x}_i, \mathbf{w})$$

Backpropagation algorithm

- For each training point (or mini-batch), a huge number of derivatives computations is required.
- Backpropagation is an efficient algorithm that accomplishes this complex task
 - Forward propagation step: data is passed through a network to determine the cost function
 - Backward propagation step: adjust the model's parameters to reduce the cost function

Forward propagation

• Computational graph for the forward pass



[https://towardsdatascience.com/neural-networks-backpropagation-by-dr-lihi-gur-arie-27be67d8fdce]

1)
$$\mathbf{Z}^{(1)} = \mathbf{W}^{(1)}\mathbf{X} + \mathbf{b}^{(1)}$$

2) $\mathbf{A}^{(1)} = \sigma \left(\mathbf{Z}^{(1)}\right)$
3) $\mathbf{Z}^{(2)} = \mathbf{W}^{(2)}\mathbf{A}^{(1)} + \mathbf{b}^{(2)}$
4) $\hat{\mathbf{Y}} = \sigma \left(\mathbf{Z}^{(2)}\right)$
5) $\mathcal{C} = \frac{1}{2}(\hat{\mathbf{Y}} - Y)^2 \rightarrow \mathcal{C}\left(\mathbf{b}^{(1)}, \mathbf{b}^{(2)}, \mathbf{W}^{(1)}, \mathbf{W}^{(2)}\right)$

Márcio Ferreira (CFisUC)

Backward propagation



[https://towardsdatascience.com/neural-networks-backpropagation-by-dr-lihi-gur-arie-27be67d8fdce]

$$\frac{\partial \mathcal{C}}{\partial \mathbf{W}^{(2)}} = \frac{\partial \mathcal{C}}{\partial \mathbf{A}^{(2)}} \frac{\partial \mathbf{A}^{(2)}}{\partial \mathbf{Z}^{(2)}} \frac{\partial \mathbf{Z}^{(2)}}{\partial \mathbf{W}^{(2)}}$$

- The total derivative is determined from the product of local derivatives
- All information required to compute these local gradients was saved during the forward pass

Backward propagation



[https://towardsdatascience.com/neural-networks-backpropagation-by-dr-lihi-gur-arie-27be67d8fdce]

$$\frac{\partial \mathcal{C}}{\partial \mathbf{b}^{(1)}} = \frac{\partial \mathcal{C}}{\partial \mathbf{A}^{(2)}} \frac{\partial \mathbf{A}^{(2)}}{\partial \mathbf{Z}^{(2)}} \frac{\partial \mathbf{Z}^{(2)}}{\partial \mathbf{A}^{(1)}} \frac{\partial \mathbf{A}^{(1)}}{\partial \mathbf{Z}^{(1)}} \frac{\partial \mathbf{Z}^{(1)}}{\partial \mathbf{b}^{(1)}}$$

After applying the procedure for each network's weight

$$\begin{pmatrix} \mathbf{W}_{t+1}^{(1)} \\ \mathbf{W}_{t+1}^{(2)} \\ \mathbf{b}_{t+1}^{(1)} \\ \mathbf{b}_{t+1}^{(2)} \end{pmatrix} = \begin{pmatrix} \mathbf{W}_t^{(1)} \\ \mathbf{W}_t^{(2)} \\ \mathbf{b}_t^{(1)} \\ \mathbf{b}_t^{(2)} \end{pmatrix} - \eta_t \begin{pmatrix} \partial \mathcal{C} / \partial \mathbf{W}_t^{(1)} \\ \partial \mathcal{C} / \partial \mathbf{W}_t^{(2)} \\ \partial \mathcal{C} / \partial \mathbf{b}_t^{(1)} \\ \partial \mathcal{C} / \partial \mathbf{b}_t^{(2)} \end{pmatrix}$$

Márcio Ferreira (CFisUC)

Machine learning: Neural Networks

Notation: batch size, mini-batch, epoch, and iteration

- Consider a dataset with N = 100 points
- By setting a **batch size** of 20, we are creating 5 **mini-batches**:

 $\{B_1, B_2, B_3, B_4, B_5\}$

- An **epoch** refers to one cycle through the full dataset
- Having 5 mini-batches, an epoch is completed in 5 iterations.



