

# Classification non supervisée (clustering)

## La librairie pytorch

*PyTorch* est un framework d'apprentissage profond (deep learning) open source principalement développé par le groupe de recherche en IA de Facebook. Pour installer *PyTorch* dans votre environnement conda, rendez-vous sur <https://pytorch.org> et sélectionnez votre système d'exploitation, conda, Python, CPU (ceux qui possèdent une GPU peuvent installer la version CUDA mais ça peut déclencher plein de bugs à gérer pour installer la bonne version de cuda). Copiez le texte de la boîte "Run this command :". Ouvrez maintenant un terminal. Collez le texte et exécutez. Une fois l'installation terminée, vous devez redémarrer Jupyter.

## Jeu de données MNIST

Pour ce travail, nous allons faire appel au jeu de données MNIST. C'est un jeu de données d'images de chiffres manuscrits de 0 à 9. Chaque image est au format  $28 \times 28$  pixels qui prennent chacun 256 valeurs dans une gamme allant du blanc ( $= 0$ ) au noir ( $= 1$ ). Les étiquettes appartenant aux images sont également incluses. Heureusement, *PyTorch* est équipé d'un chargeur de données MNIST. La première fois que vous exécutez le code de la cellule ci-dessous, elle télécharge automatiquement le jeu de données MNIST. Cela peut prendre quelques minutes. Les principaux types de données dans *PyTorch* sont des tenseurs. Dans un premier temps, nous allons convertir ces tenseurs en tableaux *numpy*. Ensuite, nous allons utiliser le module *torch* pour travailler directement avec les tenseurs *PyTorch*.

```
import torch
from torchvision import datasets, transforms

train_dataset = datasets.MNIST('../data', train=True, download=True,
                             transform=transforms.Compose([
                                 transforms.ToTensor(),
                                 transforms.Normalize((0.1307,), (0.3081,)))
                             ]))

train_labels = train_dataset.train_labels.numpy()
train_data = train_dataset.train_data.numpy()
# Nous allons implémenter l'algorithme EM (on a besoin de vecteur)
train_data = train_data.reshape(train_data.shape[0], -1)
```

## Algorithme EM

Nous allons implémenter l'algorithme EM (Expectation Maximisation) pour la reconnaissance des chiffres manuscrits dans le jeu de données MNIST. La loi de probabilité des pixels d'une image est un modèle de mélange de Bernoulli (Voir le livre de Bishop partie §9.3.3):

$$p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\pi}) = \sum_{k=1}^K \pi_k \prod_{j=1}^D \mu_{kj}^{x_j} (1 - \mu_{kj})^{(1-x_j)}$$

où  $x_j$  est la valeur du  $j$  ème pixel d'une image, le paramètre  $\mu_{kj}$  représente la probabilité que le  $j$  ème pixel de la classe  $k$  soit noir, et  $\{\pi_1, \dots, \pi_K\}$  sont les proportions du mélange. On souhaite utiliser ce jeu de données pour construire une méthode de classification de nouvelles images de chiffres manuscrits.

### Binarisation des données

La transformation des tenseurs *PyTorch* en vecteurs *numpy* a permis d'aplatir les images en vecteurs. Comme nous allons utiliser un modèle de mélange de Bernoulli, nous avons besoin de binariser les pixels d'un vecteur.

#### Question 1

- Écrire une fonction `binarize` pour convertir chaque vecteur de pixels représentant une image MNIST en vecteur de pixels binaires  $x_i \in \{0, 1\}$ , en utilisant un seuil (à choisir).
- Écrire une fonction qui prend en entrées un chiffre allant de 0 à 9 et un nombre de figures correspondantes à ce chiffre à échantillonner et à afficher à côté de la version binarisée de chaque figure. Échantillonner quelques images des chiffres 2, 5 et 9 pour tester la fonction. Mettre comme titre de l'image, son étiquette.

