

# AN INTRODUCTION TO STATISTICAL LEARNING

## Chapitre 1 - Introduction

L'apprentissage statistique se réfère à des outils pour comprendre les données.

On peut les classer en deux catégories : apprentissage supervisé et non supervisé

L'apprentissage supervisé a pour but de construire un modèle pour estimer, prédire une sortie basée sur des entrées.

L'apprentissage non supervisé quand à lui, il possède des sorties mais pas de sortie. On apprend les relations et structures des données

## Chapitre 2 – apprentissage statistique

I Qu'est-ce que l'apprentissage statistique ?

De façon générale, on suppose observer une réponse quantitative et  $p$  différents prédicteurs  $X_1, X_2, \dots, X_p$ .

On suppose que la relation entre  $Y$  et  $X = (X_1, \dots, X_p)$  peut être écrite sous la forme

$$Y = f(X) + \varepsilon$$

où  $f$  est une fonction inconnue de  $X_1, \dots, X_p$  et  $\varepsilon$  est une erreur aléatoire indépendante de  $X$  et de moyenne nulle.

L'apprentissage statistique consiste à s'approcher d'une estimation de  $f$

Il y a deux raisons pour que l'on puisse désirer d'estimer  $f$  : prédiction et inférence

a) prédiction

On peut prédire  $Y$  en utilisant la moyenne  $\hat{Y} = \hat{f}(X)$

où  $\hat{f}$  représente notre estimation de  $f$  et  $\hat{Y}$  représente le résultat de prédiction de  $Y$

$\hat{f}$  est souvent considéré comme une boîte noire

La précision de  $\hat{Y}$  comme prédiction de  $Y$  dépend de deux quantités : erreur réductible et l'irréductible erreur. En général  $\hat{f}$  ne sera pas une estimation parfaite de  $f$ , et son imprécision introduira quelques erreurs. Cette erreur est réductible car on peut potentiellement améliorer la précision de  $\hat{f}$  en utilisant une technique statistique plus appropriée pour estimer  $f$ .

Cependant même s'il est possible de se rapprocher d'une estimation parfaite de  $f$ , notre prédiction aura encore des erreurs. Cela est dû au fait que  $Y$  est aussi une fonction de  $\varepsilon$  qui ne peut être prédite en utilisant  $X$ . Donc  $\varepsilon$  peut affecter notre prédiction. Cela est connu sous le nom d'erreur irréductible.

Considérons une estimation de  $\hat{f}$  et un jeu de prédicteurs  $X$ , qui implique la prédiction  $\hat{Y} = \hat{f}(X)$ . On suppose que  $\hat{f}$  et  $X$  sont fixés. On peut donc aisément montrer que

$$E(Y - \hat{Y})^2 = E[f(X) + \varepsilon - \hat{f}(X)]^2$$

$$E(Y - \hat{Y})^2 = [f(X) - \hat{f}(X)]^2 + \text{Var}(\varepsilon)$$

où la première partie est réductible mais la seconde irréductible

Où  $E(Y - \hat{Y})^2$  représente la moyenne ou la valeur attendue du carré de la différence entre la valeur prédite et la valeur actuelle de  $Y$

Démonstration:

$$\text{Var}[X] = E[X^2] - E[X]^2$$

$$E[X^2] = \text{Var}[X] + E[X]^2$$

$$E[f] = f$$

$$y = f + \varepsilon$$

$$E[\varepsilon] = 0$$

$$E[y] = E[f + \varepsilon] = E[f] = f$$

$$\text{Var}[y] = E[(y - E[y])^2] = E[(y - f)^2]$$

$$E[(y - \hat{f})^2] = E[y^2 + \hat{f}^2 - 2y\hat{f}]$$

$$E[(y - \hat{f})^2] = E[y^2] + E[\hat{f}^2] - E[2y\hat{f}]$$

$$E[(y - \hat{f})^2] = \text{Var}[y] + E[y]^2 + \text{Var}[\hat{f}] + E[\hat{f}]^2 - 2fE[\hat{f}]$$

$$E[(y - \hat{f})^2] = \text{Var}[y] + \text{Var}[\hat{f}] + (f - E[\hat{f}])^2$$

$$E[(y - \hat{f})^2] = \text{Var}[y] + \text{Var}[\hat{f}] + E[(f - \hat{f})^2]$$

$$E[(y - \hat{f})^2] = \text{Var}[\varepsilon] + \text{Var}[\hat{f}] + \text{Biais}[\hat{f}]^2$$

Ce livre se concentre sur les techniques d'estimations de  $f$  avec le but de minimiser l'erreur reductible.

## b) inférence

On est souvent intéressé à comprendre le fait que  $Y$  soit affecté par  $X_1, \dots, X_p$  lorsqu'il change.

Dans cette situation on veut estimer  $f$  mais notre objectif n'est pas nécessairement de faire des prédictions de  $Y$ . Au lieu de cela on veut comprendre la relation entre  $X$  et  $Y$  et plus spécifiquement comprendre comment  $Y$  change en fonction de  $X_1, \dots, X_p$ . Maintenant  $\hat{f}$  ne peut être considéré comme une boîte noire.

On s'intéressera à différentes questions:

Quel prédicteur est associé à la réponse ?

Quel est la relation entre la réponse et chaque prédicteur ?

Est-ce que la relation entre  $Y$  et chaque prédicteur puisse être résumée en utilisant une équation linéaire, ou plus compliquée ?

Suivant la complexité du modèle on aura un compromis à avoir (plus le modèle est compliqué, moins on peut effectuer d'inférence) et inversement (plus le modèle est simple plus on peut effectuer des inférences).

Le modèle linéaire par exemple est relativement simple et interprétable mais la précision des prédictions peut être faible.

Tout au long de ce livre, on étudiera des modèles linéaires et non-linéaires pour estimer  $f$ .

On supposera toujours que nous observons un jeu de  $n$  données différentes

Notre but est d'appliquer les méthodes statistiques au jeu d'entraînement pour estimer la fonction inconnue  $f$ .

En d'autres termes on veut trouver une fonction  $\hat{f}$  comme  $Y \approx \hat{f}(X)$  pour toute observation  $(X, Y)$ .

La plupart des modèles d'apprentissage statistique pour cette tâche peuvent être caractérisés comme paramétrique ou non-paramétrique.

#### a) les méthodes paramétriques

Les méthodes paramétriques sont basées sur deux étapes.

1-/ premièrement, on fait une supposition de la forme et de la dimension de la fonction  $f$ .

Par exemple, une très simple supposition est que  $f$  est linéaire en  $X$ :

$$f(X) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p$$

Ceci est un modèle linéaire que l'on étudiera plus dans le chapitre 3

Une fois les suppositions faites, au lieu d'avoir à estimer une fonction entière de dimension  $p$ , on peut se retrouver simplement avec quelques coefficients (par exemple  $p+1$  coefficient pour le modèle linéaire).

2-/Après que le modèle soit sélectionné, on a besoin d'une procédure qui utilise les données d'entraînement pour adapter le modèle. Pour l'exemple du modèle linéaire, on utilise souvent les moindres carrés (à voir dans le chapitre 3). Les moindres carrés est une des multitudes de méthodes pour ajuster le modèle linéaire.

L'approche du modèle décrit ci-dessus est considéré comme paramétrique: on réduit le problème de l'estimation de  $f$  en estimant un certain nombre de paramètres.

Le principal désavantage d'une approche paramétrique est que le modèle choisit ne corresponde pas avec la vraie fonction inconnue  $f$ .

Si notre modèle est trop loin de la réalité nos estimations seront mauvaises.

On pourra essayer d'essayer des modèles flexibles qui peuvent ajuster sur différents modèles possibles de  $f$ .

Mais en général l'ajustement d'un modèle plus flexible requiert d'estimer plus de paramètres. Ces modèles plus flexibles peuvent amener à du surajustement.

Avec un trop petit nombre de données, l'ajustement linéaire peut s'avérer être le plus pratique et le plus proche.

#### b) les méthodes non-paramétriques

Les méthodes non-paramétriques ne créent pas des suppositions explicites sur la forme de la fonction  $f$ .

Au lieu de cela, elles recherchent une estimation de  $f$  qui devient aussi proche des données que possible sans être trop rugueuse ou ondulée. Ce genre d'approche peut avoir un avantage majeur sur les approches paramétriques en évitant une supposition d'une forme de fonction de  $f$  particulière, elles ont le potentiel de précision plus large que possible des dimensions de  $f$ .

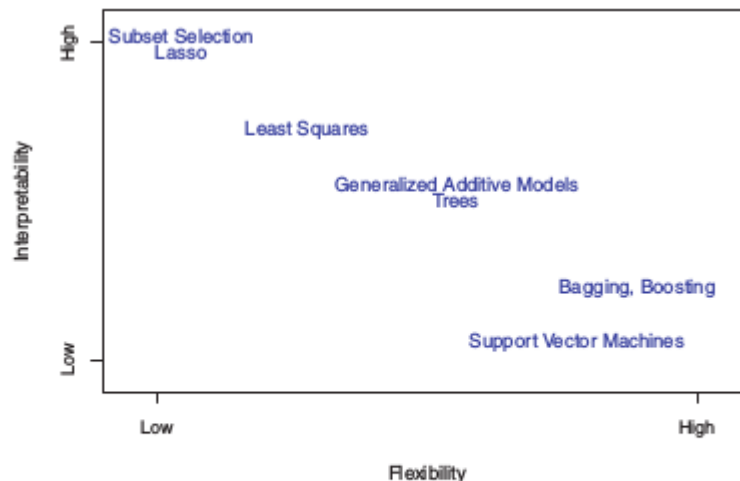
Mais l'approche non paramétrique souffre d'un désavantage majeur: elle ne réduit pas le problème d'estimation de  $f$  à un petit nombre de paramètres. Un plus grand nombre de données est requise dans le but d'obtenir une estimation précise de  $f$ .

Pour obtenir d'obtenir un ajustement proche de  $f$ , on doit utiliser un plus petit niveau de souplesse.

Comme nous avons pu le voir, il y a des avantages et des inconvénients des deux côtés d'approche paramétrique et non-paramétrique. On explorera les deux méthodes dans ce livre.

Comme nous verrons dans ce livre, nous examinerons que certaines sont moins flexibles ou plus restrictives. Par exemple, la régression linéaire est relativement inflexible, parce que elle peut générer uniquement des fonctions linéaires.

Voici un tableau de certaines techniques de l'interprétation possible en fonction de la flexibilité.



**FIGURE 2.7.** A representation of the tradeoff between flexibility and interpretability, using different statistical learning methods. In general, as the flexibility of a method increases, its interpretability decreases.

D'autres méthodes, comme thin-plate spline sont considérablement plus flexibles car elle génèrent un plus large champ de possibilités pour estimer  $f$ .

Pourquoi choisirions nous d'utiliser des méthodes plus restrictives au lieu d'approches plus flexibles ?

Il y a plusieurs raisons qui nous ferais preferer un modèle plus restrictif. Si nous sommes intéressés par l'inférence, alors les modèles restrictifs sont plus interprétables. Par exemple quand inférer est le but, le modèle linéaire peut être un bon choix puisqu'il serait plus facile de comprendre les relations entre  $Y$  et  $X_1, \dots, X_p$ . D'une autre manière, des approches plus flexibles, comme les splines et les méthodes de boosting peuvent amener à des estimations compliquées de  $f$  qui serait difficile de comprendre.

