

# Lectures on Machine Learning

## Lecture 2: from non-linear models to hyperparameter tune

---

Stefano Carrazza

TAE2019, 8-21 September 2019

University of Milan and INFN Milan (UNIMI)



## Lecture 1 (yesterday)

- Artificial intelligence
- Machine learning
- Model representation
- Metrics
- Parameter learning

## Lecture 2 (today)

- Non-linear models
- Beyond neural networks
- Clustering
- Cross-validation
- Hyperparameter tune

# Artificial neural networks

---

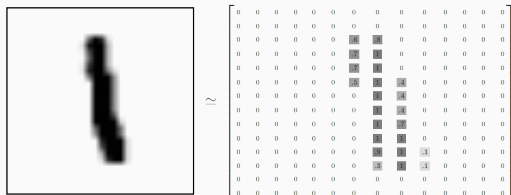
# Limitations of linear models

Why not linear models everywhere?

# Limitations of linear models

Why not linear models everywhere?

**Example:** consider 1 image from the MNIST database:



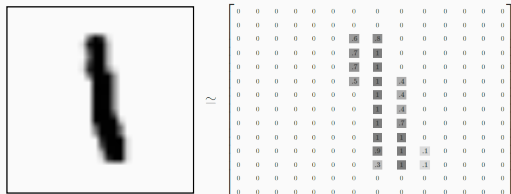
Each image has  $28 \times 28$  pixels = 785 features (x3 if including RGB colors).

If consider quadratic function  $\mathcal{O}(n^2)$  so linear models are impractical.

# Limitations of linear models

Why not linear models everywhere?

**Example:** consider 1 image from the MNIST database:

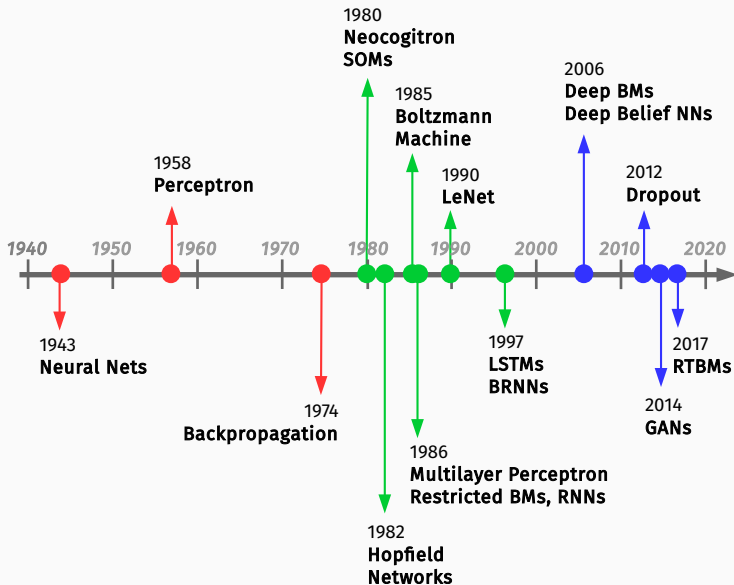


Each image has  $28 \times 28$  pixels = 785 features (x3 if including RGB colors).

If consider quadratic function  $\mathcal{O}(n^2)$  so linear models are impractical.

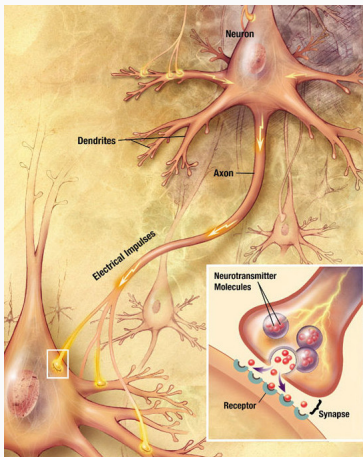
**Solution:** use non-linear models.

# Non-linear models timeline



# Neural networks

Artificial neural networks are computer systems inspired by the biological neural networks in the brain.

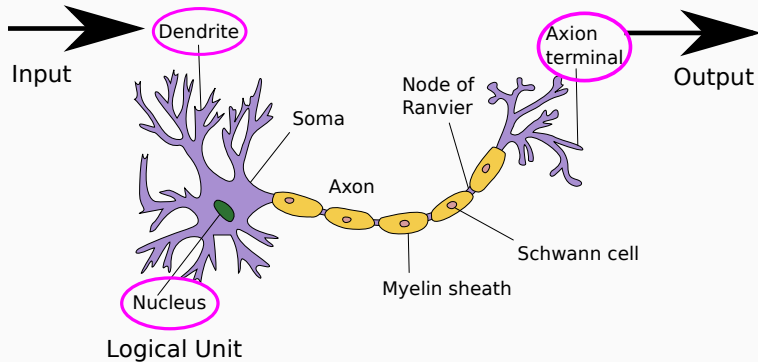


Currently the state-of-the-art technique for several ML applications.



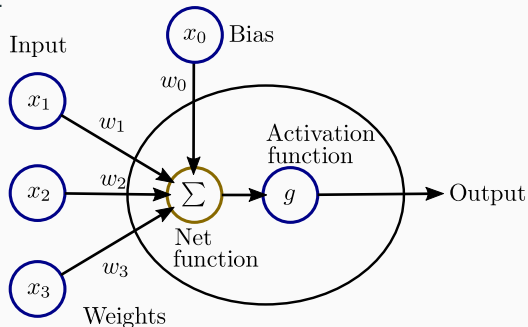
# Neuron model

We can imagine the following data communication pattern:



# Neuron model

Schematically:



where

- each **node** has an associate weights and bias  $w$  and inputs  $x$ ,
- the output is modulated by an **activation function**,  $g$ .

