

## Objectif

Montrer le fonctionnement des SVM (Support Vector Machine – Machines à vecteurs de support) dans TANAGRA.

Comparaison avec l'analyse discriminante linéaire, une méthode qui induit un séparateur linéaire dans l'espace de représentation.

L'implémentation est une traduction littérale en DELPHI du code source JAVA issu de la version 3-4 de WEKA (classe SMO.JAVA).

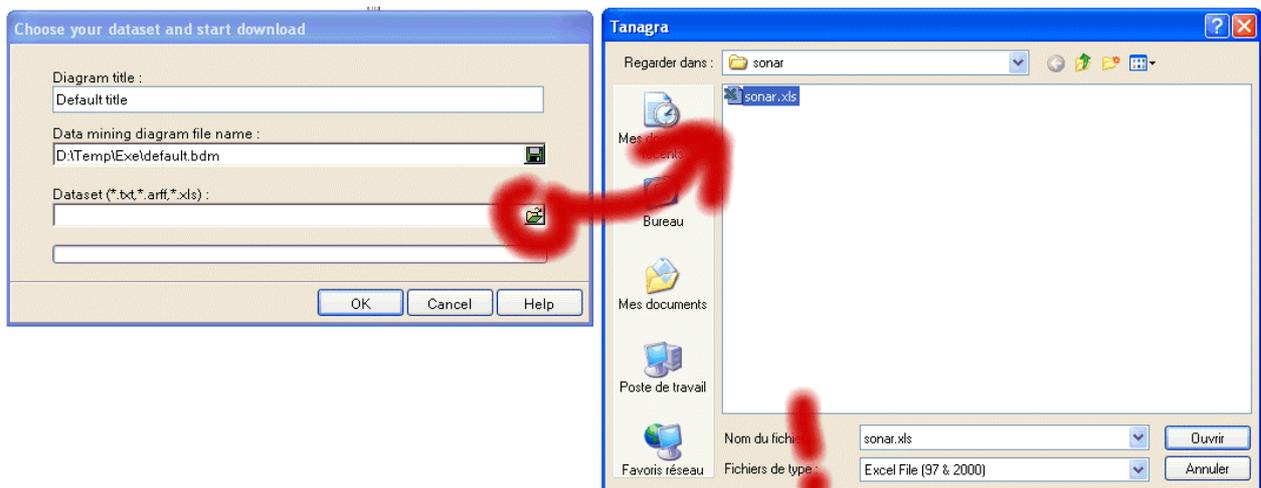
## Fichier

Le fichier SONAR : prédire la nature d'un objet (roche ou mine) à partir des informations récupérées par un détecteur de mines. L'attribut classe possède deux modalités (le seul cas traité par notre implémentation des SVM pour l'instant), il y a 60 descripteurs, tous continus.

## SVM

### Charger le fichier de données

Importez le fichier SONAR.XLS dans TANAGRA en suivant la procédure habituelle « File / New ». Il faut sélectionner le type de fichier EXCEL dans la boîte de sélection.



Les données sont importées, vérifiez que vous obtenez les mêmes spécifications que ci-dessous.

The screenshot shows a software interface with a left sidebar containing a tree view with 'Dataset (sonar.xls)'. The main area is titled 'Dataset (sonar.xls)' and contains sections for 'Parameters', 'Results', 'Download information', 'Workbook information', 'Datasource processing', and 'Dataset description'. A large blue exclamation mark is overlaid on the 'Dataset description' section.

**Dataset (sonar.xls)**

**Parameters**

Database : D:\DataMining\Databases\_for\_mining\benchmark\_datasets\sonar\sonar.xls

**Results**

**Download information**

**Workbook information**

Number of sheets	2
Selected sheet	sonar
Sheet size	209 x 61
Dataset size	209 x 61

**Datasource processing**

Computation time	62 ms
Allocated memory	118 KB

**Dataset description**

61 attribute(s)  
208 example(s)

### Analyse discriminante linéaire

Dans un premier temps, nous réalisons une analyse discriminante linéaire. Comme son nom l'indique, cette méthode induit un séparateur linéaire dans l'espace de représentation, ses qualités et ses défauts sont bien connus maintenant, nous allons nous en servir comme méthode de référence.

Le diagramme de traitements correspondant est le suivant, nous avons mis en TARGET l'attribut « CLASS » (discret et binaire), en INPUT les autres attributs (continus). Le taux d'erreur en resubstitution est de 10%, ce qui semble satisfaisant par rapport au taux d'erreur du classifieur par défaut qui est de 46% (97/208).

