

Choix du modèle

Frédéric Bertrand & Myriam Maumy-Bertrand¹

¹IRMA, Université de Strasbourg
France

Master 1 01-02-2012

Remarque

Dans tous les modèles que nous allons considérer, la constante sera par défaut automatiquement présente. Ainsi, si le modèle comporte $p - 1$ variables explicatives, il y aura au total p paramètres de régression notés $\beta_0, \dots, \beta_{p-1}$.

Critères de choix

- Il existe plusieurs critères pour sélectionner $p - 1$ variables explicatives parmi k variables explicatives disponibles, avec $k \geq p - 1$.
- Le critère du R^2 se révèle le plus simple à définir.

Ce cours s'appuie sur l'ouvrage suivant

- « Analyse de régression appliquée »
Yadolah Dodge
Dunod, 1999

Le lecteur intéressé par ce sujet pourra consulter l'ouvrage suivant

- « Le modèle linéaire par l'exemple »
J.-M. Azais et J.-M. Bardet
Dunod, 2005

Inconvénient majeur du R^2

Il augmente de façon monotone avec l'introduction de nouvelles variables même si celles-ci sont peu corrélées avec la variable expliquée Y .

Pour parler à cet inconvénient

Il est conseillé de se tourner vers l'utilisation des alternatives suivantes

- 1 le R^2 ajusté
- 2 le C_p de Mallows qui est un autre critère relatif au biais
- 3 le critère AIC
- 4 le critère AIC_c
- 5 le critère BIC .

Remarque

Ces six critères vont être maintenant présentés.

Coefficient de détermination multiple : R^2

- Il a déjà été introduit dans le cours portant sur la régression linéaire simple.
- C'est une mesure qui permet d'évaluer le degré d'adéquation du modèle.
- Lors de l'introduction du test de Fisher partiel : accroissement de R^2 au fur et à mesure de l'introduction de variables dans le modèle. Il atteint son maximum lorsque toutes les variables disponibles au départ sont incluses.

Pour comparer deux modèles ayant le même nombre de variables explicatives

Comparer les R^2 obtenus et choisir le modèle pour lequel R^2 est le plus grand.

Pour comparer un sous-modèle avec $p - 1$ variables d'un modèle avec $(p - 1 + r)$ variables

Utiliser le test de Fisher partiel. Que nous dit ce test ?
Ce test dit si l'introduction des variables supplémentaires augmente suffisamment le R^2 ou non.

Coefficient de détermination multiple ajusté R^2_{aj}

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

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

- Estimer R^2_{pop} par :

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

Propriétés sur R^2_{aj}

- 1 $R^2_{aj} < R^2$ dès que $p \geq 2$.
- 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.

Solution

Que faisons-nous dans la pratique ?

- 1 Nous estimons σ^2 par le s^2 du modèle qui fait intervenir toutes les k variables explicatives du modèle à disposition. Pour ce modèle, qui a $k + 1$ paramètres, nous avons toujours : $C_{k+1} = k + 1$. Et pour les autres ? C_p prendra d'autres valeurs que p .
- 2 **Critère du C_p de Mallows**
Nous choisissons parmi les modèles le modèle où le C_p de Mallows est le plus proche de p .

Statistique du C_p de Mallows

La statistique du C_p de Mallows est définie par :

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

Remarque : problème

Il y a un problème ! Nous ne pouvons pas 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.

Critère d'information d'Akaike (AIC)

Le critère d'information AIC s'applique aux modèles estimés par **une méthode du maximum de vraisemblance** : les analyses de variance, les régressions linéaires multiples, les régressions logistiques et de Poisson peuvent rentrer dans ce cadre.

Le critère AIC est défini par :

$$AIC = -2 \log(\tilde{L}) + 2\tilde{k}$$

où \tilde{L} est la vraisemblance maximisée et \tilde{k} le nombre de paramètres libres dans le modèle.

Avec ce critère, la déviance du modèle $-2 \log(\tilde{L})$ est pénalisée par 2 fois le nombre de paramètres libres.

Interprétation

Le critère **AIC** représente donc un **compromis** entre le **biais**, diminuant avec le nombre de paramètres libres, et la **parcimonie**, volonté de décrire les données avec le plus petit nombre de paramètres possibles.

Expression dans le cadre du modèle linéaire gaussien

Le critère AIC devient

$$AIC = 2\tilde{k} + n \left[\ln \left(\frac{2\pi \times SC_{res}}{n} \right) + 1 \right].$$

Critère d'information d'Akaike corrigé (AIC_c)

Lorsque le **nombre de paramètres libres \tilde{k} est grand** par rapport au nombre d'observations n , c'est-à-dire si $n/\tilde{k} < 40$, il est recommandé d'utiliser l'**AIC corrigé**.
 Le critère d'information d'Akaike corrigé, AIC_c, est défini par :

$$AIC_c = AIC + \frac{2\tilde{k}(\tilde{k} + 1)}{n - \tilde{k} - 1}.$$

Référence : « Model selection for extended quasi-likelihood models in small samples. »
 Hurvich, C. M. and Tsai, C.-L., 1995. *Biometrics* **51** : 1077-1084.

