

Présentation du didacticiel

Dans ce didacticiel, nous montrons les étapes à suivre pour effectuer une analyse de type « régression ».

Le fichier que nous utilisons contient les caractéristiques de moteurs d'automobiles, et leur consommation. C'est la liaison entre cette donnée « consommation » et les autres variables du fichier que nous allons modéliser par régression.

Vous apprendrez à utiliser les opérateurs suivants :

Onglet	Opérateur	Commentaire
Data visualization	View dataset	Visualisation du fichier de données
Feature selection	Define status	Précision des variables à utiliser
Regression	Multiple linear regression	

Charger les données dans Tanagra

➤ Ouvrir un diagramme existant

1 – Choisissez *File/Open...* dans le menu de Tanagra.

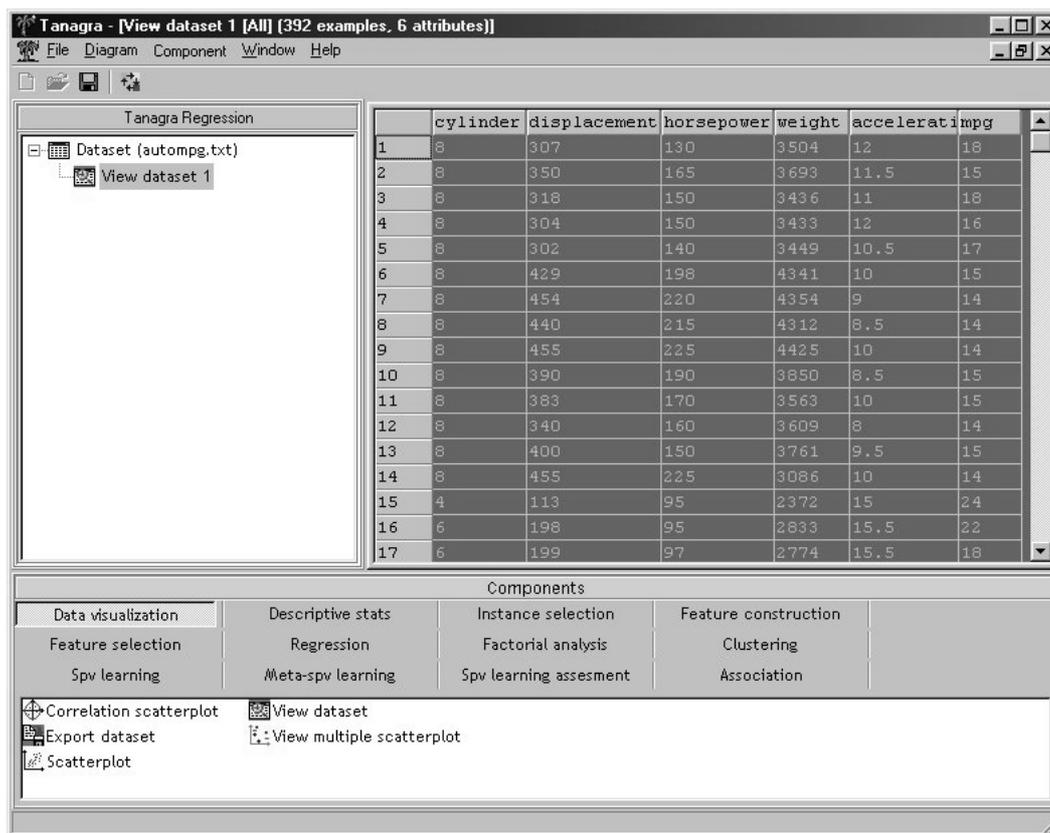
2 – Allez chercher le fichier « autompg.bdm » situé dans le sous-répertoire « Dataset » de Tanagra.

➤ Ajouter un opérateur au diagramme pour visualiser les données

1 – Ajoutez un opérateur **View dataset** au diagramme. Pour cela, cliquez sur l'onglet DATA VISUALIZATION de la palette des opérateurs.

Positionnez la souris sur l'opérateur **View dataset** et, en maintenant le bouton gauche de la souris enfoncé, amenez-le sur le diagramme. Relâchez quand vous êtes au-dessus du nœud "Dataset" (il doit apparaître sélectionné comme ci-dessous).

