

Machine learning and physical modelling-2

julien.brajard@nersc.no

October 2022

NERSC

<https://github.com/brajard/MAT330>

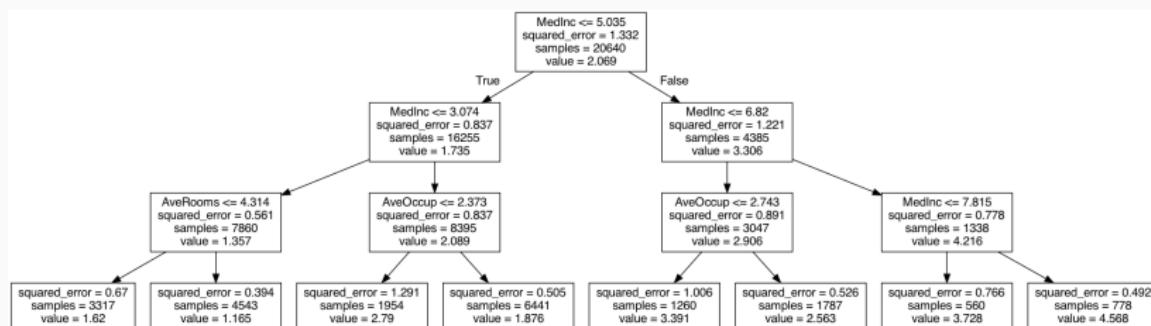
Table of contents

1. A standard Machine learning model: Random Forests
2. Neural Networks
3. A quick typology of few neural nets
4. Optimization using gradient descent
5. Gradient backpropagation
6. Optimizing a machine learning (gradient method)
7. Other regularization techniques

A standard Machine learning model: Random Forests

A decision tree

Predict house price (in \$100,000) from 8 features:



MedInc

median income in block group

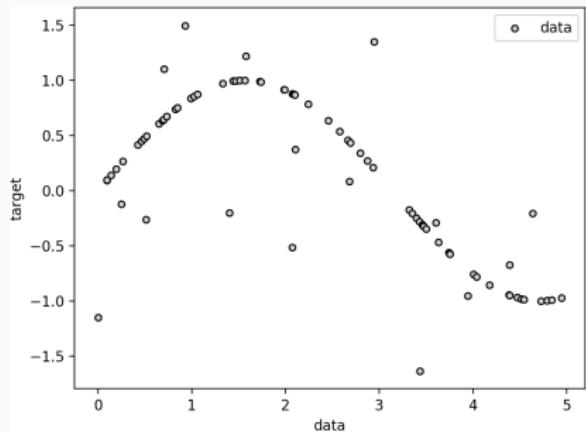
AveRooms

average number of rooms per household

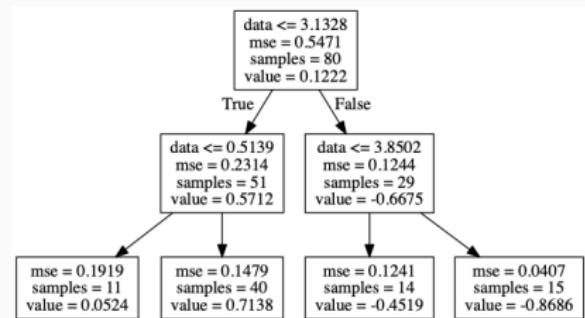
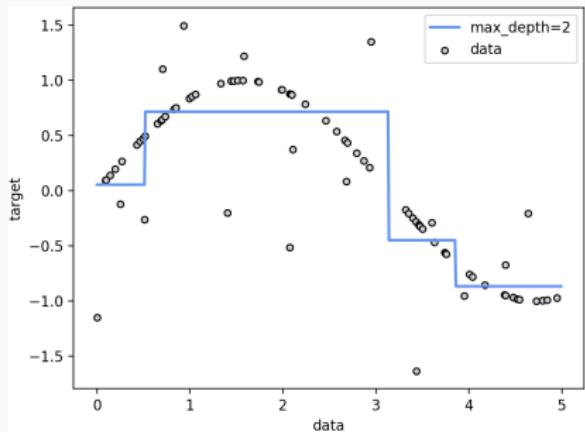
AveOccup

average number of household members

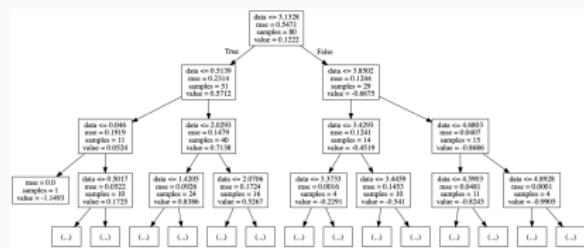
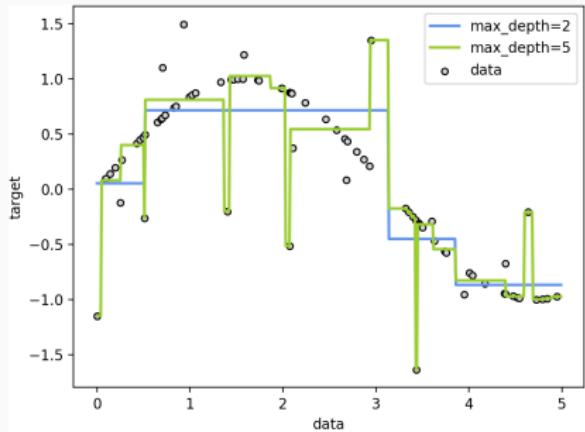
Uni-variate example



Uni-variate example



Uni-variate example



From tree to forest

Disadvantages of regression tree:

- Can overfit the data



From tree to forest

Disadvantages of regression tree:

- Can overfit the data

One extension of Regression Tree: **Random Forest**

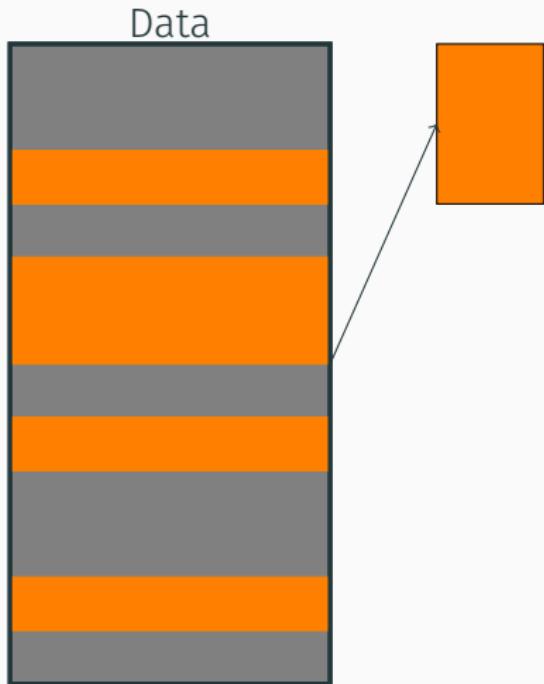


The (over simplified) principle of Random Forest

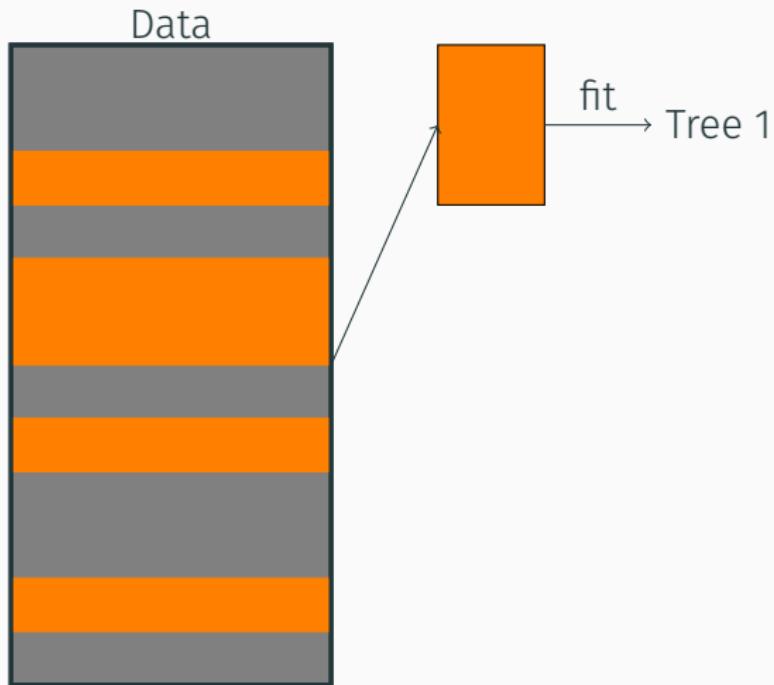
Data



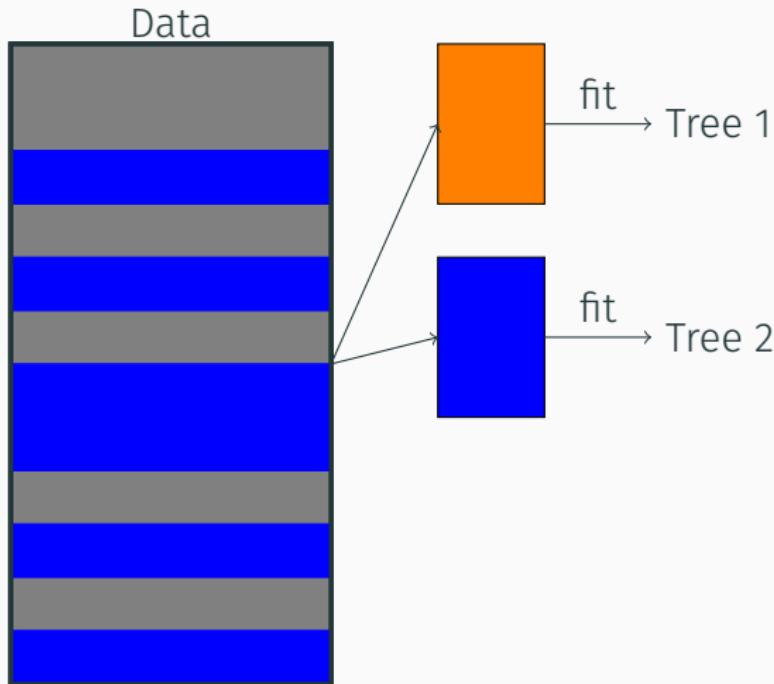
The (over simplified) principle of Random Forest



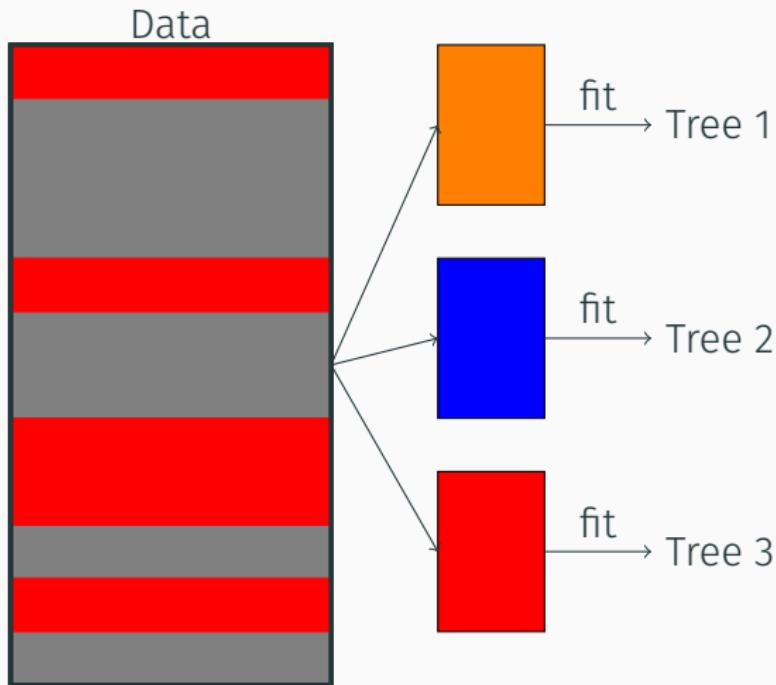
The (over simplified) principle of Random Forest



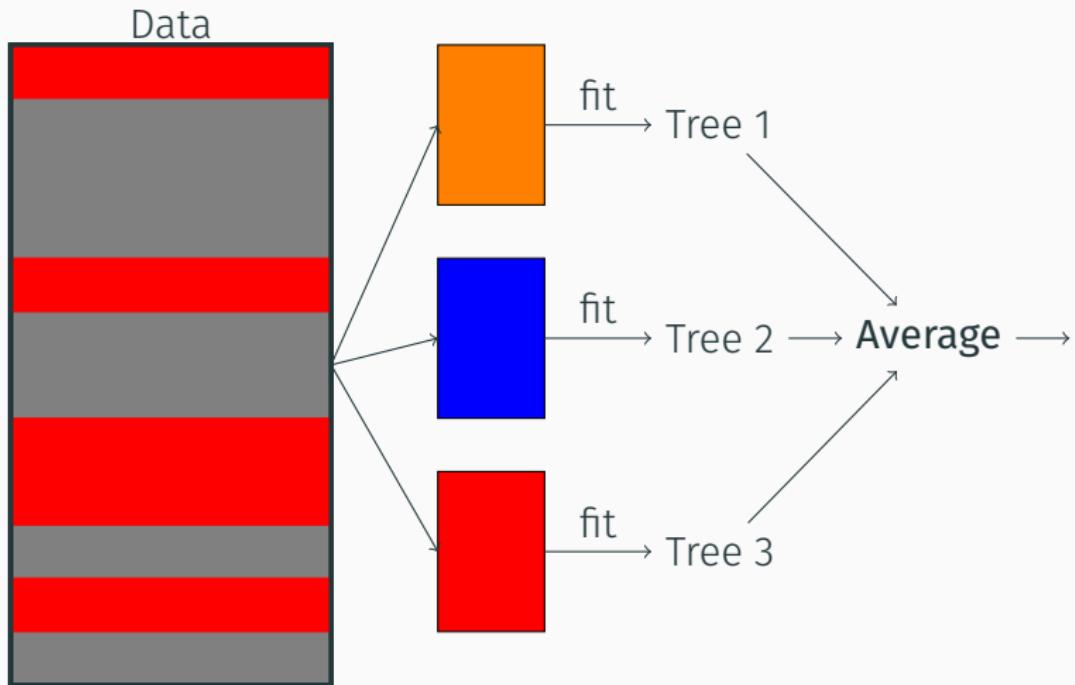
The (over simplified) principle of Random Forest



