

Deep Learning in Practice

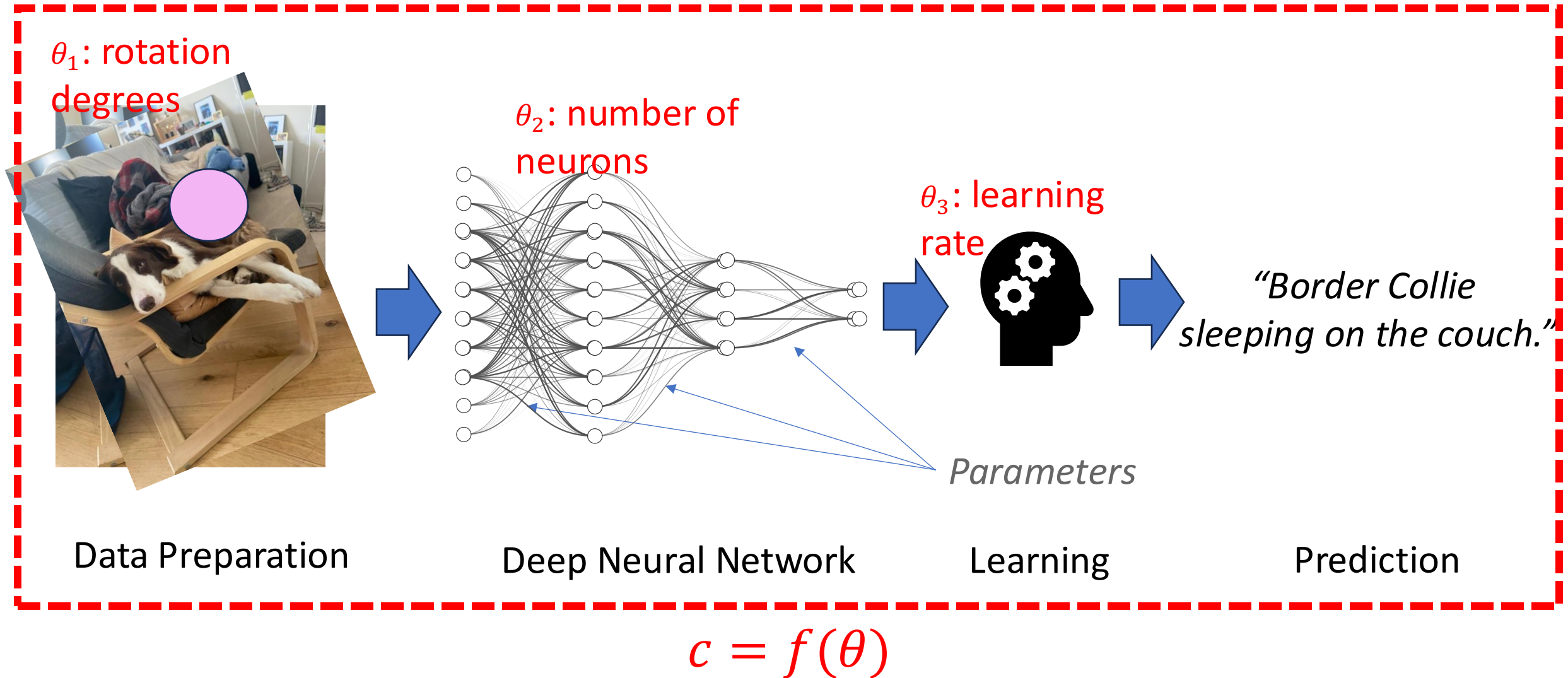
Optimization of Learning Workflows

Hyperparameter Optimization, Neural Architecture Search, Ensemble

Romain Egele

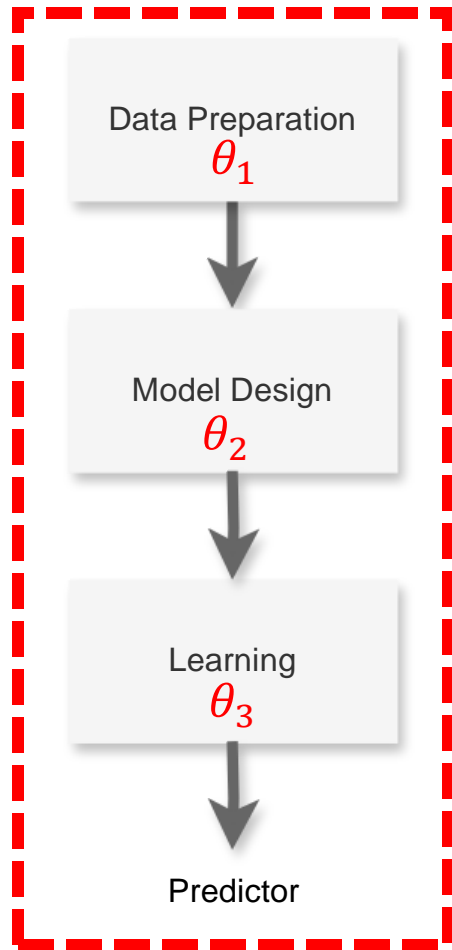


An Example of Learning Workflow



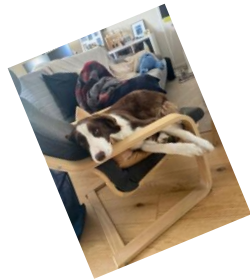
Parametrization of Learning Workflows

Categorical, Discrete, **Real**



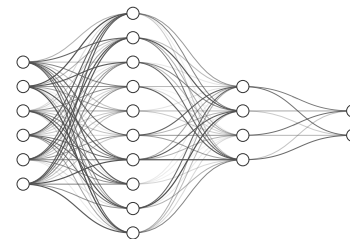
$$c = f(\theta)$$

Data Preparation



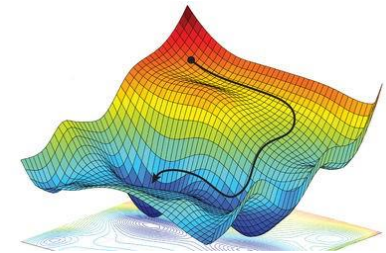
θ_1 {
Rotation Degrees
Augmentation Type
Noise Rate
...

Model Design



θ_2 {
Number of layers
Layer Type
Dropout rate
...

