

TP1 – Prise en main de WEKA  
 {Yann.Esposito}@lif.univ-mrs.fr  
 22 septembre 2005

☞ CE TP SERA DÉDIÉ À LA PRISE EN MAIN DE WEKA QUI PERMETTRA L'INTRODUCTION DE DIVERSES NOTIONS ESSENTIELLES À L'APPRENTISSAGE AUTOMATIQUE.

## 0 Ouverture

Dans le cas où le logiciel WEKA n'est pas déjà installé, veuillez l'installer en allant sur le site <http://www.cs.waikato.ac.nz/~ml/weka/>. Téléchargez alors le fichier d'installation avec java 1.4 (fichier `weka-3-4-5jre.exe`). Exécutez-le. Une fois l'installation terminée, lancez WEKA.

☞ LORS DE SON OUVERTURE WEKA OUVRE UNE PETITE FENÊTRE AVEC QUATRE BOUTONS. CHOISISSEZ `Explorer`. LORSQUE LA FENÊTRE `Weka Explorer` S'OUVRE, VOUS POUVEZ Y VOIR 6 ONGLETS : `Preprocess`, `Classify`, `Cluster`, `Associate`, `Select attributes`, `Visualize`. LE MODE `Preprocess` PERMET DE CHARGER UN FICHIER DE DONNÉES D'EN AVOIR UN APERÇU RAPIDE ET D'APPLIQUER DES FILTRES À CES DONNÉES. LES MODES `Classify`, `Cluster` ET `Associate` PERMETTENT DE FAIRE DE L'APPRENTISSAGE PROPREMENT DIT. `Classify` PERMET DE FAIRE DE L'APPRENTISSAGE SUPERVISÉ (LES DONNÉES SONT ÉTIQUETÉES), `Cluster` PERMET D'UTILISER DES ALGORITHMES D'APPRENTISSAGE NON SUPERVISÉS ET `Associate` PERMET DE GÉNÉRER DES RÈGLES D'ASSOCIATIONS DES DONNÉES. DANS CE TP, NOUS NE NOUS INTÉRESSERONS PAS AUX MODES `Cluster` ET `Associate`. LE MODE `Select attributes` PERMET DE SÉLECTIONNER LES ATTRIBUTS DE MANIÈRE FINE AVEC DES ALGORITHMES SUPERVISÉS. LE MODE `Visualize` PERMET DE VISUALISER LES DONNÉES CE QUI PERMET D'AMÉLIORER L'IDÉE QUE L'ON SE FAIT DE LEUR ORGANISATION.

## 1 Les données

☞ WEKA UTILISE LE FORMAT DE FICHIER ARFF POUR ENREGISTRER LES DONNÉES. IL S'AGIT D'UNE LISTE D'EXEMPLE AUXQUELS SONT ASSOCIÉES DES VALEURS D'ATTRIBUTS.

Ouvrez avec un éditeur quelconque (notepad par exemple) le fichier `contactlenses.arff` qui se trouve dans le répertoire `data` du répertoire d'installation de WEKA. Tout d'abord remarquons que l'on peut insérer des commentaires dans ces fichiers en commençant la ligne avec le caractère `%`. La première ligne qui n'est pas un commentaire est `@relation contact-lenses`. Il s'agit de donner le nom `contact-lenses` au fichier de données.

Il y a ensuite une liste de lignes commençant par `@attribute`. Ces lignes définissent un *attribut*. La ligne

```
@attribute age {young, pre-presbyopic, presbyopic}
```

définit l'attribut `age` qui peut avoir les valeurs `young` ou `pre-presbyopic` ou `presbyopic`.

Une fois que les attributs ont tous été définis, la ligne `@data` indique le début des données. Il s'agit d'une liste de valeurs d'attributs ordonnés. Chaque ligne correspond à un exemple. Ainsi la ligne `young,myope,no,reduced,none` correspond à l'instance

age	spectacle-prescrip	astigmatism	tear-prod-rate	contact-lenses
young	myope	no	reduced	none

1.1. Lancez WEKA. Cliquez sur **Explorer**, une autre fenêtre s’ouvre. Sur cette dernière cliquez sur **Open File**, allez dans le répertoire **data** puis ouvrez le fichier **iris.arff**.

(a) Combien y a-t-il d’instances? d’attributs?

(b) Quels sont les noms des attributs et leur numéro?

(c) À quoi correspond la zone **Selected attribute** à droite juste sous le bouton **Apply**? Vérifiez ce qui se produit lorsque vous cliquez sur différents attributs. À quoi correspondent les valeurs; **Name**, **Type**, **Missing**, **Distinct** et **Unique**?