Márcio Ferreira (CFisUC)

Machine learning: Neural Networks

[https://doi.org/10.3348/kjr.2019.0312] vorks 27 - 30 June (2022) 16 / 38

Activation functions

- There are different activation functions with specific computational properties
- Their derivatives' properties have a crucial impact on training NN
 - Sigmoid/tanh have vanishing gradients for $|x| \gg 1$



[https://medium.com/@shrutijadon10104776/survey-on-activation-functions-for-deep-learning-9689331ba092]

Márcio Ferreira (CFisUC)

Machine learning: Neural Networks

Using DNN: classification problem

MNIST dataset

 $\mathcal{D}\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), ..., (\mathbf{x}_N, y_N)\},\$

where \mathbf{x}_i is a 2D array of gray-scale and y is the label



- Training set: 60000 images
- Testing set: 10000 images

Márcio Ferreira (CFisUC)

Machine learning: Neural Networks

MNIST dataset

• The whole dataset is represented by a 3D array (N imes X imes Y)

 $x_{ijk} \in X_{\text{train}} \rightarrow [1:60000, 1:28, 1:28]$ (rank 3 tensor)

• 2D arrays (slices) represent samples $(X \times Y)$

 $x_{1jk} \rightarrow [1, 1:28, 1:28]$ (rank 2 tensor)

 $x_{5jk} \rightarrow [5, 1:28, 1:28]$

• The outut values are in a 1D array (N)

 $y_{\text{train}} = \{y_1, y_2, ..., y_{60000}\} \rightarrow [1:60000] \text{ (rank 1 tensor)}$

• Their values indicate the class:

$$y_1 = [1] = 5, \quad y_2 = [2] = 0, \quad y_3 = [3] = 4, \quad \dots$$

• y_i are *slices* (scalars, rank 0 tensors) from the vector y_{train}

• $x_{1jk} \rightarrow [28 \times 28]$ (2D array with the gray-scale of each pixel)

^ V:	L °	V2	° V3	° V4	° v	5 0	v6 °	V7 -	V8	V9 :	V10 :	V11 :	V12 0	V13 °	V14 :	V15 °	V16 :	V17 :	V18 °	V19 °	V20 °	v21 °	V22 ÷	V23 :	V24 °	V25 0	V26 °	V27	V28	
1	0		0	0	0	0	0	0		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	-	3	0
2	0		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		J	0
3	0		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		3	0
4	0		0	0	0	0	0	0		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		3	0
5	0		0	0	0	0	0	0		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		3	0
6	0		0	0	0	0	0	0		0	0	0	0	3	18	18	18	126	136	175	26	166	255	247	127	0	0		3	0
7	0		0	0	0	0	0	0	0	30	36	94	154	170	253	253	253	253	253	225	172	253	242	195	64	0	0		3	0
8	0		0	0	0	0	0	0	45	238	253	253	253	253	253	253	253	253	251	93	82	82	56	39	0	0	0		3	0
9	0		0	0	0	0	0	0	18	219	253	253	253	253	253	198	182	247	241	0	0	0	0	0	0	0	0		3	0
10	0		0	0	0	0	0	0	0	80	156	107	253	253	205	11	0	43	154	0	0	0	0	0	0	0	0		J	0
11	0		0	0	0	0	0	0	0	0	14	1	154	253	90	0	0	0	0	0	0	0	0	0	0	0	0		J	0
12	0		0	0	0	0	0	0	0	0	0	0	139	253	190	2	0	0	0	0	0	0	0	0	0	0	0		J	0
13	0		0	0	0	0	0	0	0	0	0	0	11	190	253	70	0	0	0	0	0	0	0	0	0	0	0		J	0
14	0		0	0	0	0	0	0	0	0	0	0	0	35	241	225	160	108	1	0	0	0	0	0	0	0	0		J	0
15	0		0	0	0	0	0	0	0	0	0	0	0	0	81	240	253	253	119	25	0	0	0	0	0	0	0		3	0
16	0		0	0	0	0	0	0	0	0	0	0	0	0	0	45	186	253	253	150	27	0	0	0	0	0	0		3	0
17	0		0	0	0	0	0	0	0	0	0	0	0	0	0	0	16	93	252	253	187	0	0	0	0	0	0		3	0
18	0		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	249	253	249	64	0	0	0	0	0		3	0
19	0		0	0	0	0	0	0	0	0	0	0	0	0	0	46	130	183	253	253	207	2	0	0	0	0	0		د	0
20	0		0	0	0	0	0	0	0	0	0	0	0	39	148	229	253	253	253	250	182	0	0	0	0	0	0		د	0
21	0		0	0	0	0	0	0	0	0	0	24	114	221	253	253	253	253	201	78	0	0	0	0	0	0	0		2	0
22	0		0	0	0	0	0	0	0	23	66	213	253	253	253	253	198	81	2	0	0	0	0	0	0	0	0		3	0
23	0		0	0	0	0	0	18	171	219	253	253	253	253	195	80	9	0	0	0	0	0	0	0	0	0	0		3	0
24	0		0	0	0	55	172	226	253	253	253	253	244	133	11	0	0	0	0	0	0	0	0	0	0	0	0	-	3	0
25	0		0	0	0	136	253	253	253	212	135	132	16	0	0	0	0	0	0	0	0	0	0	0	0	0	0	-	2	0
26	0		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0)	0
27	0		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0)	0
28	0		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		2	0