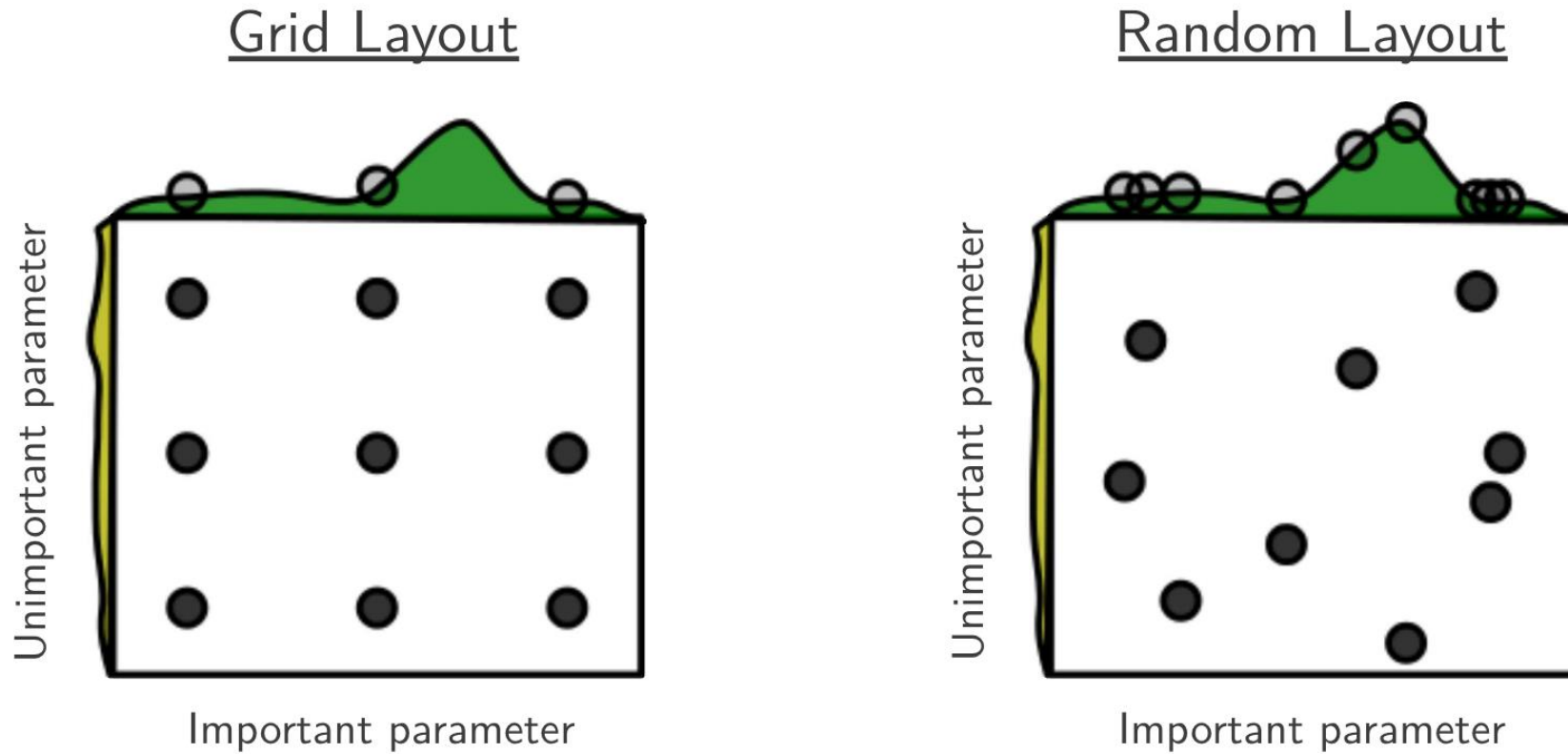
Learning



θ_3 {
Optimizer
Learning rate
Batch size
...

[Egelé and Jacques Junior, 2024]

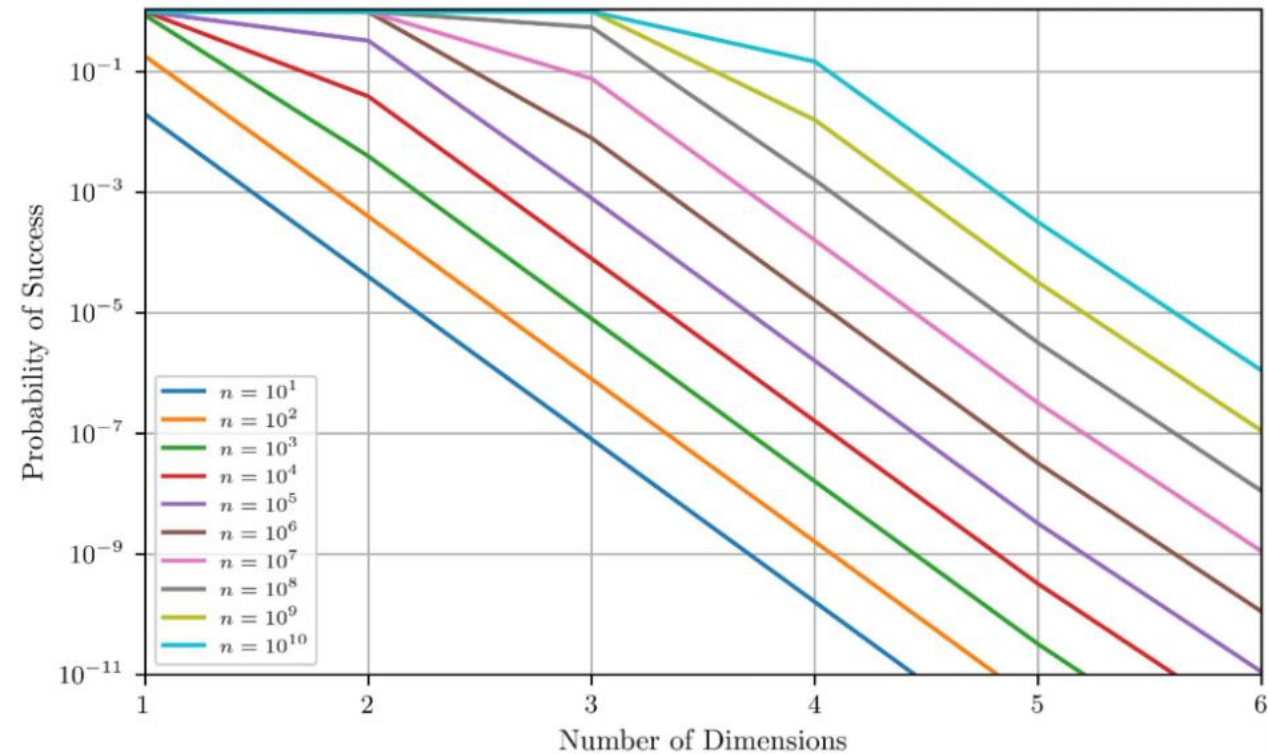
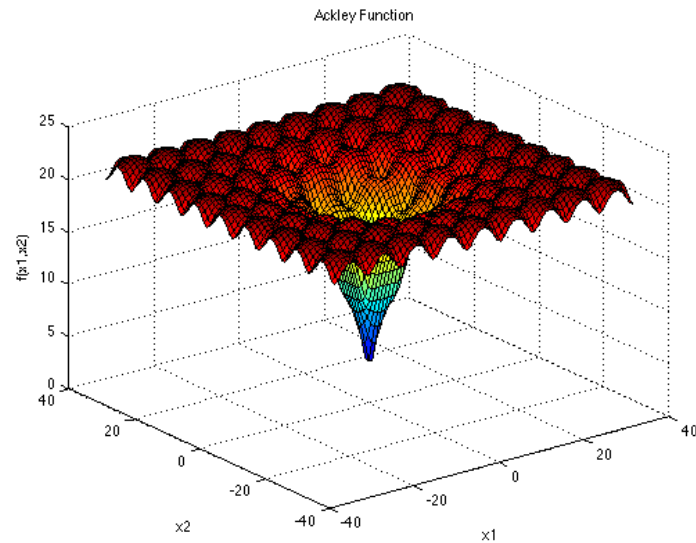
Grid Search



(Bergstra et Bengio, 2012)

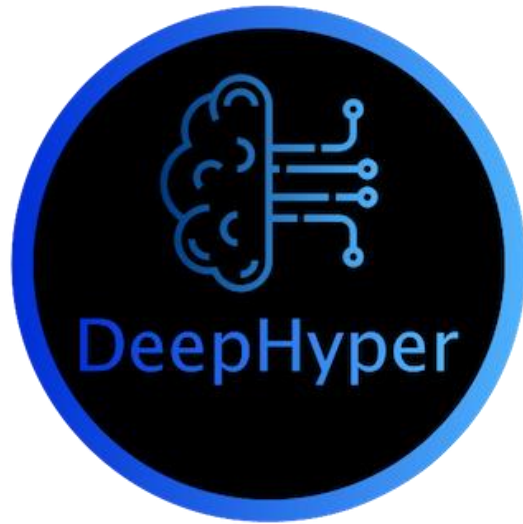
Random Search

$$p_{\epsilon} := P(|\theta - \theta^*| \leq \epsilon)$$



(b) Increasing the Number of Samples

Example with DeepHyper



<http://tinyurl.com/deephyper-grid>

Black-Box Optimization

$$\Theta^* = \arg \min_{\theta \in \Theta} f(\theta)$$

In General

The **input space** can be a mixed of Real, Discrete, Categorical parameters.

