

Reconnaissance d'images avec des réseaux denses 2/2

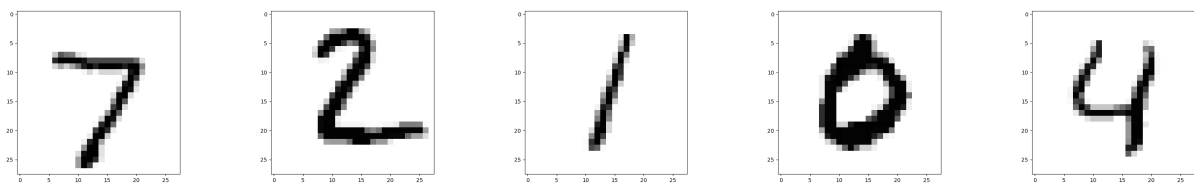
Objectifs ^a :

- Utilisation d'un réseau dense pour la reconnaissance d'image
- Comprendre les limitations de la méthode

a. Version 2025 adapté BUT 3 Informatique (Orléans), inspiré du livre d'Arnaud Bodin et François Recher

1. Données et composition du réseau

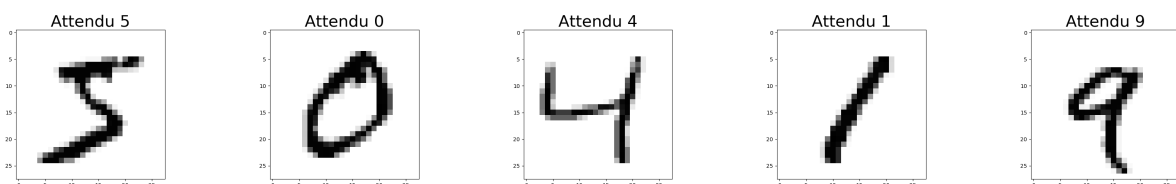
Il s'agit de reconnaître de façon automatique des chiffres écrits à la main.



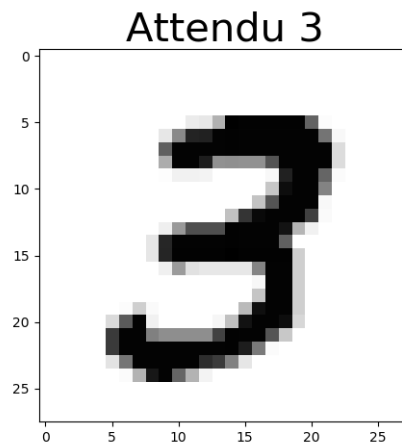
C'est l'un des succès historiques des réseaux de neurones qui permet par exemple le tri automatique du courrier par lecture du code postal.

1.1. Base MNIST

Un élément essentiel de l'apprentissage automatique est de disposer de données d'apprentissage nombreuses et de qualité. Une telle base est la base MNIST dont voici les premières images.



Une donnée est constituée d'une image et du chiffre attendu.



Plus en détails :

- La base est formée de 60 000 données d'apprentissage et de 10 000 données de test.
- Chaque donnée est de la forme : [une image, le chiffre attendu].
- Chaque image est de taille 28×28 pixels, chaque pixel contenant un des 256 niveaux de gris (numérotés de 0 à 255).

Ces données sont accessibles très simplement avec Keras :

```
from keras.datasets import mnist
(X_train_data, Y_train_data), (X_test_data, Y_test_data) = mnist.load_data()
```

Il faut passer un peu de temps à comprendre et à manipuler les données. `X_train_data[i]` correspond à une image et `Y_train_data[i]` au chiffre attendu pour cette image. Nous parlerons plus tard des données de test. On renvoie au fichier `tf2_chiffres_data.py` pour une exploration de la base. Noter que l'on peut afficher une image à l'aide de *matplotlib* par la commande :