The (over simplified) principle of Random Forest

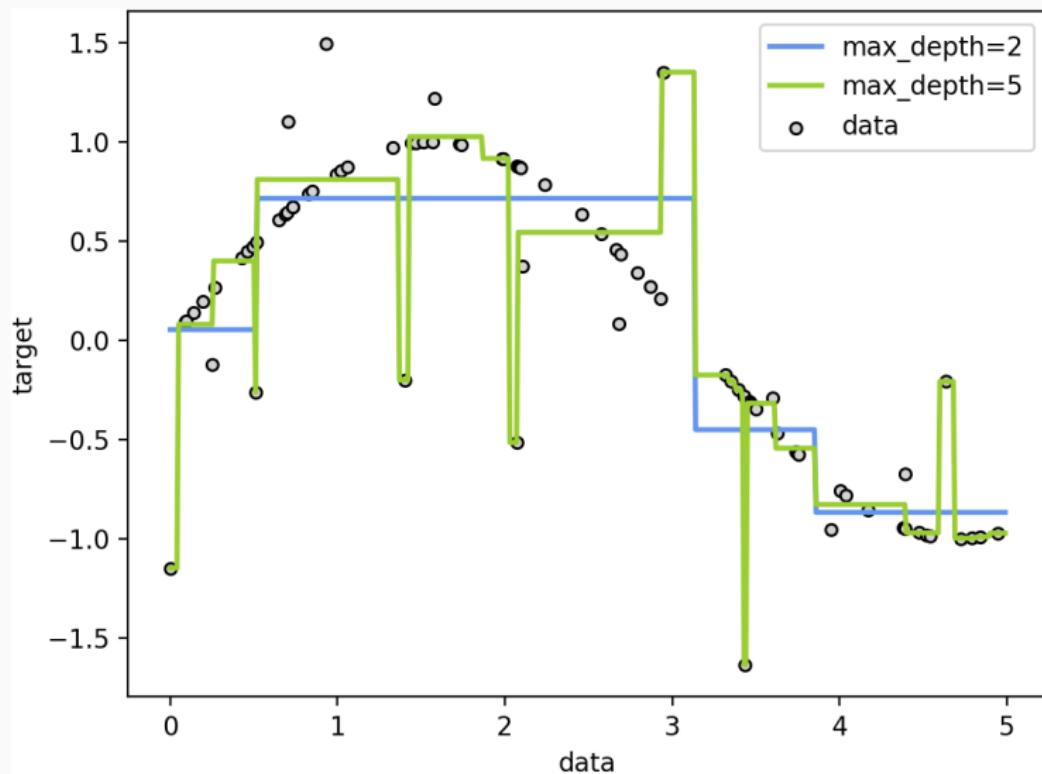


The (over simplified) principle of Random Forest



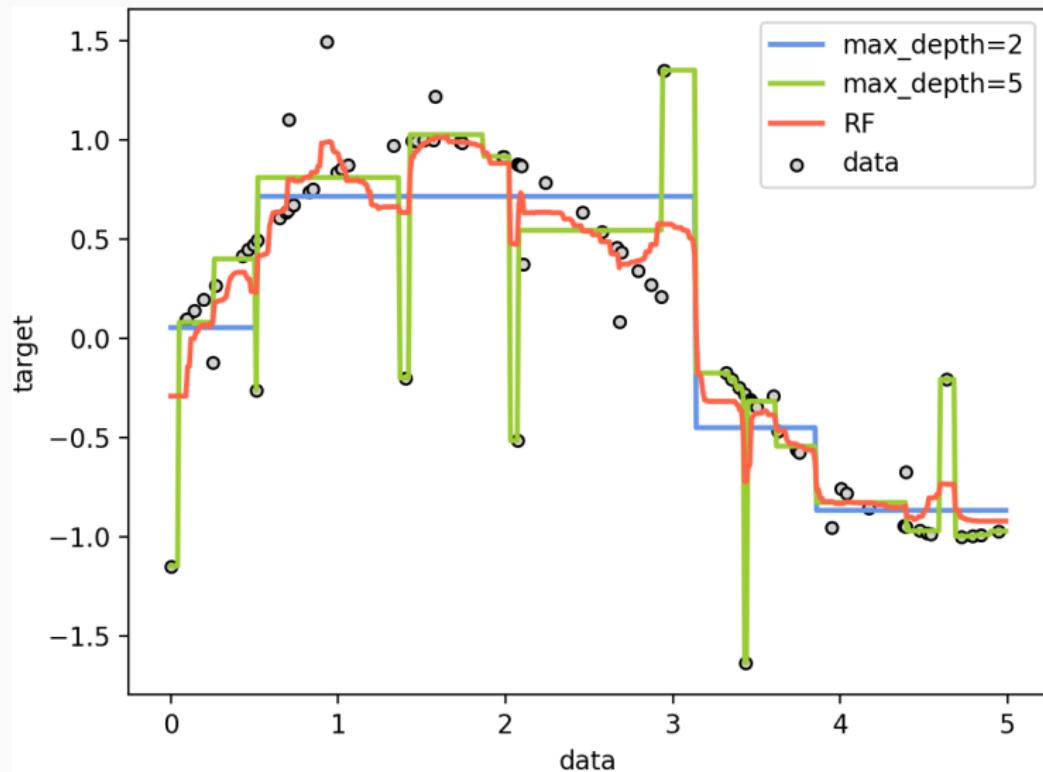
Results on the univariate experiment

Prediction of Randoms trees



Results on the univariate experiment

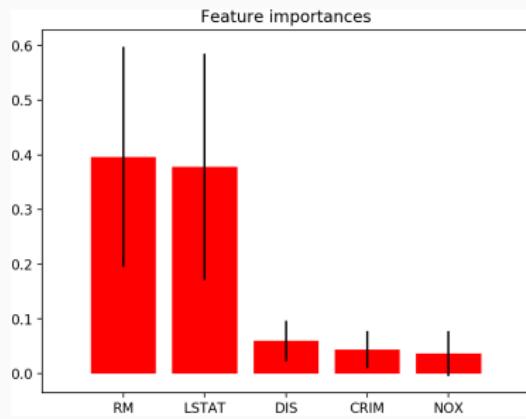
Prediction of a Random Forest



Feature importance

```
rf = RandomForestRegressor(n_estimators=1000,  
    max_features=10,random_state=10)  
rf.fit(X,y)  
importances = rf.feature_importances_
```

Indicates the impact of a feature in predicting the target.



CRIM	per capita crime rate by town
NOX	nitric oxides concentration
RM	average number of rooms per dwelling
DIS	distance to employment centres
LSTAT	lower status of the population

Some key parameters

```
from sklearn.ensemble import RandomForestClassifier  
  
rf = RandomForestRegressor(n_estimators=n, max_features=  
    maxf, min_samples_split=min_split, ...)
```

- **n_estimators**: number of trees (generally the larger is the better)

Some key parameters

```
from sklearn.ensemble import RandomForestClassifier  
  
rf = RandomForestRegressor(n_estimators=n, max_features=  
    maxf, min_samples_split=min_split, ...)
```

- **n_estimators**: number of trees (generally the larger is the better)
- **max_features**: number of features to consider at each split. The default number is the total number of features. A larger value makes provides a smaller bias (accuracy) but a bigger variance (risk of overfitting)

Some key parameters

```
from sklearn.ensemble import RandomForestClassifier  
  
rf = RandomForestRegressor(n_estimators=n, max_features=  
    maxf, min_samples_split=min_split, ...)
```

- **n_estimators**: number of trees (generally the larger is the better)
- **max_features**: number of features to consider at each split. The default number is the total number of features. A larger value makes provides a smaller bias (accuracy) but a bigger variance (risk of overfitting)
- **min_samples_fit**: number of features to consider at each split. The minimum value of 2 means that the tree is fully developed (small bias but great variance).

Determination of the hyperparameters

- Parameters that are not optimized during the training are called **hyperparameters**.

Determination of the hyperparameters

- Parameters that are not optimized during the training are called **hyperparameters**.
- They can be determined using a score on the **validation** dataset or using a **cross-validation** procedure.

Determination of the hyperparameters

- Parameters that are not optimized during the training are called **hyperparameters**.
- They can be determined using a score on the **validation** dataset or using a **cross-validation** procedure.

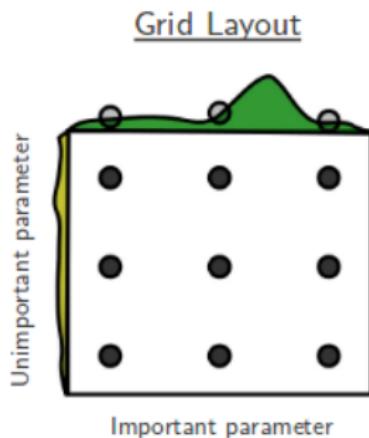
Determination of the hyperparameters

- Parameters that are not optimized during the training are called **hyperparameters**.
- They can be determined using a score on the **validation dataset** or using a **cross-validation** procedure.



Stir the pile: The gridsearch

1. Specify a list of hyperparameters to be tested.
2. For each of the parameters, specify a set of values to test
3. Train a model for each of the possible combinations of hyperparameters
4. Retain the best model (using, e.g., cross-validation)



<https://medium.com/@senapati.dipak97/grid-search-vs-random-search-d34c92946318>

Remarks on the gridsearch procedure

- It make an **exhaustive** search of the hyperparameters

Remarks on the gridsearch procedure

- It make an **exhaustive** search of the hyperparameters
- The procedure is easy to **parallelized**.

Remarks on the gridsearch procedure

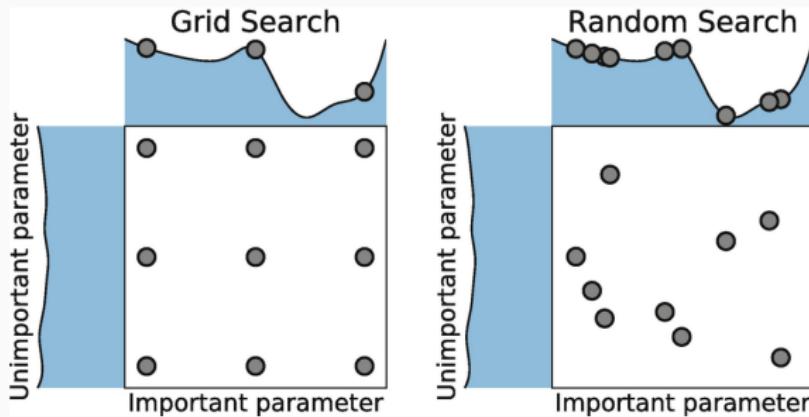
- It make an **exhaustive** search of the hyperparameters
- The procedure is easy to **parallelized**.
- it is not naturally adapted for quantitative hyperparameters.

Remarks on the gridsearch procedure

- It makes an **exhaustive** search of the hyperparameters
- The procedure is easy to **parallelized**.
- it is not naturally adapted for quantitative hyperparameters.
- it can become **very costly**. (e.g. 8 hyperparameters with 8 values each to test = $8^8 = 16,777,216$ trainings.)

Random search

1. Specify a list of hyperparameters to be tested.
2. For each of the parameters, specify a set of values to test or a law to draw a random value.
3. Draw n combinations of the hyperparameters.
4. Train a model for each of the combinations.
5. Retain the best model (using, e.g., cross-validation)



Remarks on the random search procedure

- It **does not make** an **exhaustive** search of the hyperparameters
- The procedure is easy to **parallelized**.
- The cost is predictable (number of draw).

Remarks on the random search procedure

- It **does not make** an **exhaustive** search of the hyperparameters
- The procedure is easy to **parallelized**.
- The cost is predictable (number of draw).

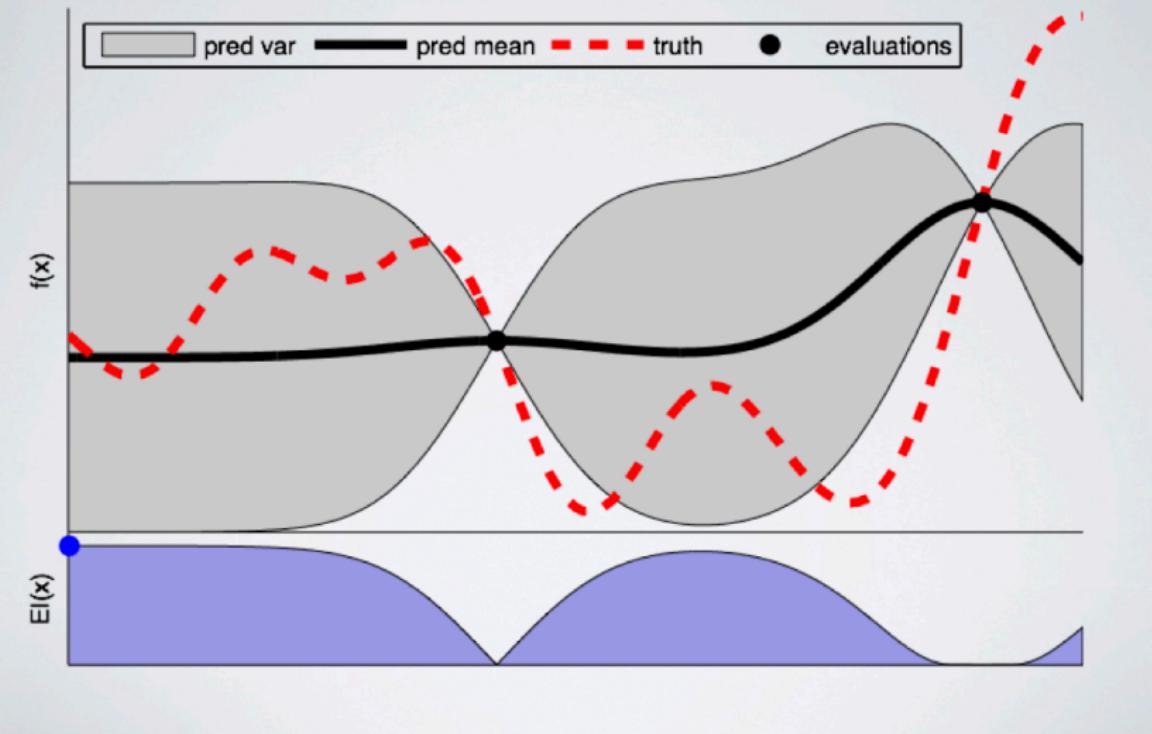
Both gridsearch and random search are implemented and easy to use in scikit-learn.

A more "advance" method is the **Bayesian** optimization.

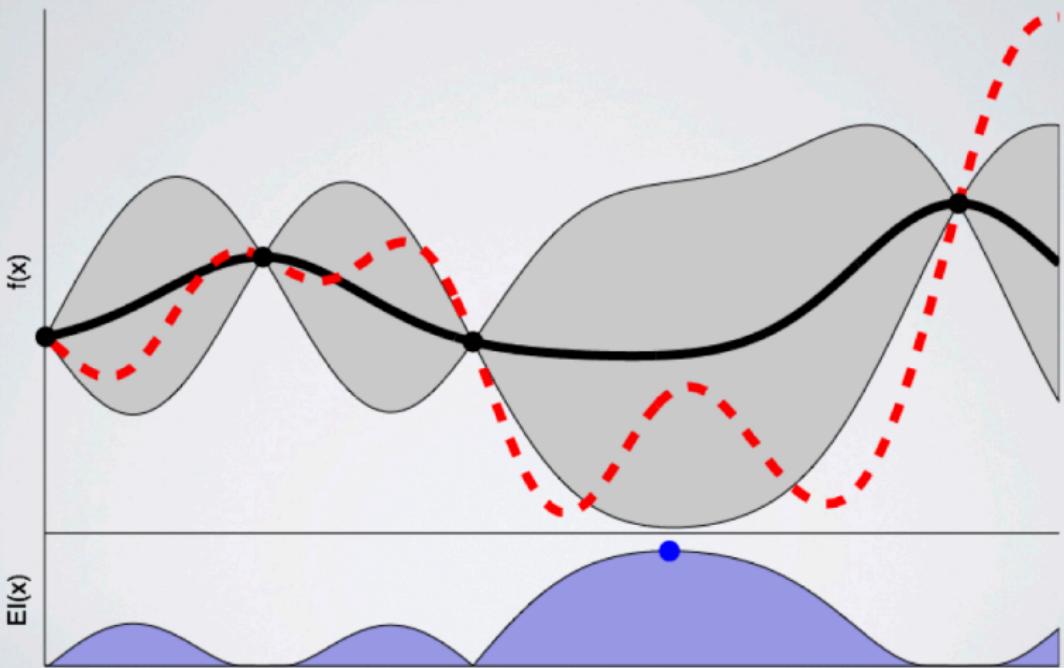
The principle is to search hyperparameters where it is more likely to improve the model (based on previous attempts)

Illustration taken from Ryan P. Adams: "A Tutorial on Bayesian Optimization for Machine Learning"