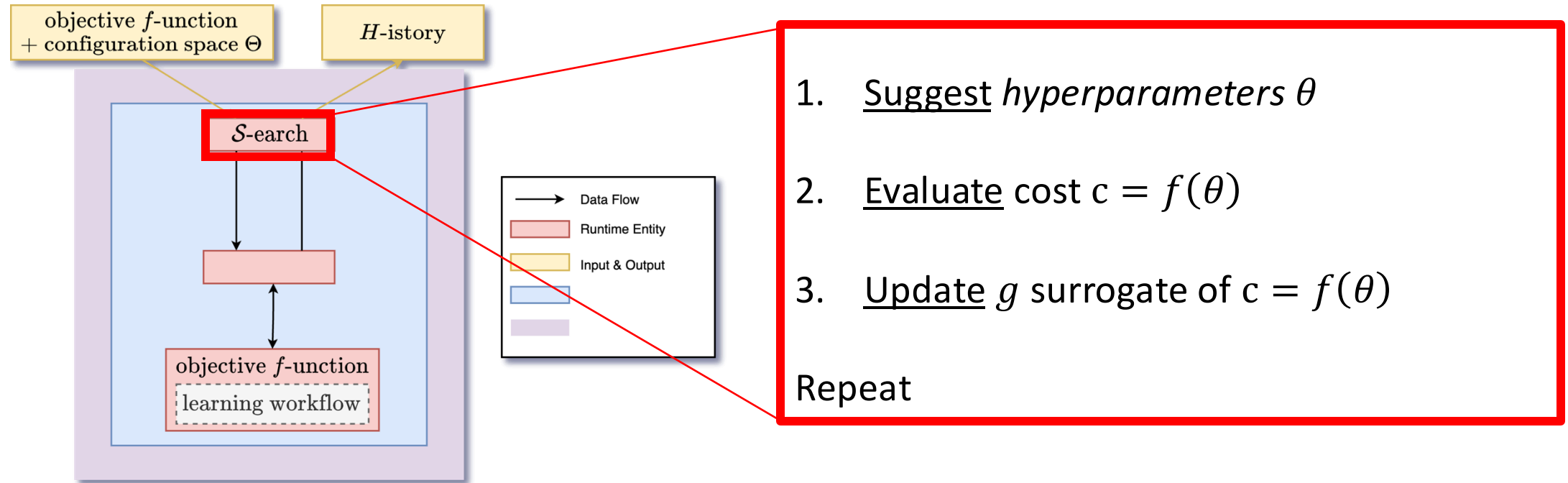
The **output space** is often mapped to real.

The **function** can be continuous, noisy, derivable.

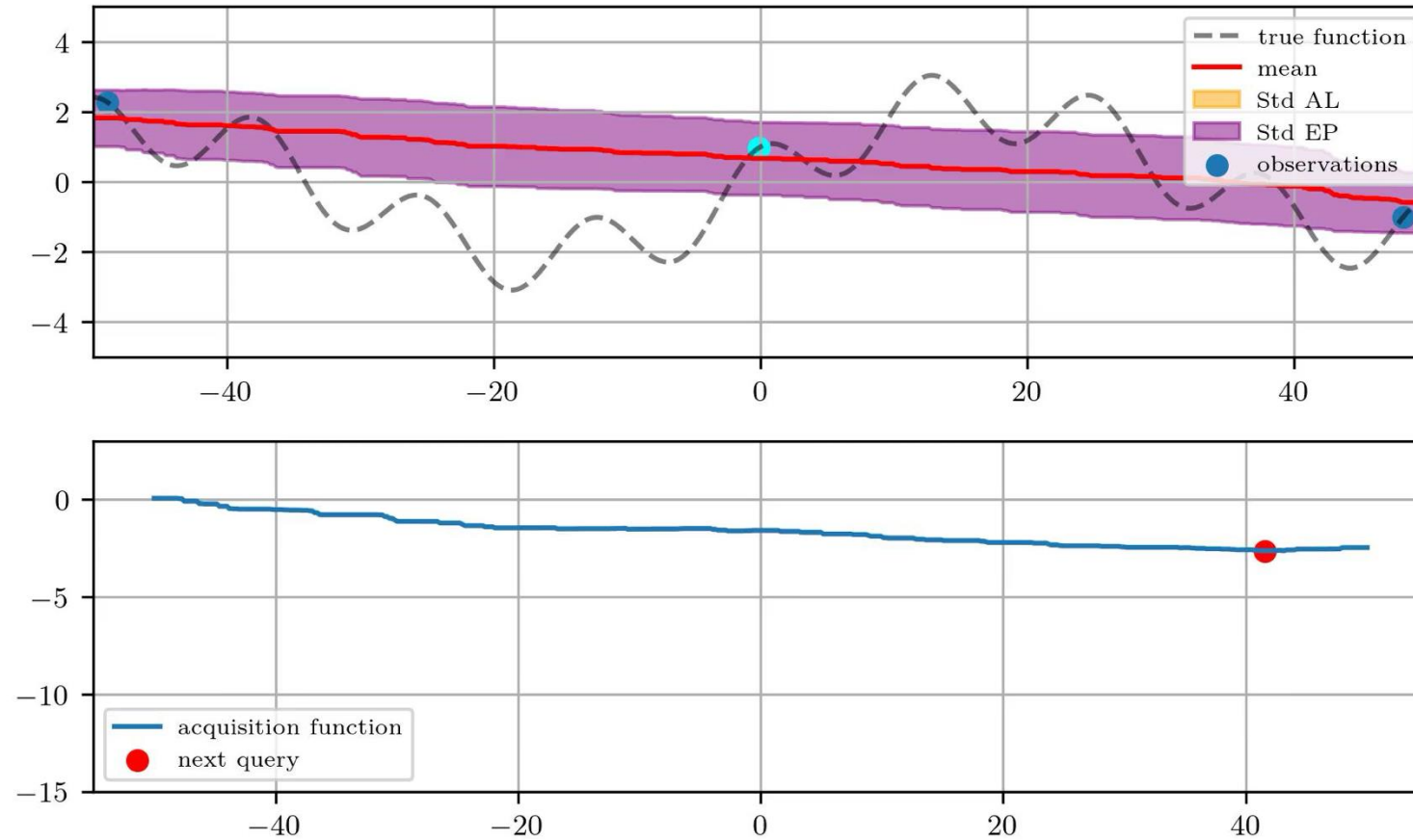
Bayesian Optimization Framework

Find hyperparameters θ to improve cost $c = f(\theta)$
of the learning workflow

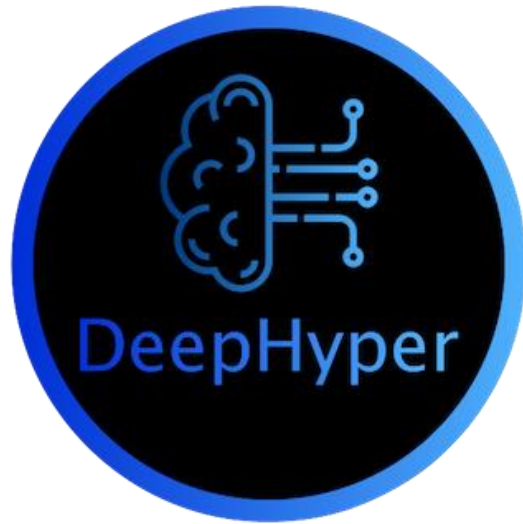
$$c^* = \min_{\theta} f(\theta)$$



An Example of Bayesian Optimization



Example with DeepHyper



<http://tinyurl.com/deephyper-bbo>

Sequential Bayesian Optimization Algorithm

Algorithm 1: Bayesian Optimization (*a.k.a.*, Efficient Global Optimization (EGO))

