

1. Objectif

Estimation du taux d'erreur des méthodes supervisées par des techniques de ré échantillonnage.

L'évaluation des classifieurs est une question récurrente en apprentissage supervisé. Parmi les différents indicateurs existants, la performance en prédiction calculée à l'aide du taux d'erreur (ou son complémentaire à 1, le taux de bon classement) est un critère privilégié. Du moins dans les publications scientifiques car, dans les études réelles, d'autres considérations sont au moins aussi importantes : l'évaluation des performances en intégrant les coûts de mauvaise affectation, l'interprétation des résultats, les possibilités de mise en production, etc.

Le taux d'erreur théorique est défini comme la probabilité de mal classer un individu dans la population. Bien entendu, il est impossible de le calculer directement, essentiellement parce qu'il n'est pas possible d'accéder à toute la population. Nous devons produire une estimation. Qui dit estimation dit utilisation d'un échantillon, un estimateur de bonne qualité doit être le moins biaisé possible (en moyenne, on tombe sur la bonne valeur du taux d'erreur théorique), et le plus précis possible (la variabilité autour de la vraie valeur est petite).

L'estimateur trivial est le taux d'erreur empirique que l'on appelle également taux d'erreur en resubstitution. Il s'agit de ré appliquer le modèle sur l'échantillon de données qui a servi à le construire. Tous les logiciels produisent cette estimation accompagnée d'un tableau de contingence, dite matrice de confusion, qui croise, pour l'ensemble des individus de l'échantillon, la vraie modalité prise par la variable à prédire et la modalité affectée par le modèle de classement. Le principal reproche que l'on peut adresser à l'erreur en resubstitution est qu'elle est fortement biaisée, on parle de « biais d'optimisme ». En effet, elle sous estime souvent le taux d'erreur théorique. La raison est simple, le fichier de données est à la fois « partie », il a servi à construire le modèle, et « juge », il est utilisé pour savoir si le modèle classe correctement. On montre que plus une observation pèse sur sa prédiction, plus l'optimisme sera important. L'exemple extrême est la méthode du plus proche voisin (1-ppv), le plus proche voisin d'un point est lui-même, le taux d'erreur en resubstitution est mécaniquement égal à zéro (dans un espace de description continu, c.-à-d. lorsque la probabilité que deux observations aient la même description est nulle). De manière générale, il y a un fort optimisme lorsque les techniques « collent » exagérément aux données (ex. un Perceptron avec trop de neurones dans la couche cachée, un arbre de décision trop grand) ou lorsque la dimensionnalité est trop importante au regard du nombre d'observations.

Pour se dégager de cet écueil, on conseille souvent de subdiviser l'échantillon en 2 parties : une première partie, dit fichier d'apprentissage, utilisée pour construire le modèle ; une seconde partie, dit fichier test, utilisée pour évaluer les performances du modèle. L'erreur ainsi mesurée est appelée « erreur en test ». Elle estime de manière non biaisée l'erreur théorique. Tout serait donc parfait si nous ne sommes pas confronté à un nouveau problème : quelle proportion des données devons nous consacrer à l'apprentissage ? La pratique veut que l'on réserve entre 60% et 70% pour l'apprentissage. Mais au delà de cette règle empirique, nous devons arbitrer entre deux exigences contradictoires, d'autant plus crucial que l'échantillon est de petite taille : plus nous réservons des

données pour l'apprentissage, moins l'estimation de l'erreur en test sera précise ; si nous favorisons la partie test, nous pénalisons l'apprentissage, nous retirons de l'information qui peut s'avérer déterminante pour la construction d'un modèle efficace.

Peut-on quantifier « échantillon de petite taille » si nous devons travailler sur un problème réel ? En deçà du millier d'observations, on pourrait considérer que la base est de petite taille. En réalité, il faut surtout appréhender le problème sous l'angle du rapport entre la complexité du modèle (ou le nombre de variables s'il s'agit par exemple d'un modèle linéaire) et le nombre d'observations.

Dans un contexte où les observations sont (relativement) rares, comment estimer au mieux le taux d'erreur théorique si nous souhaitons consacrer l'ensemble du fichier à la construction du modèle ?

Les techniques de ré échantillonnage permettent de répondre à cette question. Nous étudierons plus particulièrement la **validation croisée**, le **leave one out**, et le **bootstrap**. Il s'agit de répéter plusieurs fois, sous des configurations pré définies, le schéma apprentissage test. Attention, il s'agit bien d'une estimation de l'erreur du modèle construit sur l'ensemble des données. Les modèles intermédiaires, élaborés lors des apprentissages répétés, servent uniquement à l'évaluation de l'erreur. Ils ne sont pas accessibles à l'utilisateur, ils n'ont pas d'utilité intrinsèque.

Ce didacticiel vise à compléter le cours décrivant les techniques de ré échantillonnage accessible à l'adresse suivante : http://eric.univ-lyon2.fr/~ricco/cours/slides/resampling_evaluation.pdf. Nous utiliserons l'analyse discriminante linéaire (LDA) pour illustrer notre propos.

2. Données

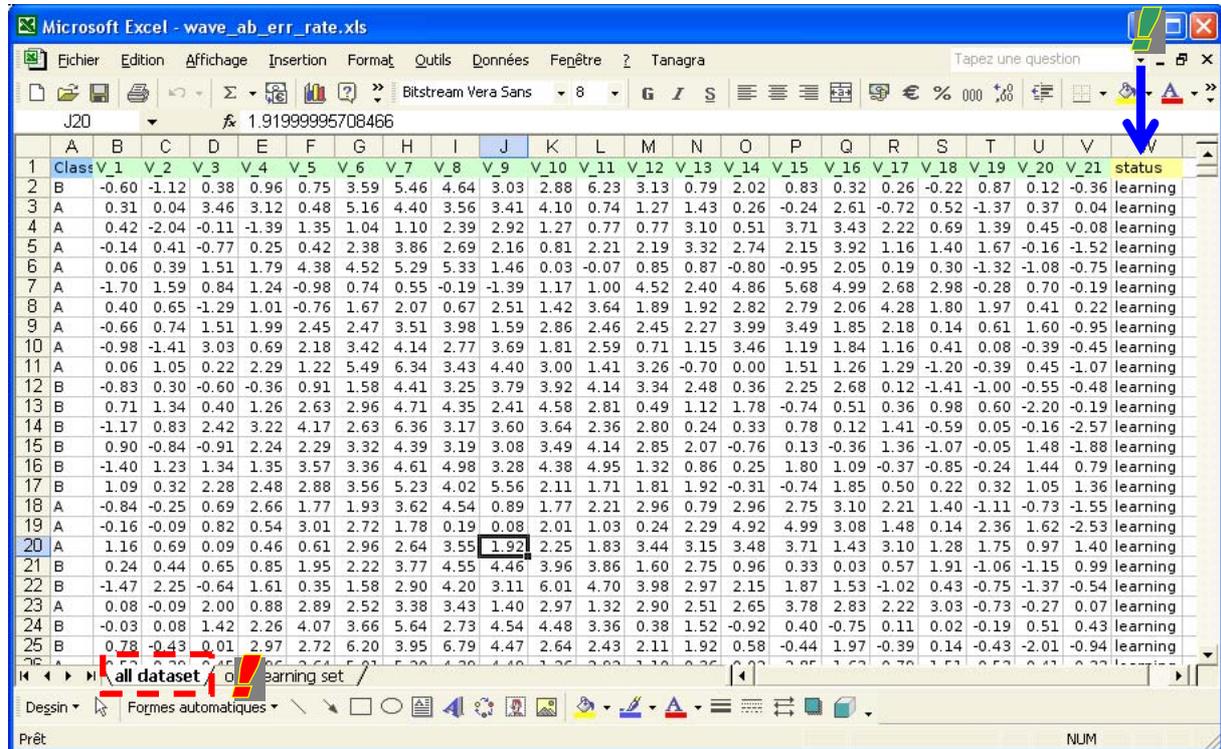
Nous utilisons les ondes de Breiman (Breiman et al., 1984) dans ce didacticiel. L'objectif est de prédire la catégorie d'un objet à partir de 21 mesures. Le fichier originel comporte 3 catégories d'objets, nous nous contenterons de la reconnaissance des 2 premières classes dans ce didacticiel.

