

TP3 – Sélection de modèle
{Yann.Esposito}@lif.univ-mrs.fr
7 octobre 2005

☞ IL VA ÊTRE QUESTION DE COMPARER LES DIFFÉRENTES MÉTHODES DE SÉLECTION DE MODÈLES AFIN D'ÉTABLIR UN PROTOCOLE GÉNÉRAL.

1 Importation des données (que du plaisir)

Allez sur le site : <http://www.lif.univ-mrs.fr/~esposito/> Dans la section enseignement, téléchargez tous les fichiers nécessaires au TP n°3 :

- wave-train.csv
- wave-test.csv

1. Lancez SAS ;
2. Dans le menu fichier choisissez importer ;
3. Choisir le format csv ;
4. Donnez le bon chemin de fichier ;
5. Choisissez Library : SASUSER et Member : nom du fichier sans tiret, par exemple Hugetrain.
6. Cliquez sur terminer.

Recommencez cette opération pour le fichier wave-test.csv.

2 Découpage

☞ LA FAÇON LA PLUS IMMÉDIATE DE CHOISIR UN MODÈLE EST D'ÉVALUER SA QUALITÉ DIRECTEMENT AVEC SA QUALITÉ SUR UN ENSEMBLE DE TEST.

2.1. Dans tout le TP les tests de qualité se feront par rapport au jeu de données test.csv. Dans un premier temps créez les nœuds **Input data source** correspondant à chacun des fichiers de données. N'oubliez de sélectionner la variable **class** et de lui attribuer le rôle de cible ainsi que le "measurment" interval. De plus attribuez le rôle **TEST** au nœud correspondant à l'échantillon test. Reliez le nœud de l'ensemble d'entraînement à un nœud **Sampling** puis reliez celui-ci à un nœud **Data Partition**. Dans un premier temps mettez 100% de l'échantillon en **Train**. Reliez le nœud Data partition à un nœud de régression. Connectez le nœud correspondant à l'ensemble de test à l'ensemble de régression.

(a) Exécutez la régression linéaire en changeant à chaque étape le pourcentage du **Sampling** : 5%, 10%, 25% et 100%.

ATTENTION ne modifiez pas la valeur du **seed** tout au long du TP sauf si cela vous est demandé.

Évaluez la qualité par rapport en l'ensemble de test. Par exemple en utilisant le critère **Mean square error** obtenu en cliquant sur les résultats de la régression puis en sélectionnant l'onglet **Statistics**.

(b) Maintenant séparez chacun des échantillons d'apprentissage en deux sous-échantillon : train + validation. Puis dans la fenêtre de régression. Choisissez l'onglet Selection Method et choisissez Validation Error. Y-a-t'il eu une améliorations ?

(c) modifiez l'ordre des instances afin de modifier les différents ensembles train, validation. Faites cette opération au moins 3 fois. Constatez-vous des disparités dans les résultats ?

2.2. Décrivez les avantages et les inconvénients du découpage.

Avantages :

- facile à mettre en œuvre
- temps de calcul réduit

Inconvénients :

- on utilise pas toutes les données, en nécessite donc beaucoup :
 - deux ensembles pour l'évaluation d'un modèles
 - trois pour la sélection et l'évaluation
- sensible au découpage
- réduit les données disponibles pour le modèles, mauvaise estimation des paramètres.

3 Validation croisée

3.1. toujours sur les mêmes échantillons, utilisez de la validation croisée (onglet Selection Method de la fenêtre de régression linéaire). Évaluez la sensibilité de la validation croisée par rapport à l'ordre des éléments de l'échantillon ?

3.2. Dans SAS EM, on ne peut pas modifier le nombre de blocs de la validation croisée. Mais à votre avis, est-ce un paramètre auquel sont sensible les résultats ? Justifiez votre réponse.

Bien évidemment, l'apprentissage est sensible au nombre de blocs. Exemple, avoir un seul bloc est équivalent à utiliser le découpage.

3.3. Quel sont les conséquences d'un apprentissage par validation croisée lorsque l'on apprend avec un arbre au lieu d'apprendre avec des modèles continus comme c'est le cas pour la régression linéaire ?

