

Consignes : Vous devez séparer pour chacun des exercices votre programme principal (l'équivalent du main) du reste du code (fonctions etc). Le dépôt de votre travail se fera sous la forme d'une archive (qui portera votre nom et qui contiendra également les données) à déposer sur ecampus (exclusivement) le 14/12/2017 avant 16h00, heure stricte.

Barème : Le barème est donné à titre indicatif. Les exercices pour lesquels les codes ne s'exécutent pas, seront notés sur la moitié des points du barème.

Ex1 : Analyse supervisée (10 pts) :

On considère un échantillon indépendant $((\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n))$ d'individus \mathbf{X}_i décrits par d variables réelles ($\mathbf{X}_i \in \mathbb{R}^d$) d'une population de K classes hétérogènes telle que $Y_i \in \llbracket 1, K \rrbracket$ est la classe de l'individu \mathbf{X}_i . On dispose d'un échantillon observé datatrain-bis $((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n))$. L'objectif est de prédire les classes pour un échantillon de test issu de cette même population sur la base d'un modèle probabiliste paramétrique de paramètres θ . La prédiction peut ainsi s'effectuer par la règle du maximum a posteriori (MAP) qui consiste à maximiser la probabilité a posteriori pour avoir une prédiction \hat{y}_i de la classe de l'individu \mathbf{x}_i , étant donné un modèle de paramètres $\hat{\theta}$ appris à partir de l'échantillon observé :

$$\hat{y}_i = \arg \max_{y_i \in \llbracket 1, K \rrbracket} \mathbb{P}(Y_i = y_i | \mathbf{X} = \mathbf{x}_i; \theta). \tag{1}$$

Modélisation discriminative : Régression logistique : Une méthode probabiliste discriminative recherche directement un modèle de la loi conditionnelle $\mathbb{P}(Y|\mathbf{X})$ et s'intéresse donc prioritairement à séparer les classes. C'est le cas de la régression logistique qui s'appuie sur un modèle linéaire généralisé de la forme suivante pour les loi conditionnelles

$$\mathbb{P}(Y = k | \mathbf{X} = \mathbf{x}; \theta) = \frac{\exp(\beta_{k0} + \beta_k^T \mathbf{x})}{1 + \sum_{\ell=1}^{K-1} \exp(\beta_{\ell 0} + \beta_{\ell}^T \mathbf{x})} \tag{2}$$

pour $k \in \llbracket 1, K-1 \rrbracket$ et $\mathbb{P}(Y = K | \mathbf{X} = \mathbf{x}; \theta) = \frac{1}{1 + \sum_{\ell=1}^{K-1} \exp(\beta_{\ell 0} + \beta_{\ell}^T \mathbf{x})}$.

Cela correspond donc à une transformation par une fonction exponentielle normalisée (appelée aussi softmax) d'une fonction linéaire en \mathbf{x} . Cela explique en effet le fait que en régression logistique on cherche à modéliser les probabilités a posteriori des classes à travers des fonctions linéaires en \mathbf{x} puisque que l'on a $\log \left(\frac{\mathbb{P}(Y=k|\mathbf{X}=\mathbf{x})}{\mathbb{P}(Y=K|\mathbf{X}=\mathbf{x})} \right) = \beta_{k0} + \beta_k^T \mathbf{x}$ pour tout $k \in \llbracket 1, K-1 \rrbracket$.

Le vecteur paramètre du modèle défini par $\theta = (\beta_1^T, \dots, \beta_{K-1}^T)^T \in \mathbb{R}^{(K-1) \times (d+1)}$ avec $\beta_k = (\beta_{k0}, \beta_k^T)^T$ est estimé en maximisant la log-vraisemblance conditionnelle par rapport à θ :

$$\begin{aligned} \log L(\theta; \mathbf{X}, \mathbf{y}) &= \log \prod_{i=1}^n \mathbb{P}(Y_i | \mathbf{X}_i = \mathbf{x}_i; \theta) \\ &= \log \prod_{i=1}^n \prod_{k=1}^K \mathbb{P}(Y_i = k | \mathbf{X}_i = \mathbf{x}_i; \theta)^{y_{ik}} \\ &= \sum_{i=1}^n \sum_{k=1}^K y_{ik} \log \mathbb{P}(Y_i = k | \mathbf{X}_i = \mathbf{x}_i; \theta) \end{aligned} \tag{3}$$

où y_{ik} est une variable binaire telle que $y_{ik} = 1$ si et seulement si $y_i = k$ (i.e, \mathbf{x}_i appartient à la classe k). Cette log-vraisemblance ne peut pas être maximisée analytiquement et pour ce faire on utilise l'algorithme Newton-Raphson.