• The output space: $y_1 = 5$

Balanced training set

- We need to check whether the training set is balanced
- The number of observations per class shows small deviations



Transform data

• We need to reshape the input space

$$x_{ij} \in X_{\text{train}} \to [1:60000, 784]$$

- Each imagine is now represented by a vector of length 784
- Rescale the dataset (normalize the features)

$$x_{ij} \to x_{ij}/255$$

Apply the one-hot encoding to the classes

$$y = 0 \rightarrow \mathbf{y} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ y = 1 \rightarrow \mathbf{y} = \begin{pmatrix} 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ & & & & \\ y = 9 \rightarrow \mathbf{y} = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$

One-hot encoding

- Our initial y variables contained label values $\{0,1,...,9\}$ rather than numeric values
- One-hot encoding ensures that our model does not consider higher numbers to be more important

V1 V2 V3 V4 V5 V6 V7 V8 V9 V10 1: 5 1: 0 0 0 0 0 0 0 2: 2: 1 0 0 0 Θ Θ 0 0 3: 4 3: 0 0 0 0 1 0 0 0 Θ 0 4: 4: 0 1 0 0 0 0 0 0 1 Θ 0 0 0 0 0 5: 9 5: 0 0 0 0 • 1 - - -59996: 8 59996: 0 0 0 0 1 0 59997: 59997: 0 Θ Θ 0 3 0 0 0 59998: 59998: 0 0 1 0 Θ 0 5 59999: 59999: 0 Θ 0 0 0 0 1 Θ Θ 0 6 60000: 60000: Θ 0 0 0 0 0 Θ 1 0 8

• $y_{ij} \in y_{\text{train}} \rightarrow [1:60000, 10]$ (2D array, rank 2 tensor)

Define the NN model (using Keras in R)

inputs <- layer_input(shape = c(784))

```
predictions <- inputs %>%
    layer_dense(units = 256, activation = 'relu') %>%
    layer_dense(units = 128, activation = 'relu') %>%
    layer_dense(units = 10, activation = 'softmax')
```

```
model <- keras_model(inputs = inputs, outputs = predictions)</pre>
```

Layer	(type)	Output Shape	Param #
input	(InputLayer)	[(None, 784)]	0
dense	(Dense)	(None, 256)	200960
dense	(Dense)	(None, 128)	32896
dense	(Dense)	(None, 10)	1290
Total	params: 235,146		

Trainable params: 235,146 Non-trainable params: 0

• Softmax function:
$$P(y = j | \mathbf{x}) = \frac{e^{\mathbf{x} \cdot \mathbf{w}_j}}{\sum_{k=1}^{K} e^{\mathbf{x}^{\mathsf{T}} \mathbf{w}_k}}$$

Network output: $(P_0, P_1, P_2, P_3, P_4, P_5, P_6, P_7, P_8, P_9)$

т

Márcio Ferreira (CFisUC)

Machine learning: Neural Networks

Setting a batch size

• Mini-batches of size 128

```
inputs <- layer_input(batch_shape = c(128,784))</pre>
```

```
predictions <- inputs %>%
    layer_dense(units = 256, activation = 'relu') %>%
    layer_dense(units = 128, activation = 'relu') %>%
    layer_dense(units = 10, activation = 'softmax')
```

```
model <- keras_model(inputs = inputs, outputs = predictions)</pre>
```

Layer (type)	Output Shape	Param #
input_4 (InputLayer)	[(128, 784)]	0
dense_23 (Dense)	(128, 256)	200960
dense_22 (Dense)	(128, 128)	32896
dense_21 (Dense)	(128, 10)	1290
Total params: 235,146 Trainable params: 235,146 Non-trainable params: 0		

