

- **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^2_{aj} :

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

Intérêts de R^2_{aj} :

- R^2_{aj} 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^2_{aj} et choisir le modèle pour lequel R^2_{aj} est le plus grand.

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

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

- Estimer R^2_{pop} par :

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

Le critère du C_p de Mallows

est défini par :

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

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.

Que fait-on dans la pratique ?

- On estime σ^2 par le s^2 du modèle qui fait intervenir toutes les k variables explicatives du modèle à disposition.
Pour ce modèle on a : $C_p = p$. Et pour les autres ? C_p prendra d'autres valeurs que p .
- On choisit parmi les modèles le modèle où C_p de Mallows est le plus proche de p .

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 !

Méthode descendante

Que fait-on ensuite ?

- Ces équations sont réparties selon le nombre r de variables explicatives qu'elles contiennent.
- Chaque ensemble d'équations est ordonné selon le critère choisi, souvent le R^2 .
- Les meilleures équations de régression issues de ce classement sont ensuite sélectionnées pour un examen plus détaillé.

(ou élimination en arrière) est une simplification de la méthode de la recherche exhaustive.

En quoi est-elle une simplification ?

Cette méthode examine non pas toutes les régressions possibles mais uniquement une régression pour chaque nombre r de variables explicatives.

En pratique comment fait-on ?

- Calculer la régression pour le modèle incluant toutes les k variables explicatives à disposition.
- Effectuer un test de Student pour chacune des variables explicatives. **Deux cas se présentent :**
 - Les variables sont trouvées significatives. Ce modèle est alors choisi. On stoppe là notre analyse.
 - Éliminer la variable la moins significative du modèle.
- Recommencer le processus avec une variable en moins.

Le modèle final est donc un modèle au sein duquel toutes les variables sont significatives.

Méthode descendante

Conclusions :

- La méthode descendante est très satisfaisante pour l'utilisateur préférant avoir toutes les variables possibles afin de ne rien ignorer.
- C'est une procédure plus économique en terme de temps et d'interprétation
- Mais il y a un inconvénient majeur. Il n'est plus possible de réintroduire une variable une fois qu'elle a été supprimée !

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.

En pratique comment fait-on ?

- 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.

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.

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.