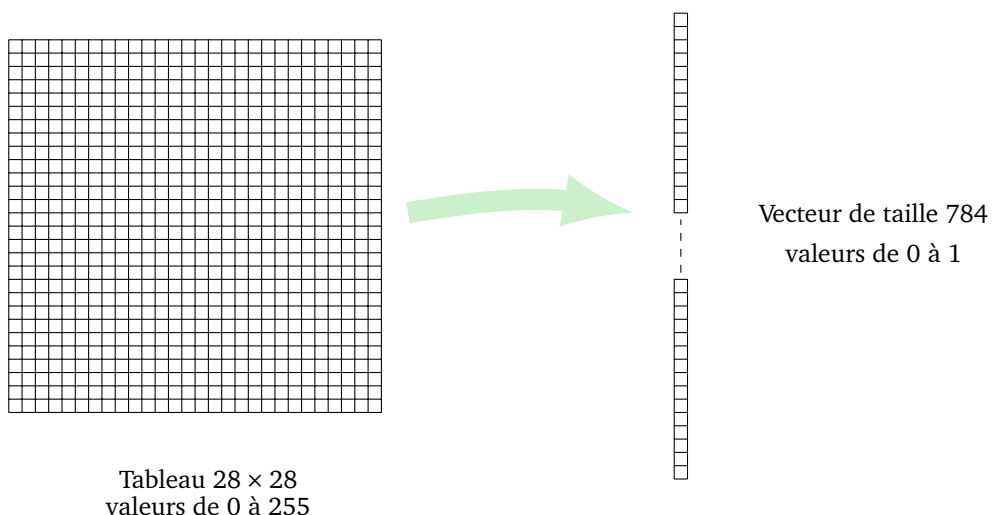
```
plt.imshow(X_train_data[i], cmap='Greys')
```

suivie de `plt.show()`.

1.2. Traitement des données

Pour une utilisation par un réseau de neurones nous devons d'abord transformer les données.

Donnée d'entrée. En entrée du réseau de neurones, nous devons avoir un vecteur. Au départ chaque image est un tableau de taille 28×28 ayant des entrées entre 0 et 255. Nous la transformons en un vecteur de taille $784 = 28^2$ et nous normalisons les données dans l'intervalle $[0, 1]$ (en divisant par 255).



Ainsi, une entrée X est un « vecteur-image », c'est-à-dire un vecteur de taille 784 représentant une image.

Donnée de sortie. Notre réseau de neurones ne va pas renvoyer le chiffre attendu, mais une liste de 10 probabilités. Ainsi chaque chiffre doit être codé par une liste de 0 et de 1.

- 0 est codé par (1, 0, 0, 0, 0, 0, 0, 0, 0, 0),
- 1 est codé par (0, 1, 0, 0, 0, 0, 0, 0, 0, 0),
- 2 est codé par (0, 0, 1, 0, 0, 0, 0, 0, 0, 0),
- ...
- 9 est codé par (0, 0, 0, 0, 0, 0, 0, 0, 0, 1).

Fonction. Nous cherchons une fonction $F : \mathbb{R}^{784} \rightarrow \mathbb{R}^{10}$, qui à un vecteur-image associe une liste de probabilités, telle que $F(X_i) \simeq Y_i$ pour nos données transformées (X_i, Y_i) , $i = 1, \dots, 60\,000$.

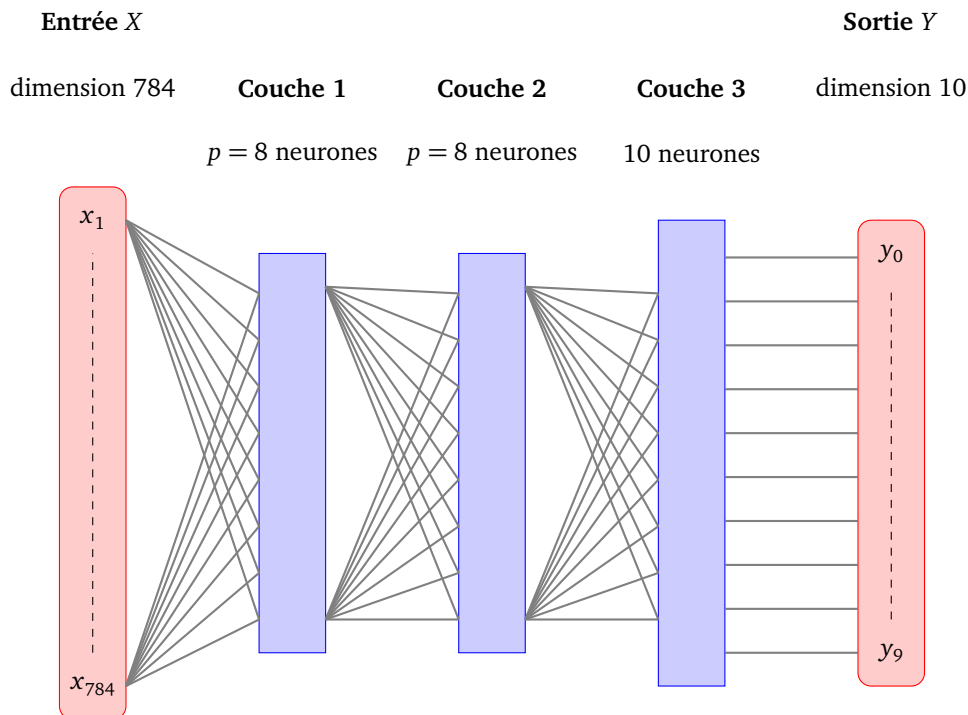
Par exemple la fonction F , évaluée sur un vecteur-image X peut renvoyer