On a donc établi que quand l'inférence est le but, il y a un avantage clair d'utiliser des simples et relativement inflexibles méthodes d'apprentissage statistiques. Dans certains cas nous sommes uniquement intéressés par la prédiction et l'interprétabilité du modèle prédictif n'est pas un enjeu.

La plupart des problèmes d'apprentissage statistique tombent dans l'une des deux catégories: apprentissage supervisé ou non-supervisé. L'apprentissage non-supervisé décrit une situation plus difficile qui pour chaque observation  $i=1, \dots, n$  on constate une mesure de vecteur  $x_i$  mais non associée à une réponse  $y_i$ . Il ne devient pas possible de faire des ajustements linéaires donc il n'y a pas de variable à prédire. Dans un certain sens on voyage à l'aveugle. On peut tenter de comprendre les relations entre les variables ou entre les observations. L'une des techniques les plus utilisées est le cluster analysis, ou encore le clustering. Le but du clustering est de vérifier sur la base des  $x_1, \dots, x_n$  si les observations créent des groupes distincts. On ne peut pas s'attendre à qu'une méthode de clustering assigne tout les points au bon groupe. Par exemple si nous avons  $p$  variables dans nos données alors  $p(p-1)/2$  différents nuages de points peuvent être réalisés. Nous ne pouvons donc pas observer dans plusieurs dimensions autant de groupes. Donc des méthodes automatisées de clustering sont importantes. Beaucoup de problèmes tombent directement dans une des deux cases apprentissage supervisé ou non-supervisé. Mais des fois ce n'est pas clair. Par exemple supposons que nous avons  $n$  observations. Pour  $m$  des observations où  $m < n$ , on a tout autant des mesures de prédictions et des réponses de mesures. Pour le reste de  $n-m$  observations, on a une mesure de prédiction mais pas de réponse de mesure. On se réfère à ce genre de problème un problème d'apprentissage semi-supervisé. Dans ce cas nous souhaitons utiliser une méthode d'apprentissage statistique qui peut prendre en compte les  $m$  observations qui pour chaque réponse de mesure soit possible autant que les  $n-m$  observations qui ne le sont pas.

Les problèmes de régression contre ceux de classification

Les variables peuvent être caractérisées soit quantitatives soit qualitatives (aussi connues sous le nom de catégoriques). Les variables quantitatives prennent des valeurs numériques. Les variables qualitatives prennent des valeurs des  $K$  différentes classes, ou catégories.

On se réfère à des problèmes de régressions pour des problèmes de réponse quantitatives et à des problèmes de classifications pour des problèmes de réponse qualitatives.

La distinction n'est pas toujours évidente. La régression linéaire avec les moindres carrés est utilisée pour des réponses quantitatives. Alors que la régression logistique est typiquement utilisée pour des qualitatives. Qui est aussi souvent utilisé comme une méthode de classification. Mais quand il estime une probabilité de classe, il peut être pensé comme une méthode de régression.

Les k plus proches voisins et le boosting peuvent autant être utilisés pour des réponses quantitatives que qualitatives.

On a tendance à choisir un modèle d'apprentissage statistique en fonction de la réponse qualitative ou quantitative. Plutôt régression linéaire si quantitatif et régression logistique quand qualitatif.

Il est souvent important de savoir si le prédicteur est qualitatif ou quantitatif.

## II Evaluation de la précision du modèle

Pourquoi ne pas présenter les meilleures méthodes que celle présentées? Parce qu'aucune méthode domine les autres sur tout type de données. Choisir la meilleure approche peut s'avérer être l'un des plus grands défis de l'apprentissage statistique.

Pour évaluer la performance d'un modèle d'apprentissage statistique sur un jeu de données, nous avons besoin d'un outil de mesure pour connaître la performance. Dans le cas de régression, l'outil de mesure le plus utilisé est le Mean Squared Error (MSE)

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \quad \text{où } \hat{f}(x_i) \text{ est la prédiction que } \hat{f} \text{ donne une } i\text{ème observation. La MSE}$$

sera petite si la prédiction sera proche de la vraie réponse et sera grande si pour une des observations, la valeur prédite diffère significativement.

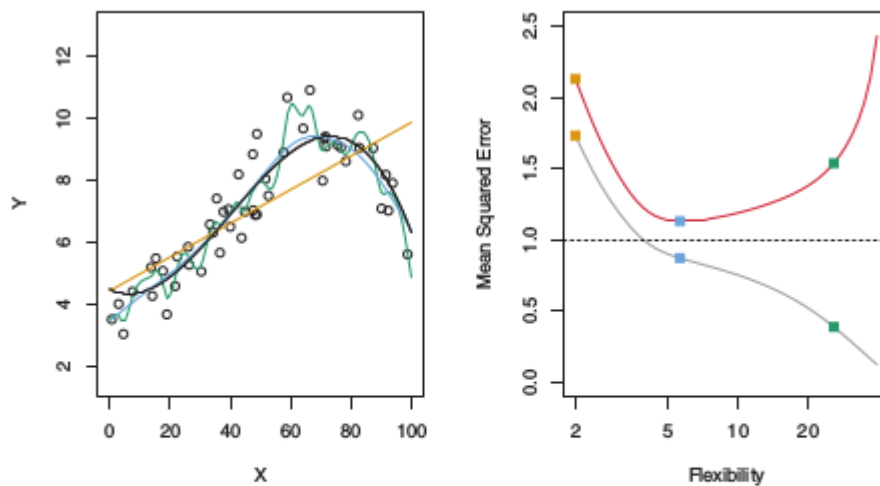
On effectue le MSE sur les données d'entraînement mais ce qui nous intéresse le plus est sur les données de test. Donc on s'intéresse à la précision de la prédiction que nous obtenons quand nous appliquons notre méthode à des données de test.

Supposons que nous faisons un ajustement de notre méthode d'apprentissage statistique sur nos observations d'entraînement  $\{(x_1, y_1), \dots, (x_n, y_n)\}$  et que nous obtenons une estimation de  $\hat{f}$ . On peut calculer  $\hat{f}(x_1), \dots, \hat{f}(x_n)$ . S'ils sont approximativement égaux à  $y_1, \dots, y_n$  alors notre MSE d'entraînement donnée par la formule plus haut est petite. Toutefois nous ne sommes plutôt pas intéressés par  $\hat{f}(x_i) \approx y_i$ , au lieu de cela on veut connaître  $\hat{f}(x_0)$  est approximativement égale à  $y_0$  où  $(x_0, y_0)$  est une observation non déjà fait par l'entraînement.

On veut choisir un modèle qui donne le plus petit test MSE contrairement au plus petit entraînement MSE. En d'autres mots si nous avons un grand nombre de données de test, on pourrait calculer

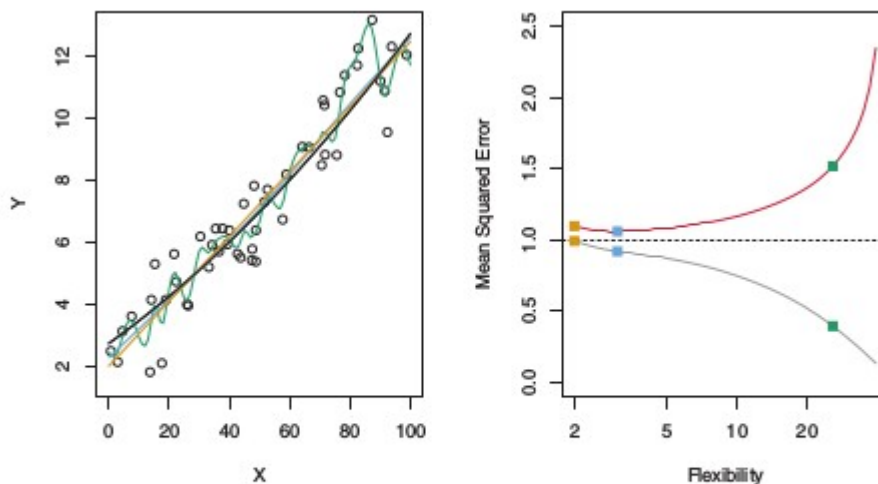
$Ave(y_0 - \hat{f}(x_0))^2$  la moyenne du carré de l'erreur de prédiction pour ce jeu de test  $(x_0, y_0)$ . Nous voudrions sélectionner le modèle pour lequel la moyenne de cette quantité – le test MSE – aussi petit que possible.

Il n'y a pas de garanties que le jeu de test aient le minimum en se basant sur le minimum du jeu d'entraînement.



**FIGURE 2.9.** Left: Data simulated from  $f$ , shown in black. Three estimates of  $f$  are shown: the linear regression line (orange curve), and two smoothing spline fits (blue and green curves). Right: Training MSE (grey curve), test MSE (red curve), and minimum possible test MSE over all methods (dashed line). Squares represent the training and test MSEs for the three fits shown in the left-hand panel.

La ligne horizontale indique  $\text{Var}(\varepsilon)$



**FIGURE 2.10.** Details are as in Figure 2.9, using a different true  $f$  that is much closer to linear. In this setting, linear regression provides a very good fit to the data.

En pratique, on peut habituellement utiliser le MSE entraînement avec relativement d'aisance. Mais estimer le MSE test est considérablement plus compliqué car souvent il n'y a pas de jeu de test. Le niveau de flexibilité correspondant au modèle avec le test MSE peut considérablement varier en pratique (en fonction des données). Dans ce livre, nous allons étudier une variété d'approches qui peuvent être utilisées en pratique pour estimer ce point minimum.

Une technique importante est la cross-validation, qui est une méthode pour estimer le test MSE en utilisant uniquement le jeu d'entraînement.

Le dilemme biais-variance

Le test MSE pour une valeur  $x_0$  donnée peut toujours être décomposé en un somme de trois quantités fondamentales:

-la variance de  $\hat{f}(x_0)$

-le carré du biais  $\hat{f}(x_0)$

-et la variance de l'erreur du terme  $\varepsilon$

$$E(y_0 - \hat{f}(x_0))^2 = \text{Var}(\hat{f}(x_0)) + [\text{Bias}(\hat{f}(x_0))]^2 + \text{Var}(\varepsilon)$$

la notation  $E(y_0 - \hat{f}(x_0))^2$  est la valeur attendue du test MSE, et se réfère à la moyenne du test MSE que l'on obtiendrait si nous répétions l'estimation de  $f$  en utilisant un grand nombre de jeu d'entraînement et testé pour  $x_0$ .

L'équation ci-dessus nous indique que pour minimiser l'erreur de test attendue, nous avons besoin de sélectionner un modèle qui accomplit la plus petite variance et le plus petit biais. A noter que la variance et le biais au carré sont deux valeurs strictement positives. Nous voyons donc que le test MSE ne peut jamais être en dessous de  $\text{Var}(\varepsilon)$  (se nommant l'erreur irréductible)

Que signifie le biais et la variance d'un modèle ?

La variance fait référence au montant que  $\hat{f}$  changerait si nous changions de jeu de données. Alors que le jeu d'entraînement est utilisé pour faire l'ajustement du modèle, différents jeux de données donneront différents  $\hat{f}$ . Si un modèle a une haute variance alors de petits changements dans les données donneront de grands changements de  $\hat{f}$ . En général les modèles les plus flexibles ont une variance élevée. On peut voir sur la figure 2,9 que la courbe verte suit les observations de manière très proche. Il a donc une haute variance car changer le jeu de données causera l'estimation de  $\hat{f}$  à changer considérablement. A l'inverse, la courbe orange est relativement inflexible et a une faible variance, car bouger n'importe quelle observation ne causera uniquement qu'un petit décalage de la position de la ligne.