Remarques

- 1 La rigueur voudrait que tous les modèles comparés dérivent tous d'un même « complet » inclus dans la liste des modèles comparés.
- 2 Il est nécessaire de vérifier que les **conditions d'utilisation** du modèle complet et de celui sélectionné sont **remplies**.
- 3 Le meilleur modèle est celui possédant l'**AIC le plus faible**.

Critère BIC

Le critère d'information bayésien BIC est défini par :

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

Remarque

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

Remarque très importante

Toutes ces procédures ne mènent pas forcément à la même solution quand elles sont appliquées au même problème.

Recherche exhaustive

Lorsque le nombre de variables explicatives, noté k , à disposition n'est pas trop élevé, il est envisageable de considérer tous les modèles possibles.

Il y a

$$C_k^r = \frac{k!}{r!(k-r)!}$$

modèles différents faisant intervenir r variables explicatives.

Recherche exhaustive - Suite et fin

Cela fait au total

$$\sum_{r=0}^k C_k^r = 2^k$$

modèles possibles à considérer.

Remarque

Nous choisissons ensuite le modèle pour lequel, par exemple, le R_{aj}^2 est le maximum.

Méthodes de type pas à pas

Les méthodes de type pas à pas consistent à considérer d'abord un modèle faisant intervenir un certain nombre de variables explicatives. Puis nous procédons par élimination ou ajout successif de variables.

- Nous parlons de la **méthode descendante** lorsque nous **éliminons des variables** (elle sera développée dans le paragraphe 4)
- Nous parlons de la **méthode ascendante** lorsque nous **ajoutons des variables** (elle sera développée dans le paragraphe 5).
- La **méthode stepwise** est une **combinaison de ces deux méthodes** (elle sera développée dans le paragraphe 6).

La recherche exhaustive

C'est une méthode fastidieuse et difficile à utiliser sans un ordinateur rapide.

Pourquoi ?

Parce qu'il faut calculer toutes les régressions possibles impliquant un sous-ensemble des k variables explicatives à disposition, soit un total de 2^k régressions.

Méthode descendante (ou élimination en arrière)

C'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.

Que faisons-nous 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 ou l' AIC .
- Les meilleures équations de régression issues de ce classement sont ensuite sélectionnées pour un examen plus détaillé.

En pratique comment faisons-nous ?

- 1 Calculer la régression pour le modèle incluant toutes les k variables explicatives à disposition.
- 2 Effectuer un test de Student pour chacune des variables explicatives.

Deux cas se présentent :

- Les variables sont trouvées significatives. Ce modèle est alors choisi. Nous arrêtons là notre analyse.
- Éliminer la variable la moins significative du modèle.
- 3 Recommencer le processus avec une variable en moins.

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

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 faisons-nous ?

- 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 nous arrêtons 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 nous arrêtons là le processus.

Sinon

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

Le processus se termine lorsque nous ne pouvons plus introduire des variable significatives dans le modèle.

Comment remédier 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.

Procédure stepwise

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 apparait **un phénomène de multicollinéarité**. La régression PLS permet de traiter ce problème.

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.

Premières remarques

Les méthodes présentées jusqu'ici :

- ne sélectionnent pas nécessairement le meilleur modèle absolu,
- mais donnent généralement un modèle acceptable.

Des procédures alternatives et combinatoires sont apparues. Nous présenterons très brièvement seulement **deux méthodes** dans ce cours.

Première variante

Effectuer d'abord la procédure stepwise.

Supposons que le modèle sélectionné contient r variables explicatives.

Effectuer alors C_k^r régressions possibles utilisant r variables.

Choisir comme modèle final celui pour lequel R^2 est maximal.
D'un point de vue pratique, **les améliorations apportées par cette procédure sont faibles...**et nécessitent beaucoup de calculs.

Deuxième variante

Utiliser deux seuils de signification différents lors de l'utilisation de la méthode stepwise :

- un certain seuil (par exemple $\alpha = 0,05$) lors de la procédure d'élimination de variables,
- un seuil plus petit (par exemple $\alpha = 0,01$) lors de la procédure d'introduction de variables.

Nous favorisons l'introduction des variables les plus significatives, tout en acceptant de les maintenir dans le modèle si elles deviennent par la suite un peu moins significatives.

C'est intéressant d'utiliser cette méthode.

Pourquoi ?

Cette méthode propose un modèle alternatif au modèle donné par **la méthode stepwise** tout en maintenant son pouvoir explicatif.

La procédure de régression stagewise

Elle 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.

Cette méthode se déroule de la façon suivante :

- 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.

Quelques conclusions

- 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.

- 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.