On considère un problème à deux classes ($Y_i \in \llbracket 0, 1 \rrbracket$, i.e., régression logistique binaire). Dans ce cas le modèle est défini par les probabilités

$$\pi(\mathbf{x}; \boldsymbol{\theta}) = \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}; \boldsymbol{\theta}) = \frac{\exp(\beta_{10} + \boldsymbol{\beta}_1^T \mathbf{x})}{1 + \exp(\beta_{10} + \boldsymbol{\beta}_1^T \mathbf{x})} \quad (4)$$

et $\mathbb{P}(Y = 0 | \mathbf{X} = \mathbf{x}; \boldsymbol{\theta}) = 1 - \pi(\mathbf{x}; \boldsymbol{\theta})$.

Le vecteur paramètre du modèle défini par $\boldsymbol{\theta} = (\beta_0, \boldsymbol{\beta}_1^T)^T \in \mathbb{R}^{d+1}$ est estimé en maximisant par rapport à $\boldsymbol{\theta}$ la log-vraisemblance conditionnelle qui dans ce cas de classification binaire prend la forme :

$$\begin{aligned} \log L(\boldsymbol{\theta}; \mathbf{X}, \mathbf{y}) &= \log \prod_{i=1}^n \mathbb{P}(Y_i = 1 | \mathbf{X}_i = \mathbf{x}_i; \boldsymbol{\theta})^{y_i} \mathbb{P}(Y_i = 0 | \mathbf{X}_i = \mathbf{x}_i; \boldsymbol{\theta})^{1-y_i} \\ &= \sum_{i=1}^n y_i \log \mathbb{P}(Y_i = 1 | \mathbf{X}_i = \mathbf{x}_i; \boldsymbol{\theta}) + (1 - y_i) \log \mathbb{P}(Y_i = 0 | \mathbf{X}_i = \mathbf{x}_i; \boldsymbol{\theta}) \\ &= \sum_{i=1}^n y_i \log \pi(\mathbf{x}_i; \boldsymbol{\theta}) + (1 - y_i) \log (1 - \pi(\mathbf{x}_i; \boldsymbol{\theta})) \\ &= \sum_{i=1}^n y_i (\beta_{10} + \boldsymbol{\beta}_1^T \mathbf{x}_i) - \log(1 + \exp(\beta_{10} + \boldsymbol{\beta}_1^T \mathbf{x}_i)). \end{aligned} \quad (5)$$

On montre que dans le cas de ce modèle, l'algorithme de Newton-Raphson consiste à partir d'un modèle initial de paramètres $\boldsymbol{\theta}^{(0)}$ et à mettre à jour, à chaque itération t , les paramètres selon l'équation de mise à jour suivante, jusqu'à ce qu'il n'y ait plus d'augmentation significative au sens d'un seuil préfixé de la log-vraisemblance conditionnelle (3) :

$$\begin{aligned} \boldsymbol{\theta}^{(t+1)} &= \boldsymbol{\theta}^{(t)} + (\mathbf{X}^T \mathbf{W}^{(t)} \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \mathbf{p}^{(t)}) \\ &= (\mathbf{X}^T \mathbf{W}^{(t)} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}^{(t)} \tilde{\mathbf{y}} \end{aligned} \quad (6)$$

où :

- \mathbf{X} est la matrice de dimensions $n \times (d + 1)$ de lignes $(1, \mathbf{x}_i^T)$ pour $i = 1, \dots, n$;
- \mathbf{y} est le vecteur colonne de dimension $n \times 1$ des labels y_i : $\mathbf{y} = (y_1, \dots, y_n)^T$;
- \mathbf{p} est le vecteur colonne de dimension $n \times 1$ des probabilités : $\mathbf{p} = (\pi(\mathbf{x}_1; \boldsymbol{\theta}), \dots, \pi(\mathbf{x}_n; \boldsymbol{\theta}))^T$ avec $\pi(\mathbf{x}_i; \boldsymbol{\theta}) = \mathbb{P}(Y_i = 1 | \mathbf{x}_i; \boldsymbol{\theta}) = \frac{\exp(\boldsymbol{\beta}_k^T \mathbf{x}_i)}{1 + \exp(\boldsymbol{\beta}_k^T \mathbf{x}_i)}$;
- \mathbf{W} est la matrice diagonale de dimension $n \times n$ d'éléments $\pi(\mathbf{x}_1; \boldsymbol{\theta}) (1 - \pi(\mathbf{x}_n; \boldsymbol{\theta}))$: $\mathbf{W} = \text{diag}(\mathbf{p})$.
- $\tilde{\mathbf{y}} = \mathbf{X} \boldsymbol{\theta}^{(t)} + (\mathbf{W}^{(t)})^{-1} (\mathbf{y} - \mathbf{p}^{(t)})$

De part la forme de (6) qui correspond à celle de l'estimateur des moindres carrés pondérés (avec une matrice de poids \mathbf{W}) on appelle l'algorithme Netwton-Raphson dans ce cas l'algorithme IRLS (Iterative Reweighted Least Squares pour moindres carrés pondérés itératifs).

