

Università di Pisa
Dipartimento di Informatica

ML 2025 Project

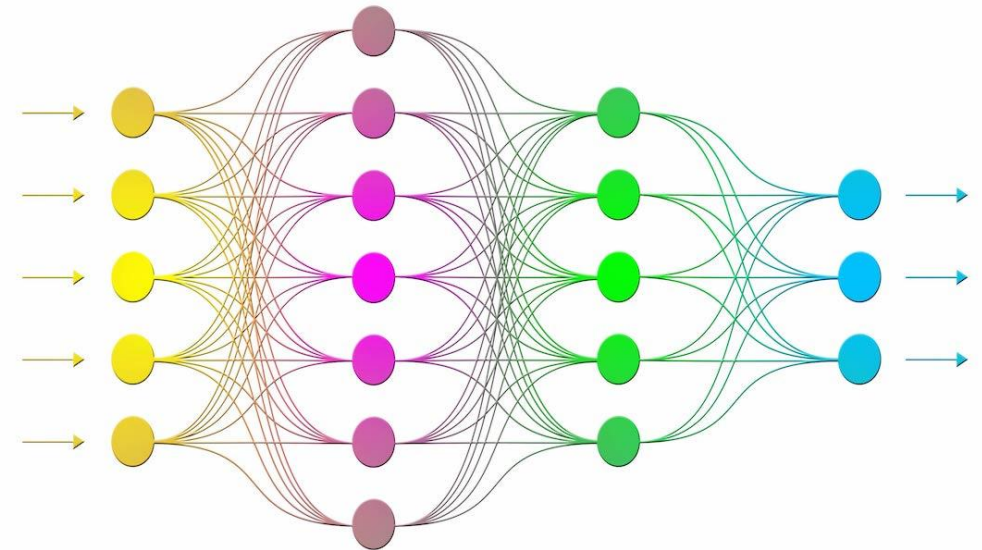
Authors: Andrea Marcheschi
(a.marcheschi6@studenti.unipi.it)
Luca Nasini
(l.nasini@studenti.unipi.it)
Simone Passera
(s.passera@studenti.unipi.it)

Team name: SimoLucaAndre

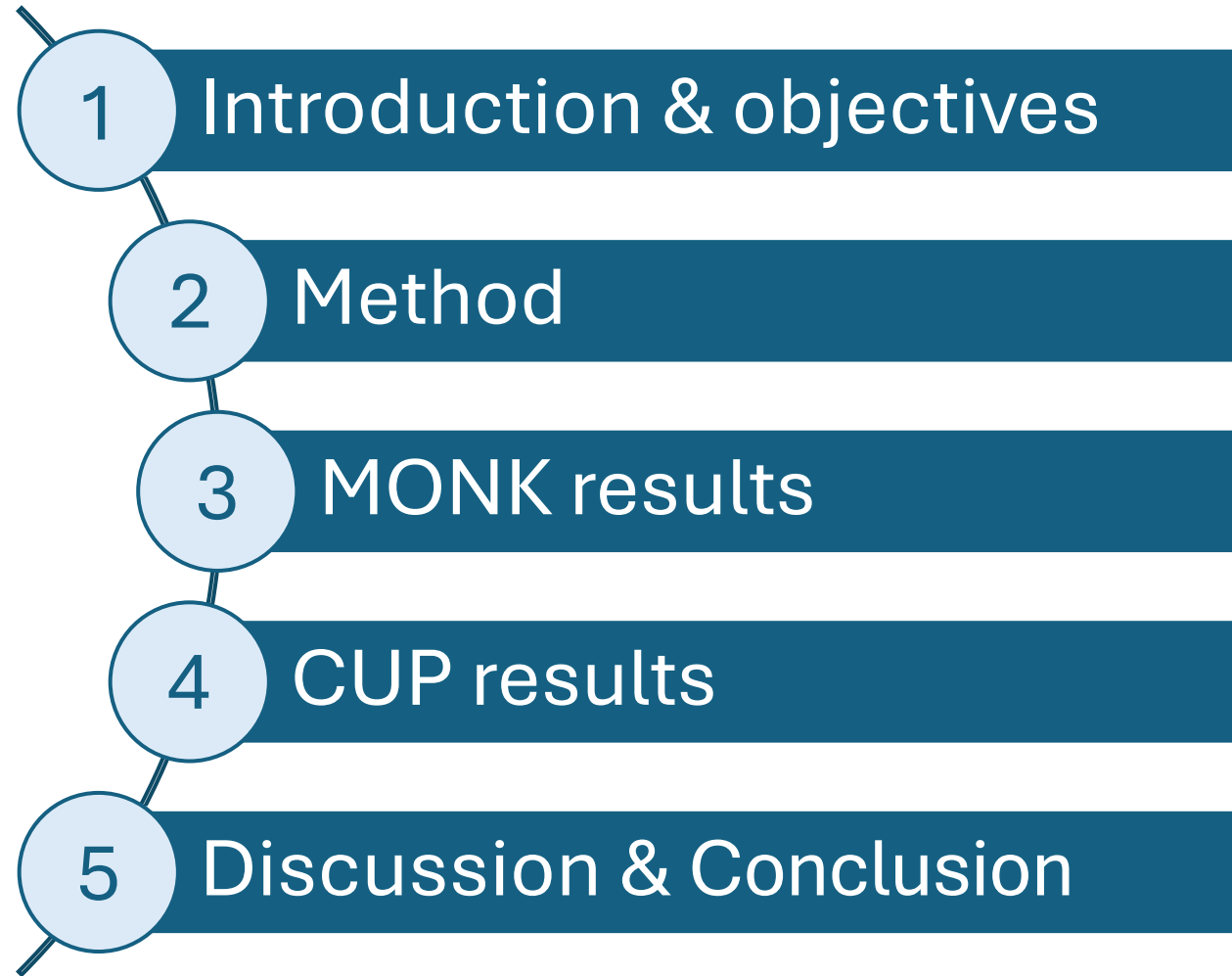
Master degree:

Date: 07/01/2026

Type of project: A



SUMMARY



Introduction & objectives



What we want to do:

Create, from scratch, a package to build a general neural network that can predict, with sufficient small error, new data target

Neural Network with:

- Early stopping, regularization, (Nesterov) momentum
- (Stochastic) Gradient Descent and Backpropagation algorithms
- 4 type of initialization (random, fan-in, xavier, he)
- Model validation and assessment (both with holdout, k-fold or leave-one-out)

Method - 1 (what is in the package)



Package scripts:

- main.py;
- model.py;
- model_selection.py;
- data_loaders.py;
- utils.py;
- activations.py;
- losses.py

External package used:

- numpy (for vectors);
- pandas (for reading .csv);
- matplotlib (for plots);
- pickle (for model savings)

Method – 2 (code description)



Philosophy of «model.py» code

Create a `NeuralNetwork` object with (dense) `NeuronLayer` (made of `Neuron`) with same activation functions. All `Neuron` have a bias and a list of weights. After the weights are initialized (and biases) all layers (hiddens and output), do the training until early stopping conditions are satisfied. Training can be batch or online (or also mini-batch). The `NeuralNetwork` class permit to plot the learning curves and save the entire neural network on a .pkl file

Philosophy of «model_selection.py» code

First select the type of search, and types of model selection/assessment. If model selection is k-fold split data in fold else respect the given proportions. Then choose the hyperparameters from the search and choose the best model from lower (average) loss. Retrain the model with `tr+vl` data. Evaluate the model with `ts` data. Then you can save the entire run (best model config, weights, biases, etc.).

Method – 3 (necessary features)



Activation functions:

- linear;
- sigmoid;
- tanh;
- (leaky)ReLU;
- softplus

Loss functions:

- MSE;
- MEE;
- binary cross-entropy

Early stopping conditions:

Monitor: «val_loss»:	reset patience if new loss is lesser than (best loss - ϵ)
Monitor: «val_accuracy»:	reset patience if new accuracy is greater than (best accuracy + ϵ)
Monitor: «train_loss»:	set patience at 0 if new loss lesser than target loss (only for retraining)

Method – 4 (necessary features)



Momentum (α):

$$\Delta w_{tu} = \eta \delta_t o_u + \alpha \Delta w_{\text{old},tu}$$

Nesterov: add α -term before computing the deltas and update weights unless this term

Regularization (λ):

$$w_{\text{new},tu} = w_{tu} + \Delta w_{tu} - \lambda w_{tu}$$

Early stopping conditions:

- | | |
|--------------------------|--|
| Monitor: «val_loss»: | reset patience if new loss is lesser than (best loss - ϵ) |
| Monitor: «val_accuracy»: | reset patience if new accuracy is greater than (best accuracy + ϵ) |
| Monitor: «train_loss»: | set patience at 0 if new loss lesser than target loss (only for retraining) |

Method – 3 (model novelties)



Nesterov momentum

Require: Learning rate ϵ , momentum α .

Require: Initial parameter θ , initial velocity v .

while stopping criterion not met **do**

Sample a minibatch of m examples from the training set $\{x^{(1)}, \dots, x^{(m)}\}$ with corresponding labels $y^{(i)}$.

Apply interim update: $\tilde{\theta} \leftarrow \theta + \alpha v$

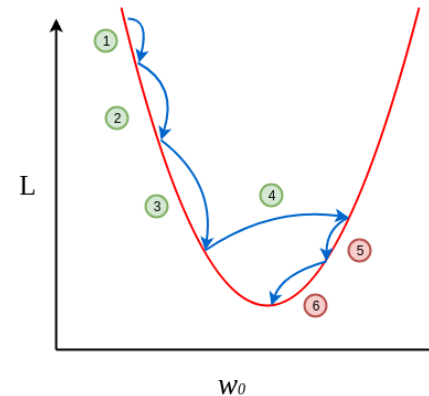
Compute gradient (at interim point):

$$g \leftarrow \frac{1}{m} \nabla_{\tilde{\theta}} \sum_i L(f(x^{(i)}; \tilde{\theta}), y^{(i)})$$

Compute velocity update: $v \leftarrow \alpha v - \epsilon g$

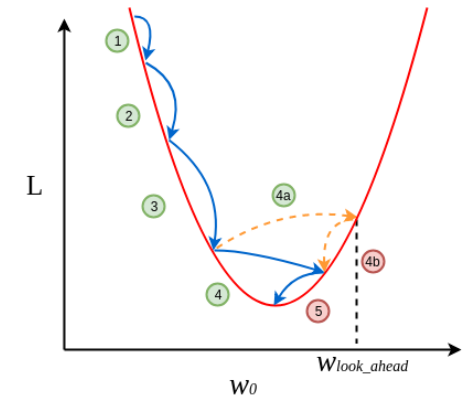
Apply update: $\theta \leftarrow \theta + v$

end while



(a) Momentum-Based Gradient Descent

$$\text{Green circle} \Rightarrow \frac{\partial L}{\partial w_0} = \frac{\text{Negative}(-)}{\text{Positive}(+)}$$



(b) Nesterov Accelerated Gradient Descent

$$\text{Red circle} \Rightarrow \frac{\partial L}{\partial w_0} = \frac{\text{Negative}(-)}{\text{Negative}(-)}$$