Son principal intérêt est qu'il s'agit de données générées. Pour évaluer un modèle construit sur un échantillon de données, nous pouvons donc produire des observations en profusion, suffisamment en tous les cas pour obtenir une estimation précise de l'erreur théorique. Nous proposons donc le schéma d'expérimentation suivant :

- 500 observations constituent le fichier d'apprentissage, nous l'utiliserons pour construire le modèle de prédiction LDA(500). Nous pouvons mesurer déjà l'erreur en resubstitution (*e-resub*).
- 42500 observations constituent le fichier test, nous l'utiliserons pour obtenir une estimation non biaisée de l'erreur théorique de LDA(500) (*e-test*). On peut penser que le nombre d'observations est suffisamment élevé pour obtenir une estimation précise.
- Lors d'une étude réelle, ces données dites de test ne sont pas disponibles (surtout en telle quantité). Nous ne pouvons utiliser que les 500 observations pour construire le modèle et en évaluer les performances : nous utiliserons pour cela, tour à tour, les techniques énumérées ci-dessus. Nous comparerons les résultats (*e-lvo*, *e-cv*, *e-boot*) avec la référence que constitue *e-test*.

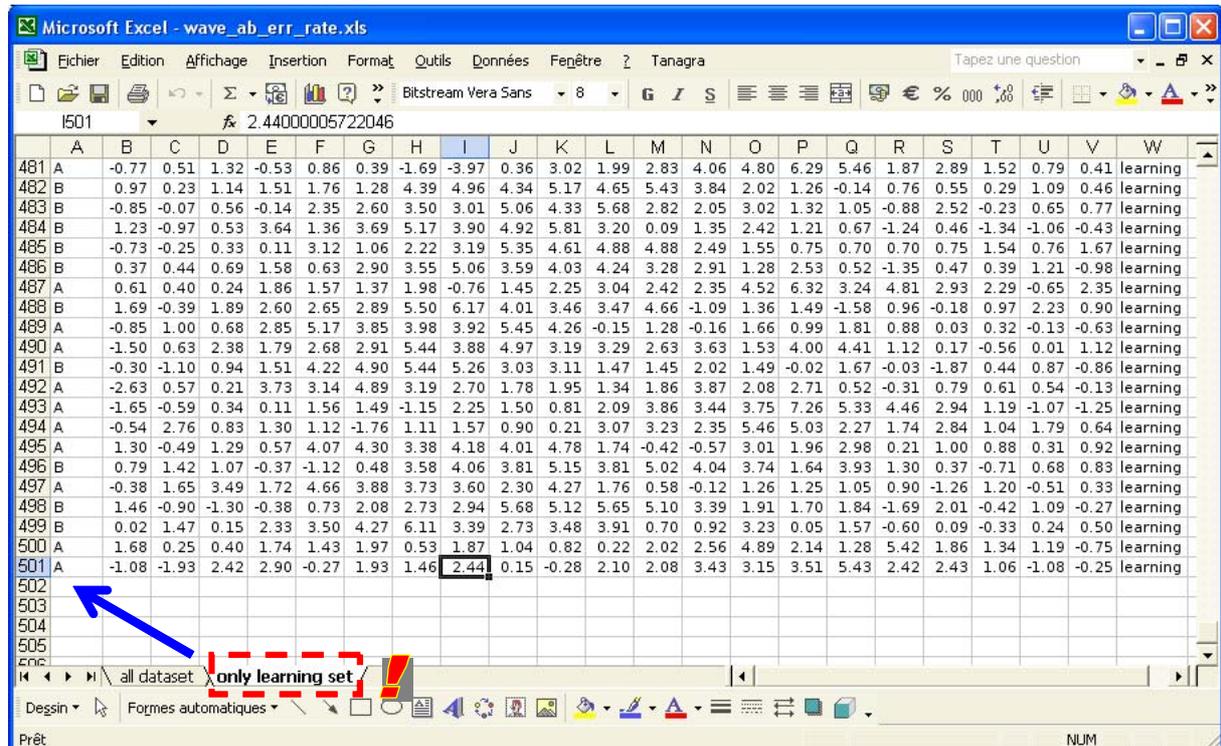
Le classeur EXCEL (wave_ab_err_rate.xls) qui accompagne ce didacticiel est subdivisé en 2 feuilles :

- « all dataset » contient 43000 observations, avec une variable supplémentaire indiquant le statut des observations (status = apprentissage ou test) ;



	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W
1	Class	V_1	V_2	V_3	V_4	V_5	V_6	V_7	V_8	V_9	V_10	V_11	V_12	V_13	V_14	V_15	V_16	V_17	V_18	V_19	V_20	V_21	status
2	B	-0.60	-1.12	0.38	0.96	0.75	3.59	5.46	4.64	3.03	2.88	6.23	3.13	0.79	2.02	0.83	0.32	0.26	-0.22	0.87	0.12	-0.36	learning
3	A	0.31	0.04	3.46	3.12	0.48	5.16	4.40	3.56	3.41	4.10	0.74	1.27	1.43	0.26	-0.24	2.61	-0.72	0.52	-1.37	0.37	0.04	learning
4	A	0.42	-2.04	-0.11	-1.39	1.35	1.04	1.10	2.39	2.92	1.27	0.77	0.77	3.10	0.51	3.71	3.43	2.22	0.69	1.39	0.45	-0.08	learning
5	A	-0.14	0.41	-0.77	0.25	0.42	2.38	3.86	2.69	2.16	0.81	2.21	2.19	3.32	2.74	2.15	3.92	1.16	1.40	1.67	-0.16	-1.52	learning
6	A	0.06	0.39	1.51	1.79	4.38	4.52	5.29	5.33	1.46	0.03	-0.07	0.85	0.87	-0.80	-0.95	2.05	0.19	0.30	-1.32	-1.08	-0.75	learning
7	A	-1.70	1.59	0.84	1.24	-0.98	0.74	0.55	-0.19	-1.39	1.17	1.00	4.52	2.40	4.86	5.68	4.99	2.68	2.98	-0.28	0.70	-0.19	learning
8	A	0.40	0.65	-1.29	1.01	-0.76	1.67	2.07	0.67	2.51	1.42	3.64	1.89	1.92	2.82	2.79	2.06	4.28	1.80	1.97	0.41	0.22	learning
9	A	-0.66	0.74	1.51	1.99	2.45	2.47	3.51	3.98	1.59	2.86	2.46	2.45	2.27	3.99	3.49	1.85	2.18	0.14	0.61	1.60	-0.95	learning
10	A	-0.98	-1.41	3.03	0.69	2.18	3.42	4.14	2.77	3.69	1.81	2.59	0.71	1.15	3.46	1.19	1.84	1.16	0.41	0.08	-0.39	-0.45	learning
11	A	0.06	1.05	0.22	2.29	1.22	5.49	6.34	3.43	4.40	3.00	1.41	3.26	-0.70	0.00	1.51	1.26	1.29	-1.20	-0.39	0.45	-1.07	learning
12	B	-0.83	0.30	-0.60	-0.36	0.91	1.58	4.41	3.25	3.79	3.92	4.14	3.34	2.48	0.36	2.25	2.68	0.12	-1.41	-1.00	-0.55	-0.48	learning
13	B	0.71	1.34	0.40	1.26	2.63	2.96	4.71	4.35	2.41	4.58	2.81	0.49	1.12	1.78	-0.74	0.51	0.36	0.98	0.60	-2.20	-0.19	learning
14	B	-1.17	0.83	2.42	3.22	4.17	2.63	6.36	3.17	3.60	3.64	2.36	2.80	0.24	0.33	0.78	0.12	1.41	-0.59	0.05	-0.16	-2.57	learning
15	B	0.90	-0.84	-0.91	2.24	2.29	3.32	4.39	3.19	3.08	3.49	4.14	2.85	2.07	-0.76	0.13	-0.36	1.36	-1.07	-0.05	1.48	-1.88	learning
16	B	-1.40	1.23	1.34	1.35	3.57	3.36	4.61	4.98	3.28	4.38	4.95	1.32	0.86	0.25	1.80	1.09	-0.37	-0.85	-0.24	1.44	0.79	learning
17	B	1.09	0.32	2.28	2.48	2.88	3.56	5.23	4.02	5.56	2.11	1.71	1.81	1.92	-0.31	-0.74	1.85	0.50	0.22	0.32	1.05	1.36	learning
18	A	-0.84	-0.25	0.69	2.66	1.77	1.93	3.62	4.54	0.89	1.77	2.21	2.96	0.79	2.96	2.75	3.10	2.21	1.40	-1.11	-0.73	-1.55	learning
19	A	-0.16	-0.09	0.82	0.54	3.01	2.72	1.78	0.19	0.08	2.01	1.03	2.24	2.29	4.92	4.99	3.08	1.48	0.14	2.36	1.62	-2.53	learning
20	A	1.16	0.69	0.09	0.46	0.61	2.96	2.64	3.55	1.92	2.25	1.83	3.44	3.15	3.48	3.71	1.43	3.10	1.28	1.75	0.97	1.40	learning
21	B	0.24	0.44	0.65	0.85	1.95	2.22	3.77	4.55	4.46	3.96	3.86	1.60	2.75	0.96	0.33	0.03	0.57	1.91	-1.06	-1.15	0.99	learning
22	B	-1.47	2.25	-0.64	1.61	0.35	1.58	2.90	4.20	3.11	6.01	4.70	3.98	2.97	2.15	1.87	1.53	-1.02	0.43	-0.75	-1.37	-0.54	learning
23	A	0.08	-0.09	2.00	0.88	2.89	2.52	3.38	3.43	1.40	2.97	1.32	2.90	2.51	2.65	3.78	2.83	2.22	3.03	-0.73	-0.27	0.07	learning
24	B	-0.03	0.08	1.42	2.26	4.07	3.66	5.64	2.73	4.54	4.48	3.36	0.38	1.52	-0.92	0.40	-0.75	0.11	0.02	-0.19	0.51	0.43	learning
25	B	0.78	-0.43	0.01	2.97	2.72	6.20	3.95	6.79	4.47	2.64	2.43	2.11	1.92	0.58	-0.44	1.97	-0.39	0.14	-0.43	-2.01	-0.94	learning