Travail demandé :

1. Implémenter le et tester le sur l'échantillon de datatest-bis issus d'une même population de K classes hétérogènes de données de dimension d . Pour cela, créer une fonction `train_LogReg.m` pour l'apprentissage et une fonction `predict_LogReg.m` pour la prédiction.
2. Calculer le taux d'erreur de prédiction pour l'échantillon de test étiqueté `datatest`.
3. Maintenant considérer l'échantillon `AllXtest` et prédire ses classes
4. Afficher, sur le même graphique, les données d'apprentissage et les données de test classées
5. Comparer vos résultats à ceux obtenus en se basant sur les fonctions `mnrfit` et `mnrval` de Matlab.
6. Comparer vos résultats à ceux obtenus en se basant sur les fonctions `glmfit` et `glmval` de Matlab.

Ex2 : Analyse non-supervisée (10 pts) :

On dispose d'un n -échantillon $(\mathbf{x}_1, \dots, \mathbf{x}_n)$ issu d'une population hétérogène de K classes d'individus binaires multivariés décrits par d variables indépendantes : $\mathbf{x}_i = (x_{i1}, \dots, x_{id})^T \in \{0, 1\}^d$ et distribués selon une loi mélange de lois de Bernoulli multivariées définie par :

$$\mathbb{P}(\mathbf{X}_i = \mathbf{x}_i; \boldsymbol{\theta}) = \sum_{k=1}^K \pi_k \mathcal{B}(\mathbf{x}_i; \mathbf{p}_k) \quad (7)$$

où $\mathcal{B}(\mathbf{x}_i; \mathbf{p}_k) = \prod_{j=1}^d \mathcal{B}(x_{ij}; p_{kj})$ est la loi de Bernoulli multivariée de la classe k avec $\mathcal{B}(x_{ij}; p_{kj})$ la loi de Bernoulli univariée de paramètre $p_{kj} \in]0, 1[$ associée à la variable x_{ij} et définie par $\mathcal{B}(x; p) = p^x (1-p)^{1-x}$.

Afin d'apprendre ce modèle en estimant les paramètres $\boldsymbol{\theta} = (\pi_1, \dots, \pi_K, \mathbf{p}_1, \dots, \mathbf{p}_K)$ à partir des données et faire la classification non-supervisée des données, on utilise l'algorithme Espérance-Maximisation (EM). On montre que pour ce modèle l'algorithme EM consiste à partir d'un modèle initial de paramètres $\boldsymbol{\theta}^{(0)}$ et alterner à chaque itération t entre les deux étapes E- et M- suivantes jusqu'à ce qu'il n'y ait plus d'augmentation significative au sens d'un seuil préfixé de la log-vraisemblance du modèle :

1. Étape E : Calculer les probabilités a posteriori : $\tau_{ik}^{(t)} = \mathbb{P}(Y_i = k | \mathbf{x}_i, \boldsymbol{\theta}^{(t)}) = \frac{\pi_k^{(t)} \mathcal{B}(\mathbf{x}_i; \mathbf{p}_k^{(t)})}{\sum_{\ell=1}^K \pi_{\ell}^{(t)} \mathcal{B}(\mathbf{x}_i; \mathbf{p}_{\ell}^{(t)})}$
2. Étape M : Mettre à jour les paramètres du modèle $\boldsymbol{\theta}^{(t+1)}$: $\pi_k^{(t+1)} = \frac{\sum_{i=1}^n \tau_{ik}^{(t)}}{n}$ et $\mathbf{p}_k^{(t+1)} = \frac{\sum_{i=1}^n \tau_{ik}^{(t)} \mathbf{x}_i}{\sum_{i=1}^n \tau_{ik}^{(t)}}$.

Travail demandé :

1. Implémenter l'algorithme EM pour estimer les paramètres $\boldsymbol{\theta}$ du modèle (7) et les classes $\mathbf{y} = (y_1, \dots, y_n)$ des données `Xbin` sur la base de ce modèle avec la règle du MAP (1). On pourra fixer la valeur de K à 3.
2. Affiche les valeurs de la log-vraisemblance enregistrée au cours de l'algorithme et les données classées. Pour ce dernier graphique, on utilisera la commande `imagesc(X(sort(y), :))` où `X` sont les données et `y` les classes obtenues.
3. Utiliser le critère BIC pour sélectionner le nombre de classes K sachant que $K \in \llbracket 1, 10 \rrbracket$