The screenshot shows a software interface with a left sidebar containing a tree view with 'Dataset (sonar.xls)', 'Define status 1', and 'Supervised Learning 1 (Linear discriminant analysis)'. The main area is titled 'Supervised Learning 1 (Linear discriminant analysis)' and contains sections for 'Parameters', 'Results', and 'Classifier performances'. The 'Classifier performances' section displays an error rate and a confusion matrix.

**Supervised Learning 1 (Linear discriminant analysis)**

**Parameters**

**Results**

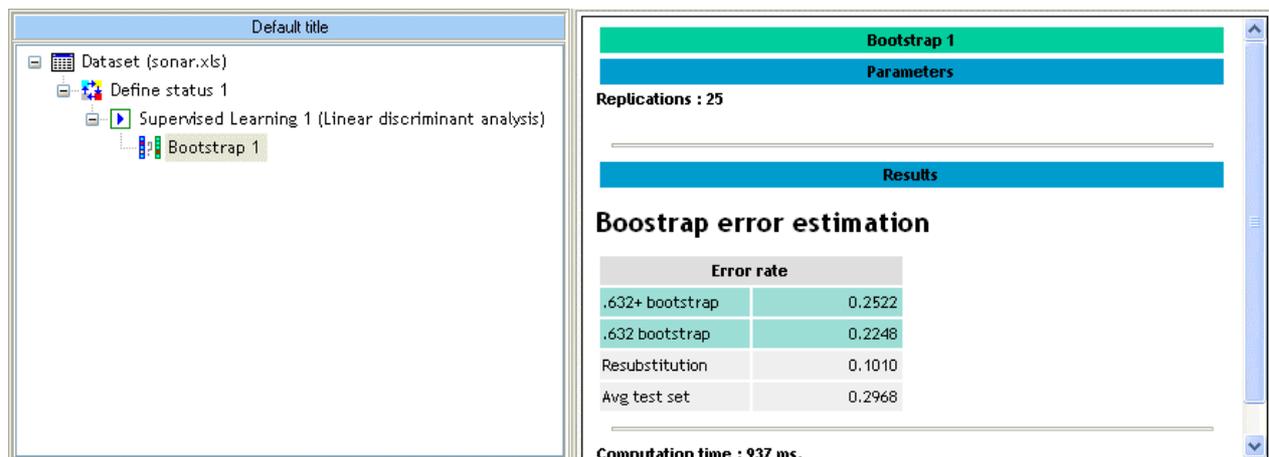
**Classifier performances**

<b>Error rate</b>	0.1010		
<b>Values prediction</b>	<b>Confusion matrix</b>		
<b>Value</b>	<b>Recall</b>	<b>1-Precision</b>	
<b>Rock</b>	0.8969	0.1122	<b>Rock</b>
<b>Mine</b>	0.9009	0.0909	<b>Mine</b>
			<b>Sum</b>
			Rock
			Mine
			Sum

Il faut être très prudent par rapport à ce résultat, nous savons que le taux d'erreur en resubstitution est un estimateur biaisé du taux d'erreur, le biais est d'autant plus fort dans

notre cas que la dimensionnalité (60 descripteurs) est élevée par rapport à la taille de l'échantillon (208 observations).

Pour calculer de manière adéquate le taux d'erreur, nous mettons en œuvre la procédure d'évaluation BOOTSTRAP. Les résultats sont édifiants, le « vrai » taux d'erreur est en réalité proche de 25%.