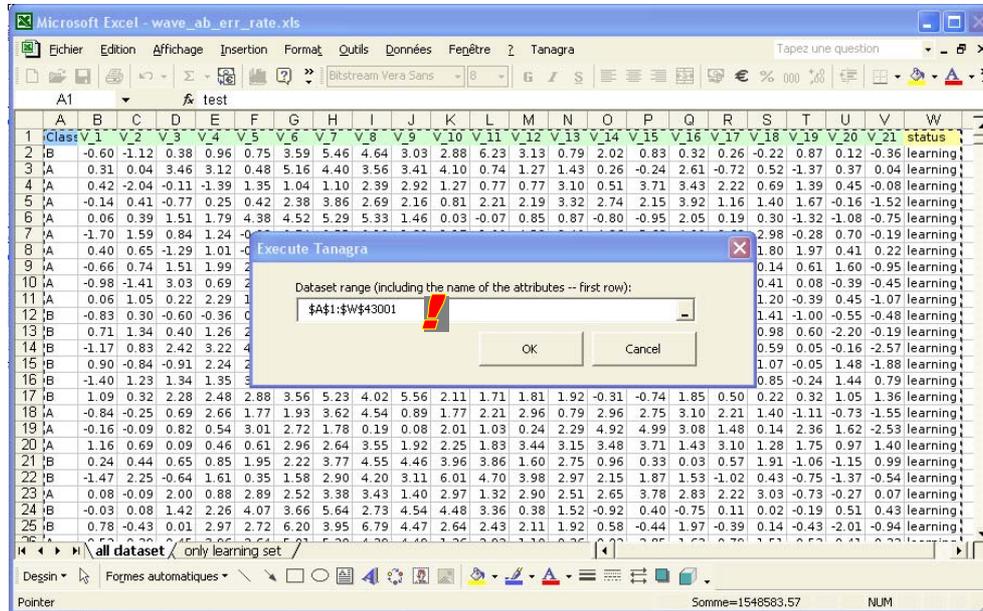
- « only learning set » ne contient que les 500 premières observations de la première feuille, correspondant à la fraction (status = apprentissage).



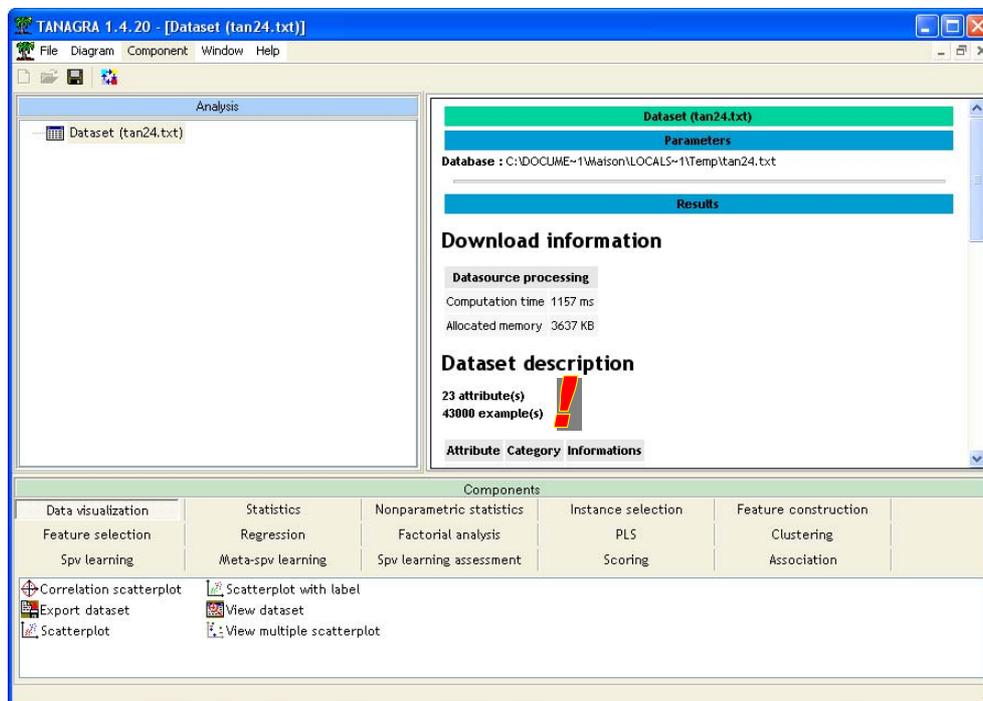
	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W
481	A	-0.77	0.51	1.32	-0.53	0.86	0.39	-1.69	-3.97	0.36	3.02	1.99	2.83	4.06	4.80	6.29	5.46	1.87	2.89	1.52	0.79	0.41	learning
482	B	0.97	0.23	1.14	1.51	1.76	1.28	4.39	4.96	4.34	5.17	4.65	5.43	3.84	2.02	1.26	-0.14	0.76	0.55	0.29	1.09	0.46	learning
483	B	-0.85	-0.07	0.56	-0.14	2.35	2.60	3.50	3.01	5.06	4.33	5.68	2.82	2.05	3.02	1.32	1.05	-0.88	2.52	-0.23	0.65	0.77	learning
484	B	1.23	-0.97	0.53	3.64	1.36	3.69	5.17	3.90	4.92	5.81	3.20	0.09	1.35	2.42	1.21	0.67	-1.24	0.46	-1.34	-1.06	-0.43	learning
485	B	-0.73	-0.25	0.33	0.11	3.12	1.06	2.22	3.19	5.35	4.61	4.88	4.88	2.49	1.55	0.75	0.70	0.70	0.75	1.54	0.76	1.67	learning
486	B	0.37	0.44	0.69	1.58	0.63	2.90	3.55	5.06	3.59	4.03	4.24	3.28	2.91	1.28	2.53	0.52	-1.35	0.47	0.39	1.21	-0.98	learning
487	A	0.61	0.40	0.24	1.86	1.57	1.37	1.98	-0.76	1.45	2.25	3.04	2.42	2.35	4.52	6.32	3.24	4.81	2.93	2.29	-0.65	2.35	learning
488	B	1.69	-0.39	1.89	2.60	2.65	2.89	5.50	6.17	4.01	3.46	3.47	4.66	-1.09	1.36	1.49	-1.58	0.96	-0.18	0.97	2.23	0.90	learning
489	A	-0.85	1.00	0.68	2.85	5.17	3.85	3.98	3.92	5.45	4.26	-0.15	1.28	-0.16	1.66	0.99	1.81	0.88	0.03	0.32	-0.13	-0.63	learning
490	A	-1.50	0.63	2.38	1.79	2.68	2.91	5.44	3.88	4.97	3.19	3.29	2.63	3.63	1.53	4.00	4.41	1.12	0.17	-0.56	0.01	1.12	learning
491	B	-0.30	-1.10	0.94	1.51	4.22	4.90	5.44	5.26	3.03	3.11	1.47	1.45	2.02	1.49	-0.02	1.67	-0.03	-1.87	0.44	0.87	-0.86	learning
492	A	-2.63	0.57	0.21	3.73	3.14	4.89	3.19	2.70	1.78	1.95	1.34	1.86	3.87	2.08	2.71	0.52	-0.31	0.79	0.61	0.54	-0.13	learning
493	A	-1.65	-0.59	0.34	0.11	1.56	1.49	-1.15	2.25	1.50	0.81	2.09	3.86	3.44	3.75	7.26	5.33	4.46	2.94	1.19	-1.07	-1.25	learning
494	A	-0.54	2.76	0.83	1.30	1.12	-1.76	1.11	1.57	0.90	0.21	3.07	3.23	2.35	5.46	5.03	2.27	1.74	2.84	1.04	1.79	0.64	learning
495	A	1.30	-0.49	1.29	0.57	4.07	4.30	3.38	4.18	4.01	4.78	1.74	-0.42	-0.57	3.01	1.96	2.98	0.21	1.00	0.88	0.31	0.92	learning
496	B	0.79	1.42	1.07	-0.37	-1.12	0.48	3.58	4.06	3.81	5.15	3.81	5.02	4.04	3.74	1.64	3.93	1.30	0.37	-0.71	0.68	0.83	learning
497	A	-0.38	1.65	3.49	1.72	4.66	3.88	3.73	3.60	2.30	4.27	1.76	0.58	-0.12	1.26	1.25	1.05	0.90	-1.26	1.20	-0.51	0.33	learning
498	B	1.46	-0.90	-1.30	-0.38	0.73	2.08	2.73	2.94	5.68	5.12	5.65	5.10	3.39	1.91	1.70	1.84	-1.69	2.01	-0.42	1.09	-0.27	learning
499	B	0.02	1.47	0.15	2.33	3.50	4.27	6.11	3.39	2.73	3.48	3.91	0.70	0.92	3.23	0.05	1.57	-0.60	0.09	-0.33	0.24	0.50	learning
500	A	1.68	0.25	0.40	1.74	1.43	1.97	0.53	1.87	1.04	0.82	0.22	2.02	2.56	4.89	2.14	1.28	5.42	1.86	1.34	1.19	-0.75	learning
501	A	-1.08	-1.93	2.42	2.90	-0.27	1.93	1.46	2.44	0.15	-0.28	2.10	2.08	3.43	3.15	3.51	5.43	2.42	2.43	1.06	-1.08	-0.25	learning