Inputs : *thetaSpace*: a configuration space

nInitial: the number of initial hyperparameter configurations

f: a function that returns the cost of the learning workflow

Output: *thetaStar* the recommended hyperparameter configuration.

```
1 thetaArray, costArray  $\leftarrow$  New empty arrays of hyperparameter configurations and costs ;
2 model  $\leftarrow$  New surrogate model ;
3 while stopping criteria not valid do
4   if Length of thetaArray < nInitial then
5     | theta  $\leftarrow$  Sample hyperparameter configuration from thetaSpace ;
6   else
7     | Update model with thetaArray, costArray ;
8     | theta  $\leftarrow$  Returns theta in thetaSpace that minimizes the acquisition function for current model ;
9   end
10  cost  $\leftarrow$  Returns the cost of learning workflow f(theta) ;
11  thetaArray, costArray  $\leftarrow$  Concatenate thetaArray with [theta] and costArray with [cost] ;
12  thetaStar  $\leftarrow$  Update recommendation ;
13 end
```

Acquisition functions

- Lower Confidence Bound

$$a_{\text{LCB}}(\theta; \kappa) := \mu(\theta) - \kappa \cdot \sigma(\theta)$$

- Probability of Improvement

$$\begin{aligned} I(\theta; \xi) &:= \max(f(\theta^*) - f(\theta) - \xi, 0) \\ &= \max(f(\theta^*) - \mu(\theta) - z \cdot \sigma(\theta) - \xi, 0) \quad \text{with } z \sim \mathcal{N}(0, 1) \end{aligned}$$

$$a_{\text{PI}}(\theta; \xi) := \mathbb{P}(0 < I(\theta; \xi))$$

- Expected Improvement

$$a_{\text{EI}}(\theta; \xi) := \mathbb{E}[I(\theta; \xi)]$$

Lower Confidence Bound

Acquisition function is based on uncertainty.

$$a_{LCB}(\theta; \kappa) = \overset{\text{Surrogate } g}{\mu(\theta)} - \kappa \cdot \sigma(\theta)$$

Lower-Confidence Bound

[Cox, 1992]

“Most optimistic outcome”

Simple, cheap, less flat

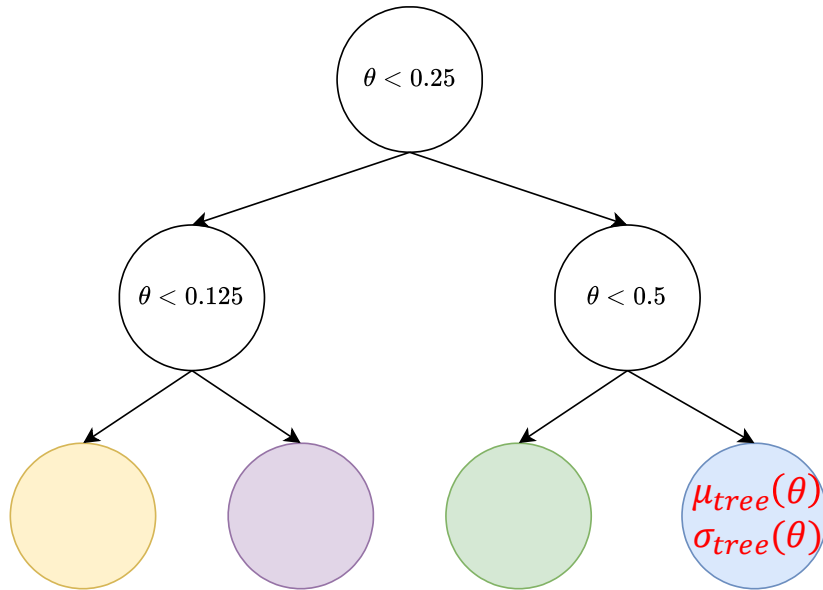
Next hyperparameters suggested have minimum LCB score

$$\theta^* = \min_{\theta} a_{LCB}(\theta; \kappa)$$

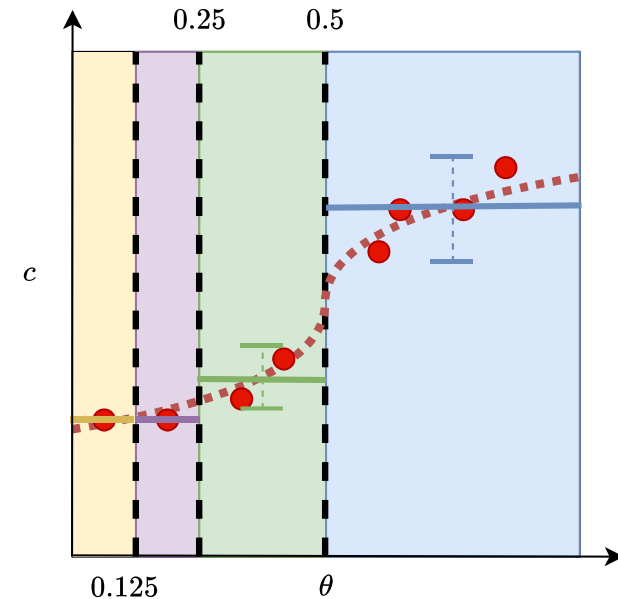
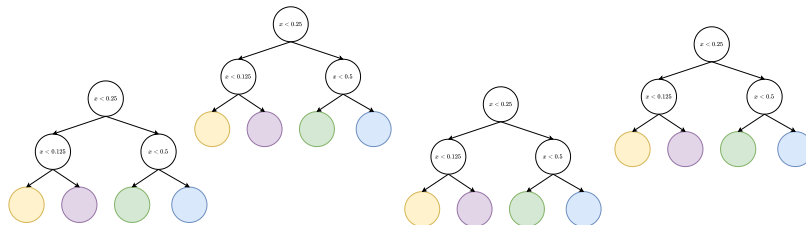
Randomized forests as model

Surrogate g is an ensemble of randomized decision trees.

Randomized
Decision Tree



Random-Forest
Surrogate g



Other Seminal Surrogates

Gaussian Process [Eriksson, 2019],
Bayesian NN [Springenberg, 2016],
Random-Forest [Hutter, 2011] (temporal
complexity and parallelizable)
TPE [Bergstra, 2013]...