$$F(X) = (0.01, 0.04, 0.03, 0.01, 0.02, 0.22, 0.61, 0.02, 0.01, 0.01).$$

Dans ce cas, le nombre le plus élevé est 0.61 au rang 6, cela signifie que notre fonction F prédit le chiffre 6 avec une probabilité de 61%, mais cela pourrait aussi être le chiffre 5 qui est prédit à 22%. Les autres chiffres sont peu probables.

1.3. Réseau

Nous allons construire un réseau de neurones qui produira une fonction $F : \mathbb{R}^{784} \rightarrow \mathbb{R}^{10}$. L'architecture est composée de 3 couches. En entrée nous avons un vecteur de taille 784. La première et la seconde couches sont composées chacune de $p = 8$ neurones. La couche de sortie est formée de 10 neurones, un pour chacun des chiffres. Pour les fonctions d'activation nous utiliserons, pour les deux premières couches, la fonction sigmoïd (σ) et pour la dernière la fonction softmax (utilisée dans un but de classification) que nous définissons dans le paragraphe suivant. Aussi nous introduirons la fonction d'erreur la plus utilisée pour la reconnaissance d'image : l'entropie croisée catégorielle.



2. Code du réseau - entraînement et test

2.1. Le code

Voici le code complet du programme qui sera commenté plus loin.

```
import numpy as np
import os
os.environ["KERAS_BACKEND"] = "torch"
import keras
import os
os.environ["KERAS_BACKEND"] = "torch"
from keras import optimizers
from Keras.keras.models import Sequential
from Keras.keras.layers import Dense

### Partie A - Les données

from Keras.keras.datasets import mnist
from Keras.keras.utils import to_categorical

# Téléchargement des données
(X_train_data, Y_train_data), (X_test_data, Y_test_data) = mnist.load_data()

N = X_train_data.shape[0] # N = 60 000 données

# Données d'apprentissage X
X_train = np.reshape(X_train_data, (N, 784)) # vecteur image
X_train = X_train/255 # normalisation

# Données d'apprentissage Y vers une liste de taille 10
Y_train = to_categorical(Y_train_data, num_classes=10)

# Données de test
X_test = np.reshape(X_test_data, (X_test_data.shape[0], 784))
X_test = X_test/255
Y_test = to_categorical(Y_test_data, num_classes=10)

### Partie B - Le réseau de neurones

p = 8
modele = Sequential()

# Première couche : p neurones (entrée de dimension 784 = 28x28)
modele.add(Dense(p, input_dim=784, activation='sigmoid'))

# Deuxième couche : p neurones
modele.add(Dense(p, activation='sigmoid'))
```

```
# Couche de sortie : 10 neurones (un par chiffre)
modele.add(Dense(10, activation='softmax'))

# La fonction d'activation 'softmax' sera décrite dans le paragraphe suivant.

# Choix de la méthode de descente de gradient
modele.compile(loss='categorical_crossentropy',
               optimizer='sgd',
               metrics=['accuracy'])

# La fonction d'erreur 'categorical_crossentropy' est décrite dans le paragraphe suivant.
# L'optimisation 'sgd' est décrite dans le paragraphe suivant.
# 'accuracy' est décrite dans le paragraphe suivant.

print(modele.summary())
```

2.2. Commentaires sur le code

Partie A - Les données

Les données d'apprentissage sont téléchargées facilement par une seule instruction. Elles regroupent $N = 60\,000$ données d'apprentissage (*train*) et $10\,000$ données de test. Les données sont de la forme (X_i, Y_i) où X_i est une image (un tableau 28×28 d'entiers de 0 à 255) et Y_i est le chiffre correspondant (de 0 à 9). Les données X_i sont transformées en un vecteur de taille 784 (fonction `reshape()`) et ses coefficients sont ramenés dans l'intervalle $[0, 1]$ par division par 255. Ainsi `X_train` est maintenant une liste *numpy* de $60\,000$ vecteurs de taille 784.