3. Apprentissage, erreur en resubstitution et erreur en test

Dans un premier temps, intéressons nous à la première feuille « all dataset ». Nous sélectionnons la plage de données et nous activons le menu TANAGRA/EXECUTE TANAGRA rendue disponible par la macro complémentaire¹ « TANAGRA.XLA ».



Après vérification des coordonnées de la plage de données, nous validons en cliquant sur le bouton OK. TANAGRA est automatiquement démarré. Les données sont chargées.

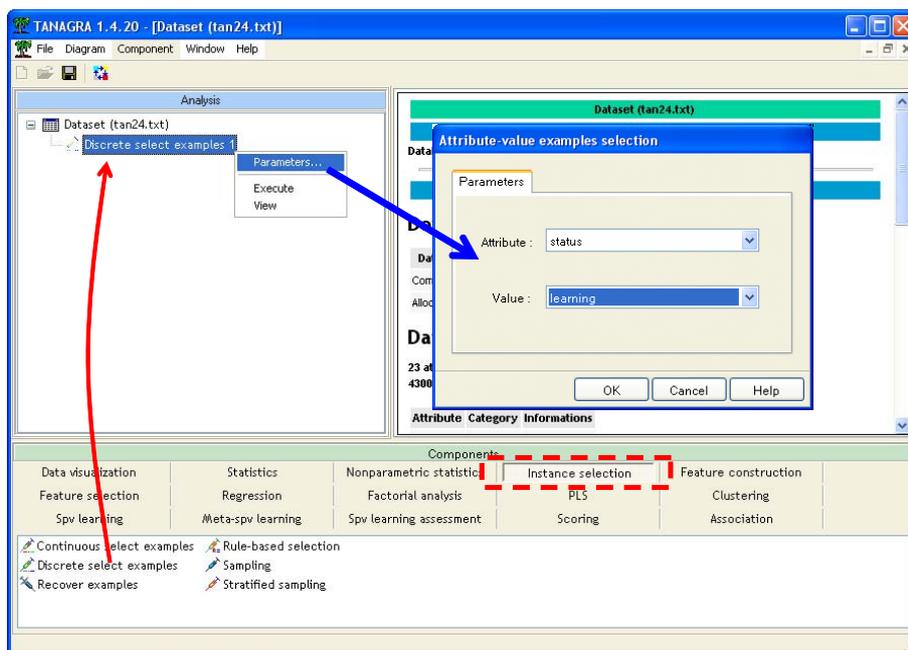


¹ Disponible depuis la version 1.4.11 de TANAGRA, voir le didacticiel http://eric.univ-lyon2.fr/~ricco/tanagra/fichiers/fr_Tanagra_Excel_AddIn.pdf pour l'installation et l'utilisation de cette macro.

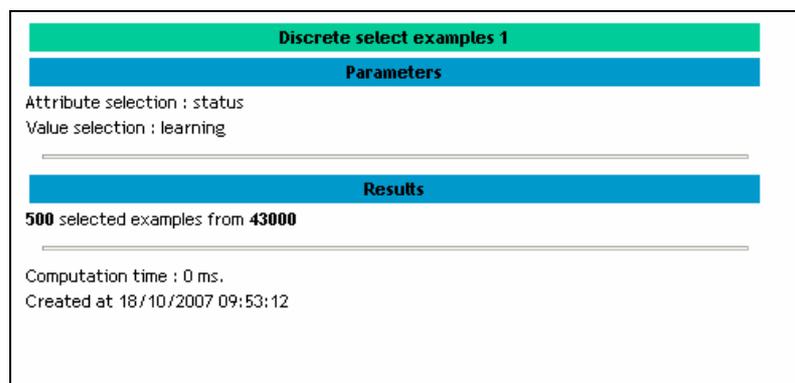
Nous devons avoir 23 variables (variable à prédire CLASS, les 21 descripteurs V1 à V21 et la variable STATUS) et 43000 observations.

Découpage apprentissage – test des données

Dans un premier temps, nous devons spécifier quelles seront les données utilisées pour la construction du modèle (données d'apprentissage) et celles qui serviront à l'évaluer (données test). Pour cela, nous utilisons le composant DISCRETE SELECT EXAMPLES (onglet INSTANCE SELECTION), il s'appuiera sur la variable supplémentaire STATUS pour subdiviser le fichier. Nous l'insérons dans le diagramme par glisser-déplacer, puis nous activons le menu contextuel PARAMETERS pour indiquer que les données à sélectionner correspondent à la modalité LEARNING de STATUS.