Le biais fait référence à l'erreur introduite par l'approximation du réel problème, qui peut être extrêmement compliqué par rapport à un modèle simple. Par exemple une approche de régression linéaire sur un problème où la fonction réelle est plus compliquée que linéaire induira un haut biais. Généralement plus le modèle est flexible moins il a de biais.

De manière générale, plus nous utilisons des modèles flexibles, plus la variance augmentera et plus le biais diminuera.

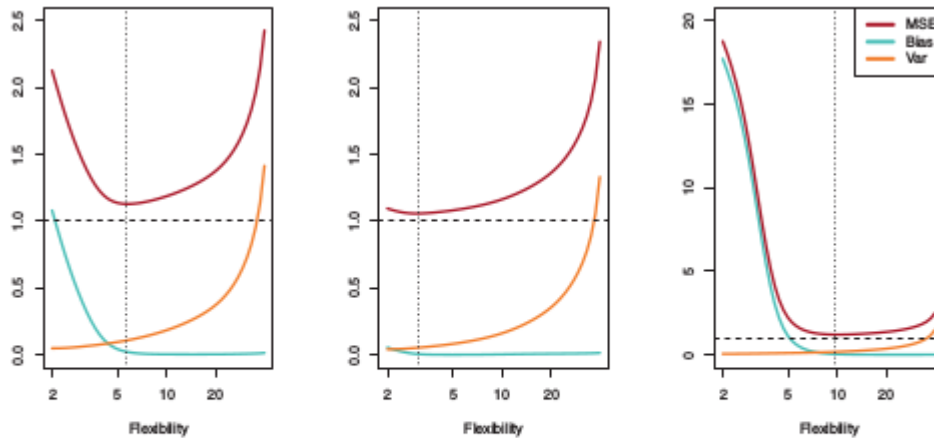
Le *biais* est l'erreur provenant d'hypothèses erronées dans l'algorithme d'apprentissage. Un biais élevé peut être lié à un algorithme qui manque de relations pertinentes entre les données en entrée et les sorties prévues (sous-apprentissage).

La *variance* est l'erreur due à la sensibilité aux petites fluctuations de l'échantillon d'apprentissage. Une variance élevée peut entraîner un surapprentissage, c'est-à-dire modéliser le bruit aléatoire des données d'apprentissage plutôt que les sorties prévues.

Donc au début le test MSE est trop sous-ajusté puis petit à petit en prenant un modèle plus flexible on voit le test MSE diminuer dû au biais qui diminue. Puis arrivé à un minimum, le test MSE réaugmente dû au surapprentissage qui est causé par l'augmentation de la variance.

VARIANCE élevée = surapprentissage

BIAIS élevé = sousapprentissage



**FIGURE 2.12.** Squared bias (blue curve), variance (orange curve),  $\text{Var}(\epsilon)$  (dashed line), and test MSE (red curve) for the three data sets in Figures 2.9–2.11. The vertical dotted line indicates the flexibility level corresponding to the smallest test MSE.

Les trois graphiques de la figure 2.12 illustrent l'équation du test MSE (voir formule séparée en plusieurs termes). On peut voir que le test MSE diffère des trois cas car le carré du biais et la variance change selon le jeu de données. Le deuxième graphique montre que  $f$  est proche d'une fonction linéaire. Pour le troisième on est clairement dans du non-linéaire.

$$\text{Biais}[\hat{f}(x)] = E[\hat{f}(x) - f(x)]$$

$$\text{Var}[\hat{f}(x)] = E[(\hat{f}(x) - E[\hat{f}(x)])^2]$$



La relation entre le biais, la variance et le test MSE donné par l'équation et affiché par la figure 2.12 est connue sous le nom de dilemme biais-variance.

Il est plus facile d'obtenir des modèles avec très peu de biais mais haute variance ou l'inverse, des modèles avec haut biais et peu de variance.

Dans la vraie vie il est rare d'avoir la fonction  $f$ . Donc il est rare de pouvoir calculer le test MSE, la variance et le biais.

La validation croisée est un bon moyen pour estimer le test MSE avec les données d'entraînement.

On a parlé presque que de régression, maintenant de classification.

Supposons que nous recherchons à estimer  $f$  sur une base de jeu de données  $\{(x_1, y_1), \dots, (x_n, y_n)\}$  où  $y_1, \dots, y_n$  sont qualitatives.

La technique la plus utilisée pour quantifier la précision de notre estimation de  $\hat{f}$  est l'error rate qui est la proportion d'erreurs faites si nous appliquons notre estimation  $\hat{f}$  aux données d'observations:

$$\frac{1}{n} \sum_{i=1}^n I(y_i \neq \hat{y}_i) \quad \text{où } \hat{y}_i \text{ est la classe prédite pour la } i\text{ème observation en utilisant } \hat{f}$$

$I(y_i \neq \hat{y}_i)$  est une variable indicatrice qui est égale à 1 si  $y_i \neq \hat{y}_i$  et 0 si  $y_i = \hat{y}_i$ .

Si  $I(y_i \neq \hat{y}_i) = 0$  alors la  $i$ ème observation a été classifiée correctement par notre modèle de classification.

Sinon il est mal classifié. Donc cette équation calcule la fraction de classification incorrectes.

Cette équation est appelée training error rate car elle calcule à partir des données d'entraînement. Comme pour la régression, on est plus intéressé par l'error rate résultant de données non d'entraînement.

Le test error rate est associé à un jeu de test de la forme  $(x_0, y_0)$  donné par

$$\text{Ave}(I(y_0 \neq \hat{y}_0))$$

Le classifieur de Bayes

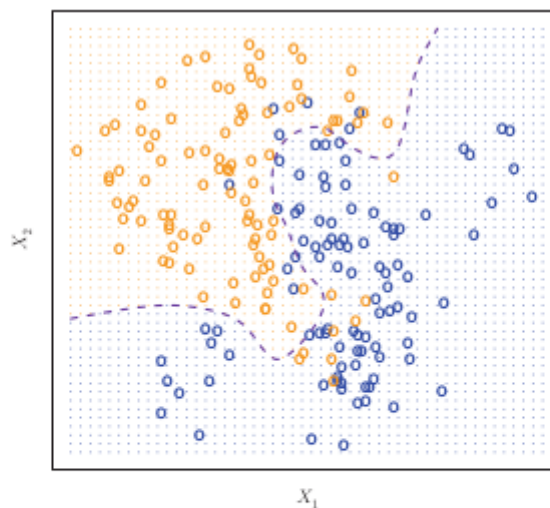
Il est possible de montrer que l'équation plus haut est minimisée en moyenne par un simple classifieur qui associe chaque observation à sa plus proche classe selon les prédictions. En d'autres mots nous souhaitons simplement assigner un test d'observation avec le prédicteur  $x_0$  de la classe  $j$  pour lequel

$Pr(Y = j | X = x_0)$  est grand. C'est une probabilité conditionnelle qui est la probabilité que  $Y = j$  selon un prédicteur de  $x_0$ . Ce très simple classifieur est appelé classifieur de Bayes.

Dans un problème à deux classes le classifieur de Bayes correspond à prédire la classe 1 si

$$Pr(Y = 1 | X = x_0) > 0.5 \quad (2.10)$$

ou la classe 2 autrement.



**FIGURE 2.13.** A simulated data set consisting of 100 observations in each of two groups, indicated in blue and in orange. The purple dashed line represents the Bayes decision boundary. The orange background grid indicates the region in which a test observation will be assigned to the orange class, and the blue background grid indicates the region in which a test observation will be assigned to the blue class.

La Figure 2.13 montre un exemple de dimension 2 ( $X_1, X_2$ ) pour un problème à deux classes. Pour chaque valeur de  $X_1$  et  $X_2$  il y a différentes probabilités de réponse en étant orange ou bleue. La limite de décision de Bayes est la ligne en pointillés.

Le classifieur de Bayes produit le plus petit test de taux d'erreur. Puisque le classifieur de Bayes choisira toujours la classe pour laquelle l'équation 2.10 est grand, le taux d'erreur à  $X=x_0$  sera

$$1 - \max_j \Pr(Y = j | X = x_0)$$

En général l'erreur globale de Bayes est donnée par  $1 - E(\max_j \Pr(Y = j | X))$  (2.11)

où l'attente moyenne la probabilité la plus probable de  $X$ . L'erreur de Bayes est analogue à l'erreur irréductible discutée plus haut.

Les k-plus-proches voisins

En théorie nous voulons toujours prédire des réponses qualitatives en utilisant le classifieur de Bayes. Mais pour des données réelles, nous ne connaissons pas la distribution conditionnelle de  $Y$  sachant  $X$ , et donc le classifieur de Bayes est impossible. Beaucoup d'approches tentent d'estimer la distribution conditionnelle de  $Y$  sachant  $X$ , et classifient une observation donnée à la classe dont l'estimation de la probabilité est la plus élevée. Un de ces modèles est le classifieur K plus proches voisins (KNN).

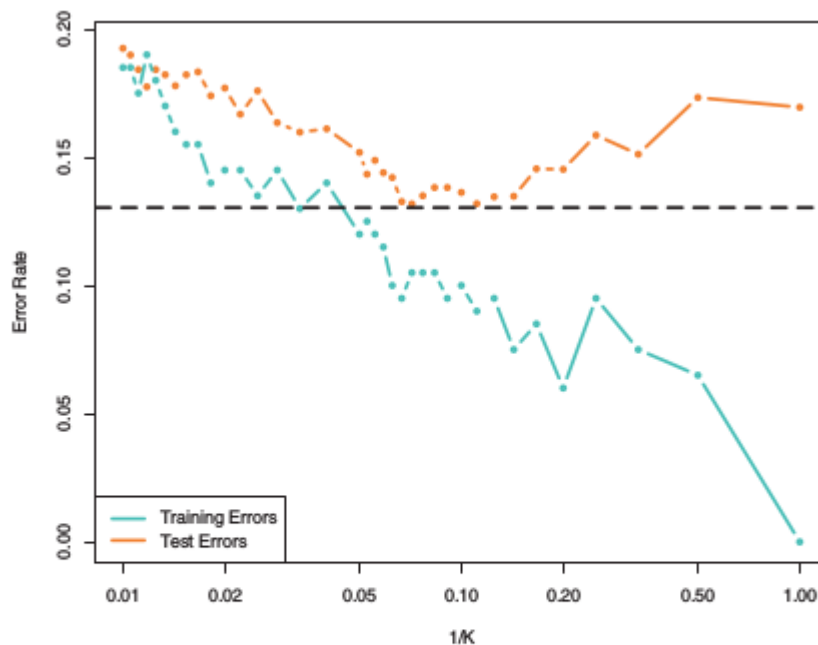
Étant donné un entier  $K$  et une observation  $x_0$ , les KNN identifient les  $K$  points dans les données d'entraînement qui sont les plus proches de  $x_0$ , représenté par  $N_0$ .

$$\Pr(Y = j | X = x_0) = \frac{1}{K} \sum_{i \in N_0} I(y_i = j)$$

En conclusion, KNN applique les règles de Bayes et classifie l'observation  $x_0$  à tester à la classe ayant la plus grande probabilité.

