

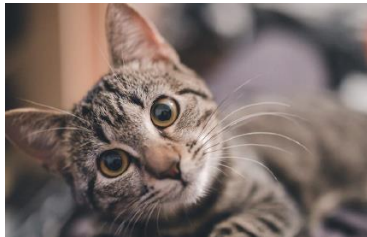
# Introduction au machine learning

ENSTA 2024

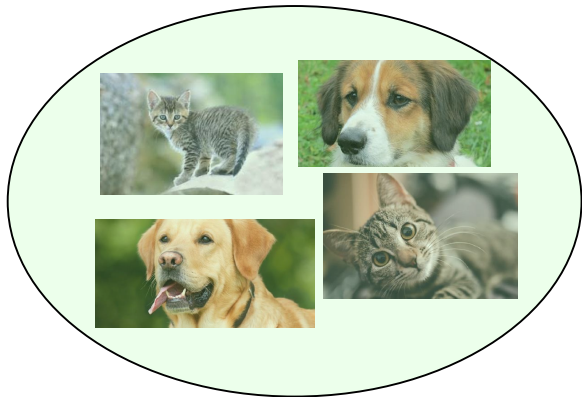
Adrien CHAN-HON-TONG

HDR ONERA

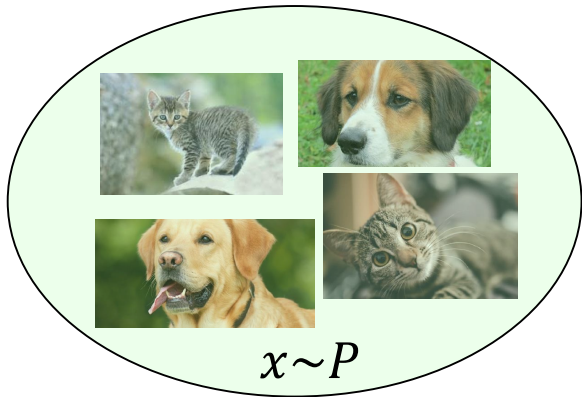
# objectif



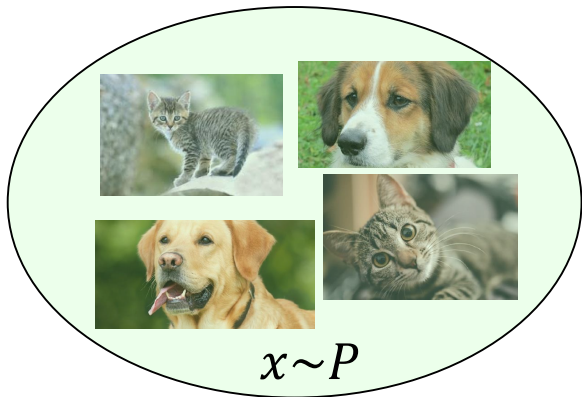
objectif



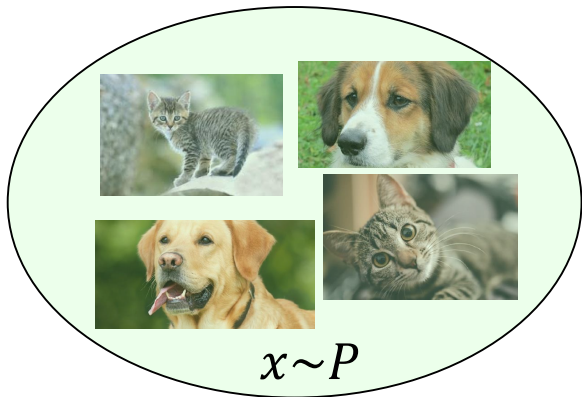
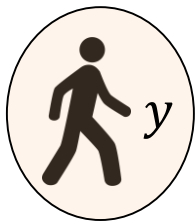
objectif



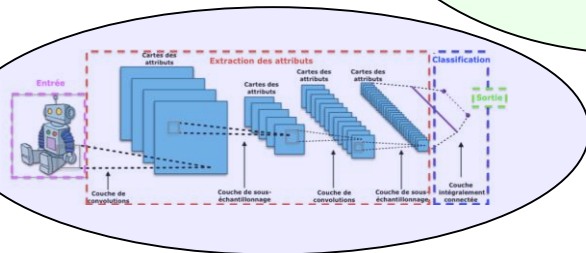
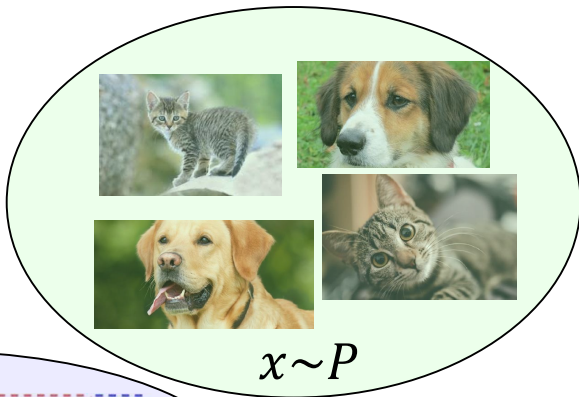
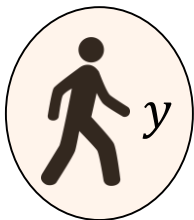
objectif



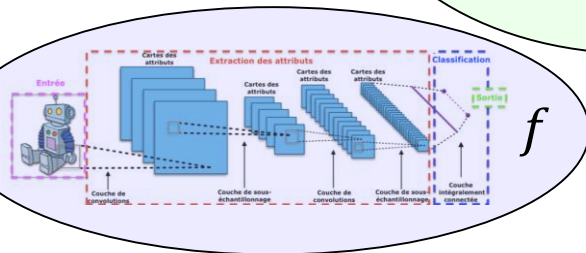
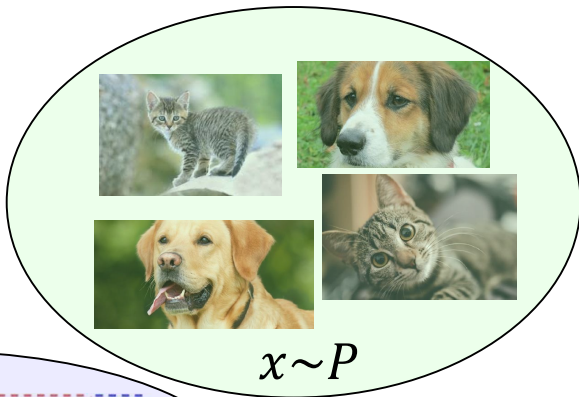
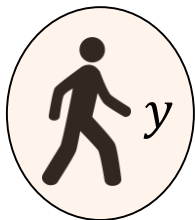
objectif



# objectif

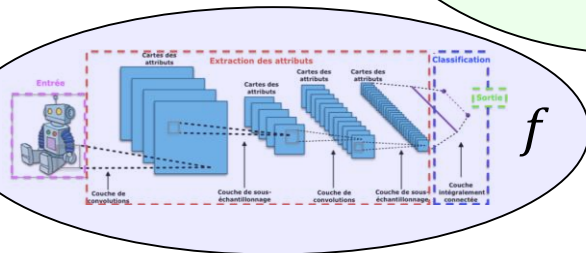
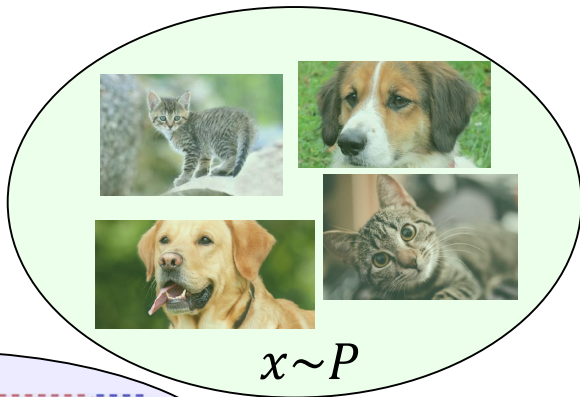
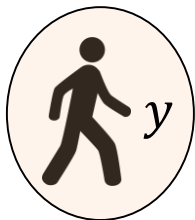


# objectif





# objectif



$$f(x) \approx y(x)$$

$$x \sim P$$

# Notion d'erreurs

$$f(x) \approx y(x)$$
$$x \sim P$$

En fait, c'est pas la valeur de  $f$  qui compte  
C'est le fait qu'on y associe la bonne valeur de  $y$ .

## Notion d'erreurs

$$f(x) \approx y(x) \\ x \sim P$$

En fait, c'est pas la valeur de  $f$  qui compte  
C'est le fait qu'on y associe la bonne valeur de  $y$ .

Dans le cas binaire  $y(x) \in \{-1, 1\}$

$$erreur = \int \mathbf{1}_{]-\infty, 0[}(f(x)y(x))P(x)dx$$

Notre objectif est de trouver  $f$  qui vérifierait

$$f(x) \approx y(x) \quad x \sim P$$

à partir de quelques exemples  $x_1, \dots, x_N \sim P$  annotés

(c'est-à-dire qu'on connaît  $y$  pour ces exemples)

(par exemple, parce qu'on les a montré à des humains)

Notre objectif est de trouver  $f$  qui vérifierait

$$\int \mathbf{1}_{]-\infty, 0[}(f(x)y(x))P(x)dx \ll 1$$

à partir de quelques exemples  $x_1, \dots, x_N \sim P$  annotés

(c'est-à-dire qu'on connaît  $y$  pour ces exemples)

(par exemple, parce qu'on les a montré à des humains)

Notre objectif est de trouver  $f$  qui vérifierait

$$\int \mathbf{1}_{]-\infty, 0[}(f(x)y(x))P(x)dx \ll 1$$

à partir de quelques exemples  $x_1, \dots, x_N \sim P$  annotés

(c'est-à-dire qu'on connaît  $y$  pour ces exemples)

(par exemple, parce qu'on les a montré à des humains)

**D'un point de vue théorique,**

**ce problème n'est PAS assez bien défini,**

**pour qu'on puisse avoir une méthode systématique !**

Notre objectif est de trouver  $f$  qui vérifierait

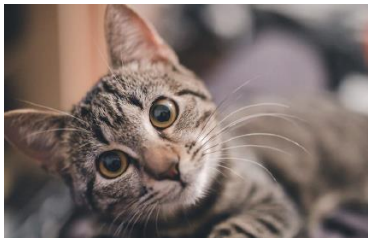
$$\int \mathbf{1}_{]-\infty, 0[}(f(x)y(x))P(x)dx \ll 1$$

à partir de quelques exemples  $x_1, \dots, x_N \sim P$  annotés

Attention. Je n'ai pas dit « on ne peut pas ».  
Je dit « il n'y a pas de méthode systematique ».