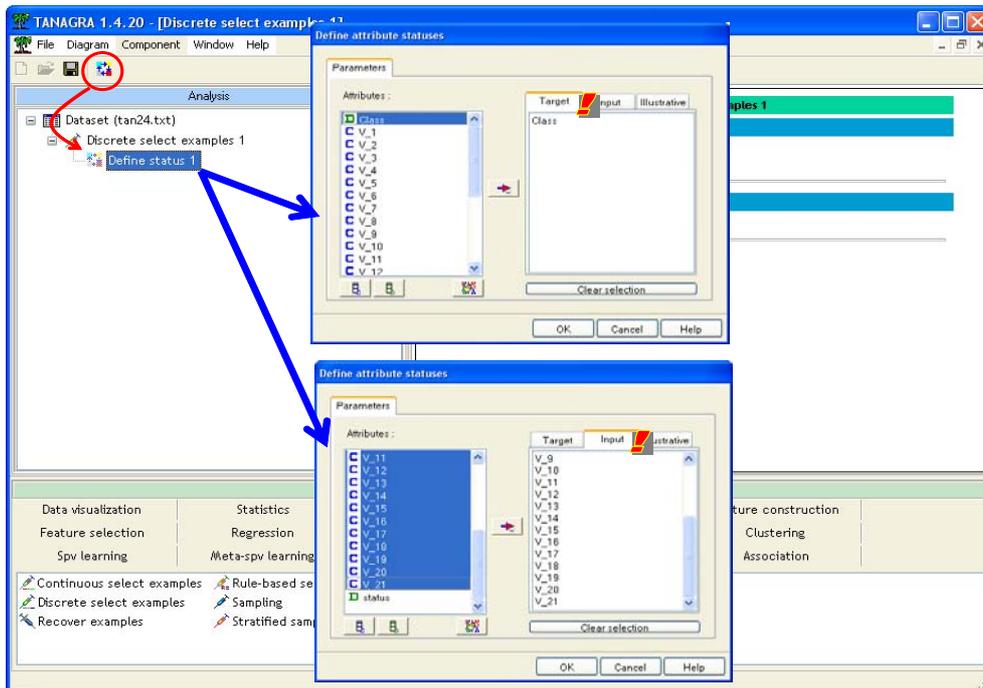
Après validation (bouton OK), nous cliquons sur le menu contextuel VIEW pour accéder aux résultats. TANAGRA nous indique bien que 500 observations sont sélectionnées pour les analyses maintenant.



Variable à prédire et variables prédictives

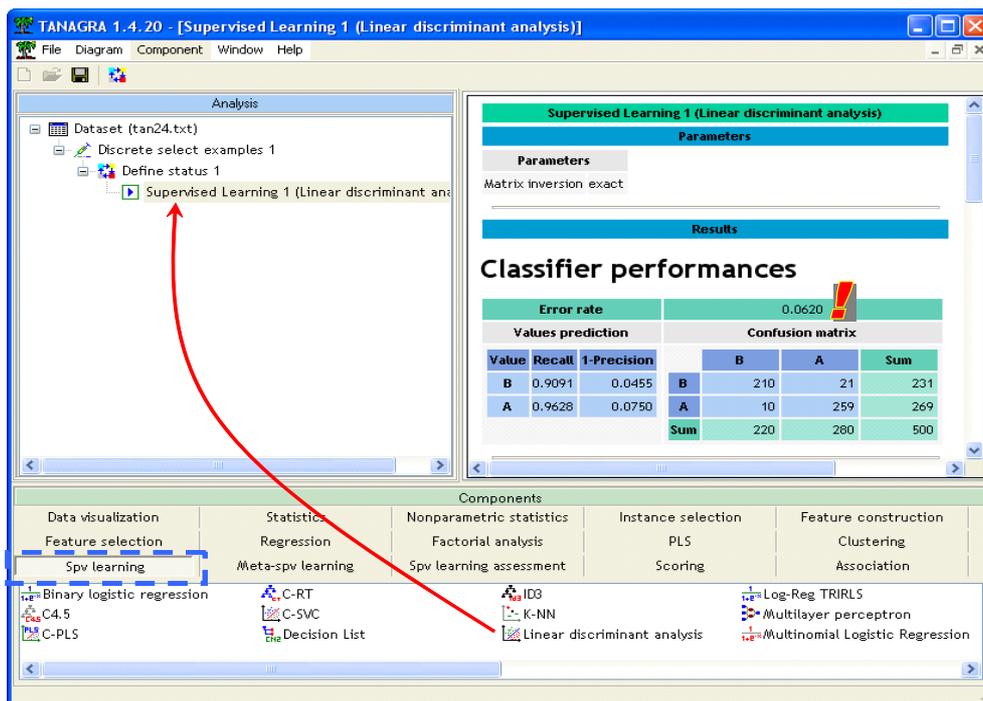
Il s'agit de préciser quelles sont les variables prédictives (INPUT) et la variable à prédire (TARGET) à l'aide du composant DEFINE STATUS. Le plus simple est de passer par le raccourci disponible dans la barre d'outils, la fenêtre de paramétrage apparaît automatiquement, nous plaçons en TARGET la

variable CLASS, en INPUT les variables V1 à V21. Bien entendu, la variable STATUS n'est pas utilisée ici.



Analyse discriminante et erreur en resubstitution

Nous pouvons construire le modèle de prédiction. Pour cela, nous plaçons dans le diagramme le composant LINEAR DISCRIMINANT ANALYSIS (onglet SPV LEARNING), toujours par glisser déplacer. Nous activons le menu VIEW pour accéder aux résultats. Qu'importe le modèle dans ce didacticiel, c'est la partie haute de la fenêtre de résultats qui nous intéresse au premier chef, il décrit la matrice de confusion et fournit l'erreur en resubstitution.

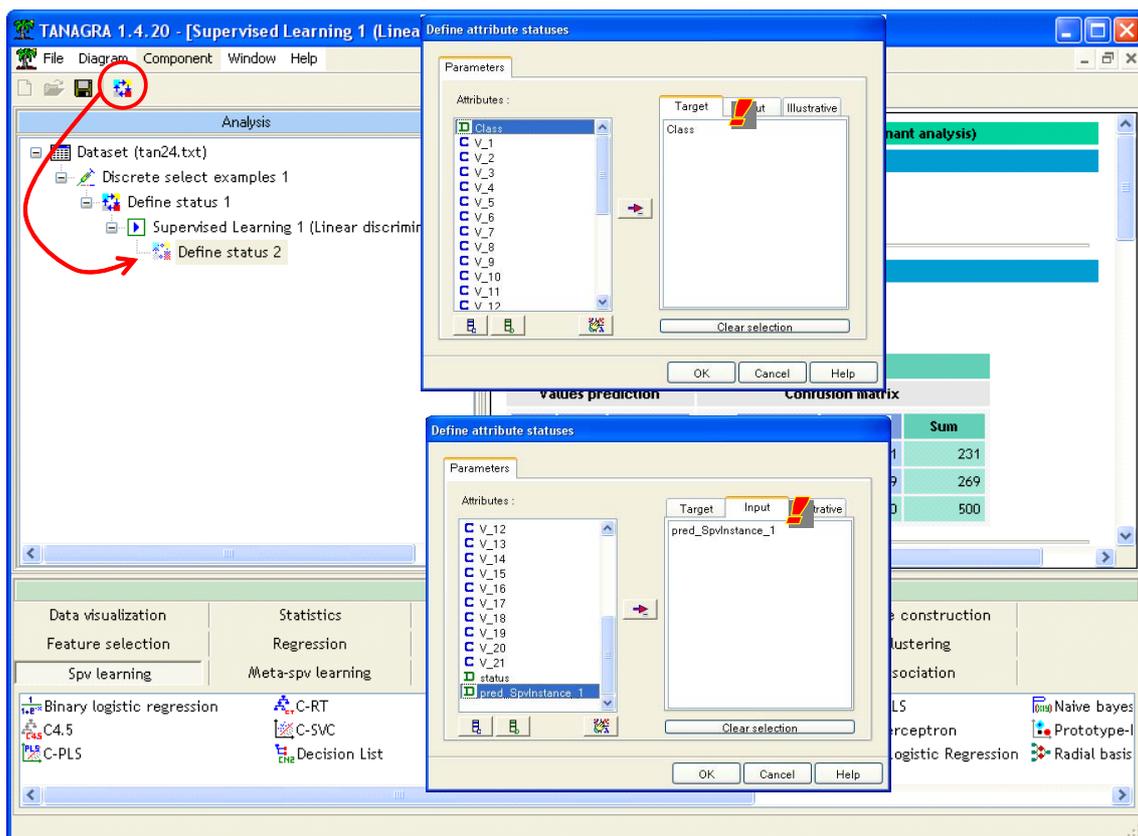


L'erreur en resubstitution est $e\text{-resub} = 6.2\%$. Si on se fie à ce chiffre, lorsque nous classons une nouvelle observation à partir du modèle fourni par la LDA, nous avons 6.2% de chances de réaliser une mauvaise affectation. En dehors de toute considération experte, nous nous bornerons à dire que sans ce modèle, la meilleure prédiction serait d'affecter systématiquement $\text{CLASS} = A$ (la modalité la plus fréquente de la variable à prédire), et dans ce cas le taux d'erreur serait $231/500 = 46.2\%$; nous lisons cela dans la marge ligne de la matrice de confusion. A priori, nous avons un modèle d'excellente qualité.