Malgré le fait que c'est une approche très simple, KNN peut souvent produire des classificateurs qui sont proches du classifieur de Bayes optimal.

Le choix de  $K$  est un grand effet sur le classifieur KNN obtenu.



**FIGURE 2.17.** The KNN training error rate (blue, 200 observations) and test error rate (orange, 5,000 observations) on the data from Figure 2.13, as the level of flexibility (assessed using  $1/K$ ) increases, or equivalently as the number of neighbors  $K$  decreases. The black dashed line indicates the Bayes error rate. The jumpiness of the curves is due to the small size of the training data set.

Pour  $K=1$ , la limite de décision est très flexible et trouve des schémas dans les données qui ne correspondent pas à la limite de Bayes.

Quand  $K$  augmente, le modèle devient de moins en moins flexible et produit une limite de décision proche

du linéaire : petite variance mais grand biais

Quand  $1/K$  augmente, le modèle devient plus flexible. Courbe en forme de U.

Dans les deux cas de classification et de régression, choisir le niveau de flexibilité est critique pour le succès de n'importe quel modèle statistique.

Le dilemme biais-variance et les résultats d'une courbe en U pour l'erreur de test, peut devenir une tâche difficile.

### Chapitre 3 – Régression linéaire

Supposons la relation linéaire  $Y \approx \beta_0 + \beta_1 X$  (3.1)

Une fois que nous avons utilisé nos données d'entraînement pour estimer  $\hat{\beta}_0$  et  $\hat{\beta}_1$  nous pouvons faire des prédictions  $\hat{y} \approx \hat{\beta}_0 + \hat{\beta}_1 x$  où  $\hat{y}$  indique la prédiction de Y sur la base de  $X=x$ .

Le chapeau ^ désigne la valeur estimée pour un paramètre inconnu ou une prédiction de la réponse.

Le plus commun est d'utiliser le critère des moindres carrés.

Posons  $\hat{y}_i \approx \hat{\beta}_0 + \hat{\beta}_1 x_i$  la prédiction de Y basée sur la ième valeur de X.

Donc  $e_i = y_i - \hat{y}_i$  représente le ième résidu.

Nous définissons les résidues de sum of squares (RSS)

$$RSS = e_1^2 + e_2^2 + \dots + e_n^2$$

$$RSS = (y_1 - \hat{\beta}_0 - \hat{\beta}_1 x_1)^2 + (y_2 - \hat{\beta}_0 - \hat{\beta}_1 x_2)^2 + \dots + (y_n - \hat{\beta}_0 - \hat{\beta}_1 x_n)^2$$

L'approche par les moindres carrés choisit  $\hat{\beta}_0$  et  $\hat{\beta}_1$  pour minimiser le RSS.

En faisant quelques calculs, on peut montrer que les minimiseurs sont

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (3.4.a)$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (3.4.b)$$

où  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$  et  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$  qui sont les moyennes

les équations (3.4) définissent les estimations des coefficients des moindres carrés pour une régression linéaire

démonstration

On cherche RSS tel que

$$f(\hat{\beta}_0, \hat{\beta}_1) = \min_{\hat{\beta}_0, \hat{\beta}_1} \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$$

on atteint le minimum pour  $\frac{\partial f(\hat{\beta}_0, \hat{\beta}_1)}{\partial \beta_k} = 0$

$$\text{donc } -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i = 0 \quad \text{et} \quad -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0$$

$$\text{donc } \sum_{i=1}^n (y_i - \hat{\beta}_1 x_i) x_i = n \bar{x} \hat{\beta}_0 \quad \text{et} \quad \bar{y} - \hat{\beta}_0 - \hat{\beta}_1 \bar{x} = 0$$

$$\text{donc } \sum_{i=1}^n (y_i x_i) - \hat{\beta}_1 \sum_{i=1}^n x_i^2 = n \bar{x} (\bar{y} - \hat{\beta}_1 \bar{x}) \quad \text{et} \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\text{donc } \sum_{i=1}^n (y_i x_i) - \hat{\beta}_1 (\sum_{i=1}^n x_i^2 - n \bar{x}^2) = n \bar{x} \bar{y} \quad \text{et} \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\text{donc } \sum_{i=1}^n (y_i x_i) - n \bar{x} \bar{y} - \hat{\beta}_1 (\sum_{i=1}^n x_i^2 - n \bar{x}^2) = 0 \quad \text{et} \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\text{or } E[(X - E[X])^2] = E[X^2] - E[X]^2 \quad \text{et} \quad E[(X - E[X])(Y - E[Y])] = E[XY] - E[X]E[Y]$$

donc 
$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \text{et} \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

On peut écrire 
$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i)}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum_{i=1}^n (y_i)(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (\text{voir cela en développant})$$

On suppose que la relation entre X et Y prend la forme  $Y = f(X) + \varepsilon$  pour une fonction f inconnue. Si f est approximé par une fonction linéaire, alors on peut écrire  $Y = \beta_0 + \beta_1 X + \varepsilon$  (3.5)  
On suppose que le terme d'erreur est indépendant de X.

Avec les coefficients de (3.4) on peut tracer la ligne des moindres carrés.

Sur un grand jeu de données nous pouvons considérer que  $\hat{\beta}_0$  et  $\hat{\beta}_1$  seront égales exactement à  $\beta_0$  et  $\beta_1$

Par exemple pour vérifier la précision de  $\hat{\mu}$  étant l'estimation de la moyenne  $\mu$  sur un jeu de données, on peut étudier l'erreur standard de  $\hat{\mu}$  écrite  $SE(\hat{\mu})$

$$Var(\hat{\mu}) = SE(\hat{\mu})^2 = \frac{\sigma^2}{n} \quad (3.7) \quad \text{où } \sigma \text{ est la déviation standard de chaque réalisations de } y_i \text{ de } Y.$$

Démonstration de  $Var(\hat{\mu}) = \frac{\sigma^2}{n}$  :

On a 
$$E(\bar{\mu}) = E\left[\frac{1}{n} \sum_{i=1}^n \mu_i\right] = \frac{1}{n} \sum_{i=1}^n E[\mu_i] = \frac{n * \mu}{n} = \mu \quad \text{car on a} \quad \bar{\mu} = \frac{1}{n} \sum_{i=1}^n \mu_i$$

$$Var(\hat{\mu}) = Var(\bar{y}) = E[(\bar{y} - E[\bar{y}])^2]$$

$$Var(\hat{\mu}) = E\left[\left(\frac{1}{n} \sum_{i=1}^n y_i - E\left[\frac{1}{n} \sum_{i=1}^n y_i\right]\right)^2\right]$$

$$Var(\hat{\mu}) = \frac{1}{n^2} E\left[\left(\sum_{i=1}^n y_i - E[y_i]\right)^2\right] \quad \text{on pose} \quad a_i = y_i - E[y_i]$$

$$Var(\hat{\mu}) = \frac{1}{n^2} E\left[\left(\sum_{i=1}^n a_i\right)^2\right]$$

$$Var(\hat{\mu}) = \frac{1}{n^2} E\left[\left(\sum_{i=1}^n a_i\right)\left(\sum_{j=1}^n a_j\right)\right]$$

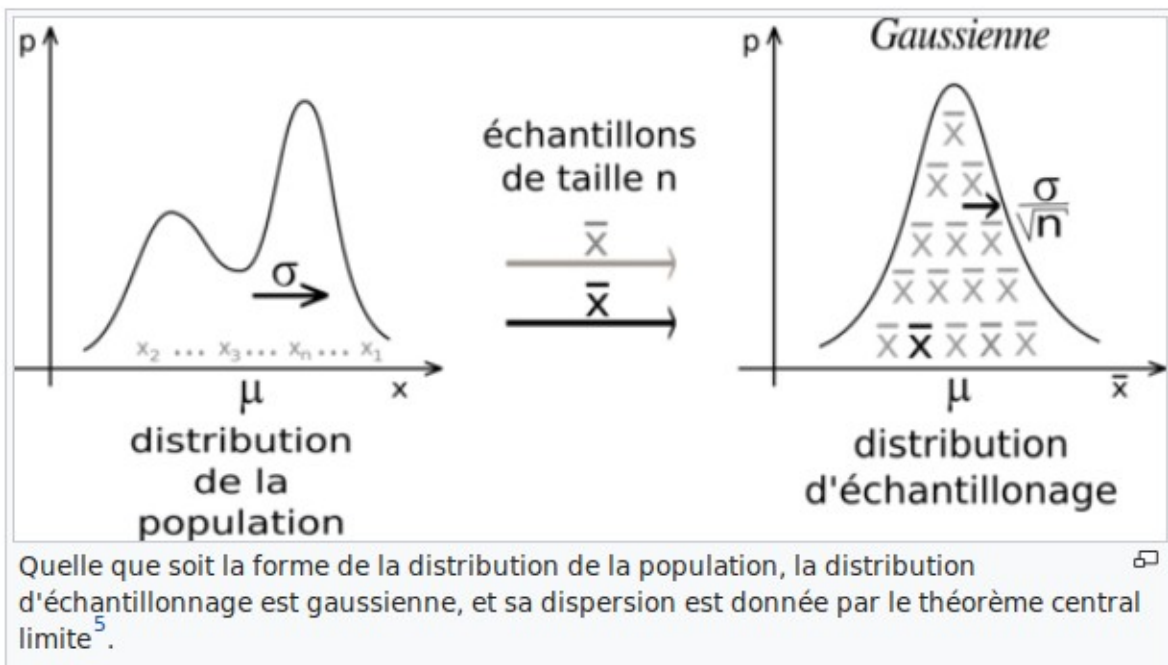
$$Var(\hat{\mu}) = \frac{1}{n^2} E\left[\sum_{i=1}^n a_i^2 + 2 \sum_{1 \leq i < j \leq n} a_i a_j\right]$$

$$Var(\hat{\mu}) = \frac{1}{n^2} \left(\sum_{i=1}^n E[a_i^2] + 2 \sum_{1 \leq i < j \leq n} E[a_i a_j]\right)$$

$$Var(\hat{\mu}) = \frac{1}{n^2} \left(\sum_{i=1}^n E[a_i^2] + 2 \sum_{1 \leq i < j \leq n} E[a_i a_j]\right) \quad \text{or} \quad E[a_i] = E[y_i - E[y_i]] = E[y_i] - E[y_i] = 0$$

$$\begin{aligned} \text{Var}(\hat{\mu}) &= \frac{1}{n^2} \sum_{i=1}^n E[(y_i - E[y_i])^2] \\ \text{Var}(\hat{\mu}) &= \frac{1}{n} E\left[\frac{1}{n} \sum_{i=1}^n (y_i - E[y_i])^2\right] \\ \text{Var}(\hat{\mu}) &= \frac{1}{n} E[\sigma^2] \\ \text{Var}(\hat{\mu}) &= \frac{\sigma^2}{n} \end{aligned}$$

On peut le voir graphiquement avec le théorème de la limite centrale



$$\begin{aligned} \text{cov}(y_i, y_j) &= \frac{1}{n^2} \sum_{i=1}^n (y_i - E[y_i]) \sum_{j=1}^n (y_j - E[y_j]) = \frac{1}{n^2} \sum_{i=1}^n (\varepsilon_i - E[\varepsilon_i]) \sum_{j=1}^n (\varepsilon_j - E[\varepsilon_j]) = \text{cov}(\varepsilon_i, \varepsilon_j) \\ \text{cov}(y_i, y_j) &= 0 \text{ si } j \neq i \text{ sinon } \text{var}(\varepsilon_i, \varepsilon_i) = \text{var}(y_i, y_j) = \sigma^2 \end{aligned}$$