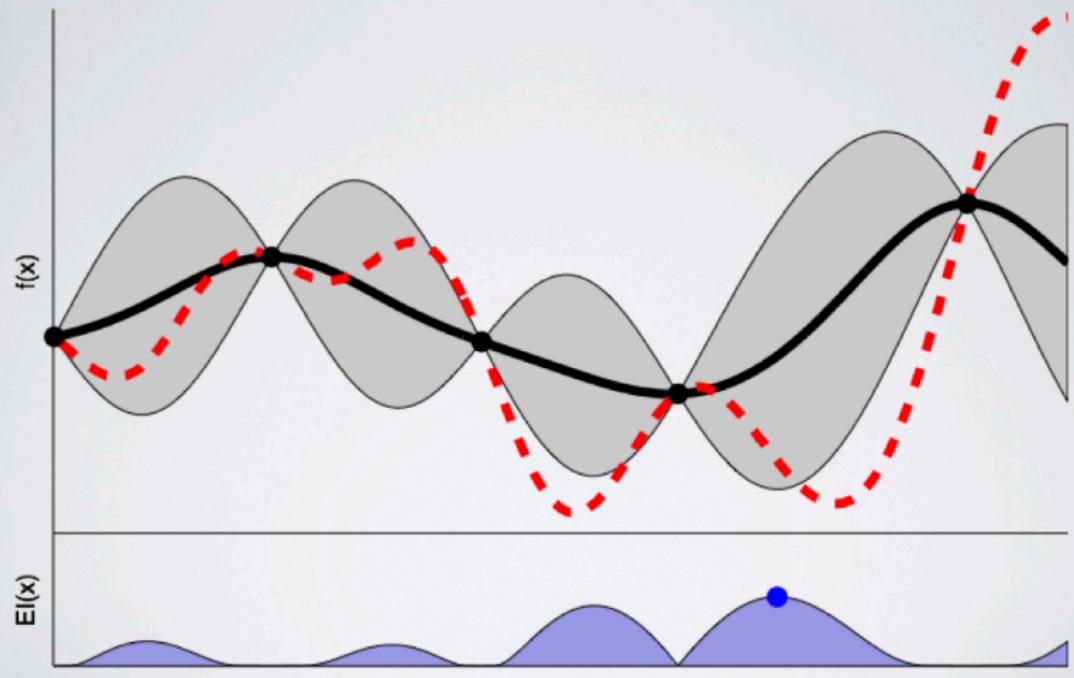
Illustrating Bayesian Optimization



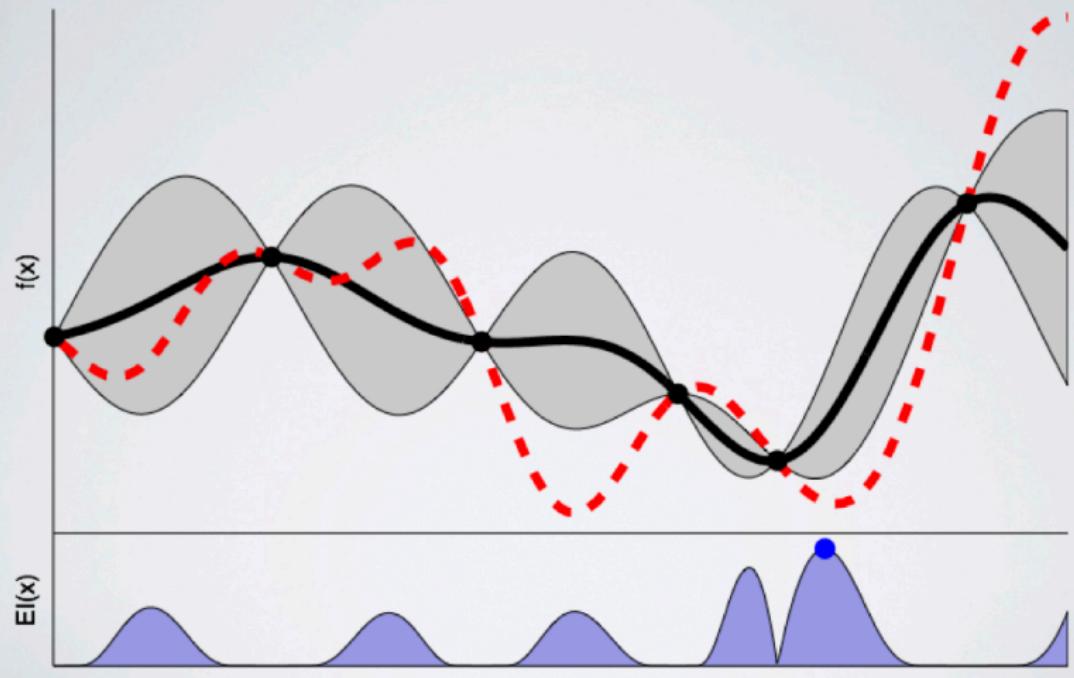
Illustrating Bayesian Optimization



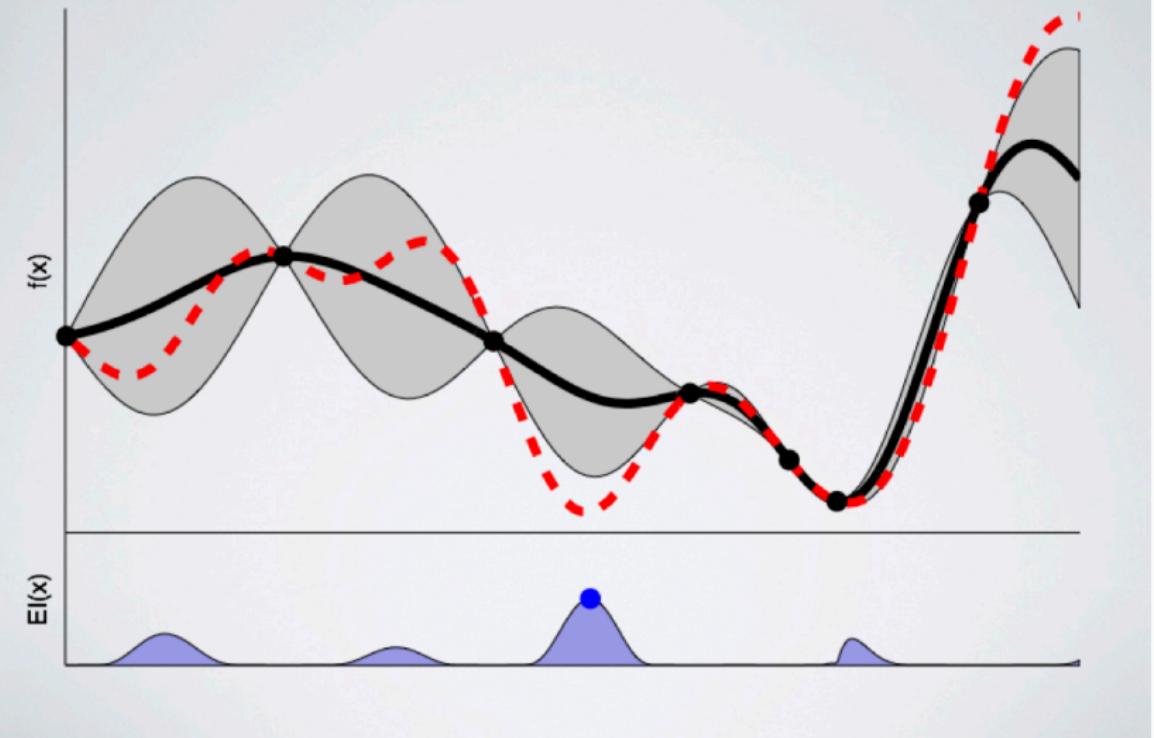
Illustrating Bayesian Optimization



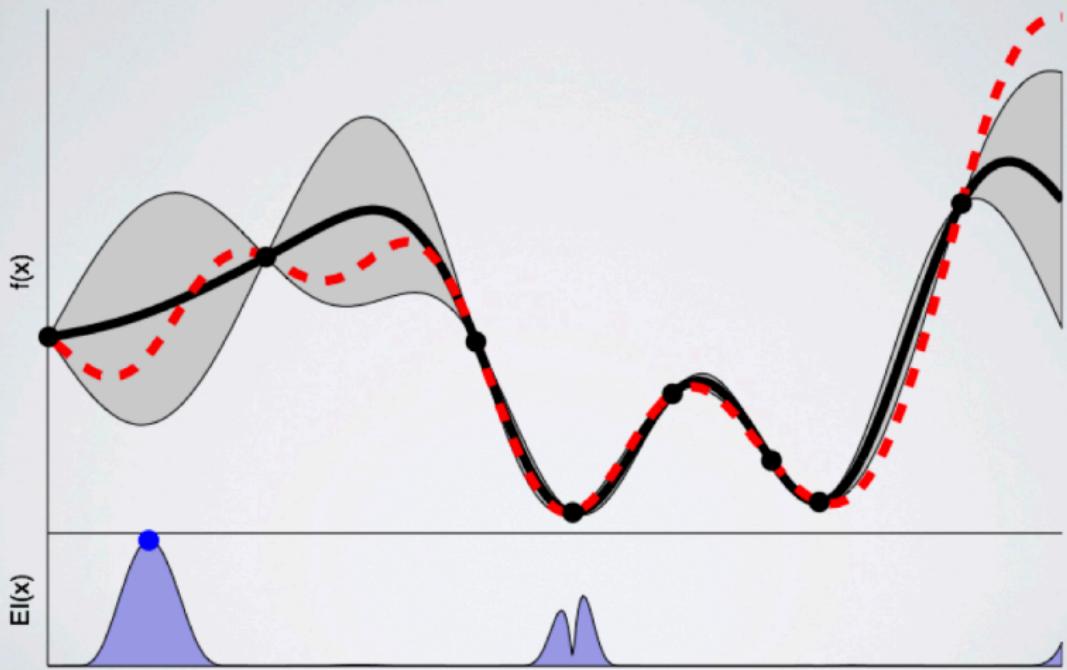
Illustrating Bayesian Optimization



Illustrating Bayesian Optimization

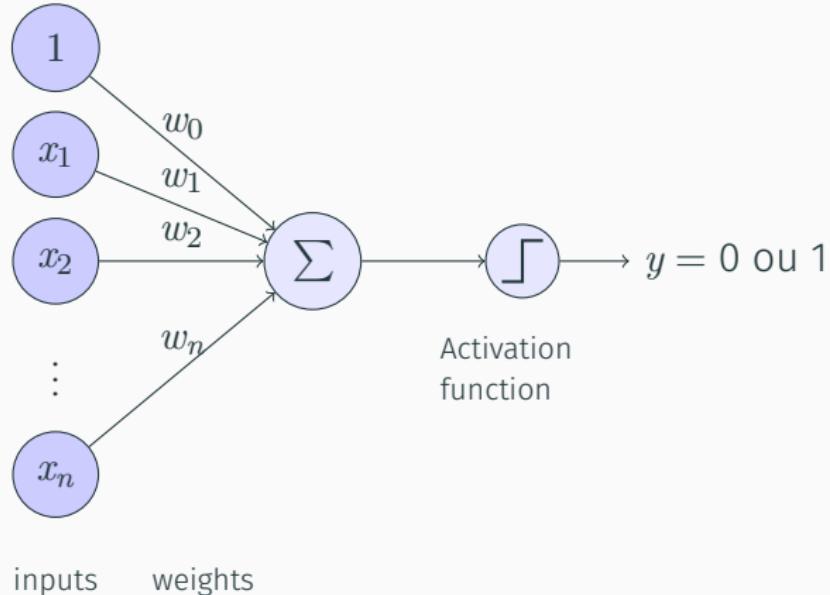


Illustrating Bayesian Optimization



Neural Networks

The perceptron : an artificial neuron



Computation

$$y = f(w_0 + w_1 \cdot x_1 + w_2 \cdot x_2 + \cdots + w_n \cdot x_n) = f(w_0 + \sum_{i=1}^n w_i \cdot x_i)$$

Some remarks

- Inputs x_i are the different features of the data

Some remarks

- Inputs x_i are the different features of the data
- Weight w_i are the parameters of the model to optimize

Some remarks

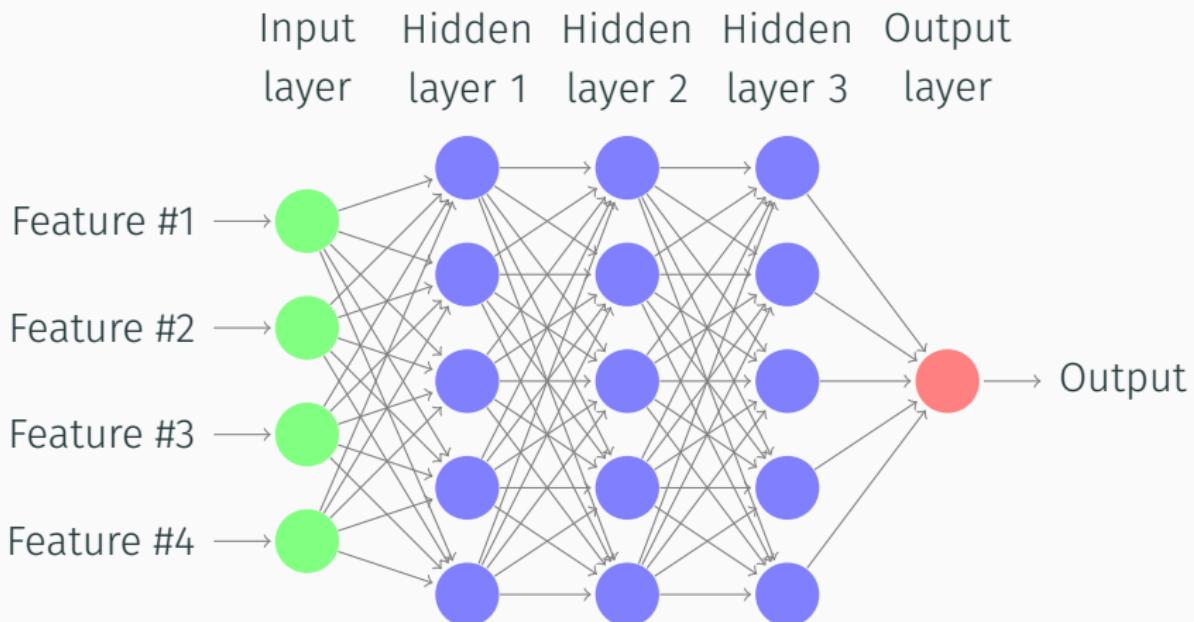
- Inputs x_i are the different features of the data
- Weight w_i are the parameters of the model to optimize
- If the activation function is identity, it is equivalent to a linear regression

Some remarks

- Inputs x_i are the different features of the data
- Weight w_i are the parameters of the model to optimize
- If the activation function is identity, it is equivalent to a linear regression

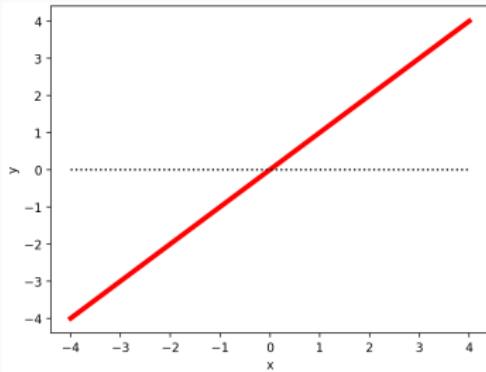
More complexe models are build by combining several perceptrons

Multi-layer perceptron (Densely connected layers)

