

# Sélection de variables

## Analyse multidimensionnelle appliquée

Léo Belzile

HEC Montréal

automne 2022

**Objectif:** bâtir un modèle pour une variable réponse  $Y$  en fonction de variables explicatives  $X_1, \dots, X_p$ .

On s'intéresse à

$$f(X_1, \dots, X_p) \quad .$$

vraie moyenne inconnue

L'analyste détermine

$$\hat{f}(X_1, \dots, X_p),$$

approximation

une fonction des variables explicatives.

On spécifie que la **moyenne** de la variable réponse  $Y$  est une fonction linéaire des variables explicatives  $X_1, \dots, X_p$ , soit

$$\underset{\text{moyenne théorique}}{E(Y \mid \mathbf{x})} = \underset{\text{somme pondérée des variables explicatives}}{\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}} .$$

en supposant que l'écart entre les observations et cette moyenne est constant,

$$\text{Va}(Y \mid \mathbf{x}) = \sigma^2.$$

Pour la  $i$ ème observation,

$$\underset{\text{réponse}}{Y_i} = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} + \underset{\text{aléa}}{\varepsilon_i}.$$

prédicteur linéaire

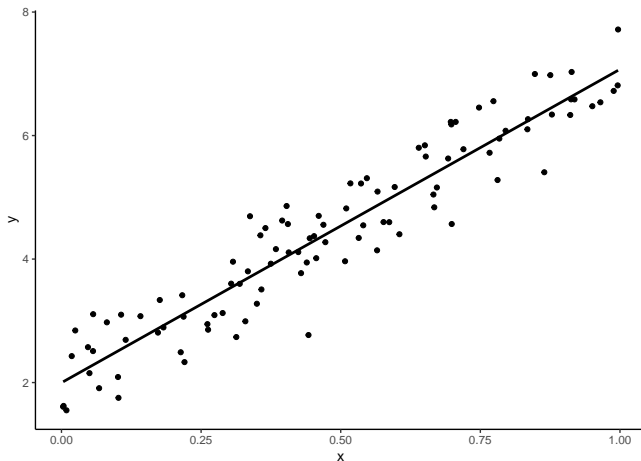
- L'aléa  $\varepsilon_i$  représente la distance **verticale** entre la vraie pente et l'observation
- Autant d'aléas que d'observations ( $n$ ), variable aléatoire inconnue...

- L'aléa  $\varepsilon_i$  représente l'erreur, soit la différence entre la valeur observée et la moyenne de la population pour les même valeurs des variables explicatives.
- On suppose que le modèle pour la moyenne est correctement spécifié: l'aléa a une moyenne théorique nulle,  $E(\varepsilon_i) = 0$ .
- On suppose que les observations sont indépendantes les unes des autres.

# Régression linéaire en deux dimensions

Si  $E(Y) = \beta_0 + \beta_1 X$ , alors

- $\beta_0$  représente l'ordonnée à l'origine (valeur quand  $X = 0$ .)
- $\beta_1$  est la pente



L'estimation des paramètres  $\hat{\beta}_0, \dots, \hat{\beta}_p$  nous donne

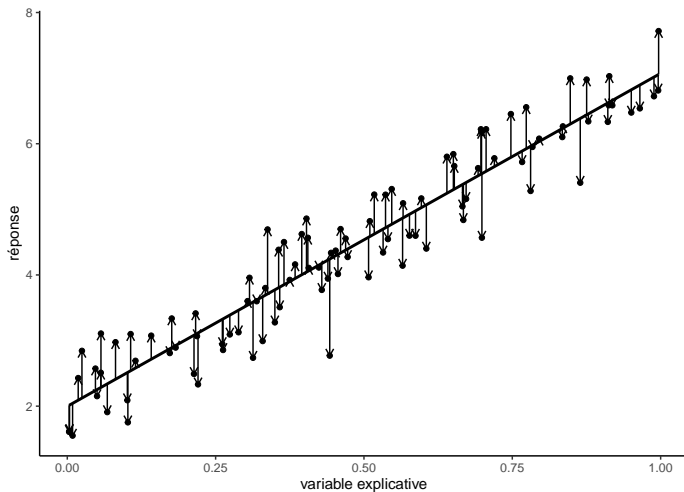
$$\underset{\text{prédiction}}{\widehat{Y}_i} = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} \dots + \hat{\beta}_p x_{ip}.$$

On peut approximer l'aléa à l'aide du **résidu ordinaires**, soit

$$\underset{\text{résidu ordinaire}}{e_i} = \underset{\text{observation}}{Y_i} - \underset{\text{prédiction}}{\widehat{Y}_i}.$$

- par construction, la moyenne des  $e_i$  est zéro.
- le résidu ordinaire est la distance verticale entre l'observation et la "droite" **ajustée**

# Illustration des résidus ordinaires





L'erreur quadratique moyenne théorique est

$$\mathbb{E} \left[ \left\{ (Y - \hat{f}(x_1, \dots, x_p)) \right\}^2 \right],$$

la moyenne de la différence au carré entre la vraie valeur de  $Y$  et la valeur prédite par le modèle.