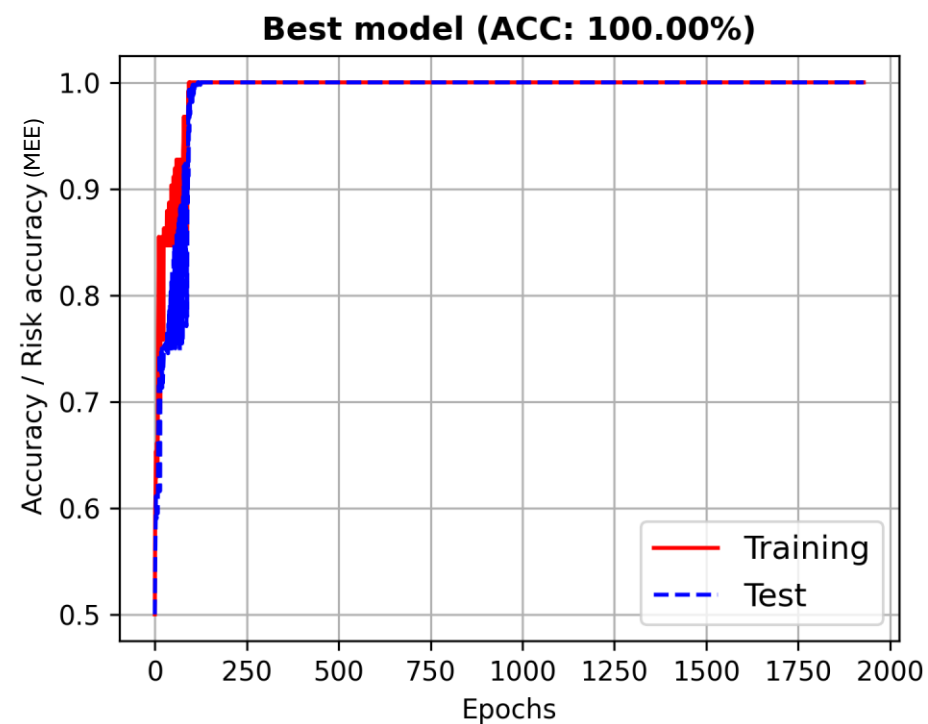
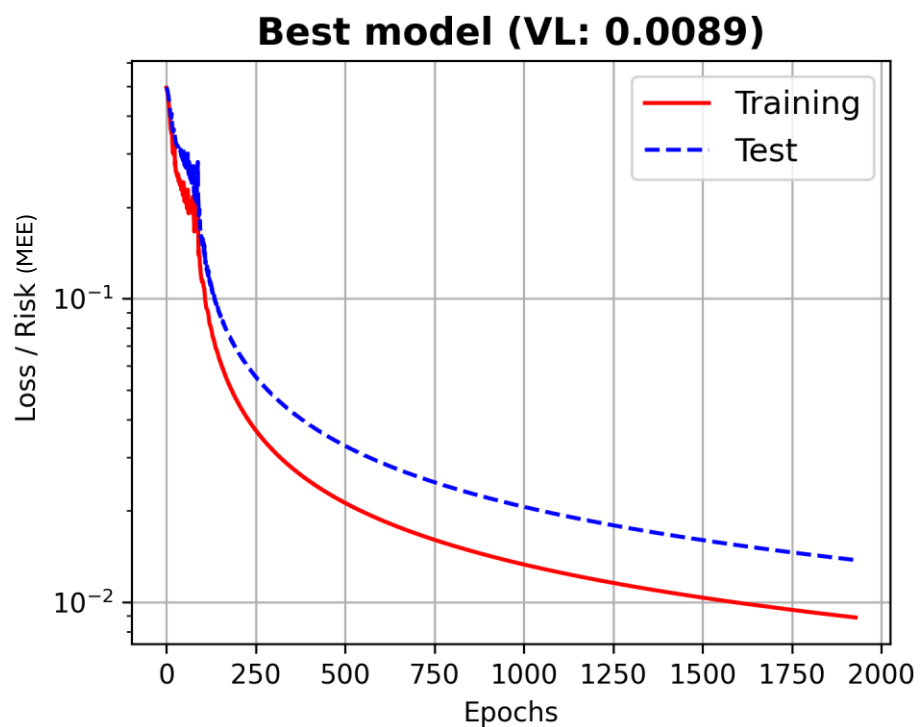
MONK results