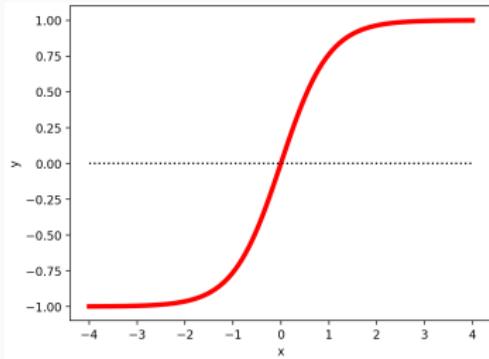


Most usual activation functions

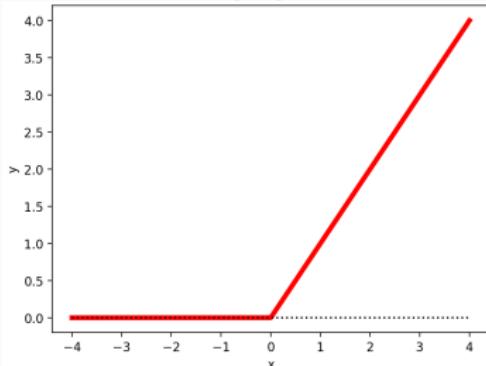
Linear



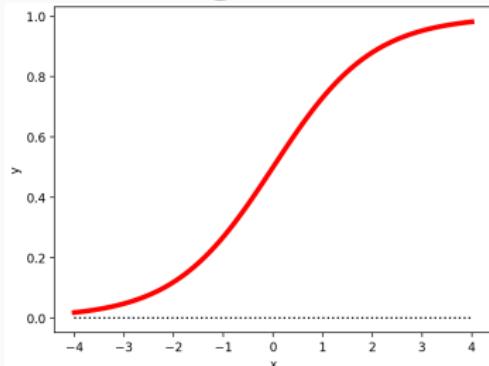
Hyperbolic tangent



ReLU

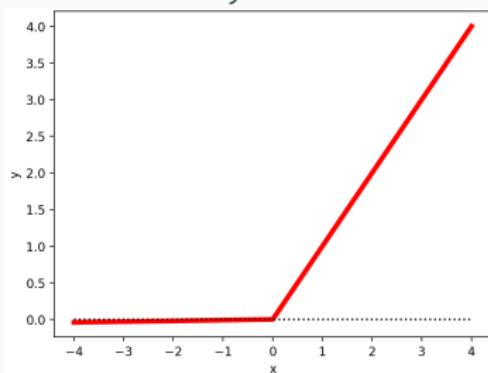


Sigmoid

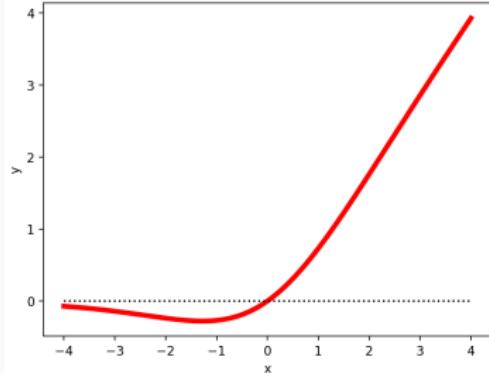


New fancy activation functions

Leaky-ReLU



Swish



Regression

- Last layer:
linear or hyperbolic
tangent
- Loss function:

$$L(\hat{y}, y) = \sum_i (\hat{y}_i - y_i)^2$$

Classification and regression loss

Regression

- Last layer:
linear or hyperbolic
tangent
- Loss function:

$$L(\hat{y}, y) = \sum_i (\hat{y}_i - y_i)^2$$

Classification

- Last layer:
Soft-max

$$p_j = f_j(\mathbf{h}) = \frac{e^{h_j}}{\sum_k e^{h_k}}$$

- Loss function:
Negative crossentropy

$$L(p, y) = - \sum_i \sum_j y_{i,j} \log p_{i,j}$$

Classification loss (binary case)

Objective: binary classification (the model determine if the input feature is in a class or not)

Exemple: Try to know if an image is a cat or not.

Dataset:

x (feature)				...
y (target)	Cat	Dog	Cat	...

How to proceed?

1. Encode the targets (1 for cat, 0 otherwise)

y (target)	Cat	Dog	Cat	...
y (encoded)	1	0	1	...

How to proceed?

1. Encode the targets (1 for cat, 0 otherwise)

y (target)	Cat	Dog	Cat	...
y (encoded)	1	0	1	...

2. Design a model with **one output** and use the **sigmoid** as output activation function of your model $f(x)$, so $0 < f(x) < 1$.

Rule: if $f(x) > \frac{1}{2}$, classify as "Cat".

$f(x)$ is interpreted as the probability of the image x to be a cat.

How to proceed?

1. Encode the targets (1 for cat, 0 otherwise)

y (target)	Cat	Dog	Cat	...
y (encoded)	1	0	1	...

2. Design a model with **one output** and use the **sigmoid** as output activation function of your model $f(x)$, so $0 < f(x) < 1$.

Rule: if $f(x) > \frac{1}{2}$, classify as "Cat".

$f(x)$ is interpreted as the probability of the image x to be a cat.

3. Loss function to minimize is binary cross entropy:

$$L = - \sum y_i \cdot \log(f(x_i)) + (1 - y_i) \cdot \log(1 - f(x_i))$$

Classification loss (multiclass)

Objective: classification (the model determine in which class the input feature (more than 2 classes))

Exemple: Try to know if an image is a cat, a dog or a duck.

Dataset:

x (feature)



y (target)

Cat



Dog



Duck

...

...

How to proceed?

1. Encode the targets using one hot encoding

y (target)	Cat	Dog	Cat	...
y (encoded)	(1, 0, 0)	(0, 1, 0)	(0, 0, 1)	...

How to proceed?

1. Encode the targets using one hot encoding

y (target)	Cat	Dog	Cat	...
y (encoded)	(1, 0, 0)	(0, 1, 0)	(0, 0, 1)	...

2. Design a model with **N outputs** (N being the number of modalities) and use the **soft-max** as output activation function

$$p_j = f_j(\mathbf{h}) = \frac{e^{h_j}}{\sum_k e^{h_k}}$$

Rule: the class is attributed to the argument of the maximum of \mathbf{p} . Ex $\mathbf{p} = (0.1, 0, 7, 0, 2)$ is classified as "dog".

p_j is interpreted as the probability of the image x to belong to the class j

How to proceed?

1. Encode the targets using one hot encoding

y (target)	Cat	Dog	Cat	...
y (encoded)	(1, 0, 0)	(0, 1, 0)	(0, 0, 1)	...

2. Design a model with **N outputs** (N being the number of modalities) and use the **soft-max** as output activation function

$$p_j = f_j(\mathbf{h}) = \frac{e^{h_j}}{\sum_k e^{h_k}}$$

Rule: the class is attributed to the argument of the maximum of \mathbf{p} . Ex $\mathbf{p} = (0.1, 0, 7, 0, 2)$ is classified as "dog".

p_j is interpreted as the probability of the image x to belong to the class j

3. Loss function to minimize is negative cross entropy

$$L = - \sum_i \sum_j y_{i,j} \cdot \log p_{i,j}$$

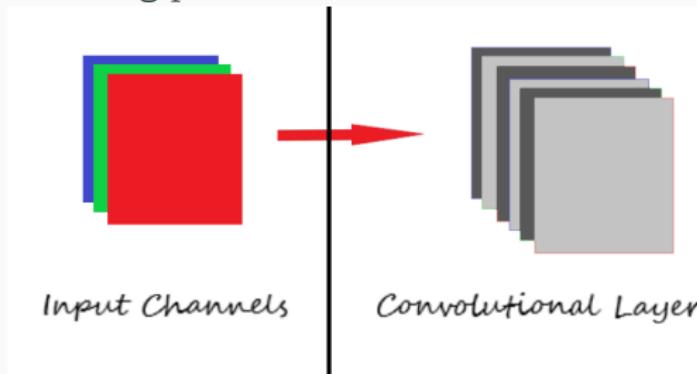
Main parameters of a convolutional layer

- Size of the filter K

Main parameters of a convolutional layer

- Size of the filter K
- Number of filters p

A convolutional layer is composed of p convolutions (size of layer) extracting p features from the data.



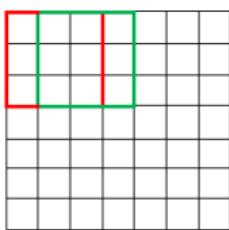
$$O = \frac{W-K+2P}{S} + 1, \text{ where } O \text{ is the output size and } W \text{ the input size.}$$

Main parameters of a convolutional layer

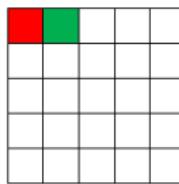
- Size of the filter K
- Number of filters p
- Strides S

$$S = 1$$

7 x 7 Input Volume

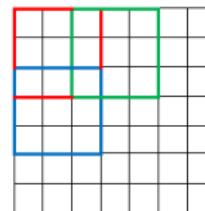


5 x 5 Output Volume

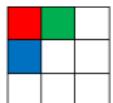


$$S = 2$$

7 x 7 Input Volume



3 x 3 Output Volume

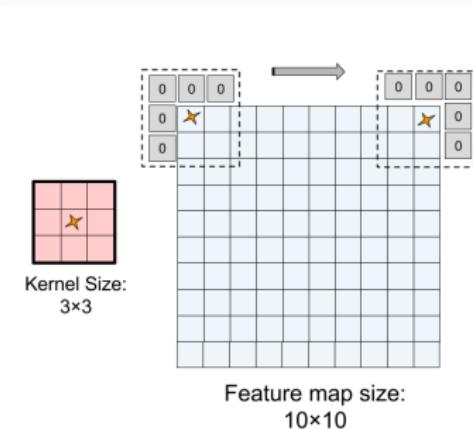


$$O = \frac{W-K+2P}{S} + 1, \text{ where } O \text{ is the output size and } W \text{ the input size.}$$

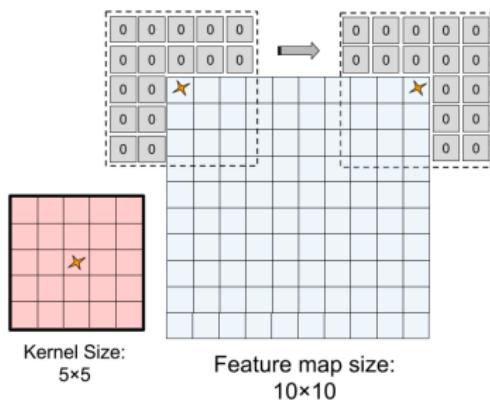
Main parameters of a convolutional layer

- Size of the filter K
- Number of filters p
- Strides S
- Padding P

$$P = 1$$



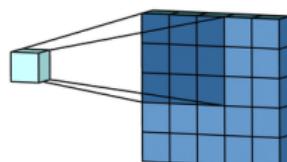
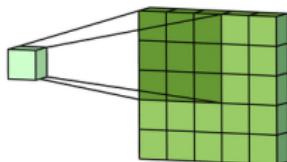
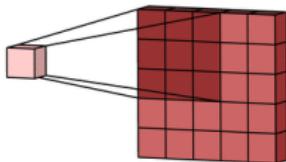
$$P = 2$$



$$O = \frac{W-K+2P}{S} + 1, \text{ where } O \text{ is the output size and } W \text{ the input size.}$$

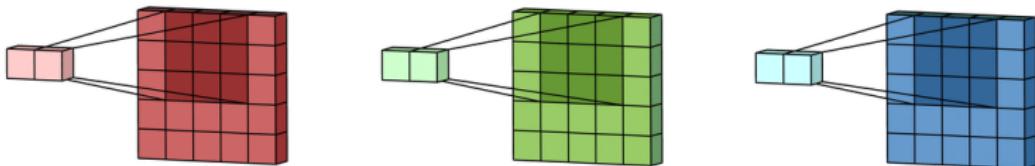
Summary of Convolutional layer steps

1. Convolution



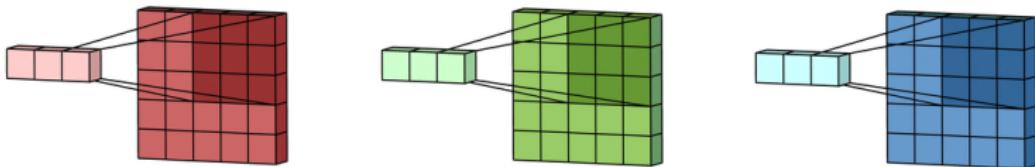
Summary of Convolutional layer steps

1. Convolution



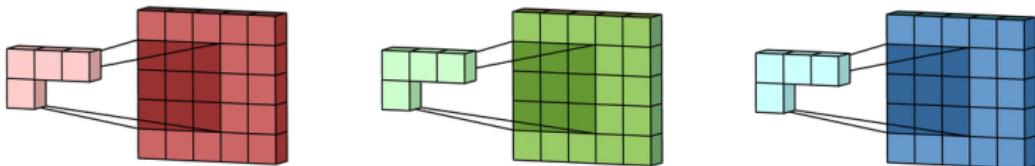
Summary of Convolutional layer steps

1. Convolution



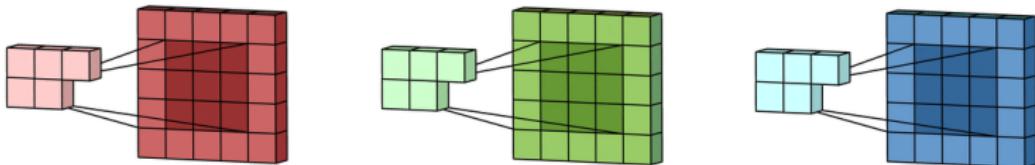
Summary of Convolutional layer steps

1. Convolution



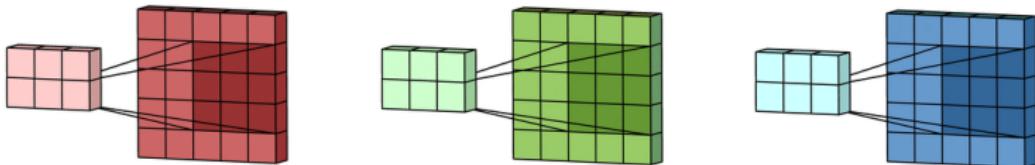
Summary of Convolutional layer steps

1. Convolution



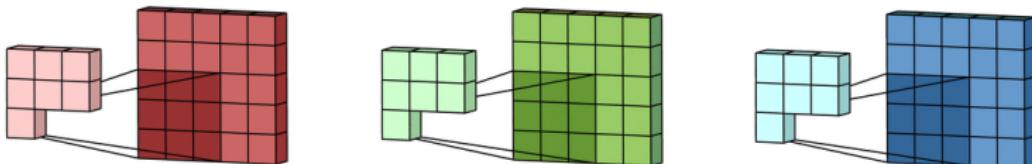
Summary of Convolutional layer steps

1. Convolution



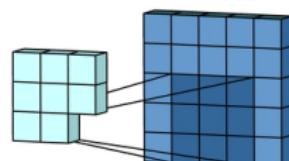
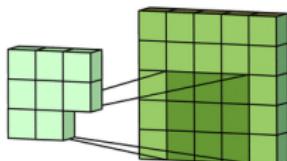
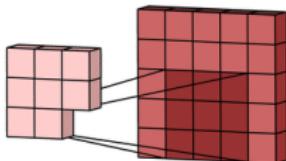
Summary of Convolutional layer steps