Chaque donnée Y_i est transformée en une liste de longueur 10 du type $[0, 0, \dots, 0, 1, 0, \dots, 0]$ avec le 1 à la place du chiffre attendu. On utilise ici la fonction `to_categorical()`.

Partie B - Le réseau de neurones

Notre réseau est composé de 3 couches. La première couche contient $p = 8$ neurones et reçoit en entrée 784 valeurs (une pour chaque pixel de l'image). La seconde couche contient aussi $p = 8$ neurones. La troisième couche contient 10 neurones (le premier pour détecter le chiffre 0, le deuxième pour le chiffre 1, ...). Pour les deux premières couches la fonction d'activation est la fonction σ . Pour la couche de sortie, la fonction d'activation est la fonction *softmax* qui est adaptée au problème. Ce qui fait que la couche de sortie renvoie une liste de 10 nombres dont la somme est 1 et qui correspond à une liste de probabilités.

La fonction d'erreur (*loss*) adaptée au problème s'appelle *categorical_crossentropy*. La méthode de minimisation de l'erreur choisie est la descente de gradient stochastique.

La valeur *accuracy* du paramètre *metrics* indique que l'on souhaite en plus mesurer la précision des résultats (cela ne change rien pour la descente de gradient qui dépend uniquement de la fonction d'erreur et pas de cette précision).

2.3. Dernières remarques

On lance le calcul des poids par la fonction `fit()` :

```
modele.fit(X_train, Y_train, batch_size=32, epochs=100)
```

- Les calculs sont effectués selon la méthode de descente de gradient choisie auparavant.
- Les données d'apprentissage utilisées sont `X_train` (valeurs de départ) et `Y_train` (valeurs d'arrivée attendues).

- L'option `batch_size` précise la taille de l'échantillon :
 - une descente de gradient stochastique pure : `batch_size=1`,
 - une descente de gradient sur la totalité des N données : `batch_size=N`,
 - ou toute valeur intermédiaire : par défaut `batch_size=32`.
- L'option `epochs` détermine le nombre d'étapes dans la méthode de gradient. Si K est la taille d'un lot (valeur de `batch_size`) et N la taille des données alors le nombre d'étapes par époque est N/K .
- Il existe une option `verbose` qui permet d'afficher plus ou moins de détails à chaque époque (0 : rien, 1 : barre de progression, 2 : numéro de l'époque).

On peut vouloir mesurer d'autres choses que la fonction d'erreur (moindre carré, entropie catégorielle). Par exemple l'option de la fonction `compile()` :

```
metrics=['accuracy']
```

va en plus mémoriser la précision du modèle (sous la forme d'un pourcentage de réussite). Par exemple, si on doit classer des images en deux catégories (chat/0 et chien/1) alors la fonction d'erreur est un outil indispensable pour notre problème, mais le résultat est mesuré concrètement par le pourcentage d'images correctement identifiées.

Pour la calculer on compare les prédictions du modèle aux données réelles (étiquettes) dans l'ensemble de données de test. Le nombre de bonnes prédictions (les prédictions correctes) est représenté par M dans la formule de l'accuracy, et N est le nombre total d'exemples dans l'ensemble de données de test. On a alors : $\text{accuracy} = M/N$.

Note : pourquoi ne pas prendre la précision comme fonction d'erreur puisque c'est ce qui nous intéresse ? Tout simplement parce que ce n'est pas une fonction différentiable, il n'y a donc pas de méthode de gradient pour la minimiser.