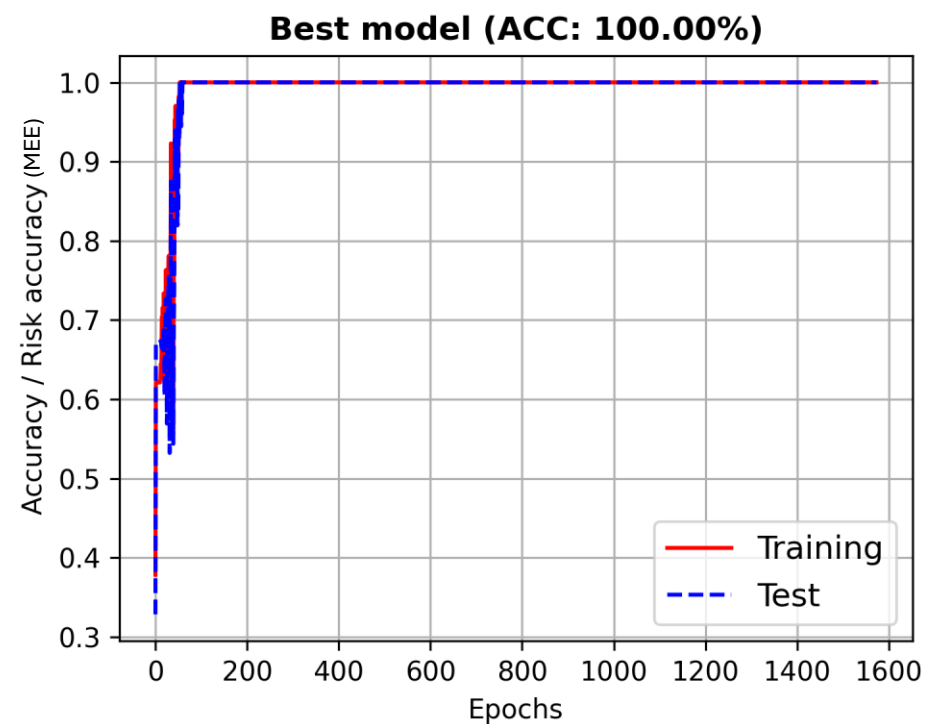
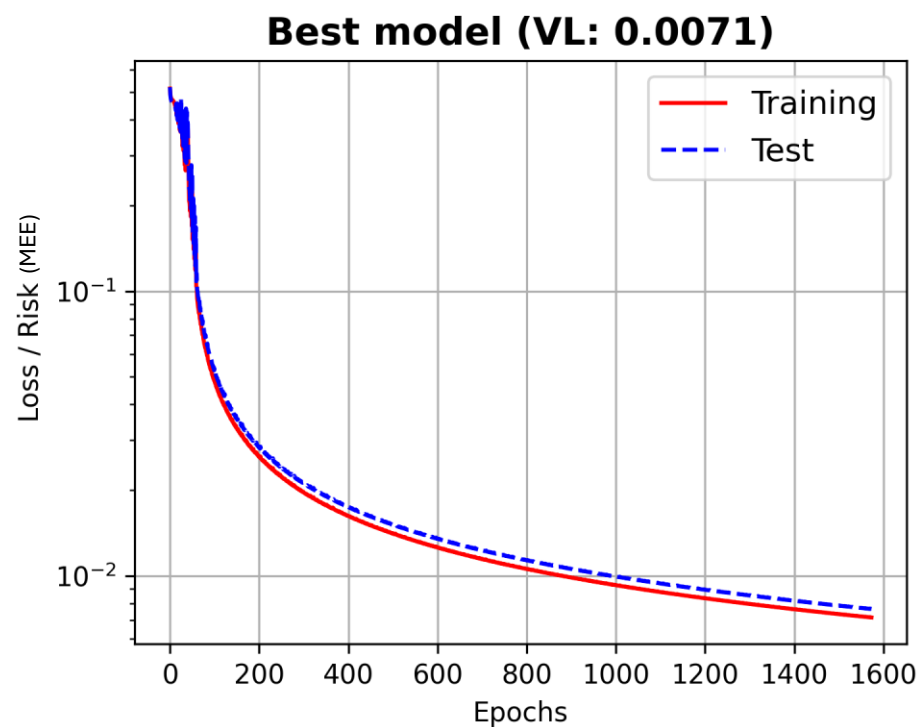
Input units: 17 (from one-hot-encoding)
Output units: 1
Activation function: ReLU (hidden layer), sigmoid (output unit)
Initialization: fan-in

Task	Units	η	λ	α	MEE (TR/TS)	Accuracy (TR/TS)
MONK1	4	0.0452	-	-	0.0089 / 0.0138	100% / 100%
MONK2	4	0.0480	-	-	0.0071 / 0.0076	100% / 100%
MONK3 (no reg.)	5	0.0143	-	-	0.0282 / 0.0779	97.54% / 92.13%
MONK3	4	0.0172	0.0011	-	0.0728 / 0.1002	99.18% / 95.37%

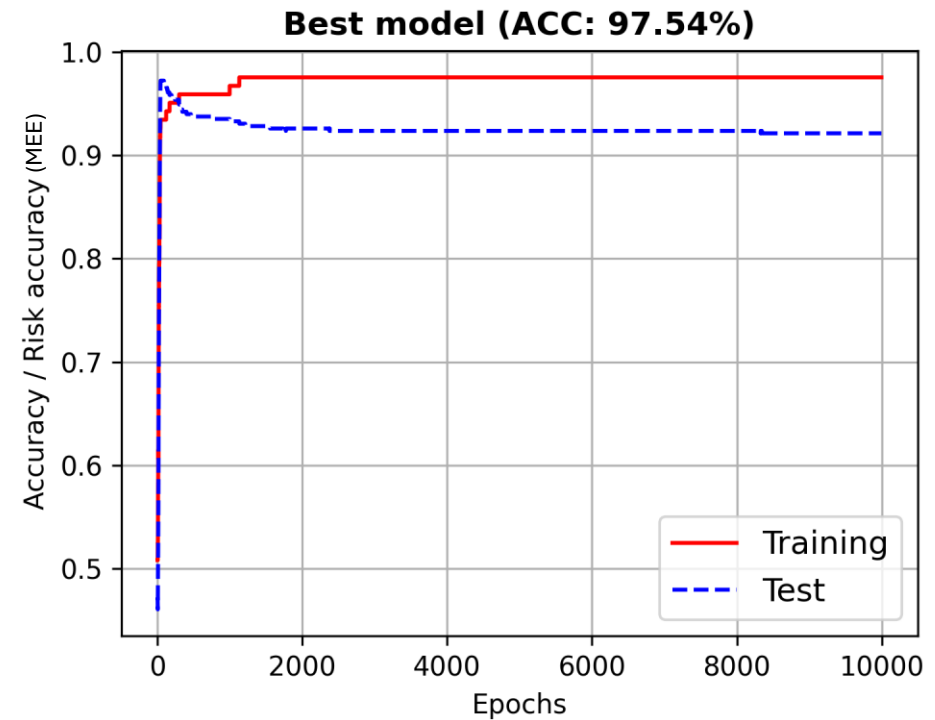
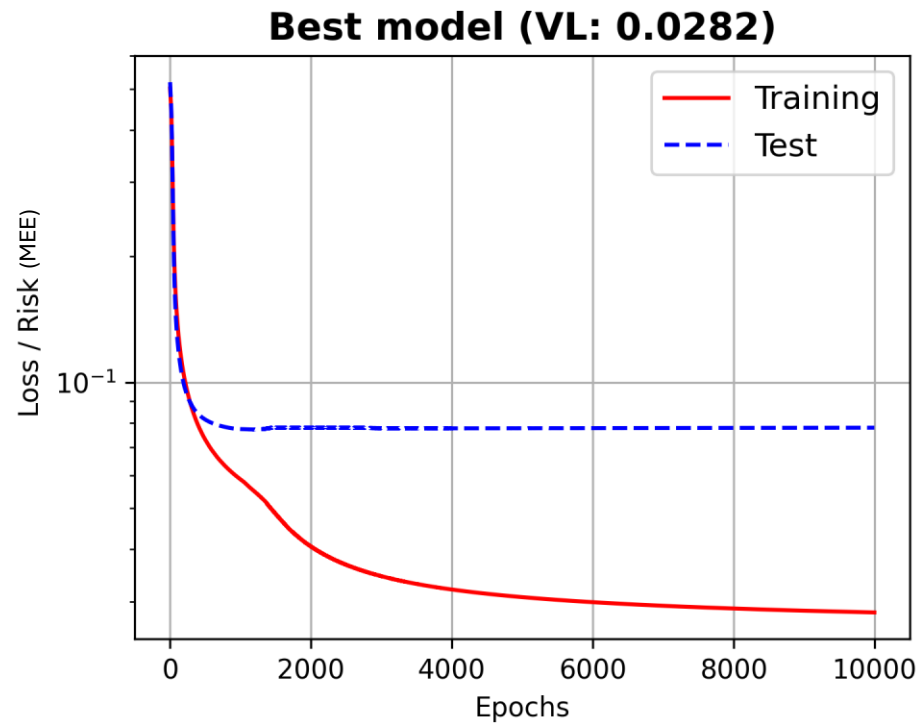
MONK 1 - plots