*The lack of a priori distinction between classifier*

# Un problème mal défini

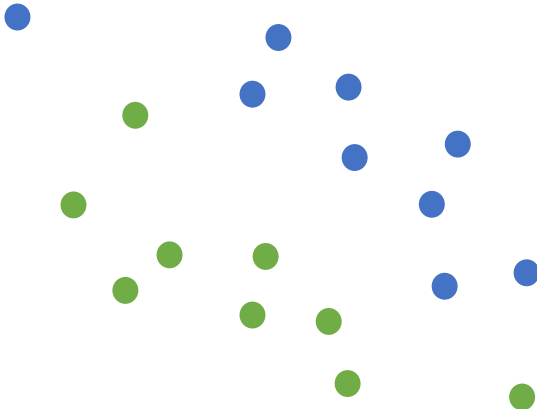




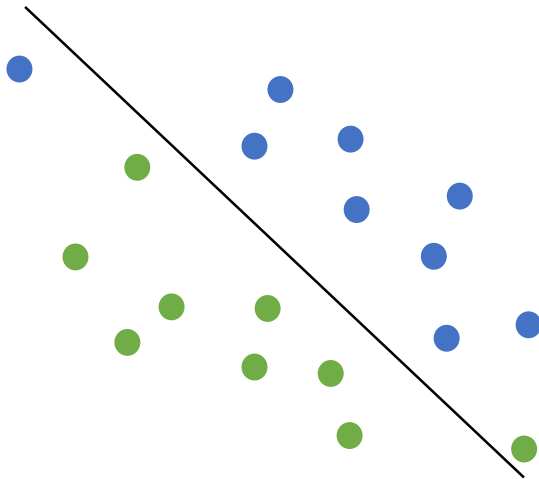
# Un problème mal défini



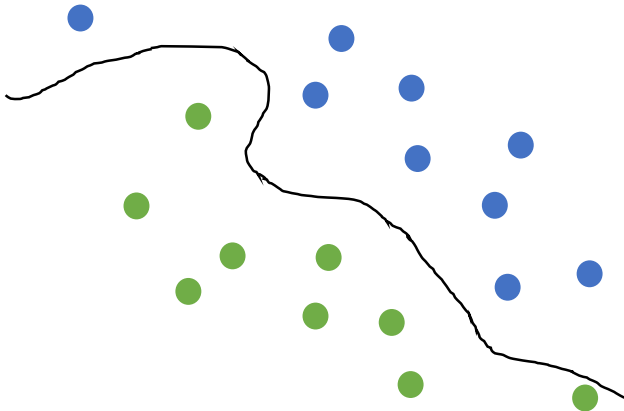
# Un problème mal défini



# Un problème mal défini



# Un problème mal défini



# Quelques résultats « théoriques »

Cas  $N = \infty$

Bornes statistiques  
pour l'évaluation  
Vapnik

# Cas $N = \infty$

( $y$  est connu pour les exemples d'apprentissage  $x_1, \dots, x_N$  )

## Le plus proche voisin

$$f(x) = y(x_i) ; ||x - x_i|| = \min_{\{x_j \in \text{Base Apprentissage}\}} ||x - x_j||$$

Cas  $N = \infty$

( $y$  est connu pour les exemples d'apprentissage  $x_1, \dots, x_N$  )

Le plus proche voisin

$$f(x) = y(x_i) ; ||x - x_i|| = \min_{\{x_j \in \text{Base Apprentissage}\}} ||x - x_j||$$

$$\int \mathbf{1}_{]-\infty, 0[}(f(x)y(x))P(x)dx \xrightarrow{N \rightarrow \infty} 0$$

Cas  $N = \infty$

( $y$  est connu pour les exemples d'apprentissage  $x_1, \dots, x_N$  )

Le plus proche voisin

$$f(x) = y(x_i) ; ||x - x_i|| = \min_{\{x_j \in \text{Base Apprentissage}\}} ||x - x_j||$$

$$\int \mathbf{1}_{]-\infty, 0[}(f(x)y(x))P(x)dx \xrightarrow{N \rightarrow \infty} 0$$

(Le résultat est plus général, dans le cas où  $y$  est une distribution, il faut utiliser un K-plus proche voisin.)



Cas  $N = \infty$

( $y$  est connu pour les exemples d'apprentissage  $x_1, \dots, x_N$  )

Le plus proche voisin

$$f(x) = y(x_i) ; ||x - x_i|| = \min_{\{x_j \in \text{Base Apprentissage}\}} ||x - x_j||$$

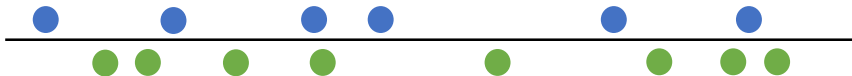
$$\int \mathbf{1}_{]-\infty, 0[}(f(x)y(x))P(x)dx \xrightarrow{N \rightarrow \infty} 0$$

Mais en pratique, ça ne marche pas bien  
hors de la base d'apprentissage ☹

# Pire cas pour le plus proche voisin



# Pire cas pour le plus proche voisin



Concept	+	-
Plus proche voisin	Erreur nulle partout quand la base d'apprentissage est infinie	Ça marche pas bien hors de la base d'apprentissage en pratique

Résultat statistique (« N » non infini)

Peut-on approximer

$$\int \mathbf{1}_{]-\infty, 0[}(f(x)y(x))P(x)dx \quad ?$$

Résultat statistique (« N » non infini)

Peut-on approximer

$$\int \mathbf{1}_{]-\infty, 0[}(f(x)y(x))P(x)dx \quad ?$$

Oui (et non – on en reparlera) :

La moyenne converge vers l'espérance

Résultat statistique (« N » non infini)

Peut-on approximer

$$\int \mathbf{1}_{]-\infty, 0[}(f(x)y(x))P(x)dx \quad ?$$

Oui (et non – on en reparlera) :

La moyenne converge vers l'espérance

Si  $\chi_1, \dots, \chi_K \sim P$  est une base de TEST

$$\frac{1}{K} \sum_k \mathbf{1}_{]-\infty, 0[}(f(\chi_k)y(\chi_k)) \rightarrow \int \mathbf{1}_{]-\infty, 0[}(f(x)y(x))P(x)dx$$

## Résultat statistique (« N » non infini)

Si  $\chi_1, \dots, \chi_K \sim P$  est une base de TEST

$$P\left(\int \mathbf{1}_{]-\infty, 0[}(f(x)y(x))P(x)dx \geq \frac{1}{K} \sum_k \mathbf{1}_{]-\infty, 0[}(f(\chi_k)y(\chi_k)) + \frac{\log(\epsilon)}{\sqrt{K}}\right) \leq \epsilon$$

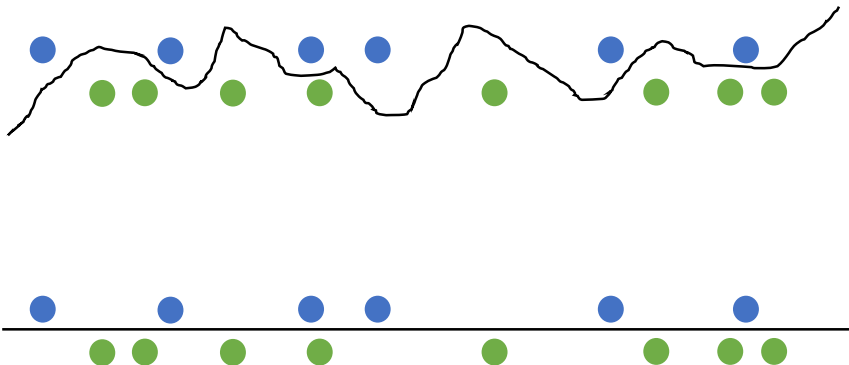


Concept	+	-
Plus proche voisin	Erreur nulle partout quand la base d'apprentissage est infinie	Ça marche pas bien hors de la base d'apprentissage en pratique
Inégalité entre erreur empirique et réelle	Permet de savoir si on a une fonction $f$ qui approxime correctement $y$	Ne dit rien sur comment choisir $f$

# Les phases d'une démarche d'apprentissage automatique

- On part d'un problème  $y, P$
- On récolte une base d'apprentissage  $x_1, \dots, x_N \sim P$
- On annote les exemples d'apprentissage (on connaît alors  $y(x_1), \dots, y(x_N)$ )
- Le cœur de l'apprentissage : à la fin on a choisi  $f$
- On récolte une base de test  $\chi_1, \dots, \chi_K \sim P$
- On annote la base de test
- $$P \left( \int \mathbf{1}_{]-\infty, 0[} (f(x)y(x)) P(x) dx \geq \frac{1}{K} \sum_k \mathbf{1}_{]-\infty, 0[} (f(\chi_k)y(\chi_k)) + \frac{\log(\epsilon)}{\sqrt{K}} \right) \leq \epsilon$$

# « Simplicité » du critère



## « Simplicité » du critère

$$P \left( \int \mathbf{1}_{]-\infty, 0[} (f(x)y(x)) P(x) dx \geq \frac{1}{N} \sum_n \mathbf{1}_{]-\infty, 0[} (f(x_n)y(x_n)) + \frac{VC(\mathcal{F}) + \log(\epsilon)}{\sqrt{K}} \right) \leq \epsilon$$

## « Simplicité » du critère

$$P \left( \int \mathbf{1}_{]-\infty, 0[} (f(x)y(x)) P(x) dx \geq \frac{1}{N} \sum_n \mathbf{1}_{]-\infty, 0[} (f(x_n)y(x_n)) + \frac{VC(\mathcal{F}) + \log(\epsilon)}{\sqrt{K}} \right) \leq \epsilon$$

La dimension de Vapnik mesure la capacité du modèle à apprendre « du bruit ».