3. Résultats

On lance maintenant la procédure d'apprentissage, puis on analyse les résultats.

```
### Partie C - Calcul des poids par descente de gradient
```

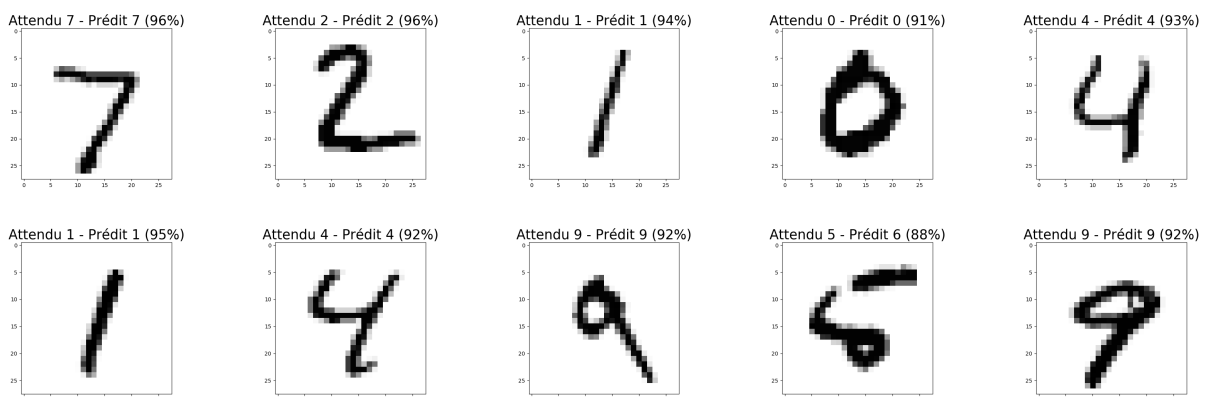
```
modele.fit(X_train, Y_train, batch_size=32, epochs=40)
```

```
### Partie D - Résultats
```

```
resultat = modele.evaluate(X_test, Y_test, verbose=0)
print('Valeur de l'erreur sur les données de test (loss):', resultat[0])
print('Précision sur les données de test (accuracy):', resultat[1])
```

Remarque : `epochs=40` signifie que l'on se sert 40 fois des mêmes données.

On estime la performance de notre réseau sur les données de test. Voici les premiers résultats.



Le réseau prédit correctement 9 valeurs sur 10. En fait, la fonction associée au réseau renvoie une liste de probabilités. Le chiffre prédit est celui qui a la plus forte probabilité. Par exemple pour la première image (en haut à gauche) la valeur renvoyée est :

$$Y_0 = (0.001, 0.000, 0.000, 0.008, 0.002, 0.005, 0.000, 0.965, 0.000, 0.020).$$

On en déduit la prédiction du chiffre 7 avec une forte probabilité de 96%.

L'avant-dernière image conduit à une mauvaise prédiction : la fonction prédit le chiffre 6 alors que le résultat attendu est le chiffre 5.

Les calculs de la descente de gradient sont faits avec une fonction d'erreur qu'il s'agit de minimiser. Par contre cette fonction d'erreur n'est pas pertinente pour évaluer la qualité de la modélisation. On préfère ici calculer la **précision** (*accuracy*) qui correspond à la proportion de chiffres détectés correctement.

Exercice 1.

Calculer l'entropie croisée catégorielle pour l'exemple ci-dessus.

Voici quelques résultats pour différentes tailles du réseau (couche 1 avec p neurones, couche 2 avec aussi p neurones, couche 3 avec 10 neurones) :

p	neurones	poids	précision
8	26	6442	90.0%
10	30	8070	91.4%
20	50	16 330	93.4%
50	110	42 310	94.4%

Sans trop d'efforts on obtient donc une précision de 95%. C'est-à-dire que 95 fois sur 100 le chiffre prédit est le chiffre correct.

Il est important de mesurer la précision sur les données de test qui sont des données qui n'ont pas été utilisées lors de l'apprentissage. Le réseau n'a donc pas appris par cœur les données d'apprentissage, mais a réussi à dégager un schéma, validé sur des données indépendantes.

3.1. Détails sur le code

Reprenons pas à pas le programme et donnons quelques explications.

Partie C - Calcul des poids par descente de gradient

La fonction `fit()` lance la descente de gradient (avec une initialisation aléatoire des poids). L'option `batch_size` détermine la taille du lot (*batch*). On indique aussi le nombre d'époques à effectuer (avec `epochs=40`, chacune des $N = 60\,000$ données sera utilisée 40 fois, mais comme chaque étape regroupe un lot de 32 données, il y a $40 \times N/32$ étapes de descente de gradient).