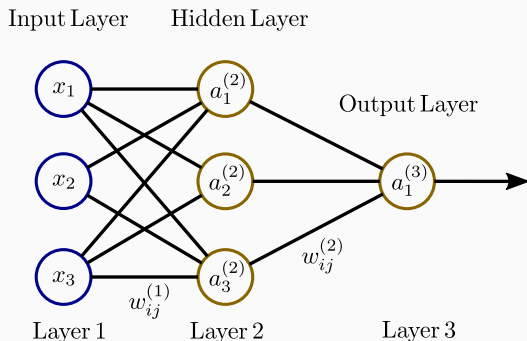
Some examples of activation functions: sigmoid, tanh, linear, ...

$$g_w(x) = \frac{1}{1 + e^{-w^T x}}, \quad \tanh(w^T x), \quad x.$$

# Neural networks

In practice, we simplify the bias term with  $x_0 = 1$ .

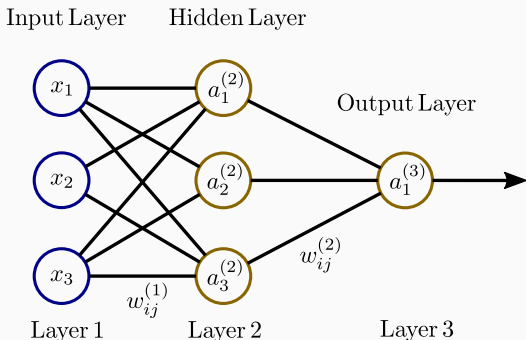
Neural network → connecting multiple units together.



where

- $a_i^{(l)}$  is the activation of unit  $i$  in layer  $l$ ,
- $w_{ij}^{(l)}$  is the weight between nodes  $i, j$  from layers  $l, l + 1$  respectively.

# Neural networks

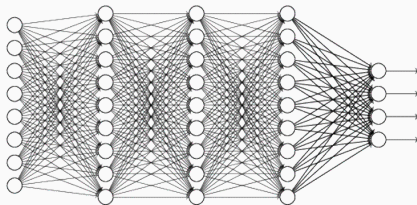


- $a_1^{(2)} = g(w_{10}^{(1)} + w_{11}^{(1)} x_1 + w_{12}^{(1)} x_2 + w_{13}^{(1)} x_3)$
- $a_2^{(2)} = g(w_{20}^{(1)} + w_{21}^{(1)} x_1 + w_{22}^{(1)} x_2 + w_{23}^{(1)} x_3)$
- $a_3^{(2)} = g(w_{30}^{(1)} + w_{31}^{(1)} x_1 + w_{32}^{(1)} x_2 + w_{33}^{(1)} x_3)$
- **Output**  $\rightarrow a_1^{(3)} = g(w_{10}^{(2)} + w_{11}^{(2)} a_1^{(2)} + w_{12}^{(2)} a_2^{(2)} + w_{13}^{(2)} a_3^{(2)})$

# Neural networks

Some useful names:

- **Feedforward neural network**: no cyclic connections between nodes from the same layer (previous example).
- **Multilayer perceptron (MLP)**: is a feedforward neural network with at least 3 layers.
- **Deep neural networks**: term referring to neural networks with more than one hidden layer.



# Training neural networks

The training NNs is usually performed with **gradient descent** methods.

Following the previous section, we have to compute the cost function gradient with respect to parameters  $w_{ij}^{(l)}$ :

$$w_{ij}^{(l)} := w_{ij}^{(l)} - \eta \nabla_{w_{ij}^{(l)}} J \quad \rightarrow \quad \nabla_{w_{ij}^{(l)}} J = \frac{\partial}{\partial w_{ij}^{(l)}} J(\mathbf{w})$$

# Training neural networks

The training NNs is usually performed with [gradient descent](#) methods.

Following the previous section, we have to compute the cost function gradient with respect to parameters  $w_{ij}^{(l)}$ :

$$w_{ij}^{(l)} := w_{ij}^{(l)} - \eta \nabla_{w_{ij}^{(l)}} J \quad \rightarrow \quad \nabla_{w_{ij}^{(l)}} J = \frac{\partial}{\partial w_{ij}^{(l)}} J(\mathbf{w})$$

Use the [backpropagation algorithm](#) to compute the gradient of a NN.

- can be used with any gradient-based optimizer, including quasi-Newton methods.
- reduces the large amount of computations thanks to chain rule
- requires the derivative of the cost function with respect to the output layer  $w_{ij}^{(l)}$  with  $l = \text{output}$ .

# Backpropagation algorithm

The backpropagation steps:

- 1: perform a **forward propagation** (calculate  $a_i^{(l)}$ )



# Backpropagation algorithm

The backpropagation steps:

- 1: perform a **forward propagation** (calculate  $a_i^{(l)}$ )
- 2: perform a **backward propagation**: evaluate for each node a “prediction error”:

$$\delta_j^{(l)} = \text{“error” of node } j \text{ in layer } l.$$

# Backpropagation algorithm

The backpropagation steps:

- 1: perform a **forward propagation** (calculate  $a_i^{(l)}$ )
- 2: perform a **backward propagation**: evaluate for each node a “prediction error”:

$$\delta_j^{(l)} = \text{“error” of node } j \text{ in layer } l.$$

- 3: calculate  $\nabla_{i,j}^{(l)} J$  using errors  $\delta_i^{(l)}$  and  $a_i^{(l)}$ .

# Backpropagation algorithm

The backpropagation steps:

- 1: perform a **forward propagation** (calculate  $a_i^{(l)}$ )
- 2: perform a **backward propagation**: evaluate for each node a “prediction error”:

$$\delta_j^{(l)} = \text{“error” of node } j \text{ in layer } l.$$

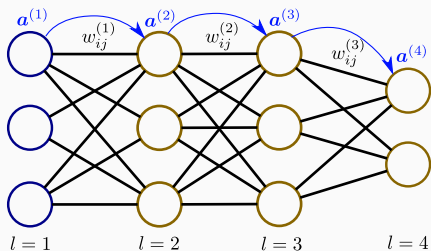
- 3: calculate  $\nabla_{ij}^{(l)} J$  using errors  $\delta_i^{(l)}$  and  $a_i^{(l)}$ .
- 4: perform **weight updates**,  $\Delta w_{ij}^{(l)}$ , via gradient descent using  $\nabla_{ij}^{(l)} J$ .