Erreur en test

Pour le vérifier, nous allons mesurer les performances du même modèle sur les 42500 observations que nous avons mis de côté. Les composants « apprentissage supervisé » génèrent automatiquement une variable supplémentaire qui correspond à la prédiction du modèle, calculé sur la totalité du fichier, sur les données d'apprentissage donc, mais aussi sur les données mises de côté. Nous allons exploiter cette propriété.

Pour cela, nous insérons de nouveau le composant DEFINE STATUS, nous plaçons en TARGET la variable CLASS, et en INPUT, la variable nouvellement générée PRED_SPV_INSTANCE_1.



Nous insérons ensuite le composant TEST (onglet SPV LEARNING ASSESSMENT) qui mesure le taux d'erreur en croisant les variables en TARGET et INPUT. Nous activons le menu VIEW pour obtenir les résultats.

Remarque 1 : Nous avons la possibilité de placer plusieurs variables en INPUT, cela permet de comparer les performances de plusieurs modèles.

Remarque 2 : Le composant est automatiquement paramétré pour que le calcul du taux d'erreur soit réalisé sur les données supplémentaires, mises de côté de la sélection des individus pour l'apprentissage. C'est ce que nous souhaitons faire ici. Nous avons néanmoins la possibilité de modifier ce paramétrage pour effectuer les calculs sur les données d'apprentissage. Dans ce cas, nous devons obtenir de nouveau l'erreur en resubstitution.

Values prediction		Confusion matrix			
Value	Recall	1-Precision	B	A	Sum
B	0.9200	0.0880	19437	1690	21127
A	0.9123	0.0798	1875	19498	21373
Sum			21312	21188	42500

Le « véritable »² taux d'erreur de notre modèle de prédiction est en réalité de 8.39%. Il y a quand même un décalage par rapport à l'erreur en resubstitution évaluée précédemment (6.2%).

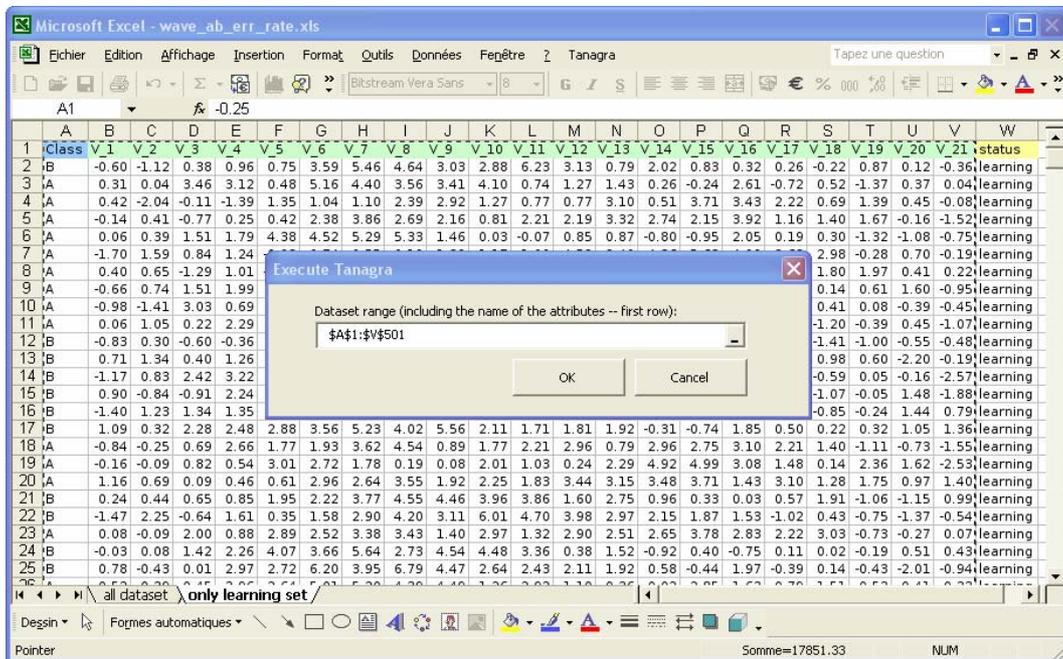
Remarque 3 : Et encore, l'analyse discriminante est une technique assez stable. Nous aurions utilisé un arbre de décision, type C4.5, qui a parfois une aptitude fâcheuse à épouser les données au plus près, les performances ne seraient pas du même ordre (e-test = 13.54%), et surtout le décalage serait plus important (e-resub = 2.2%). Nous y reviendrons à la fin de ce tutoriel.

4. Méthodes de ré échantillonnage

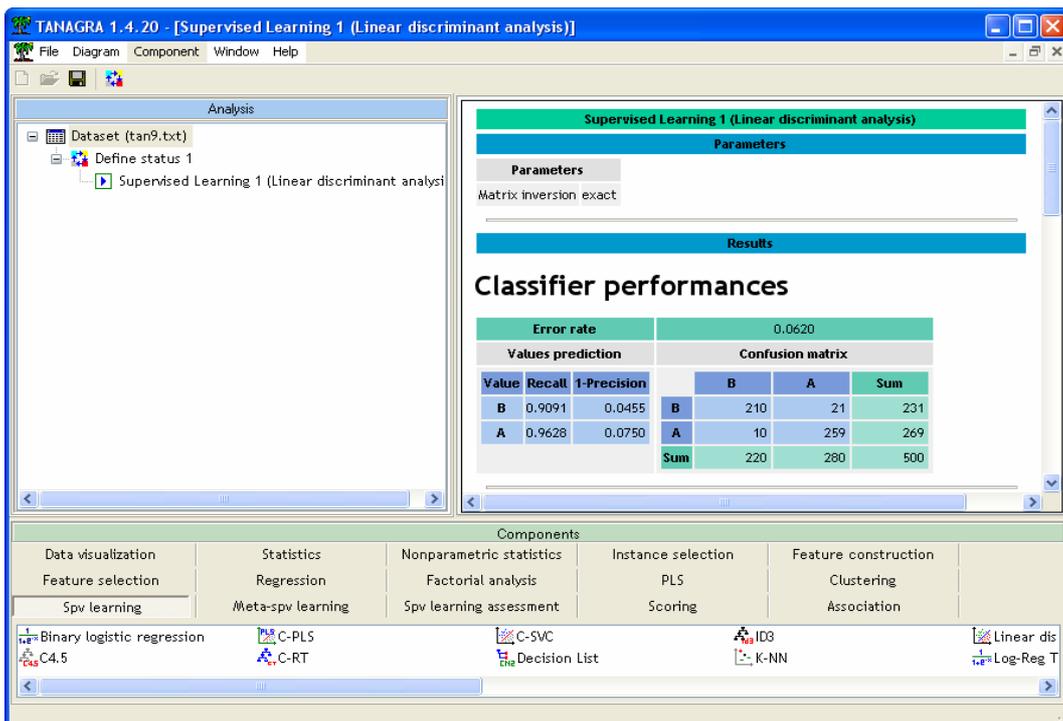
Nous pouvons clôturer cette session de travail en fermant TANAGRA. Revenons dans le classeur EXCEL, sélectionnons maintenant la seconde feuille « ONLY LEARNING SET ». En situation d'étude réelle, nous ne disposons que de ces 500 observations. Comme nous le disions plus haut, nous devons les utiliser à la fois pour construire le modèle et en mesurer les performances. Pour cela, nous allons mettre en œuvre les techniques de ré échantillonnage disponibles dans TANAGRA.

² Même si notre ensemble test comporte 42500 observations, il s'agit quand même d'un échantillon. L'erreur mesurée n'est qu'une estimation. On peut néanmoins penser qu'elle est assez précise au regard de la taille de l'échantillon test. C'est à ce titre que nous l'utiliserons comme référence pour les techniques de ré échantillonnage présentées par la suite.

Dans un premier temps, comme précédemment, nous sélectionnons la plage de données dans EXCEL, puis nous activons le menu TANAGRA/EXECUTE TANAGRA. Bien sûr, il est totalement inutile de sélectionner la dernière colonne STATUS, elle est constituée de la même valeur « LEARNING » pour toutes les lignes du tableau. Nous vérifions que la plage de cellules sélectionnée correspond bien à **\$A\$1:\$V\$501**



TANAGRA est automatiquement démarré, nous définissons les variables en TARGET (CLASS) et INPUT (V1 à V21). Puis nous insérons le composant LINEAR DISCRIMINANT ANALYSIS. Le modèle et l'erreur en resubstitution (6.2%) sont exactement les mêmes que précédemment. C'est normal. Le calcul est complètement déterministe et nous utilisons exactement les mêmes données pour l'apprentissage.