Randomized forests and uncertainty

$$a_{LCB}(\theta; \kappa) = \overset{\text{Surrogate } g}{\mu(\theta)} - \kappa \cdot \sigma(\theta)$$
$$\sigma^2(\theta) = \underbrace{E_{\text{tree}}[\sigma_{\text{tree}}^2(\theta)]}_{\text{aleatoric}} + \underbrace{V_{\text{tree}}[\mu_{\text{tree}}(\theta)]}_{\text{epistemic}}$$

Randomization Effects

Algorithms

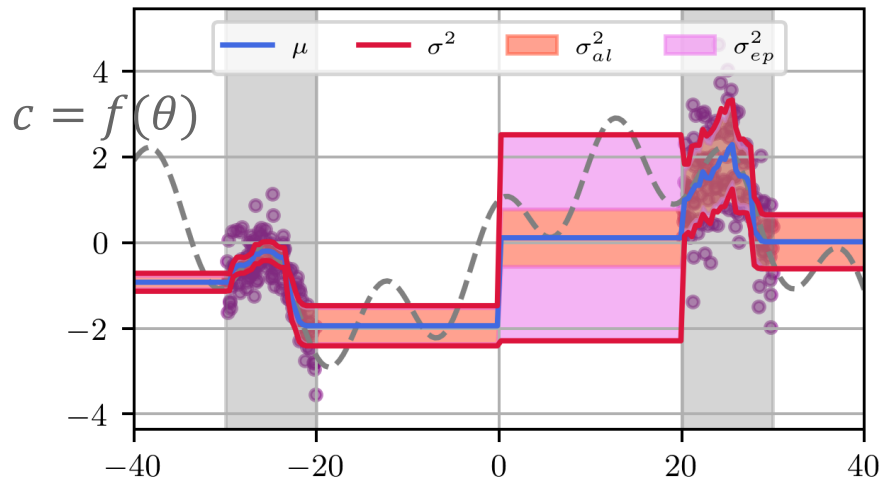
	Bootstrapping	Feature	Split
Tree Bagging (TB)	x		
Random Space (RS)		x	
Random Forest (RF)	x	x	
Extremely Randomized Trees (ET)			x
Mondrian Forest (MF)		x	x

Randomized Splits

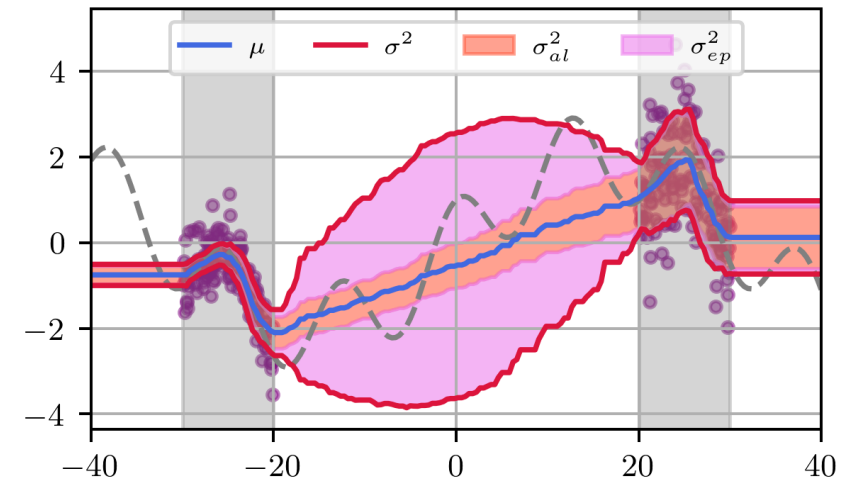
$$a_{LCB}(\theta; \kappa) = \overset{\text{Surrogate } g}{\mu(\theta)} - \kappa \cdot \sigma(\theta)$$

$$\sigma^2(\theta) = \underbrace{E_{\text{tree}}[\sigma_{\text{tree}}^2(\theta)]}_{\text{aleatoric}} + \underbrace{V_{\text{tree}}[\mu_{\text{tree}}(\theta)]}_{\text{epistemic}}$$

Best Split



Random Split



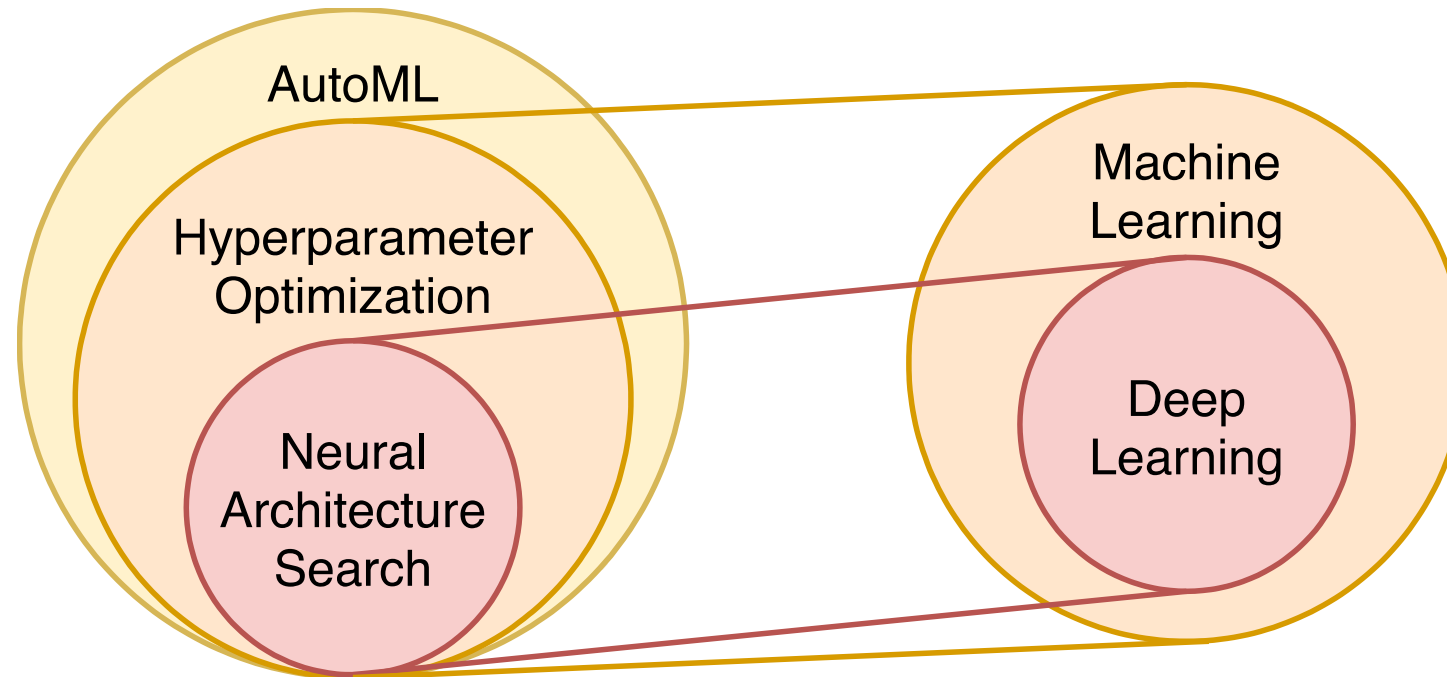
Optimization of the Acquisition Function

- Depends on the surrogate model (e.g., can we have a gradient)?
- Zero, First, Second order...