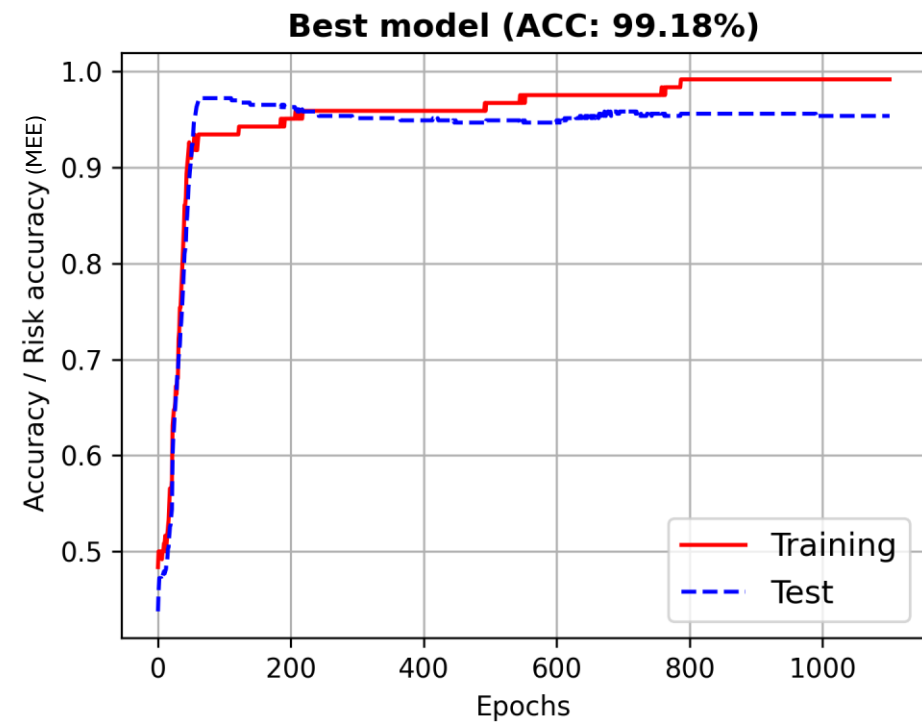
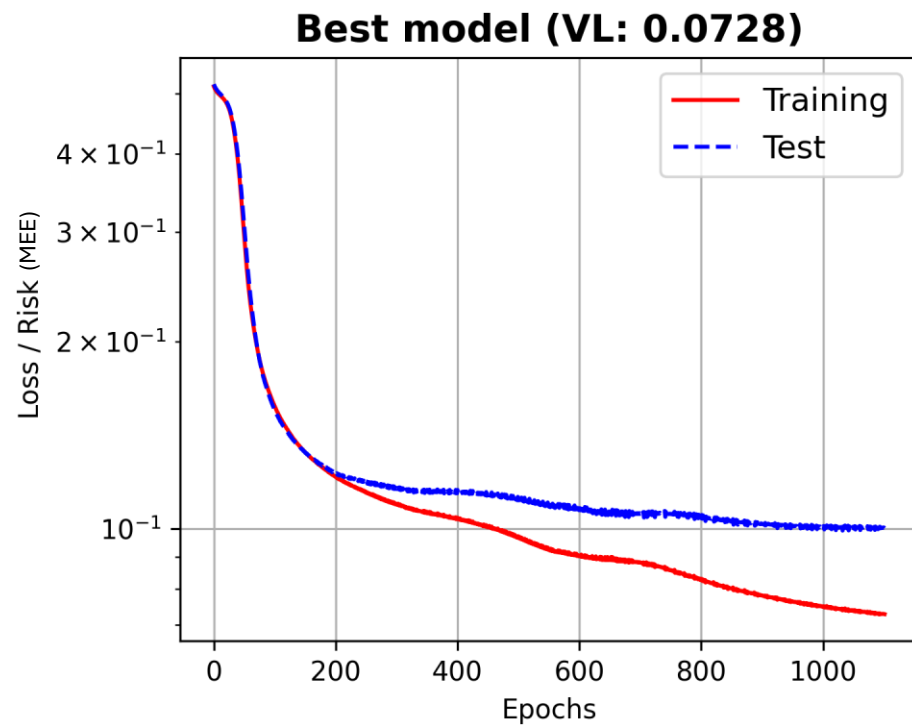
MONK 2 - plots