# Backpropagation algorithm

Suppose we have a MLP and one training example  $(\mathbf{x}, \mathbf{y})$ .

**Step 1:** We first perform a **forward propagation pass**:

- $\mathbf{a}^{(1)} = \mathbf{x}$

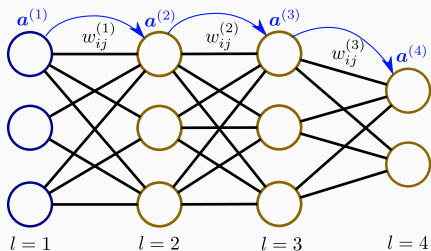


# Backpropagation algorithm

Suppose we have a MLP and one training example  $(\mathbf{x}, \mathbf{y})$ .

**Step 1:** We first perform a **forward propagation pass**:

- $\mathbf{a}^{(1)} = \mathbf{x}$
- $\mathbf{z}^{(2)} = \mathbf{w}^{(1)}\mathbf{a}^{(1)}$

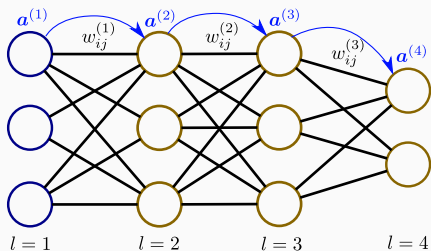


# Backpropagation algorithm

Suppose we have a MLP and one training example  $(\mathbf{x}, \mathbf{y})$ .

**Step 1:** We first perform a **forward propagation pass**:

- $\mathbf{a}^{(1)} = \mathbf{x}$
- $\mathbf{z}^{(2)} = \mathbf{w}^{(1)} \mathbf{a}^{(1)}$
- $\mathbf{a}^{(2)} = g(\mathbf{z}^{(2)})$

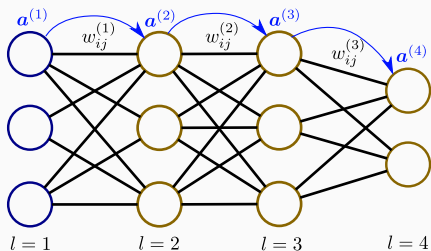


# Backpropagation algorithm

Suppose we have a MLP and one training example  $(\mathbf{x}, \mathbf{y})$ .

**Step 1:** We first perform a **forward propagation pass**:

- $\mathbf{a}^{(1)} = \mathbf{x}$
- $\mathbf{z}^{(2)} = \mathbf{w}^{(1)} \mathbf{a}^{(1)}$
- $\mathbf{a}^{(2)} = g(\mathbf{z}^{(2)})$
- $\mathbf{z}^{(3)} = \mathbf{w}^{(2)} \mathbf{a}^{(2)}$

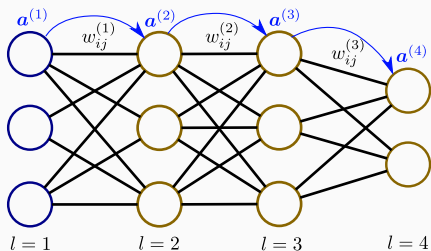


# Backpropagation algorithm

Suppose we have a MLP and one training example  $(\mathbf{x}, \mathbf{y})$ .

**Step 1:** We first perform a **forward propagation pass**:

- $\mathbf{a}^{(1)} = \mathbf{x}$
- $\mathbf{z}^{(2)} = \mathbf{w}^{(1)} \mathbf{a}^{(1)}$
- $\mathbf{a}^{(2)} = g(\mathbf{z}^{(2)})$
- $\mathbf{z}^{(3)} = \mathbf{w}^{(2)} \mathbf{a}^{(2)}$
- $\mathbf{a}^{(3)} = g(\mathbf{z}^{(3)})$



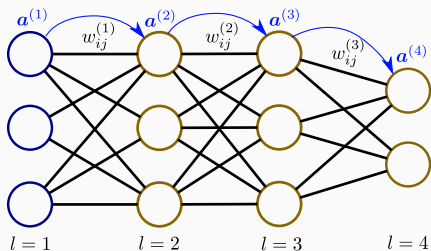


# Backpropagation algorithm

Suppose we have a MLP and one training example  $(\mathbf{x}, \mathbf{y})$ .

**Step 1:** We first perform a **forward propagation pass**:

- $\mathbf{a}^{(1)} = \mathbf{x}$
- $\mathbf{z}^{(2)} = \mathbf{w}^{(1)} \mathbf{a}^{(1)}$
- $\mathbf{a}^{(2)} = g(\mathbf{z}^{(2)})$
- $\mathbf{z}^{(3)} = \mathbf{w}^{(2)} \mathbf{a}^{(2)}$
- $\mathbf{a}^{(3)} = g(\mathbf{z}^{(3)})$
- $\mathbf{z}^{(4)} = \mathbf{w}^{(3)} \mathbf{a}^{(3)}$

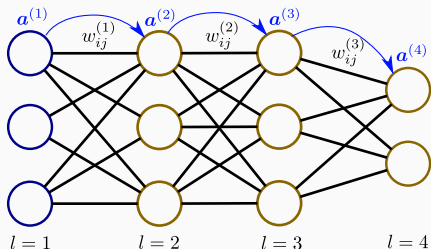


# Backpropagation algorithm

Suppose we have a MLP and one training example  $(\mathbf{x}, \mathbf{y})$ .

**Step 1:** We first perform a **forward propagation pass**:

- $\mathbf{a}^{(1)} = \mathbf{x}$
- $\mathbf{z}^{(2)} = \mathbf{w}^{(1)} \mathbf{a}^{(1)}$
- $\mathbf{a}^{(2)} = g(\mathbf{z}^{(2)})$
- $\mathbf{z}^{(3)} = \mathbf{w}^{(2)} \mathbf{a}^{(2)}$
- $\mathbf{a}^{(3)} = g(\mathbf{z}^{(3)})$
- $\mathbf{z}^{(4)} = \mathbf{w}^{(3)} \mathbf{a}^{(3)}$
- Output  $\mathbf{a}^{(4)} = g(\mathbf{z}^{(4)})$



At this step we know the output of the current MLP setup.