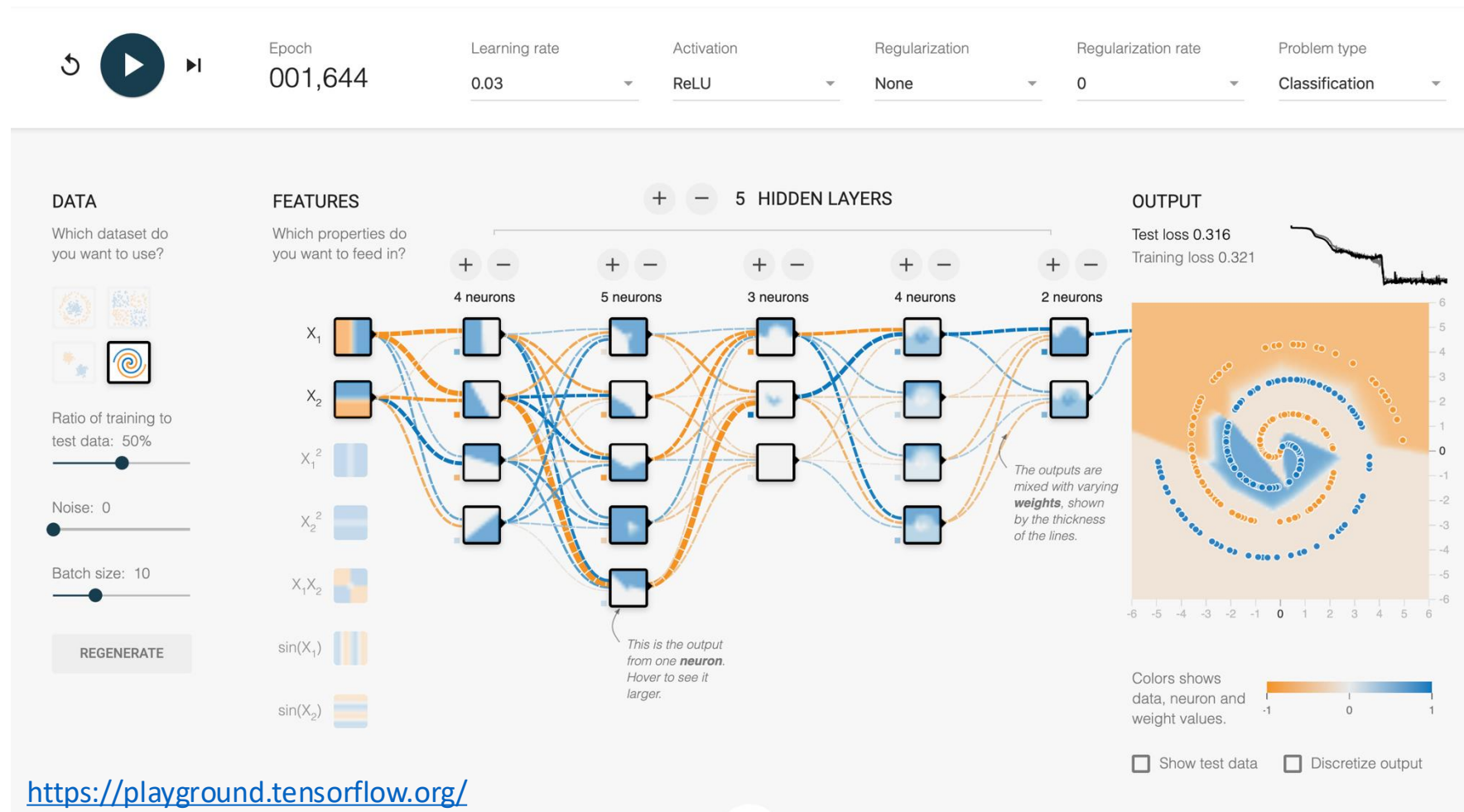
Examples:

- SGD
- LBFGS
- Genetic Algorithm
- CMAES (Covariance Matrix Adaptation Evolutionary Search)