## « Simplicité » du critère

$$P \left( \int \mathbf{1}_{]-\infty, 0[} (f(x)y(x)) P(x) dx \geq \frac{1}{N} \sum_n \mathbf{1}_{]-\infty, 0[} (f(x_n)y(x_n)) + \frac{VC(\mathcal{F}) + \log(\epsilon)}{\sqrt{K}} \right) \leq \epsilon$$

La dimension de Vapnik mesure la capacité du modèle à apprendre « du bruit ».  
Elle peut être infinie (cas plus proche voisin).

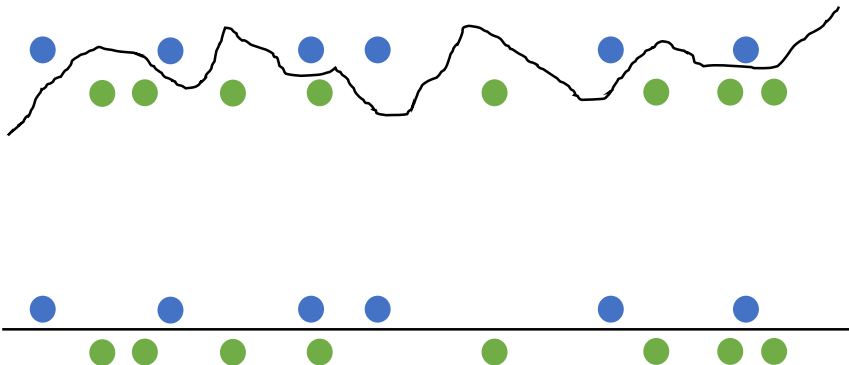
## « Simplicité » du critère

$$P \left( \int \mathbf{1}_{]-\infty, 0[} (f(x)y(x)) P(x) dx \geq \frac{1}{N} \sum_n \mathbf{1}_{]-\infty, 0[} (f(x_n)y(x_n)) + \frac{VC(\mathcal{F}) + \log(\epsilon)}{\sqrt{K}} \right) \leq \epsilon$$

La dimension de Vapnik mesure la capacité du modèle à apprendre « du bruit ». Elle peut être infinie (cas plus proche voisin).

Quand elle est faible, être capable d'avoir une erreur d'apprentissage faible conduit statistiquement à avoir une erreur réelle (et de test) faible !

# « Simplicité » du critère





# « Simplicité » du critère

$$P \left( \int \mathbf{1}_{]-\infty, 0[} (f(x)y(x)) P(x) dx \geq \frac{1}{N} \sum_n \mathbf{1}_{]-\infty, 0[} (f(x_n)y(x_n)) + \frac{VC(\mathcal{F}) + \log(\epsilon)}{\sqrt{K}} \right) \leq \epsilon$$

La dimension de Vapnik mesure la capacité du modèle à apprendre « du bruit ». Elle peut être infinie (cas plus proche voisin).

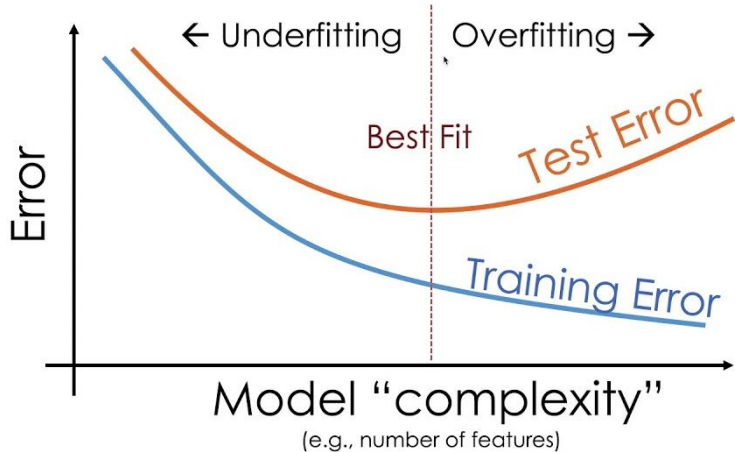
Quand elle est faible, être capable d'avoir une erreur d'apprentissage faible conduit statistiquement à avoir une erreur réelle (et de test) faible !

**En faible dimension, si un classifieur linéaire sépare vos exemples d'apprentissage en 2, il est probablement pertinent partout !**

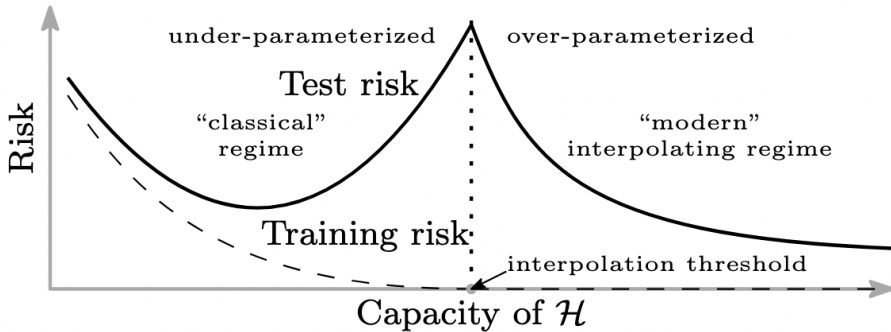


Concept	+	-
Plus proche voisin	Erreur nulle partout quand la base d'apprentissage est infinie	Ça marche pas bien hors de la base d'apprentissage en pratique
Inégalité entre erreur empirique et réelle	Permet de savoir si on a une fonction $f$ qui approxime correctement $y$	Ne dit rien sur comment choisir $f$
Classifier linéaire	Erreur d'apprentissage et réelle corrélées (en faible dimension)	Souvent mauvais même à l'apprentissage

Trouver un compromis ?  
(à N fixé)



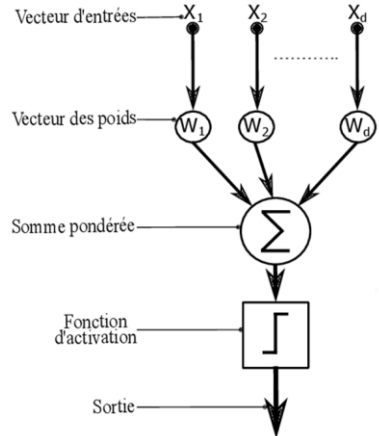
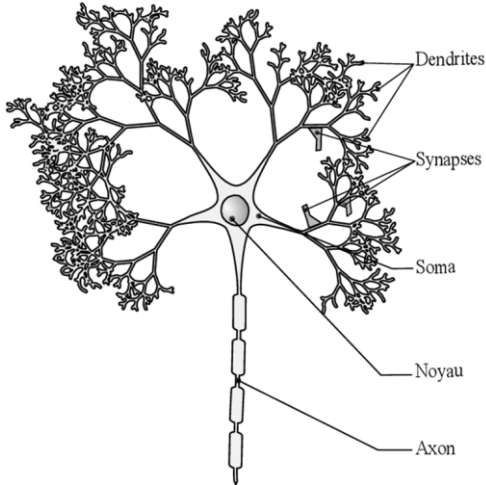
# Ce qu'on observe



# Ce qui marche le mieux aujourd'hui : les réseaux de neurones !

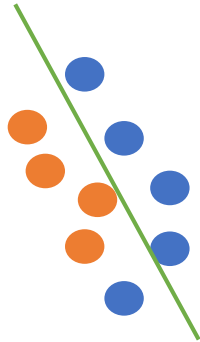
- COURS 1
  - **MLP**
  - Double descente ?
  - Apprentissage
- Cours 2
  - CNN
  - Segmentation/détection
  - Transformer et perspective

# Réseau de neurones



# La classification

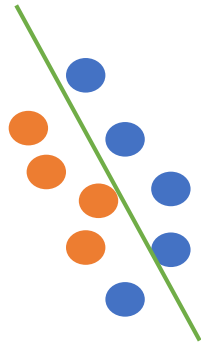
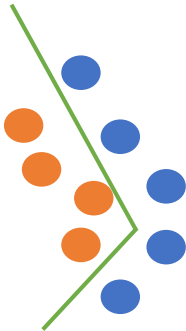
1 neurone  
ne peut être que linéaire





# La classification

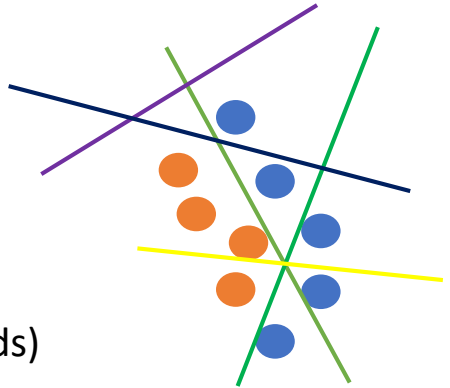
1 neurone  
ne peut être que linéaire



4 neurones + activation  
**peut** être une fonction **non** linéaire

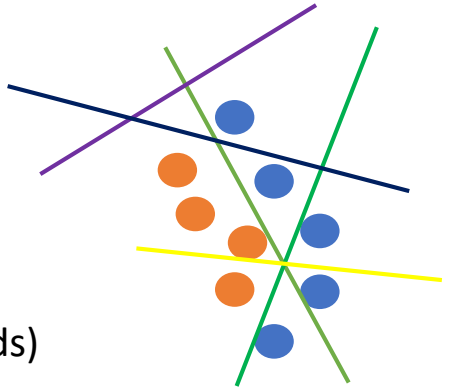
# La classification

1 neurone  
peut coder TOUTES  
ces fonctions  
(via la valeur des poids)



# La classification

1 neurone  
peut coder TOUTES  
ces fonctions  
(via la valeur des poids)



Le SVM en est 1 seule fonction  
(optimisation des poids)

# La classification



4 neurones + activation

=> On ne peut même pas dessiner  
l'ensemble des fonctions qu'on peut  
obtenir via les valeurs des poids !