(d) A quoi correspond le petit graphique en bas à droite? A quoi sert l’onglet qui se trouve juste en haut à gauche du graphique? Que fait le bouton **Vizualize all**?

## 2 Filtres sur les données

☞ IL EST SOUVENT NÉCESSAIRE DE PRÉTRAITER LES DONNÉES AVANT D’UTILISER UN ALGORITHME D’APPRENTISSAGE. ON DIT SOUVENT QUE L’ON NETTOIE LES DONNÉES. CELA PERMET NOTAMMENT DE SUPPRIMER LES INSTANCES CORRESPONDANT À DES ERREURS DE MESURES OU D’ÉLIMINER DES ATTRIBUTS SUPERFLUS. CELA PEUT AUSSI PERMETTRE D’UNIFORMISER LES DONNÉES.

Il existe un grand nombre de filtres que WEKA peut appliquer aux données. Pour les choisir, il suffit de cliquer sur le bouton **Choose** dans la zone **filter** en haut de la fenêtre.

Vous avez le choix entre des filtres supervisés (qui utilisent la classe des données) et non supervisés. Une fois ce choix accompli, vous avez le choix entre des filtres d’attributs ou d’instances.

Commencez par regarder les filtres non supervisés. Essayez de comprendre ce qu’ils font avec leur description (il suffit de cliquer sur le nom du filtre à droite du bouton **Choose**). Puis appliquez les (cliquez sur **Apply**) pour vérifier qu’il font bien ce à quoi vous vous attendez. En particulier :

2.1. Que font les filtres non supervisés d’attributs **remove**, **normalize**, **standardize**, **RandomProjection**.

2.2. Que font les filtres non supervisés d’instances **Randomize**, **RemoveMisclassified**, **resample**

2.3. Que font les filtres supervisés d’attributs ou d’instances **AttributeSelection** et **StratifiedRemoveFolds**.

2.4. Expliquez quels peuvent-être les intérêts de tous ces filtres.

## 3 Visualisation des données

☞ LA VISUALISATION DE DONNÉES PEUT PERMETTRE DE SE FAIRE UNE IDÉE DE L’ORGANISATION DE CELLES-CI.

Cliquez sur l’onglet **Visualize**. Vous voyez un tableau de graphiques étiquetés par les attributs des données. Chaque graphique correspond à un graphique où chaque point représente une instance colorée en fonction de sa classe. La position du point est donnée en abscisse par l’attribut situé en haut du tableau et en ordonnée par l’attribut situé à gauche du tableau.

3.1. Vérifiez que la diagonale du tableau ne contient que des graphiques où toutes les instances sont sur la droite  $x=y$ .

3.2. Toujours pour le jeu de données **iris.arff** pouvez vous dire s’il existe un attribut unique avec lequel on va pouvoir bien classer? Si oui lequel?

3.3. Existe-t-il deux attributs avec lesquels le classement peut être très bon? Si oui lesquels?

- ☞ EN CLIQUANT SUR UNE CROIX DANS UNE FENÊTRE DE VISUALISATION, L'ENSEMBLE DES INSTANCES CORRESPONDANT À CETTE CROIX SONT DÉCRIT DE FAÇON PLUS PRÉCISE DANS UNE AUTRE FENÊTRE. EN PARTICULIER CETTE FENÊTRE PERMET D'ACCÉDER AU NUMÉRO DES INSTANCES CLIQUÉES.

### 3.4. Qu'elle est l'utilité du Jitter ?

3.5. Expliquez comment on peut simuler un graphique tridimensionnel en changeant la coloration des instances.

## 4 Classification

- ☞ LA CLASSIFICATION EST L'APPRENTISSAGE SUPERVISÉ. LES ALGORITHMES DE CLASSIFICATION PRENNENT EN ENTRÉE UN ENSEMBLE DE DONNÉES ÉTIQUETÉES ET RENVOIENT DES MODÈLES QUI PERMETTENT DE CLASSER DE NOUVELLES DONNÉES NON ÉTIQUETÉES.