Neural Architecture Search

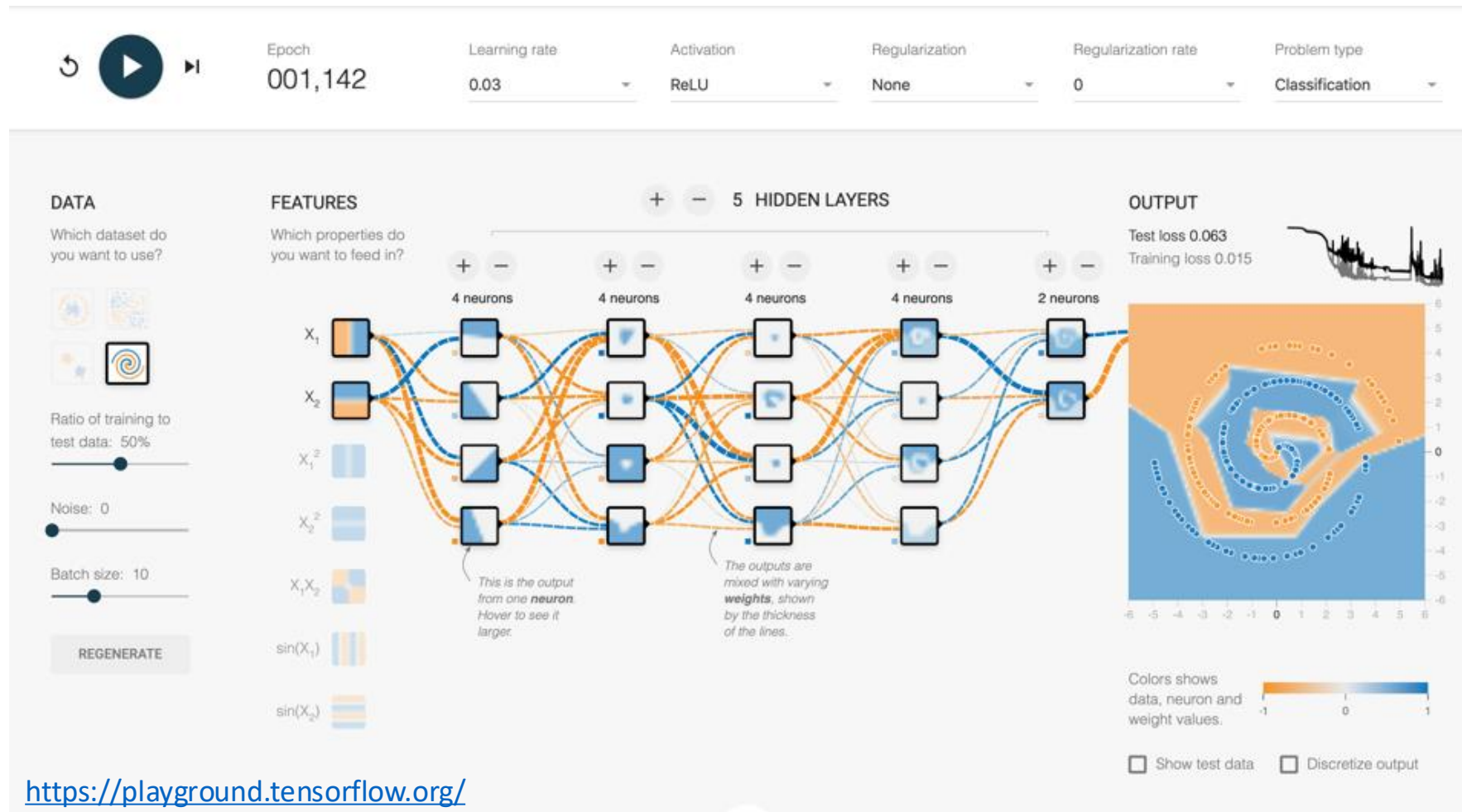


Example of sensitive neural architecture



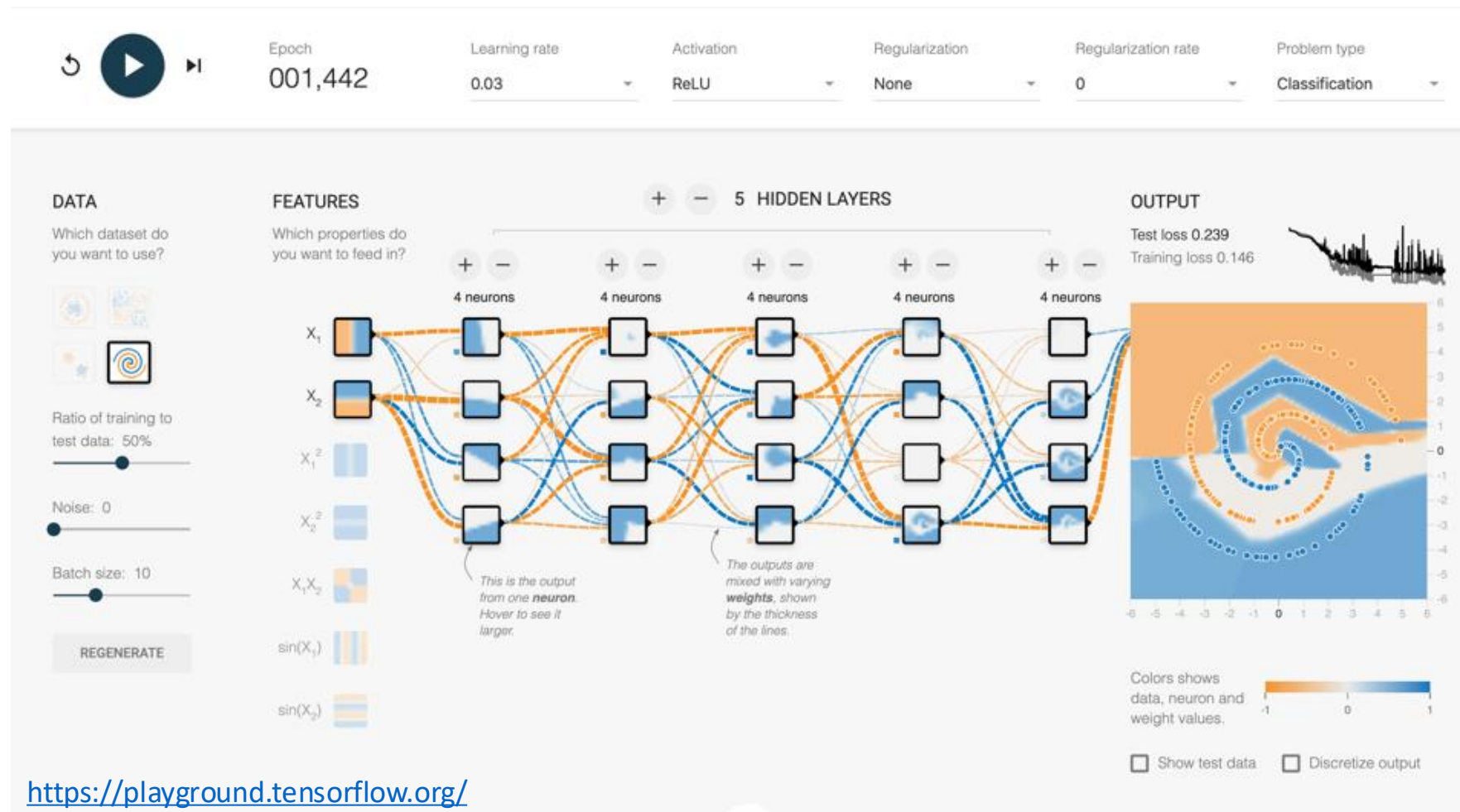
<https://playground.tensorflow.org/>

Example of sensitive neural architecture



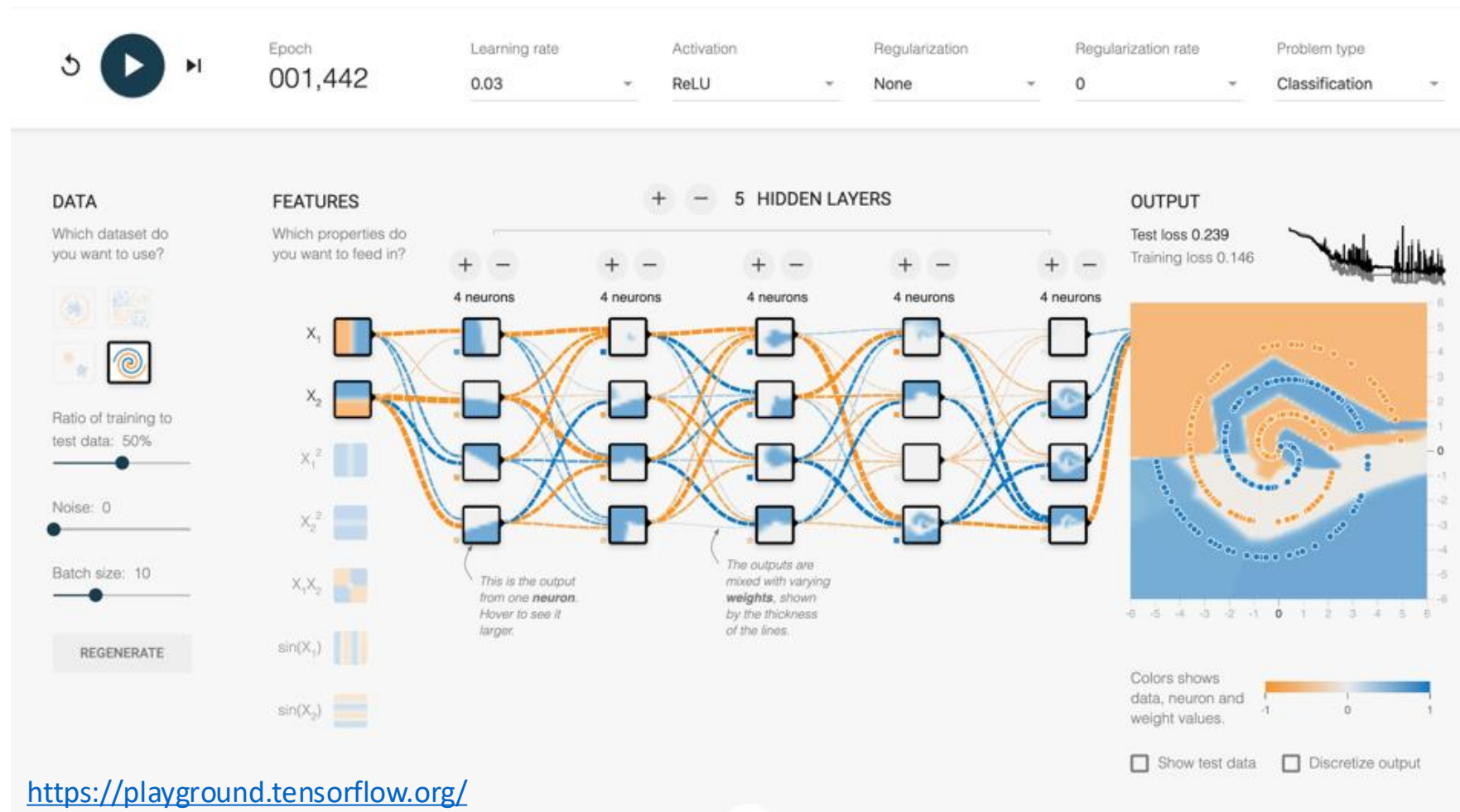
<https://playground.tensorflow.org/>

Example of sensitive neural architecture



<https://playground.tensorflow.org/>

Example of sensitive neural architecture



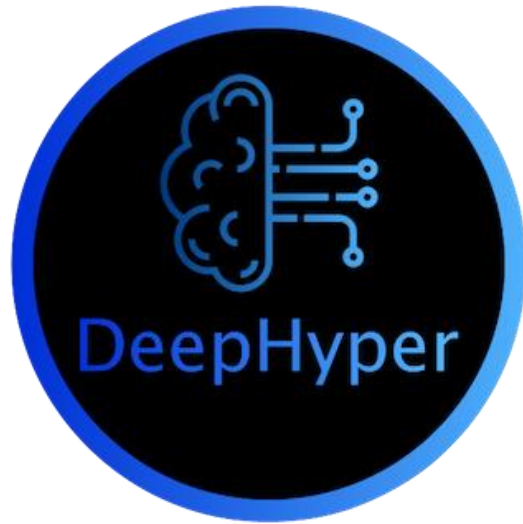
The performance of neural networks can be very sensitive with respect to the “architecture” (i.e. hyperparameters) of the neural network.

