
TP N° 2 : k -plus proches voisins, analyse
discriminante linéaire et régression logistique

Vous devez télécharger votre fichier de solution sur le site pédagogique sous le format `ipynb` avant le mercredi 27/11/2019 23h59. En plus des questions du TP, les aspects suivants seront évalués :

- aspect global de présentation : qualité de rédaction, orthographe, présentation, graphes, titres, etc...,
- aspect global du code : indentation, style, lisibilité du code, commentaires adaptés,
- absence de bugs.

Le fichier `tp_knn_source.py` est disponible sur le site pédagogique du cours. Ils contiennent le code et les fonctions utiles pour la partie sur les k -plus proches voisins.

- DÉCOUVERTE DE PYTHON -

Consulter les pages suivantes pour démarrer ou bien trouver quelques rappels de Python :

- <http://www.python.org>
- <http://scipy.org>
- <http://www.numpy.org>
- <http://scikit-learn.org/stable/index.html>
- <http://www.loria.fr/~rougier/teaching/matplotlib/matplotlib.html>
- <http://jrjohansson.github.io/>

- RAPPELS DE CLASSIFICATION -

Définitions et notations

On rappelle ici le cadre de la classification supervisée, et l'on présente les notations que l'on utilisera dans la suite. Il est à noter que pour ce TP on considère un cadre plus général qu'au TP précédent : le nombre de classes peut-être plus grand que deux.

- \mathcal{Y} est l'ensemble des étiquettes des données (*labels* en anglais). Ici on raisonne avec un nombre L quelconque de classes, et l'on choisit $\mathcal{Y} = \{1, \dots, L\}$ pour représenter les L étiquettes possibles (le cas de la classification binaire est le cas où $L = 2$).
- $\mathbf{x} = (x_1, \dots, x_p)^\top \in \mathcal{X} \subset \mathbb{R}^p$ est une observation, un exemple, un point (ou un *sample* en anglais). La j ème coordonnée de \mathbf{x} est la valeur prise par la j ème variable (*feature* en anglais).
- $\mathcal{D}_n = \{(\mathbf{x}_i, y_i), i = 1, \dots, n\}$ est l'ensemble d'apprentissage contenant les n exemples et leurs étiquettes.
- Il existe un modèle probabiliste qui gouverne la génération de nos observations selon des variables aléatoires X et $Y : \forall i \in \{1, \dots, n\}, (\mathbf{x}_i, y_i) \stackrel{i.i.d.}{\sim} (X, Y)$.
- On cherche à construire à partir de l'ensemble d'apprentissage \mathcal{D}_n une fonction appelée classifieur, $\hat{f} : \mathcal{X} \mapsto \mathcal{Y}$ qui à un nouveau point \mathbf{x}_{new} associe une étiquette $\hat{f}(\mathbf{x}_{\text{new}})$.

1. Avec la contribution majeure de Alexandre Gramfort et Joseph Salmon.

Génération artificielle de données

On considère dans cette partie des observations décrites en deux dimensions (afin de pouvoir les visualiser facilement) à savoir $p = 2$ dans le formalisme ci-dessus. On reprend en partie les jeux de données artificiels du TP sur le perceptron, avec des modifications sur les valeurs des étiquettes uniquement.

- 1) Étudiez les fonctions `rand_bi_gauss`, `rand_tri_gauss`, `rand_clown` et `rand_checkers`. Que renvoient ces fonctions ? À quoi correspond la dernière colonne ? Générez les données en utilisant chaque de ces 4 fonctions avec les paramètres du corrigé du TP1 (pour `rand_checkers` prenez 150 pour les deux premiers arguments ; pour `rand_tri_gauss` générez 50 observations dans chaque classe, centrées en $(1, 1)^\top$, $(-1, -1)^\top$ et $(1, -1)^\top$, tous les écart-types égaux à 0.9). On va nommer ces jeux de données #1, #2, #3 et #4, respectivement.
- 2) Utilisez la fonction `plot_2d` afin d'afficher les jeux de données générés avec chacune des fonctions.