L'équation 3.7 nous dit comment la déviation diminue avec n. Plus on a d'observations, plus l'erreur standard devient petite. De manière similaire on peut désirer combien  $\hat{\beta}_0$  et  $\hat{\beta}_1$  sont proches des valeurs

$\beta_0$  et  $\beta_1$ . Pour calculer l'erreur standard associé à  $\hat{\beta}_0$  et  $\hat{\beta}_1$  on utilise les formules suivantes

$$SE(\hat{\beta}_0)^2 = \sigma^2 \left[ \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right] \quad \text{et} \quad SE(\hat{\beta}_1)^2 = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

où  $\sigma^2 = \text{Var}(\varepsilon)$

Demonstration:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x}) y_i}{\sum_{i=1}^n (x_i - \bar{x})^2} = \sum_{i=1}^n w_i Y_i \quad \text{avec} \quad w_i = \frac{x_i - \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

de même  $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = \sum_{i=1}^n \left( \frac{1}{n} - \bar{x} w_i \right) y_i$

on a  $\sum_{i=1}^n w_i = \sum_{i=1}^n \frac{x_i - \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} = 0$

on a  $\sum_{i=1}^n w_i x_i = \sum_{i=1}^n \frac{(x_i - \bar{x}) x_i}{\sum_{i=1}^n (x_i - \bar{x})^2} = \sum_{i=1}^n \frac{(x_i - \bar{x})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} = 1$

on a  $\sum_{i=1}^n w_i^2 = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{\left( \sum_{i=1}^n (x_i - \bar{x})^2 \right)^2} = \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2}$

donc  $E[\hat{\beta}_1] = E\left[\sum_{i=1}^n w_i y_i\right] = E\left[\sum_{i=1}^n w_i (\beta_0 + \beta_1 x_i + \varepsilon_i)\right] = \beta_1$

$$E[\hat{\beta}_0] = E[\bar{y} - \hat{\beta}_1 \bar{x}] = E[\beta_0 \bar{y} + \beta_1 \bar{x} + \bar{\varepsilon}] - \beta_1 \bar{x} = \beta_0$$

$$Var(\hat{\beta}_0) = \sum_{i=1}^n \left( \frac{1}{n^2} - 2 \frac{\bar{x} w_i}{n} + \bar{x}^2 w_i^2 \right) \sigma^2 = \left( \frac{1}{n} - 0 + \bar{x}^2 \sum_{i=1}^n w_i^2 \right) \sigma^2 = \left( \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) \sigma^2$$

$$Var(\hat{\beta}_1) = \sum_{i=1}^n w_i^2 \sigma^2 = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

(REFAIRE LES CALCULS A LA MAIN)

$$Cov(\hat{\beta}_1, \hat{\beta}_0) = Cov\left(\sum_{i=1}^n \left(\frac{1}{n} - \bar{x} w_i\right) y_i, \sum_{i=1}^n w_i y_i\right) = \sum_{i=1}^n \sum_{j=1}^n \left(\frac{1}{n} - \bar{x} w_i\right) w_j Cov(y_i, y_j)$$

$$Cov(\hat{\beta}_1, \hat{\beta}_0) = \sum_{i=1}^n \left(\frac{1}{n} - \bar{x} w_i\right) w_i \sigma^2 = \frac{-\bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} \sigma^2$$

En général  $\sigma^2$  n'est pas connu, mais il peut être estimé depuis les données.

L'estimation de  $\sigma$  est connu sous le nom de résidu d'erreur standard et est donné par la formule

$$RSE = \sqrt{RSS/(n-2)} \quad \text{(ZT.35)}$$

démonstration:

on appelle résidus les quantités  $\hat{\varepsilon}_i = Y_i - \hat{Y}_i$   $\sum_{i=1}^n \hat{\varepsilon}_i = 0$  car  $E[\hat{\varepsilon}_i] = 0$

estimer la variance  $\sigma^2$  par  $\frac{1}{n} \sum_{i=1}^n \hat{\varepsilon}_i^2$  mais c'est un estimateur biaisé on prendra donc

$$\sigma^2 = \frac{1}{n-2} \sum_{i=1}^n \hat{\varepsilon}_i^2 \text{ qui est un estimateur non biaisé}$$

$$\hat{\varepsilon}_i = Y_i - \hat{Y}_i$$

$$\hat{\varepsilon}_i = \beta_0 + \beta_1 x_i + \varepsilon_i - \hat{\beta}_0 - \hat{\beta}_1 x_i$$

$$\hat{\varepsilon}_i = \bar{Y} - \beta_1 \bar{x} - \bar{\varepsilon} + \beta_1 x_i + \varepsilon_i - \bar{Y} + \hat{\beta}_1 \bar{x} - \hat{\beta}_1 x_i$$

$$\hat{\varepsilon}_i = (\varepsilon_i - \bar{\varepsilon}) + (\beta_1 - \hat{\beta}_1)(x_i - \bar{x})$$

$$\text{donc } \hat{\varepsilon}_i^2 = (\varepsilon_i - \bar{\varepsilon})^2 + 2(\varepsilon_i - \bar{\varepsilon})(\beta_1 - \hat{\beta}_1)(x_i - \bar{x}) + (\beta_1 - \hat{\beta}_1)^2(x_i - \bar{x})^2$$

or

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x}) y_i}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})(\beta_0 + \beta_1 x_i + \varepsilon_i)}{\sum_{i=1}^n (x_i - \bar{x})^2} = \beta_1 + \frac{\sum_{i=1}^n (x_i - \bar{x})(\varepsilon_i - \bar{\varepsilon})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\text{donc } \sum_{i=1}^n \hat{\varepsilon}_i^2 = \sum_{i=1}^n (\varepsilon_i - \bar{\varepsilon})^2 + \sum_{i=1}^n 2(\varepsilon_i - \bar{\varepsilon})(\beta_1 - \hat{\beta}_1)(x_i - \bar{x}) + \sum_{i=1}^n (\beta_1 - \hat{\beta}_1)^2(x_i - \bar{x})^2$$

$$\sum_{i=1}^n \hat{\varepsilon}_i^2 = \sum_{i=1}^n (\varepsilon_i - \bar{\varepsilon})^2 - 2 \sum_{i=1}^n (\beta_1 - \hat{\beta}_1)^2(x_i - \bar{x})^2 + \sum_{i=1}^n (\beta_1 - \hat{\beta}_1)^2(x_i - \bar{x})^2$$

$$\sum_{i=1}^n \hat{\varepsilon}_i^2 = \sum_{i=1}^n (\varepsilon_i - \bar{\varepsilon})^2 - \sum_{i=1}^n (\beta_1 - \hat{\beta}_1)^2(x_i - \bar{x})^2$$

$$E\left[\sum_{i=1}^n \hat{\varepsilon}_i^2\right] = E\left[\sum_{i=1}^n (\varepsilon_i - \bar{\varepsilon})^2\right] - E\left[\sum_{i=1}^n (\beta_1 - \hat{\beta}_1)^2(x_i - \bar{x})^2\right]$$

or

$$E\left[\sum_{i=1}^n (\varepsilon_i - \bar{\varepsilon})^2\right] = E\left[-n\bar{\varepsilon}^2 + \sum_{i=1}^n \varepsilon_i^2\right] = \frac{-n}{n^2} E\left[\left(\sum_{i=1}^n \varepsilon_i\right)^2\right] + n\sigma^2 = n\sigma^2 - \frac{1}{n} \sum_{i=1}^n E[\varepsilon_i^2] - 2 \sum_{1 \leq i < j \leq n} E[\varepsilon_i \varepsilon_j]$$

$$E\left[\sum_{i=1}^n (\varepsilon_i - \bar{\varepsilon})^2\right] = (n-1)\sigma^2 \quad (\text{A.64})$$

et

$$E\left[\sum_{i=1}^n (\beta_1 - \hat{\beta}_1)^2(x_i - \bar{x})^2\right] = E[(\beta_1 - \hat{\beta}_1)^2] E\left[\sum_{i=1}^n (x_i - \bar{x})^2\right] \quad (\text{APK.45})$$

$$E\left[\sum_{i=1}^n (\beta_1 - \hat{\beta}_1)^2(x_i - \bar{x})^2\right] = \text{var}(\hat{\beta}_1) E\left[\sum_{i=1}^n (x_i - \bar{x})^2\right] = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} E\left[\sum_{i=1}^n (x_i - \bar{x})^2\right] = \sigma^2$$

$$\text{donc } E\left[\sum_{i=1}^n \hat{\varepsilon}_i^2\right] = (n-2)\sigma^2$$



Pour les régression linéaires, il y a approximativement 95% de chances que l'intervale  $[\hat{\beta}_1 - 2 \cdot SE(\hat{\beta}_1), \hat{\beta}_1 + 2 \cdot SE(\hat{\beta}_1)]$  contiennent la vraie valeur de  $\beta_1$ . Idem pour  $\beta_0$ .

Les erreurs standards peuvent aussi être utilisées pour permettre des tests d'hypothèses sur les coefficients.

La plus connue est l'hypothèse nulle

$H_0$  : il n'y a pas de relation entre X et Y

contre l'hypothèse alternative

$H_a$  : il y a une relation entre X et Y

Mathématiquement cela correspond à tester

$H_0 : \beta_1 = 0$

contre

$H_a : \beta_1 \neq 0$

si  $\beta_1 = 0$  alors le modèle  $Y = \beta_0 + \beta_1 X + \varepsilon$  se réduit à  $Y = \beta_0 + \varepsilon$  et X est non associé à Y. Pour tester l'hypothèse nulle, nous avons besoin de déterminer si  $\hat{\beta}_1$  est suffisamment loin de zéro pour considérer que  $\beta_1$  est non nulle. De combien ? Tout dépend de la précision de  $\hat{\beta}_1$  qui dépend de  $SE(\hat{\beta}_1)$ .

En pratique on calcule une t-statistique  $t = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)}$  (3.14) qui mesure le nombre de déviation standard que

$\hat{\beta}_1$  est loin de 0.

**S'il n'y a pas de relation entre Y et X alors on s'attendra que (3.14) aura une t-distribution de n-2 degrés de liberté.**

La t-distribution a une forme en cloche et pour des valeurs de n plus grand approximativement que 30 il est quasi similaire à la distribution normale. En conséquence, il est simple de calculer la probabilité d'observer n'importe quel nombre égale à |t| ou plus grand en valeur absolue, en supposition de  $\beta_1 = 0$ . On appelle cette probabilité la p-valeur ( $p = P(X|H_0)$ ). On interprète la p-valeur comme suit: une petite p-valeur indique qu'il est peut probable d'observer une association entre le prédicteur et la réponse, en l'absence de n'importe quelle raison l'association entre le prédicteur et la réponse.

Si t est petit alors on est proche de la moyenne nulle, donc toute probabilité  $P(Y > |t|)$  doit être grande pour ne pas invalider l'hypothèse nulle.