Il va falloir en choisir 1

# MLP

Notebook illustration

+

<https://playground.tensorflow.org>

# MLP



1 neurone :  $f(x) = w^T x + b = \sum_{d=1}^D w_d x_d + b$



1 couche de 2 neurones + 1 neurone :

$$f(x) = w_3^T \begin{pmatrix} \phi(w_1^T x + b_1) \\ \phi(w_2^T x + b_2) \end{pmatrix} + b_3$$



1 couche de 2 neurones  
+ 1 couche de 3 neurones  
+ 1 neurone :

$$f(x) = w_6^T \begin{pmatrix} \phi(w_3^T (\phi(w_1^T x + b_1) + b_3)) \\ \phi(w_4^T (\phi(w_1^T x + b_1) + b_4)) \\ \phi(w_5^T (\phi(w_1^T x + b_1) + b_5)) \end{pmatrix} + b_6$$

# MLP



1 neurone :  $f(x) = Ax + b$



1 couche de 2 neurones + 1 neurone :  
 $f(x) = A_2(\phi(A_1x + b_1)) + b_2$



1 couche de 2 neurones  
+ 1 couche de 3 neurones  
+ 1 neurone :

$$f(x) = A_3(\phi(A_2(\phi(A_1x + b_1)) + b_2)) + b_3$$

# MLP

la famille de MLP à activation  $\phi$  et à  $Q$  couches de taille  $i_1, \dots, i_Q$  est la famille des fonctions qui peuvent s'écrire

$$f(x) = A_Q \left( \phi(A_{Q-1}(\dots \phi(A_2(\phi(A_1x + b_1)) + b_2) \dots) + b_{Q-1}) \right) + b_Q$$

$$\text{Avec } A_q \in \mathbb{R}^{i_q \times i_{q-1}} \text{ et } b_q \in \mathbb{R}^{i_q}$$

Convention  $i_0$  correspond à la dimension des entrées «  $x$  »

L'activation la plus classique est :  $relu(t) = \max(t, 0)$



# MLP

Notebook illustration

+

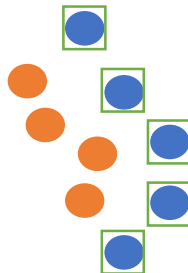
<https://playground.tensorflow.org>

# PLAN

- COURS 1
  - MLP
  - **Double descent ?**
  - Apprentissage
- Cours 2
  - CNN
  - Segmentation/détection
  - Transformer et perspective

# Théorème d'universalité

Avec 100 neurones,  
on peut faire



# Théorème d'universalité

$$\forall x \in \mathbb{R}, |x| = \text{relu}(x) + \text{relu}(-x)$$

# Théorème d'universalité

$$\forall x \in \mathbb{R}, |x| = \text{relu}(x) + \text{relu}(-x)$$

$$\forall x \in \mathbb{R}^D, \|x\|_1 = \mathbf{1}^T (\text{relu}(Ix) + \text{relu}(-Ix))$$

# Théorème d'universalité

$$\forall x \in \mathbb{R}, |x| = \text{relu}(x) + \text{relu}(-x)$$

$$\forall x \in \mathbb{R}^D, \|x\|_1 = \mathbf{1}^T (\text{relu}(Ix) + \text{relu}(-Ix))$$

$$\forall x, u \in \mathbb{R}^D, \|x - u\|_1 = \mathbf{1}^T (\text{relu}(Ix - u) + \text{relu}(-Ix + u))$$

# Théorème d'universalité

$$\forall x \in \mathbb{R}, |x| = \text{relu}(x) + \text{relu}(-x)$$

$$\forall x \in \mathbb{R}^D, \|x\|_1 = \mathbf{1}^T (\text{relu}(Ix) + \text{relu}(-Ix))$$

$$\forall x, u \in \mathbb{R}^D, \|x - u\|_1 = \mathbf{1}^T (\text{relu}(Ix - u) + \text{relu}(-Ix + u))$$

$$\forall x, u \in \mathbb{R}^D, \text{relu}(1 - \|x - u\|_1)$$

$$= \text{relu}(1 - \mathbf{1}^T (\text{relu}(Ix - u) + \text{relu}(-Ix + u)))$$

# Théorème d'universalité

$$\forall x \in \mathbb{R}, |x| = \text{relu}(x) + \text{relu}(-x)$$

$$\forall x \in \mathbb{R}^D, \|x\|_1 = \mathbf{1}^T (\text{relu}(Ix) + \text{relu}(-Ix))$$

$$\forall x, u \in \mathbb{R}^D, \|x - u\|_1 = \mathbf{1}^T (\text{relu}(Ix - u) + \text{relu}(-Ix + u))$$

$$\forall x, u \in \mathbb{R}^D, \text{relu}(1 - \|x - u\|_1)$$

$$= \text{relu} \left( 1 - \mathbf{1}^T (\text{relu}(Ix - u) + \text{relu}(-Ix + u)) \right)$$

$$\forall x, u, v \in \mathbb{R}^D, \begin{pmatrix} \text{relu}(1 - \|x - u\|_1) \\ \text{relu}(1 - \|x - v\|_1) \end{pmatrix}$$

$$= \text{relu} \left( \mathbf{1} - \begin{pmatrix} 1 & \dots & 1 & 0 & \dots & 0 \\ 0 & \dots & 0 & 1 & \dots & 1 \end{pmatrix} \text{relu} \begin{pmatrix} Ix - u \\ -Ix + u \\ Ix - v \\ -Ix + v \end{pmatrix} \right)$$

$$= \left[ \mathbf{1} - \begin{pmatrix} 1 & \dots & 1 & 0 & \dots & 0 \\ 0 & \dots & 0 & 1 & \dots & 1 \end{pmatrix} \begin{bmatrix} Ix - u \\ -Ix + u \\ Ix - v \\ -Ix + v \end{bmatrix} \right]_{+} +$$



# Théorème d'universalité

Soit la matrice  $W \in \mathbb{R}^{N \times (N \times 2D)}$  tel que  $\forall n \in \{1, \dots, N\}, i \in \{1, \dots, N \times 2D\}$ ,

$$W_{n,i} = \begin{cases} 1 & \text{ssi } 2nD \leq i < 2nD + 2D \\ 0 & \text{sinon} \end{cases}$$

alors,  $\forall x_1, \dots, x_N \in \mathbb{R}^D$

$$\left[ \mathbf{1} - \begin{pmatrix} \|x - x_1\|_1 \\ \|x - x_2\|_1 \\ \dots \\ \|x - x_N\|_1 \end{pmatrix} \right]_+ = \left[ \mathbf{1} - W \begin{bmatrix} Ix - x_1 \\ -Ix + x_1 \\ Ix - x_2 \\ -Ix + x_2 \\ \dots \\ Ix - x_N \\ -Ix + x_N \end{bmatrix} \right]_+ \Big|_+$$

# Théorème d'universalité

Soit la matrice  $W \in \mathbb{R}^{N \times (N \times 2D)}$  tel que  $\forall n \in \{1, \dots, N\}, i \in \{1, \dots, N \times 2D\}$ ,

$$W_{n,i} = \begin{cases} 1 & \text{ssi } 2nD \leq i < 2nD + 2D \\ 0 & \text{sinon} \end{cases}$$

alors,  $\forall x_1, \dots, x_N \in \mathbb{R}^D$

$$\left[ \mathbf{1} - \begin{pmatrix} \|x - x_1\|_1 \\ \|x - x_2\|_1 \\ \dots \\ \|x - x_N\|_1 \end{pmatrix} \right]_+ = \left[ \mathbf{1} - W \begin{bmatrix} Ix - x_1 \\ -Ix + x_1 \\ Ix - x_2 \\ -Ix + x_2 \\ \dots \\ Ix - x_N \\ -Ix + x_N \end{bmatrix} \right]_{++}$$

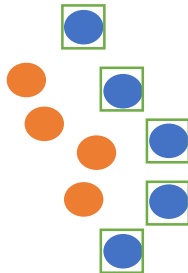
$\forall x_1, \dots, x_N \in \mathbb{Z}^D, y_1, \dots, y_N \in \{-1, 1\}$  la fonction

$$f(x) = \begin{pmatrix} y_1 \\ y_2 \\ \dots \\ y_N \end{pmatrix}^T \left[ \mathbf{1} - W \begin{bmatrix} Ix - x_1 \\ -Ix + x_1 \\ Ix - x_2 \\ -Ix + x_2 \\ \dots \\ Ix - x_N \\ -Ix + x_N \end{bmatrix} \right]_{++}$$

vérifie  $\forall n \in \{1, \dots, N\}, y_n f(x_n) > 0$ .