<https://playground.tensorflow.org/>

Hyperparameter Optimization with Constraints


- A layer is active if the number of layer is large enough.
- The parameters of a layer change depending on its type (dense, conv, batchnorm, dropout).
- Residual/skip connections can be created if input/output layers exists.

Example with DeepHyper



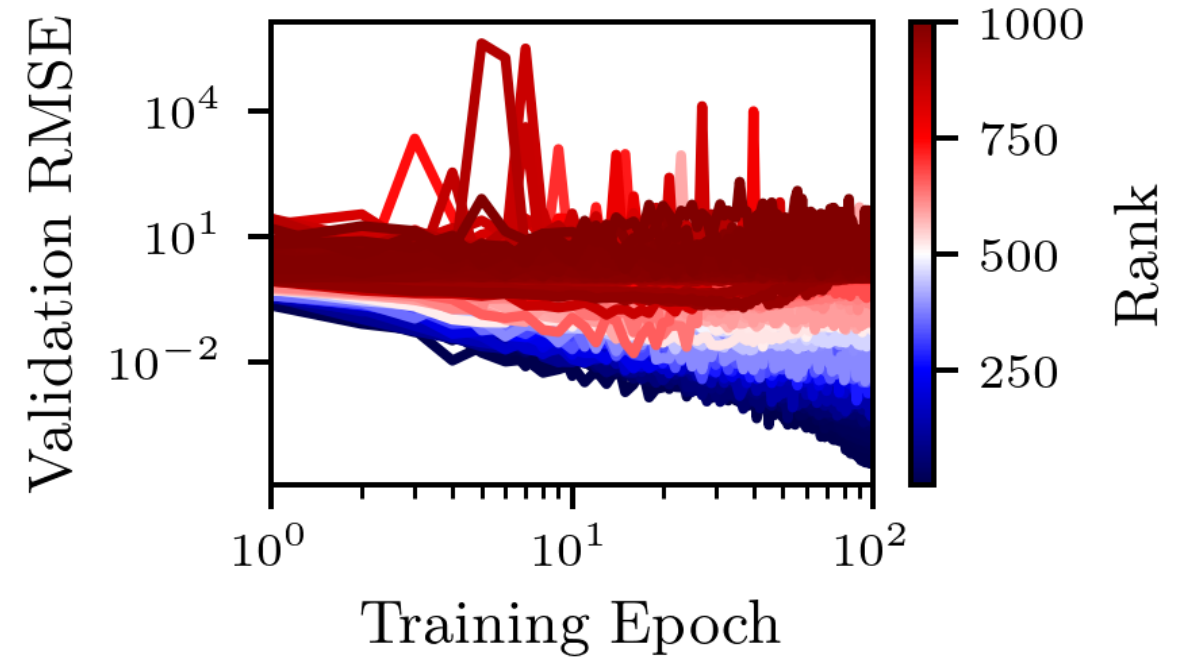
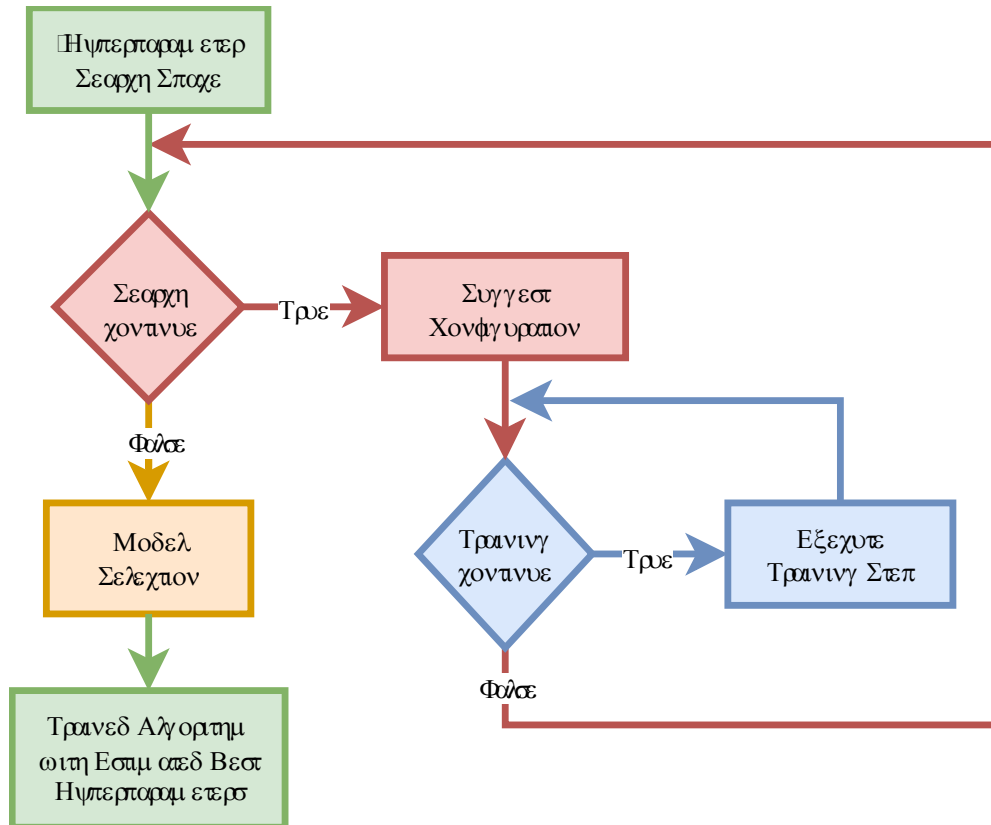
<http://tinyurl.com/deephyper-autodeuq-reg>

Overfitting in Hyperparameter Optimization

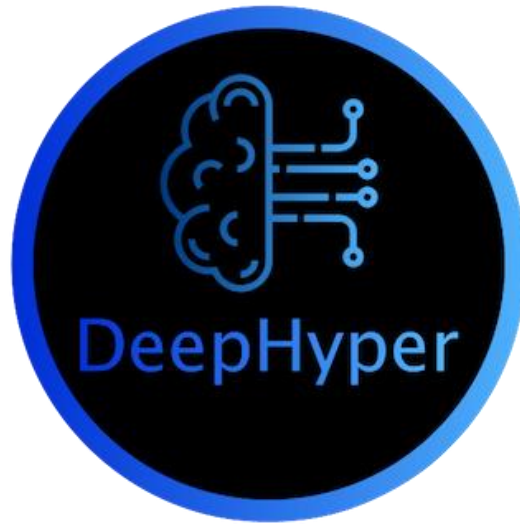
- 3-way split of the data: Training, Validation, Test
 - Training: for the weights of a neural network
 - Validation: for the hyperparameters
 - Test: generalization performance

Development Data
- Overfitting in HPO would mean that validation score improves when test score worsen.
 - Generally not observed...
 - Similar to the problem of overfitting the test set for Cifar10/ImageNet
 - Similar to the problem of development phase and final phase in machine learning competitions (Kaggle)

Early discarding



Interested to dive in?



<https://github.com/deephyper/deephyper>