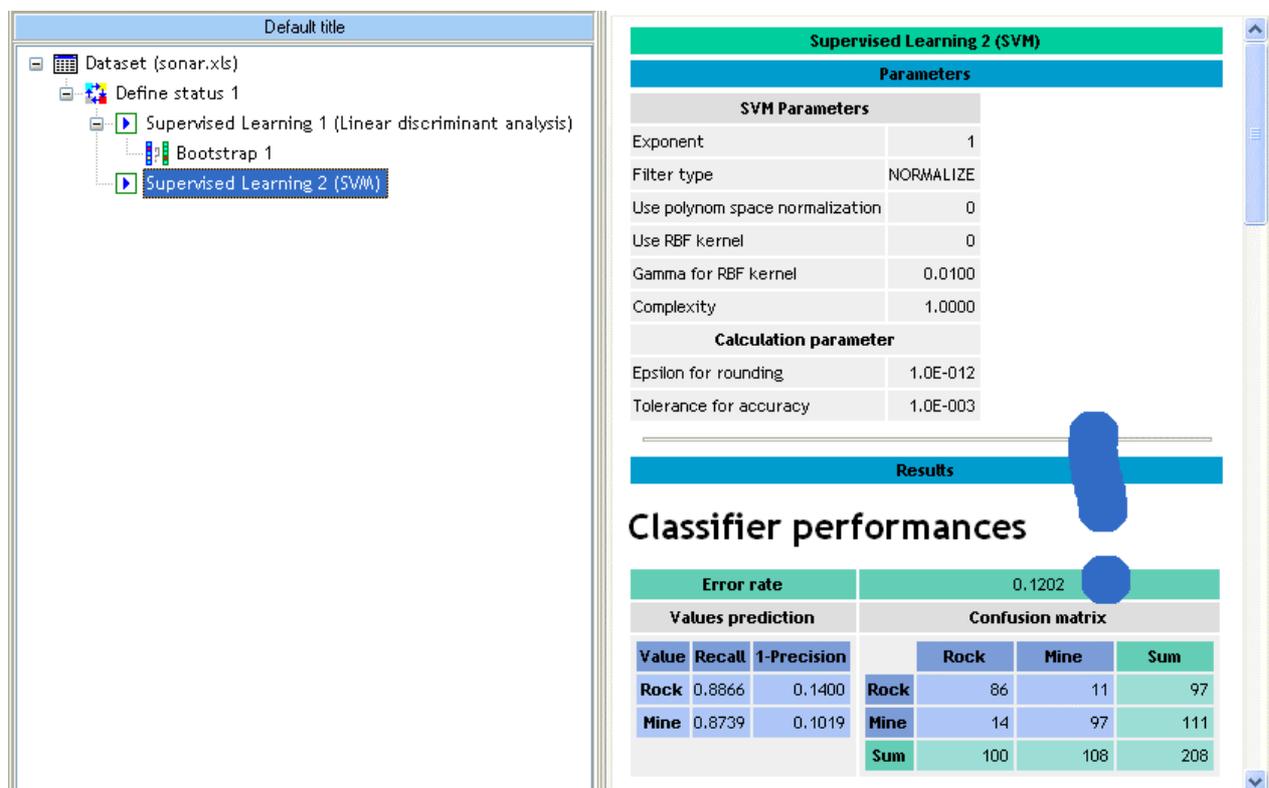
The screenshot shows a software interface with a tree view on the left and a results panel on the right. The tree view shows a dataset 'sonar.xls' with a 'Define status 1' step, followed by 'Supervised Learning 1 (Linear discriminant analysis)' and 'Bootstrap 1'. The results panel is titled 'Bootstrap 1' and shows 'Replications : 25'. Under 'Results', it displays 'Bootstrap error estimation' with the following table:

Error rate	
.632+ bootstrap	0.2522
.632 bootstrap	0.2248
Resubstitution	0.1010
Avg test set	0.2968

At the bottom, it indicates 'Computation time : 937 ms.'.

## SVM linéaire

Plaçons maintenant le composant SVM, le paramétrage par défaut permet d'obtenir un séparateur linéaire, nous pouvons comparer les résultats.



The screenshot shows the same software interface but with 'Supervised Learning 2 (SVM)' selected in the tree view. The results panel is titled 'Supervised Learning 2 (SVM)' and shows 'Parameters' for the SVM model:

SVM Parameters	
Exponent	1
Filter type	NORMALIZE
Use polynom space normalization	0
Use RBF kernel	0
Gamma for RBF kernel	0.0100
Complexity	1.0000

Below this, 'Calculation parameter' is shown:

Epsilon for rounding	1.0E-012
Tolerance for accuracy	1.0E-003

The 'Results' section is titled 'Classifier performances' and shows an 'Error rate' of 0.1202. Below this is a 'Confusion matrix' table:

Values prediction			Confusion matrix			
Value	Recall	1-Precision		Rock	Mine	Sum
Rock	0.8866	0.1400	Rock	86	11	97
Mine	0.8739	0.1019	Mine	14	97	111
			Sum	100	108	208

Les résultats semblent similaires (taux d'erreur en resubstitution de 12%). Voyons ce qu'il en est en ce qui concerne l'estimation BOOTSTRAP de l'erreur.

The screenshot shows a software interface with a project tree on the left and a results panel on the right. The project tree includes a dataset 'sonar.xls', a 'Define status 1' step, and two supervised learning tasks: 'Supervised Learning 1 (Linear discriminant analysis)' and 'Supervised Learning 2 (SVM)'. Each task has a corresponding 'Bootstrap' step. The results panel for 'Bootstrap 2' shows 'Replications : 25' and a table of error rates. The table compares the error rates of '.632+ bootstrap', '.632 bootstrap', 'Resubstitution', and 'Avg test set'. The computation time is 4328 ms, and it was created on 16/04/2005 at 09:13:48.