- LA MÉTHODE DES k -PLUS PROCHES VOISINS -

Approche intuitive

L'algorithme des k -plus proches voisins (k -nn : pour *k-nearest neighbors* en anglais) est un algorithme intuitif, aisément paramétrisable pour traiter un problème de classification avec un nombre quelconque d'étiquettes.

Le principe de l'algorithme est particulièrement simple : pour chaque nouveau point \mathbf{x} on commence par déterminer l'ensemble de ses k -plus proches voisins parmi les points d'apprentissage que l'on note $V_k(\mathbf{x})$ (bien sûr on doit choisir $1 \leq k \leq n$ pour que cela ait un sens). La classe que l'on affecte au nouveau point \mathbf{x} est alors la classe majoritaire dans l'ensemble $V_k(\mathbf{x})$. Une illustration de la méthode est donnée en Figure 1 pour le cas de trois classes.

- 3) Proposez une version adaptée de cette méthode pour la régression, *i.e.*, quand les observations y sont à valeurs réelles : $\mathcal{Y} = \mathbb{R}$.

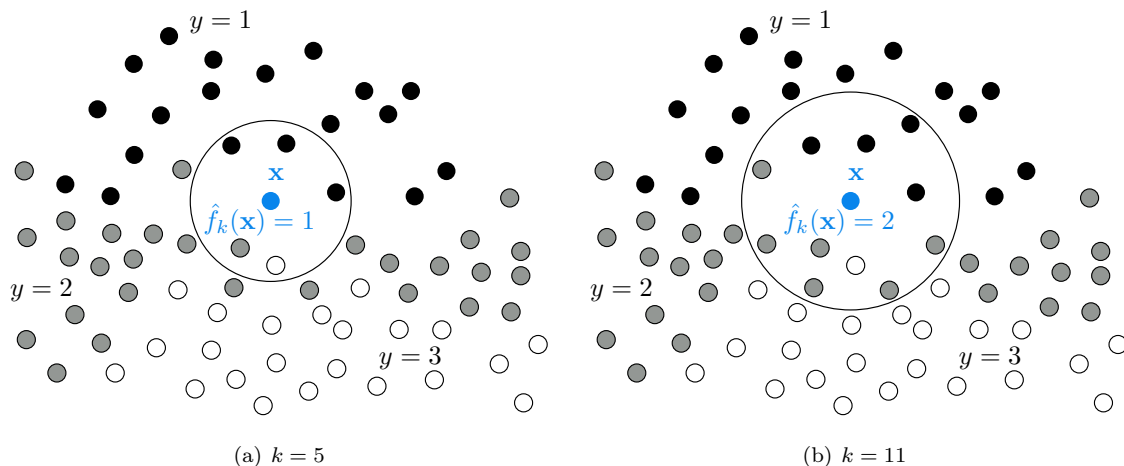


FIGURE 1 – Exemple de fonctionnement de la méthode des k -plus proches voisins pour des valeurs du paramètre $k = 5$ et $k = 11$. On considère trois classes, $L = 3$, représentées respectivement en noir ($y = 1$), en gris ($y = 2$) et en blanc ($y = 3$).

Approche formelle

Pour définir précisément la méthode, il faut commencer par choisir une distance $d : \mathbb{R}^p \times \mathbb{R}^p \mapsto \mathbb{R}$. Pour un nouveau point \mathbf{x} , on définit alors l'ensemble de ses k -plus proches voisins $V_k(\mathbf{x})$ au sens de de cette distance. On peut procéder de la manière suivante : pour chaque $\mathbf{x} \in \mathbb{R}^d$ et pour chaque $i = 1, \dots, n$,

on note $d_i(\mathbf{x})$ la distance entre \mathbf{x} et \mathbf{x}_i : $d_i(\mathbf{x}) = d(\mathbf{x}_i, \mathbf{x})$. On définit la première statistique de rang $r_1(\mathbf{x})$ comme l'indice du plus proche voisin de \mathbf{x} parmi $\mathbf{x}_1, \dots, \mathbf{x}_n$, c'est-à-dire

$$r_1(\mathbf{x}) = i^* \quad \text{si et seulement si} \quad d_{i^*}(\mathbf{x}) = \min_{1 \leq i \leq n} d_i(\mathbf{x}).$$