Par conséquent, si nous observons une petite p-valeur alors on peut inférer qu'il y a une association entre le prédicteur et la réponse. Nous rejetons l'hypothèse nulle, nous déclarons qu'il existe donc une relation entre X et Y- bien sûr si la p-valeur est suffisamment petite. De manière générale nous rejetons l'hypothèse nulle pour moins de 5 ou 1%. Quand n=30, cela correspond à une t-statistique autour de 2 et 2,75 respectivement.

Par exemple:

	Coefficient	Std. error	t-statistic	p-value
<b>Intercept</b>	7.0325	0.4578	15.36	< 0.0001
<b>TV</b>	0.0475	0.0027	17.67	< 0.0001

**TABLE 3.1.** For the **Advertising** data, coefficients of the least squares model for the regression of number of units sold on TV advertising budget. An increase of \$1,000 in the TV advertising budget is associated with an increase in sales by around 50 units (Recall that the **sales** variable is in thousands of units, and the **TV** variable is in thousands of dollars).

A noter que les coefficients sont très grands devant les erreurs standards. Donc la t-statistique est grande. Donc la p-valeur est très petite. Ce qui nous conduit à dire que la probabilité de voir ce genre de valeur si  $H_0$  est vraie est virtuellement faux. Par conséquent on peut conclure que  $\beta_1 \neq 0$  et  $\beta_0 \neq 0$

Une fois que nous avons rejeté l'hypothèse nulle (3.12) en faveur de l'hypothèse alternative, il est naturel de vouloir quantifier dans quelle mesure le modèle correspond aux données.

La qualité d'une régression linéaire est évaluée par l'utilisation de deux valeurs: la résidual standard error (RSE) et le  $R^2$ -statistique.

### L'erreur standard des résidus

un rappel du modèle  $Y = \beta_0 + \beta_1 X + \varepsilon$  qui est associé avec chaque observation un terme d'erreur  $\varepsilon$ . A cause de la présence de ces terme d'erreur, même si nous suivons la vraie ligne de régression nous ne pourrions pas être capable de prédire parfaitement Y à partir de X. La RSE est une estimation de la déviation standard de  $\varepsilon$ . C'est la moyenne des déviations des réponses provenant de la vraie ligne de régression.

$$RSE = \sqrt{\frac{1}{n-2} RSS} = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (3.15) \quad \text{avec} \quad RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (3.16)$$

Si  $\hat{y}_i \approx y_i$  pour  $i=1, \dots, n$  alors (3.15) sera petit et nous pourrions conclure que le modèle correspond bien aux données. Par contre, si  $\hat{y}_i$  est très loin de  $y_i$  pour une ou plusieurs observations, alors la RSE sera grande. Ceci indique que le modèle ne correspond pas aux données.

### $R^2$ statistique

La RSE fournit une mesure du manque d'ajustement du modèle (3.5) aux données. Mais il n'est jamais clair que la RSE est bonne ou non. La  $R^2$  statistique fournit une alternative de mesure de l'ajustement. Elle prend la forme d'une proportion et prend ses valeurs entre 0 et 1 et est indépendante des dimensions de Y (contrairement à la RSE)

$$R^2 = \frac{TSS - RSS}{TSS} = 1 - \frac{RSS}{TSS} \quad \text{où} \quad TSS = \sum (y_i - \bar{y})^2 \quad TSS = \text{Total Sum of Squares}$$

**La TSS mesure la variance totale de la réponse Y, et peut être calculée avant la régression.**

**En revanche, la RSS mesure la quantité qui reste inexpliquée après avoir effectué la régression.**

**Par conséquent, TSS – RSS mesure la quantité de variabilité de la réponse qui est expliquée (ou supprimée) en effectuant la régression, et  $R^2$  mesure la proportion de la variabilité de Y expliquée par X.**

Un  $R^2$  proche de 1 indique qu'une grande partie de la variabilité de la réponse a été expliquée par la régression. Un nombre proche de 0 indique que le contraire. Cela peut se produire si le modèle n'est pas bon, ou que  $\sigma^2$  est grand ou les deux.

La  $R^2$  statistique a une interprétabilité meilleure que la RSE car elle est entre 0 et 1.

Quoi qu'il en soit il peut être intéressant de savoir si  $R^2$  est une bonne valeur et en général cela dépendra de l'application.

La  $R^2$  statistique est une mesure de la relation linéaire entre X et Y.

$$Cor(X, Y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \quad \text{est aussi une mesure de la relation linéaire entre X et Y.}$$

Cela suppose que nous pourrions être en mesure d'utiliser  $r = Cor(X, Y)$  au lieu de  $R^2$ . En fait, on pourrait montrer que pour une simple régression linéaire  $r^2 = R^2$ . En d'autres mots, la corrélation au carré et la  $R^2$  statistique sont identiques.

## Regression lineaire multiple

En général on a plusieurs paramètres pour la regression lineaire. Ce qu'on fait c'est qu'on ajoute des coefficients

$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \varepsilon$  où  $X_j$  représente le jème predicteur et  $\beta_j$  quantifie l'association entre cette variable et la réponse. On interprète  $\beta_j$  comme l'effet moyen sur  $Y$  pour une valeur croissante en  $X_j$  en gardant les autres predicteurs fixés.

Comme précédemment, les coefficients de regression  $\beta_1, \beta_2, \dots, \beta_p$  ne sont pas connus et nous devons les estimer. Nous avons donc les estimations  $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$ . On peut faire les predictions avec la formule

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p \quad (3.21)$$

Les parametres sont estimés avec encore les moindres carrés.

On choisit donc  $\beta_1, \beta_2, \dots, \beta_p$  pour minimiser la somme des residus au carré

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$
$$RSS = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \dots - \hat{\beta}_p x_{ip})^2 \quad (3.22)$$

Les valeurs  $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$  qui minimisent (3.22) sont les multiples coefficients de la régression lineaire simple.

Quand nous utilisons la regression lineaire multiple, nous sommes souvent intéressés à répondre à ces questions.

1-/ est ce qu'au moins un des predicteurs  $X_1, \dots, X_p$  est utile pour prédire la réponse ?

2-/ Est ce que tous les predicteurs aident à expliquer  $Y$ , ou est ce qu'il y a seulement un sous ensemble utile pour predire ?

3-/ Comment le modèle s'ajuste bien par rapport aux données ?

4-/ Suivant un ensemble de valeurs prédites, quelle valeur réponse devrions nous predire et combien notre prédiction est précise ?

Première question : y a-t-il une relation entre la réponse et les predicteurs ?

Pour la régression lineaire multiple avec  $p$  predicteurs, nous avons besoin de se demander si tous les coefficients sont nuls. Donc si  $\beta_1 = \beta_2 = \dots = \beta_p = 0$ . Comme la regression lineaire simple, nous utilisons l'hypothèse de test pour répondre à la question. Nous testons l'hypothèse:

$H_0: \beta_1 = \beta_2 = \dots = \beta_p = 0$  contre l'hypothèse  $H_a$ : au moins un  $\beta_j$  est non nul

Cette hypothèse est évaluée avec la F-statistique:

$$F = \frac{(TSS - RSS)/p}{RSS/(n - p - 1)} \quad (\text{pourquoi utilise-t-on RSS en bas})$$

avec  $RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$  et  $TSS = \sum_{i=1}^n (y_i - \bar{y})^2$

Si les suppositions d'un modèle lineaire sont correctes on peut montrer que  $E[RSS/(n - p - 1)] = \sigma^2$

démonstration:

$$E[RSS] = E[\|\hat{\varepsilon}\|^2] = E[\hat{\varepsilon}'\hat{\varepsilon}] = E[\text{tr}(\hat{\varepsilon}'\hat{\varepsilon})] = E[\text{tr}(\hat{\varepsilon}\hat{\varepsilon}')] \quad \text{car } \text{tr}(ABC) = \text{tr}(BCA) = \text{tr}(CAB)$$

$$E[RSS] = \text{tr}(E[\hat{\varepsilon}\hat{\varepsilon}']) = \text{tr}(\text{Var}(\hat{\varepsilon}))$$

$$\text{or } \text{Var}(\hat{\varepsilon}) = \text{Var}(Y - \hat{Y})$$

$$\text{or } \hat{Y} = X\hat{\beta}$$

$$\text{or } \hat{\beta} = (X^T X)^{-1} X^T Y \quad (\text{point critique du lagrangien } L)$$

$$\text{donc } \hat{Y} = X(X^T X)^{-1} X^T Y$$

donc  $\text{Var}(\hat{\varepsilon}) = \text{Var}((I_n - X(X^T X)^{-1} X^T)Y) = (I_n - X(X^T X)^{-1} X^T) \text{Var}(Y) = (I_n - X(X^T X)^{-1} X^T) \sigma^2$   
donc  $\text{Var}(\hat{\varepsilon}) = \text{Var}((I_n - X(X^T X)^{-1} X^T)Y) = (I_n - X(X^T X)^{-1} X^T) \text{Var}(Y) = (I_n - X(X^T X)^{-1} X^T) \sigma^2$   
donc  $E[\text{RSS}] = \sigma^2 [\text{tr}(I_n) - \text{tr}(X^T X (X^T X)^{-1})] = \sigma^2 [\text{tr}(I_n) - \text{tr}(I_{p+1})] = \sigma^2 (n - p - 1)$   
donc  $\sigma^2 = E[\text{RSS}/(n - p - 1)]$  donc un estimateur sans biais est donné par  $\hat{\sigma}^2 = \text{RSS}/(n - p - 1)$   
dans cette démonstration on a fait l'hypothèse que l'on a effectivement un modèle linéaire

Si  $H_0$  est vrai on a  $E[(TSS - \text{RSS})/p] = \sigma^2$

démonstration:

$$E[(TSS - \text{RSS})/p] = E\left[\sum_{i=1}^n (y_i - \bar{y})^2 - \sum_{i=1}^n (y_i - \hat{y}_i)^2\right]/p$$

$$E[(TSS - \text{RSS})/p] = 1/p \left(\sum_{i=1}^n E[(\varepsilon_i - \bar{\varepsilon})^2] - \sigma^2 (n - p - 1)\right) \quad \text{si on est dans l'hypothèse } H_0$$

$$E[(TSS - \text{RSS})/p] = 1/p ((n - 1)\sigma^2 - \sigma^2 (n - p - 1)) = \sigma^2 \quad \text{voir formule (A.64)}$$

Par conséquent, lorsqu'il n'y a pas de relation entre Y et X nous nous attendons à que la F-statistique soit égale à 1.

Dans le cas où  $H_a$  est vrai alors  $E[(TSS - \text{RSS})/p] > \sigma^2$ , donc on s'attend à que la F-statistique soit supérieure à 1.

De façon évidente des termes au carré en plus donneront  $E[(TSS - \text{RSS})/p] > \sigma^2$

De combien la largeur de la F-statistique doit être avant de rejeter  $H_0$  et conclure qu'il y a une relation ? On voit bien que cela dépend des valeurs n et p. Quand n est grand, une F-statistique qui est un peu plus grand que 1 peut encore fournir l'évidence contre  $H_0$ .

Dans le cas contraire, une plus grande F-statistique est requise pour rejeter  $H_0$  si n est petit.

Quand  $H_0$  est vrai et que les erreurs  $\varepsilon_i$  sont une distribution normale, la F-statistique suit une F-distribution.