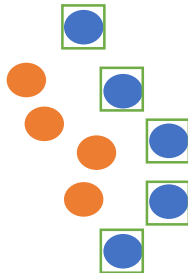
# Théorème d'universalité

Avec 100 neurones,  
on peut faire ça



# Théorème d'universalité

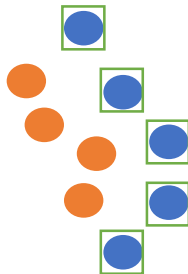
Avec 100 neurones,  
on peut faire ça



Sauf que c'est pas forcément bien  
comme avec le plus proche voisin

# Théorème d'universalité

Avec 100 neurones,  
on peut faire ça

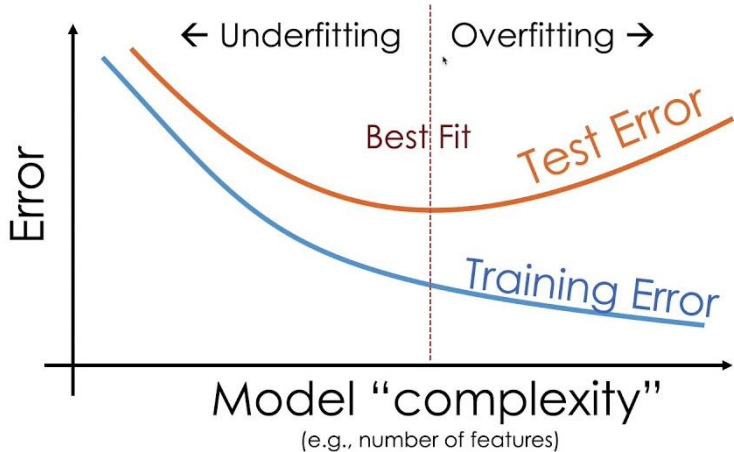


Sauf que c'est pas forcément bien  
comme avec le plus proche voisin

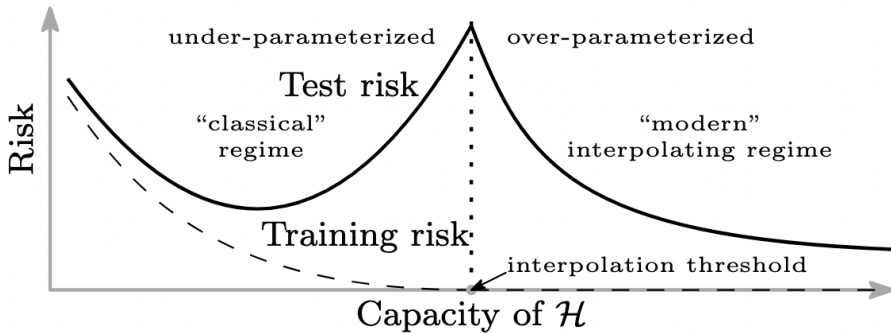


Il se trouve que ça arrive pas « trop » ???

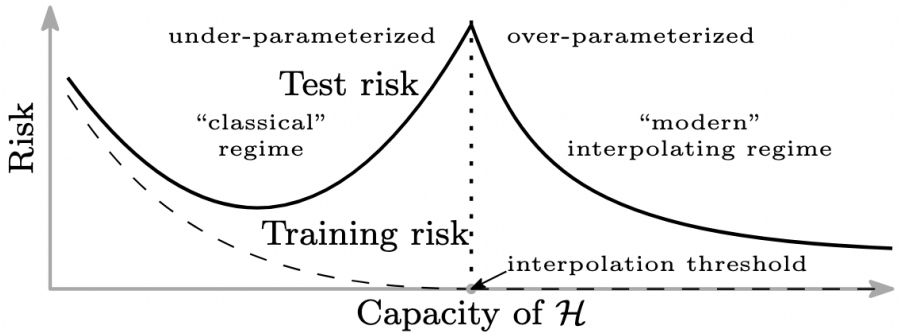
# Ancien paradigme



# Ce qu'on observe



# Ce qu'on observe



On pourrait penser que limiter le nombre de neurones  
est le moyen de limiter le sur-apprentissage.  
En pratique, c'est plus compliqué.

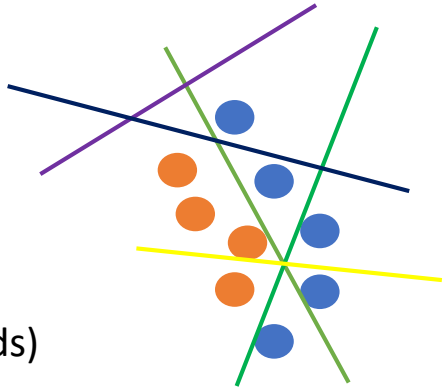


# PLAN

- COURS 1
  - MLP
  - Double descente ?
  - **Apprentissage**
- Cours 2
  - CNN
  - Segmentation/détection
  - Transformer et perspective

# Rappel

1 neurone  
peut coder TOUTES  
ces fonctions  
(via la valeur des poids)



Le SVM en est 1 seule fonction  
(optimisation des poids)

# Comment apprendre un MLP en pratique ?

- Descente de gradient
- Descente de gradient stochastique
- Fonction de perte
- Calcul du gradient

# Descente de gradient

$$J(u + \varepsilon) \approx J(u) + [\nabla_u J]^T \varepsilon + \varepsilon o(\varepsilon)$$

## Descente de gradient

$$J(u + \varepsilon) \approx J(u) + [\nabla_u J]^T \varepsilon + \varepsilon o(\varepsilon)$$

$$\begin{aligned} & J(u - \lambda \nabla_u J) \\ \approx & J(u) - \lambda [\nabla_u J]^T [\nabla_u J] + \lambda o(\lambda) \\ \approx & J(u) - \|\nabla_u J\|^2 \times \lambda + \lambda o(\lambda) \end{aligned}$$

## Descente de gradient

$$J(u + \varepsilon) \approx J(u) + [\nabla_u J]^T \varepsilon + \varepsilon o(\varepsilon)$$

$$\begin{aligned} & J(u - \lambda \nabla_u J) \\ \approx & J(u) - \lambda [\nabla_u J]^T [\nabla_u J] + \lambda o(\lambda) \\ \approx & J(u) - \|\nabla_u J\|^2 \times \lambda + \lambda o(\lambda) \end{aligned}$$

$$\begin{aligned} \nabla_u J \neq 0 &\implies \exists \lambda > 0, \\ J(u - \lambda \nabla_u J) &< J(u) \end{aligned}$$

# Descente de gradient

- Initialiser  $u$
- Si  $\|\nabla_u J\| < \delta$ , sortir
- Sinon
  - Initialiser  $\lambda$
  - Si  $J(u - \lambda \nabla_u J) < J(u)$ , alors  $u = u - \lambda \nabla_u J$
  - Sinon réessayer avec  $\lambda = \frac{\lambda}{2}$

# Descente de gradient

- Initialiser  $u$
- Si  $\|\nabla_u J\| < \delta$ , sortir
- Sinon
  - Initialiser  $\lambda$
  - Si  $J(u - \lambda \nabla_u J) < J(u)$ , alors  $u = u - \lambda \nabla_u J$
  - Sinon réessayer avec  $\lambda = \frac{\lambda}{2}$



Converge vers un point  $u$  tel que  $\|\nabla_u J\| < \delta$

Si la fonction est bornée




# Descente de gradient

$$\min_u \|Au - b\|^2 \quad \longrightarrow \quad u = A^{-1}b$$

1.  $u = 0$
2. Si  $\|2A^T(Au - b)\| < 10^{-6}$ , sortir
3.  $u = u - \lambda A^T(Au - b)$ , GOTO 2

# Descente de gradient

$$\min_u \|Au - b\|^2 \quad \begin{array}{l} \nearrow u = A^{-1}b \\ \quad O(N^\gamma L), \quad \gamma \approx 2.34 \end{array}$$

- 
1.  $u = 0$
  2. Si  $\|2A^T(Au - b)\| < 10^{-6}$ , sortir
  3.  $u = u - \lambda A^T(Au - b)$ , GOTO 2
- $O(N^2 2^L L)$

# Descente de gradient stochastique

$$J(u) = \sum_{k=1}^K j_k(u)$$

Si je tire  $k$  uniformément dans  $\{1, \dots, K\}$ , et que  $\mathcal{J} = j_k(u)$

$$\mathbf{E}[\nabla \mathcal{J}] = \frac{1}{K} \sum_{k=1}^K \nabla j_k = \nabla J$$

# Descente de gradient stochastique

$$J(u) = \sum_{k=1}^K j_k(u)$$

Si je tire  $k$  uniformément dans  $\{1, \dots, K\}$ , et que  $\mathcal{J} = j_k(u)$

$$\mathbf{E}[\nabla \mathcal{J}] = \frac{1}{K} \sum_{k=1}^K \nabla j_k = \nabla J$$

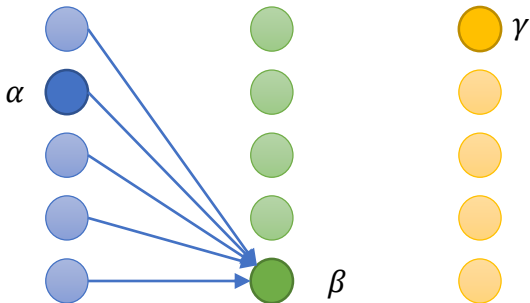


Si  $K$  est très très grand, il n'y a pas d'alternative à chercher une approximation probabiliste du gradient

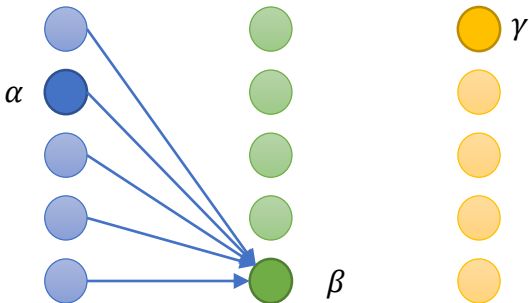
# Descente de gradient **stochastique** pour un MLP

- Choisir une fonction de perte  $J$
- Initialiser aléatoirement les poids  $w$
- Faire un certain nombre de fois
  - Sélectionner un paquet de données  $X$
  - Comparer  $f(X)$  et  $y(X)$  à travers  $J$
  - Calculer le gradient correspondant  $dw$
  - $w = w - lr * dw$

# Calcul du gradient ?



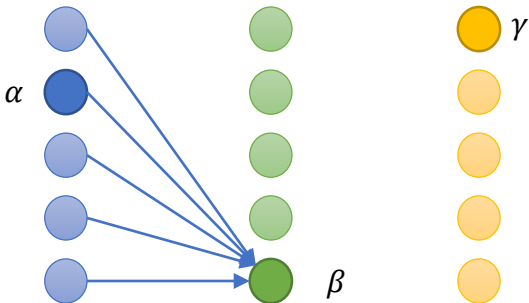
# Calcul du gradient ?