Remarque 1. *S'il y a plusieurs candidats pour la minimisation ci-dessus, on ordonne les ex-æquos de manière arbitraire (généralement aléatoirement).*

Par récurrence on peut ainsi définir le rang $r_k(\mathbf{x})$ pour tout entier $1 \leq k \leq n$:

$$r_k(\mathbf{x}) = i^* \quad \text{si et seulement si} \quad d_{i^*}(\mathbf{x}) = \min_{\substack{1 \leq i \leq n \\ i \neq r_1, \dots, r_{k-1}}} d_i(\mathbf{x}). \quad (1)$$

L'ensemble de k -plus proches voisins de \mathbf{x} s'exprime alors par $V_k(\mathbf{x}) = \{\mathbf{x}_{r_1}, \dots, \mathbf{x}_{r_k}\}$. Pour finir, la décision pour classifier le point \mathbf{x} se fait par vote majoritaire, en résolvant le problème suivant :

$$\hat{f}_k(\mathbf{x}) \in \arg \max_{y \in \mathcal{Y}} \left(\sum_{j=1}^k \mathbb{1}_{\{y_{r_j} = y\}} \right). \quad (2)$$

Le module `sklearn.neighbors` de scikit-learn, (cf. <http://scikit-learn.org/stable/modules/neighbors.html>) implémente les méthodes de classification et régression à base de k -plus proches voisins.

- 4) Écrivez votre propre classe `KNNClassifier` avec les méthodes d'apprentissage `fit` et de classification `predict`. Choisissez une stratégie de traitement des ex aequo, c'est-à-dire des points avec la même distance ou le même nombre de points de chaque classe dans $V_k(\mathbf{x})$. Vérifier la validité des résultats en les comparant à ceux de la classe `KNeighborsClassifier` de scikit-learn en utilisant le jeu de données #2. Vous proposerez votre propre méthode de comparaison (par exemple, en utilisant les observations d'indice pair pour le set d'apprentissage et celles d'indice impair pour le set de test). Vous pouvez utiliser le bloc de code ci-dessous en complétant les méthodes proposées. Pour plus d'information sur les classes on peut consulter par exemple <http://docs.python.org/3/tutorial/classes.html>.

```
from sklearn.base import BaseEstimator, ClassifierMixin

class KNNClassifier(BaseEstimator, ClassifierMixin):
    """ Homemade kNN classifier class """
    def __init__(self, n_neighbors=1):
        # Complete the method

    def fit(self, X, y):
        # Complete the method
        return self

    def predict(self, X):
        # Complete the method
```

Pour gagner en temps de calcul, vous utiliserez à partir de maintenant l'implémentation de scikit-learn.

- 5) Faites tourner sur les quatre exemples de jeu de données cet algorithme de classification, en utilisant la distance euclidienne classique $d(\mathbf{x}, \mathbf{v}) = \|\mathbf{x} - \mathbf{v}\|_2$ et $k = 5$. Visualisez les règles de classification obtenues en utilisant la fonction `frontiere_new`. (Souvent, les autres choix de distance peuvent être utiles, par exemple la distance de Mahalanobis.)
- 6) Pour les observations d'indice pair du jeu de données #2, faites varier le nombre k de voisins pris en compte : $k = 1, 2, \dots, n$. Que devient la méthode dans le cas extrême où $k = 1$? $k = n$? Afficher ces cas sur les données étudiées en utilisant la fonction `frontiere_new` et présentez les dans une forme facilement lisible. Dans quels cas la frontière est-elle complexe ? simple ?