Pour des valeurs n et p données, n'importe quel logiciel peut être utilisé pour calculer la p-valeur associée avec la F-distribution. En se basant sur la p-valeur on peut déterminer si oui ou non on rejette  $H_0$ . Des fois nous souhaitons tester un sous ensemble particulier de q coefficients nuls. Ce qui correspond à

$H_0 : \beta_{p-q+1} = \beta_{p-q+2} = \dots = \beta_p = 0$ . Dans ce cas on ajuste un second modèle qui utilise toutes les variables exceptés les derniers q. On suppose que la RSS de ce modèle est  $\text{RSS}_0$ . Donc la F-statistique appropriée est

$$F = \frac{(\text{RSS}_0 - \text{RSS})/q}{\text{RSS}/(n - p - 1)}$$

#### TEXTE A REVOIR POUR COMPRENDRE

**Given these individual p-values for each variable, why do we need to look at the overall F-statistic? After all, it seems likely that if any one of the p-values for the individual variables is very small, then at least one of the predictors is related to the response. However, this logic is flawed, especially when the number of predictors p is large.**

**For instance, consider an example in which  $p = 100$  and  $H_0 : \beta_1 = \beta_2 = \dots = \beta_p = 0$  is true, so no variable is truly associated with the response. In this situation, about 5 % of the p-values associated with each variable (of the type shown in Table 3.4) will be below 0.05 by chance. In other words, we expect to see approximately five small p-values even in the absence of any true association between the predictors and the response. In fact, we are almost guaranteed that we will observe at least one p-value below 0.05 by chance! Hence, if we use the individual t-statistics and associated p-values in order to decide whether or not there is any association between the variables and the response, there is a very high chance that we will**

**incorrectly conclude that there is a relationship. However, the F-statistic does not suffer from this problem because it adjusts for the number of predictors. Hence, if  $H_0$  is true, there is only a 5 % chance that the F-statistic will result in a p-value below 0.05, regardless of the number of predictors or the number of observations.**

L'utilisation de la F-statistic pour tester n'importe quelle association entre les prédicteurs et la réponse fonctionne quand  $p$  est relativement petit, et relativement plus petit que  $n$ .

Cependant il peut arriver que l'on ait un très grand nombre de variables. Si  $p > n$  alors il y a plus de coefficients  $\beta_j$  à estimer que d'observation pour lesquels on les estime. Dans ce cas, nous ne pouvons pas ajuster le modèle de régression linéaire en utilisant les moindres carrés, donc la F-statistique ne peut être utilisée et aucune autre technique que l'on ait pu voir de même.

Deuxième question : Decider des variables importantes

La première étape dans une régression multiple est de calculer la F-statistique et d'examiner les p-valeurs associées. Si nous concluons qu'il y a au moins un prédicteur qui est relié à la réponse alors il est naturel de savoir lesquels ne sont pas bonnes! On pourrait regarder les p-valeurs individuelles mais comme nous avons vu, si  $p$  est grand il est probable de faire de mauvaises découvertes.

Il est possible que tous les prédicteurs soient associés à la réponse mais il est plus courant que la réponse est reliée à un sous ensemble de prédicteurs. La tâche de déterminer lesquelles sont intéressantes s'appelle la sélection des variables.

Idéalement on souhaiterait effectuer une sélection de variable en essayant pleins de modèles différents, contenant chacun un ensemble de prédicteur différents (exemple pour  $p=2$  on a 4 possibilités)

Comment déterminons nous lequel des modèles est le meilleur ? Plusieurs statistiques peuvent être utilisées pour juger de la qualité du modèle. Cela inclut Mallows's  $C_p$ , Akaike information criterion (AIC), Bayesian information criterion (BIC) et  $R^2$ . Ils seront développés dans le chapitre 6.

Malheureusement il y a  $2^p$  modèles qui contiennent un sous-ensemble de  $p$  variables. Cela signifie que même pour un  $p$  modéré, essayer tous les sous-ensembles possibles de prédicteurs est impossible.

Donc si  $p$  est petit on peut considérer tous les modèles. Mais si  $p$  est grand nous aurons besoin d'une approche automatisée et efficace pour choisir un plus petit ensemble de modèles à considérer.

Il existe 3 approches classiques pour cette tâche:

- Forward selection: On commence par le modèle nul (aucun prédicteurs). Puis nous ajustons  $p$  modèles de régression simple et nous ajoutons au modèle nul la variable qui nous donne le RSS le plus petit. Puis nous ajoutons à ce modèle la variable qui a le plus petit RSS pour le nouveau modèle à deux variables. On continue cette approche jusqu'à que nous soyons satisfait des résultats.
- Backward selection : On commence avec toutes les variables, et nous supprimons la variable avec la p-valeur la plus grande – qui est la variable qui est statistiquement la moins significative. Le nouveau modèle à  $p-1$  variable est construit et la variable avec la plus grande p-valeur est supprimée. On continue jusqu'à s'arrêter au niveau d'une règle. On peut avoir comme règle de stopper quand le reste des variables ont une p-valeur sous un certain seuil.
- mixed selection: C'est une combinaison de forward et backward selection. On commence avec aucune variable en faisant comme forward selection. On continue d'en ajouter une par une. Si à un moment la p-valeur pour une des variables dans le modèle atteint un certain seuil, on supprime cette variable du modèle. On continue d'exécuter la forward et backward étapes jusqu'à que toutes les variables du modèle ont suffisamment une petite p-valeur, et que toute variable extérieure au modèle ait une grande p-valeur si on l'ajoute au modèle.

Backward selection ne peut pas être utilisé si  $p > n$  alors que forward selection peut toujours être utilisé.

Troisième question : l'ajustement du modèle

Deux des mesures numériques d'ajustement de modèle les plus connues sont RSE et  $R^2$ . Ces quantités sont calculées et interprétées de la même manière que pour une régression linéaire simple.

Un  $R^2$  proche de 1 indique que le modèle explique une grande partie de la variance dans les variables de réponse. Il faut un compromis entre la RSE et  $R^2$  pour inclure les variables ou non dans le modèle.

On a  $RSE = \sqrt{\frac{1}{n-p-1} RSS}$  pour la démonstration revoir (ZT.35) C'est à (APK.45) que ça change

$$E\left[\sum_{i=1}^n (\beta_1 + \dots + \beta_p - \hat{\beta}_1 - \dots - \hat{\beta}_p)^2 (x_i - \bar{x})^2\right] = E[(\beta_1 - \hat{\beta}_1)^2 + \dots + (\beta_p - \hat{\beta}_p)^2] E\left[\sum_{i=1}^n (x_i - \bar{x})^2\right]$$

$$E\left[\sum_{i=1}^n (\beta_1 - \hat{\beta}_1)^2 (x_i - \bar{x})^2\right] = \text{var}(\hat{\beta}_1 + \dots + \hat{\beta}_p) E\left[\sum_{i=1}^n (x_i - \bar{x})^2\right] = \sigma^2 \frac{p}{\sum_{i=1}^n (x_i - \bar{x})^2} E\left[\sum_{i=1}^n (x_i - \bar{x})^2\right] = \sigma^2 p$$

(il y a une erreur au niveau de la fraction, il faut calculer avec des  $x_1, \dots, x_p$ . mais en gros c'est ça)

Donc avec un modèle de plusieurs variables peut avoir un RSE grand si la décroissance RSS est petit devant la croissance de  $p$ .

Quatrième question : les prédictions

Une fois les modèles de régression ajustés, il est simple d'appliquer  $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p$  dans le but de prédire la réponse  $Y$  sur la base des valeurs des prédicteurs  $X_1, \dots, X_p$ .

Cependant il y a 3 sortes d'incertitudes associées à la prédiction

1-/ les estimations des coefficients  $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$  sont des estimations de  $\beta_0, \beta_1, \dots, \beta_p$ .

$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \dots + \hat{\beta}_p X_p$  est juste une estimation de  $f(X) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$

L'imprécision des estimations des coefficients est dû à l'erreur redoutable. On peut calculer l'intervalle de confiance dans le but de déterminer combien  $\hat{Y}$  est proche de  $f(X)$ .

2-/ Il y a une source additionnelle potentiellement d'erreur redoutable que l'on appelle biais du modèle. Quand nous utilisons un modèle linéaire nous estimons la meilleure approximation linéaire de la vraie surface.

Toutefois ici nous ignorons la contradiction et opérons comme si le modèle linéaire était correct

3-/ Même si nous connaissons  $f(X)$  – qui est que nous connaissons les vraies valeurs  $\beta_0, \beta_1, \dots, \beta_p$  – la réponse ne peut être prédite parfaitement à cause de l'erreur aléatoire  $\varepsilon$  du modèle. Dans le chapitre 2 nous nous référons à l'erreur irréductible. Combien  $\hat{Y}$  varie de  $Y$ ? Nous utilisons les intervalles de prédictions pour répondre à la question. Les intervalles de prédictions sont toujours plus large que l'intervalle de confiance, car ils incorporent les deux erreurs d'estimation de  $f(X)$  à savoir l'erreur irréductible et l'erreur redoutable.

Intervalle de confiance exemple pour 95% des valeurs on a  $[\hat{\beta}_1 - 2 \cdot SE(\hat{\beta}_1), \hat{\beta}_1 + 2 \cdot SE(\hat{\beta}_1)]$

intervalle de prédiction plus large que l'intervalle de confiance

Les prédicteurs qualitatifs

Il peut arriver souvent que les prédicteurs soient des variables qualitatives

Les prédicteurs avec uniquement deux niveaux sont faciles à implémenter dans le modèle de régression. On crée simplement une variable dummy (ou indicateur) qui prend deux valeurs possibles.

On peut prendre  $x_i = 1$  si homme et 0 si femme

Ce qui nous donne  $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i = \beta_0 + \beta_1 + \varepsilon_i$  si  $x_i$  est 1 (donc un homme)

$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i = \beta_0 + \varepsilon_i$  si  $x_i$  est 0 (donc une femme)

Donc  $\beta_0$  est la moyenne des femmes et  $\beta_0 + \beta_1$  est la moyenne des hommes et  $\beta_1$  l'écart moyen entre les moyennes

Inverser les dummy variables (et donc 0 pour homme et 1 pour femme) affecte les variables et les moyennes des coefficients. Pour éviter le problème on peut utiliser des dummy variables de ce genre on prend -1 si femme ou 1 si homme .

Ce qui nous donne  $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i = \beta_0 + \beta_1 + \varepsilon_i$  si  $x_i$  est 1 (donc un homme)

$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i = \beta_0 - \beta_1 + \varepsilon_i$  si  $x_i$  est -1 (donc une femme)

La différence entre ces 3 techniques et que les coefficients  $\beta_0$  et  $\beta_1$  ont différentes significations.

$\beta_0$  peut être interprété comme la moyenne de tous les hommes et femmes  $\beta_1$  représente la part moyenne supplémentaire qu'on les femmes par rapport à  $\beta_0$ .

On peut aussi se retrouver avec des prédicteurs qualitatifs ayant plus de 3 niveaux. Quand on est à plus de trois niveaux, une unique dummy variable ne suffit plus. Dans ce cas on crée plus de dummy variables.