# Backpropagation algorithm

2. evaluate for each node the error  $\delta_j^{(k)}$  for  $k = 2, 3, \dots, L$ .

## Some remarks:

It is possible to proof using derivative chain rules that:

$$\nabla_{ij}^{(l)} J = \frac{\partial J}{\partial z_i^{(l+1)}} a_j^{(l)} \equiv \delta_i^{(l+1)} a_j^{(l)},$$

for  $l = 1, \dots, L - 1$ .

# Backpropagation algorithm

2. evaluate for each node the error  $\delta_j^{(k)}$  for  $k = 2, 3, \dots, L$ .

## Some remarks:

It is possible to proof using derivative chain rules that:

$$\nabla_{ij}^{(l)} J = \frac{\partial J}{\partial z_i^{(l+1)}} a_j^{(l)} \equiv \delta_i^{(l+1)} a_j^{(l)},$$

for  $l = 1, \dots, L - 1$ .

The recursive relation for the error is:

$$\delta_i^{(l)} = \sum_k w_{ki}^{(l)} \delta_k^{(l+1)} \cdot g'(z_i^{(l)})$$

and at  $l = L$ , i.e. the highest  $l$  index:

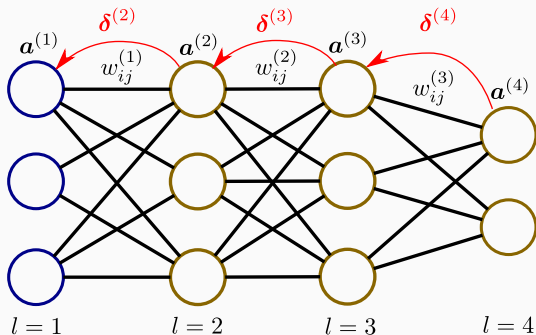
$$\delta_i^{(L)} = \frac{\partial J}{\partial a_i^{(L)}} \cdot g'(z_i^{(L)})$$

where  $g'(z_i^{(l)}) = a_i^{(l)}(1 - a_i^{(l)})$  if  $g$  is the sigmoid function.

# Backpropagation algorithm

**Example:** evaluating error  $\delta_j^{(l)}$  for a MLP with sigmoids in the hidden layers and linear activation function in the output layer:

- $\delta^{(4)} = \mathbf{a}^{(4)} - \mathbf{y}$
- $\delta^{(3)} = (\mathbf{w}^{(3)})^T \delta^{(4)} \cdot (\mathbf{a}^{(3)}(1 - \mathbf{a}^{(3)}))$
- $\delta^{(2)} = (\mathbf{w}^{(2)})^T \delta^{(3)} \cdot (\mathbf{a}^{(2)}(1 - \mathbf{a}^{(2)}))$



# Backpropagation algorithm summary

**Data:** training set  $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$  with  $i = 1, \dots, m$  examples.

**Result:** the trained neural network

Initialize network weights;

**while** *stopping criterion is not satisfied* **do**

Set all  $\Delta w_{ij}^{(l)} = 0$ .

**for**  $k = 1$  **to**  $m$  **do**

Perform forward pass and compute  $\mathbf{a}^{(l)}$  for  $l = 1, 2, 3, \dots, L$ ;

Perform backward pass and compute  $\delta^{(l)}$  for  $l = 2, \dots, L$ ;

$\Delta w_{ij}^{(l)} := \Delta w_{ij}^{(l)} + a_j^l \delta_i^{(l+1)}$

**end**

Update network weights using gradient descent;

**end**

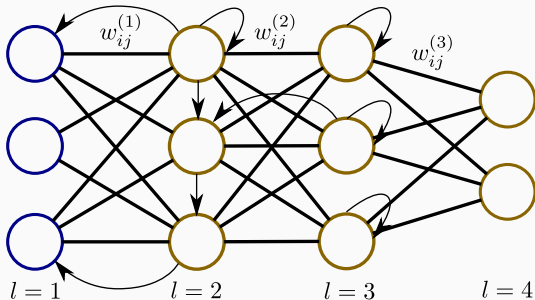
Some remarks and example of neural network initialization:

- **zero**: all weights are set to zero so all neurons perform the same calculation. The complexity of the neural network is equivalent to a single neuron.
- **random**: breaks parameter symmetry.
- **glorot/xavier**: initialize each weight with a small Gaussian value with mean zero and variance based on the in/out size of the weight.
- **he**: avoid activation function saturation. Weights are random initialized considering the size of the previous layer.

# Artificial neural networks architectures

Some examples of neural network popular architectures:

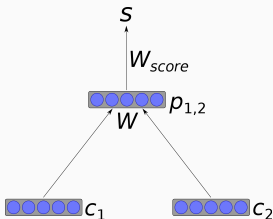
- **Recurrent neural networks:** neural networks where connections between nodes form a directed cycle.
  - built-in internal state memory
  - built-in notion of time ordering for a time sequence





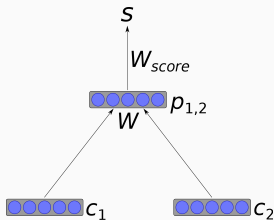
# Artificial neural networks architectures

- **Recursive neural networks**: a variation of recurrent neural network where pairs of layers or nodes are merged recursively.
  - successful applications on natural language processing.
  - some recent applications for model inference.



# Artificial neural networks architectures

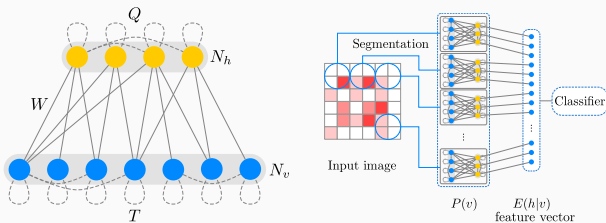
- **Recursive neural networks**: a variation of recurrent neural network where pairs of layers or nodes are merged recursively.
  - successful applications on natural language processing.
  - some recent applications for model inference.