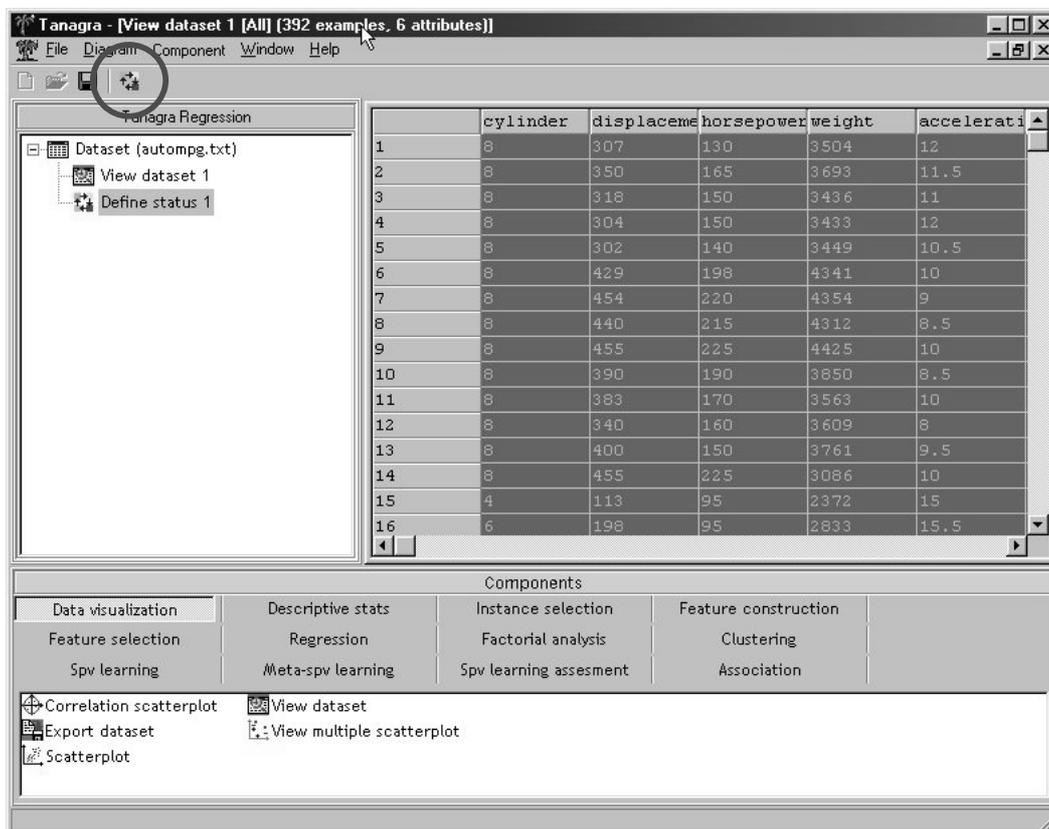
2 – Cliquez ensuite sur le nœud "View dataset" pour le sélectionner (s'il ne l'est pas déjà), et faites apparaître son menu contextuel par clic droit : choisissez la commande *View*. Les données apparaissent dans le cadre de droite.



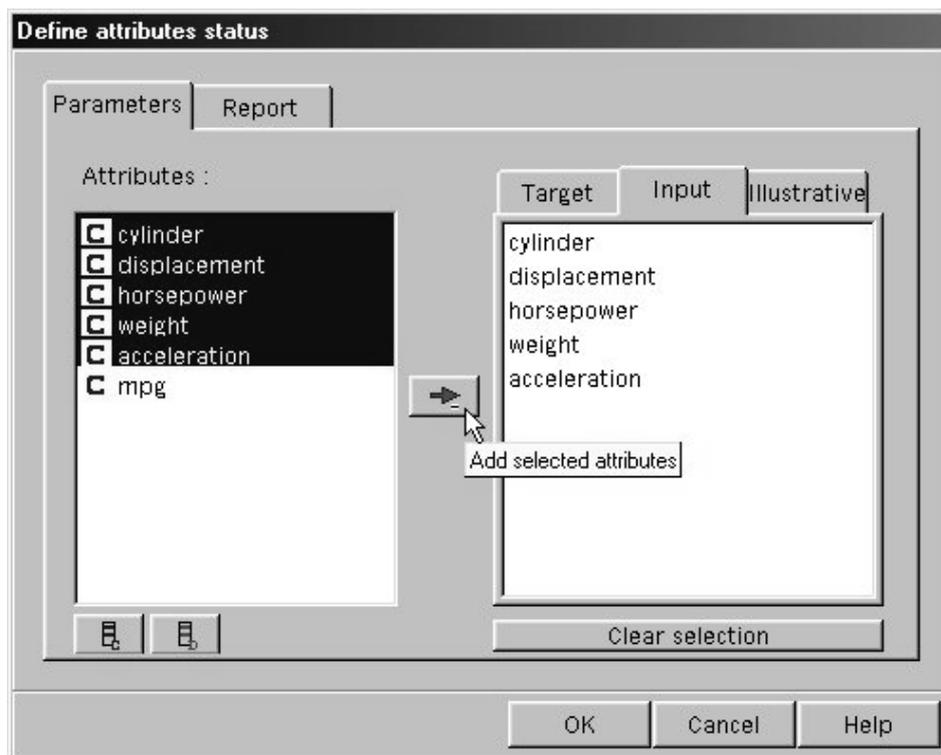
Les véhicules, au nombre de 392, sont caractérisés par 6 variables : nombre de cylindres, déplacement, puissance, poids, accélération et consommation (miles per gallon).