Error rate	
.632+ bootstrap	0.2051
.632 bootstrap	0.1953
Resubstitution	0.1202
Avg test set	0.2390

Computation time : 4328 ms.  
Created at 16/04/2005 09:13:48

Le SVM linéaire se démarque, le « vrai » taux d'erreur est autour de 20%, nettement amélioré par rapport à l'analyse discriminante.

C'est une constante remarquable de cette méthode, elle résiste très bien au sur-apprentissage même dans des dimensions très élevées (nous avons mené des tests sur les fichiers des protéines -- cf. didacticiel NIPALS – avec approximativement 7000 descripteurs pour 100 observations, la méthode tient la route alors qu'une très grande majorité des descripteurs ne sont pas pertinents).

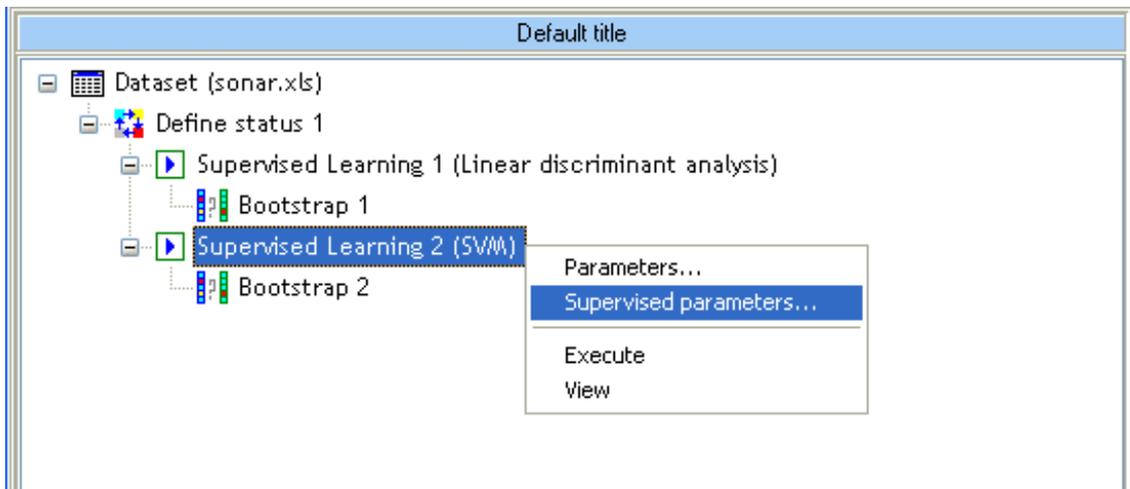
Pourquoi ? Par rapport à l'analyse discriminante, le biais de représentation est le même, nous cherchons un séparateur linéaire ; l'explication vient du biais de préférence qui est très contraignant dans les SVM, produisant ainsi un classifieur remarquablement stable, la fameuse « malédiction de la dimensionnalité » a du mal à frapper.

## SVM avec noyau polynomial de degré 2

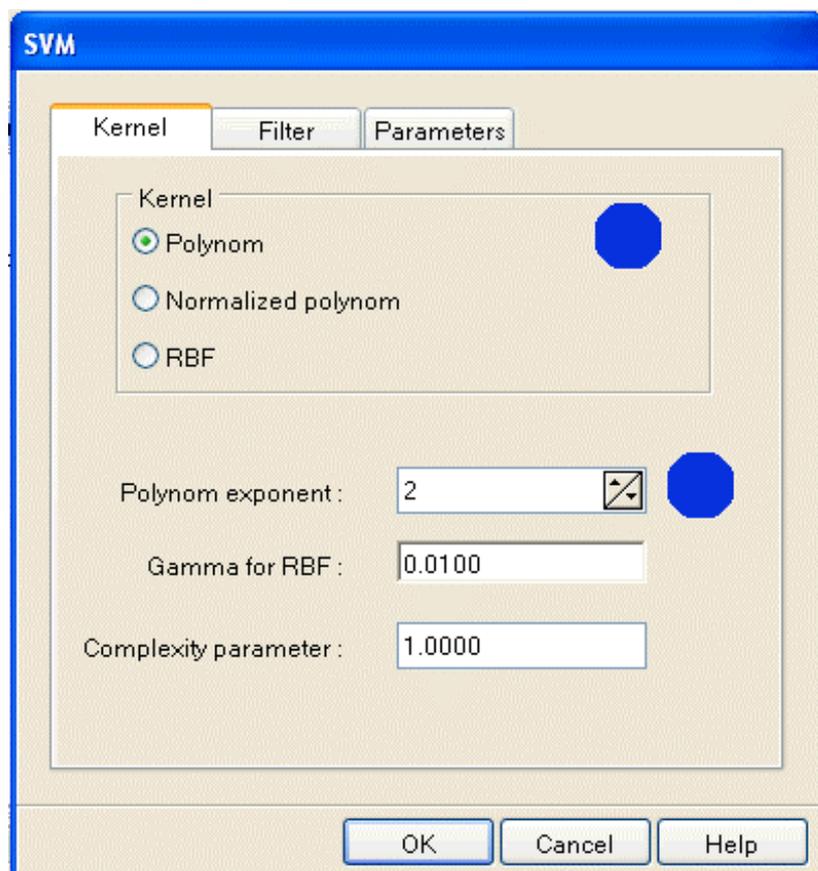
L'autre avantage des SVM est qu'il est possible de se projeter dans un espace de plus grande dimension : nous savons qu'un séparateur linéaire dans un espace de combinaisons des variables bien choisies permet d'induire un classifieur non-linéaire dans l'espace initial : **le secret des SVM réside dans l'utilisation des fameuses fonctions noyau qui permettent cette projection sans avoir à générer explicitement ces combinaisons de variables.**