1. Convolution



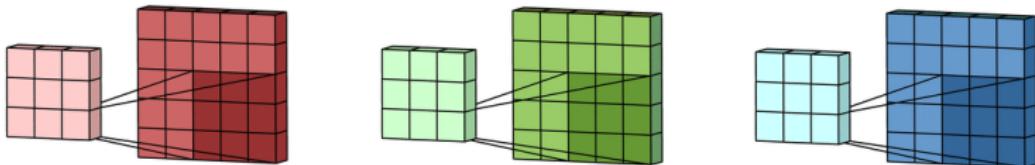
Summary of Convolutional layer steps

1. Convolution



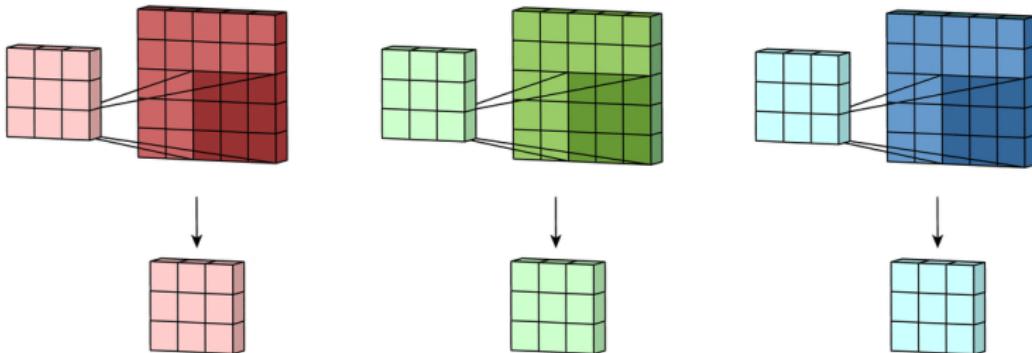
Summary of Convolutional layer steps

1. Convolution



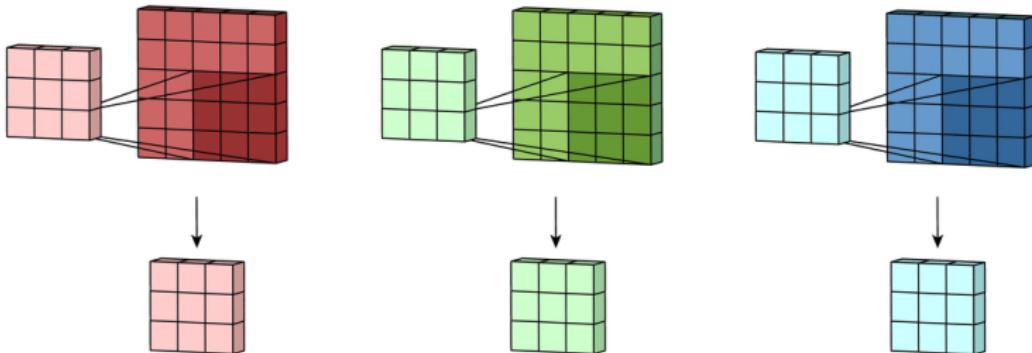
Summary of Convolutional layer steps

1. Convolution



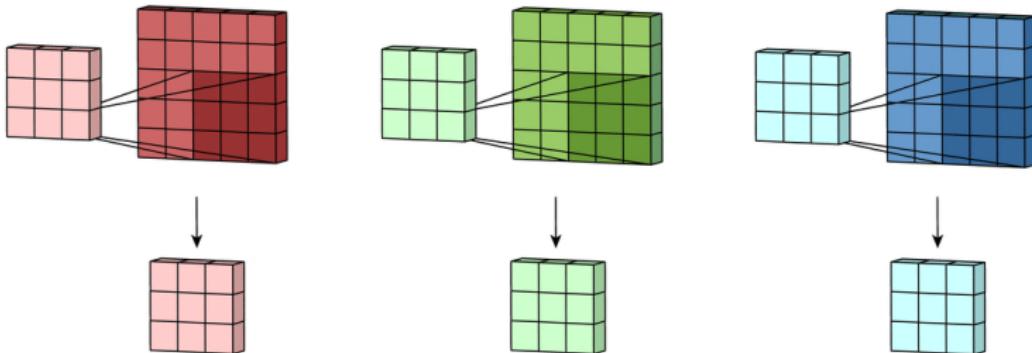
Summary of Convolutional layer steps

1. Convolution



Summary of Convolutional layer steps

1. Convolution

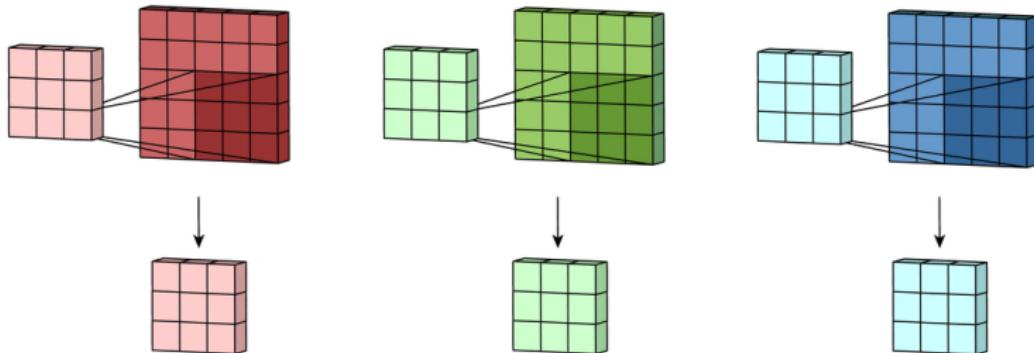


2. Addition

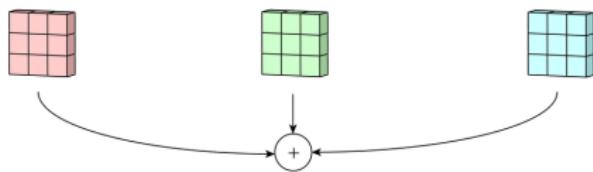


Summary of Convolutional layer steps

1. Convolution

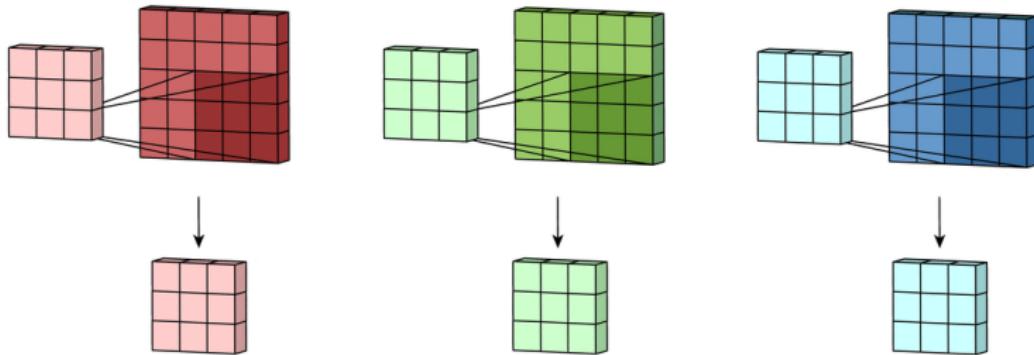


2. Addition

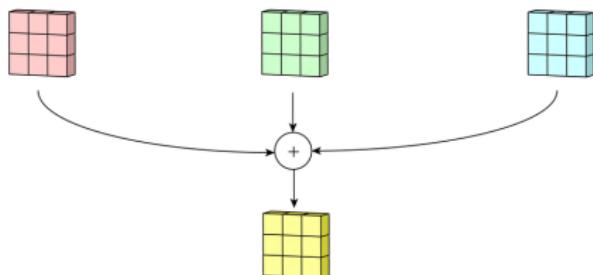


Summary of Convolutional layer steps

1. Convolution

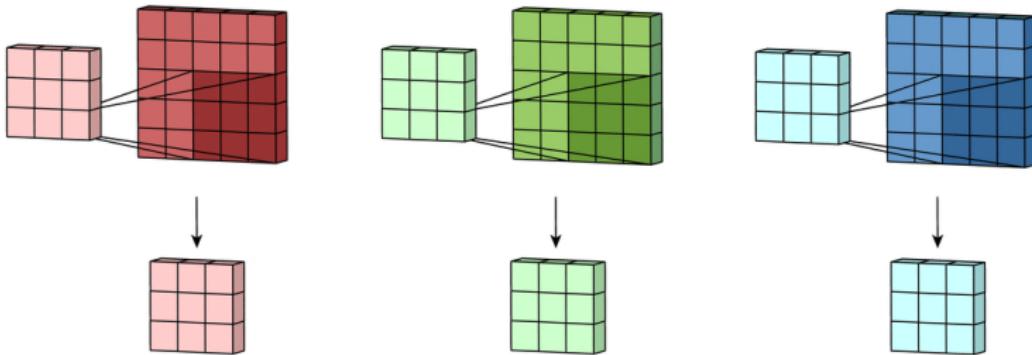


2. Addition

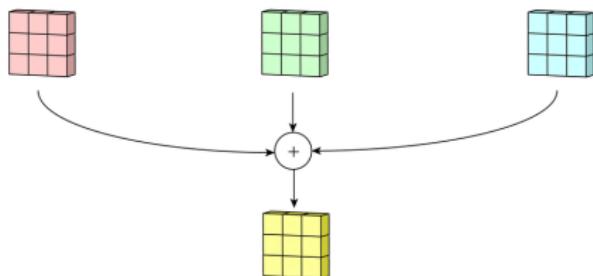


Summary of Convolutional layer steps

1. Convolution

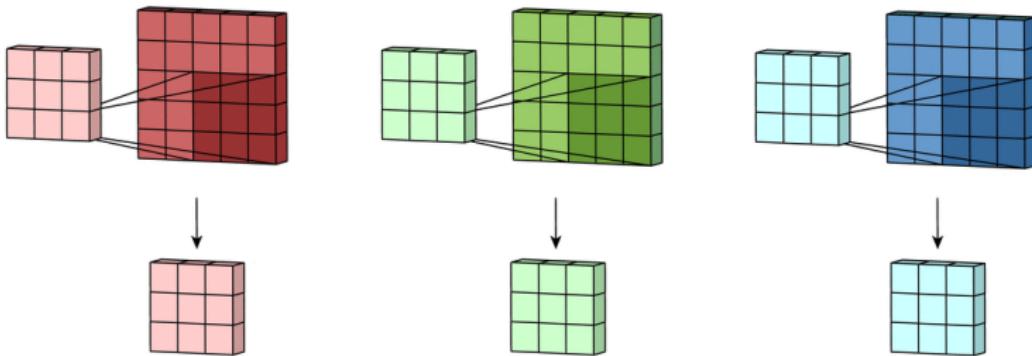


2. Addition

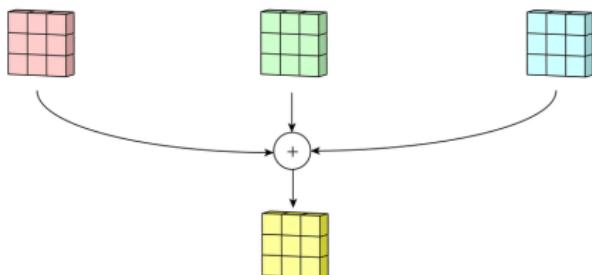


Summary of Convolutional layer steps

1. Convolution



2. Addition

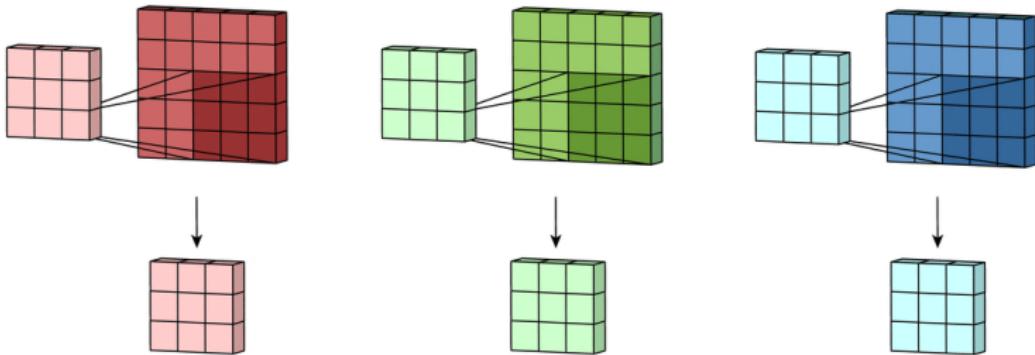


3. Bias

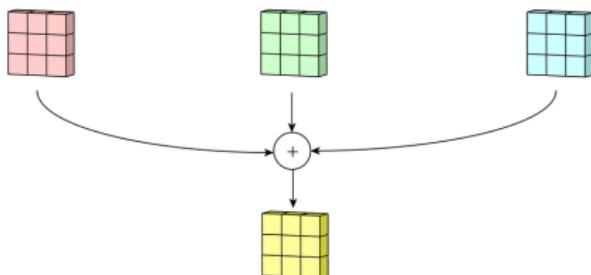


Summary of Convolutional layer steps

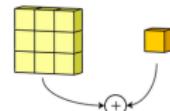
1. Convolution



2. Addition

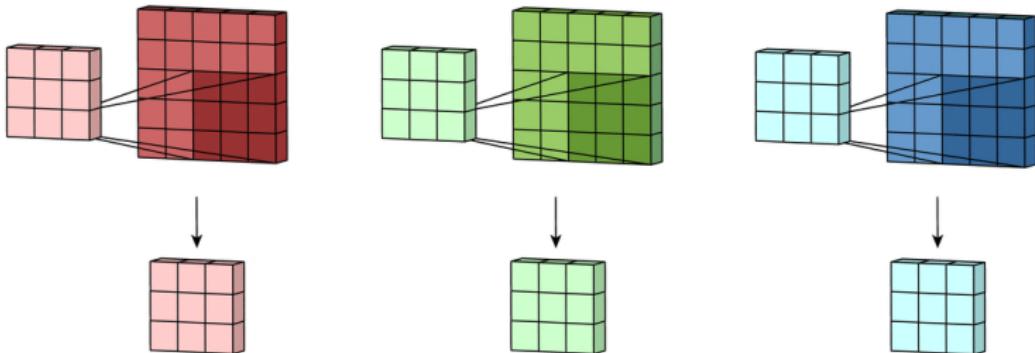


3. Bias

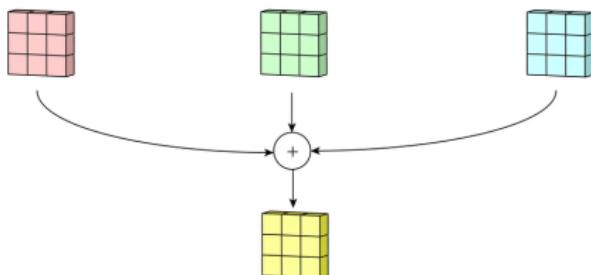


Summary of Convolutional layer steps

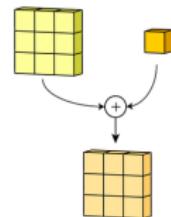
1. Convolution



2. Addition



3. Bias

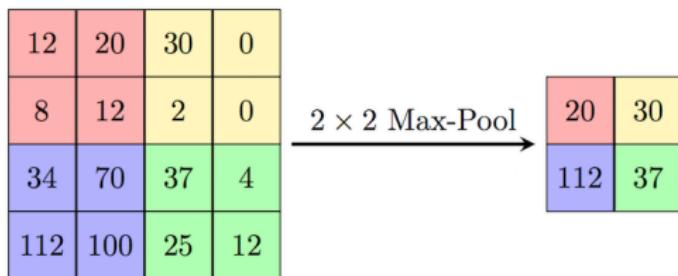


Remarks on Convolutional layers

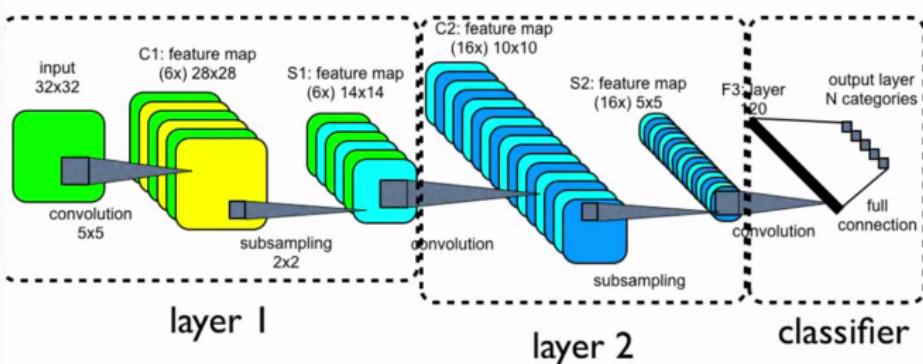
- Convolutional layers are acting locally on the image (But you can still use large scale information by adding more layers)
- Convolutions are invariant by translation (the weights do not depend on the location on the image).
- They can handle images of different sizes.