$$\beta = w_{\{\beta,\alpha\}}\alpha + \dots$$

$$\gamma = w_{\{\gamma,\beta\}}\beta + \dots$$

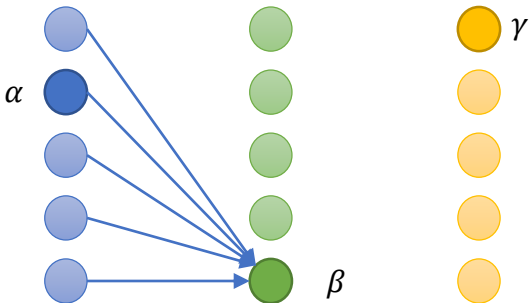
# Calcul du gradient ?



$$\begin{aligned}\beta &= w_{\{\beta,\alpha\}}\alpha + \dots & \gamma(\alpha, \dots) &= \gamma(\beta(\alpha, \dots), \dots) \\ \gamma &= w_{\{\gamma,\beta\}}\beta + \dots\end{aligned}$$



# Calcul du gradient ?

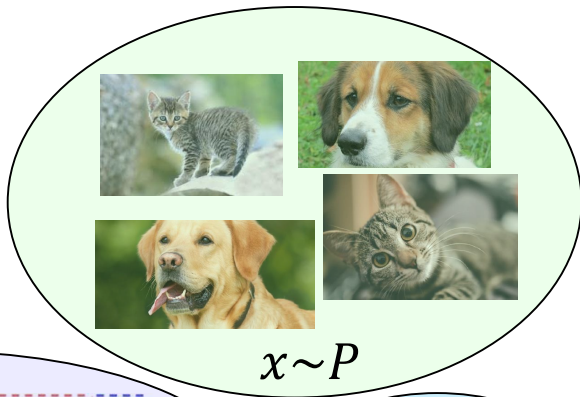
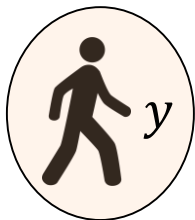


$$\begin{aligned}\beta &= w_{\{\beta,\alpha\}}\alpha + \dots & \gamma(\alpha, \dots) &= \gamma(\beta(\alpha, \dots), \dots) \\ \gamma &= w_{\{\gamma,\beta\}}\beta + \dots & \frac{\partial \gamma}{\partial \alpha} &= \frac{\partial \gamma}{\partial \beta} \frac{\partial \beta}{\partial \alpha} + \dots\end{aligned}$$

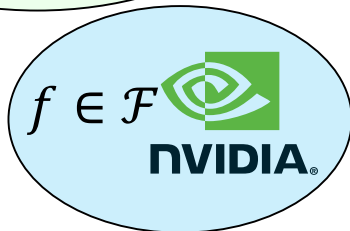
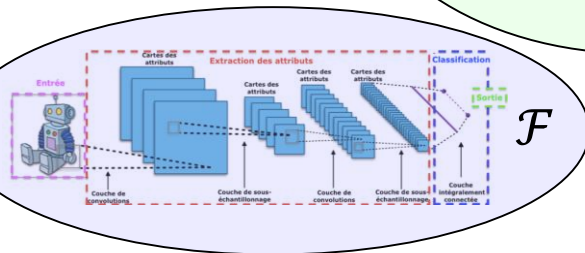
# Calcul du gradient

Et en pratique à la fin c'est pytorch qui le calcul pour vous !

# Rappel



$$x \sim P$$



# Pourquoi une fonction de perte ?

Ce qu'on veut c'est trouver  $w$  tel que

$$f_w(x) \approx y(x)$$

sur la base d'apprentissage.

# Pourquoi une fonction de perte ?

Ce qu'on veut c'est trouver  $w$  tel que

$$f_w(x) \approx y(x)$$

sur la base d'apprentissage.

➡ 1 erreur c'est  $\chi$  tel que  $y(\chi)f_w(\chi) \leq 0$

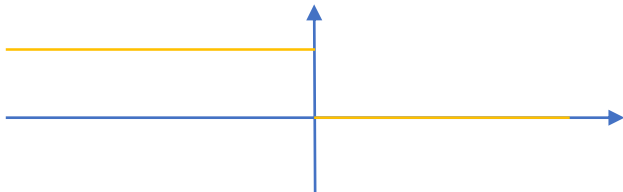
# Pourquoi une fonction de perte ?

Ce qu'on veut c'est trouver  $w$  tel que

$$f_w(x) \approx y(x)$$

sur la base d'apprentissage.

➡ 1 erreur c'est  $x$  tel que  $y(x)f_w(x) \leq 0$



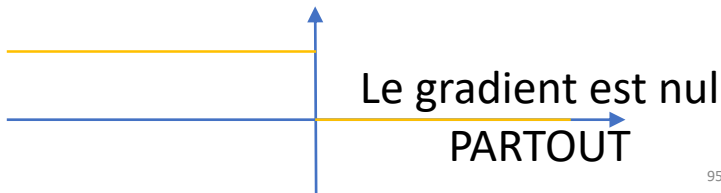
# Pourquoi une fonction de perte ?

Ce qu'on veut c'est trouver  $w$  tel que

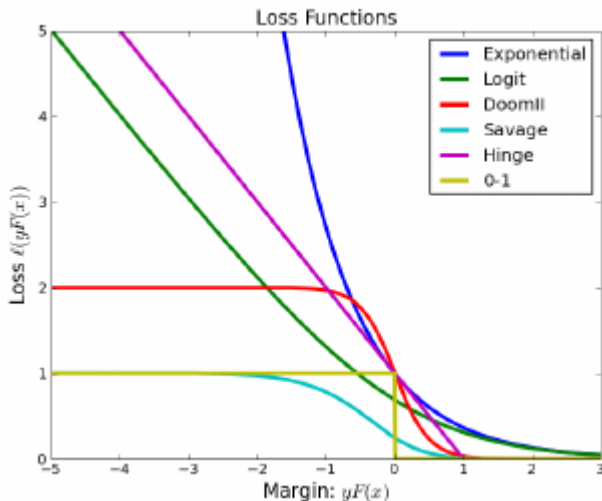
$$f_w(x) \approx y(x)$$

sur la base d'apprentissage.

➡ 1 erreur c'est  $\chi$  tel que  $y(\chi)f_w(\chi) \leq 0$

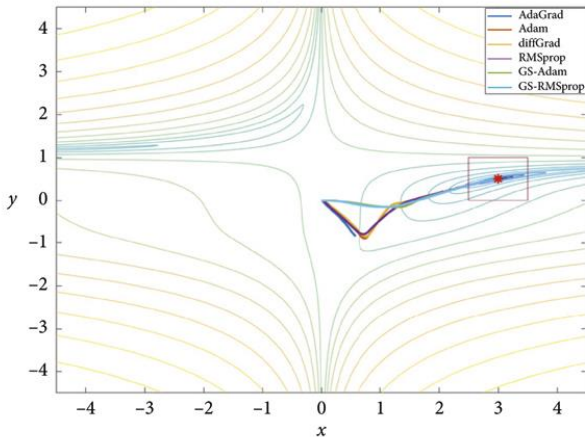


# fonction de perte





# fonction de perte



# Réseau de neurones

- Choisir une architecture
- Choisir une fonction de perte
- Choisir une variante de la SGD
- Effectuer une SGD sur les données d'apprentissage (en utilisant les gradients calculés par pytorch)



# PLAN

- COURS 1
  - MLP
  - Double descente ?
  - Apprentissage
- Cours 2
  - **CNN**
  - Segmentation/détection
  - Transformer et perspective

# Convolution



# Convolution

100	100	100	100	100
100	100	100	100	100
100	100	150	100	100
100	100	100	100	100
100	100	100	100	100

 $*$ 

0	-1	0
-1	5	-1
0	-1	0

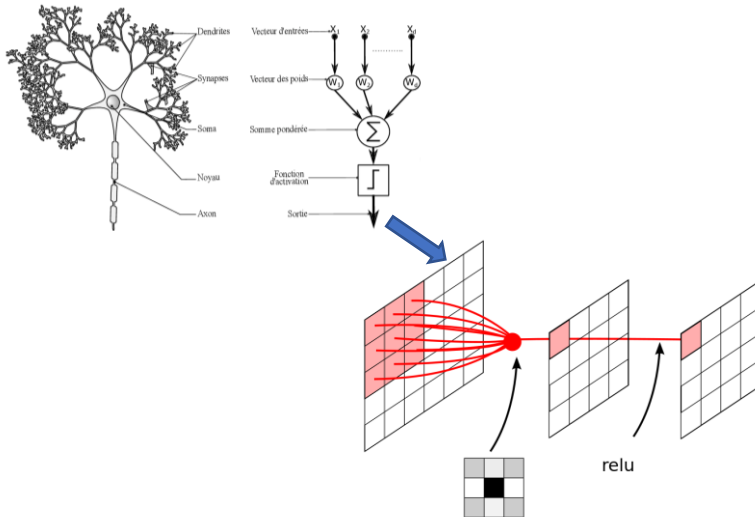
 $=$ 

100	100	100	100	100
100	100	50	100	100
100	50	350	50	100
100	100	50	100	100
100	100	100	100	100

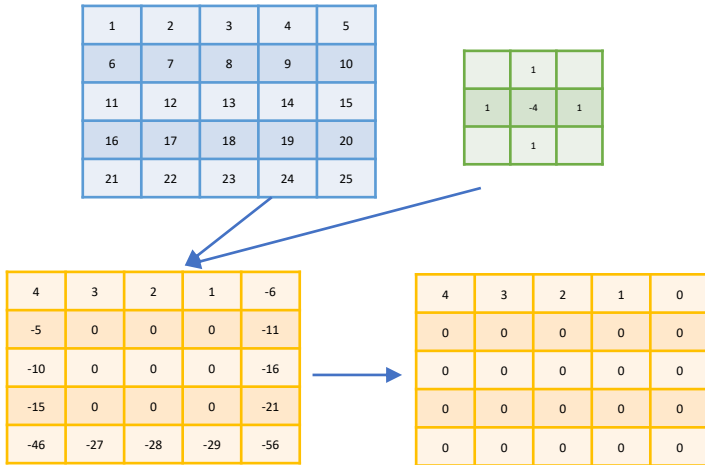
# Gabor

$$\mathbf{G}_x = \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix} * \mathbf{A} \quad \text{et} \quad \mathbf{G}_y = \begin{bmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ 1 & 2 & 1 \end{bmatrix} * \mathbf{A}$$