- **Long short-term memory**: another variation of recurrent neural networks composed by custom units cells:
  - LSTM cells have an input gate, an output gate and a forget gate.
  - powerful when making predictions based on time series data.

# Artificial neural networks architectures

- **Boltzmann Machines**: is a generative stochastic recursive artificial neural network.
  - comes with energy-based model features and advantages.
  - generalizations like RBMs can be used for pdf estimate, filtering, regression, classification and sampling.

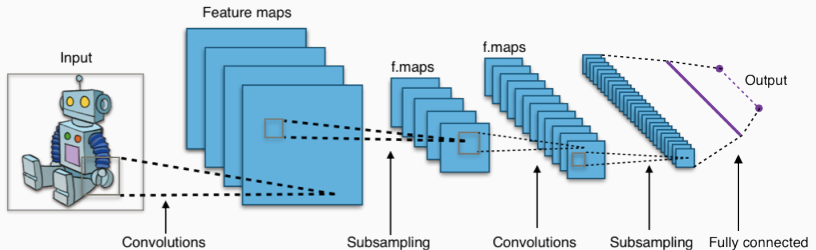


The system energy for given state vectors  $(v, h)$ :

$$E(v, h) = \frac{1}{2}v^tTv + \frac{1}{2}h^tQh + v^tWh + B_hh + Bvv$$

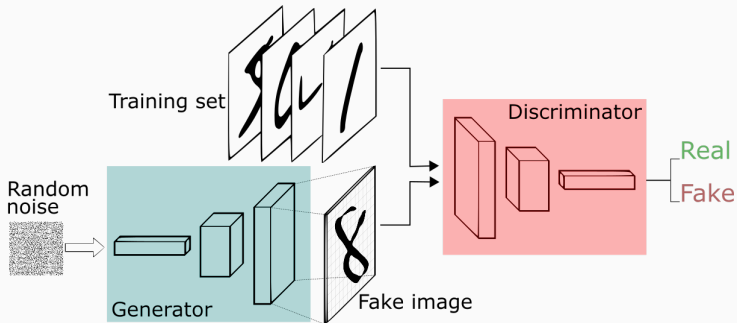
# Artificial neural networks architectures

- **Convolutional neural networks:** multilayer perceptron designed to require minimal preprocessing, *i.e.* space invariant architecture.
  - the hidden layers consist of convolutional layers, pooling layer, fully connected layers and normalization layers
  - great successful applications in image and video recognition.



# Artificial neural networks architectures

- **Generative adversarial network:** unsupervised machine learning system of two neural networks contesting with each other.
  - one network generate candidates while the other discriminates.



# Beyond neural networks

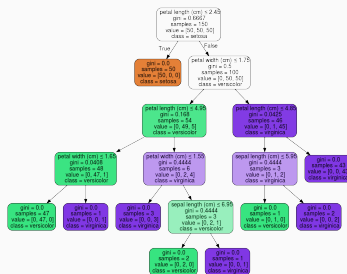
---

# Beyond neural networks

Even if neural networks are the most popular architecture nowadays employed in ML and Deep Learning, there are other models and techniques that are used frequently with great success in HEP-EXP

## Supervised learning examples:

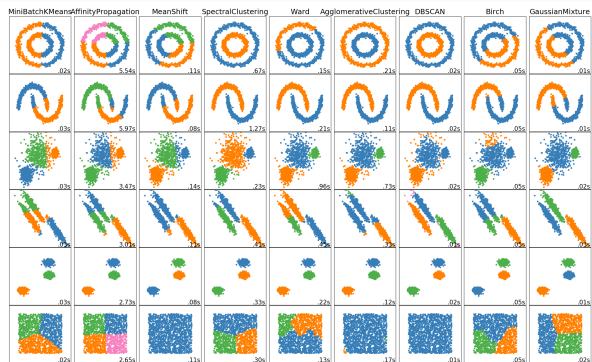
- Decision tree
- Ensemble models (random forest, bagging, boosting)
- Support Vector Machines (SVM)
- $k$ -nearest neighbors algorithm ( $k$ -NN)



# Clustering

## Unsupervised learning examples:

- $k$ -means
- Mean-shift
- Hierarchical
- Gaussian mixture models
- Density-based spatial
- Affinity propagation



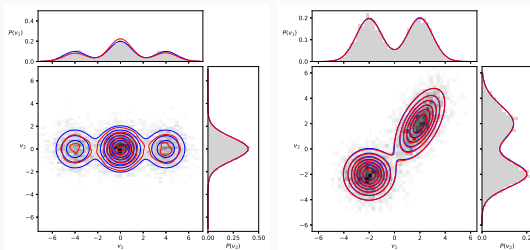


## Dimensionality reduction:

- Principal component analysis (PCA)
- Linear discriminant analysis (LDA)

## Anomaly detection:

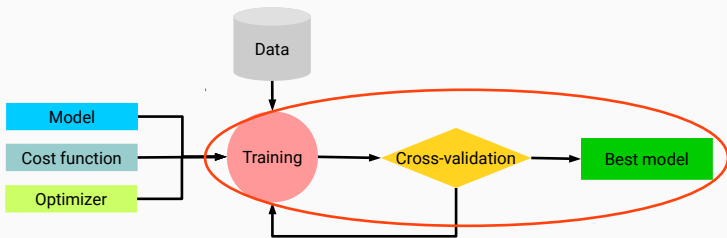
- GMM density estimate
- Kernel density estimate
- Restricted boltzmann machines
- $k$ -NN



# Hyperparameter tune

---

# Outline



# Hyperparameters summary

So far we have encountered the following hyperparameters:

- Model related:

# Hyperparameters summary

So far we have encountered the following hyperparameters:

- **Model related:**
  - model architecture / size

# Hyperparameters summary

So far we have encountered the following hyperparameters:

- **Model related:**
  - model architecture / size
  - if NN: layers, nodes, activation functions

# Hyperparameters summary

So far we have encountered the following hyperparameters:

- **Model related:**
  - model architecture / size
  - if NN: layers, nodes, activation functions
- **Regularization techniques** (to avoid overfitting):

# Hyperparameters summary

So far we have encountered the following hyperparameters:

- **Model related:**
  - model architecture / size
  - if NN: layers, nodes, activation functions
- **Regularization techniques** (to avoid overfitting):
  - weight decay (parameter  $\lambda$ )



# Hyperparameters summary

So far we have encountered the following hyperparameters:

- **Model related:**
  - model architecture / size
  - if NN: layers, nodes, activation functions
- **Regularization techniques** (to avoid overfitting):
  - weight decay (parameter  $\lambda$ )
  - if SGD: early stopping techniques

# Hyperparameters summary

So far we have encountered the following hyperparameters:

- **Model related:**
  - model architecture / size
  - if NN: layers, nodes, activation functions
- **Regularization techniques** (to avoid overfitting):
  - weight decay (parameter  $\lambda$ )
  - if SGD: early stopping techniques
  - if NN: dropout, artificial data, stochastic pooling, etc...

# Hyperparameters summary

So far we have encountered the following hyperparameters:

- **Model related:**
  - model architecture / size
  - if NN: layers, nodes, activation functions
- **Regularization techniques** (to avoid overfitting):
  - weight decay (parameter  $\lambda$ )
  - if SGD: early stopping techniques
  - if NN: dropout, artificial data, stochastic pooling, etc...
- **Training:**

# Hyperparameters summary

So far we have encountered the following hyperparameters:

- **Model related:**
  - model architecture / size
  - if NN: layers, nodes, activation functions
- **Regularization techniques** (to avoid overfitting):
  - weight decay (parameter  $\lambda$ )
  - if SGD: early stopping techniques
  - if NN: dropout, artificial data, stochastic pooling, etc...
- **Training:**
  - the SGD learning parameters  $\eta$

# Hyperparameters summary

So far we have encountered the following hyperparameters:

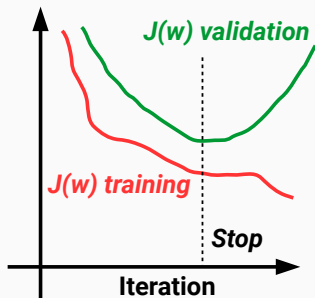
- **Model related:**
  - model architecture / size
  - if NN: layers, nodes, activation functions
- **Regularization techniques** (to avoid overfitting):
  - weight decay (parameter  $\lambda$ )
  - if SGD: early stopping techniques
  - if NN: dropout, artificial data, stochastic pooling, etc...
- **Training:**
  - the SGD learning parameters  $\eta$
  - other SGD parameters depending on the gradient descent scheme

## Other regularization techniques

**Example:** early stopping techniques are applied when the optimization is performed in an iterative procedure, .e.g. via gradient descent.

## Other regularization techniques

**Example:** early stopping techniques are applied when the optimization is performed in an iterative procedure, .e.g. via gradient descent.



These techniques monitor the cost function for the validation set and stop when this quantity has stopped improving:

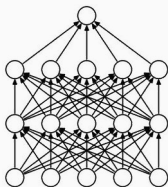
- look at the variation in a moving window
- stop at the minimum of the validation set (lookback method),

# Other regularization techniques

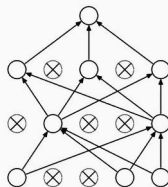
## Example: neural network dropout

At each training stage:

- individual nodes and related incoming and outgoing edges are dropped-out of the neural network with a fixed probability.
- the reduced NN is trained on the data.
- the removed nodes are reinserted in the NN with their original weights.



(a) Standard Neural Net



(b) After applying dropout.



## Other regularization techniques

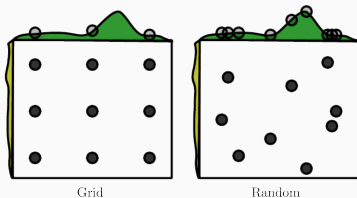
How should we proceed with hyperparameter tune?

# Other regularization techniques

## How should we proceed with hyperparameter tune?

Possible solutions:

- **grid search**: exhaustive searching through a manually subset range of hyperparameter space
- **random search**: specially powerful with small number of hyperparameters affects the final performance of the ML algorithm.

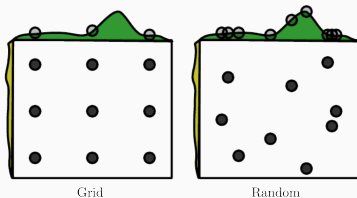


# Other regularization techniques

## How should we proceed with hyperparameter tune?

Possible solutions:

- **grid search**: exhaustive searching through a manually subset range of hyperparameter space
- **random search**: specially powerful with small number of hyperparameters affects the final performance of the ML algorithm.



Other useful methods:

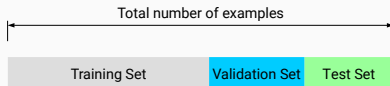
- bayesian optimization
- gradient-based optimization
- evolutionary optimization

# Cross-validation

---

# Cross-validation

The hyperparameter tune procedure still requires the training/validation/test split to choose for the best model.

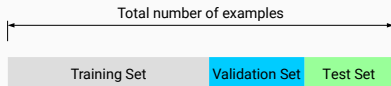


## Problems:

- how to perform the data split when the available data set is small?
- how to define a suitable split?

# Cross-validation

The hyperparameter tune procedure still requires the training/validation/test split to choose for the best model.



## Problems:

- how to perform the data split when the available data set is small?
- how to define a suitable split?

## Solution:

Use [cross-validation](#) algorithms to assess the quality of your model + hyperparameter choice.