Compile the NN model

 Categorical crossentropy (loss function): how distinguishable two discrete probability distributions are

$$L(\mathbf{W}) = -\sum_{i=0}^{9} y_i \cdot \log \,\hat{y}_i(\mathbf{W})$$

- W represents all the (235146) NN weights
- Adam Optimizer [several available: SGD, RMSprop,...]
- Metric: the fraction of the images correctly classified (accuracy)

```
model %>% compile(
    loss = "categorical_crossentropy",
    optimizer = 'adam',
    metrics = c("accuracy")
)
```

Training the NN model

- Randomly split the train dataset (60000 images):
 - $\blacksquare~80\%$ for training and 20% for validation
- 30 epochs: number of times the entire training set (48000 images) will pass through the NN
- Mini-batches of size of 128:

 $X_{in} \rightarrow \{ [1: 128, 784], [129: 257, 784], \dots \}$

- Each epoch has 48000/128 = 375 iterations
 - \blacksquare number of times the ${\bf W}$ are updated in each epoch

```
model %>% fit(
 x_train, y_train,
 epochs = 30,
 batch_size = 128,
 validation_split = 0.2
)
```

Output of the training stage

.....

Epoch 25/30
375/375 [============] - 13s 36ms/step - loss: 0.0017 - accuracy: 0.9993 - val_loss: 0.1176 - val_accuracy: 0.9791
Epoch 26/30
375/375 [===========] - 14s 37ms/step - loss: 0.0011 - accuracy: 0.9997 - val_loss: 0.1219 - val_accuracy: 0.9804
Epoch 27/30
375/375 [============] - 13s 35ms/step - loss: 0.0070 - accuracy: 0.9978 - val_loss: 0.1339 - val_accuracy: 0.9758
Epoch 28/30
375/375 [============] - 10s 26ms/step - loss: 0.0099 - accuracy: 0.9967 - val_loss: 0.1152 - val_accuracy: 0.9797
Epoch 29/30
375/375 [============] - 155 40ms/step - loss: 0.0014 - accuracy: 0.9994 - val_loss: 0.1334 - val_accuracy: 0.9793
Epoch 30/30
375/375 [===========] - 14s 39ms/step - loss: 0.0040 - accuracy: 0.9988 - val_loss: 0.1375 - val_accuracy: 0.9783

• Do we understand the meaning of all these numbers?

Márcio Ferreira (CFisUC)

Machine learning: Neural Networks

27 - 30 June (2022) 28 / 38

Deep Neural Networks with keras

• Graphical representation of the learning stage



Out-of-sample performance of the model

- The final model's performance is measured on the test set
 - 10000 images

- Accuracy of 98% (the model correctly classifies 98 out of 100 figures)
- The *evaluate()* function uses batches (default size of 32) just to speed-up evaluation (10000/32=312.5)

Model output

Model predictions

• Last layer: $(P_0, P_1, P_2, P_3, P_4, P_5, P_6, P_7, P_8, P_9)$

0 1 2 3 4 4 5 6 6 7 8 9 7 2 1 2 .243218-5 5 .204408-12 1 .0822418-68 2.795768-21 1.5747758-14 3.202079-26 9 .999995-61 6.3165798-14 5.257895-07 2 1 .3306588-18 2.8552818-14 1 .0090008-00 1.609751-16 6 .226776-28 4 .0657918-26 9 .9513944-19 5 .420499-27 3 .766663-14 1.093308-28 3 .3308088-12 9 .999995-01 4 .7822488-07 3 .4280248-13 4 .9983328-10 8 .679319-12 3 .5953268-10 2 .3224528-08 4 .7137388-08 2 .584044-15 4 .10090008-00 4 .121874-18 1 .1411989-12 2 .120390-12 1 .315718-17 3 .4485598-18 1 .0681586-12 6 .5849536-17 1 .071761-16 3 .0685388-16 5 .9.358338-14 2 .531308-14 4 .966977-12 1 .019905-17 1 .000008-00 2 .024498-20 9 .4023328-14 7 .9902718-14 3 .3898428-15 1.9173318-09 ...