MONK 3 (no reg.) - plots



MONK 3 - plots



CUP Validation Schema: data splitting



MODEL ASSESSMENT:

Hold-out technique to split the entire dataset in Selection Set (Training + Validation, 80% of Dataset) and Test Set (20% of Dataset)

MODEL SELECTION:

Hold-out technique, with 50% TR and 30% VL, (or **K-fold CV** technique, with $k=5$) to perform the model selection phase on the Selection Set

Hold-out or 5-fold-CV
TR + VL
(80%)

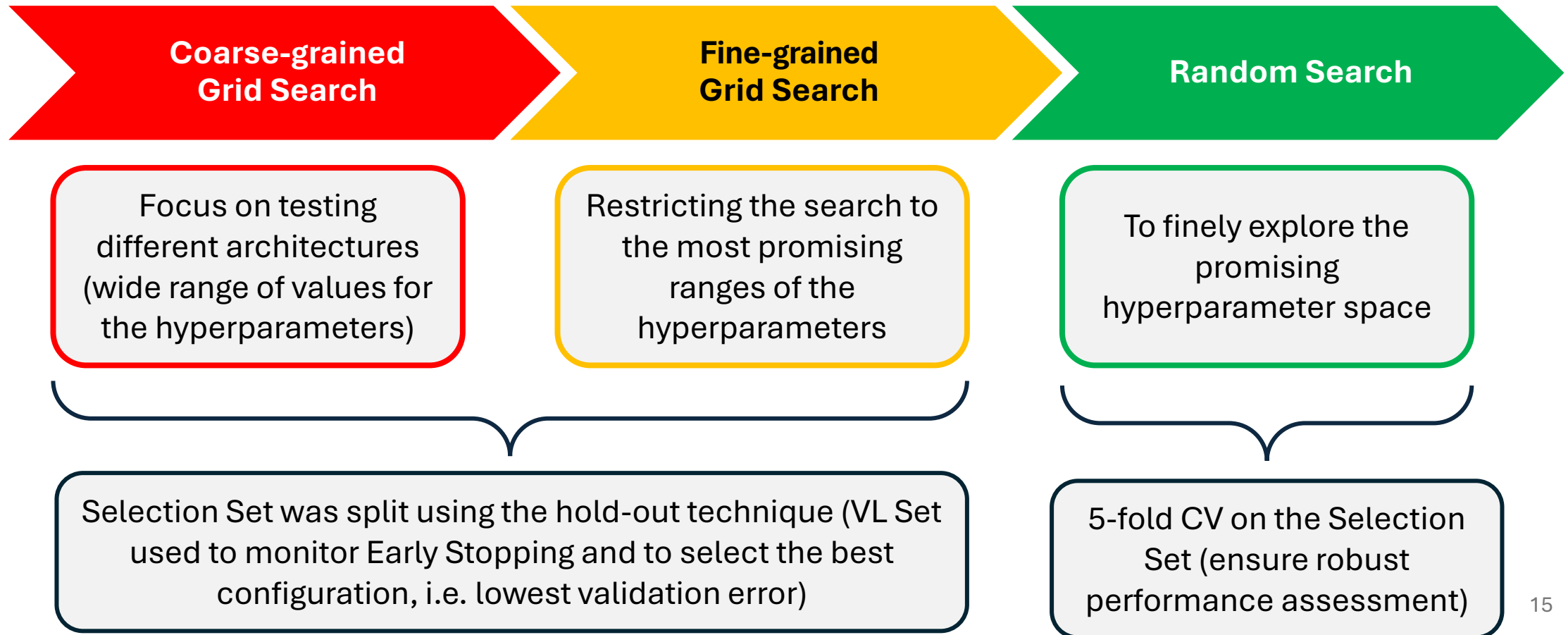
Hold-out
Internal TS
(20%)

After model selection ⇒ retraining on the entire Selection Set (Training + Validation)
After model assessment ⇒ final retraining on the entire Dataset (Training + Validation + Internal Test)

CUP Validation schema: model selection



In order to choose the best hyperparameter configuration, the model selection phase was performed:



CUP Validation schema: model selection



Coarse-grained Grid Search

Architecture	Activation functions	η	α	λ	Batch size	Initialization method	Preprocessing
12*,50,50,50,4*	Hidden: ReLU Output: Linear	0.001	0.1	1e-02	full	random (range [-0.7,0.7])	Rescaling
12*,75,75,75,4*		0.0001	0.3	1e-04			
12*,100,100,100,4*		0.00001	0.5	1e-06			
12*,128,128,128,4*			0.8 (Nesterov=False)	1e-08			

CUP Validation schema: model selection



Fine-grained
Grid Search

Architecture	Activation functions	η	α	λ	Batch size	Initialization method	Preprocessing
12*,45,45,45,4*	Hidden: ReLU Output: Linear	1e-05 3e-05 5e-05	0.1	1e-04 1e-06 1e-08	full	random (range [-0.7,0.7])	Rescaling
12*,50,50,50,4*			0.3				
12*,55,55,55,4*			0.5				
12*,68,68,68,4*			0.8				
(Nesterov=False)							

