

Classification d'images de chats et de chiens.

Gautier Appert
gautier.appert.chess@gmail.com



Soit $\mathcal{D}_n = \{x_1, \dots, x_n\}$ une base de données d'images de chat et de chien où chaque image est représentée par un vecteur de pixels $x_i \in \mathbb{R}^p$, $i \in \{1, \dots, n\}$. Notons que le nombre de pixels p est potentiellement très largement supérieur à n . On note $Y_i \in \{0, 1\}$ la variable aléatoire correspondant au label chat ou chien associé à l'image $x_i \in \mathbb{R}^p$. En pratique la base données \mathcal{D}_n stock les images en vecteur lignes

$$\mathcal{D}_n = [x_1 \mid x_2 \mid \dots \mid x_n]^\top.$$

L'objet de ce tutoriel est de prédire le label chat ou chien à l'aide d'une analyse discriminante quadratique (QDA) pour une nouvelle image en dehors de la base d'apprentissage $x \notin \mathcal{D}_n$. L'implémentation doit être faite sous le langage R. Envoyer un mail au chargé de TD afin de pouvoir récupérer les données sur la dropbox.

1. DÉCOUVERTE DE LA BASE DE DONNÉES ET RÉDUCTION DE LA DIMENSION

Deux bases de données intitulées $X_{\text{train}}.\text{RData}$ et $X_{\text{test}}.\text{RData}$ sont à disposition sur la dropbox. Ces bases contiennent des images de chiens et de chats stockées en vecteur ligne. En particulier on a $X_{\text{train}} \in \mathbb{R}^{315 \times 40000}$ et $X_{\text{test}} \in \mathbb{R}^{48 \times 40000}$. Deux autres bases de données $Y_{\text{train}}.\text{RData}$ et $Y_{\text{test}}.\text{RData}$ contiennent les labels associés aux images X_{train} et X_{test} .

QUESTION (1). Importer la base de données $X_{\text{train}}.\text{RData}$, $X_{\text{test}}.\text{RData}$, $y_{\text{train}}.\text{RData}$ et $y_{\text{test}}.\text{RData}$ à l'aide de la fonction `load`. Afficher deux à trois images des deux bases à l'aide de la fonction `image(..., col = grey(seq(0, 1, length = 256)))` en transformant préalablement les images en matrice de dimension 200×200 (fonction `matrix`). Enregistrer les images en format `pdf` ou `png` et les mettre dans votre rapport. A quoi correspond le label $y = 1$?

QUESTION (2). On souhaite réduire la dimension des données $p = 40000$. Pour cela nous allons procéder à une analyse en composante principales (ACP) des images.

(a). Concatener X_{train} et X_{test} en utilisant la fonction `rbind` et centrer les vecteurs colonnes avec la fonction `scale`. On notera $X \in \mathbb{R}^{363 \times 40000}$ la matrice résultante. Construire une ACP en utilisant une décomposition en valeur singulière (SVD) de la matrice X à l'aide la fonction `svd`. On ne retiendra que les 15 premières composantes principales. On rappelle que la décomposition en valeur singulière permet de factoriser la matrice X de la manière suivante

$$X = UDV^\top,$$

où V est la matrice des vecteurs propres. Ainsi la matrice des composantes principales est donnée par $C = XV$.

(b). Quelle est la part de variance expliquée en ne retenant que 15 composantes principales ? Désormais nous travaillerons sur les composantes principales $C \in \mathbb{R}^{363 \times 15}$ au lieu des données d'origine X . Découper la base C en $C_{\text{train}} \in \mathbb{R}^{315 \times 15}$ et $C_{\text{test}} \in \mathbb{R}^{48 \times 15}$.

2. ANALYSE DISCRIMINANTE QUADRATIQUE

On fait l'hypothèse du modèle suivant

- $Y_i \sim \mathcal{B}(\pi)$.
- $P_{c_i|Y=1} = \mathcal{N}(\mu_1, \Sigma_1)$ et $P_{c_i|Y=0} = \mathcal{N}(\mu_0, \Sigma_0)$ où c_i est le i -ième vecteur ligne de la matrice C .

