

- **Pour comparer deux modèles ayant le même nombre de variables explicatives** : comparer les R^2 obtenus et choisir le modèle pour lequel R^2 le plus grand.

- **Pour comparer un modèle avec $(p - 1)$ variables avec un modèle $(p - 1 + r)$ variables** : utiliser le test du F partiel.

Que nous dit ce test ?

Ce test dit si l'introduction des variables supplémentaires augmente suffisamment le R^2 ou non.

Propriété sur R_{aj}^2 :

- $R_{aj}^2 < R^2$ dès que $p \geq 2$.
- R_{aj}^2 peut prendre des valeurs négatives.

Intérêts de R_{aj}^2 :

- R_{aj}^2 n'augmente pas forcément lors de l'introduction de variables supplémentaires dans le modèle.
- possibilité de comparer deux modèles n'ayant pas le même nombre de variables à l'aide du R_{aj}^2 et choisir le modèle pour lequel R_{aj}^2 est le plus grand.

Le coefficient de détermination multiple ajusté

R_{aj}^2 :

- Introduire un R^2 qui concerne la population et non plus l'échantillon défini par :

$$R_{pop}^2 = 1 - \frac{\sigma^2}{\sigma^2(Y)}$$

- Estimer R_{pop}^2 par :

$$R_{aj}^2 = 1 - \frac{s^2}{s^2(Y)} = 1 - \frac{SC_{res}}{SC_{tot}} \frac{n-1}{n-p}$$

Le critère du C_p de Mallows

Le critère du C_p de Mallows est défini par :

$$C_p = \frac{SC_{res}}{\sigma^2} - (n - 2p)$$

Mais il y a un problème ! : Lequel ?

On ne peut plus estimer σ^2 par

$$s^2 = \frac{SC_{res}}{n - p}$$

Pourquoi ?

Car C_p vaudrait toujours p et alors il ne serait plus intéressant.

Le critère BIC

Le Bayesian information criterion BIC est défini par :

$$BIC = -2 \log(\tilde{L}) + k \log(n).$$

Il est plus parcimonieux que le critère AIC puisqu'il pénalise plus le nombre de variables présentent de le modèle.

Ripley, 2003, souligne que l' AIC a été introduit pour retenir des variables pertinentes lors de prévisions, et que le critère BIC vise la sélection de variables statistiquement significative dans le modèle.

Plusieurs types de procédures de sélection de variables :

- la recherche exhaustive
- les méthodes de type pas à pas.

L'efficacité de ces méthodes n'est pas démentie mais il est impossible de se fier aux résultats fournis par un programme informatique.

Exemples :

Quant à décider ou supprimer une variable dans un modèle, il faut conserver :

- une part d'intuition
- une part de déduction
- une part de synthèse.

Surtout ne pas oublier l'objectif recherché !

En pratique comment fait-on ?

Méthode ascendante (ou sélection en avant) :

- C'est également une simplification de la méthode de la recherche exhaustive.
- Cette méthode procède dans le sens inverse de la méthode descendante.
- Cette méthode examine un modèle avec une seule variable explicative puis introduction une à une d'autres variables explicatives.

Sinon

- Réitérer le processus en effectuant les $(k - 2)$ régressions possibles avec trois variables explicatives. Pour chacune d'elles, **effectuer le test de Student pour la nouvelle variable**. Retenir le modèle pour lequel la variable est la plus significative. Si aucune variable est retenue, alors on stoppe le processus.

Sinon

- Réitérer le processus en effectuant les $(k - 3)$ régressions possibles avec quatre variables explicatives...

Le processus se termine lorsqu'on ne peut plus introduire des variables significatives dans le modèle.

- Effectuer les k régressions possibles avec une seule variable explicative. Pour chacune d'elles, **effectuer le test de Student**. Retenir le modèle pour lequel la variable explicative est la plus significative.
- Effectuer les $(k - 1)$ régressions possibles avec deux variables explicatives. Pour chacune d'elles, **effectuer le test de Student pour la nouvelle variable**. Retenir le modèle pour lequel la variable est la plus significative. Si aucune variable est retenue, alors on stoppe le processus.

Parmi les méthodes présentées, **la méthode ascendante est la plus économique**.

Les avantages de la méthode ascendante :

- éviter de travailler avec plus de variables que nécessaire,
- améliorer l'équation à chaque étape.

L'inconvénient majeur de la méthode ascendante :

- une variable introduite dans le modèle ne peut plus être éliminée.

Le modèle final peut alors contenir des variables non significatives.

Ce problème est alors résolu par la procédure stepwise.

Procédure stepwise :

amélioration de la méthode descendante.

- **Pourquoi ?** À chaque étape, on réexamine toutes les variables introduites précédemment dans le modèle. En effet, une variable considérée comme la plus significative à une étape de l'algorithme peut à une étape ultérieure devenir non significative. **Pourquoi ce phénomène ?**

- En raison de ces corrélations avec d'autres variables introduites après coup dans le modèle.

La procédure stepwise

semble être la meilleure procédure de sélection de variables.

Mais

- **la procédure stepwise** peut facilement abuser l'utilisateur qui a tendance à se focaliser exclusivement sur le résultat de la sélection automatique proposé par l'outil informatique.
- En effet, il faut se méfier de certaines situations : celles où apparaît un **phénomène de multicollinéarité** que nous étudierons dans un prochain chapitre.

Comment remédie-t-on à cela ?

La procédure stepwise propose après l'introduction d'une nouvelle variable dans le modèle :

- **réexaminer les tests de Student** pour chaque variable explicative anciennement admise dans le modèle,
- après réexamen, si des variables ne sont plus significatives, alors **retirer du modèle la moins significative d'entre elles**.

Le processus continue jusqu'à ce que plus aucune variable ne puisse être introduite ni retirée du modèle.

Un exemple :

X_1 et X_2 expliquant significativement Y , sont fortement corrélées entre elles.

L'introduction de X_1 dans le modèle masquera le pouvoir explicatif de X_2 .

En effet, l'introduction de X_1 en premier va accroître le R^2 alors que l'introduction ultérieure de X_2 provoquera un faible effet.

Et réciproquement.

La procédure de régression stagewise

diffère des procédures présentées ci-dessus.

Pourquoi ?

La procédure de régression stagewise n'aboutit pas toujours à une équation obtenue par la méthode des moindres carrés.

- effectuer la régression avec la variable la plus corrélée avec Y.
- Calculer les résidus obtenus avec cette régression.
- Considérer ensuite ces résidus comme une nouvelle variable dépendante que l'on veut expliquer à l'aide des variables explicatives restantes.

Cette méthode se déroule de la façon suivante :

- Sélectionner ainsi celle d'entre elles qui est la plus corrélée avec les résidus.
- Effectuer une nouvelle régression avec les résidus de la première régression dans le rôle de la variable expliquée.
- Également calculer les résidus obtenus avec cette nouvelle régression.
- Répéter le processus avec des variables explicatives restantes, jusqu'à ce que plus aucune d'entre elles ne soit corrélée significativement avec les résidus.

- D'un point de vue théorique, **la recherche exhaustive est la meilleure.**
- **Pour les autres méthodes** : Des résultats semblables sans nécessiter autant de calculs. Ces derniers dépendent du choix des seuils de signification utilisés lors des diverses procédures.
- **Dans la pratique, la procédure stepwise et la procédure descendante** sont les plus utilisées.
- En cas de doute et si les conditions le permettent, toutes les régressions doivent être examinées.
- **Les autres méthodes** peuvent trouver leur utilisation dans des applications plus spécifiques.

Quelques conclusions