Partie D - Résultats

On insiste sur le fait que la performance du réseau avec les poids calculés doit être mesurée sur les données de test et non sur les données d'apprentissage.

La fonction `evaluate()` renvoie la valeur de la fonction d'erreur (qui est la fonction minimisée par la descente de gradient mais qui n'a pas de signification tangible). Ici la fonction renvoie aussi la précision (car on l'avait demandée en option dans la fonction `compile()`).

Un peu plus de résultats

Pour l'instant, nous avons juste mesuré l'efficacité globale du réseau. Il est intéressant de vérifier à la main les résultats. Les instructions ci-dessous calculent les prédictions pour toutes les données de test (première ligne). Ensuite, pour une donnée particulière, on compare le chiffre prédit et le chiffre attendu.

```
# Prédiction sur les données de test
Y_predict = modele.predict(X_test)

# Un exemple
i = 8 # numéro de l'image

chiffre_predit = np.argmax(Y_predict[i]) # prédiction par le réseau

print("Sortie réseau", Y_predict[i])
print("Chiffre attendu :", Y_test_data[i])
print("Chiffre prédit :", chiffre_predit)

plt.imshow(X_test_data[i], cmap='Greys')
plt.show()
```

On rappelle que `Y_predict[i]` est une liste de 10 nombres correspondant à la probabilité de chaque chiffre. Le chiffre prédit s'obtient en prenant le rang de la probabilité maximale, c'est exactement ce que fait la fonction `argmax` de *numpy*.

3.2. Un exercice pour pratiquer

1. Expérimentez le réseau précédent en modifiant son architecture :

- Modifiez le nombre de couches et de neurones.
- Explorez différentes fonctions d'activation.

2. Réalisez une analyse d'erreur :

- Identifiez des exemples mal classés par le modèle :
Que font les lignes de code suivantes ?

```
predictions = model.predict(X_test)
predicted_labels = np.argmax(predictions, axis=1)
true_labels = np.argmax(Y_test, axis=1)
```
- On utilise ensuite la méthode `where` de *numpy* comme ci-dessous :

```
misclassified_indices = np.where(predicted_labels != true_labels)[0]
```

 Que retrouve-t-on dans `misclassified_indices` ?
- Le code ci-dessous affiche les 5 premiers exemples mal classés

```
for idx in misclassified_indices[:5]:
    plt.imshow(X_test[idx].reshape(28, 28), cmap='gray')
    plt.title(f'Predicted: {predicted_labels[idx]}, True: {true_labels[idx]}')
    plt.show()
```

Auriez-vous aussi fait une erreur ?

Modifiez à nouveau la structure du réseau afin de voir si l'on peut améliorer la prédiction de ces 5 images.

4. Limites de la reconnaissance d'images avec des réseaux denses

On termine par un constat d'échec (provisoire), nos réseaux de neurones atteignent leurs limites lorsque le problème posé se complique. Nous souhaitons reconnaître des petites images et les classer selon 10 catégories.

4.1. Données

La base CIFAR-10 contient 60 000 petites images de 10 types différents.



Plus en détails :

- Il y a 50 000 images pour l'apprentissage et 10 000 pour les tests.
- Chaque image est de taille 32×32 pixels en couleur. Un pixel couleur est codé par trois entiers (r, g, b) compris entre 0 et 255. Une image est donc composée de $32 \times 32 \times 3$ nombres.
- Chaque image appartient à une des dix catégories suivantes : avion, auto, oiseau, chat, biche, chien, grenouille, cheval, bateau et camion.