# Cross-validation

Cross-validation performs a rotation estimation by:

1. **partitioning data** into **training/validation** subsets

# Cross-validation

Cross-validation performs a rotation estimation by:

1. **partitioning data** into **training/validation** subsets
2. multiple rounds of cross-validation using different partitions



# Cross-validation

Cross-validation performs a rotation estimation by:

1. **partitioning data** into **training/validation** subsets
2. multiple rounds of cross-validation using different partitions
3. results are averaged over the rounds to give an estimate of the model performance

# Cross-validation

Cross-validation performs a rotation estimation by:

1. **partitioning data** into **training/validation** subsets
2. multiple rounds of cross-validation using different partitions
3. results are averaged over the rounds to give an estimate of the model performance

Common approaches to cross-validation:

- **Exhaustive cross-validation**: test all possible ways to divide the original sample into a training and a validation set.

# Cross-validation

Cross-validation performs a rotation estimation by:

1. **partitioning data** into **training/validation** subsets
2. multiple rounds of cross-validation using different partitions
3. results are averaged over the rounds to give an estimate of the model performance

Common approaches to cross-validation:

- **Exhaustive cross-validation**: test all possible ways to divide the original sample into a training and a validation set.
  - Leave- $p$ -out: uses  $p$  observations as validation set.

# Cross-validation

Cross-validation performs a rotation estimation by:

1. **partitioning data** into **training/validation** subsets
2. multiple rounds of cross-validation using different partitions
3. results are averaged over the rounds to give an estimate of the model performance

Common approaches to cross-validation:

- **Exhaustive cross-validation**: test all possible ways to divide the original sample into a training and a validation set.
  - Leave- $p$ -out: uses  $p$  observations as validation set.
  - Leave-one-out: set  $p = 1$ .

# Cross-validation

Cross-validation performs a rotation estimation by:

1. **partitioning data** into **training/validation** subsets
2. multiple rounds of cross-validation using different partitions
3. results are averaged over the rounds to give an estimate of the model performance

Common approaches to cross-validation:

- **Exhaustive cross-validation**: test all possible ways to divide the original sample into a training and a validation set.
  - Leave- $p$ -out: uses  $p$  observations as validation set.
  - Leave-one-out: set  $p = 1$ .
- **Non-exhaustive cross-validation**: do not test all possible ways to divide the original sample but use discrete subsamples.

# Cross-validation

Cross-validation performs a rotation estimation by:

1. **partitioning data** into **training/validation** subsets
2. multiple rounds of cross-validation using different partitions
3. results are averaged over the rounds to give an estimate of the model performance

Common approaches to cross-validation:

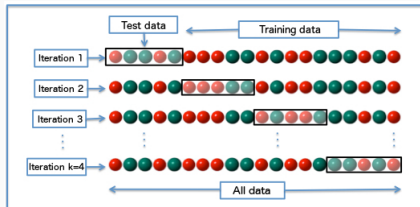
- **Exhaustive cross-validation**: test all possible ways to divide the original sample into a training and a validation set.
  - Leave- $p$ -out: uses  $p$  observations as validation set.
  - Leave-one-out: set  $p = 1$ .
- **Non-exhaustive cross-validation**: do not test all possible ways to divide the original sample but use discrete subsamples.
  - $k$ -fold cross-validation.

# Example k-fold cross-validation

## $k$ -fold cross-validation:

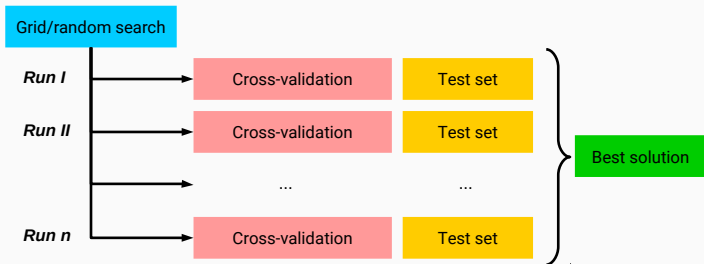
1. the original data is randomly partitioned into  $k$  equal sized subsamples.
2. from the  $k$  subsamples, a single subsample is used as validation data and the remaining  $k - 1$  subsamples are used as training data.
3. repeat the process  $k$  times by changing the validation and training partitions.
4. compute the average over the  $k$  results.

## Example of k-fold with $k = 4$ :



# Complete recipe

Perform hyperparameter tune coupled to cross-validation:



Easy parallelization at search and cross-validation stages.

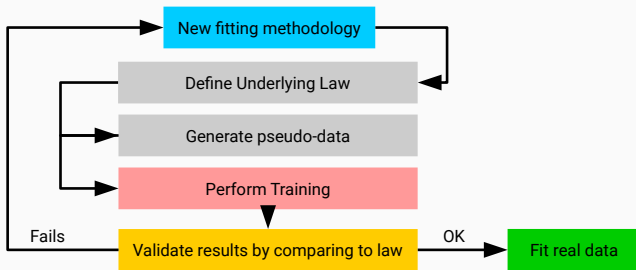


# Closure testing

---

# Closure tests

Validation and optimization of fitting strategy performed on **closure test** with known underlying law.



# ML in practice

---

# Most popular public ML frameworks

## For experimental HEP:

- TMVA: ROOT's builtin machine learning package.

# Most popular public ML frameworks

## For experimental HEP:

- TMVA: ROOT's builtin machine learning package.

## For ML applications:

- Keras: a Python deep learning library.
- Theano: a Python library for optimization.
- PyTorch: a DL framework for fast, flexible experimentation.
- Caffe: speed oriented deep learning framework.
- MXNet: deep learning framework for neural networks.
- CNTK: Microsoft Cognitive Toolkit.

# Most popular public ML frameworks

## For experimental HEP:

- TMVA: ROOT's builtin machine learning package.

## For ML applications:

- Keras: a Python deep learning library.
- Theano: a Python library for optimization.
- PyTorch: a DL framework for fast, flexible experimentation.
- Caffe: speed oriented deep learning framework.
- MXNet: deep learning framework for neural networks.
- CNTK: Microsoft Cognitive Toolkit.

## For ML and beyond:

- TensorFlow: library for numerical computation with data flow graphs.
- scikit-learn: general machine learning package.

# Most popular public ML frameworks

## For experimental HEP:

- TMVA: ROOT's builtin machine learning package.