- 7) Une variante possible très utilisée consiste à pondérer les poids du j ème voisin selon $e^{-d_j^2/h}$ (h contrôlant le niveau de pondération) : cela revient à remplacer l'Équation (2) par :

$$\hat{f}_k(\mathbf{x}) \in \arg \max_{y \in \mathcal{Y}} \left(\sum_{j=1}^k \exp(-d_j^2/h) \mathbb{1}_{\{y_{r_j}=y\}} \right). \quad (3)$$

Implémentez cette variante dans scikit-learn en passant le paramètre `weights` au constructeur de `KNeighborsClassifier`. (Une autre possibilité consiste à pondérer les variable et non seulement les observations, on le regarde pas ici.) On pourra s'inspirer de `_weight_func` de la partie test de scikit-learn : https://github.com/scikit-learn/scikit-learn/blob/master/sklearn/neighbors/tests/test_neighbors.py Pour tester l'impact du choix de h sur les frontières de classification, visualisez les règles de classification pour $k = 7$ et $h = 10^j$ avec $j = -2, -1, 0, 1, 2$; utilisez les observations d'indice pair du jeu de données #2 pour l'entraînement du classifieur (Vous pouvez utiliser la fonction `frontiere_new`.)

- 8) Quel est le taux d'erreur sur les mêmes données d'apprentissage (*i.e.*, la proportion d'erreur faite par le classifieur) lorsque $k = 1$? et sur des données de test (les observations d'indice pair) ?
- 9) Pour le jeu de données #4, en utilisant les observations d'indice pair pour l'apprentissage et les observations d'indice impair pour le test, tracez le taux d'erreur en fonction de k pour $k = 1, 2, \dots, 50$. Vous pourrez utiliser la classe fournie `ErrorCurve`.
- 10) Tracez les différentes courbes d'erreur en fonction du paramètre $k = (1, 2, \dots, 50)$ sur le jeu de données #4, pour des nombres d'échantillons d'entraînement n prenant les valeurs 100, 200, 500 à 1000. Cette fois, tirez l'ensemble d'apprentissage et l'ensemble de test indépendamment et de même taille. Quelle est la meilleure valeur de k ? Est-ce la même pour les différents datasets ? Vous pourrez utiliser la classe fournie `ErrorCurve`. Pour $n = 1000$ visualisez les données et la règle de décision sur le même graphique. (Vous pouvez utiliser la fonction `frontiere_new`.)
- 11) A votre avis, quels sont les avantages et les inconvénients de la méthode des plus proches voisins : temps de calcul ? passage à l'échelle ? interprétabilité ?
- 12) Étudiez la base DIGITS de scikit-learn. On pourra se référer à http://scikit-learn.org/stable/_downloads/plot_digits_classification.py pour le chargement et la manipulation de la base de données. Pour de plus amples informations sur la nature de la classe 'Bunch' (une sous-classe de dictionnaire, on se reportera à la documentation sur la classe 'dict' : <http://docs.python.org/2/library/stdtypes.html#mapping-types-dict>. Décrivez la nature et le format des données (précisément), affichez un exemple. Tracez l'histogramme pour des classes. Coupez l'échantillon en deux parties de même taille et utilisez la première partie pour l'apprentissage et la deuxième pour le test. Appliquez la méthode aux données issues de la base DIGITS pour un choix de $k \geq 1$ (*e.g.*, $k = 30$) et indiquez le taux d'erreur.
- 13) Estimez la matrice de confusion $(\mathbb{P}\{Y = i, C_k(X) = j\})_{i,j}$ associée au classifieur C_k ainsi obtenu et visualisez celle-ci. Pour la manipulation de telles matrices avec scikit-learn, on pourra consulter http://scikit-learn.org/stable/auto_examples/model_selection/plot_confusion_matrix.html.
- 14) Proposez une méthode pour choisir k et mettez-la en œuvre. Vous pourrez utiliser la classe fournie `L00Curve`. En utilisant toutes les données, tracez la courbe du taux d'erreur *leave-one-out* pour $k = 1, 6, 11, 16, 21, 26, 31, 36, 41, 46, 100, 200$. Pour plus d'information sur la validation croisée (*cross-validation*) on peut consulter [HTF09, Chapitre 7.10].