- ☞ LA ZONE `Test options` PERMET DE CHOISIR DE QUELLE FAÇON L'ÉVALUATION DES PERFORMANCES DU MODÈLE APPRIS SE FERA.
  - L'OPTION `Use training set` UTILISE L'ENSEMBLE D'ENTRAÎNEMENT POUR CETTE ÉVALUATION.
  - L'OPTION `Supplied test set` VA UTILISER UN AUTRE FICHIER.
  - LORSQUE L'OPTION `Cross-validation` EST SÉLECTIONNÉE, L'ENSEMBLE D'APPRENTISSAGE EST COUPÉ EN 10 (SI `Folds` VAUT 10). L'ALGORITHME VA APPRENDRE 10 FOIS SUR 9 PARTIES ET LE MODÈLE SERA ÉVALUÉ SUR LE DIXIÈME RESTANT. LES 10 ÉVALUATIONS SONT ALORS COMBINÉES.
  - AVEC L'OPTION `Percentage split`, C'EST UN POURCENTAGE DE L'ENSEMBLE D'APPRENTISSAGE QUI SERVIRA À L'APPRENTISSAGE ET L'AUTRE À L'ÉVALUATION.

Par défaut le classifieur choisit est `ZeroR`. Il s'agit du classifieur qui choisit la classe majoritaire. Choisissez le mode d'évaluation `Use training set`. Lancez le classifieur sur le jeu de données `iris.arff`.

Dans la zone `Classifier output` vous pouvez voir les informations que le classifieur renvoie.

Après la ligne

```
=== Run information ===
```

les information sont données :

- `Scheme` donne le classifieur utilisé
- `Relation` donne le nom de la relation utilisée (les données)
- `Instances` donne le nombre d'instances
- `Attributes` liste les attributs
- `Test mode` donne le mode d'évaluation du modèle : validation croisée, ensemble de test...

Après la ligne

```
=== Classifier model (full training set) ===
```

Le modèle rendu par l'algorithme est donné sous format texte.

Après la ligne

```
=== Summary ===
```

s'affiche un résumé de l'évaluation du modèle :

Correctly Classified Instances	100	66.6667 %
Incorrectly Classified Instances	50	33.3333 %
Kappa statistic	0.5	
Mean absolute error	0.2222	
Root mean squared error	0.3333	
Relative absolute error	50 %	
Root relative squared error	70.7107 %	
Total Number of Instances	150	

Tout d'abord le nombre d'instances bien classées ainsi que le pourcentage que cela représente. Puis le nombre d'instances mal classées avec la proportion associée. Puis vient la `kappa statistic` ; il s'agit de la valeur suivante :

$$\frac{P(A) - P(E)}{1 - P(E)}$$

où  $P(A)$  correspond à la proportion d'instances bien classées et  $P(E)$  correspond à l'espérance de bien classer par chance. Plus précisément  $P(E)$  est calculée à l'aide de la matrice de confusion (voir plus loin) avec la formule suivante :

$$P(E) = \frac{\sum_i (\sum_j M_{i,j}) \times (\sum_j M_{j,i})}{N^2}$$

où  $N$  est le nombre d'instances. Un coefficient de 1 signifie un modèle parfait alors que 0 correspond à un modèle qui se trompe toujours. On peut considérer qu'une valeur inférieure à 0,7 est faible et qu'une valeur supérieure à 0,8 est élevée.

Les 4 indicateurs suivants ne sont pertinents que pour la régression. `Mean absolute error` : il s'agit de l'erreur de prédiction moyenne. `Root squared error` : si cette valeur est significativement plus élevée que la `Mean absolute error` cela signifie qu'il y a des instances pour lesquelles l'erreur de prédiction est significativement plus grande que l'erreur de prédiction moyenne.

`Relative absolute error` et `Root relative squared error` correspondent aux deux même valeur que précédemment mais comparativement au classifieur de classe majoritaire.

Vient ensuite la matrice de confusion. La case à la ligne  $i$  et à la colonne  $j$  indique le nombre d'élément de classe  $i$  prédite comme faisant partie de la classe  $j$ .

**4.1.** L'algorithme OneR choisit un seul attribut et choisit de classer en fonction de celui-ci. Lancez l'algorithme OneR sur le jeu de données.

(a) Quel est le taux de réussite de l'algorithme ? L'attribut choisit est-il le même que celui que vous aviez choisi avec la visualisation ?

**4.2.** Les classifieurs sont classés en 7 classes :

- bayes – comporte notamment naïve Bayes et les réseaux Bayesiens ;
- fonctions – comporte les réseaux de neurones, les régressions linéaires... ;
- lazy – comporte IB1 (le plus proche voisin) et IBk (les k plus proches voisins) ;
- meta – comporte des algorithmes comme le Boosting notamment AdaBoost ;
- misc – Divers algorithmes exotiques ;
- trees – contient C4.5 sous le nom J48 ;
- rules – des règles d'apprentissages comme le choix de la classe majoritaire ou apprendre en utilisant un unique attribut.