4.2. Programme

```
# Partie A. Données
```

```
from Keras.keras.datasets import cifar10
```

```
(X_train_data,Y_train_data),(X_test_data,Y_test_data)=cifar10.load_data()
```

```
num_classes = 10
```

```
labels = ['airplane','automobile','bird','cat','deer',  
          'dog','frog','horse','ship','truck']
```

```
Y_train = keras.utils.to_categorical(Y_train_data, num_classes)
```

```
X_train = X_train_data.reshape(50000,32*32*3)
```

```
X_train = X_train.astype('float32')
```

```

X_train = X_train/255

# Partie B. Réseau

modele = Sequential()

p = 30
modele.add(Dense(p, input_dim=32*32*3, activation='sigmoid'))
modele.add(Dense(p, activation='sigmoid'))
modele.add(Dense(p, activation='sigmoid'))
modele.add(Dense(p, activation='sigmoid'))
modele.add(Dense(10, activation='softmax'))

modele.compile(loss='categorical_crossentropy',
               optimizer='adam', metrics=['accuracy'])

modele.summary()

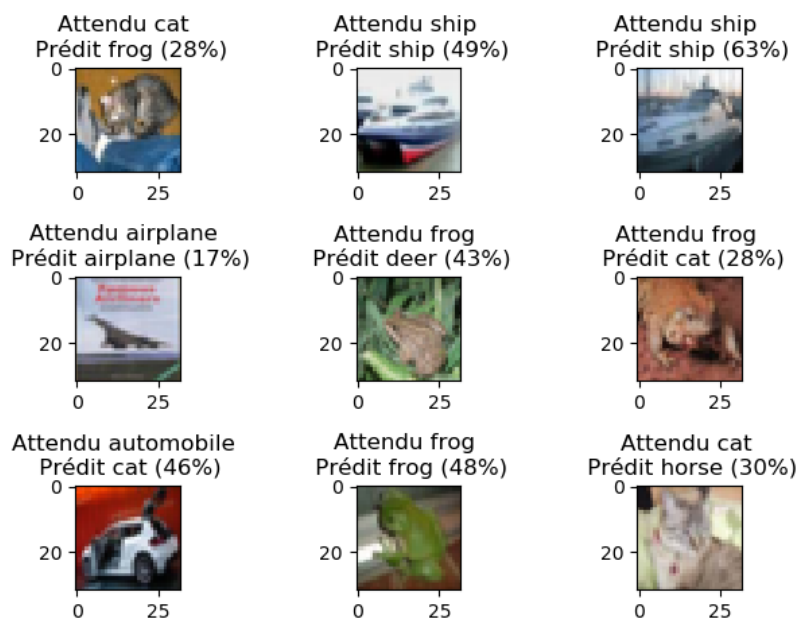
# Partie C. Apprentissage

modele.fit(X_train, Y_train, epochs=10, batch_size=32)

```

4.3. Résultats

Dans le programme ci-dessus avec $p = 30$, il y a 100 000 poids à calculer. Avec beaucoup de calculs (par exemple avec 50 époques et la descente « adam »), on obtient moins de 50% de précision. Ce qui fait que le sujet d'une image sur deux est mal prédit.



Il faut donc une nouvelle approche pour reconnaître correctement ces images !

Exercice

Indiquer si les affirmations suivantes sont **vraies** ou **fausses**.

1. Après l'instruction `X_train = X_train/255`, toutes les valeurs de `X_train` sont dans l'intervalle $[0, 1]$.
2. L'instruction `Y_train = to_categorical(Y_train_data, num_classes=10)` convertit chaque chiffre en un vecteur de 10 probabilités dont la somme vaut 1.
3. L'option `activation='softmax'` dans la couche de sortie garantit que la somme des 10 valeurs de sortie vaut toujours 1.
4. L'option `metrics=['accuracy']` modifie le calcul de la descente de gradient.
5. La fonction `model.fit(X_train, Y_train, batch_size=32, epochs=40)` effectue une mise à jour des poids après chaque lot de 32 images.
6. Lorsque `batch_size = 1`, la fonction `fit()` effectue une descente de gradient stochastique pure.
7. L'appel `modele.summary()` permet d'afficher le nombre total de poids du réseau.
8. L'appel `modele.predict(X_test)` renvoie, pour chaque image, une liste de probabilités correspondant aux 10 classes possibles.
9. L'instruction `np.argmax(Y_predict[i])` calcule l'erreur d'entropie croisée pour la donnée numéro i .
10. En utilisant le même réseau dense ($784 \rightarrow 8 \rightarrow 8 \rightarrow 10$), il est possible d'obtenir une précision de 95% sur CIFAR-10 en utilisant suffisamment d'époques.
11. La fonction d'erreur `categorical_crossentropy` compare la distribution de sortie du réseau q avec le vecteur one-hot p représentant la vraie classe.
12. Augmenter le nombre de neurones (p) dans les couches cachées augmente toujours la précision du modèle, quelle que soit la base d'images utilisée.