Pour l'exemple de Caucasiens, afro américain ou asiatique on utilise deux dummy variables

$$y_i = \beta_0 + \beta_1 x_{i,1} + \beta_2 x_{i,2} + \varepsilon_i$$

$x_{i,1} = 1$  si la  $i$ ème personne est asiatique sinon 0

$x_{i,2} = 1$  si la  $i$ ème personne est caucasienne sinon 0

Ce qui nous donne  $y_i = \beta_0 + \beta_1 + \varepsilon_i$  si la  $i$ ème personne est asiatique

$y_i = \beta_0 + \beta_2 + \varepsilon_i$  si la  $i$ ème personne est Caucasienne

$y_i = \beta_0 + \varepsilon_i$  si la  $i$ ème personne est afro américaine

Maintenant  $\beta_0$  peut être interprété comme la moyenne des afro américains,

$\beta_1$  comme la différence entre la moyenne entre asiatique et afro américain

$\beta_2$  comme la différence entre la moyenne entre Caucasien et afro américain

Conclusion : il y aura toujours un nombre  $n-1$  de dummy variable pour  $n$  catégories

## REVOIR LA F-STATISTIQUE (ce que ça signifie)

Extensions du modèle linéaire

Le modèle de régression linéaire standard permet l'interprétation des résultats et fonctionne assez bien à beaucoup de vrais problèmes. Toutefois, il suppose de très hautes hypothèses de restriction qui sont souvent violées en pratique.

Les deux suppositions les plus importantes sont la relation entre le prédicteur et la réponse (elle est additive et linéaire).

La supposition additive signifie que les effets de changement de  $X_j$  sur la réponse  $Y$  est indépendante des autres prédicteurs.

**L'hypothèse linéaire indique que la variation de la réponse  $Y$  due à une modification d'une unité de  $X_j$  est constante, quelle que soit la valeur de  $X_j$ .**

Supprimer la supposition additive

Supposons le modèle de régression linéaire suivant  $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon$

suivant ce modèle si nous augmentons  $X_1$  par une unité alors  $Y$  augmentera en moyenne d'une unité  $\beta_1$ . À noter que la présence de  $X_2$  n'altère pas l'état. Alors qu'en réalité il est possible que ça affecte  $X_2$ .

Un des moyens pour étendre ce modèle est de permettre un effet d'interaction en incluant un troisième prédicteur, appelé un terme d'interaction. Celui-ci est construit en calculant le produit de  $X_1$  et  $X_2$ .

Ce qui nous donne  $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 + \varepsilon$

$$Y = \beta_0 + (\beta_1 + \beta_3 X_2) X_1 + \beta_2 X_2 + \varepsilon$$

$$Y = \beta_0 + \tilde{\beta}_1 X_1 + \beta_2 X_2 + \varepsilon \quad \text{où} \quad \tilde{\beta}_1 = \beta_1 + \beta_3 X_2$$

Quand  $\tilde{\beta}_1$  change avec  $X_2$ , l'effet de  $X_1$  sur  $Y$  n'est plus constant. Ajuster  $X_2$  changera l'impact de  $X_1$  sur  $Y$ .

Le principe hiérarchique stipule que si nous incluons une interaction dans le modèle, nous devons aussi inclure les principaux effets et ce même si les  $p$ -valeurs associées avec leur coefficient ne sont pas significatives.

En d'autres mot si l'interaction entre  $X_1$  et  $x_2$  semble importante, nous devons inclure  $X_1$  et  $x_2$  dans le modèle.

Pour le cas des variables qualitatives, ou quantitative et qualitatives on peut procéder de la même manière. (Voir page 104 pour exemple en gros on a une dummy variable  $Beta_2$  et une  $Beta_3 * X_1$  et nulles si non étudiant)

#### Des relations non-linéaires

Comme nous avons parlé précédemment, le modèle de régression linéaire suppose une relation linéaire entre la réponse et les prédicteurs. Mais dans certains cas la relation entre la réponse et les prédicteurs peuvent être non linéaires. On peut passer à une régression polynomiale.

#### Des problèmes potentiels

Quand nous ajustons une régression linéaire à un jeu de données particulier, beaucoup de problèmes se produisent. Les plus communs sont:

- non linéarité de la relation réponse-prédicteur
- corrélation des termes d'erreur
- valeurs aberrantes
- fort levier de points
- collinéarité

En pratique, identifier et appréhender ces problèmes est un état de l'art.

##### 1- non linéarité du modèle

La régression linéaire suppose qu'il y a une stricte ligne reliant les prédicteurs et la réponse. Si la relation est très loin de linéaire, alors toutes les conclusions que nous avons établies sont mauvaises. La précision de prédiction du modèle peut être nettement réduite.

Les résidual plots sont des outils graphiques pour identifier la non-linéarité. Soit un simple modèle de régression linéaire, nous pouvons afficher les résidus  $e_i = y_i - \hat{y}_i$  contre les prédicteurs. Dans le cas d'une régression multiple on a alors plusieurs prédicteurs, au lieu de cela on affiche les résidus contre les valeurs prédites  $\hat{y}_i$ .

Si les résidual plots indiquent qu'il y a une association non linéaire dans les données, alors une approche simple est d'utiliser des transformations non-linéaires des prédicteurs comme  $\log X$ ,  $\sqrt{X}$  et  $X^2$ .

##### 2- la corrélation des termes d'erreur

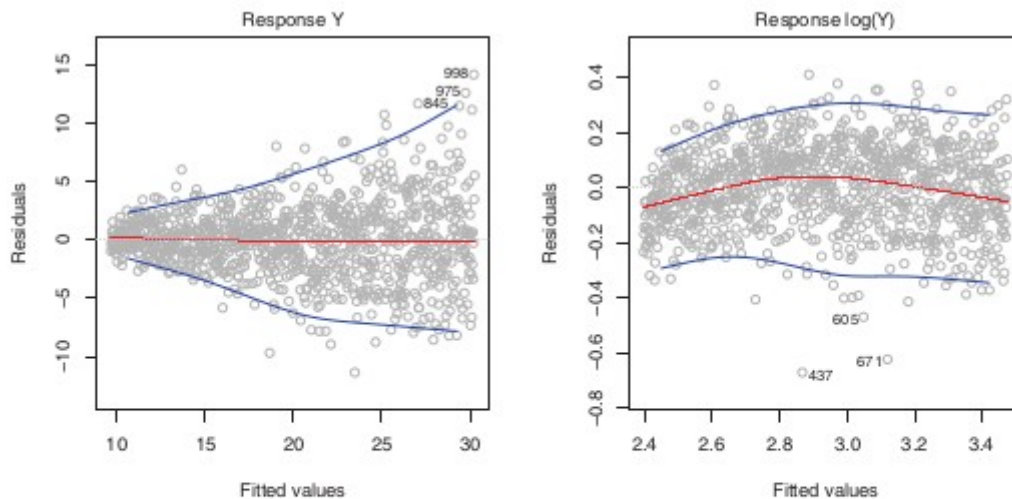
Une supposition importante du modèle de régression linéaire est que les termes d'erreur  $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$  ne sont pas corrélés. Qu'est-ce que ça signifie? Par exemple si l'erreur ne sont pas corrélés alors le fait que  $\varepsilon_i$  soit positif permet peu ou pas d'information sur le signe de  $\varepsilon_{i+1}$ . L'erreur standard qui est calculée pour calculer les coefficients de régression ou les valeurs ajustées sont basées sur la supposition de termes d'erreurs non corrélés. Si en fait il y a une corrélation parmi les termes d'erreurs, alors l'estimation des erreurs standard tendront à sous-estimer la vraie erreur standard. Les intervalles de confiance et de prédiction seraient donc trop petits qu'il ne devraient l'être: cela causerait des conclusions erronées qu'un paramètre soit significatif. En résumé, si les termes d'erreur sont corrélés, on aura un sentiment de confiance injustifié dans notre modèle.



3-/la variance du terme d'erreur non-constante

Une supposition importante de la regression linéaire est que les termes d'erreur ont une variance constante,  $\text{Var}(\varepsilon_i) = \sigma^2$

Les erreurs standards, les intervalles de confiance, et les hypothèses de tests associées a un modèle linéaire dependent de cette supposition. Malheuresmeent, il est frequent que la varaince du terme d'erreur soit non constant. Par exmple la variance du terme d'erreur pourrait augmenter avec la valeur de la reponse. On peut identifier des variance non constantes des les erreurs, ou heteroscedacité par la présence d'un entonnoir dans les residuals plots. Par exemple voir figure 3.11 en utilisant  $\log(Y)$



**FIGURE 3.11.** Residual plots. In each plot, the red line is a smooth fit to the residuals, intended to make it easier to identify a trend. The blue lines track the outer quantiles of the residuals, and emphasize patterns. Left: The funnel shape indicates heteroscedasticity. Right: The response has been log transformed, and there is now no evidence of heteroscedasticity.

avec

$\log(Y)$  les residus apparaissent avoir une variance constante, qui montre une relation non linéaire dans les données

Il arrive que l'on ait une bonne idée de la variance de chaque reponse. Parfois, nous avons une bonne idée de la variance de chaque réponse. Par exemple, la  $i$ ème réponse pourrait être une moyenne de  $n_i$  observations brutes. Si chacune de ces observations brutes n'est pas corrélée avec la variance  $\sigma^2$ , alors leur moyenne a la variance  $\sigma_i^2 = \sigma^2 / n_i$ . Dans ce cas, une solution simple consiste à ajuster notre modèle par la méthode des moindres carrés pondérés, avec des pondérations proportionnelles à la valeur inverse des variances - c'est-à-dire.  $w_i = n_i$  dans ce cas. La plupart des logiciels de régression linéaire permettent les poids d'observation.

#### 4-/ Valeurs aberrantes

Un outliers est un point pour lequel  $y_i$  est loin de la valeur prédite par le modèle. Elle peuvent apparaitre pour plusieurs raisons, comme un enregistrement incorrect durant l'observation de la collecte des donnees. Les residual plits peuvent etre utilisés pour identifier les outliers. En pratique il peut s'averer difficile de decider de combien les residus ont besoin pour etre considéré comme un outlier. Pour resoudre se probelme, au lieu d'afficher les residus, on peut afficher les residus studentisés, caluclé a partir de la division de chaque résidu par son erreur standard d'estimation. Les observation pour lesquels les residus sont plus grand que 3 en valeur absolue sont des possibles valeur aberante.

Si nous croyons qu'un outlier s'est produit dû a une erreur dans l'enregistrement des données, alors la solution est de simplement supprimer l'observation. Il faut quand meme faire attention que ce ne soit pas le modèle qui est mal choisit comme un predicteur manquant.

#### 5-/ High leverage points

On vient de voir que les outliers sont des observations pour lesquelles la réponse  $y_i$  est inhabituelle suivant un prédicteur  $x_i$ . En revanche, des observations avec high leverage ont des valeurs inhabituelles pour  $x_i$ . Des points de levier permettent de faire des régressions plus précises. Elles ne sont pas aberrantes et doivent être prises en compte pour augmenter la qualité de la régression. C'est pour cela qu'il est important d'identifier les observations de point de levier.

SUITE In a simple linear regression

Suite page 112

εβμσΠ

$\hat{f}$

ascertain : à vérifier

assessing : évaluer, évaluation

assumption : supposition

despite : malgré

discrepancy : contradiction

for instance : par exemple

however: toutefois

in contrast : en revanche

lack of: manquer de

outliers: valeurs aberrantes

overall : global

perform : permettre / exécuter

provided that : à condition que

rather : plutôt

scatter : distribution

scatterplot : nuage de points

scope: portée

seek : rechercher , tenter de

since : puisque

slope : pente

spread out : s'étendre

straightforward : simple, honnête, juste

throughout : tout au long

whether : qu'il s'agisse, si

wider : plus large