9996: 1.840926-28 1.223770e-18 1.000000e+00 4.778302e-16 8.188080e-38 1.464597e-23 2.409068e-28 1.215604e-15 7.508420e-16 2.860830e-39 9997: 1.32177e-19 3.057360e-18 4.519080e-12 9.999950e 1.2132357e-24 1.603080e-08 4.081237e-27 2.538337e-15 5.46537e-13 3.256952e-07 9996: 2.140457e-25 5.025506e-20 3.493071e-24 6.937665e-24 1.000000e+00 4.665477e-27 2.479116e-22 7.837477e-17 2.080704e-19 3.053961e-3 9996: 2.140457e-25 5.025506e-20 3.493071e-24 6.937665e-24 1.000000e+00 4.665477e-27 2.479116e-22 7.837477e-17 2.080704e-19 3.053961e-3 9996: 2.140457e-25 5.025506e-20 3.493071e-24 6.397467e-26 1.385951e-28 1.00000e+00 6.783952e-20 1.148829e-18 5.537012e-14 3.021205e-21 0.00000e+00 6.68717e-22 1.048279e-13 3.15107e-19 4.579642e-29 1.00000e+00 6.89717e-22 1.048279e-13 3.15107e-19 4.579642e-29 1.09308e-20 6.201741e-18 6.504062e-20 1.00000e+00 6.89717e-33 3.15107e-19 4.579642e-29 1.00000e+00 6.89717e-33 3.15107e-19 4.579642e-29 1.00000e+00 6.89717e-33 3.15107e-19 4.579642e-20 1.935952e-20 1.405396e-20 1.935952e-20 1.405396e-20 1.935952e-20 1.90000e+00 6.89717e-33 3.15107e-19 4.579642e-20 1.9959600e+00 6.89717e-33 3.15107e-19 4.579642e-20 1.90000e+00 6.89717e-33 3.15107e-19 4.579640e+00 6.9776400000e+00 6.89717e-33 3.15107e+000000e+00 6.9776400000e+00 6.9776400000e+00 6.89717e-33 3.15107e+000000e+00 6

Dataset values

	0	1	2	3	4	5	б	7	8	9
1:	0	0	0	0	0	0	0	1	0	0
2:	0	0	1	0	0	0	0	0	0	0
3:	0	1	0	0	0	0	0	0	0	0
4:	1	0	0	0	0	0	Θ	0	Θ	0
5:	0	0	0	0	1	0	0	0	Θ	0
9996:	0	0	1	0	0	0	0	0	0	0
9996: 9997:	0 0	0 0	1 0	0 1	0 0	0 0	0 0	0 0	0 0	0 0
9996: 9997: 9998:	0 0 0	0 0 0	1 0 0	0 1 0	0 0 1	0 0 0	0 0 0	0 0 0	0 0 0	0 0 0
9996: 9997: 9998: 9999:	0 0 0	0 0 0	1 0 0	0 1 0 0	0 0 1 0	0 0 0 1	0 0 0	0 0 0	0 0 0	0 0 0
9996: 9997: 9998: 9999: 10000:	0 0 0 0	0 0 0 0	1 0 0 0	0 1 0 0	0 0 1 0 0	0 0 0 1 0	0 0 0 0	0 0 0 0	0 0 0 0	0 0 0 0

Márcio Ferreira (CFisUC)

Machine learning: Neural Networks

27 - 30 June (2022) 31 / 38

Confusion Matrix

• Summary of prediction results



Regularization for NN: dropout

- The dropout layer **prevents overfitting** by randomly assigning neurons to 0 at each step during training with a given frequency.
- dropout rate = 0.5: half of the neurons excluded from each update iteraction.



https://jamesmccaffrey.wordpress.com/2018/05/11/neural-network-library-dropout-layers/

Adding dropout to the model

Adding dropout layers

```
inputs_m2 <- layer_input(batch_shape = c(128,784))</pre>
```

```
predictions_m2 <- inputs_m2 %>%
layer_dense(units = 256, activation = 'relu') %>%
layer_dense(units = 0.4) %>%
layer_dense(units = 128, activation = 'relu') %>%
layer_dense(units = 10, activation = 'softmax')
```

```
model_m2 <- keras_model(inputs = inputs_m2, outputs = predictions_m2)</pre>
```

Layer ((type)	Output Shape	Param #		
input	(InputLayer)	[(128, 784)]	0		
dense	(Dense)	(128, 256)	200960		
dropout	t (Dropout)	(128, 256)	0		
dense	(Dense)	(128, 128)	32896		
dropout	t (Dropout)	(128, 128)	0		
dense	(Dense)	(128, 10)	1290		
Total p Trainal Non-tra	params: 235,146 Jle params: 235,146 pinable params: 0				

Compare models (with/without dropout)



Márcio Ferreira (CFisUC)

Machine learning: Neural Networks

Models' performances

- Evaluation both models in the test set (10000 digits)
- Without dropout:

accuracy = 0.9801000

• With dropout:

accuracy = 0.98119998

• Human error rates around 2-2.5% (accuracy 97.5-98%)

Next lectures ...



37 / 38

Next lectures ...



38 / 38