- POUR ALLER PLUS LOIN -

Des détails généraux sur la méthode des k -plus proches voisins se trouvent dans [HTF09, Chapitre 13]. Pour améliorer la compréhension théorique de la méthode on peut se reporter au livre [DGL96, Chapitre 11] et les limites de la méthode quand $k = 1$ <http://certis.enpc.fr/%7Edalalyan/Download/DM1.pdf>. Enfin pour les considérations algorithmiques on pourra commencer par lire <http://scikit-learn.org/stable/modules/neighbors.html#brute-force> et les paragraphes suivants.

Le nom anglais est *Linear Discriminant Analysis* (LDA). Il est préférable de se reporter à cette dénomination en vue de trouver de l'aide en ligne pour la partie numérique. Attention toutefois à ne pas confondre avec Latent Dirichlet Allocation qui est un modèle statistique hiérarchique pour données catégorielles et qui n'a pas de lien avec l'Analyse Discriminante Linéaire.

Aspect théorique

On considère deux populations gaussiennes dans \mathbb{R}^p ayant la **même** structure de covariance. On observe des points dans le mélange de ces deux populations.

Les lois conditionnelles de X sachant $Y = +1$ (respectivement $Y = -1$) sont des gaussiennes multivariées $\mathcal{N}_p(\mu_+, \Sigma)$ (respectivement $\mathcal{N}_p(\mu_-, \Sigma)$). On notera leurs densités respectives f_+ et f_- . Les vecteurs μ_+ et μ_- sont dans \mathbb{R}^p et la matrice Σ est (symétrique) de taille $p \times p$. On note également $\pi_+ = \mathbb{P}\{Y = +1\}$.

On rappelle que la densité p -dimensionnelle de la loi $\mathcal{N}_p(\mu, \Sigma)$ est donnée par :

$$f(\mathbf{x}) = \frac{1}{(2\pi)^{p/2} \sqrt{\det(\Sigma)}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mu)^\top \Sigma^{-1} (\mathbf{x} - \mu) \right\}.$$

et que la matrice de covariance d'un vecteur aléatoire X est définie par $\Sigma = \mathbb{E}[(X - \mathbb{E}(X))(X - \mathbb{E}(X))^\top]$.

- 15) En utilisant la formule de Bayes donner la formule des probabilités a posteriori : $\mathbb{P}\{Y = +1 \mid X = \mathbf{x}\}$, $\mathbb{P}\{Y = -1 \mid X = \mathbf{x}\}$, comme fonctions de f_+ , f_- et π_+ .
- 16) Exprimer le log-ratio des deux classes :

$$\log \left(\frac{\mathbb{P}\{Y = +1 \mid X = \mathbf{x}\}}{\mathbb{P}\{Y = -1 \mid X = \mathbf{x}\}} \right)$$

en fonction de π_+ , μ_+ , μ_- et Σ .

On dispose à présent d'un échantillon de ce mélange et on suppose que π_+ , μ_+ , μ_- et Σ sont des paramètres inconnus. On suppose que l'échantillon considéré contient n observations notées $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$ et que $\sum_{i=1}^n \mathbb{1}\{y_i = +1\} = m$. On va utiliser les estimateurs sans biais $\hat{\pi}_+$, $\hat{\mu}_+$, $\hat{\mu}_-$ et $\hat{\Sigma}$ des paramètres :

$$\hat{\pi}_+ = m/n, \quad \hat{\mu}_+ = \frac{1}{m} \sum_{i=1}^n \mathbb{1}\{y_i = +1\} \mathbf{x}_i, \quad \hat{\mu}_- = \frac{1}{n-m} \sum_{i=1}^n \mathbb{1}\{y_i = -1\} \mathbf{x}_i$$