CUP Validation schema: model selection



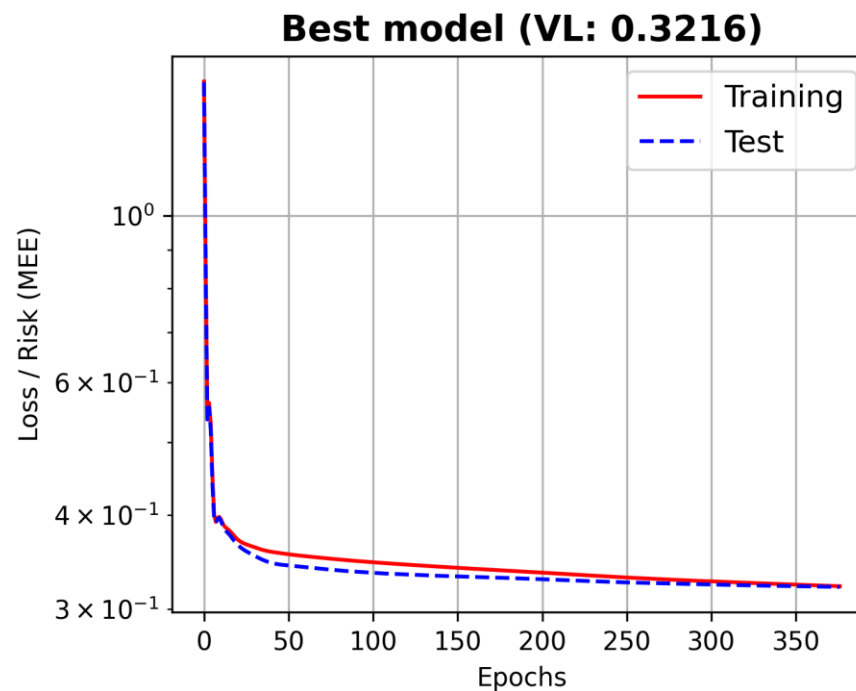
Random Search

Architecture	Activation functions	η	α	λ	Batch size	Initialization method	Preprocessing
12*,40,40,40,4* 12*,50,50,50,4*	Hidden: ReLU Output: Linear	2.5e-05 3.5e-05	0.7 0.9 (Nesterov=True)	0.5e-08 1.5e-08	full	He (range [-0.7,0.7])	Rescaling

CUP Final Model



The final model was selected based on achieving the lowest validation error combined with a stable learning curve.



Architecture	Activation functions	η	α	$\ln \lambda$	Batch size	Initialization method	Preprocessing
12*, 44, 41, 48, 4*	Hidden: ReLU Output: Linear	2.690e-05	0.895 (Nesterov=True)	-18.04	full	He (range [-0.7, 0.7])	Rescaling

CUP Final Model – results



VL	MEE TR+VL	TS
35.908436	33.987997	32.912784

```
19  "training": {
20      "epochs": 10000,
21      "search_type": "grid",
22      "number_random_trials": 50,
23      "learning_rate": {
24          "eta": 0.00002689944986613886,
25          "min_rate": 100,
26          "decay_factor": 0.0001
27      },
28      "momentum": 0.8947475896692705,
29      "nesterov": true,
30      "regularization": 1.4675420469884381e-8,
```

Discussion



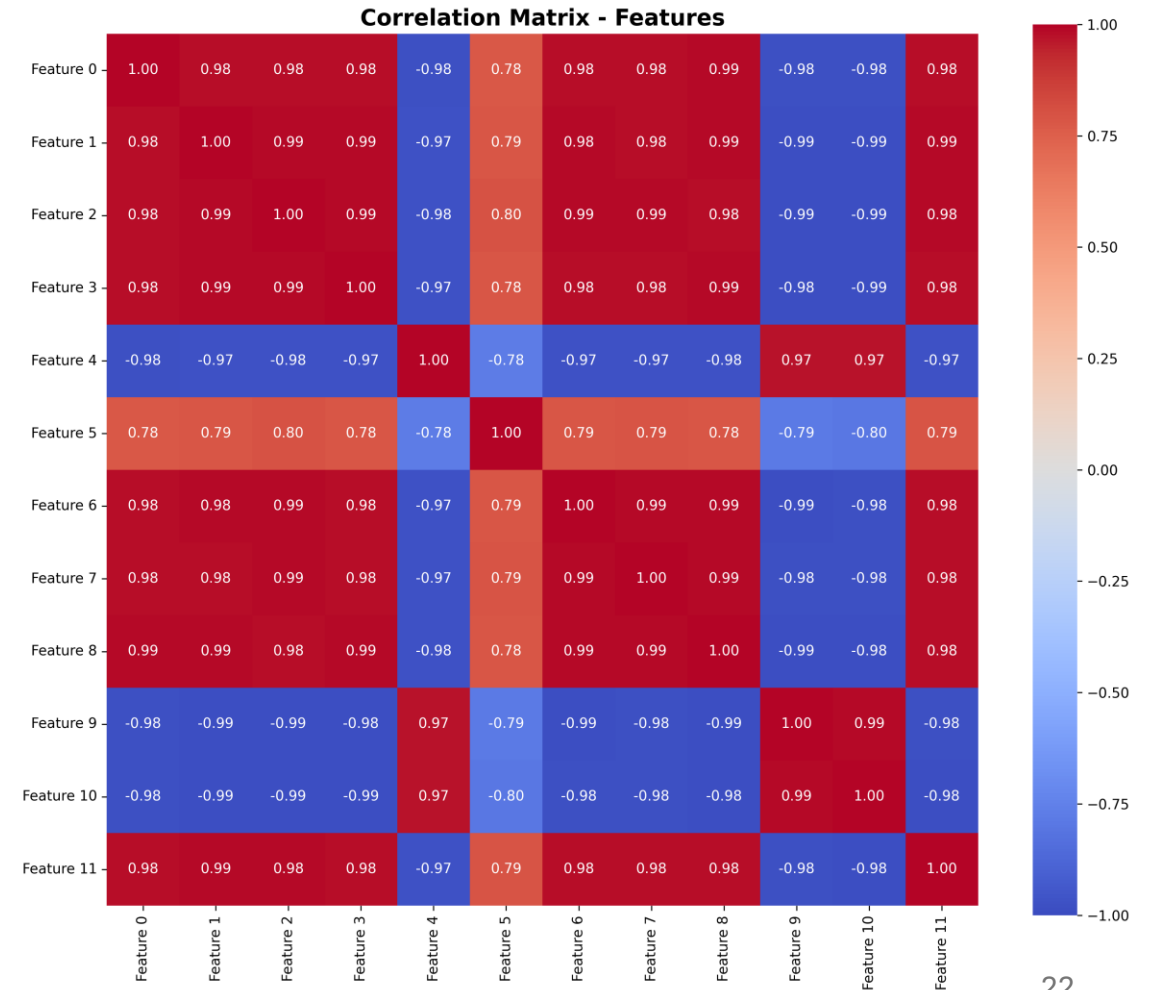
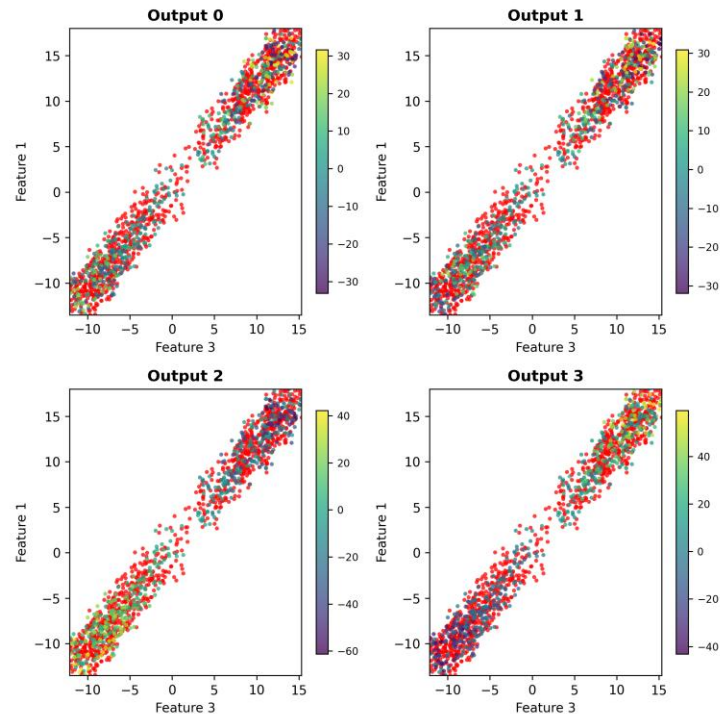
Total runs:	a few thousand (including bad ones)
Total runtime:	~3/4hrs per PC (coarse+fine search)
Most important hyperparameters:	Architecture and Learning Rate, followed by Regularization and Momentum