## For ML applications:

- Keras: a Python deep learning library.
- Theano: a Python library for optimization.
- PyTorch: a DL framework for fast, flexible experimentation.
- Caffe: speed oriented deep learning framework.
- MXNet: deep learning framework for neural networks.
- CNTK: Microsoft Cognitive Toolkit.

## For ML and beyond:

- TensorFlow: library for numerical computation with data flow graphs.
- scikit-learn: general machine learning package.

**Why use public codes?** → builtin models and automatic differentiation

**Keras** is a high-level deep learning framework in Python which runs on top of TensorFlow, CNTK or Theano.



**Keras** is a high-level deep learning framework in Python which runs on top of TensorFlow, CNTK or Theano.

**Pros:**

- fast prototyping, user friendly, common code for multiple backends.
- support several NN architectures out-of-the-box.
- runs seamlessly on CPU and GPU.

**Keras** is a high-level deep learning framework in Python which runs on top of TensorFlow, CNTK or Theano.

## Pros:

- fast prototyping, user friendly, common code for multiple backends.
- support several NN architectures out-of-the-box.
- runs seamlessly on CPU and GPU.

## Cons:

- more tricky to extend when custom ML setups are required
- runs only in Python

## Example of code using Keras:

---

```
1  model = Sequential() # allocate an empty model (MLP)
2
3  # append feed-forward layers 2-5-3-1
4  model.add(Dense(units=5, activation='sigmoid', input_dim=2))
5  model.add(Dense(units=3, activation='sigmoid', input_dim=5))
6  model.add(Dense(units=1, activation='linear', input_dim=3))
7
8  model.compile(loss='mse', optimizer='sgd') # compile the model
9
10 # train the model
11 model.fit(x_train, y_train, epochs=1000, batch_size=32)
12
13 # measure performance
14 loss_and_metrics = model.evaluate(x_test, y_test)
15
16 # generate predictions
17 classes = model.predict(x_test)
```

---

**TensorFlow** is a library for high performance numerical computation.

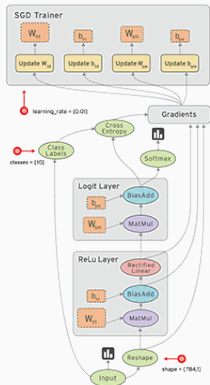
**TensorFlow** is a library for high performance numerical computation.

## Pros:

- solves optimization problems with automatic differentiation.
- can be extended in python and c/c++.
- runs seamlessly on CPU and GPU, and can uses JIT technology.

## Cons:

- do not provides builtin models from the core framework
- less automation for cross-validation and hyperparameter tune



## Example of code using TensorFlow:

---

```
1     n_input = 2
2     n_output = 1
3     n_hidden_1 = 5
4     n_hidden_2 = 3
5
6     # tf Graph input
7     X = tf.placeholder("float", [None, n_input])
8     Y = tf.placeholder("float", [None, n_output])
9
10    # Store layers weight & bias
11    weights = {
12        'h1': tf.Variable(tf.random_normal([n_input, n_hidden_1])),
13        'h2': tf.Variable(tf.random_normal([n_hidden_1, n_hidden_2])),
14        'out': tf.Variable(tf.random_normal([n_hidden_2, n_output]))
15    }
16    biases = {
17        'b1': tf.Variable(tf.random_normal([n_hidden_1])),
18        'b2': tf.Variable(tf.random_normal([n_hidden_2])),
19        'out': tf.Variable(tf.random_normal([n_output]))
20    }
```

---

## Example of code using TensorFlow:

---

```
1     ...
2
3     def MLP(x): # define the neural network
4         layer_1 = tf.add(tf.matmul(x, weights['h1']), biases['b1'])
5         layer_2 = tf.add(tf.matmul(layer_1, weights['h2']), biases['b2'])
6         return tf.matmul(layer_2, weights['out']) + biases['out']
7
8     model = MLP(X) # attach model to the input placeholder
9     loss = tf.reduce_mean(tf.square(model-Y)) # evaluate loss graph
10    train = tf.train.GradientDescentOptimizer(learning_rate).minimize(loss)
11
12    # perform training loop manually
13    ...
14    for epoch in range(1000):
15        _, cost = sess.run([train, loss], feed_dict={X: x_train, Y: y_train})
```

---

scikit-learn

Home Installation Documentation Examples

Google Custom Search

Fork me on GitHub

## scikit-learn

Machine Learning in Python

- Simple and efficient tools for data mining and data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
- Open source, commercially usable - BSD license

### Classification

Identifying to which category an object belongs to.

**Applications:** Spam detection, Image recognition.

**Algorithms:** SVM, nearest neighbors, random forest, ... — Examples

### Regression

Predicting a continuous-valued attribute associated with an object.

**Applications:** Drug response, Stock prices.

**Algorithms:** SVR, ridge regression, Lasso, ... — Examples

### Clustering

Automatic grouping of similar objects into sets.

**Applications:** Customer segmentation, Grouping experiment outcomes

**Algorithms:** k-Means, spectral clustering, mean-shift, ... — Examples

### Dimensionality reduction

Reducing the number of random variables to consider.

**Applications:** Visualization, Increased efficiency

**Algorithms:** PCA, feature selection, non-negative matrix factorization. — Examples

### Model selection

Comparing, validating and choosing parameters and models.

**Goal:** Improved accuracy via parameter tuning

**Modules:** grid search, cross validation, metrics. — Examples

### Preprocessing

Feature extraction and normalization.

**Application:** Transforming input data such as text for use with machine learning algorithms.

**Modules:** preprocessing, feature extraction. — Examples

[scikit-learn.org/stable/modules/clustering.html#mean-shift](http://scikit-learn.org/stable/modules/clustering.html#mean-shift)



Scikit-learn contains the most popular algorithms for:

- Supervised learning: neural networks, decision trees, etc.
- Unsupervised learning: density estimate, clustering, etc.
- Model selection: cross-validation, hyperparameter tune, etc.
- Dataset transformations: feature extractions, dim. reduction, etc.
- Dataset loading
- Strategies to scale computationally
- Computational performance

**Questions?**