### Implémentation

#### Question 2

Écrire une fonction `EM(X, K, max_iter)` qui implémente l'algorithme EM pour estimer les paramètres d'un modèle de mélange de Bernoulli.

- Les seuls paramètres d'entrée de la fonction doivent être :
  - `X` :: Une matrice (array *numpy*) de dimension  $n \times d$  d'images binarisées à utiliser pour l'estimation des paramètres.
  - `K` :: Le nombre de classes du mélange de Bernoulli
  - `max_iter` :: Nombre maximal d'itérations, *i.e.* chaque itération est composée d'une étape E (E-step) et une étape M (M-step)

- Initialiser les paramètres  $\mu_{kj}$  à l'aide d'une loi uniforme entre 0.15 et 0.85 et les  $\pi_k$  à  $\frac{1}{K}$ .
- L'implémentation des éléments retournés par la fonction est laissée à votre appréciation.
- Utiliser un critère d'arrêt raisonnable afin d'éviter de parcourir inutilement toutes les itérations. Penser à vectoriser les calculs en utilisant `numpy` quand c'est possible.
- Vous devez implémenter les fonctions `E_step(X, mu, pi)` et `M_step(X, gamma)` séparément pour les utiliser dans votre fonction `EM(X, K, max_iter)`.

### Test avec trois chiffres

#### Question 3

- En analogie avec le livre de Bishop partie §9.3.3, échantillonnez un ensemble de données d'apprentissage composé uniquement d'images binarisées des chiffres 2, 3, et 4. Exécutez votre algorithme EM et afficher les paramètres  $\mu_1, \mu_2, \mu_3$  sous formes d'images.
- Pouvez-vous identifier quel classe correspond à quel chiffre ? Quels sont les tailles du mélange  $\pi_1, \pi_2, \pi_3$  identifiées pour les chiffres 2, 3 et 4. Ces valeurs correspondent-elles aux vraies proportions des images des 3 chiffres dans le jeu de données ?

### Choix de modèle

La question 1.3 faisait appel à 3 classes d'images seulement (chiffres 2, 3 et 4) pour éviter un temps de calcul trop lent. La nombre de classes  $K = 3$  est une connaissance a priori car nous avons sélectionné les données nous mêmes. Souvent le nombre de classes dans un jeu de données est inconnu et son estimation est une question de choix de modèle. Ici nous allons faire appel au critère BIC *Bayesian information criterion* pour choisir le nombre de classes. À nombre de classes fixé  $K$  le critère BIC est donné par

$$BIC(K) = 2\hat{\ell}_K - p_K \ln(n),$$

où

- $\hat{\ell}_K$  est la valeur maximale de la log-vraisemblance associée au modèle à  $K$  classes.
- $p_K$  est le nombre de paramètres estimés dans le modèle (**attention :**  $\sum_{k=1}^K \pi_k = 1$ ).

#### Question 4

- Ajouter aux éléments retournés par la fonction `EM(X, K, max_iter)`, la valeur du BIC associé au modèle estimé.
- Tester les valeur  $K = 2, \dots, 5$ . Quelle est la valeur de  $K$  qui maximise le critère BIC ?

### **Images mal-classées**

Jusqu'ici, nous n'avons pas utilisé les vraies étiquettes des images dans une étape d'estimation.

#### **Question 5**

Utilisez les étiquettes des images pour identifier les cas "mal-classés" et afficher quelques exemples d'images mal-classées (afficher dans le titre la vraie étiquette et la fausse étiquette attribuée). Peut-on donner une explication du mauvais classement sur certains cas ?

### **Initialiser les paramètres $\pi_k$ avec les vraies valeurs**

#### **Question 6**

Initialisez les tailles des trois classes ( $\pi_1, \pi_2, \pi_3$ ) avec les vraies valeurs (proportions d'images correspondantes aux chiffres 2, 3 et 4 respectivement dans le jeu de données). Que remarque-t-on ?