Max-Pooling

In order to reduce the size of the feature space (and to enhance the gradients), a common operation is to perform a max-pooling.



A traditionnal CNN architecture



Example of AlexNet

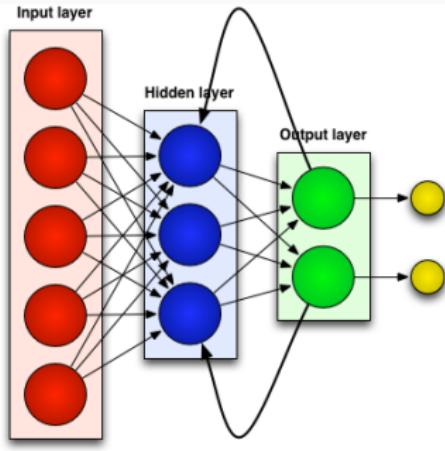
AlexNet is the first Deep architecture used on ImageNet challenge in 2012 and achieved an error of 15.3% (10% better than the previous best classifier). The paper was cited more than 34,000 times.

 Alex Krizhevsky and Geoffrey E Hinton, *ImageNet Classification with Deep Convolutional Neural Networks*, Neural Information Processing Systems (2012), 1–9.

Layer		Feature Map	Size	Kernel Size	Stride	Activation
Input	Image	1	227x227x3	-	-	-
1	Convolution	96	55 x 55 x 96	11x11	4	relu
	Max Pooling	96	27 x 27 x 96	3x3	2	relu
2	Convolution	256	27 x 27 x 256	5x5	1	relu
	Max Pooling	256	13 x 13 x 256	3x3	2	relu
3	Convolution	384	13 x 13 x 384	3x3	1	relu
4	Convolution	384	13 x 13 x 384	3x3	1	relu
5	Convolution	256	13 x 13 x 256	3x3	1	relu
	Max Pooling	256	6 x 6 x 256	3x3	2	relu
6	FC	-	9216	-	-	relu
7	FC	-	4096	-	-	relu
8	FC	-	4096	-	-	relu
Output	FC	-	1000	-	-	Softmax

A quick typology of few neural nets

Recurrent Neural Networks

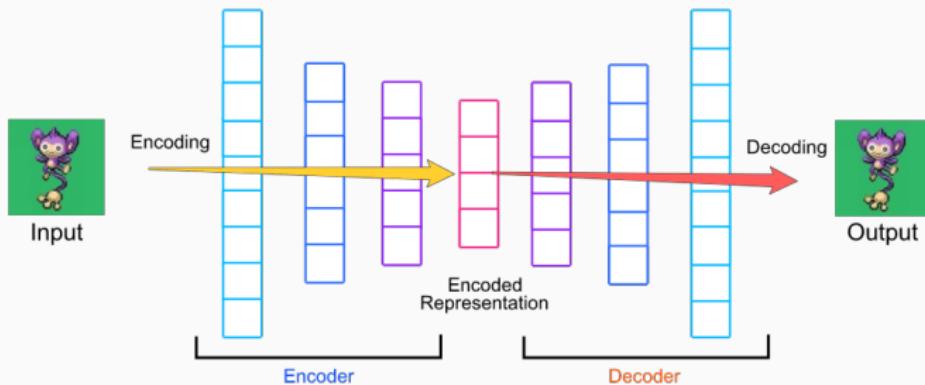


Some popular types of recurrent neural networks:

- Long short-term memory (LSTM)
- Gated Recurrent Unit (GRU)

Used in machine translation and text processing

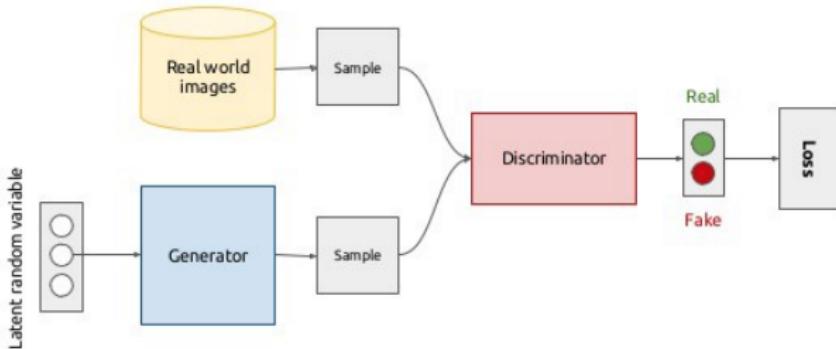
Autoencoders



Used in image denoising, compressing, generation,...

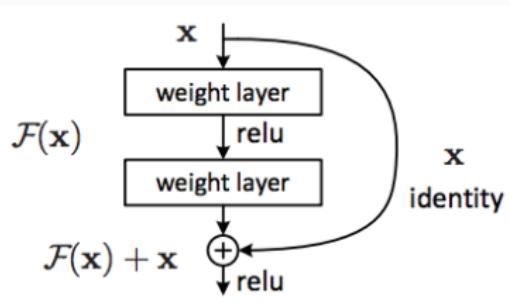
Generative adversarial networks

Generative adversarial networks (conceptual)



5

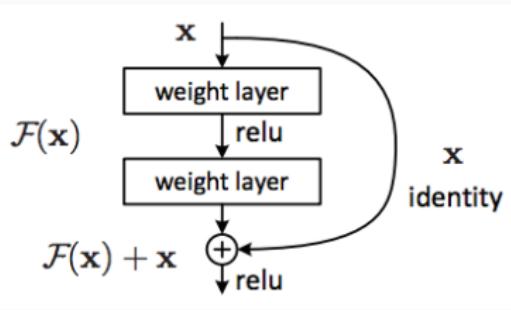
Residual Networks



x : input, y : output

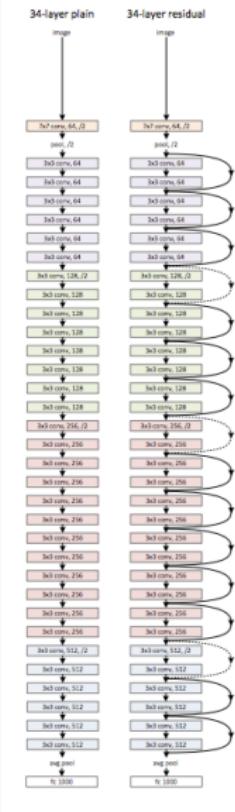
$$y = x + \mathcal{F}(x)$$

Residual Networks



x : input, y : output

$$y = x + \mathcal{F}(x)$$



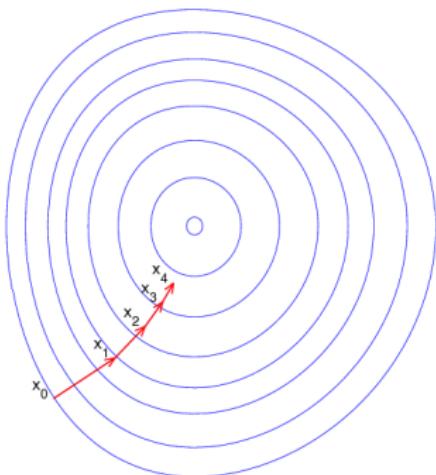
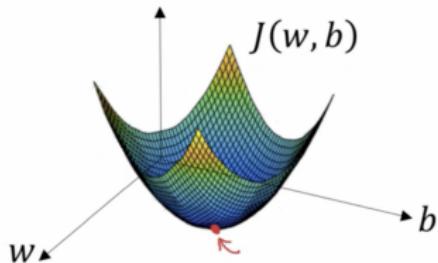
Optimization using gradient descent

What is gradient descent

Objective

Minimize the function $L(\theta)$ where θ is a vector of parameters (e.g. weights of a neural net).

Example with $\theta = (w, b)$



Iterative algorithm

1. We start with a "first guess" of the parameter θ_0

Iterative algorithm

1. We start with a "first guess" of the parameter θ_0
2. we iterate over several values of the parameters following the update rule:

$$\theta_{k+1} = \theta_k - \gamma \nabla L(\theta_k),$$

γ is called the **learning rate**.

Iterative algorithm

1. We start with a "first guess" of the parameter θ_0
2. we iterate over several values of the parameters following the update rule:

$$\theta_{k+1} = \theta_k - \gamma \nabla L(\theta_k),$$

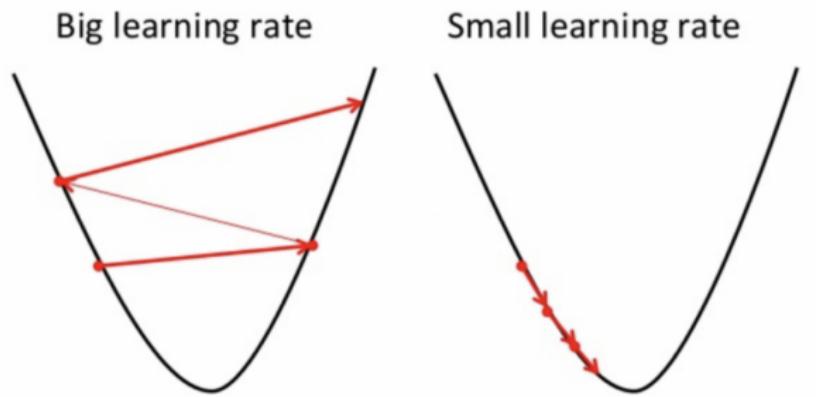
γ is called the **learning rate**.

Iterative algorithm

1. We start with a "first guess" of the parameter θ_0
2. we iterate over several values of the parameters following the update rule:

$$\theta_{k+1} = \theta_k - \gamma \nabla L(\theta_k),$$

γ is called the **learning rate**.



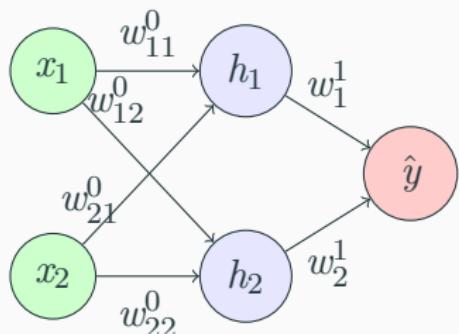
Few comments

- The learning γ is an important **hyperparameter** that needs to be tuned using, e.g random search

- The learning γ is an important **hyperparameter** that needs to be tuned using, e.g random search
- The key part of the formula is the computation of the gradient $\nabla L(\theta_k)$

Gradient backpropagation

Training a neural-net: gradient backpropagation



1. Given a couple (x, y)

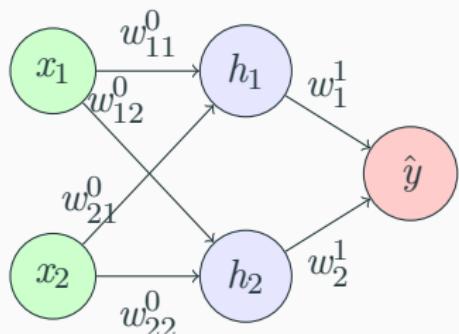
Objective

Determination of the best set of weights \mathbf{w} to minimize the Loss function

$$L(\mathbf{w}) = \|\hat{y}(\mathbf{w}) - y\|^2.$$

Calculation of $\partial L / \partial w$

Training a neural-net: gradient backpropagation



- Given a couple (x, y)
- Forward computation:

$$h_j = f_0(\sum_{i=1}^2 w_{ij}^0 \cdot x_i)$$

$$\hat{y} = f_1(\sum_{j=1}^2 w_j^1 \cdot h_j)$$

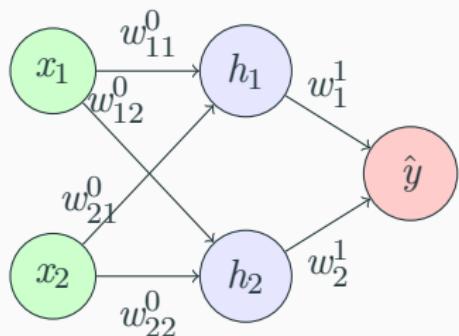
Objective

Determination of the best set of weights \mathbf{w} to minimize the Loss function

$$L(\mathbf{w}) = \|\hat{y}(\mathbf{w}) - y\|^2.$$

Calculation of $\partial L / \partial w$

Training a neural-net: gradient backpropagation



- Given a couple (x, y)
- Forward computation:**
$$h_j = f_0(\sum_{i=1}^2 w_{ij}^0 \cdot x_i)$$
$$\hat{y} = f_1(\sum_{j=1}^2 w_j^1 \cdot h_j)$$
- Compute the gradient of the loss:**

$$\partial L / \partial \hat{y}$$

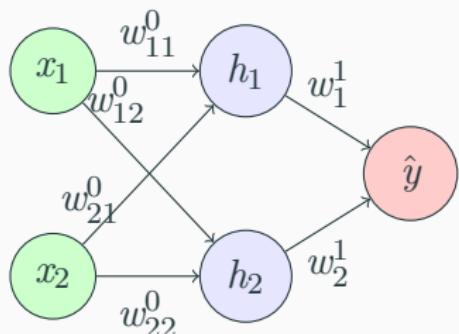
Objective

Determination of the best set of weights \mathbf{w} to minimize the Loss function

$$L(\mathbf{w}) = \|\hat{y}(\mathbf{w}) - y\|^2.$$

Calculation of $\partial L / \partial w$

Training a neural-net: gradient backpropagation



Objective

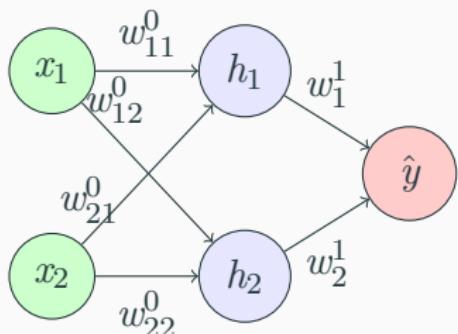
Determination of the best set of weights \mathbf{w} to minimize the Loss function

$$L(\mathbf{w}) = \|\hat{y}(\mathbf{w}) - y\|^2.$$

Calculation of $\partial L / \partial w$

1. Given a couple (x, y)
2. Forward computation:
$$h_j = f_0(\sum_{i=1}^2 w_{ij}^0 \cdot x_i)$$
$$\hat{y} = f_1(\sum_{j=1}^2 w_j^1 \cdot h_j)$$
3. Compute the gradient of the loss:
$$\boxed{\partial L / \partial \hat{y}}$$
4. Gradient Backpropagation:

Training a neural-net: gradient backpropagation



Objective

Determination of the best set of weights \mathbf{w} to minimize the Loss function

$$L(\mathbf{w}) = \|\hat{y}(\mathbf{w}) - y\|^2.$$

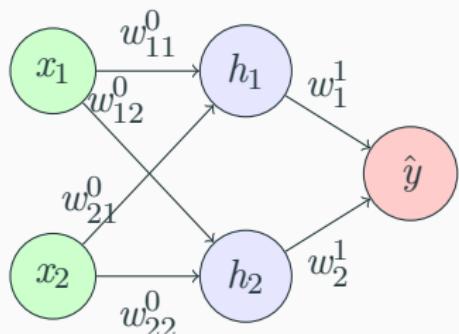
Calculation of $\partial L / \partial w$

- Given a couple (x, y)
- Forward computation:
$$h_j = f_0(\sum_{i=1}^2 w_{ij}^0 \cdot x_i)$$
$$\hat{y} = f_1(\sum_{j=1}^2 w_j^1 \cdot h_j)$$
- Compute the gradient of the loss:
$$\boxed{\partial L / \partial \hat{y}}$$
- Gradient Backpropagation:

- Layer 1

$$\partial L / \partial w_j^1 = \boxed{\partial L / \partial \hat{y}} \cdot \partial f_1 / \partial w_j^1$$
$$\boxed{\partial L / \partial h_j} = \boxed{\partial L / \partial \hat{y}} \cdot \partial f_1 / \partial h_j$$

Training a neural-net: gradient backpropagation



Objective

Determination of the best set of weights \mathbf{w} to minimize the Loss function

$$L(\mathbf{w}) = \|\hat{y}(\mathbf{w}) - y\|^2.$$

Calculation of $\partial L / \partial w$

- Given a couple (x, y)
- Forward computation:
$$h_j = f_0(\sum_{i=1}^2 w_{ij}^0 \cdot x_i)$$

$$\hat{y} = f_1(\sum_{j=1}^2 w_j^1 \cdot h_j)$$
- Compute the gradient of the loss:
$$\boxed{\partial L / \partial \hat{y}}$$
- Gradient Backpropagation:

- Layer 1

$$\partial L / \partial w_j^1 = \boxed{\partial L / \partial \hat{y}} \cdot \partial f_1 / \partial w_j^1$$

$$\boxed{\partial L / \partial h_j} = \boxed{\partial L / \partial \hat{y}} \cdot \partial f_1 / \partial h_j$$

- Layer 0

$$\partial L / \partial w_{ij}^0 = \boxed{\partial L / \partial h_j} \cdot \partial f_1 / \partial w_{ij}^0$$

Optimizing a machine learning (gradient method)

Optimizing the loss

Several loss function (depending on the problem) can be defined.