Définir les paramètres de la régression

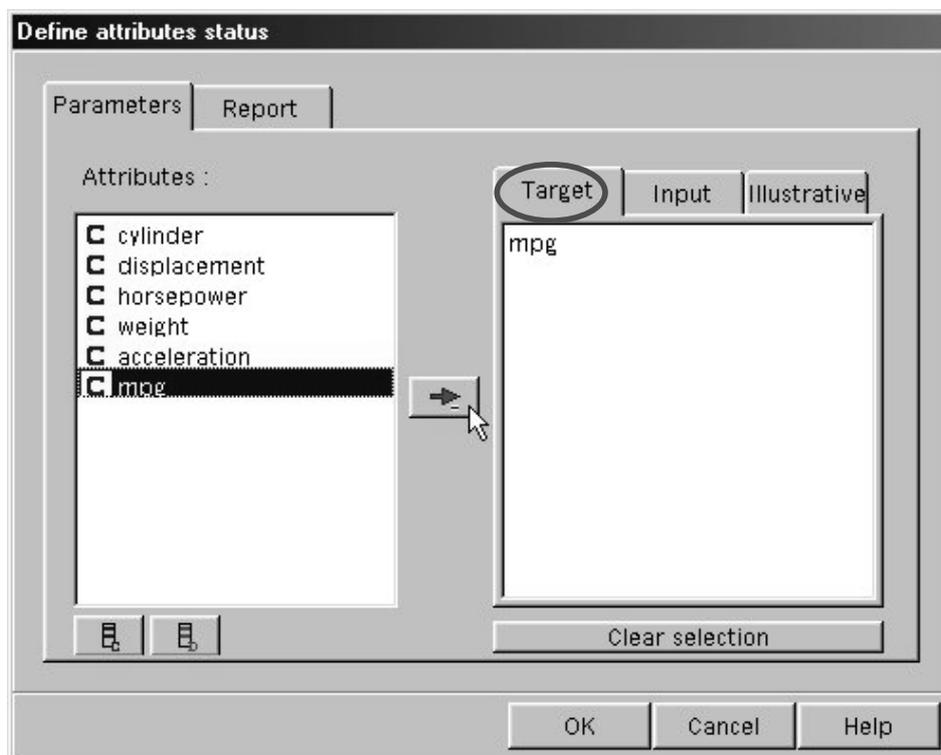
1 – Placez-vous sur le nœud « Dataset » et ajoutez un opérateur **Define Status** en cliquant sur son icône dans la barre des raccourcis. La fenêtre de dialogue permettant de définir le statut des variables apparaît automatiquement.



2 – Assurez-vous que c'est l'onglet « Input » qui est actif. Sélectionnez les cinq premières variables continues de la liste, et cliquez enfin sur le bouton flèche pour les passer dans la liste des Input.