# ConvNet : convolution

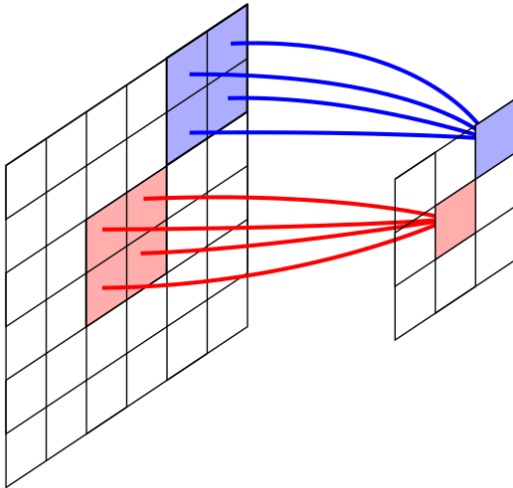


# ConvNet



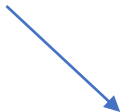


# ConvNet : pooling



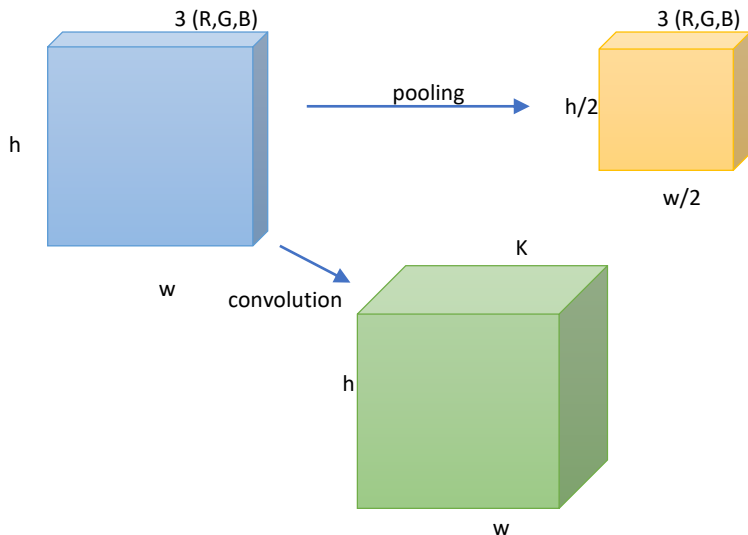
# ConvNet

1	2	3	4
6	7	8	9
11	12	13	14
16	17	18	19



7	9
17	19

# Attention : on oublie souvent la dimension spectrale



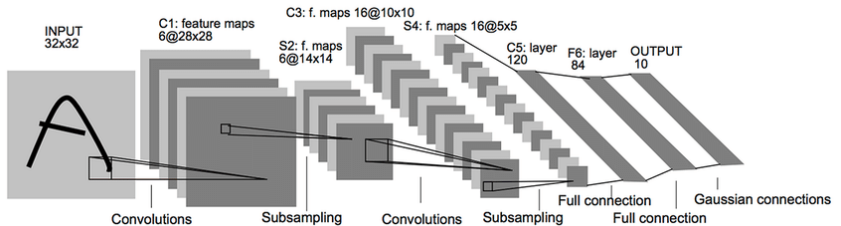
**Attention**, on manipule des paquets  
(donc 4D)

The diagram illustrates the components of the 4D tensor  $x \in \mathbb{R}^{B \times C \times H \times W}$ . Three blue arrows point from descriptive text to the dimensions in the tensor notation:

- An arrow from the word "paquet" points to the  $B$  dimension.
- An arrow from the text "dimension spectrale" points to the  $C$  dimension.
- An arrow from the text "dimensions spatiales" points to the  $H \times W$  spatial dimensions.

$x \in \mathbb{R}^{B \times C \times H \times W}$

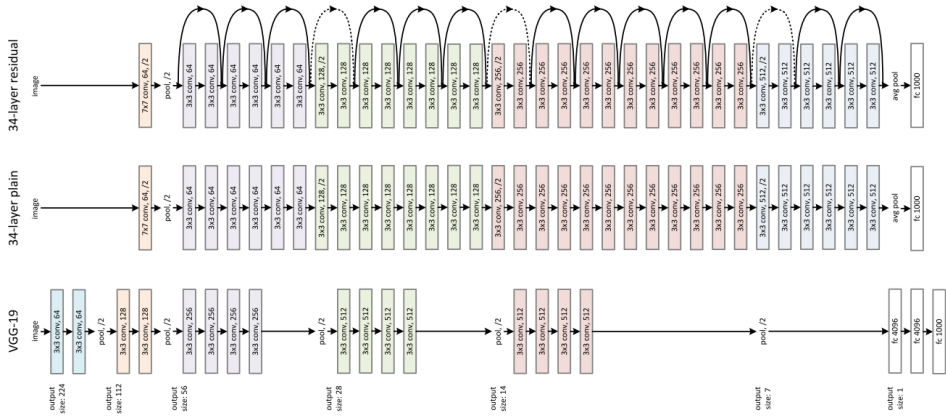
# Lenet



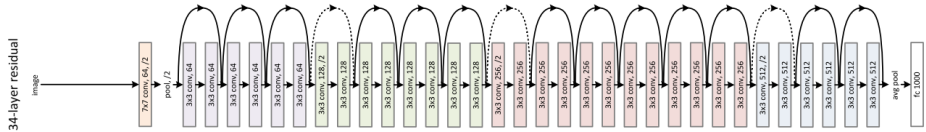
# Alexnet et VGG



# Resnet



# Resnet



$$BN(x) = \frac{x - (\sum_b x_b)}{\sqrt{\sum_b (x_b - (\sum_{b'} x_{b'}))^2 + 0.0000001}}$$

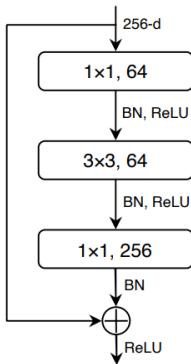


# EfficientNet, ConvNext

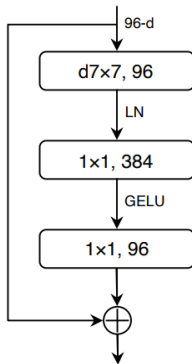
<https://pytorch.org/vision/stable/models.html>

# EfficientNet, ConvNeXt

**ResNet Block**



**ConvNeXt Block**



# EfficientNet, ConvNext

<https://pytorch.org/vision/stable/models.html>

# PLAN

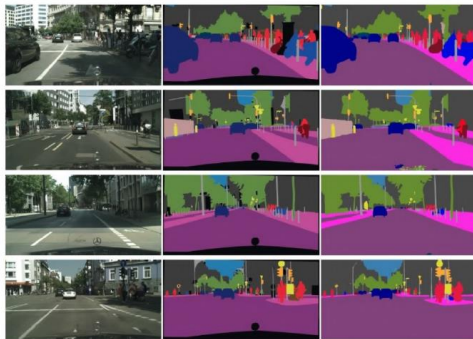
- COURS 1
  - MLP
  - Double descente ?
  - Apprentissage
- Cours 2
  - CNN
  - **Segmentation/détection**
  - Transformer et perspective

# Au-delà de la classification

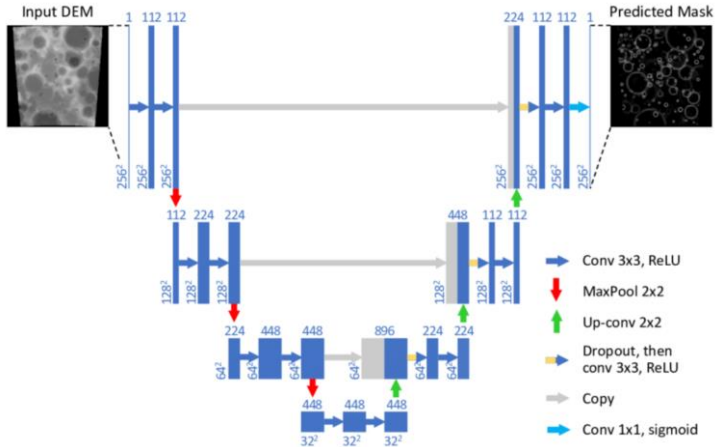
Les modèles par réseau de neurones sont flexibles.  
Ils peuvent généralement s'adapter au format du problème !

# Au-delà de la classification

Les modèles par réseau de neurones sont flexibles.  
Ils peuvent généralement s'adapter au format du problème !