For example, Mean Square Error:

Method

Find a minimum of L by adjusting the parameters (weights) \mathbf{w} given the gradient of the loss with respect to the weights $\nabla_{\mathbf{w}} L$.

Batch Vs Stochastic training

Dataset: (X, Y) with N samples denoted (\mathbf{x}_i, y_i)

Batch gradient:

Require: Learning rate(s): ν_k

Require: Initial weights: \mathbf{w}

$k \leftarrow 1$

while stopping criterion not met do

 Compute gradient:

$$\mathbf{g} \leftarrow \frac{1}{N} \sum_i^N \nabla_{\mathbf{w}} L(f(\mathbf{x}_i, y_i))$$

 Update weights: $\mathbf{w} \leftarrow \mathbf{w} - \nu_k \mathbf{g}$

$k \leftarrow k + 1$

end while

1 Update / N forwards

Batch Vs Stochastic training

Dataset: (X, Y) with N samples denoted (\mathbf{x}_i, y_i)

Batch gradient:

Require: Learning rate(s): ν_k

Require: Initial weights: \mathbf{w}

$k \leftarrow 1$

while stopping criterion not met **do**

 Compute gradient:

$$\mathbf{g} \leftarrow \frac{1}{N} \sum_i^N \nabla_{\mathbf{w}} L(f(\mathbf{x}_i, y_i))$$

 Update weights: $\mathbf{w} \leftarrow \mathbf{w} - \nu_k \mathbf{g}$

$k \leftarrow k + 1$

end while

1 Update / N forwards

Stochastic gradient:

Require: Learning rate(s): ν_k

Require: Initial weights: \mathbf{w}

$k \leftarrow 1$

while stopping criterion not met **do**

 Sample an example (\mathbf{x}, y) from (X, Y)

 Compute gradient: $\mathbf{g} \leftarrow \nabla_{\mathbf{w}} L(f(\mathbf{x}, y))$

 Update weights: $\mathbf{w} \leftarrow \mathbf{w} - \nu_k \mathbf{g}$

$k \leftarrow k + 1$

end while

1 Update / 1 forward

Mini-Batch training

Dataset: (X, y) with N samples

Mini-Batch gradient:

Require: Learning rate(s): ν_k

Require: Initial weights: \mathbf{w}

$k \leftarrow 1$

while stopping criterion not met do

 Sample m examples (\mathbf{x}_i, y_i) from (X, y)

 Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \sum_i^m \nabla_{\mathbf{w}} L(f(\mathbf{x}_i, y_i))$

 Update weights: $\mathbf{w} \leftarrow \mathbf{w} - \nu_k \mathbf{g}$

$k \leftarrow k + 1$

end while

Mini-Batch training

Dataset: (X, y) with N samples

Mini-Batch gradient:

Require: Learning rate(s): ν_k

Require: Initial weights: \mathbf{w}

$k \leftarrow 1$

while stopping criterion not met do

 Sample m examples (\mathbf{x}_i, y_i) from (X, y)

 Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \sum_i^m \nabla_{\mathbf{w}} L(f(\mathbf{x}_i, y_i))$

 Update weights: $\mathbf{w} \leftarrow \mathbf{w} - \nu_k \mathbf{g}$

$k \leftarrow k + 1$

end while

1 Update / m forward

$m = 1$: Pure stochastic gradient.

$m = N$: Batch gradient

Let's have a break

<https://playground.tensorflow.org>

Other regularization techniques

Regularization

Definition:

Regularization refers to the set of techniques that constraints the optimization. It is generally used to avoid overfitting, but can also be used to inject prior knowledge during the training phase (e.g. set known limits to parameters)

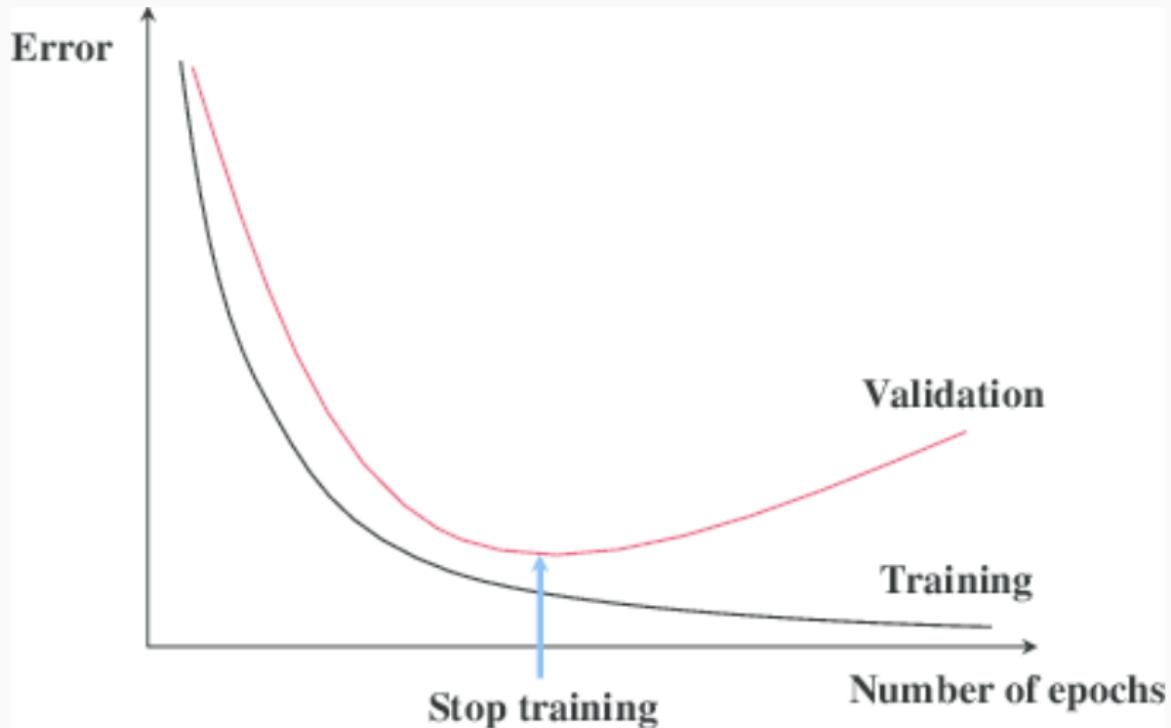
Regularization

Definition:

Regularization refers to the set of techniques that constraints the optimization. It is generally used to avoid overfitting, but can also be used to inject prior knowledge during the training phase (e.g. set known limits to parameters)

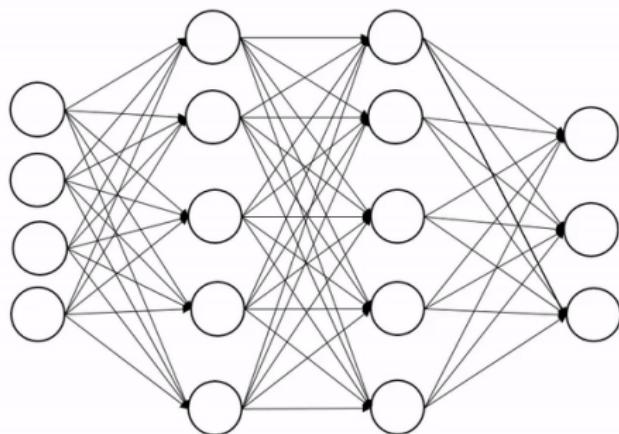
- Stochastic mini-batch gradient is a regularization technique

Early Stopping



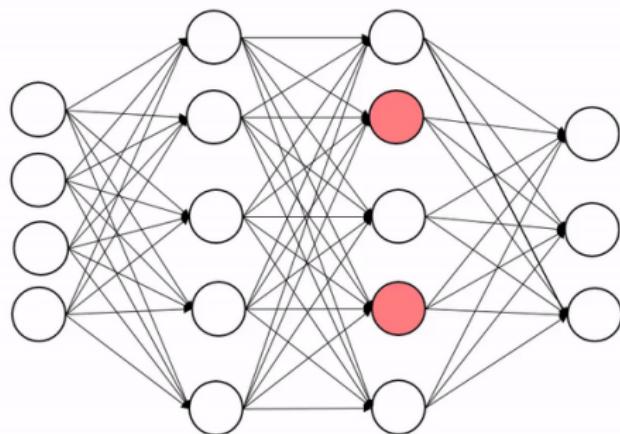
Dropout

During training, randomly remove neurons on a layer with probability p .



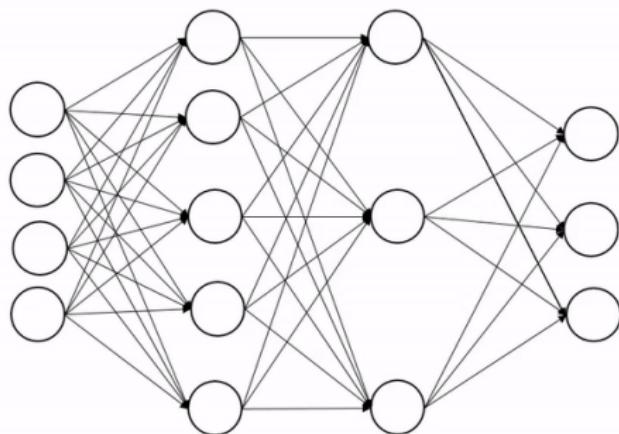
Dropout

During training, randomly remove neurons on a layer with probability p .



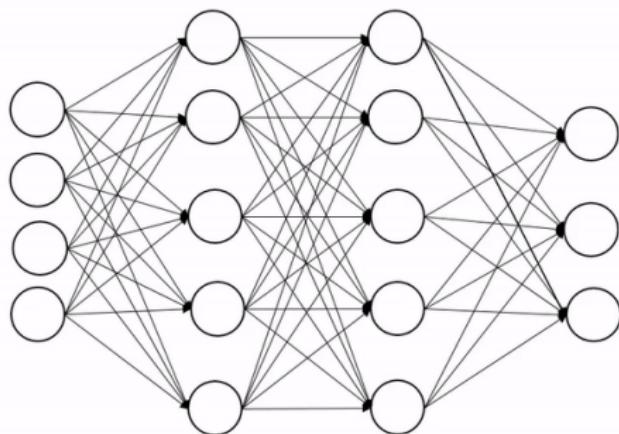
Dropout

During training, randomly remove neurons on a layer with probability p .



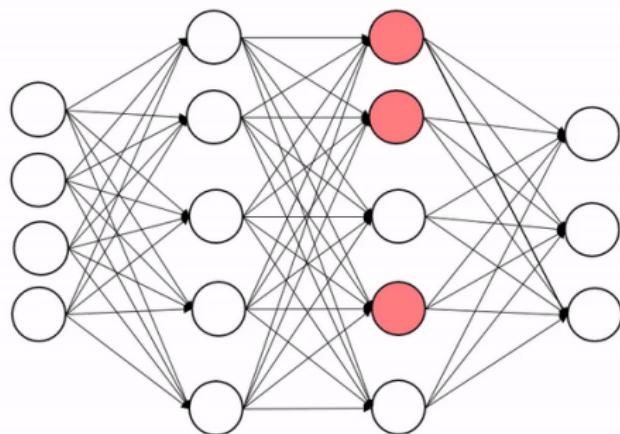
Dropout

During training, randomly remove neurons on a layer with probability p .



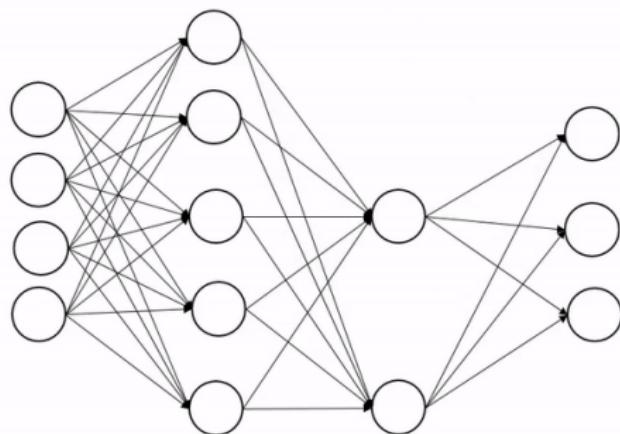
Dropout

During training, randomly remove neurons on a layer with probability p .



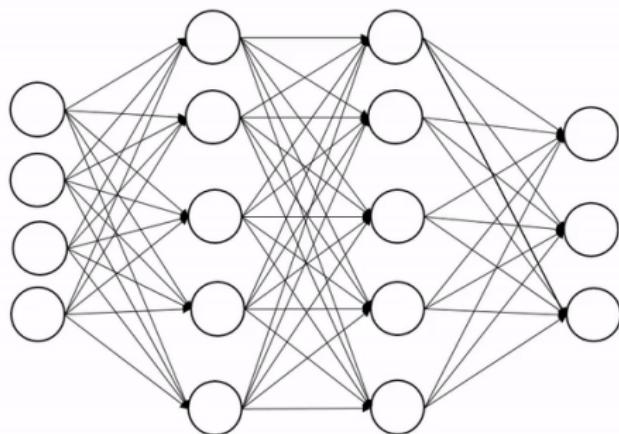
Dropout

During training, randomly remove neurons on a layer with probability p .



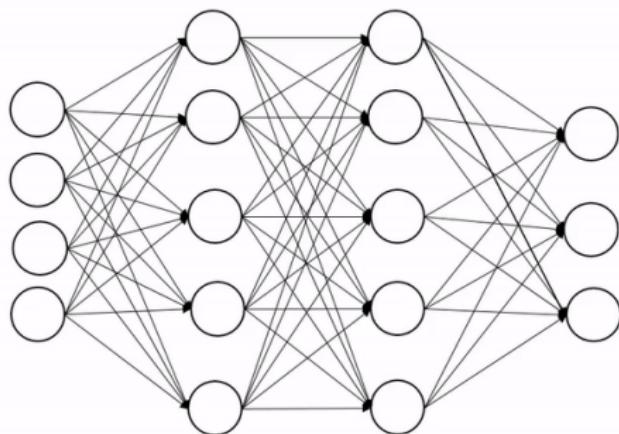
Dropout

During training, randomly remove neurons on a layer with probability p .