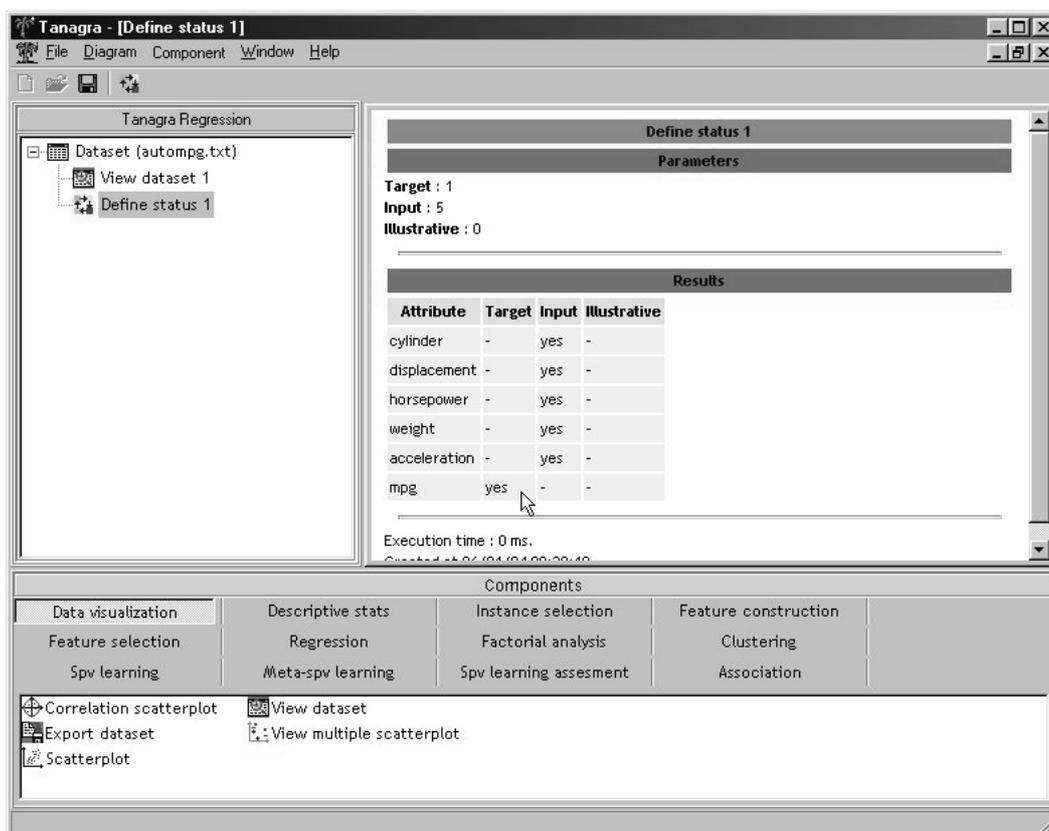


3 – Toujours en restant dans la fenêtre de dialogue, activez l'onglet Target. Cliquez sur la variable « mpg » pour la sélectionner, puis sur le bouton flèche.



4 – Vous venez de définir la variable à prédire (« mpg » = Target) et les variables explicatives (les autres = Input). Appuyez sur OK pour valider et fermer cette fenêtre.

5 – Double-cliquez sur l'opérateur **Define status** que vous venez d'ajouter, vous pouvez ainsi vérifier le statut de vos variables dans le cadre de droite.



Effectuer la régression

1 – Ajoutez un opérateur **Multiple linear regression** (onglet REGRESSION) au diagramme, sous le nœud «Define status 1».

2 – Comme il n'y a rien à paramétrer de plus au niveau de cet opérateur, vous pouvez activer directement la commande *View* du menu contextuel. Les résultats de la régression apparaissent dans le cadre de droite.

Multiple linear regression 1

Parameters

Regression parameters

Include intercept yes

Results

Global results

Predicted attribute	mpg
Number of observations	392
Residual error	4.24705542
Coefficient of determination	0.70769263
Adjusted coef of determination	0.70390627
Variance ratio (explained/resid.)	F= 186.90556 with Prob(>F) = 0.0000

Coefficients

Parameter	Est.value	Std.dev.	t Student	Prob(> t)
constant	46.2643	2.6694	17.3313	0.0000
cylinder	-0.3979	0.4105	-0.9693	0.3330
displacement	-0.0001	0.0091	-0.0092	0.9927
horsepower	-0.0453	0.0167	-2.7162	0.0069
weight	-0.0052	0.0008	-6.3515	0.0000
acceleration	-0.0291	0.1258	-0.2314	0.8171

Components: Data visualization, Descriptive stats, Instance selection, Feature construction, Feature selection, **Regression**, Factorial analysis, Clustering

On constate que la régression est d'assez bonne qualité (le coefficient de détermination est de 0.70). Les variables les plus significatives sont la puissance (horse power) et le poids (weight). Il faudrait cependant tenir compte des problèmes de colinéarité dans cet exemple.