Le paramètre inconnu est $\theta = (\pi, \mu_0, \mu_1, \Sigma_1, \Sigma_0)$ où $\pi \in]0, 1[$, $(\mu_0, \mu_1) \in \mathbb{R}^{15} \times \mathbb{R}^{15}$ et $(\Sigma_0, \Sigma_1) \in \mathbb{R}^{15 \times 15} \times \mathbb{R}^{15 \times 15}$ sont des matrices définies positives. On définit $P_\theta = P_{c,Y}$ et on dispose d'un échantillon $(c_1, y_1), \dots, (c_n, y_n) \stackrel{\text{i.i.d}}{\sim} P_\theta$.

QUESTION (3). Ecrire le modèle statistique associé aux observations $(c_1, y_1), \dots, (c_n, y_n)$.

QUESTION (4). On pose $N_1 = \sum_{i=1}^n y_i$ et $N_2 = n - N_1$. En utilisant le fait que $f_{c,Y}(c, y) = f_{c|Y=y}(c) f_Y(y)$, montrer que la log vraisemblance $\ell((c_1, y_1), \dots, (c_n, y_n); \theta)$ s'écrit

$$\ell((c_1, y_1), \dots, (c_n, y_n); \theta) = N_1 \log(\pi) + N_2 \log(1 - \pi) - \frac{N_1}{2} \log(\det(\Sigma_1)) - \frac{1}{2} \sum_{i:y_i=1} (c_i - \mu_1)^\top \Sigma_1^{-1} (c_i - \mu_1) - \frac{N_2}{2} \log(\det(\Sigma_0)) - \frac{1}{2} \sum_{i:y_i=0} (c_i - \mu_0)^\top \Sigma_0^{-1} (c_i - \mu_0).$$

QUESTION (5). En utilisant les formules $\nabla_\Sigma \log(\det(\Sigma)) = \Sigma^{-1}$ et $\nabla_\Sigma (a^\top \Sigma^{-1} b) = -\Sigma^{-1} a b^\top \Sigma^{-1}$, écrire l'équation du premier ordre pour le maximum de vraisemblance et montrer que l'on obtient les estimateurs

$$\hat{\pi} = \frac{N_1}{n} \quad \hat{\mu}_1 = \frac{1}{N_1} \sum_{i:y_i=1} c_i \quad \hat{\mu}_0 = \frac{1}{N_0} \sum_{i:y_i=0} c_i$$

$$\hat{\Sigma}_1 = \frac{1}{N_1} \sum_{i:y_i=1} (c_i - \hat{\mu}_1)(c_i - \hat{\mu}_1)^\top \quad \hat{\Sigma}_0 = \frac{1}{N_0} \sum_{i:y_i=0} (c_i - \hat{\mu}_0)(c_i - \hat{\mu}_0)^\top.$$

QUESTION (6). Montrer que la sous Hessienne $\nabla_{\pi, \mu_1, \mu_0}^2 \ell(\theta)$ est bien définie négative. On ne regardera pas les conditions du second ordre avec Σ_1 et Σ_0 .

QUESTION (7). Montrer que $\hat{\pi}$ est sans biais et montrer que $\hat{\mu}_1$ et $\hat{\mu}_0$ sont sans biais (conditionner par rapport à l'échantillon $\{y_1, \dots, y_n\}$ via la loi des espérances itérées.)

QUESTION (8). Montrer que les estimateurs issus de la méthode des moments coïncident avec les estimateurs du maximum de vraisemblance. (on pourra utiliser la définition de l'espérance conditionnelle sachant un événement $\mathbb{E}[C|Y=y] = \frac{\mathbb{E}[C \mathbb{1}(Y=y)]}{\mathbb{E}[\mathbb{1}(Y=y)]}$).

QUESTION (9). Coder une fonction sous R intitulée `computeML(C, Y)` prenant en argument une matrice C et un vecteur Y , et qui renvoie sous forme de liste les estimateurs du maximum de vraisemblance $\hat{\pi}, \hat{\mu}_1, \hat{\mu}_0, \hat{\Sigma}_1, \hat{\Sigma}_0$. Lancer la fonction `computeML` sur `Ctrain, Ytrain`. Comparer les estimateurs obtenus avec la fonction `qda(Ctrain, Ytrain)` du package MASS. (La fonction `qda` ne fournit pas les estimateurs concernant les matrices de variances covariances mais fournit le log du déterminant).