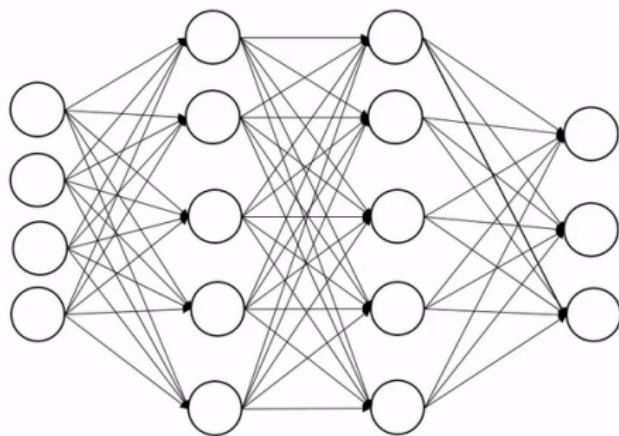
Dropout

During training, randomly remove neurons on a layer with probability p .



Dropout

During training, randomly remove neurons on a layer with probability p .



Practical remarks:

- Avoid Dropout on convolutive layer
- Avoid Dropout on the last layer.

Batch normalization

From Ioffe et al. 2015, Batch normalization...

Batch Normalization is a new type of layer.

If we use mini-batch training with a minibatch of size m :

Mini Batch Layer:

Input: Values of $\mathbf{x}_1 \dots \mathbf{x}_m$

Input: Initial parameters to be optimized: γ, β

Output: $\mathbf{z}_i = BN_{\gamma, \beta}(\mathbf{x}_i)$

$$\mu \leftarrow \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i \quad \triangleright \text{mini-batch mean}$$

$$\sigma^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i - \mu)^2 \quad \triangleright \text{(mini-batch variance)}$$

$$\hat{\mathbf{x}}_i \leftarrow (\mathbf{x}_i - \mu) / \sqrt{\sigma^2 + \epsilon} \quad \triangleright \text{normalize}$$

$$\mathbf{z}_i = \gamma \hat{\mathbf{x}}_i + \beta \quad \triangleright \text{Scale and shift}$$

return \mathbf{z}_i

μ and σ^2 are **non-trainable parameters**. They are fixed for inferring new result (in test/validation).

Link with data assimilation

Example of **BLUE**: Best Linear Unbiased Estimator

Given a state vector $\mathbf{x} \in \mathbb{R}^n$ and a data vector $\mathbf{d} \in \mathbb{R}^m$:

$$\begin{aligned}\mathbf{x}^f &= \mathbf{x}^t + \mathbf{p}, & \bar{\mathbf{p}} &= 0, & \overline{\mathbf{p}\mathbf{p}^T} &= \mathbf{C}_{xx}. \\ \mathbf{d} &= \mathbf{H}\mathbf{x}^t + \boldsymbol{\epsilon}, & \bar{\boldsymbol{\epsilon}} &= 0, & \overline{\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T} &= \mathbf{C}_{\epsilon\epsilon}.\end{aligned}$$

Data Assimilation

Example of **BLUE**: Best Linear Unbiased Estimator

Given a state vector $\mathbf{x} \in \mathbb{R}^n$ and a data vector $\mathbf{d} \in \mathbb{R}^m$:

$$\mathbf{x}^f = \mathbf{x}^t + \mathbf{p}, \quad \bar{\mathbf{p}} = 0, \quad \overline{\mathbf{p}\mathbf{p}^T} = \mathbf{C}_{xx}.$$

$$\mathbf{d} = \mathbf{Hx}^t + \boldsymbol{\epsilon}, \quad \bar{\boldsymbol{\epsilon}} = 0, \quad \overline{\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T} = \mathbf{C}_{\epsilon\epsilon}. \quad \text{Estimating}$$

\mathbf{x}^t by minimizing the estimation error leads to minizing the following function:

$$\mathcal{J}(\mathbf{x}) = (\mathbf{d} - \mathbf{Hx})^T \mathbf{C}_{\epsilon\epsilon}^{-1} (\mathbf{d} - \mathbf{Hx}) + (\mathbf{x} - \mathbf{x}^f)^T \mathbf{C}_{xx}^{-1} (\mathbf{x} - \mathbf{x}^f)$$

Data Assimilation

Example of **BLUE**: Best Linear Unbiased Estimator

Given a state vector $\mathbf{x} \in \mathbb{R}^n$ and a data vector $\mathbf{d} \in \mathbb{R}^m$:

$$\mathbf{x}^f = \mathbf{x}^t + \mathbf{p}, \quad \bar{\mathbf{p}} = 0, \quad \overline{\mathbf{p}\mathbf{p}^T} = \mathbf{C}_{xx}.$$

$$\mathbf{d} = \mathbf{Hx}^t + \boldsymbol{\epsilon}, \quad \bar{\boldsymbol{\epsilon}} = 0, \quad \overline{\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T} = \mathbf{C}_{\epsilon\epsilon}. \quad \text{Estimating}$$

\mathbf{x}^t by minimizing the estimation error leads to minizing the following function:

$$\mathcal{J}(\mathbf{x}) = (\mathbf{d} - \mathbf{Hx})^T \mathbf{C}_{\epsilon\epsilon}^{-1} (\mathbf{d} - \mathbf{Hx}) + (\mathbf{x} - \mathbf{x}^f)^T \mathbf{C}_{xx}^{-1} (\mathbf{x} - \mathbf{x}^f)$$

This is data assimilation!

We correct a forecast \mathbf{x}^f given some observational data \mathbf{d}

Is it machine learning?

- Ridge regression: $J(\boldsymbol{\theta}) = (\mathbf{y} - h_{\boldsymbol{\theta}}(\mathbf{x}))^T(\mathbf{y} - h_{\boldsymbol{\theta}}(\mathbf{x})) + \alpha\boldsymbol{\theta}^T\boldsymbol{\theta}$
- BLUE: $\mathcal{J}(\mathbf{x}) = (\mathbf{d} - \mathbf{Hx})^T\mathbf{C}_{\epsilon\epsilon}^{-1}(\mathbf{d} - \mathbf{Hx}) + (\mathbf{x} - \mathbf{x}^f)^T\mathbf{C}_{xx}^{-1}(\mathbf{x} - \mathbf{x}^f)$

Is it machine learning?

- Ridge regression: $J(\boldsymbol{\theta}) = (\mathbf{y} - h_{\boldsymbol{\theta}}(\mathbf{x}))^T(\mathbf{y} - h_{\boldsymbol{\theta}}(\mathbf{x})) + \alpha\boldsymbol{\theta}^T\boldsymbol{\theta}$
- BLUE: $\mathcal{J}(\mathbf{x}) = (\mathbf{d} - \mathbf{H}\mathbf{x})^T\mathbf{C}_{\epsilon\epsilon}^{-1}(\mathbf{d} - \mathbf{H}\mathbf{x}) + (\mathbf{x} - \mathbf{x}^f)^T\mathbf{C}_{xx}^{-1}(\mathbf{x} - \mathbf{x}^f)$

BLUE	Ridge
data \mathbf{d}	target \mathbf{y}
Observation operator \mathbf{H}	feature \mathbf{x}
State \mathbf{x}	parameters $\boldsymbol{\theta}$
$\mathbf{C}_{\epsilon\epsilon}\mathbf{C}_{xx}^{-1}$	α

Is it machine learning?

- Ridge regression: $J(\boldsymbol{\theta}) = (\mathbf{y} - h_{\boldsymbol{\theta}}(\mathbf{x}))^T(\mathbf{y} - h_{\boldsymbol{\theta}}(\mathbf{x})) + \alpha\boldsymbol{\theta}^T\boldsymbol{\theta}$
- BLUE: $\mathcal{J}(\mathbf{x}) = (\mathbf{d} - \mathbf{H}\mathbf{x})^T\mathbf{C}_{\epsilon\epsilon}^{-1}(\mathbf{d} - \mathbf{H}\mathbf{x}) + (\mathbf{x} - \mathbf{x}^f)^T\mathbf{C}_{xx}^{-1}(\mathbf{x} - \mathbf{x}^f)$

BLUE	Ridge
data \mathbf{d}	target \mathbf{y}
Observation operator \mathbf{H}	feature \mathbf{x}
State \mathbf{x}	parameters $\boldsymbol{\theta}$
$\mathbf{C}_{\epsilon\epsilon}\mathbf{C}_{xx}^{-1}$	α

Further links...

If a numerical model is integrated over several time steps, it can be related to successive layers of a Neural Network.

Probabilistic interpretation

Maximum likelihood estimator and loss

We can assume that the observation y follows a Gaussian law:

$$p(y/x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp -\frac{(y - \mu(x))^2}{2\sigma^2},$$

where x is observed and $\mu(x)$ is a function of x .

Given a set of samples $(x_k, y_k)_{1:K}$, the negative log-likelihood is defined by

$$L = \sum_{k=1}^K \left(\frac{\log 2\pi\sigma^2}{2} + \frac{(y_k - \mu(x_k))^2}{2\sigma^2} \right)$$

Minimizing L is maximizing the probability of having the observations y_k given x_k .

Loss function of a neural net

First case: σ is constant

$\mu(x)$ is parametrized by a neural net $G_\mu(x, \theta_\mu)$

The Maximum likelihood estimator is found by minizming:

$$L(\theta_\mu) = \sum_{k=1}^K (y_k - G_\mu(x, \theta_\mu))^2$$

which is exactly the regression loss already introduced.

Loss function of a neural net

First case: σ is constant

$\mu(x)$ is parametrized by a neural net $G_\mu(x, \theta_\mu)$

The Maximum likelihood estimator is found by minimizing:

$$L(\theta_\mu) = \sum_{k=1}^K (y_k - G_\mu(x, \theta_\mu))^2$$

which is exactly the regression loss already introduced.

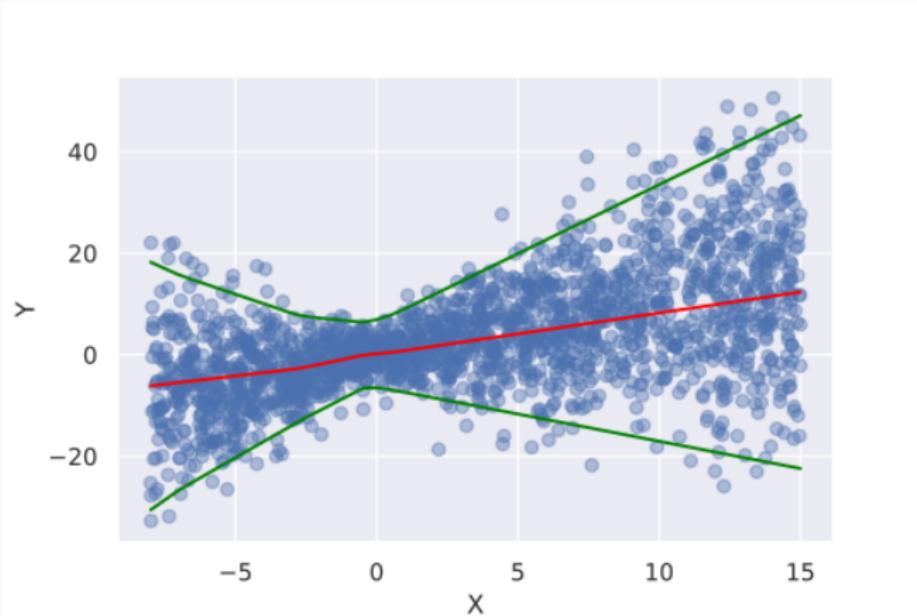
Second case: $\sigma(x)$ is a function of x .

In addition to $G_\mu(x, \theta_\mu)$, $\sigma(x)$ is parametrized by $G_\sigma(x, \theta_\sigma)$. The loss to minimize is then:

$$L(\theta_\mu, \theta_\sigma) = \sum_{k=1}^K \left(\frac{\log 2\pi G_\sigma(x_k, \theta_\sigma)^2}{2} + \frac{(y_k - G_\mu(x_k, \theta_\mu))^2}{2G_\sigma(x_k, \theta_\sigma)^2} \right)$$

The neural net $G = (G_\mu, G_\sigma)$ gives also the uncertainty of its estimation in the form of the standard deviation.

Illustration

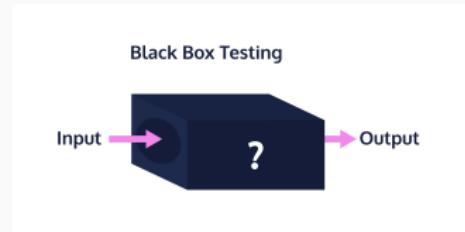


In red the estimation of the mean $G_\mu(x, \theta_\mu)$
In green the confidence interval $G_\mu(x, \theta_\mu) \pm \sigma(x_k, \theta_\sigma)$

A black box?

The black box paradigm

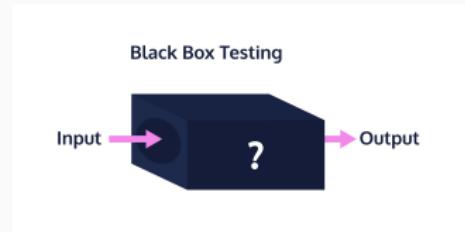
The machine-learning based model is a **black box**. It gives some results but we don't understand how.



Do we need to understand the model?

The black box paradigm

The machine-learning based model is a **black box**. It gives some results but we don't understand how.



Do we need to understand the model?

"Every time I fire a linguist, the performance of our speech recognition system goes up."

F. Jelinek, 1988



Frederick Jelinek 1932-2010

Motivation

- Build models that can be trusted
- Use other source of knowledge (e.g. physical properties) when data are not sufficient.

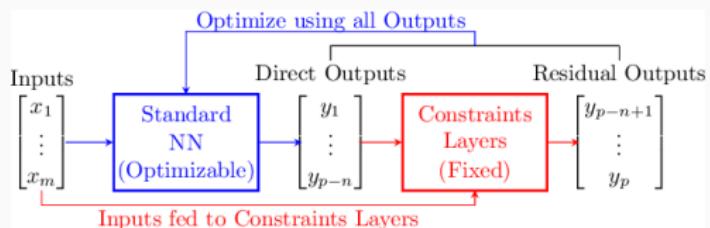
Two directions:

- Add physical constraints to ML models
- Explainable/Transparent ML.

Add physical constraints to ML models

- **Simple example:** enforce the positivity of some quantities (e.g. concentration)
- **More complex:** enforce conservation laws

Beucler, T., Pritchard, M., Rasp, S., Ott, J., Baldi, P. and Gentine, P., 2021. Enforcing analytic constraints in neural networks emulating physical systems. *Physical Review Letters*, 126(9), p.098302.



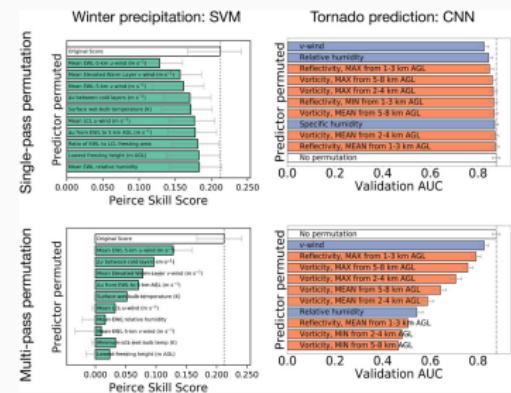
Beucler et al.

Explainable/Transparent ML.

The objective is to understand how the machine learning makes a prediction (e.g. which feature is important for the prediction).

McGovern et al., 2019, Making the black box more transparent: Understanding the physical implications of machine learning. *Bulletin of the American Meteorological Society*

Sonnewald et al., 2021.
Bridging observation, theory and numerical simulation of the ocean using Machine Learning. *Environ. Res. Lett.*



McGovern et al.

Questions addressed in this lecture

- What is gradient backpropagation? [GBC16,6]
- What are the different types of gradient descent techniques? [GBC16,8]
- What are some other types of regularizations? [GBC16,7]
- How machine learning can be related to data assimilation?

Refs

[Van16,*n*]: Jake VanderPlas, *Python Data Science Handbook*, section *n*

[GBC16,*n*]: Goodfellow et al., Deep Learning, chapter *n*