**4.3.** Nous allons nous intéresser aux algorithmes **ZeroR**, **OneR**, **Naïve Bayes**, **IB1**, **IBk** et **J48**. Essayez chacun de ces algorithmes en explorant les paramètres modifiables.

**4.4.** Évaluez l'impact des différentes méthodes d'évaluations sur chaque algorithme.

- utilisez l'ensemble d'apprentissage comme ensemble test ;
- utilisez un ensemble donné ;
- utilisez la validation croisée ;
- utilisez une découpe de l'échantillon d'apprentissage.

(a) Expliquez les problèmes liés à l'utilisation de l'ensemble d'apprentissage pour l'évaluation du modèle.

(b) Expliquez pourquoi il vaut parfois mieux utiliser la découpe de l'ensemble d'apprentissage que de la validation croisée.

(c) Expliquez l'intérêt de faire de la validation croisée où le nombre de "Folds" est égal au nombre d'instances ?

☞ LORSQUE DES CLASSIFICATIONS SONT EFFECTUÉES, DANS LA ZONE **Result list**, EN CLIQUANT AVEC LE BOUTON DROIT, UNE LISTE D'OPTIONS APPARAÎSENT. EN PARTICULIER IL Y A L'OPTION **visualize classifier errors** QUI PERMET D'ENTER EN MODE VISUALISATION ET QUI ENCADRE LES INSTANCES ERRONÉES.

UNE AUTRE FAÇON DE CONSTATER CE QUI S'EST PRODUIT EST DE CLIQUER SUR LE BOUTON **More options...** UNE FENÊTRE D'OPTION S'OUVRE DANS LAQUELLE ON PEUT CHOISIR L'OPTION **Output predictions**

**4.5.** Est-ce que les instances qui posent problèmes sont souvent les mêmes ? Par exemple entre **j48** et **NaiveBayes**. Préférez le mode **percentage split** pour faire cette vérification.

## 5 Étude approfondie des filtres de prétraitement

**5.6.** Dans un premier temps nous allons étudier l'influence de certains filtres sur **Naive Bayes**. Il faut appliquer les filtres sur tous les attributs à l'exception de la classe. Pour cela sélectionnez tous les attribut en cliquant sur **All** dans la zone **Attributes**. Désélectionnez ensuite **class**. De plus toutes les évaluations se feront avec **Percentage split** à 66%.

(a) Quelle est l'influence du filtre **Normalize** ? Expliquez.

(b) Quelle est l'influence du filtre **Standardize** ? Expliquez.

(c) Quelle est l'influence de **RandomProjection** en utilisant 3 dimensions ? En utilisant 1 dimension (valeur de **seed 42**) ? Que vous inspirent ces résultats ?

## 6 choix d'attributs

☞ IL FAUT SOUVENT LIMITER LE NOMBRE D'ATTRIBUTS POUR AMÉLIORER LES RÉSULTATS D'APPRENTISSAGE AINSI QUE DU TEMPS DE CALCUL. LA QUESTION VA ÊTRE DE SAVOIR QUELLES INFORMATIONS DOIVENT-ÊTRE ÉCARTÉES.

Allez dans **Select attributes** et choisissez comme évaluateur d'attribut **CfsSubsetEval** et comme méthode de recherche **BestFirst**. Cliquez sur **start**.

**6.7.** Interprétez ce qu'il y a dans la zone **Attribute selection output**. Ce résultat correspond-t-il au choix que vous aviez fait lors de la visualisation des données ?

**6.8.** Apprenez en utilisant seulement ces deux attributs. La qualité de l'apprentissage est-elle meilleure ou moins bonne ? Expliquez pourquoi, en particulier pour J48.

## 7 Recherche du meilleur modèle

☞ IL VA ÊTRE QUESTION DE TROUVER LE MEILLEUR MODÈLE D'APPRENTISSAGE.

**7.9.** Téléchargez le fichier `learn.arff` à l'adresse <http://www.lif.univ-mrs.fr/~esposito/pub/learn.arff>.

Le but du jeu va être de trouver le meilleur modèle possible sachant qu'il existe un autre jeu de données test qui servira une fois vos choix de modèles accomplis.

Pensez à utiliser le choix d'attributs, des filtres, faites des visualisations. Utilisez plusieurs classificateurs. Certains d'entre eux peuvent être inutilisables. Pensez à trouver des moyens pour les utiliser malgré tout.

En fait, tous les moyens sont bons pour trouver la solution. Lorsque vous estimerez avoir trouvé le meilleur modèle, envoyez-le moi par email (en précisant le titre WEKA). Je testerai ce modèle avec un fichier test. Bonne chance.