En pratique, on remplace la moyenne théorique par une moyenne empirique obtenue à partir d'un échantillon aléatoire.

Comment estimer les paramètres  $\beta_0, \dots, \beta_p$ ?

**Optimisation:** trouver les valeurs qui minimisent l'erreur quadratique moyenne **empirique** avec l'échantillon des  $n$  observations, soit

$$\frac{e_1^2 + \dots + e_n^2}{n}$$

Solution explicite  $\hat{\beta} = (\mathbf{x}^\top \mathbf{x})^{-1} \mathbf{x}^\top \mathbf{Y}$ !

La fonction `lm` calcule l'ajustement du modèle linéaire.

Arguments:

- `formula`: formule de type `reponse ~ variables explicatives`, où les variables explicatives sont séparées par un signe `+`
- `data`: base de données

```
1 modlin <- lm(mpg ~ hp + wt,  
2             data = mtcars)  
3 summary(modlin)
```

# Sortie de summary

Call:

```
lm(formula = mpg ~ hp + wt, data = mtcars)
```

Residuals:

Min	1Q	Median	3Q	Max
-3.941	-1.600	-0.182	1.050	5.854

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	37.22727	1.59879	23.285	< 2e-16 ***
hp	-0.03177	0.00903	-3.519	0.00145 **
wt	-3.87783	0.63273	-6.129	1.12e-06 ***

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2.593 on 29 degrees of freedom

Multiple R-squared: 0.8268, Adjusted R-squared: 0.8148

F-statistic: 69.21 on 2 and 29 DF, p-value: 9.109e-12

- Formule de l'appel
- Statistiques descriptives des résidus ordinaires  $e_1, \dots, e_n$ .
- Tableau des estimations
  - Coefficients  $\hat{\beta}_j$
  - Erreurs-types,  $\text{se}(\hat{\beta}_j)$
  - Statistique du test- $t$  pour  $\mathcal{H}_0 : \beta_j = 0$ , soit  $t = \hat{\beta}_j / \text{se}(\hat{\beta}_j)$
  - Valeur- $p$  selon loi nulle  $\text{St}(n - p - 1)$
- Estimation de l'écart-type  $\hat{\sigma}$  et degrés de liberté  $n - p - 1$
- Estimations du coefficient de détermination,  $R^2$  et  $R^2$  ajusté
- Statistique  $F$  d'ajustement global et valeur- $p$  de  $F(p, n - p - 1)$ 
  - $\mathcal{H}_a$ : modèle linéaire
  - $\mathcal{H}_0$ : modèle avec uniquement ordonnée à l'origine (chaque observation prédite par la moyenne des réponses,  $\bar{Y}$ )

- `resid` pour les résidus ordinaires  $e_i$
- `fitted` pour les valeurs ajustées  $\widehat{Y}_i$
- `coef` pour les estimations des paramètres  $\hat{\beta}_0, \dots, \hat{\beta}_p$
- `plot` pour des diagnostics graphiques d'ajustement
- `anova` pour la comparaison de modèles emboîtés
- `predict` pour les prédictions (avec nouvelles données)
- `confint` pour intervalles de confiance pour les paramètres.

- Les facteurs (<factor>) sont traités adéquatement par **R**.
- Si la variable a  $K$  valeurs possibles (niveaux), le modèle inclut  $K - 1$  indicatrices 0/1.
- Par défaut dans **\$**, la catégorie de référence est la plus petite en ordre alphanumérique.

Considérons une variable catégorielle `cat` avec niveaux 1, 2, et 3.

cat	cat2	cat3
1	0	0
2	1	0
3	0	1

La catégorie de référence est associée à l'ordonnée à l'origine (quand `cat2=0` et `cat3=0`).



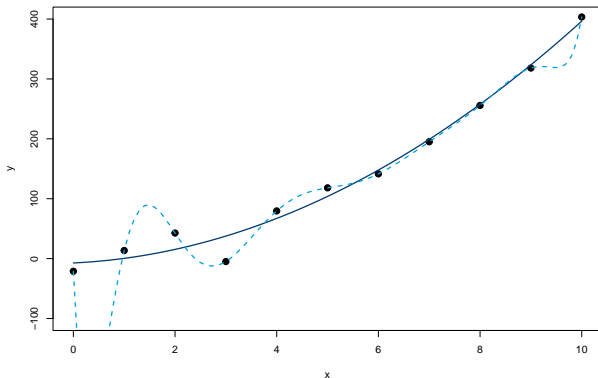


- Comment choisir quelles variables inclure?
- Quel est la spécification adéquate pour  $f(x_1, \dots, x_p)$ ?
  - régression, réseaux de neurone, forêts aléatoires, etc.
  - transformations de variables,  $\text{age}^2$ ,  $\ln(\text{age})$ , etc.

Notre but sera de sélectionner un **bon** modèle, selon les objectifs de l'étude



# Illustration du surajustement



## Séparation des données

Ne pas utiliser les données employés pour ajuster un modèle pour **prédire la performance**

échantillons d'apprentissage/validation/test (fixes)

validation croisée (avec  $K = 5, 10$  groupes), mais *résultat aléatoire*