Dans ce qui suit, nous montrons comment utiliser les composants d'évaluation de l'erreur en ré échantillonnage dans TANAGRA. Nous vérifierons également si l'erreur calculée dans ce cas se rapproche de la « vraie » erreur (avec les réserves émises plus haut) 8.39% estimée sur le fichier de 42500 observations.

Erreur « leave one out »

L'idée de l'erreur *leave one out* est de réitérer l'opération suivante pour chaque observation i : construire le modèle sur les $(n-1)$ observations en excluant l'observation numéro i , l'appliquer sur cet individu, et vérifier qu'il est bien classé ($d_i = 1$ s'il est mal classé, $d_i = 0$ sinon). L'erreur $e-lvo$ s'écrit

$$e-lvo = \frac{1}{n} \sum_i d_i$$

Dans TANAGRA, il s'agit de glisser le composant LEAVE-ONE-OUT (onglet SPV LEARNING ASSESSMENT) à la suite de l'analyse discriminante et d'activer le menu VIEW pour obtenir les résultats. TANAGRA effectue donc 500 fois l'opération « apprentissage sur 499 observations et test sur 1 observation ». Ca semble considérable mais compte tenu des caractéristiques de notre étude (taille du fichier + méthode), le temps d'attente est raisonnable (#8 secondes).

The screenshot shows the TANAGRA 1.4.20 interface. The 'Analysis' pane on the left shows a tree view with 'Leave-One-Out 1' selected. The 'Results' pane on the right displays the following information:

Leave-One-Out 1

Parameters

Results

Overall cross-validation error rate: 0.0800

Values prediction

Value	Recall	1-Precision
B	0.8961	0.0717
A	0.9405	0.0866

Confusion matrix

	B	A	Sum
B	207	24	231
A	16	253	269
Sum	223	277	500

Computation time : 8547 ms.
Created at 20/10/2007 06:48:20

Le taux d'erreur $e-lvo$ est égal à 8%. L'estimation semble satisfaisante. De manière générale, la technique *leave one out* permet d'estimer correctement l'erreur, avec un biais faible, mais une variance plus élevée que les autres techniques. Elle est totalement inappropriée en revanche dans les cas de fort sur apprentissage (par exemple, ratio nombre de descripteurs – nombre d'observations très élevé associé à la méthode du plus proche voisin).

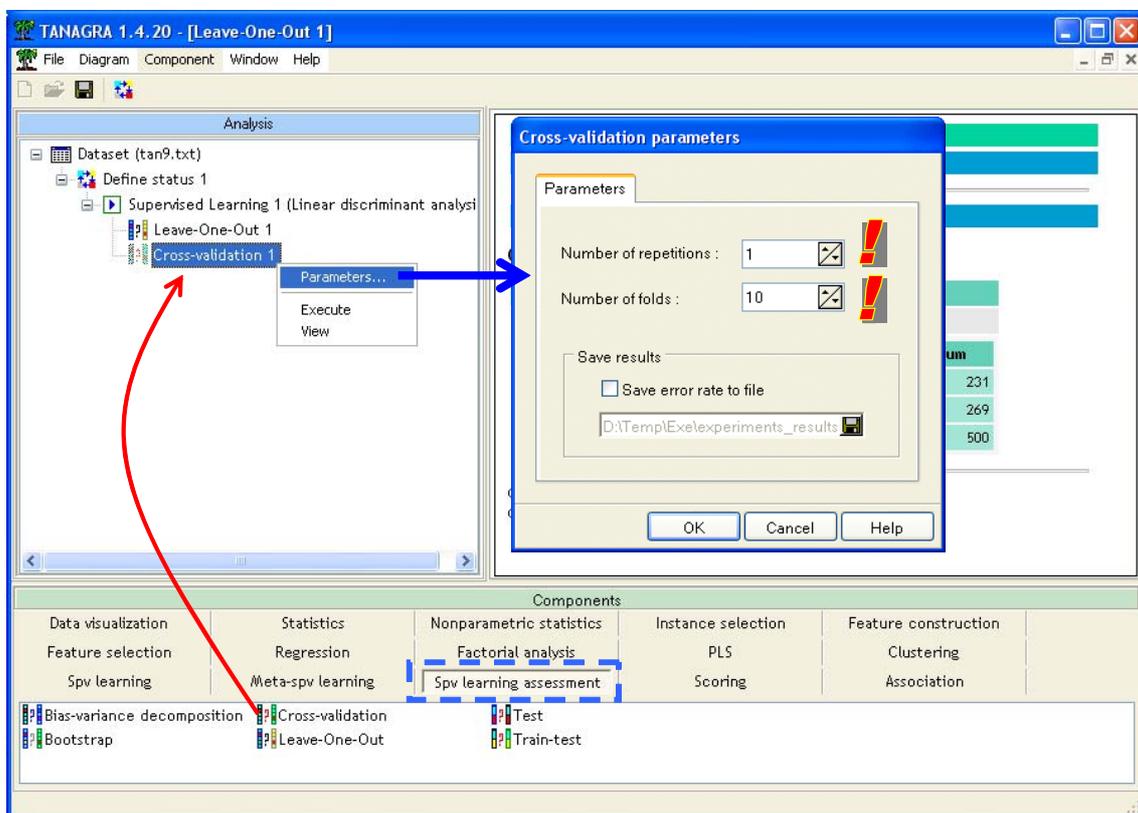
Erreur « validation croisée »

La validation croisée est une sorte de simplification du leave one out où l'on subdivise les données en K blocs. On répète alors K fois le processus suivant : apprentissage sur (K-1) blocs, test sur le K-ième bloc, nous mesurons l'erreur en test e_k . L'erreur e-cv s'écrit

$$e - cv = \frac{1}{K} \sum_k e_k$$

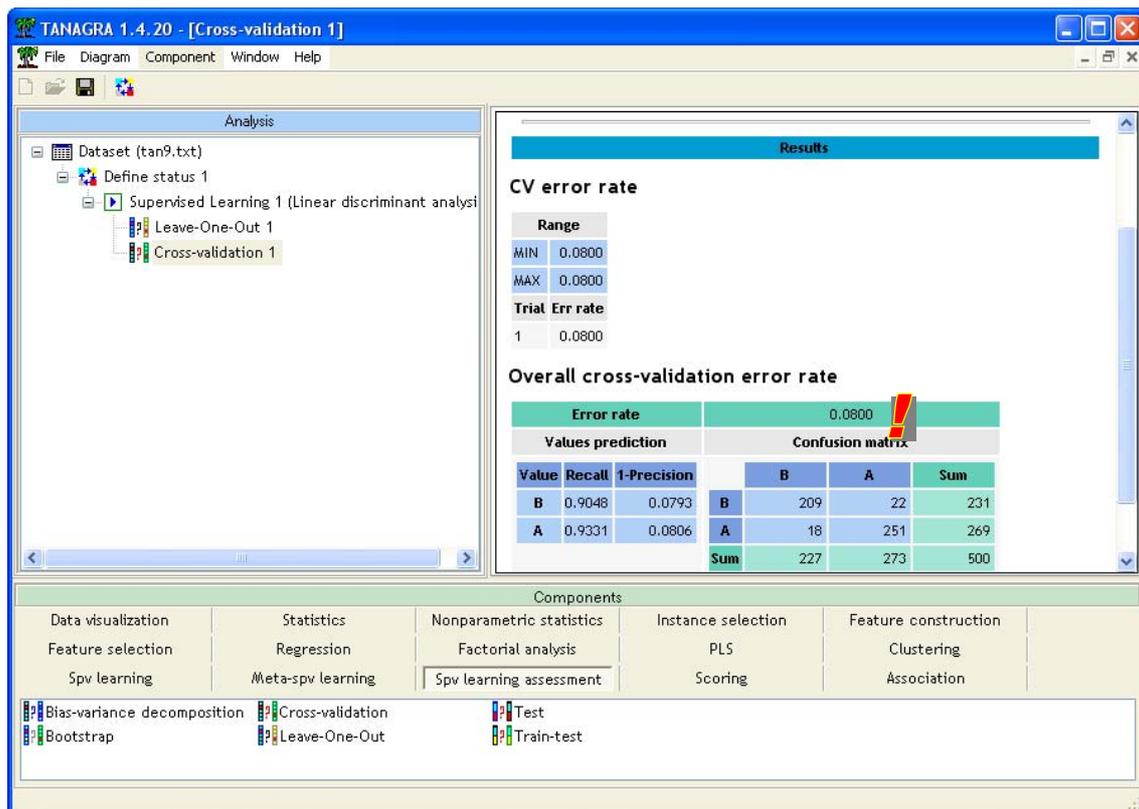
Son principal atout est de réduire la quantité de calcul par rapport au leave one out, sans détériorer la précision de l'estimation. Elle peut même s'avérer meilleure dans certaines situations. Les études empiriques montrent que K=10 semble un bon compromis.