et

$$\hat{\Sigma} = \frac{1}{n-2} ((m-1)\hat{\Sigma}_+ + (n-m-1)\hat{\Sigma}_-)$$

avec

$$\hat{\Sigma}_+ = \frac{1}{m-1} \sum_{i=1}^n \mathbb{1}\{y_i = +1\} (\mathbf{x}_i - \hat{\mu}_+)(\mathbf{x}_i - \hat{\mu}_+)^\top, \quad \hat{\Sigma}_- = \frac{1}{n-m-1} \sum_{i=1}^n \mathbb{1}\{y_i = -1\} (\mathbf{x}_i - \hat{\mu}_-)(\mathbf{x}_i - \hat{\mu}_-)^\top.$$

- 17) Justifier le choix du classifieur suivant :

$$\begin{cases} 1 & \text{si } \mathbf{x}^\top \hat{\Sigma}^{-1} (\hat{\mu}_+ - \hat{\mu}_-) > \frac{1}{2} \hat{\mu}_+^\top \hat{\Sigma}^{-1} \hat{\mu}_+ - \frac{1}{2} \hat{\mu}_-^\top \hat{\Sigma}^{-1} \hat{\mu}_- + \log(1 - m/n) - \log(m/n), \\ -1 & \text{sinon.} \end{cases}$$

Mise en oeuvre

- 18) Écrivez votre propre classe `LDAClassifier` avec les méthodes d'apprentissage `fit` et de classification `predict`.
- 19) Importez le module `sklearn.discriminant_analysis` qui contient en particulier la classe `LinearDiscriminantAnalysis` qui nous servira dans la suite.

```
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
```

Vérifiez la validité des résultats obtenus avec votre méthode en les comparant à ceux de la classe `LinearDiscriminantAnalysis` de scikit-learn en utilisant un jeu de données simulé. Vous proposerez votre propre méthode de comparaison (l'échantillon d'apprentissage doit être petit et l'échantillon de validation doit être assez grand; on considère par ailleurs que m est différent de $n - m$, *i.e.*, les deux classes ne sont pas représentées par le même nombre d'échantillons observés). Indiquez les taux d'erreur de `LDAClassifier` et de `LinearDiscriminantAnalysis` et le nombre de cas où les prédictions coïncident. En utilisant votre classe `LDAClassifier`, visualisez la règle de classification. (Vous pouvez utiliser la fonction `frontiere_new`.)

- 20) En utilisant votre classe `LDAClassifier`, visualisez la règle de classification pour les jeux de données #1 et #3. (Vous pouvez utiliser la fonction `frontiere_new`.) Discutez l'efficacité de la méthode dans ces deux cas.

- RÉGRESSION LOGISTIQUE -

Méthode discriminative avec régression logistique

Importer le module `sklearn.linear_model` qui contient en particulier la classe `LogisticRegression` qui nous servira dans la suite.

```
from sklearn import linear_model
```

- 21) Appliquez la classification par régression logistique sur les données `rand_bi_gauss`. Comparer les résultats avec la LDA, notamment lorsque une classe est beaucoup plus petite que l'autre (a beaucoup moins d'observations). On parle alors de classes déséquilibrées.
- 22) À quoi correspond la variable `coef_` du modèle? `intercept_`?
- 23) Utiliser la fonction `frontiere_new` pour visualiser la frontière de décision.
- 24) Appliquez la classification par régression logistique à des données issues de la base DIGITS. Comme précédemment, coupez l'échantillon en deux parties de même taille et utilisez la première partie pour l'apprentissage et la deuxième pour tester. Indiquez le taux d'erreur.

Références

- [DGL96] L. Devroye, L. Györfi, and G. Lugosi. *A probabilistic theory of pattern recognition*, volume 31 of *Applications of Mathematics (New York)*. Springer-Verlag, New York, 1996. 4
- [HTF09] T. Hastie, R. Tibshirani, and J. Friedman. *The elements of statistical learning*. Springer Series in Statistics. Springer, New York, second edition, 2009. <http://www-stat.stanford.edu/~tibs/ElemStatLearn/>. 4