# Segmentation



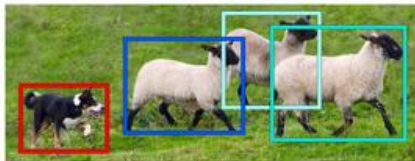
# Détection



Image Recognition



Semantic Segmentation



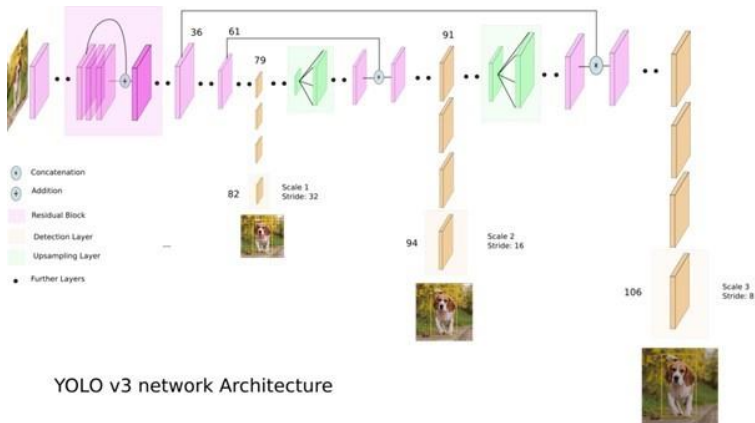
Object Detection



Instance Segmentation



# Détection



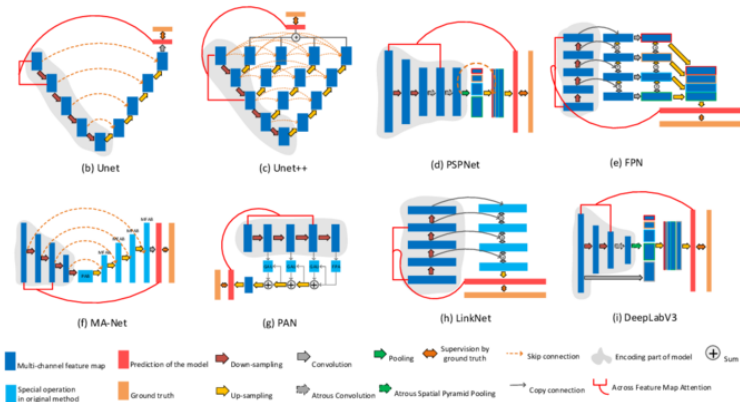
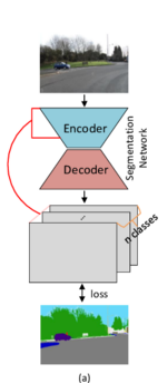
YOLO v3 network Architecture

# Segmentation et détection

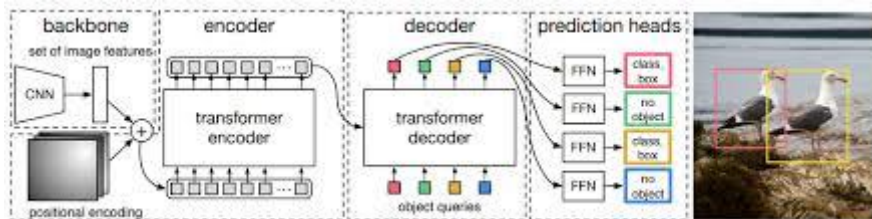
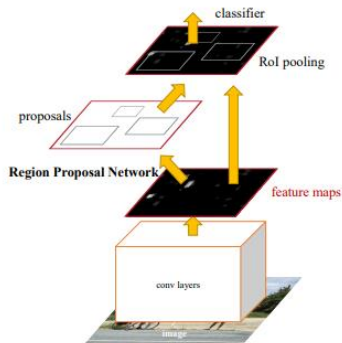
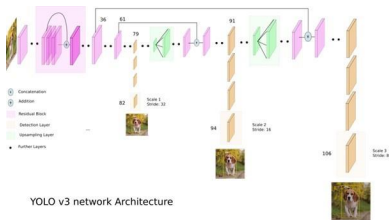
Le framework classique

- Un encodeur
  - qui part de l'image, et
  - produit des cartes de features multirésolutions
- Une tête qui
  - fusionne des cartes de features multirésolutions
  - produit la sortie attendue

# Segmentation



# Détection



# Segmentation et détection

Le framework classique

- Un encodeur
  - qui part de l'image, et
  - produit des cartes de features multirésolutions
- Une tête qui
  - fusionne des cartes de features multirésolutions
  - produit la sortie attendue



Grosse unification en cours  
via le mécanisme d'attention !

# PLAN

- COURS 1
  - MLP
  - Double descente ?
  - Apprentissage
- Cours 2
  - CNN
  - Segmentation/détection
  - **Transformer et perspective**

# Transformer

Si  $H = h_1, \dots, h_R$  sont  $R$  vecteurs de dimension  $D$ , la self attention s'exprime comme

$$S(H) = \text{softmax} \left( \frac{(QH)(KH)^T}{\sqrt{L}} \right) (VH)$$

avec  $Q, K \in \mathbb{R}^{L \times D}$

# Transformer

Si  $H = h_1, \dots, h_R$  sont  $R$  vecteurs de dimension  $D$ , la self attention s'exprime comme

$$S(H) = \text{softmax} \left( \frac{(QH)(KH)^T}{\sqrt{L}} \right) (VH)$$

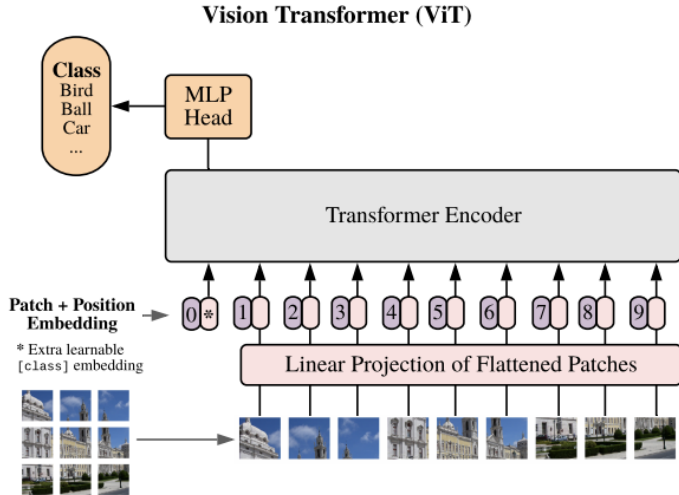
avec  $Q, K \in \mathbb{R}^{L \times D}$



Le nombre de poids ne dépend pas de  $R$  !

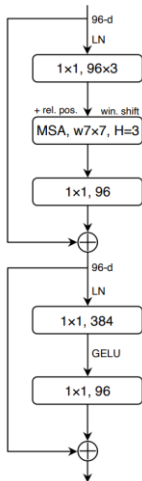


# Visual transformer

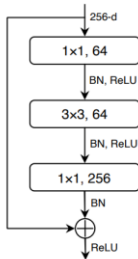


# ViT vs ConvNet

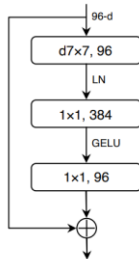
**Swin Transformer Block**



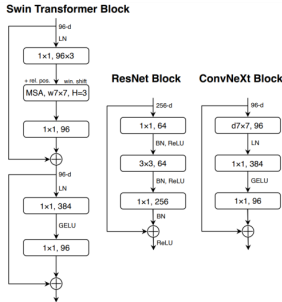
**ResNet Block**



**ConvNeXt Block**



# ViT vs ConvNet



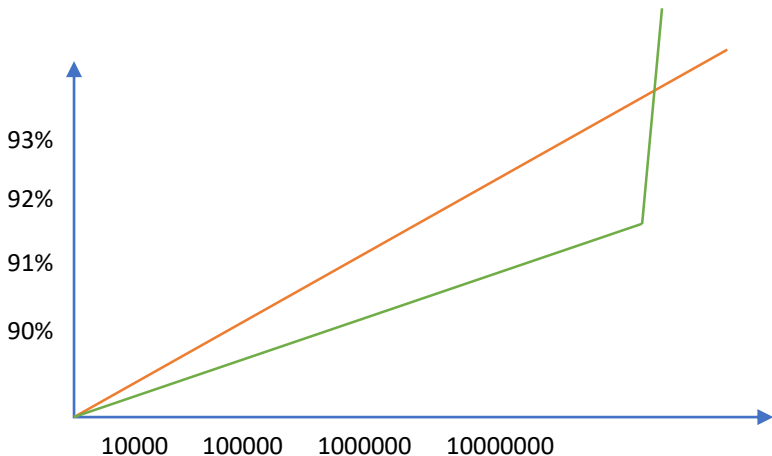
$$S(H) = \text{softmax} \left( \frac{(QH)(KH)^T}{\sqrt{L}} \right) (VH)$$

Augmenter brutalement la taille n'apporte pas tant que cela sur des architectures convolutives, alors que l'augmentation de performance semble constante avec des architectures Transformer !

# L'état de l'art de la vision par ordinateur

<https://pytorch.org/vision/stable/models.html>

# Performance vs tailles-données-modèles



??

# Fondation models

EfficientNet

100 Mo

SAM (Segment Everything by Meta)

2Go de poids

21Go pour faire un batch de 2 sur le petit modèle

LLAMA (un concurrent de chatGPT)

**128 Go** de poids !!!!!!!!!!!!!!!!!!!!!!!

8 x 32Go pour le faire tourner

200000€

# Fondation models



Sam Scarpino  
@svscarpino

...

"Most good ideas still come from academia." 100 100 100

Prof @ylecun on the role of academia in AI and the importance of good ideas (even if you don't have access to 50k gpus for compute). Fireside chat @Northeastern with @Experiential\_AI's @usamaf.

[Traduire le Tweet](#)



8:32 PM · 24 mai 2023 depuis Boston, MA · 25,5 k vues

# Fondation models

- Lenet AT&T
- Alexnet toronto / google
- VGG Oxford
- Resnet Microsoft
- Gpipe Google
- Efficientnet Google
- ViT Google
- SWIM Microsoft
- ConvNext Meta



Sam Scarpino  
@svscarpino

"Most good ideas still come from academia." 100 100 100

Prof @ylecun on the role of academia in AI and the importance of good ideas (even if you don't have access to 50k gpus for compute). Fireside chat @Northeastern with @Experiential\_AI's @usamaf.

[Traduire le Tweet](#)



8:32 PM · 24 mai 2023 depuis Boston, MA · 25,5 k vues



# Perspectives

Il existe de très bon modèles convolutifs (EfficientNet, ConvNext) mais leurs performances saturent

L'utilisation de couche Transformer  
est autant prometteuse que couteuse !

Le passage à l'échelle ne se fera pas dans les labos,  
car les ordres de grandeur mis en jeux ont explosé !

Merci pour votre attention