Nous insérons le composant CROSS-VALIDATION à la suite de l'analyse discriminante. Nous activons le menu PARAMETERS pour définir les paramètres de la méthode.



Nous spécifions bien K=10 (NUMBER OF FOLDS), nous notons qu'il est possible de réitérer plusieurs fois le processus (NUMBER OF REPETITIONS), cela permet d'apprécier la variabilité des résultats, au prix certes d'une quantité de calculs plus élevée. Dans notre cas, nous laissons NUMBER OF REPETITIONS = 1.

Après avoir validé, et activer menu VIEW, nous obtenons les résultats suivants :



Le taux d'erreur e-cv est 8%, identique à e-lvo. C'est un hasard, nous observons d'ailleurs que la matrice de confusion est un peu différente. Il reste néanmoins que les taux d'erreur mesurés en validation croisée et par leave one out sont souvent assez proches.

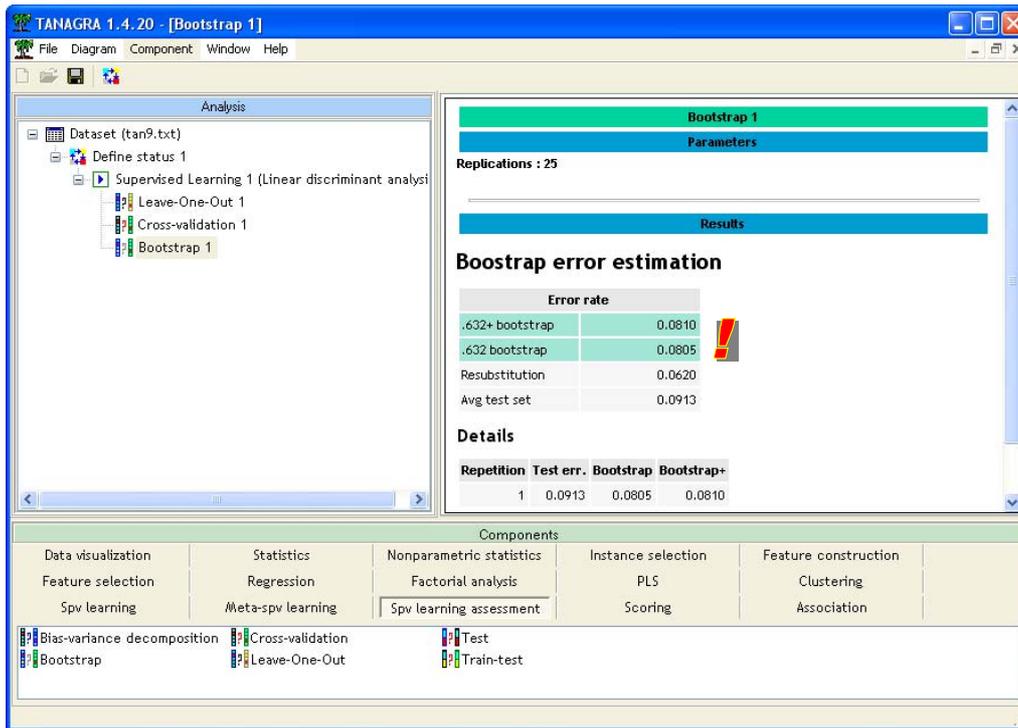
Erreur « bootstrap »

Contrairement aux deux techniques précédentes, le bootstrap ne cherche pas à estimer directement l'erreur, mais plutôt le biais d'optimisme. La présentation est un peu complexe, pour avoir plus de détails, mieux vaut consulter le document en ligne http://eric.univ-lyon2.fr/~ricco/cours/slides/resampling_evaluation.pdf

Deux estimateurs bootstrap sont disponibles dans TANAGRA, le standard 0.632 bootstrap et l'indicateur modifié 0.632 bootstrap+, censé tenir compte des spécificités de la technique d'apprentissage. Dans certaines situations, la correction peut être importante.

Le seul paramétrage possible est le nombre de réplifications, dans la majorité des cas, 25 répétitions suffisent largement pour obtenir une estimation satisfaisante. En augmentant le nombre de réplifications, nous réduisons la variance de l'estimateur. Le biais d'estimation lui reste stable. On convient généralement que le bootstrap est légèrement biaisé par rapport à la validation croisée, il présente en revanche une variance nettement plus faible à effort de calcul égal.

Nous insérons le composant BOOTSTRAP dans le diagramme. Nous obtenons les résultats suivants :



Le taux d’erreur $e_{-boot+} = 8.1\%$.

Nous résumons les résultats dans le tableau suivant

Erreur	LDA	Ecart
Resubstitution	6.2 %	- 2.19
« Vraie » (test)	8.39 %	x
LVO	8 %	- 0.39
CV (10)	8 %	- 0.39
Bootstrap+ (25)	8.1 %	- 0.29

Rappelons que ce que nous qualifions de « vraie » erreur est aussi une estimation dans ce didacticiel puisqu’il est calculé sur un échantillon. Nous pouvons néanmoins lui faire crédit car elle est calculée sur un ensemble de données à part extraite de la population, qui n’a en aucune manière participé à la construction des modèles, et qui possède un effectif suffisamment élevé.

Concernant les techniques d’estimation de l’erreur, celui qui se rapproche le plus du « vrai » taux d’erreur est l’estimateur « Bootstrap » avec un écart de -0.29 %. Mais les différences entre les techniques restent minimes, il ne faut surtout pas en tirer des conclusions hâtives. **La précision des estimateurs dépend des caractéristiques des données analysées et des méthodes d’apprentissage utilisées.**

Enfin, le fait que les techniques sous-estiment toutes la véritable erreur dans cet exemple est fortuit. Il arrive parfois que la technique de ré échantillonnage la surestime comme nous le verrons plus bas.

5. Conclusion

Dans ce didacticiel, nous avons mis en œuvre plusieurs techniques de ré échantillonnage pour évaluer l'erreur de prédiction sans disposer d'un échantillon test. Sur des données simulées, nous constatons que l'erreur ainsi évaluée est assez proche du véritable taux d'erreur que nous avons calculé par ailleurs sur un échantillon test de taille virtuellement infinie (si tant est que 42500 est suffisamment grand pour obtenir une précision satisfaisante).

Il reste que nous nous sommes placés dans un cadre plutôt favorable en combinant une méthode relativement stable (l'analyse discriminante) et un problème somme toute facile (500 observations pour 21 variables avec une distribution monomodale des classes, le test basé sur Lambda de Wilks confirme clairement cette appréciation). Nous avons voulu vérifier le comportement de ces techniques pour une méthode avec des caractéristiques assez différentes (arbre de décision C4.5), connue pour sa propension au sur apprentissage, nous obtenons les résultats suivants :

Erreur	C4.5	Ecart
Resubstitution	2.2%	- 11.34
« Vraie » (test)	13.54%	x
LVO	16%	+ 2.46
CV (10)	16.2%	+ 2.66
Bootstrap+ (25)	12.24%	- 1.3

Manifestement, C4.5 est nettement moins performant que la LDA sur ces données.

Concernant les estimations, l'erreur en resubstitution est clairement mauvaise, trop optimiste. Cela est dû en grande partie aux caractéristiques de C4.5 qui a parfois tendance à trop coller aux données avec un arbre sur dimensionné.

De manière générale, la précision des estimations par ré échantillonnage est moindre. Nous constatons que dans certains cas (LVO, CV), l'erreur est surestimée.