- Only a small range of lr is good (otherwise divergence/instability/very slow learning)
- Random search allows to explore this range quite well
- Small value of Learning Rate due to absence of $1/\text{batch_size}$ prefactor

Possible improvements - 1



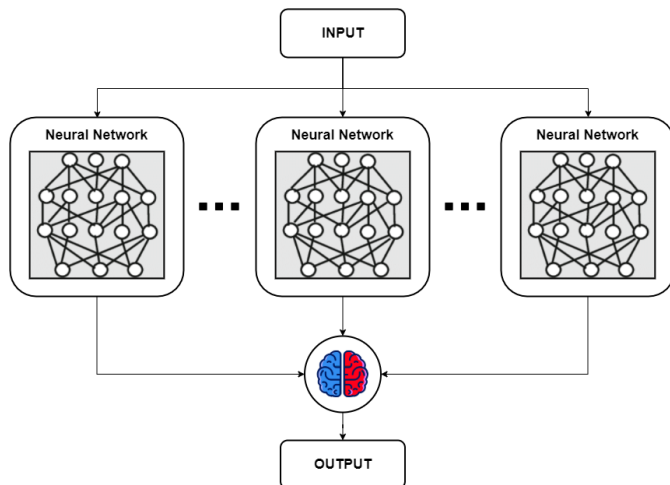
Some of the features or target CUP data are highly (anti)correlated with other features.
It would be usefull to develop by scratch some form of PCA to do a better preprocess.



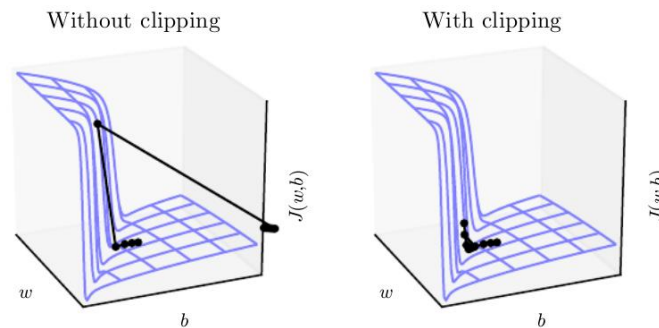
Possible improvements - 2



Ensembling



Gradient clipping



Multithreading scripts



Conclusion



What we learnt

- Build a NN from scratch using only numerical libraries (i.e. Numpy) without higher level ones (i.e. Scikit, Pytorch, ...)
- Cooperate (also between people from different curriculum)
- Read books and papers, don't trust the internet!

References



- Datta, Leonid. «A Survey on Activation Functions and their relation with Xavier and He Normal Initialization». arXiv, 2020. *DOI.org (Datacite)*, <https://doi.org/10.48550/ARXIV.2004.06632>.
- Shlens, Jonathon. «A Tutorial on Principal Component Analysis». arXiv, 2014. *DOI.org (Datacite)*, <https://doi.org/10.48550/ARXIV.1404.1100>.
- Ian Goodfellow, Yoshua Bengio, & Aaron Courville (2016). *Deep Learning*. MIT Press.



Thank you for your time! 😊