Dans ce qui suit, nous utilisons un noyau polynomial de degré 2, ce qui revient à « construire » un espace composé du carré des variables et des produits croisés entre les variables.

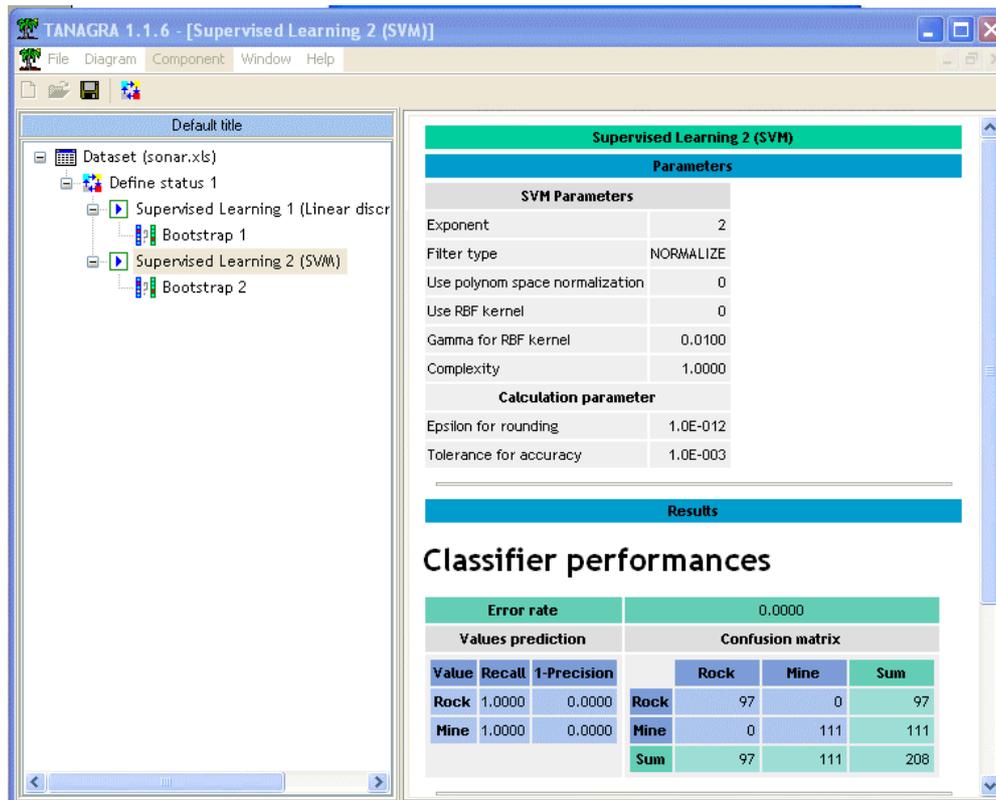
Le paramétrage de la méthode est réalisé via un clic droit sur le composant SVM



Puis la modification du degré du noyau polynomial.



Le taux d'erreur en resubstitution (0%) est très trompeur, dans un tel espace il est usuel qu'en apprentissage nous trouvons une courbe de séparation parfaite.



Le taux d'erreur du SVM à noyau polynomial de degré 2 estimé à l'aide du BOOTSTRAP est d'à peu près 15% sur le fichier SONAR (par comparaison, l'analyse discriminante propose un taux d'erreur de 25% et un SVM linéaire 20%).