On peut difficilement choisir le modèle "moyen" entre tous. Il faut alors faire des choix qui peuvent s'avérer destructeurs. Par exemple, il est possible que le résultat soit aussi mauvais que celui de n'avoir fait qu'un seul des apprentissages de la validation croisée.

3.4. Donnez les avantages et les inconvénients de la validation croisée comme mode de sélection de modèle.

Avantages :

- facile à mettre en œuvre
- utilise toutes les données

Inconvénients :

- sensible au découpage :
 - choix du nombre de blocs
 - choix des blocs eux-mêmes
- temps de calculs élevé
- difficile de savoir quel modèle choisir lorsqu'il s'agit d'arbre ou plus généralement de modèles discontinus.

4 Contrôle de complexité

☞ IDÉE DE BASE ÉTUDIER LA COMBINAISON

$$\mathcal{E} + \mathcal{C}$$

OÙ \mathcal{E} EST L'ERREUR ET \mathcal{C} LA COMPLEXITÉ DU MODÈLE.

CRITÈRE DE MALLOW'S :

$$E + 2 \frac{W}{N} \sigma^2$$

où

- E DÉSIGNE L'ERREUR QUADRATIQUE MOYENNE SUR L'ENSEMBLE D'APPRENTISSAGE
- W LE NOMBRE DE PARAMÈTRES DU MODÈLE LINÉAIRE
- N LE NOMBRE DE DONNÉES ET
- σ^2 UNE ESTIMATION DE LA VARIANCE DU BRUIT

DANS SAS VOUS AVEZ ACCÈS À DEUX CRITÈRES QUI GÉNÉRALISENT LE CRITÈRE DE MALLOW'S : AIC ET SBC. POUR PLUS DE PRÉCISIONS ALLEZ DANS L'AIDE ET CHOISISSEZ EM REFERENCE. PUIS CHOISISSEZ REGRESSION, ENFIN FAITE UNE RECHERCHE SUR LE TERME AIC.

4.1. Comparez les résultats avec ceux précédemment obtenus avec les autres méthodes de sélection de modèle.

4.2. Donnez les avantages et les inconvénient de la méthode de contrôle de complexité.

Avantages :

- relativement facile à mettre en œuvre
- utilise toute les données
- temps de calcul additionnel négligeable

Inconvénients :

- AIC sélectionne des modèles trop complexes avec N grand
- comportement parfois décevant quand N est 'raisonnable'
- il faut estimer le bruit.

5 Conclusion

☞ IL EXISTES D'AUTRES MÉTHODES DE SÉLECTIONS DE MODÈLES. LE CONTRÔLE DE COMPLEXITÉ ;

CETTE MÉTHODE ÉVALUE LA VALEUR : ERREUR + COMPLEXITÉ. ASSEZ FACILE À METTRE EN ŒUVRE, MAIS SE LIMITE GÉNÉRALEMENT AUX MODÈLES LINÉAIRES. LE BOOTSTRAP ; CETTE MÉ-

THODE PERMET D'AVOIR DES BORNES SUR L'ERREUR MAIS COÛTE TRÈS CHER EN CALCULS ET EN PRATIQUE EST AUSSI PERFORMANT QUE LA VALIDATION CROISÉE. UTILISER LA DIMENSION DE

VAPNIK-CHEVONENKIS ; IL S'AGIT D'UNE DES THÉORIE LES PLUS AVANCÉES POUR L'APPRENTISSAGE AUTOMATIQUE. EN PRATIQUE POURTANT IL EST SOUVENT TRÈS DIFFICILE DE CALCULER LA VC-DIM. DE PLUS CETTE MÉTHODE DONE DES BORNES PESSIMISTES, EN GRANDE PARTIE CAR AUCUNE HYPOTHÈSE N'EST FAITE SUR LES DONNÉES.

VOUS AVEZ ÉTUDIÉ LES MÉTHODES LES PLUS UTILISABLES EN PRATIQUES. IL S'AGIT MAINTENANT DE SYNTHÉTISER CE QUE VOUS AVEZ APPRIS.

5.1. établissez un protocole général qui vous permettra de faire le bon choix parmi les différentes méthodes de sélection de modèles.

Dans le cas général, il existe d'autres façon de faire de la sélection de modèle : En pratique :

- si le temps de calcul est acceptable : validation croisée
- sinon contrôle de complexité
- sinon découpage (Train, validation, test)

*