3. PRÉDICTION DES LABELS SUR LA BASE TEST

On souhaite dans cette partie prédire les labels correspondant aux données C_{test} à l'aide de l'analyse discriminante quadratique dont les paramètres ont été estimés sur l'échantillon d'apprentissage (C_{train} , Y_{train}). En effet, l'analyse discriminante quadratique permet de modéliser les probabilités $\mathbb{P}(Y = 1|c)$ et $\mathbb{P}(Y = 0|c)$. C'est pourquoi nous prendrons la règle de prédiction suivante $\hat{y} = \arg \max_{y \in \{0,1\}} \mathbb{P}(Y = y|c)$.

QUESTION (10). A l'aide de la formule de Bayes, montrer que

$$\mathbb{P}(Y = 1|c) = \frac{\pi \varphi(c; \mu_1, \Sigma_1)}{\pi \varphi(c; \mu_1, \Sigma_1) + (1 - \pi) \varphi(c; \mu_0, \Sigma_0)}$$

où $\varphi(c; \mu, \Sigma)$ représente la densité de la Gaussienne multivariée $\mathcal{N}(\mu, \Sigma)$. Calculer de la même manière $\mathbb{P}(Y = 0|c)$.

QUESTION (11). En déduire que

$$\begin{aligned} \log \left(\frac{\mathbb{P}(Y = 1|c)}{\mathbb{P}(Y = 0|c)} \right) &= -\frac{1}{2} \log(\det(\Sigma_1)) - \frac{1}{2} (c - \mu_1)^\top \Sigma_1^{-1} (c - \mu_1) + \log(\pi) \\ &\quad + \frac{1}{2} \log(\det(\Sigma_0)) + \frac{1}{2} (c - \mu_0)^\top \Sigma_0^{-1} (c - \mu_0) - \log(1 - \pi). \end{aligned}$$

QUESTION (12). Montrer que

$$\mathbb{1} \left(\log \left(\frac{\mathbb{P}(Y = 1|c)}{\mathbb{P}(Y = 0|c)} \right) > 0 \right) = \arg \max_{y \in \{0,1\}} \mathbb{P}(Y = y|c).$$

Ainsi, on utilisera la règle de prédiction suivante (méthode Plug-in): pour toutes lignes $c \in C_{\text{test}}$

$$\begin{aligned} \hat{y} &= \mathbb{1} \left(-\frac{1}{2} \log(\det(\hat{\Sigma}_1)) - \frac{1}{2} (c - \hat{\mu}_1)^\top \hat{\Sigma}_1^{-1} (c - \hat{\mu}_1) + \log(\hat{\pi}) \right. \\ &\quad \left. + \frac{1}{2} \log(\det(\hat{\Sigma}_0)) + \frac{1}{2} (c - \hat{\mu}_0)^\top \hat{\Sigma}_0^{-1} (c - \hat{\mu}_0) - \log(1 - \hat{\pi}) > 0 \right). \end{aligned}$$

QUESTION (13). En utilisant le fait que $(y, \hat{y}) \in \{0, 1\}^2$, montrer la double égalité

$$\mathbb{E}[(y - \hat{y})^2] = \mathbb{E}[|y - \hat{y}|] = \mathbb{P}(\hat{y} \neq y).$$

Empiriquement on prendra

$$\hat{\mathbb{E}}[|y - \hat{y}|] = \frac{1}{n} \sum_{i=1}^n |\hat{y}_i - y_i|,$$

ce qui correspond à l'erreur de classification.

QUESTION (14). Ecrire une fonction R intitulée `computeLogRatio(c, pi, mu1, mu0, Sigma1, Sigma0)`

prenant en argument un vecteur $c \in C_{\text{test}}$ et le paramètre θ , et qui renvoie la quantité: $\log \left(\frac{\mathbb{P}(Y=1|c)}{\mathbb{P}(Y=0|c)} \right)$.

Puis coder une fonction `computePred(C, pi, mu1, mu0, Sigma1, Sigma0)` prenant en argument une matrice C et le paramètre θ et qui renvoie la prédiction des labels pour chaque ligne de la matrice C .

QUESTION (15). Prédire les labels de la base de données test C_{test} avec `computePred` et donner l'erreur de classification à l'aide de Y_{test} . Comparer avec la prédiction en utilisant l'estimation du modèle faite avec la fonction `qda` de **R**. La prédiction est-elle meilleur que le prédicteur aléatoire ?

RÉFÉRENCE

- *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, Trevor Hastie, Robert Tibshirani, Jerome Friedman.